

Reinforcement Learning Policy Recommendation for Interbank Network Stability

Author names are hided for blind review purposes

Abstract

In this paper, we analyze the effect of a policy recommendation on the performance of an artificial interbank market. Financial institutions stipulate lending agreements following a public recommendation and their individual information. The former is modeled by a reinforcement learning optimal policy that maximizes the system's fitness and gathers information on the economic environment. The policy recommendation directs economic actors to create credit relationships through the optimal choice between a low interest rate or a high liquidity supply. The latter, based on the agents' balance sheet, allows to determine the liquidity supply and interest rate that the banks optimally offer **their clients within the market**. Thanks to the combination between the public and the private signal, financial institutions create or cut their credit connections over time via a preferential attachment evolving procedure able to generate a dynamic network. Our results show that the emergence of a core-periphery interbank network, combined with a certain level of homogeneity in the size of lenders and borrowers, is essential to ensure the system's resilience. Moreover, the optimal policy recommendation obtained through reinforcement learning is crucial in mitigating systemic risk.

1 Introduction

At the height of the sovereign debt crisis, the former president of the European Central Bank, Trichet, declared: “When the crisis came, the serious limitations of existing economic and financial models immediately became apparent. . . As a policy-maker during the crisis, I found the available models of limited help. In fact, I would go further: in the face of the crisis, we felt abandoned by conventional tools. . . The key lesson I would draw from our experience is the danger of relying on a single tool, methodology or paradigm. Policy-makers need to have input from various theoretical perspectives and from a range of empirical approaches. . . In this context, I would very much welcome inspiration from other disciplines: physics, engineering, psychology, biology. Bringing experts from these fields together with economists and central bankers is potentially very creative and valuable. . .” (see Trichet (2010)).

Inspired by the words of Trichet, welcoming new and multidisciplinary policy tools, in this paper, we are explicitly interested in understanding the effect that an unconventional and environmentally dependent policy recommendation has on the stability of the interbank system. From the point of view of the functioning of the interbank market, our work follows Berardi and Tedeschi (2017), where financial institutions establish preferential lending arrangements to insure themselves against the unexpected withdrawal of deposits. Financial connections might change over time via a preferential attachment evolving procedure (see Barabási and Albert (1999)) such that each agent can enter into a lending relationship with others with a probability proportional to a fitness measure. Specifically, the attractiveness of agents is based either on their high supply of liquidity or their low interest rate. The authors show how the implementation of one or the other strategy generates different architectures of the credit network, which dissimilarly impact the spread of systemic risk.

The originality of this work with respect to the one mentioned above concerns the mechanism that drives banks to choose between the two strategies. Where in Berardi and Tedeschi (2017) the choice is exogenous and fixed, here we introduce a time-dependent policy recommendation based on a reinforcement learning approach that directs banks to optimize the entire banking system’s long-term fitness. Specifically, the regulator directs the interbank system towards an optimal strategy that chooses between favoring a high liquidity supply rather than a low interest rate, by collecting information from the environment. Once the policy recommendation is made public, each bank signals **to her counterparty**

within the interbank market her optimal liquidity supply or interest rate level, which are used to establish credit agreements via the above-mentioned preferential attachment mechanism. In a nutshell, we might think that the central bank directs the interbank system to choose between interest rate and liquidity supply by announcing the interest rate corridor that it publishes periodically. The corridor dynamics determines the choice of interbank interest rates (such as EONIA) and thus directs financial institutions to choose their strategies.

Compared to Berardi and Tedeschi (2017), therefore, the reinforcement learning mechanism allows us both to endogenize and identify the optimal strategy and to model a policy recommendation useful to tame systemic risk. Although this tool is helpful for modeling reward-seeking behavior of agents in complex systems¹ (see(Osoba et al.; 2020)), to the best of our knowledge, it is barely employed in the agent-based framework. Interesting exceptions are Liu et al. (2018), and Lozano et al. (2007), which use reinforcement learning to model the credit allocation strategy of financial institutions in the interbank market. Apart from the modeling differences omitted here that distinguish us from those works, it is important to point out the methodological distinction. Where these works use a tabular reinforcement learning algorithm, as proposed by (Watkins and Dayan; 1992), we use a state-of-the-art reinforcement learning algorithm with neural network approximators (Schulman et al.; 2017), which describes the complex reward-seeking behavior. While the advantages and disadvantages of these algorithms are well documented and concern issues such as the computational efficiency, the curse of dimensionality, and the convergence (Bellman; 1956), the better performance of the neural network-powered algorithms emerge. These models are beneficial when solving complex problems where the underlying environment changes rapidly and is also defined by the different forces that relate and compete with each other. These capabilities have already demonstrated effective in solving complex financial and economics problems (see (Du et al.; 2020; Jiang et al.; 2017; Lin and Beling; 2020; Zhang et al.; 2020)).

Without delving into technical details, some clarifications on how the proposed algorithm works should be done. The selected reinforcement learning algorithm optimizes an objective function that, in our context, corresponds to the aggregate fitness of the interbank system. The optimization is carried out by training a

¹ We refer the reader to Charpentier et al. (2021) and Mosavi et al. (2020) for comprehensive reviews of different use cases of reinforcement learning in financial and economic contexts.

neural network model. The neural network receives input variables concerning the economic conditions of the interbank system and returns as output the strategy, i.e., the policy recommendation directing the system towards competing on liquidity supply rather than on the interest rate.

This family of algorithms is often criticized regarding the interpretability of inputs' impact on the results. The output, in fact, often appears as a black box whose determinants remain hidden from the user. To avoid this problem, we act in the following way. Firstly, we limit the choice of inputs to variables readily available to the regulator. To this end, we use aggregate systemic variables such as the interbank system's minimum, maximum, and average interest rate and liquidity supply. The choice of a limited set of input variables allows us to understand their effects in determining the output and to model a system with incomplete and asymmetric information (see Bernanke et al. (1999)). Secondly, we directly study each input's impact on the output's determination through the SHapley Additive exPlanation (SHAP) framework (Lundberg and Lee; 2017).

The introduction of the reinforcement learning framework into the interbank market model proposed by Berardi and Tedeschi (2017) allows us to draw some important conclusions about the systemic stability of the system and to determine some policy interventions capable of curbing contagion. Firstly, the proposed algorithm fully endogenizes the evolution of the interbank network, whose architecture, therefore, changes over time. In this way, we can identify that the topology that emerges when the policy recommendation suggests a high supply of liquidity is more resilient in the face of exogenous shocks. Also, at the individual level, this policy produces better microeconomic performance. In this circumstance, banks are less heterogeneous, which generates a uniform risk exposure among counterparties able to favor the system's resiliency. **Although not unequivocally accepted (see, for instance, Haldane and May (2011)),** the negative impact of heterogeneity on systemic stability is in line with other theoretical and empirical studies (see Caccioli et al. (2012), Iori et al. (2006) and Tedeschi et al. (2012)). On the other hand, the worse performance of a system dominated by low interest rates reflects the empirical evidence. Indeed, it is well documented that a credit market dominated by "low-for-long" interest rates adversely affects both the banks' and the economy's stability. For financial institutions, low rates might reduce resilience by lowering profitability and thus their ability to replenish capital after a negative shock. This strategy would encourage risk-taking for the system, undermining systemic stability (see Bindseil

(2018), for a general overview on the topic). Finally, our results suggest that the policy recommendation implemented via reinforcement learning can more mitigate systemic risk than alternative tools.

Related literature

The increasingly recurrent and impactful crises affecting the socio-economic system have called for a deep rethinking of the economic theory. Firstly, the literature has made an effort to understand and include in the economic models the sources of contagion. Regardless of the modeling approach used, which ranges from New Keynesian models solved globally or using reduced functional form (see, for instance, Boissay et al. (2016), Gertler et al. (2020), Svensson (2017)) to agent-based models and the most recent network-oriented approaches (see Battiston et al. (2012a,b), Georg (2013), Haldane and May (2011), Upper (2011), Capponi et al. (2020), Calice et al. (2020)), there is a general agreement that identifies interaction and heterogeneity as the drivers of endogenous crises. Moreover, the post-Lehman studies have placed particular emphasis on the propagation of contagion, determining the direction of the attack from financial to real markets and its fuse in the portfolio structure of financial institutions (see Brunnermeier et al. (2012)). Many interesting studies, for example, have identified the source of contagion in the asset or liability side of banks' balance sheets. Among them, the effect of the fire-sale price and the (re)payment system between creditors and debtors have proven to be particularly important in generating financial instability (see Acharya and Yorulmazer (2008a), Angelini et al. (1996), Dasgupta (2004) Rochet and Tirole (1996)). In this vein, maturity transformation, sharing risk, herding behavior, and interbank linkages are just some of the various components able to trigger instability or collapse in financial markets (see Acharya and Yorulmazer (2008b), Allen and Gale (2000) and Tedeschi et al. (2021), among the many).

Once the origin of the disease and the channels through which it spreads have been identified, the literature has turned to treatment, that is, identifying the best tools to mitigate financial contagion. The scientific community has heavily focused on developing new tools to overcome systemic instability. In this regard, several conventional and non-conventional monetary policies and other alternative tools have been proposed. However, their effects on financial stability are controversial and depend on the overall economic condition (see Goldberg et al. (2020), and Altavilla et al. (2021)). A strand of literature, for example,

has emphasized the importance of a strict, rule-based, and predictable monetary policy to tame systemic risk (see Jiménez et al. (2014) and Taylor (2011)). On another side, instead, different studies have bet on alternative rules, compatible with the underlying economic conditions (see Boissay et al. (2021), De Grauwe (2011) and Galí (2015)). Unfortunately, the weak empirical evidence, due to the fairly recent development of these alternative techniques, which also include the so-called macro-prudential policies, makes it difficult to prove the supremacy of one approach over the other. While the empirical facts are still uncertain, recent theoretical models have attempted to resolve this "certamen". An interesting contribution in this direction is the model of Boissay et al. (2021). The authors use a globally solved New Keynesian model with heterogeneous agents to generate endogenous crises. The paper compares two monetary policy instruments, one that follows a strict inflation targeting rule and the other that allows the central bank to curb financial booms and busts. The authors show how the policies that mitigate output fluctuations help prevent financial crises by acting on agents' expectations. In support of cyclical policies determined by the economic background, there are also many agent-based models (see, Cincotti et al. (2012), Giri et al. (2019) and Riccetti et al. (2018), among the many). Generating complex dynamics in evolving systems is an ideal environment for testing the effect of (un)conventional policies/measures on financial stability.

The rest of the work is organized as follows. In Section 2 we present the functioning of the interbank market, placing particular emphasis on the evolution of the credit network and the implementation of the reinforcement learning algorithm. In Section 3 we show the results. Specifically, we follow three steps: firstly, we verify the performance and robustness of the reinforcement learning algorithm; secondly, we investigate its implication on the interbank network morphology and the performances of the financial institutions; thirdly, we present the effect on the interbank systemic stability of the policy recommendation. Finally, Section 4 concludes with some remarks on the achieved results and the provided contribution.

2 Model

This section describes the formation and evolution of credit relationships between financial institutions. Due to unexpected future movements of deposits, banks enter into preferential lending agreements to have a potential credit channel when

needed. These lending agreements are fast lanes created before use, and their set defines a *potential interbank network*. In order to build their preferential lending agreements, banks **report their credit conditions to their customers** through an attractiveness measure. We model bank fitness as a combination of a policy recommendation and private information. The first ingredient is a signal obtained via a reinforcement learning mechanism, through which the regulator directs banks to choose the best strategy given the underlying environmental conditions. In particular, the regulator recommends the weight to assign to high liquidity supply rather than to low interest rates, thus directing the competition. The second ingredient is a private signal, based on the bank's capital structure, consisting of the actual interest rate and credit provision offered. Potential credit relationships might change over time via a preferential attachment evolving procedure that depends on bank fitness. As the deposit shock materializes, financial institutions face liquidity surpluses or shortages, which induce them to exploit their preferential lending agreements and enter the interbank market as lenders or borrowers. At this point, the previously potential network becomes an *active credit network*. Only the potential links of the banks facing a liquidity shortage are activated and correspond to a very sparse network.

2.1 The interbank market microstructure

We consider a sequential economy operating in discrete time, which is denoted by $t = \{0, 1, 2, \dots, T\}$. At any time t , the system is populated by a large number N of active banks $i, j \in \Omega = \{1, \dots, N\}$. Financial institutions interact with each other through credit relationships represented by the set V_t , whose elements are ordered pairs of different banks. Banks (nodes or vertices) and their connections (edges or links) form the interbank network $G_t = (\Omega, V_t)$. The daily balance sheet structure of each bank is defined as

$$L_t^i + C_t^i + R_t^i = D_t^i + E_t^i, \quad (1)$$

where assets are on the left-hand side and liabilities are on the right-hand one. In particular, L , C , and R represent long-term assets, liquidity, and reserves, while D and E deposits and equity of bank i at time t . Reserves are a portion of deposits, $R_t^i = \hat{r}D_t^i$, where the required reserve rate², \hat{r} , meets the legal requirement of 2%.

² This rate replicates a central bank regulation that sets the minimum amount that a commercial bank must hold in liquid assets and is com-

At every time t , deposits are exogenously shocked, and the balance sheet in Eq. 1 modifies accordingly. Specifically, deposits evolve as

$$D_t^i = D_{t-1}^i(\mu + \omega U(0, 1)), \quad (2)$$

with $U(0, 1)$ a uniformly distributed noise between 0 and 1 and μ and ω modeling the expected number of negative shocks and thus different market conditions. On the one hand, financial institutions with a negative change in deposits and subject to a complete erosion of their liquidity become potential debtors in the interbank market. On the other hand, banks that suffer a small negative shock or an increase in deposits become potential creditors to the system. Consequently, the respective demand d_t^i and supply s_t^i of liquidity of potential borrowers and lenders are given by

$$\begin{aligned} &\text{borrower if: } \Delta D_t^i + C_t^i \leq 0, \text{ with demand of liquidity } d_t^i = |\Delta D_t^i + C_t^i| \\ &\text{lender if: } \Delta D_t^i + C_t^i > 0, \text{ with supply of liquidity } s_t^i = \Delta D_t^i + C_t^i. \end{aligned}$$

Since we do not assume a Walrasian tâtonnement mechanism, the system may endogenously generate a mismatch between credit supply and demand. Moreover, since the interbank network is not fully connected, even at a micro level, the demand for liquidity of a borrower bank might not match the credit supply offered by the lender banks connected to it. Specifically, we define the granted loan from a generic lender i to a generic borrower j as $l_t^{i,j} = \min(s_t^i, d_t^j)$. Borrowing banks rationed in the interbank market can sell their long-term assets at a fire-sale price as a method of last resort. The amount of loan the borrower has to sell for covering its residual liquidity need is equal to $\Delta L_t^j = \frac{d_t^j - s_t^i}{\rho}$, where ρ is the 'fire-sale' price

At the beginning of the next day, the repayment round takes place. Financial institutions encounter a new deposit movement that increases or decreases their liquidity. On the one hand, lending banks facing a positive (negative) change in deposits remain potential creditors (became potential debtors). On the other hand, borrowing banks face different scenarios depending on whether the deposit shock is positive or negative. Specifically, in the case of a positive shock, it can

monly referred to as the reserve ratio. The central bank determines this minimum amount based on a specified proportion of bank deposit liabilities.

happen that: i) the change in deposits is sufficient to repay the principal and the interest, or ii) the deposit variation is insufficient to cope with the loan. In the first case, the debtor can quickly meet her obligations, but in the second case, she must sell a number of long-term assets sufficient to repay the creditor at a fire-sale price fully. On the other hand, in the case of a negative shock, banks must sell their long-term assets to pay for previous interbank borrowings and meet the new liquidity needs. All institutions that do not raise enough liquidity to meet their obligations via the fire sale fail, thus creating a bad debt on the lender. The creditor's loss, $B_t^{i,j}$, is equal to the granted loan after liquidating the debtor assets. Hence the equity of the bank i obeys the following law of motion:

$$E_t^i = E_{t-1}^i + \sum_j l_{t-1}^{i,j} r_{t-1}^{i,j} - \sum_{j \in \theta_t^i} B_t^{i,j} - (1 - \rho) \hat{L}_t^j, \quad (3)$$

where the second term on the right-hand side is the repayment, at the agent-specific interest rate $r^{i,j}$, of the granted loan $l^{i,j}$, and the third term is the bad debt of the subset of the bank i clients, θ_t^i , unable to repay their debts because they go bankrupt and the last term represents fire sales. If the bank has not fulfilled the loan requirements (i.e., if she cannot repay the principal and interest in full), the lender no longer provides credit, forcing her to exit the market. Thus, the borrower exits the market when assets fall short of liabilities, that is $E_t^i < 0$. The failed banks leave the market. The banks exiting in t are replaced in $t + 1$ by new entrants, which are, on average smaller than incumbents. So, entrants' size is drawn from a uniform distribution centered around the mode of the size distribution of incumbent banks (see Bartelsman et al. (2005)).

2.2 Banks microfoundations: the dynamics of lending agreements and trading strategies

In order to meet their liquidity needs, at the beginning of each day, agents meet in the interbank market and sign bilateral potential lending agreements representing the directed links $(i, j) \in V_t$. These agreements can be interpreted as credit lines, which are valid during t , and can be used at the request of the borrower j in case of the lender i available liquidity. The set of all potential lending agreements reproduces the potential interbank network topology³.

Let us now explain in detail the mechanism that governs the formation/evolution

³ The creation of these links predates the deposit shock, which is why they are potential. These credit lanes, common in interbank markets, can be

of credit relationships between financial institutions. We assume banks are risk-neutral agents operating in a perfect competition environment to optimize their expected profit. The bank i expected profit for a loan provided to j is given by

$$\mathbb{E}[\Pi_t^{i,j}] = p_t^j(r_t^{i,j} c_t^{i,j}) + (1 - p_t^j)(\xi A_t^j - c_t^{i,j}) + \phi A_t^j - \chi A_t^i, \quad (4)$$

where p_t^j is the probability that the borrower does not fail, $r_t^{i,j}$ the interest rate asked by the lender i to the borrower j , $c_t^{i,j}$ the maximum amount i is willing to lend to j . Moreover, ξ is the liquidation cost of assets, A_t^j , pledged as collateral, and ϕ and χ the screening costs of creating a credit link that decrease with the debtor dimension and increase with the creditor size (see Dell’Ariccia and Marquez (2004), and Maudos and De Guevara (2004), for empirical evidence). Specifically, Eq. 4 captures the lender’s expected revenue if the borrower does or does not meet her obligations (the first and the second term on the right side, respectively), and the opportunity cost of the agreement (last two variables in Eq. 4). Moreover, we apply a heuristic rule to model a proxy for the debtor’s j survival probability. Recalling that the borrower fails if her equity becomes negative, $E_t^j < 0$, the probability of surviving is given by the closeness between j ’s equity and the highest net-worth in the system, i.e.

$$p_t^j = \frac{E_t^j}{E_t^{\max}}. \quad (5)$$

The bank’s probability of surviving is connected to financial fragility. A financial institution leaves the system if her net worth is so low that an adverse shock makes it negative or if she suffers a loss so huge as to deplete all the net worth accumulated in the past (see Greenwald and Stiglitz (1993)). The Eq. 5 can also be interpreted as a rule of thumb for determining the risk premium that lenders charge to a borrower⁴. Finally, the maximum amount that the lender i is willing to lend to j , that is,

interpreted as mutual ’promises’ of help between financial institutions in case of liquidity needs.

⁴ We acknowledged that the simulation results are robust even when implementing a survival probability where E_t^{\max} does not change over time. Specifically by using in the denominator of Eq. 5 the average maximum equity over all the timestep.

the lending capacity, $c_t^{i,j}$, in Eq. 4 is defined as

$$\begin{cases} c_t^{i,j} = (1 - h_t^j)A_t^j > 0, \text{ if } (i, j) \in V_t, \\ c_t^{i,j} = 0 \text{ otherwise,} \end{cases}$$

with $h_t^j \in (0, h_t^{max})$ to be the borrower haircut, defined as the j 's leverage, λ_t^j , with respect to the maximum one. Hence $h_t^j = \frac{\lambda_t^j}{\lambda_t^{max}}$, with $\lambda_t^j = \frac{L_t^j}{E_t^j}$. By setting Eq. 4 equal to zero and rewriting it as a function of $r_t^{i,j}$, we get the interbank rate that guarantees zero expected profit:

$$r_t^{i,j} = \frac{\chi A_t^i - \phi A_t^j - (1 - p_t^j)(\xi A_t^j - c_t^{i,j})}{p_t^j c_t^{i,j}}. \quad (6)$$

In line with the assumption of asymmetric information and costly state verification (see Bernanke et al. (1999)), the lender applies an interest rate that increases with her size⁵ (that is, her assets) and the financial vulnerability of the borrower (that is j 's leverage). This last implication derives from the budget identity (see Eq. 1) from which we can derive that $A_t^j = \frac{L_t^j}{\lambda_t^j} + D_t^j$, where $\lambda_t^j = \frac{L_t^j}{E_t^j}$. In addition, the interest rate in Eq. 6 is not linearly related to the bank's probability of surviving and capacity.

We now have all the elements to describe how traders select their counterparts in the interbank system, i.e., how lending arrangements are formed and evolve. We develop a measure of agent attractiveness to generate an endogenous preferential attachment mechanism. **Specifically, banks signal themselves to their pool of clients based on their low interest rate or abundant supply of liquidity. The dichotomy between these two strategies is microfounded and stems from the expected profit of banks (see Eq. 4), where the screening costs of creating a credit link increase with the creditor size. This implies, as shown in Eq. 6, that lenders attractive in terms**

⁵ The relationship between screening costs and the interest rate has been widely explored in the economic literature and often associated with the imperfect information paradigm (see Aleem (1990); Bester (1985); Hoff and Stiglitz (1990)). Following this interpretation, the explanation for the high interest rate lies in the problem of asymmetric information. Specifically, lenders having less information than borrowers about the latter's ability and willingness to repay a loan have to screen applicants and charge the cost of this operation to borrowers. However, it is infrequent to find evidence about the costs associated with screening and, more generally, about the effect of imperfect information on the behavior of credit market participants.

of higher liquidity supply offer higher interest rates. Symmetrically, banks offering low interest rates are necessarily less liquid⁶.

Although all agents start from the same initial conditions, financial institutions are characterized by heterogeneous levels of their agent-specific variables as time goes by. In line with this, the fitness of each agent μ_t^i is a combination of her liquidity relative to the highest liquidity provided in the market, C_t^{\max} , and her interest rate compared to the cheapest one, r_t^{\min} , i.e.

$$\mu_t^i = \eta_t \left(\frac{C_t^i}{C_t^{\max}} \right) + (1 - \eta_t) \left(\frac{r_t^{\min}}{r_t^i} \right). \quad (7)$$

The parameter η_t reflects a policy recommendation at time t , addressing the choice of the banking sector towards one of two possible strategies. On the one hand, η approaching zero identifies an interbank system moving towards the cheapest interest rates. On the other hand, η close to one highlights a liquidity-based system. **The signal disseminated by the regulator that directs the system toward the optimal strategy can be interpreted as the central bank's announcement of the interest rate corridor. This corridor conditions the interbank interest rate and, consequently, the choice of each financial institution on her credit condition (see Giannone et al. (2011)).** We refer the reader to Subsection 3.1 for a detailed explanation of the policy recommendation evolution. One of the main contributions of our work is to assume η_t endogenously evolving through a reinforcement learning mechanism, modeling the regulator's will to address the banking system toward the best credit strategy for system stability. **It is worth emphasizing that, although in Eq.7 the public signal is homogeneous in the baseline model, banks' attractiveness remains highly heterogeneous as the private signals on the liquidity, C_t^i , and interest rate, r_t^i , are agent-specific. Let us assume, for example, that the system is directed towards a low-interest rate, $\eta = 0$. Since interest rates in the fitness measure are bank-specific, interest rates applied by lenders to their clients are different. Further-**

⁶ Assuming screening costs that increase with borrower's dimension and decrease with the lender's dimension implies an inverse relationship between the lender's size and the interest rate the financial institution offers on the interbank market: the most liquid lender provides the best conditions in terms of interest rate. In this circumstance, the two banks' strategies collapse into the same. Since the banks' strategies go in the same direction, their impact on the simulated dynamics is similar, and the reinforcement learning mechanism achieves precisely the same effects as a random choice, given the perfect overlap of the two tactics.

more, the liquidity supply of those lenders chosen to grant credit is also agent-specific, which ensures heterogeneity in granting credit. A similar dynamic applies to the case where the signal directs toward a high liquidity supply, i.e., $\eta = 1$. In other words, the only element of homogeneity is the public signal that directs the system toward the optimally selected strategy⁷.

Regarding our interbank network model, credit links are directional because they are created and deleted by the agent j who looks for a loan and points to the agent i that provides credit. The information on credit conditions (and then loan) flows opposite. **It is worth noting that credit terms are bilateral (between creditor and debtor) and, therefore, not available from other market members.**

In general local interaction models, the agent interacts directly with a finite number of counter-parties in the population. The set of nodes with which a single node is linked is called its neighborhoods. In our model, the number of outgoing links is constrained to be a small number \hat{d} . Thus borrowers can only get loans from \hat{d} lenders. With this assumption of network sparsity, the topology is always locally tree-like, avoiding loops that would preclude us from fully understanding the network architecture's impact on economic dynamics, such as systemic risk, failures, and liquidity diffusion.

At the time $t = 0$, each bank j starts having \hat{d} random outgoing links (i.e., potential borrowing positions) and possibly with some incoming links from other agents (i.e., potential lending position). At the beginning of each period, links are rewired in the following way. For any outgoing link i , each borrower j randomly selects a new bank k . Comparing the fitness of the new financial institution with the one of its previous lender i , the borrower j cuts her old link with i and creates a new one with k according to the probability

$$P_t^j = \frac{1}{1 + e^{-\beta(\mu_t^k - \mu_t^i)}}, \quad (8)$$

or keep its previous link with probability $1 - P_t^j$. The proposed mechanism for reviewing credit agreements ensures that the most attractive lenders get the highest number of borrowers (i.e., incoming links) and earn the highest profits. Nevertheless, the degree of randomness in the algorithm guarantees that some

⁷ This assumption is modified in Sec 3.3, where heterogeneity is also introduced in the public signal.

links with very high-performing agents may be cut in favor of less attractive creditors. The amount of randomness is regulated by β and has a double purpose: from a practical point of view, it prevents the system from being centralized around a single financial hub; from a theoretical perspective, it allows us to model incomplete information and bounded rationality.

The evolution of the banking system: determining the policy recommendation

As anticipated in the previous section, we use the reinforcement learning paradigm to move the parameter η_t and obtain an optimal policy recommendation in the described banking system. Reinforcement learning aims to solve a decision-making problem in which the timing of costs and benefits is relevant. In an interbank market that follows the specified dynamics for the creation of lending agreements, reinforcement learning can help determine the policy recommendation that better identifies the optimal attachment strategy to follow in Eq. 7, even when partial information about the system is provided. Hereafter, we refer to the reinforcement learning algorithm as the learning algorithm. A Markov Decision Process (MDP) is the mathematical formalism under which the reinforcement learning problem is usually defined. A MDP comprises of a set of possible states $S_t \in \mathcal{S}$, a set of possible actions $A_t \in \mathcal{A}$ and a transition probability $P[S_{t+1} = s' \mid S_t = s, A_t = a]$. At each time t , a learning agent that is in state S_t , takes an action A_t and receives a reward $R_{t+1}(S_t, A_t, S_{t+1}) \in \mathbb{R}$ from the environment before moving to the next state S_{t+1} . We define the agent strategy $\pi : \mathcal{S} \times \mathcal{A} \mapsto [0, 1]$ as the conditional probability $\pi(a \mid s)$ of taking the action $A_t = a$ being in the state $S_t = s$. The reinforcement learning problem is the stochastic control problem of maximizing the expected discounted cumulative reward

$$\mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t R_{t+1}(S_t, A_t, S_{t+1}) \right], \quad (9)$$

where $\gamma \in [0, 1)$ is a discount factor, and the expectation is w.r.t. the sequence of states and actions reached following the strategy π .

In our MDP, the sequential economy in which the banking system operates plays the role of the environment. Banks interact with the environment by changing their credit lines: each day, they can adapt their attachment strategy between liquidity supply and interest rate discount, which is regulated through Eq. 7, with the choice of η_t , playing the role of the action A_t . We assume

the agent is the system as a whole rather than the single bank and that the optimal strategy is realized at the system level, i.e., that the regulator directs financial institutions towards the correct combination of the two strategies. This assumption has a twofold purpose. On the one hand, it helps us to model a system with incomplete/asymmetric information, where the central bank has richer information set than the single economic actor (see, for instance, Hoff and Stiglitz (1990), and Thakor (2020)). On the other hand, it allows us to incorporate economic policy, seen as the optimal indication that the regulator gives to the system to reduce the interbank market vulnerability (see Trichet (2010), for a global overview)⁸. **As shown above, the central bank’s recommendation is made through the optimal interest rate corridor announcement, which conditions interest rates and the liquidity supplied by financial institutions.**

The state S_t includes information on both the liquidity C_t and the interest rate r_t distributions of the banking system. Specifically, the state space is defined as

$$S_t = (C_t^{\max}, C_t^{\min}, r_t^{\max}, C_t^{\text{avg}}, r_t^{\min}, r_t^{\text{avg}}),$$

where $x_t^{\max} = \max_{i \in \Omega} x_t^i$, $x_t^{\min} = \min_{i \in \Omega} x_t^i$, $x_t^{\text{avg}} = \sum_{i=1}^N x_t^i / N$, being x the variable of interest. We believe that this state-space setting is realistic enough to model the partial information of the regulator about the banking system: it would be difficult and costly to retrieve detailed and specific data on all the banks included in the system at each time step. It is easier to gather information about the best and the worst liquidity provider in the interbank network as much as average estimates of the entire market.

Finally, the reward function we consider is the system’s total fitness

$$R_t(S_t, A_t, S_{t+1}) = \sum_{i=1}^N \mu_t^i \quad (10)$$

Moreover, the problem in Eq. 9 becomes a maximization of the discounted cumulative banks’ total fitness. From the definition of bank fitness, this means guaranteeing a better flow of liquidity through the banking system and an efficient

⁸ Considering η as a system variable allows us to reduce the problem’s mathematical and computational complexity and study the banking system’s behavior as a whole. Making η bank specific leads towards multi-agent reinforcement learning applications (Buşoniu et al. (2010)), which consider agents that compete with each other and are an out-of-the scope of the present paper.

allocation at a more convenient interest rate. **We recall here that maximizing the fitness of financial institutions corresponds to optimizing their expected profit. The motivation behind this modeling assumption is twofold. Firstly, for the recommendation to be followed by the banks, it must have a goal of interest to the banks themselves, namely their profit. Second, the regulator, by maximizing the fitness of the system, succeeds ex-post in safeguarding the resilience of the financial system, given the inverse relationship between expected profits and failures of financial institutions.**

The learning algorithm operates in a model-free setting because it only receives partial information on the relevant variables of the system. At the same time, it has no knowledge of the internal dynamics (i.e., transition probability) with which the banks' balance sheets moves and lending agreements are generated. This information has to be inferred through the sequence of states, actions, and rewards during the learning process.

2.3 The optimization algorithm: Proximal Policy Optimization

The optimization problem in Eq. 9 can be solved using a policy gradient algorithm like the Proximal Policy Optimization (PPO) (Schulman et al.; 2017). A policy gradient algorithm directly parametrizes the optimal strategy $\pi_\theta = \pi(a | s; \theta)$, for example using a multilayer neural network with parameters θ . The optimization problem is approximately solved by computing the gradient of the cumulative fitness of the system $J(\theta) = \sum_{t=0}^{\infty} \gamma^t R_{t+1}(S_t, A_t, S_{t+1}; \pi_\theta)$ and then carrying out gradient ascent updates according to

$$\theta_{t+1} = \theta_t + \alpha \nabla_\theta J(\theta_t), \quad (11)$$

where α is a scalar learning rate. The policy gradient theorem (Marbach and Tsitsiklis (2001); Sutton et al. (2000)) provides an analytical expression for the gradient of $J(\theta)$ as

$$\begin{aligned} \nabla_\theta J(\theta) &= \mathbb{E}_{\pi_\theta} \left[\frac{\nabla_\theta \pi(A_t | S_t; \theta)}{\pi(A_t | S_t; \theta)} Q_{\pi_\theta}(S_t, A_t) \right] \\ &= \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi(A_t | S_t; \theta) Q_{\pi_\theta}(S_t, A_t)], \end{aligned} \quad (12)$$

where the expectation, with respect to (S_t, A_t) , is taken along a trajectory (episode) that occurs adopting the strategy π_θ and the action-value function

$$Q_\pi(s, a) \equiv \mathbb{E} \left[\sum_{k=0}^{\infty} \rho^k R_{t+1+k} \mid S_t = s, A_t = a, \pi \right], \quad (13)$$

represents the long-term reward associated with the action a taken in the state s if the strategy π is followed hereafter. It can be proven that it is possible to modify the action value function $Q_\pi(s, a)$ in (12) by subtracting a baseline that reduces the variance of the empirical average along the episode while keeping the mean unchanged. A popular baseline choice is the state-value function

$$V_\pi(s) \equiv \mathbb{E} \left[\sum_{k=0}^{\infty} \rho^k R_{t+1+k} \mid S_t = s, \pi \right], \quad (14)$$

which reflects the long-term reward starting from the state s if the strategy π is adopted onwards. The gradient thus can be rewritten as

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi(A_t \mid S_t; \theta_t) \mathbb{A}_{\pi_\theta}(S_t, A_t)] \quad (15)$$

where

$$\mathbb{A}_\pi(s, a) \equiv Q_\pi(s, a) - V_\pi(s), \quad (16)$$

is called advantage function and can be interpreted as the gain obtained by choosing a specific value of a in a given state with respect to its average value for the strategy π .

Different policy gradient algorithms derive from the way the advantage function is estimated. In PPO the advantage estimator $\mathbb{A}(s, a; \psi)$ is parameterized by another neural network with parameters ψ . This approach is known as actor-critic: the actor is represented by the policy estimator $\pi(a|s; \theta)$ that outputs a probability for each possible value of $a \in \mathcal{A}$, which the learning algorithm uses to sample actions, while the critic is the advantage function estimator $\mathbb{A}(s, a; \psi)$ whose output is a single scalar value. The two neural networks interact during the learning process: the critic drives the updates of the actor, which successively collects new sample sequences that will be used to update the critic and again evaluated by it for new updates. The extended objective function can therefore describe the PPO algorithm

$$J^{\text{PPO}}(\theta, \psi) = J(\theta) - c_1 L^{\text{AF}}(\psi) + c_2 H(\pi(a \mid s; \theta)) \quad (17)$$

where the second term is a loss between the advantage function estimator $\mathbb{A}(s, a; \psi)$ and a target \mathbb{A}^{targ} , represented by the cumulative sum of discounted reward, needed to train the critic neural network. The last term represents an entropy bonus to guarantee an adequate level of exploration. Details about the specific choice of the target, together with additional information about the general algorithm implementation, are given in the App B. In what follows, PPO can be generally referred to as the learning algorithm.

3 Simulation Results

In this section, we perform numerical experiments to test the capability of the learning algorithm to identify an optimal strategy for selecting η and trading off the two competing ways of establishing credit relationships. In this respect, we analyze the effects of the η dynamics on agents' economic performances, the interbank network topology, and its resilience in the face of exogenous shocks. Finally, we study the effect of the policy recommendation obtained through reinforcement learning in controlling credit crunch phenomena and mitigating systemic risk.

The results provided in the following subsections are obtained from simulated tests, which shares some choices for the parameter involved in the dynamic simulation of the system. The number of Monte Carlo simulations performed is $M = 200$, and each simulation is $T = 1000$ periods long. We simulate a system with $N = 50$ banks whose out-degree is $\hat{d} = 1$, so each bank can obtain at most one outgoing link at each time step while can have many possible incoming links. Each bank is subjected to an initial probability of being isolated, set at 0.25. The parameters of the screening costs χ and ϕ that enters in Eq. 6 are set respectively at 0.015 and 0.025, while the liquidation cost of collateral ξ is 0.3. The parameters μ and ω shifting the uniformly distributed noise that shocks the bank deposits are set at 0.7 and 0.55. All the banks starts with the same initial interest rate equal to 2% and are endowed with the same initial balance sheet $C_0 = 30$, $L_0 = 120$, $D_0 = 135$ and $E_0 = 15$. The price of fire sale $\rho = 0.3$ and the intensity for breaking the connection between banks $\beta = 5$ in Eq. 8 are other parameters common to all the agents in the network. In the App. A we check the robustness of our qualitative results by changing some key parameters. Specifically, we vary the intensity of choice, β , from 0 to 40 with steps of 2; the fire-sale price, ρ , from 0.1 to 0.5 with steps of 0.1 and, finally, the parameter ω regarding the volatility shock on bank deposit. We have then studied the

moments of the distributions of the statistics of interest. Results confirm that our findings are robust to some variations of the banking system simulation.

The PPO algorithm parametrizes a discrete strategy function so that the learning algorithm can choose the value of η among a finite set of actions $\mathcal{A} = \{0, 0.5, 1\}$ ⁹

3.1 Training the PPO algorithm

As the first step in our numerical analysis, we evaluate the performance of the strategy learned by the PPO algorithm. We train four PPO instances on $E_{in} = 1000$ consecutive episodes, which are independent simulations of the banking system. The PPO instances differ for the random seed used to initialize the neural networks and to train them using a stochastic gradient descent approach. Multiple concurrent training of different instances is needed to provide an average performance together with a confidence interval that highlights the robustness of the learning process. Each training episode consists of a simulation of the banking system for T periods that allow the learning algorithm to collect samples of data with which it can perform updates of the model parameters. During the learning phase, we evaluate the learning progress of each instance at several intermediate steps. We fix the weights of the neural networks that parametrize the η public signal and perform $E_{out} = 5$ out-of-sample test episodes before carrying on the training process to assess the learned behavior up to that point. We refer to the App. B for the technical difference between an in-sample and an out-of-sample test episode.

After training the PPO algorithm, the reinforcement learning agents tends to select only the extremes of the set $\mathcal{A} = \{0, 0.5, 1\}$, which corresponds to an interest rate strategy ($\eta = 0.0$) or a liquidity strategy ($\eta = 1.0$). For this reason, we highlight such a dichotomy in the baseline model since it is the pattern that emerges when all the banks in the

⁹ Under the same setting, training PPO instances that are allowed to pick fine-grained discrete values between 0 and 1 as a possible action is computationally expensive because the algorithm needs to explore a broader set of possible state-action pairs. Such an implementation would let the algorithm runtime grow and would not necessarily improve the results because the algorithm would not be able to alias between consecutive actions. A fine-grained action space \mathcal{A} would make the η -strategy less interpretable. Hence in our analysis, we decided to distinguish three specific scenarios, which are the two extreme cases ($\eta = 0.0$ and $\eta = 1.0$) and the middle case ($\eta = 0.5$).

system follow the policy recommendation. Considering the emerging dichotomy in the selected action, we compare the PPO performance with respect to a dynamic random baseline that picks the value of η according to a Bernoulli distribution with a parameter equal to 0.5. This random policy that chooses between 0 and 1 with equal probability represents a meaningful benchmark, as we observe in the left-hand side of Fig. 1, where the values of η in both scenarios are identically distributed over the 200 performed Monte Carlo simulations. The Kolmogorov-Smirnoff test statistically confirms up to the 1% confidence level that the distribution of the η values generated by the selected¹⁰ PPO instance is not significantly different from the one of the random baseline. The right-hand side of Fig. 1 summarizes the learning process results where the system’s average cumulative fitness in Eq. 10 is represented on the y -axis. Every PPO instance is tested $E_{out} = 5$ times using Monte Carlo simulations of length T . We notice that the performance metric is always greater for PPO than for the random recommendation, signaling that the banks in the system generated by the PPO signal tend to be more attractive for the borrowers by exhibiting a higher aggregated fitness through time. Moving η randomly causes banks to be less attractive to the borrowers in their interbank market. This result implies that the PPO instances learn to choose the value of η by leveraging the information available about the system without changing the distribution of the values with respect to the random case. The learning procedure allows us to discover when it is convenient to pick a side in this trade-off. **A further comparison with some fixed signals and a decentralized mechanism for the η dynamic are provided in Sec.3.4.** However, fixing the η for all time steps has an evolutionary impact on the system, which has already been studied Berardi and Tedeschi (2017) and is not centered on studying the effect of an η that changes through time.

¹⁰ It is common in reinforcement learning applications to train different instances of the same algorithm and then select the best performing one over some out-of-sample tests (Andrychowicz et al.; 2020)

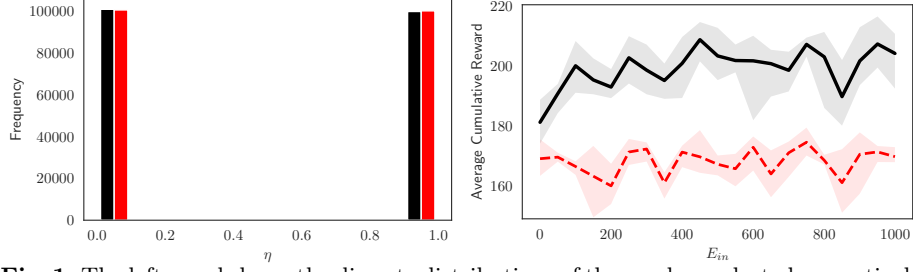


Fig. 1: The left panel shows the discrete distributions of the η values selected respectively by PPO (in black) and by a Bernoulli distribution (in red) with a parameter equal to 0.5 over 200 Monte Carlo simulations of the system. The right panel shows the average cumulative fitness of the system as a function of the number of training episodes for the trained PPO instances (in solid black) and the Bernoulli distribution of η (in dashed red) with the corresponding confidence intervals.

In order to shed light on the decisions taken by the best performing trained PPO instance, we use the SHapley Additive exPlanation (SHAP)¹¹ framework (see Lundberg and Lee (2017), Shapley (2016)). This approach explains a complex nonlinear model like a neural network by shedding light on the contribution of each input feature to the output formation. For each input vector $x \in \mathbb{R}^K$ and a model f , the SHAP value $\phi_i(f, x)$, $i = 1, \dots, K$ quantifies the effect (in a sense, the importance) on the output $f(x)$ of the i -th feature. To compute this effect one measures, for any subset $S \subseteq \{1, \dots, K\}$, the effect of adding/removing the i -th feature to the set, i.e. $f_{S \cup \{i\}}(x) - f_S(x)$. The SHAP value is defined as the weighted average

$$\phi_i(f, x) = \sum_{S \subseteq \{1, \dots, K\} \setminus \{i\}} \frac{|S|! (K - |S| - 1)!}{K!} [f_{S \cup \{i\}}(x) - f_S(x)], \quad (18)$$

where the weights ensure that $\sum_i \phi_i = f(x)$.

Figure 2 shows the magnitude of the Shapley values for the policy recommendation learned by the best performing PPO instance referred to the two possible outcomes $\eta = 0$ and $\eta = 1$. The left-hand side shows that high values for the maximum liquidity available in the system tend to favor the choice of an η based on the interest rate. Also, a low average interest rate and a high maximum interest rate point to the choice of $\eta = 0$. The right-hand side shows an opposite input relevance with a dominant role for high values of the average interest rate

¹¹ For the implementation, we use the Python package linked to Lundberg and Lee (2017)

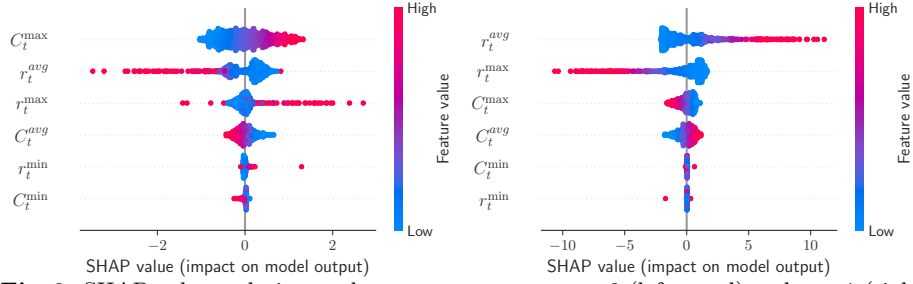


Fig. 2: SHAP values relative to the strategy outputs $\eta = 0$ (left panel) and $\eta = 1$ (right panel). The cloud of colored dots for each input variable expresses the importance and the correlation with respect to the model output. Features are ordered on the y -axis by relevance, so the first on the top influences the model output the most.

and low values of the maximum interest rate. The two figures show that the trained learning algorithm chooses one of the two signals by looking at the main characteristics of the opposite one. When it chooses $\eta = 0$, it is more interesting to know if there are participants in the network who are large. In contrast, when it chooses $\eta = 1$, it looks for homogeneity of interest rate, a common feature obtained by always playing towards the interest rate. The learning algorithm suggests a switch towards the other competing recommendation to avoid extreme cases in which a disadvantage of one or the other choice exacerbates. For instance, a huge financial institution that gathers all the borrowers' demand when $\eta = 1$ could not be sustainable in the long term, so the algorithms suggest switching to the other option. On the other hand, most medium-size banks offer medium rates when $\eta = 0$ could not gather enough liquidity to deal with deposit shocks, and it would be better to resort to the opposite signal.

3.2 Micro and macro consequences of the policy recommendation

In this subsection, we deal with the implications that the dynamics of the η parameter have on the interbank network morphology and the resulting performances of the financial institutions. Finally, we study the effects of the emerging network topology on the market's stability. **All network-related results presented in the following Sessions refer to the active credit network.**

Topology and evolution of the interbank network

Before starting the analysis, it is worth remembering the dynamics of η , that appears in the banks' fitness (see Eq. 7), determines the probability of creating credit links in the system as shown in Eq. 8. Therefore, it is appropriate to begin the analysis by describing the topology of the interbank network.

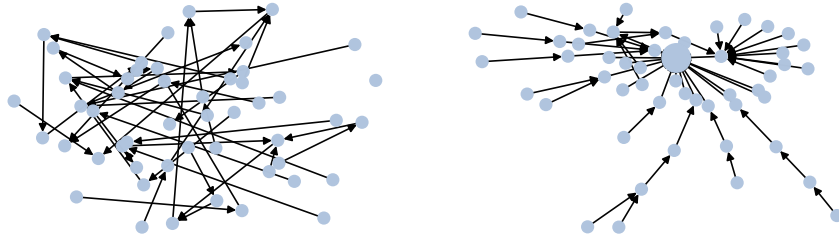


Fig. 3: Network configuration at time $t=0$ (left side), and $t=800$ (right side).

In Fig. 3, we plot the configuration of the endogenous interbank network at two different time steps of a single simulation of the system. As the reader can appreciate, the market configuration goes through different phases ranging from a random topology with isolated agents to a highly centralized architecture where a few hubs compete in credit supply. A more detailed analysis of the evolution of the interbank network architecture over time can be found on the left-hand side of Fig. 4, where we show the time series of network degree centrality

$$C_t^{\text{Net}} = \frac{\sum_i (k_t^{\text{max}} - k_t^i)}{N(N-1) - |V_t|}, \quad (19)$$

where N is the number of banks, $|V_t|$ is the total number of incoming links in the system, k_t^i is the number of incoming links for the i -th bank, and k_t^{max} is the number of incoming links holds by the hub of the network.

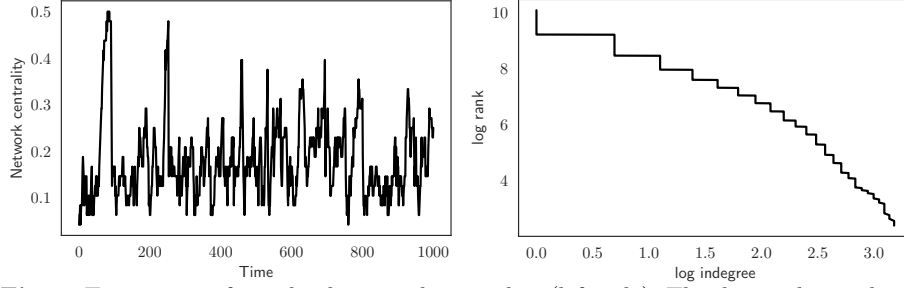
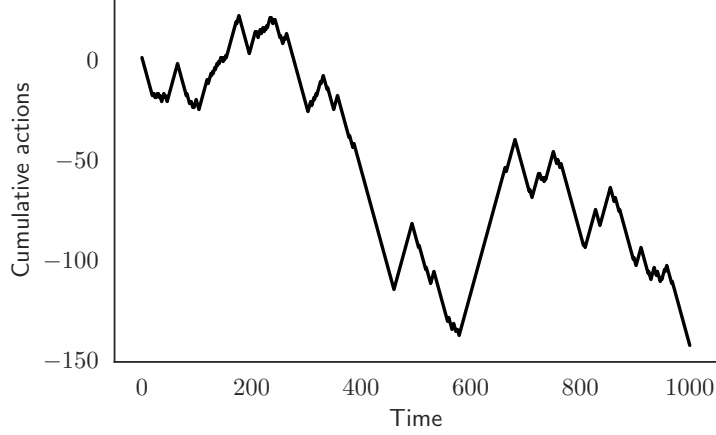


Fig. 4: Time series of interbank network centrality (left side). The decumulative distribution (DDF) of the in-degree (right side).

The dynamics of network centrality show how the morphology of the credit market evolves, going from periods in which the network is decentralized and made of many small components to periods in which more than 45% of banks are connected to a single hub. In addition, the topology of the emerging network is different from that of the random graph, where the in-degree distribution decays exponentially. Similar to real credit networks, in our system, some banks are found to have a disproportionately large number of incoming links. In contrast, others have very few (see Iori and Mantegna (2018), for a survey of the relevant literature). This result is shown in the right-hand side of Fig. 4 where we plot the decumulative distribution function of the in-degree. As the reader can observe, this distribution is in keeping with scale-free networks and displays a 'fat tail.'



y_t	b_0	b_1
Centrality	0.1830*** (599.06)	-0.0047*** (-11.8013)
Density	0.1042*** (317.08)	-0.0070*** (-16.01)
Diameter	10.14*** (1104.02)	0.22*** (17.10)
Components	1.48*** (656.45)	0.0095*** (3.04)
Avg nodes per components	39.50*** (940.20)	-0.21*** (-3.57)

*** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$

Fig. 5: Top Panel: Time series of η cumulative values over the simulation. Bottom Panel: Estimated results with the respective T-test in brackets for Eq.20. b_0 is the estimated mean value of y when $\eta = 1$ and b_1 the deviation from this mean value when $\eta = 0$. Data are obtained through 200 Monte Carlo simulations of the system.

To conclude the analysis of the interbank market architecture, we deal with the effect of the η parameter on the credit network topology. In the top panel of Fig. 5, we plot a single realization of the cumulative value of η over time. The figure shows how the reinforcement learning algorithm generates a time evolution in the choice of policy recommendations. Precisely, increasing (decreasing) values in the curve correspond to a signal that directs the system toward a high liquidity supply (low interest rate), i.e., $\eta = 1$ (i.e., $\eta = 0$).

The effect of the signal in shaping the topology of the interbank network is, instead, shown in the lower panel of Fig. 5, where we estimate a categorical regression model

$$y_t = b_0 + b_1(1 - \eta_t), \quad (20)$$

where b_0 is the estimated mean value assumed by the dependent variable y when $\eta = 1$ and $b_0 + b_1$ is the mean when $\eta = 0$. As shown in the bottom panel of Fig. 5, when the system selects low interest rates, the interbank network is less centralized, more sparse, and with a larger diameter. Moreover, the graph is fragmented into many scarcely-populated islands.

Having described the architecture of the interbank network, let us now examine its evolution over time. It is worth remembering that banks signal in the market their attractiveness μ according to the recommendation from the regulator, i.e., whether to compete more on low interest rates, $\eta = 0$, or on high liquidity supply, $\eta = 1$. While the regulator's signal is market-specific, liquidity supplies and interest rates (based on Eq.6) are bank-specific variables. This mechanism creates competition among financial institutions for credit allocation. The war in granting credit, modeled through the possibility of redefining lending agreements via Eq. 8 is shown in Fig. 6.

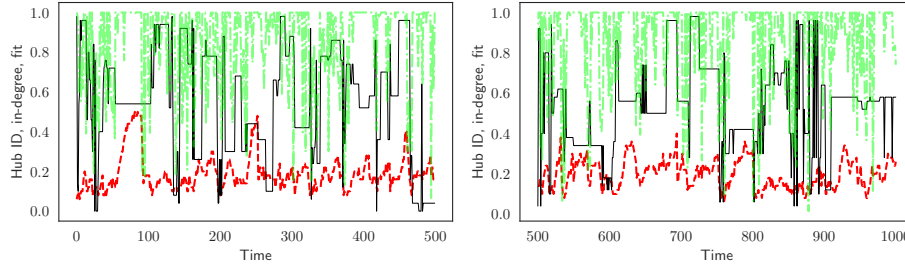


Fig. 6: Time series of the most connected lender (hub) evolution along the time T . The solid black line identifies the normalized hub id, the red dashed line her number of clients (incoming links), and the green dotted line the hub' fitness. Colors are available on the website version.

The black solid, red dashed, and green dotted lines represent the normalized id of the lender with the highest number of clients (i.e., the hub), her incoming links (i.e., number of clients), and her fitness, respectively. As the reader can appreciate, the simulation presents periods of hub stability and periods of alternation and competition between different hubs. When the hub stands out from

her competitors and signals a significantly higher fitness (i.e., the green dotted line approaches the unit), she can attract numerous clients, as shown by her high number of incoming links. However, the attractiveness of the hub may work against her. A large portfolio of customers increases the likelihood that some of them may fail. This either decreases the attractiveness of the hub herself or even causes her failure. **The reduction of the hub’s fitness due to one of her clients’ failure works in the following way. On the one hand, when the fitness uses a strategy based on a low interest rate, the client’s approach to the bankruptcy threshold increases the borrower’s financial fragility and probability of bankruptcy. Both these effects increase the lending interest rate, making the hub less attractive (see Eq.6). On the other hand, when μ moves towards a high liquidity supply, the borrower’s bad debt is absorbed by the lender’s net worth. The fall in the latter causes a parallel reduction in the hub liquidity, as shown by the balance sheet identity (see Eq. 1). Interestingly, reducing the hub net-worth could reduce liquidity higher than proportionally, given the Basel rules on maximum capital and leverage ratio. In any case, the drop in the agent’s fitness gradually reduces her number of clients and makes other lenders more attractive. These agents can replace the unsuccessful hub and so become, in turn, the most appealing lenders.**

Micro consequences of the reinforcement learning policy

In this subsection, we investigate how the dynamics of η affect the hub’s performance and other financial institutions.

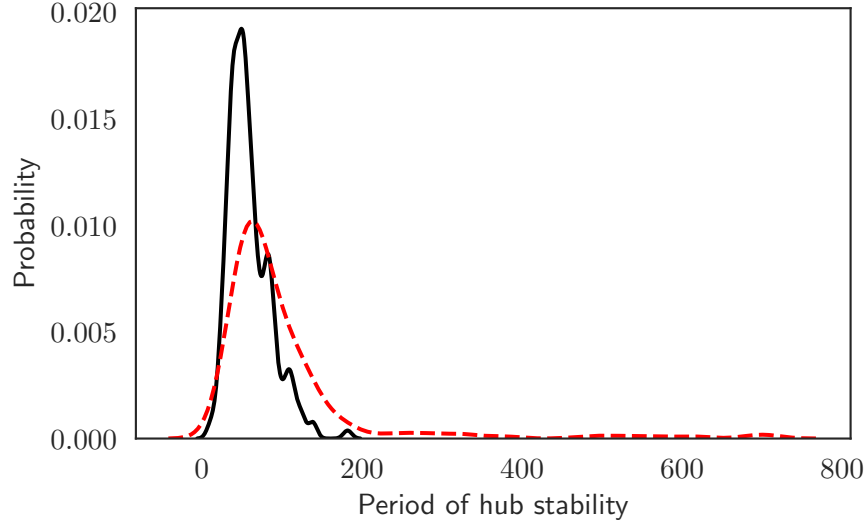


Fig. 7: Density distributions over 200 Monte Carlo simulations of the maximum period of hub stability in which the strategy does not change. The black solid and red dashed lines show $\eta = 0$ and $\eta = 1$, respectively.

In Fig. 7, we show how the choice between a low interest rate and a high liquidity supply strategy affects the hub's longevity. The figure shows the distribution, over 200 simulations, of the maximum period of hub stability in which the strategy does not change, respectively, for $\eta = 0$ (black) or $\eta = 1$ (red). The figure shows that the hub is generally more stable if the regulator recommends a high liquidity supply (red dashed line in Fig 7). Moreover, also at a micro level, we show that $\eta = 1$ seems to produce better individual performances. This result is shown in the top panel of Fig. 8, where we report the effect of the two possible values of η on some key individual variables.

y_t	b_0	b_1
Liquidity	2960.34*** (1062.34)	331.42*** (82.59)
Equity	888.96*** (1171.31)	-110.72*** (-135.79)
Leverage	0.01673*** (435.43)	0.000175*** (3.23)
Rationing	0.33*** (71.80)	0.28*** (38.31)
Bad debt	36.05*** (446.46)	2.19*** (19.49)
Failed banks	3.14*** (520.17)	0.35*** (41.00)

*** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$

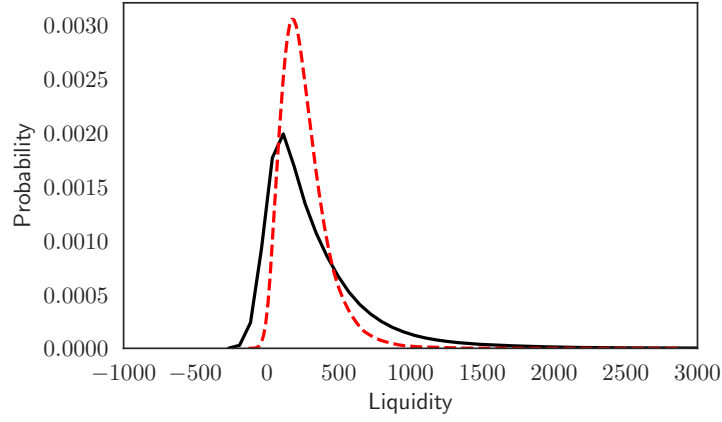


Fig. 8: Top Panel: Estimated results with the respective T-test in brackets for Eq.20. b_0 is the estimated mean value of y when $\eta = 1$ and b_1 the deviation from this mean value when $\eta = 0$. Data are obtained through 200 Monte Carlo simulations of the system. Bottom Panel: Density distributions of aggregated liquidity over times over 200 Monte Carlo simulations. The black solid and red dashed lines show $\eta = 0$ and $\eta = 1$, respectively.

Specifically, our results, estimated via the categorical regression model in Eq. 20, show that a signal that directs the system toward an abundant supply of liquidity (i.e., $\eta = 1$) produces better results in controlling leverage, rationing, bad debt, and bankruptcies. Moreover, according to the hypothesis that banks fail as net-worth falls below a minimum threshold, the equity is higher in the case of $\eta = 1$.

The result on the liquidity is, however, less intuitive. The system that competes on the interest rate level is significantly more liquid than the one adopting high liquidity, with an average liquidity value of 3291 in the case of $\eta = 0$ and 2960 in the opposite case. The reason for the apparent better performance on liquidity in the case of $\eta = 0$ lies in the competition among banks using interest rates. As clarified by Eq. 6, the financial institutions applying the lowest interest rates are the smallest ones. This implies that the biggest banks are less attractive to borrowers because they charge higher rates. Therefore, the system excludes these economic agents from trading while encouraging small institutions to provide liquidity. This mechanism of selection has a twofold effect. On the one hand, it generates a substantial heterogeneity between lenders and borrowers. Creditors, much smaller than debtors, are overwhelmed in the event of their clients' bankruptcy. On the other hand, the exclusion from the exchanges of the largest institutions leaves a consistent level of unallocated liquidity in the system. The first effect, i.e., agents' heterogeneity, determines the worst performances under $\eta = 0$, while the second effect, i.e., exclusion, determines the highest level of unallocated liquidity in the system. In contrast, a signal that directs the system towards an abundant liquidity supply produces a more homogeneous distribution among banks' liquidity, as shown in the bottom panel of Fig. 8. This homogeneity between economic agents generates a uniform risk exposure among counterparties, favoring the system's resiliency in front of shocks. **This result, although not unanimously shared (see Haldane and May (2011)), is in line with other studies showing that agents' heterogeneity is a leading force in generating propagation of systematic failure (see, for instance, Caccioli et al. (2012), Berardi and Tedeschi (2017), Iori et al. (2006), Lenzu and Tedeschi (2012) and Tedeschi et al. (2012)).**

Systemic impact of the network

To conclude the section, we combine the results on network topology and individual performance as a function of η to capture the interbank architecture's overall

Indep. Variable	Dep. Variable		
	Rationing	Failed banks	Leverage
Net centrality	-0.25*** (-6.04)	-2.09*** (-41.82)	-0.016*** (-55.69)
Density	-1.32*** (-52.64)	-9.14*** (-254.08)	-0.051*** (-235.33)
Diameter	0.011*** (8.95)	0.032*** (21.69)	0.0002*** (27.02)
Components	0.029*** (5.28)	0.020*** (3.09)	0.0004*** (10.57)
Avg nodes per comp	-0.0011*** (-4.18)	-0.0022*** (-6.97)	-0.00002*** (-13.19)

*** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$

Table 1: Regression results between indicators of the interbank stability and network measures. T-stats for each coefficient are provided in parentheses. Data are obtained through 200 Monte Carlo simulations of the system.

effect on systemic stability. To this end, in Tab. 1, we report the results of a linear regression estimated through ordinary least squares where the independent variables are some measures of the interbank network topology and dependent variables are some indicators of systemic market stability. In line with what has been observed so far, when the network tends to be centralized, i.e., denser towards the hub and with a smaller diameter, the risk of contagion decreases, i.e., bankruptcies, rationing, and leverage is reduced. This architecture corresponds to a graph composed of a few highly populated components. It is worth noting that this topology emerges when the interbank system is oriented towards an abundant supply of liquidity, which generates a certain homogeneity among agents able to compensate for the imbalance between lenders and borrowers present in the case of $\eta = 0$. In this respect, clarification is essential: $\eta = 1$ is not the absolute best signal. This is the best strategy given the individual and aggregate conditions of the system at the time of the choice. The algorithm is designed to identify one recommendation as optimal based on the underlying environmental conditions. The robustness of this observation is shown in Sec.3.3 and Sec. 3.4. In the former, we show that the system governed by a regulator that directs the choice via the implemented reinforcement learning algorithm outperforms a system based on a random selection between the two signals. In

the latter, we demonstrate the better performances of the reinforcement learning rather than keeping constant the two values of η **or modeling an η evolving with decentralized dynamics..**

3.3 The reinforcement learning based recommendation for taming systemic risk

In this subsection, we study the effect on the interbank systemic stability of the policy recommendation obtained through the reinforcement learning mechanism solved by the PPO algorithm.

Specifically, we answer the following question: how would the interbank system perform in terms of aggregate resiliency when the regulator directs financial institutions to choose the optimal strategy between competing on the low interest rate, $\eta = 0$, or on high liquidity, $\eta = 1$? Again we compare the effects of the learned strategy on the market stability with those of a random strategy. **Finally, the last part of this session is devoted to understanding the effects of herding on systemic stability. Specifically, we study the market performance as the percentage of banks that follow the policy recommendation changes.**

A common finding in several theoretical and empirical works is that the interbank market works better when credit flows efficiently through the system, thus ensuring it against liquidity shocks (see, for instance, Allen and Gale (2000); Carlin et al. (2007); Freixas et al. (2000)).

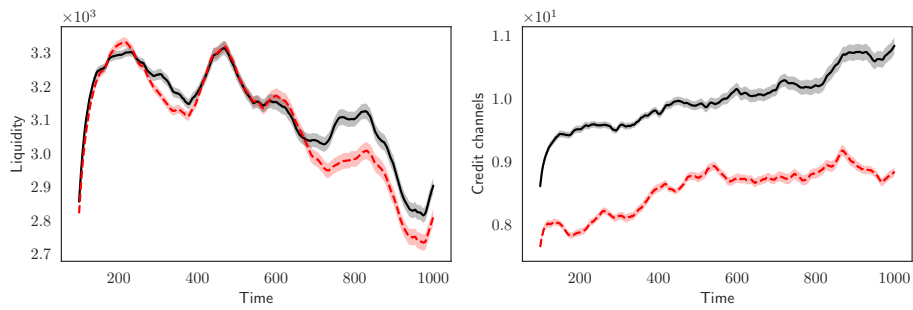


Fig. 9: Liquidity of the system (left panel) and the number of credit channels (right panel). Black solid and red dashed lines refer to the best-performing reinforcement learning optimal and random strategies, respectively. The curves reproduce the mean and the standard deviation over **200** simulations of the system and a rolling window of 100 timesteps.

Starting from this consideration and recalling the severity of liquidity crises, we show in Fig. 9 (left side) the effectiveness of the implemented reinforcement learning strategy in spreading liquidity through the system. In the figure, once the best performing learned strategy is selected, as shown on the right-hand side of Fig. 1, the aggregated average liquidity of **200** simulations over a rolling window of 100 timesteps is shown through time. Although the learned strategy strongly competes with the random one in some periods, its supremacy becomes evident from step 700 onwards. In addition, the average liquidity, over all periods and simulations, of the learned strategy is statistically higher than the one obtained with the random strategy (i.e., 3129.98 (std. 1.5128) vs. 3091.51 (std. 4.4258), respectively).

A possible explanation for this phenomenon can be seen in the right-hand side of Fig. 9, where we plot the active credit links in the two frameworks¹². As the reader can appreciate, the number of activated credit channels is higher when the system follows the learned strategy with respect to the case of random strategy, and this guarantees a higher circulation of liquidity in the system. In detail, the average number of credit channels, over time and simulations, in the first scenario is 9.9823 (std. 0.4321), while in the second case is 8.5464 (std. 0.3596). On the whole, this result reveals the ability of the reinforcement learning optimal policy to design an interbank network architecture promoting an efficient credit allocation and, therefore, reducing liquidity shortage phenomena. As a consequence, the emerging topology of the credit network effectively controls rationing and avoids failures due to credit crunch phenomena, as shown in Fig. 10, left and right panel, respectively.

¹² By the terms credit channels and credit links we refer to the linkages through which the liquidity needed by borrowers due to the deposit shock flows. These are, therefore, the credit lines used in the active credit network.

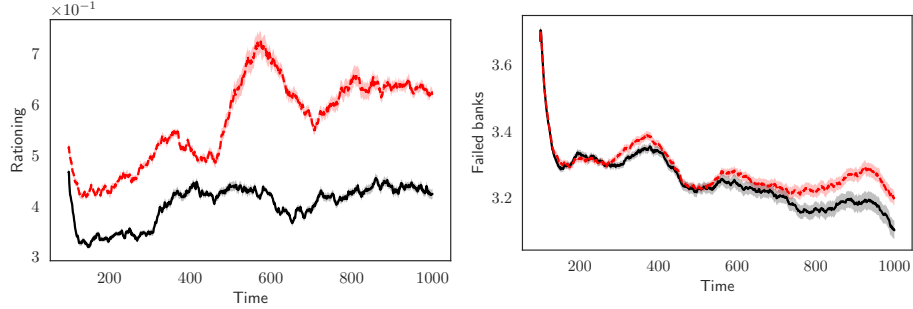


Fig. 10: Rationing of the system (left panel) and the number of failed banks (right panel). Black solid and red dashed lines refer to the best-performing reinforcement learning optimal and random strategies, respectively. The curves reproduce the mean and the standard deviation over **200** simulations of the system and a rolling window of 100 timesteps.

The robustness of the two latter results is confirmed by the average values of these variables over all timesteps and simulations. Specifically, the mean and standard deviation of the rationing in the case of the learned strategy (resp. random strategy) are 0.4024 and 0.0375 (resp. 0.5671 and 0.08465), while the mean and standard deviation of the number of failed banks in the case of the learned policy (resp. random policy) are 3.2101 and 0.0410 (resp. 3.2931 and 0.0423).

It is essential to note the ability of the reinforcement learning mechanism to generate an interbank network whose architecture is resilient in the face of financial attacks. This characteristic provides, on the one hand, an additional monetary policy tool that can be implemented in times of economic adversity and, on the other hand, enriches the vast literature that emphasizes the importance of credit network architecture in dealing with systemic shocks (see Grilli et al. (2017), for a survey of the relevant literature).

We conclude this section by analyzing the effect of the reinforcement learning optimal policy on the market’s financial (in)stability. The approach followed here in explaining the materialization of financial frictions is very close in spirit to the Minskyan financial instability hypothesis and therefore uses banks’ leverage as the leading indicator (see Minsky (1964)). In our stylized market, leverage and systemic instability are connected through a specific structure. Given our naive banks’ balance sheet (see Eq.1), leverage is defined as assets on equity. Moreover, credit costs (i.e., interest rates) are strongly positively affected by the leverage (see Eq.6). When a lender grants a loan to a bank with a low

probability of surviving (i.e., an over-leveraged borrower), she charges a higher interest rate via the financial accelerator. This, in turn, exacerbates the financial condition of the borrower herself pushing her towards a bankruptcy state. If one or more borrowers cannot pay back their loans, even the lenders' equity is affected by bad debts. Therefore, lenders decrease their credit supply and increase the borrowers' rationing. In this way, the profit margin of borrowers decreases, and a new round of failures may occur. The leverage dynamics when the system follows the reinforcement learning recommended policy and in the random case are shown on the left-hand side of Fig. 11. The figure highlights two important features. First, the recommended learned policy keeps the leverage below the values obtained with the random policy. Specifically, the average leverage in the first scenario, over time and simulations, is 1.59 (std. 0.042), while in the second case is 1.69 (std. 0.031). Second, the leverage fluctuates over time, thus recalling the different phases of lending suggested by Minsky. There are periods when financial institutions grant more loans without considering the overall financial fragility. However, banks can underestimate their credit risk, making the system more vulnerable when default materializes. This ambiguous effect of the leverage, first positive and then negative, on interbank stability, is clearly shown in the right-hand side of Fig. 11, where the correlation wave between bankruptcies and agents' leverage first decreases from lag $\tau = -21$ up to $\tau = -11$, then increases from $\tau = -8$ up to $\tau = 9$, and finally, returns to decrease from $\tau = 15$.

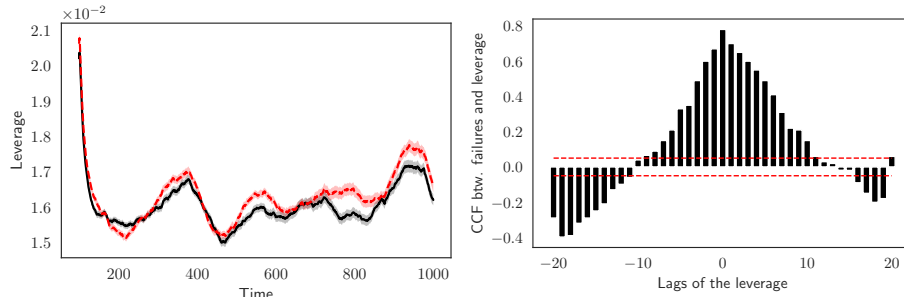


Fig. 11: Left side: Leverage of the system. Black solid and red dashed lines refer to the best performing reinforcement learning optimal strategy and to the random strategy, respectively. The curves reproduce the mean and the standard deviation over **200** simulations of the system and a rolling window of 100 time steps. Right side: Average correlation between number of bankruptcy and lagged leverage, at a 1% confidence level.

Should they herd or should they do not? An exercise on the herding effect

Let us introduce an additional element of heterogeneity concerning the signal itself. Whereas in the previous experiment, all the banks followed the signal on the optimal strategy, here we modify this assumption. We simulate a system where different percentages of banks follow the signal while the others randomly go to different possible strategies. This experiment allows us, on the one hand, to introduce an additional element of differentiation and, on the other, to understand what is the minimum threshold of herding required by the system with the RL-generated signal to be more resilient than the one with a random strategy. To this end, we fix a percentage κ of banks that follow the reinforcement learning strategy, while the remaining $N(1 - \kappa)$ banks randomly sample the strategy in the set $\{0, 0.5, 1\}$. We train the reinforcement learning algorithm following the same procedure of the previous subsections, letting the mixture parameter vary on a discrete range of values. Every time we change the value of κ , a new algorithm is trained. Several system simulations are performed to evaluate the effect of such heterogeneity in the strategy followed by the banks.

Before studying the impact on the interbank systemic stability of the different percentages of financial institutions applying the policy recommendation obtained through the reinforcement learning mechanism and comparing it with the random strategy, an important consideration is necessary. Fig. 12 shows that as herding rates vary, the PPO algorithm selects different categories of strategies. For example, when only 10% of the banks follow the optimal signal, the most common strategy steers the banks towards a low interest rate (see solid black line). However, in this scenario, even if with low probability, a mixed strategy (i.e., $\eta = 0.5$) or a high liquidity supply strategy (i.e., $\eta = 1.0$) can emerge (see the brown and yellow lines, respectively). This competition among different optimally selected strategies varies as the herding varies. However, in the case of total herding, i.e., when all banks follow the policy recommendation, the system stabilizes, with equiprobability, on the two extremes, i.e., $\eta = 0$ and $\eta = 1$. When $\kappa = 0.1$ (i.e., herding of 10%), the probability of a low interest rate signal is

93%, while the probability of a mixed strategy is 2.51%. Instead, the probability of a signal pointing to a high supply of liquidity is 4.38%. Moving towards herding of 50%, the selected strategies vary. Specifically with $\kappa = 0.5$ the probability of $\eta = 0.0$ is 57%, that of $\eta = 0.5$ is 40% and finally $\eta = 1$ is 3%. Tab.2 shows the portion of the chosen optimal strategy for each herding percentage.

	Herding Percentage									
	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0
Strategy $\eta = 0.0$	0.93	0.68	0.22	0.18	0.57	0.31	0.66	0.38	0.24	0.5
$\eta = 0.5$	0.03	0.21	0.66	0.66	0.40	0.69	0.32	0.62	0.76	0.0
$\eta = 1.0$	0.04	0.11	0.12	0.16	0.03	0.00	0.01	0.00	0.00	0.5

Table 2: Percentage of the chosen optimal strategy ($\eta = 0$; $\eta = 0.5$ and $\eta = 1$) by varying the herding parameter κ from 1% to 100%.

Let us now analyze how the interbank system performs in terms of aggregate resiliency when the regulator convinces different percentages of banks to follow the optimal signal. As in the first part of this subsection, the results obtained with the optimized strategy are compared with those obtained from a random choice of strategy. As in the baseline case, the dynamic random scenario picks the value of η according to the probabilities shown in Tab 2.

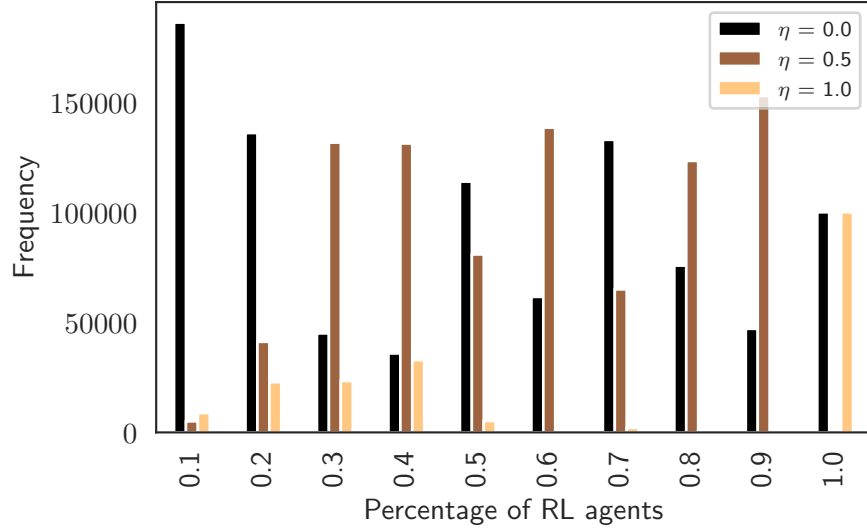


Fig. 12: Discrete distributions of the η values selected by PPO over 200 Monte Carlo simulations of the system.

We report the aggregated results at the macroeconomic level for some of the critical systemic variables. In each panel of the Figures 13 and 14, we show the variation of the aggregated measure obtained averaging through 200 simulations and through the timesteps of the simulations (1000). The aggregated measure is displayed on the y -axis, while the herding parameter κ varies on the x -axis.

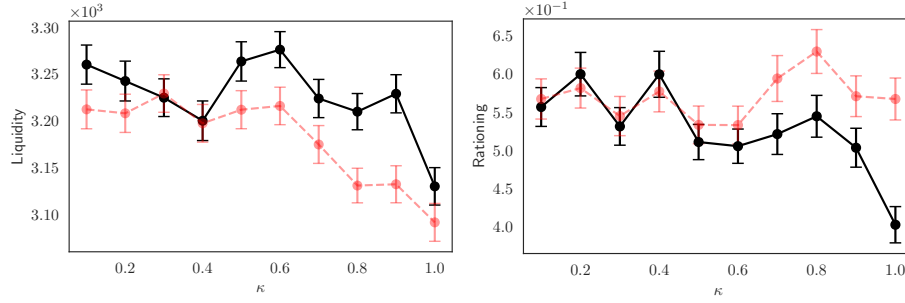


Fig. 13: Average liquidity (left panel) and rationing of the system (right panel) as a function of the herding percentage, κ . Black solid and red dashed lines refer to the best-performing reinforcement learning optimal and random strategies, respectively. The curves reproduce the mean and the standard deviation over 200 simulations of the system.

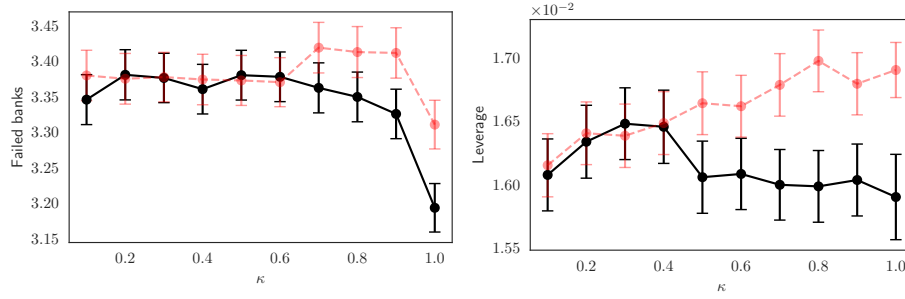


Fig. 14: Average number of failed banks (left panel) and their leverage (right panel) as a function of the herding percentage, κ . Black solid and red dashed lines refer to the best-performing reinforcement learning optimal and random strategies, respectively. The curves reproduce the mean and the standard deviation over 200 simulations.

When the regulator cannot convince a sufficient percentage of banks to follow the policy recommendation, the system generated with the optimal signal obtained via the reinforcement learning algorithm (black solid) does not significantly differ from that generated with the random signal (red dashed lines). This holds for all the considered variables, such as the liquidity and rationing of the system (see Fig. 13) and the number of failures and leverage of financial institutions (see Fig. 14). Instead, when the regulator can convince a share of banks equal to/greater than 60%, higher systemic stability is observed in the model using the optimized signal than in the ran-

dom one. In fact, above this percentage, the optimized system, on the one hand, generate higher liquidity and lower rationing, on the other hand, fewer bankruptcies and less leverage for financial institutions.

3.4 A competition among different behavioural strategies

This session compares the aggregate performances obtained by the reinforcement learning strategy with other possible tactics. Whereas in the previous sessions, the comparison is made only by considering a strategy that randomly selected the value of η , which, however, follows the same distribution as the η obtained with the reinforcement learning algorithm, here we develop two other possible dynamics for this parameter. In the first experiment, we focus on the comparison with exogenous and fixed η . In the second study, we implement the parameter in a dynamic and decentralized way so that each bank has her plan of action.

In all these experiments, we run our model 200 times for different values of the initial seed generating the pseudo-random numbers over a time span of $T = 1000$ periods. Moreover, all the agents' initialization parameters, except for the variations studied here, coincide with those presented in Sec 3 with a percentage of herding $\kappa = 1$.

Let us start with the first scenario where the η obtained with the reinforcement learning algorithm is compared with two fixed and constant η values, i.e. $\eta = 0$ and $\eta = 1$. This experiment allows us to verify the resilience of our simulated system with respect to the one obtained by implementing a fixed mechanism of parameter choice as in Berardi and Tedeschi (2017).

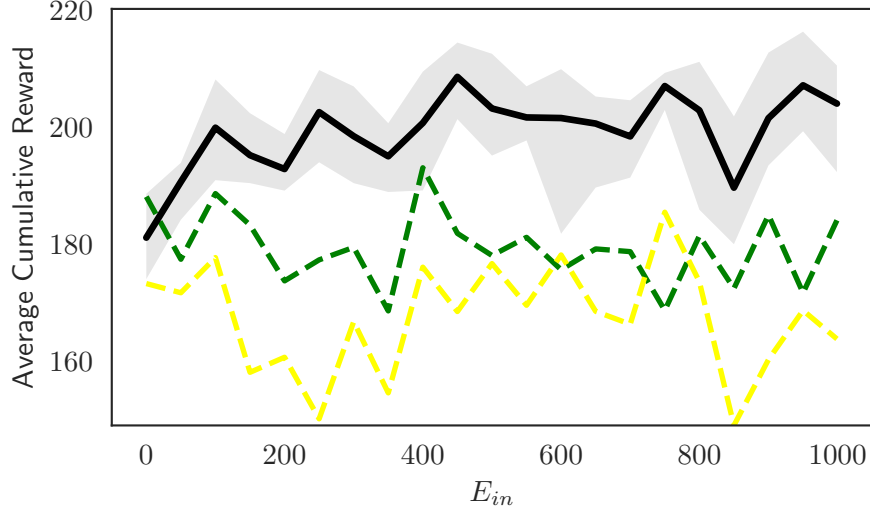


Fig. 15: Average cumulative fitness of the system as a function of the number of training episodes for the trained PPO instances with the corresponding confidence intervals (in solid black), and the fixed strategies $\eta = 0$ (in dashed green) and $\eta = 1$ (in dashed yellow).

In Fig. 15, we show that the reinforcement learning algorithm outperforms the two fixed signals in terms of aggregated fitness of the system. This result indicates that a dynamic selection of the η by looking at the available information allows for more attractive banks than maintaining a fixed η . Finally, in Tab. 3, we show the results of the economic performance of the three different signals.

It is worth noting that when the regulator adopts an η evolving through reinforcement learning, the system is more liquid and absorbs shocks better than in the other two cases (as shown by the lower leverage and lower rationing and the number of failures associated with endogenous η). Furthermore, in line with what has been described in Sec. 3.2, the $\eta = 1$ strategy consistently outperforms the $\eta = 0$ one, with the only known exception for liquidity.

Let us now investigate what happens to the system if strategies are decentralized and left in the hands of individual financial institutions. This experiment, although still very preliminary, gives interesting insights into the comparison between a centralized signal and a decentralized one where agents compete with each other without following any public recommendation from the regulator.

Ave. value	$\eta = 0.0$	$\eta = 1.0$	$\eta = RL$
Liquidity	302405.63 (173.84)	280389.40 (146.69)	312998.02 (151.28)
Leverage	1.75 (0.054)	1.72 (0.064)	1.59 (0.042)
Rationing	57.40 (2.90)	48.51 (2.27)	40.24 (3.75)
Failed banks	332.49 (4.86)	325.21 (4.81)	321.01 (4.10)
Credit channels	1032.21 (18.66)	1260.51 (120.67)	998.23 (43.21)

Table 3: Average values with standard deviations in parentheses, over times and all 200 Montecarlo simulations, of the aggregated economic variables obtained for $\eta = 0$, $\eta = 1$ and η evolving via the reinforcement learning algorithm, i.e., $\eta = RL$.

Let us begin by describing the dynamics of the η_t^i parameter, which now becomes dynamic and agent-specific. Specifically, denoting η_t^i as the weight that bank i gives to the liquidity or the interest rate in the fitness function, it becomes a function of the recent performance of the agent in terms of attractiveness. Namely, if $\mu_t^i - \mu_{t-1}^i \geq 0$, the agent i intensifies the strategy she is already pursuing, then

$$\eta_{t+1}^i = \begin{cases} \eta_t^i + a & \text{if } \eta_t^i \geq 0.5 \\ \eta_t^i - a & \text{if } \eta_t^i < 0.5 \end{cases} \quad (21)$$

On the other hand, if $\mu_t^i - \mu_{t-1}^i < 0$, the bank i weakens the strategy she is pursuing, intensifying the opposite one

$$\eta_{t+1}^i = \begin{cases} \eta_t^i - a & \text{if } \eta_t^i \geq 0.5 \\ \eta_t^i + a & \text{if } \eta_t^i < 0.5 \end{cases}, \quad (22)$$

where a is a scalar parameter defining the movement step towards liquidity or interest rate.

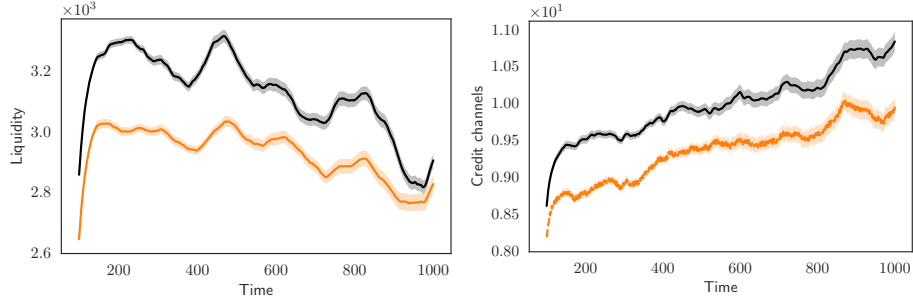


Fig. 16: Liquidity of the system (left panel) and the number of credit channels (right panel). Black solid and red dashed lines refer to the best performing reinforcement learning optimal strategy and the decentralized strategy, respectively. The curves reproduce the mean and the standard deviation over 200 simulations of the system and a rolling window of 100 timesteps.

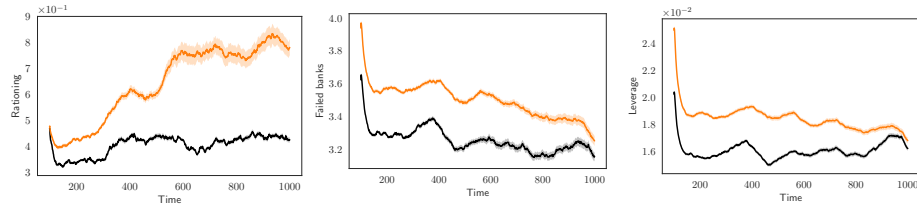


Fig. 17: Rationing of the system (left panel), number of failed banks (middle panel), and Leverage of the system (right panel). Black solid and red dashed lines refer to the best performing reinforcement learning optimal strategy and the decentralized strategy, respectively. The curves reproduce the mean and the standard deviation over 200 simulations of the system and a rolling window of 100 timesteps.

As the reader can easily grasp from Figs. 16 and 17, where the well-known systemic dynamics presented in Sec. 3.3 are reproduced, once again the RL-based strategy (black solid line) generates more desirable systemic patterns than the new decentralised strategy (red dashed line). Although these latest results are only sketchy, they further highlight the validity of the reinforcement learning mechanism for the stability of the interbank system.

4 Concluding remarks

This work shows the effects of a policy recommendation obtained through a reinforcement learning mechanism in an artificial interbank market. Specifically,

we assume that the financial institutions receive a signal from the regulator regarding the best strategy to adopt for the creation of their lending agreements. Depending on the underlying economic conditions, the signal directs the system towards providing a high liquidity supply or a low interest rate. Using a reinforcement learning approach to provide this public signal has proven effective since the method exploits the available information and redirects the system towards an efficient flow of liquidity compared to other different static and dynamic **behavioral tactics**. Moreover, through the use of the SHAP framework, which dissects the contribution of each piece of information to the recommended policy, we have been able to interpret what is the primary input that drives the choice of the policy. We have acknowledged that the occurrence of one circumstance (liquidity vs. interest rate) generates significant consequences affecting the agents' performances and the topology and resiliency of the interbank network. Specifically, when the signal directs the system toward an abundant liquidity provision, the interbank network, composed of a few populated communities, is more centralized and dense towards hub banks than in the low interest rate scenario. This network architecture is accompanied by better individual performances and higher system resilience in the face of exogenous shocks. Our results have shown that the better general conditions underlying this signal are due to the homogeneity between lenders and borrowers, which generates a uniform risk exposure among counterparties able to favor the system's resiliency.

Leaving aside the results of the comparison between the two signals, we have analyzed the general effect of the policy recommendation implemented via the reinforcement learning procedure in the second part of the paper. Our results have shown how systemic risk is mitigated by such a tool and how this outperforms other alternatives **behavioral strategies**.

Acknowledgments

This research was supported by grants from the Spanish Ministerio de Ciencia, Innovacion y Universidades (grant RTI2018-096927-B-100)

Bibliography

Acharya, V. V. and Yorulmazer, T. (2008a). Cash-in-the-market pricing and optimal resolution of bank failures, *The Review of Financial Studies* **21**(6): 2705–2742.

- Acharya, V. V. and Yorulmazer, T. (2008b). Information contagion and bank herding, *Journal of money, credit and Banking* **40**(1): 215–231.
- Aleem, I. (1990). Imperfect information, screening, and the costs of informal lending: A study of a rural credit market in pakistan, *The World Bank Economic Review* **4**(3): 329–349.
- URL:** <http://www.jstor.org/stable/3989880>
- Allen, F. and Gale, D. (2000). Financial contagion, *Journal of political economy* **108**(1): 1–33.
- Altavilla, C., Lemke, W., Linzert, T., Tapking, J. and von Landesberger, J. (2021). Assessing the efficacy, efficiency and potential side effects of the ecb’s monetary policy instruments since 2014.
- Andrychowicz, M., Raichuk, A., Stańczyk, P., Orsini, M., Girgin, S., Marinier, R., Hussenot, L., Geist, M., Pietquin, O., Michalski, M. et al. (2020). What matters in on-policy reinforcement learning? a large-scale empirical study, *arXiv preprint arXiv:2006.05990*.
- Angelini, P., Maresca, G. and Russo, D. (1996). Systemic risk in the netting system, *Journal of Banking & Finance* **20**(5): 853–868.
- Barabási, A.-L. and Albert, R. (1999). Emergence of scaling in random networks, *science* **286**(5439): 509–512.
- Bartelsman, E., Scarpetta, S. and Schivardi, F. (2005). Comparative analysis of firm demographics and survival: evidence from micro-level sources in oecd countries, *Industrial and corporate change* **14**(3): 365–391.
- Battiston, S., Gatti, D. D., Gallegati, M., Greenwald, B. and Stiglitz, J. E. (2012a). Default cascades: When does risk diversification increase stability?, *Journal of Financial Stability* **8**(3): 138–149.
- Battiston, S., Gatti, D. D., Gallegati, M., Greenwald, B. and Stiglitz, J. E. (2012b). Liaisons dangereuses: Increasing connectivity, risk sharing, and systemic risk, *Journal of economic dynamics and control* **36**(8): 1121–1141.
- Bellman, R. (1956). Dynamic programming and lagrange multipliers, *Proceedings of the National Academy of Sciences of the United States of America* **42**(10): 767.
- Berardi, S. and Tedeschi, G. (2017). From banks’ strategies to financial (in) stability, *International Review of Economics & Finance* **47**: 255–272.
- Bernanke, B. S., Gertler, M. and Gilchrist, S. (1999). The financial accelerator in a quantitative business cycle framework, *Handbook of macroeconomics* **1**: 1341–1393.

- Bester, H. (1985). Screening vs. rationing in credit markets with imperfect information, *The American Economic Review* **75**(4): 850–855.
URL: <http://www.jstor.org/stable/1821362>
- Bindseil, U. (2018). *Financial stability implications of a prolonged period of low interest rates*, BIS.
- Boissay, F., Collard, F., Gali, J. and Manea, C. (2021). Monetary policy and endogenous financial crises.
- Boissay, F., Collard, F. and Smets, F. (2016). Booms and banking crises, *Journal of Political Economy* **124**(2): 489–538.
- Brunnermeier, M. K., Eisenbach, T. M. and Sannikov, Y. (2012). Macroeconomics with financial frictions: A survey.
- Buşoniu, L., Babuška, R. and De Schutter, B. (2010). Multi-agent reinforcement learning: An overview, *Innovations in multi-agent systems and applications-1* pp. 183–221.
- Caccioli, F., Catanach, T. A. and Farmer, J. D. (2012). Heterogeneity, correlations and financial contagion, *Advances in Complex Systems* **15**(supp02): 1250058.
- Calice, G., Sala, C. and Tantari, D. (2020). Contingent convertible bonds in financial networks, *arXiv preprint arXiv:2009.00062*.
- Capponi, A., Sun, X. and Yao, D. D. (2020). A dynamic network model of interbank lending—systemic risk and liquidity provisioning, *Mathematics of Operations Research* **45**(3): 1127–1152.
- Carlin, B. I., Lobo, M. S. and Viswanathan, S. (2007). Episodic liquidity crises: Cooperative and predatory trading, *The Journal of Finance* **62**(5): 2235–2274.
- Charpentier, A., Elie, R. and Remlinger, C. (2021). Reinforcement learning in economics and finance, *Computational Economics* pp. 1–38.
- Cincotti, S., Raberto, M. and Teglioni, A. (2012). Macroprudential policies in an agent-based artificial economy, *Revue de l'OFCE* (5): 205–234.
- Dasgupta, A. (2004). Financial contagion through capital connections: A model of the origin and spread of bank panics, *Journal of the European Economic Association* **2**(6): 1049–1084.
- De Grauwe, P. (2011). The banking crisis: causes, consequences and remedies, *Systemic Implications of Transatlantic Regulatory Cooperation and Competition*, World Scientific, pp. 23–46.
- Dell’Ariccia, G. and Marquez, R. (2004). Information and bank credit allocation, *Journal of financial Economics* **72**(1): 185–214.

- Du, J., Jin, M., Kolm, P. N., Ritter, G., Wang, Y. and Zhang, B. (2020). Deep reinforcement learning for option replication and hedging, *The Journal of Financial Data Science* **2**(4): 44–57.
- Freixas, X., Parigi, B. M. and Rochet, J.-C. (2000). Systemic risk, interbank relations, and liquidity provision by the central bank, *Journal of money, credit and banking* pp. 611–638.
- Galí, J. (2015). *Monetary policy, inflation, and the business cycle: an introduction to the new Keynesian framework and its applications*, Princeton University Press.
- Georg, C.-P. (2013). The effect of the interbank network structure on contagion and common shocks, *Journal of Banking & Finance* **37**(7): 2216–2228.
- Gertler, M., Kiyotaki, N. and Prestipino, A. (2020). A macroeconomic model with financial panics, *The Review of Economic Studies* **87**(1): 240–288.
- Giannone, D., Lenza, M., Pill, H. and Reichlin, L. (2011). *Non-standard monetary policy measures and monetary developments*, Macroeconomic Policy Making, Cambridge University Press, p. 195–221.
- Giri, F., Riccetti, L., Russo, A. and Gallegati, M. (2019). Monetary policy and large crises in a financial accelerator agent-based model, *Journal of Economic Behavior & Organization* **157**: 42–58.
- Goldberg, J. E., Klee, E., Prescott, E. S. and Wood, P. R. (2020). Monetary policy strategies and tools: Financial stability considerations.
- Greenwald, B. C. and Stiglitz, J. E. (1993). Financial market imperfections and business cycles, *The Quarterly Journal of Economics* **108**(1): 77–114.
- Grilli, R., Iori, G., Stamboglis, N. and Tedeschi, G. (2017). A networked economy: A survey on the effect of interaction in credit markets, *Introduction to agent-based economics*, Elsevier, pp. 229–252.
- Grilli, R., Tedeschi, G. and Gallegati, M. (2014). Network approach for detecting macroeconomic instability, *2014 Tenth International Conference on Signal-Image Technology and Internet-Based Systems*, IEEE, pp. 440–446.
- Haldane, A. G. and May, R. M. (2011). Systemic risk in banking ecosystems, *Nature* **469**(7330): 351–355.
- Hoff, K. and Stiglitz, J. E. (1990). Introduction: Imperfect information and rural credit markets: Puzzles and policy perspectives, *The world bank economic review* **4**(3): 235–250.
- Ioffe, S. and Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift, *arXiv preprint arXiv:1502.03167*.

- Iori, G., Jafarey, S. and Padilla, F. G. (2006). Systemic risk on the interbank market, *Journal of Economic Behavior & Organization* **61**(4): 525–542.
- Iori, G. and Mantegna, R. N. (2018). Empirical analyses of networks in finance, *Handbook of Computational Economics*, Vol. 4, Elsevier, pp. 637–685.
- Jiang, Z., Xu, D. and Liang, J. (2017). A deep reinforcement learning framework for the financial portfolio management problem, *arXiv preprint arXiv:1706.10059*.
- Jiménez, G., Ongena, S., Peydró, J.-L. and Saurina, J. (2014). Hazardous times for monetary policy: What do twenty-three million bank loans say about the effects of monetary policy on credit risk-taking?, *Econometrica* **82**(2): 463–505.
- Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimization, *arXiv preprint arXiv:1412.6980*.
- Lenzu, S. and Tedeschi, G. (2012). Systemic risk on different interbank network topologies, *Physica A: Statistical Mechanics and its Applications* **391**(18): 4331–4341.
- Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D. and Wierstra, D. (2015). Continuous control with deep reinforcement learning, *arXiv preprint arXiv:1509.02971*.
- Lin, S. and Beling, P. A. (2020). An end-to-end optimal trade execution framework based on proximal policy optimization., *IJCAI*, pp. 4548–4554.
- Liu, A., Mo, C. Y. J., Paddrik, M. E. and Yang, S. Y. (2018). An agent-based approach to interbank market lending decisions and risk implications, *Information* **9**(6): 132.
- Lozano, F., Lozano, J. and García Molina, M. (2007). An artificial economy based on reinforcement learning and agent based modeling, *Documentos de Trabajo, Facultad de Economía, Universidad del Rosario* (18).
- Lundberg, S. M. and Lee, S.-I. (2017). A unified approach to interpreting model predictions, *Proceedings of the 31st international conference on neural information processing systems*, pp. 4768–4777.
- Marbach, P. and Tsitsiklis, J. N. (2001). Simulation-based optimization of markov reward processes, *IEEE Transactions on Automatic Control* **46**(2): 191–209.
- Maudos, J. and De Guevara, J. F. (2004). Factors explaining the interest margin in the banking sectors of the european union, *Journal of Banking & Finance* **28**(9): 2259–2281.
- Minsky, H. P. (1964). Longer waves in financial relations: financial factors in the more severe depressions, *The American Economic Review* **54**(3): 324–335.

- Mnih, V., Badia, A. P., Mirza, M., Graves, A., Lillicrap, T., Harley, T., Silver, D. and Kavukcuoglu, K. (2016). Asynchronous methods for deep reinforcement learning, *International conference on machine learning*, PMLR, pp. 1928–1937.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G. et al. (2015). Human-level control through deep reinforcement learning, *nature* **518**(7540): 529–533.
- Mosavi, A., Faghan, Y., Ghamisi, P., Duan, P., Ardabili, S. F., Salwana, E. and Band, S. S. (2020). Comprehensive review of deep reinforcement learning methods and applications in economics, *Mathematics* **8**(10): 1640.
- Osoba, O. A., Vardavas, R., Grana, J., Zutshi, R. and Jaycocks, A. (2020). Policy-focused agent-based modeling using rl behavioral models, *arXiv preprint arXiv:2006.05048*.
- Ricchetti, L., Russo, A. and Gallegati, M. (2018). Financial regulation and endogenous macroeconomic crises, *Macroeconomic Dynamics* **22**(4): 896–930.
- Rochet, J.-C. and Tirole, J. (1996). Interbank lending and systemic risk, *Why Are there So Many Banking Crises?* p. 140.
- Schulman, J., Levine, S., Abbeel, P., Jordan, M. and Moritz, P. (2015). Trust region policy optimization, *International conference on machine learning*, PMLR, pp. 1889–1897.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A. and Klimov, O. (2017). Proximal policy optimization algorithms, *arXiv preprint arXiv:1707.06347*.
- Shapley, L. S. (2016). *17. A value for n-person games*, Princeton University Press.
- Sutton, R. S., McAllester, D. A., Singh, S. P. and Mansour, Y. (2000). Policy gradient methods for reinforcement learning with function approximation, *Advances in neural information processing systems*, pp. 1057–1063.
- Svensson, L. E. (2017). Cost-benefit analysis of leaning against the wind, *Journal of Monetary Economics* **90**: 193–213.
- Taylor, J. B. (2011). Macroeconomic lessons from the great deviation, *NBER macroeconomics annual* **25**(1): 387–395.
- Tedeschi, G., Mazlounian, A., Gallegati, M. and Helbing, D. (2012). Bankruptcy cascades in interbank markets, *PloS one* **7**(12): e52749.
- Tedeschi, G., Vidal-Tomás, D., Delli-Gatti, D. and Gallegati, M. (2021). The macroeconomic effects of default and debt restructuring: An agent based exploration, *International Review of Economics & Finance* **76**: 1146–1163.
- Thakor, A. V. (2020). Fintech and banking: What do we know?, *Journal of Financial Intermediation* **41**: 100833.

- Trichet, J.-C. (2010). Reflections on the nature of monetary policy non-standard measures and finance theory, opening address at the ecb central banking conference. frankfurt, germany, 18.
- Upper, C. (2011). Simulation methods to assess the danger of contagion in interbank markets, *Journal of Financial Stability* **7**(3): 111–125.
- Watkins, C. J. and Dayan, P. (1992). Q-learning, *Machine learning* **8**(3-4): 279–292.
- Zhang, Z., Zohren, S. and Roberts, S. (2020). Deep reinforcement learning for trading, *The Journal of Financial Data Science* **2**(2): 25–40.

A A sensitivity analysis on model parameters

In this appendix, we investigate the performances of the learning algorithm by varying some key parameters. The first investigated parameter, β , governs the network topology (see Grilli et al. (2014), for a mathematical explanation). As the intensity of choice increases, the interbank architecture ranges from a random configuration to a star one. The effect of the network topology on the interbank system is studied by changing β from 0 to 40 with steps of 2. The second parameter we consider is fire sale price ρ . An increase in ρ impacts both lenders and borrowers. On the one hand, it compensates the losses that lenders incur due to the failure of their clients (see Eq. 4). On the other hand, a higher fire-sale increases the likelihood that the borrower, rationed in the interbank market, can face deposit repayments. Here we vary the fire-sale price, ρ , from 0.1 to 0.5 with steps of 0.1. Thirdly, we modify the skewness of the distribution of the random shock affecting the bank deposit at the beginning of each period. Recalling the equation for the deposit movements as $D_t^i = D_{t-1}^i(\mu + \omega U(0, 1))$, we remark that it allows us to reproduce bearish and bullish market periods. The uniformly distributed noise component can be shifted towards more negative or positive shocks at convenience to represent different market situations. Having fixed $\mu = 0.7$ in our simulations, we let ω vary from 0.52 to 0.6 with steps of 0.02, corresponding to a highly negatively skewed and perfectly symmetrical shock distribution.

The role of μ and ω is critical to regulate the magnitude of the aggregated shock that affects the interbank system. Precisely, μ and ω determine the probability of the sign of the deposit's shock. When $\mu = 0.7$ and $\omega = 0.6$, the likelihood of a negative shock is equal to that of a positive one. This parameter configuration corresponds to a consistent stock-flow model, where on average, the other half of the market participants recover what is eroded by the adverse market condition. Even though we do not necessarily respect all the requirements to guarantee the stock-flow consistency, we have checked that when $\mu = 0.7$ and $\omega = 0.6$, the total number of assets for each bank matches the total number of liabilities, hence the aggregated balance sheet of the system sum to 0. In our baseline model, we used $\mu = 0.7$ and $\omega = 0.55$ to favor more adverse shocks and, therefore, to have more interbank market activity to cover such needs.

The last part of this appendix is dedicated to investigating the effects of a change in the reserve ratio, \hat{r} , previously set at 2%. This analysis has a twofold value. On the one hand, it is a further experiment on the robustness of the model by changing the parameters space. On the other, it corresponds to a conventional monetary policy.

In all these experiments, we run our model 100 times for different values of the initial seed generating the pseudo-random numbers over a time span of $T = 1000$ periods. Moreover, all the agents' initialization parameters, except for the variations studied here, coincide with those presented in Sec. 3.

Let us begin the analysis by focusing on the three-parameter variations' implications on the model's results. Each parameter variation represents a different configuration of the banking system, which is used to test the different strategies over 100 simulations. The cumulative reward of these simulations is then averaged to obtain the mean values and the respective confidence interval for the reinforcement learning strategy and the random strategy. Fig. 18 shows the average cumulative reward over the 100 simulations as a function of a single parameter variation. We notice that the performance of the reinforcement learning algorithm solved with the PPO procedure is still superior with respect to the random strategy for all three sensitivity cases presented. Therefore, we can conclude that the effect analysis in the main paper still holds if one modifies some characteristics of the underlying financial system.

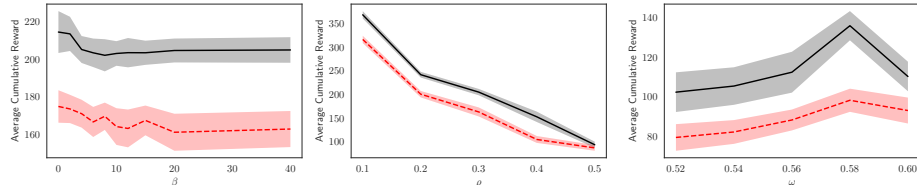


Fig. 18: Average cumulative fitness of the system as a function of changes in β , ρ , and ω in the first, second, and third panels, respectively. The reinforcement learning algorithm is in solid black, while the random strategy is in dashed red.

In Fig. 19 we show the sensitivity of the average values, over all the 100 simulations and a rolling window of 100 timesteps, of relevant quantities at the systemic level with respect to the three parameters described above¹³. **Before going into**

¹³ We refer the reader to Sec. 3.3 for a detailed explanation on the implementation of Fig. 19

the details concerning the systemic impact of the single parameters, we can observe that the reinforcement learning strategy (solid black in Fig.19) consistently outperforms the random strategy (dashed red line in Fig.19) over all parameters and variables considered¹⁴. The system generated with the reinforcement learning algorithm produces, on the one hand, higher liquidity and more credit channels and, on the other hand, lower rationing, bankruptcies, and leverage than the one with the random algorithm.

Let us now turn to the analysis that variations in each parameter have on the system's stability. In the first column of Fig. 19, we show the effects that the intensity of choice, β , has on the systemic variables. When β increases from 0 to 40, the liquidity and the credit channels increase to $\beta = 10$ and stabilize. **This pattern occurs in both scenarios (i.e., with optimal and random strategies).** The underlying reason for this dynamic is as follows: a β value greater than or equal to 10 generates a stable topology in the interbank network, which makes the investigated values insensitive to further changes in the parameter. Similar to the trend of the previous variables are the leverage dynamics, which increase with β but at a decreasing rate, **which is confirmed for both the adopted strategies.** Indeed, the more liquidity is available in the system, the more exchange of loans between banks happens. Finally, an increasing β causes the amount of rationing of the system to decrease **in both the considered scenarios, while the failures of the agent happen to be stable over the period under the optimal strategy or increase under the random scenario.**

In the second column of Fig. 19, we focus on the effects produced by a variation in the fire-sale price. An increase of ρ protects both lenders and borrowers from losses, and it is beneficial when looking at the liquidity up to $\rho = 0.3$. From that level, borrowers do not enter the interbank market frequently because they can cover their needs by selling their long-term assets at a satisfactory price. This is also reflected in the amount of rationing and failures that decrease when ρ is above 0.3. The leverage immediately decreases with ρ because the increase

¹⁴ To appreciate the statistical significance of the reinforcement learning strategy with respect to the random strategy, we performed a series of T-tests for each variable in the figures presented. The results show a statistically significant difference between each pair of curves at least the 5% level. We omitted here the table, including the p-values that are available under requests, as well as the results of the sensitivity analysis that we performed on the parameters \hat{d} , χ , ϕ and ξ .

in the system's liquidity is more than compensated by the increase in equity since lenders are usually repaid by borrowers and do not lose parts of their equity. **The dynamics produced by the fire-sale price variation are valid when observing the system with the optimal signal and the one with the random signal.** Finally, in the last column of the figure, the impact of the deposit's motion is investigated. The increase of the ω parameter causes an increase in liquidity since the shocks become gradually less and less harmful. This also explains the decrease in the leverage and the rationing because banks are less negatively impacted by the deposit shock and, consequently, need to gather less money from the market. For the same reason, the amount of credit channels decreases with a more symmetric shock distribution. In contrast, the failures are substantially stable, except for a higher variability when ω describes a highly asymmetric shock. **Also, for this last parameter, the system dynamics produced with the optimal signal follow the same trend as those obtained with the random signal.**

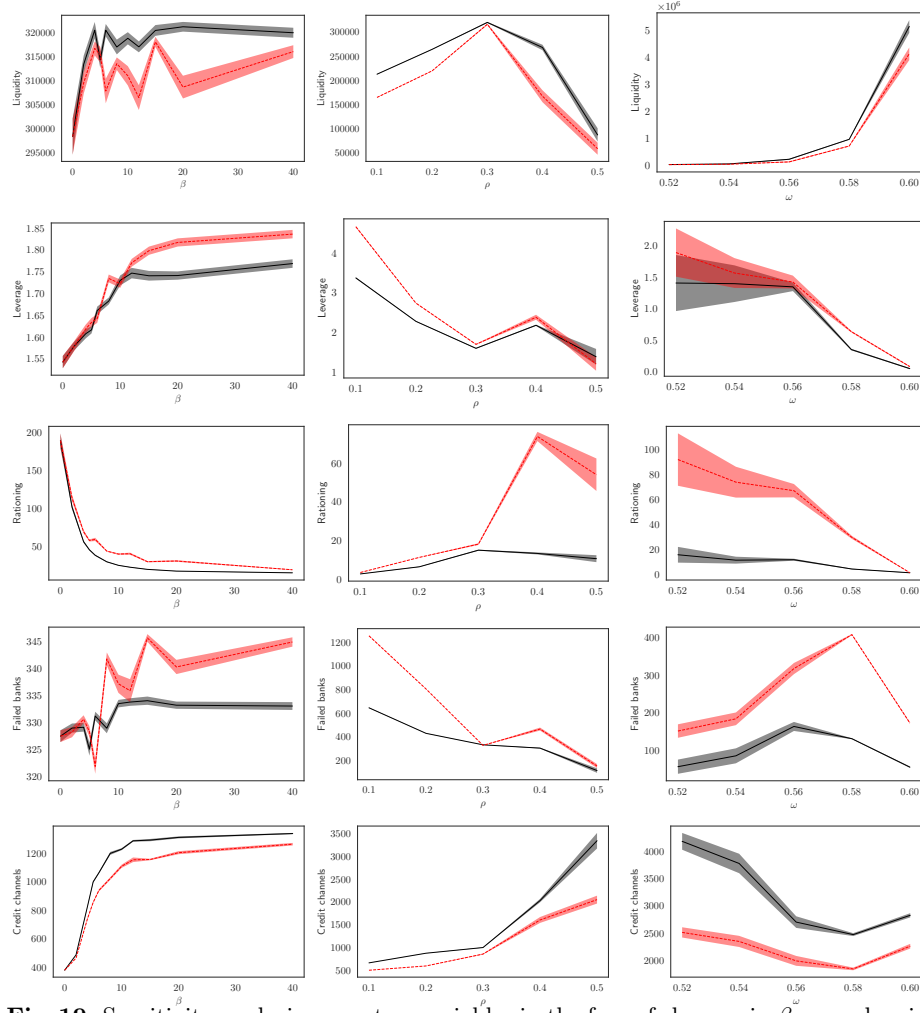


Fig. 19: Sensitivity analysis on system variables in the face of changes in β , ρ , and ω , in the first, second and third columns, respectively. The reinforcement learning algorithm is in solid black, while the random strategy is in dashed red.

In the last part of our analysis, we study how the system's resilience varies as the reserve requirement ratio varies from 1% to 10%. Fig20 shows the sensitivity of the average values over all the 100 simulations and a rolling window of 100 timesteps of relevant quantities at the systemic level with respect to the variation of \hat{r} . Before describing the effects of the contractionary monetary policy on market stability, it is worth noting that the system obtained through the optimized η (solid black line) consistently outperforms that with the random signal (dashed red line). The former always generates higher liquidity, lower leverage, and several failures. If we now observe the systemic effects of the increase in reserve ratio in the framework with the optimized signal, we can see an inverted U-shaped trend in liquidity. For \hat{r} -values between 1% and 5% , liquidity increases, showing that a non-excessively high reserve ratio promotes interbank stability by decreasing the number of failures. However, when the central bank imposes a reserve ratio above 5%, the contractionary effect of the policy takes over. The system becomes less liquid, and this causes a spike in failures as banks can no longer cope with their adverse deposit shocks. Finally, the behavior of the leverage, always in the context of the optimal signal, is timidly monotonically increasing with \hat{r} (see black line in the right-hand panel of Fig.20). For values of \hat{r} up to 5%, the leverage increases due to the rise in the granting of a loan. Above this threshold, the increase in leverage is mainly caused by the higher number of bankruptcies, which negatively impacts the net worth of financial institutions.

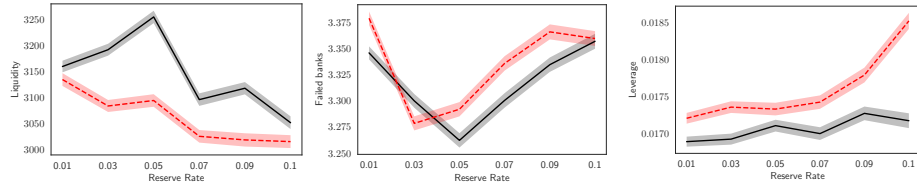


Fig. 20: Sensitivity analysis on the liquidity, number of bankruptcies, and leverage in the face of changes in reserve ratio \hat{r} in the first, second, and third columns, respectively. The reinforcement learning algorithm is in solid black, while the random strategy is in dashed red.

B Algorithms and Hyperparameters

The PPO algorithm is easier to implement than a trust-region method (Schulman et al.; 2015) and easier to tune with respect to of Deep-Q network (DQN) Mnih et al. (2015) or its continuous counterpart (Lillicrap et al.; 2015). Our implementation of PPO follows Andrychowicz et al. (2020), which performs an extensive empirical study of the effect of implementation and parameter choices on PPO performances. Even if we use the algorithm in a different context than their testbed, we follow the direction of their results in order to tune our hyperparameters.

As described in the main, we implement PPO in an actor-critic setting without shared architectures. When used to parametrize discrete strategies, policy gradient methods like PPO output a normalized set of logits to get the corresponding probabilities. Then, a greedy strategy selects the action which obtains the maximum probability. The entropy bonus guarantees exploration during training in the objective function.

The on-policy feature of PPO makes the training process episodic so that experience is collected by interacting with the environment and then discarded immediately once the strategy has been updated. In principle, on-policy learning appears a more obvious learning setup, even if it comes with some caveats. It makes the training less sample efficient and computationally expensive since a new sequence of experiences must be collected after each update step. In this process, the advantage function is computed before the optimization steps, when the discounted sum of returns over the episode can be computed. In order to increase the training efficiency, after one sweep through the collected samples, we compute the advantage estimator again and perform another sweep through the same experience. This trick reduces the computational expense of recollecting experiences and increases the sample efficiency of the training process. Usually, we do at most three sweeps (epochs) over a set of collected experiences before moving on and collecting a new set.

The gradient descent optimizer is Adam (Kingma and Ba (2014)), which performs a batch update of size 100 with a learning rate of 0.005. Since the data are not all available in a reinforcement learning setting at the beginning of the training, we can not normalize our input variables as usual in the preprocessing step of a supervised learning context. Hence, we add a Batch Normalization layer (Ioffe and Szegedy (2015)) before the first hidden layer to normalize the inputs batch by batch and obtain the same effect.

Maximizing the objective function that returns the gradient in Eq. 15 is unstable since updates are not bounded and can move the strategy too far from the local optimum. Similarly to TRPO (Schulman et al. (2015)), PPO optimizes an alternative objective to mitigate the instability

$$J^{\text{CLIP}}(\theta, \psi) = \mathbb{E}_{\pi_{\theta}} \left[\min \left(r(\theta) \hat{\mathbb{A}}(s, a; \psi), \text{clip}(r(\theta), 1 - \epsilon, 1 + \epsilon) \hat{\mathbb{A}}(s, a; \psi) \right) \right] \quad (23)$$

where $r(\theta) = \frac{\pi(A_t|S_t;\theta)}{\pi(A_t|S_t;\theta_{\text{old}})}$ is a ratio indicating the relative probability of an action under the current strategy with respect to the old one. Instead of introducing a hard constraint as in TRPO, the ratio is bounded according to a tolerance level ϵ to limit the magnitude of the updates. The combined objective function in Eq. 17 can be easily optimized by the PyTorch’s automatic differentiation engine, which quickly computes the gradients with respect to the two sets of parameters θ and ψ . The implemented advantage estimator depends on the parameterized value function V_{ψ} and is a truncated version of the one introduced by (Mnih et al. (2016)) for a rollout trajectory (episode) of length T :

$$\hat{\mathbb{A}}_t = \delta_t + (\gamma\tau)\delta_{t+1} + \dots + \dots + (\gamma\tau)^{T-t+1}\delta_{T-1} \quad (24)$$

where $\delta_t = r_t + \gamma V_{\psi}(s_{t+1}) - V_{\psi}(s_t)$, γ is a discount rate with the same role of ρ in DQN and τ is the exponential weight discount which controls the bias-variance trade-off in the advantage estimation. The generalized advantage estimator (GAE) uses a discounted sum of temporal difference residuals.