

The AI Economist: Optimal Economic Policy Design via Two-level Deep Reinforcement Learning

Stephan Zheng^{*,†,1}, Alexander Trott^{*,1}, Sunil Srinivasa¹, David C. Parkes^{1,2}, and Richard Socher³

¹Salesforce Research

²Harvard University

³You.com

August 24, 2021

AI and reinforcement learning (RL) have improved many areas, but are not yet widely adopted in economic policy design, mechanism design, or economics at large. At the same time, current economic methodology is limited by a lack of counterfactual data, simplistic behavioral models, and limited opportunities to experiment with policies and evaluate behavioral responses. Here we show that machine-learning-based economic simulation is a powerful policy and mechanism design framework to overcome these limitations. The AI Economist is a two-level, deep RL framework that trains both agents and a social planner who co-adapt, providing a tractable solution to the highly unstable and novel two-level RL challenge. From a simple specification of an economy, we learn rational agent behaviors that adapt to learned planner policies and vice versa. We demonstrate the efficacy of the AI Economist on the problem of optimal taxation. In simple one-step economies, the AI Economist recovers the optimal tax policy of economic theory. In complex, dynamic economies, the AI Economist substantially improves both utilitarian social welfare and the trade-off between equality and productivity over baselines. It does so despite emergent tax-gaming strategies, while accounting for agent interactions and behavioral change more accurately than economic theory. These results demonstrate for the first time that two-level, deep RL can be used for understanding and as a complement to theory for economic design, unlocking a new computational learning-based approach to understanding economic policy.

*: equal contribution. †: Correspondence to: stephan.zheng@salesforce.com.

1 Introduction

Economic policies need to be optimized to tackle critical global socio-economic issues and achieve social objectives. For example, tax policy needs to balance equality and productivity, as large inequality gaps cause loss of economic opportunity (1) and adverse health effects (2). However, the problem of optimal policy design is very challenging, even when the policy objectives can be agreed upon.

Policy optimization poses a *mechanism design* (3) problem: the government (*social planner*) aims to find a policy under which the (boundedly) rational behaviors of affected economic agents yield the desired social outcome. Theoretical approaches to policy design are limited by analytical tractability and thus fail to capture the complexity of the real world. Empirical studies are challenged by the lack of counterfactual data and face the *Lucas critique* (4) that historical data do not capture behavioral responses to policy behavior. Furthermore, opportunities for rigorous, real-world experimentation are limited and come with ethical questions (5).

Computational and machine learning techniques for *automated* mechanism design (6–10) show promise towards overcoming existing limitations, but a general computational framework for policy design remains lacking. The challenge with policy design comes from needing to solve a highly non-stationary, *two-level*, sequential decision-making problem where all actors (the agents and the government) are *learning*: economic agents learn rational, utility-maximizing behaviors and the government learns to optimize its own objective via policy choices.

A New Machine Learning Challenge. Using deep reinforcement learning (RL) with multiple agents has been underexplored as a solution framework for mechanism design. Recent advances in deep RL have mostly studied the single-level setting; for example, state-of-the-art deep RL systems such as AlphaGo (11) and AlphaStar (12) optimized actors under fixed reward functions. In contrast, in the two-level setting agents’ effective reward functions depend on (changes in) the planner’s policy, which leads to a highly unstable learning and co-adaptation problem.

Significant advances in multi-agent RL have focused on cooperative problems (12, 13), and social dilemmas with fixed reward functions (14), but dynamical systems of heterogenous self-interested agents with changing incentives have been little studied at scale.

As such, few tractable computational learning approaches to mechanism design exist that scale to sequential settings with high-dimensional feature spaces. Consequently, machine learning has so far not been widely applied to economic policy design. In fact, more generally, economics as a field has not seen wide adoption of deep RL or related AI methods.

A.T. and S.Z. contributed equally. R.S. and S.Z. conceived and directed the project; S.Z., A.T., and D.P. developed the theoretical framework; A.T., S.S., and S.Z. developed the economic simulator, implemented the reinforcement learning platform, and performed experiments; A.T., S.Z., and D.P. processed and analyzed experiments with AI agents; S.Z., A.T., and D.P. drafted the manuscript; R.S. planned and advised the work, and analyzed all results; All authors discussed the results and commented on the manuscript.

We only consider rational behaviors in this work, although our framework can be extended to include boundedly rational actors.

The AI Economist. Here we introduce the *AI Economist*, a new and powerful framework that combines machine learning and AI-driven economic simulations to overcome the limitations faced by existing approaches. Specifically, the AI Economist shows the efficacy and viability of using 1) AI-driven economic simulations and 2) two-level RL as a new paradigm for economic policy design.

AI-driven Simulations. We show that AI-driven simulations capture features of real-world economies without the need for hand-crafted behavioral rules or simplifications to ensure analytic tractability. We use both a single-step economy and a multi-step, micro-founded economic simulation, *Gather-Trade-Build*. Gather-Trade-Build features multiple heterogeneous economic agents in a two-dimensional spatial environment. Productivity and income elasticity emerge as the result of the strategic behavior of multiple agents, rather than from statistical assumptions. Moreover, Gather-Trade-Build includes trading between agents and simulates the economy over extended periods of time, i.e., spanning 10 tax periods, each of 100 days. As such, the dynamics of Gather-Trade-Build are more complex than those considered in traditional tax frameworks and serve as a rich testbed for AI-driven policy design.

AI-driven Policy Design with Two-level, Deep RL. The AI Economist uses two-level, deep RL to learn optimal policies: at the level of individual agents within the economy and at the level of the social planner. Both the agents and the social planner use deep neural networks to implement their policy model. Two-level RL compares the performance of billions of economic designs, making use of agents whose behaviors are learned along with the optimal planner policy.

Two-level RL is natural in many contexts, e.g., mechanism design, the principal-agent problem, or regulating systems with (adversarial) agents with misaligned or unethical incentives. However, it poses a highly unstable learning problem, as agents need to continuously adapt to changing incentives. The AI Economist solves the two-level problem through the use of *learning curricula* (15) and *entropy-based regularization* (16), providing a tractable and scalable solution. Our approach stabilizes training using two key insights: (1) agents should not face significant utility costs that discourage exploration early during learning, and (2) the agents and social planner should be encouraged to gradually explore and co-adapt.

The AI Economist framework provides numerous advantages.

- It does not suffer from the Lucas critique. By design, it considers actors who co-adapt with economic policy.
- Nor does it suffer the problems from using simulated agents with *ad hoc* behavioral rules; rather, the use of RL provides rational agent behavior.
- The simulation framework is flexible, supporting a configurable number of agents and various choices in regard to economic processes.

- The designer is free to choose any policy objective and this does not have to be analytically tractable or differentiable.
- The use of RL requires only observational data and does not require prior knowledge about the simulation or economic theory.

Optimal Tax Policy. We demonstrate the efficacy of the AI Economist on the problem of *optimal tax policy design* (17–19), which aims to improve social welfare objectives, for example finding the right balance of equality and productivity. In brief, tax revenue can be used to redistribute wealth, invest in infrastructure, or fund social programs. At the same time, tax rates that are too high may disincentivize work and elicit strategic responses by tax-payers.

Theory-driven approaches to tax policy design have needed to make simplifications in the interest of analytical tractability (20). For example, typical models use static, one-step economies (21, 22) and make use of assumptions about people’s sensitivity to tax changes (elasticity). Although work in *New Dynamic Public Finance (NDPF)* (23, 24) seeks to model multi-step economies, these models quickly become intractable to study analytically. Concrete results are only available for two-step economies (25). These theoretical models also lack interactions between agents, such as market-based trading, and consider simple, inter-temporal dynamics.

Previous simulation work that makes use of agent-based modeling (ABM) (26–31) avoids problems of analytical tractability but uses complex and *ad hoc* behavioral rules to study emergent behavior, this complicating the interpretation of results. Moreover, the behavior of ABM agents is often rigid and lacking in strategic or adaptive behavior.

Experimental Validation. We provide extensive proof that the AI Economist provides a sound, effective, and viable approach to understanding, evaluating, and designing economic policy design. We study optimal tax design in a single-step economy and the multi-step Gather-Trade-Build environment, which implements a dynamic economy of heterogeneous, interacting agents that is more complex than the economic environments assumed in state-of-the-art tax models. We show that the use of RL yields emergent agent behaviors that align well with economic intuition, such as specialization and tax gaming, phenomena that are not captured through analytical approaches to tax policy design. This happens even with a small number of agents (4 and 10 agents in our experiments).

We show that policy models using two-level RL are effective, flexible, and robust to strategic agent behaviors through substantial quantitative and qualitative results:

- In one-step economies, the AI Economist recovers the theoretically optimal tax policy derived by Saez (21). This demonstrates the use of two-level RL is sound.
- In Gather-Trade-Build economies, tax policies discovered by the AI Economist provide a substantial improvement in social welfare for two different definitions of social welfare and

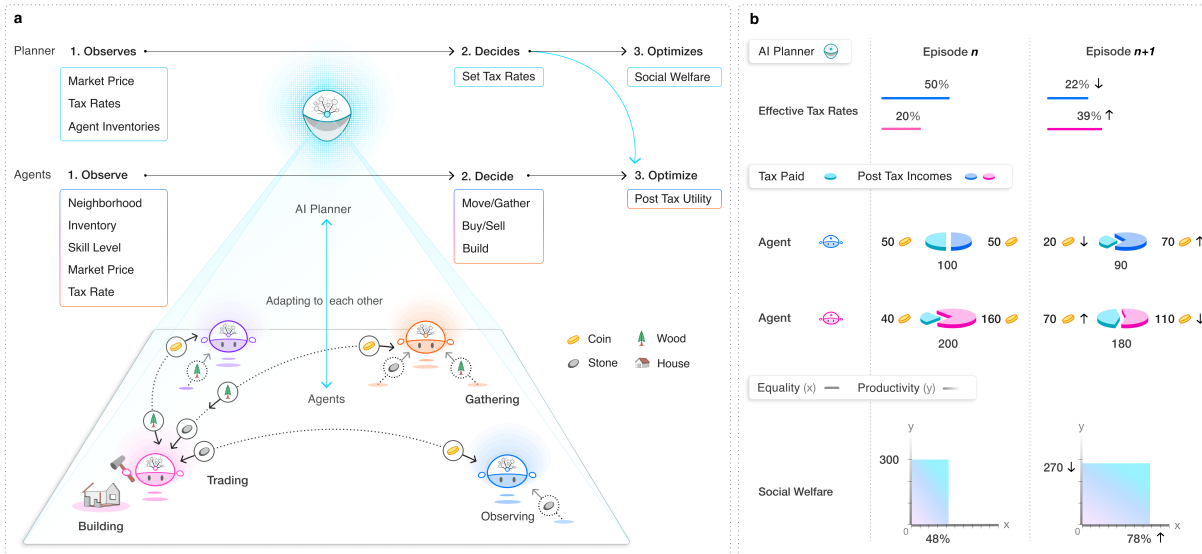


Figure 1: AI-driven economic simulations and two-level reinforcement learning (RL). **a**, An AI social planner optimizes social welfare by setting income tax rates in an economic simulation with AI agents. The agents optimize their individual post-tax utility by deciding how to perform labor and earn income. Both the planner and agents use RL to co-adapt and optimize their behavior. Agents need to optimize behavior in a non-stationary environment, as the planner’s tax decisions change the reward that agents experience. **b**, Illustration of co-adaptation and two-level learning in an economy with two agents. Simulations proceed in episodes that last for 10 tax years, with 100 timesteps in each simulated year. During learning, between any episodes n and $n + 1$, the planner changes tax rates, which, after behavioral changes, leads to higher social welfare, here defined as the product of productivity and equality.

in various spatial world layouts; e.g., in the Open-Quadrant world with four agents, *utilitarian social welfare* increases by 8%, and the *trade-off between equality and productivity* increases by 12% over the prominent Saez tax framework (21).

- In particular, AI social planners improve social welfare despite strategic behavior by AI agents seeking to lower their tax burden.
- AI-driven tax policies improve social welfare by using different kinds of tax schedules than baseline policies from economic theory. This demonstrates that analytical methods fail to account for all of the relevant aspects of an economy, while AI techniques do not require simplifying assumptions.
- Our work gives new economic insights: it shows that the well-established Saez tax model, while optimal in a static economy, is suboptimal in more realistic dynamic economies where it fails to account for interactions between agents. Our framework enables us to precisely quantify behavioral responses and agent interactions.

Ethical Disclaimer. As a point of caution, while the Gather-Trade-Build environments provide a rich testbed for demonstrating the potential of AI-driven simulation, they do not articulate the full range of economic opportunities, costs, and decisions faced by *real-world individuals*, nor their distribution of relevant attributes. More realistic AI-driven simulations are needed to support real-world policymaking, and defining the criteria for sufficient realism will require widespread consultation. By extension, any conclusions drawn from experiments in these environments face the same limitations and, therefore, are not meant to be applied to any specific real-world economies. See Section 9 for an extensive discussion on ethical risk.

2 AI-driven Economic Simulations

The AI Economist framework applies RL in two key ways: (1) to describe how *rational* agents respond to alternative policy choices, and (2) to *optimize* these policy choices in a principled economic simulation. Specifically, economic simulations need to capture the relevant economic drivers that define rational behavior. As such, a key strength of this framework is that finding rational behaviors along with an optimal policy remains tractable even with complex specifications of economic incentives and dynamics.

Simulation Dynamics. We apply the AI Economist to the problem of optimal taxation (Figure 1). The set-up follows the Mirrleesian framework of non-linear optimal taxation subject to incentive constraints (18). Here, the incentive constraints are represented through the rational behavior of agents, who optimize behavior subject to income tax and income redistribution.

Our simulations run for a finite number of timesteps H and capture several key features of the Mirrleesian framework: that agents perform *labor* l in order to earn *income* z , where *skill* determines how much income an agent earns for a given amount of labor; that an agent’s utility increases with its *post-tax* income and decreases with its labor; and that agents are heterogeneously skilled.

The simulation captures these concepts through its dynamics, i.e. the actions available to the actors and how those actions a_t influence the world state s_t at timestep t . For example, agents may move spatially to collect resources, trade with one another, or spend resources to build houses; each such action accrues *labor* but may generate *income*, with higher *skill* ν leading to higher incomes for the same actions.

Taxation. Agents pay *taxes* on the income they earn according to a tax schedule $T(z, \tau)$, which determines taxes owed as a function of income and a set of *bracketed marginal tax rates* τ . The planner controls these tax rates, with all agents facing the same tax schedule, where this schedule can change at the start of each tax year. Collected taxes are evenly redistributed back to agents. For simplicity, we use fixed bracket intervals, and the planner only sets the marginal rates.

Behavioral Models. Each actor (whether agent or planner) uses a deep neural network to encode its behavior as a probability distribution $\pi(a_t|o_t)$ over available actions, given observation o_t . Following economic theory, each actor observes only a portion of the full world state s_t . For instance, the planner can observe trade activity but not an agent’s skill level. Actors’ objectives, i.e. post-tax utility for agents and social welfare for the planner, are captured in the *reward function* used to train each behavioral policy π . In this way, the AI Economist uses RL to describe rational agent behavior and optimize policy choices in complex, sequential economies beyond the reach of traditional analysis.

3 Two-Level Reinforcement Learning

Under the AI Economist framework, all actors (i.e. the AI agents and the AI planner) learn and adapt using RL (33), see Algorithm 1. Each actor learns a behavioral policy π to maximize its objective (expected sum of future rewards). Each actor also learns a *value function*, which estimates this expectation given observation o_t . Actors iteratively *explore* actions by sampling from their current behavioral model, and improve this model across episodes by training on experiential data. RL agents can be optimized for any reward function and this does not have to be analytical.

An agent i maximizes *expected total discounted utility*:

$$\max_{\pi_i} \mathbb{E}_{a_i \sim \pi_i, a_{-i} \sim \pi_{-i}, s' \sim \mathcal{P}} \left[\sum_{t=1}^H \gamma^t r_{i,t} + u_{i,0} \middle| \tau \right], \quad r_{i,t} = u_{i,t} - u_{i,t-1}, \quad (1)$$

given tax rates τ , discount factor γ , and utility $u_{i,t}$. Here s' is the state following s , and \mathcal{P} represents the simulation dynamics. We use *isoelastic utility* (34):

$$u_{i,t} = \frac{C_{i,t}^{1-\eta} - 1}{1-\eta} - L_{i,t}, \quad \eta > 0, \quad (2)$$

which models diminishing marginal utility over money endowment $C_{i,t}$, controlled by $\eta > 0$, and the linear disutility of total labor $L_{i,t}$. The planner maximizes expected social welfare:

$$\max_{\pi_p} \mathbb{E}_{\tau \sim \pi_p, a \sim \pi, s' \sim \mathcal{P}} \left[\sum_{t=1}^H \gamma^t r_{p,t} + \text{swf}_0 \right], \quad r_{p,t} = \text{swf}_t - \text{swf}_{t-1}, \quad (3)$$

where swf_t is social welfare at time t . We take swf as a utilitarian objective (an average of all agent utilities weighted by their inverse pre-tax income), or alternatively as the product of equality and productivity (representing a balance between equality and productivity). For details, see Methods.

Agents need to adapt to policies set by the planner, and vice versa (Figure 1). This is a challenging non-stationary learning problem. While learning, the planner in effect adjusts agent

reward functions because taxes influence the post-tax income that agents receive as a result of payments and redistributions. As the tax schedule changes, the optimal behavior for agents changes. This instability is exacerbated by mutual exploration.

These challenging learning dynamics reflect the nested optimization problem that two-level RL attempts to solve. That is, we aim to find the tax rates that maximize social welfare, subject to the constraint that agents’ behaviors maximize their own utility given the tax rates. Planner learning (the outer level) serves to maximizing social welfare, whereas agent learning (the inner level) serves to ensure that the constraint is satisfied. Our approach to two-level RL follows from the intuition that instability depends on how well the agent-optimality constraint is satisfied during learning.

To stabilize learning, our approach combines two key ideas: *curriculum learning* (15) and *entropy regularization* (16). This effectively *stagger*s agent and planner learning such that agents are well-adapted to a wide range of tax settings before the learning of the planner begins. In particular, we use the early portion of training to gradually introduce labor costs and, later, taxes. These curricula are based on the key intuition that suboptimal agent strategies may incur punitively high cost of labor and taxes, while earning insufficient income to yield positive utility, and this may discourage RL agents from continuing to learn. We schedule the entropy regularization applied to π_p such that agents are initially exposed to highly random taxes. Random taxes provide the training experience needed for agent policies to appropriately condition actions on the observed tax rates, for a wide range of possible taxes. As described above, this is an important precondition for stably introducing planner optimization. Lastly, the entropy of policy models are strongly regularized to encourage exploration and gradual co-adaptation between the agents and social planner throughout the remainder of training. For details, see Methods.

We note that unlike previous strategies for overcoming instability in multi-agent RL (35, 36), ours is tailored to the nested optimization intrinsic to the two-level setting.

4 Validation in a One-Step Economy

The most prominent solution for optimal taxation is the analytical framework developed by Saez (21). This framework analyzes a simplified model where both the planner and the agents each make a single decision: the planner setting taxes and the agents choosing labor. This analysis describes the welfare impact of a tax rate change via its mechanical effect on redistribution and its behavioral effect on the underlying income distribution. The resulting formula computes theoretically optimal tax rates as a function of the income distribution and the *elasticity* of income with respect to the marginal tax rate.

We first validate our approach in these simplified one-step economies. Each agent chooses an amount of labor that optimizes its post-tax utility, and this optimal labor depends on its skill and the tax rates, and it does not depend on the labor choices of other agents. Before the agents

In practice, these income elasticities typically need to be estimated from empirical data, which is a non-trivial task (37).

act, the planner sets the marginal tax rates in order to optimize social welfare, taken here to be utilitarian.

We compare the economy under the *Saez tax*, and the AI Economist. In both cases, AI agents learn to optimize their own utility given their tax setting. The Saez tax baseline computes tax rates based on our implementation of the Saez formula (induced through an optimal elasticity parameter found via grid-search, as detailed in the Methods), and the AI Economist learns tax rates via two-level RL. We include two additional baseline tax models here and throughout this work: the *free market* (no taxes) and a stylized version of the *US Federal* progressive tax schedule (see Methods for details). There is no *a priori* expectation that either of the additional baselines should maximize social welfare; rather, they provide useful comparison and help to characterize behavioral responses to different tax choices. The AI Economist and the Saez tax schedule produce highly consistent tax schedules and social welfare, as shown in Figure 3a-b. In comparison, the free market and US Federal achieve substantially worse social welfare. These results show that the AI Economist can reproduce optimal tax rates in economies that satisfy the simplifying assumptions of optimal tax theory and validate the soundness of our learning-based approach.

5 Gather-Trade-Build: a Dynamic Economy

We study the *Gather-Trade-Build* economy, a two-dimensional, spatiotemporal economy with agents who move, gather resources (stone and wood), trade, and build houses. Gather-Trade-Build captures the fundamental trade-off between equality and productivity intrinsic to optimal tax design (see below), and is a rich testbed to demonstrate the advantages of AI-driven policy design.

Each simulation simulates 10 tax years. Each tax year lasts 100 timesteps (so that $H = 1000$), with the agents acting each timestep, and the planner setting and changing tax rates at the start of each tax year. The Gather-Trade-Build environment is depicted in Figure 1. For details, see Methods.

AI-driven Simulations Capture Macro-Economic Phenomena. A key advantage is that AI-driven simulations capture macro-level features of real economies that are emergent purely through learned rational behavior and without being manually implemented. To illustrate this, we showcase three examples of AI-driven emergent behavior.

Example 1: Emergent Specialization. Each agent varies in its skill level. We instantiate this in our simulation as *build-skill*, which sets how much income an agent receives from building a house. Build-skill is distributed according to a Pareto distribution. As a result, we observe that utility-maximizing agents *learn to specialize* their behavior based on their build-skill, see Figure 2. Agents with low build-skill become *gatherers*: they earn income by gathering and selling resources. Agents with high build-skill become *builders*: they learn that it is more profitable to

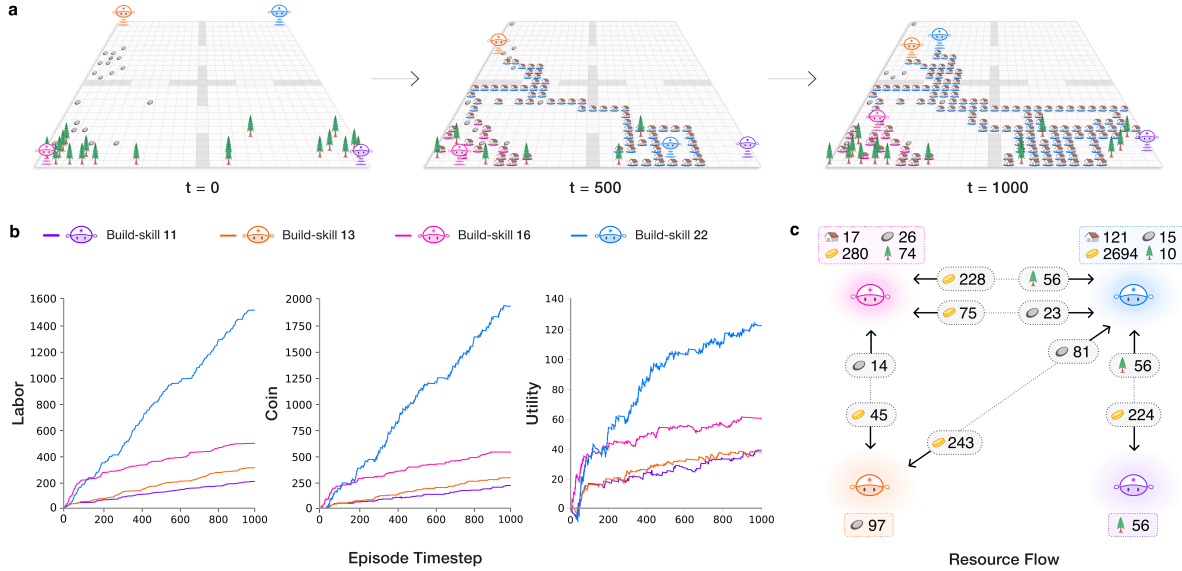


Figure 2: **Emergent phenomena in AI-driven economic simulations under the free market.** **a**, Visualization of the spatial state of the world at $t = 0$, 500, and 1000 of an example episode in the 4-agent Open-Quadrant Gather-Trade-Build scenario. Agents specialize as *builders* (blue agent) or *gatherers* (others) depending on their build-skill. **b**, Labor, income, and utility over the course of the episode for all agents. Each quantity increases with build-skill in this setting. The highest build-skill (blue) agent chooses to do the most work, and earns larger income and ultimately experience the most utility. **c**, Net resource flow between agents during the episode. The box adjacent to each agent show the resources it gathered and the coin it earned from building. Arrows between agents denote coin and resources exchanged through trading.

buy resources and then build houses. This emergent behavior is entirely due to *heterogeneous* their experienced utility for different economic activity, and not due to fixed behavioral rules as in most traditional agent-based modeling.

Example 2: Equality-Productivity Trade-off. Our AI simulations capture the trade-off between equality and productivity: as tax rates increase, equality increases through wealth transfers, but productivity falls as agents are less incentivized to work due to lower post-tax incomes (Figure 3 and Figure 4). As a demonstration of this, the free market (no tax) baseline always yields the highest productivity and lowest equality compared to the alternative tax models. Unlike standard theoretical models that rely on elasticity assumptions to capture this trade-off, we observe it as an emergent consequence of rational behavior.

Example 3: AI Tax Gaming Strategies. Our AI simulations yield emergent strategic behaviors. High-income agents learn to avoid taxes by moving labor and thus income between tax

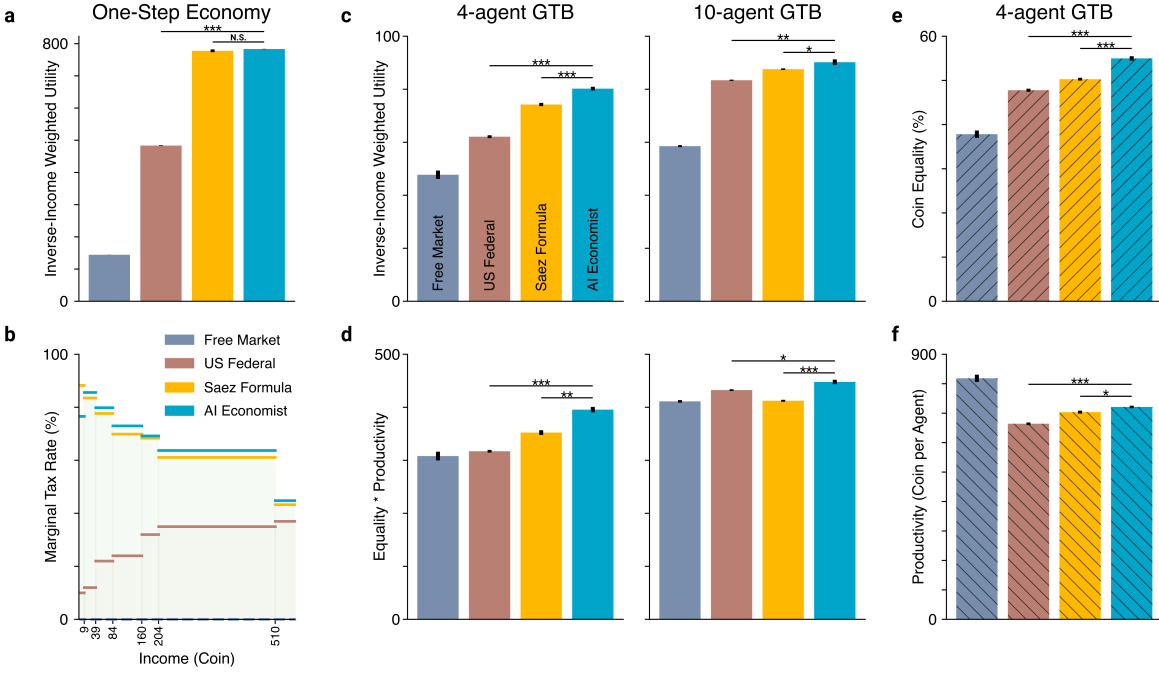


Figure 3: **Quantitative results in a one-step economy and the Open-Quadrant Gather-Trade-Build environment.** **a-b**, The results of the AI Economist and the Saez tax are highly consistent in the one-step economy, both in terms of utilitarian social welfare (a) and the tax schedule (b). **c-d**, In the Gather-Trade-Build environment (GTB) with 4 and 10 agents, the AI Economist outperforms baselines when optimizing the utilitarian social welfare objective (c) and when optimizing the equality-times-productivity objective (d). **e-f**, Overall coin equality (e) and average productivity (f) achieved by each tax model in the 4-agent Open Quadrant scenario. Each bar represents the average end-of-training metrics over 10 random seeds (5 for the one-step economy), with error bars denoting standard error. Asterisks indicate a statistically significant difference at an α level of 0.05 (*), 0.001 (**), or 0.00001 (***). N.S. denotes not statistically significant ($p > 0.05$). All social welfare, productivity, and equality differences between the AI Economist and baselines are statistically significant, except for the difference in social welfare between the AI Economist and the Saez tax in the one-step economy (a).

years in order to move more income to low-rate brackets. This can reduce the overall tax paid in comparison to earning a constant amount each year (Figure 6c). Given the complexity of Gather-Trade-Build and similar dynamic economic environments, it is prohibitively complex for theory-driven methods to derive such temporal behavioral strategies.

Together, these examples show that AI-driven simulations capture features of real-world economies, purely through RL. Hence, AI-driven simulations provide a rich class of environments for policy design, unconstrained by analytic tractability.

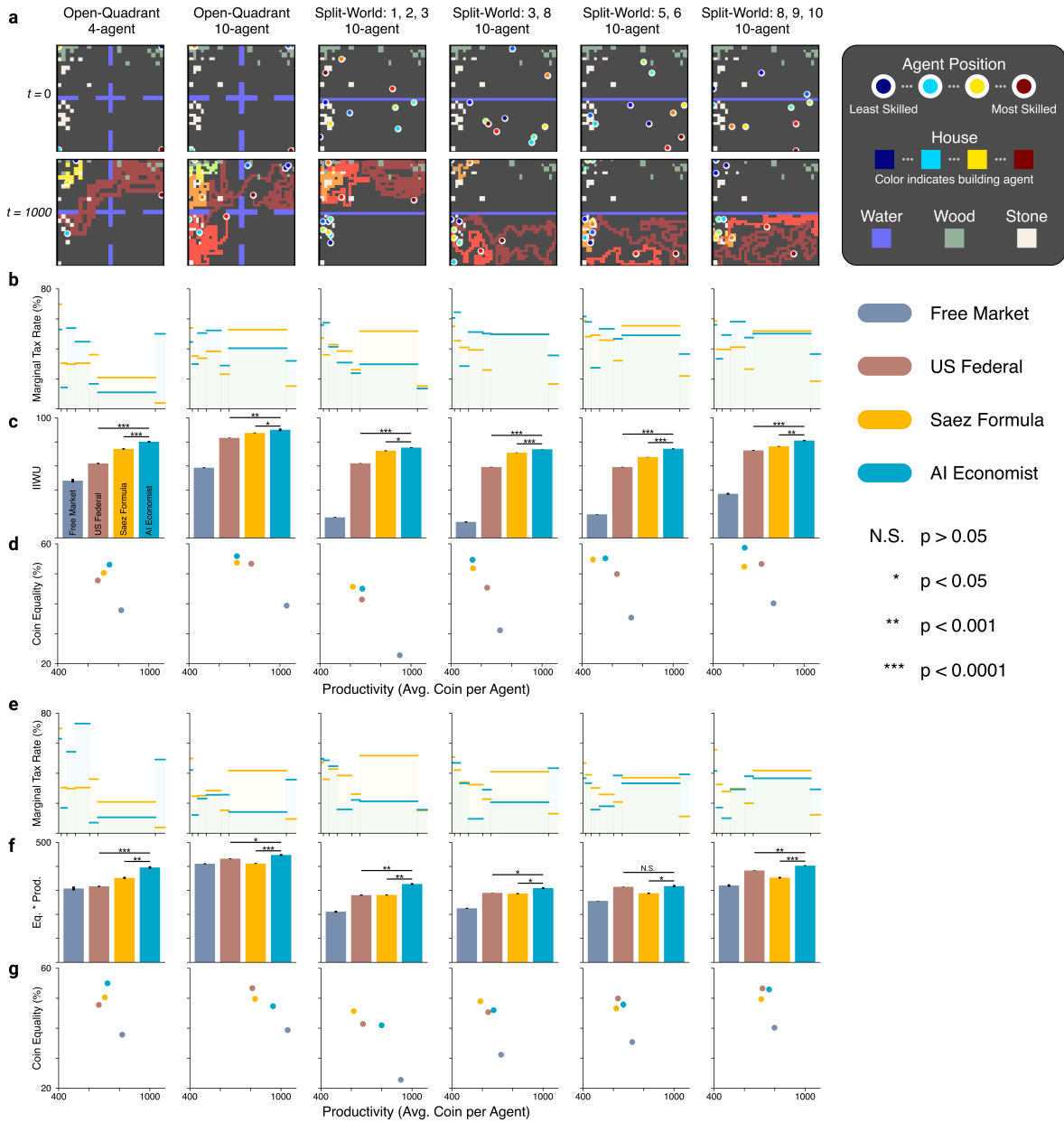


Figure 4: Description on the next page.

6 AI-Driven Optimal Taxation

We evaluate the AI Economist across different Gather-Trade-Build economies to validate that AI-driven policy design is effective, can be applied to different economic environments, and adapts to strategic behavior more successfully than baseline tax policies.

Figure 4: **Comprehensive quantitative results in the Gather-Trade-Build environment with the utilitarian or equality-times-productivity planner objective, across all settings: Open-Quadrant and 4 Split-World scenarios; 4 and 10 agents.** The AI Economist achieves significantly higher social welfare than all baselines. **a**, Spatial layouts of the Open-Quadrant and Split-World scenarios at the start ($t = 0$) and end ($t = 1000$) of example episodes. **b**, Tax schedules for the Saez tax (yellow) and the AI Economist (blue). **c**, Utilitarian social welfare objective (inverse-income weighted utility, labeled “IIWU”) for all planners. **d**, Equality and productivity for all planners. For the data in b-d, the AI Economist is trained to maximize the utilitarian social welfare objective, and the Saez taxes use the best-performing elasticity for the utilitarian objective. **e-g**, As b-d, but for the data in e-g the AI Economist is trained to maximize the equality-times-productivity social welfare objective, and the Saez taxes use the best-performing elasticity for this objective, which is shown in f. Bars and dots represent the average end-of-training metrics over 10 (5) random seeds for the Open-Quadrant (Split-World) scenarios, with error bars denoting standard error. Asterisks indicate a statistically significant difference at an α level of 0.05 (*), 0.001 (**), or 0.00001 (***). N.S. denotes not statistically significant ($p > 0.05$). All social welfare differences between the AI Economist and baselines are statistically significant, except for the difference in equality-times-productivity (f) between the AI Economist and the US Federal tax in the *Split-World-5,6* scenario.

Settings. We use two spatial layouts, *Open-Quadrant* and *Split-World*, each with different physical barrier placements and different agent starting positions. *Open-Quadrant* features four areas laid out in a 2×2 pattern, each area having a connection with its neighbor to allow agents to move between areas. *Split-World* features two halves, separated by an impassable water barrier. This prevents agents from moving between the top and bottom halves of the map, which blocks agents from directly accessing certain resources.

We consider four Split-World scenarios, each with 10 agents but differing in the subset of agents assigned to the resource-rich half. We consider two Open-Quadrant scenarios, with 4 agents in one version and 10 agents in the other. All 6 scenarios are illustrated in Figure 4a. For ease of exposition, we focus our fine-grained analyses on results in the 4-agent Open-Quadrant scenario.

Improved Social Welfare. As with the one-step economy, we compare the AI Economist against the *free market*, *US Federal*, and *Saez tax* baselines across all of these settings (see Methods). The AI Economist achieves the highest social welfare throughout. The combined results of these experiments are presented in Figure 4. In the *Open-Quadrant* layout with four (ten) agents (Figure 3), AI-driven taxes improve the utilitarian objective by over 8% (2%) and the product of equality and productivity by over 12% (8.6%) over the Saez tax.

We observe that the relative performance of the baselines depends on the choice of social welfare objective: the utilitarian objective is always higher when using the Saez tax compared to the US Federal tax; however, the opposite is often true for the equality-times-productivity

objective (especially in settings with 10 agents). In contrast, the AI Economist is not tailored towards a particular definition of social welfare and flexibly adjusts its tax schedule to optimize the chosen objective, yielding the best social welfare throughout.

These results show the AI Economist is flexible, maintains performance with more agents, can be successfully optimized for distinct objectives, and works well in the face of adaptive, strategic behavior.

Adaptation During Training. During training, the AI Economist increases rates on the first (incomes of 0 to 9), third (39 to 84), and fourth (84 to 160) brackets, maintaining low rates otherwise, see Figure 5. This does not significantly shift the pre-tax income distribution, while the post-tax income distribution becomes more equal. The resulting tax schedule is distinctly different from the baselines, which use either increasing (progressive) or decreasing (regressive) schedules (Figure 5a). The AI Economist is neither: on average, it sets the highest marginal rates for incomes between 39 and 160 coins and the lowest rates for the adjacent brackets (9 to 39 and 160 to 510 coins). Under the AI Economist, the low build-skill agents earn 9% more from trading (Figure 6b), wealth transfers from the highest build-skill agent to others are 46% larger (Figure 5d), income equality is at least 9% higher (Figure 3e), the number of incomes in the second-to-highest bracket (204 to 510 coins) is at least 64% higher, and, 92% smaller for the top bracket, compared to baselines (Figure 5b). These numbers are measured over the last 400 episodes within each experiment group, which amounts to 4000 total tax periods and 16000 total incomes per group.

Behavior of Learned AI Tax Policies. The AI Economist adapts to different environments: Figure 4 shows that the best-performing AI taxes behave differently across scenarios.

For instance, in the Open-Quadrant, the AI tax schedules are similar when optimizing for the two different social welfare objectives with 4 agents but this pattern changes with 10 agents, where objective-specific tax schedules emerge. Tax rates for the brackets between 9 and 160 coins follow different patterns, for example, and overall tax rates are lower when optimizing for equality times productivity.

Furthermore, in the Split-World, the AI tax schedule depends on which agents are in the resource-rich top half of the environment. As an example, when optimizing for equality times productivity, when the two agents with the highest build-skill (Agents 1, 2) are (not) in the top half, taxes in the 204 to 510 bracket are lower (higher) than those in the 0 to 84 range.

Owing to the complexity of these environments, it is not possible to provide an intuitive explanation of these AI tax schedules. Nevertheless, it is not surprising that differences between scenarios reflect in the optimal tax rates, as the various combinations of skill and resource access promote difference economic forces and resulting equilibria. Such is demonstrated even in the range of free market social outcomes across these scenarios (Figure 4d,g). Considering that the AI tax schedules maximize social welfare within their respective scenarios, we view their scenario-specific idiosyncrasies as evidence of the adaptability of the AI Economist framework.

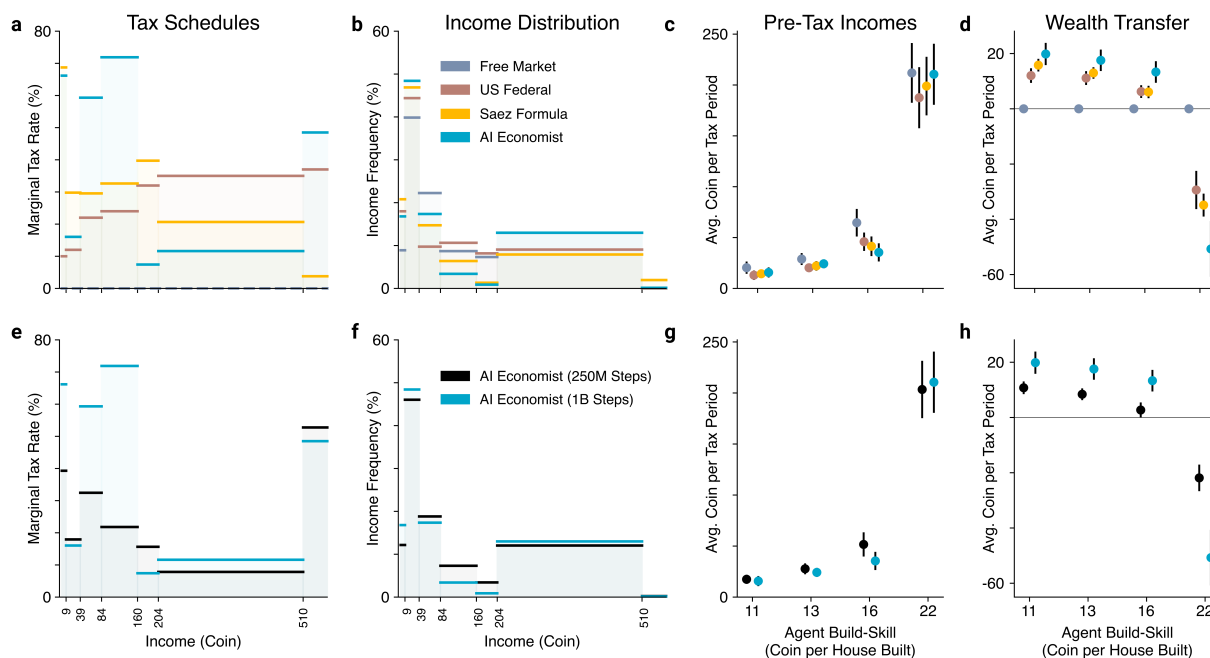


Figure 5: **Comparison of tax policies in the 4-agent Open-Quadrant Gather-Trade-Build environment.** **a**, Average marginal tax rates within each tax bracket. **b**, Frequency with which agent incomes fall within each bracket. **c**, Average pre-tax income of each agent (sorted by build-skill) under each of the tax models. **d**, Average wealth transfer resulting from taxation and redistribution. **e-h**, Same as a-d, comparing the AI Economist from early during training (250 million training samples) versus at the end of training (1 billion training samples). Dots denote averages and error bars denote standard deviation across episodes.

7 Policy Design Beyond Independence Assumptions

Micro-founded AI-driven simulations such as Gather-Trade-Build enable optimal tax policy design in multi-step economies with coupled agent behaviors and interactions, through two-level RL. In contrast, analytical solutions are not available for these kinds of environments: traditional methods fail to account for interactions and thus only achieve suboptimal social welfare.

To illustrate the effect of interactions, Figure 6a-b shows that the income of the two agents with the lowest build-skill depends on the second-to-highest bracket tax rate, even though this income bracket only directly applies to the agent with the highest build-skill. As this tax rate increases, the agent with the highest build-skill buys fewer resources. In turn, the average resource price as well as the trade volume decreases, reducing the incomes of the low build-skill agents. Hence, a behavioral change of one agent can change the optimal policy of another agent.

However, the Saez analysis uses assumptions and a standard definition of elasticity that fail to account for interactions that arise in multi-step (real-world) economies, these interactions arising through trading for example. The Saez analysis assumes that behavioral changes of agents are

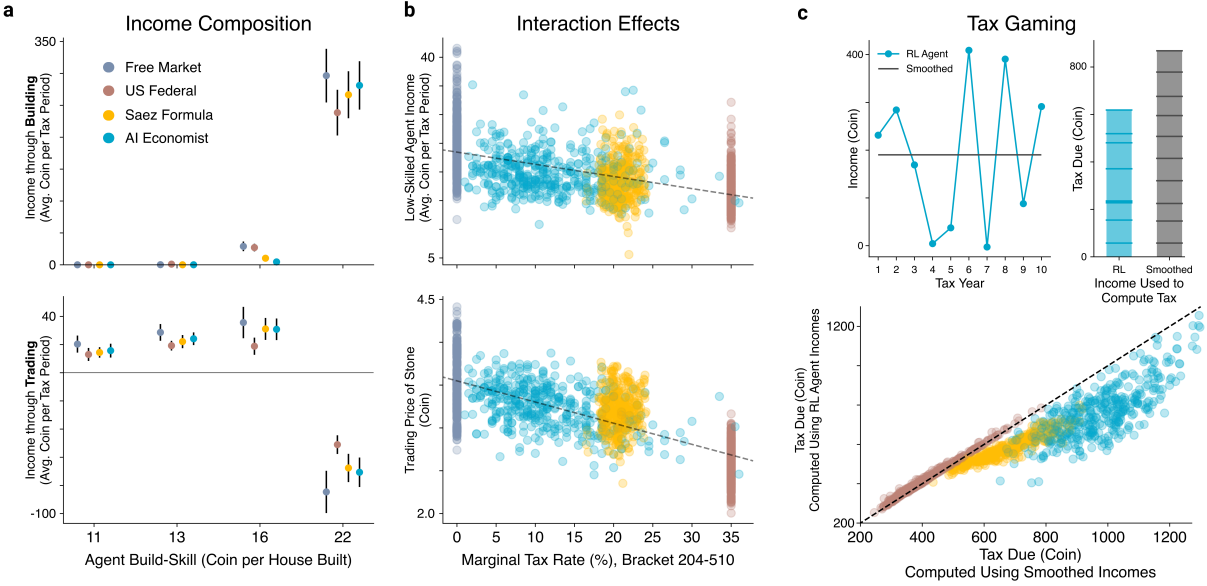


Figure 6: Specialization, interactions, and tax gaming in the 4-agent Open-Quadrant Gather-Trade-Build environment. **a**, Average net income from building (a, top) and trading (a, bottom) of each agent. Negative values denote net expenditure. **b**, The income of the two lowest build-skill agents (b, top) and average trading price (b, bottom) decrease as the tax rate in the higher 204-510 tax bracket increases, even though the agents’ incomes are below the cutoff for this bracket. Hence, the trading behavior of high-skilled agents affects the income of the low-skilled agents. The standard definition of elasticity does not capture this interaction effect. **c**, RL agents learn to strategize across each of the 10 tax years, lowering their total payable tax compared to a smoothed strategy that earns the same, average income in each year: the top panels illustrate this for a single episode; the bottom panel shows the saving relative to a smoothed income across all episodes used in the analysis. We do not observe this tax gaming under the progressive US Federal tax schedule.

independent and do not affect each other. This limitation results in suboptimal policy and lost social welfare under the Saez tax, when applied to the Gather-Trade-Build environment.

To illustrate this, for the four agent, *Open Quadrant* scenario, a typical regression of observed taxes paid and reported incomes would estimate elasticity at around 0.87, see Methods for details. However, by evaluating the Saez tax over a wide range of elasticity values, we find that an assumed elasticity of around 3 optimizes social welfare when used in Saez’s framework. This mismatch between offline estimates and imputed optimal values for agent elasticity is in significant part due to interactions between agents.

8 Discussion

The AI Economist demonstrates for the first time that economic policy design using RL, together with principled economic simulation, is sound, viable, flexible, and effective. It suggests an exciting research agenda: using AI to enable a new approach to economic design. The AI Economist framework can be used to study different policy goals and constraints, and, as AI-driven simulations grow in sophistication, may help to address the modern economic divide. In particular, AI-driven simulations enable economic policies to be tested in more realistic environments than those available to analytical methods, and show promise in validating assumptions in policy proposals and evaluating ideas coming from economic theory.

However, these results are a first step and are not ready to be implemented as real-world policy. Future research should scale up AI-driven simulations and calibrate them to real-world data, along with learning AI policies that are explainable and robust to simulation-to-reality gaps. Also, designing simulations to incorporate different societal values and be representative of different parts of society will be an important direction for future work.

AI-driven policy design could democratize policymaking, for instance, through easily accessible open-source code releases that enable a broad multidisciplinary audience to inspect, debate, and build future policymaking frameworks. As such, we hope the potential of AI-driven policy design will motivate building fair and inclusive data, computation, and governance structures that ultimately improve the social good.

9 Ethics

While the current version of the AI Economist provides only a limited representation of the real world, we recognize that it could be possible to manipulate future, large-scale iterations of the AI Economist to increase inequality and hide this action behind the results of an AI system.

Furthermore, either out of ignorance or malice, bad training data may result in biased policy recommendations, particularly in cases where users will train the tool using their own data. For instance, the under-representation of communities and segments of the work-force in training data might lead to bias in AI-driven tax models. This work also opens up the possibility of using richer, observational data to set individual taxation, an area where we anticipate a strong need for robust debate.

Economic simulation enables studying a wide range of economic incentives and their consequences, including models of stakeholder capitalism. However, the simulation used in this work is not an actual tool that can be currently used with malintent to reconfigure tax policy. We encourage anyone utilizing the AI Economist to publish a model card and data sheet that describes the ethical considerations of trained AI-driven tax models to increase transparency, and by extension, trust, in the system. Furthermore, we believe any future application or policy built on economic simulations should be built on inspectable code and subject to full transparency.

In order to responsibly publish this research, we have taken the following measures:

- To ensure accountability on our part, we have consulted academic experts on safe release of code and ensured we are in compliance with their guidance. We shared the paper and an assessment of the ethical risks, mitigation strategies, and assessment of safety to publish with the following external reviewers: Dr. Simon Chesterman, Provost’s Chair and Dean of the National University of Singapore Faculty of Law, and Lofred Madzou, AI Project Lead at the World Economic Forum’s Center for the Fourth Industrial Revolution. None of the reviewers identified additional ethical concerns or mitigation strategies that should be employed. All affirmed that the research is safe to publish.
- To increase transparency, we are also publishing a summary of this work as a blog post, thereby allowing robust debate and broad multidisciplinary discussion of our work.
- To further promote transparency, we will release an open-source version of our environment and sample training code for the simulation. This does not prevent future misuse, but we believe, at the current level of fidelity, transparency is key to promote grounded discussion and future research.

With these mitigation strategies and other considerations in place, we believe this research is safe to publish. Furthermore, this research was not conducted with any corporate or commercial applications in mind.

10 Methods

Dataset Use and Availability. No independent, third-party datasets were used in this work. All results were obtained through the use of simulation. The data and code used to visualize the results is available for all Figures, specifically:

- [Figure 2](#)
- [Figure 3](#)
- [Figure 5](#)
- [Figure 6](#)
- [Figure 4](#)
- [Figure 8](#)

Code Availability. All code for the economic simulations, reinforcement learning algorithms, and analysis are available upon request from the corresponding author.

One-Step Economy. We trained the AI Economist in a stylized, one-step economy with $N = 100$ agents, indexed by i , that each choose how many hours of labor l_i to perform. Each agent i has a *skill level* ν_i , which is a private value that represents its hourly wage. Based on labor, each agent i earns a pre-tax income $z_i = l_i \cdot \nu_i$. Each agent i also pays income tax $T(z_i)$ which is evenly redistributed back to the agents. As such, the post-tax income is defined as $\tilde{z}_i = z_i - T(z_i) + \frac{1}{N} \sum_{j=1}^N T(z_j)$. As a result, each agent i experiences a *utility* $u(\tilde{z}_i, l_i) = \tilde{z}_i - c \cdot l_i^\delta$, which increases linearly with post-tax income \tilde{z}_i and decreases exponentially with labor l_i , with exponent $\delta > 0$ and constant $c > 0$ (for exact values used, see Table 1).

Parameter		Value
Number of agents	N	100
Minimum skill value		1.24
Maximum skill value		159.1
Maximum labor choice		100
Labor disutility coefficient	c	0.0005
Labor disutility exponent	δ	3.5
Min bracket rate		0%
Max bracket rate		100%
Rate discretization (AI Economist)		5%

Table 1: Hyperparameters for the One-Step Economy environment.

Gather-Trade-Build Simulation. *Gather-Trade-Build* simulates a multi-step trading economy in a two-dimensional grid-world. Table 2 provides details regarding the simulation hyperparameters. Agents can gather resources, earn coins by using the resources of stone and wood to build houses, and trade with other agents to exchange resources for coins. Agents start at different initial locations in the world and are parameterized by different skill levels (described below). Simulations are run in episodes of 1000 timesteps, which is subdivided into 10 tax periods, each lasting 100 timesteps.

The state of the world is represented as an $n_h \times n_w \times n_c$ tensor, where n_h and n_w are the size of the world and n_c is the number of unique entities that may occupy a cell, and the value of a given element indicates which entity is occupying the associated location.

The action space of the agents includes 4 movement actions: up, down, left, and right. Agents are restricted from moving onto cells that are occupied by another agent, a water tile, or another agent’s house.

Stone and wood stochastically spawn on special resource regeneration cells. Agents can gather these resources by moving to populated resource cells. After harvesting, resource cells remain empty until new resources spawn. By default, agents collect 1 resource unit, with the possibility of a bonus unit also being collected, the probability of which is determined by the agent’s *gather-skill*. Resources and coins are accounted for in each agent’s *endowment* x , which represents how many coins, stone, and wood each agent owns.

Parameter		Value
Episode length	H	1000
World height	n_h	25
World width	n_w	25
Resource respawn probability		0.01
Max resource health		1
Starting agent coin	$C_{i,0}$	0
Iso-elastic utility exponent	η	0.23
Move labor		0.21
Gather labor		0.21
Trade labor		0.05
Build labor		2.1
Minimum build payout		10
Build payment max skill multiplier		3
Max bid/ask price		10
Max bid/ask order duration		50
Max number of open orders per resource		5
Tax period duration	\mathcal{T}	100
Min bracket rate		0%
Max bracket rate		100%
Rate discretization (AI Economist)		5%

Table 2: Hyperparameters for the Gather-Trade-Build environment.

Agent observations include the state of their own endowment (wood, stone, and coin), their own skill levels, and a view of the world state tensor within an egocentric spatial window (see Figure 7).

Our experiments use a world of size 25-by-25 (40-by-40) for four agent (ten agent) environments, where agent spatial observations have size 11-by-11 and are padded as needed when the observation window extends beyond the world grid.

The planner observations include each agent’s endowment but not skill levels (see Figure 7). We do not include the spatial state in the planner’s observations (in pilot experiments, we observed that this choice did not affect performance).

Trading. Agents can buy and sell resources from one another through a *continuous double auction*. Agents can submit *asks* (the number of coins they are willing to accept) or *bids* (how much they are willing to pay) in exchange for one unit of wood or stone.

The action space of the agents includes 44 actions for trading, representing the combination of 11 price levels (0, . . . , 10 coin), 2 directions (bids and asks), and 2 resources (wood and stone). Each trade action maps to a single order (i.e. bid 3 coins for 1 wood, ask for 5 coins in exchange for 1 stone, etc.). Once an order is submitted, it remains open until either it is matched (in which

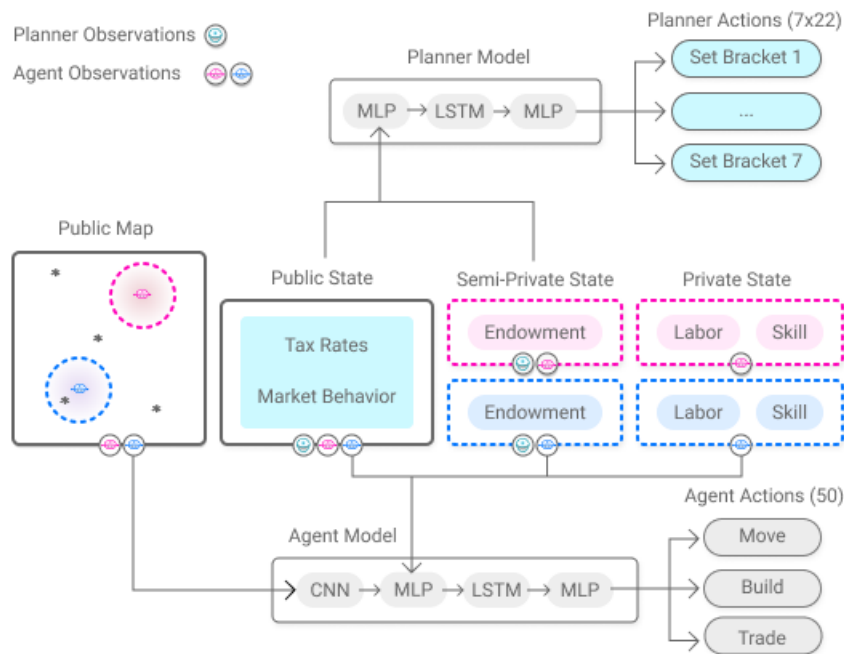


Figure 7: **Observation and action spaces for economic agents and the social planner.** The agents and the planner observe different subsets of the world state. Agents observe their spatial neighborhood, market prices, tax rates, inventories, and skill level. Agents can decide to move (and therefore gather if moving onto a resource), buy, sell, or build. There are 50 unique actions available to the agents. The planner observes market prices, tax rates, and agent inventories. The planner decides how to set tax rates, choosing one of 22 settings for each of the 7 tax brackets.

case a trade occurs) or it expires (after 50 timesteps). Agents are restricted from having more than 5 open orders for each resource, and are restricted from placing orders that they cannot complete (they cannot bid with more coins than they possess and cannot submit asks for resources that they do not have).

A bid/ask pair forms a valid trade if they are for the same resource and the bid price matches or exceeds the ask price. When a new order is received it is compared against complementary orders to identify potential valid trades. When a single bid (ask) could be paired with multiple existing asks (bids), priority is given to the ask (bid) with the lowest (highest) price; in the event of ties, priority then is given to the earliest order and then at random. Once a match is identified, the trade is executed using the price of whichever order was placed first.

For example, if the market receives a new bid that offers 8 coins for 1 stone and the market has two open asks offering 1 stone for 3 coins and 1 stone for 7 coins, received in that order, the market would pair the bid with the first ask and a trade would be executed for 1 stone at a price of 3 coins. The bidder loses 3 coins and gains 1 stone; the asker loses 1 stone and gains 3 coins. Once a bid and ask are paired and the trade is executed, both orders are removed.

The state of the market is captured by the number of outstanding bids and asks at each price level for each resource. Agents observe these counts both for their own bids/asks as well as the cumulative bids/asks of other agents. The planner observes the cumulative bids/asks of all agents. In addition, both the agents and the planner observe historical information from the market: the average trading price for each resource, as well as the number of trades at each price level.

Building. Agents can choose to spend one unit of wood and one unit of stone to build a house, and this places a house tile at the agent’s current location and earns the agent some number of coins. Agents are restricted from building on source cells as well as locations where a house already exists. The number of coins earned per house is identical to an agent’s *build-skill*, a numeric value between 10 and 30. As such, agents can earn between 10 and 30 coins per house built. Skill is heterogeneous across agents and does not change during an episode. Each agent’s action space includes 1 action for building.

Labor. Over the course of an episode of 1000 timesteps, agents accumulate labor cost, which reflects the amount of effort associated with their actions. Each type of action (moving, gathering, trading, and building) is associated with a specific labor cost. All agents experience the same labor costs.

Taxation Mechanism. Taxation is implemented using income brackets and bracket tax rates. All taxation is anonymous: tax rates and brackets do not depend on the identity of taxpayers. The payable tax for income z is computed as follows:

$$T(z) = \sum_{j=1}^B \tau_j \cdot ((b_{j+1} - b_j) \mathbf{1}[z > b_{j+1}] + (z - b_j) \mathbf{1}[b_j < z \leq b_{j+1}]), \quad (4)$$

where B is the number of brackets, and the τ_j and b_j are marginal tax rates and income boundaries of the brackets, respectively.

Each simulation episode has 10 tax years. On the first time step of each tax year, marginal tax rates are set that will be used to collect taxes when the tax year ends. For baseline models, tax rates are set either formulaically or fixed. For taxes controlled by a deep neural network, the action space of the planner is divided into seven action subspaces, one for each tax bracket: $(0, 0.05, 0.10, \dots, 1.0)^7$. Each subspace denotes the set of discretized marginal tax rates available to the planner. Discretization of tax rates only applies to deep learning networks, enabling standard techniques for RL with discrete actions.

Each agent observes the current tax rates, indicators of the temporal progress of the current tax year, and the set of sorted and anonymized incomes the agents reported in the previous tax year. In addition to this global tax information, each agent also observes the marginal rate at the level of income it has earned within the current tax year so far. The planner also observes this global tax information, as well as the non-anonymized incomes and marginal tax rate (at these incomes) of each agent in the previous tax year.

Redistribution Mechanism. An agent’s pretax income z_i for a given tax year is defined simply as the change in its coin endowment C_i since the start of the year. Accordingly, taxes are collected at the end of each tax year by subtracting $T(z_i)$ from C_i .

Taxes are used to redistribute wealth: the total tax revenue is evenly redistributed back to the agents. In total, at the end of each tax year, the coin endowment for agent i changes according to $\Delta C_i = -T(z_i) + \frac{1}{N} \sum_j^N T(z_j)$, where N is the number of agents. Through this mechanism, agents may gain coin when they receive more through redistribution than they pay in taxes.

Gather-Trade-Build Scenarios. We considered two spatial layouts: *Open-Quadrant* and *Split-World*, see Figure 4.

Open-Quadrant features four regions delineated by impassable water with passageways connecting each quadrant. Quadrants contain different combinations of resources: both stone and wood, only stone, only wood, or nothing. Agents can freely access all quadrants, if not blocked by objects or other agents.

Split-World features two disconnected regions: the top contains stone and wood, while the bottom only has stone. Water tiles prevent agents from moving from one region to the other.

All scenarios use a fixed set of *build-skills* based on a clipped Pareto distribution (sampled skills are clipped to the maximum skill value) and determine each agent’s starting location based on its assigned build-skill. The *Open-Quadrant* scenario assigns agents to a particular corner of the map, with similarly skilled agents being placed in the same starting quadrant. (Agents in the lowest build-skill quartile start in the wood quadrant; those in the second quartile start in the stone quadrant; those in the third quartile start in the quadrant with both resources; and agents in the highest build-skill quartile start in the empty quadrant.) The *Split-World* scenario allows control over which agents have access to both wood and stone versus access to only stone. We consider 4 *Split-World* variations, each with ten agents. Each variation gives stone and wood access to a specific subset of the ten agents, as determined by their build-skill rank. For example: *Split-World-1,2,3* places the 3 highest-skilled agents in the top, *Split-World-8,9,10* places the 3 lowest-skilled agents in the top, and *Split-World-5,6* places the 2 middle-skilled agents in the top.

Agent Utility. Following optimal taxation theory, agent utilities depend positively on accumulated coin $C_{i,t}$, which only depends on post-tax income $\tilde{z} = z - T(z)$. In contrast, the utility for agent i depends negatively on accumulated labor $L_{i,t} = \sum_{k=0}^t l_{i,k}$ at timestep t . The utility for an agent i is:

$$u_{i,t} = \frac{C_{i,t}^{1-\eta} - 1}{1 - \eta} - L_{i,t}. \quad (5)$$

Agents learn behaviors that maximize their expected total discounted utility for an episode. We found that build-skill is a significant determinant of behavior; agents’ gather-skill empirically does not affect optimal behavior in our settings.

All of our experiments use a fixed set of build-skills, which, along with labor costs, are roughly calibrated so that (1) agents need to be strategic in how they choose to earn income, and

(2) the shape of the resulting income distribution roughly matches that of the 2018 US economy with trained optimal agent behaviors.

Social Planner. The simulation environment includes a *social planner* who uses tax policy and lump-sum redistribution to influence social outcomes. Each episode is divided into 10 tax years. At the start of each tax year, the planner chooses a tax schedule $T(z)$ that determines the amount of taxes each agent will owe as a function of its income z earned during the tax year and redistributes tax revenue.

We compare four kinds of planners: (1) *Free Market*: a fixed-rate planner where all tax rates are 0%; (2) *US Federal*: a fixed-rate planner where bracketed marginal tax rates follow a progressive scheme adapted from the 2018 US federal single-filer income tax schedule; (3) *Saez tax*: an adaptive planner that computes theoretically optimal marginal rates using the empirical income distribution and elasticity of income with respect to taxation; and (4) *AI Economist*: a deep neural network, adaptive planner that maps a set of planner observations to bracketed marginal tax rates, which is trained via reinforcement learning (RL) to maximize social welfare.

Two-level Deep Reinforcement Learning. RL provides a flexible way to simultaneously optimize and model the behavioral effects of tax policies. We instantiate RL at two levels, that is, for two types of actors: training agent behavioral policy models and a taxation policy model for the social planner.

We train each actor’s behavioral policy using deep reinforcement learning, which learns the weights θ_i of a neural network $\pi(a_{i,t}|o_{i,t}; \theta_i)$ that maps an actor’s observations to actions. Network weights are trained to maximize the expected total discounted reward of the output actions.

Specifically, for an agent i using a behavioral policy $\pi_i(a_t|o_t; \theta_i)$, the RL training objective is (omitting the tax policy π_p):

$$\max_{\pi_i} \mathbb{E}_{a_1 \sim \pi_1, \dots, a_N \sim \pi_N, s' \sim \mathcal{P}} \left[\sum_{t=0}^H \gamma^t r_t \right], \quad (6)$$

where s' is the next state and \mathcal{P} denotes the dynamics of the environment. The objective for the planner policy π_p is similar. Standard model-free policy gradient methods update the policy weights θ_i using

$$\Delta \theta_i \propto \mathbb{E}_{a_1 \sim \pi_1, \dots, a_N \sim \pi_N, s' \sim \mathcal{P}} \left[\sum_{t=0}^H \gamma^t r_t \nabla_{\theta_i} \log \pi_i(a_{i,t}|o_{i,t}; \theta_i) \right]. \quad (7)$$

In our work, we use *proximal policy gradients* (PPO) (32), an extension of Formula 7 to train all actors (both agents and planner).

To improve learning efficiency, we train a single agent policy network $\pi(a_{i,t}|o_{i,t}; \theta)$ whose weights are shared by all agents, that is, $\theta_i = \theta$. This network is still able to embed diverse, agent-specific behaviors by conditioning on agent-specific observations.

At each timestep t , each agent observes: its nearby spatial surroundings; its current endowment (stone, wood, and coin); private characteristics, such as its building skill; the state of the markets for trading resources; and a description of the current tax rates. These observations form the inputs to the policy network, which uses a combination of convolutional, fully connected, and recurrent layers to represent spatial, non-spatial, and historical information, respectively. For recurrent components, each agent maintains its own hidden state. This is visualized in Figure 7. For the detailed model architecture and training hyperparameters, see Tables 3 and 4. The

Parameter		Value
Number of parallel environment replicas		30
Sampling horizon (steps per replica)	\mathcal{H}	200
Agent SGD minibatch size (# agents = 4)		600
Agent SGD minibatch size (# agents = 10)		1500
Planner SGD minibatch size		1500
SGD sequence length		25
Policy updates per horizon (agent)		40
Policy updates per horizon (planner)		4
CPU		15
Learning rate (agent)		0.0003
Learning rate (planner)		0.0001
Entropy regularization coefficient (agent)		0.025
Entropy regularization coefficient (planner)		0.125
Discount factor	γ	0.998
Generalized Advantage Estimation discount parameter	λ	0.98
Gradient clipping norm threshold		10
Value function loss coefficient		0.05
Phase <i>one</i> training duration		25M steps
Phase <i>two</i> training duration		1B steps
Phase <i>two</i> initial max τ		10%
Phase <i>two</i> tax annealing duration		27M steps
Phase <i>two</i> entropy regularization annealing duration		50M steps

Table 3: Hyperparameters for two-level reinforcement learning (RL), which trains multiple agents and a social planner. The base RL algorithm is the proximal policy gradient algorithm (32).

policy network for the social planner follows a similar construction, but differs somewhat in the information it observes. Specifically, at each timestep, the planner policy observes: the current inventories of each agent; the state of the resource markets; and a description of the current tax rates. The planner cannot directly observe private information such as an agent’s skill level.

Training Objectives. Rational economic agents train their policy π_i to optimize their total discounted utility over time, while experiencing tax rates τ set by the planner’s policy π_p . The

Parameter	Value
Number of convolutional layers	2
Number of fully-connected layers	2
Fully-connected layer dimension (agent)	128
Fully-connected layer dimension (planner)	256
LSTM cell size (agent)	128
LSTM cell size (planner)	256
Agent spatial observation box half-width	5

Table 4: Hyperparameters for the neural networks implementing the agent and planner policy models.

agent training objective is:

$$\forall i : \max_{\pi_i} \mathbb{E}_{\tau \sim \pi_p, \mathbf{a}_i \sim \pi_i, \mathbf{a}_{-i} \sim \pi_{-i}, s' \sim \mathcal{P}} \left[\sum_{t=1}^H \gamma^t r_{i,t} + u_{i,0} \right], \quad r_{i,t} = u_{i,t} - u_{i,t-1}, \quad (8)$$

where the instantaneous reward $r_{i,t}$ is the marginal utility for agent i at timestep t , and we use the isoelastic utility u_t as defined in Equation 5. Bold-faced quantities denote vectors, and the subscript “ $-i$ ” denotes quantities for all agents except for i .

For an agent population with monetary endowments $\mathbf{C}_t = (C_{1,t}, \dots, C_{N,t})$, we define equality $\text{eq}(\mathbf{C}_t)$ as:

$$\text{eq}(\mathbf{C}_t) = 1 - \frac{N}{N-1} \text{gini}(\mathbf{C}_t), \quad 0 \leq \text{eq}(\mathbf{C}_t) \leq 1, \quad (9)$$

where the Gini index is defined as

$$\text{gini}(\mathbf{C}_t) = \frac{\sum_{i=1}^N \sum_{j=1}^N |C_{i,t} - C_{j,t}|}{2N \sum_{i=1}^N C_{i,t}}, \quad 0 \leq \text{gini}(\mathbf{C}_t) \leq \frac{N-1}{N}. \quad (10)$$

We also define productivity as the sum of all incomes:

$$\text{prod}(\mathbf{C}_t) = \sum_i C_{i,t}. \quad (11)$$

Note that we assume the economy is closed: subsidies are always redistributed evenly among agents, no tax money leaves the system. Hence, the sum of pre-tax and post-tax incomes is the same. The planner trains its policy π_p to optimize social welfare:

$$\max_{\pi_p} \mathbb{E}_{\tau \sim \pi_p, \mathbf{a} \sim \pi, s' \sim \mathcal{P}} \left[\sum_{t=1}^H \gamma^t r_{p,t} + \text{swf}_0 \right], \quad r_{p,t} = \text{swf}_t - \text{swf}_{t-1}. \quad (12)$$

The *utilitarian* social welfare objective is the family of linear-weighted sums of agent utilities, defined for weights $\omega_i \geq 0$:

$$\text{swf}_t = \sum_{i=1}^N \omega_i \cdot u_{i,t}. \quad (13)$$

We use inverse-income as the weights: $\omega_i \propto \frac{1}{C_i}$, normalized to sum to 1. We also adopt an objective that optimizes a trade-off between equality and productivity, defined as the product of equality and productivity:

$$\text{swf}_t = \text{eq}(\mathbf{C}_t) \cdot \text{prod}(\mathbf{C}_t). \quad (14)$$

As agent incomes z_i depend on skill and access to resources, the heterogeneity in initial locations and build-skill are the main drivers of both economic inequality and specialization in Gather-Trade-Build.

Training Strategies. Two-level RL can be unstable, as the planner’s actions (setting tax rates) affect agent rewards (marginal utility depending on post-tax income).

We employ three learning curricula and two training phases to stabilize two-level RL. In phase one, agent policies are trained from scratch in a free-market (no-tax) environment for 50 million steps. In phase two, agents continue to learn in the presence of taxes for another 1 billion steps.

The first learning curriculum occurs during phase one: agents use a curriculum in phase one that anneals the utility cost associated with labor. The reason is that many actions cost labor, but few yield income. Hence, if exploring without a curriculum, a suboptimal policy can experience too much labor cost and converge to doing nothing.

The second learning curriculum occurs during phase two: we anneal the maximum marginal tax to prevent planners from setting extremely high taxes during exploration that reduce post-tax income to zero and discourage agents from improving their behaviors.

We also carefully balance entropy regularization, to prevent agent and planner policies from prematurely converging and promote the co-adaption of agent and planner policies. The entropy of the policy π for agent i , given an observation o_i , is defined as:

$$\text{entropy}(\pi) = -\mathbb{E}_{a \sim \pi} [\log \pi(a|o_i; \theta_i)]. \quad (15)$$

When training the AI Economist planner, we introduce the third learning curriculum by annealing the level of planner policy entropy regularization. Enforcing highly entropic planner policies during the early portion of phase two allows the agents to learn appropriate responses to a wide range of tax levels before the planner is able to optimize its policy.

Training Procedure. For training, we use proximal policy gradients (PPO) on mini-batches of experience collected from 30 parallel replicas of the simulation environment. Each environment replica runs for 200 steps during a training iteration. Hence, for each training iteration, 6,000 transitions are sampled for the planner and $N \cdot 6,000$ transitions are sampled for the agents, where N is the number of agents in the scenario, using the latest policy parameters.

The planner policy model is updated using transition mini-batches of size 1500, with one PPO update per minibatch (4 updates per iteration). The agent policy model is updated using transition mini-batches of size 400 (1500) for 4 (10) agent scenarios (40 updates per iteration). Table 4 provides details regarding the training hyperparameters. Algorithm 1 describes the full training procedure.

Action Spaces and Masks. Both agents and planners use discrete action spaces. We use action masking to prevent invalid actions, e.g., when agents cannot move across water, and to implement learning curricula. Masks control which actions can be sampled at a given time by assigning zero probability to restricted actions.

In addition, we include a no-operation action (NO-OP) in each action space. For the planner, each of the 7 action subspaces includes a NO-OP action. The NO-OP action allows agents to idle and the planner to leave a bracket’s tax rates unchanged between periods.

Action masks allow the planner to observe every timestep while only acting at the start of each new tax year. After the first timestep of a tax year, action masks enforce that only NO-OP planner actions are sampled.

Saez Tax. The Saez tax computes tax rates using an analytical formula (21) for a one-step economy with income distribution $f(z)$ and cumulative distribution $F(z)$. These rates maximize a weighted average $\sum_i w_i u_i$ of agent utilities, where the weights w_i reflect the redistributive preferences of the planner, and are optimal in an idealized one-step economy. The Saez tax computes marginal rates as:

$$\tau(z) = \frac{1 - G(z)}{1 - G(z) + a(z)e(z)}, \quad (16)$$

where z is pre-tax income, $G(z)$ is an income-dependent social welfare weight and $a(z)$ is the local Pareto parameter.

Specifically, let $\alpha(z)$ denote the *marginal average income at income z* , normalized by the fraction of incomes above z , i.e.,

$$\alpha(z) = \frac{z \cdot f(z)}{1 - F(z)}. \quad (17)$$

Let $G(z)$ denote the *normalized, reverse cumulative Pareto weight* over incomes above a threshold z , i.e.,

$$G(z) = \frac{1}{1 - F(z)} \int_{z'=z}^{\infty} p(z')g(z')dz'. \quad (18)$$

where $g(z)$ is the normalized social marginal welfare weight of an agent earning income z , and $1 - F(z)$ is the fraction of incomes above income z . In this way, $G(z)$ represents how much the social welfare function weights the incomes above z . Let *elasticity* $e(z)$ denote the *sensitivity of an agent’s income to changes in the tax rate when that agent’s income is z* , defined as

$$e(z) = \frac{1 - \tau(z)}{z} \frac{dz}{d(1 - \tau(z))}. \quad (19)$$

Both $G(z)$ and $a(z)$ can be computed directly from the (empirical) income distribution, but typically $e(z)$ needs to be estimated (which is challenging).

We set the social welfare weights $w_i \propto \frac{1}{z_i}$, normalized so the sum over all individuals is 1. This choice encodes a welfare focus on low-income agents.

Empirical Income Distribution. To apply the Saez tax, we use rollout data from a temporal window of episodes to estimate the empirical income distribution and compute $G(z)$ and $a(z)$. We aggregate reported incomes over a look-back window. We maintain a buffer of recent incomes reported by the agents, where each data point in this buffer represents the income reported by a single agent during a single tax year. Each simulation episode includes 10 tax years. As such, a single agent may report incomes in multiple different brackets in a single episode.

To compute $G(z)$ and $a(z)$, we first discretize the empirical income distribution and compute $\tau(z)$ within each of the resulting income bins. To get the average tax rate τ for each tax bracket, we take the average of the binned rates over the bracket’s income interval. Following the Saez analysis (21), when computing the top bracket rate, $G(z)$ is the total social welfare weight of the incomes in the top bracket, and $a(z)$ is computed as $\frac{m}{m - z^+}$ where m is the average income of those in the top bracket, and z^+ is the income cutoff for the top bracket (510 in our implementation, see Figure 5).

Estimating Elasticity. The most substantial obstacle to implementing the Saez tax is correctly identifying the elasticity $e(z)$, defined as in Equation 19. Owing to the complexity of the Gather-Trade-Build economy and agent learning dynamics, it is challenging to reliably measure local elasticities $e(z)$ as a function of income z . The large variance in empirical incomes caused large variance in the estimated local elasticity, leading to unstable two-level RL.

Therefore, we used a *global* elasticity estimate e , which assumes that elasticity is the same at all income levels. Empirically, we observe that the elasticity does not vary greatly across income ranges, hence justifying using a global elasticity.

For comparison, we also estimated the elasticity $e(z)$ using classic techniques, which use regression on observed incomes and marginal tax-rates obtained from agents trained under varying fixed flat-tax systems (37). Using a global constant elasticity for all agents, we instantiate this method by regression on K tuples $[(Z_k, \tau_k)]_{k=1}^K$ of observed total income $Z = \sum_i z_i$ and manually fixed flat tax rates τ in the simulation. Specifically, we use a linear model:

$$\log(Z) = \hat{e} \cdot \log(1 - \tau) + \log(\hat{Z}^0), \quad (20)$$

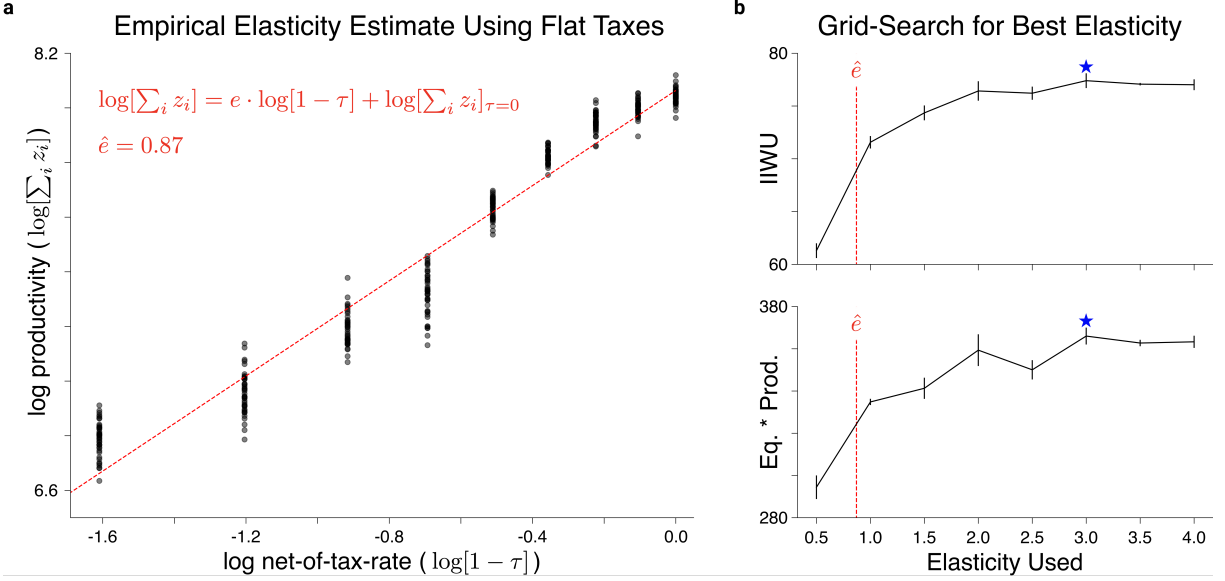


Figure 8: **Estimating elasticity for the Saez tax in the 4-agent Open-Quadrant scenario.** **a**, Regression on income and marginal tax-rate data yields elasticity estimates \hat{e} of approximately 0.87 (slope of the red dotted line). The net-of-tax-rate ($1 - \tau$) is the fraction of income agents retain after paying taxes. Productivity ($\sum_i z_i$) is the total pre-tax income earned by the agents. Each dot represents a $(\sum_i z_i, \tau)$ pair observed from a sweep over flat tax rates (see Methods). **b**, Social welfare with agents trained to convergence under the Saez tax, using a grid-search over elasticity parameters. Social welfare is highest under the Saez tax when the used elasticity parameter is approximately 3 (blue star), for both the inverse-income-weighted-utility objective (top) and the equality-times-productivity objective (bottom). Error bars denote standard error across the 3 random seeds used for each elasticity value.

where \hat{Z}^0 is a bias unit. Using a flat tax rate ensures agents always face the same tax rate during episodes, allowing for more consistent estimates. To generate data, we sweep over a range of values for τ , and collect observed total income data Z . This yields an estimate of $\hat{e} \sim 1$, which produces suboptimal social welfare. See Figure 8.

To provide the best possible performance for the Saez framework, we optimize the Saez tax using a grid search over possible e values. For each scenario, we separately conduct experiments involving sweeps over a range of potential values of e and select the best-performing one for each social welfare objective to use as a fixed elasticity estimate. This yields optimal elasticity estimate $e \sim 3$ in the 4-agent Open-Quadrant scenario, substantially higher than that estimated through regression techniques. See Figure 8.

Quantification and Statistical Significance. All experiments in the Open-Quadrant Gather-Trade-Build scenarios were repeated with 10 random seeds; experiments in the Split-World

Gather-Trade-Build scenarios and the One-Step Economy were repeated with 5 random seeds.

For a given repetition, we compute each performance metric, e.g. equality or social welfare, as its average value over the last 3000 episodes of training (the last 100 episodes for each of the 30 parallel environments). We report the average and standard error of these metrics across the 5 or 10 random seeds within a particular experiment group (Figure 3, Figure 4). Statistical significance is computed using a two-sample t-test.

In other analyses (Figures 5 and 6), we compute agent-wise statistics, e.g. pre-tax income and wealth transfer, using agent-specific statistics for each of the 10 tax periods in the episode. We conduct our analyses using the 40 most recent episodes (prior to the end of training, or prior to 250 million training steps where noted) for each repetition. For these analyses, we report the averages and standard deviations across the 400 associated episodes within each group of experiments.

References

1. United Nations, *Inequality Matters: Report of the World Social Situation 2013* (Department of Economic and Social Affairs, 2013).
2. S. v. Subramanian, I. Kawachi, *Epidemiologic Reviews* **26**, 78–91 (2004).
3. R. B. Myerson, *Mathematics of Operations Research* **6**, 58–73 (1981).
4. R. E. Lucas Jr, presented at the Carnegie-Rochester Conference Series on Public Policy, vol. 1, pp. 19–46.
5. A. M. Rivlin, P. M. Timpane, *Ethical and legal issues of social experimentation* (Brookings Institution Washington, DC, 1975), vol. 4.
6. V. Conitzer, T. Sandholm, presented at the Proceedings of the 18th Conference on Uncertainty in Artificial Intelligence, pp. 103–110.
7. T. Sandholm, presented at the International Conference on Principles and Practice of Constraint Programming, pp. 19–36.
8. H. Narasimhan, S. Agarwal, D. C. Parkes, presented at the Proceedings of the 25th International Joint Conference on Artificial Intelligence, pp. 433–439.
9. T. Baumann, T. Graepel, J. Shawe-Taylor, *arXiv:1806.04067 [Cs]*, arXiv: 1806.04067, (2018; <http://arxiv.org/abs/1806.04067>) (June 2018).
10. P. Dütting, Z. Feng, H. Narasimhan, D. C. Parkes, S. S. Ravindranath, presented at the Proc. 36th Int. Conf. On Machine Learning, pp. 1706–1715.
11. D. Silver *et al.*, *Nature* **550**, 354 (2017).
12. O. Vinyals *et al.*, *Nature* **575**, 350–354 (2019).
13. OpenAI, *OpenAI Five*, <https://blog.openai.com/openai-five/>, 2018.

14. J. Z. Leibo, V. Zambaldi, M. Lanctot, J. Marecki, T. Graepel, *arXiv:1702.03037 [Cs]*, arXiv: 1702.03037, (2018; <http://arxiv.org/abs/1702.03037>) (Feb. 2017).
15. Y. Bengio, J. Louradour, R. Collobert, J. Weston, presented at the ICML.
16. R. J. Williams, J. Peng, *Connection Science* **3**, 241–268 (1991).
17. P. A. Diamond, J. A. Mirrlees, *The American Economic Review* **61**, 8–27 (1971).
18. J. A. Mirrlees, *Journal of Public Economics* **6**, 327–358, ISSN: 0047-2727, (2019; <http://www.sciencedirect.com/science/article/pii/0047272776900475>) (Nov. 1976).
19. N. G. Mankiw, M. Weinzierl, D. Yagan, en, *Journal of Economic Perspectives* **23**, 147–174, ISSN: 0895-3309, (2019; <https://www.aeaweb.org/articles?id=10.1257/jep.23.4.147>) (Dec. 2009).
20. A. Auerbach, J. Hines, in *Handbook of Public Economics*, ed. by A. J. Auerbach, M. Feldstein (Elsevier, ed. 1, 2002), vol. 3, chap. 21, pp. 1347–1421, (<https://EconPapers.repec.org/RePEc:eee:pubchp:3-21>).
21. E. Saez, *The Review of Economic Studies* **68**, 205–229 (2001).
22. E. Saez, S. Stantcheva, *American Economic Review* **106**, 24–45 (2016).
23. N. R. Kocherlakota, *The New Dynamic Public Finance* (Princeton University Press, STU - Student edition, 2010), ISBN: 978-0-691-13915-9, (2019; www.jstor.org/stable/j.ctt7s9rn).
24. S. Stantcheva, *Annual Review of Economics*, (https://www.dropbox.com/s/xca67zq04v03zqr/Stantcheva_Dynamic_Taxation_Final.pdf?dl=0) (2020).
25. S. Albanesi, C. Sleet, *The Review of Economic Studies* **73**, 1–30 (2006).
26. E. Bonabeau, en, *Proceedings of the National Academy of Sciences* **99**, 7280–7287, ISSN: 0027-8424, 1091-6490, (2019; https://www.pnas.org/content/99/suppl_3/7280) (May 2002).
27. K. Bloomquist, *Public Finance Review* **39**, 25–49, (2019; https://econpapers.repec.org/article/saepubfin/v_3a39_3ay_3a2011_3ai_3a1_3ap_3a25-49.htm) (2011).
28. F. J. Miguel, J. A. Noguera, T. Llacer, E. Tapia, presented at the Ecms.
29. N. Garrido, L. Mittone, en, *The Journal of Socio-Economics* **42**, 24–30, ISSN: 1053-5357, (2019; <http://www.sciencedirect.com/science/article/pii/S105353571200114X>) (Feb. 2013).
30. *Penn Wharton Budget Model* (<https://budgetmodel.wharton.upenn.edu/tax-policy-1>).

31. J. Gokhale, (<https://budgetmodel.wharton.upenn.edu/issues/2018/2/6/w2018-1>) (2018).
32. J. Schulman, F. Wolski, P. Dhariwal, A. Radford, O. Klimov, *arXiv:1707.06347* (2017).
33. R. S. Sutton, A. G. Barto, *Reinforcement Learning: An Introduction*, en, Google-Books-ID: uWV0DwAAQBAJ (MIT Press, Oct. 2018), ISBN: 978-0-262-35270-3.
34. K. J. Arrow, *Essays in the theory of risk-bearing*, 90–120 (1971).
35. J. N. Foerster *et al.*, *arXiv:1709.04326 [Cs]*, arXiv: 1709.04326, (2019; <http://arxiv.org/abs/1709.04326>) (Sept. 2017).
36. R. Lowe *et al.*, *arXiv:1706.02275 [Cs]*, arXiv: 1706.02275, (2018; <http://arxiv.org/abs/1706.02275>) (June 2017).
37. J. Gruber, E. Saez, *Journal of Public Economics* **84**, 1–32 (2002).

11 End Notes

- **Acknowledgements.** We thank Kathy Baxter for the ethical review. We thank Nikhil Naik, Lofred Madzou, Simon Chesterman, Rob Reich, Mia de Kuijper, Scott Kominers, Gabriel Kriendler, Stefanie Stantcheva, Stefania Albanesi, and Thomas Piketty for valuable discussions.
- **Author Contributions.** A.T. and S.Z. contributed equally. R.S. and S.Z. conceived and directed the project; S.Z., A.T., and D.P. developed the theoretical framework; A.T., S.S., and S.Z. developed the economic simulator, implemented the reinforcement learning platform, and performed experiments; A.T., S.Z., and D.P. processed and analyzed experiments with AI agents; S.Z., A.T., and D.P. drafted the manuscript; R.S. planned and advised the work, and analyzed all results; All authors discussed the results and commented on the manuscript.
- Source code for the economic simulation is available at <https://www.github.com/salesforce/ai-economist>.
- The authors declare no competing interests.
- All data needed to evaluate the conclusions in the paper are present in the paper and/or the Supplementary Materials.
- The data can be provided by Stephan Zheng pending scientific review and a completed material transfer agreement. Requests for the data should be submitted to: stephan.zheng@salesforce.com.
- The authors acknowledge that they received no funding in support for this research.

Algorithm 1 Two-level Reinforcement Learning. Agents and social planner learn simultaneously. Bold-faced symbols indicate quantities for multiple agents. Note that agents share weights.

Input

- \mathcal{H} Sampling horizon
- \mathcal{T} Tax period length
- \mathbb{A} On-policy learning algorithm (in this work, PPO (32))
- \mathcal{C} Stopping criterion (for instance, agent and planner rewards have not improved)

Output

- θ Trained agent policy weights
- ϕ Trained planner policy weights

$s, \mathbf{o}, o_p, \mathbf{h}, h_p \leftarrow s_0, \mathbf{o}_0, o_{p,0}, \mathbf{h}_0, h_{p,0}$ \triangleright Reset episode: initialize world state s , observation o , hidden states h

$\theta, \phi \leftarrow \theta_0, \phi_0$ \triangleright Initial agent and planner policy weights

$\mathcal{D}, \mathcal{D}_p \leftarrow \{\}, \{\}$ \triangleright Reset agent and planner transition buffers

while training **do**

for $t = 1, \dots, \mathcal{H}$ **do**

$\mathbf{a}, \mathbf{h} \leftarrow \pi(\cdot | \mathbf{o}, \mathbf{h}, \theta)$ \triangleright Sample agent actions; update hidden state

if $t \bmod \mathcal{T} = 0$ **then** \triangleright First timestep of tax period

$\tau, h_p \leftarrow \pi_p(\cdot | o_p, h_p, \phi)$ \triangleright Sample marginal tax rates; update planner hidden state

else

$\text{no-op}, h_p \leftarrow \pi_p(\cdot | o_p, h_p, \phi)$ \triangleright Only update planner hidden state

end if

$s', \mathbf{o}', o'_p, \mathbf{r}, r_p \leftarrow \text{Env.step}(s, \mathbf{a}, \tau)$ \triangleright Next state / observations, pre-tax reward,
planner reward

if $t \bmod \mathcal{T} = \mathcal{T} - 1$ **then** \triangleright Last timestep of tax period

$s', \mathbf{o}', o'_p, \mathbf{r}, r_p \leftarrow \text{Env.tax}(s', \tau)$ \triangleright Apply taxes; compute post-tax rewards

end if

$\mathcal{D} \leftarrow \mathcal{D} \cup \{(\mathbf{o}, \mathbf{a}, \mathbf{r}, \mathbf{o}')\}$ \triangleright Update agent transition buffer

$\mathcal{D}_p \leftarrow \mathcal{D}_p \cup \{(o_p, \tau, r_p, o'_p)\}$ \triangleright Update planner transition buffer

$s, \mathbf{o}, o_p \leftarrow s', \mathbf{o}', o'_p$

end for

 Update θ, ϕ using data in $\mathcal{D}, \mathcal{D}_p$ and \mathbb{A} .

$\mathcal{D}, \mathcal{D}_p \leftarrow \{\}, \{\}$ \triangleright Reset agent and planner transition buffers

if episode is completed **then**

$s, \mathbf{o}, o_p, \mathbf{h}, h_p \leftarrow s_0, \mathbf{o}_0, o_{p,0}, \mathbf{h}_0, h_{p,0}$ \triangleright Reset episode

end if

if criterion \mathcal{C} is met **then return** θ, ϕ

end if

end while
