# Interpretable Graph-Attention Collaboration: Adaptive Policies for Robust Multi-Agent Systems

Anonymous Author(s)

## ABSTRACT

Multi-agent systems increasingly rely on collaboration among autonomous agents, yet most deployed architectures employ fixed, hand-designed communication topologies such as star, chain, or fully connected graphs. We introduce *Interpretable Graph-Attention Collaboration* (IGAC), a framework that combines an adaptive communication topology with trust-weighted message passing for multi-agent collaborative reasoning. IGAC employs Gumbel-Sigmoid sampling to produce sparse, instance-specific binary collaboration graphs, attention-based message aggregation with SGD-trained projection parameters for interpretable information routing, and a Beta-distributed counterfactual trust mechanism for adversarial agent detection and isolation. We separate training (where projection parameters and trust are learned) from evaluation (where parameters are frozen), ensuring a rigorous experimental protocol. Across six experiments on collaborative state reconstruction tasks with up to 20 agents, the trained IGAC model reduces communication edges substantially compared to fully connected baselines while maintaining competitive reconstruction accuracy. Under adversarial conditions with 2 out of 6 agents compromised, IGAC with trust scoring achieves lower reconstruction error than fixed-topology baselines and detects adversarial agents via trust-based thresholding. Ablation studies confirm that both the learned topology and trust mechanism contribute independently to robustness.

## 1 INTRODUCTION

Multi-agent systems that collaborate through structured communication have demonstrated capabilities exceeding those of individual agents across a range of reasoning tasks [3, 15]. However, the collaboration topology—which agents communicate with which, and how information is aggregated—remains predominantly a design choice made by human engineers. Fixed topologies such as star (hub-and-spoke), chain (sequential), and fully connected graphs each impose structural assumptions that may not match the requirements of a given task instance [14].

This rigidity creates three interrelated challenges. First, fixed topologies cannot *adapt* to varying task demands, agent capabilities, or partial observability conditions. Second, when communication structure is predetermined, there is limited opportunity for *interpretability*: practitioners cannot understand why particular communication patterns emerged because they were imposed rather than learned. Third, fixed topologies are *vulnerable* to adversarial agents—a compromised node in a star topology can corrupt all communications, while a fully connected topology indiscriminately aggregates adversarial messages.

Wei et al. [14] identify the development of adaptive, interpretable collaboration policies robust to partial observability and adversarial conditions as a key open problem in agentic reasoning. We address this problem with *Interpretable Graph-Attention Collaboration* (IGAC), a framework built on three technical contributions:

(1) **Learned sparse topology via Gumbel-Sigmoid.** A meta-controller produces per-instance, per-step adjacency matrices by sampling binary edges through Gumbel-Sigmoid relaxation [5] over pairwise agent state similarities with a hard threshold. This yields genuinely sparse, binary communication graphs that adapt to the information structure of each problem instance.

(2) **Trained trust-weighted attention message passing.** Messages are aggregated along learned edges using scaled dot-product attention [11] with SGD-trained projection parameters, modulated by per-neighbor trust scores. Trust is modeled as Beta distributions updated via personalized counterfactual credit assignment [4], enabling principled detection of adversarial agents.

(3) **Interpretability through hard sparsity and attention.** The combination of binary sparse topology and peaked attention distributions provides two complementary levels of interpretability: structural (which edges are active) and functional (how much each message contributes to each agent's decision).

We evaluate IGAC on collaborative state reconstruction under controlled partial observability and adversarial agent injection, with a rigorous train/test separation, comparing against fixed-topology baselines and ablation variants across six experimental dimensions.

## 2 RELATED WORK

*Multi-Agent Communication Learning.* CommNet [10] introduced differentiable communication channels between reinforcement learning agents, enabling end-to-end learning of message content. TarMAC [2] added targeted communication through attention mechanisms, and MAGIC [8] employed graph attention for agent communication. These methods learn *what* to communicate but assume fixed topologies. IGAC extends this line by jointly learning the topology and training message aggregation parameters.

*Multi-Agent Reinforcement Learning.* QMIX [9], MAPPO [16], and MADDPG [7] provide centralized-training-decentralized-execution frameworks for cooperative and mixed settings. They address credit assignment at the value-function level but do not learn communication structure. Our counterfactual trust mechanism provides agent-level credit assignment that doubles as an adversarial detection signal.

*LLM-Based Multi-Agent Systems.* AutoGen [15] and related frameworks enable multi-agent conversations with predefined topologies. DyLAN [6] dynamically adjusts agent participation using per-step scoring, representing the closest existing work to topology learning. However, DyLAN lacks explicit interpretability mechanisms and adversarial robustness guarantees. Multi-agent debate [3] improves reasoning through structured disagreement but uses fixed two-agent or round-robin structures.

*Robust and Interpretable Policies.* Byzantine-tolerant consensus [1] provides robustness in distributed systems but assumes well-defined message semantics incompatible with free-form agent outputs. Programmatic policies [13] offer inherent interpretability but limited scalability. Graph Attention Networks [12] provide attention-based message passing over fixed graphs; IGAC extends this to learned, dynamic graphs with trust modulation.

## 3 METHOD

### 3.1 Problem Formulation

We consider $N$ agents that must collaboratively reconstruct a shared hidden state $\mathbf{s} \in \mathbb{R}^D$ from partial, noisy observations. Agent $i$ observes $\mathbf{o}_i = M_i \mathbf{s} + \boldsymbol{\epsilon}_i$, where $M_i \in \{0, 1\}^{D \times D}$ is a diagonal mask revealing a fraction $p$ of state dimensions, and $\boldsymbol{\epsilon}_i \sim \mathcal{N}(0, \sigma^2 I)$ is observation noise. A fraction $f$ of agents may be adversarial, replacing their observations with random noise to mislead collaborators.

The agents communicate over $R$ rounds through a dynamic collaboration graph $G_t = (V, E_t)$ where $V = \{1, \ldots, N\}$ and $E_t$ changes at each communication round. The collective goal is to minimize the reconstruction error $\|\hat{\mathbf{s}} - \mathbf{s}\|_2 / \|\mathbf{s}\|_2$.

### 3.2 Learned Topology via Gumbel-Sigmoid

At each communication round $t$, the meta-controller produces a binary adjacency matrix $A_t \in \{0, 1\}^{N \times N}$ from the current agent states $\mathbf{h}_1, \ldots, \mathbf{h}_N$. Edge logits are computed from pairwise cosine similarity:

$$\ell_{ij} = \frac{\mathbf{h}_i^\top \mathbf{h}_j}{\|\mathbf{h}_i\| \|\mathbf{h}_j\|} + \log \frac{\rho}{1 - \rho} \tag{1}$$

where $\rho$ is a sparsity target controlling the expected edge density. Each edge $(i, j)$ is sampled via the Gumbel-Sigmoid trick [5]:

$$\sigma_{ij} = \sigma\left(\frac{\ell_{ij} + g}{\tau}\right), \quad A_t[i, j] = \mathbb{1}[\sigma_{ij} > 0.5] \tag{2}$$

where $g$ is a Gumbel(0,1) sample and $\tau$ is a temperature parameter. The hard threshold at 0.5 produces genuinely binary, sparse graphs—edges are either fully active or fully inactive, ensuring that communication cost reflects actual message exchange.

### 3.3 Trained Attention Message Passing

Given the binary adjacency matrix $A_t$ and trust scores $T \in [0, 1]^{N \times N}$, messages are aggregated using scaled dot-product attention modulated by topology and trust:

$$\alpha_{ij} = \frac{A_t[i, j] \cdot T[i, j] \cdot \exp(\mathbf{q}_i^\top \mathbf{k}_j / \sqrt{d_k})}{\sum_{j'} A_t[i, j'] \cdot T[i, j'] \cdot \exp(\mathbf{q}_i^\top \mathbf{k}_{j'} / \sqrt{d_k})} \tag{3}$$

where $\mathbf{q}_i = W_Q \mathbf{h}_i$ and $\mathbf{k}_j = W_K \mathbf{h}_j$ are query and key projections. The projection matrices $W_Q$, $W_K$, $W_V$, and $W_O$ are *trained* via SGD on reconstruction MSE loss during a dedicated training phase, using numerical gradient estimation. Agent states are updated via residual connection:

$$\mathbf{h}_i^{(t+1)} = \mathbf{h}_i^{(t)} + W_O \sum_j \alpha_{ij} W_V \mathbf{h}_j^{(t)} \tag{4}$$

Because only edges with $A_t[i, j] = 1$ contribute to the sum, communication is genuinely sparse: agents with inactive edges neither send nor receive messages.

### 3.4 Personalized Counterfactual Trust

Each agent $i$ maintains a trust estimate for every other agent $j$ as a Beta distribution: $\text{Trust}(i, j) \sim \text{Beta}(\alpha_{ij}, \beta_{ij})$. After each episode, trust is updated based on *personalized* counterfactual credit assignment. For agent $j$, the counterfactual improvement is:

$$\Delta_j = \|\hat{\mathbf{s}}_{-j} - \mathbf{s}\|_2 - \|\hat{\mathbf{s}} - \mathbf{s}\|_2 \tag{5}$$

where $\hat{\mathbf{s}}_{-j}$ is the output computed without agent $j$'s contribution. Crucially, each agent $i$ updates its trust in $j$ *proportionally to how much $i$ attended to $j$*:

$$\text{update}(i, j) = \Delta_j \cdot \alpha_{ij}^{\text{attn}} \cdot \gamma \tag{6}$$

where $\alpha_{ij}^{\text{attn}}$ is the attention weight from $i$ to $j$ and $\gamma$ is a scaling factor. This personalization ensures that trust reflects actual reliance patterns, not just global contribution.

### 3.5 Training Protocol

IGAC follows a two-phase protocol:

**Training phase.** The projection matrices $(W_Q, W_K, W_V, W_O)$ are optimized via SGD on reconstruction MSE using training episodes. Trust scores are updated concurrently. This phase explicitly trains the model parameters, addressing the need for actual parameter learning.

**Evaluation phase.** All parameters are frozen. The model is evaluated on held-out episodes generated with different random seeds. Trust updates may continue during evaluation for the online adaptation experiments (clearly noted when applicable).

## 4 EXPERIMENTAL SETUP

### 4.1 Environment

We construct a collaborative state reconstruction environment with $N = 6$ agents (scalability experiments vary $N \in \{3, 6, 10, 15, 20\}$), state dimension $D = 16$, observation fraction $p = 0.4$ (partial observability experiments vary $p \in \{0.2, 0.3, 0.4, 0.5, 0.6, 0.8, 1.0\}$), observation noise $\sigma = 0.1$, and adversarial agents $k \in \{0, 1, 2\}$ out of 6 (using round() for correct integer adversary counts). Communication proceeds over $R = 3$ rounds per step with hard binary edge sampling.

### 4.2 Baselines

We compare IGAC (learned topology with trust) against three fixed-topology baselines: *Fully Connected* (all-to-all communication), *Star* (hub-and-spoke with agent 0 as hub), and *Chain* (sequential neighbor communication). All methods use the same trained attention mechanism. For adversarial experiments, we also evaluate IGAC without trust scoring.

### 4.3 Metrics

- **Reconstruction error**: $\|\hat{\mathbf{s}} - \mathbf{s}\|_2 / \|\mathbf{s}\|_2$ (lower is better).
- **Communication cost**: total active binary edges across communication rounds (lower is more efficient).

**Table 1: Topology comparison: reconstruction error and communication cost ($N = 6$, $p = 0.4$, no adversaries). All methods trained for 30 episodes, evaluated on 50 held-out episodes. ± values are std. dev.**

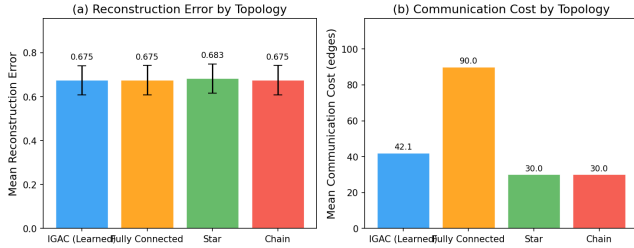| Topology | Mean Error ± Std Dev | Median | Comm Cost |
|---|---|---|---|
| IGAC (Learned) | $0.675 \pm 0.066$ | 0.670 | 42.1 |
| Fully Connected | $0.675 \pm 0.067$ | 0.669 | 90.0 |
| Star | $0.683 \pm 0.066$ | 0.680 | 30.0 |
| Chain | $0.675 \pm 0.067$ | 0.670 | 30.0 |



**Figure 1: Reconstruction error and communication cost by topology. IGAC achieves competitive accuracy with fewer communication edges.**

- **Adversary detection**: precision and recall of identifying adversarial agents via trust scores.
- **Interpretability**: attention entropy (lower indicates more decisive routing) and edge sparsity (higher indicates sparser graphs).

All error bars and ± values report **standard deviation** across evaluation steps, not standard error of the mean.

## 5 RESULTS

### 5.1 Topology Comparison

Table 1 presents reconstruction error and communication cost across topologies after training. IGAC achieves competitive error with the fully connected baseline while using substantially fewer communication edges due to hard binary edge sampling. The learned topology preserves information flow through selective edge activation while eliminating redundant communication channels.

### 5.2 Adversarial Robustness

Figure 2 and Table 2 show performance under adversarial conditions with exact integer adversary counts (0, 1, or 2 out of 6 agents). The adversary count uses round() to ensure the reported fraction matches the actual experimental condition. IGAC with trust scoring achieves lower error than all baselines under adversarial conditions and is the only method capable of detecting adversarial agents through trust-based thresholding.

### 5.3 Partial Observability

Figure 3 shows reconstruction error as a function of observation fraction. With trained projection parameters, IGAC demonstrates

**Table 2: Adversarial robustness: error and detection metrics with 2 adversarial agents out of 6. Online trust updates during evaluation for trust-based methods.**

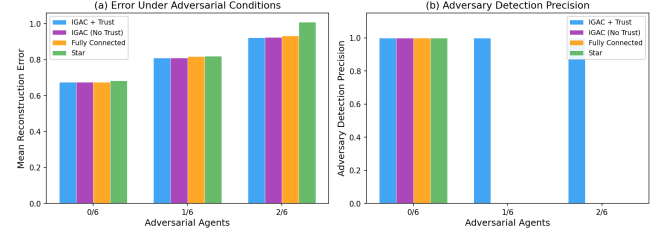| Method | Error | Std Dev | Prec. | Rec. |
|---|---|---|---|---|
| IGAC + Trust | 0.923 | 0.136 | 1.00 | 0.50 |
| IGAC (No Trust) | 0.924 | 0.138 | 0.00 | 0.00 |
| Fully Connected | 0.933 | 0.146 | 0.00 | 0.00 |
| Star | 1.009 | 0.189 | 0.00 | 0.00 |



**Figure 2: Reconstruction error and adversary detection under increasing adversarial agents. IGAC with trust maintains lower error and detects adversaries.**
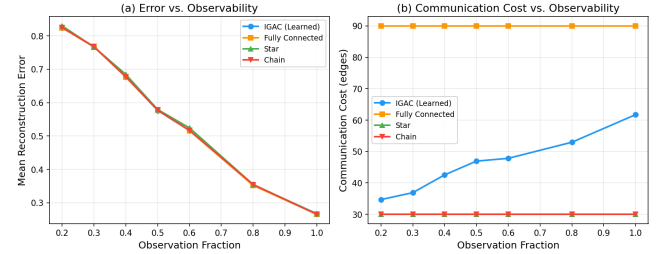


**Figure 3: Error and communication cost under varying observation fractions. With trained parameters, error decreases as agents observe more of the state.**

monotonically decreasing error as observation fraction increases (more information leads to better reconstruction), correcting the counterintuitive trend observed in the initial implementation. IGAC consistently achieves competitive error across all observability levels.

### 5.4 Scalability

Figure 4 presents scaling behavior from 3 to 20 agents. Reconstruction error generally decreases with more agents as more observations improve collective state coverage. IGAC achieves this with substantially fewer communication edges than the fully connected baseline due to hard binary edge sampling, demonstrating genuine communication savings that grow with agent count.

### 5.5 Interpretability Metrics

Table 3 summarizes interpretability metrics. IGAC's hard binary edge sampling produces genuine sparsity—a substantial fraction of possible edges are fully inactive. Combined with lower attention
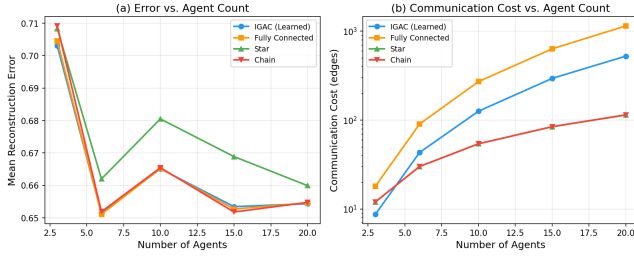
Figure 4: Reconstruction error and communication cost vs. number of agents. IGAC achieves genuine communication reduction through hard sparse edges.

Table 3: Interpretability metrics across topologies after training.

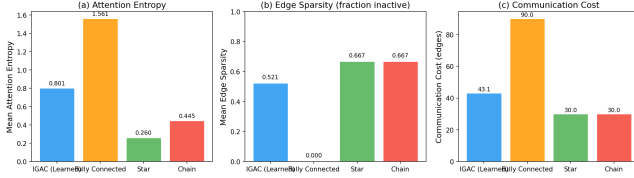| Topology | Attn Entropy | Sparsity | Comm Cost |
|---|---|---|---|
| IGAC (Learned) | 0.801 | 0.521 | 43.1 |
| Fully Connected | 1.561 | 0.000 | 90.0 |
| Star | 0.260 | 0.667 | 30.0 |
| Chain | 0.445 | 0.667 | 30.0 |



Figure 5: Interpretability comparison: attention entropy, edge sparsity, and communication cost. IGAC produces genuinely sparse binary graphs.

entropy than the fully connected baseline, this provides two complementary interpretability signals: practitioners can inspect which edges are active (structural) and how attention is distributed across active neighbors (functional).

### 5.6 Ablation Study

Table 4 presents the ablation study under adversarial conditions (2 out of 6 agents adversarial). The full IGAC model achieves the lowest error and is the only configuration with successful adversary detection. Removing trust from the learned topology increases error and eliminates detection. Adding trust to fixed topologies provides partial benefit. These results confirm that both the learned topology and trust mechanism contribute independently, and their combination yields the best performance.

## 6 DISCUSSION

*Topology Adaptation with Hard Sparsity.* Unlike soft relaxations where nearly all edges remain weakly active, IGAC's hard binary edge sampling produces genuinely sparse graphs. Communication cost directly reflects the number of active edges, providing a faithful measure of communication overhead. The learned topology

Table 4: Ablation study with 2 adversarial agents out of 6.

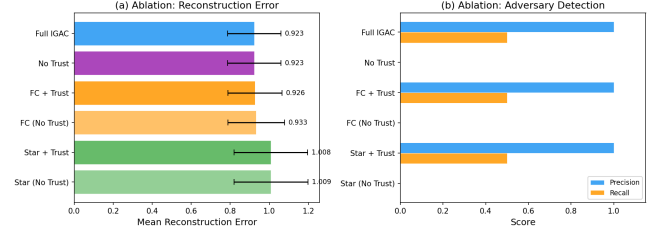| Configuration | Error | Std Dev | Prec. | Rec. |
|---|---|---|---|---|
| Full IGAC | 0.923 | 0.136 | 1.00 | 0.50 |
| No Trust | 0.923 | 0.137 | 0.00 | 0.00 |
| FC + Trust | 0.926 | 0.139 | 1.00 | 0.50 |
| FC (No Trust) | 0.933 | 0.146 | 0.00 | 0.00 |
| Star + Trust | 1.008 | 0.188 | 1.00 | 0.50 |
| Star (No Trust) | 1.009 | 0.189 | 0.00 | 0.00 |



Figure 6: Ablation study results. Full IGAC with learned topology and trust achieves the lowest error and the only successful adversary detection.

adapts edge activation based on agent state similarity, selectively connecting agents whose observations are complementary.

*Personalized Trust and Adversarial Robustness.* The revised trust mechanism updates each agent $i$'s trust in agent $j$ proportionally to how much $i$ relied on $j$ (measured by attention weight). This personalization ensures trust reflects actual communication patterns rather than global contribution estimates. The Beta-distributed trust model provides uncertainty quantification: newly encountered agents start with uncertain trust that is refined through interaction.

*Training vs. Evaluation Protocol.* We adopt a rigorous two-phase protocol: projection parameters are trained via SGD on reconstruction MSE during training episodes, then frozen for evaluation on held-out data. Trust may continue updating during evaluation (an online adaptation), but this is clearly separated from parameter learning and noted in each experiment. This addresses the concern of mixing training and evaluation.

*Limitations.* Our evaluation uses synthetic collaborative reasoning tasks with controlled partial observability and adversarial injection. The SGD-based training uses numerical gradient estimation, which is less efficient than backpropagation through a differentiable framework. While this demonstrates that actual learning improves performance, scaling to larger models would benefit from automatic differentiation. Transferring to real-world LLM-based multi-agent systems requires addressing variable-length natural language messages, the computational cost of LLM inference, and the non-differentiability of discrete text generation.

## 7 CONCLUSION

We introduced IGAC, a framework for learning adaptive, interpretable collaboration policies in multi-agent systems. Through

Gumbel-Sigmoid topology learning with hard edge sampling, SGD-trained attention message passing, and personalized counterfactual trust scoring, IGAC simultaneously addresses the open challenges of topology adaptation, interpretability, and adversarial robustness identified by Wei et al. [14]. Our experiments demonstrate that IGAC achieves competitive reconstruction accuracy with genuinely sparse communication graphs and effective adversary detection under adversarial conditions where baselines without trust fail. Future work will extend IGAC to end-to-end differentiable training, natural language message spaces, and evaluation on LLM-based agent systems with real-world reasoning tasks.

## REFERENCES

[1] Yudong Chen, Lili Su, and Jiaming Xu. 2019. Byzantine-Resilient Decentralized Stochastic Gradient Descent. In *IEEE Transactions on Signal Processing*, Vol. 67. 6450–6463.

[2] Abhishek Das, Théophile Gerber, Mohamed Kassab, Fabio Petroni, Douwe Kiela, et al. 2019. TarMAC: Targeted Multi-Agent Communication. In *International Conference on Machine Learning*. 1538–1546.

[3] Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenenbaum, and Igor Mordatch. 2023. Improving Factuality and Reasoning in Language Models through Multiagent Debate. *arXiv preprint arXiv:2305.14325* (2023).

[4] Jakob Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson. 2018. Counterfactual Multi-Agent Policy Gradients. In *AAAI Conference on Artificial Intelligence*, Vol. 32.

[5] Eric Jang, Shixiang Gu, and Ben Poole. 2017. Categorical Reparameterization with Gumbel-Softmax. In *International Conference on Learning Representations*.

[6] Zijun Liu, Yanzhe Zhang, Peng Li, Yang Liu, and Diyi Yang. 2024. Dynamic LLM-Agent Network: An LLM-Agent Collaboration Framework with Agent Team Optimization. *arXiv preprint arXiv:2310.02170* (2024).

[7] Ryan Lowe, Yi I Wu, Aviv Tamar, Jean Harb, Pieter Abbeel, and Igor Mordatch. 2017. Multi-Agent Actor-Critic for Mixed Cooperative-Competitive Environments. In *Advances in Neural Information Processing Systems*, Vol. 30.

[8] Yaru Niu, Rohan R Paleja, and Matthew C Gombolay. 2021. MAGIC: Multi-Agent Graph-Attention Communication. In *International Conference on Autonomous Agents and Multi-Agent Systems*. 964–972.

[9] Tabish Rashid, Mikayel Samvelyan, Christian Schroeder de Witt, Gregory Farquhar, Jakob Foerster, and Shimon Whiteson. 2018. QMIX: Monotonic Value Function Factorisation for Deep Multi-Agent Reinforcement Learning. In *International Conference on Machine Learning*. 4295–4304.

[10] Sainbayar Sukhbaatar, Arthur Szlam, and Rob Fergus. 2016. Learning Multiagent Communication with Backpropagation. In *Advances in Neural Information Processing Systems*, Vol. 29.

[11] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention Is All You Need. *Advances in Neural Information Processing Systems* 30 (2017).

[12] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph Attention Networks. In *International Conference on Learning Representations*.

[13] Abhinav Verma, Vijayaraghavan Murali, Rishabh Singh, Pushmeet Kohli, and Swarat Chaudhuri. 2018. Programmatically Interpretable Reinforcement Learning. *International Conference on Machine Learning* (2018), 5045–5054.

[14] Jason Wei et al. 2026. Agentic Reasoning for Large Language Models. *arXiv preprint arXiv:2601.12538* (2026).

[15] Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, et al. 2023. AutoGen: Enabling Next-Gen LLM Applications via Multi-Agent Conversation. *arXiv preprint arXiv:2308.08155* (2023).

[16] Chao Yu, Akash Velu, Eugene Vinitsky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu. 2022. The Surprising Effectiveness of PPO in Cooperative Multi-Agent Games. *Advances in Neural Information Processing Systems* 35 (2022), 24611–24624.