

---

# DSPG: Finding Heterogeneous-Agent Equilibria with Distribution-based Differentiable Policy Gradients

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We extend the Structural Reinforcement Learning (SRL[19]) framework by em-  
2 ploying a U-Net to efficiently process the full distributional and structural infor-  
3 mation of a macroeconomic environment, solving in minutes a Heterogeneous-  
4 Agent Model that features both a non-trivial market-clearing condition and aggre-  
5 gate uncertainty—a class of macroeconomic models widely regarded as extremely  
6 challenging to solve. We formally show that our algorithm, DSPG (Distribution-  
7 based Structural Policy Gradient), obtains global nonlinear solutions without  
8 numerically relaxing the market-clearing condition or approximating aggregate  
9 shocks via perturbation. Additional experiments demonstrate that DSPG handles  
10 not only general-equilibrium but also partial-equilibrium models, and significantly  
11 outperforms conventional Deep Reinforcement Learning (DRL) baselines, includ-  
12 ing PPO, SAC, and DDPG. To our knowledge, DSPG is the first algorithm to  
13 obtain an accurate global solution of the Huggett [8] economy. Our results also  
14 confirm the validity of the SPG[19] algorithm under its dimensionality reduction  
15 and condition relaxation, whose solutions closely match ours. In the equilibrium,  
16 the simulation reproduces a negative relationship between interest rates and TFP  
17 shocks, capturing the general-equilibrium tension in a small open economy with-  
18 out domestic production, where aggregate shocks are absorbed entirely through  
19 interest-rate adjustments that clear the bond market by reconciling the competing  
20 responses of heterogeneous agents. Code is available at GitHub<sup>1</sup>.

## 21 1 Introduction

22 Modern macroeconomists use heterogeneous-agent models to study important problems such as  
23 inequality, business cycles, and optimal policy<sup>2</sup>. However, as models become more complex, the  
24 difficulty of solving them increases substantially. Two features that are computationally challenging  
25 are *aggregate uncertainty* and *non-trivial market clearing*. Huggett [8] was the first to incorporate  
26 a non-trivial market-clearing condition into a heterogeneous-agent model and solved it using binary  
27 search. Krusell and Smith [10] introduced aggregate risk into an Aiyagari [1]-style model and pro-  
28 posed a numerical solution. However, each of these seminal models incorporates only one of the two  
29 challenging features. Economists have long found it difficult to solve models that combine aggregate  
30 uncertainty with non-trivial market clearing. As Moll [12] observes, this technical bottleneck has  
31 constrained the development of macroeconomic research; the natural aspiration is that researchers  
32 should be able to specify rich, complex models without being blocked by computational solvability.

---

<sup>1</sup>The entire codebase is written in JAX, natively supporting GPU-accelerated parallel computation.

<sup>2</sup>Heterogeneous-agent models are also used to study housing, health, the environment, and many other topics. See Sargent et al. [14], Huggett [8], Aiyagari [1], Krusell and Smith [10], Kaplan et al. [9], Dietz and Venmans [4], Sun [17].

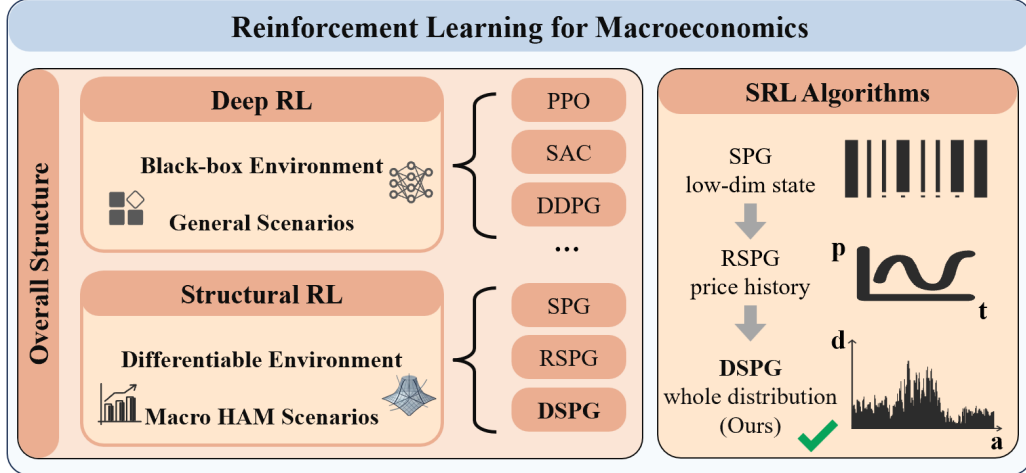


Figure 1: A general framework for reinforcement learning in macroeconomics. Conventional DRL assumes that the structural information of the environment is unknown, which facilitates generalization across diverse scenarios; in contrast, SRL exploits structural economic information in the HAM environment and improves solution efficiency by constructing a differentiable environment. The differences among SRL-style methods are as follows: (1) SPG uses a low-dimensional state space to simplify the problem; (2) RSPG augments the state with price history to strengthen the agent’s decision-making capacity; (3) DSPG uses the full distribution as the state representation for transition dynamics, preserving both computational efficiency and mathematical rigor.

33 Recently, several groups have sought to tackle complex heterogeneous-agent models. Han et al. [7]  
 34 solved the Krusell and Smith [10] model within an hour using deep neural networks and agent-based  
 35 simulation. Auclert et al. [2] solved a complex HANK [9] model with aggregate risk using first-order  
 36 perturbation to handle aggregate uncertainty. Azinovic et al. [3] proposed an equilibrium neural  
 37 network training method that exploits full distribution information to solve heterogeneous-agent  
 38 models. Gabriele et al. [5] applied several classical reinforcement learning algorithms to the Krusell  
 39 and Smith [10] model and compared their performance. Yang et al. [19] relaxed rational expectations  
 40 and proposed Structural Reinforcement Learning (SRL), solving several complex models within  
 41 minutes. SRL is efficient and easy to use; however, the structural policy gradient (SPG) algorithm in  
 42 the original paper relies on simplifying assumptions that shrink the state space and thereby constrain  
 43 agent rationality. This simplification introduces a discrepancy between the model’s theoretical state  
 44 space and the one actually used in computation, which may compromise solution accuracy in models  
 45 with richer dynamics. Wibault et al. [18] propose RSPG (RNN+SPG), which encodes the entire  
 46 price history in the state, to address this issue in the Krusell and Smith [10] setting. Nevertheless,  
 47 it remains an open question whether a model combining non-trivial market clearing with aggregate  
 48 risk can be solved with both full state-space consistency and computational efficiency.

49 Therefore, in this paper we propose DSPG, a novel method to solve the Huggett [8] economy with  
 50 aggregate risk. Despite its importance for understanding the interaction between macroeconomic un-  
 51 certainty (i.e. business cycles and TFP shocks) and non-trivial market clearing (supply matching de-  
 52 mand), this class of models remains difficult to solve with existing computational techniques. DSPG  
 53 extends SPG by incorporating the full wealth distribution into the state representation, retaining com-  
 54 putational efficiency while substantially improving equilibrium accuracy. Unlike DEQN [3], which  
 55 also conditions on distribution information but relies on a different training paradigm, DSPG embeds  
 56 the entire distribution in the state space within a differentiable policy-gradient framework, remaining  
 57 computationally tractable. To handle the resulting high-dimensional inputs efficiently, we employ a  
 58 one-dimensional U-Net—first introduced by Ronneberger et al. [13] for image segmentation—to pa-  
 59 rameterize agents’ policy functions. In the Huggett [8] economy where non-trivial market clearing  
 60 and aggregate risk coexist, DSPG recovers the optimal consumption–saving policy within minutes  
 61 and achieves market-clearing errors below  $10^{-7}$  on a server with a single NVIDIA RTX 3090 GPU.

62 Further experiments show that DSPG is consistent with SPG while reducing average market-clearing  
 63 error significantly. To our knowledge, DSPG is the first method to obtain an accurate global non-

64 linear solution to the Huggett [8] economy in the simultaneous presence of aggregate uncertainty  
65 and non-trivial market clearing. Compared with classical DRL baselines<sup>3</sup> on a simplified bench-  
66 mark environment, DSPG exhibits strong stability and efficiency<sup>4</sup>. Figure 1 summarizes how DSPG  
67 relates to existing DRL and SRL algorithms. Together with this paper, we release two benchmark  
68 models<sup>5</sup> for validating reinforcement learning algorithms in this setting. We believe that, owing  
69 to its efficiency and accuracy, DSPG can serve as an effective computational paradigm for solving  
70 heterogeneous-agent macroeconomic models.

## 71 2 Preliminaries

72 **Structural Reinforcement Learning** First introduced by Yang et al. [19] and later extended by  
73 Wibault et al. [18], SRL solves heterogeneous-agent models by combining RL with structural knowl-  
74 edge of the economic environment. Consider agents with individual state  $s_t \in \mathcal{S}$  who choose actions  
75  $a_t \in \mathcal{A}^6$  according to a policy  $\pi_\theta(a_t|s_t, p_t)$ . The agent receives per-period payoff

$$\mathcal{R}_t = \mathcal{R}(s_t, a_t, p_t, z_t),$$

76 where  $z_t$  denotes aggregate shocks and  $p_t$  equilibrium prices. Individual states evolve according to  
77 known structural dynamics, usually a budget constraint,

$$s_{t+1} = f(s_t, a_t, z_t, p_t),$$

78 and agents maximize expected discounted returns<sup>7</sup>

$$J(\theta) = \lim_{T \rightarrow \infty} \mathbb{E} \left[ \sum_{t=0}^T \beta^t \mathcal{R}(s_t, a_t, p_t, z_t) \right].$$

79 In heterogeneous-agent models, prices depend on the cross-sectional distribution of agents, which  
80 makes the environment high-dimensional. SRL addresses this by simulating equilibrium trajectories  
81 and directly optimizing the policy parameters  $\theta$ . A key feature of SRL is that the individual transi-  
82 tion function  $f(\cdot)$  is known and differentiable. This makes the environment partially differentiable,  
83 allowing gradients of  $J(\theta)$  to be computed by differentiating through the transition dynamics. As a  
84 result, policy parameters can be optimized using gradient-based methods rather than relying solely  
85 on likelihood-ratio policy gradients.

## 86 3 A Huggett Economy with Aggregate Uncertainty

87 **Households** There is a continuum of households indexed by  $i \in [0, 1]$ . Each household earns  
88 income  $z_t e_{i,t}$  at time  $t$ , consisting of an aggregate component  $z_t$  and an idiosyncratic component  
89  $e_{i,t}$ , both following Markov processes

$$z_{t+1} \sim \mathcal{T}_z(\cdot | z_t) \quad \text{and} \quad e_{i,t+1} \sim \mathcal{T}_e(\cdot | e_{i,t}). \quad (1)$$

90 The wealth  $a_{i,t}$  of household  $i$  at time  $t$  evolves according to

$$a_{i,t+1} + c_{i,t} = (1 + r_t) a_{i,t} + z_t e_{i,t}, \quad (2)$$

91 where  $c_{i,t}$  is consumption and  $r_t$  is the interest rate. Households also face a borrowing constraint

$$a_{i,t} \geq \underline{a}, \quad (3)$$

92 where  $\underline{a} \leq 0$ . Each household chooses consumption  $c_{i,t}$  to maximize lifetime expected utility

$$\max_{c_{i,t}} \lim_{T \rightarrow \infty} \mathbb{E} \left[ \sum_{t=0}^T \beta^t u(c_{i,t}) \right], \quad (4)$$

93 subject to (1), (2), and (3), taking as given the interest rate  $r_t$ .

<sup>3</sup>Including PPO [15], SAC [6], and DDPG [16]. A more detailed comparative study of classical DRL methods on the Krusell and Smith [10] HAM appears in TaxAI [11].

<sup>4</sup>These gains stem largely from the differentiable environment and the convexity of the optimization problem—features common to SRL methods. Note that although DSPG, RSPG and SPG differ algorithmically, they are all SRL methods.

<sup>5</sup>GE and PE versions of the dynamic Huggett [8] economy.

<sup>6</sup>Here  $a_t$  denotes a generic action variable, not to be confused with asset holdings  $a_{i,t}$  defined in Section 3.

<sup>7</sup>In practice, we set a relatively large  $T$  cutoff so that the program can complete running.

94 **Market** We consider a closed economy with exogenous foreign demand for asset in which the  
 95 domestic bond market interacts with an exogenous foreign sector. Each period, the foreign sector  
 96 demands a fixed volume  $B > 0$  of bonds from domestic households<sup>8</sup>—this can be interpreted as  
 97 external sovereign or institutional borrowing, reflecting a constant foreign appetite for domestic safe  
 98 assets—and repays principal plus interest at the equilibrium rate  $r_t$ . Here  $r_t$  is determined by a non-  
 99 trivial market-clearing condition following Huggett [8]: it adjusts until aggregate domestic saving  
 100 absorbs the foreign borrowing requirement,

$$\int_0^1 a_{i,t+1} di = B, \quad \forall t \geq 0. \quad (5)$$

101 Domestic households are price-takers with respect to both foreign demand and the interest rate.  
 102 Since  $B$  is fixed and exogenous, the equilibrium  $r_t$  is fully pinned down by the cross-sectional distri-  
 103 bution of domestic wealth—the hallmark of non-trivial market clearing in the Huggett [8] tradition.

104 **Equilibrium** Since the aggregate shock  $z_t$  renders the cross-sectional distribution non-degenerate  
 105 as a payoff-relevant state, we adopt a recursive formulation. Let  $\mu$  denote the joint distribution of  
 106 individual states  $(a, e)$  across households, so that  $(z, \mu)$  constitutes the aggregate state. A *recursive*  
 107 *competitive equilibrium* consists of household policies  $c(a, e; z, \mu)$ ,  $a'(a, e; z, \mu)$ , a pricing func-  
 108 tion  $r(z, \mu)$ , and a law of motion  $\mu' = H(\mu, z, z')$  such that: (i) given  $r(\cdot)$  and  $H(\cdot)$ , the policies  
 109 solve the household problem (4) subject to (2), (3), and (1)—equivalently, the value function satis-  
 110 fies the Bellman optimality condition  $V(a, e; z, \mu) = \max_{a'} \{u(c) + \beta \mathbb{E}[V(a', e'; z', \mu') \mid e, z]\}$ ;  
 111 (ii) price  $r(z, \mu)$  clears the bond market (5); and (iii)  $H$  is induced by  $a'(\cdot)$  and (1); for measurable  
 112 subsets  $\mathcal{B}, \mathcal{E}'$  of the asset and efficiency state spaces,  $\mu'(\mathcal{B} \times \mathcal{E}') = \int \mathcal{T}_e(\mathcal{E}' \mid e) \mathbf{1}\{a'(a, e; z, \mu) \in$   
 113  $\mathcal{B}\} d\mu(a, e)$ , under which the joint process  $(z_t, \mu_t)$  admits a unique ergodic distribution.<sup>9</sup>

114 What distinguishes this model from the original Huggett [8] economy is the presence of aggregate  
 115 risk  $z_t$ , which renders the equilibrium interest rate  $r_t$  time-varying and only *approximately Markov*<sup>10</sup>.  
 116 The economy settings ensure that no domestic production side absorbs the aggregate shock; instead,  
 117  $r_t$  must adjust endogenously to clear the bond market given the entire cross-sectional distribution,  
 118 which is the central computational challenge.<sup>11</sup>

## 119 4 Our Method

### 120 4.1 Economic Dynamics Reformulation

121 **Global solution** Since aggregate shock  $z_t$  affects income level  $z_t e_{i,t}$  through (2), the optimal  
 122 policy of households depends on the cross-sectional distribution and the current aggregate state. In  
 123 general, the optimal policy may depend on the entire history of distributions; however, when the  
 124 aggregate shock  $z_t$  is Markov and the distribution  $\mathbf{g}_t$  is a sufficient statistic for the cross-sectional  
 125 state, we define the policy as a function of  $(\mathbf{g}_t, z_t)$  alone. Since  $\mathbf{g}_t$  lives on a finite grid of dimension  
 126  $n_a \times n_e$ , this yields a well-defined (though high-dimensional) state representation

$$c_t^* = \pi(\mathbf{g}_t, z_t) \quad \text{where} \quad \mathbf{g}_t = \begin{bmatrix} d(a_1, e_1) & \cdots & d(a_1, e_{n_e}) \\ \cdots & & \cdots \\ d(a_{n_a}, e_1) & \cdots & d(a_{n_a}, e_{n_e}) \end{bmatrix} \quad \text{and} \quad \sum_{i=1}^{n_a} \sum_{j=1}^{n_e} d(a_i, e_j) = 1, \quad (6)$$

127 where  $c_t^*$  denotes the optimal consumption policy evaluated on the full grid (yielding the matrix  
 128  $\mathbf{c}_t \in \mathbb{R}^{n_a \times n_e}$  defined in (10)),  $\mathbf{g}_t$  the distribution,  $d(\cdot) \in (0, 1)$  the density of distribution,  $\mathbf{a} =$   
 129  $\{a_1, \dots, a_{n_a}\}$  the wealth grid and  $\mathbf{e} = \{e_1, \dots, e_{n_e}\}$  the income grid. See Appendix A.2 for more  
 130 details of discretization. (6) requests distribution  $\mathbf{g}_t$  and aggregate shock  $z_t$  both in the households'  
 131 state space to get the optimal global solution. Therefore, the decision variable of household is high-  
 132 dimensional which makes the problem hard to solve.

<sup>8</sup>Note that this differs from the original Huggett setup where  $B = 0$ . This is to make the model definition cleaner; however, these two setups are entirely equivalent. When  $\underline{a} < 0$ , one can redefine  $\tilde{a} = a - \underline{a}$  as the distance from the borrowing constraint, and then the model would revert to our definition.

<sup>9</sup>Our method targets the policy and pricing functions evaluated along this ergodic measure. In the notation of Section 2, the individual state is  $s = (a, e)$ , the action is  $c$ , the equilibrium price is  $p \equiv r$ , and the per-period payoff is  $\mathcal{R} = u(c)$ .

<sup>10</sup>In the Krusell and Smith [10] economy, prices are approximately AR(1), since tomorrow's prices  $(r_{t+1}, w_{t+1})$  can be predicted accurately from current aggregate capital  $K_t$  and aggregate state  $z_t$  alone; see Appendix B.1 for the model and Appendix B.2 for the solution method.

<sup>11</sup>See Appendix A.1 for the model parameters.

133 **Non-trivial market clearing** From (2) and (5), we know that the aggregate saving is

$$S_t(\mathbf{g}_t, z_t, r_t) = \int_0^1 a_{i,t+1} di = \int_0^1 [(1+r_t)a_{i,t} + z_t e_{i,t} - c_{i,t}] di = B, \quad (7)$$

134 where  $S_t$  is the saving supply curve. To express (7) in vector form, we define the following no-  
 135 tation on the discretized economy: let  $\mathbf{a} \in \mathbb{R}^{n_a}$  denote the wealth grid vector,  $\mathbf{g}_t \in \mathbb{R}^{n_a \times n_e}$   
 136 the distribution matrix, and  $\mathbf{c}_t \in \mathbb{R}^{n_a \times n_e}$  the consumption policy matrix. We further write  
 137  $\mathbf{a} \odot \mathbf{g}_t \equiv \sum_{i=1}^{n_a} \sum_{j=1}^{n_e} a_i \cdot d(a_i, e_j)$  for the cross-sectional mean of wealth (element-wise multipli-  
 138 cation followed by summation, broadcasting  $\mathbf{a}$  across income states);  $\mathbf{c}_t \odot \mathbf{g}_t$  is defined analogously  
 139 for consumption. We adopt this notation throughout the paper. Then we have

$$\mathbf{a} \odot \mathbf{g}_{t+1} = S_t(\mathbf{g}_t, z_t, r_t) = (1+r_t) \cdot \mathbf{a} \odot \mathbf{g}_t + z_t \cdot \bar{e} - \mathbf{c}_t \odot \mathbf{g}_t = B, \quad (8)$$

140 where  $\bar{e}$  is the ergodic mean of idiosyncratic income  $e_{i,t}$  and  $\mathbf{c}_t$  is the consumption policy vector  
 141 over distribution. Then, we can directly get the interest rate that clears the market

$$1+r_t = \frac{B - z_t \cdot \bar{e} + \mathbf{c}_t \odot \mathbf{g}_t}{\mathbf{a} \odot \mathbf{g}_t}. \quad (9)$$

142 Equation (9) shows that the market-clearing interest rate  $r_t$  can be obtained in closed form from  
 143 a single formula<sup>12</sup>. This is possible because the full cross-sectional distribution  $\mathbf{g}_t$  is included in  
 144 the state space: given  $\mathbf{g}_t$  and the policy  $\mathbf{c}_t$ , all aggregate quantities required for market clearing are  
 145 directly computable, so  $r_t$  is determined without any iterative root-finding or fixed-point loop.

## 146 4.2 Distribution-based Structural Policy Gradient (DSPG)

147 **U-net policy function** From (6), the state space (input) is a distribution of dimension  $n_a \times n_e$  and  
 148 the policy (output) is the consumption value at each grid point. We parameterize this mapping with  
 149 a 1D U-Net for three reasons: (i) the encoder–decoder structure captures both local interactions  
 150 between neighboring wealth grid points and global distributional features through progressively  
 151 coarser representations; (ii) skip connections preserve fine-grained grid-level information that would  
 152 otherwise be lost in deep networks, which is critical for satisfying the borrowing constraint at spe-  
 153 cific grid points; (iii) the convolutional architecture exploits the ordered structure of the wealth grid,  
 154 sharing parameters across grid points—a fully connected (MLP) network would require  $\mathcal{O}((n_a n_e)^2)$   
 155 parameters in its first layer alone, making it computationally prohibitive for high-dimensional grids,  
 156 whereas the U-Net scales linearly in  $n_a n_e$ . The policy network is thus

$$\mathbf{c}_t = U\text{-net}(\mathbf{g}_t, z_t; \Theta) = \begin{bmatrix} c(a_1, e_1) & \cdots & c(a_1, e_{n_e}) \\ \cdots & & \cdots \\ c(a_{n_a}, e_1) & \cdots & c(a_{n_a}, e_{n_e}) \end{bmatrix}, \quad (10)$$

157 where  $\Theta$  is the parameters of the 1D U-net and the detailed structure is in Appendix E.

158 **Differentiable learning** Unlike classical reinforcement learning algorithms, we create a differen-  
 159 tiable environment where agents know the analytical policy gradient of the cumulative utility (4),  
 160 and directly optimize the objective function by using the loss function

$$\mathbb{J}(\Theta) = - \int_0^1 \mathbb{E} \left[ \sum_{t=0}^{\infty} \beta^t u(c_{i,t}) \right] di \approx - \frac{1}{n} \sum_{m=1}^n \sum_{t=0}^T \beta^t \int_0^1 u(c_{i,t}) di, \quad (11)$$

161 where the approximation truncates the infinite horizon to  $T$  periods and replaces the expectation  
 162 with an empirical average over  $n$  sampled trajectories. (11) can also be written in vector form

$$\mathbb{J}(\Theta) \approx - \frac{1}{n} \sum_{m=1}^n \sum_{t=0}^T \beta^t u(\mathbf{c}_t) \odot \mathbf{g}_t, \quad (12)$$

163 which is efficient when using JAX and GPU acceleration in practice. Empirically, we observe that  
 164 the loss landscape of (12) behaves as a convex optimization problem: training consistently converges

<sup>12</sup>Note that this formula is a price update formula, not a price definition formula. It represents the clearing price for the current period. Of course, under equilibrium, it is the equilibrium price.

165 to the same solution from different random initializations, which is consistent with the concavity of  
 166  $u(\cdot)$  and the linearity of the budget constraint.

167 Every step of the DSPG computation graph—U-Net evaluation (10), closed-form market clearing  
 168 (9), the budget-constraint wealth update, and the distributional transition—is a differentiable func-  
 169 tion of its inputs. Conditional on an exogenous shock realization  $\{z_t\}$ , reverse-mode automatic  
 170 differentiation therefore yields  $\nabla_{\Theta} \mathbb{J}(\Theta)$  *exactly*, with no score-function (REINFORCE) term (Ap-  
 171 pendix F, Proposition F.3). The only stochasticity in our estimator comes from averaging over  $n$   
 172 independent shock trajectories, so its variance decays at the standard  $O(1/n)$  rate (Proposition F.5);  
 173 truncating the horizon at  $T$  introduces a bias bounded by  $M\beta^{T+1}/(1-\beta)$  (equation (36)), which is  
 174 negligible for the horizons used in our experiments.

175 **Monotonicity-preserving parameterization** In consumption-saving models, the optimal con-  
 176 sumption policy is monotonically increasing in wealth. To enforce this structural property and  
 177 prevent the network architecture from producing non-monotone policies, we parameterize the first-  
 178 order difference rather than the level of consumption:

$$U\text{-net}(\mathbf{g}_t, z_t; \Theta) = \begin{bmatrix} c(a_1, e_1) & \cdots & c(a_1, e_{n_e}) \\ \delta(a_2, e_1) & \cdots & \delta(a_2, e_{n_e}) \\ \vdots & \ddots & \vdots \\ \delta(a_{n_a}, e_1) & \cdots & \delta(a_{n_a}, e_{n_e}) \end{bmatrix} \quad \text{and} \quad c(a_i, e_j) = c(a_1, e_j) + \sum_{k=2}^i \delta(a_k, e_j), \quad (13)$$

179 where  $\delta(a_i, e_j) = c(a_i, e_j) - c(a_{i-1}, e_j) \geq 0$  with  $i \geq 2$  and  $j \geq 1$  is the first-step difference  
 180 of the policy vector along wealth grids. Note that (13) is equivalent to (10). In order to ensure  
 181 stability at the beginning of the training, we firstly solve a Huggett [8] model without aggregate  
 182 uncertainty (usually called steady state and easy to solve) and pre-train the U-net to fit the optimal  
 183 policy function of the steady state. Figure 2 shows the steady state results, where the policy is used  
 184 to initialize the network and the distribution is used to initialize the simulation.

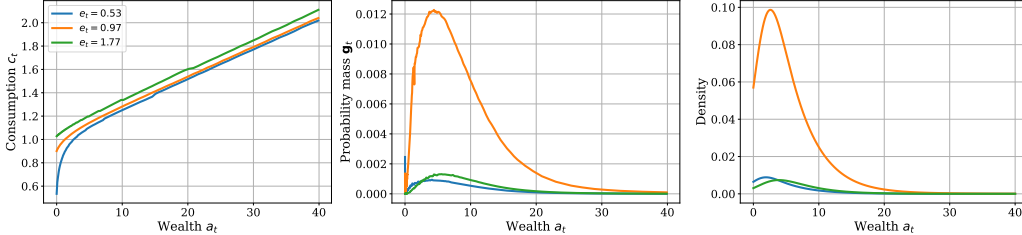


Figure 2: Steady-state solution of the Huggett [8] economy. **Left:** optimal consumption policy as a function of wealth for each idiosyncratic income level. **Middle:** probability mass function of the ergodic wealth distribution. **Right:** kernel density estimate (KDE) of the same distribution.

185 **Relationships to SPG and RSPG** All three algorithms—SPG, RSPG, and DSPG—belong to the  
 186 SRL family and share a differentiable environment for variance-free gradient computation. They  
 187 differ in the state representation fed to the policy function:

- 188 • **SPG** [19] uses a low-dimensional state space  $(a_t, e_t, r_t, z_t)$ , which is a summary of the  
 189 distribution  $\mu_t$ , being computational friendly.
- 190 • **RSPG** [18] augments the state space with a history of prices  $\{r_{t-\ell}, \dots, r_t\}$  encoded by an  
 191 RNN, partially recovering distributional information through price signals.
- 192 • **DSPG** directly conditions on the full distribution  $\mu_t$ , preserving all cross-sectional informa-  
 193 tion at the cost of a higher-dimensional input, which we handle via the U-Net architecture.

194 As shown in Section 5, DSPG retains SRL-level efficiency—solving the model within minutes—  
 195 while targeting the exact global equilibrium; the close agreement with SPG solutions in turn vali-  
 196 dates SPG’s approximation. Beyond SRL, DSPG produces global nonlinear solutions (unlike per-  
 197 turbation methods [2]), optimizes through a differentiable simulation rather than an equilibrium net-  
 198 work (unlike DEQN [3]), and eliminates score-function variance relative to classical DRL baselines  
 199 (Appendix F). A limitation is the fixed-grid cost  $O(n_a n_e)$ : for very high-dimensional distributions,

200 backpropagation through  $T$  steps may exhaust GPU memory regardless of the network architecture,  
 201 motivating future work on adaptive distributional representations.

### 202 4.3 Algorithm Summary

---

#### Algorithm 1 DSPG: Distribution-based Structural Policy Gradient

---

- 1: **Phase 1: Steady-state pre-training**
  - 2: Solve the steady-state Huggett [8] economy (no aggregate risk) via Value Iteration
  - 3: Pre-train U-Net parameters  $\Theta_0$  to fit the steady-state consumption policy
  - 4: Initialize distribution  $\mathbf{g}_0$  from the steady-state distribution
  - 5: **Phase 2: Policy optimization**
  - 6: **for** iteration  $k = 1, 2, \dots, K$  **do**
  - 7:     **for** trajectory  $m = 1, \dots, n$  **do**
  - 8:         Sample aggregate shock sequence  $\{z_t\}_{t=0}^T$  from  $\mathcal{T}_z$
  - 9:         **for**  $t = 0, 1, \dots, T$  **do**
  - 10:             Compute consumption policy:  $\mathbf{c}_t = \text{U-Net}(\mathbf{g}_t, z_t; \Theta_k)$  ▷ Eq. (10)
  - 11:             Compute equilibrium interest rate  $r_t$  ▷ Eq. (9)
  - 12:             Update distribution  $\mathbf{g}_{t+1}$  via budget constraint and idiosyncratic transition  $\mathcal{T}_e$
  - 13:         **end for**
  - 14:     **end for**
  - 15:     Compute loss  $\mathbb{J}(\Theta_k)$  by averaging over  $n$  trajectories ▷ Eq. (12)
  - 16:     Update  $\Theta_{k+1} \leftarrow \Theta_k - \alpha \nabla_{\Theta} \mathbb{J}(\Theta_k)$  via backpropagation through differentiable dynamics
  - 17: **end for**
  - 18: **return** Trained policy parameters  $\Theta_K$
- 

203 **Implementation details** The entire pipeline is implemented in JAX for GPU-accelerated differentiable  
 204 computation. All training hyperparameters—including learning rate, number of trajectories  $n$ ,  
 205 and truncation length  $T$ —are reported in Appendix D.1. The 1D U-Net architecture details (number  
 206 of levels, channel widths) are provided in Appendix E; model calibration parameters are in Appen-  
 207 dices A.1. Each forward pass evaluates the U-Net and computes  $r_t$  in  $\mathcal{O}(n_a n_e)$ ; the backward  
 208 pass through  $T$  time steps has complexity  $\mathcal{O}(T \cdot n_a n_e \cdot |\Theta|)$ . Memory scales linearly with  $T$  due to  
 209 storing intermediate activations for backpropagation.

## 210 5 Experiments

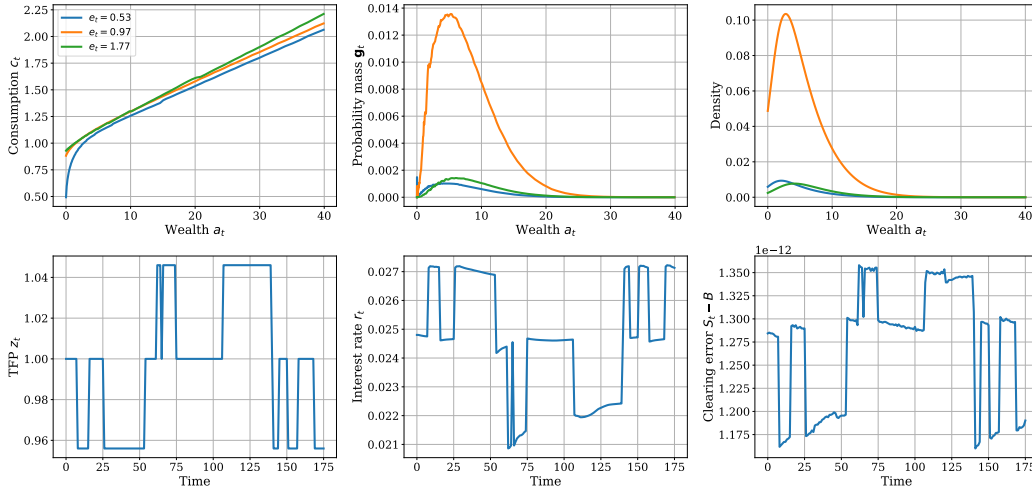


Figure 3: DSPG solution of the Huggett [8] economy with aggregate risk. **Top row:** optimal consumption policy, probability mass, and KDE density of the ergodic distribution. **Bottom row:** a simulated path of aggregate shock  $z_t$ , equilibrium interest rate  $r_t$ , and market clearing error  $S_t - B$ .

211 **5.1 Optimal Policy**

212 Figure 2 shows the steady-state solution used for initialization. Figure 3 presents the DSPG solution  
 213 of the full model with aggregate risk: the policy and ergodic distribution closely match those of the  
 214 steady-state economy under the same calibration, while producing even smoother curves than value  
 215 iteration. The simulation path reveals a negative relationship between the equilibrium interest rate  $r_t$   
 216 and TFP  $z_t$ , and aggregate saving residuals remain below  $10^{-10}$  throughout—both consistent with  
 217 economic theory.

218 We also ablate the number of trajectories  $n$  per update (Figure 4). Increasing  $n$  from 32 to 512  
 219 progressively reduces training variance at the cost of longer wall-clock time;  $n = 128$  offers a  
 220 practical balance. All three settings converge to comparable final loss values. Other hyperparameters  
 221 are in Appendix D.1.

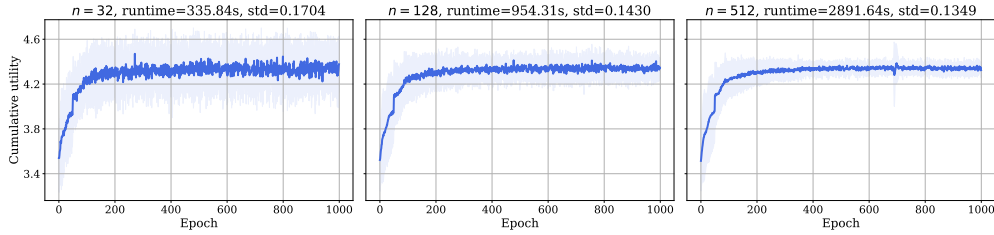


Figure 4: Ablation on the number of trajectories  $n \in \{32, 128, 512\}$  (10 independent runs each).  
**Left:** training loss curves (solid line: mean; shaded:  $\pm 1$  std). **Middle:** wall-clock time per run.  
**Right:** final training variance. Larger  $n$  reduces variance but increases computation.

222 Together, these results demonstrate that DSPG solves the Huggett [8] economy with aggregate risk  
 223 and positive bond supply in a stable, accurate, and robust manner.

224 **5.2 Economic Findings**

225 Table 1 reports the equilibrium interest rate  $r_t$  conditional on each TFP level  $z_t$ , revealing a clear  
 226 negative relationship. The mechanism is intuitive: classify households with  $a_{i,t} > B$  as net savers  
 227 and those with  $a_{i,t} < B$  as net borrowers. When TFP is high, all households’ income rises, increas-  
 228 ing the aggregate supply of savings relative to borrowing demand; the interest rate falls to restore  
 229 market clearing. Conversely, when TFP is low, saving supply contracts and the interest rate must rise  
 230 to attract sufficient lending. In short,  $r_t$  adjusts to reconcile the shifting balance between aggregate  
 231 saving and borrowing, generating the observed negative co-movement with  $z_t$ .

Table 1: Equilibrium interest rate  $r_t$  vs TFP shock  $z_t$

TFP shock $z_t$	0.91	0.96	1.00	1.05	1.09
AVG $r_t$	0.0283	0.0267	0.0250	0.0222	0.0188
STD $r_t$	0.0007	0.0006	0.0008	0.0011	0.0016

232 **5.3 Comparison to SPG**

233 Figure 5 compares the mean ergodic distributions produced by DSPG and SPG<sup>13</sup>. The DSPG distri-  
 234 bution is very close to that of SPG yet noticeably smoother, reflecting the continuity of the U-Net  
 235 policy in contrast to SPG’s grid-based interpolation. Figure 6 shows that the aggregate consumption  
 236 paths  $C_t$  also agree closely. Quantitatively, DSPG achieves average market-clearing residuals below  
 237  $10^{-11}$ , compared with roughly  $10^{-5}$  for SPG—an improvement of six orders of magnitude. To our  
 238 knowledge, this is the first high-accuracy global solution of this economy.

<sup>13</sup>SPG [19] was previously the only effective algorithm for this model and yields solutions close to the rational-expectations equilibrium under its bounded-rationality assumptions, making it a natural reference.

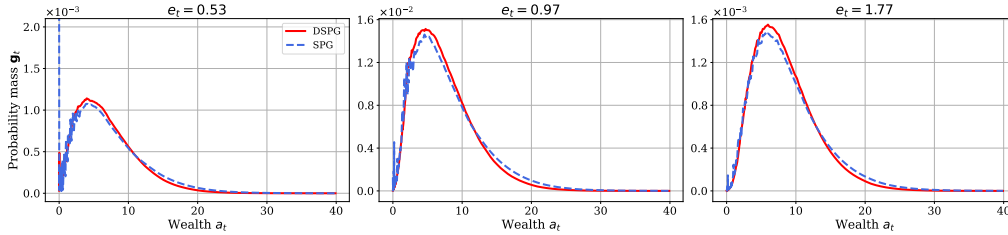


Figure 5: Mean ergodic distribution  $\mathbf{g}_t$ : DSPG (ours) vs. SPG. The two distributions are nearly indistinguishable, with DSPG producing smoother curves.

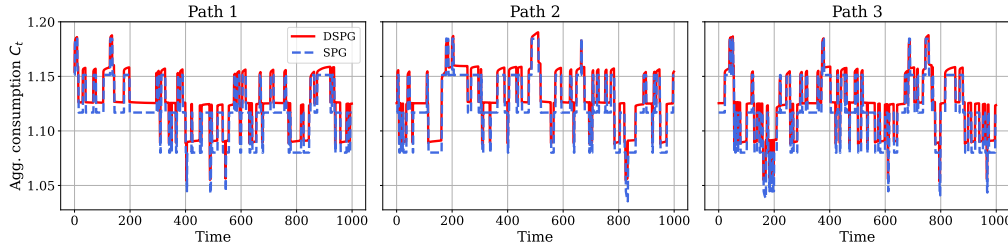


Figure 6: Aggregate consumption  $C_t$ : DSPG (ours) vs. SPG. Under three fixed representative TFP shocks, the total consumption paths from both methods closely agree.

## 239 5.4 Validation

240 We compare DSPG with conventional DRL algorithms—PPO, SAC, and DDPG—on the partial-  
 241 equilibrium (PE) version of our economy, where a value-function iteration (VFI) solution serves as  
 242 the ground-truth baseline (see Appendix C.1 for the PE setup, Appendix C.2 for the VFI solution,  
 243 and Appendix C.3 for the rationale of using the PE benchmark). Table 2 and Figure 7 show that  
 244 DSPG converges faster, with lower variance and closer final performance to VFI than all three DRL  
 245 baselines. Hyperparameters for each method are reported in Appendices D.2–D.5.

Table 2: Summary of performance

Algorithm	Last eval		Best eval	
	AVG	STD	AVG	STD
<b>DSPG</b>	<b>5.661</b>	<b>0.003</b>	5.732	<b>0.001</b>
PPO	5.128	0.132	<b>5.910</b>	0.067
SAC	4.266	0.060	4.716	0.019
DDPG	4.502	0.501	5.225	0.221
VFI	5.653	—	5.653	—

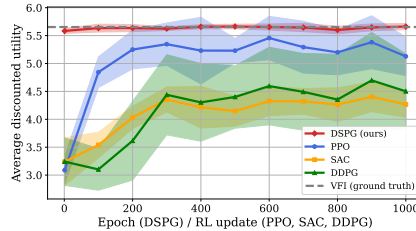


Figure 7: Training curves on the PE model (10 runs; shaded:  $\pm 1.96$  std).

## 246 6 Conclusion

247 We introduced DSPG, an efficient distribution-based structural policy gradient method that solves  
 248 heterogeneous-agent models with both aggregate uncertainty and non-trivial market clearing. By  
 249 conditioning a U-Net policy on the full cross-sectional distribution and computing exact path-  
 250 wise gradients through a differentiable simulation, DSPG obtains global nonlinear equilibria within  
 251 minutes—without linearizing around a steady state or relaxing the market-clearing condition. On the  
 252 dynamic Huggett [8] economy with aggregate risk, DSPG achieves market-clearing residuals below  
 253  $1e-7$ , improving upon prior SRL methods by two orders of magnitude. To our knowledge, this is the  
 254 first high-accuracy global solution of this economy. Future work includes extending DSPG to richer  
 255 environments such as HANK models, improving scalability via adaptive or mesh-free distributional  
 256 representations, and establishing formal convergence guarantees under NN parameterization.

257 **References**

- 258 [1] S Rao Aiyagari. Uninsured idiosyncratic risk and aggregate saving. *The Quarterly Journal of*  
259 *Economics*, 109(3):659–684, 1994.
- 260 [2] Adrien Auclert, Bence Bardóczy, Matthew Rognlie, and Ludwig Straub. Using the sequence-  
261 space Jacobian to solve and estimate heterogeneous-agent models. *Econometrica*, 89(5):2375–  
262 2408, 2021.
- 263 [3] Marlon Azinovic, Luca Gaegauf, and Simon Scheidegger. Deep equilibrium nets. *International*  
264 *Economic Review*, 2022.
- 265 [4] Simon Dietz and Frank Venmans. Cumulative carbon emissions and economic policy: in  
266 search of general principles. *Journal of Environmental Economics and Management*, 96:108–  
267 129, 2019.
- 268 [5] Federico Gabriele, Aldo Glielmo, and Marco Taboga. Heterogeneous rbcs via deep multi-agent  
269 reinforcement learning. *arXiv preprint arXiv:2510.12272*, 2025.
- 270 [6] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-  
271 policy maximum entropy deep reinforcement learning with a stochastic actor. In *International*  
272 *conference on machine learning*, pages 1861–1870. Pmlr, 2018.
- 273 [7] Jiequn Han, Yucheng Yang, and Weinan E. DeepHAM: A global solution method for hetero-  
274 geneous agent models with aggregate shocks. *arXiv preprint arXiv:2112.14377*, 2021.
- 275 [8] Mark Huggett. The risk-free rate in heterogeneous-agent incomplete-insurance economies.  
276 *Journal of Economic Dynamics and Control*, 17(5-6):953–969, 1993.
- 277 [9] Greg Kaplan, Benjamin Moll, and Giovanni L Violante. Monetary policy according to HANK.  
278 *American Economic Review*, 108(3):697–743, 2018.
- 279 [10] Per Krusell and Anthony A Smith, Jr. Income and wealth heterogeneity in the macroeconomy.  
280 *Journal of Political Economy*, 106(5):867–896, 1998.
- 281 [11] Qirui Mi, Siyu Xia, Yan Song, Haifeng Zhang, Shenghao Zhu, and Jun Wang. Taxai: A  
282 dynamic economic simulator and benchmark for multi-agent reinforcement learning. *arXiv*  
283 *preprint arXiv:2309.16307*, 2023.
- 284 [12] Benjamin Moll. The trouble with rational expectations in heterogeneous agent models: A  
285 challenge for macroeconomics. *Economic Journal Lecture, Royal Economic Society*, 2024.
- 286 [13] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for  
287 biomedical image segmentation. In *International Conference on Medical image computing*  
288 *and computer-assisted intervention*, pages 234–241. Springer, 2015.
- 289 [14] Thomas J Sargent, Christopher A Sims, et al. Business cycle modeling without pretending to  
290 have too much a priori economic theory. *New methods in business cycle research*, 1:145–168,  
291 1977.
- 292 [15] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal  
293 policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 294 [16] David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Ried-  
295 miller. Deterministic policy gradient algorithms. In *International conference on machine*  
296 *learning*, pages 387–395. Pmlr, 2014.
- 297 [17] Jeffrey Sun. The distributional consequences of climate change: The role of housing wealth,  
298 expectations, and uncertainty. Technical report, Princeton University, 2023.
- 299 [18] Clarisse Wibault, Johannes Forkel, Sebastian Towers, Tiphaine Wibault, Juan Duque, George  
300 Whittle, Andreas Schaab, Yucheng Yang, Chiyuan Wang, Michael Osborne, et al. Recur-  
301 rent structural policy gradient for partially observable mean field games. *arXiv preprint*  
302 *arXiv:2602.20141*, 2026.
- 303 [19] Yucheng Yang, Chiyuan Wang, Andreas Schaab, and Benjamin Moll. Structural reinforcement  
304 learning for heterogeneous agent macroeconomics, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2512.18892)  
305 [2512.18892](https://arxiv.org/abs/2512.18892).

306 **Appendix**

307 **A Details of the Huggett Economy**

308 **A.1 Calibration**

309 Following the literature, we set the household discount factor  $\beta = 0.975$  and utility function  $u(c) =$   
 310  $\frac{c^{1-\sigma}}{1-\sigma}$  with  $\sigma = 1$ . Then, we set borrowing constraint  $\underline{a} = 0$  and total bond supply  $B = 5$ , a setting  
 311 that is different from the original [8] model. The idiosyncratic shock follows a three-state Markov  
 312 chain to approximate the log AR(1) process

$$\log(e_{i,t+1}) = \rho_e \cdot \log(e_{i,t}) + \nu_e \cdot \sqrt{1 - \rho_e^2} \cdot \epsilon_{i,t}, \quad \epsilon_{i,t} \sim N(0, 1), \quad (14)$$

313 where the persistence  $\rho_e = 0.6$  and the variance  $\nu_e = 0.2$ , same as [1] model. The grids  $\mathbf{e}$  and  
 314 transition matrix  $\mathbf{P}_e$  of idiosyncratic shock are

$$\mathbf{e} = \begin{bmatrix} 0.5343 \\ 0.9735 \\ 1.7739 \end{bmatrix} \quad \text{and} \quad \mathbf{P}_e = \begin{bmatrix} 0.6460 & 0.3539 & 0.0001 \\ 0.0304 & 0.9392 & 0.0304 \\ 0.0001 & 0.3539 & 0.6460 \end{bmatrix}, \quad (15)$$

315 with  $\mathbf{e}$  and  $\mathbf{P}_e$  rounded to 4 decimals and the ergodic mean average  $\bar{e}$  is 1.

316 The aggregate shock follows a five-state Markov chain to approximate the log AR(1) process similar  
 317 to (14)

$$\log(z_{t+1}) = \rho_z \cdot \log(z_t) + \nu_z \cdot \sqrt{1 - \rho_z^2} \cdot \epsilon_t, \quad \epsilon_t \sim N(0, 1), \quad (16)$$

318 where the persistence  $\rho_z = 0.9$  and  $\nu_z = 0.03$  to make the aggregate shock more steady and  
 319 controllable, thereby reducing the difficulty of solving. The grids  $\mathbf{z}$  and transition matrix  $\mathbf{P}_z$  of  
 320 aggregate shock are

$$\mathbf{z} = \begin{bmatrix} 0.9139 \\ 0.9560 \\ 1.0000 \\ 1.0460 \\ 1.0942 \end{bmatrix} \quad \text{and} \quad \mathbf{P}_z = \begin{bmatrix} 0.8478 & 0.1522 & 0.0000 & 0.0000 & 0.0000 \\ 0.0195 & 0.8962 & 0.0843 & 0.0000 & 0.0000 \\ 0.0000 & 0.0427 & 0.9147 & 0.0427 & 0.0000 \\ 0.0000 & 0.0000 & 0.0843 & 0.8962 & 0.0195 \\ 0.0000 & 0.0000 & 0.0000 & 0.1522 & 0.8478 \end{bmatrix}, \quad (17)$$

321 with  $\mathbf{z}$  and  $\mathbf{P}_z$  rounded to 4 decimals.

322 **A.2 Discretization**

323 There are two heterogeneous attributes, wealth  $a_{i,t}$  and income  $e_{i,t}$ . The latter has already been  
 324 discretized to  $n_e = 3$  states. We discretize wealth to  $n_a = 200$  grid points with  $a_{\min} = \underline{a} = 0$  and  
 325  $a_{\max} = 150$ . So the distribution  $\mathbf{g}_t$  is a vector with  $n_a \times n_e = 600$  dimensions. There are  $n_z = 5$   
 326 different aggregate states.

327 **B Krusell and Smith [10] model**

328 **B.1 Model Setup**

329 **Households** There's a continuum of households indexed by  $i \in [0, 1]$  and each household has  
 330 productivity level  $e_{i,t}$  and wealth  $a_{i,t}$ . Productivity level (idiosyncratic shock)  $e_{i,t}$  follows a three-  
 331 state Markov process which is the same as (15).

332 The aggregate shock  $z_t$  follows a two-state Markov process

$$\mathbf{z} = [z_b \quad z_g] = [0.99 \quad 1.01] \quad \text{and} \quad \mathbf{P}_z = \begin{bmatrix} 0.875 & 0.125 \\ 0.125 & 0.875 \end{bmatrix}, \quad (18)$$

333 where  $\mathbf{z}$  is the grid of aggregate shock and  $\mathbf{P}_z$  is the transition matrix.

334 Wealth  $a_{i,t}$  evolves according to the budget constraint

$$a_{i,t+1} + c_{i,t} = (1 + r_t)a_{i,t} + w_t e_{i,t}, \quad (19)$$

335 where  $r_t$  is interest rate and  $w_t$  is wage. Households take prices as given and live permanently to  
 336 optimize cumulative utility

$$\max_{c_{i,t}} \mathbb{E}_0 \left[ \sum_{t=0}^{\infty} \beta^t u(c_{i,t}) \right], \quad (20)$$

337 where  $\beta = 0.975$  is the discount factor and  $u(c) = \frac{c^{1-\sigma}}{1-\sigma}$  with  $\sigma = 1$  is the utility function.

338 **Firm** There is a representative firm who borrows capital from households and produces goods  
 339 according to Cobb-Douglas function

$$Y_t = z_t K_t^\alpha L_t^{1-\alpha}, \quad (21)$$

340 where  $K_t$  is aggregate capital,  $L_t$  is aggregate labor supply and  $\alpha = 0.36$  is capital share. Firm use  
 341 FOC condition to determine interest rate and wage

$$r_t = \frac{\partial Y_t}{\partial K_t} - \delta = \alpha z_t k_t^{\alpha-1} - \delta, \quad w_t = \frac{\partial Y_t}{\partial L_t} = (1-\alpha) z_t k_t^\alpha, \quad (22)$$

342 where  $\delta = 0.025$  is the capital depreciation rate and  $k_t = \frac{K_t}{L_t}$  is capital per capita.

343 **Equilibrium** Households maximize cumulative utility (20) subject to the budget constraint (19)  
 344 and taking as given prices according to firm's optimality condition (22).

## 345 B.2 KS Method

346 From (22), we know that as long as we get the knowledge of the law of motion of aggregate capital  
 347  $K_t$ , we get the knowledge of interest rate  $r_t$  and wage  $w_t$ . So in [10] method, we assume the  
 348 aggregate capital  $K_t$  evolves according to

$$\log(K_{t+1}) = \begin{cases} x_g \cdot \log(K_t) + y_g, & z_t = z_g \\ x_b \cdot \log(K_t) + y_b, & z_t = z_b \end{cases}, \quad (23)$$

349 where  $\theta = \begin{bmatrix} x_g & y_g \\ x_b & y_b \end{bmatrix}$  is unknown parameter. Then, there is a fixed point problem between unknown  
 350 parameter  $\theta$  and optimal consumption policy  $\pi^*$ .

351 Given parameter  $\theta \in \mathbb{R}^{2 \times 2}$ , we can find the corresponding optimal policy  $\pi^*(\theta)$  using Value Iteration,  
 352 for all transition dynamics are known. And at the same time, given a policy function  $\pi$ , we can  
 353 simulate the economy to get massive samples of  $(K_t, z_t, K_{t+1})$  and then run regression of (23) to  
 354 get optimal parameter  $\theta(\pi)$  that satisfy (23). So as long as we can find

$$\theta = \theta(\pi^*), \quad \pi^* = \pi^*(\theta), \quad (24)$$

355 we get the solution of Krusell and Smith [10] model.

356 KS method can be extended to solve nearly all heterogeneous agent models. However, not all models  
 357 exist a fixed point including what we solved in this paper. Economists have tried many function  
 358 forms of (23) to solve the Huggett [8] model with aggregate risk, yet all failed. A potential reason  
 359 why there does not exist a fixed point is the sensitiveness of the market clearing price.

## 360 C A PE Model of the Huggett [8] Economy

### 361 C.1 Model Setup

362 **Households** There's a representative household with wealth  $a_t$  and productivity level  $e_t$  at period  
 363  $t$ . There are two prices, interest rate  $r_t$  and wage  $w_t$ , that evolves exogenously. Productivity level  $e_t$ ,  
 364 interest rate  $r_t$  and wage  $w_t$  all follow a Markov process,

$$e_{t+1} \sim \mathcal{T}_e(\cdot | e_t), \quad r_{t+1} \sim \mathcal{T}_r(\cdot | r_t) \quad \text{and} \quad w_{t+1} \sim \mathcal{T}_w(\cdot | w_t), \quad (25)$$

365 where we discretize  $\mathcal{T}_e$  to three states,  $\mathcal{T}_r$  to five states and  $\mathcal{T}_w$  to seven states.

366 The discrete values and transition matrix of  $e_t$  are

$$\mathbf{e} = \begin{bmatrix} 0.5343 \\ 0.9735 \\ 1.7739 \end{bmatrix} \quad \text{and} \quad \mathbf{P}_e = \begin{bmatrix} 0.6460 & 0.3539 & 0.0001 \\ 0.0304 & 0.9392 & 0.0304 \\ 0.0001 & 0.3539 & 0.6460 \end{bmatrix}, \quad (26)$$

367 where we use the same process of our original model.

368 The discrete values and transition matrix of  $r_t$  are

$$\mathbf{r} = \begin{bmatrix} 0.0193 \\ 0.0218 \\ 0.0247 \\ 0.0279 \\ 0.0315 \end{bmatrix} \quad \text{and} \quad \mathbf{P}_r = \begin{bmatrix} 0.5655 & 0.4172 & 0.0173 & 0.0000 & 0.0000 \\ 0.0920 & 0.6079 & 0.2966 & 0.0035 & 0.0000 \\ 0.0034 & 0.1423 & 0.7032 & 0.1507 & 0.0004 \\ 0.0000 & 0.0056 & 0.2287 & 0.7131 & 0.0526 \\ 0.0000 & 0.0000 & 0.0081 & 0.3438 & 0.6480 \end{bmatrix}, \quad (27)$$

369 with the ergodic interest rate  $\bar{r} = 2.5\%$ .

370 The discrete values and transition matrix of  $w_t$  are

$$\mathbf{w} = \begin{bmatrix} 0.9418 \\ 0.9608 \\ 0.9802 \\ 1.0000 \\ 1.0202 \\ 1.0408 \\ 1.0618 \end{bmatrix} \quad \text{and} \quad \mathbf{P}_w = \begin{bmatrix} 0.6657 & 0.3312 & 0.0031 & 0.0000 & 0.0000 & 0.0000 & 0.0000 \\ 0.0541 & 0.7002 & 0.2442 & 0.0014 & 0.0000 & 0.0000 & 0.0000 \\ 0.0001 & 0.0842 & 0.7363 & 0.1787 & 0.0007 & 0.0000 & 0.0000 \\ 0.0000 & 0.0003 & 0.1254 & 0.7487 & 0.1254 & 0.0003 & 0.0000 \\ 0.0000 & 0.0000 & 0.0007 & 0.1787 & 0.7363 & 0.0842 & 0.0001 \\ 0.0000 & 0.0000 & 0.0000 & 0.0014 & 0.2442 & 0.7002 & 0.0541 \\ 0.0000 & 0.0000 & 0.0000 & 0.0000 & 0.0031 & 0.3312 & 0.6657 \end{bmatrix}, \quad (28)$$

371 with the ergodic wage  $\bar{w} = 1$ .

372 Wealth  $a_t$  evolves according to the budget constraint

$$a_{t+1} + c_t = (1 + r_t)a_t + w_t e_t, \quad (29)$$

373 where we use the natural borrowing constraint  $a_t \geq 0$  for all  $t$ . And the representative household  
374 optimizes the cumulative utility

$$\max_{c_t} \mathbb{E}_0 \left[ \sum_{t=0}^{\infty} \beta^t u(c_t) \right], \quad (30)$$

375 where we use the same parameters of  $\beta = 0.975$  and  $u(c) = \frac{c^{1-\sigma}}{1-\sigma}$  with  $\sigma = 1$ .

## 376 C.2 Numerical Solution (Value Iteration)

377 Since all transition dynamics are known, we use Value Iteration to solve it. Given the current state  
378  $(a_t, w_t, r_t, e_t)$ , we know the transition probability of  $(e_t, r_t, w_t)$  follows (25) and the transition dy-  
379 namic of  $a_t$  follows (29). Then, we can write the Bellman equation

$$\begin{aligned} \max_{c_t} \quad & u(c_t) + \beta \mathbb{E}[V(a_{t+1}, e_{t+1}, r_{t+1}, w_{t+1}) \mid (a_t, e_t, r_t, w_t)] \\ \text{s.t.} \quad & a_{t+1} = (1 + r_t)a_t + w_t e_t - c_t \\ & e_{t+1} \sim \mathcal{T}_e(\cdot \mid e_t) \\ & r_{t+1} \sim \mathcal{T}_r(\cdot \mid r_t) \\ & w_{t+1} \sim \mathcal{T}_w(\cdot \mid w_t) \end{aligned}, \quad (31)$$

380 therefor we can solve the optimal control by using Value Iteration.

## 381 C.3 Why PE the baseline?

382 The reason why we use PE as the baseline is we want to compare our DSPG algorithm with con-  
383 ventional deep reinforcement learning algorithms, including PPO, SAC and DDPG. PE is a real  
384 low-dimensional single-agent problem which can be easily extended as a reinforcement learning en-  
385 vironment and solved by existing open-source RL repos. Though we discussed in previous section  
386 that our complex environment is also equivalent to a single-agent problem, it's high-dimensional.  
387 The dimension of inputs and outputs are several hundreds, which is difficult for current open-source  
388 DRL repos to handle. So instead we use PE as the baseline.

389 **D Hyper-parameters**

390 **D.1 USPG (Huggett Economy)**

Symbol	Description	Value
$n_a$	Number of asset grid points	200
$a_{\min}, a_{\max}$	Borrowing limit and asset upper bound	0, 150
$n_e, n_z$	Labor productivity states, TFP states	3, 5
$\beta$	Discount factor	0.975
$\sigma$	Coefficient of relative risk aversion	1 (log utility)
$B$	Total bond supply (market clearing)	5
$\varepsilon_{\text{trunc}}$	Truncation error for infinite horizon	$10^{-2}$
$T$	Training rollout length	$\lfloor \ln(\varepsilon_{\text{trunc}}) / \ln \beta \rfloor$ (as in script)
$B_{\text{train}}$	Batch size	64
$N_{\text{epoch}}$	Training epochs	1000
$\eta_0$	Adam initial learning rate	$2 \times 10^{-3}$
LR schedule	<code>optax.exponential_decay</code> (multiplicative per epoch)	Factor 0.5
Optimizer		Adam
Warm-up	Initial distribution $g_0$ for early epochs	Epochs 1–50: <code>steady_dist</code> ; thereafter previous <code>final_g</code>
U-Net width	Output channels per <code>DoubleConv</code> block	4
Convolution	<code>Conv1D / Conv1DTranspose</code>	Kernel 3; <code>SAME</code> padding
Activation	Hidden layers	<code>LeakyReLU</code>
Downsampling	Pooling	<code>AvgPool</code> , window 2, stride 2
Output head	Consumption parameterization	<code>sigmoid</code> then cumulative sum for $c(a, e, z)$
Repeats	Outer loop in <code>ablation_study.py</code>	10

391 **D.2 USPG (PE Economy)**

Symbol	Description	Value
$N_a$ (NA)	Number of asset grid points	200
$a_{\min}, a_{\max}$	Asset bounds ( <code>PEEnv</code> )	0, 100
$n_e, n_r, n_w$	States for productivity, interest rate, wage factor	3, 5, 7
$\beta$	Discount factor	0.975
$\sigma$	Coefficient of relative risk aversion	1
$c_{\min}$	Minimum consumption (feasibility)	$10^{-3}$
$T$	Training and evaluation horizon	Default <code>PEEnv.T</code> ( $\lceil \ln(10^{-2}) / \ln \beta \rceil$ , $\approx 182$ ); override with <code>-horizon</code>
$B_{\text{train}}$	Batch size	64
$N_{\text{epoch}}$	Training epochs	1000
$N_{\text{rep}}$	Independent runs	10
$\eta_0$	Adam initial learning rate	$2 \times 10^{-3}$
LR schedule	<code>optax.exponential_decay</code>	Factor 0.5 per epoch
Optimizer		Adam
$N_{\text{warm}}$	Epochs with ergodic $g_0$ warm-start	50 ( <code>WARMUP_EPOCHS</code> )
$N_{\text{eval}}$	Interval for fixed- $g_0$ curve evaluation	100 ( <code>eval_every</code> ; first and last epoch always evaluated)
U-Net width	Channels per <code>DoubleConv</code>	4
Conv. / pooling		Same pattern as Huggett USPG (kernel 3, <code>LeakyReLU</code> , <code>AvgPool</code> )
Policy output	Relative to VFI upper bound	Consumption over $n_e$ ; $(r, w)$ enter through wealth, no separate head

392 **D.3 PPO (PE Economy)**

Symbol	Description	Value
$N_{\text{update}}$	Total PPO updates	500
$N_{\text{env}}$	Parallel environments	32
$H$	Rollout length per env per update	64
$d_h$	Shared trunk hidden width (MLP)	128 (two <code>tanh</code> layers)
$\eta$	Learning rate	$3 \times 10^{-4}$
$\gamma$	Discount factor	Default $\beta$ ( <code>PEEnv.beta</code> , 0.975)
$\lambda_{\text{GAE}}$	GAE parameter	0.95
$\varepsilon_{\text{clip}}$	PPO clip range	0.2
$E_{\text{ppo}}$	PPO epochs per update	10
$ \mathcal{M} $	Minibatch size	256
$c_v$	Value loss coefficient	0.5
$c_H$	Entropy bonus coefficient	0.01
–	Obs. normalization scale ( <code>normalize_obs</code> )	$[a_{\text{max}} - a_{\text{min}}, 2.5, 0.035, 0.15]$
–	Obs. normalization offset ( <code>normalize_obs</code> )	$[a_{\text{min}}, 0, r_{\text{min}} - 0.01, w_{\text{min}} - 0.02]$
Eval paths	Post-training Monte Carlo evaluation	4096 ( <code>-eval_paths</code> )
Ergodic eval paths	Ergodic eval when logging	256 ( <code>-log_ergodic_eval_paths</code> )

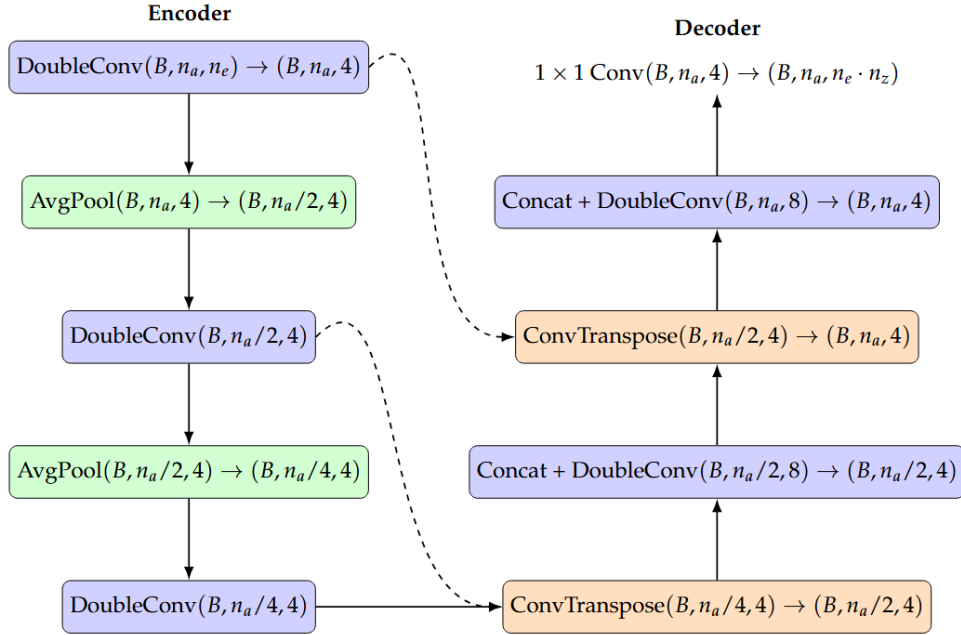
393 **D.4 SAC (PE Economy)**

Symbol	Description	Value
$N_{\text{update}}$	Total gradient-update rounds	500
$N_{\text{env}}$	Parallel environments	32
$H$	Rollout-related horizon	64
$d_h$	Actor / critic MLP hidden width	128 (actor: two ReLU layers)
$\eta$	Learning rate (Q and actor)	$3 \times 10^{-4}$
$\gamma$	Discount factor	Default $\beta = 0.975$
$\tau$	Target soft-update coefficient	0.005
$\alpha_{\text{SAC}}$	Entropy temperature (fixed)	0.2
$ \mathcal{B} $	Replay buffer capacity	500 000
$B_{\text{train}}$	Batch size per gradient step	256
$N_{\text{grad}}$	Gradient steps per environment step (after warm-up)	80
$N_{\text{start}}$	Minimum transitions before learning	5000
$\log \sigma_{\text{clip}}$	Stability	$[-20, 2]$ ( <code>LOG_STD_MIN/MAX</code> )
Eval paths	Post-training MC evaluation	4096 ( <code>-eval_paths</code> )
Ergodic eval paths	Ergodic eval when logging	256 ( <code>-log_ergodic_eval_paths</code> )

394 **D.5 DDPG (PE Economy)**

Symbol	Description	Value
$N_{\text{update}}$	Total update rounds	500
$N_{\text{env}}$	Parallel environments	32
$H$	Horizon	64
$d_h$	Actor / Q hidden width	128 (two ReLU layers)
$\eta$	Learning rate	$3 \times 10^{-4}$
$\gamma$	Discount factor	Default $\beta = 0.975$
$\tau$	Target soft-update coefficient	0.005
$\sigma_{\text{noise}}^{(0)}$	Initial exploration noise std (on share)	0.1
$\sigma_{\text{noise}}^{(\text{end})}$	Final exploration noise std	0.02
$ \mathcal{B} $	Replay buffer capacity	500 000
$B_{\text{train}}$	Batch size	256
$N_{\text{grad}}$	Gradient steps per update	80
$N_{\text{start}}$	Minimum samples before learning	5000
Eval paths	Post-training MC evaluation	4096 (-eval_paths)
Ergodic eval paths	Ergodic eval when logging	256 (log_ergodic_eval_paths)

395 **E Overall Structure of the 1D U-Net**



396 We employ a 1D U-Net architecture operating on the asset grid. The input tensor has shape  
 397  $(B, n_a, n_e)$ , where  $n_a$  denotes the asset grid and  $n_e$  the productivity states. Convolutions are applied  
 398 along the asset dimension while treating the productivity states as channels. The network consists of  
 399 two down-sampling layers followed by two symmetric up-sampling layers with skip connections.

## 400 F Theoretical Analysis of DSPG

401 This appendix analyzes the gradient computation and variance properties of DSPG using the notation  
 402 of Sections 3–4:  $\Theta$  denotes the U-Net parameters,  $\mathbf{g}_t \in \mathbb{R}_+^{n_a \times n_e}$  the discretized wealth–productivity  
 403 distribution with entries  $g_{k\ell,t} = d(a_k, e_\ell)$ ,  $\mathbf{c}_t \in \mathbb{R}^{n_a \times n_e}$  the consumption matrix,  $z_t$  the aggregate  
 404 shock, and  $r_t$  the equilibrium interest rate. We write  $u(\mathbf{c}_t) \odot \mathbf{g}_t := \sum_{k=1}^{n_a} \sum_{\ell=1}^{n_e} g_{k\ell,t} u(c_{k\ell,t})$  for  
 405 the aggregate welfare at date  $t$ , following the convention of Section 4.

### 406 F.1 The DSPG Computation Graph

407 At each period  $t$ , DSPG executes four operations:

- 408 (i) **Policy.**  $\mathbf{c}_t = \text{U-Net}(\mathbf{g}_t, z_t; \Theta)$  via (10).
- 409 (ii) **Market clearing.**  $r_t$  is obtained in closed form from the aggregated budget constraint via  
 410 (9):  $1 + r_t = (B - z_t \bar{e} + \mathbf{c}_t \odot \mathbf{g}_t) / (\mathbf{a} \odot \mathbf{g}_t)$ , which involves only weighted sums and one  
 411 scalar division.
- 412 (iii) **Wealth update.** At each grid point  $(k, \ell)$ , next-period wealth is  $a'_{k\ell,t} = (1 + r_t)a_k +$   
 413  $z_t e_\ell - c_{k\ell,t}$ , projected onto  $[a, \infty)$  to enforce the borrowing constraint (3).
- 414 (iv) **Distributional update.**  $\mathbf{g}_{t+1} = \Phi(\mathbf{g}_t, \mathbf{c}_t, r_t, z_t, z_{t+1})$ , where  $\Phi$  redistributes mass across  
 415 wealth grid points (via linear-interpolation weights that are piecewise linear in  $a'_{k\ell,t}$ ) and  
 416 composes with the exogenous transition matrix for  $e$ .

417 Composing steps (i)–(iv) over  $T$  periods yields the graph

$$\Theta \xrightarrow{\text{U-Net}} \mathbf{c}_0 \xrightarrow{(9)} r_0 \longrightarrow \mathbf{g}_1 \xrightarrow{\text{U-Net}} \mathbf{c}_1 \longrightarrow \cdots \longrightarrow \mathbf{g}_T, \quad (32)$$

418 every arrow of which is a (piecewise) differentiable map. Aggregate shocks  $\{z_t\}$  and idiosyncratic  
 419 realizations  $\{e_{i,t}\}$  are exogenous and do not depend on  $\Theta$ , so they act as fixed random inputs for a  
 420 given sample path.

### 421 F.2 Assumptions

422 [Differentiability] The utility  $u : \mathbb{R}_{++} \rightarrow \mathbb{R}$  is twice continuously differentiable, strictly increasing,  
 423 and strictly concave. The U-Net map  $\Theta \mapsto \text{U-Net}(\mathbf{g}, z; \Theta)$  is continuously differentiable. The redis-  
 424 tribution map  $\Phi$  is piecewise continuously differentiable in  $(\mathbf{g}_t, \mathbf{c}_t, r_t)$ . Non-smoothness is confined  
 425 to (a) the borrowing-constraint projection in step (iii) and (b) the boundaries between neighboring  
 426 grid cells in the interpolation of  $\Phi$ .

427 [Bounded rewards] Consumption is bounded along any trajectory:  $c_{k\ell,t} \in [\underline{c}, \bar{c}]$  with  $0 < \underline{c} \leq \bar{c} <$   
 428  $\infty$ , so that  $|u(c_{k\ell,t})| \leq M$  for some  $M < \infty$ .

429 [Positive aggregate wealth] The wealth grid contains at least one strictly positive point and  $\mathbf{g}_t$  assigns  
 430 positive mass to it, so that  $\mathbf{a} \odot \mathbf{g}_t > 0$  for all  $t$ ; with  $B > 0$ , the denominator in (9) is bounded away  
 431 from zero.

### 432 F.3 DSPG Objective and Gradient

433 Recall from (12) that the DSPG sample objective is

$$\mathbb{J}(\Theta) \approx -\frac{1}{n} \sum_{m=1}^n \sum_{t=0}^T \beta^t \underbrace{u(\mathbf{c}_t^{(m)}) \odot \mathbf{g}_t^{(m)}}_{=: U_t^{(m)}}, \quad (33)$$

434 where  $m$  indexes independent shock trajectories. For a fixed  $m$ ,  $U_t^{(m)}$  depends on  $\Theta$  through the  
 435 graph (32).

436 [Exact pathwise gradient] Under Assumptions F.2–F.2, for each trajectory  $m$  and almost every real-  
 437 ization of exogenous shocks, the map  $\Theta \mapsto \sum_{t=0}^T \beta^t U_t^{(m)}$  is piecewise differentiable, and

$$\nabla_{\Theta} \mathbb{J}(\Theta) = -\frac{1}{n} \sum_{m=1}^n \sum_{t=0}^T \beta^t \frac{dU_t^{(m)}}{d\Theta}, \quad (34)$$

438 where  $dU_t^{(m)}/d\Theta$  is computed by reverse-mode automatic differentiation through the graph (32).  
 439 No score-function (REINFORCE) estimator is used.

440 Fix a trajectory  $m$  and its exogenous shock sequence. Given these, the sequence  
 441  $(\mathbf{g}_0, \mathbf{c}_0, r_0, \mathbf{g}_1, \mathbf{c}_1, r_1, \dots, \mathbf{g}_T)$  is a deterministic, piecewise differentiable function of  $\Theta$ :

- 442 •  $\mathbf{g}_0$  is an exogenous initial distribution (independent of  $\Theta$ ).
- 443 •  $\mathbf{c}_t = \text{U-Net}(\mathbf{g}_t, z_t; \Theta)$  is  $C^1$  in  $\Theta$  by Assumption F.2.
- 444 •  $r_t$  is a ratio of affine functions of  $(\mathbf{g}_t, \mathbf{c}_t)$  via (9); by Assumption F.2 the denominator is  
 445 bounded away from zero, so  $r_t$  is  $C^1$  in  $(\mathbf{g}_t, \mathbf{c}_t)$ .
- 446 •  $\mathbf{g}_{t+1} = \Phi(\mathbf{g}_t, \mathbf{c}_t, r_t, z_t, z_{t+1})$  is piecewise  $C^1$  by Assumption F.2.

447 By the chain rule,  $U_t^{(m)}$  is piecewise differentiable in  $\Theta$  for every  $t$ , and so is the finite sum  
 448  $\sum_{t=0}^T \beta^t U_t^{(m)}$ . The gradient (34) follows from applying the chain rule through the entire graph,  
 449 which is exactly what reverse-mode automatic differentiation computes.

450 Because  $\{z_t\}, \{e_{i,t}\}$  enter as fixed exogenous inputs (not functions of  $\Theta$ ), no likelihood-ratio term  
 451 appears. Expectation and differentiation commute by dominated convergence: the integrand is  
 452 bounded by  $M \sum_{t=0}^T \beta^t < M/(1 - \beta)$  under Assumption F.2.

453 **Gradient decomposition.** The total derivative of  $U_t$  with respect to  $\Theta$  admits the natural decom-  
 454 position

$$\frac{dU_t}{d\Theta} = \underbrace{\frac{\partial U_t}{\partial \mathbf{c}_t} \frac{\partial \mathbf{c}_t}{\partial \Theta}}_{\text{direct effect}} + \underbrace{\frac{\partial U_t}{\partial \mathbf{g}_t} \frac{d\mathbf{g}_t}{d\Theta}}_{\text{distributional effect}}, \quad (35)$$

455 where  $d\mathbf{g}_t/d\Theta$  depends recursively on all prior periods through  $\Phi$ ,  $r_s$ , and  $\mathbf{c}_s$  for  $s < t$ . The  
 456 first term captures the immediate impact of  $\Theta$  on current consumption; the second captures the  
 457 cumulative effect of  $\Theta$  on the wealth distribution, which feeds back through market clearing (9) into  
 458  $r_t$ . Both are computed automatically by backpropagation.

#### 459 F.4 Truncation Bias

460 Let  $\mathbb{J}^\infty(\Theta) = -\mathbb{E}[\sum_{t=0}^\infty \beta^t U_t]$  be the infinite-horizon objective. DSPG uses a finite truncation  
 461  $T < \infty$ . Under Assumption F.2,

$$|\mathbb{J}^\infty(\Theta) - \mathbb{J}(\Theta)| \leq \frac{M \beta^{T+1}}{1 - \beta}, \quad (36)$$

462 i.e., the bias decays geometrically in  $T$ . The same bound applies to gradients, with  $M$  replaced by  
 463 the product of  $M$  and the Lipschitz constant of the graph (32). For  $\beta \approx 0.96$ , truncation at  $T$  on the  
 464 order of hundreds makes the bias numerically negligible.

#### 465 F.5 Variance Reduction over REINFORCE

466 [Zero conditional variance] For a fixed exogenous shock sequence, the DSPG gradient (34) is a  
 467 deterministic function of  $\Theta$ ; it therefore has zero variance conditional on the shocks. In contrast,  
 468 a REINFORCE estimator applied to the same trajectory carries strictly positive variance from the  
 469 score-function weighting  $\nabla_\Theta \log \pi_\Theta \cdot R$ .

470 The residual variance of (34) arises only from averaging over  $n$  independent shock trajectories and  
 471 decreases at the standard  $O(1/n)$  rate. DSPG thus eliminates one of the two variance sources carried  
 472 by REINFORCE-style estimators; empirically, this permits training with far fewer trajectories than  
 473 classical DRL baselines.

#### 474 F.6 Convexity of the Inner Problem

475 The per-period objective  $U_t = \sum_{k,\ell} g_{k\ell,t} u(c_{k\ell,t})$  is concave in  $\mathbf{c}_t$  (non-negative weighted sum of  
 476 concave functions), and the budget constraint (2) is linear in  $c_{k\ell,t}$ . The inner problem of choosing  $\mathbf{c}_t$   
 477 given a fixed distribution and price is therefore convex in the minimization sense of  $\mathbb{J}$ . This structure

478 does *not* imply that  $\mathbb{J}(\Theta)$  is globally convex in the U-Net parameters  $\Theta$ , but it yields a benign loss  
 479 landscape: empirically (Section 5), gradient descent converges to the same solution from diverse  
 480 random initializations, consistent with the observation in Section 4.

## 481 F.7 Remarks on Non-smoothness

482 Proposition F.3 asserts *piecewise* differentiability rather than global smoothness. Two sources of  
 483 non-smoothness arise in the DSPG graph:

484 **Borrowing constraint.** Step (iii) projects next-period wealth onto  $[\underline{a}, \infty)$  via  $a'_{k\ell,t} \mapsto$   
 485  $\max(a'_{k\ell,t}, \underline{a})$ . The max operator is non-differentiable at  $a'_{k\ell,t} = \underline{a}$ . However, for a given  $\Theta$  the  
 486 set of parameters at which some grid point exactly hits the kink is a measure-zero subset of the  
 487 parameter space. Away from this set, the projection is locally linear ( $a'_{k\ell,t}$  or  $\underline{a}$ , whichever is active)  
 488 and poses no obstacle to automatic differentiation. In practice, JAX computes a valid subgradient at  
 489 the kink, which is sufficient for SGD convergence [? ].

490 **Grid interpolation.** Step (iv) redistributes mass using linear-interpolation weights that are piece-  
 491 wise linear in the continuous wealth value  $a'_{k\ell,t}$ . The interpolation weights change their functional  
 492 form when  $a'_{k\ell,t}$  crosses a grid boundary. As with the borrowing constraint, the set of  $\Theta$  values  
 493 where  $a'_{k\ell,t}$  lies exactly on a grid boundary has measure zero, and the interpolation is smooth al-  
 494 most everywhere. Furthermore, because linear interpolation is a convex combination, the resulting  
 495 gradients are bounded, so the occasional subgradient at a grid boundary does not destabilize training.

496 In summary, the non-smooth points form a measure-zero set in  $\Theta$ -space at each time step. The  
 497 composition over  $T$  steps remains piecewise  $C^1$  almost everywhere, and gradient descent with sub-  
 498 gradients at non-smooth points is well studied and known to converge under mild conditions [? ].

## 499 F.8 Lipschitz Regularity of the Gradient

500 The truncation-bias bound on gradients in Section F.4 (equation (36)) involves the Lipschitz constant  
 501 of the computational graph. We argue here that this constant is finite under our assumptions.

502 [Finite Lipschitz constant] Under Assumptions F.2–F.2, let  $L_{\text{net}}$  denote the Lipschitz constant of the  
 503 U-Net map  $(\mathbf{g}, z) \mapsto \text{U-Net}(\mathbf{g}, z; \Theta)$  with respect to its inputs (for a fixed  $\Theta$ ), and let  $L_{\Phi}$  denote  
 504 the Lipschitz constant of the redistribution map  $\Phi$  with respect to  $(\mathbf{g}_t, \mathbf{c}_t, r_t)$ . Then the one-step  
 505 Jacobian of the graph (32) satisfies

$$\left\| \frac{d\mathbf{g}_{t+1}}{d\mathbf{g}_t} \right\| \leq L_{\Phi} (1 + L_{\text{net}})(1 + L_r), \quad (37)$$

506 where  $L_r$  is the Lipschitz constant of  $r_t$  with respect to  $(\mathbf{g}_t, \mathbf{c}_t)$ , which is finite by Assumption F.2  
 507 (the denominator in (9) is bounded away from zero, and the numerator is Lipschitz in its arguments).

508 From the graph,  $\mathbf{g}_{t+1}$  depends on  $\mathbf{g}_t$  through three channels: (a) directly via  $\Phi$ , (b) via  $\mathbf{c}_t =$   
 509  $\text{U-Net}(\mathbf{g}_t, z_t; \Theta)$ , and (c) via  $r_t$  which depends on  $(\mathbf{g}_t, \mathbf{c}_t)$ . By the chain rule and the triangle  
 510 inequality:

$$\left\| \frac{d\mathbf{g}_{t+1}}{d\mathbf{g}_t} \right\| \leq \left\| \frac{\partial\Phi}{\partial\mathbf{g}_t} \right\| + \left\| \frac{\partial\Phi}{\partial\mathbf{c}_t} \right\| \left\| \frac{\partial\mathbf{c}_t}{\partial\mathbf{g}_t} \right\| + \left\| \frac{\partial\Phi}{\partial r_t} \right\| \left\| \frac{\partial r_t}{\partial\mathbf{g}_t} \right\| + \left\| \frac{\partial\Phi}{\partial r_t} \right\| \left\| \frac{\partial r_t}{\partial\mathbf{c}_t} \right\| \left\| \frac{\partial\mathbf{c}_t}{\partial\mathbf{g}_t} \right\|.$$

511 Each factor is finite:  $\|\partial\Phi/\partial\cdot\| \leq L_{\Phi}$  by assumption;  $\|\partial\mathbf{c}_t/\partial\mathbf{g}_t\| \leq L_{\text{net}}$ ;  $\|\partial r_t/\partial\cdot\| \leq L_r$  since  $r_t$   
 512 is a ratio with bounded-away-from-zero denominator. Collecting terms yields (37).

513 By induction, the  $T$ -step Jacobian is bounded by  $(L_{\Phi}(1 + L_{\text{net}})(1 + L_r))^T$ . When multiplied by the  
 514 discount  $\beta^t$ , the contribution of period- $t$  gradients to the total gradient is bounded by  $M \cdot (\beta \cdot L_{\Phi}(1 +$   
 515  $L_{\text{net}})(1 + L_r))^t$ . In the regime  $\beta \cdot L_{\Phi}(1 + L_{\text{net}})(1 + L_r) < 1$  (which holds when  $\beta$  is sufficiently small  
 516 relative to the per-step expansion), the gradient series converges absolutely and the truncation bias  
 517 on the gradient also decays geometrically. In practice, the effective per-step expansion is moderate  
 518 because the market-clearing mechanism acts as a stabilizing feedback: if households over-consume,  
 519  $r_t$  rises to restore equilibrium, dampening the distributional response.

520 **F.9 Pre-training and Initialization**

521 The steady-state pre-training described in Section 4 produces an initial parameter vector  $\Theta_0$  by fitting  
522 the value-iteration policy of a Huggett [8] economy without aggregate risk. This step is a separate  
523 optimization and does not affect the theoretical properties of DSPG’s main training phase:  $\Theta_0$  is  
524 treated as a fixed starting point, and all gradient computations in (34) are with respect to the current  
525  $\Theta$  only. The initial distribution  $\mathbf{g}_0$  is likewise set to the steady-state distribution and is independent  
526 of  $\Theta$ . Pre-training serves only to place  $\Theta_0$  in a region where the policy produces economically  
527 reasonable consumption values from the outset, ensuring numerical stability (in particular, avoiding  
528 negative or exploding consumption that would violate Assumption F.2).

529 **Summary.** DSPG assembles a fully differentiable simulation—U-Net policy, closed-form market  
530 clearing (9), wealth update, and distributional transition—and obtains exact pathwise gradients via  
531 backpropagation. Compared with REINFORCE, it eliminates score-function variance entirely; the  
532 only remaining stochasticity arises from sampling over exogenous shock trajectories and decays at  
533 the  $O(1/n)$  rate. Non-smooth points (borrowing constraint, grid boundaries) form a measure-zero  
534 set and do not impede SGD convergence. The per-step Lipschitz constant of the computational  
535 graph is finite, ensuring that both truncation bias and gradient norms remain controlled. Horizon  
536 truncation contributes a geometrically decaying bias that is negligible in practice.

537 **NeurIPS Paper Checklist**

538 **1. Claims**

539 Question: Do the main claims made in the abstract and introduction accurately reflect the  
540 paper’s contributions and scope?

541 Answer: [Yes]

542 Justification: The abstract and introduction clearly state that DSPG solves the Huggett [8]  
543 economy with aggregate risk and non-trivial market clearing, achieving market-clearing er-  
544 rors below  $10^{-11}$ . These claims are supported by the theoretical analysis (Appendix F) and  
545 experimental results (Section 5). Limitations regarding grid scalability are also discussed  
546 in Section 4.

547 Guidelines:

- 548 • The answer [N/A] means that the abstract and introduction do not include the claims  
549 made in the paper.
- 550 • The abstract and/or introduction should clearly state the claims made, including the  
551 contributions made in the paper and important assumptions and limitations. A [No] or  
552 [N/A] answer to this question will not be perceived well by the reviewers.
- 553 • The claims made should match theoretical and experimental results, and reflect how  
554 much the results can be expected to generalize to other settings.
- 555 • It is fine to include aspirational goals as motivation as long as it is clear that these  
556 goals are not attained by the paper.

557 **2. Limitations**

558 Question: Does the paper discuss the limitations of the work performed by the authors?

559 Answer: [Yes]

560 Justification: Section 4 explicitly discusses that DSPG relies on a fixed grid whose cost  
561 scales with  $n_a \times n_e$ , and that backpropagation through  $T$  time steps may exhaust GPU  
562 memory for very high-dimensional distributions regardless of network architecture. Ap-  
563 pendix F discusses non-smoothness at borrowing constraints and grid boundaries, and the  
564 truncation bias introduced by finite horizons.

565 Guidelines:

- 566 • The answer [N/A] means that the paper has no limitation while the answer [No] means  
567 that the paper has limitations, but those are not discussed in the paper.
- 568 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 569 • The paper should point out any strong assumptions and how robust the results are to  
570 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
571 model well-specification, asymptotic approximations only holding locally). The au-  
572 thors should reflect on how these assumptions might be violated in practice and what  
573 the implications would be.
- 574 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
575 only tested on a few datasets or with a few runs. In general, empirical results often  
576 depend on implicit assumptions, which should be articulated.
- 577 • The authors should reflect on the factors that influence the performance of the ap-  
578 proach. For example, a facial recognition algorithm may perform poorly when image  
579 resolution is low or images are taken in low lighting. Or a speech-to-text system might  
580 not be used reliably to provide closed captions for online lectures because it fails to  
581 handle technical jargon.
- 582 • The authors should discuss the computational efficiency of the proposed algorithms  
583 and how they scale with dataset size.
- 584 • If applicable, the authors should discuss possible limitations of their approach to ad-  
585 dress problems of privacy and fairness.
- 586 • While the authors might fear that complete honesty about limitations might be used by  
587 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
588 limitations that aren’t acknowledged in the paper. The authors should use their best

589 judgment and recognize that individual actions in favor of transparency play an impor-  
590 tant role in developing norms that preserve the integrity of the community. Reviewers  
591 will be specifically instructed to not penalize honesty concerning limitations.

### 592 3. Theory assumptions and proofs

593 Question: For each theoretical result, does the paper provide the full set of assumptions and  
594 a complete (and correct) proof?

595 Answer: [Yes]

596 Justification: Appendix F states three explicit assumptions (Assumptions F.2–F.2) and pro-  
597 vides complete proofs for Proposition F.3 (exact pathwise gradient), Proposition F.5 (zero  
598 conditional variance), and Proposition F.8 (finite Lipschitz constant). The truncation bias  
599 bound is also derived with an explicit formula.

600 Guidelines:

- 601 • The answer [N/A] means that the paper does not include theoretical results.
- 602 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
603 referenced.
- 604 • All assumptions should be clearly stated or referenced in the statement of any theo-  
605 rems.
- 606 • The proofs can either appear in the main paper or the supplemental material, but if  
607 they appear in the supplemental material, the authors are encouraged to provide a  
608 short proof sketch to provide intuition.
- 609 • Inversely, any informal proof provided in the core of the paper should be comple-  
610 mented by formal proofs provided in appendix or supplemental material.
- 611 • Theorems and Lemmas that the proof relies upon should be properly referenced.

### 612 4. Experimental result reproducibility

613 Question: Does the paper fully disclose all the information needed to reproduce the main  
614 experimental results of the paper to the extent that it affects the main claims and/or conclu-  
615 sions of the paper (regardless of whether the code and data are provided or not)?

616 Answer: [Yes]

617 Justification: The model specification is fully described in Section 3. The algorithm is  
618 presented as pseudocode in Algorithm 1. All hyperparameters (learning rate, number of  
619 trajectories, truncation length, etc.) are reported in the appendices. The U-Net architecture  
620 is detailed in Appendix E and model calibration in Appendix A.1. Hardware specifications  
621 are provided in the introduction.

622 Guidelines:

- 623 • The answer [N/A] means that the paper does not include experiments.
- 624 • If the paper includes experiments, a [No] answer to this question will not be per-  
625 ceived well by the reviewers: Making the paper reproducible is important, regardless  
626 of whether the code and data are provided or not.
- 627 • If the contribution is a dataset and/or model, the authors should describe the steps  
628 taken to make their results reproducible or verifiable.
- 629 • Depending on the contribution, reproducibility can be accomplished in various ways.  
630 For example, if the contribution is a novel architecture, describing the architecture  
631 fully might suffice, or if the contribution is a specific model and empirical evaluation,  
632 it may be necessary to either make it possible for others to replicate the model with  
633 the same dataset, or provide access to the model. In general, releasing code and data  
634 is often one good way to accomplish this, but reproducibility can also be provided via  
635 detailed instructions for how to replicate the results, access to a hosted model (e.g., in  
636 the case of a large language model), releasing of a model checkpoint, or other means  
637 that are appropriate to the research performed.
- 638 • While NeurIPS does not require releasing code, the conference does require all sub-  
639 missions to provide some reasonable avenue for reproducibility, which may depend  
640 on the nature of the contribution. For example  
641 (a) If the contribution is primarily a new algorithm, the paper should make it clear  
642 how to reproduce that algorithm.

- 643 (b) If the contribution is primarily a new model architecture, the paper should describe  
644 the architecture clearly and fully.
- 645 (c) If the contribution is a new model (e.g., a large language model), then there should  
646 either be a way to access this model for reproducing the results or a way to re-  
647 produce the model (e.g., with an open-source dataset or instructions for how to  
648 construct the dataset).
- 649 (d) We recognize that reproducibility may be tricky in some cases, in which case au-  
650 thors are welcome to describe the particular way they provide for reproducibility.  
651 In the case of closed-source models, it may be that access to the model is limited in  
652 some way (e.g., to registered users), but it should be possible for other researchers  
653 to have some path to reproducing or verifying the results.

## 654 5. Open access to data and code

655 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
656 tions to faithfully reproduce the main experimental results, as described in supplemental  
657 material?

658 Answer: [Yes]

659 Justification: The complete codebase (written entirely in JAX) is publicly available at  
660 <https://github.com/leafDancer/DSPG>. The repository includes code for both GE  
661 and PE benchmark models, training scripts, and steady-state value-iteration solvers.

662 Guidelines:

- 663 • The answer [N/A] means that paper does not include experiments requiring code.
- 664 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/  
665 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 666 • While we encourage the release of code and data, we understand that this might not  
667 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not  
668 including code, unless this is central to the contribution (e.g., for a new open-source  
669 benchmark).
- 670 • The instructions should contain the exact command and environment needed to run to  
671 reproduce the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 672 • The authors should provide instructions on data access and preparation, including how  
673 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 674 • The authors should provide scripts to reproduce all experimental results for the new  
675 proposed method and baselines. If only a subset of experiments are reproducible, they  
676 should state which ones are omitted from the script and why.
- 677 • At submission time, to preserve anonymity, the authors should release anonymized  
678 versions (if applicable).
- 679 • Providing as much information as possible in supplemental material (appended to the  
680 paper) is recommended, but including URLs to data and code is permitted.

## 681 6. Experimental setting/details

682 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-  
683 rameters, how they were chosen, type of optimizer) necessary to understand the results?

684 Answer: [Yes]

685 Justification: The optimizer (gradient descent via JAX autodiff), learning rate, number  
686 of trajectories  $n$ , truncation length  $T$ , and all model calibration parameters are reported  
687 in the appendices (Appendix D.1 for DSPG on the GE model, Appendices D.2–D.5 for  
688 all methods on the PE benchmark). The U-Net architecture is specified in Appendix E.  
689 Hardware details (RTX 3090 GPUs) are stated.

690 Guidelines:

- 691 • The answer [N/A] means that the paper does not include experiments.
- 692 • The experimental setting should be presented in the core of the paper to a level of  
693 detail that is necessary to appreciate the results and make sense of them.
- 694 • The full details can be provided either with the code, in appendix, or as supplemental  
695 material.
- 696

697 **7. Experiment statistical significance**

698 Question: Does the paper report error bars suitably and correctly defined or other appropri-  
699 ate information about the statistical significance of the experiments?

700 Answer: [Yes]

701 Justification: All main experiments are repeated over 10 independent runs. Table 2 reports  
702 both mean and standard deviation for each method. Training curves in Figure 7 and Fig-  
703 ure 4 show  $\pm 1.96$  standard deviation shading across runs. The variability source (different  
704 random seeds for aggregate shock trajectories) is stated.

705 Guidelines:

- 706 • The answer [N/A] means that the paper does not include experiments.
- 707 • The authors should answer [Yes] if the results are accompanied by error bars, confi-  
708 dence intervals, or statistical significance tests, at least for the experiments that support  
709 the main claims of the paper.
- 710 • The factors of variability that the error bars are capturing should be clearly stated (for  
711 example, train/test split, initialization, random drawing of some parameter, or overall  
712 run with given experimental conditions).
- 713 • The method for calculating the error bars should be explained (closed form formula,  
714 call to a library function, bootstrap, etc.)
- 715 • The assumptions made should be given (e.g., Normally distributed errors).
- 716 • It should be clear whether the error bar is the standard deviation or the standard error  
717 of the mean.
- 718 • It is OK to report 1-sigma error bars, but one should state it. The authors should prefer-  
719 ably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of  
720 Normality of errors is not verified.
- 721 • For asymmetric distributions, the authors should be careful not to show in tables or fig-  
722 ures symmetric error bars that would yield results that are out of range (e.g., negative  
723 error rates).
- 724 • If error bars are reported in tables or plots, the authors should explain in the text how  
725 they were calculated and reference the corresponding figures or tables in the text.

726 **8. Experiments compute resources**

727 Question: For each experiment, does the paper provide sufficient information on the com-  
728 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
729 the experiments?

730 Answer: [Yes]

731 Justification: The introduction specifies the hardware used (a server with NVIDIA RTX  
732 3090 GPUs and 32 AMD CPU cores). Section 4 reports per-step computational complexity.  
733 The ablation study (Figure 4) reports wall-clock training time under different trajectory  
734 counts. The introduction states that the model is solved within minutes.

735 Guidelines:

- 736 • The answer [N/A] means that the paper does not include experiments.
- 737 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
738 or cloud provider, including relevant memory and storage.
- 739 • The paper should provide the amount of compute required for each of the individual  
740 experimental runs as well as estimate the total compute.
- 741 • The paper should disclose whether the full research project required more compute  
742 than the experiments reported in the paper (e.g., preliminary or failed experiments  
743 that didn't make it into the paper).

744 **9. Code of ethics**

745 Question: Does the research conducted in the paper conform, in every respect, with the  
746 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

747 Answer: [Yes]

748 Justification: This work develops a computational method for solving macroeconomic mod-  
749 els. It does not involve human subjects, personal data, or any activity that could conflict  
750 with the NeurIPS Code of Ethics.

751 Guidelines:

- 752 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of  
753 Ethics.
- 754 • If the authors answer [No], they should explain the special circumstances that require  
755 a deviation from the Code of Ethics.
- 756 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
757 eration due to laws or regulations in their jurisdiction).

## 758 10. Broader impacts

759 Question: Does the paper discuss both potential positive societal impacts and negative  
760 societal impacts of the work performed?

761 Answer: [N/A]

762 Justification: DSPG is a computational method for solving academic macroeconomic mod-  
763 els. It does not directly enable deployment decisions, surveillance, or other applications  
764 with foreseeable negative societal impact. The positive impact is enabling economists to  
765 study richer models of inequality and policy, which is discussed in the introduction.

766 Guidelines:

- 767 • The answer [N/A] means that there is no societal impact of the work performed.
- 768 • If the authors answer [N/A] or [No], they should explain why their work has no soci-  
769 etal impact or why the paper does not address societal impact.
- 770 • Examples of negative societal impacts include potential malicious or unintended uses  
771 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
772 (e.g., deployment of technologies that could make decisions that unfairly impact spe-  
773 cific groups), privacy considerations, and security considerations.
- 774 • The conference expects that many papers will be foundational research and not tied  
775 to particular applications, let alone deployments. However, if there is a direct path to  
776 any negative applications, the authors should point it out. For example, it is legitimate  
777 to point out that an improvement in the quality of generative models could be used to  
778 generate Deepfakes for disinformation. On the other hand, it is not needed to point out  
779 that a generic algorithm for optimizing neural networks could enable people to train  
780 models that generate Deepfakes faster.
- 781 • The authors should consider possible harms that could arise when the technology is  
782 being used as intended and functioning correctly, harms that could arise when the  
783 technology is being used as intended but gives incorrect results, and harms following  
784 from (intentional or unintentional) misuse of the technology.
- 785 • If there are negative societal impacts, the authors could also discuss possible mitiga-  
786 tion strategies (e.g., gated release of models, providing defenses in addition to attacks,  
787 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from  
788 feedback over time, improving the efficiency and accessibility of ML).

## 789 11. Safeguards

790 Question: Does the paper describe safeguards that have been put in place for responsible  
791 release of data or models that have a high risk for misuse (e.g., pre-trained language models,  
792 image generators, or scraped datasets)?

793 Answer: [N/A]

794 Justification: The released code solves stylized economic models and does not pose risks  
795 of misuse comparable to generative models or scraped datasets.

796 Guidelines:

- 797 • The answer [N/A] means that the paper poses no such risks.
- 798 • Released models that have a high risk for misuse or dual-use should be released with  
799 necessary safeguards to allow for controlled use of the model, for example by re-  
800 quiring that users adhere to usage guidelines or restrictions to access the model or  
801 implementing safety filters.

- 802 • Datasets that have been scraped from the Internet could pose safety risks. The authors  
803 should describe how they avoided releasing unsafe images.  
804 • We recognize that providing effective safeguards is challenging, and many papers do  
805 not require this, but we encourage authors to take this into account and make a best  
806 faith effort.

## 807 12. Licenses for existing assets

808 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
809 the paper, properly credited and are the license and terms of use explicitly mentioned and  
810 properly respected?

811 Answer: [N/A]

812 Justification: The paper does not use pre-existing code packages, datasets, or pre-trained  
813 models from external sources. All code is written from scratch by the authors. The eco-  
814 nomic models referenced are theoretical constructs from the academic literature and are  
815 properly cited.

816 Guidelines:

- 817 • The answer [N/A] means that the paper does not use existing assets.
- 818 • The authors should cite the original paper that produced the code package or dataset.
- 819 • The authors should state which version of the asset is used and, if possible, include a  
820 URL.
- 821 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 822 • For scraped data from a particular source (e.g., website), the copyright and terms of  
823 service of that source should be provided.
- 824 • If assets are released, the license, copyright information, and terms of use in the pack-  
825 age should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has  
826 curated licenses for some datasets. Their licensing guide can help determine the li-  
827 cense of a dataset.
- 828 • For existing datasets that are re-packaged, both the original license and the license of  
829 the derived asset (if it has changed) should be provided.
- 830 • If this information is not available online, the authors are encouraged to reach out to  
831 the asset’s creators.

## 832 13. New assets

833 Question: Are new assets introduced in the paper well documented and is the documenta-  
834 tion provided alongside the assets?

835 Answer: [Yes]

836 Justification: We release two benchmark models (GE and PE versions of the dy-  
837 namic Huggett [8] economy) along with the DSPG codebase at <https://github.com/leafDancer/DSPG>. The code is written entirely in JAX and includes documentation for  
838 reproducing all experiments.  
839

840 Guidelines:

- 841 • The answer [N/A] means that the paper does not release new assets.
- 842 • Researchers should communicate the details of the dataset/code/model as part of their  
843 submissions via structured templates. This includes details about training, license,  
844 limitations, etc.
- 845 • The paper should discuss whether and how consent was obtained from people whose  
846 asset is used.
- 847 • At submission time, remember to anonymize your assets (if applicable). You can  
848 either create an anonymized URL or include an anonymized zip file.

## 849 14. Crowdsourcing and research with human subjects

850 Question: For crowdsourcing experiments and research with human subjects, does the pa-  
851 per include the full text of instructions given to participants and screenshots, if applicable,  
852 as well as details about compensation (if any)?

853 Answer: [N/A]

854 Justification: This work does not involve crowdsourcing or human subjects.

855 Guidelines:

- 856 • The answer [N/A] means that the paper does not involve crowdsourcing nor research  
857 with human subjects.
- 858 • Including this information in the supplemental material is fine, but if the main contri-  
859 bution of the paper involves human subjects, then as much detail as possible should  
860 be included in the main paper.
- 861 • According to the NeurIPS Code of Ethics, workers involved in data collection, cura-  
862 tion, or other labor should be paid at least the minimum wage in the country of the  
863 data collector.

864 **15. Institutional review board (IRB) approvals or equivalent for research with human**  
865 **subjects**

866 Question: Does the paper describe potential risks incurred by study participants, whether  
867 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
868 approvals (or an equivalent approval/review based on the requirements of your country or  
869 institution) were obtained?

870 Answer: [N/A]

871 Justification: This work does not involve human subjects research.

872 Guidelines:

- 873 • The answer [N/A] means that the paper does not involve crowdsourcing nor research  
874 with human subjects.
- 875 • Depending on the country in which research is conducted, IRB approval (or equiva-  
876 lent) may be required for any human subjects research. If you obtained IRB approval,  
877 you should clearly state this in the paper.
- 878 • We recognize that the procedures for this may vary significantly between institutions  
879 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
880 guidelines for their institution.
- 881 • For initial submissions, do not include any information that would break anonymity  
882 (if applicable), such as the institution conducting the review.

883 **16. Declaration of LLM usage**

884 Question: Does the paper describe the usage of LLMs if it is an important, original, or  
885 non-standard component of the core methods in this research? Note that if the LLM is used  
886 only for writing, editing, or formatting purposes and does *not* impact the core methodology,  
887 scientific rigor, or originality of the research, declaration is not required.

888 Answer: [N/A]

889 Justification: LLMs are not used as a component of the DSPG algorithm or any part of the  
890 core methodology. Any use of LLMs was limited to writing assistance and does not affect  
891 the scientific content.

892 Guidelines:

- 893 • The answer [N/A] means that the core method development in this research does not  
894 involve LLMs as any important, original, or non-standard components.
- 895 • Please refer to our LLM policy in the NeurIPS handbook for what should or should  
896 not be described.