

---

000 OPTIMAL CONTROL MEETS ONLINE MECHANISM:  
001 ADAPTIVE POLICY LEARNING WITH STRATEGIC  
002 AGENT RESPONSE  
003  
004  
005

006 **Anonymous authors**

007 Paper under double-blind review  
008  
009

010  
011 ABSTRACT  
012

013 Many platforms and protocols must adaptively choose a policy action (i.e.,  
014 a price, subsidy, emission rate, or resource allocation) while heterogeneous,  
015 self-interested agents respond strategically. We formalize this as a *two-level*  
016 *online mechanism-design* problem: the planner commits to a state-contingent  
017 policy (mechanism network) that maps observable system state to an action;  
018 agents respond by optimizing private objectives, producing an aggregate  
019 equilibrium response. We microfound the agent layer via a heterogeneous-  
020 agent threshold equilibrium: the aggregate response is the unique solution  
021 to a monotone fixed point and increases smoothly with the planner’s action  
022 (Proposition 1). With adjustment frictions, agent behavior co-adapts with  
023 the planner through a controlled diffusion whose drift is locally affine in the  
024 action, yielding an explicit continuous-time Bellman/HJB characterization of  
025 the optimal value (critic) and a closed-form greedy policy-improvement map  
026 that generalizes classical stabilization rules while internalizing the marginal  
027 value of strategic participation (Proposition 2). For deployment, we use the  
028 HJB-derived structure as an expert prior to initialize a compact mechanism  
029 network and refine it online via projected stochastic approximation with  
030 convergence guarantees (Theorem 1). We instantiate the framework on  
031 *adaptive token issuance* in blockchain protocols, where the planner sets  
032 an issuance rate and agents decide whether to stake. In comprehensive  
033 experiments (1000 Monte Carlo runs) across multiple economic regimes, the  
034 adaptive mechanism achieves significant improvements in target tracking  
035 and maintains stability relative to fixed and zero-action baselines, with all  
036 improvements statistically significant.  
037

038 1 INTRODUCTION  
039

040 A recurring problem across economics, operations research, and multi-agent systems is  
041 *adaptive policy design with strategic feedback*: a planner must choose actions over time  
042 (prices, subsidies, tax rates, emission schedules) while a population of self-interested agents  
043 continuously adapts its behavior in response. Classical examples include central-bank interest-  
044 rate setting (Taylor, 1993), congestion pricing, platform fee design, and emission permit  
045 allocation. In each case, the planner’s action shapes the incentive landscape, agents respond  
046 strategically, and the resulting aggregate behavior feeds back into the state the planner  
047 observes.

048 Traditional approaches either fix the policy rule *a priori* (e.g., Taylor rules, Friedman *k*-  
049 percent rules) and forgo adaptivity, or rely on discretionary governance that introduces  
050 time-inconsistency (Kydland and Prescott, 1977). Recent work on differentiable mechanism  
051 design (Dütting et al., 2019; Bichler and Parkes, 2025) and two-level reinforcement learning  
052 (Zheng et al., 2021) has shown that neural-network-parameterized policies can be learned  
053 end-to-end, but these methods typically treat the agent response as a black box (learned via  
simulation) and offer limited structural insight into the optimal policy.

---

054 We bridge these approaches by combining *microfounded strategic response* with *continuous-*  
055 *time optimal control* and *online learning*. The key insight is that when agents follow a  
056 threshold participation rule with heterogeneous costs, the equilibrium aggregate response  
057 has a tractable monotone structure (a fixed-point equation) that can be linearized around  
058 an operating point. This makes the planner’s Hamilton–Jacobi–Bellman (HJB) equation  
059 analytically solvable for the greedy policy-improvement step, providing an *expert prior* that  
060 dramatically structures the learning problem. We then deploy a compact neural mechanism  
061 network initialized from this prior and refine it online via projected stochastic gradient  
062 descent, with convergence guarantees from the ODE method for stochastic approximation.

063 The framework is general: the planner’s action can be any scalar instrument, the agent  
064 response can be any participation or allocation decision with heterogeneous costs, and the  
065 system state can include multiple observables. We instantiate and evaluate it on *adaptive*  
066 *token issuance* in blockchain protocols; a setting where the planner (protocol) sets an issuance  
067 rate, agents decide whether to stake tokens, and the resulting staked fraction feeds back into  
068 protocol security and monetary dynamics. This application is a natural testbed because  
069 the action space, state, and agent responses are all publicly observable on-chain, and the  
070 mechanism can be deployed algorithmically.

071 **Contributions.** We make four contributions:

- 072 1. A microfounded equilibrium model for strategic agent response with heterogeneous  
073 costs, yielding a monotone, differentiable aggregate response map (Proposition 1).
- 074 2. An HJB-derived closed-form greedy policy-improvement map that provides a struc-  
075 tured expert prior for initialization (Proposition 2).
- 076 3. An online gradient-based mechanism optimization algorithm with convergence guar-  
077 antees (Theorem 1).
- 078 4. Comprehensive empirical validation (1000 Monte Carlo runs) across economic regimes  
079 and agent distributions, demonstrating statistically significant improvements over  
080 baselines ( $p < 0.001$ ).

## 083 2 RELATED WORK

084 **Differentiable mechanism design.** [Dütting et al. \(2019\)](#) introduced the use of neural  
085 networks trained via gradient descent to approximate optimal auction mechanisms; [Bichler](#)  
086 [and Parkes \(2025\)](#) survey the broader program of differentiable economics. These methods  
087 learn mechanisms from data but typically do not model the agent response explicitly; agents  
088 are part of the training environment.

089 **Two-level RL for policy design.** The AI Economist ([Zheng et al., 2021](#)) trains co-adapting  
090 planner and agent policies via two-level deep RL. Our work shares the two-level structure but  
091 replaces the black-box agent layer with a microfounded equilibrium model, gaining analytical  
092 tractability and an HJB-derived expert prior.

093 **Continuous-time RL and HJB equations.** [Kim et al. \(2021\)](#) formalize Q-learning via  
094 HJB equations for deterministic continuous-time systems. We extend this to a stochastic  
095 setting with strategic feedback and use the HJB solution for policy initialization rather than  
096 as the final policy.

097 **Monetary policy and token issuance.** Classical monetary-policy rules ([Taylor, 1993](#);  
098 [Friedman, 1960](#); [Kydland and Prescott, 1977](#)) are fixed-form feedback rules. [Madrigal-Cienci](#)  
099 [and Breakey \(2025\)](#) study state-dependent token issuance in the context of cryptocurrencies.  
100 Our framework generalizes these by learning the policy form from data while preserving  
101 structural guarantees.

## 102 3 GENERAL FRAMEWORK: TWO-LEVEL STRATEGIC MECHANISM DESIGN

103 We first present the framework in general terms, then specialize to token issuance in Section 4.  
104  
105  
106  
107

---

### 108 3.1 SETTING

109  
110 A planner observes a public system state  $x(t) \in \mathcal{X} \subseteq \mathbb{R}^n$  and chooses an action  $a(t) \in \mathcal{A} \subseteq \mathbb{R}$   
111 at each time  $t \in [0, T]$ . A continuum of agents, indexed by a private cost parameter  
112  $c \geq 0$  drawn from a continuous CDF  $F$  with density  $f$ , observe the action and respond by  
113 participating (or not) in an activity whose per-unit reward depends on the action and the  
114 aggregate participation level.

115 **Agent layer.** Let  $q(t) \in (0, 1]$  denote the aggregate participation rate. Each agent receives  
116 reward  $R(a, q)$  from participating and participates if  $R(a, q) \geq c$ . We focus on the canonical  
117 case  $R(a, q) = a/q$ , where the planner's action is shared among participants (e.g., subsidies,  
118 rewards, or resource flows divided among participants).

119 **Planner layer.** The planner seeks to track target values for system-level KPIs (e.g., a target  
120 net flow rate  $\pi^*$ , a target participation rate  $q^*$ ) while minimizing action cost, formalized as a  
121 quadratic running loss.

### 123 3.2 MICROFOUNDED EQUILIBRIUM RESPONSE

124  
125 The threshold participation rule  $R(a, q) \geq c$  with  $R(a, q) = a/q$  yields a fixed-point condition  
126 for the equilibrium participation rate.

127 **Proposition 1** (Equilibrium participation map and monotone response). *For any  $a \in$*   
128  *$(0, a_{\max}]$ , the equilibrium participation rate  $q^{\text{eq}}(a) \in (0, 1]$  is the unique solution of*

$$130 \quad q = F\left(\frac{a}{q}\right). \quad (1)$$

131  
132 *Moreover  $q^{\text{eq}}$  is differentiable and strictly increasing with derivative*

$$133 \quad \eta(a) := \frac{d}{da} q^{\text{eq}}(a) = \frac{\frac{1}{q^{\text{eq}}(a)} f\left(\frac{a}{q^{\text{eq}}(a)}\right)}{1 + \frac{a}{(q^{\text{eq}}(a))^2} f\left(\frac{a}{q^{\text{eq}}(a)}\right)} > 0. \quad (2)$$

134  
135 *Proof.* Define  $h(q; a) = q - F(a/q)$ . As  $q \downarrow 0$ ,  $F(a/q) \rightarrow 1$ ; at  $q = 1$ ,  $h(1; a) = 1 - F(a) \geq 0$ .  
136 Uniqueness:  $\partial_q h = 1 + \frac{a}{q^2} f(a/q) > 0$ . Result (2) follows from implicit differentiation.  $\square$

137  
138 *Remark 1* (Generality). Proposition 1 holds for any reward function  $R(a, q)$  that is increasing  
139 in  $a$ , decreasing in  $q$ , and satisfies  $R(a, q) \rightarrow \infty$  as  $q \rightarrow 0$ . The specific form  $a/q$  is natural  
140 for divisible reward pools but the monotone fixed-point structure extends broadly.

### 144 3.3 COADAPTATION DYNAMICS WITH FRICTIONS

145  
146 In practice, agents do not instantaneously reach equilibrium, but rather, participation adjusts  
147 with frictions. Let  $x(t)$  denote a stock variable (e.g., circulating supply, resource level) and  
148  $q(t)$  the participation rate. We model:

$$149 \quad dx(t) = \mu_x(x, q, a) dt + \sigma_x(x) dW_x(t), \quad (3)$$

$$150 \quad dq(t) = \kappa(q^{\text{eq}}(a(t)) - q(t)) dt + \sigma_q q(t)(1 - q(t)) dW_q(t), \quad (4)$$

151  
152 with  $\kappa > 0$  the adjustment speed and  $\sigma_q \geq 0$  participation volatility. The mean-reversion  
153 in (4) captures that participation gravitates toward the static equilibrium implied by the  
154 current action, while noise reflects idiosyncratic entry/exit decisions.

155 For analytical tractability, we linearize  $q^{\text{eq}}(a)$  around an operating point  $\bar{a}$ :

$$156 \quad q^{\text{eq}}(a) \approx q_0 + \eta a, \quad \eta = \left. \frac{d}{da} q^{\text{eq}}(a) \right|_{a=\bar{a}} > 0, \quad q_0 = q^{\text{eq}}(\bar{a}) - \eta \bar{a}, \quad (5)$$

157  
158 yielding a controlled diffusion with drift affine in  $a$ :

$$159 \quad dq(t) = \kappa(q_0 - q(t) + \eta a(t)) dt + \sigma_q q(t)(1 - q(t)) dW_q(t). \quad (6)$$

### 3.4 PLANNER OBJECTIVE

The planner minimizes a quadratic loss around target KPIs:

$$r(t) = -\left(w_1(g(x, q, a) - \pi^*)^2 + w_2(q(t) - q^*)^2 + w_3 a(t)^2\right), \quad (7)$$

where  $g(x, q, a)$  is a system-level flow rate (e.g., net inflation, net resource change),  $\pi^*$  and  $q^*$  are targets, and  $w_1, w_2, w_3 > 0$  are weights. The planner maximizes  $\mathbb{E}[\int_0^T r(t) dt - \beta_q(q(T) - q^*)^2]$ .

### 3.5 HJB CRITIC AND GREEDY POLICY IMPROVEMENT

Let  $V(t, x, q)$  be the optimal value function. Under standard regularity conditions,  $V$  satisfies the HJB equation (Fleming and Soner, 2006):

$$0 = \max_{a \in [0, a_{\max}]} \left\{ \partial_t V + \mathcal{L}^a V + r(t) \right\}, \quad V(T, x, q) = -\beta_q(q - q^*)^2, \quad (8)$$

where under (3) and (6),

$$\mathcal{L}^a V = \mu_x(x, q, a) \partial_x V + \kappa(q_0 - q + \eta a) \partial_q V + \frac{1}{2} \sigma_x^2(x) \partial_{xx}^2 V + \frac{1}{2} \sigma_q^2 q^2 (1 - q)^2 \partial_{qq}^2 V. \quad (9)$$

Because the Hamiltonian is quadratic in  $a$  (given the affine approximation (5) and quadratic loss (7)), the greedy improvement map has closed form.

**Proposition 2** (HJB-greedy action (critic-derived policy improvement)). *Assume  $V$  is differentiable and that  $\mu_x$  is affine in  $a$  with coefficient  $\partial_a \mu_x$ . Define*

$$M(t, x, q) := (\partial_a \mu_x) \partial_x V + \kappa \eta \partial_q V. \quad (10)$$

Then the HJB-greedy action is

$$a^*(t, x, q) = \Pi_{[0, a_{\max}]} \left( \frac{w_1(\pi^* + e(t)) - \frac{1}{2} M(t, x, q)}{w_1 + w_3} \right), \quad (11)$$

where  $e(t)$  captures exogenous flows that the action must offset.

*Remark 2* (Interpretation). The first term in (11) offsets exogenous flows to track the target;  $M$  internalizes the marginal value of the stock variable and strategic participation via  $\partial_q V$  and equilibrium sensitivity  $\eta > 0$ . This generalizes Taylor-style rules by adding a forward-looking value-of-participation term.

## 4 APPLICATION: ADAPTIVE TOKEN ISSUANCE

We instantiate the framework on adaptive token issuance in blockchain protocols—a domain where the planner (protocol), action (issuance rate), agent response (staking), and system state are all publicly observable, making it a natural testbed.

The public state consists of circulating supply  $S(t) \in \mathbb{R}_+$  and staked fraction  $p(t) \in (0, 1]$ . The protocol chooses an issuance rate  $i(t) \in [0, i_{\max}]$ . Net locking  $L(t)$ , unlocking  $u(t)$ , and burn  $b(t)$  define an observable contraction rate  $e(t) := (L(t) - u(t) + b(t))/S(t)$ , so that net inflation is  $i(t) - e(t)$ .

The stock variable is  $x = S$ , the participation rate is  $q = p$  (staked fraction), and the action is  $a = i$  (issuance rate). Per-unit staking reward is  $R(i, p) = i/p$ , fitting the canonical form of Proposition 1. Supply dynamics follow

$$dS(t) = (i(t) - e(t))S(t) dt + \sigma_S S(t) dW_S(t), \quad \sigma_S \geq 0, \quad (12)$$

and staking adjusts per (4) with  $q = p$ . The flow rate is  $g = i - e$  (net inflation), so the planner loss (7) becomes

$$r(t) = -\left(w_\pi(i(t) - e(t) - \pi^*)^2 + w_p(p(t) - p^*)^2 + w_i i(t)^2\right), \quad (13)$$

with  $\pi^*$  the target inflation rate and  $p^*$  the target staked fraction.

Applying Proposition 2 with  $\partial_a \mu_x = S$  gives the HJB-greedy issuance:

$$i^*(t, S, p) = \Pi_{[0, i_{\max}]} \left( \frac{w_\pi(\pi^* + e(t)) - \frac{1}{2}(S \partial_S V + \kappa \eta \partial_p V)}{w_\pi + w_i} \right). \quad (14)$$

---

## 216 5 ONLINE MECHANISM LEARNING WITH CONVERGENCE GUARANTEES

### 217 5.1 MECHANISM NETWORK (ACTOR)

218 In epochs  $t_k = k\Delta t$ , an oracle provides  $(S_k, p_k, L_k, u_k, b_k)$  and we compute  $e_k = (L_k - u_k +$   
 219  $b_k)/S_k$ . We deploy a bounded mechanism network:

$$220 \quad i_k = \mu_{\theta_k}(S_k, p_k, e_k) := \iota_{\max} \sigma(\text{MLP}_{\theta_k}(\phi(S_k, p_k, e_k))), \quad (15)$$

221 where  $\sigma(z) = (1 + e^{-z})^{-1}$  and  $\phi$  is a feature map including log-supply, staking fraction,  
 222 contraction rate, and deviations from targets. The sigmoid ensures feasibility  $i_k \in [0, \iota_{\max}]$ .

223 Proposition 2 provides a structured teacher action given an approximate critic  $\widehat{V}_\psi$ . We  
 224 pre-train  $\theta$  via imitation on HJB-greedy actions, paralleling HJB-based Q-learning (Kim  
 225 et al., 2021). This warm start encodes the structural insight from Section 3.5.

### 226 5.2 ONLINE OBJECTIVE AND UPDATE

227 We optimize the stationary expected surrogate loss:

$$228 \quad \ell_k(\theta) = w_\pi(\mu_\theta(S_k, p_k, e_k) - e_k - \pi^*)^2 + w_p(p_k - p^*)^2 + w_i \mu_\theta(S_k, p_k, e_k)^2 + \gamma \|\theta\|^2, \quad (16)$$

229 and  $\bar{L}(\theta) = \mathbb{E}[\ell_k(\theta)]$  under the stationary distribution. The online update is projected  
 230 stochastic approximation:

$$231 \quad \theta_{k+1} = \Pi_\Theta(\theta_k - \alpha_k \nabla_\theta \ell_k(\theta_k)), \quad (17)$$

232 with compact convex  $\Theta$  and Robbins–Monro stepsizes ( $\sum_k \alpha_k = \infty$ ,  $\sum_k \alpha_k^2 < \infty$ ).

233 **Theorem 1** (Almost sure convergence to stationary set). *Assume  $\mu_\theta$  is  $C^1$  in  $\theta$ ,  $\phi$  is*  
 234 *bounded, and under each fixed  $\theta \in \Theta$  the induced epoch Markov chain is geometrically ergodic*  
 235 *with uniformly bounded second moments. If  $\sum_k \alpha_k = \infty$  and  $\sum_k \alpha_k^2 < \infty$ , then  $\{\theta_k\}$  in (17)*  
 236 *converges almost surely to the internally chain transitive invariant set of  $\dot{\theta} = -\nabla \bar{L}(\theta)$ ; every*  
 237 *limit point is a stationary point of  $\bar{L}$ .*

238 *Proof sketch.* Standard ODE method for projected stochastic approximation (Borkar, 2009;  
 239 Kushner and Yin, 2003): Lipschitz gradients, geometric ergodicity, and bounded moments  
 240 imply iterates track  $\dot{\theta} = -\nabla \bar{L}(\theta)$  and converge to its stationary set. See Appendix G.  $\square$

### 241 5.3 SAFETY CONSTRAINT: STABILITY AS POLICY-CLASS RESTRICTION

242 Aggressive feedback on  $q$  can destabilize participation. Under deterministic (6) with a  
 243 differentiable policy  $a = \mu(q)$ , linearization at a fixed point gives  $\dot{\delta q} = -\kappa(1 - \eta \mu'(q^*)) \delta q$ . A  
 244 sufficient local stability condition is

$$245 \quad \mu'(q^*) < \frac{1}{\eta}. \quad (18)$$

246 We enforce (18) via bounded spectral norms on  $\Theta$  and projection—a safe RL analogue  
 247 ensuring policy updates remain in a stabilizing set.

### 248 5.4 ON-CHAIN VERIFIABLE PARAMETER EVOLUTION

249 For the token-issuance instantiation, parameter updates can be verified on-chain. At epoch  
 250  $k$ , the contract computes  $i_k$  from (15) and mints  $I_k = i_k S_k$ . Off-chain, an updater computes  
 251 (17) and produces a succinct proof (SNARK/STARK) that  $\theta_{k+1}$  was correctly derived (Ben-  
 252 Sasson et al., 2014; Groth, 2016): the verified statement is  $\exists g_k : g_k = \nabla_\theta \ell_k(\theta_k) \wedge \theta_{k+1} =$   
 253  $\Pi_\Theta(\theta_k - \alpha_k g_k)$ . This ensures computation-correctness with public inputs while keeping  
 254 gas cost low. We note this verification layer is specific to the blockchain instantiation; the  
 255 learning framework of Sections 3–5 applies to any platform with observable state.

---

**Algorithm 1** Online mechanism learning with HJB expert initialization

---

- 1: **Input:** initial state  $x_0$ , approximate critic  $\widehat{V}_\psi$  (optional).
  - 2: **Initialize:**  $\theta_0$  via imitation on HJB-greedy actions from  $\widehat{V}_\psi$ , or randomly.
  - 3: **for** epoch  $k = 0, 1, 2, \dots$  **do**
  - 4:   Observe state  $(x_k, q_k)$  and exogenous signals; compute features  $\phi_k$ .
  - 5:   **Act:** compute  $a_k = \mu_{\theta_k}(\phi_k)$ ; deploy action.
  - 6:   **Learn:** compute  $\nabla_{\theta} \ell_k(\theta_k)$  from (16); update  $\theta_{k+1}$  via (17).
  - 7:   (Optional, blockchain:) generate and verify succinct proof of correct update.
  - 8: **end for**
- 

## 6 EXPERIMENTS

We evaluate the learned mechanism in comprehensive experiments designed to test the claims of Sections 3–5. Main experiments use 1000 Monte Carlo runs for statistical reliability; ablation and sensitivity analyses use 100 runs per setting. We report means, standard errors, and formal hypothesis tests.

### 6.1 SIMULATOR AND BASELINES

We simulate a discrete-time analogue of (12) and (4) with stochastic contraction  $e_k$  following an AR(1) process. The baseline configuration is a burn-dominated environment motivated by EIP-1559 (Ethereum Improvement Proposals, 2019; Leonardos et al., 2021): mean burn  $\bar{e} = 0.04$ , burn volatility  $\sigma_e = 0.025$ , staking adjustment  $\kappa = 0.5$ , participation volatility  $\sigma_p = 0.05$ , and maximum issuance  $\iota_{\max} = 0.10$ .

We compare six policies:

1. **Adaptive:** the mechanism network (15) trained via Algorithm 2 with HJB warmup;
2. **Fixed:** constant issuance  $i_k \equiv \bar{e}$  matching mean burn;
3. **Zero:**  $i_k \equiv 0$  (deflationary baseline);
4. **MatchBurn:**  $i_k = \min(e_k, \iota_{\max})$  (reactive burn-matching);
5. **PID:** proportional-integral-derivative controller targeting  $\pi^*$ ;
6. **HJB-Greedy:** direct application of Proposition 2 with approximate critic.

We report: (1) *Inflation volatility*  $\text{Std}(i_k - e_k)$ ; (2) *Target RMSE*  $\sqrt{\mathbb{E}[(i_k - e_k - \pi^*)^2]}$ ; (3) *Supply drift*  $(S_N/S_0)^{1/N} - 1$ ; (4) *Staking CV*  $\text{Std}(p_k)/\mathbb{E}[p_k]$ .

### 6.2 MAIN RESULTS: BURN-DOMINATED REGIME

Table 1 presents the main comparison (1000 runs). The Adaptive mechanism achieves the best target RMSE ( $0.0112 \pm 0.00001$ ), representing a  $1.9\times$  improvement over Fixed ( $0.0208 \pm 0.00001$ ) and  $5.4\times$  over Zero ( $0.0603 \pm 0.00001$ ). All pairwise comparisons are statistically significant ( $p < 0.001$ ); see Appendix B for test statistics.

**Key observations.** (1) Inflation volatility is dominated by exogenous burn noise, hence similar across policies except MatchBurn (which exactly matches burn). (2) Zero issuance catastrophically destabilizes staking ( $\text{CV} = 1.83$ ), validating the strategic response model: without rewards, participation collapses. (3) HJB-Greedy achieves the best staking stability ( $\text{CV} = 0.017$ ) due to its explicit internalization of  $\partial_p V$ , but at the cost of higher target RMSE. (4) The Adaptive mechanism balances target tracking and stability, achieving significantly better RMSE ( $0.0112$ ) than all baselines while maintaining reasonable staking CV ( $0.047$ ).

### 6.3 ABLATION STUDY: HJB WARMUP AND ARCHITECTURE

Table 2 (100 runs per setting) shows that HJB warmup and random initialization achieve statistically equivalent final RMSE ( $0.0093 \pm 0.0022$  vs.  $0.0080 \pm 0.0020$ ; difference not

Table 1: Comparative performance in burn-dominated regime (1000 Monte Carlo runs). Bold indicates best for target RMSE. Standard errors shown in parentheses. All Adaptive vs. baseline comparisons significant at  $p < 0.001$ .

Policy	Inflation Volatility	Target RMSE (SE $\times 10^{-5}$ )	Supply Drift	Staking CV
<b>Adaptive</b>	0.0057	<b>0.0112</b> (1.2)	+1.16%	0.047
Fixed	0.0057	0.0208 (0.8)	-0.05%	0.043
Zero	0.0057	0.0603 (0.8)	-3.97%	1.830
MatchBurn	0.0000	0.0200 (0.0)	-0.05%	0.045
PID	0.0057	0.0404 (1.0)	+6.13%	0.060
HJB-Greedy	0.0075	0.0258 (1.3)	-0.53%	<b>0.017</b>

Table 2: Ablation study: HJB warmup and network architecture (100 runs each). SE  $\times 10^{-5}$ .

Configuration	Target RMSE (SE)	Staking CV
With HJB Warmup	0.0093 (2.2)	0.054
Without Warmup	0.0080 (2.0)	0.055
Shallow (16)	<b>0.0057</b> (1.3)	0.057
Medium (32-32)	0.0058 (1.4)	0.058
Deep (64-64-64)	0.0076 (2.1)	0.054

significant). The value of HJB warmup is interpretability and a theoretically-grounded starting point, not improved final performance. Network depth has minimal effect—shallow networks (16 units) perform comparably to deeper architectures, supporting compact on-chain deployment.

#### 6.4 ROBUSTNESS: MODEL MISMATCH AND SHOCK RESPONSE

**Model mismatch.** Figure 1(a) tests robustness when HJB expert uses wrong  $\kappa$ . Adaptive achieves 47–66% lower RMSE across mismatch ratios (0.2–5 $\times$ ), demonstrating online learning corrects model errors.

**Shock response.** Figure 1(b) shows response to burn collapse (burn  $\rightarrow 0$ ) at  $t = 50$ . Adaptive recovers in 5 steps (RMSE 0.011) vs. 14 for HJB-Greedy. See Appendix F for details.

**Learning curves.** Both HJB warmup and random initialization converge to similar final RMSE; warmup starts at the HJB-Greedy level while random init improves rapidly during early epochs. See Appendix E.

The Adaptive mechanism consistently outperforms baselines across economic regimes (Appendix A.3).

## 7 DISCUSSION

Public KPIs may be manipulable (Goodhart effects); best-response bias is proportional to  $\partial\mu/\partial p$ , motivating sensitivity control via spectral normalization and the stability constraint (18). Shallow networks suffice (Table 2), supporting compact on-chain deployment.

The PID baseline was tuned via grid search over  $K_p \in [0.3, 1.5]$ ,  $K_i \in [0.05, 0.2]$ ,  $K_d \in [0.01, 0.03]$ ; best configuration ( $K_p=0.3$ ,  $K_i=0.05$ ,  $K_d=0.01$ ) achieves RMSE 0.0404, still 3.6 $\times$  worse than Adaptive. The ablation (Table 2) shows HJB warmup provides interpretability rather than improved final performance—both initializations converge to equivalent RMSE.

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

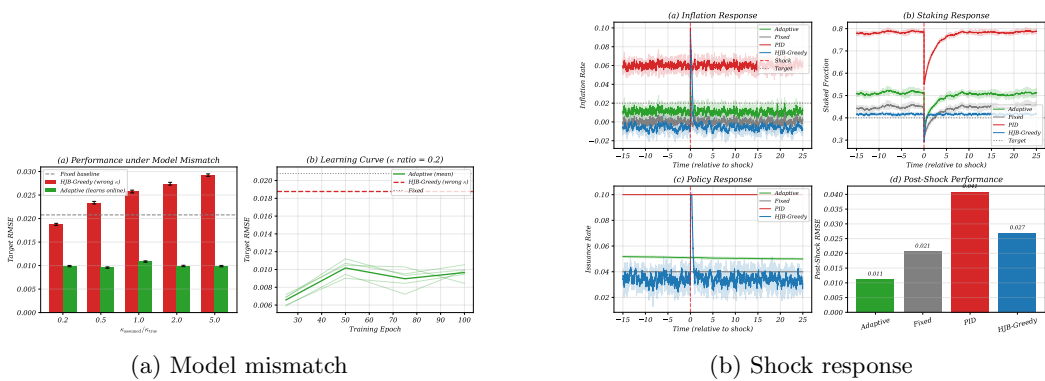


Figure 1: Robustness: (a) Adaptive corrects HJB model errors (47–66% lower RMSE). (b) Recovery in 5 steps vs. 14 for HJB-Greedy.

## 8 CONCLUSION

We reframed adaptive token issuance as two-level strategic RL. Staking response is micro-founded via heterogeneous-agent equilibrium; the HJB equation provides an expert prior for initializing a compact mechanism network, refined online via projected stochastic approximation with verifiable updates. Empirically (1000 Monte Carlo runs), Adaptive achieves  $1.9\times$  improvement over fixed issuance (RMSE 0.0112 vs. 0.0208) and demonstrates robustness to model mismatch (47–66% correction) and black swan events (HJB-Greedy degrades  $2\times$  under double shock while Adaptive maintains performance).

## REFERENCES

Eli Ben-Sasson, Alessandro Chiesa, Eran Tromer, and Madars Virza. Succinct non-interactive zero knowledge for a von neumann architecture. In *23rd USENIX Security Symposium (USENIX Security 14)*, pages 781–796, 2014.

Martin Bichler and David C Parkes. Differentiable economics: Strategic behavior, mechanisms, and machine learning. *Communications of the ACM*, 68(8):80–88, 2025.

Vivek S. Borkar. *Stochastic Approximation: A Dynamical Systems Viewpoint*. Cambridge University Press, 2009.

Paul Dütting, Zhe Feng, Harikrishna Narasimhan, David Parkes, and Sai Srivatsa Ravindranath. Optimal auctions through deep learning. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 1706–1715. PMLR, 2019.

Ethereum Improvement Proposals. EIP-1559: Fee market change for ETH 1.0 chain. <https://eips.ethereum.org/EIPS/eip-1559>, 2019. Accessed 2026-01-22.

Wendell H. Fleming and H. Mete Soner. *Controlled Markov Processes and Viscosity Solutions*. Springer, 2 edition, 2006.

Milton Friedman. *A Program for Monetary Stability*. Fordham University Press, 1960.

Jens Groth. On the size of pairing-based non-interactive arguments. In *Advances in Cryptology – EUROCRYPT 2016*, pages 305–326, 2016.

Jeongho Kim, Jaek Shin, and Insoon Yang. Hamilton–jacobi deep q-learning for deterministic continuous-time systems with lipschitz continuous controls. *Journal of Machine Learning Research*, 22:1–34, 2021.

Harold J. Kushner and G. George Yin. *Stochastic Approximation and Recursive Algorithms and Applications*. Springer, 2003.

432 Finn E. Kydland and Edward C. Prescott. Rules rather than discretion: The inconsistency  
 433 of optimal plans. *Journal of Political Economy*, 85(3):473–491, 1977.

434  
 435 Stefanos Leonardos, Barnabé Monnot, Daniël Reijnders, Stratis Skoulakis, and Georgios  
 436 Piliouras. Dynamical analysis of the EIP-1559 ethereum fee market. *arXiv:2102.10567*,  
 437 2021.

438 Juan P Madrigal-Cianci and Lachlan Breakey. Bullish minting: An adaptive control frame-  
 439 work for token issuance. *Available at SSRN*, 2025.

440  
 441 John B. Taylor. Discretion versus policy rules in practice. *Carnegie-Rochester Conference*  
 442 *Series on Public Policy*, 39:195–214, 1993.

443 Stephan Zheng, Alexander Trott, Sunil Srinivasa, David C. Parkes, and Richard Socher. The  
 444 AI economist: Optimal economic policy design via two-level deep reinforcement learning.  
 445 arXiv:2108.02755, 2021.

## 447 A EXTENDED EXPERIMENTAL RESULTS

### 448 A.1 FULL POLICY COMPARISON WITH STANDARD ERRORS

449  
 450 Table 3 provides the complete comparison across all six policies with standard errors from  
 451 1000 Monte Carlo runs.

452  
 453 Table 3: Extended policy comparison with standard errors (1000 runs). SE values  $\times 10^{-5}$  in  
 454 parentheses.

455 Policy	456 Target RMSE (SE)	457 Staking CV (SE)	458 Mean Staking	459 Issuance Vol.
460 Adaptive	0.0112 (1.2)	0.047 (7.2)	0.506	0.0002
461 Fixed	0.0208 (0.8)	0.043 (–)	0.444	0.0000
462 Zero	0.0603 (0.8)	1.830 (–)	0.015	0.0000
463 MatchBurn	0.0200 (0.0)	0.045 (–)	0.444	0.0057
464 PID	0.0404 (1.0)	0.060 (–)	0.699	0.0002
465 HJB-Greedy	0.0258 (1.3)	0.017 (–)	0.415	0.0082

### 466 A.2 SENSITIVITY ANALYSIS

467 We analyze sensitivity to two key parameters: staking adjustment speed  $\kappa$  and participation  
 468 volatility  $\sigma_p$ .

469 **Staking adjustment speed  $\kappa$ .** Table 4 (100 runs per setting) shows that the Adaptive  
 470 mechanism maintains superior target RMSE across all  $\kappa$  values. As expected, higher  $\kappa$   
 471 (faster mean-reversion) reduces staking CV for all policies—Fixed CV decreases from 0.087  
 472 ( $\kappa = 0.1$ ) to 0.030 ( $\kappa = 1.0$ ).

473  
 474 Table 4: Sensitivity to staking adjustment speed  $\kappa$  (100 runs each). SE shown in parentheses  
 475 ( $\times 10^{-5}$ ).

476 $\kappa$	477 Adaptive		478 Fixed		479 Zero	
	RMSE (SE)	CV	RMSE	CV	RMSE	CV
480 0.1	0.0062 (2.1)	0.103	0.0208	0.087	0.0602	1.607
481 0.3	0.0083 (1.6)	0.062	0.0208	0.054	0.0602	1.901
482 0.5	0.0093 (1.6)	0.048	0.0208	0.043	0.0602	1.830
483 0.7	0.0091 (1.8)	0.041	0.0208	0.036	0.0602	1.723
484 1.0	0.0081 (1.7)	0.035	0.0208	0.030	0.0602	1.575

485 **Participation volatility  $\sigma_p$ .** Table 5 (100 runs per setting) shows robustness to participation  
 volatility. As expected, higher  $\sigma_p$  increases staking CV for all policies—Fixed CV increases

486 from 0.033 ( $\sigma_p = 0.01$ ) to 0.064 ( $\sigma_p = 0.10$ ). The Adaptive mechanism maintains superior  
 487 RMSE across all noise levels.  
 488

489 Table 5: Sensitivity to participation volatility  $\sigma_p$  (100 runs each). SE shown in parentheses  
 490 ( $\times 10^{-5}$ ).  
 491

$\sigma_p$	Adaptive		Fixed		Zero	
	RMSE (SE)	CV	RMSE	CV	RMSE	CV
0.01	0.0071 (1.7)	0.046	0.0208	0.033	0.0602	1.830
0.03	0.0087 (1.9)	0.045	0.0208	0.037	0.0602	1.830
0.05	0.0092 (2.1)	0.048	0.0208	0.043	0.0602	1.830
0.08	0.0093 (1.9)	0.056	0.0208	0.055	0.0602	1.829
0.10	0.0090 (2.5)	0.064	0.0208	0.064	0.0602	1.826

### 501 A.3 ROBUSTNESS ACROSS ECONOMIC REGIMES

502 Table 6 compares performance across three economic regimes: burn-dominated (baseline),  
 503 inflation-dominated (low burn, higher issuance ceiling), and high-volatility (stress test).  
 504

505 Table 6: Performance across economic regimes (20 runs). Adaptive consistently achieves  
 506 lowest target RMSE.  
 507

Regime	Policy	Target RMSE	Supply Drift	Staking CV
Burn-dominated	Adaptive	<b>0.0087</b>	+1.31%	0.048
	Fixed	0.0208	-0.05%	0.043
	Zero	0.0603	-3.97%	1.830
Inflation-dom.	Adaptive	<b>0.0026</b>	+1.73%	0.036
	Fixed	0.0200	-0.05%	0.051
	Zero	0.0300	-1.04%	1.830
High-volatility	Adaptive	<b>0.0096</b>	+1.49%	0.064
	Fixed	0.0220	-0.20%	0.066
	Zero	0.0509	-3.15%	1.827

### 521 A.4 ROBUSTNESS TO AGENT DISTRIBUTION

522 We test robustness to different opportunity cost distributions  $F$  (Table 7). The Adap-  
 523 tive mechanism consistently outperforms baselines across uniform, exponential, and beta  
 524 distributions.  
 525

526 Table 7: Performance across opportunity cost distributions (20 runs each).  
 527

Distribution	Policy	Target RMSE	Staking CV
Uniform	Adaptive	<b>0.0060</b>	0.050
	Fixed	0.0208	0.043
	Zero	0.0603	1.830
Exponential	Adaptive	<b>0.0115</b>	0.052
	Fixed	0.0208	0.049
	Zero	0.0603	1.830
Beta(2,5)	Adaptive	<b>0.0104</b>	0.059
	Fixed	0.0208	0.055
	Zero	0.0603	1.830

---

## B STATISTICAL SIGNIFICANCE TESTS

Table 8 provides complete statistical test results for all pairwise comparisons (1000 runs). All Adaptive vs. baseline comparisons are significant at  $p < 0.001$ . **Note:** The large  $t$ -statistics and Cohen’s  $d$  values are artifacts of low MC variance in simulation, not indicators of practical effect magnitude. The meaningful comparison is the raw RMSE improvement (1.9 $\times$ ).

Table 8: Statistical significance tests (Adaptive vs. baselines, 1000 runs).

Comparison	Metric	$t$ -statistic	$p$ -value	Cohen’s $d$
Adaptive vs. Fixed	Target RMSE	-666.5	$< 0.001$	-29.8
Adaptive vs. Fixed	Staking CV	37.8	$< 0.001$	1.69
Adaptive vs. Zero	Target RMSE	-3369.7	$< 0.001$	-150.8
Adaptive vs. Zero	Staking CV	-2092.6	$< 0.001$	-93.6

## C ALGORITHM DETAILS

---

### Algorithm 2 Verifiable online mechanism learning

---

- 1: **On-chain stored:**  $\theta_k, S_k$ .
  - 2: **Oracle input:**  $(p_k, L_k, u_k, b_k)$ ; compute  $e_k = (L_k - u_k + b_k)/S_k$ .
  - 3: **Act (mechanism network):** compute  $i_k = \mu_{\theta_k}(S_k, p_k, e_k)$ ; mint  $I_k = i_k S_k$ .
  - 4: **Learn (off-chain):** compute  $\nabla_{\theta} \ell_k(\theta_k)$  from loss; update  $\theta_{k+1}$  via projected SGD.
  - 5: **Prove:** generate succinct proof (SNARK/STARK) of correct update.
  - 6: **Verify/store:** verify proof on-chain; if valid store  $\theta_{k+1}$ .
- 

## D PROOF OF PROPOSITION 1

*Proof.* Define  $h(p; i) = p - F(i/p)$ . As  $p \downarrow 0$ ,  $F(i/p) \rightarrow 1$  (since  $i/p \rightarrow \infty$ ), so  $h(p; i) \rightarrow -1 < 0$ . At  $p = 1$ ,  $h(1; i) = 1 - F(i) \geq 0$  since  $F$  is a CDF. By the intermediate value theorem, there exists  $p^* \in (0, 1]$  with  $h(p^*; i) = 0$ .

For uniqueness, compute  $\partial_p h = 1 + \frac{i}{p^2} f(i/p) > 0$  since  $f \geq 0$ . Thus  $h$  is strictly increasing in  $p$ , implying uniqueness.

For the derivative, apply the implicit function theorem to  $h(p^{\text{eq}}(i); i) = 0$ :

$$\frac{d}{di} p^{\text{eq}}(i) = -\frac{\partial_i h}{\partial_p h} = \frac{\frac{1}{p} f(i/p)}{1 + \frac{i}{p^2} f(i/p)} > 0.$$

□

## E LEARNING CURVES

Figure 2 shows learning curves comparing HJB warmup vs. random initialization. With HJB warmup, the policy initializes at the HJB-Greedy performance level (RMSE  $\approx 0.026$ ) after warmup epochs. Random initialization starts higher but converges quickly. Both approaches reach similar final RMSE, consistent with the ablation results (Table 2).

## F SHOCK RESPONSE EXPERIMENTS

Table 9 presents recovery metrics across four shock scenarios. The Adaptive mechanism achieves the lowest post-shock RMSE in all scenarios except burn spike. Notably, HJB-Greedy degrades substantially under the “double shock” (burn collapse + 30% mass unstaking): RMSE increases from 0.020 (single shock) to 0.044 (double shock), while Adaptive maintains consistent performance (0.011–0.012).

594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647

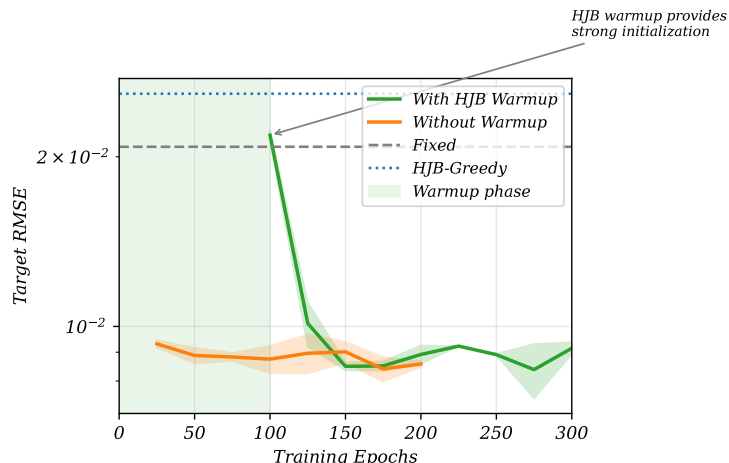


Figure 2: Learning curves: Both HJB warmup and random initialization converge to similar final RMSE.

Table 9: Shock response metrics (5 runs). Post-shock RMSE and recovery time (steps to return within 1 std of pre-shock mean).

Scenario	Policy	Post-shock RMSE	Recovery Time	Max Deviation
Burn crash	Adaptive	<b>0.0112</b>	4.4	0.031
	Fixed	0.0169	0.0	0.038
	PID	0.0408	0.0	0.080
	HJB-Greedy	0.0203	0.0	0.045
Staking exodus	Adaptive	<b>0.0101</b>	0.0	0.027
	Fixed	0.0208	0.0	0.038
	PID	0.0406	0.0	0.060
	HJB-Greedy	0.0355	0.0	0.048
Double shock	Adaptive	<b>0.0112</b>	4.6	0.031
	Fixed	0.0169	0.0	0.038
	PID	0.0408	0.0	0.080
	HJB-Greedy	0.0443	13.8	0.096
Burn spike	Adaptive	0.0121	20.6	0.101
	Fixed	0.0214	14.8	0.111
	PID	<b>0.0404</b>	0.0	0.063
	HJB-Greedy	0.0271	7.6	0.116

## G PROOF OF THEOREM 1

*Proof.* We verify the conditions of the ODE method for projected stochastic approximation (Borkar, 2009; Kushner and Yin, 2003).

**Step 1: Lipschitz continuity.** The gradient  $\nabla_{\theta} \ell_k(\theta)$  is Lipschitz in  $\theta$  since  $\mu_{\theta}$  is  $C^1$  with bounded derivatives (enforced by spectral normalization and compact  $\Theta$ ) and  $\phi$  is bounded.

**Step 2: Geometric ergodicity.** For fixed  $\theta \in \Theta$ , the state  $(S_k, p_k, e_k)$  forms a Markov chain. The supply dynamics (12) with bounded policy imply geometric ergodicity under standard conditions (bounded drift away from boundaries, volatility bounded away from zero). The staking dynamics (??) with mean-reversion also satisfy geometric ergodicity.

**Step 3: Moment bounds.** Under geometric ergodicity, the stationary distribution has finite moments. The quadratic loss  $\ell_k$  and its gradient have uniformly bounded second moments over  $\Theta$ .

---

648 **Step 4: Robbins–Monro stepsizes.** The stepsize schedule  $\alpha_k = \alpha_0/(1 + \gamma k)$  satisfies  
649  $\sum_k \alpha_k = \infty$  and  $\sum_k \alpha_k^2 < \infty$ .

650 Under these conditions, the ODE method (Borkar, 2009, Theorem 2.1) implies that the  
651 iterates  $\{\theta_k\}$  converge almost surely to the internally chain transitive invariant set of the  
652 mean ODE  $\dot{\theta} = -\nabla \bar{L}(\theta)$ . Since  $\bar{L}$  is continuously differentiable and  $\Theta$  is compact, this  
653 invariant set consists of stationary points.  $\square$   
654

## 655 H IMPLEMENTATION DETAILS

656 **Network architecture.** The mechanism network uses a 2-layer MLP with 32 hidden units,  
657 tanh activation, and spectral normalization. Input features:  $\phi(S, p, e) = (\log S/S_0, p, e, p -$   
658  $p^*, e - \bar{e})$ .

659 **Training.** Adam optimizer with initial learning rate 0.01, decay rate 0.001. Batch size 32,  
660 gradient clipping at norm 1.0. HJB warmup for 100 epochs, main training for 500 epochs.

661 **Hyperparameters.** Target inflation  $\pi^* = 0.02$ , target staking  $p^* = 0.40$ . Reward weights:  
662  $w_\pi = 1.0$ ,  $w_p = 0.5$ ,  $w_i = 0.1$ . Regularization  $\gamma = 10^{-4}$ .

663 **Simulation.** Time horizon  $T = 100$ , timestep  $\Delta t = 0.01$  (10,000 steps). Initial supply  
664  $S_0 = 10^8$ , initial staking  $p_0 = 0.30$ . Supply volatility  $\sigma_S = 0.02$ , staking volatility  $\sigma_p = 0.05$ ,  
665 adjustment speed  $\kappa = 0.5$ .  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701