

Emergent Coordination and Phase Structure in Independent Multi-Agent Reinforcement Learning

Azusa Yamaguchi
School of Physics and Astronomy
University of Edinburgh
ayamaguc@staffmail.ed.ac.uk

December 1, 2025

Abstract

A clearer understanding of when coordination emerges, fluctuates, or collapses in decentralized multi-agent reinforcement learning (MARL) is increasingly sought in order to characterize the dynamics of multi-agent learning systems. We revisit fully independent Q-learning (IQL) as a minimal decentralized testbed and run large-scale experiments across environment size L and agent density ρ . We construct a phase map using two axes—cooperative success rate (CSR) and a stability index derived from TD-error variance—revealing three distinct regimes: a coordinated and stable phase, a fragile transition region, and a jammed/disordered phase. A sharp *double* Instability Ridge separates these regimes and corresponds to persistent kernel drift, the time-varying shift of each agent’s effective transition kernel induced by others’ policy updates. Synchronization analysis further shows that temporal alignment is required for sustained cooperation, and that drift–synchronization competition generates the fragile regime. Ablation without agent identifiers removes drift entirely and collapses the three-phase structure, demonstrating that small inter-agent asymmetries are a necessary driver of drift. Overall, our results show that decentralized MARL exhibits a coherent phase structure governed by the interaction between scale, density, and kernel drift, suggesting that emergent coordination behaves as a distribution–interaction driven phase phenomenon.

1 Introduction

A central aim of artificial intelligence is to develop agents that can perceive, adapt, and interact with others in complex environments. Reinforcement learning (RL) provides a basic framework for learning through experience, with many of its successes demonstrated in single-agent domains [1, 10].

In multi-agent reinforcement learning (MARL), however, several agents update their policies concurrently, altering each other’s transition dynamics and creating an inherently non-stationary learning problem [7]. Such non-stationarity is frequently associated with instability, divergence, and failures of coordination.

A number of influential approaches mitigate these challenges by introducing centralized critics or structural biases. CTDE methods such as MADDPG [11] and COMA [5] stabilize learning through centralized information, while value-decomposition approaches such as VDN [17] and QMIX [14] impose additive or monotonic structures to promote cooperation. These techniques are highly effective but rely on explicit architectural assumptions that constrain the joint value landscape [12, 15]. Consequently, they offer limited insight into how coordination might emerge—or collapse—*without* centralized bias.

Emergent coordination has also been examined from perspectives including evolutionary games [19, 3], social dilemmas [9], and collective multi-agent systems [8, 4, 16]. These studies suggest that MARL behaviors can resemble physical collective phenomena. However, because coordination is actively enforced in these frameworks, they offer limited insight into how coordination might arise—or fail spontaneously—in fully decentralized settings *without centralized critics, structural assumptions, or explicit coordination mechanisms*.

In this work, we revisit Independent Q-Learning (IQL), which provides a clean testbed for examining emergent coordination because it operates *without imposing centralized inductive biases or decompositional structures*.

Although IQL is often regarded as unstable or limited [18, 6], it is precisely its lack of centralized critics, structural decomposition, or enforced cooperation that makes it an ideal setting for studying spontaneous coordination. By scanning a broad range of (L, ρ) conditions, we show that this minimal MARL system exhibits a rich phase structure driven by distributional non-stationarity.

Our contributions are as follows:

1. **A phase map of coordination and stability.** Combining the coordination success rate (CSR) with a TD-error–variance stability index S , we identify coordinated, fragile, and jammed/disordered regimes separated by a *double* instability ridge in the (L, ρ) plane.
2. **Kernel drift as a mechanism for MARL non-stationarity.** Temporal analysis of TD-error variance and gradient-norm variance shows that instability correlates strongly with kernel drift—the time-varying drift of the effective transition kernel induced by others’ policy updates.
3. **Synchronization as a requirement for sustained coordination.** Arrival-time spread and co-reach statistics reveal that temporal synchronization is necessary for maintaining coordination and that insufficient synchronization characterizes the fragile regime.
4. **Spontaneous coordination in decentralized MARL.** Even without centralized bias, certain scale–density combinations yield stable coordi-

nated behavior, suggesting that decentralized MARL can exhibit phase-transition-like phenomena.

Together, these results provide a unified perspective in which emergent coordination, fluctuation, and collapse arise from interactions among scale, density, and kernel drift, forming a coherent phase structure.

2 Methods

This section describes the environment, learning setup, and evaluation metrics used to analyze emergent coordination, non-stationarity, and kernel drift in fully decentralized Independent Q-Learning (IQL).

2.1 Environment: Grid-Based Multi-Agent Navigation

We use an $L \times L$ grid world with $L \in \{8, 16, 24, 32\}$, containing N agents and a single goal placed without overlap. Agent density is defined as

$$\rho_{\text{agents}} = N/L^2, \quad \rho \in \{0.03125, 0.0625, 0.125, 0.25, 0.5\},$$

excluding ($L = 32, \rho = 0.5$) due to computational cost.

Agents choose from

$$A = \{\text{stay}, \text{up}, \text{down}, \text{left}, \text{right}\},$$

and receive -0.005 per step, $+1$ upon reaching the goal, and a **hold** action thereafter which no longer affects the transition kernel.

Episodes terminate once the accumulated reward reaches the target score $0.8N$, or when the maximum horizon of $8L$ steps is reached. Each condition is trained for up to 1500 episodes. A single goal ensures that changes in (L, ρ) naturally induce sparsity and congestion.

2.2 Learning Algorithm: Parameter-Shared Double DQN

All agents share parameters and learn via Double DQN. Given online parameters θ and target parameters θ^- , the TD target is

$$a' = \arg \max_a Q_\theta(\tilde{s}, a), \quad y = r + \gamma(1 - d) Q_{\theta^-}(\tilde{s}, a').$$

We optimize the Huber loss

$$L(\theta) = \mathbb{E}[\ell_{\text{Huber}}(y - Q_\theta(s, a))],$$

and update the target network via Polyak averaging:

$$\theta^- \leftarrow (1 - \tau)\theta^- + \tau\theta.$$

Exploration follows exponentially decaying ϵ -greedy, and evaluation uses $\epsilon = 0$. Parameter sharing improves sample efficiency while retaining decentralized execution.

2.3 Network, Optimization, and Replay Buffer

We use a two-layer MLP (128 units each, ReLU), gradient clipping, and the Adam optimizer (PyTorch defaults $\beta_1 = 0.9, \beta_2 = 0.999$). Full hyperparameters appear in Appendix 1.

A shared replay buffer \mathcal{D} stores 10^5 transitions; training updates start once it contains at least 1,500 transitions, with a mini-batch size of 64. To permit spontaneous role differentiation despite parameter sharing, each transition augments observations with a one-hot agent identifier:

$$o_{\text{ID}}^{(i)} = [o^{(i)} \parallel \text{ID}_i].$$

2.4 Phase Structure: Coordination and Drift Metrics

We characterize each (L, ρ) condition using two axes:

- **Cooperative Success Rate (CSR)**,
- **Stability Index S** based on TD-error variance.

2.4.1 Cooperative Success Rate (CSR)

During evaluation episodes K_{eval} (with $\epsilon = 0$),

$$\text{CSR}(L, \rho) = \frac{1}{K_{\text{eval}}} \sum_{k=1}^{K_{\text{eval}}} \mathbf{1}\{\text{all agents reached}\}.$$

2.4.2 Stability Index S : TD-Error Variance as Kernel Drift Proxy

Non-stationarity arises from temporal drift in the effective transition kernel P_i^t caused by policy updates of other agents. We estimate this effect using the per-episode variance of the TD error

$$\delta_t = y_t - Q(s_t, a_t),$$

denoted $v(L, \rho)$. Normalizing by the maximum variance across all conditions,

$$S(L, \rho) = 1 - \frac{v(L, \rho)}{v_{\text{max}}},$$

where high S indicates weak drift (stable) and low S indicates strong drift (unstable). Gradient-norm variance is analyzed in the Appendix.

2.4.3 Reference Point and Phase Distance

Thresholds are defined as the 60th percentiles of CSR and S :

$$\tau_{\text{CSR}} = \text{Perc}_{60}(\text{CSR}), \quad \tau_S = \text{Perc}_{60}(S).$$

The phase distance is

$$d_{\text{phase}}(L, \rho) = \sqrt{(\text{CSR} - \tau_{\text{CSR}})^2 + (S - \tau_S)^2}. \quad (1)$$

This yields three regimes:

- **Stable Coordinated Phase:** high CSR, high S ; kernel drift suppressed.
- **Fragile Transitional Region:** near the minimum of d_{phase} ; coordination fluctuates.
- **Jammed / Disordered Phase:** low CSR, low S ; congestion-induced stagnation.

2.5 Experimental Setup

2.5.1 Training

- 1,500 episodes,
- Maximum steps = $8L$,
- 50 random seeds (mean and 95% confidence interval),
- Learning rate 1.5×10^{-4} , discount $\gamma = 0.95$, Polyak $\tau = 10^{-3}$.

2.5.2 Evaluation Conditions

CSR is computed during a separate evaluation phase using greedy policies ($\epsilon = 0$). In contrast, the stability index S is derived from TD-error statistics recorded *during training*, under the standard ϵ -greedy exploration schedule. These two quantities are then combined to construct the phase maps.

2.5.3 Compute Resources

Experiments were conducted on the University of Edinburgh’s *Eddie* HPC cluster (Rocky Linux 9, Altair GridEngine).

3 Results

3.1 Global Structure of Coordination and Non-coordination: CSR–Stability Phase Map

Fig. 1 shows CSR and the stability index S (TD-error–variance-based) for 19 conditions ($L \in \{8, 16, 24, 32\}$ and $\rho \in \{0.03125, \dots, 0.5\}$), computed from the last 25% of training episodes.

At small scales and low densities, CSR and S are simultaneously high, indicating that coordinated behavior can emerge even under fully decentralized IQL.

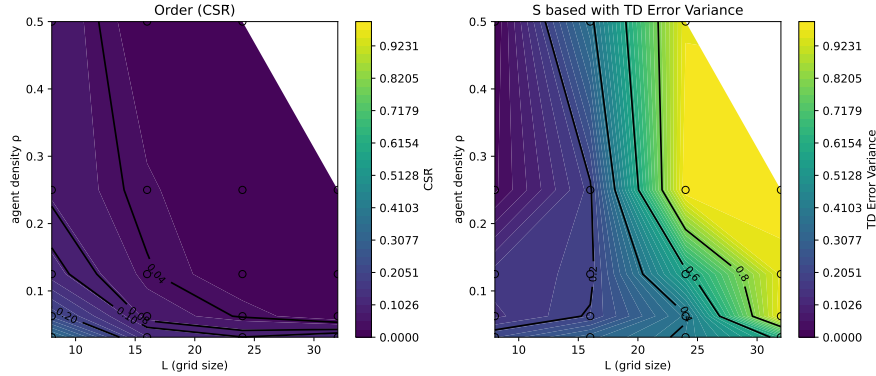


Figure 1: Cooperative success rate (CSR, left) and stability index S (right) for all (L, ρ) conditions, computed from the last 25% of episodes. Coordination and stability occur only at small scales and low densities, while both metrics collapse sharply as scale or density increases. The low- S region marks strong non-stationarity and forms an Instability Ridge.

As density increases, CSR drops sharply toward zero across all scales, while TD-error variance increases. This pattern suggests that congestion-induced exploration difficulty amplifies kernel drift, making the policy updates increasingly unstable.

At fixed density, increasing scale also degrades both CSR and S monotonically, revealing that coordination becomes progressively fragile as exploration cost and agent interactions intensify. Taken together, the CSR- S axes provide a clear separation between coordinated, partially coordinated, and non-coordinated regimes, and illustrate that coordination in independent MARL is strongly dependent on environmental scale and density, collapsing rapidly with growing non-stationarity.

3.2 Phase Geometry via Distance from the Coordinated Attractor

Fig. 2 visualizes the normalized distance d_{phase} from Eq. 1. A notable feature is the presence of two contour lines around $d_{\text{phase}} \approx 0.4$, forming a double Instability Ridge that separates coordinated and non-coordinated behavior.

This structure partitions the (L, ρ) space into three regimes:

- **Coordinated & Stable Phase** ($d_{\text{phase}} > 0.4$ in the low- L , low- ρ region): High CSR and high S , with weak drift and stable synchronization. Convergence toward a coordinated attractor is consistently observed.
- **Fragile Region** ($d_{\text{phase}} < 0.4$, between the two ridges): Synchronization and collapse alternate, producing non-monotonic and transitional dynamics. Kernel drift competes with synchronization, yielding only temporary

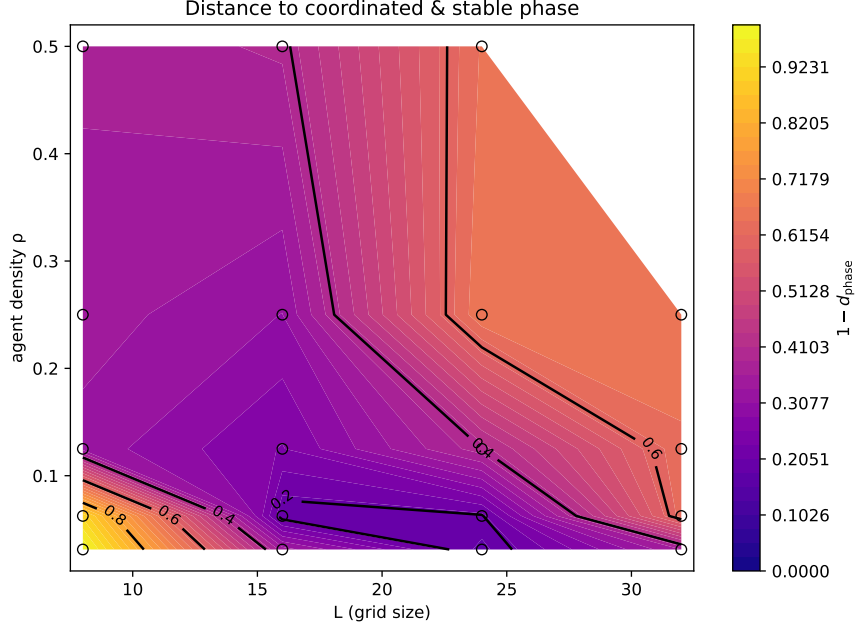


Figure 2: Phase geometry based on the normalized distance d_{phase} . Two contour lines near $d_{\text{phase}} \approx 0.4$ form a double Instability Ridge. The low- L , low- ρ region corresponds to coordinated and stable behavior, the region between the ridges to a fragile transitional regime, and larger scales to jammed/disordered outcomes.

or partial coordination.

- **Jammed/Disordered Phase** ($d_{\text{phase}} > 0.4$ at larger L): CSR is near zero and S remains low. Although kernel drift weakens, gradient noise later dominates, leading to irreversible jammed/disordered behavior.

This three-part geometry indicates that IQL exhibits bifurcating dynamics between a coordinated attractor and drift-driven instability, resembling a phase transition.

3.3 Synchronization and Partial Coordination Near the Instability Ridge

Fig. 3 compares synchronization (arrival-time spread) and coordination (co-reach) dynamics for two ridge-adjacent conditions.

In the pre-ridge case, drift remains weak, spread collapses quickly, and co-reach increases smoothly—reflecting stable convergence to a coordinated attractor.

On the ridge, spread and co-reach repeatedly undergo collapse–recovery cycles. This oscillatory pattern reflects competition between synchronization and

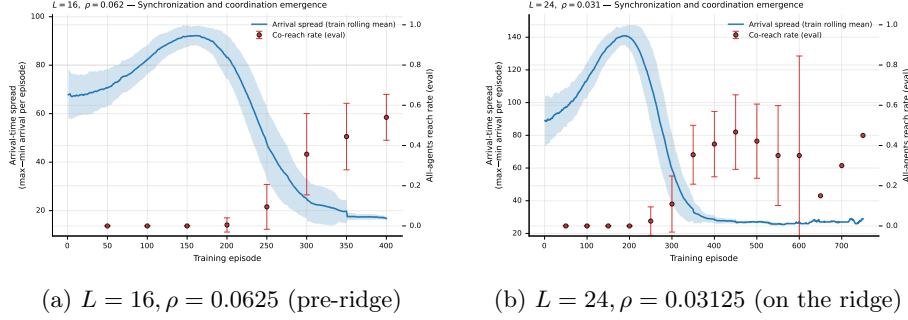


Figure 3: Temporal profiles of arrival-time spread (synchronization) and co-reach rate for two representative conditions near the Instability Ridge. Pre-ridge ($L=16, \rho=0.0625$) shows rapid convergence toward coordination, whereas on-ridge ($L=24, \rho=0.03125$) exhibits alternating collapse–recovery cycles.

kernel drift, giving rise to fragile coordination that does not persist.

These observations support the view that synchronization is a necessary condition for coordination, and that drift near the ridge disrupts this mechanism. (High-density cases transition into jammed or oscillatory behavior; see Appendix Fig. 9.)

3.4 Kernel Drift vs. Gradient Noise: Two Forms of Instability

Fig. 4 shows that the dominant source of instability changes with density:

- **Near-ridge** ($\rho = 0.0625$): TD variance grows throughout training, and gradient variance increases later. Drift does not settle, preventing convergence to a coordinated fixed point.
- **On-ridge** ($\rho = 0.03125$): TD variance remains high, and gradient variance increases gradually, producing fragile oscillatory dynamics.
- **Outside-ridge** ($\rho = 0.125$): Both variances saturate early, reflecting premature lock-in to incomplete patterns rather than coordination.
- **High densities** ($\rho \geq 0.25$): Drift weakens but gradient variance diverges in late training, causing loss of update coherence and jammed/disordered outcomes.

These trends indicate that the Instability Ridge corresponds to a transition from drift-dominated to gradient-noise-dominated instability, offering a unified perspective on coordination, partial coordination, and collapse. Complete time-series statistics of TD error, gradient norms, and their variances across all conditions are provided in Appendix Fig. 10, with further discussion in Appendix A.3.

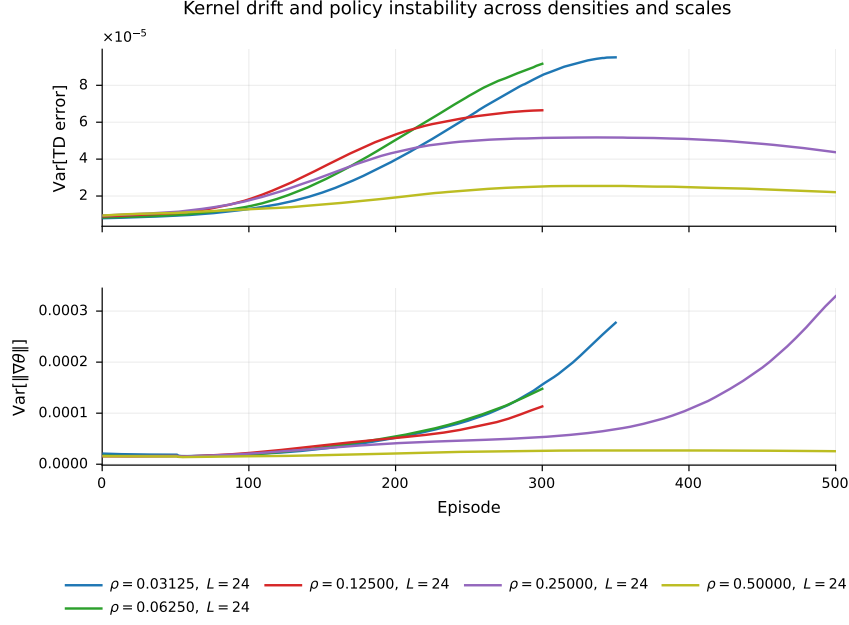


Figure 4: TD-error variance (top) and gradient-norm variance (bottom) for $L = 24$ across densities. The dominant source of instability differs by density: persistent variance growth near the Ridge, oscillatory behavior on the Ridge, early saturation outside the Ridge, and late-stage divergence at high densities. These four densities were selected because they represent pre-ridge, on-ridge, post-ridge, and high-density regimes, respectively.

3.5 Macroscopic Statistics of the Coordination Transition: Synchronization and Density Thresholds

Fig. 5a shows that only conditions with small spread achieve high co-reach. The absence of points with large spread and high co-reach indicates that non-synchronized coordination is not observed. Large fluctuations near the ridge again reflect drift-synchronization competition.

Fig. 5b shows the effective coordinated throughput $\rho_{\text{eff}} = \rho_{\text{agents}} \cdot \text{CSR}$, which reveals scale-dependent critical densities $\rho_{\text{crit}}(L)$. At medium and large scales, throughput peaks at intermediate densities and drops sharply at high densities. This behavior mirrors the synchronization and coordination collapse observed in Fig. 3, appearing as a macroscopic phase transition.

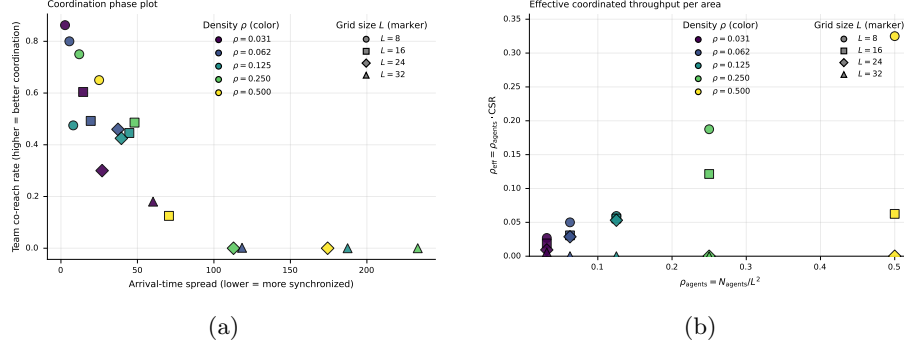


Figure 5: (Left) Relationship between arrival-time spread and co-reach across all conditions. High co-reach occurs only when spread is small; large fluctuations appear near the Instability Ridge. (Right) Effective coordinated throughput $\rho_{\text{eff}} = \rho_{\text{agents}} \text{CSR}$, showing scale-dependent critical densities where throughput peaks and then declines.

4 Discussion

These findings position kernel drift not merely as a correlate of instability, but as a mechanistic driver whose growth rate predicts the onset of fragile and jammed regimes.

Non-stationarity is a central challenge in multi-agent reinforcement learning (MARL), arising from the fact that each agent’s policy update modifies the effective transition dynamics faced by all others. Prior work has examined how replay-induced drift [6], non-convex–non-concave game structures [20], and regularization or learning-rate control [13] contribute to instability. This interpretation complements prior taxonomies of non-stationarity Hernandez-Leal et al. (2019), but identifies kernel drift as a unifying mechanism driving the observed phase structure.

Our results suggest that these diverse forms of non-stationarity can be viewed through a single unifying mechanism: *kernel drift*, the distributional drift in each agent’s effective transition kernel. For agent i , the kernel

$$P_i^t(s'|s, a_i) = \sum_{a_{-i}} P(s'|s, a_i, a_{-i}) \prod_{j \neq i} \pi_j^t(a_j|s) \quad (2)$$

fluctuates as other agents update, producing a drift term ΔP_i^t . These fluctuations amplify the variance of TD targets, create a mismatch between replayed and current dynamics, and push policy updates away from stable fixed points. This mechanism is qualitatively consistent with observations in mean-field MARL, replicator dynamics, and actor–critic oscillations.

The resulting phase diagram reveals three regimes shaped by the interaction between kernel drift and gradient noise:

- **Stable coordinated phase:** At low density and small scale, kernel drift is weak, synchronization holds, and learning consistently converges toward a coordinated attractor.
- **Fragile transitional phase:** Near the Instability Ridge, kernel drift grows critically and competes with synchronization, producing alternating collapse–recovery cycles and highly non-monotonic dynamics.
- **Failure phase (jammed/disordered):** Outside the ridge and at higher densities, kernel drift weakens but gradient noise becomes dominant, leading to early lock-in or jammed/disordered behavior from which coordination does not recover.

Notably, with a single shared goal, increasing the scale L reduces collision-induced asymmetries, weakening kernel-drift sources and causing the dynamics at $L = 24$ and $L = 32$ to converge toward similar regimes—an observation that further motivates the ID-removal ablation examined below.

The ablation in Appendix A.4 further clarifies the mechanism: removing the agent ID restores full input symmetry, yielding identical transition kernels for all agents, $\Delta P_{\text{brk}}^t \equiv 0$. Kernel drift is therefore canceled over time. TD-error and gradient-norm variances remain low and flat, CSR becomes density-independent, and neither coordinated nor transitional phases emerge. This provides direct evidence that the observed phase structure is not an artifact of algorithmic bias, but rather arises from small symmetry-breaking differences that accumulate through distributional interaction.

These observations align with prior work on role emergence or symmetry breaking in MARL [2], evolutionary games [19, 3], and collective robotics [4]. The phase-transition-like structure observed here appears to be a spontaneous phenomenon rooted in distributional interaction rather than an explicit coordination mechanism.

Overall, kernel-drift dynamics highlight how MARL learning trajectories may be understood as collective processes shaped by distribution-level feedback. A rigorous theoretical account remains open, but our empirical findings point toward promising connections with distributional or mean-field stability analyses.

5 Conclusion

This work examined how independent multi-agent reinforcement learning (IQL), despite having no centralized critic or structural coordination bias, exhibits systematic patterns of emergent coordination, fragility, and failure across environment scale L and agent density ρ . By combining cooperative success (CSR) with a stability index based on TD-error variance, we constructed phase diagrams that reveal three regimes:

- **Coordinated & stable:** synchronization holds and learning converges toward a coordinated attractor.

- **Fragile coordination:** kernel drift competes with synchronization, producing collapse–recovery oscillations.
- **Failure / jammed–disordered:** kernel drift weakens but gradient noise dominates, preventing recovery of coordination.

The Instability Ridge emerges as a sharp transition boundary where kernel drift amplifies and the learning dynamics switch from drift-dominated to gradient-noise-dominated behavior. These observations indicate that MARL learning trajectories may exhibit phase-transition-like structure shaped by distributional interaction.

Ablation experiments further showed that removing agent IDs eliminates the entire phase structure: TD-error variance stays uniformly low, CSR becomes flat across densities, and neither coordination nor fragility nor collapse emerges. This supports the interpretation that the observed dynamics arise from spontaneous symmetry breaking and distributional interaction, rather than from explicit coordination mechanisms.

Taken together, our results suggest that coordination in MARL may be understood as a spontaneous, distribution-driven phenomenon governed by physical-style parameters such as scale, density, kernel drift, and gradient noise. Treating kernel drift as a distributional fluctuation in effective transition dynamics provides a unifying view on non-stationarity and offers a foundation for future stability analyses based on distributional or mean-field models.

Implications for Multi-Agent Economics and Financial Markets

The structure observed here—emergence, fluctuation, and breakdown of coordination driven by distributional interactions—has natural parallels in multi-agent economics and market microstructure. Many-agent systems in which actions modify the distributional environment that, in turn, shapes optimal behavior share the same feedback structure as MARL.

The phenomena observed in this work— emergence of coordination (price stability), transient or partial synchronization, critical-density collapse, and drift-induced deviation from equilibrium— suggest that scale, density, and distributional drift may provide useful explanatory principles for when markets stabilize and when they destabilize.

These connections point toward future work on the stability of many-agent financial models, critical phenomena in order-book dynamics, and MARL-based market simulations. Understanding coordination and its breakdown through the lens of phase geometry may offer a promising direction for integrating ideas from AI, economics, and statistical physics.

6 Acknowledgments

A.Y. received support from a sponsored research agreement with Intel Corporation, which provided partial funding for this research.

References

- [1] Kai Arulkumaran et al. “Deep Reinforcement Learning: A Brief Survey”. In: *IEEE Signal Processing Magazine* 34.6 (2017), pp. 26–38.
- [2] Bowen Baker et al. *Emergent Tool Use From Multi-Agent Autocurricula*. 2020. arXiv: 1909.07528 [cs.LG].
- [3] Daan Bloembergen et al. “Evolutionary dynamics of multi-agent learning: a survey”. In: *J. Artif. Int. Res.* 53.1 (May 2015), pp. 659–697.
- [4] Dmitry Bratsun and Kirill Kostarev. “Phase Transition in a Dense Swarm of Self-Propelled Bots”. In: *Fluid Dynamics & Materials Processing* 20.8 (2024), pp. 1785–1798.
- [5] Jakob Foerster et al. “Counterfactual Multi-Agent Policy Gradients”. In: *AAAI Conference on Artificial Intelligence*. 2018.
- [6] Jakob Foerster et al. *Stabilising Experience Replay for Deep Multi-Agent Reinforcement Learning*. 2018. arXiv: 1702.08887 [[cs.AI] (<http://cs.ai/>)].
- [7] Pablo Hernandez-Leal et al. *A Survey of Learning in Multiagent Environments: Dealing with Non-Stationarity*. 2019. arXiv: 1707.09183.
- [8] Jason Hindes et al. “Critical transition for colliding swarms”. In: *Physical Review E* 103.6 (June 2021).
- [9] Joel Z. Leibo et al. *Multi-agent Reinforcement Learning in Sequential Social Dilemmas*. 2017. arXiv: 1702.03037 [cs.MA].
- [10] Yuxi Li. “Deep Reinforcement Learning: An Overview”. In: *arXiv preprint arXiv:1701.07274* (2017).
- [11] Ryan Lowe et al. “Multi-Agent Actor-Critic for Mixed Cooperative-Competitive Environments”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2017.
- [12] Anuj Mahajan et al. *MAVEN: Multi-Agent Variational Exploration*. 2020. arXiv: 1910.07483 [cs.LG].
- [13] Georgios Papoudakis et al. *Dealing with Non-Stationarity in Multi-Agent Deep Reinforcement Learning*. 2019. arXiv: 1906.04737 [cs.LG].
- [14] Tabish Rashid et al. *QMIX: Monotonic Value Function Factorisation for Deep Multi-Agent Reinforcement Learning*. 2018. arXiv: 1803.11485 [cs.LG].
- [15] Tabish Rashid et al. *Weighted QMIX: Expanding Monotonic Value Function Factorisation for Deep Multi-Agent Reinforcement Learning*. 2020. arXiv: 2006.10800 [cs.LG].

- [16] Michael Rubenstein, Alejandro Cornejo, and Radhika Nagpal. “Programmable self-assembly in a thousand-robot swarm”. In: *Science* 345.6198 (2014), pp. 795–799.
- [17] Peter Sunehag et al. *Value-Decomposition Networks For Cooperative Multi-Agent Learning*. 2017. arXiv: 1706.05296 [[cs.AI](http://cs.ai/)].
- [18] Ming Tan. “Multi-Agent Reinforcement Learning: Independent versus Cooperative Agents”. In: *Proceedings of the Tenth International Conference on Machine Learning (ICML 1993)*. Morgan Kaufmann, 1993, pp. 330–337.
- [19] K.P. Tuyls and A. Nowé. “Evolutionary Game Theory and Multi-Agent Reinforcement Learning”. English. In: *Knowledge Engineering Review* 20(01) (Jan. 2005), pp. 63–90.
- [20] Kaiqing Zhang, Zhuoran Yang, and Tamer Başar. *Multi-Agent Reinforcement Learning: A Selective Overview of Theories and Algorithms*. 2021. arXiv: 1911.10635 [cs.LG].

A Appendix

This appendix provides additional analyses and ablation studies that complement the main text. In particular, we present: (i) full hyperparameter specifications, (ii) extended phase maps, (iii) additional examples of synchronization–coordination dynamics, (iv) scale-dependent properties of kernel-drift indicators, and (v) ablations in which agent identifiers are removed to examine the collapse of the phase structure.

A.1 Hyperparameter Details

A.2 Stability Index Based on Gradient-Norm Variance and the Corresponding Phase Maps

In addition to the TD-error–variance stability index S_{TD} , we also construct a gradient-norm–based stability index S_{∇} and visualize the corresponding phase maps in Figs. 6 and 7.

At the global level, S_{∇} recovers the same broad structure as S_{TD} : a coordinated and stable phase at small scale and low density, and a jammed/disordered phase at high density. However, the phase boundaries become noticeably less sharp under S_{∇} . In particular, at $L = 24$ and $\rho = 0.25$, S_{∇} exhibits a pronounced local drop (Fig. 8a), whereas S_{TD} shows a clear peak around $\rho \approx 0.06$ – 0.08 (Fig. 8b), precisely identifying the Instability Ridge.

This discrepancy arises because gradient-norm variance is sensitive not only to kernel drift—the source of distributional non-stationarity—but also to curvature of the local loss landscape and optimizer momentum. For instance, at $(L, \rho) = (24, 0.25)$, CSR and arrival-time spread remain in a metastable semi-

Table 1: Hyperparameters used in all experiments.

| Item | Symbol | Value / Setting |
|--|---------------------------|---|
| Number of episodes (max) | – | 1500 |
| Max steps per episode | T_{\max} | $8L$ |
| Number of random seeds per condition | – | 50 |
| Learning rate | α | 1.5×10^{-4} |
| Discount factor | γ | 0.95 |
| Polyak coefficient (soft target update) | τ | 10^{-3} |
| Optimizer | – | Adam ($\beta_1 = 0.9$, $\beta_2 = 0.999$)* |
| Batch size | B | 64 |
| Replay buffer capacity | $ \mathcal{D} $ | 10^5 |
| Warm-up transitions | – | 1500 |
| Initial exploration rate | ϵ_{start} | 1.0 |
| Minimum exploration rate | ϵ_{min} | 0.01 |
| Exploration decay (per episode, exponential) | λ | 0.98 |
| Exploration during evaluation | ϵ_{eval} | 0 (greedy) |
| Network architecture | – | Input-FC(128, ReLU)-FC(128, ReLU)-Output |
| Activation function | – | ReLU |
| Gradient clipping | – | $\ \nabla\ _2 \leq 1.0$ |
| Loss function | – | Huber loss (smooth L1) |

*PyTorch’s default hyperparameters for Adam were used.

jammed configuration, yet gradient norms fluctuate due to fine-scale pattern adjustments, which blurs the phase boundary at the map level.

Taken together, these observations indicate that S_{TD} is the more reliable primary indicator for detecting the growth of kernel drift and for identifying the Instability Ridge. We therefore present the S_{∇} -based maps only as supplementary analysis.

A.3 Kernel Drift Dynamics Across Scales

Fig. 10 reports the TD error, gradient norm, and their variances across all (L, ρ) conditions (50 seeds; 95% CI). The main trends are as follows.

1. TD-error mean Across all settings the mean TD error quickly approaches zero, indicating that the Bellman update itself stabilizes early.
2. TD-error variance $\text{Var}[\text{TD}]$ depends primarily on scale L , not density. For small grids ($L = 8, 16$) the variance is larger and more dynamic, whereas for large grids ($L = 24, 32$) it remains small and often nearly flat—especially in jammed/disordered regimes such as $(L = 24, \rho > 0.125)$ and $(L = 32, \rho > 0.0625)$.

3. **Gradient norm and its variance** Gradient statistics also exhibit strong scale dependence. Large grids show late-episode divergence at specific densities (e.g., $L = 24, \rho = 0.25$ and $L = 32, \rho = 0.125$), indicating gradient-noise-dominated instability even when drift is weak.
4. **Density effects** Density modulates the timing and magnitude of fluctuations but does not produce a consistent monotonic trend across scales. Scale L and drift-synchronization competition near the Instability Ridge (discussed in Fig. 4) dominate the global behavior.
5. **Overall structure** These patterns complement the phase diagram in Fig. 2, showing that:
 - drift suppression yields stable coordination,
 - drift-synchronization competition produces fragile dynamics, and
 - weak drift with strong gradient noise leads to jammed/disordered outcomes.

A.4 Ablation: Removing Agent ID and Symmetry Breaking

To isolate the role of symmetry breaking in the emergence of kernel drift and the phase structure reported in the main text, we conduct ablations in which agent identifiers (IDs) are removed. Without IDs, all agents share identical observation and policy spaces; hence their policy mappings become fully symmetric.

Following the definition of the effective transition kernel in Eq. 2, removing agent IDs makes all agents share identical observation and policy spaces. As a result, each agent i satisfies

$$P_i^t(s'|s, a_i) \approx \bar{P}, \quad \Delta P_i^t := P_i^t - \bar{P} \approx 0,$$

and the asymmetry-driven component ΔP_{brk}^t in the decomposition $\Delta P_i^t = \Delta P_{\text{sym}}^t + \Delta P_{\text{brk}}^t$ vanishes. Kernel drift therefore cannot accumulate.

Learning dynamics without IDs. The TD-error mean decreases rapidly at the beginning, similar to the ID-present case, but then exhibits a secondary rise followed by a slow, noise-like decay. The TD-error variance shows an early peak comparable to the ID-present setting, yet does not settle; instead, it undergoes repeated small-amplitude fluctuations. These patterns indicate that distributional drift is suppressed, and that the remaining non-stationarity is dominated by weak exploration noise arising from policy symmetry.

Gradient-norm statistics. Removing IDs produces a sharper distinction: gradient magnitudes become substantially larger—often several times those observed with IDs—and the corresponding variance rises steeply before decaying slowly. When all agents share an identical policy, update directions lose diversity, so even small perturbations are coherently amplified, creating persistent and exaggerated update noise.

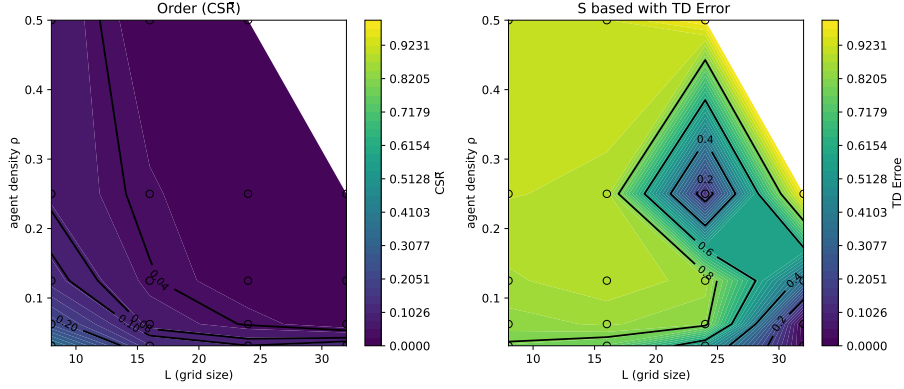


Figure 6: Phase map using the gradient-norm-based stability index S_{∇} . Left: CSR. Right: S_{∇} . Global trends resemble the TD-based map, but ridge boundaries become noticeably less distinct due to optimizer- and curvature-induced fluctuations.

Consequences for phase behavior. Overall, removing IDs yields a characteristic structure in which (i) TD-error variance remains suppressed, confirming the disappearance of kernel drift, while (ii) enlarged gradient norms and heightened gradient-norm variance show that update noise becomes the dominant driver of instability. This produces superficially similar “non-coordination” outcomes to the jammed/disordered regime observed with IDs, but the underlying mechanism is fundamentally different: the instability arises not from drift, but from homogeneous, noise-amplifying updates due to complete policy symmetry.

These findings demonstrate that the phase transitions in the main text are driven by

$$\text{kernel drift} \propto \Delta P_{\text{brk}}^t,$$

i.e., by small but persistent inter-agent asymmetries. Removing such asymmetries suppresses kernel drift and eliminates both coordinated and fragile regimes. This result is consistent with prior work on symmetry breaking in learning dynamics [19, 3], role emergence in MARL [2], and self-organized differentiation in collective robotic systems [4]. It supports the interpretation that the coordination transition observed in decentralized MARL arises as a spontaneous distribution–interaction phenomenon rather than an architectural artifact.

This appendix is self-contained and provides supplementary analyses and ablations that support and extend the results presented in the main text.

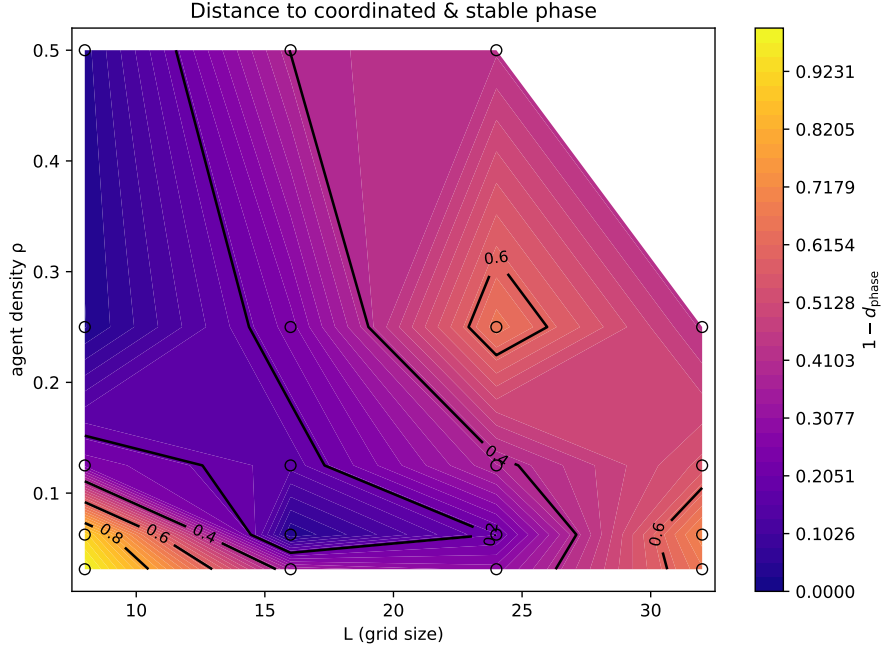


Figure 7: Phase-distance map constructed from the gradient-norm-based stability index. A ridge structure remains visible, but is substantially blurred compared to the TD-based counterpart.

B Appendix Figures

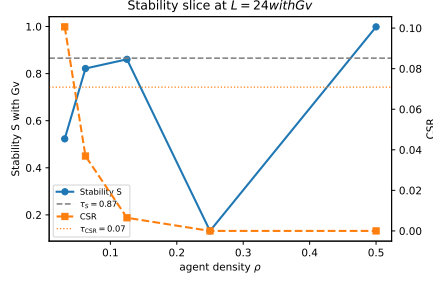
B.1 Phase Maps Based on Gradient-Norm Variance

B.2 Stability Index Slices at Fixed Scale

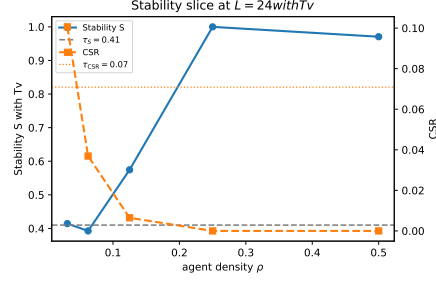
B.3 Synchronization and Coordination Dynamics

B.4 Kernel Drift and Gradient-Noise Regimes

B.5 Effect of Removing Agent IDs

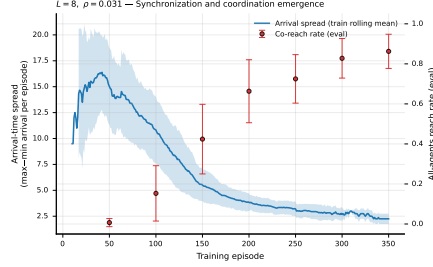


(a) S_{∇} at $L = 24$.

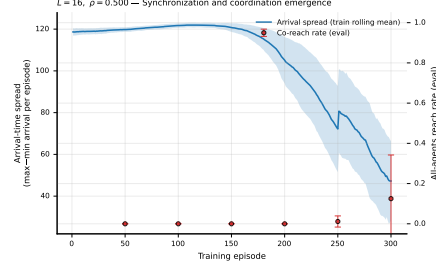


(b) S_{TD} at $L = 24$.

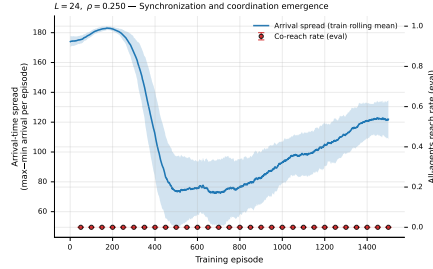
Figure 8: Stability-index slices at $L = 24$. Left: S_{∇} shows a local collapse at $\rho = 0.25$, which obscures the phase boundary. Right: S_{TD} exhibits a distinct peak near $\rho \approx 0.06$ – 0.08 , sharply locating the Instability Ridge.



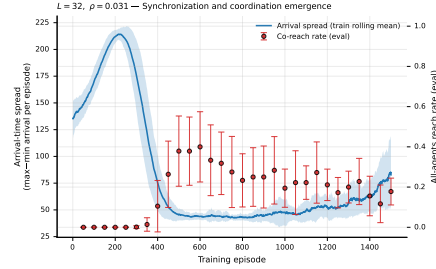
(a) $L = 8$, $\rho = 0.03125$



(b) $L = 16$, $\rho = 0.50$



(c) $L = 24$, $\rho = 0.25$



(d) $L = 32$, $\rho = 0.03125$

Figure 9: Time-series dynamics of arrival-time spread (synchronization) and co-reach rate (coordination) across four representative conditions. (a) $L = 8$, $\rho = 0.03125$: rapid synchronization and stable coordinated behavior. (b) $L = 16$, $\rho = 0.50$: approaching the Instability Ridge from the high-density side, where synchronization is hindered and co-reach exhibits large fluctuations. (c) $L = 24$, $\rho = 0.25$: a jammed/disordered regime outside the Ridge, with persistently large spread and low co-reach. (d) $L = 32$, $\rho = 0.03125$: approaching the Ridge from the low-density, large-scale side, showing partial synchronization with drift-induced fluctuations.

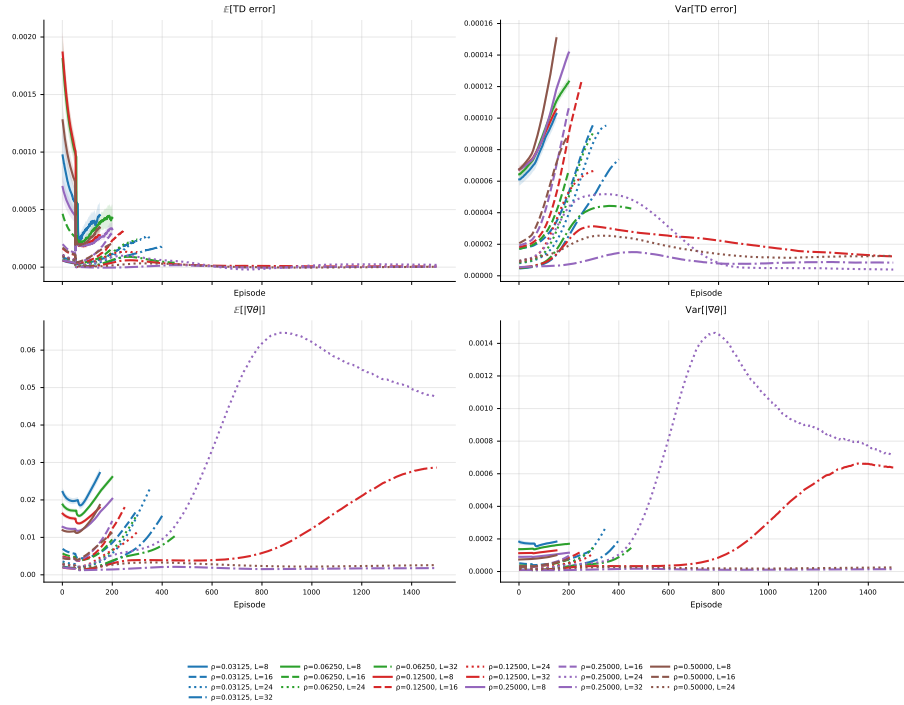


Figure 10: TD-error and gradient-norm statistics across all (L, ρ) conditions. Both metrics exhibit strong scale dependence: larger grids ($L = 24, 32$) suppress TD-error variance except near the Instability Ridge, whereas gradient-norm variance shows late-episode divergence only in jammed/disordered regimes (e.g., $L = 24, \rho = 0.25$; $L = 32, \rho = 0.125$). These patterns reflect a transition from drift-dominated to gradient-noise-dominated instability. Averaged over 50 seeds; 95% CIs shown.

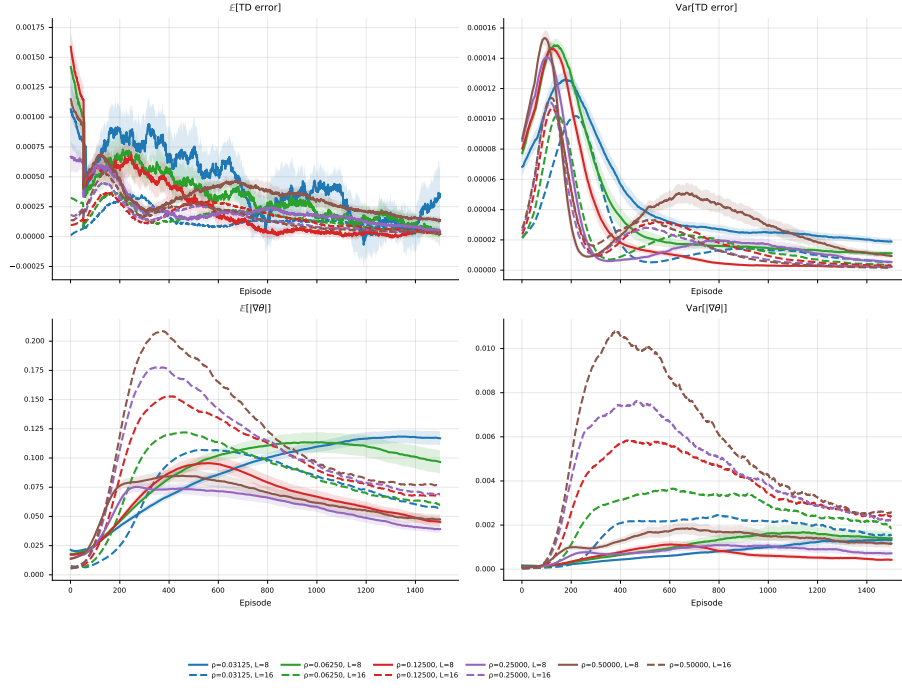


Figure 11: Mean and variance of the TD error (top) and gradient norm (bottom) for all (L, ρ) conditions in the ID-removed setting. Removing agent IDs suppresses the asymmetry-driven component of kernel drift, resulting in low and non-accumulating TD-error variance. In contrast, gradient norms remain substantially larger and exhibit pronounced early-episode fluctuations, reflecting update-noise amplification under complete policy symmetry. Curves show averages over 25 seeds with 95% confidence intervals for $L \in \{8, 16\}$ and $\rho \in \{0.03125, \dots, 0.5\}$.

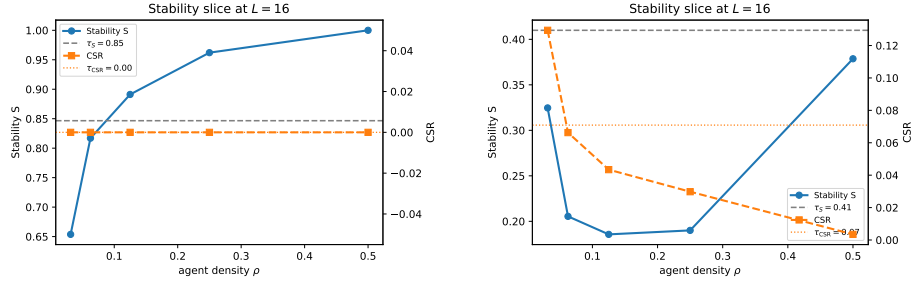


Figure 12: Left: **Without agent IDs**—CSR and the stability index S remain flat across densities; no coordinated, fragile, or jammed phases emerge. Right: **With IDs**—the coordinated, fragile, and jammed regimes re-emerge, confirming that symmetry breaking is required for sustaining kernel drift and for producing the phase structure observed in the main text.