

Learning Macroeconomic Policies based on Microfoundations: A Stackelberg Mean Field Game Approach

Qirui Mi^{1,2}, Zhiyu Zhao^{1,2}, Siyu Xia^{1,2}, Yan Song¹, Jun Wang³, Haifeng Zhang^{1,2,4}

¹ Institute of Automation, Chinese Academy of Sciences ² School of Artificial Intelligence, University of Chinese Academy of Sciences ³ University College London

⁴ Nanjing Artificial Intelligence Research of IA
miqirui2021@ia.ac.cn, haifeng.zhang@ia.ac.cn

Abstract

The Lucas critique emphasizes the importance of considering microfoundations—how micro-agents (i.e., households) respond to policy changes—in macroeconomic policymaking. However, due to the vast scale and complex dynamics among micro-agents, predicting microfoundations is challenging. Consequently, this paper introduces a Stackelberg Mean Field Game (SMFG) approach that models macroeconomic policymaking based on microfoundations, with the government as the leader and micro-agents as dynamic followers. This approach treats large-scale micro-agents as a population, to optimize macroeconomic policies by learning the dynamic response of this micro-population. Our experimental results indicate that the SMFG approach outperforms real-world macroeconomic policies, existing AI-based and economic methods, enabling the learned macroeconomic policy to achieve the highest performance while guiding large-scale micro-agents toward maximal social welfare. Additionally, when extended to real-world scenarios, households that do not adopt the SMFG policy experience lower utility and wealth than adopters, thereby increasing the attractiveness of our policy. In summary, this paper contributes to the field of AI for economics by offering an effective tool for modeling and solving macroeconomic policymaking issues.

1 Introduction

Macroeconomic policymaking is fundamental to the sustained development of an economy (Schneider and Frey 1988; Persson and Tabellini 1999). Governments influence economic production, wealth distribution, social stability, and welfare through economic policies such as interest rates, taxation, and fiscal spending (Sachs and Warner 1995). Therefore, accurately modeling and solving optimal macroeconomic policies, and evaluating the effects of their implementation are critical issues. According to the Lucas critique (Lucas Jr 1976), predicting the impact of macroeconomic policies effectively requires modeling individual micro-level behaviors—termed “*microfoundations*”—to forecast micro-level responses to policy shifts and aggregate these into macroeconomic effects. Nonetheless, the vast number of micro-agents (i.e., households) and the complex dynamic interactions among them render predictions based on microfoundations exceedingly challenging.

To address the challenges outlined, this paper introduces a Stackelberg Mean Field Game (SMFG) approach: initially

proposing a Dynamic Stackelberg Mean Field Game (Dynamic SMFG) framework to model macroeconomic policymaking based on microfoundations (see Figure 1). Within this framework, we design a model-free reinforcement learning algorithm, Stackelberg Mean-Field Reinforcement Learning (SMFRL), aimed at optimizing macroeconomic policies by learning the dynamic responses of a large-scale population of micro-agents. In the Dynamic SMFG framework, the government acts as the leader, iteratively adjusting macroeconomic policies, while a vast array of micro-agents, serving as followers, react accordingly. This process entails sequential decision-making, with both the leader and the followers striving to optimize their cumulative long-term benefits. Given the leader’s policy, the SMFG is reduced to a mean field game, treating the large-scale micro-agents as a population to study their dynamic response to policy shifts. This dynamic SMFG framework responds to the Lucas critique. Existing models often simplify Dynamic SMFGs into single-step decisions (Pawlick and Zhu 2017), repeated games (Alcantara-Jimenez and Clempner 2020), or scenarios with myopic followers (Zhong et al. 2021), which fail to capture the complexities of macroeconomic policymaking adequately. Therefore, we introduce the SMFRL algorithm to optimize long-term utility for both the macro-agent and numerous micro-agents under the dynamic SMFG framework. In the SMFRL algorithm, we utilize mean field approximation to reduce the dimensionality of the micro-agents’ high-dimensional joint states and actions and employ a leader-follower update to learn the micro-population’s best response to macroeconomic policy adjustments, as well as to optimize macroeconomic policies considering dynamic feedback from the micro-population. This learning method enhances the stability and convergence of policies without requiring transition information or prior knowledge of the environment.

Our experiments, conducted on the TaxAI (Mi et al. 2023) environment, firstly demonstrate SMFG approach’s superiority over AI Economist (Zheng et al.), the Saez tax (Saez 2001), the 2022 U.S. Federal tax, and free market policy (Backhouse 2005) from perspectives of macroeconomic policy evaluation and dynamic responses to economic shocks. Further, through ablation studies, we verify the significant roles of mean field approximation and leader-follower update in terms of macroeconomic policy effects, convergence, and maximal social welfare. When applied in real-world scenar-

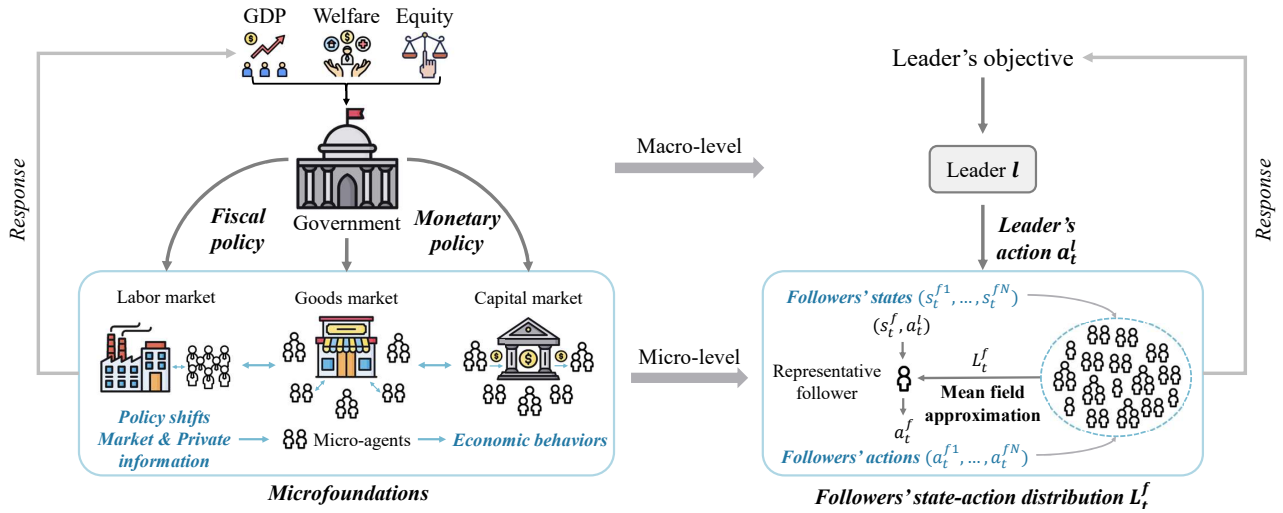


Figure 1: The mapping from macroeconomic policymaking (left) to Dynamic Stackelberg Mean Field Games (right). On the left, macro-level government actions, including fiscal or monetary policies, aim to optimize GDP, welfare, or equity. The microfoundations reflect the economic behaviors of large-scale micro-agents responding to shifts in macroeconomic policy. In Dynamic SMFGs, the government acts as the leader, and micro-agents as followers, respond to the leader’s policy.

ios, the SMFG approach shows that households refusing the policy have lower wealth and utility than adopters, which enhances the SMFG policy’s attractiveness.

In summary, this paper’s contributions are threefold:

1. We introduce the Dynamic Stackelberg Mean Field Games framework, a novel approach for modeling macroeconomic policymaking based on microfoundations, responding to the Lucas critique.

2. We propose a model-free Stackelberg Mean Field Reinforcement Learning algorithm, designed to optimize macroeconomic policies by learning the dynamic responses of large-scale micro-agents within the Dynamic SMFG framework.

3. Experimental results demonstrate the superiority of our method in modeling and solving macroeconomic policymaking, and further explore its effectiveness in real-world settings.

Our code is available in an anonymous GitHub repository — https://anonymous.4open.science/r/SMFG_macro_7740

2 Related work

Economic Methods

In the field of economics, research on economic policy typically relies on theoretical model studies, empirical analysis, and econometric methods. Classic theoretical models, such as the IS-LM model (Hicks 1980; Gali 1992), are applicable for analyzing short-term policy effects but overlook long-term factors and price fluctuations. The aggregate demand-aggregate supply (AD-AS) model (Dutt 2006; Lee, Pesaran, and Smith 1997) focuses on the relationship between aggregate demand and supply, integrating short-term and long-term factors, but simplifies macroeconomic dynamics. The Solow growth model (Brock and Taylor 2010) provides a framework for economic growth but does not consider market imperfections and externalities. New Keynesian models (Blanchard and Galí 2007; Gabaix 2020) emphasize the stickiness of

prices and wages, suitable for explaining economic fluctuations and policy interventions, but are based on strong assumptions. DSGE models (An and Schorfheide 2007; Smets and Wouters 2007), built on micro-foundations for macroeconomic forecasting, offer a consistent framework for policy analysis but are limited in handling nonlinearities and market imperfections. Empirical analysis (Ramesh et al. 2010; Vedung 2017) primarily evaluates policy effects through historical data but faces limitations in data timeliness and accuracy. Econometric methods (Johnston and DiNardo 1963; Davidson, MacKinnon et al. 2004), such as regression and time series analysis, can quantify policy effects but are constrained by model settings and data quality. Additionally, the Saez tax model (Saez 2001) uses earnings elasticities to optimize tax rates, enhancing social welfare without deterring economic activity. This model continues to guide current fiscal policy through effective taxation strategies.

AI for Economics

Artificial intelligence (AI), such as reinforcement learning, may offer new perspectives for solving complex economic problems (Tilbury 2022). In terms of macroeconomics, AI Economist (Zheng et al.) employs curriculum learning on tax policy design, (Trott et al. 2021) studies the collaboration between central and local governments under COVID-19, as well as monetary policy (Hinterlang and Tänzer; Chen et al. 2023), international trade (Sch 2021), market pricing (Danassis et al. 2023). However, these studies consider the government as a single agent, neglecting the dynamic response of the households to policies. (Koster et al.) studies democratic AI by human voting models, (Yaman et al. 2023) examines the impact of social sanction rules on labor division. However, these works are based on simplified settings and remain a gap in real-world economic policy. In terms of microeconomics, (Shi 2021; Rui and Shi 2022; Atashbar and Aruhan Shi 2023) study the optimal saving and consumption problems of micro-agents, and other researchers

explore rational expectation equilibrium (Kuriksha 2021; Hill, Bardoscia, and Turrell 2021), multiple equilibria under real-business-cycle (Curry et al. 2022), the emergence of barter behavior (Johanson et al. 2022), optimal asset allocation and savings strategies (Ozhamaratli and Barucca 2022).

Stackelberg Mean Field Game

Stackelberg Mean Field Games (SMFGs) model the dynamic interactions between a leader and numerous homogeneous followers, enabling the development of macroeconomic policies grounded in microfoundations. The methods for SMFGs are divided into model-based and model-free algorithms. Model-based methods, such as (Bensoussan, Chau, and Yam 2015; Fu and Horst 2020; Dayanikli and Lauriere 2023; Bergault, Cardaliaguet, and Rainer 2023) employ forward-backward stochastic functional differential equations to ascertain the followers' mean field equilibrium prior to computing the optimal strategy for the leader. Additionally, studies such as (Bensoussan et al. 2017; Lin, Jiang, and Zhang 2018; Moon and Başar 2018; Du, Wu et al. 2019; Huang and Yang 2020) simplify scenarios through linear-quadratic formulations. (Guo, Hu, and Zhang 2022), convert SMFGs into min-max optimization problems. However, these model-based methods generally depend on simplified representations of SMFGs and explicit environmental transition information. In contrast, model-free approaches leverage interaction data with the environment to learn strategies without requiring detailed knowledge of the transition dynamics and reward functions. (Pawlick and Zhu 2017) explores the resolution of SMFGs in single-step decision scenarios. (Campbell et al. 2021) utilizes a modified deep backward stochastic differential equation framework to calculate the followers' equilibrium despite constraints on the leader's decisions. (Miao et al. 2024) investigates locally optimal defensive strategies given predefined attacker trajectories. Furthermore, (Li et al. 2024) learns followers' transition functions from empirical data, thus aiding the solution of SMFGs through the empirical Fokker-Planck equation. Nevertheless, the assumptions and constraints inherent in these algorithms often do not meet the complex needs of macroeconomic policymaking.

3 Stackelberg Mean Field Game Modeling

In this section, we provide details about macroeconomic policy-making problems and propose a modeling framework based on Dynamic Stackelberg Mean Field Games (Dynamic SMFGs), with the relevant mappings are shown in Figure 1.

Macroeconomic Policy-Making

To foster economic growth, enhance social welfare, or maintain social equity, governments implement economic policies such as fiscal and monetary policies. These policies directly influence micro-agents (i.e., households) and the markets they constitute, altering their expectations and behaviors. The aggregation of micro-level behaviors from numerous households manifests as macroeconomic phenomena. Therefore, in formulating policies, governments should model micro-agents' behaviors and predict their responses to new policies, which is so-called *microfoundations* in the Lucas critique.

However, in reality, micro-agents form a large and dynamically interacting group, each seeking to maximize personal benefits. This complexity significantly complicates predicting their responses to macroeconomic policy shifts.

To address these challenges, we propose using Dynamic SMFGs to model macroeconomic policy-making problems:

1. For **microfoundations**, we model the government as the leader who first sets policies, with micro-agents as followers who adjust their behaviors according to these policies. The interaction between the leader and followers forms a Stackelberg game, where the leader considers the followers' responses to make informed decisions.

2. For **large scale micro-agents**, we employ a Mean Field Game to model micro-interactions among a large number of households, treating them as a population and predicting their collective response to macroeconomic policy changes.

3. For **dynamic decision-making**, the formulation of macroeconomic policies is inherently a multi-step decision-making process, where both the government and households make decisions at each step to optimize long-term benefits, thereby constituting a dynamic game.

Dynamic Stackelberg Mean Field Games

For macroeconomic policy-making issues, we consider a dynamic SMFG with a leader and N homogeneous follower agents. At any time step $t \in \{0, \dots, T\}$, the leader selects an action $a_t^l \in \mathcal{A}^l$ based on the leader's state $s_t^l \in \mathcal{S}^l$ and policy $\pi^l : \mathcal{S}^l \rightarrow \mathcal{A}^l$. Subsequently, followers determine their actions based on the leader's action a_t^l and their private states $s_t^f \in \mathcal{S}^f$. The representative follower's action $a_t^f \in \mathcal{A}^f$ is derived from a shared policy $\pi^f : \mathcal{S}^f \times \mathcal{A}^l \rightarrow \mathcal{A}^f$. The sequences $\{\pi_t^l\}_{t=0}^T$ and $\{\pi_t^f\}_{t=0}^T$ are denoted as π^l and π^f , respectively. We denote the population state-action distribution of followers under leader action a_t^l and follower policy π^f as $L_t(s_t^f, a_t^f; \pi^f, a_t^l)$, and abbreviate it as $L_t(s_t^f, a_t^f)$.

$L_t(s_t^f, a_t^f; \pi^f, a_t^l) \in \mathcal{P}(\mathcal{S}^f \times \mathcal{A}^f)$, where $a_t^f = \pi^f(s_t^f, a_t^l)$

Followers' Mean Field Game Given the leader's policy π^l , the dynamic SMFG is simplified into a mean field game for the followers. At time $t \in \{0, \dots, T-1\}$, given $(a_t^l, s_t^f, a_t^f, L_t)$, the representative follower then receives a reward $r^f(s_t, a_t^l, a_t^f, L_t)$ and transitions to the next state s_{t+1}^f according to the transition probability $s_{t+1}^f \sim P(\cdot | s_t, a_t^l, a_t^f, L_t)$, where joint state $\mathbf{s}_t = \{s_t^l, s_t^f\}$. Each follower aims to find the optimal policy π^f that maximizes his cumulative reward over the time horizon:

$$J^f(\pi^l, \pi^f, L) = \mathbb{E}_{s_0^f \sim \mu_0^f, s_{t+1} \sim P} \left[\sum_{t=0}^T r^f(s_t, a_t^l, a_t^f, L_t) \right]$$

where $a_t^l = \pi^l(s_t^l)$, $a_t^f = \pi^f(s_t^f, a_t^l)$ and $L_t = L_t(s_t^f, a_t^f)$.

Definition 1 (Followers' Best Response for Leader's Policy). *Given any leader's policy $\pi^l \in \Pi^l$ and the followers' state-action distributions $L = \{L_t\}_{t=0}^T$, the followers' best response policy $\pi^{f*}(\pi^l, L)$ is defined as*

$$\pi^{f*}(\pi^l, L) \in \arg \max_{\pi^f} J^f(\pi^l, \pi^f, L).$$

Leader's Game At time $t \in \{0, \dots, T-1\}$, given the leader's state s_t^l , action a_t^l , the followers' policy π^f and state-action distribution L_t , the leader will receive a reward $r^l(s_t, a_t^l, L_t)$, and move to the next state s_{t+1}^l according to the transition probabilities $P(\cdot | s_t, a_t^l, L_t)$. The leader aims to find the optimal policy π^l to maximize the total expected reward:

$$J^l(\pi^l, \pi^f, L) = \mathbb{E}_{s_0^l \sim \mu_0^l, s_{t+1} \sim P} \left[\sum_{t=1}^T r^l(s_t, a_t^l, L_t) \right], \quad (1)$$

where $a_t^l = \pi^l(s_t^l)$, $a_t^f = \pi^f(s_t^f, a_t^l)$ and $L_t = L_t(s_t^f, a_t^f)$.

Definition 2 (Leader's Optimal Policy in Dynamic Stackelberg Mean Field Games). *Considering the followers' best response (π^f, L) to the leader's policy, which satisfies Definition 1, learning the leader's optimal policy π^{l*} in dynamic Stackelberg mean field games is equivalent to solving the following fixed-point problem given initial condition (μ_0^l, μ_0^f) :*

$$\begin{aligned} \pi^{l*} &\in \arg \max_{\pi^{l'}} J^l(\pi^{l'}, \pi^f, L) \\ \text{s.t. } \pi^f &\in \arg \max_{\pi^f} J^f(\pi^{l*}, \pi^f, L) \end{aligned}$$

where the followers' state-action distribution L_t satisfies the following McKean-Vlasov equation:

$$\begin{aligned} L_{t+1}(s_{t+1}^f, a_{t+1}^f) &= \sum_{s_t^l, a_t^l, s_t^f, a_t^f} L_t(s_t^f, a_t^f) \pi_t^l(a_t^l | s_t^l) \mu_t^l(s_t^l) \\ &\quad P(s_{t+1}^f | s_t^f, a_t^f, a_t^l, L_t) \pi_{t+1}^f(a_{t+1}^f | s_{t+1}^f), \\ \mu_{t+1}^l(s_{t+1}^l) &= \sum_{s_t^l, a_t^l} \mu_t^l(s_t^l) \pi_t^l(a_t^l | s_t^l) P(s_{t+1}^l | s_t^l, a_t^l, L_t). \end{aligned}$$

In conclusion, to study the issue of macroeconomic policy-making, we model it as a dynamic SMFG. Under this modeling framework, we can optimize the macroeconomic policy of the leader, namely the government, while considering the followers' best response over discrete timesteps (Definition 2).

4 Stackelberg Mean Field Reinforcement Learning

To solve macroeconomic policymaking under Dynamic SMFG framework, we propose the Stackelberg Mean Field Reinforcement Learning (SMFRL) algorithm, which is a centralized training with decentralized execution (CTDE) algorithm. We introduce the leader's policy $\pi^l(s^l)$ and a central critic Q-function $Q^l(s^l, a^l, s^f, \mathbf{a}^f)$ for the leader agent, and the shared policy $\pi^f(s^f, a^l)$ and $Q^f(s^l, a^l, s^f, \mathbf{a}^f)$ for each follower agent. SMFRL algorithm comprises two key components: mean field approximation and leader-follower update (see Figure 2).

Mean Field Approximation

Due to the dimensions of joint state s^f and action \mathbf{a}^f scaling with the number of follower agents, the standard Q-function $Q^l(s^l, a^l, s^f, \mathbf{a}^f)$ becomes infeasible to learn when dealing with a large number of followers. Therefore, we factorize

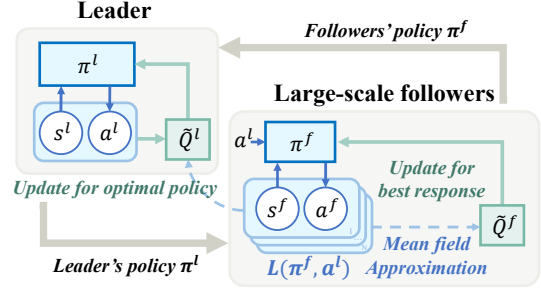


Figure 2: The architecture of SMFRL algorithm.

the standard Q-function using pairwise local interactions between the leader and each follower (Yang et al. 2018):

$$\begin{aligned} Q^l(s^l, a^l, s^f, \mathbf{a}^f) &= \frac{1}{N} \sum_i \tilde{Q}^l(s^l, a^l, s^{fi}, a^{fi}) \\ &\approx \mathbb{E}_{(s^f, a^f) \sim L} \tilde{Q}^l(s^l, a^l, s^f, a^f) \end{aligned} \quad (2)$$

Similarly, the i -th follower's Q-function $Q^f(s^l, a^l, s^f, \mathbf{a}^f)$ can be factorized using local interactions between the leader, the i -th follower, and each of the remaining followers:

$$Q^f(s^l, a^l, s^f, \mathbf{a}^f) \approx \mathbb{E}_{(s^f, a^f) \sim L} \tilde{Q}^f(s^l, a^l, s^{fi}, a^{fi}, s^f, a^f)$$

where (s^{fi}, a^{fi}) and (s^f, a^f) are the state-action pair for the i -th follower and another follower agent, respectively. We use $\hat{L}(\pi^f, a^l)$, a finite set of state-action pairs sampled under the follower's policy π^f and the leader's action a^l , to construct an empirical distribution that approximates the state-action distribution $L(s^f, a^f)$.

$$\hat{L}(\pi^f, a^l) = \{(s^{fi}, a^{fi})\}_{i=1}^N, \text{ where } a^{fi} = \pi^f(s^{fi}, a^l).$$

Leader-follower Update

To enhance the convergence and stability, we propose the leader-follower update: first, by fixing the leader's policy, we train the followers' shared policy and Q-networks towards the best response; subsequently, based on the followers' policies, we optimize the leader's policy and Q net, alternating these steps until convergence. We measure the distance between the agents' policies and their best responses by using exploitability.

Based on mean field approximation, we will train these networks π_{θ^l} and \tilde{Q}_{ϕ^l} for the leader agent, shared π_{θ^f} and \tilde{Q}_{ϕ^f} for the follower agents, with parameters $\theta^l, \phi^l, \theta^f$, and ϕ^f . To ensure training stability, we introduce target networks with parameters $\theta_-^l, \phi_-^l, \theta_-^f$, and ϕ_-^f . At any step $t \in \{0, \dots, T\}$, the tuple $(s_t^l, a_t^l, s_t^f, \mathbf{a}_t^f, s_{t+1}^l, s_{t+1}^f, r_t^l, \mathbf{r}_t^f)$ is stored in the replay buffer \mathcal{D} for training. When the number of followers is particularly large, a finite subset of state-action pairs $\{(s_t^{fi}, a_t^{fi})\}_{i=1}^N$ can be selected for storage.

Followers' Update for Best Response Given leader's policy π_{θ^l} , the followers' policy network π_{θ^f} is updated using a deterministic policy gradient approach (Silver et al. 2014).

The followers' policy gradient is estimated:

$$\nabla_{\theta^f} J \approx \mathbb{E}_{\mathbf{s}_t, a_t^l \sim \mathcal{D}, (s_t^f, a_t^f) \sim \hat{\mathcal{L}}_{\mathcal{D}}} \left[\nabla_{\theta^f} \pi_{\theta^f}(s_t^{fi}, a_t^l) \nabla_{a_-^i} \tilde{Q}_{\phi^f}(s_t^l, a_t^l, s_t^{fi}, a_t^{fi}, s_t^f, a_t^f) \right]$$

where $a_-^{fi} = \pi_{\theta^f}(s_t^{fi}, a_t^l)$, $\hat{\mathcal{L}}_{\mathcal{D}}$ denotes state-action pairs sampled from replay buffer \mathcal{D} . The action-value function \tilde{Q}_{ϕ^f} is updated by minimizing the mean squared error loss:

$$\mathcal{L}(\phi^f) = \mathbb{E}_{\substack{\mathbf{s}_t, a_t^l, \mathbf{s}_{t+1} \sim \mathcal{D} \\ (s_t^f, a_t^f) \sim \hat{\mathcal{L}}_{\mathcal{D}}}} \left[\left(y_t^{fi} - \tilde{Q}_{\phi^f}(s_t^l, a_t^l, s_t^{fi}, a_t^{fi}, s_t^f, a_t^f) \right)^2 \right]$$

$$y_t^{fi} = r_t^{fi} + \gamma \mathbb{E}_{(s_{t+1}^f, a_{t+1}^f) \sim \hat{\mathcal{L}}(\pi_{\theta^f}, a_{t+1}^l)} \tilde{Q}_{\phi^f}(s_{t+1}^l, a_{t+1}^l, s_{t+1}^{fi}, a_{t+1}^{fi}, s_{t+1}^f, a_{t+1}^f)$$

where $a_{t+1}^l = \pi_{\theta^l}(s_{t+1}^l)$, $a_{t+1}^{fi} = \pi_{\theta_-^f}(s_{t+1}^{fi}, a_{t+1}^l)$. The gradient of the loss function $\mathcal{L}(\phi^f)$ is derived as:

$$\nabla_{\phi^f} \mathcal{L}(\phi^f) = \mathbb{E}_{\substack{\mathbf{s}_t, a_t^l, \mathbf{s}_{t+1} \sim \mathcal{D} \\ (s_t^f, a_t^f) \sim \hat{\mathcal{L}}_{\mathcal{D}}}} \left[\left(y_t^{fi} - \tilde{Q}_{\phi^f}(s_t^l, a_t^l, s_t^{fi}, a_t^{fi}, s_t^f, a_t^f) \right) \nabla_{\phi^f} \tilde{Q}_{\phi^f}(s_t^l, a_t^l, s_t^{fi}, a_t^{fi}, s_t^f, a_t^f) \right].$$

Leader's Update for Optimal Policy Given followers' policy π_{θ^f} , the leader's policy π_{θ^l} is optimized by DPG approach, and the leader's policy gradient is estimated:

$$\nabla_{\theta^l} J \approx \mathbb{E}_{\mathbf{s}_t \sim \mathcal{D}, (s^f, a^f) \sim \hat{\mathcal{L}}(\pi_{\theta^f}, a_-^l)} \left[\nabla_{\theta^l} \pi_{\theta^l}(s_t^l) \nabla_{a_-^l} \tilde{Q}_{\phi^l}(s_t^l, a_-^l, s^f, a^f) \right], \text{ where } a_-^l = \pi_{\theta^l}(s_t^l).$$

This network is periodically updated to minimize the loss:

$$\mathcal{L}(\phi^l) = \mathbb{E}_{\substack{\mathbf{s}_t, a_t^l, \mathbf{s}_{t+1} \sim \mathcal{D} \\ (s_t^f, a_t^f) \sim \hat{\mathcal{L}}_{\mathcal{D}}}} \left[\left(y_t^l - \tilde{Q}_{\phi^l}(s_t^l, a_t^l, s_t^f, a_t^f) \right)^2 \right]$$

The target value y_t^l is given by:

$$y_t^l = r_t^l + \gamma \mathbb{E}_{(s_{t+1}^f, a_{t+1}^f) \sim \hat{\mathcal{L}}(\pi_{\theta^f}, a_{t+1}^l)} \tilde{Q}_{\phi^l}(s_{t+1}^l, a_{t+1}^l, s_{t+1}^f, a_{t+1}^f)$$

where $a_{t+1}^l = \pi_{\theta_-^l}(s_{t+1}^l)$, and γ is the discount factor. Differentiating the loss function $\mathcal{L}(\phi^l)$ yields the gradient utilized for training:

$$\nabla_{\phi^l} \mathcal{L}(\phi^l) = \mathbb{E}_{\mathbf{s}_t, a_t^l, \mathbf{s}_{t+1} \sim \mathcal{D}, (s_t^f, a_t^f) \sim \hat{\mathcal{L}}_{\mathcal{D}}} \left[\left(y_t^l - \tilde{Q}_{\phi^l}(s_t^l, a_t^l, s_t^f, a_t^f) \right) \nabla_{\phi^l} \tilde{Q}_{\phi^l}(s_t^l, a_t^l, s_t^f, a_t^f) \right].$$

The pseudocode for the SMFRL algorithm 1 is presented in Appendix A.

5 Experiment

In this section, we conduct experiments of the SMFG method, which is modeled by dynamic SMFGs and solved by the SMFRL algorithm, within the TaxAI (Mi et al. 2023) environment to answer the following questions:

1. For macroeconomic policy-making issues, how does our SMFG method compare to real-world policies, classic economic methods, and existing AI-based policy?

2. Are both the leader-follower update and mean-field approximation modules of the SMFG method essential?

3. Given the likely real-world scenario where some households do not adopt the SMFG policy, does the SMFG approach remain effective?

TaxAI Environment TaxAI is a simulation environment based on real data designed to study optimal taxation policies. It simulates the dynamic interactions between the government and large-scale households, to promote economic growth (i.e., GDP) as the objective.

Comparison of Macroeconomic Policies

This section compares the tax policy derived from SMFG with existing taxation policies in two aspects: (1) **Macroeconomic Policy Evaluation**: We assess the performance of different policies using key macroeconomic indicators. (2) **Dynamic Response**: At step 100, we simulate an economic shock caused by a financial crisis to evaluate the dynamic response capabilities of the different tax policies.

Baselines and Metrics We select four different types of tax policy baselines, including the **Free Market** (Backhouse 2005) policy, which entails no government intervention; the **Saez Tax** (Saez 2001), a well-known economic approach often advocated for specific tax reforms in practice; the **U.S. Federal Tax**, which is the actual policy derived from OECD data <https://data.oecd.org/>; and the **AI Economist** (Zheng et al.), a popular AI-based policy utilizing a two-level Proximal Policy Optimization (PPO) framework. Additionally, the **AI Economist-BC** method employs behavior cloning on real-world data for pretraining. In the TaxAI environment, we evaluate the performance of these policies using key macroeconomic indicators derived from real-world data. *Per Capita GDP* reflects the level of economic development, while *Income Gini* and *Wealth Gini* measure inequality in household income and wealth, respectively—a lower Gini index indicates greater social equality. The *Years* metric represents the sustainable duration of an economy, with a maximum cap of 300 years. *Average Wealth*, *Income*, and *Consumption* are crucial assessment metrics related to financial crises.

Macroeconomic Policy Evaluation In the context of optimizing economic growth, our SMFG method achieves the highest per capita GDP, social welfare, and the longest duration of outcomes across both $N = 100$ and $N = 1000$ scenarios. Table 1 shows the average values of macroeconomic indicators across 10 test runs for different algorithms under varying household numbers. Although the Gini index for SMFG is higher than that of the Saez Tax when $N = 100$, this is acceptable when optimizing GDP, as SMFG's per capita GDP is nearly 4 times that of the Saez Tax. When $N = 1000$, SMFG outperforms all other metrics, indicating that the advantages of our SMFG method increase with the number of households. Figure 7 in the Appendix provides the corresponding training curves.

Dynamic Response We simulate an economic shock similar to a financial crisis in the TaxAI environment: at step 100, all households lose 50% of their wealth. Due to the short duration (Years) under the Free Market and AI economist

Tax Policies	Per Capita GDP		Social Welfare		Income Gini		Wealth Gini		Years	
	100	1000	100	1000	100	1000	100	1000	100	1000
Free Market	1.37e+05	1.41e+05	32.97	334.79	0.89	0.90	0.92	0.93	1.10	1.00
Saez Tax	2.34e+12	6.35e+11	73.82	498.88	0.21	0.68	0.38	0.73	300.00	100.58
US Federal Tax	4.88e+11	1.41e+05	94.19	351.17	0.40	0.89	0.40	0.93	289.55	1.00
AI Economist-BC	4.24e+12	N/A	97.24	N/A	0.54	N/A	0.52	N/A	299.55	N/A
AI Economist	1.26e+05	N/A	72.81	N/A	0.88	N/A	0.91	N/A	1.00	N/A
SMFG	9.59e+12	1.10e+13	96.87	968.94	0.52	0.54	0.51	0.53	300.00	300.00

Table 1: Assessment results of 6 macroeconomic policies for different numbers of households $N=100$ and $N=1000$.

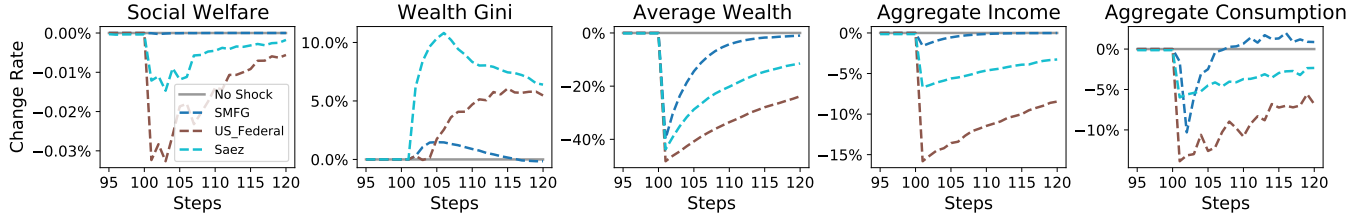


Figure 3: Dynamic response curves of 3 macroeconomic policies to economic shocks at step 100. The SMFG policy (dark blue line) exhibits the least fluctuation and the fastest recovery in key indicators, indicating superior dynamic response capabilities.

strategies, we focus on comparing the Saez Tax, US Federal Tax, and SMFG policy.

From Figure 3, we find that the SMFG policy (dark blue line) demonstrates the fastest dynamic response across all indicators: Social Welfare remains largely unaffected, while Average Income and Average Consumption recover to their original levels within about 5 steps. Wealth Gini and Average Wealth return to their previous levels within about 15 steps. This level of recovery is unmatched by any other tax policy.

Ablation Studies of SMFG

We conduct ablation studies to analyze the impact of the leader-follower update and mean field approximation modules in the SMFG method on macroeconomic indicators, as well as on convergence and social welfare.

Baselines We conduct comparisons based on the SMFG method by excluding the leader-follower update (denoted as **S**) and the mean field approximation (denoted as **MF**) individually:

1. **SMFG-S**: SMFG without leader-follower update. The government agent uses the independent DDPG (de Witt et al. 2020), while interactions among large-scale households are modeled as a mean field game and trained using the mean-field MARL algorithm (Yang et al. 2018).

2. **SMFG-MF**: SMFG without mean field approximation. Both the government and household agents are trained using a variant of MADDPG algorithm (Lowe et al. 2017), with the government making decisions first, followed by households.

3. **SMFG-S-MF**: SMFG without both leader-follower update and mean field approximation. In this case, both the government and household agents use independent DDPG.

Macroeconomic Indicators Analysis Figure 4 presents the results of the ablation studies on various macroeconomic indicators, evaluated across five seeds. The height of each bar represents the mean of the macroeconomic indicators, while the error bars indicate the variance. Our findings reveal

that: (1) Excluding either module leads to a significant performance drop compared to the full SMFG method, with this effect becoming more pronounced as the number of agents N increases. This underscores the critical importance of both the **S** and **MF** modules. (2) The SMFG-S variant outperforms SMFG-MF and SMFG-S-MF in terms of higher Per Capita GDP, social welfare, and Years, as well as lower Gini index. This suggests the dimensionality reduction of the **MF** module is crucial for the SMFG method’s overall effectiveness. (3) The low variance of the SMFG method indicates its stability.

Convergence and Social Welfare Table 2 presents the performance of different baselines across three game theory-related metrics: leader’s payoff, social welfare, and exploitability. Leader’s payoff reflects the performance of the leader agent (i.e., the government), social welfare corresponds to the total utility of all households, and exploitability measures the distance between current policy and Nash equilibrium. The methods for calculating the metrics and optimal values are provided in Appendix D.

From Table 2, we find: (1) SMFG achieves the highest leader’s payoff, the highest social welfare, and near-zero exploitability, indicating that the SMFG policy is closest to both the equilibrium and the maximal social welfare. (2) The high exploitability exhibited by SMFG-S indicates that the leader-follower update module facilitates the algorithm’s convergence toward equilibrium. (3) SMFG-MF outperforms SMFG-S-MF in terms of the leader’s payoff, suggesting that the leader-follower update module aids in optimizing the leader’s objective. However, it is challenging for SMFG-MF to optimize the followers’ social welfare.

Effectiveness Analysis of SMFG

Problem Description The above experiments have shown the effectiveness of the SMFG method in macroeconomic policymaking. Assuming the government is willing to implement macroeconomic policies learned by the SMFG method,

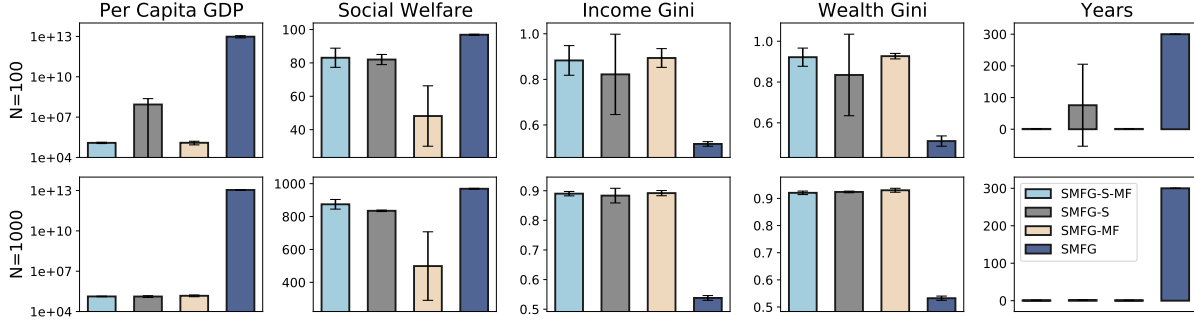


Figure 4: Ablation studies of SMFG method based on macroeconomic indicators for households' numbers $N=100$ and $N=1000$.

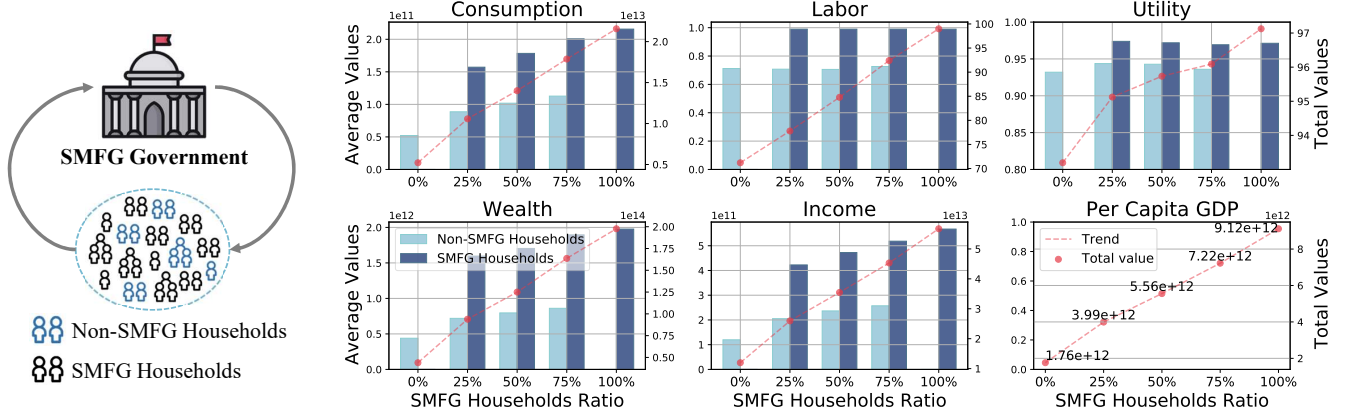


Figure 5: The left subfigure illustrates a scenario with heterogeneous households adopting various policies. The right subfigure presents test results for households' decisions (consumption, labor), microeconomic indicators (utility, wealth, income), and macroeconomic indicators (GDP) when there are different proportions of SMFG households ($N = 100$).

Methods\Num	Leader's Payoff		Exploitability		Social Welfare	
	100	1000	100	1000	100	1000
Optimal Value	\	\	0.	0.	100	1000
SMFG-S-MF	-774	-716	0.652	1.023	83	859
SMFG-S	-448	-856	3.161	1.213	82	782
SMFG-MF	-535	-441	0.782	0.725	54	499
SMFG (ours)	3294	3376	0.002	0.023	98	971

Table 2: Ablation studies of the SMFG method based on game theory metrics for $N=100$ and $N=1000$. Optimal values are provided for reference.

we aim to explore the impact of non-adoption by some households. In our experiments, we introduce **Non-SMFG Households** whose policies are trained using a behavior cloning algorithm based on real data, reflecting the decision-making styles of real-world households. Details of these behavior cloning experiments are detailed in Appendix E. Additionally, we denote governments and households that adopt the SMFG policy as **SMFG Government** and **SMFG Households**, respectively. We set different proportions of SMFG households (0%, 25%, 50%, 75%, and 100%) and evaluate our SMFG method using microeconomic and macroeconomic indicators. The results is shown in Figure 5, each subplot displays average values for two types of households (bars on the left Y-axis) and their aggregate values (dots on the right Y-axis).

GDP, a macroeconomic indicator, is shown only in aggregate, not for individual households.

Results From Figure 5, we observe several intriguing findings: (1) SMFG Households outperform Non-SMFG Households across microeconomic indicators such as wealth, income, and utility, suggesting that the policy of SMFG households is more attractive and likely to encourage broader adoption; (2) The presence of Non-SMFG Households does not prevent SMFG Households from achieving near-optimal utility, which is capped at 1 in the TaxAI environment; (3) As the proportion of SMFG Households increases, the per capita GDP also rises, indicating that the government should encourage households to adopt the SMFG policy.

6 Conclusion

Addressing issues in macroeconomic policymaking, this paper innovatively proposes the SMFG method, which includes dynamic SMFG for modeling and SMFRL algorithm for solving. The SMFG method optimizes government's policy by considering the dynamic responses of large-scale households to policy shifts, ultimately deriving optimal policies for both government and households. Experimental results demonstrate the SMFG method's superior performance, stability, and convergence over existing approaches. When extended to real-world scenarios, even if some households do not adopt the SMFG policy, their wealth and utility are not as high

as those who do, encouraging more households to adopt the SMFG policy. In conclusion, this paper contributes an effective approach to modeling and solving optimal macroeconomic policies in the fields of AI for economics and AI for social impact.

References

- Alcantara-Jimenez, G.; and Clempner, J. B. 2020. Repeated Stackelberg security games: Learning with incomplete state information. *Reliability Engineering & System Safety*, 195: 106695.
- An, S.; and Schorfheide, F. 2007. Bayesian analysis of DSGE models. *Econometric reviews*, 26(2-4): 113–172.
- Arrow, K. J.; and Debreu, G. 1954. Existence of an equilibrium for a competitive economy. *Econometrica: Journal of the Econometric Society*, 265–290.
- Atashbar, T.; and Aruhan Shi, R. 2023. AI and Macroeconomic Modeling: Deep Reinforcement Learning in an RBC Model.
- Backhouse, R. E. 2005. The rise of free market economics: Economists and the role of the state since 1970. *History of political economy*, 37(Suppl.1): 355–392.
- Bensoussan, A.; Chau, M.; Lai, Y.; and Yam, S. C. P. 2017. Linear-quadratic mean field Stackelberg games with state and control delays. *SIAM Journal on Control and Optimization*, 55(4): 2748–2781.
- Bensoussan, A.; Chau, M. H.; and Yam, S. C. P. 2015. Mean field Stackelberg games: Aggregation of delayed instructions. *SIAM Journal on Control and Optimization*, 53(4): 2237–2266.
- Bergault, P.; Cardaliaguet, P.; and Rainer, C. 2023. Mean Field Games in a Stackelberg problem with an informed major player. *arXiv preprint arXiv:2311.05229*.
- Blanchard, O.; and Galí, J. 2007. Real wage rigidities and the New Keynesian model. *Journal of money, credit and banking*, 39: 35–65.
- Brock, W. A.; and Taylor, M. S. 2010. The green Solow model. *Journal of Economic Growth*, 15: 127–153.
- Campbell, S.; Chen, Y.; Shrivats, A.; and Jaimungal, S. 2021. Deep learning for principal-agent mean field games. *arXiv preprint arXiv:2110.01127*.
- Chen, M.; Joseph, A.; Kumhof, M.; Pan, X.; and Zhou, X. 2023. Deep Reinforcement Learning in a Monetary Model. *arxiv:2104.09368*.
- Curry, M.; Trott, A.; Phade, S.; Bai, Y.; and Zheng, S. 2022. Analyzing Micro-Founded General Equilibrium Models with Many Agents Using Deep Reinforcement Learning. *arxiv:2201.01163*.
- Danassis, P.; Filos-Ratsikas, A.; Chen, H.; Tambe, M.; and Faltings, B. 2023. AI-driven Prices for Externalities and Sustainability in Production Markets. *arxiv:2106.06060*.
- Davidson, R.; MacKinnon, J. G.; et al. 2004. *Econometric theory and methods*, volume 5. Oxford University Press New York.
- Dayanikli, G.; and Lauriere, M. 2023. A Machine Learning Method for Stackelberg Mean Field Games. *arXiv preprint arXiv:2302.10440*.
- de Witt, C. S.; Gupta, T.; Makoviichuk, D.; Makoviyshuk, V.; Torr, P. H.; Sun, M.; and Whiteson, S. 2020. Is independent learning all you need in the starcraft multi-agent challenge? *arXiv preprint arXiv:2011.09533*.
- Du, K.; Wu, Z.; et al. 2019. Linear-quadratic Stackelberg game for mean-field backward stochastic differential system and application. *Mathematical Problems in Engineering*, 2019.
- Dutt, A. K. 2006. Aggregate demand, aggregate supply and economic growth. *International review of applied economics*, 20(3): 319–336.
- Fu, G.; and Horst, U. 2020. Mean-field leader-follower games with terminal state constraint. *SIAM Journal on Control and Optimization*, 58(4): 2078–2113.
- Gabaix, X. 2020. A behavioral New Keynesian model. *American Economic Review*, 110(8): 2271–2327.
- Gali, J. 1992. How well does the IS-LM model fit postwar US data? *The Quarterly Journal of Economics*, 107(2): 709–738.
- Guo, X.; Hu, A.; and Zhang, J. 2022. Optimization frameworks and sensitivity analysis of Stackelberg mean-field games. *arXiv preprint arXiv:2210.04110*.
- Hicks, J. 1980. IS-LM: an explanation. *Journal of post Keynesian economics*, 3(2): 139–154.
- Hill, E.; Bardoscia, M.; and Turrell, A. 2021. Solving Heterogeneous General Equilibrium Economic Models with Deep Reinforcement Learning. *arXiv:2103.16977*.
- Hinterlang, N.; and Tänzer, A. ??? Optimal Monetary Policy Using Reinforcement Learning.
- Huang, M.; and Yang, X. 2020. Mean field Stackelberg games: State feedback equilibrium. *IFAC-PapersOnLine*, 53(2): 2237–2242.
- Johanson, M. B.; Hughes, E.; Timbers, F.; and Leibo, J. Z. 2022. Emergent Bartering Behaviour in Multi-Agent Reinforcement Learning. *arxiv:2205.06760*.
- Johnston, J.; and DiNardo, J. 1963. Econometric methods.
- Koster, R.; Balaguer, J.; Tacchetti, A.; Weinstein, A.; Zhu, T.; Hauser, O.; Williams, D.; Campbell-Gillingham, L.; Thacker, P.; Botvinick, M.; and Summerfield, C. ??? Human-Centred Mechanism Design with Democratic AI. 6(10): 1398–1407.
- Kuriksha, A. 2021. An Economy of Neural Networks: Learning from Heterogeneous Experiences. *arxiv:2110.11582*.
- Lee, K.; Pesaran, M. H.; and Smith, R. 1997. Growth and convergence in a multi-country empirical stochastic Solow model. *Journal of applied Econometrics*, 12(4): 357–392.
- Li, P.; Yu, R.; Wang, X.; and An, B. 2024. Transition-Informed Reinforcement Learning for Large-Scale Stackelberg Mean-Field Games.
- Lin, Y.; Jiang, X.; and Zhang, W. 2018. An open-loop Stackelberg strategy for the linear quadratic mean-field stochastic differential game. *IEEE Transactions on Automatic Control*, 64(1): 97–110.

Lowe, R.; Wu, Y. I.; Tamar, A.; Harb, J.; Pieter Abbeel, O.; and Mordatch, I. 2017. Multi-agent actor-critic for mixed cooperative-competitive environments. *Advances in neural information processing systems*, 30.

Lucas Jr, R. E. 1976. Econometric policy evaluation: A critique. In *Carnegie-Rochester conference series on public policy*, volume 1, 19–46. North-Holland.

Malul, M.; and Bar-El, R. 2009. The gap between free market and social optimum in the location decision of economic activity. *Urban Studies*, 46(10): 2045–2059.

Mi, Q.; Xia, S.; Song, Y.; Zhang, H.; Zhu, S.; and Wang, J. 2023. TaxAI: A Dynamic Economic Simulator and Benchmark for Multi-Agent Reinforcement Learning. *arXiv preprint arXiv:2309.16307*.

Miao, L.; Li, S.; Wu, X.; and Liu, B. 2024. Mean-Field Stackelberg Game-Based Security Defense and Resource Optimization in Edge Computing. *Applied Sciences*, 14(9): 3538.

Moon, J.; and Başar, T. 2018. Linear quadratic mean field Stackelberg differential games. *Automatica*, 97: 200–213.

Ozhamaratli, F.; and Barucca, P. 2022. Deep Reinforcement Learning for Optimal Investment and Saving Strategy Selection in Heterogeneous Profiles: Intelligent Agents Working towards Retirement. *arxiv:2206.05835*.

Pawlick, J.; and Zhu, Q. 2017. A mean-field stackelberg game approach for obfuscation adoption in empirical risk minimization. In *2017 IEEE Global Conference on Signal and Information Processing (GlobalSIP)*, 518–522. IEEE.

Persson, T.; and Tabellini, G. 1999. Political economics and macroeconomic policy. *Handbook of macroeconomics*, 1: 1397–1482.

Ramesh, M.; Wu, X.; Howlett, M.; and Fritzen, S. 2010. *The public policy primer: Managing the policy process*.

Rui; and Shi. 2022. Learning from Zero: How to Make Consumption-Saving Decisions in a Stochastic Environment with an AI Algorithm. *arxiv:2105.10099*.

Sachs, J. D.; and Warner, A. 1995. Economic convergence and economic policies.

Saez, E. 2001. Using elasticities to derive optimal income tax rates. *The review of economic studies*, 68(1): 205–229.

Sch, A. A. O. 2021. Intelligence in the Economy: Emergent Behaviour in International Trade Modelling with Reinforcement Learning.

Schneider, F.; and Frey, B. S. 1988. Politico-economic models of macroeconomic policy: A review of the empirical evidence. *Political business cycles*, 239–275.

Shi, R. A. 2021. Can an AI Agent Hit a Moving Target. *arXiv preprint arXiv*, 2110.

Silver, D.; Lever, G.; Heess, N.; Degris, T.; Wierstra, D.; and Riedmiller, M. 2014. Deterministic policy gradient algorithms. In *International conference on machine learning*, 387–395. Pmlr.

Smets, F.; and Wouters, R. 2007. Shocks and frictions in US business cycles: A Bayesian DSGE approach. *American economic review*, 97(3): 586–606.

Tilbury, C. 2022. Reinforcement Learning in Macroeconomic Policy Design: A New Frontier? *arXiv preprint arXiv:2206.08781*.

Trott, A.; Srinivasa, S.; van der Wal, D.; Haneuse, S.; and Zheng, S. 2021. Building a Foundation for Data-Driven, Interpretable, and Robust Policy Design using the AI Economist. *arXiv:2108.02904*.

Vedung, E. 2017. *Public policy and program evaluation*. Routledge.

Yaman, A.; Leibo, J. Z.; Iacca, G.; and Wan Lee, S. 2023. The emergence of division of labour through decentralized social sanctioning. *Proceedings of the Royal Society B: Biological Sciences*, 290(2009).

Yang, Y.; Luo, R.; Li, M.; Zhou, M.; Zhang, W.; and Wang, J. 2018. Mean field multi-agent reinforcement learning. In *International conference on machine learning*, 5571–5580. PMLR.

Zheng, S.; Trott, A.; Srinivasa, S.; Parkes, D. C.; and Socher, R. ????. The AI Economist: Taxation Policy Design via Two-Level Deep Multiagent Reinforcement Learning. 8(18): eabk2607.

Zhong, H.; Yang, Z.; Wang, Z.; and Jordan, M. I. 2021. Can Reinforcement Learning Find Stackelberg-Nash Equilibria in General-Sum Markov Games with Myopic Followers? *arxiv:2112.13521*.

Reproducibility Checklist

1. This paper:

- Includes a conceptual outline and/or pseudocode description of AI methods introduced (yes/partial/no/NA) yes
- Clearly delineates statements that are opinions, hypothesis, and speculation from objective facts and results (yes/no) yes
- Provides well marked pedagogical references for less-familare readers to gain background necessary to replicate the paper (yes/no) yes

2. Does this paper make theoretical contributions? (yes/no) no

If yes, please complete the list below.

- All assumptions and restrictions are stated clearly and formally. (yes/partial/no)
- All novel claims are stated formally (e.g., in theorem statements). (yes/partial/no)
- Proofs of all novel claims are included. (yes/partial/no)
- Proof sketches or intuitions are given for complex and/or novel results. (yes/partial/no)
- Appropriate citations to theoretical tools used are given. (yes/partial/no)
- All theoretical claims are demonstrated empirically to hold. (yes/partial/no/NA)
- All experimental code used to eliminate or disprove claims is included. (yes/no/NA)

3. Does this paper rely on one or more datasets? (yes/no)

yes

If yes, please complete the list below.

- A motivation is given for why the experiments are conducted on the selected datasets (yes/partial/no/NA) yes
- All novel datasets introduced in this paper are included in a data appendix. (yes/partial/no/NA) no
- All novel datasets introduced in this paper will be made publicly available upon publication of the paper with a license that allows free usage for research purposes. (yes/partial/no/NA) NA
- All datasets drawn from the existing literature (potentially including authors' own previously published work) are accompanied by appropriate citations. (yes/no/NA) yes
- All datasets drawn from the existing literature (potentially including authors' own previously published work) are publicly available. (yes/partial/no/NA) yes
- All datasets that are not publicly available are described in detail, with explanation why publicly available alternatives are not scientifically satisfying. (yes/partial/no/NA) NA

4. Does this paper include computational experiments?

(yes/no) yes

If yes, please complete the list below.

- Any code required for pre-processing data is included in the appendix. (yes/partial/no) yes
- All source code required for conducting and analyzing the experiments is included in a code appendix. (yes/partial/no) yes
- All source code required for conducting and analyzing the experiments will be made publicly available upon publication of the paper with a license that allows free usage for research purposes. (yes/partial/no) yes
- All source code implementing new methods have comments detailing the implementation, with references to the paper where each step comes from (yes/partial/no) partial
- If an algorithm depends on randomness, then the method used for setting seeds is described in a way sufficient to allow replication of results. (yes/partial/no/NA) NA
- This paper specifies the computing infrastructure used for running experiments (hardware and software), including GPU/CPU models; amount of memory; operating system; names and versions of relevant software libraries and frameworks. (yes/partial/no) yes
- This paper formally describes evaluation metrics used and explains the motivation for choosing these metrics. (yes/partial/no) yes
- This paper states the number of algorithm runs used to compute each reported result. (yes/no) yes
- Analysis of experiments goes beyond single-dimensional summaries of performance (e.g., average; median) to include measures of variation, confidence, or other distributional information. (yes/no) no

- The significance of any improvement or decrease in performance is judged using appropriate statistical tests (e.g., Wilcoxon signed-rank). (yes/partial/no) yes
- This paper lists all final (hyper-)parameters used for each model/algorithm in the paper's experiments. (yes/partial/no/NA) yes
- This paper states the number and range of values tried per (hyper-) parameter during development of the paper, along with the criterion used for selecting the final parameter setting. (yes/partial/no/NA) yes

Algorithm 1: Stackelberg Mean Field Reinforcement Learning (SMFRL)

Initialize $\tilde{Q}_{\phi^l}, \tilde{Q}_{\phi_-^l}, \tilde{Q}_{\phi^f}, \tilde{Q}_{\phi_-^f}, \pi_{\theta^l}, \pi_{\theta_-^l}, \pi_{\theta^f}, \pi_{\theta_-^f}$, and replay buffer \mathcal{D} .

for epoch = 1 to M **do**

Receive initial state $\mathbf{s}_t = \{s_t^l, \mathbf{s}_t^f\}$.

for $t = 1$ to max-epoch-length **do**

For the leader agent, select action $a_t^l = \pi_{\theta^l}(s_t^l) + \mathcal{N}_t$; for each follower agent i , select action $a_t^{fi} = \pi_{\theta^f}(s_t^{fi}, a_t^l) + \mathcal{N}_t$.

Execute actions $\mathbf{a}_t = \{a_t^l, \mathbf{a}_t^f\}$ and observe reward $\mathbf{r}_t = \{r_t^l, \mathbf{r}_t^f\}$ and next state \mathbf{s}_{t+1} .

Store $(\mathbf{s}_t, \mathbf{s}_{t+1}, \mathbf{a}_t, \mathbf{r}_t)$ in \mathcal{D} .

$\mathbf{s}_t \leftarrow \mathbf{s}_{t+1}$.

for $j = 1$ to update-cycles **do**

Sample a random minibatch of N_b samples from \mathcal{D} .

Update follower agents by executing inner updates:

for $k = 1$ to inner-update-cycles **do**

Compute target for each follower:

$$y_t^{fi} = r_t^{fi} + \gamma \mathbb{E}_{(\mathbf{s}_{t+1}^f, \mathbf{a}_{t+1}^f) \sim \hat{L}(\pi_{\theta_-^f}, a_{t+1}^l)} \tilde{Q}_{\phi_-^f}(s_{t+1}^l, a_{t+1}^l, s_{t+1}^{fi}, a_{t+1}^{fi}, s_{t+1}^f, a_{t+1}^f),$$

Update follower's critic by minimizing the loss:

$$\mathcal{L}(\phi^f) = \mathbb{E}_{\mathbf{s}_t, a_t^l, \mathbf{s}_{t+1} \sim \mathcal{D}, (s_t^f, a_t^f) \sim \hat{L}_{\mathcal{D}}} \left[\left(y_t^{fi} - \tilde{Q}_{\phi^f}(s_t^l, a_t^l, s_t^{fi}, a_t^{fi}, s_t^f, a_t^f) \right)^2 \right],$$

Update follower's actor by policy gradient:

$$\nabla_{\theta^f} J \approx \mathbb{E}_{\mathbf{s}_t, a_t^l \sim \mathcal{D}, (s_t^f, a_t^f) \sim \hat{L}_{\mathcal{D}}} \left[\nabla_{\theta^f} \pi_{\theta^f}(s_t^{fi}, a_t^l) \nabla_{a_-^f} \tilde{Q}_{\phi^f}(s_t^l, a_t^l, s_t^{fi}, a_-^f, s_t^f, a_t^f) \right].$$

end for

Update the leader's networks:

Compute target for leader:

$$y_t^l = r_t^l + \gamma \mathbb{E}_{(\mathbf{s}_{t+1}^f, \mathbf{a}_{t+1}^f) \sim \hat{L}(\pi_{\theta^f}, a_{t+1}^l)} \tilde{Q}_{\phi_-^l}(s_{t+1}^l, a_{t+1}^l, s_{t+1}^f, a_{t+1}^f),$$

Update leader's critic by minimizing the loss:

$$\mathcal{L}(\phi^l) = \mathbb{E}_{\mathbf{s}_t, a_t^l, \mathbf{s}_{t+1} \sim \mathcal{D}, (s_t^f, a_t^f) \sim \hat{L}_{\mathcal{D}}} \left[\left(y_t^l - \tilde{Q}_{\phi^l}(s_t^l, a_t^l, s_t^f, a_t^f) \right)^2 \right],$$

Update leader's actor by sampled policy gradient:

$$\nabla_{\theta^l} J \approx \mathbb{E}_{\mathbf{s}_t \sim \mathcal{D}, (s^f, a^f) \sim \hat{L}(\pi_{\theta^f}, a_-^l)} \left[\nabla_{\theta^l} \pi_{\theta^l}(s_t^l) \nabla_{a_-^l} \tilde{Q}_{\phi^l}(s_t^l, a_-^l, s^f, a^f) \right], \text{ where } a_-^l = \pi_{\theta_-^l}(s_t^l).$$

end for

Update target network parameters periodically:

$$\phi_-^l \leftarrow \tau_{\phi} \phi^l + (1 - \tau_{\phi}) \phi_-^l,$$

$$\theta_-^l \leftarrow \tau_{\theta} \theta^l + (1 - \tau_{\theta}) \theta_-^l,$$

$$\phi_-^f \leftarrow \tau_{\phi} \phi^f + (1 - \tau_{\phi}) \phi_-^f,$$

$$\theta_-^f \leftarrow \tau_{\theta} \theta^f + (1 - \tau_{\theta}) \theta_-^f.$$

end for
end for

A SMFRL algorithm pseudocode

B Assumptions and Limitations

Assumptions This paper models the problem of macroeconomic policy-making as a Dynamic Stackelberg Mean Field Game, based on the following assumptions: (1) Homogeneous followers: We assume that a large-scale group of households is homogeneous. They can use different characteristics as observations to influence decisions, but there are commonalities in human behavioral strategies. (2) Rational Expectations: We assume that both macro and micro agents engage in rational decision-making, adjusting their future expectations based on observed information. However, in reality, the level of rationality varies among different households. Most households exhibit bounded rationality, and their expectations and preferences differ accordingly. (3) Experimental environment: We validate our approach through experiments in the TaxAI environment, based on the assumption that results within this environment can provide insights applicable to real-world scenarios. Addressing and potentially relaxing these assumptions will be a primary focus of our future research.

Limitations The limitations of our SMFG method will be thoroughly investigated in future work: (1) We plan to consider dynamic games involving multiple leaders and large-scale followers to explore policy coordination across various macroeconomic sectors. (2) We will continue to develop theoretical proofs for the equilibrium solutions in Stackelberg mean field games. Currently, our approach is empirically demonstrated by showing that followers converge toward their best responses and that the leader achieves higher performance compared to other baselines. (3) We intend to examine dynamic games between a leader and a large, heterogeneous group of followers, including scenarios where followers dynamically alter their strategies, to determine the optimal leader policy. Addressing these limitations will provide further insights applicable to real-world scenarios.

C Saez tax

The Saez tax policy is often considered a suggestion for specific tax reforms in the real world. The specific calculation method is as follows (Saez 2001). The Saez tax utilizes income distribution $f(z)$ and cumulative distribution $F(z)$ to get the tax rates. The marginal tax rates denoted as $\tau(z)$, are expressed as a function of pretax income z , incorporating elements such as the income-dependent social welfare weight $G(z)$ and the local Pareto parameter $\alpha(z)$.

$$\tau(z) = \frac{1 - G(z)}{1 - G(z) + \alpha(z)e(z)}$$

To further elaborate, the marginal average income at a given income level z , normalized by the fraction of incomes above z , is denoted as $\alpha(z)$.

$$\alpha(z) = \frac{zf(z)}{1 - F(z)}$$

The reverse cumulative Pareto weight over incomes above z is represented by $G(z)$.

$$G(z) = \frac{1}{1 - F(z)} \int_{z'=z}^{\infty} p(z') g(z') dz'$$

From the above calculation formula, we can calculate $G(z)$ and $\alpha(z)$ by income distribution. We obtain the data of income and marginal tax rate through the interaction between the agent and environment and store them in the buffer. It is worth noting that the amount of buffer is fixed.

To simplify the environment, we discretize the continuous income distribution, by dividing income into several brackets and calculating a marginal tax rate $\tau(z)$ for each income range. Within each tax bracket, we determine the tax rate for that bracket by averaging the income ranges in that bracket. In other words, income levels falling within the income range are calculated as the average of that range. In particular, when calculating the top bracket rate, it is not convenient to calculate the average because its upper limit is infinite. So here $G(z)$ represents the total social welfare weight of incomes in the top bracket, when calculating $\alpha(z)$, we take the average income of the top income bracket as the average of the interval.

Elasticity $e(z)$ shows the sensitivity of the agent's income z to changes in tax rates. Estimating elasticity is very difficult in the process of calculating tax rates, here we estimate the elasticity $e(z)$ using a regression method through income and marginal tax rates under varying fixed flat-tax systems, which produces an estimate equal to approximately 1.

$$e(z) = \frac{1 - \tau(z)}{z} \frac{dz}{d(1 - \tau(z))}$$

$$\log(Z) = \hat{e} \cdot \log(1 - \tau) + \log(\hat{Z}^0)$$

where $Z = \sum_i z_i$ when tax rates is τ .

D Game Theory Metrics

We will utilize the following metrics related to game theory to evaluate the effectiveness of the leader and follower policies: (1) The leader's payoff, which indicates the performance of the leader's policy in optimizing the leader's objective; (2) Exploitability, which measures the deviation of the agent's policy from the best response; (3) Social welfare, which assesses the deviation of the current state from the social optimum.

Leader's Payoff We define the leader's payoff using the long-term expected rewards of the leader's policy π^l over discrete timesteps, as detailed in Equation 1.

Exploitability Exploitability is a critical metric in evaluating the convergence of policies and quantifying the divergence from the best response strategy in game theory. For a follower, exploitability $\mathcal{E}^f(\pi^f; \pi^l)$ is defined as the difference in payoffs between the follower's actual policy π^f and its optimal response π^{f*} , given the leader's policy π^l . Formally, it is represented as:

$$\mathcal{E}^f(\pi^f; \pi^l) = J^f(\pi^l, \pi^{f*}, L) - J^f(\pi^l, \pi^f, L),$$

where J^f denotes the cumulative reward for the follower, defined in Section 3.

Similarly, the leader’s exploitability $\mathcal{E}^l(\pi^l; \pi^f)$ measures the payoff difference between the leader’s policy π^l and its best response π^{l*} , given the followers’ response policy π^f and state-action distribution L . This is given by:

$$\mathcal{E}^l(\pi^l; \pi^f) = J^l(\pi^{l*}, \pi^f, L) - J^l(\pi^l, \pi^f, L),$$

with J^l representing the cumulative reward for the leader, and $L = \{L_t\}_{t=0}^T$ detailing the state-action distribution for followers over time (see Section 3).

The overall exploitability, which measures the discrepancy from Nash equilibrium for both the leader and the followers, is defined as:

$$\mathcal{E}(\pi^l, \pi^f) = \mathcal{E}^f(\pi^f; \pi^l) + \mathcal{E}^l(\pi^l; \pi^f),$$

A near-zero value of $\mathcal{E}(\pi^l, \pi^f)$ indicates that the policies of both the leader and the followers are approaching their respective optimal strategies π^{l*} and π^{f*} , signifying an equilibrium state.

Social Optimum and Social Welfare In economic theory, the *Social Optimum* describes a state in which the allocation of resources achieves maximum efficiency, as measured by social welfare (Arrow and Debreu 1954; Malul and Bar-El 2009). Given the leader’s policy π^l and the representative follower’s policy π^f among large-scale followers, social welfare $\mathcal{SW}(\pi^l, \pi^f)$ is approximately calculated as the sum of the utility functions defined in Section 3 of the N followers:

$$\begin{aligned} \mathcal{SW}(\pi^l, \pi^f) &= \sum_{i=1}^N J^{fi}(\pi^l, \pi^f, L) \\ &= \mathbb{E}_{s_0^f \sim \mu_0^f, s_{t+1} \sim P} \left[\sum_{t=0}^T \sum_{i=1}^N r^{fi}(s_t, a_t^l, a_t^{fi}, L_t) \right] \end{aligned}$$

E Additional Results

Compute Resources

All experiments are run on 2 workstations: A 64-bit server with dual AMD EPYC 7742 64-Core Processors @2.25 GHz, 256 cores, 512 threads, 503GB RAM, and 2 NVIDIA A100-PCIE-40GB GPU. A 64-bit workstation with Intel Core i9-10920X CPU @ 3.50GHz, 24 cores, 48 threads, 125 GB RAM, and 2 NVIDIA RTX2080 Ti GPUs. The following Table 3 shows the approximate training times for several algorithms.

Experimental Results of Ablation Studies

In this section, we present additional experimental results from ablation studies, including training curves 6 and a test table 4, as well as experiments incorporating the use of behavior cloning as a pre-training strategy for follower agents. We find that the SMFG method without behavior cloning as pre-training still surpasses other baselines that utilize behavior cloning. More specifically, we compared SMFG with 5 baselines across 4 different experimental setups: without behavior cloning as pre-training for follower agents at $N=100$

Algorithm	N=100 (hours)	N=1000 (hours)
SMFG	4	14
SMFG-MF	3	16
SMFG-S	4	9
SMFG-S-MF	2	6
Free Market	0.25	2
Saez Tax	4	23
AI Economist	6.5	N/A

Table 3: The approximate training times for baselines in our experiments.

and $N=1000$ (marked as $N=100$ without BC and $N=1000$ without BC); with BC-based pre-training for follower agents at $N=100$ and $N=1000$ ($N=100$ -BC; $N=1000$ -BC). Figure 6 illustrates the training curves of 4 key macroeconomic indicators under these four settings. The solid line represents the average value of the metrics across the 5 random seeds, while the shaded area represents the standard deviation. Each row corresponds to one setting, and each column to a macroeconomic indicator, including per capita GDP, social welfare, income Gini, and wealth Gini. A rise in per capita GDP indicates economic growth, an increase in social welfare implies happier households and a lower Gini index indicates a fairer society. Each subplot’s Y-axis represents the indicators’ values, and the X-axis represents the training steps. Table 4 displays the test results of the 6 algorithms across 5 indicators, with each column corresponding to an experimental setting.

Figure 6 and Table 4 present two experimental findings: (1) Using BC as a pre-training method for the follower’s policy enhances the algorithms’ stability and performance. Comparing settings with and without BC (the first two rows), our method, SMFG, shows similar convergence outcomes; however, the performance of other algorithms significantly improves across all four indicators with BC-based pre-training. Furthermore, the training curves of each algorithm are more stable. (2) The SMFG method substantially outperforms other algorithms in solving SMFGs, both in large-scale followers and without pre-training scenarios. In the setting of $N=100$ -BC, SMFG achieved a significantly higher per capita GDP compared to other algorithms, while its social welfare and Gini index are similar to others, essentially reaching the upper limit. Besides, in $N=100$ without BC and $N=1000$ -BC, SMFG consistently obtains the most optimal solutions across all indicators.

Training Curves for Various Tax Policies

We compare the performance of 6 policies across four economic indicators under two settings: with $N=100$ and $N=1000$ households. Figure 7 displays the training curves and Table 1 shows the test results. Both Figure 7 and Table 1 indicate that the SMFG method significantly surpasses other policies in the task of optimizing GDP, and achieves the highest social welfare. When $N=100$, the Saez tax achieves the lowest income and wealth Gini coefficients, suggesting greater fairness. However, at $N=1000$, SMFG performs optimally

Settings	Algorithm	Per Capita GDP	Social Welfare	Income Gini	Wealth Gini	Years
N=100 without BC	Rule-based+Random	1.41e+05	69.27	0.89	0.92	1.00
	DDPG+Random	1.41e+05	70.91	0.88	0.92	1.00
	SMFG-MF	1.23e+05	48.17	0.89	0.93	1.00
	SMFG-S-MF	1.21e+05	83.09	0.88	0.92	1.00
	SMFG-S	8.66e+07	82.02	0.82	0.83	75.75
	SMFG	9.59e+12	96.87	0.52	0.51	300.00
N=1000 without BC	Rule-based+Random	1.42e+05	705.85	0.90	0.93	1.00
	DDPG+Random	1.42e+05	654.37	0.90	0.93	1.00
	SMFG-MF	1.48e+05	499.01	0.89	0.93	1.02
	SMFG-S-MF	1.33e+05	874.53	0.89	0.92	1.00
	SMFG-S	1.31e+05	834.93	0.88	0.92	1.50
	SMFG	1.10e+13	968.94	0.54	0.53	300.00
N=100 with BC	Rule-based+Real	3.66e+11	79.23	0.52	0.53	217.45
	DDPG+Real	2.03e+12	94.50	0.46	0.48	299.85
	SMFG-MF	6.38e+12	93.89	0.57	0.58	268.53
	SMFG-S-MF	7.41e+12	98.16	0.53	0.55	300.00
	SMFG-S	5.44e+12	98.21	0.50	0.52	300.00
	SMFG	1.01e+13	96.90	0.51	0.53	299.89
N=1000 with BC	Rule-based+Real	1.41e+05	334.79	0.90	0.93	1.00
	DDPG+Real	4.92e+11	527.09	0.75	0.79	100.68
	SMFG-MF	6.82e+12	954.88	0.56	0.62	278.50
	SMFG-S-MF	2.79e+12	512.19	0.77	0.81	100.68
	SMFG-S	1.13e+05	440.00	0.90	0.93	1.00
	SMFG	9.68e+12	975.15	0.52	0.51	300.00

Table 4: Performance metrics of six algorithms across different settings on key macroeconomic indicators.

across all economic indicators, while the effectiveness of other policies noticeably diminishes as the number of households increases. The Saez tax also reduces the Gini index, but not as effectively or stably as the SMFG.

Efficiency-Equity Tradeoff of Policies

In economics, the Efficiency-Equity Tradeoff is a highly debated issue. We find that our SMFG method is optimal in balancing efficiency-equity, except in cases of extreme concern for social fairness. In our study, we depict the economic efficiency (Per capita GDP) on the Y-axis and equity (wealth Gini) on the X-axis of Figure 8(a) for various policies. Different policies are represented by circles of different colors, with their sizes proportional to social welfare. Different circles of the same color correspond to different seeds. Figure 8 (a) shows that the wealth Gini indices for SMFG and AI Economist-BC are similar, but SMFG has a higher GDP, suggesting its superiority over AI Economist-BC. SMFG significantly outperforms the free market policy and AI Economist due to its higher GDPs and lower wealth Ginis. However, comparing SMFG with the Saez tax and the U.S. Federal tax policy in terms of both economic efficiency (GDP) and social equity (Gini) is challenging. Therefore, we introduce Figure 8 (b) to demonstrate the performance of different policies under various weights in a multi-objective assessment.

In Figure 8 (b), the Y-axis shows the weighted values of the multi-objective function $Y = \log(\text{per capita GDP}) + \alpha(\text{wealth Gini})$, and the X-axis represents the weight of the

wealth Gini index. For each weight α , we compute the multi-objective weighted values for those policies, represented as circles of different colors. Due to the logarithmic treatment of GDP in (b), when $\alpha = 10$, the overall objective focuses solely on social fairness; when $\alpha = 0$, the overall objective is concerned only with efficiency. Our findings in (b) reveal that only when $\alpha \geq 8$, which indicates a substantial emphasis on social equity, does the Saez tax outperform SMFG. However, SMFG consistently proves to be the most effective policy under a wide range of preference settings.

Behavior Cloning Experiments

We conduct behavior cloning based on real data to simulate the behavior strategies of households in realistic scenarios, which are then used in Experiment 5 to compare with SMFG Households. We collect the statistical data from the 2022 Survey of Consumer Finances (SCF) (<https://www.federalreserve.gov/econres/scfindex.htm>) as the real data buffer \mathcal{D}_{real} .

Based on real data, we fetch a large number of followers' state-action pairs $\{s^f, a^f\}$ from a real-data buffer \mathcal{D}_{real} for behavior cloning. For different settings of network structures, we have chosen two types of loss: when the neural network outputs a probability distribution of actions, we use the negative log-likelihood loss (NLL loss); when the neural network outputs action values, we employ the mean square error loss (MSE loss). Our goal is to find the optimal parameters θ as the follower's policy network π_θ initialization, thereby

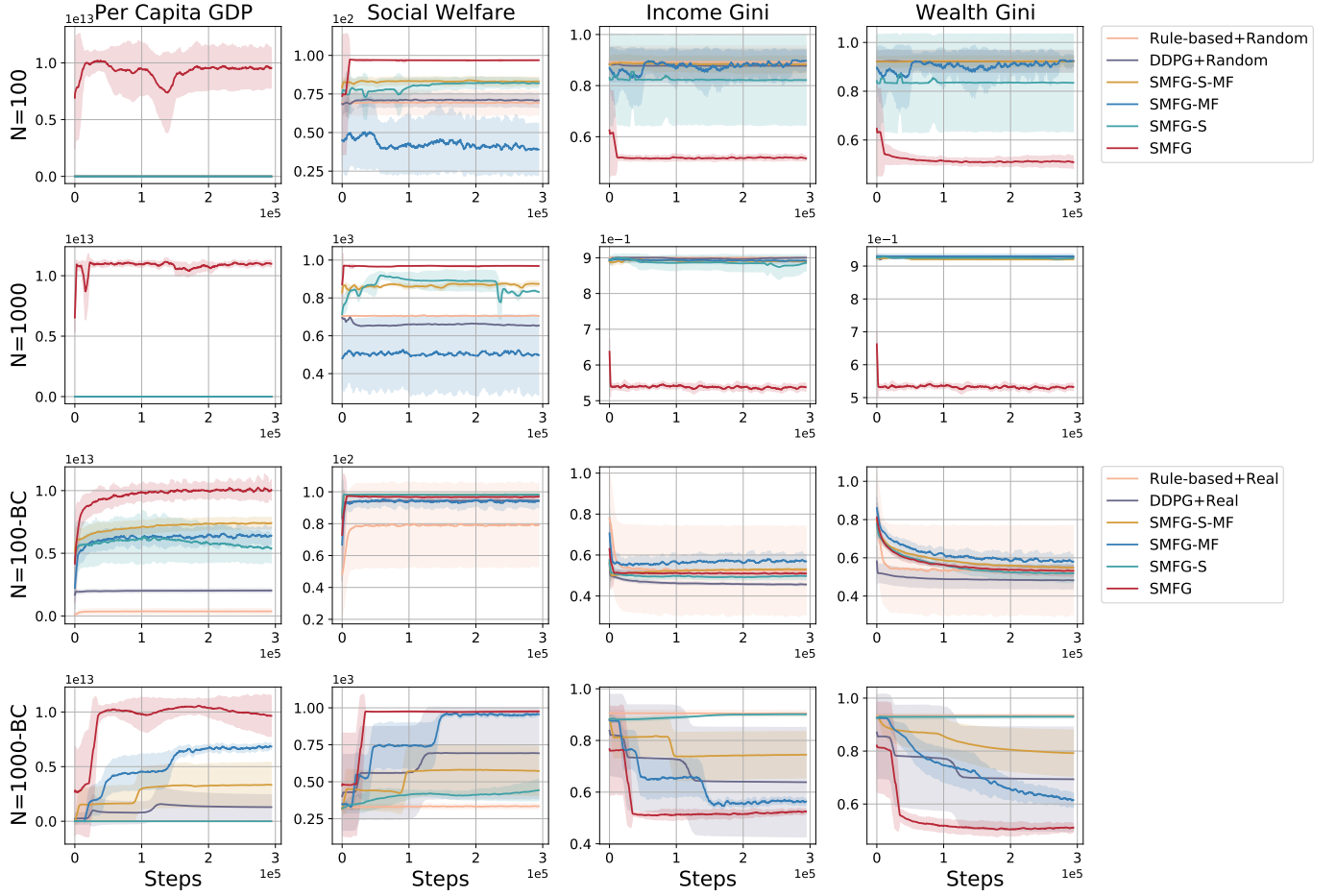


Figure 6: The training curves for 6 algorithms on 4 macroeconomic indicators, comparing settings without behavior cloning as pre-train (N=100 & N=1000) and with behavior cloning (N=100-BC & N=1000-BC).

minimizing the loss to its lowest convergence.

$$\begin{aligned} \min_{\theta} \mathcal{L}_{NLL} &= -\mathbb{E}_{s^f, a^f \sim \mathcal{D}} \log \pi_{\theta}(a^f | s^f), \\ \min_{\theta} \mathcal{L}_{MSE} &= \mathbb{E}_{s^f, a^f \sim \mathcal{D}} (a^f - a)^2 |_{a=\pi_{\theta}(s^f)}. \end{aligned}$$

This experiment conducts behavior cloning on networks for four different household policies: Multilayer Perceptron (MLP), AI economist’s network (MLP+LSTM+MLP), SMFG-S, and SMFG-MF network. The first two, as their network outputs, are probability distributions, use negative log-likelihood loss (Figure 9 left); the latter two’s networks employ deterministic policies, hence they use mean square error loss against real data (Figure 9 right). The loss convergence curve of behavior cloning is shown in Figure 9. It can be observed that the AI economist’s network, due to its complexity, struggles to converge to near -1 like MLP. The losses corresponding to MFRL and SMFG-MF can converge to below 0.1.

F Economic Model Details

Economic activities among households aggregate into labor markets, capital markets, goods markets, etc. In the labor market, households are the providers of labor, with the aggregate supply $S(W_t) = \sum_i^N e_t^i h_t^i$, and firms are the demanders of labor, with the aggregate demand $D(W_t) = \mathcal{L}_t$. When supply equals demand in the labor market, there exists an equilibrium price W_t^* that satisfies:

$$S(W_t^*) = D(W_t^*), \mathcal{L}_t = \sum_i^N e_t^i h_t^i.$$

In the capital market, financial intermediaries play a crucial role, lending the total deposits of households $A_{t+1} = \sum_i^N a_{t+1}$ to firms as production capital K_{t+1} , and purchasing government bonds B_{t+1} at the interest rate r_t . The capital market clears when supply equals demand:

$$K_{t+1} + B_{t+1} - A_{t+1} = (r_t + 1)(K_t + B_t - A_t)$$

In the goods market, firms produce and supply goods, while all households, the government, and physical capital invest-

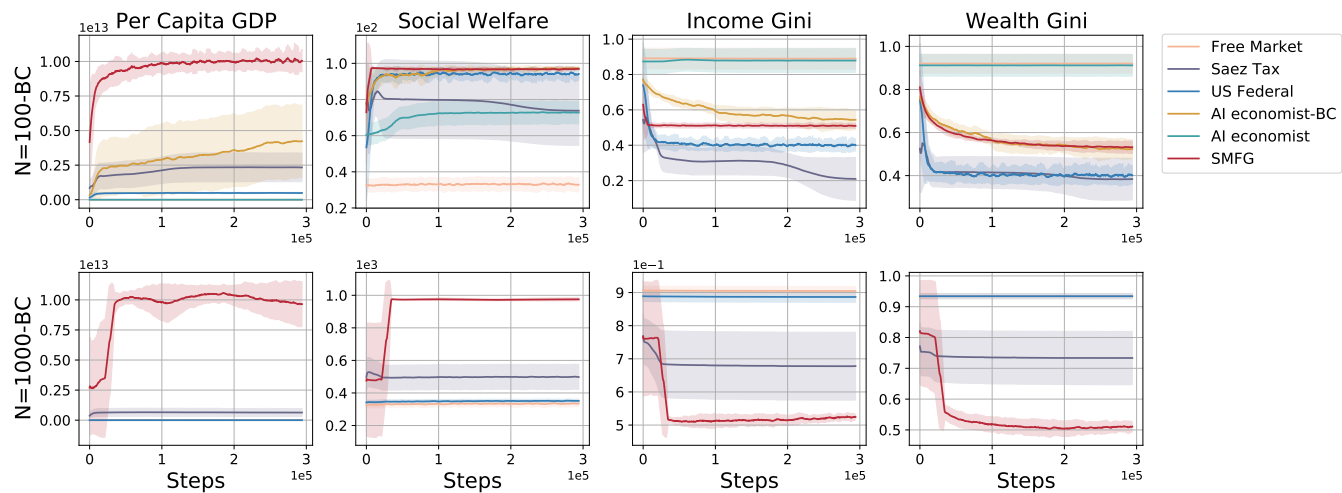


Figure 7: The training curves for 6 tax policies on 5 macroeconomic indicators ($N=100$ & $N=1000$ with BC)

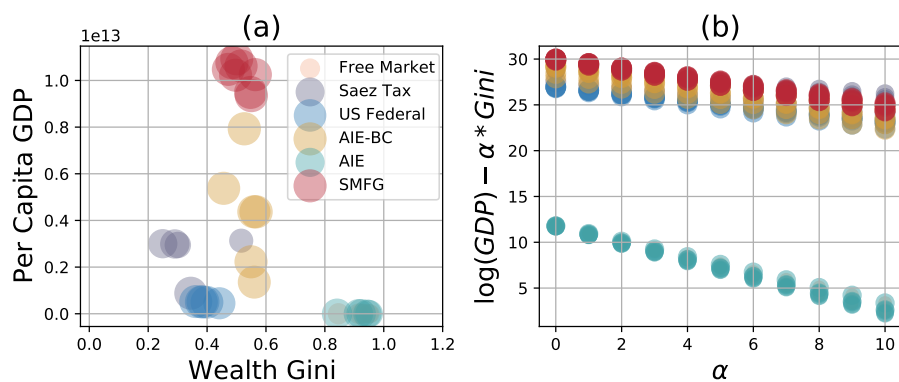


Figure 8: Comparative performance of various policies under multi-objective assessment (Efficiency-Equity).

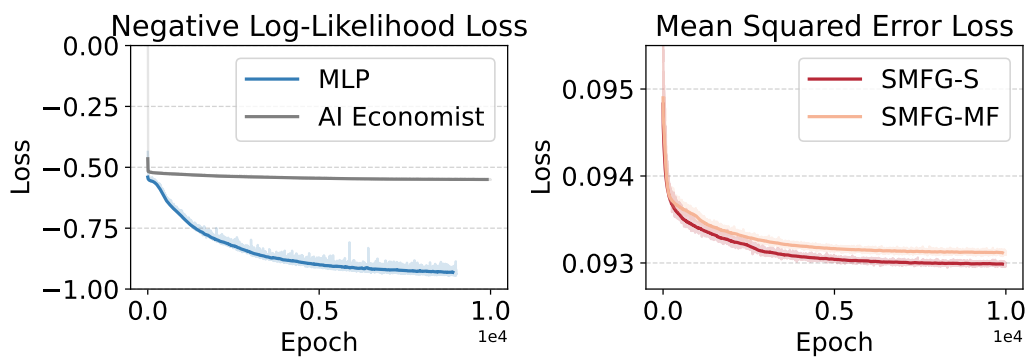


Figure 9: The behavior cloning loss for 4 networks in two loss types.

ments X_t demand them. The goods market clears when:

$$Y_t = C_t + G_t + X_t$$

where $C_t = \sum_i^N c_t^i$ represents the total consumption of consumers, and G_t is government spending. The supply, demand, and price represent the states of the market.

Economic Shocks

In Experiment 5, we simulate economic shocks analogous to a financial crisis: at the 100-th step, the wealth of all households is reduced by 50%. In our economic model, this scenario is mathematically represented as follows: for each household member, the wealth a_t^i at time t is updated according to the rule

$$a_t^i = 0.5a_{t-1}^i, \quad \forall i \in \{1, \dots, N\}$$

where N denotes the total number of household members.

G Hyperparameters

Hyperparameter	Value
Discount factor γ	0.975
Replay buffer size	1e6
Num of epochs	1000
Epoch length	300
Batch size	128
Adam epsilon	1e-5
Update cycles	100
Evaluation epochs	10
Hidden size	128
Tau	0.95
Critic initial learning rate	3e-4
Actor initial learning rate	3e-4
Learning rate adjustment	$0.95^{(\text{epoch}/35)}$

Table 5: Hyperparameters of SMFG methods and its variants.

Hyperparameter	Value
Noise rate	0.01
Epsilon start	0.1
Epsilon end	0.05
Epsilon decay	1e-5

Table 6: Hyperparameters of SMFG-MF algorithm different from SMFG method.

Ethical Statement

This research introduces a novel SMFG method, designed to optimize macroeconomic policies by modeling complex interactions at the micro level. The potential impact of this work extends across several domains:

Hyperparameter	Value
Tau τ	5e-3
Gamma γ	0.95
Eps ϵ	1e-5
Clip	0.1
Vloss coef	0.5
Ent coef	0.01
Government’s initial learning rate	3e-4
Learning rate adjustment	0, epoch < 10 $0.97^{(\text{epoch}/35)}$, epoch ≥ 10
Households’ initial learning rate	1e-6
Learning rate adjustment	$0.97^{(\text{epoch}/35)}$

Table 7: Hyperparameters of AI Economist Algorithm different from SMFG approach.

Academic Contributions The framework and algorithm proposed represent significant advancements in AI for economics and AI for social impact field, potentially serving as foundational tools for future research in macroeconomic policy making. By addressing the Lucas critique through dynamic modeling of individual agents within a mean field game, this work encourages more accurate and robust economic predictions and policy evaluations.

Policy Making and Societal Impact By enabling the optimization of macroeconomic policies through real-time, dynamic responses of micro-agents, this model provides policymakers with a powerful tool for assessing the impact of different economic strategies, leading to more informed decisions that maximize social welfare and economic stability, particularly in response to economic shocks. The application of this model can have profound implications for wealth distribution and social equity, helping ensure that economic policies are beneficial to a broader section of the population, potentially reducing inequality and enhancing societal well-being.

Ethical Considerations While the model aims to improve economic outcomes, the manipulation of macroeconomic policies must be approached with caution to avoid unintended negative consequences such as increased inequality or destabilization of economic sectors. Further, the reliance on AI-based decisions necessitates continuous scrutiny to ensure that the model accurately represents all population segments.

Limitations and Risks The complexity of the models also introduces risks related to the oversimplification of real-world dynamics and potential biases in the simulation of economic responses. Continuous validation against empirical data and diverse economic scenarios is essential to ensure the reliability and ethical application of the proposed methods.

In summary, the proposed SMFG framework and SMFRL algorithm hold the potential to significantly impact both academic research and practical policy making, offering a new perspective on dynamic economic modeling that prioritizes realistic, individual-level responses within large-scale economic systems.