

Monetary Policy Operations: Theory, Evidence, and Tools for Quantitative Analysis*

Ricardo Lagos
New York University

Gastón Navarro
Federal Reserve Board

November 20, 2023

Abstract

We formulate a quantitative dynamic equilibrium theory of trade in the fed funds market, calibrate it to fit a comprehensive set of marketwide and micro-level cross-sectional observations, and use it to make two contributions to the operational side of monetary policy implementation. First, we produce global structural estimates of the aggregate demand for reserves—a crucial decision-making input for modern central banks. Second, we propose diagnostic tools to gauge the central bank’s ability to track a given fed funds target, and the heterogeneous incidence of policy actions on the shadow cost of funding across banks.

Keywords: monetary policy, fed funds, central-bank reserves, OTC markets
JEL classification: G1, C78, D83, E44

*We thank Joshua Herman, Patrick Molligo, Maddie Penn, Charlotte Singer, and Kenji Wada for superb research assistance. We also thank Michele Cavallo, Jeff Huther, and Cindy Vojtech for useful comments and discussions. We are grateful to Heather Ford and Gina Sellito for helping us navigate the compliance requirements to access some of the data. The views expressed in the paper are those of the authors and are not necessarily reflective of views at the Federal Reserve Board or the Federal Reserve System.

1 Introduction

The Federal Reserve uses the fed funds rate to communicate and implement its monetary-policy stance. In each of the eight regularly scheduled meetings during the year, the Federal Open Market Committee (FOMC) chooses a fed funds rate target and issues implementation notes specifying the policy instruments that will be used to create market conditions for fed funds to trade at rates near the target. Two kinds of instruments are typically used to achieve this goal. The first, open-market operations, affect the market price of fed funds by changing the quantity of reserves. The second, a set of administered rates offered by standing facilities, such as the Discount-Window rate (DWR), the interest rate paid on reserves (IOR), and the offering rate on overnight reverse repurchase agreements (ONRRP), affect the market price of fed funds by changing banks' return from holding (or borrowing) reserves at the central bank. An *operating framework* for implementing monetary policy is a consistent usage of these instruments to implement the fed funds rate target. For example, an operating framework may rely mostly on managing *quantities* of reserves, and another on managing administered *rates*.

Figure 1 shows the stylized theoretical demand-and-supply model that policymakers use to think about how different operating frameworks can achieve a fed funds rate target.¹ Before the Great Financial Crisis of 2007-2008 (GFC), aggregate reserves were scarce, e.g., around a relatively low level such as Q_0 in the first panel of Figure 1. In this context, target rates like r_0^* were achieved by changing the quantity of reserves (represented by the vertical line in the figure) through open-market operations. This operating framework is known as a *corridor system* since it can implement any target rate inside the *corridor* defined by a *ceiling rate*, typically the Discount-Window rate, ι_w (possibly plus a stigma premium or other costs associated with borrowing from the central bank), and a *floor rate*, such as the interest that the central bank pays on bank reserves, ι_r (or a lower rate if not all fed funds participants can earn interest on reserves held at the central bank).²

¹An individual bank's demand for reserves is typically thought of in terms of Poole (1968), who derives it from the static decision problem of an individual bank that chooses how much of its beginning-of-day reserves to lend (to earn a given market interest), and how much to hold (to insure against an otherwise uninsurable exogenous reduction in reserves, which would force the bank to engage in end-of-day borrowing at a central-bank discount rate higher than the market rate). This "Poole model" is the go-to framework in policy circles, see, e.g., Ennis and Keister (2008), Keister et al. (2008), Keister (2012), Afonso et al. (2020b), and Åberg et al. (2021).

²In a corridor system the target rate and the administered rates are typically chosen so that the target rate lies in the middle of the corridor, and the central bank does not use the standing facilities as instruments to manage the fed funds rate (although they are available for fed funds participants who wish to borrow from, or lend reserves to the central bank).

During, and in the aftermath of the GFC, the Federal Reserve undertook a series of large-scale asset purchase programs that increased banks’ reserve balances to a very high level such as Q_1 in the top-right panel of Figure 1.³ With a such a large supply of reserves, the Fed can no longer rely on routine open-market operations (which entail relatively small changes in the quantity of reserves) to instrument changes in the fed funds rate target. In this context, target rates like r_1^* are achieved by changing the administered rates, i.e., the DWR, the IOR, and the ONRRP. This operating framework is known as a *floor system*.⁴

To describe the levels of reserves compatible with these operating frameworks, policymakers often use “scarce reserves” to refer to the range for which the slope of the aggregate demand for reserves is “steep”, “ample reserves” for the range for which it is “gentle”, and “abundant reserves” for the range for which it is “flat”, as illustrated in the top-right panel of Figure 1.⁵ The Federal Reserve intends to continue operating a floor system in which an “ample” supply of reserves ensures that the fed funds rate is controlled by the administered rates, and in which “active management of the supply of reserves is not required.”⁶ In terms of the schematic in Figure 1, this operating framework seems easy to manage: the Fed just needs to ensure the supply of reserves remains “ample”, i.e., at a level close to Q_1 in the top-right panel.

In practice, however, it is difficult to quantify how large the supply of reserves needs to be for it to be “ample”. A supply of reserves larger than Q_1 in the top-right panel of Figure 1 would still allow the Fed to operate a floor system, but would entail unnecessary costs.⁷ Conversely, a supply of reserves smaller than Q' would imply operating on the steep part of the aggregate demand for reserves, i.e., *de facto* abandoning the preferred operating framework for a corridor system that requires active day-to-day management of the supply of reserves to achieve the fed

³Total reserves were about \$40 bn before the GFC through mid 2008, and reached \$2.8 tn in 2014 (see Figure 16).

⁴In a floor system, the central bank actively operates two *standing facilities*: a lending facility that lends reserves to banks, and a deposit facility that enables qualifying institutions to lend reserves to the central bank. In the United States, for example, the Fed now sets three administered rates: DWR, IOR, and ONRRP. The IOR, which is regarded as the primary policy instrument, is the rate banks can earn by holding reserves in the deposit facility. The ONRRP, which is typically set lower than the IOR, is the rate that a broader set of financial institutions (including banks but also GSEs and money-market funds) can earn by holding reserves in the deposit facility. The logic is that the IOR acts a reservation price for lending banks, and the ONRRP acts as a reservation price for other (non-bank) lending institutions, so depending on the composition of trades, one of the two rates should act as a floor for negotiated rates.

⁵See Afonso et al. (2020b) and Afonso et al. (2022). The term “ample” has become standard language in FOMC press releases (see, e.g., Federal Reserve Board (2019c)).

⁶See, e.g., Federal Reserve Board (2019b).

⁷See Bernanke and Kohn (2016) and Ireland (2018) for discussions of the economic and political risks associated with an operating framework that involves paying interest on a large stock of reserves.

funds rate target. Notice that the operational difficulty goes beyond resolving the arbitrariness involved in specifying numerical thresholds for the slope of the reserve demand to be considered “steep”, “gentle”, or “flat”. Even if we agreed on a precise definition of “gentle slope”, the main difficulty is to find the associated quantity of reserves, which requires reliable estimates of the slope of the aggregate demand for reserves for a wide range of the aggregate quantity of reserves. In other words, running a floor system like the one the Federal Reserve has adopted, requires *global* estimates of the slope of the aggregate demand for reserves. This presents a significant challenge because existing state-of-the-art empirical methods only deliver *local* estimates of the slope of the demand for reserves, i.e., estimates obtained based on instrumented variation around a relatively narrow range of the aggregate supply of reserves.⁸

In this paper we develop a quantitative model of the fed funds market, calibrate it to match key micro and macro statistics that characterize fed funds trading in the United States—including available empirical estimates of the *local* slope of the aggregate demand for reserves—and use it to bridge the *local-global gap*. Specifically, we use the relationship between the aggregate supply of reserves and the equilibrium fed funds rate implied by the theory, to infer the *global* shape of the actual aggregate demand for reserves.

The theory incorporates search and bilateral bargaining to represent the well-documented over-the-counter microstructure of the fed funds market. The theory also accounts for relevant institutional considerations, such as the differential regulatory treatment of the reserve balances held by Government Sponsored Enterprises (GSEs) vis á vis depository institutions, and incorporates the array of policy instruments and regulations that affect participants’ demands for reserves, such as the administered policy rates (DWR, IOR, ONRRP), the regulatory requirements on reserve holdings, and the aggregate quantity of reserves supplied to the system. The theory also accommodates the large degree of heterogeneity among fed funds participants across several dimensions, such as: market power in bilateral negotiations, frequency and size distribution of idiosyncratic payment shocks originated by forces outside the fed funds trading motives, measures of trading activity (frequency of trade, number of counterparties, participation rate in aggregate volume of trade), and degree of centrality in the endogenous market-making activity that reallocates reserves across the trading network.

⁸The empirical challenge is illustrated in the bottom panels of Figure 1, which show situations in which structural parameters are Π_i at the time the quantity-price pair (Q_i, r_i^*) is observed, for $i \in \{0, 1\}$. Without theoretical guidance to identify the structural parameters whose variation, e.g., from Π_0 to Π_1 , shift the demand for reserves, one may be led to believe that the observations $\{(Q_i, r_i^*)\}_{i \in \{0, 1\}}$ lie on a single demand curve, and therefore overestimate (bottom-left panel) or underestimate (bottom-right panel) the relevant slope.

We calibrate the parameters of the theory that govern the heterogeneity in payment and trading activities using high-frequency micro-level transaction data from Fedwire, and find that the model is able to fit the targeted observations well, e.g., as in the data, a small number of very active banks account for the majority of loans, and carry out most of the intermediation. The calibration strategy also ensures that—at the current level of total reserves—the magnitude of the variation in the equilibrium fed funds rate induced by exogenous variation in the supply of reserves is in line with standard reduced-form empirical estimates of the “liquidity effect”. The calibrated model is also broadly consistent with empirical observations not targeted in the calibration, such as the cross-sectional distribution of bilateral interest rates, the distribution of bid-ask spreads, and the intraday flow of reserves and supporting interest rates between pairs of banks in different positions on the trading network.

We use the quantitative theory to make two practical contributions to the operational side of monetary policy implementation. First, we use the calibrated model—disciplined and validated by micro data—to deliver global structural estimates of the aggregate demand for reserves, which should be useful to central banks that wish to operate floor systems. Second, we use the calibrated model as the basis for two “navigational instruments” for monetary policy implementation. The first, which we term *Monetary Confidence Band* (MCB), is a hybrid of theory and data: it is a simple procedure that uses the empirical distribution of daily reserve-draining shocks to construct a confidence band around the aggregate demand for reserves that, for each outstanding quantity of reserves, will contain the equilibrium fed funds rate with a desired degree of confidence, e.g., 95%. The second is the cross-sectional distribution of banks’ shadow cost of procuring funding in the fed funds market implied by the theory.

This paper contributes to the large empirical and theoretical literature that studies the fed funds market, e.g., Poole (1968), Hamilton (1996), Furfine (1999), Carpenter and Demiralp (2006), Ashcraft and Duffie (2007), Bech and Atalay (2010), Afonso et al. (2011), Bech and Klee (2011), Afonso and Lagos (2012, 2015a,b), Ennis and Weinberg (2013), Armenter and Lester (2017), Afonso et al. (2019), Ennis (2019), Chiu et al. (2020), Beltran et al. (2021), Copeland et al. (2021), and Afonso et al. (2022). Methodologically, we build on the strand of the finance microstructure literature that uses search theory to model over-the-counter markets, e.g., Duffie et al. (2005), Lagos and Rocheteau (2007, 2009), Weill (2007), Lagos et al. (2011), Üslü (2019), and Hugonnier et al. (2020). Specifically, our model builds on Afonso and Lagos (2015b), which we generalize along several dimensions to make it a serviceable quantitative tool for monetary

policy implementation.

The rest of the paper is organized as follows. Section 2 presents the model, discusses the main assumptions, and defines equilibrium. Section 3 documents the key statistics that will guide the quantitative implementation of the theory, e.g., bank-level measures of fed funds trading activity (participation and intermediation), frequency and size distribution of micro-level intraday payments between banks, and the typical beginning-of-day cross-sectional distribution of reserve balances (net of predictable payments and regulatory requirements). Section 3 also reports empirical estimates of the distribution of aggregate daily reserve-draining shocks for the fed funds market since the GFC, and of the “liquidity effect” associated with exogenous variation in the aggregate supply of reserves. Section 4 discusses the calibration strategy. Section 5 reports how well the quantitative model fits empirical price and quantity observations not targeted by the calibration. Section 6 analyzes the aggregate demand for reserves generated by the model, and uses it as the basis for a quantitative-theoretic estimation of the aggregate demand for reserves in the United States. Section 7 proposes two “navigational instruments” to guide the routine monetary policy operations necessary to implement a fed-funds rate target. Section 8 uses the quantitative theory to rationalize the well-known fed-funds rate spikes of September 2019. Section 9 concludes. The appendices contain supplementary material.

2 Theory

There is a unit measure of *banks* that are heterogeneous along several dimensions. We represent this heterogeneity with a finite set \mathbb{N} of bank types, and let $n_i \in [0, 1]$ represent the proportion of banks of type $i \in \mathbb{N}$, with $\sum_{i \in \mathbb{N}} n_i = 1$. Banks hold an asset we interpret as (claims to) *reserve balances* that can be traded with other banks during the time interval $\mathbb{T} = [0, T]$. The reserve balance that a bank holds at a given time is represented by a real number, e.g., $a \in \mathbb{R}$. The cumulative distribution function of reserve balances across all banks at time $t \in \mathbb{T}$ is denoted $F_t(a) = \sum_{i \in \mathbb{N}} n_i F_t^i(a)$, where $F_t^i(a) : \mathbb{R} \times \mathbb{T} \rightarrow [0, 1]$ is the cumulative distribution of balances across banks of type i at time t . The initial distribution, $\{F_0^i(\cdot)\}_{i \in \mathbb{N}}$, is given, and so is the aggregate supply of reserve balances throughout the trading session, denoted $Q \equiv \int a dF_0(a)$.

Banks trade reserves with other banks in a bilateral over-the-counter market where a bank of type $i \in \mathbb{N}$ contacts another bank at random times generated by a Poisson process with arrival rate $\beta_i \in \mathbb{R}_+$. Conditional on a meeting, the counterparty is a random (uniform) draw from the population of banks. Once two banks have made contact, they bargain over the size

of the loan and the quantity of reserve balances to be repaid by the borrower. The bargaining outcome is determined by Nash bargaining. When a bank of type $i \in \mathbb{N}$ negotiates with a bank of type $j \in \mathbb{N}$, we assume the bargaining power of the former is $\theta_{ij} = 1 - \theta_{ji} \in [0, 1]$. After the transaction, the banks part ways.

We assume all loans are settled at time $\bar{T} > T$, and that banks value reserve balances linearly at that time. Specifically, if $c \in \mathbb{R}$ is a bank's net credit position to be settled at \bar{T} that has resulted from a certain history of trades, then $e^{-r(\bar{T}-t)}c$ is the bank's payoff from this credit balance at time $t \in [0, T]$, where $r \in \mathbb{R}_+$ is the discount rate common to all banks.

Banks receive payment shocks that cause reallocations of reserve balances among pairs of banks. Specifically, with Poisson rate $\lambda_i \in \mathbb{R}_+$, a bank of type $i \in \mathbb{N}$ is forced to make an immediate transfer of reserves to a counterparty that is drawn randomly (uniformly) from the population of banks. This process for the arrival of payment shocks is independent across banks and independent of the processes that generate bilateral trading opportunities. Conditional on the arrival of a payment shock, the quantity of reserves that the bank of type i sends the bank of type j is modeled as a random variable with cumulative distribution function $G_{ij} : \mathbb{Z} \rightarrow [0, 1]$, where $\mathbb{Z} \subseteq \mathbb{R}$ is the support of G_{ij} , and $dG_{ij}(z) = dG_{ji}(-z)$, which captures the notion that these payments are *transfers* between pairs of bank types.

For each $i \in \mathbb{N}$, define the function $U_i : \mathbb{R} \rightarrow \mathbb{R}$, where $U_i(a)$ represents the payoff to a bank of type i from holding reserve balance $a \in \mathbb{R}$ at the end of the trading session. Similarly, for each $i \in \mathbb{N}$, define the function $u_i : \mathbb{R} \rightarrow \mathbb{R}$, where $u_i(a)$ represents the flow payoff to a bank of type i from holding reserve balance $a \in \mathbb{R}$ during the trading session. The type of a bank is defined by a collection of primitives, i.e., type $i \in \mathbb{N}$ is defined by $(n_i, \beta_i, \lambda_i, \{\theta_{ij}, G_{ij}\}_{j \in \mathbb{N}}, u_i, U_i)$, in the sense that each of the n_i banks of type i has trading frequency β_i , bargaining powers $\{\theta_{ij}\}_{j \in \mathbb{N}}$, payment frequency λ_i , probability distributions $\{G_{ij}\}_{j \in \mathbb{N}}$ of payment sizes, intraday payoff function u_i , and end-of-day payoff function U_i .

2.1 Discussion

The market for federal funds is a market for unsecured loans of reserve balances at the Federal Reserve Banks. These unsecured loans, commonly referred to as *federal funds* (or *fed funds*) are delivered on the same day, and their maturity is typically overnight. Most fed funds transactions and interbank payments are conducted through *Fedwire Funds Services* (Fedwire), a large-value real-time gross settlement system operated by the Federal Reserve Banks. Participants in the

fed funds market are institutions that hold reserve balances in accounts at the Federal Reserve, which include commercial banks, savings banks, thrift institutions, credit unions, agencies and branches of foreign banks in the United States, government securities dealers, government agencies such as federal or state governments, and Government Sponsored Enterprises (GSEs, e.g., Freddie Mac, Fannie Mae, and Federal Home Loan Banks). The fed funds market is *over the counter*: in order to trade, a financial institution must first find a willing counterparty and then bilaterally negotiate the size and rate of the loan. The fed funds market is the epicenter of monetary policy implementation in the sense that the *effective fed funds rate* (EFFR)—the policy rate that the Federal Reserve uses to communicate and implement monetary policy—is a daily volume-weighted average of the bilateral interest rates negotiated by fed funds participants.

We use a search-based model with ex post bargaining to represent the bilateral over-the-counter nature of the fed funds market. Search captures three layers of randomness in trading activity in our model. First, the time it takes a bank of type $i \in \mathbb{N}$ to contact a counterparty is an exponentially distributed random variable with mean $1/\beta_i$. Second, conditional on having contacted a counterparty, the type of the counterparty is a uniform random draw. Third, conditional on having met a counterparty of type $j \in \mathbb{N}$ at time t , the current reserve balance of the counterparty is a random variable with cumulative distribution function $\{F_t^j(\cdot)\}_{j \in \mathbb{N}}$. We use the generalized Nash bargaining solution to represent the outcome of the bilateral negotiations between counterparties in actual fed funds trades.

The motives for trading fed funds may vary across participants and their specific circumstances on any given day, but there are two main reasons in general. First, some participants may regard fed funds as an investment vehicle—an interest-yielding asset that can be used to deposit balances overnight. Also, some institutions such as commercial banks use the fed funds market to offset the effects of random payment shocks (resulting from transactions initiated by their clients or by profit centers within the institutions themselves) that would otherwise leave them with a reserve position deemed too low relative to regulatory requirements. In the theory, the Poisson rate λ_i represents the frequency of these payment shocks for a bank of type $i \in \mathbb{N}$, and G_{ij} represents the size distribution of payment shocks between two banks of types i and j . In our model, all payoff-relevant policy and regulatory considerations are captured by the intraday and end-of-day payoff functions, $\{u_i(\cdot), U_i(\cdot)\}_{i \in \mathbb{N}}$.

Fedwire and the fed funds market operate 21.5 hours each business day, from 9:00 pm eastern standard time (EST) on the preceding calendar day to 6:30 pm EST. Although there is

occasionally some activity between 9:00 pm and 9:00 am, the bulk of the fed funds transactions and interbank payments take place between 9:00 am and 6:30 pm. Thus, in the theory, we think of $t = 0$ as standing in for 9:00 am and use the initial condition $\{F_0^i(\cdot)\}_{i \in \mathbb{N}}$ to represent the distribution of reserve balances at this time.

2.2 Equilibrium

Let $J_t^i(a, c) : \mathbb{N} \times \mathbb{T} \times \mathbb{R}^2 \rightarrow \mathbb{R}$ be the maximum attainable payoff to a bank of type i that at time $t \in \mathbb{T}$ has reserve balance $a \in \mathbb{R}$ and net credit position $c \in \mathbb{R}$. In Appendix A (Lemma 1) we show that $J_t^i(a, c) = V_t^i(a) + e^{-r(\bar{T}-t)}c$, where $V_t^i(a) : \mathbb{N} \times \mathbb{T} \times \mathbb{R} \rightarrow \mathbb{R}$ is the maximum expected discounted payoff a bank of type $i \in \mathbb{N}$ can obtain when holding $a \in \mathbb{R}$ reserve balances at time $t \in \mathbb{T}$. Whenever two banks contact each other during the trading session, they bargain over the size of the loan and the repayment. Consider a bank of type i with reserve balance a that contacts a bank of type j with reserve balance \tilde{a} . The pair $(b_t^{ij}(a, \tilde{a}), R_t^{ji}(\tilde{a}, a))$ denotes the bilateral terms of trade negotiated by these banks at time t , where $b_t^{ij}(a, \tilde{a})$ is the quantity of reserves that the bank of type i with balance a lends to the bank of type j with balance \tilde{a} , and $R_t^{ji}(\tilde{a}, a)$ is the quantity of balances that the latter commits to repay the former at time \bar{T} . For any $(a, \tilde{a}, t) \in \mathbb{R}^2 \times \mathbb{T}$, $(b_t^{ij}(a, \tilde{a}), R_t^{ji}(\tilde{a}, a))$ is the solution to

$$\max_{(b, R) \in \bar{\mathbb{R}} \times \mathbb{R}} \left[V_t^i(a - b) + e^{-r(\bar{T}-t)}R - V_t^i(a) \right]^{\theta_{ij}} \left[V_t^j(\tilde{a} + b) - e^{-r(\bar{T}-t)}R - V_t^j(\tilde{a}) \right]^{\theta_{ji}},$$

with $\bar{\mathbb{R}} \equiv [-\bar{b}, \bar{b}]$, where $\bar{b} \in \mathbb{R}_+ \cup \{\infty\}$ represents a limit on bilateral credit exposures (there is no borrowing limit if $\bar{b} = \infty$). The first-order conditions for this problem imply

$$b_t^{ij}(a, \tilde{a}) \in \arg \max_{b \in \bar{\mathbb{R}}} S_t^{ij}(a, \tilde{a}, b) \tag{1}$$

$$e^{-r(\bar{T}-t)}R_t^{ji}(\tilde{a}, a) = \theta_{ij} \left[V_t^j(\tilde{a} + b_t^{ij}(a, \tilde{a})) - V_t^j(\tilde{a}) \right] + \theta_{ji} \left[V_t^i(a) - V_t^i(a - b_t^{ij}(a, \tilde{a})) \right], \tag{2}$$

where

$$S_t^{ij}(a, \tilde{a}, b) \equiv V_t^i(a - b) + V_t^j(\tilde{a} + b) - V_t^i(a) - V_t^j(\tilde{a}).$$

Condition (1) characterizes the loan size, and (2) gives the repayment given the loan size. The implied gross interest rate on this loan is

$$1 + \rho_t^{ji}(\tilde{a}, a) = \frac{R_t^{ji}(\tilde{a}, a)}{b_t^{ij}(a, \tilde{a})}.$$

In Appendix A (Lemma 2) we show that the value function $V_t^i(a)$ satisfies

$$\begin{aligned} rV_t^i(a) - \dot{V}_t^i(a) &= u_i(a) + \lambda_i \sum_{j \in \mathbb{N}} \pi_j \int [V_t^i(a-z) - V_t^i(a)] dG_{ij}(z) \\ &+ \beta_i \sum_{j \in \mathbb{N}} \sigma_j \theta_{ij} \int \max_{b \in \mathbb{R}} S_t^{ij}(a, \tilde{a}, b) dF_t^j(\tilde{a}), \end{aligned} \quad (3)$$

with boundary condition $V_T^i(a) = U_i(a)$, where

$$\pi_j \equiv \frac{\lambda_j n_j}{\sum_{i \in \mathbb{N}} \lambda_i n_i}$$

is the probability the counterparty in a bilateral payment is of type j , and

$$\sigma_j \equiv \frac{\beta_j n_j}{\sum_{k \in \mathbb{N}} \beta_k n_k}$$

is the probability the counterparty in a bilateral trade is of type j .

Let $f_t^i \equiv dF_t^i$ denote the probability density function of reserve holdings among banks of type i at time t . This density follows the law of motion

$$\begin{aligned} \dot{f}_t^i(a) + (\beta_i + \lambda_i) f_t^i(a) &= \lambda_i \sum_{j \in \mathbb{N}} \pi_j \int \int \mathbb{I}_{\{x-z=a\}} dG_{ij}(z) dF_t^j(x) \\ &+ \beta_i \sum_{j \in \mathbb{N}} \sigma_j \int \int \mathbb{I}_{\{x-b_t^{ij}(x, \tilde{a})=a\}} dF_t^j(\tilde{a}) dF_t^i(x). \end{aligned} \quad (4)$$

Hereafter, let $\mathbf{U}(\cdot) = \{U_i(\cdot)\}_{i \in \mathbb{N}}$, $\mathbf{V}_t(\cdot) = \{V_t^i(\cdot)\}_{i \in \mathbb{N}}$, $\mathbf{b}_t(\cdot, \cdot) = \{b_t^{ij}(\cdot, \cdot)\}_{i, j \in \mathbb{N}}$, $\mathbf{R}_t(\cdot, \cdot) = \{R_t^{ij}(\cdot, \cdot)\}_{i, j \in \mathbb{N}}$, and $\mathbf{F}_t(\cdot) = \{F_t^i(\cdot)\}_{i \in \mathbb{N}}$.

Definition 1 *An equilibrium is a time-path $\{\mathbf{b}_t(\cdot, \cdot), \mathbf{R}_t(\cdot, \cdot), \mathbf{V}_t(\cdot), \mathbf{F}_t(\cdot)\}_{t \in \mathbb{T}}$ that satisfies (1), (2), (3), and (4), given the initial condition \mathbf{F}_0 and the terminal condition $\mathbf{V}_T = \mathbf{U}$.*

3 Data

In this section we document the fed funds market facts that will guide the quantitative implementation of the theory. Section 3.1 presents the joint distribution of two bank-level measures of fed funds trading activity: a bank's *participation rate* in marketwide trade volume, and a *reallocation index* that quantifies the degree to which a bank is a net borrower or lender of funds. Section 3.2 reports estimates of the frequency and size distribution of micro-level intraday payments between banks. Section 3.3 presents estimates of a typical beginning-of-day

cross-sectional distribution of reserve balances. Section 3.4 reports estimates of the distribution of aggregate daily reserve-draining shocks for the fed funds market since the GFC of 2007-2008. Section 3.5 presents empirical estimates of the slope of the aggregate demand for reserve balances. Finally, Section 3.6 describes an empirical interpolation procedure to map changes in the aggregate quantity of reserves into changes in the cross-sectional distributions of reserves that is consistent with available observations, and will be used in our quantitative analysis.

Since some of the regulations introduced in the wake of the GFC are likely to have affected trading incentives in the fed funds market, we report facts separately for the period before, and after these regulations had been implemented.⁹ In this section we use the years 2006 and 2019 as typical pre- and post-GFC-regulation periods, respectively. However, since some of our quantitative exercises will require sample variation in the aggregate quantity of reserves while keeping regulation constant, we will also report facts for the years 2014 and 2017.¹⁰

We use transaction data from the Fedwire Funds Service (*Fedwire*). Our typical Fedwire participant, which we call a *bank*, corresponds to a *bank holding company*. Our sample consists of 754 Fedwire participants for the year 2006, 404 for the year 2014, 395 for the year 2017, and 412 for the year 2019.¹¹ We use a modified version of the *Furfine algorithm* to identify overnight loans of reserves from the universe of Fedwire transfers; and we regard the remaining transactions as payments (presumably unrelated to loan issuance or repayment).¹² We focus on transactions that occur between 9:00 am and 6:30 pm EST.

⁹Some of these regulations increased the shadow value of liquid assets (including reserves), or introduced leverage constraints that increased the shadow cost of borrowing funds (including overnight fed funds purchases). Two prominent examples of such regulations are the *Liquidity Coverage Ratio* (LCR) and the *Supplementary Leverage Ratio* (SLR) requirements. We discuss these regulations in Appendix B.

¹⁰The LCR was phased in between January 2015 and January 2017. Non-foreign bank organizations began reporting SLR to U.S. regulators in July of 2013, SLR disclosures become mandatory in January of 2015, and SLR compliance became mandatory in January of 2018. We regard 2006 and 2014 as pre-GFC-regulation years, and the 2017 and 2019 as post-GFC-regulation years. In terms exploiting sample variability in the quantity of outstanding reserves in the system, the years 2006, 2014, 2017, and 2019 are natural benchmarks for the following reasons. The year 2006 is a typical pre-GFC period with excess reserves close to zero, and the year 2014 is a post-GFC but pre-GFC-regulation period with very high level of excess reserves (close to the pre-2020 historical peak). The year 2017 is a post-GFC-regulation period with very high level of excess reserves (again, close to the pre-2020 historical peak), while the year 2019 has the lowest level of excess reserves in the post-GFC-regulation era.

¹¹In Appendix D (Section D.1.2) we explain our sample selection criteria, and how we assigned Fedwire transactions to bank holding companies.

¹²The algorithm, which is based on Furfine (1999), was made available to us by the Money Market Analysis Section at the Monetary Affairs Division of the Federal Reserve Board.

3.1 Fed funds trading network

Let \mathbb{B} denote the collection of banks in our sample in a given year, and \mathbb{Y} denote the collection of all *trading periods* in that year.¹³ Let v_{nd}^e be the dollar value of all loans extended by bank $n \in \mathbb{B}$ in period $d \in \mathbb{Y}$, and use $v_d \equiv \sum_{n \in \mathbb{B}} v_{nd}^e$ to denote the dollar value of all the loans traded in period d . Also, let v_{nd}^r be the dollar value of all loans received by bank $n \in \mathbb{B}$ in period $d \in \mathbb{Y}$. For each bank n and period d , define

$$\begin{aligned}\mathcal{P}_{nd} &\equiv \frac{v_{nd}^e + v_{nd}^r}{2v_d} \\ \mathcal{R}_{nd} &\equiv \frac{v_{nd}^e - v_{nd}^r}{v_{nd}^e + v_{nd}^r}.\end{aligned}$$

We refer to \mathcal{P}_{nd} as bank n 's *participation rate* during period d , since it measures the share of the total period- d trade volume that is accounted for by bank n 's trading activity. For any given bank n in period d , $\mathcal{P}_{nd} \in [0, 1/2]$, with $\mathcal{P}_{nd} = 0$ corresponding to a bank that did not trade, and $\mathcal{P}_{nd} = 1/2$ corresponding to a bank that acted as a counterparty in every trade. In general, if a bank n participated as a counterparty in $x\%$ of the dollar value of all the loans traded in period d , then $2\mathcal{P}_{nd} = x/100$. We refer to \mathcal{R}_{nd} as bank n 's *reallocation index* during period d , since it is an index of the degree to which a bank is a net borrower or lender of funds. For any given bank n in period d , $\mathcal{R}_{nd} \in [-1, 1]$, with $\mathcal{R}_{nd} = -1$ corresponding to a bank that only borrowed, $\mathcal{R}_{nd} = 1$ corresponding to a bank that only lent, and $\mathcal{R}_{nd} = 0$ corresponding to a bank whose trading activity in period d consisted of pure intermediation. A typical bank n will have either $\mathcal{R}_{nd} \in (-1, 0)$, meaning it is a net borrower that engaged in some intermediation, or $\mathcal{R}_{nd} \in (0, 1)$, meaning it is a net lender that engaged in some intermediation.¹⁴ To provide a parsimonious description of the typical trading activity for each bank, we construct a bank-level

¹³In our empirical work, a *trading period* will correspond either to a *trading day*, or to a typical 14-day (*reserve*) *maintenance period* used to calculate a bank's reserve requirement. Our convention is to use \mathbb{Y} to denote a generic set of trading periods in a year, \mathbb{D} to denote the set of trading days in a year, and \mathbb{H} to denote the set of maintenance periods in a year. See Section B.1 in Appendix B for institutional details on reserve requirements and maintenance periods.

¹⁴Notice that $\mathcal{X}_{nd} \equiv 1 - |\mathcal{R}_{nd}|$ is a measure of the proportion of the total volume of funds traded by bank n in period d that the bank *intermediated* during that period, and $(v_{nd}^e + v_{nd}^r)\mathcal{X}_{nd}$ is what Afonso and Lagos (2015b) call *excess funds reallocation* (a measure of the volume of funds that an individual bank trades over and above what is required to accommodate its daily net trade).

participation rate and reallocation index averaged over all trading periods in a given year, i.e.,

$$\begin{aligned}\mathcal{P}_n &= \frac{1}{N_Y} \sum_{d \in \mathbb{Y}} \mathcal{P}_{nd} \\ \mathcal{R}_n &= \frac{1}{N_Y} \sum_{d \in \mathbb{Y}} \mathcal{R}_{nd},\end{aligned}$$

where $N_Y \equiv \sum_{d \in \mathbb{Y}} \mathbb{I}_{\{d \in \mathbb{Y}\}}$ is the number of trading periods in the year, and each trading period corresponds to one of bank n 's (reserve) maintenance periods during the year.¹⁵

Figure 2 shows the empirical cumulative distribution function (ECDF) of participation rates for the banks that are in our sample in the year 2006 (the circles) and the banks that are in our sample in the year 2019 (the crosses). We use the bank-level participation rate to sort each bank into one of three *groups*, denoted S , M , and F , depending on whether the bank's participation rate is low, medium, or high, respectively.¹⁶ Specifically, in each year we label the 4 banks with highest participation rate as, F ; the banks outside the top 4 that have participation rate at least as large as 0.5%, as M ; and all other banks, as S . Notice that individually, each of the top four most active banks that compose group F participated as a counterparty roughly in at least 10% of the total volume of loans traded in an average reserve maintenance period. And together, these four banks accounted for a large share of the aggregate trade volume: 45.6% in 2006, and 43.1% in 2019. In contrast, the large majority of banks, which belong to group S , have extremely low participation rates. We regard this large skewness in loan trading activity across banks as a key empirical regularity of the fed funds market structure.

Among the institutions assigned to group S based on the ECDF in Figure 2, there is a subgroup of non-bank Fedwire participants typically referred to as *Government Sponsored Enterprises* (GSEs), which includes the Federal Home Loan Banks, the Federal National Mortgage Association (Fannie Mae), and the Federal Home Loan Mortgage Corporation (Freddie Mac). Even though on the basis of their trading activity GSEs would belong in group S , in what follows we consider them a different type of participant because their business model and regulatory treatment make their payoffs from holding reserves different from the rest of the participating institutions.¹⁷ To offer a parsimonious representation of the data, we will sort institutions into

¹⁵See Appendix B (Section B.1) for institutional information on maintenance periods.

¹⁶The mnemonic is that banks of type S , M , and F , are slow, medium, and fast, at contacting counterparties.

¹⁷In contrast to banks, GSEs have very predictable cashflows (so payment shocks are not relevant for their day-to-day trading motives), and for most of our sample they did not earn interest on reserves—although nowadays they may lend reserves in the Federal Reserve's overnight reverse repo (ONRRP) facility.

four *types*, i.e., $\mathbb{N} = \{F, M, S, G\}$. Types F , M , and S correspond to the F , M , and S groups defined above, but excluding GSEs, and type G is composed exclusively of GSEs.¹⁸

Figure 3 shows the location of each bank type $i \in \{F, M, S, G\}$ in the coordinate axes defined by the reallocation index, \mathcal{R}_i , and the participation rate, \mathcal{P}_i , in the years 2006, 2014, 2017, and 2019. The figure shows an empirical trading network that conveys information on the distribution of trading activity across bank types, the flows of reserves implied by the fed funds lending among the four types of banks, and the average interest rates on the underlying loans. The participation, reallocation, and loans measures are all computed at the bank-type level.¹⁹ Each node represents the set of banks assigned to a particular type, labeled accordingly as F , M , S , or G . The arrows from one node to another represent loans extended from banks of that type to the other. The position of each node indicates how active the corresponding bank type is in the fed funds market and whether banks of that type are, on average, net borrowers, net lenders, or intermediaries. The size of each node is proportional to the volume of trade between banks of the that type. The width of each arrow is proportional to the volume of trade between the bank types connected by the arrow. The colors of the arrows and nodes are: light blue, dark blue, light red, or dark red, if the volume-weighted average interest rate on the loans between the two types of banks, expressed as a spread over the EFFR, falls in the first, second, third, or fourth quartile, respectively.²⁰

While specifics vary somewhat across years, several stable trading patterns emerge from Figure 3. Banks of type F account for about 1/2 of aggregate trade volume (i.e., $\mathcal{P}_F \approx 1/2$) and intermediate a large share of what they trade, with a tendency to act as net lenders. Banks of type M and banks of type S tend to be net borrowers; the former account for more than 1/4 of aggregate trade volume, and the latter much less (e.g., less than a quarter in 2006, and

¹⁸Our sample for 2006 consists of 4 banks of type F , 22 banks of type M , 716 banks of type S , and 12 GSEs. Our sample for 2019 consists of 4 banks of type F , 18 banks of type M , 379 banks of type S , and 11 GSEs. If we apply the same classification criteria for the years 2014 and 2017, we find that our sample for 2014 consists of 4 banks of type F , 15 banks of type M , 373 banks of type S , and 12 GSEs, while our sample for 2017 consists of 4 banks of type F , 18 banks of type M , 362 banks of type S , and 11 GSEs.

¹⁹The participation rate for each bank type $i \in \{F, M, S, G\}$ on a given year was calculated as follows. For each maintenance period, we summed the participation rates of all the banks of a given type, and then averaged across all maintenance periods in the year. The reallocation index for each bank type is calculated as follows. For each maintenance period, we summed all the loans sent, and all the loans received, by banks of a given type, and used these aggregate measures of loans sent and received by the type to calculate the reallocation index for that bank type in that given maintenance period, and then averaged across all maintenance periods in the year. We followed the same aggregation procedure to calculate volume-weighted interest rates across groups. See Appendix D.2 (Section D.2.2) for more details.

²⁰Arrow widths and node sizes are defined relative to trades within a year; thus not comparable across years.

less than 1/8 in later years). GSEs account for about a 1/8 of aggregate trade volume, and participate almost exclusively as lenders.

3.2 Interbank payments

In the previous section we analyzed transfers of reserves associated with overnight borrowing and lending between banks. In this section we focus on transfers that are unrelated to loan issuance or repayment. We regard these transfers as *payments*, which may reflect transactions originated by the banks' clientele, or by sections of the bank other than the ones in charge of actively managing reserve balances.

We identify as *payments* all Fedwire transfers that are not flagged as loans or repayments by the Furfine algorithm. These payments are likely to have a predictable component, but also a random component, which we refer to as *payment shocks*. Since these components affect trading incentives differently in the theory, we construct a measure of the predictable component, and estimate a process for the payment shocks of a typical bank of type F , M , or S .²¹ As in the theory, we model payment shocks as a compound process with a parameter that determines the frequency with which a bank of type i receives a payment shock (i.e., λ_i in the theory), and a conditional probability distribution for the payment size, which is allowed to depend on the types of the banks sending and receiving the payment (i.e., G_{ij} in the theory). Next, we describe our procedure to estimate the process for high-frequency interbank payment shocks.

Let \mathbf{T} denote the set of all one-second time intervals in a trading day $d \in \mathbb{D}$. For every pair of banks $m, n \in \mathbb{B}$, let $s_{mn}(t, d) \in \mathbb{R}$ denote the dollar value of all payments from bank m to bank n in the one-second time interval $t \in \mathbf{T}$ during trading day $d \in \mathbb{D}$.²² Let \bar{s}_{mn} denote the value of the average payment between banks m and n in a given year, and define $\tilde{s}_{mn}(t, d) \equiv s_{mn}(t, d) - \bar{s}_{mn}$ for all $(t, d) \in \mathbf{T} \times \mathbb{D}$. In this way, we decompose every high-frequency payment $s_{mn}(t, d)$ between a pair of banks into a *predictable component*, \bar{s}_{mn} , and a *payment shock*, $\tilde{s}_{mn}(t, d)$. For each pair of bank types $i, j \in \mathbb{N}$, we pool all payment shocks to form the data set

$$\tilde{\mathbb{S}}^{ij} = \{\tilde{s}_{mn}(t, d) : m \in \mathbb{B}_i, n \in \mathbb{B}_j \text{ for all } (t, d) \in \mathbf{T} \times \mathbb{D}\},$$

where \mathbb{B}_i is the set of banks of type $i \in \mathbb{N}$. We then use the data set $\tilde{\mathbb{S}}^{ij}$ to estimate a

²¹The business model of a GSE makes its reserve balances unlikely to be subject to unexpected payment shocks of significant magnitude, so we regard all GSE payments as predictable.

²²The bilateral payment credits bank n 's account if $0 < s_{mn}(t, d)$, and bank m 's account if $s_{mn}(t, d) < 0$.

Gaussian kernel density that we regard as the size distribution of payment shocks between each pair of bank types i and j , i.e., the empirical counterpart of the probability density function corresponding to G_{ij} in the theory.²³ Figures 4 and 5 display the empirical histogram along with the corresponding estimated kernel of payment shocks for each pair of bank types using data from the years 2006, and 2019, respectively.

For each bank type $i \in \mathbb{N}$ we estimate the empirical counterpart of λ_i in our theory, as the average number of payment shocks that a typical bank of type i receives in a one-second time interval, $t \in \mathbf{T}$, during a trading day, d , in year \mathbb{Y} . Let $f_m(t, d)$ denote the number of payment shocks between a bank $m \in \mathbb{B}$ and any other bank during the one-second time interval t in trading day d , i.e., $f_m(t, d) = \sum_{n \in \mathbb{B} \setminus \{m\}} \mathbb{I}_{\{s_{mn}(t,d) \neq 0\}}$. The corresponding average across seconds in a trading day, and trading days in the year is $\bar{f}_m = \frac{1}{N_D} \sum_{d \in \mathbb{D}} \left[\frac{1}{N_T} \sum_{t \in \mathbf{T}} f_m(t, d) \right]$, where $N_T \equiv \sum_{t \in \mathbf{T}} \mathbb{I}_{\{t \in \mathbf{T}\}}$ is the number of seconds in a trading day, and $N_D \equiv \sum_{d \in \mathbb{D}} \mathbb{I}_{\{d \in \mathbb{D}\}}$ is the number of trading days in a year. We use these bank-level empirical frequencies of payment shocks to estimate the probability that an average bank of type $i \in \{F, M, S\}$ receives a payment shock in a typical one-second time period, i.e., we set $\lambda_i = \frac{1}{N_i} \sum_{m \in \mathbb{B}_i} \bar{f}_m$, where $N_i \equiv \sum_{m \in \mathbb{B}} \mathbb{I}_{\{m \in \mathbb{B}_i\}}$ denotes the number of banks of type i in our sample. The estimates for $\{\lambda_i\}_{i \in \{F, M, S\}}$ for the years 2019 and 2006 are reported in Table 1 and Table 3, respectively.²⁴

3.3 Distribution of reserve balances

In this section we estimate beginning-of-day distributions of reserve balances (for each bank type) that are the empirical counterparts of the beginning-of-day distributions in the theory, i.e., $\{F_0^i\}$. Our calculations begin with a primitive bank-level quantity of reserves, and involve constructing a notion of *unencumbered* excess reserves by subtracting regulatory reserve requirements, and netting predictable Fedwire transfers (both, outright payments, and fed funds repayments).

For each bank in our sample, the Monetary Policy Operations and Analysis (MPOA) section at the Monetary Affairs Division at the Federal Reserve Board calculates the daily reserve balance at 6:30 pm. We devise an algorithm that uses this end-of-day balance to calculate the “basic” beginning-of-day balance at 9:00 am on the following day for each bank. Specifically, the algorithm starts with the bank’s end-of-day balance for day $d - 1$ provided by MPOA, adds

²³See Appendix D (Section D.2.3) estimation details.

²⁴We set $\lambda_G = 0$ for every year, for the reasons explained in footnote 21.

all fed funds repayments received during day d , and subtracts all fed funds repayments sent during day d that correspond to fed funds loans originated during day $d-1$.²⁵ For each bank m , and each reserve maintenance period, h , that belongs to the set \mathbb{H} of all maintenance periods in a given year, we calculate the average beginning-of-day balance across trading days in the maintenance period, which we denote $a_m(h)$.²⁶ We make two additional adjustments to this average “basic” measure of beginning-of-day balance at the bank level.

The first adjustment consists of subtracting the quantity of *required reserves*, i.e., the minimum level of reserves that the bank must hold during the maintenance period in order to comply with Regulation D and the minimum *Liquidity Coverage Ratio* requirement (LCR).²⁷ Specifically, for each individual bank m , we compute the average beginning-of-day *excess reserves* during a maintenance period h , as $x_m(h) = a_m(h) - \underline{a}_m^D(h) - \underline{a}_m^L(h)$, where $\underline{a}_m^D(h)$ and $\underline{a}_m^L(h)$ denote the Regulation D and LCR reserve requirements, respectively.²⁸

²⁵Repayments are identified using the send-receive matching from the Furfine algorithm. The rationale for netting the *predictable transfers*, which include the *repayments of fed funds* borrowed in the previous trading day, as well as the *predictable component of payments* (discussed below), is that through the lens of our theory, the beginning-of-day balance that is relevant for a bank’s incentives to trade reserves during the day ought to be net of *anticipated* transfers that the bank knows will receive or have to make during the trading day. The beginning-of-day- d balance for each GSE is constructed by taking the GSE’s end-of-day balance for day $d-1$ provided by MPOA, and netting all repayments of fed funds loans traded during day $d-1$ (between the GSE and any other bank that meets the sample selection criteria described in Section D.1.2), as well as *payments* sent or received during trading day d (and that involve *any* bank, not only those that meet the sample selection criteria described in Section D.1.2). The rationale for netting all transfers that will occur during day d to obtain the GSE’s balance at the beginning of day d is that a GSE’s business model generates very predictable cashflows, so through the lens of our theory, we regard the GSE as being able to predict all its intraday Fedwire transfers at the beginning of the trading day.

²⁶What motivates our focus on beginning-of-day balances *averaged over all trading days in a reserve maintenance period* is the fact that the reserve requirement regulations that influence banks’ payoffs from holding reserves must be met not on a daily basis, but on average over all days in the maintenance period. See Section B.1 in Appendix B for details on the reserve requirements stipulated by Regulation D.

²⁷Appendix B gives an overview of the relevant regulation. Our motivation for estimating reserves net of regulatory requirements is that this notion of *excess reserves* will play an important role in our quantitative theoretical exercises, e.g., it will be a key input to determine whether the central bank is implementing a monetary policy framework with “ample reserves” or a “corridor system”. For this reason, in the quantitative implementation of the theory we specify banks’ end-of-day payoffs in terms of excess reserves.

²⁸The bank-level data for Regulation D requirements are provided by MPOA. The LCR regulation requires a bank to maintain (typically on a daily basis) a quantity of *High Quality Liquid Assets* (HQLA) at least as large as a measure of total net cash outflows in a 30-day standardized stress scenario. Specifically, if we let $H_m(d)$ denote the quantity of qualifying HQLA held by bank m in a trading period d , and $L_m(d)$ denote the corresponding measure of outflows in the stress scenario, the LCR regulation requires $L_m(d) \leq H_m(d)$. Both these quantities are publicly available for each bank at a quarterly frequency (see Section D.1.3 in Appendix D for details). The set of qualifying HQLA includes reserves in excess of Regulation D, as well as securities issued or guaranteed by the U.S. Treasury (and also other securities, but subject to caps and haircuts). The fact that the LCR regulation allows banks to meet the requirement with assets other than reserves presents a challenge when trying to identify the quantity of reserves that bank m treats as “required” to satisfy the LCR constraint in period

The second adjustment to the average basic measure of a bank’s beginning-of-day balance consists of subtracting the *predictable component of payments*. Specifically, for each bank m we compute $\hat{s}_m = \sum_{n \in \mathbb{B}} \hat{s}_{mn}$, where $\hat{s}_{mn} \equiv \sum_{d \in \mathbb{D}} \frac{1}{N_D} \sum_{t \in \mathcal{T}} s_{mn}(t, d)$ is the average (over the set \mathbb{D} of N_D trading days in the year) net daily payment from bank m to bank n . Then, for each bank m and reserve maintenance period h , we construct $q_m(h) = x_m(h) - \hat{s}_m$, which is a bank’s average (across days in the maintenance period) beginning-of-day measure of *unencumbered reserves*.²⁹

For each bank type $i \in \mathbb{N}$, define

$$\mathbb{Q}^i = \{q_m(h) : m \in \mathbb{B}_i \text{ for all } h \in \mathbb{H}\}.$$

We pool the data in the set \mathbb{Q}^i and use it to estimate a Gaussian kernel density that we regard as the empirical counterpart of the beginning-of-day distribution of reserves, F_0^i , in the theory.³⁰

Figures 6-9 show the kernel density estimates of the distributions of reserves for each bank type $i \in \mathbb{N}$ for the years 2006, 2014, 2017, and 2019, respectively. In every year, the distribution of unencumbered reserves across banks of type S is fairly concentrated around zero. In 2006 (a typical year before the GFC), about 60% of bank-period observations for type S have beginning-of-day reserves close to zero, with dispersion in both directions. In 2014, 2017, and 2019 (the post-GFC period with very high level of total reserves), the pattern for banks of type S is similar: about 60% of bank-period observations have beginning-of-day reserves close to zero, with some bank-period observations with positive reserves, and almost no bank-period observations with negative reserves. The distributions of beginning-of-day reserves for banks of type F and M , on the other hand, exhibit significant dispersion. For type M there are virtually no bank-period observations with negative reserves for the years 2014, 2017, and 2019, and the dispersion over positive holdings is sizeable. For type F there is significant dispersion of reserves around zero in the years 2017 and 2019, largely due to the predictable component of payments.

d , i.e., $\underline{a}_m^L(d)$. Our strategy to tackle this identification problem is to set $\underline{a}_m^L(d) = \max(0, L_m(d) - A_m(d))$, where $A_m(d) \equiv H_m(d) - \max(0, a_m(d) - \underline{a}_m^D(d))$ is the quantity of qualifying HQLA in excess of (i.e., other than) reserves net of the Regulation D requirement. Notice that the resulting measure of excess reserves, $x_m(d)$, selects the largest level of excess reserves net of the Regulation D requirement that is consistent with the LCR constraint. (Section B.2.1 in Appendix B discusses our strategy to identify the quantity of required reserves induced by the LCR regulation.) For banks that are not subject to LCR regulation, we set $\underline{a}_m^L(d) = 0$. Since GSEs are not subject to Regulation D or LCR regulation, we set $\underline{a}_m^D(d) = \underline{a}_m^L(d) = 0$ for $m \in \mathbb{B}_G$.

²⁹Unless otherwise specified, whenever we refer to “beginning-of-day reserves”, we will be alluding to *unencumbered reserves*, i.e., the quantity of reserves in excess of Regulation D and LCR requirements, and net of predictable interbank Fedwire transfers.

³⁰See Appendix D (Section D.2.3) estimation details.

3.4 Reserve-draining shocks

The aggregate demand for reserves is determined by the decisions of individual banks, who demand reserve balances as payment instruments, as safe short-term investment vehicles, and to meet regulatory requirements. The aggregate supply of reserves, on the other hand, is largely determined by the central bank’s actions. But the central bank does not have complete control over the supply of reserves: The supply of reserves available to private banks also depends on transactions for which the Federal Reserve is not a counterparty, such as those that involve private-sector bank accounts and the account that the U.S. Treasury holds at the Federal Reserve. We will term the changes in the aggregate quantity of reserves resulting from the actions of entities other than the Federal Reserve, *exogenous supply shocks*. For example, whenever corporations or households pay taxes or purchase issuances of treasury securities, reserves are transferred from private banks to the Treasury’s account at the Federal Reserve, which from the perspective of domestic banks, amounts to an aggregate contractionary (reserve-draining) supply shock. Conversely, expansionary (reserve-augmenting) supply shocks take place whenever the Treasury makes payments to the private sector (e.g., when redeeming outstanding debt instruments).³¹ In this section we use daily data for the 2011-2019 sample period to estimate the size distribution of exogenous shocks to the supply of reserves.

Reserves were relatively scarce before 2007, and the Open Market Trading Desk (“the Desk”) at the Federal Reserve Bank of New York (FRB-NY) routinely conducted open-market operations to offset the effects of exogenous supply shocks on the fed funds rate. These systematic policy responses make it challenging to identify exogenous shifts in the supply of reserves in the pre-2007 period. The sharp increase in excess reserves and the very low fed funds rate target that followed the GFC made it unnecessary for the Desk to actively respond to daily market conditions in order to implement the target. In fact, post-2008, the Federal Reserve interventions that affected the stock of reserves were driven by longer-term objectives (e.g., implementation of quantitative easing policies) rather than by day-to-day managing of the fed funds rate in response to high-frequency exogenous supply shocks to the quantity of reserves. Thus, in the post-GFC era we can identify exogenous supply shocks using high-frequency (e.g., daily) changes in the aggregate quantity of reserves held by financial institutions. The middle

³¹Three other common sources of reserve-draining or reserve-augmenting shocks are: foreign official reverse repurchase agreements, changes in the quantity of currency in circulation (which imply swaps of currency for reserves or vice versa), and Federal Reserve “float” that is caused by the mismatch in timing between the debiting of reserves from a paying bank and the crediting of reserves to a receiving bank.

panel of Figure 10 shows that the variation in total reserves has been much larger since 2008, which is in line with our identifying assumption that the Desk did not react to exogenous supply shocks to the stock of reserves in the post-GFC period.

We estimate the distribution of reserve-draining shocks as follows. For each trading day d in the set \mathbb{D} of all trading days in a given year, let A_d denote the aggregate quantity of reserves held by all banks at the end of day d , and define the corresponding 40-day (two-sided) moving average, $\bar{A}_d \equiv \frac{1}{41} \sum_{k=-20}^{20} A_{d+k}$.³² The top panel of Figure 10 shows the time series $\{A_d, \bar{A}_d\}$ between the years 2001 and 2019. The middle panel of Figure 10 shows the deviations between total reserves and its own moving average, i.e., $\{z_d\}$, with $z_d \equiv A_d - \bar{A}_d$.³³ In time periods when the Federal Reserve does not react systematically to exogenous shocks to the supply of reserves, $\{z_d\}$ can be interpreted as a measure of the supply shocks themselves. Define the set $\mathbb{Z} = \{z_d : d \in \mathbb{D}\}$, where \mathbb{D} denotes the collection of trading days during the sample period January 2011–July 2019. We use the pooled data in the set \mathbb{Z} to estimate a Gaussian kernel density for the distribution of shocks to the aggregate quantity of reserves.³⁴ The bottom panel of Figure 10 displays the empirical histogram based on the daily observations in \mathbb{Z} , along with its kernel estimate. The figure also depicts the intervals that contain the daily realization of the “aggregate supply shock” with 99% or 95% probabilities, i.e., $[-\$279 \text{ bn}, \$130 \text{ bn}]$, and $[-\$115 \text{ bn}, \$99 \text{ bn}]$, respectively.

To assess the plausibility of our estimates, consider Anbil et al. (2020), who in the context of the market events of September 16–17, 2019, estimate a reserve-draining shock of \$120 bn, and remark “it is not uncommon for reserves to fall about \$100 bn over a day or two” (p. 5). Our estimates imply that the probability of a reserve-draining shock of \$110 bn or larger is about 2.5%.

3.5 Liquidity effect

In this section we present empirical estimates of the change in the fed funds rate in response to an exogenous change in the aggregate quantity of reserves—the so-called *liquidity effect*.³⁵ It is customary to think of the fed funds rate as resulting from the intersection of a vertical supply

³²For the purpose of these calculations we include *all* banks, not only those that meet the sample selection criteria based on fed funds trading activity described in Section D.1.2.

³³Since the daily time series cannot be made public, the top and middle panels of Figure 10 show the *weekly* versions. But we use the daily time series for the purposes of the kernel estimation discussed below.

³⁴See Appendix D (Section D.2.3) estimation details.

³⁵See Carpenter and Demiralp (2006) for a review, and Afonso et al. (2022) for more recent references.

and a downward sloping demand for reserves (e.g., as in Poole (1968)). Framed in this way, the slope of the demand for reserves is the key determinant of the liquidity effect. Traditionally the main challenge for estimation has been to identify *exogenous* shifts in the supply of reserves.³⁶

In an influential paper, Hamilton (1997) proposed a proxy for exogenous shifts in the aggregate quantity of reserves, and Carpenter and Demiralp (2006) subsequently proposed another.³⁷ The range of estimates obtained by Hamilton (1997) (for the period 1989/04/06–1991/11/27) and Carpenter and Demiralp (2006) (for the period 1989/05/19–2003/06/27) is similar: the estimated increase in the fed funds rate in response to an unexpected, temporary (one-day) \$1 bn aggregate reserve-draining shock, ranges between 1 and 2 basis points (and can be as high as 3 basis points on “settlement Wednesdays”).³⁸

To estimate the liquidity effect for the post-GFC sample period with large excess reserves that were not actively managed by the central bank, we run the following regression:

$$s_t - s_{t-1} = \gamma_0 + \gamma(Q_t - Q_{t-1}) + \varepsilon_t, \quad (5)$$

where s_t denotes the spread between the effective fed funds rate (EFFR, published by the FRB-NY) and the administered interest rate on reserves (IOR) on day t , Q_t denotes the aggregate quantity of reserves at the end of day t (provided by the Monetary Affairs Division at the Federal Reserve Board), ε_t is an error term, and γ is the coefficient of interest.

We estimate regression (5) at daily frequency for the sample period 2019/05/02–2019/09/13.

³⁶This was the main estimation challenge for the pre-GFC regime in which the Desk was actively conducting open-market operations reacting to market conditions in order to manage the fed funds rate. The challenges are different for the post-GFC era (e.g., until mid September 2019), when reserves were not actively managed by the Desk. For example, within the post-GFC period when the Federal Reserve began managing the fed funds rate by setting administered rates rather than the quantity of reserves, our theory prescribes controlling for the spreads between the administered rates (see Section 6.1).

³⁷Hamilton (1997) proposed the deviations between the actual end-of-day balance of the Treasury’s Fed account and an empirical forecast of the end-of-day balance of the Treasury’s Fed account as a proxy for unexpected changes in the quantity of reserves. Carpenter and Demiralp (2006) build on the work of Hamilton (1997) by replacing his measure of unexpected changes in the Treasury’s Fed account with a more accurate and comprehensive measure: the difference between the realized quantity of reserves on a given day, and the forecast for the quantity of reserves for that day that is used by the Desk (or the FRB) to perform its daily accommodative open-market operations. Relative to Hamilton’s, the Carpenter-Demiralp measure of unexpected changes in reserves is more comprehensive because it contemplates all possible sources of variation in the supply of reserves (not only fluctuations in the Treasury’s Fed account), and it is more accurate because, by definition, these daily “forecast misses” are changes in the quantity of reserves that the Desk did not accommodate.

³⁸Since this range of estimates was obtained from time series during a period in which reserves in excess of Regulation D were very close to zero, and post-GFC regulation had not been introduced, we will use it in our historical calibration exercise (Appendix F) to discipline the parameters that determine the magnitude of the liquidity effect in our quantitative theory *locally*, i.e., around the equilibrium point that results when excess reserves are close to zero.

We base our estimation on the year 2019 because it is the baseline year we will use to calibrate our theory in Section 4. Our identifying assumption is that the daily changes in the aggregate quantity of reserves can be regarded as exogenous because, as discussed in Section 3.4, the Federal Reserve was not actively managing the quantity of reserves in response to developments in the fed funds market during the post-GFC sample periods that we consider for this regression.³⁹ The estimate is $\gamma = -0.0119$ (significant at the 1% level), with 95% confidence interval $[-0.0187, -0.0052]$. Since the independent variable is measured in billions of dollars and the dependent variable in basis points, these estimates mean that a \$1 bn increase in the quantity of reserves decreases the EFRR-IOR spread by 0.01 basis points (i.e., about one hundred times smaller than the estimates obtained by Hamilton (1997) and Carpenter and Demiralp (2006) for the pre-GFC corridor system with scarce reserves).

The sample period that we use in our estimation is chosen so that the spread between the (primary credit) Discount-Window rate (DWR) and the overnight reverse repo rate (ONRRP), and the spread between the IOR and the ONRRP, are constant (and in particular, equal to 75 and 10 bps per annum, respectively, as in our baseline calibration of Section 4).⁴⁰ This is important because, as we show in Section 6.1, our theory predicts that changes in these spreads shift the aggregate demand for reserves. To illustrate the perils of not controlling for these spreads, we run regression (5) at daily frequency for an extended sample period: 2019/01/01–2019/09/13. This sample period consists of two subperiods with different spreads between administered rates: a first subperiod (from 2019/01/01 to 2019/05/01) with IOR-ONRRP spread equal to 15 bps, and a second subperiod (starting on 2019/05/02) with IOR-ONRRP spread equal to 10 bps. (The DWR-ONRRP spread is equal to 75 bps throughout.) The resulting estimate is $\gamma = -0.0062$ (significant at the 1% level), with 95% confidence interval $[-0.00975, -0.00264]$. Since the independent variable is measured in billions of dollars and the

³⁹The sample goes up to mid-September 2019, when the overnight money market rates exhibited unusual spikes and exhibited significant volatility. This sample includes 2019/09/13 (Friday) and deliberately stops there because on 2019/09/16 (Monday), in response to the fed funds rate printing at the upper limit the target range, the Desk announced an overnight repo operation to be conducted at 9:30 AM on 2019/09/17 (Tuesday), offering up to \$75 billion against Treasury, agency, and agency MBS collateral. This operation, which injected \$53 billion in additional reserves and led to an immediate decline in rates, was the first time since the GFC that the Desk conducted an open-market operation to manage the fed funds rate. The sample we use to estimate γ ought to end before this policy response since it would clearly violate our identifying assumption. See Afonso et al. (2022) for a more comprehensive estimation exercise under different identifying assumptions. See Anbil et al. (2020) for a detailed narrative of the money-market rate spikes of mid-September 2019, and Section 8 below for a quantitative theoretical analysis of this episode.

⁴⁰The time series of the three administered rates (DWR, IOR, ONRRP) are displayed in Figure 16.

dependent variable in basis points, this estimate means that a \$1 bn increase in the quantity of reserves decreases the EFFR-IOR spread by about 0.006 basis points.⁴¹

To validate our estimates, we can compare them with those from Afonso et al. (2022), who provide time-varying estimates for the period 2009-2021 of the slope of the aggregate demand for reserves using an instrumental variable approach combined with a time-varying vector autoregressive model of the joint dynamics of reserves and federal fund rates. The slope of the aggregate demand for reserves for the year 2019 estimated by Afonso et al. (2022) implies that a 1 percentage point increase in the ratio of total reserves to total assets held by commercial banks leads to a 1 basis point reduction in the EFFR-IOR spread (see the entry in panel (a), row 1 of the column labeled “2019” in their Table 1). Since the value of total assets held by commercial banks was about \$17,000 bn in 2019, a 1 percentage point daily increase in the ratio of total reserves to total assets held by commercial banks corresponds roughly to a \$170 bn increase in total reserves. Thus, the estimate for 2019 that Afonso et al. (2022) report in Table 1 means that a \$1 bn increase in the quantity of reserves decreases the EFFR-IOR spread by about 0.00588 basis points, which is essentially the same as the estimate we obtain from regression (5) when we do not control for variation in the IOR-ONRRP spread.

Estimates of the liquidity-effect coefficient (e.g., γ in our regression equation (5), or the analogous estimates from Hamilton (1997), Carpenter and Demiralp (2006), and Afonso et al. (2022)) are to be interpreted as *local* estimates of the slope of the aggregate demand for reserves, since they can be thought of as the empirical counterparts of the slope of the demand for reserves in the Poole (1968) model—calculated using a relatively narrow range of variation in the aggregate supply of reserves. Unlike the Poole (1968) model, our theory does not have a primitive demand for reserves. But as we change the exogenous quantity of reserves, the model traces out a series of equilibrium interest rates, which together with the respective quantities of reserves, can be regarded as a model-generated “demand for reserves”.

3.6 An interpolation procedure for counterfactual experiments

Several of the counterfactual and policy experiments that we conduct below involve changes in the aggregate quantity of reserves. Since our theory features *ex ante* heterogeneity in reserve

⁴¹Since in the quantitative implementation of the theory we focus on a subset of fed funds participants (see Section D.1.2 in Appendix D for our sample selection criteria), we have also run a version of (5) where the loan rate used to compute the spread s_t is the volume weighted average of loans in our sample, and the quantity of reserves Q_t is the aggregate level of reserves held by all banks in our sample. The estimate is $\gamma = -0.0057$, which is within the 95% confidence interval of the estimate reported above.

balances, changing the aggregate supply of reserves requires us to specify the underlying change in the distributions of reserve balances across banks. For example, in order to implement a \$1 bn decrease in the aggregate quantity of reserves in the model, we must specify the associated changes in the beginning-of-day distributions of reserve balances of the four bank types. How is the \$1 bn being drained exactly? Only from fast banks? Only from slow banks? Uniformly from all banks? We tackle this issue with a simple interpolation procedure that allows us to map changes in the aggregate quantity of reserves into changes in the cross-sectional distributions of reserves that is consistent with available observations.⁴² The procedure is as follows.

Let \bar{n}_Y^i denote the proportion of banks of type i in our sample for the year Y , and let \bar{F}_Y^i denote the empirical beginning-of-day distribution of reserve balances across banks of type i , estimated from all trading days in year Y (as described in Section 3.3). Let Y_0 and Y_1 denote two sample years for which we have estimates of $\{\bar{F}_{Y_0}^i, \bar{F}_{Y_1}^i\}_{i \in \mathbb{N}}$. For each $i \in \mathbb{N}$, and each $Y \in \{Y_0, Y_1\}$, discretize the continuous cumulative distribution function \bar{F}_Y^i with N quantiles, denoted $\{x_Y^i(p_n)\}_{n=1}^N$, where $\{p_n\}_{n=0}^{N+1}$ is a sequence that satisfies $p_{N+1} = 1 - p_0 = 1$, with $p_n < p_{n+1}$ for all $n \in \{0, \dots, N\}$, and $x_Y^i(p_n)$ is the number that satisfies $\bar{F}_Y^i(x_Y^i(p_n)) = p_n$ for each $n \in \{1, \dots, N\}$.⁴³ For each $i \in \mathbb{N}$, $Y \in \{Y_0, Y_1\}$, $n \in \{1, \dots, N\}$, and $\omega \in \mathbb{R}$, use the pair of quantiles $\{x_{Y_0}^i(p_n), x_{Y_1}^i(p_n)\}$ to define the *synthetic quantile*,

$$x_{Y_\omega}^i(p_n) \equiv \omega x_{Y_1}^i(p_n) + (1 - \omega) x_{Y_0}^i(p_n). \quad (6)$$

We then use $\omega \in \mathbb{R}$ to define a family of economies indexed by the following distribution of

⁴²Empirical studies (e.g., those that estimate the liquidity effect discussed in Section 3.5) typically abstract from how reserve-draining or reserve-augmenting shocks are distributed in the cross section of banks. The theoretical challenge of having to specify a path for the distribution of reserve balances associated with a certain path for the aggregate quantity of reserves (which is the variable we usually regard as being under direct control of the central bank) is common to all existing micro-based models of the fed funds market that allow for heterogeneity in reserve holdings across banks. Afonso and Lagos (2015b), for example, parametrize the beginning-of-day distribution of reserves with a Gaussian mixture with two components, and implement changes in the aggregate quantity of reserves by draining reserves from the two components in a way that their variances and the ratio of their means remain constant (see footnote 26, and Section C.2 in the Supplemental Material of Afonso and Lagos (2015b) for details). Afonso et al. (2019), whose main quantitative experiment involves draining a large quantity of aggregate reserves, assume a two-stage draining scheme: Reserves are drained exclusively from the banks with the largest initial holdings until their reserves become low enough; and are drained proportionately from all banks thereafter.

⁴³See Appendix C (Section C.1) for more details on the grids that we use in our quantitative implementation.

banks across types and distributions of reserves for each bank type $i \in \mathbb{N}$:

$$\bar{n}_{Y_\omega}^i \equiv \omega \bar{n}_{Y_1}^i + (1 - \omega) \bar{n}_{Y_0}^i \quad (7)$$

$$\bar{F}_{Y_\omega}^i(a) \equiv \sum_{n \in \{1, \dots, N\}: x_{Y_\omega}(p_n) \leq a} (p_n - p_{n-1}), \quad (8)$$

so the corresponding aggregate quantity of reserves is

$$Q_{Y_\omega} \equiv \sum_{i \in \mathbb{N}} \bar{n}_{Y_\omega}^i \int a d\bar{F}_{Y_\omega}^i(a). \quad (9)$$

Notice that for $\omega = 1$ the distribution of banks across types and the distributions of reserves for each bank type are as in the base year Y_1 , and for $\omega = 0$ they are as in the base year Y_0 . Thus, by varying ω on $[0, 1]$ we can use (9) to span any aggregate level of reserves between Q_{Y_1} (the aggregate supply of reserves held by all banks in our sample in base year Y_1) and Q_{Y_0} (the aggregate supply of reserves held by all banks in our sample in base year Y_0). Conversely, for any aggregate quantity of reserves, Q , between Q_{Y_1} and Q_{Y_0} , there is an $\omega \in [0, 1]$ implied by (9), denoted $\omega(Q)$, that decomposes Q into a particular distribution of banks across types and distributions of reserves for each bank type, namely $\left\{ \bar{n}_{Y_{\omega(Q)}}^i, \bar{F}_{Y_{\omega(Q)}}^i \right\}_{i \in \mathbb{N}}$ implied by (7) and (8). For any $\omega \in [0, 1]$ our procedure produces a distribution of banks across types and a set of distributions of reserves for each bank type that are linear interpolations of the corresponding distributions for the base years. We will use this procedure to conduct counterfactual and policy experiments in our quantitative model.⁴⁴

4 Calibration

In this section we calibrate the model to match the key statistics that describe fed funds trading activity in the year 2019.⁴⁵ The model primitives are: trading session, $[0, T]$, discount rate, r ,

⁴⁴The procedure also allows for linear extrapolations, e.g., corresponding to parametrizations with $\omega < 0$ or $\omega > 1$. An alternative to our empirical interpolation/extrapolation procedure would be to integrate a fully specified capital-structure theory of the bank into our dynamic stochastic heterogeneous-bank fed-funds trading model in order to establish a theoretical link between market conditions (e.g., policy choices of administered rates and aggregate supply of reserves) and the cross section of the composition of banks' assets, and in particular, their choices of reserve balances. The main challenge would then be to ensure that the endogenous portfolio choices implied by the theory are quantitatively consistent with the empirical paths for the cross-sectional distributions of reserves that have accompanied the observed long- and medium-term changes the aggregate supply of reserves. An attractive feature of our empirical interpolation procedure is that, by construction, it ensures that this is the case (at least for moderate deviations in the aggregate supply of reserves from those prevailing the base years). We think that integrating fed funds microstructure theory with a macroeconomic theory of the capital structure of the banking sector is a promising avenue of research (see Bianchi and Bigio (2022) for work along these lines).

⁴⁵We use 2019 as the baseline year for our calibration because it has the lowest level of total reserves of the current post-GFC-regulation policy regime (see Figure 16). As we explain in Section 6, this allows us to test

set of bank types, \mathbb{N} , population shares of bank types, $\{n_i\}_{i \in \mathbb{N}}$, beginning-of-day distributions of reserve balances, $\{F_0^i(a)\}_{i \in \mathbb{N}}$, payment shock frequencies, $\{\lambda_i\}_{i \in \mathbb{N}}$, conditional size distributions of payment shocks, $\{G_{ij}\}_{i,j \in \mathbb{N}}$, bargaining powers, $\{\theta_{ij}\}_{i,j \in \mathbb{N}}$, intraday payoffs, $\{u_i\}_{i \in \mathbb{N}}$, end-of-day payoffs, $\{U_i\}_{i \in \mathbb{N}}$, and trading frequencies, $\{\beta_i\}_{i \in \mathbb{N}}$. In the quantitative implementation it is useful to augment the model to include proportional borrowing costs, $\{\kappa_i\}_{i \in \mathbb{N}}$, that proxy for institutional and regulatory considerations that affect banks' incentives to borrow in the fed funds market. In Appendix A (Section A.2) we derive a generalization of (1), (2), and (3) to the case of proportional borrowing costs.

Our calibration strategy is as follows. We regard the trading session in the model as an average trading day in a typical 14-day reserve maintenance period. As discussed in Section 3, there is little trading activity between 9:00 pm on day $h - 1$ and 9:00 am on day h , so we think of $[0, T]$ as corresponding to the time interval that starts at 9:00 am and ends at 6:30 pm EST on an actual trading day. In the quantitative implementation of the theory we discretize the time interval $[0, T]$ into 800 periods, so each period in the model corresponds approximately to a 42-second interval of the trading day.⁴⁶ With a model period this short, we abstract from pure discounting, and set $r = 0$. As described in Section 3.1, we sort institutions into four types, i.e., $\mathbb{N} = \{F, M, S, G\}$, based on their participation rates in the volume of fed funds trade, business model, and regulatory treatment. We set $n_i = N_i / \sum_{j \in \mathbb{N}} N_j$, where N_i denotes the number of banks of type $i \in \mathbb{N}$ in the base year. We interpret reserve balances in the theory as

the quantitative predictions of the theory by varying the level of excess reserves from their lowest level in the post-GFC-regulation era, to the level they reached in the year 2017 (a post-GFC-regulation year with a level of excess reserves that is close to the pre-2020 historical peak).

⁴⁶A model period corresponding to 42 seconds is short enough to approximate the empirical frequency of loans even for the most active banks. Payment shocks, however, are much more frequent than loans: In Section 3.2 we had to use a period length of 1 second in order to get a good approximation to the empirical frequency of payment shocks (especially for fast banks, which typically experience several payment shocks per minute, and sometimes even more than one payment shock per second). In order to allow for such high frequency of payment shocks, we could simply discretize $[0, T]$ into 34,200 periods, each corresponding to 1 second. With so many periods, however, the computational burden would increase significantly, so we took a different approach. Payment shocks, although very frequent, are computationally cheap since they involve no optimization (they are just "forced" transfers between banks). Loans on the other hand, are computationally more expensive (they involve maximization of the joint surplus), but are also significantly less frequent than payment shocks in the data. In the quantitative implementation of the model, we balance these considerations as follows. We regard each model period as being composed of 42 *subperiods*, each corresponding to 1 second in the actual trading day. We then treat the first 41 subperiods as "payment-shock rounds" (in each of these rounds, there are only bilateral payment shocks among banks), and treat the 42nd subperiod as a "loan round" in which banks get bilateral opportunities to negotiate loans. In sum, this allows us to have payment shocks that are as frequent as 1 per second, and loans that are as frequent as one every 42 seconds, while economizing on computation time. See Appendix C for a more detailed discussion of computational issues.

a bank's *unencumbered reserves* in the data, and therefore set the theoretical beginning-of-day distributions, $\{F_0^i(\cdot)\}_{i \in \mathbb{N}}$, equal to the corresponding empirical kernel estimates reported in Section 3.3. The frequencies of payment shocks, $\{\lambda_i\}_{i \in \mathbb{N}}$, are calibrated to match the empirical one-second frequencies of payment shocks reported in Section 3.2.⁴⁷ The size distributions of payment shocks, $\{G_{ij}\}_{i,j \in \mathbb{N}}$, are set equal to the corresponding empirical kernel estimates reported in Section 3.2. We set $\theta_{ij} = \underline{\theta}$ if $i \in \{G\}$ and $j \in \mathbb{N} \setminus \{G\}$, and $\theta_{ij} = 1/2$ otherwise (i.e., unless one of the parties in the trade is a GSE, we abstract from differences in relative market power purely driven by a bank's *type*). We set $u_i(a) = 0$ for all $(a, i) \in \mathbb{R} \times \mathbb{N}$ (i.e., we abstract from banks' *intraday* payoffs from holding reserves, such as the regulatory costs associated with running an intraday overdraft with the Federal Reserve).

End-of-day payoffs are parametrized by

$$U_i(a) = (1 + \mathbb{I}_{\{0 \leq a\}} \bar{\iota}_r + \mathbb{I}_{\{a < 0\}} \bar{\iota}_w) a, \quad (10)$$

for any $(a, i) \in \mathbb{R} \times \{F, M, S\}$, where $\mathbb{I}_{\{a < 0\}}$ is an indicator function that equals 1 if $a < 0$ and 0 otherwise, $\bar{\iota}_r \equiv \iota_r + \iota_\ell$, $\bar{\iota}_w \equiv \iota_w + \iota_\ell + \iota_s$, and a denotes end-of-day balance in excess of reserve requirements.⁴⁸ We use ι_r to denote the interest rate that a bank earns from the Federal Reserve per dollar of end-of-day reserves (IOR), and ι_ℓ to represent a *liquidity return* that proxies for a bank's benefits from holding reserves that are not captured by the administered rates.⁴⁹ We use ι_w to denote the (primary credit) Discount-Window rate (DWR) that the Federal Reserve charges a bank that needs to borrow to make up an end-of-day shortfall of reserves relative to the required level, and ι_s to represent the additional costs associated with borrowing from the Discount Window.⁵⁰ For GSEs, the end-of-day payoff is $U_G(a) = (1 + \mathbb{I}_{\{0 \leq a\}} \bar{\iota}_o + \mathbb{I}_{\{a < 0\}} \bar{\iota}_w) a$, with $\bar{\iota}_o \equiv \iota_o + \iota_\ell$, where ι_o denotes the interest rate that the Federal Reserve offers on the

⁴⁷In the discrete-time approximation that we compute, λ_i corresponds to the probability that a bank of type i receives a payment shock in a one-second time interval (see footnote 46).

⁴⁸Since our calibration strategy maps beginning-of-day reserve balances in the theory to *unencumbered reserves* in the data, which are reserves in excess of reserve requirements (and net of predictable payments), we specify a bank's end-of-day payoff as a function of its *excess reserves*. This allows us to have end-of-day payoff functions that are *type specific* but not *bank specific*, despite the fact that in the data, two banks will typically have different reserve requirements even if they are of the same type, $i \in \mathbb{N}$. To see this, let $\mathcal{U}_i(b, \underline{b})$ be the end-of-day payoff of a bank of type i with reserve requirement \underline{b} , and reserve balance b (gross of the reserve requirement). We would parametrize this function as $\mathcal{U}_i(b, \underline{b}) = b + \bar{\iota}_r \underline{b} + (\mathbb{I}_{\{0 \leq b - \underline{b}\}} \bar{\iota}_r + \mathbb{I}_{\{b - \underline{b} < 0\}} \bar{\iota}_w) (b - \underline{b})$, which is equivalent to $U_i(a)$ in the sense that they only differ by a constant, i.e., $\mathcal{U}_i(b, \underline{b}) = U_i(a) + (1 + \bar{\iota}_r) \underline{b}$, where $a \equiv b - \underline{b}$ denotes excess reserves, as in (10).

⁴⁹For example, ι_ℓ may stand in for the additional return associated with the use of reserves as means of payment, or for the additional return resulting from lending reserves outside the fed funds market (e.g., in repo markets, or as loans to corporate or retail bank customers).

⁵⁰It is a well documented phenomenon that banks often borrow from other banks at a premium over the

overnight reverse repo facility.⁵¹ The administered rates, i.e., ι_r , ι_w , and ι_o , are set equal to their empirical counterparts in the base year.

The remaining eleven parameters, $\underline{\theta}$, ι_ℓ , ι_s , and $\{\beta_i, \kappa_i\}_{i \in \mathbb{N}}$, are calibrated so that the equilibrium of the model matches the following eleven empirical moments: (1) average value-weighted fed funds rate; (2) average value-weighted fed funds rate for loans with rates lower than the IOR; (3) regression estimates of the “liquidity effect” (at the average level of aggregate reserves outstanding in the base year, as reported in Section 3.5); (4) ratio of the average number of loans traded by banks of type F relative to the average number of loans traded by all banks; (5)-(8) reallocation indices $\{\mathcal{R}_i\}_{i \in \mathbb{N}}$ (as defined in Section 3.1); (9)-(11) participation rates $\{\mathcal{P}_i\}_{i \in \mathbb{N} \setminus \{F\}}$ (as defined in Section 3.1).⁵²

Table 1 reports the parameter values, the targeted moments, and the corresponding theoretical moments for the 2019 calibration. Banks of type F , M , and S , accounted for about 1%, 4%, and 92%, of all the institutions that were active in the fed funds market in 2019, respectively. To interpret the frequencies of payment shocks, $\{\lambda_i\}_{i \in \mathbb{N}}$, recall that λ_i represents the probability that a bank of type i receives a payment shock in a one-second time interval, so for example, $\lambda_M = 0.257$ implies a bank of type M receives a payment shock approximately every 4 seconds, on average. Similarly, $\lambda_F = 0.920$ implies a bank of type F receives approximately a payment shock per second, and $\lambda_S = 0.011$ implies a bank of type S receives a payment shock approximately every 90 seconds, on average. The values of ι_w (DWR), ι_r (IOR), and ι_o (ON-RRP) are set to 3.00%, 2.35%, and 2.25% per annum, respectively, which were the administered policy rates in effect from May through July of 2019.

The calibration strategy delivers a liquidity return (ι_ℓ) of 4.9 bps per annum, an additional cost associated to Discount-Window borrowing (ι_s) of 75.8 bps per annum (i.e., about one quarter of the DWR), and $\underline{\theta} = 1/20$, which means that a GSE reaps 5% of the total gains from

DWR. The most common explanation is stigma associated with Discount-Window borrowing (see e.g., Artuç and Demiralp (2010), Ennis and Weinberg (2013), Armantier et al. (2015), Ennis (2019), and Klee et al. (2021)). Another reason why fed funds may trade above the DWR is that Discount-Window loans must be collateralized, and Reserve Banks require a perfected security interest in all collateral pledged to secure these loans, which entails costs for the borrower. Assets accepted as collateral are assigned a lendable value deemed appropriate by the Reserve Bank that issues the loan (a market value or an internally-modelled fair market value estimate multiplied by a margin, possibly adjusted as a function of the financial condition of the borrowing institution). For details, see <https://www.frbdiscountwindow.org/Pages/General-Information/The-Discount-Window>.

⁵¹We use ι_o rather than ι_r in the payoff for GSEs because regulation prevents them from earning interest on reserves. We use ι_r rather than ι_o in the payoffs of other bank types because $\iota_o \leq \iota_r$ throughout our sample.

⁵²The participation rate of type F banks is not an explicit calibration target because it is implied by the participation rates of the other three bank types, since $\sum_{i \in \mathbb{N}} \mathcal{P}_i = 1$.

lending to a non-GSE. The frequency of trade, β_i , is interpreted as the probability that a bank of type i contacts a trading partner during a 42-second time interval. Thus, the calibrated values $\{\beta_i\}_{i \in \mathbb{N}}$ imply that banks of type F , M , S , and G , trade fed funds approximately every 23 minutes, 4.86 hours, 16.7 hours, and 3.24 hours, respectively. The calibration also ensures that the magnitude of the “liquidity effect” in the theory is in line with the range of empirical estimates for the year 2019 reported in Section 3.5.⁵³ The borrowing costs needed to match the calibration targets, $\{\kappa_i\}_{i \in \mathbb{N}}$, which proxy for institutional and regulatory considerations that affect banks’ incentives to buy fed funds, are positive for banks of type F and S , and zero for banks of type M .⁵⁴

5 Validation

In this section we report the model fit of empirical price and quantity observations not targeted in the calibration. We organize the material in four sections. The first focuses on the cross-sectional distribution of loan rates for all transactions. The second, on the distribution of loan rates for transactions with rates higher than the DWR. The third, on the distribution of borrowing and lending rates for each bank type. The fourth, on the trading network.

5.1 Distribution of loan rates

Figure 12 shows the empirical and theoretical cumulative distribution functions of bilateral negotiated fed funds rates in the year 2019, along with the administered rates prevailing in the sample period (all expressed in percent per annum).⁵⁵ The model delivers a reasonable fit for the distribution of bilateral fed funds rate, which was not targeted in the calibration.⁵⁶

⁵³Figure 11 shows the magnitude of the liquidity effect in the calibrated model along with the confidence bands for the regression estimates for the sample period 2019/05/02–2019/09/13 presented in Section 3.5. In the model, the liquidity effect is computed by extracting \$100 bn reserves (approximately 2 standard deviations of the size distribution of reserve-draining shocks) using the procedure described in Section 3.6. The figure shows that the model-generated liquidity effect is within the 95% confidence bands of the empirical estimate.

⁵⁴The value of κ_G is set large enough to match the observation that GSEs essentially do not borrow in the fed funds market, but its exact value is inconsequential.

⁵⁵Data are for every trading day in the period 2019/06/06–2019/07/31, which covers eight reserve maintenance periods during which the policy rate remained constant and the administered rates (DWR, IOR, ONRRP) were as in our baseline calibration. To obtain the equilibrium rates for 2019, the model is calibrated as in Table 1.

⁵⁶The model, however, does not generate enough dispersion of rates relative to the data. This is the case for loans that trade above the IOR, but also for loans that trade below the IOR. One way to match the larger empirical dispersion of loans with above-IOR rates would be to allow for heterogeneity in bargaining powers across banks of types, i.e., to let θ_{ij} differ in trades between two non-GSEs. Notice that a significant part of the large dispersion for below-IOR trades in the data comes from trades with rates lower than the ONRRP. This

5.2 Conditional distribution of loan rates in excess of DWR

In the model, as in the data, banks sometimes trade at rates higher than the DWR. In the model this is possible because ι_s is calibrated to a positive value. In this section we compare the theoretical and empirical distributions of traded rates conditional on the rate being higher than the DWR, which were not targeted in our calibration. During the sample period 2019/06/06–2019/07/31 the DWR was set at 3%; the 10th percentile, mean, and 90th percentile were 3%, 3.1%, and 3.3%, respectively, both in the data and in the calibrated model. The maximum loan rate observed in our data sample was 3.45%, and the maximum possible rate a bank is willing to pay in the equilibrium of the model is $\iota_w + \iota_\ell + \iota_s = 0.038$.

5.3 Bid-ask spreads

Each of the panels on the right side of Figure 13 shows an empirical cumulative distribution function of borrowed reserves over borrowing rates, denoted \mathcal{H}_i^B (represented by a solid line), and an empirical cumulative distribution function of lent reserves over lending rates, denoted \mathcal{H}_i^L (represented by a dashed line), for $i \in \{F, M, S\}$. In words, $\mathcal{H}_i^B(\iota)$ is the proportion of reserves borrowed by banks of type i that bear interest rates lower than ι , and $\mathcal{H}_i^L(\iota)$ is the proportion of reserves lent by banks of type i that bear interest rates lower than ι .

Each of the panels on the left side of Figure 13 shows the theoretical counterpart of the adjoining right-side panel. The top-left and middle-left panels show the theory predicts $\mathcal{H}_i^B(\iota) \leq \mathcal{H}_i^L(\iota)$ for $i \in \{F, M\}$. That is, banks of type F and M tend to borrow at lower rates than they earn when they lend. This theoretical prediction also holds in the data, as long as we focus on loans with rates that are not lower than the IOR (2.35%).⁵⁷ In contrast, according to the bottom-left panel, the theory predicts $\mathcal{H}_S^B(\iota) \leq \mathcal{H}_S^L(\iota)$, i.e., banks of type S tend to borrow at higher rates than they earn when they lend.⁵⁸ This theoretical prediction also holds in the data, and the fit is remarkably good for loans with rates that are not lower than the IOR.

observation is difficult to rationalize through the lens of the theory, and may be indicative of some repo loans being misclassified as fed funds in our dataset (e.g., as suggested by Armantier and Copeland (2015)).

⁵⁷As mentioned in footnote 56, rates below the IOR are likely to correspond to repo loans that are misclassified as fed funds by the Furfine algorithm.

⁵⁸The model counterparts of $\mathcal{H}_S^B(\iota)$ and $\mathcal{H}_S^L(\iota)$ are constructed excluding loans between a G and a bank of type S . The rationale is that our model abstracts from the institutional details that make these trades very rare in the data. For example, there was only one loan of this kind in our sample period.

5.4 Distributions of loan rates between pairs of bank types

Each of the panels on the right side of Figure 14 shows an empirical cumulative distribution of rates for loans extended from bank type $i \in \{F, M, S\}$ to bank type $j \in \{F, M, S\}$. For example, for each interest rate ι on the horizontal axis, the height of the curve labeled “ S ” in the top right panel represents the fraction of the total volume of loans extended from banks of type “ F ” to banks of type “ S ” with interest rate less than or equal to ι . Each of the panels on the left side of Figure 14 shows the theoretical counterparts of the adjoining right-side panel. The theory predicts that, regardless of lender type, banks of type “ S ” tend to borrow at higher rates than banks of other types, and this is true in the data.⁵⁹

5.5 Fed funds trading network

Figure 15 shows the empirical fed funds trading network for the year 2019 (bottom panel) and the corresponding trading network generated by the model (top panel). As explained in Section 3.1, these network plots show the location of the four bank types in the coordinate axes defined by the reallocation index, \mathcal{R}_i , and the participation rate, \mathcal{P}_i , and convey information on the sizes of the flows of reserves associated with fed funds lending across and within bank types, as well as on the average interest rates on the underlying loans.⁶⁰

The theoretical network matches several characteristics of the empirical one. For example, the model replicates quite well the direction and volume of the loans between and within bank types (represented by the direction of the arrows, their width, and the sizes of the nodes). In this regard, one difference is that the model predicts a significant volume of loans from GSEs to banks of type S that is not present in the data.⁶¹ The model predicts that GSEs tend to lend at lower rates than banks, and in particular, lower than the rates that banks of type F tend to charge banks of type M . However, the opposite is true in the data.⁶²

⁵⁹The theory also predicts that banks of type “ M ” tend to borrow at higher rates than banks of type “ F ”, but this is not as evident in the data.

⁶⁰In comparing the top and bottom panels of Figure 15, notice that while the positions of the four nodes in \mathcal{R}_i - \mathcal{P}_i space have been used as calibration targets, the remaining collection of statistics that shape these network representations were not targeted. This includes the size and color of each node, and the direction, color, and width of each arrow.

⁶¹This discrepancy is likely due to the fact that our theory abstracts from the real-world institutional details that cause GSEs to lend reserves only to a relatively small subset of counterparties, which tend to be big banks that are very active in fed funds trading.

⁶²This discrepancy may be due to the fact that the Furfine algorithm, used to identify overnight uncollateralized loans from the universe of Fedwire transfers, may pick up some overnight *collateralized* loans (i.e., repos) that trade at lower rates. Armantier and Copeland (2015) provide some evidence consistent with this interpretation.

6 Aggregate demand for reserves

It is customary to think of the fed funds rate in the context of a static, perfectly competitive loan market, i.e., as being determined by the intersection of a vertical supply of reserves controlled by the central bank, and an aggregate downward-sloping demand for reserves implied by the solution to a reserve-management problem faced by individual banks.⁶³ Our over-the-counter theory does not involve a bank-level reserve demand.⁶⁴ However, as we vary the aggregate quantity of reserves, Q , e.g., by changing the beginning-of-day distributions of reserves, our theory can generate a negative relationship between Q and a volume-weighted average of all the equilibrium bilateral loan rates, ι^* , which can be interpreted as the *aggregate demand for reserves* implied by the theory. We can write this relationship as $\iota^* = \mathcal{D}(Q; \Pi)$, where $Q = \sum_{i \in \mathbb{N}} n_i \int adF_0^i(a)$ is the supply of reserves, and $\Pi \equiv \{\beta_i, \lambda_i, \{\theta_{ij}, G_{ij}\}_{j \in \mathbb{N}}, u_i, U_i\}_{i \in \mathbb{N}}$ is the full set of model primitives, which makes clear that the mapping depends on all the structural parameters of the model.⁶⁵

Consider the model calibrated to the year 2019, as described in Table 1. Then, using the notation introduced in Section 3.6, let $Y_0 = 2017$ and $Y_1 = 2019$, i.e., Y_0 and Y_1 represent the years 2017, and 2019, respectively, with \bar{n}_{2017}^i and \bar{F}_{2017}^i given by the estimates reported in Section 3.3. Construct a grid, $\mathbb{G} \subset \mathbb{R}$ for ω , and for each $\omega \in \mathbb{G}$, use the interpolation procedure described by (7) and (8) to generate the sample $\{(\bar{n}_{Y\omega}^i, \bar{F}_{Y\omega}^i)\}_{(i,\omega) \in \mathbb{N} \times \mathbb{G}}$. For each pair $(\bar{n}_{Y\omega}^i, \bar{F}_{Y\omega}^i)$, compute the equilibrium value-weighted fed funds rate, which we denote $\iota_{Y\omega}^*$, and let $Q_{Y\omega} \equiv \sum_{i \in \mathbb{N}} \bar{n}_{Y\omega}^i \int ad\bar{F}_{Y\omega}^i(a)$. This procedure delivers a collection of pairs, $\{(Q_{Y\omega}, \iota_{Y\omega}^*)\}_{\omega \in \mathbb{G}}$, that define the mapping $\iota_{Y\omega}^* = \mathcal{D}(Q_{Y\omega}; \Pi)$. This mapping—the aggregate demand for reserves generated by the theory—is the curve labeled “Benchmark” in all panels of Figure 17.⁶⁶

⁶³E.g., as in the “Poole model” described in footnote 1.

⁶⁴This is because in our dynamic over-the-counter theory, the end-of-day reserve holding of an individual bank is a random variable that depends on the bank’s beginning-of-day balance, the number of counterparties it encounters throughout the trading session, and the bilateral bargaining outcomes that in turn depend on the counterparties’ individual characteristics (such as their reserve holdings at the time of the trade, bargaining powers, and abilities to contact other counterparties). In other words, the fact that our theory is distinctively non-Walrasian implies there is no natural or useful counterpart to the notion of an optimal quantity of reserves chosen by an individual bank who can borrow and lend frictionlessly at a given market interest rate.

⁶⁵For example, equation (16) in Afonso and Lagos (2015b) gives an explicit formula for this mapping for a special case of our model that allows an analytical solution (identical banks and heterogeneity in reserve balances restricted to the set $\{0, 1, 2\}$).

⁶⁶We use 2017 and 2019 as endpoints for our interpolation procedure because this choice maximizes the sample variation in total reserves during the post-GFC-regulation era (prior to the large reserve injection that took place in response to the COVID shock in the year 2020). Specifically, as illustrated in Figure 16, 2017 is the post-GFC-regulation year with highest level of total reserves (\$2,254.27 bn, which is roughly the pre-2020 historical peak),

We have calibrated the model using empirical beginning-of-day distributions of reserves that are net of predictable transfers, net of Regulation D and LCR requirements, and that only include banks that had at least one fed funds transaction in the baseline year. Hence, the measure “ Q ” of aggregate reserves reported in the primary horizontal axis of any figure that displays the aggregate demand for reserves generated by the theory (e.g., Figure 17, Figure 19, Figure 21, or the right panels of Figure 18) is the sum of *excess* reserves (net of Regulation D and LCR requirements) across *active* banks (in the sense of having had at least one fed funds transaction in the baseline year). We term this measure of aggregate reserves *active excess reserves*, to distinguish it from *total reserves*, which is gross of reserve requirements, and includes *all* institutions that hold reserve balances at the Federal Reserve Banks.⁶⁷

The notion of active excess reserves arises naturally in our theory, since reserve requirements determine incentives to hold reserves, and reserve balances at banks that are inactive in the fed funds market are inconsequential. However, we want to establish a mapping between our notion of active excess reserves and the notion of total reserves for two reasons. First, doing so will make our results easier to interpret, since the latter is a well known and readily available measure of aggregate reserves.⁶⁸ Second, for some of our quantitative exercises (e.g., the demand estimation illustrated in the top-right panel of Figure 18) we will want to overlay empirical observations for total reserves, which we may denote Q_t^D , on the theoretical demand for reserves, which is computed as a function of active excess reserves, which we may denote Q_t^M . For these two reasons, in Appendix D (Section D.2.5) we show how to “translate” the value of Q_t^D into a value of Q_t^M using a mapping that preserves the variation in the relevant sample $\{Q_t^D\}$ and that is consistent with the observed relationship between the sample mean of $\{Q_t^D\}$ and the sample mean of $\{Q_t^M\}$ in the two base years that we use to derive the theoretical demand for reserves (i.e., 2017 and 2019). Whenever a figure shows active excess reserves on the primary horizontal axis, we often include a secondary horizontal axis (above the figure) that shows the corresponding values for total reserves, to facilitate the translation between these units.⁶⁹

while the year 2019 has the lowest level of total reserves in the post-GFC-regulation era (roughly \$1,568.26 bn).

⁶⁷Fed funds transfers approximately sum to zero in our sample of active banks, so the fact that the beginning-of-day distributions that we feed the theory are net of predictable transfers does not contribute much to the difference between *total reserves* and *active excess reserves*.

⁶⁸E.g., Total reserves at weekly frequency is published in *Federal Reserve Balance Sheet: Factors Affecting Reserve Balances - H.4.1* (shown in Figure 16), and available at monthly frequency as “TOTRESNS” at <https://fred.stlouisfed.org>.

⁶⁹See, e.g., Figure 19, Figure 21, and the right panels of Figure 18. The average quantity of active excess reserves was about \$1,150.86 bn in 2017, and \$910.73 bn in 2019. Thus, by varying ω on $[0, 1]$ and using (9), our

In Section 6.1 we study how changes in structural parameters affect the position and shape of the theoretical aggregate demand for reserves. In Section 6.2 we use these insights and our quantitative theory to tackle the well-known empirical challenges involved in obtaining global estimates of the slope of the aggregate demand for reserves.

6.1 Reserve demand counterfactuals

In this section we study how the theoretical aggregate demand responds to changes in the administered rates and key marketstructure parameters. The results are reported in Figure 17.

In all panels of Figure 17, the curve labeled “Benchmark” is the theoretical aggregate demand $\iota_{Y_\omega}^* = \mathcal{D}(Q_{Y_\omega}; \Pi)$ for the model calibrated as in Table 1, and with $\iota_{Y_\omega}^*$ and Q_{Y_ω} computed with the interpolation procedure described in Section 3.6, for $Y_0 = 2017$ and $Y_1 = 2019$. We wish to make two observations about this demand for reserves generated by our theory. First, it exhibits the kind of logistic sigmoid shape that is characteristic of the popular “Poole model” (see, e.g., p. 784 in Poole (1968)). Second, for the baseline calibration, the demand lies within the DWR-IOR corridor. This means that despite there being GSEs that earn a lower interest on reserves than banks (i.e., the ONRRP rather than the IOR), the average equilibrium fed funds rate is above the IOR for all levels of reserves in the baseline calibration.⁷⁰

The top-left panel of Figure 17 shows two experiments. In the first, the DWR is increased by 50 bps (so that it is equal to the ONRRP plus 125 bps, rather than equal to the ONRRP plus 75 bps as in the baseline calibration). This shifts the demand up, with the size of the shift being decreasing in the quantity of reserves. Intuitively, the DWR has little effect on the equilibrium average interest rate when reserves are abundant, but a stronger effect when reserves are scarce. The second experiment consists of increasing the IOR by 15 bps (so that it is equal to the ONRRP plus 25 bps, rather than equal to the ONRRP plus 10 bps as in the baseline calibration). This policy change increases the equilibrium average rate when the quantity of reserves is relatively large, and it also implies that—if reserves are abundant enough—the equilibrium average fed funds rate lies between the IOR and the ONRRP.⁷¹ The

choice of endpoints ($Y_0 = 2017$ and $Y_1 = 2019$) allow us to interpolate any level of active excess reserves between \$1,150.86 bn and \$910.73 bn. The corresponding quantities of total reserves for 2017 and 2019 are \$2,254.27 bn and \$1,568.26 bn, respectively.

⁷⁰This observation is in line with the data, since the EFFR was consistently above the IOR during the 2019 sample period that we used as baseline for our calibration (see Figure 16).

⁷¹This observation is in line with the data, since the EFFR was consistently between the IOR and the ONRRP during most of the post-GFC period ranging from 2008 until 2018 when, as in this experiment, the IOR was set 25 bps above the ONRRP.

top-right panel of Figure 17 shows that increasing all administered rates (DWR, IOR, and ONRRP) by 75 bps simply causes a parallel upward shift in the aggregate demand for reserves.

The bottom-left panel of Figure 17 shows three experiments. The first, is to multiply the trading probabilities of all bank types by a factor of 10, which makes the marketstructure more competitive (i.e., “less OTC”), reducing rates (and increasing the slope of the demand for reserves) when the quantity of reserves is low to moderate. The second marketstructure experiment is to set $\beta_F = 0$, which effectively excludes all banks of type F from fed funds trading, and causes the aggregate demand for reserves to rotate clockwise around an intermediate quantity of aggregate reserves (about \$700 bn). This experiment causes the average fed funds rate to rise for relatively low levels of reserves, and to fall for relatively high levels of reserves. This rotation reflects the intermediation role that banks of type F play in the equilibrium: When reserves are scarce there are many banks with deficient reserve balances who, absent type- F counterparties, find it more difficult to meet a counterparty eager to lend, which reduces their effective market power thus leading to higher negotiated loan rates on average. When reserves are abundant there are many banks with excess reserves who, absent type- F counterparties, find it more difficult to meet a counterparty eager to borrow, thus leading to lower average negotiated rates. The third experiment is to eliminate the proportional borrowing costs from the baseline calibration. This shifts up the aggregate demand for reserves, reflecting that the borrowing costs stifle individual banks’ incentives to borrow.

The bottom-right panel of Figure 17 shows two experiments involving payment risk: One where we eliminate payment shocks for all banks, i.e., we set $\lambda_i = 0$ for all $i \in \mathbb{N}$, and another where we set $\lambda_i = \lambda_F$ for all $i \in \mathbb{N}$, i.e., we assume all bank types experience the same—very high—frequency of payment shocks as banks of type F . In both cases the result is an upward shift in the demand for reserves. In the second experiment the demand shifts up due to a heightened precautionary motive for holding reserves. In the first experiment the upward shift occurs because of a compositional effect: In an equilibrium with payment shocks, there are banks that borrow because their balances are deficient, and banks with moderate positive reserves that borrow to self-insure against payment shocks. The former values reserves more, and thus are willing to pay higher rates than the latter. The precautionary motive for borrowing disappears when $\lambda_i = 0$ for all $i \in \mathbb{N}$, and therefore the average negotiated rate increases.⁷²

⁷²The fact that the size of the upward shift that results from setting $\lambda_i = 0$ for all $i \in \mathbb{N}$ is decreasing in the quantity of reserves is consistent with this intuition.

6.2 Quantitative-theoretic reserve demand estimation

As discussed in Section 1, the floor system that the Federal Reserve has chosen as operating framework for monetary policy implementation relies on the ability to ascertain what level of reserves is “ample enough” so that active management of the supply of reserves is not required to instrument the fed funds rate target. In other words, operating a floor system requires *global* estimates of the aggregate demand for reserves, and in particular, reliable estimates of its slope for wide ranges of the aggregate supply of reserves. This presents two empirical challenges.

The first challenge is the potential endogeneity of the supply of reserves, which complicates the estimation of the demand equation. In terms of the simple demand-and-supply picture in the first panel of Figure 1, the issue is to identify *exogenous* variation in the quantity of reserves that allow to estimate the slope of the demand. This problem is well-understood, and has been addressed by the empirical literature that studies the *liquidity effect*.⁷³

The second challenge is to obtain *global* estimates for the slope of the demand; i.e., to identify the slope of the demand for a range of values of the supply of reserves that is wide enough to span the “abundant”, “ample”, and “scarce” segments of the demand curve, as illustrated in the top-right panel of Figure 1. The issue is that, empirically, spanning substantial variation in the supply of reserves usually entails spanning a substantial period of time during which the demand for reserves itself is likely to have shifted due to structural changes, e.g., in the marketstructure of the fed funds market, or in banks’ incentives to hold reserves (due to changes in policy, regulation, or portfolio allocation frameworks within banks).⁷⁴

This low-frequency demand-shift identification problem has not been overcome by the empirical literature on *liquidity effects*—possibly due to limited theoretical guidance on the key structural variables that determine the shape and position of the aggregate demand for reserves. Available empirical estimates of liquidity effects tend to be *local*, i.e., estimated from daily time-series variation in the quantity of reserves over relatively short sample periods during

⁷³We discussed these identification issues in Section 3.5, where we also reported estimates of the slope of the reserve demand for different sample periods based on the identification strategies of Hamilton (1997), Carpenter and Demiralp (2006), and Afonso et al. (2022).

⁷⁴These are the kinds of shifts in the demand for reserves that we studied in Section 6.1. The bottom panels of Figure 1 show situations in which structural parameters are Π_i at the time the quantity-price pair (Q_i, r_i^*) is observed, for $i \in \{0, 1\}$. The problem is that, without controlling for the structural change from Π_0 to Π_1 , other considerations may lead one to assume the observations $\{(Q_i, r_i^*)\}_{i \in \{0, 1\}}$ lie on a single demand curve, and therefore overestimate (in the example in the bottom-left panel) or underestimate (in the example in the bottom-right panel) the (absolute value of the) slope.

which the average quantity of reserves remains relatively stable.⁷⁵ We will use our quantitative model to bridge this *local-global gap*. The idea is to use the structure imposed by the theory, i.e., the equilibrium *aggregate demand relationship*, $\iota^* = \mathcal{D}(Q; \Pi)$, with the microstructure and policy parameters, Π , calibrated to match the key micro-level and market-level moments that describe the fed funds market—and in particular the available *local* estimates of the liquidity effect in the base year—to estimate the *global* shape of the aggregate demand for reserves.⁷⁶

To motivate and illustrate our quantitative-theoretic identification approach, consider Figure 18. The top-left panel displays pairs of empirical observations of the total quantity of reserves, and the corresponding EFFR-IOR spread for every trading day in the sample period 2017/01/20–2019/09/13. Through the lens of standard theory (e.g., Poole (1968)), each of these observations depicts the intersection point of the supply and demand for reserves on a given day. To inform monetary policy operations, one needs to estimate the liquidity effect for each level of reserves over a wide range of reserves, which can be done by estimating an aggregate demand for reserves. A natural approach is to posit a flexible reduced-form model of the demand for reserves, e.g., $s_t = D(Q_t)$, where s_t denotes the EFFR-IOR spread on day t and Q_t denotes the aggregate quantity of reserves at the end of day t , with

$$D(Q_t) \equiv \underline{s} + \frac{\bar{s} - \underline{s}}{1 + e^{(Q_t - Q_0)\xi}}, \quad (11)$$

and estimate the parameters $(\underline{s}, \bar{s}, \xi, Q_0)$. The top-left panel of Figure 18 displays the fitted demand curve that results from estimating (11) on the full sample (2017/01/20–2019/09/13) by nonlinear least squares (NLS).⁷⁷ This estimation presumes all observations in the sample lie on a single demand curve.⁷⁸ The estimated slope evaluated at the mean quantity of total

⁷⁵Hamilton (1997), Carpenter and Demiralp (2006), and our estimate of the coefficient γ in (5) are examples of this standard methodology. Afonso et al. (2022) follow an alternative methodology that involves estimating a time-varying vector autoregressive model at daily frequency (with an instrumental variable approach to address endogeneity of the supply of reserves) to obtain a 10-year time series of daily estimates of the elasticity of the fed funds rate to (instrumented) variation in the aggregate quantity of reserves (from 2010 until 2020). Their estimation, however, cannot recover the whole demand function. The reason is that without information on whether structural factors have shifted the demand schedule during the sample period, it is not possible to infer the global shape of the reserve demand from a sequence of (local, linear, daily) estimates of the sensitivity of the fed funds rate to (instrumented) changes in aggregate reserves. Having said this, below we will find that the reduced-form estimates from Afonso et al. (2022) can be a useful guide once complemented with our quantitative theory, which can help identify the structural shifts in the demand for reserves.

⁷⁶Alvarez and Argente (2023) use a similar strategy to extrapolate a demand for cash-paid Uber rides in Mexico using relatively narrow empirical variation in prices.

⁷⁷See Appendix D (Section D.2.4) for details.

⁷⁸In a similar estimation exercise, Afonso et al. (2022, Sec. 6) justify this particular identifying assumption by splitting their sample period (2010–2021/03/29) according to the different low-frequency cycles of expansion

reserves for the full sample (about \$1,974.69 bn) is -0.016 , which means a \$1 bn decrease in total reserves increases the EFR by 0.016 bps when the supply of total reserves is around \$2 tn. This local estimate (at about \$2 tn) is similar to the linear estimates in Section 3.5.

In Section 6.1 we showed that keeping the DWR-ONRRP spread constant, changes in the IOR-ONRRP spread shift the demand for reserves. This minimal theoretical insight implies that the data points from the sample period 2017/01/20–2019/09/13 plotted in the top-left panel of Figure 18 do not all lie on the same demand curve (contrary to what we implicitly assumed when running (11) on the full sample). The bottom-left panel of Figure 18 displays the same data points as the top-left panel, but partitioned into four subsamples, each determined by the size of the IOR-ONRRP spread: 10 bps (2019/05/02–2019/09/13), 15 bps (2018/12/20–2019/05/01), 20 bps (2018/06/14–2018/12/19), or 25 bps (2017/01/20–2018/06/13).⁷⁹ The bottom-left panel also displays the four fitted demand curves that result from estimating (11) on each subsample by nonlinear least-squares.

To illustrate the perils associated with the atheoretical demand estimation in the top-left panel of Figure 18, let’s focus on the demand estimation for the policy regime with IOR-ONRRP equal to 10 bps in the bottom-left panel, and highlight two discrepancies with the top-left panel. First, the liquidity effect at about \$1,974.69 bn (the mean of *total reserves* for the full sample) is -0.0001 bps; but it was estimated to be -0.016 bps in the full sample—much bigger in absolute value.⁸⁰ Second, suppose we want to use the estimated demand to identify the quantity of reserves that determines the end of the “ample” and the beginning of the “abundant” range for reserves, i.e., we want to estimate a quantity such as the Q_1 illustrated in the top-right panel of Figure 1. For practical purposes we adopt the convention that a supply of reserves, Q , is considered “abundant” if reducing Q by \$1 bn increases the EFR by no more than on hundredth of a basis point. Given this definition of “abundant”, the demand estimated for the subsample with IOR-ONRRP equal to 10 bps implies $Q_1 = \$1,300$ bn, while the demand estimated on the full sample implies $Q_1 = \$2,943$ bn. Discrepancies this large

and contraction of the Federal Reserve balance sheet. Specifically, they split it into three periods: the initial post-GFC expansionary period (2010–2014), the subsequent post-GFC and pre-COVID contractionary period (2015–2020/3/13), and the most recent post-COVID expansionary period (2020/03/16–2021/03/29). Thus, all the data points displayed in our Figure 18 belong to their pre-COVID contractionary period, which Afonso et al. (2022) fit with a single reduced-form demand curve (the gray curve in their Figure 9, p. 30), like we do in the top panel of Figure 18.

⁷⁹The DWR-ONRRP spread was constant (equal to 75 bps) throughout the full sample (see Figure 16).

⁸⁰The slope of the demand estimated for the subsample with IOR-ONRRP equal to 10 bps evaluated at the mean for the *subsample* (about \$1,521.48 bn of total reserves) is -0.0186 bps.

make one wary of relying on these kinds of estimations to guide monetary policy operations.

A shortcoming of the atheoretical reduced-form approach to estimating a global aggregate demand curve for reserves is that the extrapolations the empirical model makes for ranges of Q for which there are not many observations (e.g., very low values of Q) can be very sensitive to our ability to identify the structural parameters that shift the aggregate demand being estimated. It seems sensible to try to control for these “policy regimes”, and the sample split in the bottom-left panel of Figure 18 is an attempt to do so; but is this the right way to split the sample? Can variation in other policy or microstructure parameters shift or rotate the aggregate demand for reserves? To tackle these questions, we propose a quantitative theory-based approach.

The top-right panel of Figure 18 depicts several theoretical demands, $\iota_{y_\omega}^* = \mathcal{D}(Q_{y_\omega}; \Pi)$. The curve labeled “IOR-ONRRP = 10 bps” is the demand generated by the baseline calibration.⁸¹ The curves labeled “IOR-ONRRP = 15 bps”, “IOR-ONRRP = 20 bps”, and “IOR-ONRRP = 25 bps” are the theoretical demands corresponding to $\iota_r - \iota_o = 0.0015/360$, $\iota_r - \iota_o = 0.0020/360$, and $\iota_r - \iota_o = 0.0025/360$, respectively, with all other parameters as in the baseline calibration. The top-right panel of Figure 18 also displays pairs of empirical observations of the quantity of active excess reserves, and the corresponding EFFR-IOR spread for every trading day in the sample period 2017/01/20–2019/09/13. As before, the sample is partitioned into four subsamples, each determined by the size of the IOR-ONRRP spread: 10 bps (2019/05/02–2019/09/13), 15 bps (2018/12/20–2019/05/01), 20 bps (2018/06/14–2018/12/19), or 25 bps (2017/01/20–2018/06/13). The bottom-right panel of Figure 18 displays the same data points as the top-right panel, along with the four fitted demand curves that result from estimating (11) on each subsample by nonlinear least-squares.

There are two takeaways from comparing the top-right and bottom-right panels of Figure 18. First, the quantitative-theoretic demands fit the data reasonably well.⁸² Second, *locally*, i.e., for range of Q for which there is available data, the theoretical demand in the top-right panel and the reduced-form demands in the bottom-right panel fit about as well.⁸³ However, as

⁸¹That is, with the parameter values reported in Table 1, and $\iota_{y_\omega}^*$ and Q_{y_ω} computed with the interpolation procedure described in Section 3.6, for $y_0 = 2017$ and $y_1 = 2019$.

⁸²The height and slope of the demand curve labeled “IOR-ONRRP = 10 bps” were calibrated to match the average EFFR-IOR spread and the local liquidity effect for the corresponding subsample, but the other subsamples were not targeted. The theoretical demand labeled “IOR-ONRRP = 25 bps” predicts an EFFR-IOR spread that is somewhat high, but it only takes a 2 bp reduction in the liquidity return parameter, ι_ℓ , to bring the theoretical EFFR-IOR spread in line with the data.

⁸³In terms of local fit, the reduced-form specification is, as expected, no worse than the theory since it is more flexible. E.g., it allows us to choose *four* parameters, i.e., $(\underline{s}, \bar{s}, \xi, Q_0)$, to match the data corresponding to each

is evident from the figure, their predictions for lower levels of Q are quite different. To illustrate this point, focus on the subsample with IOR-ONRRP spread equal to 10 bps. The reduced-form model in the bottom-right panel estimates the steepest point on the corresponding demand at about \$934 bn (or about \$1,637 bn of total reserves), and predicts that reducing the supply of reserves below \$800 (or about \$1,255 bn of total reserves) would have essentially no effect on the equilibrium EFFR-IOR spread.⁸⁴ In contrast, our theory estimates the steepest point on the demand at about \$500 bn (or about \$662 bn of total reserves), and predicts that reductions in the supply of reserves start to cause significant increases in the EFFR-IOR spread for levels of reserves roughly below \$700 (below \$1,064 bn of total reserves). For the reduced-form approach, the extrapolation to out-of-sample levels of Q is essentially driven by the assumed functional form. In contrast, the theoretical extrapolation is based on the explicit equilibrium borrowing-and-lending activity that underlie the equilibrium *aggregate demand relationship*, $\iota^* = \mathcal{D}(Q; \Pi)$, with the microstructure and policy parameters, Π , calibrated to match the key micro-level and market-level moments that describe the fed funds market (as documented in Section 3).

6.2.1 Lessons for the practice of estimating reserve demands

We draw two conclusions from the quantitative-theoretic and the reduced-form estimations in this section. In Appendix E we show that these conclusions generalize to a wider range of reduced-form estimation strategies.

First, our theory identifies a set of structural “shifters” of the aggregate demand relationship that can help with the identification problems that pervade all reduced-form econometric estimations of the aggregate demand for reserves. For example, the theoretical counterfactuals in Section 6.1 show that the set of shifters include: the widths of the spreads between the administered policy rates; the parameters that regulate the trading frequencies and bargaining powers of the different types of fed funds participants; the bank-level idiosyncratic payment-shock processes; and balance-sheet borrowing costs induced by regulation.

Second, our quantitative-theoretic approach delivers estimates of the demand for reserves

subsample, while the theoretical demands corresponding to each subsample are generated by changing only one parameter, i.e., the policy spread $\iota_r - \iota_o$.

⁸⁴The reduced-form demand curves estimated using *active excess reserves* (reported in the bottom-right panel of Figure 18) are essentially identical to the ones estimated using *total reserves* (reported in the bottom-left panel). For example, the steepest point on the reduced-form demand estimated on the subsample with IOR-ONRRP spread equal to 10 bps in the bottom-left panel is at \$1,637 bn of total reserves; the slope at that point is -0.0002 bps, which is the same slope that the reduced-form demand estimated using active excess reserves achieves at \$934 bn.

that fit available data as well as the reduced-form approaches, but these approaches have very different out-of-sample predictions. In other words, all the estimated demand relationships are similar *locally*, i.e., for the range of reserve balances that have been observed since 2017, but are very different *globally*, i.e., for levels of reserves that the fed funds market has not visited in the past fifteen years.⁸⁵

Since the global estimate of the aggregate reserve demand based on our quantitative-theoretic approach is different from the global estimates based on reduced-form econometric approaches, and in turn the global estimates based on different but seemingly reasonable reduced-form econometric specifications are themselves different, a natural question arises: Which estimation approach should we favor?

We think our quantitative-theoretic approach has a clear advantage over the reduced-form econometric approaches whenever the global estimation entails large extrapolations from observed data. The advantage is that in our quantitative-theoretic approach, the global (out-of-sample) shape of the reserve demand is determined by the choice of “deep” microstructure parameters that can be disciplined with micro data.⁸⁶ In contrast, the out-of-sample shape of the reserve demand from reduced-form econometric specifications depends on the ad hoc functional-form specification that is assumed. And moreover, as we show in Appendix E, very reasonable specifications give very different out-of-sample predictions.

Another advantage of our approach is that the micro-structural foundations for the aggregate reserve demand allow us to run counterfactuals. For example, experiments involving changes in policy parameters, such as the spreads between the administered policy rates, or the regulatory costs of leverage. Or experiments involving changes in marketstructure parameters, such as those that regulate the trading frequencies and bargaining powers of the different types of fed funds participants, or changes in the bank-level idiosyncratic payment-shock processes.

⁸⁵For example, the slope of the quantitative-theoretic demand becomes virtually flat for total reserves in excess of \$1.3 tn, while the slope implied by reasonable reduced-form econometric estimates remains positive even for total reserves as large as \$2.5 tn. For relatively low levels of reserves, the model-generated demand becomes quite steep at about \$600 bn of total reserves and flattens for levels lower than \$340 bn. In contrast, the slope of implied by reasonable reduced-form estimates increases (often exponentially) as total reserves decrease, and becomes unreasonably large at low levels of reserves (e.g., at pre-GFC levels). See Appendix E for details.

⁸⁶As discussed in Section 4, our way of disciplining the shape of the demand is to calibrate the microstructure parameters so that the model replicates a wide array of high-frequency micro-level loan and payment data from Fedwire. As an additional source of validation (discussed in Section 5), recall that the calibrated model is also consistent with micro-level empirical observations not targeted in the calibration, such as the cross-sectional distribution of bilateral interest rates, the distribution of bid-ask spreads, and the intraday flow of reserves and supporting interest rates between pairs of banks in different positions on the trading network.

7 Navigational instruments for central banks

In this section we propose two diagnostic tools, or “navigational instruments” to aid monetary policy operations: (i) the *Monetary Confidence Band* (MCB), and (ii) the theory-based cross-sectional distribution of banks’ shadow cost of procuring funding in the fed funds market.

7.1 Confidence bands for monetary policy implementation

Monetary policy implementation changed drastically during the last decade as the supply of reserve balances increased to unprecedented levels, effectively turning the Fed’s operating framework into a *floor system*. There are unanswered operational questions about this system. The most elementary is: what is the smallest quantity of outstanding aggregate reserves needed to ensure that plausible market shocks do not cause significant deviations of the fed funds rate from its policy target? In this section we use the estimated quantitative theory to frame our answer in terms of a new policy-evaluation instrument: the *Monetary Confidence Band* (MCB).

Let $\iota = \mathcal{D}(Q)$ denote an aggregate-demand relationship between the equilibrium fed funds rate, ι , and the aggregate supply of reserves, Q . The mapping $\mathcal{D}(\cdot)$ could be obtained from the equilibrium of the theory, as in Section 6.1, or from another procedure (e.g., by estimating something like (11)). Let Z_p denote the p^{th} percentile of the empirical distribution of reserve-draining shocks estimated in Section 3.4. We define the “ $p\%$ MCB” as a pair of functions, $(\underline{\iota}(Q), \bar{\iota}(Q))$ with $\underline{\iota}(Q) \equiv \mathcal{D}\left(Q + Z_{\frac{100+p}{2}}\right)$ and $\bar{\iota}(Q) \equiv \mathcal{D}\left(Q + Z_{\frac{100-p}{2}}\right)$. The idea is that the reserve-augmenting or reserve-draining shocks induce randomness in the supply of reserves, which in turn induces randomness in the fed funds rate. For example, for a given beginning-of-day supply of reserves, Q , the equilibrium fed funds rate lies inside the 95% MCB, $(\mathcal{D}(Q + Z_{97.5}), \mathcal{D}(Q + Z_{2.5}))$, with 95% probability. Figure 19 presents several examples of MCBs where $\mathcal{D}(\cdot)$ is the aggregate demand for reserves derived from our theory.

The top-left panel in Figure 19 displays the 95% and 99% MCBs around the aggregate demand corresponding to the baseline calibration, which is labeled “Mean (volume weighted) rate”. There are two ways to use the MCB. First, for a given beginning-of-day supply of reserves, we can use the MCB to estimate the probability that the fed funds rate will be in a certain range. For example, for a typical day in the sample period targeted by this baseline calibration, the beginning-of-day quantity of active excess reserves, Q , was about \$900 bn, and the IOR was 235 bps. Under these conditions, the MCB indicates that the Desk should be able

to implement any target rate in the range IOR-IOR+25bps with certainty. Second, for any target range for the fed funds rate, the MCB yields the minimum quantity of reserves needed to meet the target with a desired degree of confidence. For example, if the Desk wanted the fed funds rate to be within the IOR-IOR+25bps range with 95% confidence, it would have to supply the market at least about \$670 bn in beginning-of-day active excess reserves. A 99% degree of confidence would instead require at least about \$850 bn.

The other panels in Figure 19 report the MCBs for calibrations that differ in one parameter from the baseline calibration. The top-right panel sets $\beta_F = 0$ (the baseline has $\beta_F = 0.03$). This could be interpreted as a day in which all banks of type F withdraw from the fed funds market. Under these conditions, the Desk would have to supply the market at least about \$700 bn beginning-of-day active excess reserves to keep the fed funds rate within the IOR-IOR+25bps with 95% confidence (about \$30 bn more than in the baseline). The bottom-left panel increases the IOR by 15 bps (from ONRRP + 10 bps in the baseline, to ONRRP + 25 bps). Notice that in this case the Desk would have to make beginning-of-day active excess reserves very scarce—less than \$500 bn—to ensure the FFR is higher than the IOR with 95% confidence. This is in contrast with the baseline calibration, which guarantees the FFR will be higher than the IOR with certainty for any level of reserves.

The bottom-right panel assumes $u_i(a) = \iota_d \mathbb{I}_{\{a < 0\}}$ for all i , with $\iota_d = \frac{0.2}{800} \iota_w$ (the baseline has $u_i(a) = 0$ for all $(a, i) \in \mathbb{R} \times \mathbb{N}$). The parameter ι_d captures the regulatory, reputational, or other costs associated with running an *intraday overdraft* (defined as a negative intraday excess reserve balance), which gained notoriety after the spikes in money-market rates of September 2019.⁸⁷ The main takeaway from this exercise is that even modest costs of not meeting the LCR and Regulation-D thresholds *on an intraday basis* can cause a significant upward shift in the demand for reserves. For example, from the bottom-right panel of Figure 19 we see that, with a level of beginning-of-day active excess reserves of about \$800 bn, the Desk cannot keep the FFR within the IOR-IOR+25bps range with 99% confidence (but the Desk can do so if $\iota_d = 0$, as in the baseline of the top-left panel).

⁸⁷See, e.g., Copeland et al. (2021, Section 4). With no available evidence on the value of ι_d , for illustrative purposes, here we have chosen it so that a bank that incurs intraday overdraft for a whole trading day (composed of 800 model periods) suffers a per-dollar cost equal to 20% of the DWR. In Section 8 we use our theory augmented with $0 < \iota_d$ to rationalize the spikes in the EFR of September 16 and 17, 2019. In Appendix D (Section D.3) we give a more detailed account of the events that took place during September 13–20, 2019 (reserve-draining shocks, associated rate spikes, and ensuing policy interventions).

7.2 Distribution of shadow price of reserves

When analyzing commercial banks’ decisions to lend to households, corporations, or money-market participants, the fed funds rate is usually regarded as a measure of the (opportunity) cost of the loanable funds. The logic is that a bank that is long in reserves could lend in the fed funds market rather than to a client, and a bank that is short may borrow in the fed funds market and lend elsewhere. Thus, in a competitive marketstructure the (opportunity) cost of funds for *all* banks is summarized by a single statistic—the equilibrium fed funds rate. But in an over-the-counter marketstructure where loans are negotiated bilaterally and sequentially over time as in the actual fed funds market, each bank faces different borrowing and lending rates depending on their own and their counterparties’ characteristics, such as their reserve balance at the time of the trade, degree of market power (e.g., θ_i), ability to find counterparties (e.g., β_i), and regulatory treatment (e.g., the administered rates they earn for holding reserves or pay for overdrafts).

In a dynamic OTC marketstructure like the fed funds market, each participating bank of type $i \in \mathbb{N}$ with reserve balance $a \in \mathbb{R}$ at time t has its own opportunity cost, or “shadow price” of reserves, which is summarized by $\frac{\partial V_t^i(a)}{\partial a}$. In this context, at any point in time the opportunity cost of loanable funds is characterized by a whole cross-bank distribution rather than by a single number, which may be more or less representative of the majority of banks. Below, we show that according to our baseline calibration, neither the EFFR nor the distribution of traded rates are representative of the distribution of shadow prices of reserves of the majority of the banks that participate in the fed funds market.

While outside the scope of our model, one could envision a more general model in which banks make lending decisions to outside clients at a first stage knowing they will later participate of a fed funds trading stage like the one we have modeled above.⁸⁸ In this setup, the relevant opportunity cost of loanable funds in the first stage for a bank of type i is given by $\mu_i(a) \equiv \frac{\partial V_0^i(a)}{\partial a} - 1$, where a is the bank’s residual balance after having made loans to outside clients in the first stage. We summarize this heterogeneity with a cumulative distribution function $\mathcal{M}_i(\iota) \equiv \int \mathbb{I}_{\{a: \mu_i(a) \leq \iota\}} dF_0^i(a)$, i.e., $\mathcal{M}_i(\iota)$ is the proportion of banks of type $i \in \mathbb{N}$ whose shadow

⁸⁸This setup would be a natural way to incorporate a repo market into the theory, since the majority of repo transactions are executed early in the business day. Copeland et al. (2021), for example, report that a large fraction of interdealer repo trades are conducted between 7:00 am and 7:20 am, EST, and use this fact to argue that when intermediating the Treasury repo market, the marginal value to a dealer bank of holding balances at the Fed is sensitive to anticipated intraday payment stresses on these balances.

price of reserves at the beginning of the fed funds market trading day is lower than $\iota \in \mathbb{R}$.

The top-left panel of Figure 20 shows

$$\mathcal{M}(\iota) \equiv \sum_{i \in \{F, M, S\}} \frac{n_i \mathcal{M}_i(\iota)}{\sum_{i \in \{F, M, S\}} n_i},$$

along with the cumulative distribution function of *all* bilateral loan rates negotiated throughout the day, denoted \mathcal{H} (both calculated for the baseline calibration). Intuitively, $\mathcal{H}(\iota)$ is the proportion of reserves traded at rates below ι . The dashed vertical line labeled “EFFR” denotes the volume-weighted average fed funds rate on *all* trades implied by the theory. The IOR and DWR are denoted by solid vertical lines. Notice that \mathcal{H} is very concentrated around the EFFR (about 60% of the funds are traded at the EFFR), so although there is heterogeneity in negotiated loan rates, the EFFR is quite representative of the overall distribution of traded rates. On the other hand, neither the distribution of traded rates nor the EFFR are representative of the distribution of shadow prices of reserves across all banks, represented by \mathcal{M} . For example, 80% of banks have a shadow price of reserves higher than the EFFR, but only about 10% of reserves are traded at rates higher the EFFR. The reason is that banks of type S , which constitute more than 90% of the population of banks, account for a small share of trades, and are therefore underrepresented in the statistics computed on actual trades, such as the EFFR and the distribution \mathcal{H} .

The remaining three panels of Figure 20 display the beginning-of-day cumulative distribution function of shadow prices for banks of type i , denoted \mathcal{M}_i , and the cumulative distribution function of all loan rates paid or received by banks of type i , denoted \mathcal{H}_i . These panels show that the EFFR and the distribution of traded rates, \mathcal{H}_i , are fairly representative of the distribution of shadow prices of reserves across banks, \mathcal{M}_i , only for types $i \in \{F, M\}$, but not for type S . This means that for about 90% of banks that participate in the fed funds market, the EFFR does not adequately capture the shadow cost of procuring funding, and is therefore not the relevant cost of lending in the retail and corporate loan markets.

8 Tuesday, September 17, 2019

On Tuesday September 17, the EFFR printed at 230 bps, exceeding the upper limit of the FOMC’s target range by 5 bps.⁸⁹ This event garnered the attention of market analysts and

⁸⁹The 99th percentile of the distribution of fed funds rates reached about 400 bps on September 17. Repo markets also experienced rate spikes, e.g., the secured overnight financing rate (SOFR) printed at 243 bps on

policymakers for two reasons. First, it was the first upward deviation from target in the 11 years since the FOMC began announcing a target range for the EFFR in December 2008. Second, it seemed inconsistent with the widespread view that the \$1.3 tn of reserves in excess of Regulation D outstanding at the time ought to be “ample enough” to run a floor system in which the Federal Reserve can implement its EFFR target without having to micromanage the supply of reserves.

To frame the discussion, consider the top panel of Figure 21, which displays a data scatterplot with the EFFR-IOR spread on the vertical axis (in percent per annum), and the quantity of reserves on the horizontal axis (in billions of dollars).⁹⁰ The data points labeled “IOR-ONRRP = 10 bps” are all trading days in the sample 2019/05/02–2019/09/13 (the period we used to estimate the liquidity effect in Section 3.5). The six darkest data points labeled “Sept 13-20 2019” are September 13, 16, 17, 18, 19, and 20. The dashed lines labeled “Target Upper Limit” and “Target Lower Limit” are the top and bottom of the fed funds target range minus the IOR for the period 2019/05/02–2019/09/18. On the scatterplot we have overlaid the MCB implied by the baseline calibration of the model (this is the same MCB displayed in the top-left panel of Figure 19, but this time with the EFFR-IOR spread on the vertical axis).

Friday, September 13 is the dark dot that sits on the demand for reserves generated by the theory—well within the EFFR target range. Monday, September 16 is the rightmost dark dot that sits on the upper limit of the target range for the EFFR-IOR spread, and September 17 is the uppermost dark dot, with an EFFR-IOR spread of 20 bps (5 bps higher than the spread between the upper limit of the EFFR target range and the IOR). Wednesday, September 18 is the leftmost dark dot that sits on the upper limit of the target range for the EFFR-IOR spread.⁹¹ The most cited culprits for the rate spikes of September 16 and 17 are two anticipated reserve-draining shocks that reduced the supply of reserves by about \$120 bn over two business days.⁹² From the top panel of Figure 21, we see that the EFFR-IOR spreads for September 16–18

Monday September 16 (13 bps higher than the previous business day), and exceeded 500 bps on September 17. See Afonso et al. (2020a) and Anbil et al. (2020) for detailed accounts of these money-market events.

⁹⁰As in Figure 19, the primary horizontal axis represents *active excess reserves* (as defined in Section 6), and the secondary horizontal axis translates them into *total reserves* (as explained in Section D.2.5 of Appendix D).

⁹¹September 19 and 20 are the dark dots with an EFFR-IOR spread of 10 bps.

⁹²The first was a quarterly corporate tax payment transferred from corporations’ bank and money market mutual fund accounts to the Treasury’s account. The second, a \$54 bn settlement of Treasury debt paid by primary dealers into the Treasury’s account on September 16. In Section D.3 we give a more detailed account of the reserve-draining shocks, associated rate spikes, and ensuing policy interventions that took place during September 13–20, 2019. Table 2 summarizes the main facts.

lie outside the 99% MCB. This means that (under our baseline calibration) our quantitative model cannot rationalize these observations as resulting from a “typical” daily reserve-draining shock—even if we define a “typical” shock as one with probability larger than 1%.

These events raised several questions: In a context with \$1.3 tn of excess reserves in the banking system, how could an anticipated \$120 bn reserve-draining shock cause such large spikes in money-market rates? Why didn’t banks lend some of their excess reserves to exploit the high overnight rates? In response to these questions during an earnings call on October 15, 2019, Jamie Dimon (Chairman and CEO of JPMorgan Chase) famously alluded to internal reserve management practices to ensure compliance with liquidity regulations:

As I said, we have \$120 bn in our checking account at the Fed, and it goes down to \$60 bn and then back to \$120 bn during the average day. But we believe the requirement under CLAR (Comprehensive Liquidity Analysis and Review) and resolution and recovery is that we need enough in that account, so if there’s extreme stress during the course of the day, it doesn’t go below zero. If you go back to before the crisis, you’d go below zero all the time during the day. So the question is, how hard is that as a red line? Was the intent of regulators between CLAR and resolution to lock up that much of reserves in the account with Fed? And that’ll be up to regulators to decide. But right now, we have to meet those rules and we don’t want to violate anything we’ve told them we’re going to do.⁹³

To explore this hypothesis, the middle and bottom panels of Figure 21 overlay, on the same data scatterplot of the top panel, the MCB implied by the baseline calibration of the model, but with $u_i(a) = \iota_d a \mathbb{I}_{\{a < 0\}}$ for all i , with $\iota_d = \frac{x}{800} \iota_w$. The middle panel has $x = 0.1$, and the bottom panel, $x = 0.2$.⁹⁴ The parameter ι_d stands in for a bank’s perceived penalty from going below Dimon’s “red line” (e.g., associated to the possible loss of reputation with regulatory supervisors for failing to maintain prudent liquidity buffers, as suggested by Copeland et al.

⁹³For a full transcript of the call, see: <https://tinyurl.com/29scwszt>. There is other evidence that the introduction of post-GFC liquidity regulations and associated supervisory programs have changed banks’ liquidity risk management practices. Afonso et al. (2020a), for example, point to a recent survey conducted by the Federal Reserve in which the majority of bank respondents identified “meeting routine intraday payments flows and satisfying internal liquidity stress metrics as the main drivers of their demand for reserves”. See, e.g., the August 2019 Senior Financial Officer Survey, <https://www.federalreserve.gov/data/sfos/sfos.htm>.

⁹⁴The baseline calibration in the top panel corresponds to the special case with $x = 0$. The case with $x = 0.2$ was one of the counterfactual exercises considered in Section 7.1. Recall that $x = 0.1$, for example, implies a bank that incurs intraday overdraft for a whole trading day suffers a per-dollar cost equal to 10% of the DWR.

(2021)). The middle panel of Figure 21 then shows that a shadow cost of intraday overdraft equal to 10% of the DWR, e.g., caused by precautionary reserve internal management practices designed to ensure compliance with liquidity regulations, is enough for the model to rationalize the September 16-18 EFRR-IOR observations as resulting from “typical” daily reserve-draining shocks, in the sense of being within the 99% MCB. The bottom panel of Figure 21 shows that a shadow cost of intraday overdraft equal to 20% is enough to put these observations within the 95% MCB (marginally within, in the case of September 16).

In another segment of the JPM earnings call of October 15, 2019, Dimon also alluded to the LCR requirement as a relevant determinant of banks’ demand for reserves:

We have a checking account at the Fed with a certain amount of cash in it. That cash, we believe, is required under resolution and recovery and liquidity stress testing. And therefore, we could not redeploy it into repo market, which we would’ve been happy to do. [...] You’re also going to hit a red line in LCR, like HQLA, which cannot be redeployed either.

Throughout the paper we have defined a bank’s “excess reserves” as its reserve balance net of the Regulation D reserve requirement, and net of the minimum quantity of reserves necessary to meet the LCR requirement given the bank’s holdings of other qualifying HQLA.⁹⁵ In other words, our baseline calculation of excess reserves corresponds to a world in which banks have a preference for satisfying the LCR requirement with non-reserve assets to the extent possible.⁹⁶ In reality, however, there are anecdotal accounts that banks appear to have a preference for meeting LCR requirements with reserves rather than with other HQLA. By incorporating a very modest preference for complying with the LCR requirement with reserves, e.g., something that reduces our measure of aggregate beginning-of-day excess reserves by as little as \$50 bn, the baseline calibration is able to rationalize the events of September 16–17, 2019 (in the sense that the compliance preference would shift all the dark dots in the top panel of Figure 21 leftwise by \$50 bn, and into the 99% MCB).

⁹⁵See footnote 28 and the more comprehensive discussion in Section B.2.1.

⁹⁶We have adopted this identifying assumption for our baseline because we regard it the most conservative option in the sense that—even to this day—the common definition of “excess reserves” considers the pre-GFC Regulation D requirement but not the post-GFC LCR requirement.

9 Conclusion

In this paper we have taken several steps toward developing models of the fed funds market with explicit over-the-counter microstructures into useful tools to guide monetary policy operations. Our framework incorporates the main microstructure ingredients of the fed funds market, accounts for the most salient institutional features, and includes the collection of policy instruments and regulations that shape participants' demands for reserves. The model also incorporates the large degree of heterogeneity among participants across several dimensions, such as: market power in bilateral loans, frequency and size distribution of payment shocks, and degree of centrality in market-making.

We documented a comprehensive set of novel marketwide and micro-level observations that describe the market dynamics, and showed that the quantitative model is flexible enough to match these observations. We then used the quantitative theory to deliver structural estimates of the aggregate demand for reserves, and developed two policy instruments to assess the cross-bank inequality in the shadow cost of procuring funding, and the central bank's ability to implement a given fed funds target.

While we think we have made significant progress, we also realize we have touched upon some questions and ideas that would be worth studying further in future work. First, we have allowed for heterogeneity in contact rates across bank types to capture the core-periphery structure of the fed funds market, but we have treated these contact rates as parameters. While the exogeneity of contact rates may be a reasonable assumption during periods when regulation and the deeper marketstructure parameters are relatively constant, it is not difficult to imagine settings or questions where it would be desirable to endogenize search intensity, (e.g., perhaps along the lines of Farboodi et al. (2023)). A similar point can be made about the beginning-of-day distributions of reserves, which for many applications would be best derived from an explicit portfolio problem of banks that takes place *prior* to the fed-funds trading stage that we have focused on.

Finally, a monetary-policy operating framework consists of two parts: an *operating target* (e.g., the fed funds rate), and *policy instruments* (e.g., standing facilities, open-market operations). Monetary models in the macro tradition focus on the macroeconomic effects of choosing different values (or rules) for the operating target, and leave operational implementation considerations outside the scope of their analysis. Here we have instead focused on the operational

side of the monetary policymaking process, and have left macro considerations outside the scope of our analysis. We think that exploring the macroeconomic implications of the microstructure of interbank lending and payments is a promising avenue of research (examples of work along these lines include Arce et al. (2020), Bianchi and Bigio (2022), De Fiore et al. (2018), Li and Li (2021), and Piazzesi and Schneider (2018)).

References

- ÅBERG, P., M. CORSI, V. GROSSMANN-WIRTH, T. HUDEPOHL, Y. MUDDE, T. ROSOLIN, AND F. SCHOBERT (2021): “Demand for Central Bank Reserves and Monetary Policy Implementation Frameworks: The Case of the Eurosystem,” *European Central Bank Occasional Paper*.
- AFONSO, G., R. ARMENTER, AND B. LESTER (2019): “A Model of the Federal Funds Market: Yesterday, Today, and Tomorrow,” *Review of Economic Dynamics*, 33, 177–204.
- AFONSO, G., M. CIPRIANI, A. M. COPELAND, A. KOVNER, G. LA SPADA, AND A. MARTIN (2020a): “The Market Events of Mid-September 2019,” *FRB of New York Staff Report*.
- AFONSO, G., D. GIANNONE, G. LA SPADA, AND J. C. WILLIAMS (2022): “Scarce, Abundant, or Ample? A Time-Varying Model of the Reserve Demand Curve,” Working Paper 1019, Federal Reserve Bank of New York.
- AFONSO, G., K. KIM, A. MARTIN, E. NOSAL, S. POTTER, AND S. SCHULHOFER-WOHL (2020b): “Monetary Policy Implementation with an Ample Supply of Reserves,” *Finance and Economics Discussion Series 2020-020*. Washington: Board of Governors of the Federal Reserve System.
- AFONSO, G., A. KOVNER, AND A. SCHOAR (2011): “Stressed, Not Frozen: The Federal Funds Market in the Financial Crisis,” *Journal of Finance*, 66, 1109–1139.
- AFONSO, G. AND R. LAGOS (2012): “An Empirical Study of Trade Dynamics in the Fed Funds Market,” *FRB of New York Staff Report*.
- (2015a): “The Over-the-Counter Theory of the Fed Funds Market: A Primer,” *Journal of Money, Credit and Banking*, 47, 127–154.
- (2015b): “Trade Dynamics in the Market for Federal Funds,” *Econometrica*, 83, 263–313.
- ALVAREZ, F. AND D. ARGENTE (2023): “Consumer Surplus of Alternative Payment Methods: Paying Uber with Cash,” Manuscript, University of Chicago.

- ANBIL, S., A. G. ANDERSON, AND Z. SENYUZ (2020): “What Happened in Money Markets in September 2019?” *FEDS Notes*, 27.
- ARCE, O., G. NUNO, D. THALER, AND C. THOMAS (2020): “A Large Central Bank Balance Sheet? Floor vs Corridor Systems in a New Keynesian Environment,” *Journal of Monetary Economics*, 114, 350–367.
- ARMANTIER, O. AND A. M. COPELAND (2015): “Challenges in Identifying Interbank Loans,” *Federal Reserve Bank of New York, Economic Policy Review*, 1–17.
- ARMANTIER, O., E. GHYSELS, A. SARKAR, AND J. SHRADER (2015): “Discount Window Stigma During the 2007–2008 Financial Crisis,” *Journal of Financial Economics*, 118, 317–335.
- ARMENTER, R. AND B. LESTER (2017): “Excess Reserves and Monetary Policy Implementation,” *Review of Economic Dynamics*, 23, 212–235.
- ARTUÇ, E. AND S. DEMIRALP (2010): “Discount Window Borrowing After 2003: The Explicit Reduction in Implicit Costs,” *Journal of Banking & Finance*, 34, 825–833.
- ASHCRAFT, A. B. AND D. DUFFIE (2007): “Systemic Illiquidity in the Federal Funds Market,” *American Economic Review*, 97, 221–225.
- BASEL COMMITTEE ON BANKING SUPERVISION (2010): *Basel III: International Framework for Liquidity Risk Measurement, Standards and Monitoring*, Bank for International Settlements.
- BECH, M. L. AND E. ATALAY (2010): “The Topology of the Federal Funds Market,” *Physica A: Statistical Mechanics and its Applications*, 389, 5223–5246.
- BECH, M. L. AND E. KLEE (2011): “The Mechanics of a Graceful Exit: Interest on Reserves and Segmentation in the Federal Funds Market,” *Journal of Monetary Economics*, 58, 415–431.
- BELTRAN, D. O., V. BOLOTNYI, AND E. KLEE (2021): “The Federal Funds Network and Monetary Policy Transmission: Evidence from the 2007–2009 Financial Crisis,” *Journal of Monetary Economics*, 117, 187–202.

- BERNANKE, B. S. AND D. KOHN (2016): “The Fed’s Interest Payments to Banks,” *Brookings Institution blog*.
- BIANCHI, J. AND S. BIGIO (2022): “Banks, Liquidity Management, and Monetary Policy,” *Econometrica*, 90, 391–454.
- BOTEV, Z. I., J. F. GROTH, D. P. KROESE, ET AL. (2010): “Kernel Density Estimation via Diffusion,” *The Annals of Statistics*, 38, 2916–2957.
- CARPENTER, S. AND S. DEMIRALP (2006): “The Liquidity Effect in the Federal Funds Market: Evidence from Daily Open Market Operations,” *Journal of Money, Credit and Banking*, 38, 901–920.
- CHIU, J., J. EISENSCHMIDT, AND C. MONNET (2020): “Relationships in the Interbank Market,” *Review of Economic Dynamics*, 35, 170–191.
- COPELAND, A., D. DUFFIE, AND Y. YANG (2021): “Reserves Were Not So Ample After All,” Working Paper 29090, National Bureau of Economic Research.
- DE FIORE, F., M. HOEROVA, AND H. UHLIG (2018): “Money Markets, Collateral and Monetary Policy,” Working Paper 25319, National Bureau of Economic Research.
- DUFFIE, D., N. GÂRLEANU, AND L. H. PEDERSEN (2005): “Over-the-Counter Markets,” *Econometrica*, 73, 1815–1847.
- ENNIS, H. M. (2019): “Interventions in Markets with Adverse Selection: Implications for Discount Window Stigma,” *Journal of Money, Credit and Banking*, 51, 1737–1764.
- ENNIS, H. M. AND T. KEISTER (2008): “Understanding Monetary Policy Implementation,” *FRB Richmond Economic Quarterly*, 94, 235–263.
- ENNIS, H. M. AND J. A. WEINBERG (2013): “Over-the-Counter Loans, Adverse Selection, and Stigma in the Interbank Market,” *Review of Economic Dynamics*, 16, 601–616.
- FARBOODI, M., G. JAROSCH, AND R. SHIMER (2023): “The Emergence of Market Structure,” *Review of Economic Studies*, 90, 261–292.
- FEDERAL RESERVE BOARD (2019a): *Reserve Maintenance Manual*, Board of Governors of the Federal Reserve System.

- (2019b): “Statement Regarding Monetary Policy Implementation,” *Press Release, Washington, DC, October 11, 11:00 a.m. EDT.*
- (2019c): “Statement Regarding Monetary Policy Implementation and Balance Sheet Normalization,” *Press Release, Washington, DC, January 30, 2:00 p.m. EST.*
- FEINMAN, J. N. (1993): “Reserve Requirements: History, Current Practice, and Potential Reform,” *Fed. Res. Bull.*, 79, 569.
- FURFINE, C. H. (1999): “The Microstructure of the Federal Funds Market,” *Financial Markets, Institutions & Instruments*, 8, 24–44.
- HAMILTON, J. D. (1996): “The Daily Market for Federal Funds,” *Journal of Political Economy*, 104, 26–56.
- (1997): “Measuring the Liquidity Effect,” *American Economic Review*, 87, 80–97.
- HUGONNIER, J., B. LESTER, AND P.-O. WEILL (2020): “Frictional Intermediation in Over-the-Counter Markets,” *Review of Economic Studies*, 87, 1432–1469.
- IRELAND, P. (2018): “Fed Should Stop Paying Interest on Reserves,” *Economics 21 blog*.
- KEISTER, T. (2012): “Corridors and Floors in Monetary Policy,” *Liberty Street Economics*.
- KEISTER, T., A. MARTIN, AND J. MCANDREWS (2008): “Divorcing Money from Monetary Policy,” *Economic Policy Review*, 14.
- KLEE, E. ET AL. (2021): “The First Line of Defense: The Discount Window during the Early Stages of the Financial Crisis,” *International Journal of Central Banking*, 17, 143–190.
- LAGOS, R. AND G. ROCHETEAU (2007): “Search in Asset Markets: Market Structure, Liquidity, and Welfare,” *American Economic Review*, 97, 198–202.
- (2009): “Liquidity in Asset Markets With Search Frictions,” *Econometrica*, 77, 403–426.
- LAGOS, R., G. ROCHETEAU, AND P.-O. WEILL (2011): “Crises and Liquidity in Over-the-Counter Markets,” *Journal of Economic Theory*, 146, 2169–2205.
- LI, Y. AND Y. LI (2021): “Payment Risk and Bank Lending,” Working Paper 2021-03, Ohio State University Fisher College of Business.

- LÓPEZ-SALIDO, D. AND A. VISSING-JORGENSEN (2023): “Reserve Demand, Interest Rate Control, and Quantitative Tightening,” .
- PIAZZESI, M. AND M. SCHNEIDER (2018): “Payments, Credit and Asset Prices,” Manuscript, Stanford University.
- POOLE, W. (1968): “Commercial Bank Reserve Management in a Stochastic Model: Implications for Monetary Policy,” *The Journal of Finance*, 23, 769–791.
- ÜSLÜ, S. (2019): “Pricing and Liquidity in Decentralized Asset Markets,” *Econometrica*, 87, 2079–2140.
- WEILL, P.-O. (2007): “Leaning Against the Wind,” *Review of Economic Studies*, 74, 1329–1354.

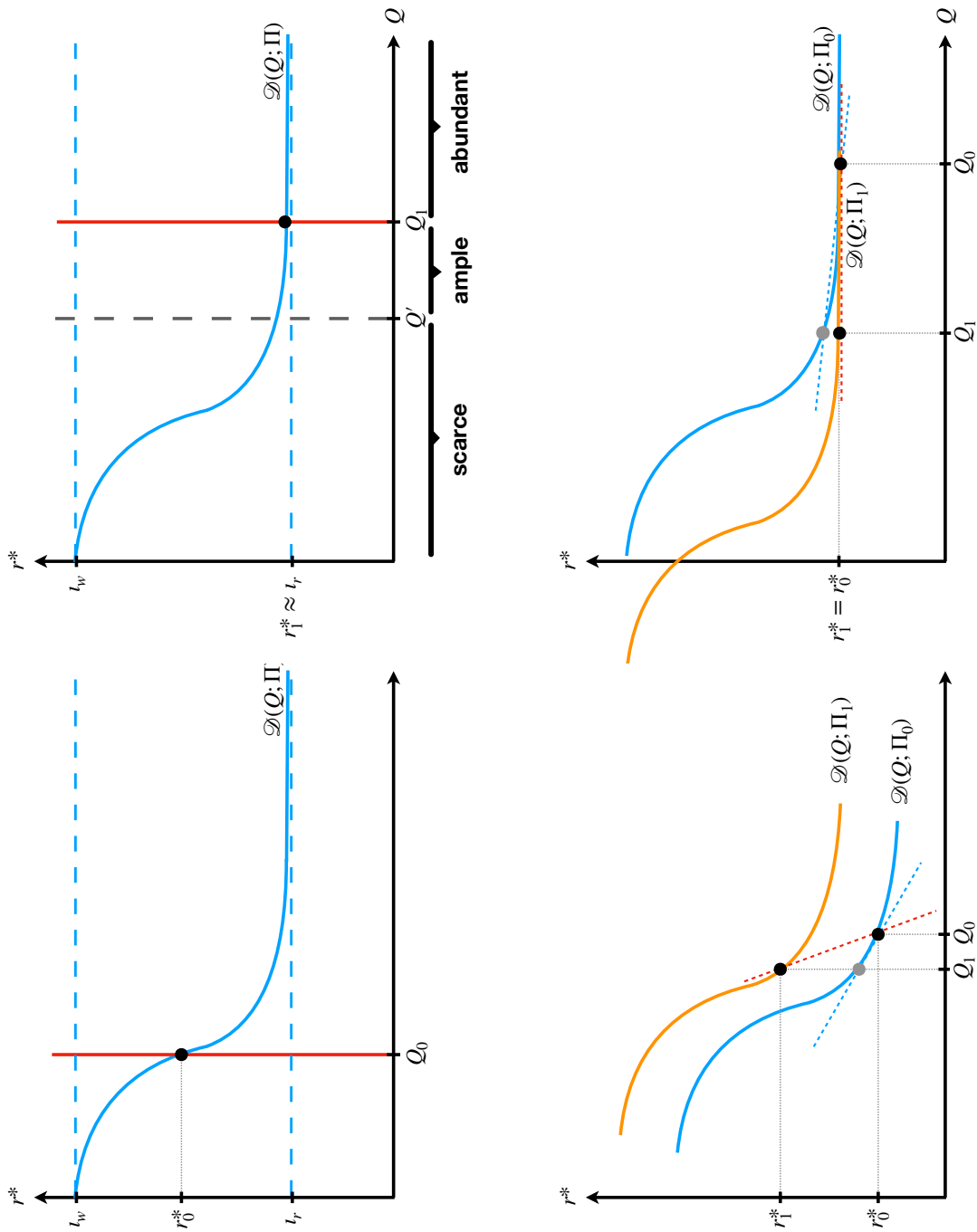


Figure 1: Stylized model of the determination of the fed funds rate.

Notes: In this figure, Q denotes the aggregate quantity of reserves, r^* the fed funds rate, and Π a set of parameters that determine the position of the aggregate demand for reserves, \mathcal{D} . The administered rates in the lending and deposit facility are denoted i_w and i_r , respectively.

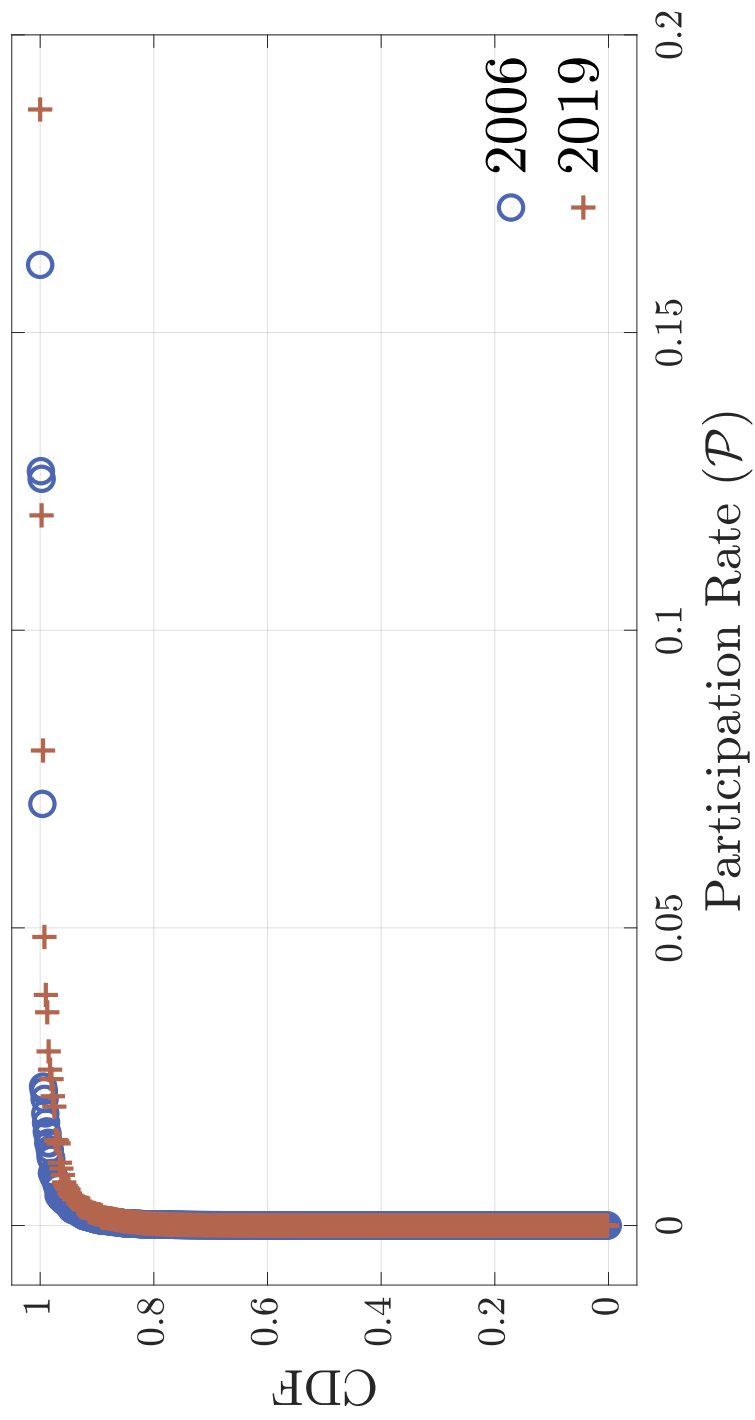


Figure 2: Empirical cumulative distribution function of bank-level participation rates for 2006 and 2019.

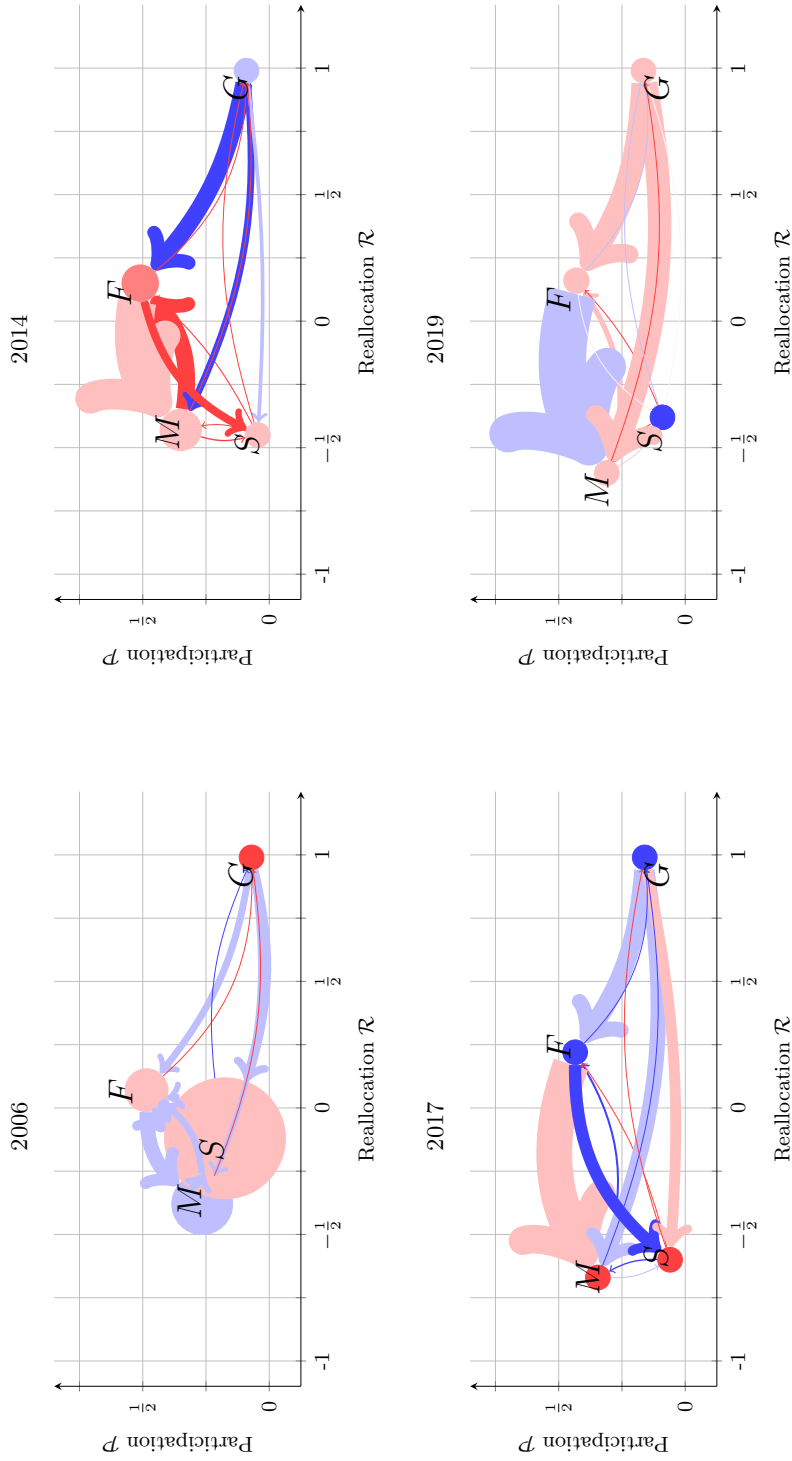


Figure 3: Fed funds trading networks.

Notes: Each node corresponds to one of the four bank types, and labeled accordingly as F , M , S , or G . An arrow from a type to another represents loans extended from banks in the former to banks in the latter, with the size of the arrow proportional to the volume of loans. The size of each node is proportional to the volume of loans extended by banks of that type to other banks of that same type. The colors of the arrows and nodes are: light blue, dark blue, light red, or dark red, if the volume-weighted average interest rate on the loans between the two types of banks, expressed as a spread over the effective fed funds rate (EFFR), falls in the first, second, third, or fourth quartile, respectively.

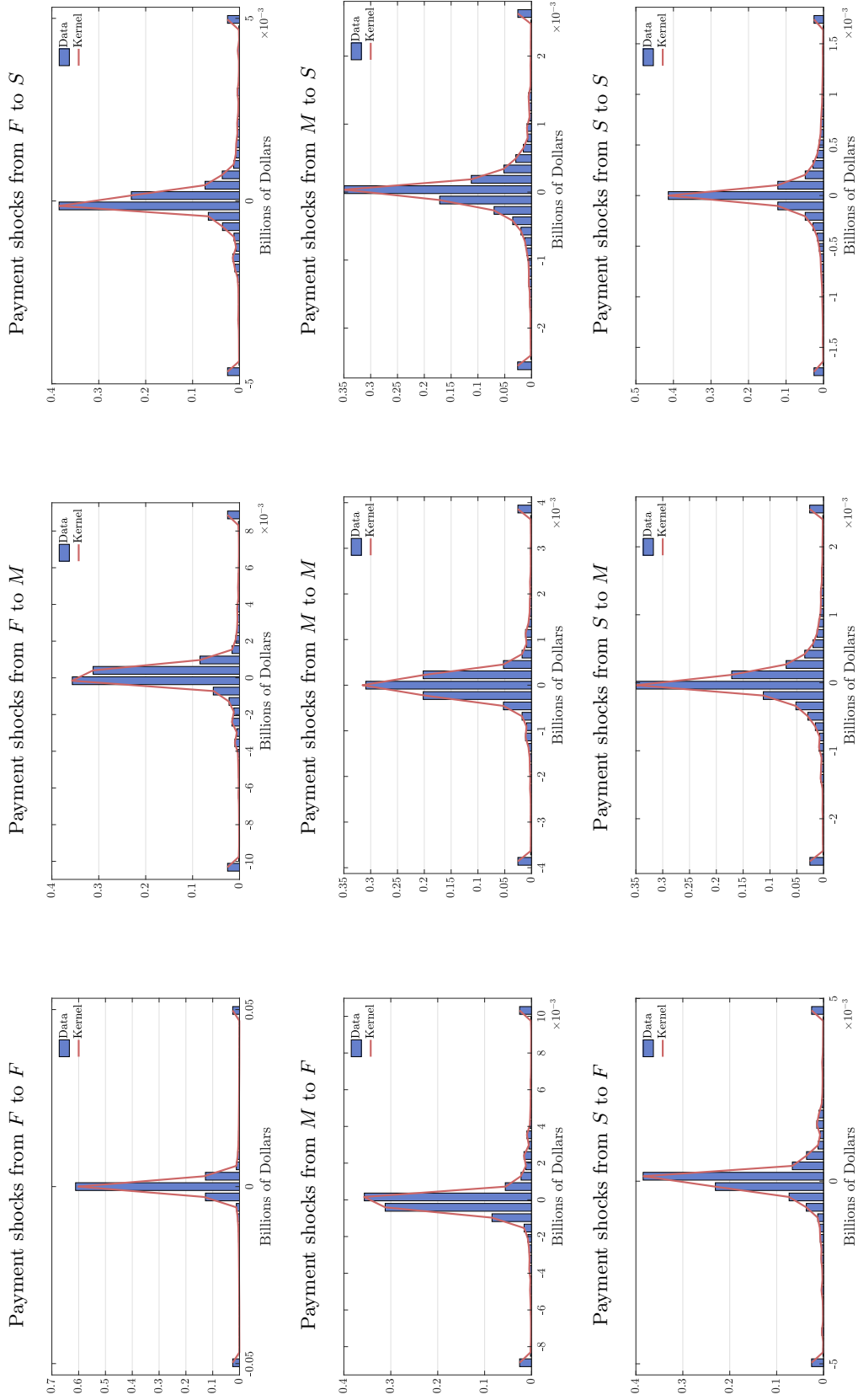


Figure 4: High-frequency payment shocks between pairs of bank types in 2006.

Notes: Empirical size distributions of high-frequency payment shocks between banks of each type (blue histograms) and their corresponding Gaussian kernel density estimates (red curves).

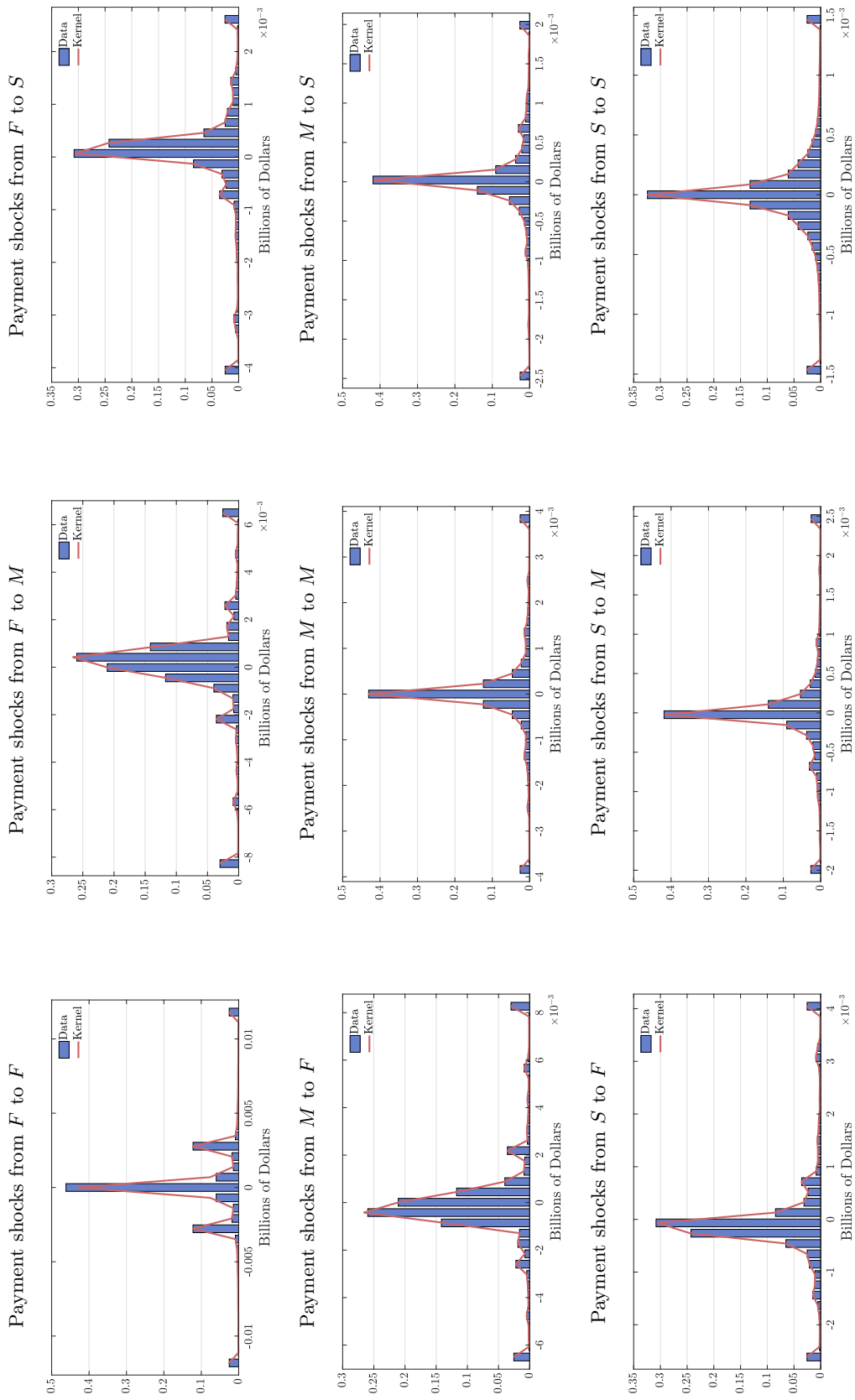


Figure 5: High-frequency payment shocks between pairs of bank types in 2019.

Notes: Empirical size distributions of high-frequency payment shocks between banks of each type (blue histograms) and their corresponding Gaussian kernel density estimates (red curves).

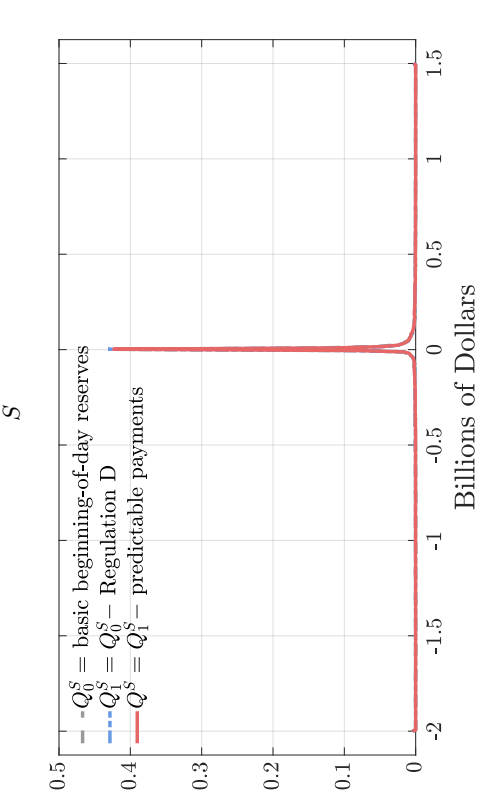
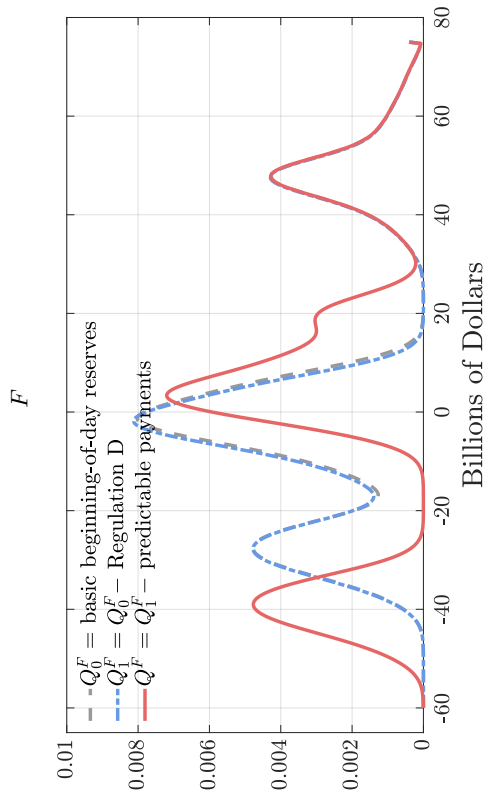
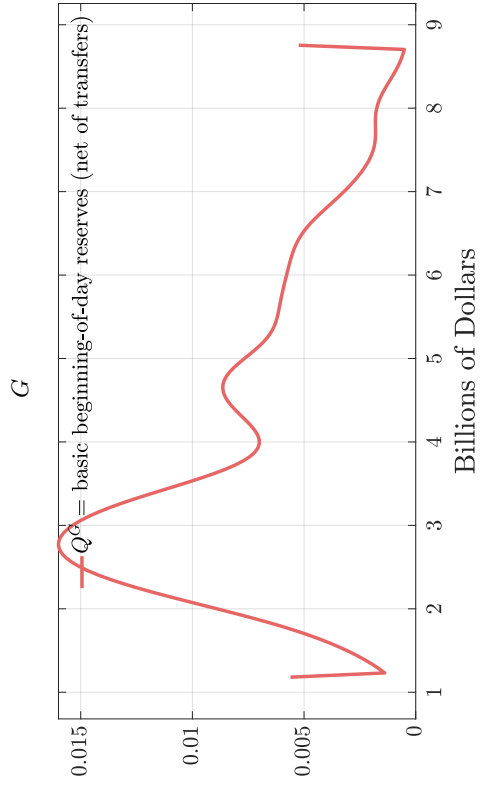
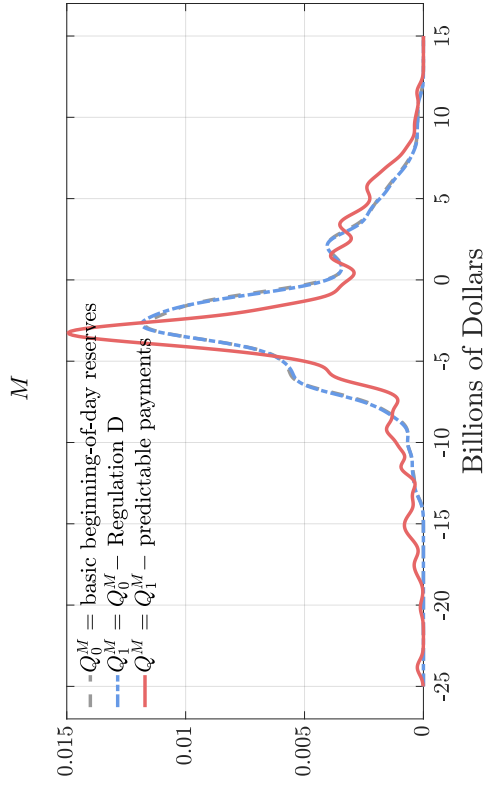


Figure 6: Estimated beginning-of-day distributions of reserves by bank type for the year 2006.

Notes: For every bank type $i \in \mathbb{N}$, the curve labeled Q^i is the (Gaussian kernel density) estimate of the empirical distribution of beginning-of-day excess reserves net of predictable transfers. For bank type $i \in \{F, M, S\}$, the curve labeled Q_0^i is the estimate of the distribution of “basic” beginning-of-day reserves (net of fed funds repayments), and the curve labeled Q_1^i is the estimate of the distribution of “basic” beginning-of-day reserves net of the Regulation D reserve requirement.

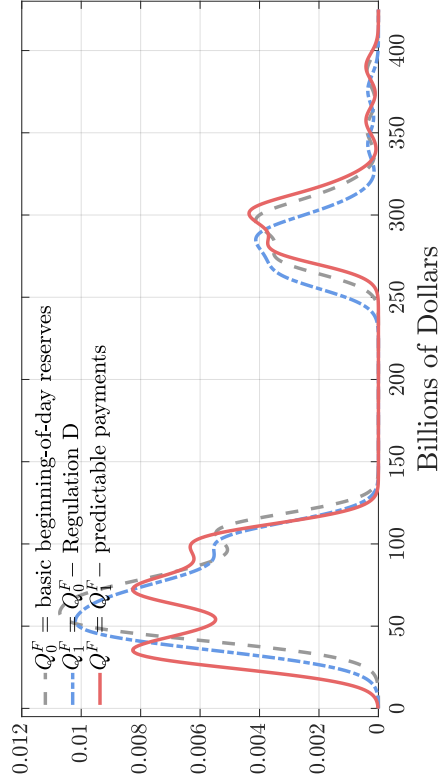
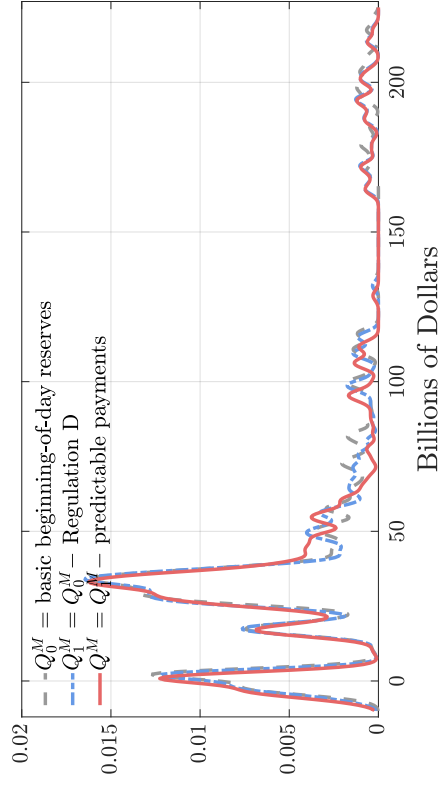
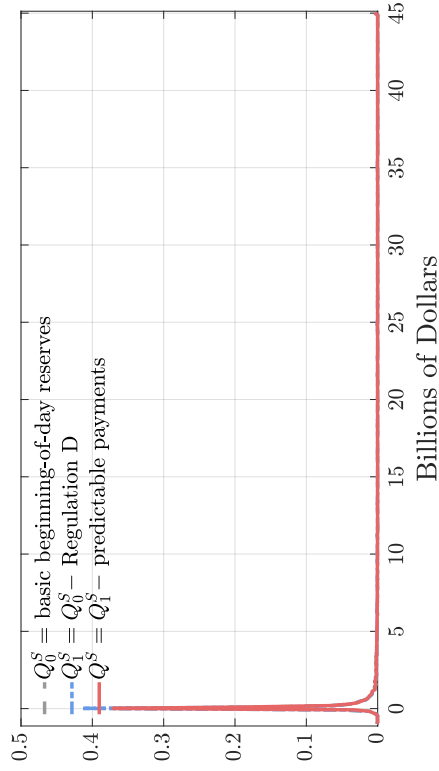
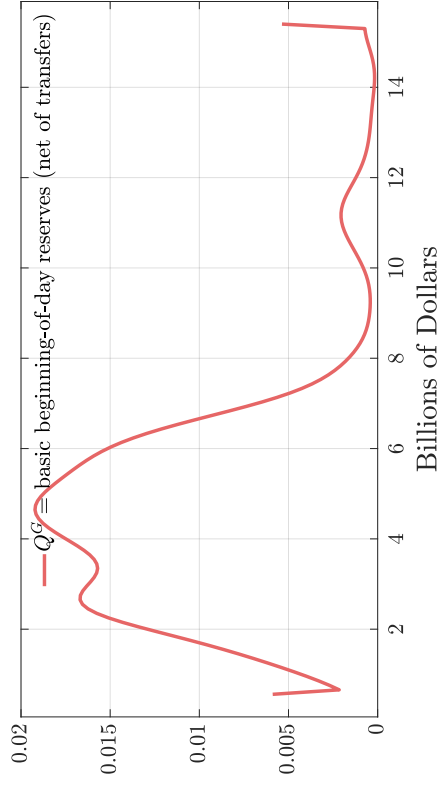
F  M  S  G 

Figure 7: Estimated beginning-of-day distributions of reserves by bank type for the year 2014.

Notes: For every bank type $i \in \mathbb{N}$, the curve labeled Q^i is the (Gaussian kernel density) estimate of the empirical distribution of beginning-of-day excess reserves net of predictable transfers. For bank type $i \in \{F, M, S\}$, the curve labeled Q_0^i is the estimate of the distribution of “basic” beginning-of-day reserves (net of fed funds repayments), and the curve labeled Q_1^i is the estimate of the distribution of “basic” beginning-of-day reserves net of the Regulation D reserve requirement.

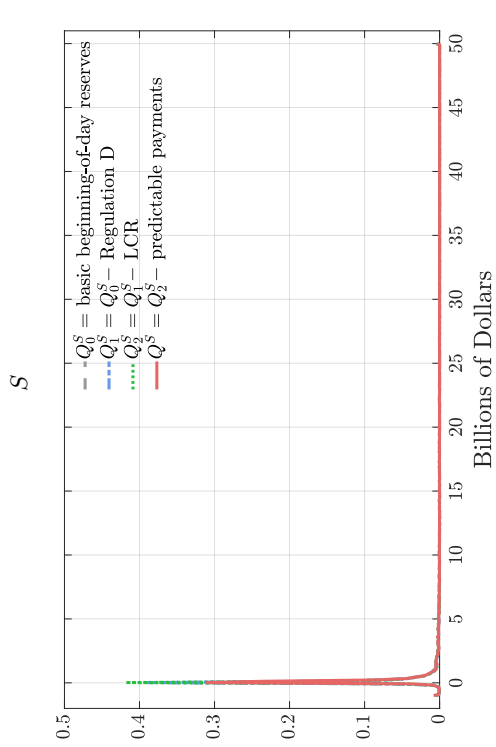
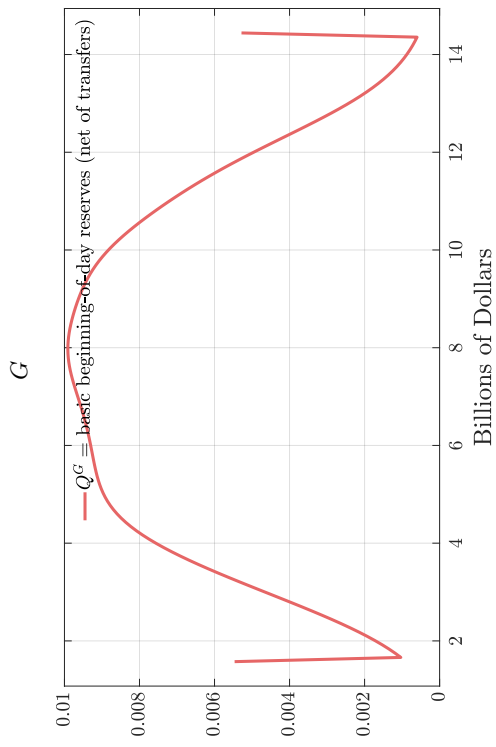
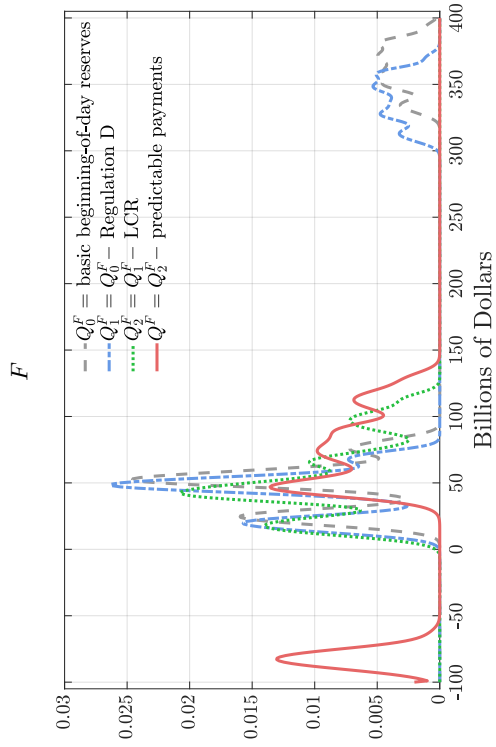
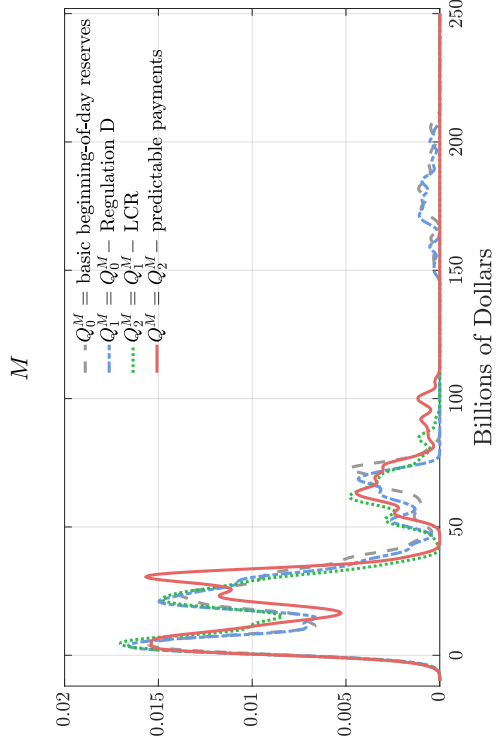


Figure 8: Estimated beginning-of-day distributions of reserves by bank type for the year 2017.

Notes: For every bank type $i \in \mathbb{N}$, the curve labeled Q^i is the (Gaussian kernel density) estimate of the empirical distribution of beginning-of-day excess reserves net of predictable transfers. For bank type $i \in \{F, M, S\}$, the curve labeled Q_0^i is the estimate of the distribution of “basic” beginning-of-day reserves (i.e., net of fed funds repayments), the curve labeled Q_1^i is the estimate of the distribution of “basic” beginning-of-day reserves net of the Regulation D reserve requirement, and the curve labeled Q_2^i is the estimate of the distribution of “basic” beginning-of-day reserves net of the Regulation D and LCR requirements.

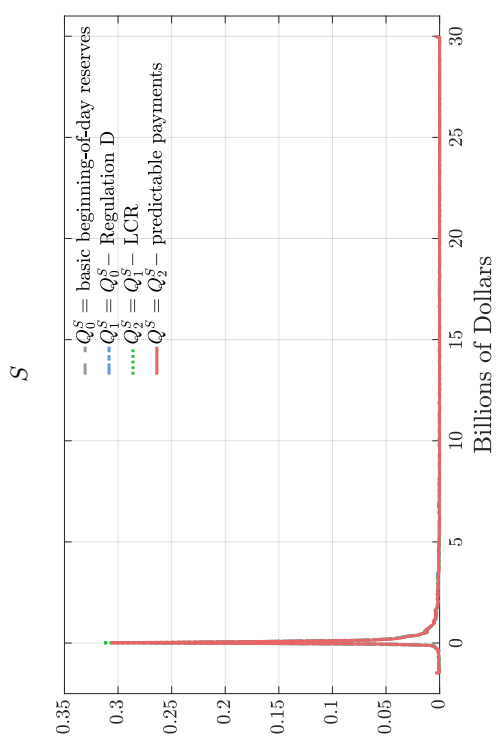
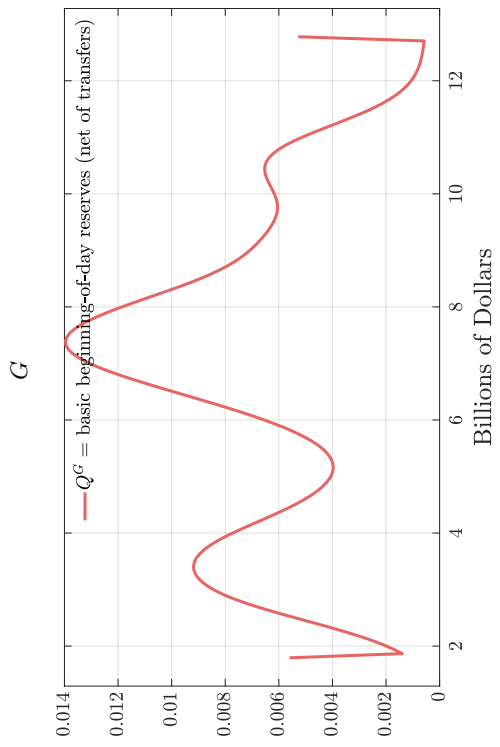
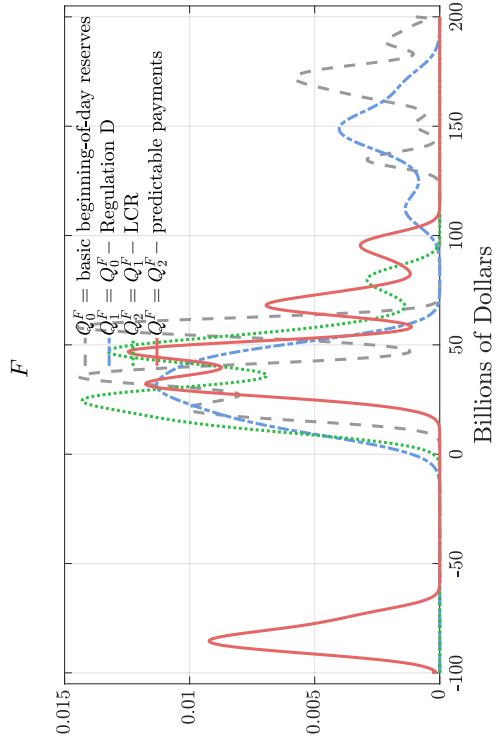
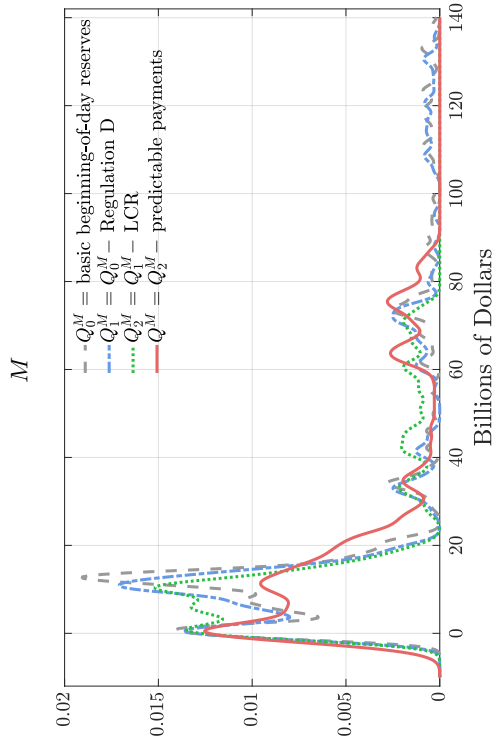


Figure 9: Estimated beginning-of-day distributions of reserves by bank type for the year 2019.

Notes: For every bank type $i \in \mathbb{N}$, the curve labeled Q^i is the (Gaussian kernel density) estimate of the empirical distribution of beginning-of-day excess reserves net of predictable transfers. For bank type $i \in \{F, M, S\}$, the curve labeled Q_0^i is the estimate of the distribution of “basic” beginning-of-day reserves (i.e., net of fed funds repayments), the curve labeled Q_1^i is the estimate of the distribution of “basic” beginning-of-day reserves net of the Regulation D reserve requirement, and the curve labeled Q_2^i is the estimate of the distribution of “basic” beginning-of-day reserves net of the Regulation D and LCR requirements.

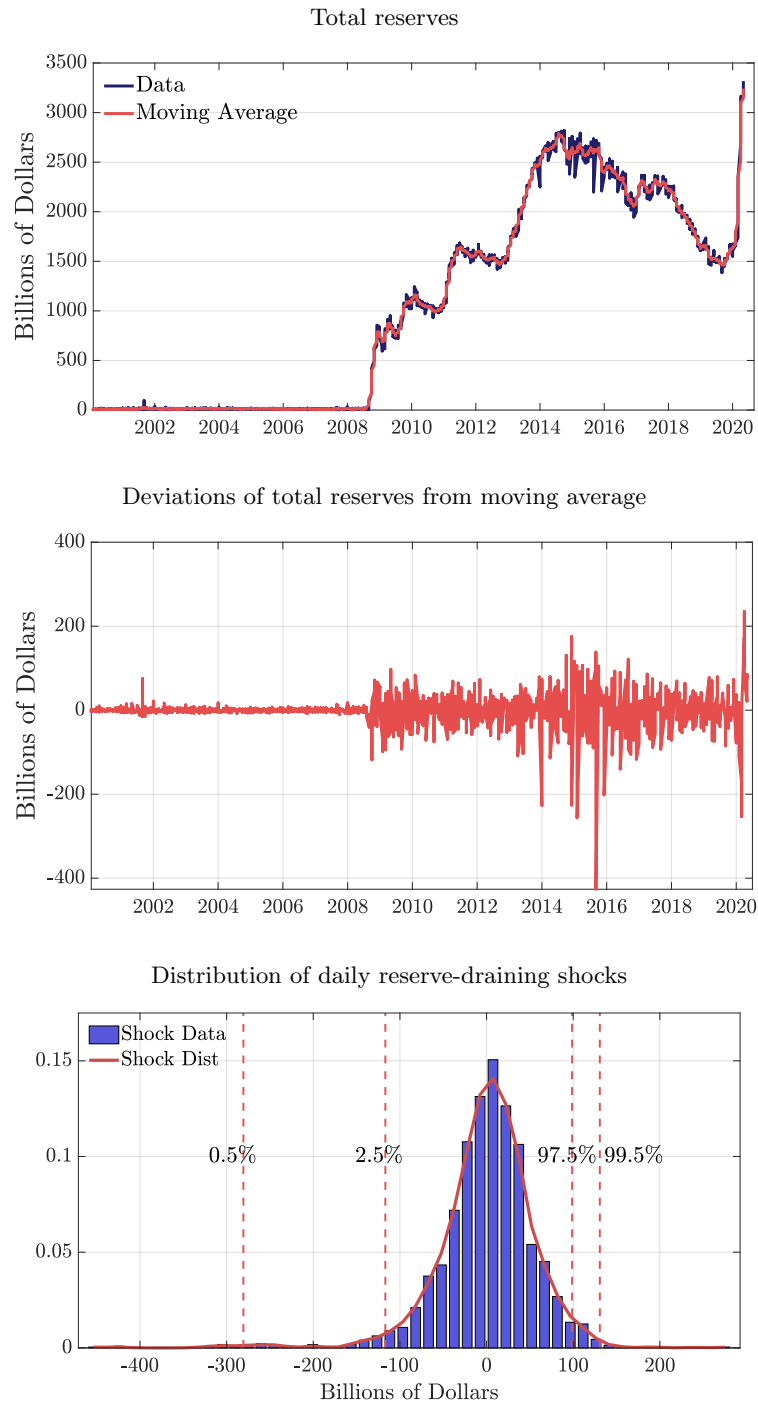


Figure 10: Aggregate supply of reserves and reserve-draining shocks.

Notes: Top panel: weekly time series of aggregate quantity of reserves and corresponding 40-day two-sided moving average. Middle panel: difference between the two time series in the top panel. Bottom panel: empirical histogram of daily deviations of the aggregate quantity of reserves from its 40-day two-sided moving average (January 2011-July 2019), and the corresponding Gaussian kernel estimate.

Parameter	Target	Moment	
		Data	Model
$n_F = 0.010$	proportion of financial institutions of type F	4/412	0.010
$n_M = 0.044$	proportion of financial institutions of type M	18/412	0.044
$n_S = 0.920$	proportion of financial institutions of type S	379/412	0.920
$n_G = 0.026$	proportion of financial institutions of type G	11/412	0.026
$\lambda_F = 0.951$	bank-level share of unexpected payments per second for type F	0.951	0.951
$\lambda_M = 0.257$	bank-level share of unexpected payments per second for type M	0.257	0.257
$\lambda_S = 0.011$	bank-level share of unexpected payments per second for type S	0.011	0.011
$\lambda_G = 0$	bank-level share of unexpected payments per second for type G	0	0
$\iota_w = 0.0300/360$	DWR (3.00% per annum, primary credit)	0.0300/360	0.0300/360
$\iota_r = 0.0235/360$	IOR (2.35% per annum)	0.0235/360	0.0235/360
$\iota_o = 0.0225/360$	ONRRP (2.25% per annum)	0.0225/360	0.0225/360
$\iota_\ell = 0.00049/360$	average value-weighted fed funds rate	0.0239/360	0.0239/360
$\iota_s = 0.00758/360$	estimated liquidity effect for 2019 (bps per \$1 bn decrease in reserves)	$\in [-0.019, -0.005]$	-0.0073
$\theta = 1/20$	conditional (below the IOR) average value-weighted fed funds rate	0.0229/360	0.0231/360
$\beta_F = 0.0300$	number of loans of financial institutions of type F relative to average	24	25
$\beta_M = 0.0024$	participation rate of financial institutions of type M (i.e., \mathcal{P}_M)	0.31	0.27
$\beta_S = 0.0007$	participation rate of financial institutions of type S (i.e., \mathcal{P}_S)	0.09	0.08
$\beta_G = 0.0036$	participation rate of financial institutions of type G (i.e., \mathcal{P}_G)	0.17	0.14
$\kappa_F = 0.039e-3$	reallocation index of financial institutions of type F (i.e., \mathcal{R}_F)	0.16	0.13
$\kappa_M = 0$	reallocation index of financial institutions of type M (i.e., \mathcal{R}_M)	-0.61	-0.64
$\kappa_S = 0.003e-3$	reallocation index of financial institutions of type S (i.e., \mathcal{R}_S)	-0.38	-0.37
$\kappa_G = 1.25e-3$	reallocation index of financial institutions of type G (i.e., \mathcal{P}_G)	1	1

Table 1: Calibration for the year 2019.

Notes: Each non-shaded parameter is calibrated externally (i.e., to match a corresponding target moment, independently of the model and other parameters). Shaded parameters are calibrated internally (i.e., jointly, to match the set of shaded target moments, using the equilibrium conditions of the model, and given the values of the parameters calibrated externally). The calibration assumes a model period corresponding to approximately to 42 seconds in a trading day, $r = 0$, $\mathbb{N} = \{F, M, S, G\}$ (as discussed in Section 3.1), $\theta_{i,j} = 1/2$ for all $i, j \in \mathbb{N} \setminus \{G\}$, $\theta_{i,j} = \theta$ if $i \in \{G\}$ and $j \in \mathbb{N} \setminus \{G\}$, $\{G_{i,j}\}_{i,j \in \mathbb{N}}$ are estimated as described in Section 3.2, $\{F_0^i\}_{i \in \mathbb{N}}$ are estimated as described in Section 3.3, $u_i = 0$ for all $i \in \mathbb{N}$, and $\{U_i\}_{i \in \mathbb{N}}$ are as in Section 4. The liquidity effect in the model is computed by extracting \$100 bn reserves using the procedure described in Section 3.6.

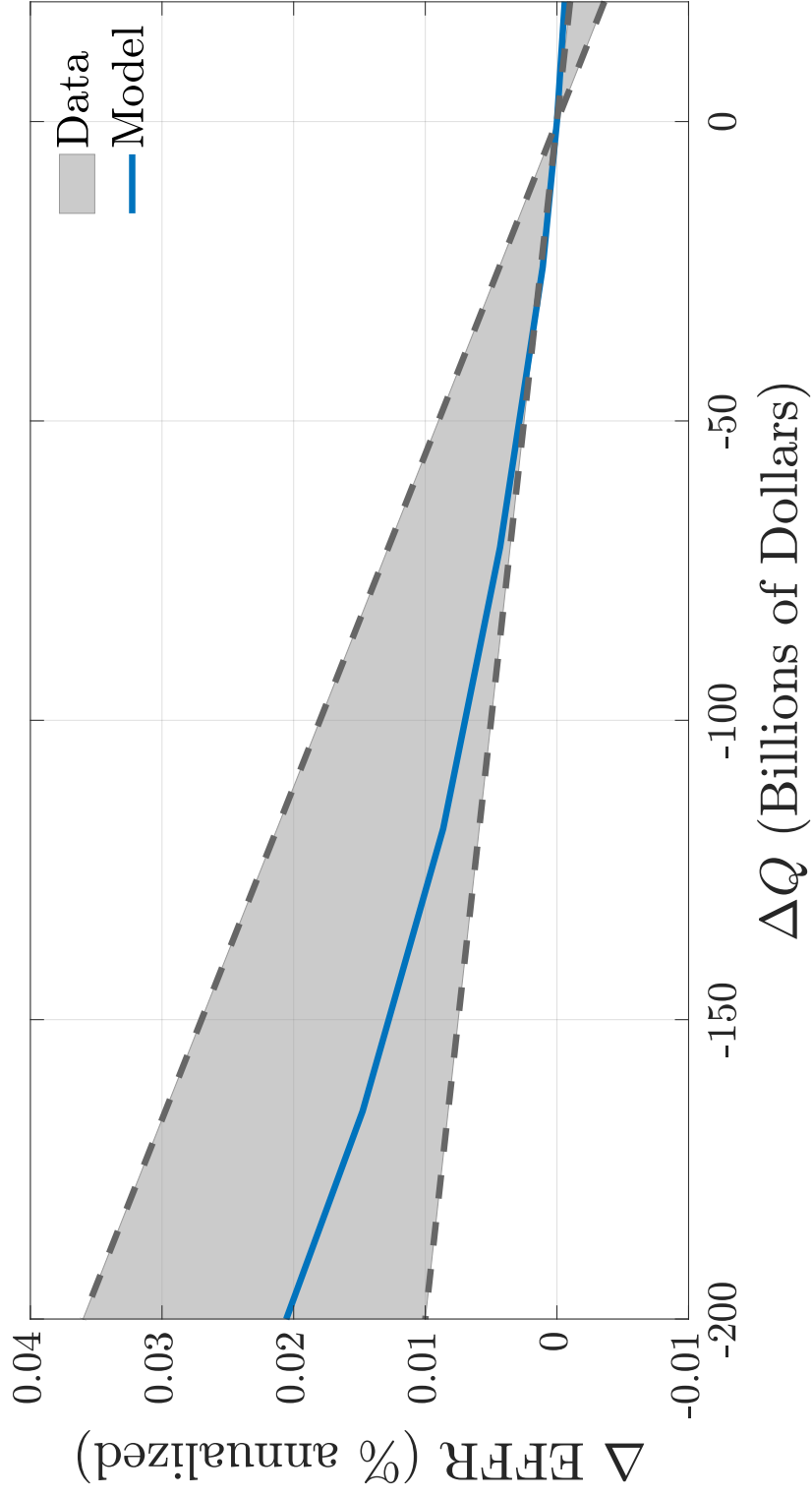


Figure 11: Liquidity effect: model and empirical estimates for the year 2019.

Notes. Rates in the vertical axis are in percent per annum. The shaded area represents the 95% confidence interval for the point estimates of the liquidity effect from specification (5). The solid line is the change in the equilibrium fed funds rate implied by the theory in response to changes in the total quantity of reserves (starting from the quantity of reserves corresponding to the 2019 calibration, and extracting reserves using the procedure described in Section 3.6.)

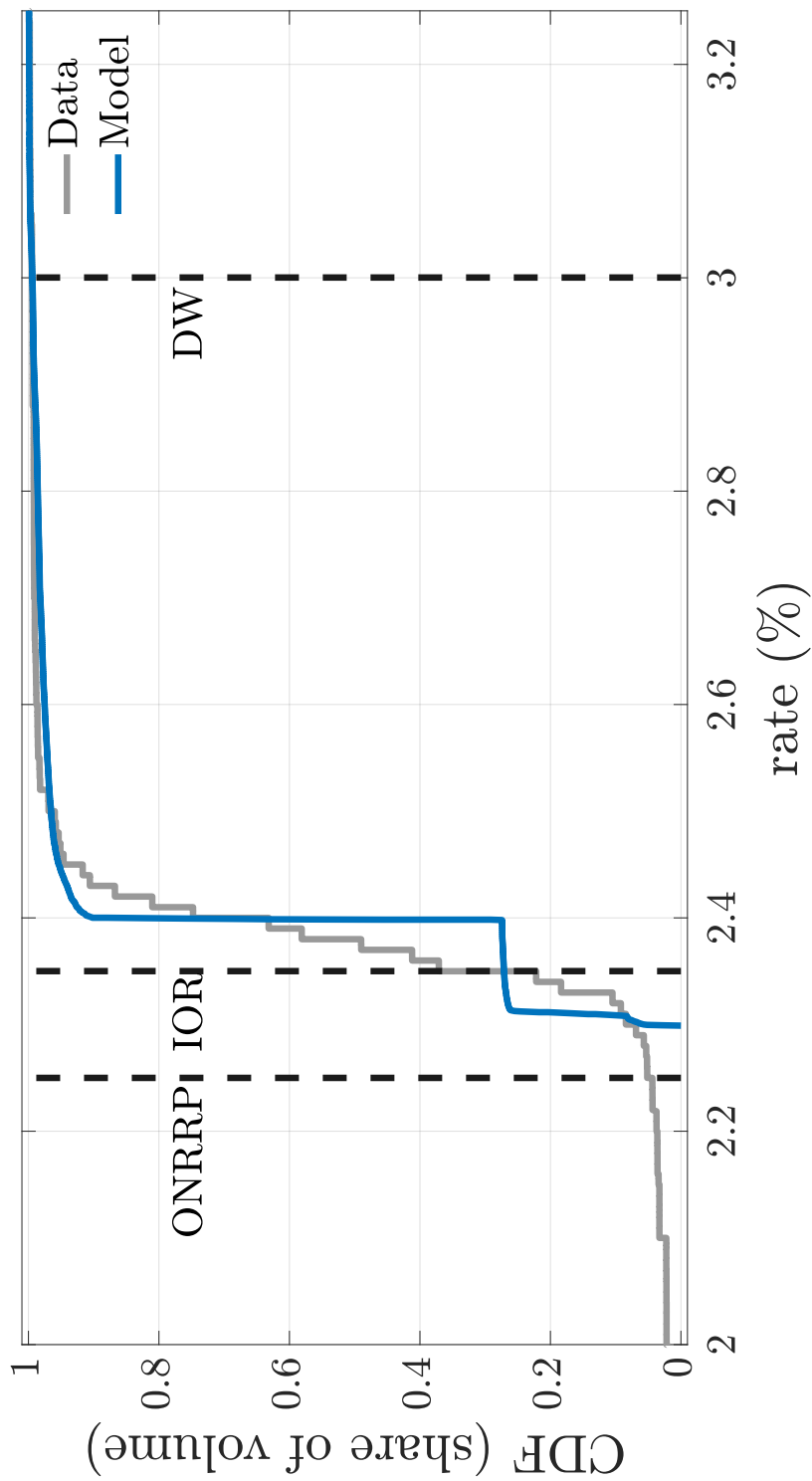


Figure 12: Empirical and theoretical cumulative distribution functions of bilateral fed funds rates for the year 2019. *Notes:* For each loan rate ι , the curve labeled “Data” (“Model”) gives the fraction of total loan volume traded at rates lower than ι in the data (model). Data are for every trading day in the period 2019/06/06–2019/07/31. The model calibrated as in Table 1. Rates are in percent per annum.

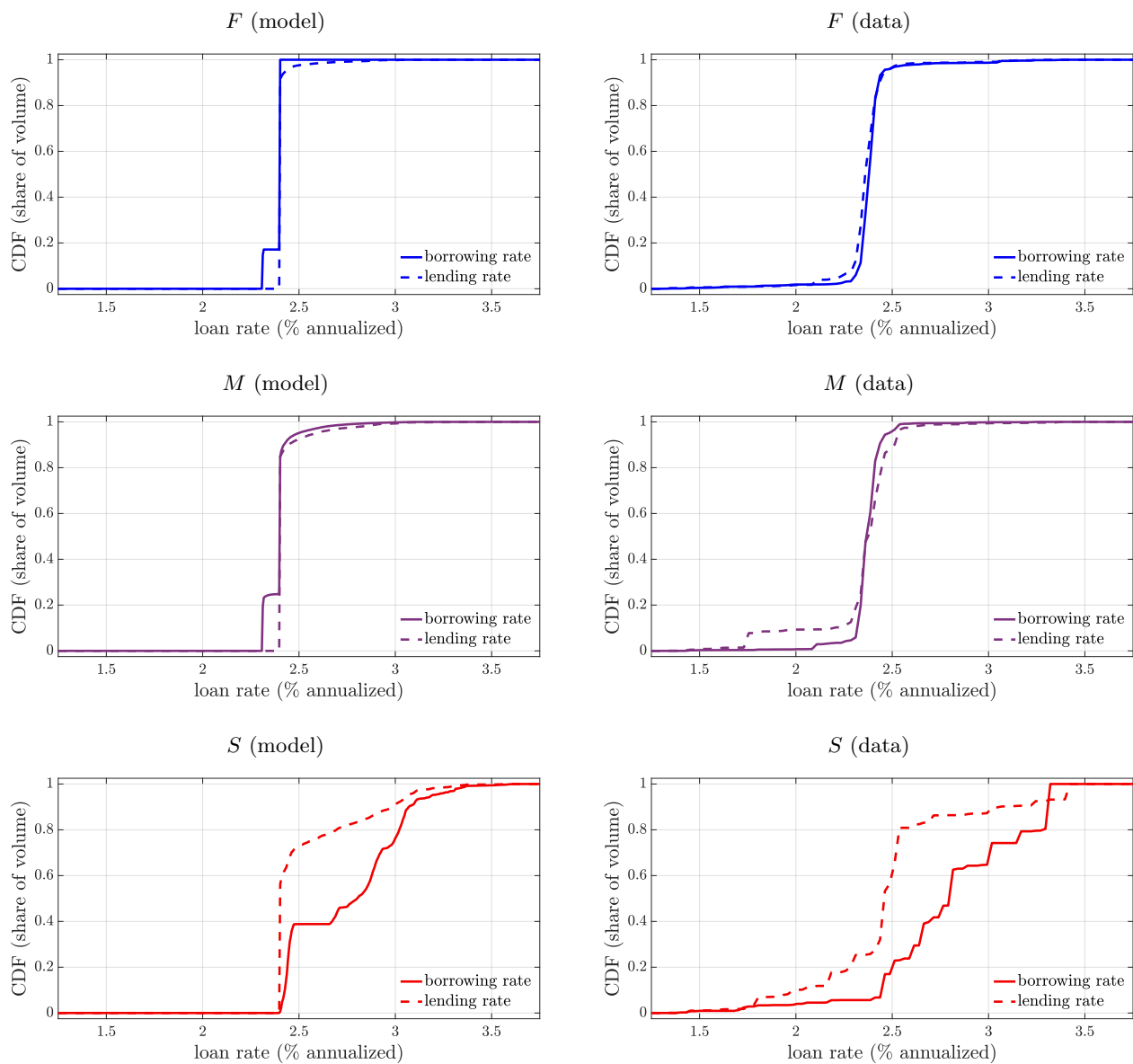


Figure 13: Cumulative distributions of borrowing and lending rates by bank type.

Notes: For each loan rate, the curve labeled “borrowing rate” (“lending rate”) gives the fraction of total reserves borrowed (lent) by banks of the type indicated in the panel heading, at rates lower than that rate. The panels on the left are for the model calibrated as in Table 1. The panels on the right are from data, for every trading day in the period 2019/06/06–2019/07/31. Rates are in percent per annum.

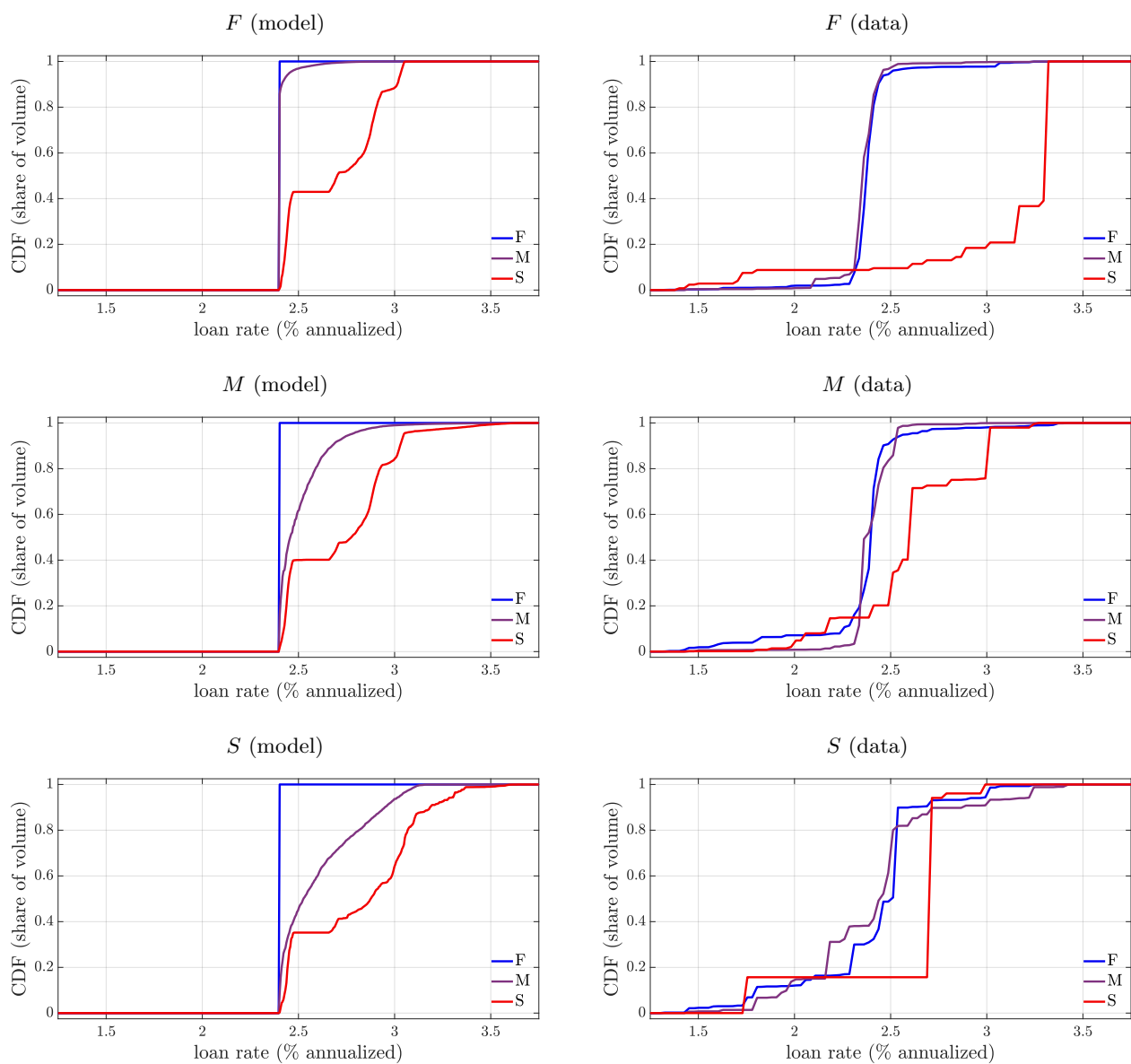


Figure 14: Cumulative distributions of loan rates between pairs of bank types.

Notes: For each loan rate, the curve labeled “ i ” (for $i \in \{F, M, S\}$) gives the fraction of total reserves borrowed by banks of type i from the bank types indicated in the panel heading, at rates lower than that rate. The panels on the left are for the model calibrated as in Table 1. The panels on the right are from data, for every trading day in the period 2019/06/06–2019/07/31. Rates are in percent per annum.

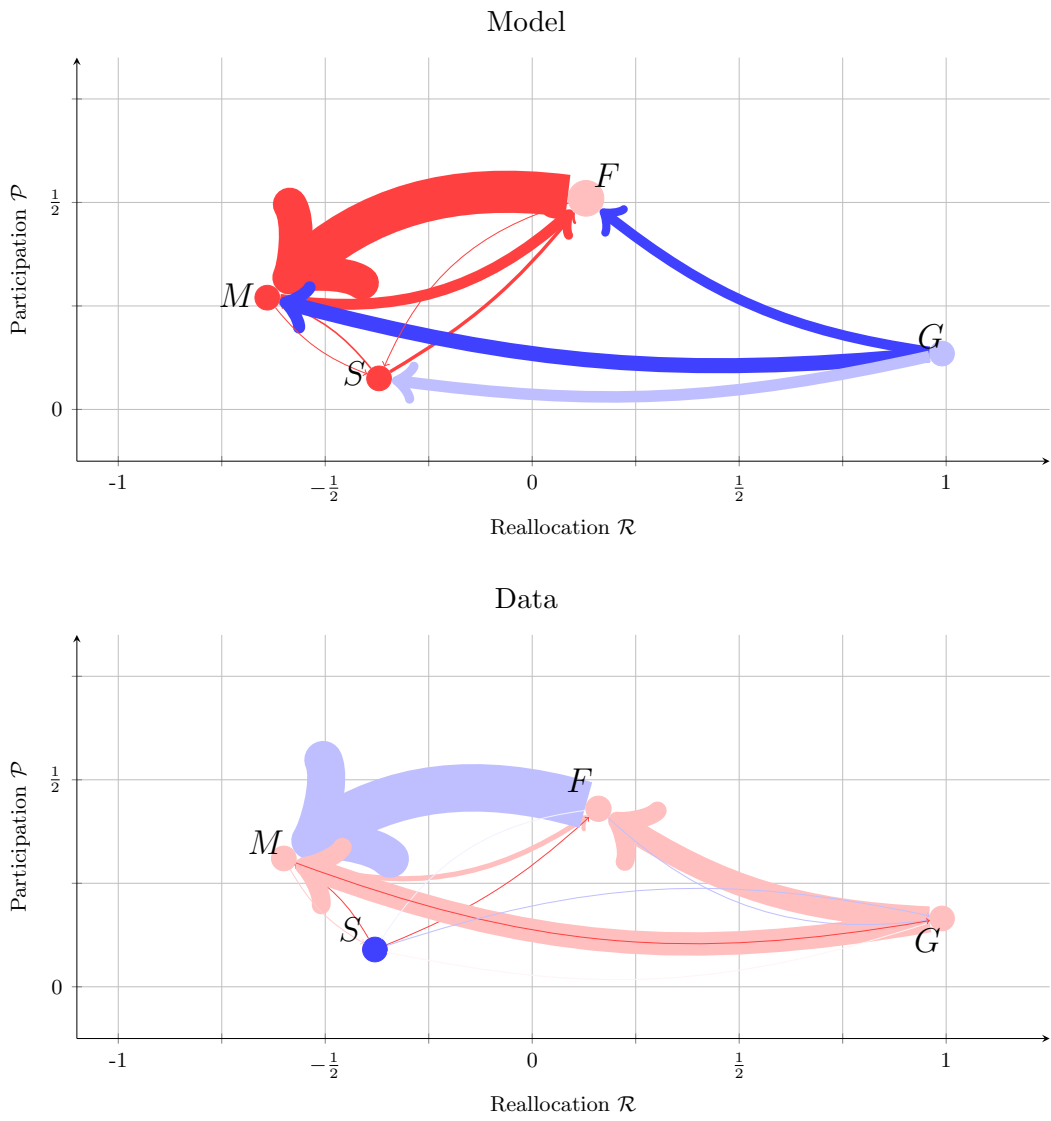


Figure 15: Theoretical and empirical fed funds trading networks for 2019.

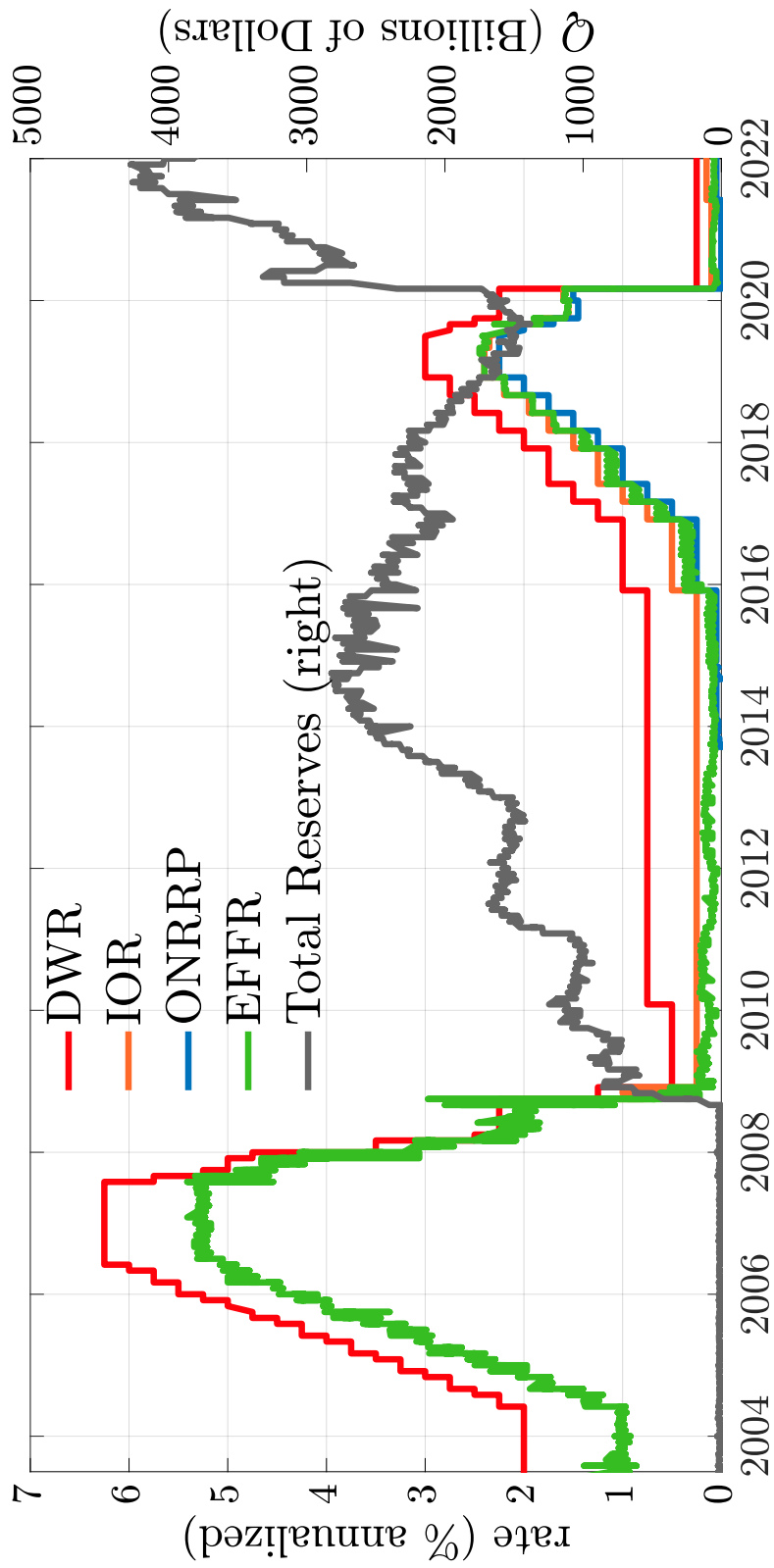


Figure 16: Time series of Total Reserves, and administered rates: Discount-Window rate (DWR), interest on reserves (IOR), and overnight reverse repo rate (ONRRP). Reserves are in billions of dollars. Rates are in percent per annum.

Notes: Total Reserves is "Reserve balances with Federal Reserve Banks: Wednesday level" from *Federal Reserve Balance Sheet: Factors Affecting Reserve Balances - H.4.1*. Administered rates are from <https://fred.stlouisfed.org>. DWR is "DPCREDIT"; IOR is "IOER" (until 2021/07) and "IORB" since (2021/08); ONRRP is "RRPONTSYAWARD"; EFFR is "DFF".

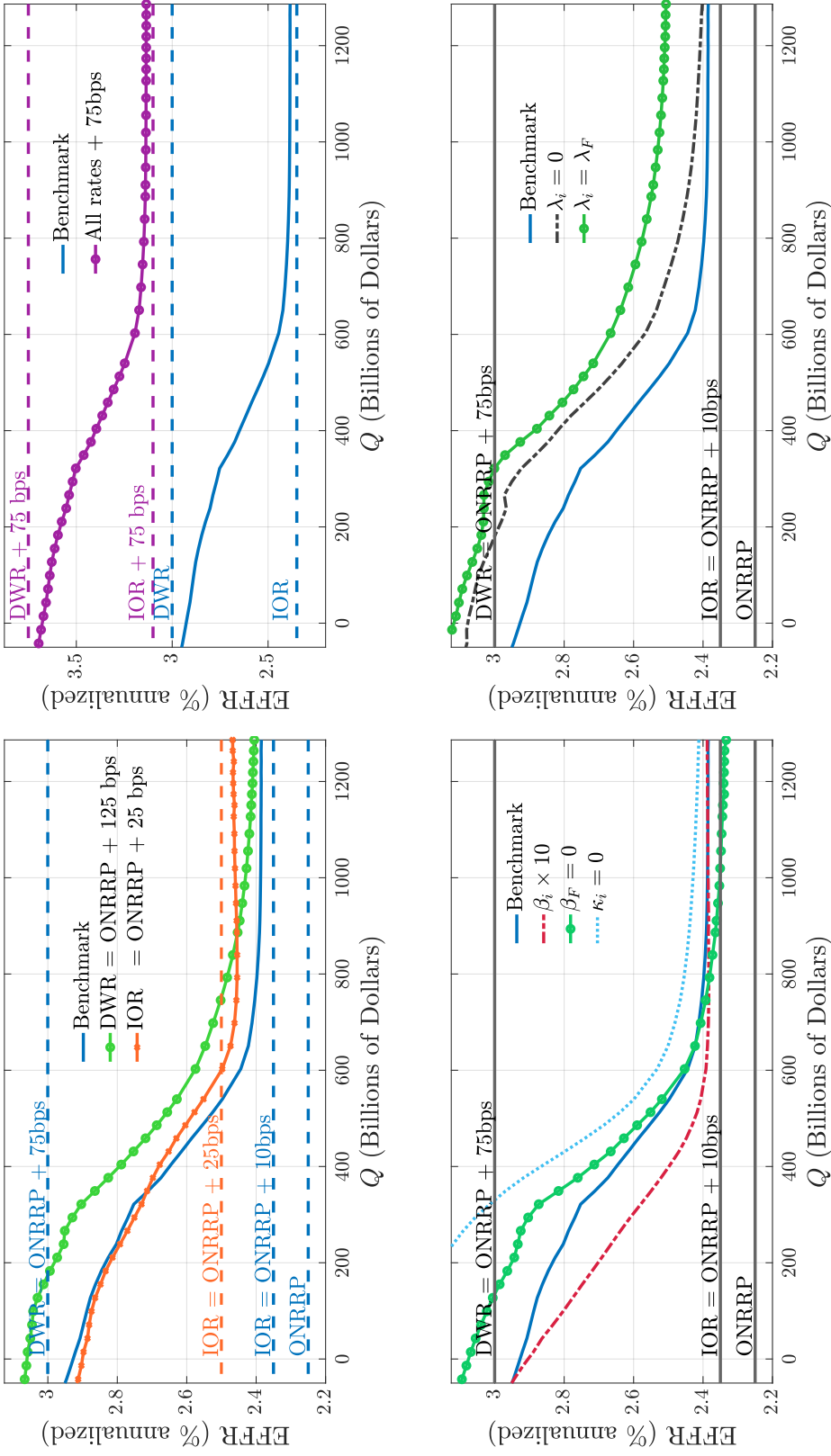


Figure 17: Theoretical aggregate demand for reserves: shifts and rotations.

Notes: In all panels, the curve labeled “Benchmark” is the theoretical aggregate demand $\iota_{v,w}^* = \mathcal{D}(Q_{v,w}; \Pi)$ for the model calibrated as in Table 1, and with $\iota_{v,w}^*$ and $Q_{v,w}$ computed with the interpolation procedure described in Section 3.6, for $\gamma_0 = 2017$ and $\gamma_1 = 2019$. Top-left panel: benchmark aggregate demand, and aggregate demands resulting from two experiments: (i) increase IOR by 15 bps. Top-right panel: benchmark aggregate demand, and aggregate demands resulting from increasing all administered rates (i.e., DWR, IOR, and ONRRP) by 75 bps. Bottom-left panel: benchmark aggregate demand, and aggregate demands resulting from three experiments: (i) multiply $\{\beta_i\}_{i \in \mathbb{N}}$ by 10; (ii) set $\beta_F = 0$; (iii) set $\kappa_F = \kappa_S = 0$. Bottom-right panel: benchmark aggregate demand, and aggregate demands resulting from two experiments: (i) set $\lambda_i = 0$ for all $i \in \mathbb{N}$; (ii) set $\lambda_i = \lambda_F$ for all $i \in \mathbb{N}$.

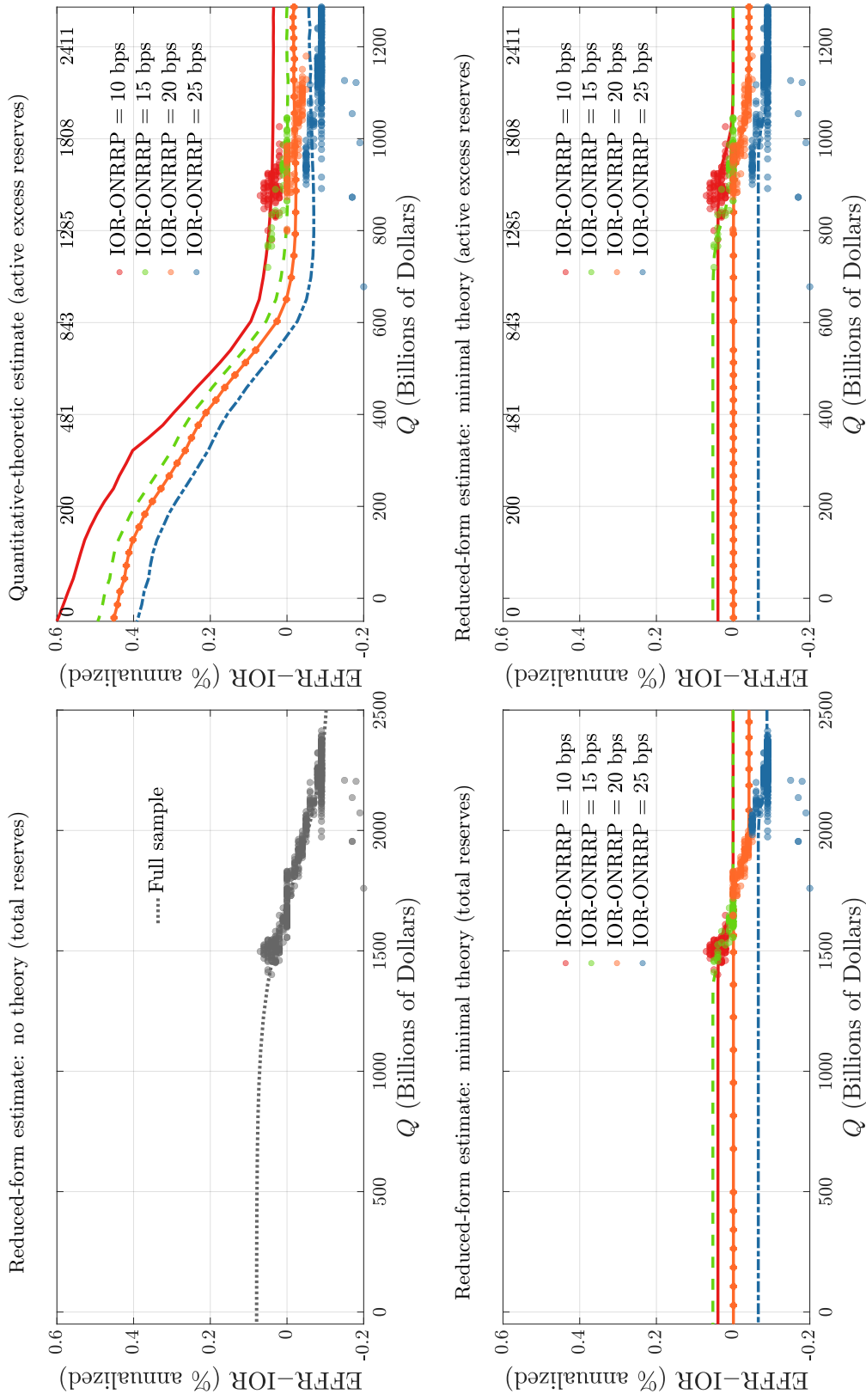


Figure 18: Aggregate demand for reserves: estimation.

Notes: In each panel: vertical axis is EFFR-IOR spread; horizontal axis is *total reserves* or *active excess reserves* as defined in Section 6. Full sample: all trading days in 2017/01/20–2019/09/13. Top-left panel: total reserves and NLS fit of (11) on full sample. Bottom-left panel: total reserves and NLS fits of (11) on each subsample defined by IOR-ONRRP spread. Top-right panel: active excess reserves and aggregate demands implied by the theory for each IOR-ONRRP spread (all other parameters as in the baseline calibration). Bottom-right panel: active excess reserves and NLS fits of (11) on each subsample. Secondary horizontal axis (above top-left and bottom-left panels) translates active excess reserves into total reserves.

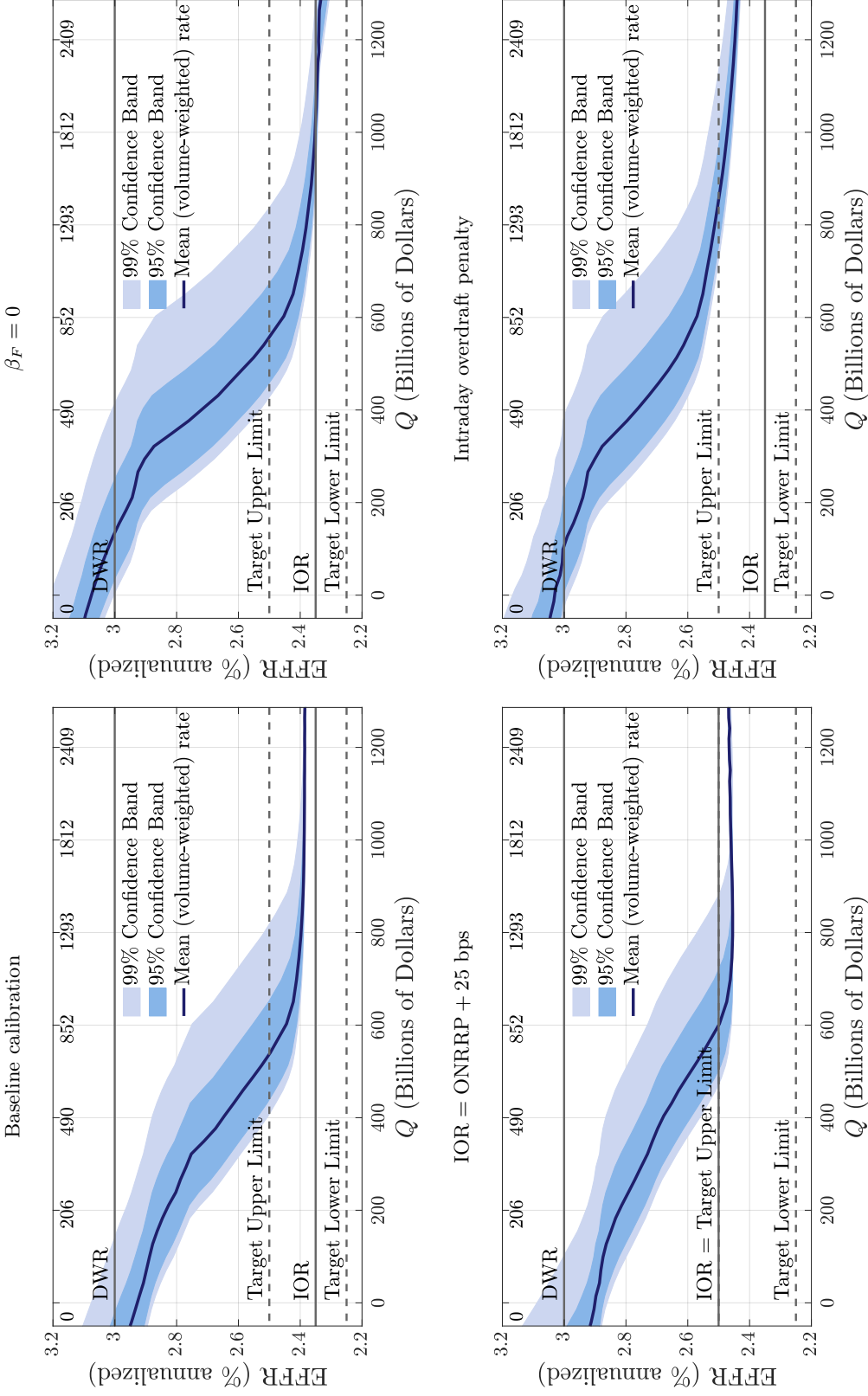


Figure 19: Monetary confidence bands.

Notes: In each panel, the curve labeled “Mean (volume-weighted) rate” is the theoretical money demand, $\mathcal{D}(Q)$. The lower and upper boundaries of the shaded area labeled “99% Confidence Band” are $\mathcal{D}(Q + Z_{99.5})$ and $\mathcal{D}(Q + Z_{0.5})$, respectively, where Z_p is the p^{th} percentile of the empirical distribution of reserve-draining shocks. The lower and upper boundaries of the shaded area labeled “95% Confidence Band” are $\mathcal{D}(Q + Z_{97.5})$ and $\mathcal{D}(Q + Z_{2.5})$, respectively. The top-left panel corresponds to the baseline calibration. The other panels are for calibrations that differ in one parameter from the baseline calibration. The top-right panel sets $\beta_F = 0$ (the baseline has $\beta_F = 0.03$). The bottom-left panel increases the IOR by 15 bps (from ONRRP + 10 bps in the baseline, to ONRRP + 25 bps). The bottom-right panel assumes $u_i(a) = \iota_d a \mathbb{1}_{\{a < 0\}}$ for all i , with $\iota_d = \frac{0.2}{800} \iota_w$ (the baseline has $\iota_d = 0$).

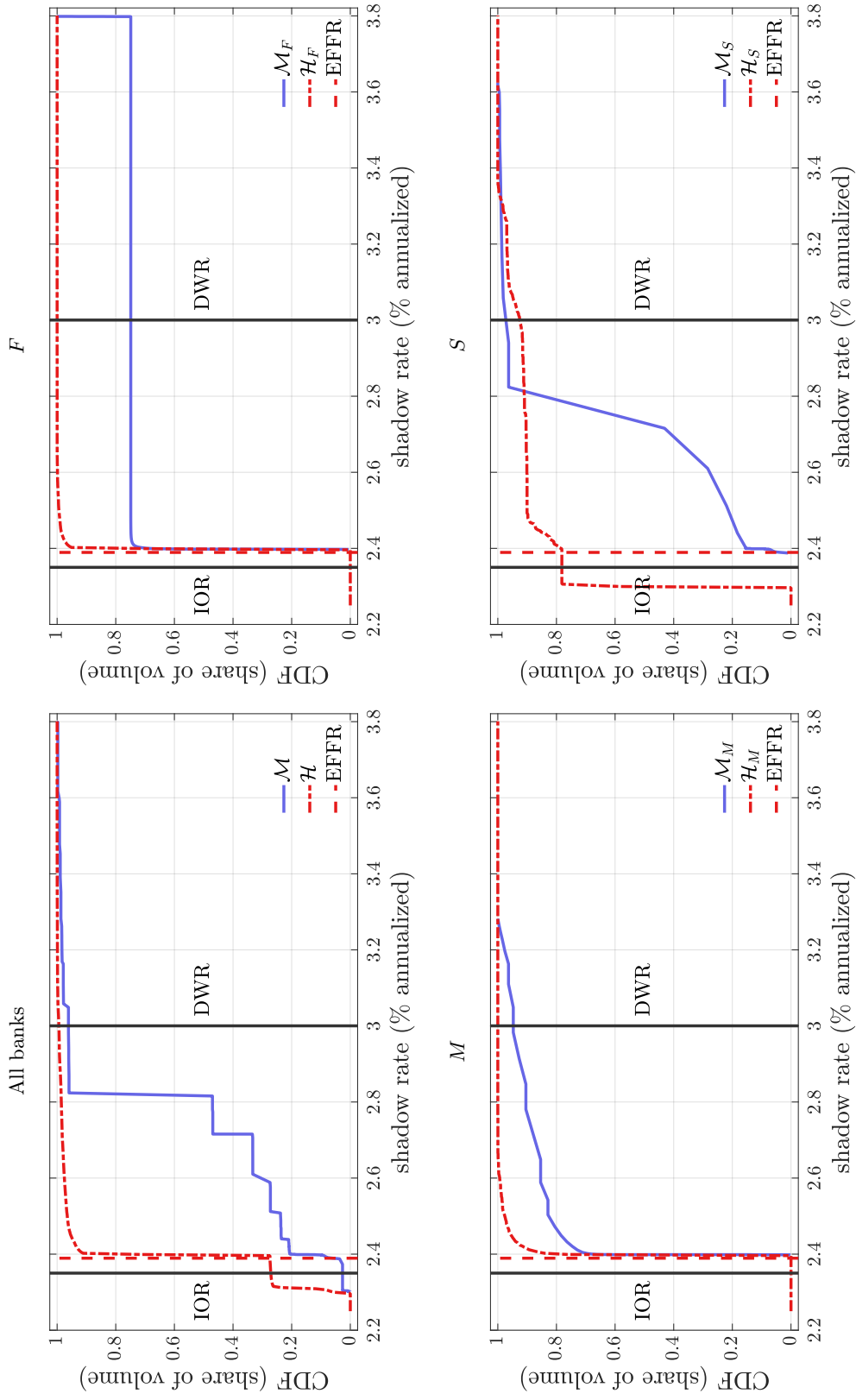


Figure 20: Cross-bank distributions of the shadow price of reserves.

Notes: All panels are constructed using data generated from the model under the baseline calibration. The beginning-of-day cumulative distribution function of shadow prices is denoted \mathcal{M}_i for banks of type i , and \mathcal{M} for all banks. The cumulative distribution function of all the bilateral loan rates negotiated throughout the day is denoted \mathcal{H} . The cumulative distribution function of all loan rates paid or received by banks of type i is denoted \mathcal{H}_i . The dashed vertical line labeled “EFFR” denotes the volume-weighted average fed funds rate on *all trades* implied by the theory. The IOR and DWR are denoted by solid vertical lines. All rates are in percent per annum.

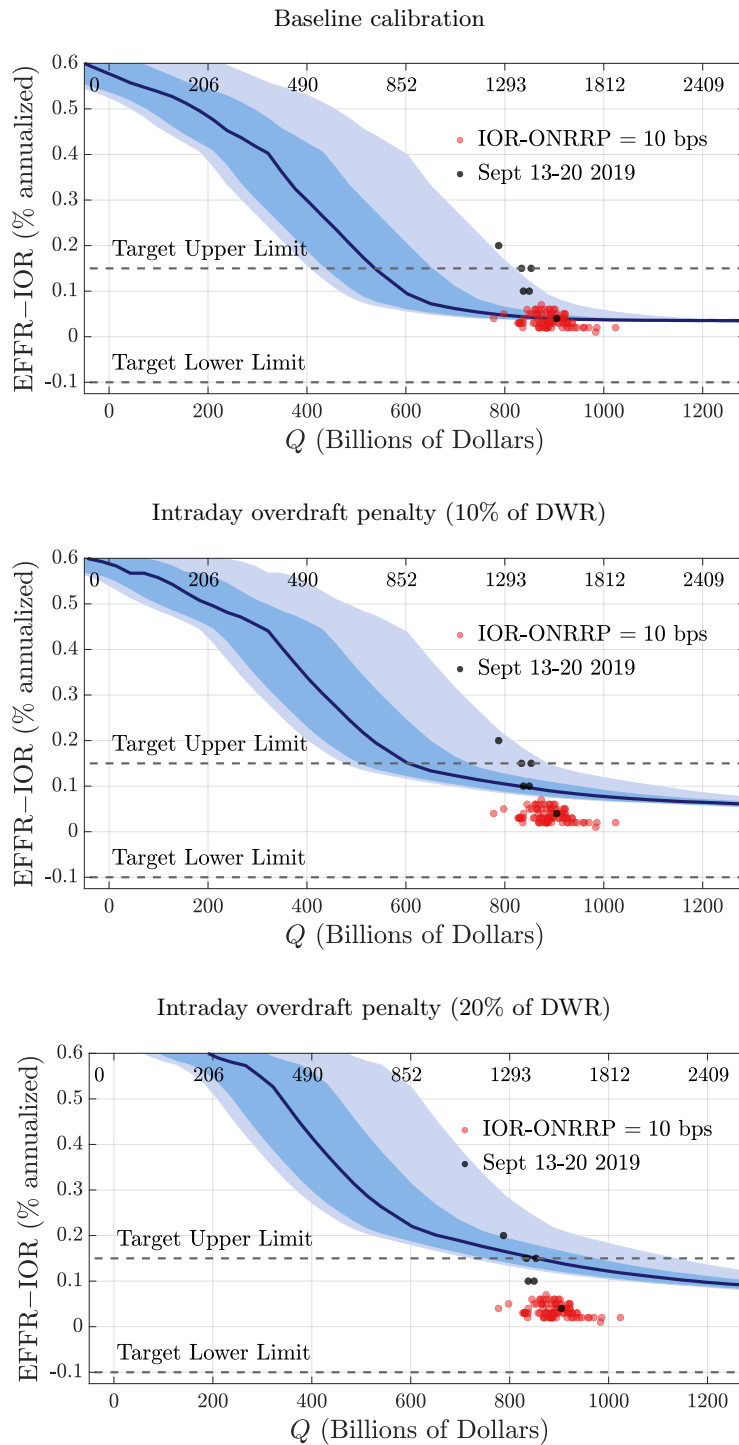


Figure 21: The events of September 13–20, 2019.

Notes: Each panel shows an MCB with the EFR-IO spread on the vertical axis (in percent per annum). The MCBs assume $u_i(a) = \iota_d a \mathbb{I}_{\{a < 0\}}$ for all i , with $\iota_d = \frac{x}{800} \iota_w$; the top panel has $x = 0$ (the baseline calibration), the middle panel $x = 0.1$, and the bottom panel $x = 0.2$. The data points labeled “IOR-ONRRP=10 bps” are for the period 2019/05/02–2019/09/13. The dashed lines labeled “Target Upper Limit” and “Target Lower Limit” are the top and bottom of the fed funds target range minus the IOR for the period 2019/05/02–2019/09/18.

Day	Administered Rates						FFFR	EFFR-IOR	Fed Repo
	ONRRP	IOR	DWR	TRL	TRU				
September 13 (Friday)	2.00	2.10	2.75	2.00	2.25	2.14	0.04	0	
September 16 (Monday)	2.00	2.10	2.75	2.00	2.25	2.25	0.15	0	
September 17 (Tuesday)	2.00	2.10	2.75	2.00	2.25	2.30	0.20	53	
September 18 (Wednesday)	2.00	2.10	2.75	2.00	2.25	2.25	0.15	75	
September 19 (Thursday)	1.70	1.80	2.50	1.75	2.00	1.90	0.10	75	
September 20 (Friday)	1.70	1.80	2.50	1.75	2.00	1.90	0.10	75	

Table 2: The events of September 13–20, 2019.

Notes: ONRRP, IOR, and DWR, denote the overnight reverse repo rate, the interest rate paid on reserves, and the discount-window rate, respectively. The lower limit of the fed funds target range is denoted TRL, and the upper limit is denoted TRU. The effective fed funds rate is denoted EFFR, and EFFR-IOR denotes the EFFR-IOR spread. The column labeled “Fed Repo” reports the quantity of reserves injected by the Federal Reserve during day t through overnight repo operations. All rates are in percent per annum. All quantities in billions of dollars.

A Theory: supplementary results

A.1 Value function

Let $J_t^i(a, c) : \mathbb{N} \times \mathbb{T} \times \mathbb{R}^2 \rightarrow \mathbb{R}$ denote the maximum attainable payoff to a bank of type i that at time $t \in \mathbb{T}$ has reserve balance $a \in \mathbb{R}$ and net credit position $c \in \mathbb{R}$. Then, $J_t^i(a, c)$ satisfies

$$\begin{aligned}
& J_t^i(a, c) \\
&= \mathbb{E}_t \left\{ \mathbb{I}_{\{T-t \leq \min[\tau(\beta_i), \tau(\lambda_i)]\}} \left[\int_t^T e^{-r(s-t)} u_i(a) ds + e^{-r(T-t)} \left[U_i(a) + e^{-r(\bar{T}-T)} c \right] \right] \right. \\
&+ \mathbb{I}_{\{\tau(\lambda_i) < \min[\tau(\beta_i), T-t]\}} \left[\int_t^{t+\tau(\lambda_i)} e^{-r(s-t)} u_i(a) ds \right. \\
&+ \left. \left. e^{-r\tau(\lambda_i)} \sum_{j \in \mathbb{N}} \pi_j \int J_{t+\tau(\lambda_i)}^i(a-z, c) dG_{ij}(z) \right] \right. \\
&+ \mathbb{I}_{\{\tau(\beta_i) < \min[\tau(\lambda_i), T-t]\}} \left[\int_t^{t+\tau(\beta_i)} e^{-r(s-t)} u_i(a) ds \right. \\
&+ \left. \left. e^{-r\tau(\beta_i)} \sum_{j \in \mathbb{N}} \sigma_j \int J_{t+\tau(\beta_i)}^i \left[a - b_{t+\tau(\beta_i)}^{ij}(a, \tilde{a}), c + R_{t+\tau(\beta_i)}^{ji}(\tilde{a}, a) \right] dF_{t+\tau(\beta_i)}^j(\tilde{a}) \right] \right\}, \quad (12)
\end{aligned}$$

where $\tau(\zeta)$ denotes the exponentially distributed first passage time of the Poisson process with arrival rate ζ ,

$$\begin{aligned}
\pi_j &\equiv \frac{\lambda_j n_j}{\sum_{i \in \mathbb{N}} \lambda_i n_i} \\
\sigma_j &\equiv \frac{\beta_j n_j}{\sum_{k \in \mathbb{N}} \beta_k n_k},
\end{aligned}$$

and

$$\begin{aligned}
& (b_t^{ij}(a, \tilde{a}), R_t^{ji}(\tilde{a}, a)) \\
&= \arg \max_{(b, R) \in \mathbb{R}^2} \left[J_t^i(a-b, c+R) - J_t^i(a, c) \right]^{\theta_{ij}} \left[J_t^j(\tilde{a}+b, c-R) - J_t^j(\tilde{a}, c) \right]^{\theta_{ji}}. \quad (13)
\end{aligned}$$

Lemma 1 *The function*

$$J_t^i(a, c) = V_t^i(a) + e^{-r(\bar{T}-t)} c \quad (14)$$

satisfies (12) if and only if $V_t^i(a)$ satisfies

$$\begin{aligned}
V_t^i(a) = & \mathbb{E}_t \left\{ \mathbb{I}_{\{T-t \leq \min[\tau(\beta_i), \tau(\lambda_i)]\}} \left[\int_t^T e^{-r(s-t)} u_i(a) ds + e^{-r(T-t)} U_i(a) \right] \right. \\
& + \mathbb{I}_{\{\tau(\lambda_i) < \min[\tau(\beta_i), T-t]\}} \left[\int_t^{t+\tau(\lambda_i)} e^{-r(s-t)} u_i(a) ds \right. \\
& + \left. e^{-r\tau(\lambda_i)} \sum_{j \in \mathbb{N}} \pi_j \int V_{t+\tau(\lambda_i)}^i(a-z) dG_{ij}(z) \right] \\
& + \mathbb{I}_{\{\tau(\beta_i) < \min[\tau(\lambda_i), T-t]\}} \left[\int_t^{t+\tau(\beta_i)} e^{-r(s-t)} u_i(a) ds \right. \\
& + \left. e^{-r\tau(\beta_i)} \sum_{j \in \mathbb{N}} \sigma_j \int \left[V_{t+\tau(\beta_i)}^i(a - b_{t+\tau(\beta_i)}^{ij}(a, \tilde{a})) + \bar{R}_{t+\tau(\beta_i)}^{ji}(\tilde{a}, a) \right] dF_{t+\tau(\beta_i)}^j(\tilde{a}) \right] \left. \right\}, \quad (15)
\end{aligned}$$

with

$$\bar{R}_{t+\tau(\beta_i)}^{ji}(\tilde{a}, a) \equiv e^{-r\{\bar{T}-[t+\tau(\beta_i)]\}} R_{t+\tau(\beta_i)}^{ji}(\tilde{a}, a),$$

and $(b_t^{ij}(a, \tilde{a}), R_t^{ji}(\tilde{a}, a))$ given by (1) and (2).

Proof. With (14), (13) becomes equivalent to (1) and (2). Substitute (14) into (12) to get

$$\begin{aligned}
& V_t^i(a) + e^{-r(\bar{T}-t)} c \\
= & \mathbb{E}_t \left\{ \mathbb{I}_{\{T-t \leq \min[\tau(\beta_i), \tau(\lambda_i)]\}} \left[\int_t^T e^{-r(s-t)} u_i(a) ds + e^{-r(T-t)} \left[U_i(a) + e^{-r(\bar{T}-T)} c \right] \right] \right. \\
& + \mathbb{I}_{\{\tau(\lambda_i) < \min[\tau(\beta_i), T-t]\}} \left[\int_t^{t+\tau(\lambda_i)} e^{-r(s-t)} u_i(a) ds \right. \\
& + \left. e^{-r\tau(\lambda_i)} \sum_{j \in \mathbb{N}} \pi_j \int \left[V_{t+\tau(\lambda_i)}^i(a-z) + e^{-r\{\bar{T}-[t+\tau(\lambda_i)]\}} c \right] dG_{ij}(z) \right] \\
& + \mathbb{I}_{\{\tau(\beta_i) < \min[\tau(\lambda_i), T-t]\}} \left[\int_t^{t+\tau(\beta_i)} e^{-r(s-t)} u_i(a) ds + e^{-r\tau(\beta_i)} \sum_{j \in \mathbb{N}} \sigma_j \int \left[V_{t+\tau(\beta_i)}^i(a - b_{t+\tau(\beta_i)}^{ij}(a, \tilde{a})) \right. \right. \\
& + \left. \left. e^{-r\{\bar{T}-[t+\tau(\beta_i)]\}} [c + R_{t+\tau(\beta_i)}^{ji}(\tilde{a}, a)] \right] dF_{t+\tau(\beta_i)}^j(\tilde{a}) \right] \left. \right\},
\end{aligned}$$

which after cancelling the terms proportional to c , becomes identical to (15). ■

Lemma 2 *The Bellman equation (15) can be written as*

$$\begin{aligned}
V_t^i(a) &= \left[1 - e^{-(r+\beta_i+\lambda_i)(T-t)} \right] \frac{u_i(a)}{r + \beta_i + \lambda_i} + e^{-(r+\beta_i+\lambda_i)(T-t)} U_i(a) \\
&+ \lambda_i \int_t^T e^{-(r+\beta_i+\lambda_i)(\tau-t)} \left[\sum_{j \in \mathbb{N}} \pi_j \int V_\tau^i(a-z) dG_{ij}(z) \right] d\tau \\
&+ \beta_i \int_t^T e^{-(r+\beta_i+\lambda_i)(\tau-t)} \left[V_\tau^i(a) + \sum_{j \in \mathbb{N}} \sigma_j \theta_{ij} \int \max_{b \in \mathbb{R}} S_\tau^{ij}(a, \tilde{a}, b) dF_\tau^j(\tilde{a}) \right] d\tau \quad (16)
\end{aligned}$$

or equivalently, as (3) with boundary condition $V_T^i(a) = U_i(a)$.

Proof. With the bargaining outcomes (1) and (2), (15) can be rewritten as

$$\begin{aligned}
V_t^i(a) &= \mathbb{E}_t \left\{ \mathbb{I}_{\{T-t \leq \min[\tau(\beta_i), \tau(\lambda_i)]\}} \left[\int_t^T e^{-r(s-t)} u_i(a) ds + e^{-r(T-t)} U_i(a) \right] \right. \\
&+ \mathbb{I}_{\{\tau(\lambda_i) < \min[\tau(\beta_i), T-t]\}} \left[\int_t^{t+\tau(\lambda_i)} e^{-r(s-t)} u_i(a) ds + e^{-r\tau(\lambda_i)} \sum_{j \in \mathbb{N}} \pi_j \int V_{t+\tau(\lambda_i)}^i(a-z) dG_{ij}(z) \right] \\
&+ \mathbb{I}_{\{\tau(\beta_i) < \min[\tau(\lambda_i), T-t]\}} \left[\int_t^{t+\tau(\beta_i)} e^{-r(s-t)} u_i(a) ds \right. \\
&\left. \left. + e^{-r\tau(\beta_i)} \sum_{j \in \mathbb{N}} \sigma_j \int \left[V_{t+\tau(\beta_i)}^i(a) + \theta_{ij} \max_{b \in \mathbb{R}} S_{t+\tau(\beta_i)}^{ij}(a, \tilde{a}, b) \right] dF_{t+\tau(\beta_i)}^j(\tilde{a}) \right] \right\},
\end{aligned}$$

where

$$S_t^{ij}(a, \tilde{a}, b) \equiv V_t^i(a-b) + V_t^j(\tilde{a}+b) - V_t^i(a) - V_t^j(\tilde{a}).$$

The first term on the right side of $V_t^i(a)$ can be written as

$$\begin{aligned}
&\mathbb{E}_t \left\{ \mathbb{I}_{\{T-t \leq \min[\tau(\beta_i), \tau(\lambda_i)]\}} \left[\int_t^T e^{-r(s-t)} u_i(a) ds + e^{-r(T-t)} U_i(a) \right] \right\} \\
&= e^{-(\beta_i+\lambda_i)(T-t)} \left\{ \left[1 - e^{-r(T-t)} \right] \frac{u_i(a)}{r} + e^{-r(T-t)} U_i(a) \right\}.
\end{aligned}$$

The second term on the right side of $V_t^i(a)$ can be written as

$$\begin{aligned} & \mathbb{E}_t \left\{ \mathbb{I}_{\{\tau(\lambda_i) < \min[\tau(\beta_i), T-t]\}} \left[\int_t^{t+\tau(\lambda_i)} e^{-r(s-t)} u_i(a) ds + e^{-r\tau(\lambda_i)} \sum_{j \in \mathbb{N}} \pi_j \int V_{t+\tau(\lambda_i)}^i(a-z) dG_{ij}(z) \right] \right\} \\ &= \frac{\lambda_i}{\beta_i + \lambda_i} \frac{r [1 - e^{-(\beta_i + \lambda_i)(T-t)}] - (\beta_i + \lambda_i) e^{-(\beta_i + \lambda_i)(T-t)} [1 - e^{-r(T-t)}]}{r + \beta_i + \lambda_i} \frac{u_i(a)}{r} \\ &+ \int_0^{T-t} \lambda_i e^{-(r+\beta_i+\lambda_i)y} \left[\sum_{j \in \mathbb{N}} \pi_j \int V_{t+y}^i(a-z) dG_{ij}(z) \right] dy. \end{aligned}$$

The third term on the right side of $V_t^i(a)$ can be written as

$$\begin{aligned} V_t^i(a) &= \mathbb{E}_t \left\{ \mathbb{I}_{\{\tau(\beta_i) < \tau(\lambda_i)\}} \mathbb{I}_{\{\tau(\beta_i) < T-t\}} \left[\int_t^{t+\tau(\beta_i)} e^{-r(s-t)} u_i(a) ds \right. \right. \\ &\quad \left. \left. + e^{-r\tau(\beta_i)} \sum_{j \in \mathbb{N}} \sigma_j \int \left[V_{t+\tau(\beta_i)}^i(a) + \theta_{ij} \max_{b \in \mathbb{R}} S_{t+\tau(\beta_i)}^{ij}(a, \tilde{a}, b) \right] dF_t^j(\tilde{a}) \right] \right\} \\ &= \frac{\beta_i}{\beta_i + \lambda_i} \frac{r [1 - e^{-(\beta_i + \lambda_i)(T-t)}] - (\beta_i + \lambda_i) e^{-(\beta_i + \lambda_i)(T-t)} [1 - e^{-r(T-t)}]}{r + \beta_i + \lambda_i} \frac{u_i(a)}{r} \\ &+ \int_0^{T-t} \beta_i e^{-(r+\beta_i+\lambda_i)z} \left[\sum_{j \in \mathbb{N}} \sigma_j \int \left[V_{t+z}^i(a) + \theta_{ij} \max_{b \in \mathbb{R}} S_{t+z}^{ij}(a, \tilde{a}, b) \right] dF_{t+z}^j(\tilde{a}) \right] dz. \end{aligned}$$

Thus, we can write

$$\begin{aligned} V_t^i(a) &= \left[1 - e^{-(r+\beta_i+\lambda_i)(T-t)} \right] \frac{u_i(a)}{r + \beta_i + \lambda_i} + e^{-(r+\beta_i+\lambda_i)(T-t)} U_i(a) \\ &+ \lambda_i \int_0^{T-t} e^{-(r+\beta_i+\lambda_i)y} \left[\sum_{j \in \mathbb{N}} \pi_j \int V_{t+y}^i(a-s) dG_{ij}(s) \right] dy \\ &+ \beta_i \int_0^{T-t} e^{-(r+\beta_i+\lambda_i)z} \left[\sum_{j \in \mathbb{N}} \sigma_j \int \left[V_{t+z}^i(a) + \theta_{ij} \max_{b \in \mathbb{R}} S_{t+z}^{ij}(a, \tilde{a}, b) \right] dF_{t+z}^j(\tilde{a}) \right] dz. \end{aligned}$$

With a change of variables in the integrals with respect to time,

$$\begin{aligned} V_t^i(a) &= \left[1 - e^{-(r+\beta_i+\lambda_i)(T-t)} \right] \frac{u_i(a)}{r + \beta_i + \lambda_i} + e^{-(r+\beta_i+\lambda_i)(T-t)} U_i(a) \\ &+ \lambda_i \int_t^T e^{-(r+\beta_i+\lambda_i)(\tau-t)} \left[\sum_{j \in \mathbb{N}} \pi_j \int V_\tau^i(a-z) dG_{ij}(z) \right] d\tau \\ &+ \beta_i \int_t^T e^{-(r+\beta_i+\lambda_i)(\tau-t)} \left[V_\tau^i(a) + \sum_{j \in \mathbb{N}} \sigma_j \theta_{ij} \int \max_{b \in \mathbb{R}} S_\tau^{ij}(a, \tilde{a}, b) dF_\tau^j(\tilde{a}) \right] d\tau. \quad (17) \end{aligned}$$

To obtain (3), simply differentiate (16) with respect to t . ■

A.2 Extension: regulatory borrowing costs

In this section we generalize the theory to allow for proportional borrowing costs to proxy for the effects of regulatory constraints that affect banks' incentives to buy fed funds. Let

$$\Gamma_t^i(a, b, R) \equiv V_t^i(a - b) - V_t^i(a) + [1 + \mathbb{I}_{\{b < 0\}} \kappa_i] e^{-r(\bar{T}-t)} R \quad (18)$$

denote payoff of a bank of type $i \in \mathbb{N}$, with pre-trade balance a , that at time t sells a loan of size b in exchange for a repayment of size R delivered at time \bar{T} , with $\kappa_i \in \mathbb{R}_+$. Intuitively, if $b, R \in \mathbb{R}_+$, then the bank is “selling fed funds” (i.e., lending) and the gain from trade is as in Section 2.2. Conversely, if $b, R \in \mathbb{R}_-$, then the bank is “buying fed funds” (i.e., *borrowing*), and κ_i captures the effects of policies that increase the shadow cost of the bank's liabilities. In all our calibrations we set κ_G large enough to make our theory consistent with the fact that the business model of a GSE consists of lending, but not borrowing in the fed funds market. In our 2019 calibration we use κ_i for $i \in \{F, M, S\}$ to capture the effects of the prudential liquidity regulations discussed in Appendix B (Section B.2). With borrowing costs, the bargaining outcome at time t between two banks of type i and j , with respective balances a and \tilde{a} , denoted $(b_t^{ij}(a, \tilde{a}), R_t^{ji}(\tilde{a}, a))$, is the solution to

$$\max_{(b, R) \in \mathbb{R} \times \mathbb{R}} \Gamma_t^i(a, b, R)^{\theta_{ij}} \Gamma_t^j(\tilde{a}, -b, -R)^{\theta_{ji}}. \quad (19)$$

The correspondig first-order condition with respect to R is

$$\theta_{ij} [1 + \mathbb{I}_{\{b < 0\}} \kappa_i] \Gamma_t^j(\tilde{a}, -b, -R) = \theta_{ji} [1 + \mathbb{I}_{\{0 < b\}} \kappa_j] \Gamma_t^i(a, b, R),$$

which implies $R_t^{ji}(\tilde{a}, a)$ is given by

$$\begin{aligned} e^{-r(\bar{T}-t)} R_t^{ji}(\tilde{a}, a) &= \frac{\theta_{ij}}{1 + \mathbb{I}_{\{0 < b_t^{ij}(a, \tilde{a})\}} \kappa_j} [V_t^j(\tilde{a} + b_t^{ij}(a, \tilde{a})) - V_t^j(\tilde{a})] \\ &\quad + \frac{\theta_{ji}}{1 + \mathbb{I}_{\{b_t^{ij}(a, \tilde{a}) < 0\}} \kappa_i} [V_t^i(a) - V_t^i(a - b_t^{ij}(a, \tilde{a}))], \end{aligned} \quad (20)$$

and

$$b_t^{ij}(a, \tilde{a}) \in \arg \max_{b \in \mathbb{R}} \hat{S}_t^{ij}(a, \tilde{a}, b), \quad (21)$$

where

$$\hat{S}_t^{ij}(a, \tilde{a}, b) \equiv \hat{\Gamma}_t^{ij}(a, \tilde{a}, b)^{\theta_{ij}} \hat{\Gamma}_t^{ji}(\tilde{a}, a, -b)^{\theta_{ji}},$$

with

$$\begin{aligned}\hat{\Gamma}_t^{ij}(a, \tilde{a}, b) &\equiv \theta_{ij} \left\{ S_t^{ij}(a, \tilde{a}, b) - \frac{\mathbb{I}_{\{0 < b\}} \kappa_j - \mathbb{I}_{\{b < 0\}} \kappa_i}{1 + \mathbb{I}_{\{0 < b\}} \kappa_j} [V_t^j(\tilde{a} + b) - V_t^j(\tilde{a})] \right\} \\ \hat{\Gamma}_t^{ji}(\tilde{a}, a, -b) &\equiv \theta_{ji} \left\{ S_t^{ij}(a, \tilde{a}, b) - \frac{\mathbb{I}_{\{0 < b\}} \kappa_j - \mathbb{I}_{\{b < 0\}} \kappa_i}{1 + \mathbb{I}_{\{b < 0\}} \kappa_i} [V_t^i(a) - V_t^i(a - b)] \right\}.\end{aligned}$$

In summary, the bargaining solution, $(b_t^{ij}(a, \tilde{a}), R_t^{ji}(\tilde{a}, a))$, is given by (21) and (20), and the value function $V_t^i(a)$ now satisfies

$$\begin{aligned}rV_t^i(a) - \dot{V}_t^i(a) &= u_i(a) + \lambda_i \sum_{j \in \mathbb{N}} \pi_j \int [V_t^i(a - z) - V_t^i(a)] dG_t^{ij}(z) \\ &\quad + \beta_i \sum_{j \in \mathbb{N}} \sigma_j \int \Gamma_t^i(a, b_t^{ij}(a, \tilde{a}), R_t^{ji}(\tilde{a}, a)) dF_t^j(\tilde{a}),\end{aligned}\tag{22}$$

with Γ_t^i as defined in (18). Notice that (21), (20), and (22) generalize (1), (2), and (3), respectively (and the former reduce to the latter if $\kappa_i = 0$ for all $i \in \mathbb{N}$).

B Institutional background and regulation

In this section we review three financial regulations that affect banks' incentives to borrow and lend in the fed funds market. Two of them directly increase a bank's shadow value of holding reserves by imposing regulatory balance-sheet constraints that can be satisfied with reserve balances (traditional reserve requirements, discussed in Section B.1, and the *Liquidity Coverage Ratio*, discussed in Section B.2.1). The third, is a leverage constraint that increases a bank's shadow cost of all borrowing, including fed funds purchases (the *Supplementary Leverage Ratio*, discussed in Section B.2.2).

B.1 Traditional reserve requirements (Regulation D)

Reserve requirements have been a part of the financial landscape in the United States since before the Federal Reserve Act of 1913 that created the system of Reserve Banks.⁹⁷ Regulation D ("Reserve Requirements for Depository Institutions") is the Federal Reserve regulation

⁹⁷Reserve requirements at the national level were first established with the passage of the National Bank Act in 1863. In the original Federal Reserve Act of 1913, for example, banks were required to hold in reserve different percentages of their demand deposits, depending on whether they were classified as *central reserve city banks* (18 percent), *reserve city banks* (15 percent), or *country banks* (12 percent). See Feinman (1993) for more background and references on the history of reserve requirements in the United States.

that stipulates reserve requirements for depository institutions (i.e., commercial banks, savings banks, thrift institutions, credit unions, and agencies and branches of foreign banks located in the United States).

Until March 2020, Regulation D required depository institutions to keep a minimum amount of reserves against their transaction accounts (such as demand deposits).⁹⁸ This reserve requirement was 0%, 3%, or 10% of transaction account deposits depending on the size of the bank's reservable liabilities.⁹⁹ Institutions had to satisfy reserve requirements by holding cash in their vaults or as a balance in the institution's account at the Federal Reserve Bank in the Federal Reserve District in which the institution is located (either an account of the institution or an account of the institution's Federal Reserve pass-through correspondent).

Reserve requirements were calculated based on a bank's deposit accounts during *computation periods* that depended on the frequency (either weekly or quarterly) with which an institution files an FR 2900 report.¹⁰⁰ Each reserve computation period was used to calculate the reserve requirement that a bank had to satisfy on a lagged basis, i.e., during a 14-day (*reserve*) *maintenance period* in the future.

For institutions that file the FR 2900 report weekly, a (*FR 2900*) *reporting period* is one week long, covering the seven consecutive calendar days beginning on a Tuesday and ending on the following Monday. The *computation period* for weekly reporters consisted of two *reporting periods*, i.e., 14 consecutive days beginning on a Tuesday and ending on the second Monday thereafter. A *maintenance period* consisted of 14 consecutive days beginning on a Thursday and ending on the second Wednesday thereafter. Each reserve computation period was used to calculate the reserve requirement that a bank had to satisfy on a lagged basis: The reserve balance requirement that had to be satisfied during a maintenance period was based on the average level of net transaction accounts and vault cash held during the computation period that had ended 17 days earlier.¹⁰¹

Federal Reserve Banks were authorized to assess charges for deficiencies at a rate of 1 percentage point per year above the primary credit rate in effect for borrowings from the Federal Reserve Bank on the first day of the calendar month in which the deficiencies occurred. Charges were assessed on the basis of daily average deficiencies during each maintenance period.

⁹⁸There was an explicit exemption from Regulation D for bank obligations in nondeposit form to another bank, which included "federal funds purchased".

⁹⁹The Federal Reserve Board reduced all reserve requirement ratios to 0% effective March 26, 2020.

¹⁰⁰This report collects information on select deposits and vault cash from depository institutions.

¹⁰¹See Federal Reserve Board (2019a) for details.

B.2 Post-GFC regulation

In the years following the Great Financial Crisis (GFC), the Federal Reserve Board (FRB), the Federal Deposit Insurance Corporation (FDIC), and the Office of the Comptroller of the Currency (OCC) implemented versions of two regulations agreed to by the Basel Committee on Banking Supervision (BCBS), and consistent with the Dodd-Frank Wall Street Reform and Consumer Protection Act: The *Liquidity Coverage Ratio* (LCR), a prudential liquidity standard, and the *Supplementary Leverage Ratio* (SLR), a prudential leverage standard. Both affect banks' payoffs from trading in the fed funds market. We discuss each in turn.

B.2.1 Liquidity Coverage Ratio (LCR)

The first objective of the Basel III accord agreed upon by the members of the Basel Committee on Banking Supervision (BCBS) is to promote the short-term resilience of the liquidity risk profile of banks. The BCBS developed the LCR to achieve this objective.¹⁰² Specifically, the LCR is designed to ensure that a bank maintains an adequate level of unencumbered, *High Quality Liquid Assets* (HQLA) that can be converted into cash to meet its liquidity needs for a 30-calendar-day time horizon under a liquidity stress scenario specified by supervisors.

The LCR is defined as

$$LCR \equiv \frac{H}{L}, \quad (23)$$

where H denotes HQLA, and L is a measure of total net cash outflows in a 30-day standardized stress scenario. The HQLA consist of Level 1 assets and Level 2 assets. Level 1 assets, which are not subject to haircuts or quantitative caps, include reserves in excess of Regulation D held at a Federal Reserve Bank, as well as securities issued or guaranteed by the U.S. Treasury. Level 2 assets are subject to prescribed haircuts and are capped at no more than 40% of a banking organization's total HQLA.¹⁰³ For our purposes, we can think of H as consisting of two components: (i) reserves, denoted Q_0 , minus Regulation D required reserves, denoted R_D ;

¹⁰²See Basel Committee on Banking Supervision (2010) for more details on the rationale for the regulation.

¹⁰³Level 2 assets are further divided into Level 2A and Level 2B assets. Level 2A assets, which are subject to a 15% haircut, include claims on or guaranteed by a U.S. government-sponsored enterprise (GSE) such as Fannie Mae and Freddie Mac. Level 2B assets, which are subject to a 50% haircut and are capped at no more than 15% of a banking organization's total HQLA, include certain corporate debt securities issued by non-financial companies, and certain publicly traded common equities issued by non-financial companies that are included in the Russell 1000 Index or a foreign equivalent index for shares held in foreign jurisdictions.

and (ii) the value (net of haircut) of all other assets that qualify as HQLA, denoted A , i.e.,

$$H = A + M,$$

where

$$M \equiv \max(Q_1, 0) \tag{24}$$

and $Q_1 \equiv Q_0 - R_D$ denotes the quantity of reserves in excess of the Regulation D requirement. The “max” in (24) reflects that only reserves in excess of Regulation D qualify as HQLA.

Banks report H and L , and these reports are publicly available at a quarterly frequency.¹⁰⁴ Given H , since we have independent information on Q_0 and R_D (and therefore M), we can infer A . The LCR regulation requires

$$1 \leq LCR \tag{25}$$

daily, or monthly, depending on the size and other characteristics of the bank.¹⁰⁵ For our purposes, the key implication of the policy constraint (25) is that it may cause a bank to treat certain holdings of HQLA as required to comply with the LCR regulation. By this we mean that the LCR constraint may cause the bank to impute an additional shadow cost of reducing its holdings of HQLA on a typical day—including reserve balances. In the specific case of reserve balances, the bank may impute an additional shadow cost of selling fed funds, since this may drive the bank’s reserves (net of the Regulation D requirement) below the level of reserves that the bank routinely allocates to comply with the LCR regulation. Thus, in practice, banks may regard some of the reserves in excess of the Regulation D requirement as being “required” to satisfy the LCR constraint. The fact that the LCR regulation allows for substitutability among the HQLA in the numerator of the left side of (25) presents us with an identification challenge when trying to estimate the share of a bank’s reserve balances in excess of the Regulation D requirement that the bank treats as “required” to satisfy the LCR constraint. Next, we formalize this identification problem, and describe how we address it.

For each bank, we observe H , M , and A . We want to express M as the sum of a component, \hat{M}^R , that represents the quantity of reserves (in excess of the Regulation D requirement) that

¹⁰⁴E.g., from the S&P Global Capital IQ database. See Appendix D (Section D.1.3) for details.

¹⁰⁵Relatively large institutions regulated by the FRB must calculate and maintain a liquidity coverage ratio that is equal to or greater than 1 on each business day (or, in the case of a smaller FRB-regulated institutions, on the last business day of the applicable month). The LCR rule is codified at 12 CFR part 50 (OCC), 12 CFR part 249 (FRB), and 12 CFR part 329 (FDIC).

the bank relies on to comply with the LCR regulation, and a component, \hat{M}^E , that represents reserves in excess of the Regulation D *and* the LCR requirements. Similarly, a bank may hold HQLA (other than reserves) in excess of what would be necessary to meet the LCR requirement for reasons other than having to comply with the LCR regulation, so we can also decompose A into two (unobserved) components: \hat{A}^R , which represents the value of HQLA (other than reserves in excess of the Regulation D requirement) that the bank regards as being necessary to comply with the LCR regulation, and \hat{A}^E , which represents the value of HQLA (other than reserves in excess of the Regulation D requirement) that the bank regards as being in excess of what is required to meet the LCR regulation. In summary, $\{\hat{M}^j, \hat{A}^j\}_{j \in \{R, E\}}$ satisfy:

$$M = \hat{M}^R + \hat{M}^E \quad (26)$$

$$A = \hat{A}^R + \hat{A}^E \quad (27)$$

$$\hat{A}^R + \hat{M}^R \leq L, \text{ with “=” if } L \leq A + M \quad (28)$$

$$\hat{A}^E + \hat{M}^E = 0, \text{ if } A + M < L \quad (29)$$

$$\hat{M}^j, \hat{A}^j \in \mathbb{R}_+ \text{ for } j \in \{R, E\}. \quad (30)$$

We are interested in using the policy constraint (25) along with data on M , A , and L , and (26)-(30), to estimate bank-level bounds for \hat{M}^R .

There are three special cases in which the constraint (25) together with knowledge of M , A , and L , and the definitions (26)-(30) are sufficient to identify \hat{M}^R and \hat{A}^R . First, if a bank has $LCR \leq 1$ (i.e., if it is not complying with the LCR regulation in a given sample period), then the bank is clearly holding no excess HQLA of any type, so $\hat{M}^R = M$, $\hat{A}^R = A$, and $\hat{M}^E = \hat{A}^E = 0$, as implied by (26), (27), (29), and (30). Second, if $LCR \geq 1$ and $Q_1 \leq 0$, then $M = 0$, so the LCR requirement, L , is being satisfied exclusively with HQLA other than reserves, i.e., $\hat{A}^R = L$ and $\hat{A}^E = A - L$, with $\hat{M}^R = \hat{M}^E = 0$, as implied by (26), (27), (28), and (30). Third, if $LCR \geq 1$ and $A = 0$, then the LCR requirement, L , is being satisfied exclusively with reserves, M , i.e., $\hat{M}^R = L$ and $\hat{M}^E = M - L$, with $\hat{A}^R = \hat{A}^E = 0$, as implied by (26), (27), (28), and (30).

In practice, most banks satisfy the LCR constraint (25) with $\min(M, A) \geq 0$, and for such banks it is not obvious how to decompose the level of *required* HQLA, i.e., L , into the two unobserved components, \hat{M}^R and \hat{A}^R . However, notice that conditions (26)-(30) imply \hat{M}^R

must satisfy the following bounds:

$$\hat{M}^R \begin{cases} = M & \text{if } A + M < L \\ \in [\max(0, L - A), \min(L, M)] & \text{if } L \leq A + M. \end{cases} \quad (31)$$

We can write (31) as

$$\hat{M}^R = \begin{cases} M & \text{if } A + M < L \\ \rho \min(L, M) + (1 - \rho) \max(0, L - A) & \text{if } L \leq A + M, \end{cases} \quad (32)$$

for some $\rho \in [0, 1]$. For a given ρ , (26)-(30), and (32) imply

$$\hat{A}^R = \begin{cases} A & \text{if } A + M < L \\ (1 - \rho) \min(L, A) + \rho \max(0, L - M) & \text{if } L \leq A + M, \end{cases}$$

and given \hat{M}^R and \hat{A}^R , \hat{M}^E and \hat{A}^E are implied by (26) and (27).

The parameter $\rho \in [0, 1]$ represents the bank's (unobserved) preference for satisfying the LCR requirement, L , with reserves (rather than with other HQLA). For example, if $\rho = 1$, the bank has a strong preference for satisfying the LCR with reserves, and this will reduce the bank's willingness to lend reserves in the fed funds market. If $\rho = 0$, the bank has a strong preference for satisfying the LCR with HQLA *other* than reserves, and will be less constrained by its reserve balance when trading in the fed funds market.

According to elementary theory, the quantity of reserves in excess of regulatory reserve requirements is a key determinant of a bank's "fundamental" incentive to borrow and lend in the fed funds market. For example, a bank whose reserve balance is lower than the minimum regulatory requirement, has a fundamental incentive to borrow (at a rate no larger than the shadow cost of violating the regulatory requirement). Conversely, a bank whose reserve balance is higher than the regulatory requirement, would have, all else equal, an incentive to lend (e.g., to banks with negative excess reserves, at a rate between the lender's and the borrower's respective shadow prices of reserves). For this reason, it is important to impute an accurate notion of "excess reserves" in any empirical implementation of a theory of interbank loans.

The traditional definition of "excess reserves", which only subtracts the Regulation D requirement from the bank's reserve balance is not an adequate notion of excess reserves for institutions that must comply with the LCR regulation.¹⁰⁶ In our empirical and quantitative work we use a more comprehensive notion of "required reserves" that includes not only the level

¹⁰⁶The LCR regulation applies to bank holding companies (BHCs) and savings and loans holding (SLHCs) with at least \$50 bn in total consolidated assets.

of reserves that a bank is required to hold to comply with Regulation D, but also the level of reserves that the bank holds toward meeting the LCR requirement. Specifically, our benchmark definition of “excess reserves” for any bank that is subject to, and satisfies the LCR constraint (25), is $Q_2 \equiv Q_1 - R_L$, where $R_L \equiv \max(0, L - A)$. In other words, to construct our preferred notion of excess reserves, we start from the traditional notion of reserves in excess of the Regulation D requirement, Q_1 , and subtract the minimum level of reserves needed to comply with the LCR requirement, i.e., R_L .¹⁰⁷ Notice that our measure of excess reserves coincides with the traditional measure for a bank that has enough HQLA other than reserves to meet the LCR requirement, i.e., $Q_1 - Q_2 = R_L = 0$ if $L \leq A$. But our measure of excess reserves is lower than the traditional measure for a bank whose holdings of HQLA other than reserves are insufficient to meet the LCR requirement, i.e., if $A < L$, then $0 < Q_1 - Q_2 = R_L = L - A$.

B.2.2 Supplementary Leverage Ratio (SLR)

The SLR is the U.S. banking agencies’ implementation of the “Basel III Tier 1 Leverage Ratio”, which is defined as

$$SLR \equiv \frac{\text{Tier 1 Capital}}{\text{Total Leverage Exposure}}. \quad (33)$$

The numerator (defined in U.S. Basel III) includes common stock and retained earnings. The denominator is a comprehensive measure of assets, composed of four elements: (1) on-balance sheet assets, (2) derivative exposures, (3) repo-style transaction exposures, and (4) other off-balance sheet exposures. The SLR regulation requires a bank to maintain an SLR above a threshold; specifically, either $SLR \geq 0.03$, or $SLR \geq 0.05$.¹⁰⁸

B.2.3 Resolution Planning

In the aftermath of the GFC, regulatory authorities started requiring large “systemically important” financial institutions (e.g., BHCs with total consolidated assets of \$50 bn or more) to periodically submit a resolution plan (also known as “living will”) to the Federal Reserve

¹⁰⁷From (32), we see that R_L is the same as \hat{M}^R when $\rho = 0$ (in the empirically relevant case with $L \leq A + M$). In this sense, our preferred notion of excess reserves selects the largest level of excess reserves that is consistent with the LCR constraint, (25).

¹⁰⁸The threshold equals 3% for *advanced approaches firms*, which include state banks, savings associations, bank holding companies (BHCs), and saving and loan holding companies (SLHCs) with more than \$250 bn in total consolidated assets, or more than \$10 bn of on-balance sheet foreign exposures. The threshold equals 5% for the 8 US bank-holding companies that have been identified by the Financial Stability Board as global systemically important banks (and their U.S. insured depository institution subsidiaries).

and the Federal Deposit Insurance Corporation. A resolution plan describes in some detail the company’s strategy for rapid and orderly resolution in the event of material financial distress or failure of the company.

B.2.4 Effects of LCR and SLR regulation on fed funds trading incentives

In this section we discuss the effects of the LCR and SLR regulation on banks’ incentives to borrow and lend in the fed funds market.

First, we consider the effect of LCR regulation on banks’ incentives to borrow and lend in the fed funds market. Reserves appear (with weight =1) in the numerator of the LCR in (23), and overnight fed fund liabilities appear in the denominator (also with weight =1). Consider a bank that borrows ℓ in the fed funds market. The LCR before the trade is $\frac{H}{L}$ and after the trade it is $\frac{H+\ell}{L+\ell}$. Since

$$\frac{\partial}{\partial \ell} \left(\frac{H + \ell}{L + \ell} \right) = \frac{L - H}{(L + \ell)^2},$$

it follows that the trade does not affect the LCR if the bank is satisfying it exactly pre-trade (i.e., if $LCR \equiv \frac{H}{L} = 1$), increases the LCR if the borrowing bank is below the LCR target pre trade (i.e., if $LCR \equiv \frac{H}{L} < 1$), and decreases the LCR if the borrowing bank is above the LCR target pre trade (i.e., if $LCR \equiv \frac{H}{L} > 1$). For a bank that lends ℓ in the fed funds market, the LCR before the trade is $\frac{H}{L}$ and after the trade it is $\frac{H-\ell}{L}$.¹⁰⁹ Hence, selling fed funds unambiguously reduces the LCR. To summarize, LCR regulation increases the shadow cost of selling fed funds (because lending reserves tightens the LCR constraints of lenders), and increases the shadow cost of borrowing for banks whose LCR constraints are slack at the time of the trade (because borrowing reserves tightens the LCR constraints of such banks).

Second, we consider the effect of SLR regulation on banks’ incentives to borrow and lend in the fed funds market. Let \mathcal{A} denote assets, \mathcal{L} denote liabilities, and $\mathcal{C} \equiv \mathcal{A} - \mathcal{L}$ denote capital. Then, we can write (33) as

$$SLR \equiv \frac{\mathcal{C}}{\mathcal{A}} = \frac{\mathcal{A} - \mathcal{L}}{\mathcal{A}}. \quad (34)$$

Notice that *lending* in the fed funds market does not change the SLR because the bank that acts as a lender is just exchanging reserves for an overnight credit of reserves, which leaves both \mathcal{L} and \mathcal{A} unchanged. However, *borrowing* in the fed funds market reduces the SLR, since

¹⁰⁹The quantity of reserves sold, ℓ , is subtracted from the HQLA of the lender, but the corresponding fed funds credit is not added to the total of HQLA of the lender because fed funds not qualify as a HQLA.

borrowing ℓ dollars worth of reserves increases liabilities from \mathcal{L} to $\mathcal{L} + \ell$, and increases assets from \mathcal{A} to $\mathcal{A} + \ell$, and therefore the *SLR* is *reduced* from $\frac{\mathcal{A}-\mathcal{L}}{\mathcal{A}}$ to $\frac{\mathcal{A}-\mathcal{L}}{\mathcal{A}+\ell}$. To summarize, SLR regulation has no effect on the shadow cost of lending fed funds (because lending reserves does not alter the SLR constraint), but increases the shadow cost of buying fed funds (because borrowing reserves tightens the SLR constraint of solvent banks).

C Computation

In this section we discuss computational issues. Section C.1 outlines the solution algorithm. Section C.2 explains how we compute, in the quantitative theory, the statistics that we compare with their empirical counterparts.

C.1 Solution algorithm

The steps we use to solve for the equilibrium of the model are as follows.

Step 0: Set grids. We think of the time interval $[0, T]$ as corresponding to a trading day in the fed-funds market, which consists of 9.5 hours (from 9.00 AM to 5.30 PM). We divide the interval $[0, T]$ into $N_T + 1$ periods, denoted $t = 0, 1, \dots, N_T$, and set $N_T = 799$. As we have 800 periods, each period represents approximately 42 seconds (i.e., $\frac{9.5 \times 60 \times 60}{800} = 42.75$ seconds).

For each bank type $i \in \mathbb{N}$, we construct an equally spaced grid for reserve balances, $\mathbb{A}^i = \{a_1^i, a_2^i, \dots, a_{N_a}^i\}$, with $N_a = 150$. We interpret each unit of reserves in the model as corresponding to \$10 bn in the data. For the benchmark years 2017 and 2019, we set a_1^i and $a_{N_a}^i$ equal to the 0.5th and 99.5th percentiles of the kernel estimate of the beginning-of-day distributions, respectively (see Section 3.3). We use the interpolation procedure explained in Section 3.6 to construct grids whenever we change the total quantity of balances. In all cases we add 5 additional points to the grid, $\{-0.2, -0.1, 0, 0.1, 0.2\}$.¹¹⁰

For each pair of bank types, $i, j \in \mathbb{N}$, we construct a grid for payment sizes, $\mathbb{Z}^{ij} = \{z_1^{ij}, z_2^{ij}, \dots, z_{N_z}^{ij}\}$, with $N_z = 35$. The probability mass function for payment sizes, $\{G_{ij}(z)\}_{z \in \mathbb{Z}^{ij}}$, is constructed with the procedure described in Section 3.2.

¹¹⁰We add these grid points because the value functions are numerically close to having a kink around $a = 0$ towards the end of the trading day (i.e., as t gets closer to N_T).

Step 1: Guess the distribution of balances. For each $a \in \mathbb{A}^i$, let $f_t^i(a)$ be the fraction of banks of type $i \in \mathbb{N}$ that hold a quantity of reserves equal to a at the beginning of period t , with $\sum_{a \in \mathbb{A}^i} f_t^i(a) = 1$. The beginning-of-day distribution, $f_0^i(\cdot)$, is given since $F_0^i(a) \equiv \sum_{x \in \mathbb{A}^i} f_0^i(x) \mathbb{I}_{\{x \leq a\}}$ is estimated from the data with the procedure described in Section 3.3. Guess the distributions $\{f_t^i(a)\}_{a \in \mathbb{A}^i, i \in \mathbb{N}}$ for each $t \in \{1, 2, \dots, N_T\}$.

Step 2: Compute the value function. Since for each $i \in \mathbb{N}$ and $a \in \mathbb{A}^i$ we have the terminal condition, $V_{N_T}^i(a) = U_i(a)$, where $U_i(\cdot)$ is the exogenous end-of-day payoff function, we can then solve backwards for the value function, $\{V_t^i(a)\}_{a \in \mathbb{A}^i, i \in \mathbb{N}, t \in \{0, \dots, N_T-1\}}$. Each of these backward iterations between period $t \in \{N_T, \dots, 1\}$ and period $t-1$ consists of two steps. In the first step, for each pair of bank types $i, j \in \mathbb{N}$, we compute the bargaining outcomes, $\{b_t^{ij}(a^i, a^j), R_t^{ji}(a^j, a^i)\}_{a^i \in \mathbb{A}^i, a^j \in \mathbb{A}^j}$, taking $\{V_t^i(a)\}_{a \in \mathbb{A}^i, i \in \mathbb{N}}$ as given. In the second step we solve for the value function backwards, i.e., we solve for $\{V_{t-1}^i(a)\}_{a \in \mathbb{A}^i, i \in \mathbb{N}}$ given the one-period-ahead bargaining outcomes and values, i.e., $\{b_t^{ij}(a^i, a^j), R_t^{ji}(a^j, a^i), V_t^i(a^i)\}_{a^i \in \mathbb{A}^i, a^j \in \mathbb{A}^j, (i,j) \in \mathbb{N}^2}$. Next, we explain these two steps in detail.

Step 2.1: Solve for $b_t^{ij}(\cdot, \cdot)$ and $R_t^{ji}(\cdot, \cdot)$. Given the values $\{V_t^i(\cdot)\}_{i \in \mathbb{N}}$, we compute the bargaining outcome for the loan size, $b_t^{ij}(a, \tilde{a})$, as in (21), which can be written as:

$$b_t^{ij}(a, \tilde{a}) = \arg \max_b \left\{ \frac{1}{1 + \mathbb{I}_{\{b < 0\}} \kappa_i} V_t^i(a - b) + \frac{1}{1 + \mathbb{I}_{\{0 < b\}} \kappa_j} V_t^j(\tilde{a} + b) - \epsilon |b| \right\}, \quad (35)$$

where $\epsilon = 1e-9$ is a small trading cost introduced to rule out loans with negligible gains from trade. Since a unit of reserves in the model corresponds to \$10 bn in the data, the value of ϵ implies a trading cost of \$10 for a loan of size \$1 bn. We use a Golden search routine to solve for $b_t^{ij}(a, \tilde{a})$ in (35), for each $a \in \mathbb{A}^i$, $\tilde{a} \in \mathbb{A}^j$, and $(i, j) \in \mathbb{N} \times \mathbb{N}$. Given the bargained loan sizes, $\{b_t^{ij}(a^i, a^j)\}_{a^i \in \mathbb{A}^i, a^j \in \mathbb{A}^j, (i,j) \in \mathbb{N}^2}$, we can compute the associated repayments, $R_t^{ji}(a^j, a^i)$, as in (20), and the gain from trade to the bank of type i and balance a^i , i.e., $\Gamma_t^i(a, b_t^{ij}(a^i, a^j), R_t^{ji}(a^j, a^i))$, as in (18).

Step 2.2: Solve for $V_t^i(a)$ backwards. We divide each period into two stages. Random payments between pairs of banks take place in the first stage. Trade between pairs of banks takes place in the second stage. The first stage is divided further into N_S subperiods, each indexed by $s \in \{1, 2, \dots, N_S\}$ with $N_S = 42$, so each of these subperiods corresponds to 1 second (since each full model period corresponds to approximately 42 seconds). We solve for the value function within each model period backwards: we start by solving for the value of

trade decisions in the second stage, and then integrate the value of the payment shocks in the 42 subperiods of the first stage.

Let $\hat{V}_t^i(a)$ be the value of a bank of type $i \in \mathbb{N}$ with balance $a \in \mathbb{A}^i$ at the beginning of the second stage of period t . This value satisfies

$$(1 + \Delta r)\hat{V}_t^i(a) = \Delta u_i(a) + \Delta \beta_i \sum_{j \in \mathbb{N}} \sigma_j \sum_{\tilde{a} \in \mathbb{A}_j} \bar{\Gamma}_{t+1}^{ij}(a, \tilde{a}) f_{t+1}^j(\tilde{a}) + V_{t+1}^i(a), \quad (36)$$

where $\bar{\Gamma}_t^{ij}(a, \tilde{a}) \equiv \Gamma_t^i(a, b_t^{ij}(a, \tilde{a}), R_t^{jj}(\tilde{a}, a))$, and $\Delta = 1/[N_S(N_T + 1)]$ is the size of the time interval (including all trade and payment periods in the day). Let $\tilde{V}_{t,s}^i(a)$ be the value of a bank of type $i \in \mathbb{N}$ with balance $a \in \mathbb{A}^i$ at the beginning of subperiod s of the first stage of period t . This value satisfies

$$(1 + \Delta r)\tilde{V}_{t,s}^i(a) = \Delta u_i(a) + \Delta \lambda_i \sum_{j \in \mathbb{N}} \pi_j \sum_{z \in \mathbb{Z}^{ij}} [V_t^i(a - z) - V_t^i(a)] G_{ij}(z) + \tilde{V}_{t,s+1}^i(a), \quad (37)$$

for $s = 1, \dots, N_S$, with boundary conditions $\tilde{V}_{t,N_S+1}^i(a) = \hat{V}_t^i(a)$, and $\tilde{V}_{t,1}^i(a) = V_t^i(a)$. Equations (36) and (37) are the discrete-time approximations to the Bellman equation (22).

We solve (36)-(37) backwards, as follows. Given $V_{t+1}^i(\cdot)$ (recall the terminal condition $V_{N_T}^i(\cdot) = U_i(\cdot)$), we compute $\bar{\Gamma}_{t+1}^{ij}(\cdot, \cdot)$, and given our guess of $\{f_{t+1}^i(\cdot)\}_{i \in \mathbb{N}}$, we compute $\hat{V}_t^i(\cdot)$ using (36). We then compute $\{\tilde{V}_{t,s}^i(\cdot)\}_{s \in \{1, 2, \dots, N_S\}}$ using (37) and the terminal condition $\tilde{V}_{t,N_S+1}^i(\cdot) = \hat{V}_t^i(\cdot)$ by iterating backwards, which delivers $V_t^i(\cdot) = \tilde{V}_{t,1}^i(\cdot)$.

Step 3: Compute the implied distribution of balances Given the negotiated loan sizes, $b_t^{ij}(\cdot, \cdot)$ and the distribution of random payments, we can solve for the distribution of balances forward from an initial condition, $f_0^i(a)$. As in **step 2**, we need to compute the evolution of balances for the two within-period stages (the payments stage, and the trading stage). Since we are solving for the distribution of reserves forward, we start with the first stage and then move to the second stage.

Let $\tilde{f}_{t,s}^{i,\text{new}}(a_m)$ be the fraction of banks of type $i \in \mathbb{N}$ with balance $a_m \in \mathbb{A}^i$, at the beginning of subperiod s of the first stage of period t . We use the superscript “new” to emphasize that this is the new distribution implied by the bargaining outcomes in **step 2** (rather than the distribution that was used to derive those outcomes). Then,

$$\tilde{f}_{t,s}^{i,\text{new}}(a_m) = (1 - \Delta \lambda_i) \tilde{f}_{t,s-1}^{i,\text{new}}(a_m) + \Delta \lambda_i \sum_{j \in \mathbb{N}} \sum_{a \in \mathbb{A}^i} \sum_{z \in \mathbb{Z}^{ij}} \pi_j \mathbb{L}(a_m, a - z) G_{ij}(z) \tilde{f}_{t,s-1}^{i,\text{new}}(a) \quad (38)$$

for $s = 1, \dots, N_S$, with initial condition $\tilde{f}_{t,0}^{i,\text{new}}(a_m) = f_t^{i,\text{new}}(a_m)$, and where

$$\mathbb{L}(a_m, x) \equiv \mathbb{I}_{\{x \in (a_{m-1}, a_m]\}} \frac{x - a_{m-1}}{a_m - a_{m-1}}$$

implements a linear interpolation. The recursion (38) is initialized with the exogenous time-0 distribution of balances, i.e., $f_0^{i,\text{new}}(a_m) = f_0^i(a_m)$.

Let $\hat{f}_t^{i,\text{new}}(a_m)$ be the fraction of banks of type $i \in \mathbb{N}$ with balance $a_m \in \mathbb{A}^i$ after the trades in the second stage of period t ; it is given by

$$\begin{aligned} \hat{f}_t^{i,\text{new}}(a_m) &= (1 - \Delta\beta_i) \tilde{f}_{t,N_S}^i(a_m) \\ &+ \Delta\beta_i \sum_{j \in \mathbb{N}} \sum_{a \in \mathbb{A}^j} \sum_{\tilde{a} \in \mathbb{A}^j} \sigma_j \mathbb{L}\left(a_m, a - b_t^{ij}(\hat{a}, \tilde{a})\right) \tilde{f}_{t,N_S}^{i,\text{new}}(a) \tilde{f}_{t,N_S}^{j,\text{new}}(\tilde{a}). \end{aligned}$$

Having solved for $\hat{f}_t^{i,\text{new}}(\cdot)$, set $f_{t+1}^{i,\text{new}}(\cdot) = \tilde{f}_{t+1,0}^{i,\text{new}}(\cdot) = \hat{f}_t^{i,\text{new}}(\cdot)$, and move to next period.

Step 4: Check for convergence. We use two criteria for convergence.

Criterion 1. We determine that the algorithm has converged if the probability distribution in **step 1** is close enough to the probability distribution obtained after **step 4**. Specifically, we consider the algorithm has converged if $\mathcal{E}(f) \equiv \max_{a,i,t} |f_t^i(a) - f_t^{i,\text{new}}(a)| < 1e-4$.

Criterion 2. We determine that the algorithm has converged if some key theoretical moments have stabilized across iterations. In particular, we look at convergence in the distribution of interest rates and measures of trading activity.¹¹¹ Specifically, let ρ_t^p denote the p -percentile of the (volume weighted) distribution of interest rates at time t . Every 10 iterations of the algorithm, we compute the rate percentiles ρ_t^p for $p \in \{0.05, 0.10, 0.30, 0.50, 0.70, 0.90, 0.95\}$, and then compute the error $\mathcal{E}(\rho) \equiv \max_{p,t} |\rho_t^p - \rho_t^{p,\text{new}}|$. Every 10 iterations, we also compute: the effective fed-funds rate (EFFR), the participation for each bank, \mathcal{P}_i , and the reallocation for each bank, \mathcal{R}_i . We consider the algorithm has converged if after 10 iterations, we have: (i) $\mathcal{E}(f) < 1e-3$, (ii) $\mathcal{E}(\rho) < 1e-4$, and (iii) the errors for the EFFR, \mathcal{P}_i , and \mathcal{R}_i all below $1e-4$. For all these error computations, we check errors comparing results 10 iterations apart (e.g.: the EFFR this iteration compared with the EFFR 10 iterations ago), which ensure that results are stable across algorithm iterations.

The reason we sometimes use **Criterion 2** is that, despite our using the trading cost ϵ in equation (35), we sometimes observe loans that entail very small gains from trade, but still

¹¹¹See Section C.2 for details on computation of theoretical moments.

affect the distributions $\{f_t^i(\cdot)\}$. These small-surplus trades may keep the error $\mathcal{E}(F)$ above the convergence tolerance, but have no significant effect on the distribution of rates nor in the relevant measures of trading activity. To ensure the algorithm has stabilized, we only start implementing **Criterion 2** after 25 iterations, and check errors $\mathcal{E}(\rho)$, EFFR, \mathcal{P}_i , and \mathcal{R}_i every 10 iterations. We have found that using **Criterion 1** exclusively has no significant effect on the main results, but it typically takes longer for the algorithm to converge.

C.2 Computation of theoretical moments

Many of the statistics that we compute from model output are volume-weighted, which is the standard way many official statistics are computed (e.g., the EFFR). In this section we provide more details on how to perform these calculations in the theory.

Let $\omega_t^{ij}(a, \tilde{a})$ be the share of loans between banks type $i \in \mathbb{N}$ and $j \in \mathbb{N}$ with balances $a \in \mathbb{A}^i$ and $\tilde{a} \in \mathbb{A}_j$ at time t , relative to the total volume of loans in the trading-day, v . That is,

$$\omega_t^{ij}(a, \tilde{a}) = \frac{\tilde{v}_t^{ij}(a, \tilde{a})}{v}$$

where

$$\tilde{v}_t^{ij}(a, \tilde{a}) = (\Delta\beta_i) n_i \sigma_j F_t^i(a) F_t^j(\tilde{a}) |b_t^{ij}(a, \tilde{a})|,$$

and

$$v = \sum_t \sum_{i,j \in \mathbb{N}} \sum_{(a, \tilde{a}) \in \mathbb{A}^i \times \mathbb{A}^j} \tilde{v}_t^{ij}(a, \tilde{a})$$

is the total volume of loans in the trading day.

The EFFR is the volume-weighted mean of all daily traded rates, i.e.,

$$\text{EFFR} = \sum_t \sum_{i,j \in \mathbb{N}} \sum_{(a, \tilde{a}) \in \mathbb{A}^i \times \mathbb{A}^j} \omega_t^{ij}(a, \tilde{a}) \rho_t^{ij}(a, \tilde{a}). \quad (39)$$

Let v_i^e and v_i^r denote the values of all the loans that were extended and received, respectively, throughout the trading day by all banks of type $i \in \mathbb{N}$, i.e.,

$$\begin{aligned} v_i^e &= \sum_t \sum_{i,j \in \mathbb{N}} \sum_{(a, \tilde{a}) \in \mathbb{A}^i \times \mathbb{A}^j} \omega_t^{ij}(a, \tilde{a}) b_t^{ij}(a, \tilde{a}) \mathbb{I}_{\{b_t^{ij}(a, \tilde{a}) > 0\}} \\ v_i^r &= \sum_t \sum_{i,j \in \mathbb{N}} \sum_{(a, \tilde{a}) \in \mathbb{A}^i \times \mathbb{A}^j} \omega_t^{ij}(a, \tilde{a}) b_t^{ij}(a, \tilde{a}) \mathbb{I}_{\{b_t^{ij}(a, \tilde{a}) < 0\}}. \end{aligned}$$

The participation and reallocation measures are $\mathcal{P}_i = (v_i^e + v_i^r)/2v$, and $\mathcal{R}_i = (v_i^e - v_i^r)/(v_i^e + v_i^r)$, respectively.

D Data

This appendix discusses the data we use in the paper. Section D.1 describes the data sources, how we merged them, and our sample selection procedure. Section D.2 describes our own calculations of statistics that we used in the empirical and quantitative sections of the paper. Section D.3 gives a detailed account of the market events of September 13–20, 2019.

D.1 Reserve balances, transfers, and regulatory requirements

We used three databases for bank-level data on: (1) reserve balances and Regulation-D requirements, (2) high-frequency reserve transfers, and (3) Liquidity Coverage Ratio (LCR) requirements. We discuss each below.

D.1.1 Reserve balances and Regulation D

Bank-level end-of day balances at daily frequency were provided by the Monetary Policy Operations and Analysis (MPOA) section of the Monetary Affairs Division at the Federal Reserve Board of Governors. MPOA also supplied us the bank-level data on Regulation-D reserve requirements for each two-week maintenance period. Reserve balances and Regulation-D requirements are reported at the level of the bank holding company (and we used the bank holding company as the relevant unit of observation throughout). MPOA reports end-of-day balances as of 6:30 pm EST. We imputed next-day beginning-of-day balances as of 9:00 am EST with the procedure explained in Section D.2.1.

D.1.2 Reserve transfers

Fedwire Funds Services (*Fedwire*) is an electronic large-value real-time gross settlement system operated by the Federal Reserve Banks. Fedwire participants include commercial banks, savings banks, thrift institutions, credit unions, agencies and branches of foreign banks in the United States, government securities dealers, government agencies such as federal or state governments, and Government Sponsored Enterprises (GSEs, e.g., Freddie Mac, Fannie Mae, and Federal Home Loan Banks). These institutions hold reserve balances in accounts at the Federal Reserve, and use Fedwire to transfer reserves to other participants, e.g., to settle payments, or to lend or repay loans of reserve balances.

Every Fedwire participant is identified by a *Fedwire account number*. Whenever an institution uses multiple Fedwire account numbers, we followed the guidelines from the Reserve Bank

Operations and Payment Systems Division at the Federal Reserve Board for linking those Fedwire account numbers to a single *bank ID*. Whenever institutions with different bank IDs belong to the same bank-holding company, we aggregated them into a single entity (since regulations, e.g., reserve requirements, LCR and SLR requirements, and interest-on-reserves calculations, etc., typically apply at the level of the bank-holding-company level). In a few instances, a bank ID could not be matched to a bank-holding company. Those account numbers were excluded from the sample. We also excluded any bank ID that did not have any fed funds trading activity in a given year. Our sample consists of 754 Fedwire participants for the year 2006, 404 for the year 2014, 395 for the year 2017, and 412 for the year 2019.

Having mapped Fedwire account numbers to bank-holding companies, we assigned the identity of each Fedwire sender or receiver to a bank holding company. We used the output of the Furfine algorithm to identify the set of overnight loans from the universe of Fedwire transfers, and treated the remaining transfers as *payments* unrelated to overnight borrowing and lending. All individual payments with value lower than \$10,000 between a pair of banks during a trading day are consolidated into a single payment.

D.1.3 Liquidity Coverage Ratio

LCR regulation requires a bank to maintain (typically on a daily basis) a quantity of *High Quality Liquid Assets* (HQLA) at least as large as a measure of total net cash outflows in a 30-day standardized stress scenario. If we let $H_m(d)$ denote the quantity of qualifying HQLA held by bank m in a trading period d , and $L_m(d)$ denote the corresponding measure of outflows in the stress scenario, the LCR regulation requires $L_m(d) \leq H_m(d)$.¹¹²

Both, $L_m(d)$ and $H_m(d)$ are made public by each bank at a quarterly frequency. We obtain data on the ratio of these quantities from S&P Global Capital IQ database.¹¹³ We used SNL Classic Data and run a Companies (Classic) screener to search for our data. We extracted quarterly LCR (LIQUIDITY_COV_RATIO) data from 1990Q1 to 2021Q2. For some banks, LCR data were missing in some quarters. For these cases, we obtained the LCR data directly

¹¹²Appendix B (Section B.2.1) describes the LCR regulation in greater detail.

¹¹³The S&P database can be accessed at: <https://www.spglobal.com/marketintelligence/en/>.

from the bank’s website.¹¹⁴

We merged our balances data from MPOA (described in Section D.1.1) with the S&P LCR data using the Replication Server System Database (RSSD) ID. (The balances data from MPOA contains the RSSD of each bank holding company.) We then created a manual cross-walk to match RSSDs to parent company names in the S&P database, using the National Information Center repository from the Federal Financial Institutions Examination Council (<https://www.ffeec.gov/NPW>). We always matched RSSDs to the parent bank holding company to which the LCR regulation applies. In general, this procedure implies matching the RSSDs to the highest level parent company in the corporate structure, except for cases in which the parent company is a sovereign government and the LCR constraint applies to the second highest parent company level.

D.2 Empirical computations

D.2.1 Balances: beginning-of-day imputation

This section provides further details about the construction of beginning-of-day (BOD) balances that we discussed in Section 3.3. The BOD balances used in the paper were obtained from the following three-step procedure for each bank:

- Step 1. We started with the end-of-day (EOD) balance for trading day $d - 1$ obtained from MPOA, and calculated a “basic” measure of the BOD balance for trading day d , by adding (subtracting) the repayments received (sent) corresponding to loans extended (received) during trading day d .
- Step 2. From the “basic” measure of BOD balance calculated in Step 1, we calculated an “adjusted” measure of BOD balance by subtracting the quantity of required reserves, i.e., the minimum level of reserves that the bank must hold during the maintenance period in order to comply with Regulation D and the minimum LCR requirement.

¹¹⁴This was the case for the following three banks:

- Credit Agricole Group (<https://www.credit-agricole.com/en/pdfPreview/186985>)
- DNB ASA (<https://ir.dnb.no/capital-framework>)
- State Street Corporation (<https://investors.statestreet.com/filings-and-reports/u-s-liquidity-coverage-ratio-disclosures/default.aspx>).

- Step 3. From the “adjusted” measure of BOD balance calculated in Step 2, we calculated a measure of “unencumbered” BOD reserve balance for trading day d , by the netting predictable payments that take place during trading day d .

Next, we discuss each step in more detail.

Step 1: Netting repayments of previous-day loans. For each bank holding company m in our sample, we obtained the EOD balance as of 6:30 pm EST of day d from MPOA (see Section D.1.1), which we denote $a_m^{\text{eod}}(d)$. For each bank m , we used the output of the Furfine algorithm to compute the repayments to be sent and received on day d corresponding to loans originated during day $d - 1$. Let $\text{receive}_m(d)$ and $\text{send}_m(d)$ denote the amounts of reserves that bank m will receive or send, respectively, on day d , and define the net repayment corresponding to loans originated during day $d - 1$, as $\text{net}_m(d) \equiv \text{receive}_m(d) - \text{send}_m(d)$. We then computed $a_m(d) = a_m^{\text{eod}}(d - 1) + \text{net}_m(d)$, which is our “basic” measure of BOD balance for bank m on day d . Finally, we computed the BOD “basic” balance for the maintenance period h as the average of $a_m(d)$ for days $d \in h$: $a_m(h) = \frac{1}{N_h} \sum_{d \in h} a_m(d)$, where N_h is the number of trading days in a maintenance period h .

As mentioned in Section 3.2 (footnote 21), for the purpose of calculating the “basic” BOD balance, we treated GSEs differently than banks. In the case of a GSE, we did not only net out the repayments corresponding to loans issued on day $d - 1$ (i.e., $\text{net}_m(d)$), but *all transfers* sent or received during trading day d —involving *any* counterparty, not only those that meet the sample selection criteria described Section D.1.2. The rationale for netting all transfers that will occur during day d to obtain the GSE’s “basic” BOD balance for day d is that a GSE’s business model generates very predictable cashflows, so through the lens of our theory, we regard the GSE as being able to predict all its intraday Fedwire transfers at the beginning of the trading day.

Step 2: Subtracting reserve requirements. For each bank m in maintenance period h , we computed “adjusted” (excess) reserves as $x_m(h) \equiv a_m(h) - \underline{a}_m^D(h) - \underline{a}_m^L(h)$, where $\underline{a}_m^D(h)$ and $\underline{a}_m^L(h)$ denote the Regulation D and LCR reserve requirements, respectively. The bank-level Regulation-D requirement, $\underline{a}_m^D(h)$, was provided by MPOA. The reserve requirement implied by the LCR regulation is less straightforward, as we discuss next.

As explained in Appendix B (Section B.2.1), the LCR regulation requires a bank to maintain

(on a daily or monthly basis) a quantity of *High Quality Liquid Assets* (HQLA) at least as large as a measure of total net cash outflows in a 30-day standardized stress scenario. Specifically, if we let $H_m(d)$ denote the quantity of qualifying HQLA held by bank m in a trading period d (a day or a month, depending on the type of institution, see footnote 105) and $L_m(d)$ denote the corresponding measure of outflows in the stress scenario, the LCR regulation requires $L_m(d) \leq H_m(d)$. The set of qualifying HQLA includes reserves in excess of Regulation D, as well as securities issued or guaranteed by the U.S. Treasury (and also other securities, but subject to caps and haircuts). The fact that the LCR regulation allows banks to meet the requirement with assets other than reserves presents a challenge when trying to identify the quantity of reserves that bank m treats as “required” to satisfy the LCR constraint in period d , i.e., $\underline{a}_m^L(d)$. Our strategy to tackle this identification problem is to set $\underline{a}_m^L(d) = \max(0, L_m(d) - A_m(d))$, where $A_m(d) \equiv H_m(d) - \max(0, a_m(d) - \underline{a}_m^D(d))$ is the quantity of qualifying HQLA in excess of (i.e., other than) reserves net of the Regulation D requirement.¹¹⁵ Our proposed measure of excess reserves, $x_m(h) \equiv a_m(h) - \underline{a}_m^D(h) - \underline{a}_m^L(h)$, selects the largest level of excess reserves net of the Regulation D requirement that is consistent with the LCR constraint.

For banks that are not subject to LCR regulation (such as banks with assets below \$50 bn in our sample period), we set $\underline{a}_m^L(h) = 0$. Since GSEs are not subject to Regulation D or LCR regulation, we set $\underline{a}_m^D(h) = \underline{a}_m^L(h) = 0$ for $m \in \mathbb{B}_G$. Finally, since we only have quarterly LCR observations (see Section D.1.3), we imputed the same LCR-induced reserve requirement for all maintenance periods within the quarter.

Step 3: Netting predictable payments. To go from the bank-level “adjusted” measure of BOD balance calculated in Step 2, to the bank-level measure of “unencumbered” BOD reserve balance for period h , we netted (at the individual bank level) all *predictable payments* that take place during period h , as explained in Section 3.3.

D.2.2 Network statistics

In this section we describe the calculations of the network statistics reported in Figure 3.¹¹⁶ We begin by introducing some notation. Let v_{md}^e denote the dollar value of all loans extended, and v_{md}^r denote the dollar value of loans received, by bank m on day d . Let v_{mh}^e and v_{mh}^r

¹¹⁵See Section B.2.1 in Appendix B for a more detailed explanation of our strategy to identify the quantity of required reserves induced by the LCR regulation.

¹¹⁶The theoretical counterparts of these computations are discussed in Appendix C (Section C.2).

denote the dollar values of loans extended and received, respectively, during the maintenance period h , i.e., $v_{mh}^e = \sum_{d \in h} v_{md}^e$ and $v_{mh}^r = \sum_{d \in h} v_{md}^r$. Finally, let $v_h = \sum_m v_{mh}^e$ denote the total dollar value of loans extended in maintenance period h . We compute the participation and reallocation values by *bank type*, $i \in \{F, M, S, G\}$, as follows.

Participation rate by bank type. We computed the participation rate for bank type $i \in \{F, M, S, G\}$ during maintenance period h as $\mathcal{P}_{ih} = \sum_{m \in i} \frac{v_{mh}^e + v_{mh}^r}{2v_h}$. We then computed the participation rate of bank type $i \in \{F, M, S, G\}$ in a given year as $\mathcal{P}_i = \frac{1}{N_h} \sum_h \mathcal{P}_{ih}$, where N_h denotes the number of maintenance periods in the year.

Reallocation index by bank type. We computed the reallocation index for bank type $i \in \{F, M, S, G\}$ during maintenance period h as $\mathcal{R}_{ih} = \frac{\sum_{m \in i} v_{mh}^e - \sum_{m \in i} v_{mh}^r}{\sum_{m \in i} v_{mh}^e + \sum_{m \in i} v_{mh}^r}$. We then computed the reallocation index of group i in a given year as $\mathcal{R}_i = \frac{1}{N_h} \sum_h \mathcal{R}_{ih}$, where N_h denotes the number of maintenance periods in the year.

As explained in Section 3.1, the arrows from one node to another in Figure 3 represent loans extended from banks of that type to the other. The arrow width is proportional to the volume of trade between the bank types connected by the arrow. The node size is proportional to the volume of trade between banks of a given type. The arrow widths and node sizes are defined relative to the trades within a year, so they are not comparable across years. The colors of the arrows and nodes are: light blue, dark blue, light red, or dark red, depending on whether the volume-weighted average interest rate on the loans between the two types of banks, expressed as a spread over the EFFR, falls in the first, second, third, or fourth quartile, respectively.

D.2.3 Kernel density estimations

We use Gaussian kernel densities to estimate the distributions of payment shocks, beginning-of-day reserve balances, and aggregate reserve-draining shocks. For the distributions of payment shocks, and the distribution of reserve-draining shocks, we set the smoothing parameter, h , using a standard “rule of thumb”, namely $h = 0.9 \min\left(\hat{\sigma}, \frac{\text{IQR}}{1.34}\right) n^{-1/5}$, where n , $\hat{\sigma}$, and IQR denote the number of observations, standard deviation, and interquartile range of the sample, respectively. For the distributions of beginning-of-day reserves we use the [iterative] methodology described in Botev et al. (2010) to set the smoothing parameter (since they may be multimodal, as seen in Figures 6-9).

D.2.4 Reduced-form estimation of the reserve demand (11)

As in Section 6.2, let s_t denote the EFFR-IOR spread on day t , and Q_t denote the aggregate quantity of reserves at the end of day t . We estimated equation (11) using a nonlinear least-squares procedure. For each sample period, we estimated the vector of parameters, $\nu \equiv \{\underline{s}, \tilde{s}, \xi, Q_0\}$, with $\tilde{s} \equiv \bar{s} - \underline{s}$, to solve

$$\Xi = \min_{\nu} \left\{ \sum_t (s_t - D(Q_t))^2 \right\} \text{ s.t. } 0 \leq \tilde{s}, 0 \leq \xi, \quad (40)$$

where $D(\cdot)$ is as defined in equation (11). We found the solution to (40) by following two steps. In the first step, we did a thorough grid search: we set equally spaced grids for each parameter in ν , computed the hypercube combining all these grids, and then evaluated Ξ for each entry in this hypercube.¹¹⁷ Let ν_{grid} be the vector of parameters that delivered the lowest value of Ξ . In the second step, we used a Nelder-Mead optimization starting from ν_{grid} .

D.2.5 A mapping between reserves of all banks and reserves of active banks

Let Q_t^D denote the quantity of *total reserves* on day t in the sample of all banks in the data (e.g., the quantity of reserves shown in Figure 16). Let Q_t^M denote the quantity of *active excess reserves* on day t that we use to calibrate our model to the year 2019 (and in the interpolation procedure described in Section 3.6, which also uses the year 2017 as an endpoint).¹¹⁸ Let \mathbb{T} denote a subset of trading days and let \mathbf{t} be the cardinality of this set. For any sample $\{Q_t^D\}_{t \in \mathbb{T}}$, define $\bar{Q}_{\mathbb{T}}^D \equiv \frac{1}{\mathbf{t}} \sum_{t \in \mathbb{T}} Q_t^D$. Similarly, for any sample $\{Q_t^M\}_{t \in \mathbb{T}}$, define $\bar{Q}_{\mathbb{T}}^M \equiv \frac{1}{\mathbf{t}} \sum_{t \in \mathbb{T}} Q_t^M$.

Our model output, e.g., the aggregate demand for reserves, is computed using a quantity of reserves $Q \in \mathbb{R}$ constructed with the interpolation procedure described in Section 3.6, which uses \bar{Q}_{2017}^M and \bar{Q}_{2019}^M , i.e., the average quantity of reserves in excess of LCR and Regulation D *in our subsample of active banks* for the two base years. For some exercises (e.g., the top-right panel of Figure 18) we want to show—in the same graph—the model output along with actual daily data observations of total reserves and interest rates, but the observations that we have

¹¹⁷We used grid sizes of: 50 points for \underline{s} , 50 points for \tilde{s} , 123 points for Q_0 , and 63 points for ξ . This gave a combination of 19,372,500 values for ν . The bounds for each grid were: -0.50 and 0.01 for \underline{s} , 0.00 and 1.00 for \tilde{s} , $-3 \times Q_{2019}$ and $3 \times Q_{2019}$ for Q_0 , and $1e-6$ and 0.10 for ξ . We always found the the optimal value for ν well within our bounds.

¹¹⁸That is, $\{Q_t^M\}$ is the time series for the aggregate quantity of reserves for the subsample of banks that were active in fed funds trading during the years under study, net of Regulation D and LCR requirements, as explained in Section 3.3.

available at a daily frequency are $\{Q_t^D\}_{t \in \mathbb{T}}$, not $\{Q_t^M\}_{t \in \mathbb{T}}$. So we need a way to “transform” each daily observation, Q_t^D , into an estimate of Q_t^M .

We adopt a transformation \mathcal{G} , such that $Q_t^M = \mathcal{G}(Q_t^D; \mathbb{T})$ for all $t \in \mathbb{T}$, which satisfies two properties for any sample \mathbb{T} : (1) daily variation in reserves in the full sample of banks is the same as daily variation in reserves in the subsample of banks, i.e., $Q_{t+1}^M - Q_t^M = Q_{t+1}^D - Q_t^D$ for all $t \in \mathbb{T}$ (this is consistent with our strategy of calibrating the slope of our model-generated reserve demand to match the liquidity effect associated with variation in the quantity of reserves of the full sample of banks); and (2) $\bar{Q}_{\mathbb{T}}^M = \mathcal{F}(\bar{Q}_{\mathbb{T}}^D)$, where \mathcal{F} is a linear function that satisfies $\mathcal{F}(\bar{Q}_{2017}^D) = \bar{Q}_{2017}^M$ and $\mathcal{F}(\bar{Q}_{2019}^D) = \bar{Q}_{2019}^M$ (the subscript “2017” denotes the sample of all trading days in the year 2017, and the subscript “2019” denotes the sample of all trading days in the year 2019). For any sample \mathbb{T} of trading days, we posit

$$\begin{aligned} Q_t^M &= \mathcal{G}(Q_t^D; \mathbb{T}) \\ &\equiv Q_t^D - \bar{Q}_{\mathbb{T}}^D + \bar{Q}_{\mathbb{T}}^M, \end{aligned} \quad (41)$$

where

$$\bar{Q}_{\mathbb{T}}^M \equiv \omega_{\mathbb{T}}^D \bar{Q}_{2019}^M + (1 - \omega_{\mathbb{T}}^D) \bar{Q}_{2017}^M, \quad (42)$$

with

$$\omega_{\mathbb{T}}^D \equiv \frac{\bar{Q}_{\mathbb{T}}^D - \bar{Q}_{2017}^D}{\bar{Q}_{2019}^D - \bar{Q}_{2017}^D}. \quad (43)$$

For each day $t \in \mathbb{T}$, the transformation (41) constructs Q_t^M from Q_t^D by first subtracting from Q_t^D its sample mean, $\bar{Q}_{\mathbb{T}}^D$, and then recentering the resulting quantity by adding an *imputed* sample mean, $\bar{Q}_{\mathbb{T}}^M$, corresponding to the *subset of active banks*. The imputed sample mean $\bar{Q}_{\mathbb{T}}^M$ is defined by (42) and (43) as a convex combination of \bar{Q}_{2017}^M and \bar{Q}_{2019}^M (the observed sample means for the subset of active banks in the baseline years 2017 and 2019).

Next, we verify that the mappings \mathcal{G} and \mathcal{F} defined by (41)-(43) satisfy the desired properties. First, notice that for any sample \mathbb{T} , (41) implies $Q_{t+1}^M - Q_t^M = Q_{t+1}^D - Q_t^D$ for all $t \in \mathbb{T}$, so property (1) is satisfied. Second, notice that (42)-(43) define a linear transformation, \mathcal{F} , such that $\bar{Q}_{\mathbb{T}}^M = \mathcal{F}(\bar{Q}_{\mathbb{T}}^D)$, with

$$\mathcal{F}(\bar{Q}_{\mathbb{T}}^D) \equiv \frac{\bar{Q}_{2017}^D \bar{Q}_{2019}^M - \bar{Q}_{2019}^D \bar{Q}_{2017}^M}{\bar{Q}_{2017}^D - \bar{Q}_{2019}^D} + \frac{\bar{Q}_{2017}^M - \bar{Q}_{2019}^M}{\bar{Q}_{2017}^D - \bar{Q}_{2019}^D} \bar{Q}_{\mathbb{T}}^D,$$

which satisfies the desired property (2), i.e., $\mathcal{F}(\bar{Q}_{2017}^D) = \bar{Q}_{2017}^M$ and $\mathcal{F}(\bar{Q}_{2019}^D) = \bar{Q}_{2019}^M$.

The linear mapping $\bar{Q}^D = \mathcal{F}^{-1}(\bar{Q}^M)$ from the quantity of *active excess reserves* to the quantity of *total reserves* is a reasonable approximation for relatively narrow ranges of reserve balances (e.g., for quantities of reserves between \bar{Q}_{2019}^M and \bar{Q}_{2017}^M). However, for some of our quantitative exercises (e.g., Figure 19, Figure 21, and the right panels of Figure 18) we want a mapping to transform values of \bar{Q}^M into values of \bar{Q}^D that performs well *globally* (i.e., for quantities of reserves balances that are far from \bar{Q}_{2019}^M and \bar{Q}_{2017}^M). For this reason, whenever a figure includes a secondary horizontal axis for *total reserves* (i.e., $\bar{Q}_{\mathbb{T}}^D$) that “translates” the *active excess reserves* (i.e., $\bar{Q}_{\mathbb{M}}^D$) on the primary axis, we obtain $\bar{Q}_{\mathbb{T}}^D$ from the following *quadratic* mapping:

$$\bar{Q}_{\mathbb{T}}^D = \mathcal{T}(\bar{Q}_{\mathbb{T}}^M) \equiv A\bar{Q}_{\mathbb{T}}^M + B(\bar{Q}_{\mathbb{T}}^M)^2,$$

with

$$A \equiv \frac{(\bar{Q}_{2017}^M)^2 \bar{Q}_{2019}^D - (\bar{Q}_{2019}^M)^2 \bar{Q}_{2017}^D}{(\bar{Q}_{2017}^M - \bar{Q}_{2019}^M) \bar{Q}_{2017}^M \bar{Q}_{2019}^M}$$

$$B \equiv \frac{\bar{Q}_{2019}^M \bar{Q}_{2017}^D - \bar{Q}_{2017}^M \bar{Q}_{2019}^D}{(\bar{Q}_{2017}^M - \bar{Q}_{2019}^M) \bar{Q}_{2017}^M \bar{Q}_{2019}^M}.$$

The mapping \mathcal{T} satisfies $\bar{Q}_{2017}^D = \mathcal{T}(\bar{Q}_{2017}^M)$, $\bar{Q}_{2019}^D = \mathcal{T}(\bar{Q}_{2019}^M)$, and $\mathcal{T}(0) = 0$, and it is consistent with the linear mapping \mathcal{F}^{-1} (as defined by (42) and (43)) in the sense that for all practical purposes, the difference between the quadratic mapping $Q^D = \mathcal{T}(Q^M)$ and the linear mapping $Q^D = \mathcal{F}^{-1}(Q^M)$ is very small for all $Q^M \in [Q_{2019}^M, Q_{2017}^M]$.¹¹⁹

D.3 Events of September 13–20, 2019

In Section 8 we use our quantitative theory to analyze the fed-funds rate spikes of September 2019. In this section we give more background on the associated reserve-draining shocks, and

¹¹⁹Notice that $Q^D = \mathcal{F}^{-1}(Q^M)$ is the secant line to the quadratic mapping $Q^D = \mathcal{T}(Q^M)$ through the points (Q_{2019}^M, Q_{2019}^D) and (Q_{2017}^M, Q_{2017}^D) . To see that the values of the mappings $\mathcal{F}^{-1}(Q^M)$ and $\mathcal{T}(Q^M)$ are indeed very close for $Q^M \in [Q_{2019}^M, Q_{2017}^M]$, we verify that

$$\arg \max_{Q \in [Q_{2019}^M, Q_{2017}^M]} [\mathcal{F}^{-1}(Q) - \mathcal{T}(Q)] = \frac{Q_{2019}^M + Q_{2017}^M}{2} \equiv Q_M^*,$$

and

$$\frac{\mathcal{F}^{-1}(Q_M^*) - \mathcal{T}(Q_M^*)}{\mathcal{T}(Q_M^*)} = \frac{14.21}{1897.05} \approx 0.0075.$$

the monetary-policy interventions that followed these rate spikes.¹²⁰

The events unfolded as follows. On Friday, September 13 the beginning-of-day supply of reserves was about \$1.5 tn, and the EFFR printed at 214 bps. In the top panel of Figure 21, September 13 is the dark dot that sits on the demand for reserves generated by the theory—well within the FFR target range. On Monday, September 16 the beginning-of-day supply of reserves was \$51.5 bn lower than on the previous trading day (due to reserve-draining shocks that occurred throughout Friday, September 13), and the EFFR printed at 225 bps (the upper limit of the target range). In the top panel of Figure 21, September 16 is the rightmost dark dot that sits on the upper limit of the target range for the EFFR. On Tuesday, September 17 the beginning-of-day supply of reserves was \$65.72 bn lower than on the previous trading day (due to reserve-draining shocks that occurred throughout Monday, September 16), and the EFFR printed at 230 bps (5 bps above the upper limit of the target range). In the top panel of Figure 21, September 17 is the uppermost dark dot. Following an overnight repo operation that injected \$53 billion on Tuesday, September 17, the beginning-of-day supply of reserves on September 18 was \$46.3 bn higher than on the previous day, and the EFFR fell to 225 bps.¹²¹

The morning of Tuesday, September 17 was the first time since the GFC that the Desk conducted an open-market operation to manage the fed funds rate. That Tuesday afternoon the Desk announced it would conduct an overnight operation at 8:15 a.m. on Wednesday, September 18. This operation injected \$75 bn, which contributed to the beginning-of-day supply of reserves on Thursday, September 19, being \$3.67 bn higher than the previous day. Similar operations were used to inject \$75 bn every day until the end of the week. The EFFR printed at 190 bps on September 19 and September 20.¹²²

¹²⁰Table 2 summarizes the main facts. Most of the events we describe in this section are based on the detailed accounts provided Afonso et al. (2020a) and Anbil et al. (2020).

¹²¹On Monday afternoon (2019/09/16), in response to the observed upward pressure on the EFFR, the Desk announced an overnight repo operation to be conducted at 9:30 AM on Tuesday (2019/09/17), offering up to \$75 billion against Treasury, agency, and agency MBS collateral, of which only \$53 bn were subscribed.

¹²²On Thursday, September 19, the Federal Reserve also made adjustments to administered rates and the FFR target range. The ONRRP was reduced from 200 bps to 170 bps, the IOR from 210 bps to 180 bps, and the DWR from 275 bps to 250 bps. The lower limit of the FFR target range was reduced from 200 bps to 175 bps, and the upper limit was reduced from 225 bps to 200 bps. On the morning of Friday, September 20, the Desk announced a series of operations over the quarter-end, which included three two-week operations covering the quarter-end and daily overnight operations of \$75 billion through October 10. The September 16–17 event seem to have had lasting an impact on the conduct of monetary policy. As Afonso et al. (2020a, p. 24) recount:

On October 11, 2019, the FOMC announced its intention to maintain an ample supply of reserve balances at or above the level that prevailed in early September. The FOMC instructed the Desk to purchase Treasury bills at least into the second quarter of 2020 (and to continue repo operations) in order to supply reserves and mitigate money market pressures that might impede policy imple-

D.3.1 JPM earnings call for the period ending September 30, 2019

In this section we report the key excerpts of the earnings call of October 15, 2019, in which Jamie Dimon, Chairman and CEO of JPMorgan Chase (JPM), answered questions about why JPM was not more active lending in money markets during the week of September 16, 2019.

Question: Glenn Schorr (Analyst, Evercore ISI)

Curious your take on everything that went on in the repo markets during the quarter, and I would love it if you could put it in the context of maybe the fourth quarter of last year. If I remember correctly, you stepped in in the fourth quarter, saw higher rates, threw money at it, made some more money, and it calmed the markets down. I'm curious what's different this quarter that did not happen, and curious if you think we need changes in the structure of the market to function better on a go-forward basis.

Answer: Jamie Dimon (Chairman and CEO, JPM)

So, if I remember correctly, you got to look at the concept of – we have a checking account at the Fed with a certain amount of cash in it. Last year we had more cash than we needed for regulatory requirements. So when repo rates went up, we went from the checking account, which [ph] was paying (00:14:10) IOER into repo. Obviously makes sense, you make more money. But now the cash in the account, which is still huge – it's \$120 billion in the morning and goes down to \$60 billion during the course of the day and back to \$120 billion at the end of the day – that cash, we believe, is required under resolution and recovery and liquidity stress testing. And therefore, we could not redeploy it into repo market, which we would have been happy to do. And I think it's up to the regulators to decide they want to recalibrate the kind of liquidity they expect us to keep in that account. Again, I look at this as technical; a lot of reasons why those balances dropped to where they were. I think a lot of banks were in the same position, by the way. But I think the real issue, when you think about it, is what does that mean if we ever have bad markets? Because that's kind of hitting the red line in the Fed checking account, you're also going to

mentation. The goal of the bill purchases was to ensure the smooth functioning of money markets at the current monetary policy stance, not to change the monetary policy stance.

For details, see <https://www.federalreserve.gov/newsevents/pressreleases/monetary20191011a.htm>.

hit a red line in LCR, like HQLA, which cannot be redeployed either. So, to me, that will be the issue when the time comes. And it's not about JPMorgan. JPMorgan will be fine in any event. It's about how the regulators want to manage the system and who they want to intermediate when the time comes.

Question: Erika Najarian (Analyst, Bank of America Merrill Lynch)

Yes, good morning. My first question is a follow-up to Glenn's question. As we think about the crosscurrents of resolution planning, LCR, and liquidity stress testing, could you help us – what is the level of excess deployable cash at JPMorgan?

Answer: Jamie Dimon (Chairman and CEO, JPM)

As I said, we have \$120 billion in our checking account at the Fed, and it goes down to \$60 billion and then back to \$120 billion during the average day. But we believe the requirement under CLAR and resolution and recovery is that we need enough in that account, so if there's extreme stress during the course of the day, it doesn't go below zero. If you go back to before the crisis, you'd go below zero all the time during the day. So the question is, how hard is that as a red line? Was the intent of regulators between CLAR and resolution to lock up that much of reserves in the account with Fed? And that'll be up to regulators to decide. But right now, we have to meet those rules and we don't want to violate anything we've told them we're going to do.

For a full transcript of the call, visit: <https://tinyurl.com/29scwszt>.

E Aggregate demand for reserves: alternative estimations

Our baseline reduced-form demand estimation in Section 6.2 consisted of estimating the parameters $(\underline{g}, \bar{s}, \xi, Q_0)$ in (11) by nonlinear least squares. The estimated demand fits the data well, but performs poorly for out-of-sample levels of reserves: Notably, the estimated demand predicts the EFRR-IOR spread would remain unchanged if total reserves were drained from \$1 tn to zero. In this section we consider several alternative econometric specifications of the reduced-form estimation of the aggregate demand for reserves. Section E.1 tries to improve on the empirical model of Section 6.2 by imposing theoretically-motivated constraints on the

estimation. Section E.2 considers a semi-log specification that is common among practitioners. Section E.3 considers a variant of the semi-log specification.

All the alternative reduced-form estimation strategies we consider support the two main lessons of Section 6.2.1. First, our theory identifies a set of structural “shifters” of the aggregate demand relationship that can help with the identification problems that pervade all reduced-form econometric estimations of the aggregate demand for reserves. Second, our quantitative-theoretic approach delivers estimates of the demand for reserves that fit available data as well as the reduced-form approaches, but these approaches have very different out-of-sample predictions. Specifically:

- (a) For large levels of reserves, the slope of the quantitative-theoretic demand becomes virtually flat (e.g., \$1.3 tn of total reserves), while the slopes of the reduced-form econometric estimates tend to remain positive even for very large reserves (e.g., even for total reserves in excess of \$2.5 tn).
- (b) For relatively low levels of reserves, the model-generated demand becomes quite steep for total reserves between \$600 bn and \$340 bn, and flattens for levels lower than \$340 bn. In contrast, in the specifications of Sections E.2 and E.3, the slopes of the reduced-form demand estimates increase exponentially as total reserves decrease, and become unreasonably large at low (e.g., pre-GFC) levels of reserves.

E.1 NLS estimation of (11) with constraints motivated by theory

In this section we try to improve the reduced-form empirical model we used in Section 6.2 by imposing two constraints on the estimation that are grounded on elementary theory. Specifically, we redo the NLS estimation of (11) but imposing that \bar{s} should equal the largest value of $\bar{\iota}_w - \iota_r$ in the relevant sample, and that \underline{s} should equal the lowest value of $\bar{\iota}_o - \iota_r$ in the relevant sample (with $\bar{\iota}_w$, $\bar{\iota}_o$, and ι_r as defined in Section 4). The results are reported in Figure 22, which is analogous of Figure 18. The global fit of this reduced-form approach looks somewhat more credible than the unconstrained version in Figure 18; at least the EFFR-IOR spread now rises as the quantity of reserves falls below \$1 tn. However, even after we control by IOR-ONRRP regime, as we do in the bottom-left and bottom-right panels of the figure, the estimated demands are still quite different from our quantitative-theoretic estimates. To illustrate, compare the quantitative-theoretic estimate in the top-right panel with the reduced-form estimate in the

bottom-right panel for the sample with IOR-ONRRP spread equal to 10 bps. In the former, the slope of the demand becomes flat somewhere above \$1.3 tn of total reserves, while the slope of the latter remains positive even if total reserves exceed \$2.5 tn. The behavior is also quite different for relatively low levels of reserves: the model-generated demand becomes quite steep at about \$600 bn of total reserves, while the slope of the reduced-form estimate does not vary much with the quantity of reserves (even as the quantity of total reserves approaches zero).

E.2 A semi-log specification

In this section we consider the following semi-log specification for the demand for reserves:

$$s_t = a + b \ln(Q_t), \quad (44)$$

where s_t denotes the EFFR-IOR spread on day t and Q_t denotes the aggregate quantity of reserves at the end of day t . We estimate the parameters a and b by ordinary least squares (OLS); the results are reported in Figure 23.

The main points we made in Section 6.2 still hold. First, comparing the top-left and bottom-left panels of Figure 23 we see that incorporating the minimal theoretical insight that changes in the IOR-ONRRP spread act like demand shifters, makes a big difference for the global estimates of the demand for reserves. Second, even after we control by IOR-ONRRP regime, as we do in the bottom-left and bottom-right panels of the figure, the estimated demands are still quite different from our quantitative-theoretic estimates. To illustrate, compare the quantitative-theoretic estimate in the top-right panel with the reduced-form estimate in the bottom-right panel for the sample with IOR-ONRRP spread equal to 10 bps. According to the former, the slope of the demand becomes flat somewhere above \$1.3 tn of total reserves, while the slope of the latter remains positive even if total reserves exceed \$2.5 tn. The behavior is also quite different for relatively low levels of reserves: the model-generated demand becomes quite steep at about \$600 bn of total reserves, but then flattens at about \$340 bn. In contrast, the slope of the reduced-form estimate increases exponentially as total reserves decrease, and eventually becomes unreasonably large.

E.3 The López-Salido and Vissing-Jorgensen (2023) specification

In this section we consider the reduced-form specification for the demand for reserves proposed in López-Salido and Vissing-Jorgensen (2023), who assume

$$s_t = a + b \ln(Q_t) + c \ln(D_t), \quad (45)$$

where s_t denotes the EFRR-IOR spread in period t , Q_t denotes the aggregate quantity of reserves in period t , and D_t is a measure of bank deposits in period t .¹²³ We estimate the parameters a , b , and c by OLS; the results are reported in Figure 24. The demand estimates are very similar to the ones we obtained in Section E.2 and reported in Figure 23, so the main points we made about the specification (44) also hold for (45).

F Quantitative analysis for the pre-GFC-regulation regime

Our quantitative analysis in the body of the paper focuses on the current post-GFC monetary-policy framework. For completeness, and because the pre-GFC period is of historical interest, in this section we also study the pre-GFC framework. The pre-GFC and post-GFC frameworks differ in two ways. First, the quantity of excess reserves was close to zero in the former, but is very large in the latter. Second, as discussed in Section 3, regulations introduced after the GFC have affected banks' payoffs from fed funds trading. For this reason, in this section we recalibrate the model for a base year before the GFC, which we choose to be 2006.¹²⁴

F.1 Calibration

We set ι_w to match the prevailing DWR, and $\iota_o = 0$ (since there was no ONRRP facility in 2006). The remaining nine parameters, ι_r and $\{\beta_i, \kappa_i\}_{i \in \mathbb{N}}$, are calibrated so that the equilibrium

¹²³The baseline measure of D_t in López-Salido and Vissing-Jorgensen (2023) is “DPSACBW027SBOG” (*Deposits, All Commercial Banks*) from <https://fred.stlouisfed.org>.

¹²⁴Our main motive for recalibrating the model is that the trading network, which in our theory is represented by the parameters $\{\beta_i, \kappa_i\}_{i \in \mathbb{N}}$, may not be stable across policy regimes. For example, it is reasonable to imagine that the trading patterns represented by the type-specific meeting rates may change in response to regulatory constraints, in particular those post-GFC regulations that increased the cost of borrowing, and therefore the cost of intermediating fed funds. We use 2006 as the baseline year for the pre-GFC period for two reasons. First, policy rates and total reserves remained stable for most of that year, and it was the last “normal” year before the GFC that spurred the policy interventions that changed the landscape of the fed funds market. Second, the 2006 calibration will allow us to assess the model fit in a pre-GFC-regulation environment, and it will also allow us to test the quantitative predictions of the theory as we vary the level of aggregate excess reserves from near zero (the level they had during 2006) to \$2.689 tn (the level they reached for all the banks in our sample in 2014, which was the last pre-GFC-regulation year).

of the model matches the following nine empirical moments: (i) effective fed funds rate¹²⁵; (ii)-(v) reallocation indices $\{\mathcal{R}_i\}_{i \in \mathbb{N}}$ (as defined in Section 3.1); (vi)-(viii) participation rates $\{\mathcal{P}_i\}_{i \in \mathbb{N} \setminus \{F\}}$ (as defined in Section 3.1)¹²⁶; (ix) empirical estimates of the “liquidity effect” (at the average level of aggregate reserves outstanding in the base year, as reported in Section 3.5).

Table 3 reports the parameter values, empirical targeted moments, and the corresponding theoretical moments for the 2006 calibration. Banks of type F , M , and S , accounted for about 0.5%, 3%, and 95%, of all the institutions that were active in the fed funds market in 2006, respectively.¹²⁷ To interpret the frequencies of payment shocks, $\{\lambda_i\}_{i \in \mathbb{N}}$, recall that λ_i represents the probability that a bank of type i receives a payment shock in a one-second time interval, so for example, $\lambda_F = 0.901$ implies a bank of type F receives a payment shock approximately every 1.1 seconds, on average. Similarly, $\lambda_M = 0.402$ implies a bank of type M receives a payment shock approximately every 2.5 seconds, and $\lambda_S = 0.007$ implies a bank of type S receives a payment shock approximately every 2.38 minutes, on average. The rate ι_w corresponds to a DWR equal to 6.25% per annum, which was in effect in the second half of 2006. The calibrated value of ι_r is 4.81% per annum.¹²⁸

The frequency of trade, β_i , is the probability a bank of type i contacts a trading partner during a 42-second time interval. Thus, the calibrated values $\{\beta_i\}_{i \in \mathbb{N}}$ for 2006 imply that banks of type F , M , S , and G trade fed funds approximately every 1.75 minutes, 8 minutes, 20 minutes, and 3.5 minutes, respectively. The calibration also ensures that, when computed in the neighborhood of zero excess reserves, the magnitude of the “liquidity effect” in the theory is within the range of the empirical estimates reported in Section 3.5 (i.e., about a 1.7 bp

¹²⁵Our calibration strategy uses the effective fed funds rate as a calibration target unless the Federal Reserve pays interest on reserves (IOR) in the base year, in which case we simply set ι_r to match the IOR. For example, the IOR was 2.35% per annum in May–July 2019, so we set $\iota_r = 0.0235/360$ in the 2019 calibration. The Federal Reserve did not pay IOR before October 9, 2008, so in the 2006 calibration we regard ι_r as a proxy for a bank’s unmodelled opportunity cost of lending reserves in the fed funds market, and calibrate it internally so that the average (volume-weighted) interest rate in the model is equal to 5.25% per annum, which was the effective fed-funds rate prevailing during the second half of 2006.

¹²⁶The participation rate of type F banks is not an explicit calibration target because it is implied by the participation rates of the other three types, since $\sum_{i \in \mathbb{N}} \mathcal{P}_i = 1$.

¹²⁷The main change in the bank population between 2006 and 2019 is the reduction in the total number of active banks in our sample, mostly due to the fact that almost half of the banks of type S that were fed funds market participants in 2006 did not trade fed funds during 2019.

¹²⁸For comparison, 4.81% per annum is the 0.5 percentile of the volume-weighted distribution of rates observed in the second half of 2006. That is, only half of one percent of the fed funds traded in the second semester of 2006 had a rate below 4.81%, so we regard 4.81% as a reasonable proxy for the unmodelled opportunity cost of an alternative use of reserves. We focus on July–December because in that period the administered rates (i.e., the Discount-Window rate and the fed funds rate target) were constant and equal to the rates targeted in the 2006 calibration (the administered rates had been gradually increasing in the first half of 2006).

increase in the fed funds rate per \$1bn reduction in the aggregate quantity of reserves).¹²⁹ The borrowing costs $\{\kappa_i\}_{i \in \mathbb{N}}$, which proxy for institutional and regulatory considerations that affect banks' incentives to buy fed funds, are null for banks of type F , M , and S in the 2006 calibration.¹³⁰

F.2 Validation

In this section we report the model fit of empirical observations that were not targeted in the calibration. We organize the material in two sections: the first focuses on the cross-sectional distribution of loan rates, and the second on the main features of the fed funds trading network.

F.2.1 Distribution of interest rates

Figure 26 shows the empirical and theoretical cumulative distribution functions of bilateral fed funds rates in the year 2006 (expressed in percent per annum).¹³¹ The model delivers a reasonable fit for the distribution of bilateral fed funds rate, which was not a calibration target.

F.2.2 Fed funds trading network

Figure 27 shows the empirical fed funds trading network for the year 2006 (top panel) and the corresponding trading network generated by the model for the 2006 calibration (bottom panel). As explained in Section 3.1, these network plots show the location of the four bank types in the coordinate axes defined by the reallocation index, \mathcal{R}_i , and the participation rate, \mathcal{P}_i , and convey information on the sizes of the flows of reserves associated with fed funds lending across and within bank types, as well as on the average interest rates on the underlying loans.¹³²

¹²⁹Figure 25 shows the magnitude of the liquidity effect in the model calibrated to 2006 (extracting reserves using the procedure described in Section 3.6), as well as the confidence bands for the estimates from Carpenter and Demiralp (2006) reported in Section 3.5. The model-generated liquidity effect is within the range of empirical estimates.

¹³⁰In every calibration the value of κ_G is set large enough to match the observation that GSEs essentially do not borrow in the fed funds market.

¹³¹The empirical interest rates for 2006 are from the sample period July–December because throughout that period the Discount-Window rate and the fed funds rate target were constant and equal to the rates targeted in the 2006 calibration. To obtain the equilibrium rates for 2006, the model is calibrated as in Table 3.

¹³²In comparing the top and bottom panels of Figure 27, notice that while the positions of the four nodes (each of which represents the set of banks assigned to a particular type) in \mathcal{R}_i - \mathcal{P}_i space have been used as calibration targets, the remaining collection of statistics that shape these network representations were not targeted. This includes the node sizes (each of which is proportional to the volume of trade between banks of a given type), the direction of each arrow (which indicates which bank type lends), the width of each arrow (which is proportional to the volume of trade between the bank types connected by the arrow), the colors of the arrows and nodes (which are light blue, dark blue, light red, or dark red, if the volume-weighted average interest rate on the loans

The theoretical network matches several characteristics of the empirical one. For example, it replicates quite well the direction and volume of the loans between bank types (represented by the direction and width of the arrows between the nodes). In this regard, one shortcoming of the model is that it underpredicts the volume of loans within bank types S and M . The model is consistent with the empirical facts that banks of type S lend to each other at relatively high rates, while banks of type F can borrow at relatively low rates from banks of type M , S , and GSEs. In terms of shortcomings, the model predicts that banks of type S borrow at relatively high rates from GSEs, that loans between banks of type F carry relatively low rates, and that loans between banks of type M carry relatively high rates, as do loans from type F to type M , but these predictions do not match the empirical patterns.

F.3 Aggregate demand for reserves

Consider the model calibrated to the year 2006, as described in Table 3, but with $\iota_w = 0.0075/365$ and $\iota_r = 0.0025/365$, to match the DWR and IOR in the year 2014. Then, using the notation introduced in Section 3.6, let $Y_0 = 2006$ and $Y_1 = 2014$, i.e., Y_0 and Y_1 represent the years 2006, and 2014, respectively, with \bar{n}_{2014}^i and \bar{F}_{2014}^i given by the estimates reported in Section 3.3. Construct a grid, $\mathbb{G} \subset \mathbb{R}$ for ω , and for each $\omega \in \mathbb{G}$, use the interpolation procedure described by (7) and (8) to generate the sample $\{(\bar{n}_{Y_\omega}^i, \bar{F}_{Y_\omega}^i)\}_{(i,\omega) \in \mathbb{N} \times \mathbb{G}}$. For each pair $(\bar{n}_{Y_\omega}^i, \bar{F}_{Y_\omega}^i)$ of elements of this sample, use the model to compute the corresponding equilibrium value-weighted fed funds rate, which we denote $\iota_{Y_\omega}^*$, and let $Q_{Y_\omega} \equiv \sum_{i \in \mathbb{N}} \bar{n}_{Y_\omega}^i \int ad \bar{F}_{Y_\omega}^i(a)$. This procedure delivers a sample of pairs, $\{(Q_{Y_\omega}, \iota_{Y_\omega}^*)\}_{\omega \in \mathbb{G}}$, which we represent with the mapping $\iota_{Y_\omega}^* = \mathcal{D}(Q_{Y_\omega}; \Pi)$. This mapping, which we interpret as the aggregate demand for reserves generated by the theory, is shown in Figure 28.¹³³

between the two types of banks, expressed as a spread over the EFFR, falls in the first, second, third, or fourth quartile, respectively).

¹³³We use 2006 as one endpoint for our interpolation procedure since it was the last year of the scarce-reserve regime that prevailed until the GFC. We use 2014 as the other endpoint because it is the year when the quantity of reserves achieved its maximum historical level of the pre-2020 era. By varying ω on $[0, 1]$ we can use (9) to span any aggregate level of excess reserves between 0 (roughly the pre-GFC level prevailing in 2006) and \$2.7 tn (roughly the level achieved in 2014).

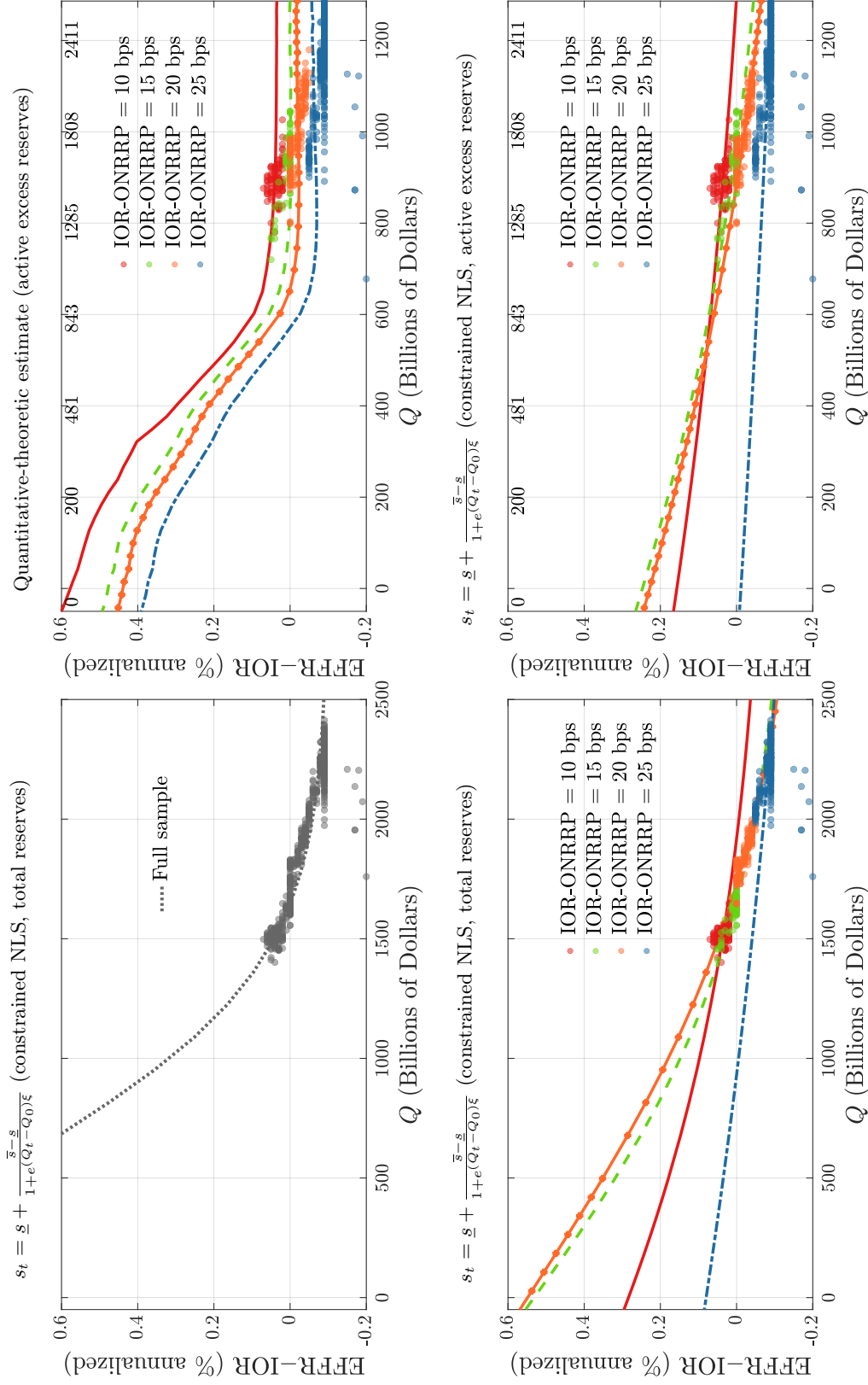


Figure 22: Reserve demand estimation: model vs. NLS fit of (11) with theoretically motivated constraints on \bar{s} and \underline{s} .

Notes: In each panel: vertical axis is EFR-IOR spread; horizontal axis is *total reserves* or *active excess reserves* as defined in Section 6. Full sample: all trading days in 2017/01/20–2019/09/13. Top-left panel: total reserves and constrained NLS fit of (11) on full sample. Bottom-left panel: total reserves and constrained NLS fits of (11) on each subsample defined by IOR-ONRRP spread. Top-right panel: active excess reserves and aggregate demands implied by the theory for each IOR-ONRRP spread (all other parameters as in the baseline calibration). Bottom-right panel: active excess reserves and constrained NLS fits of (11) on each subsample. Secondary horizontal axis (above top-left and bottom-left panels) translates active excess reserves into total reserves.

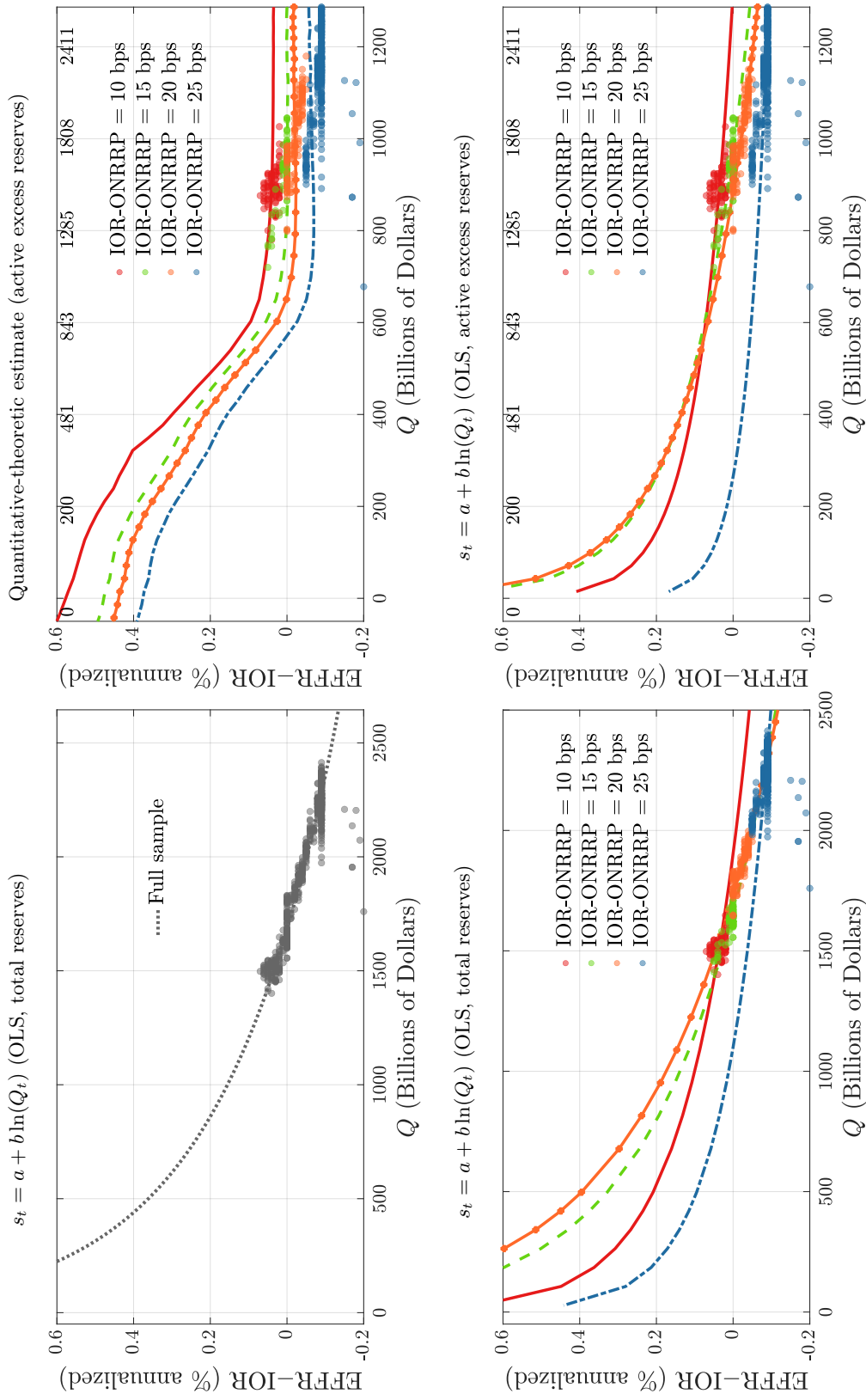


Figure 23: Reserve demand estimation: model vs. OLS fit of (44).

Notes: In each panel: vertical axis is EFRR-IOR spread; horizontal axis is *total reserves* or *active excess reserves* as defined in Section 6. Full sample: all trading days in 2017/01/20–2019/09/13. Top-left panel: total reserves OLS fit of (44) on full sample. Bottom-left panel: total reserves and OLS fits of (44) on each subsample defined by IOR-ONRRP spread. Top-right panel: active excess reserves and aggregate demands implied by the theory for each IOR-ONRRP spread (all other parameters as in the baseline calibration). Bottom-right panel: active excess reserves and OLS fits of (44) on each subsample. Secondary horizontal axis (above top-left and bottom-left panels) translates active excess reserves into total reserves.

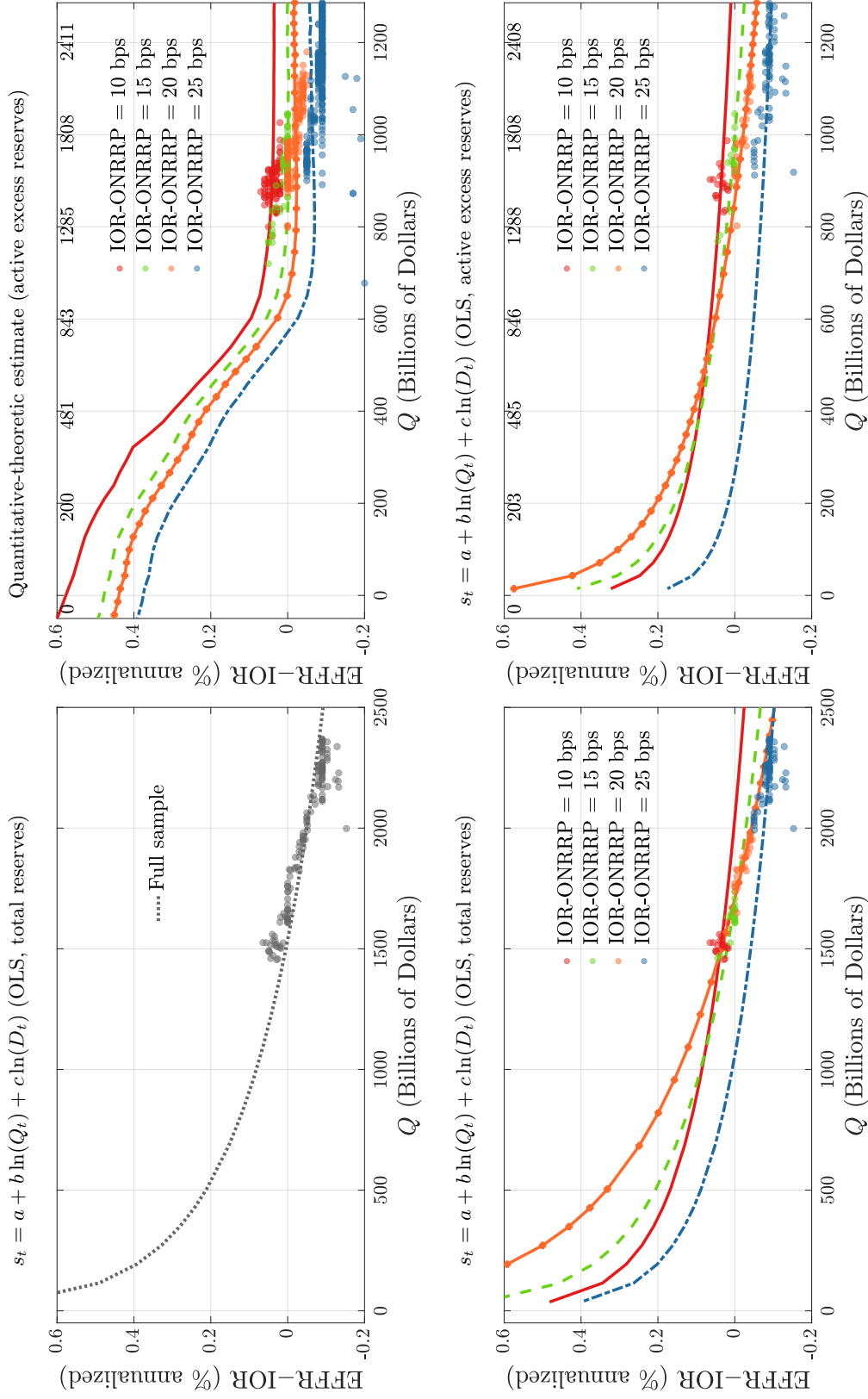


Figure 24: Reserve demand estimation: model vs. OLS fit of (45).

Notes: In each panel: vertical axis is EFFR-IOR spread; horizontal axis is *total reserves* or *active excess reserves* as defined in Section 6. Full sample: all trading days in 2017/01/20–2019/09/13. Top-left panel: total reserves OLS fit of (45) on full sample. Bottom-left panel: total reserves and OLS fits of (45) on each subsample defined by IOR-ONRRP spread. Top-right panel: active excess reserves and aggregate demands implied by the theory for each IOR-ONRRP spread (all other parameters as in the baseline calibration). Bottom-right panel: active excess reserves and OLS fits of (45) on each subsample. Secondary horizontal axis (above top-left and bottom-left panels) translates active excess reserves into total reserves.

Parameter	Target	Moment	
		Data	Model
$n_F = 0.005$	proportion of financial institutions of type F	4/754	0.005
$n_M = 0.029$	proportion of financial institutions of type M	22/754	0.029
$n_S = 0.950$	proportion of financial institutions of type S	716/754	0.950
$n_G = 0.016$	proportion of financial institutions of type G	12/754	0.016
$\lambda_F = 0.901$	bank-level share of unexpected payments per second for type F	0.901	0.901
$\lambda_M = 0.402$	bank-level share of unexpected payments per second for type M	0.402	0.402
$\lambda_S = 0.007$	bank-level share of unexpected payments per second for type S	0.007	0.007
$\lambda_G = 0$	bank-level share of unexpected payments per second for type G	0	0
$\iota_w = 0.0625/360$	Discount-Window rate (primary credit, 6.25% per annum)	0.0625/360	0.0625/360
$\iota_r = 0.0481/360$	effective fed funds rate (5.25% per annum)	0.0525/360	0.0525/360
$\beta_F = 0.401$	estimated liquidity effect around zero excess reserves (bps per \$1 bn)	$\in [1, 3]$	1.7
$\beta_M = 0.089$	participation rate of financial institutions of type M (i.e., \mathcal{P}_M)	0.27	0.27
$\beta_S = 0.033$	participation rate of financial institutions of type S (i.e., \mathcal{P}_S)	0.18	0.19
$\beta_G = 0.200$	participation rate of financial institutions of type G (i.e., \mathcal{P}_G)	0.07	0.07
$\kappa_F = 0$	reallocation index of financial institutions of type F (i.e., \mathcal{R}_F)	0.068	0.064
$\kappa_M = 0$	reallocation index of financial institutions of type M (i.e., \mathcal{R}_M)	-0.385	-0.268
$\kappa_S = 0$	reallocation index of financial institutions of type S (i.e., \mathcal{R}_S)	-0.268	-0.126
$\kappa_G = 1.25e-3$	reallocation index of financial institutions of type G (i.e., \mathcal{P}_G)	0.995	1

Table 3: Calibration for the year 2006.

Notes: Each non-shaded parameter is calibrated externally (i.e., to match a corresponding target moment, independently of the model and other parameters). Shaded parameters are calibrated internally (i.e., jointly, to match the set of shaded target moments, using the equilibrium conditions of the model, and given the values of the parameters calibrated externally). The calibration assumes a model period corresponding to approximately to 42 seconds in a trading day, $r = 0$, $\mathbb{N} = \{F, M, S, G\}$ (as discussed in Section 3.1), $\theta_i = 1/2$ for all $i \in \mathbb{N}$, $\{G_{ij}\}_{i,j \in \mathbb{N}}$ are estimated as described in Section 3.2, $\{F_0^i\}_{i \in \mathbb{N}}$ are estimated as described in Section 3.3, $u_i = 0$ for all $i \in \mathbb{N}$, $\{U_i^j\}_{i \in \mathbb{N}}$ are as in Section 4).

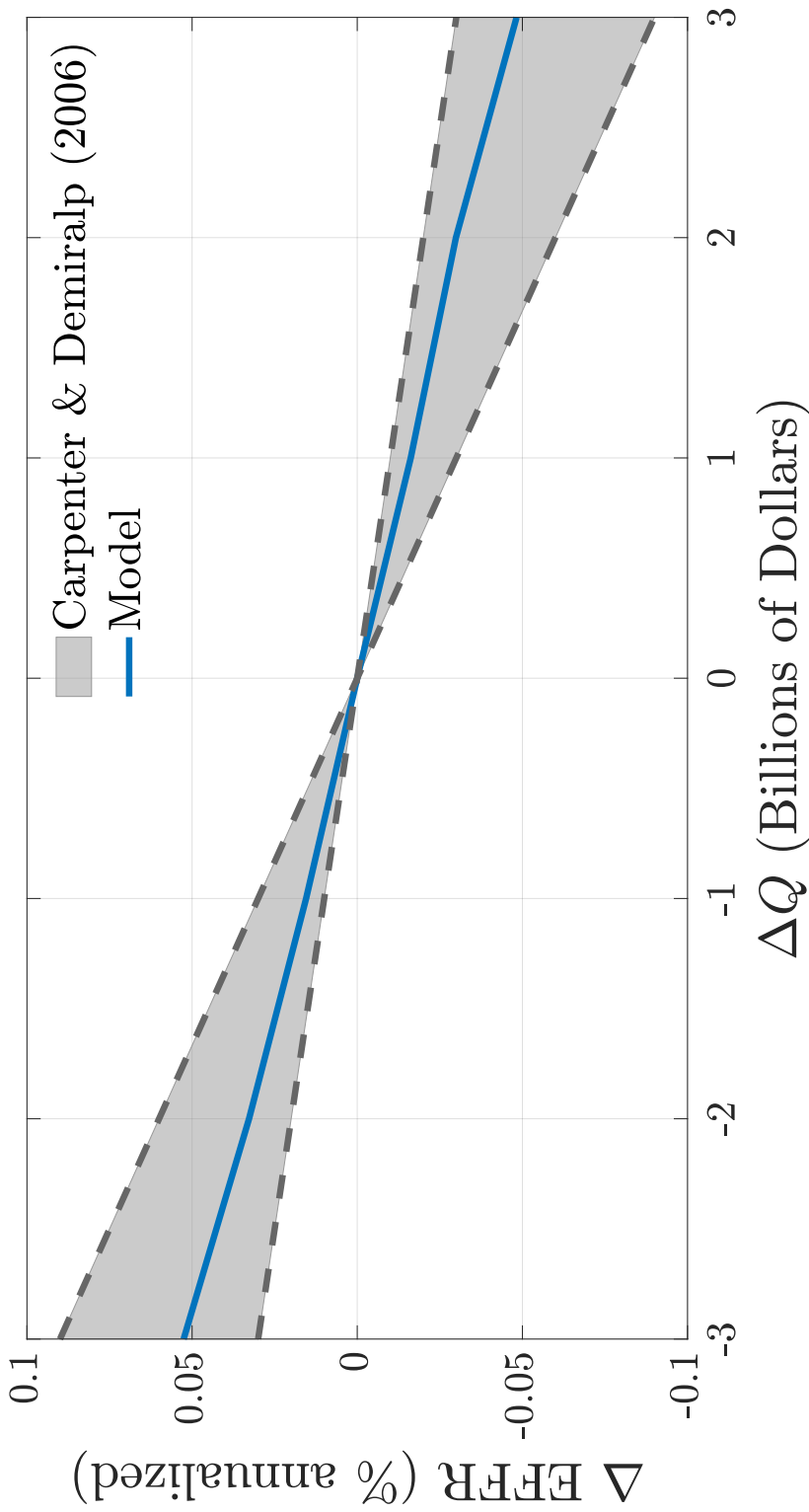


Figure 25: Liquidity effect: model and empirical estimates for the year 2006.

Notes: Rates in the vertical axis are in percent per annum. The shaded area represents the 95% confidence interval for the point estimate of the liquidity effect from Carpenter and Demiralp (2006). The solid line is the change in the equilibrium fed funds rate implied by the theory in response to changes in the total quantity of reserves (starting from the quantity of reserves corresponding to the 2006 calibration, and extracting reserves using the procedure described in Section 3.6.)

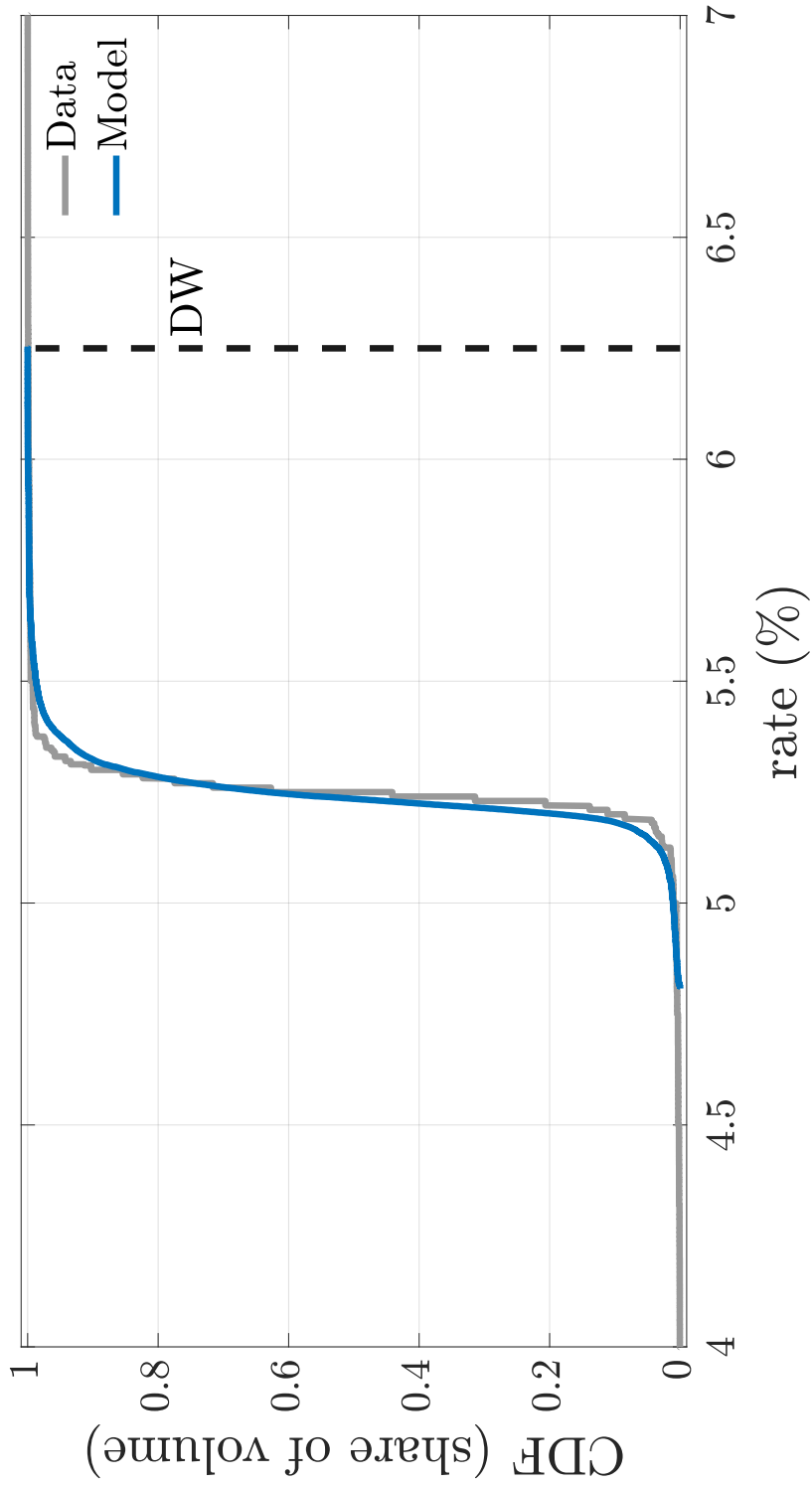


Figure 26: Empirical and theoretical cumulative distribution functions of bilateral fed funds rates for the year 2006.

Notes: Rates in the horizontal axis are in percent per annum.

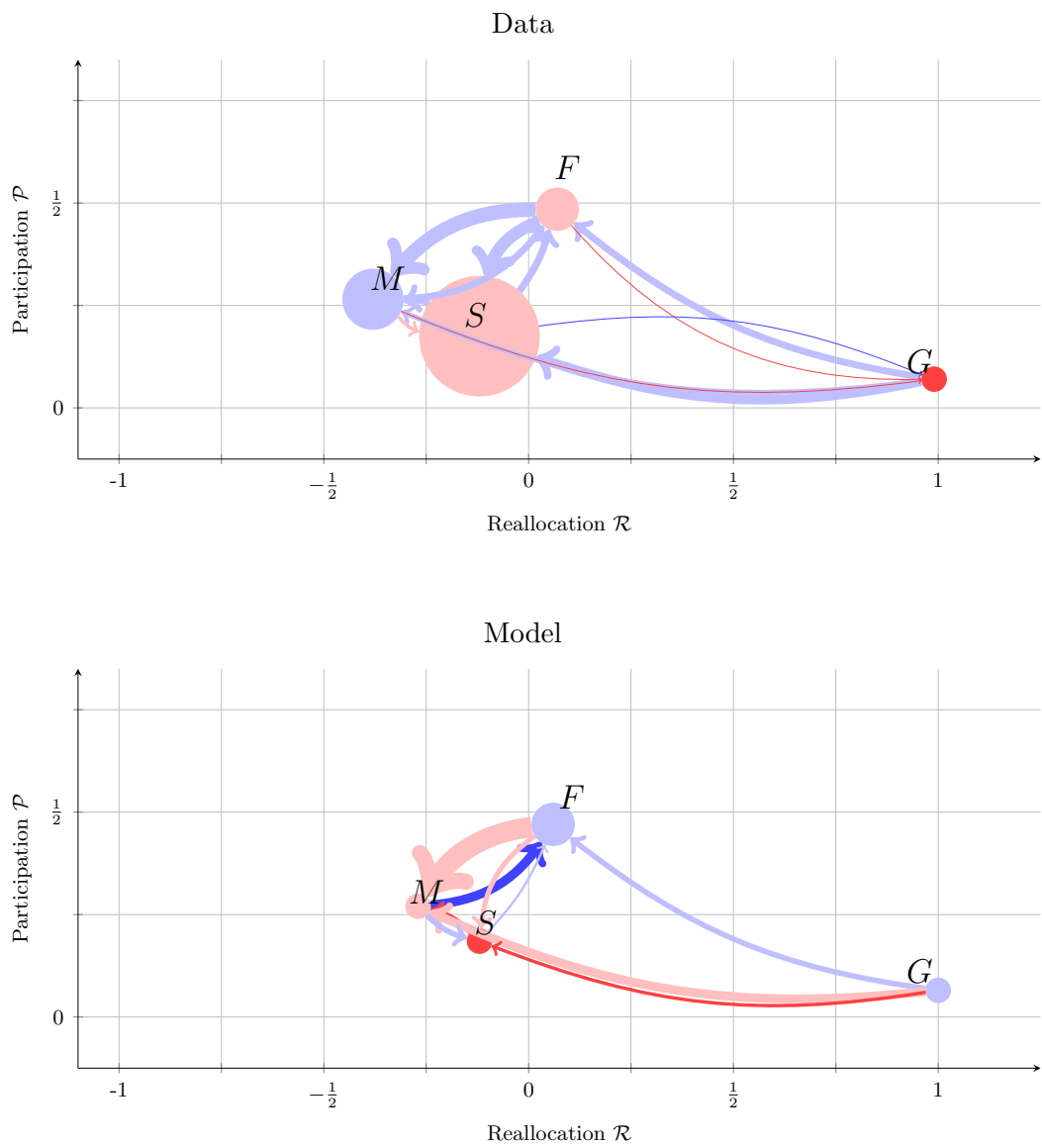


Figure 27: Empirical and theoretical fed funds trading networks for 2006.

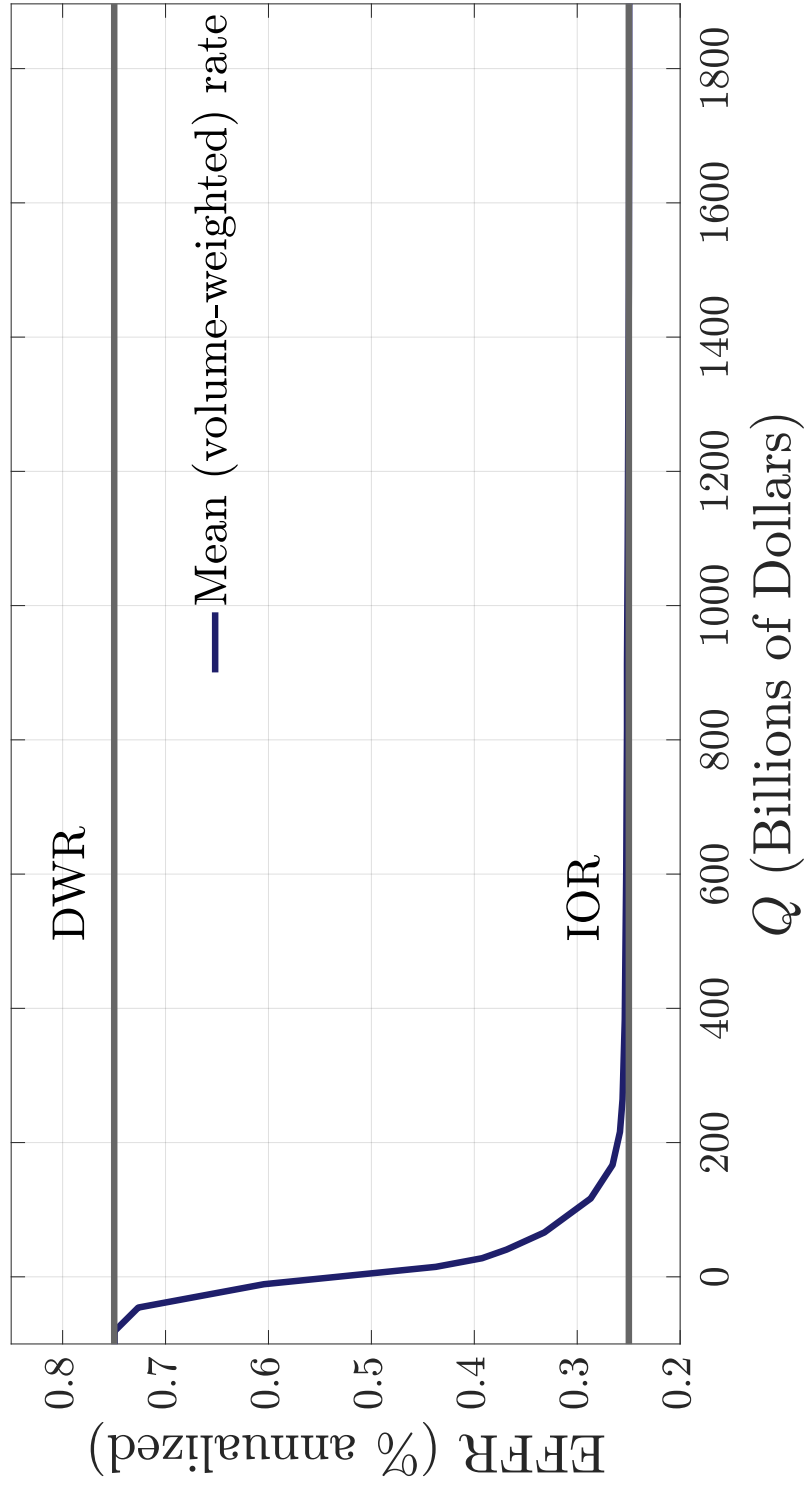


Figure 28: Theoretical aggregate demand for reserves for the year 2006 calibration.

Notes: Theoretical aggregate demand $\ell_{Y,\omega}^* = \mathcal{D}(Q_{Y,\omega}; \Pi)$ for the model calibrated as in Table 3, and with $\ell_{Y,\omega}^*$ and $Q_{Y,\omega}$ computed with the interpolation procedure described in Section 3.6, for $Y_0 = 2006$ and $Y_1 = 2014$.