# Opponent Modeling with In-context Search

**Yuheng Jing**[1,2]**, Bingyun Liu**[1,2]**, Kai Li**[1,2,†]**, Yifan Zang**[1,2]**,**
**Haobo Fu**[6]**, Qiang Fu**[6]**, Junliang Xing**[5]**, Jian Cheng**[1,3,4,†]

1 Institute of Automation, Chinese Academy of Sciences
2 School of Artificial Intelligence, University of Chinese Academy of Sciences
3 School of Future Technology, University of Chinese Academy of Sciences
4 AiRiA     5 Tsinghua University     6 Tencent AI Lab

# Background

**Opponent Modeling (OM)** enhances decision-making in multi-agent environments by modeling the behaviors, goals, and etc. of ***other agents*** *(including adversaries and teammates, collectively termed as **opponents**).*

Existing OM approaches fall into two categories:

**1. Pretraining-Focused Approaches (PFAs)**: Train models to generalize across diverse opponent policies.

**2. Testing-Focused Approaches (TFAs)**: Update pretrained models at testing time to adapt to unknown opponents through finetuning.

# Background

**Opponent Modeling (OM)** enhances decision-making in multi-agent environments by modeling the behaviors, goals, and etc. of **other agents** *(including adversaries and teammates, collectively termed as **opponents**).*

Existing OM approaches fall into two categories:

**1. Pretraining-Focused Approaches (PFAs)**: Train models to generalize across diverse opponent policies.

**2. Testing-Focused Approaches (TFAs)**: Update pretrained models at testing time to adapt to unknown opponents through finetuning.

**Problems with Current Approaches**:

- **PFAs**: **Limited Generalization** to unknown opponents due to a lack of theoretical analyses and minimal testing-stage adjustments.

- **TFAs**: **Performance Instability** from reliance on small datasets for testing-time gradient updates and sensitivity to hyperparameters.

A unified approach is needed to tackle these challenges.

# Contributions

We introduce **Opponent Modeling with In-context Search (OMIS)**, a novel OM approach which leverages **In-Context Learning (ICL)** and **Decision-Time Search (DTS)** to tackle the existing challenges.
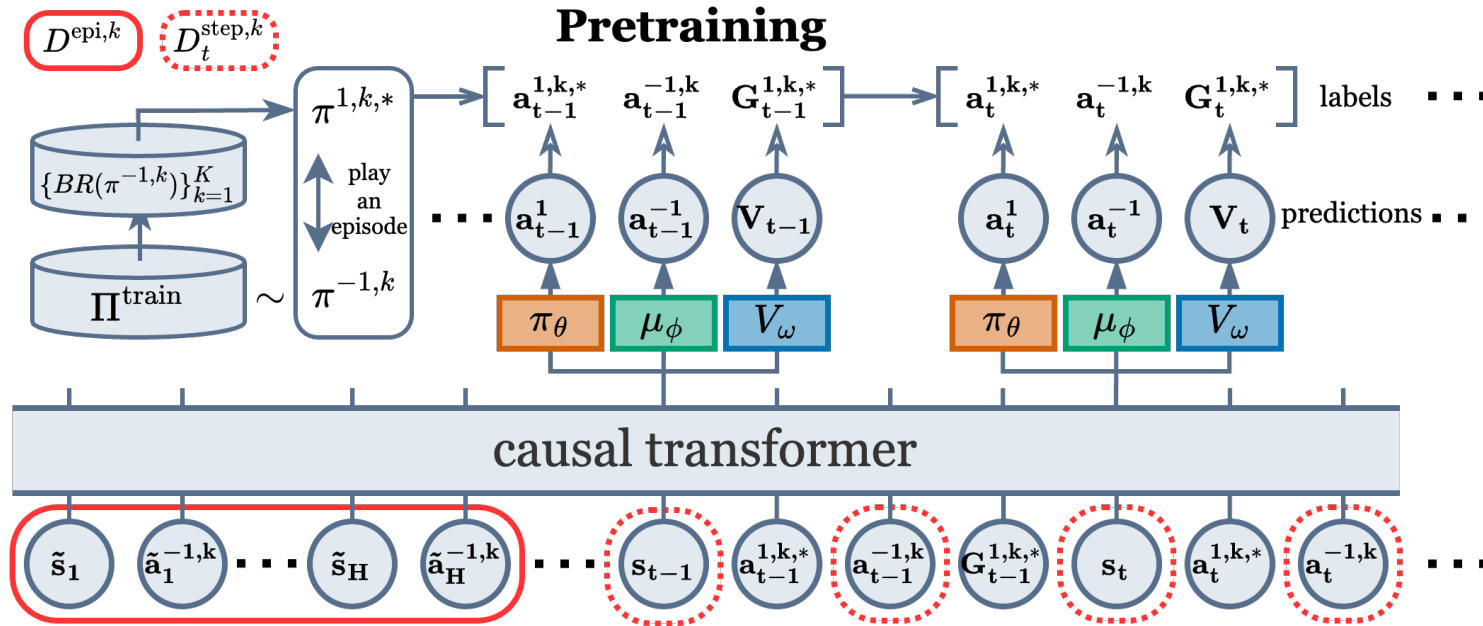
**The Core Motivation of OMIS**:

- We encode the idea of '*think before you act*' through DTS.
- We make DTS *possible* and essentially *adaptive* through ICL.

# Contributions

We introduce **Opponent Modeling with In-context Search (OMIS)**, a novel OM approach which leverages **In-Context Learning (ICL)** and **Decision-Time Search (DTS)** to tackle the existing challenges.

**The Core Motivation of OMIS**:

- We encode the idea of '*think before you act*' through DTS.
- We make DTS *possible* and essentially *adaptive* through ICL.
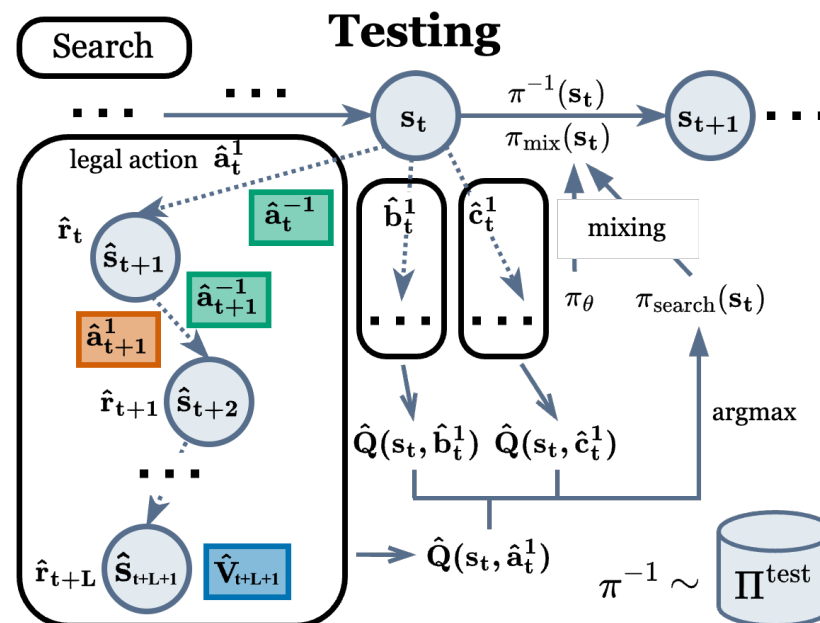
**The Core Innovations of OMIS**:

- A *theoretically grounded* **DTS mechanism** to refine the original policy in real-time rather than gradient updates ***to avoid performance instability*** *(main problems of TFAs)*.
- A *theoretically grounded* **ICL-pretrained Transformer** with three in-context components ***to enable DTS and enhance generalization*** *(main problems of PFAs)*:

  (1) **Actor**: Learn to respond appropriately to the current opponent.

  (2) **Opponent Imitator**: Predict the actions of the current opponent.

  (3) **Critic**: Evaluate the values of the states during DTS.
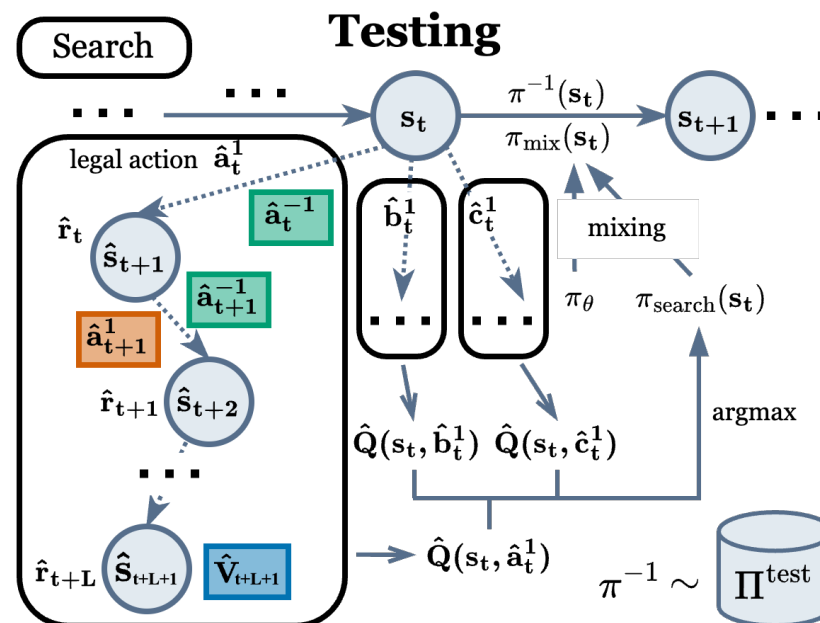
# Methodology



**The Pretraining Procedure and Architecture of OMIS**. The pretraining steps are as follows: (1) Train *Best Responses* (BRs) against all policies in the *training set of opponent policies* $\Pi^{\text{train}}$. (2) Continuously sample opponent policy from $\Pi^{\text{train}}$ and collect training data by playing against it using its BR. (3) Train a Transformer model using ICL-based supervised learning, where the model consists of *three in-context components*.

# Methodology



**The Testing Procedure of OMIS**. During testing, OMIS refines the *actor* $\pi_\theta$ through DTS at each timestep. The DTS steps are as follows: (1) Do multiple $L$-step rollouts for each legal action, where the *actor* $\pi_\theta$ and the *opponent imitator* $\mu_\phi$ are used to simulate actions for the self-agent and opponent, respectively. The *critic* $V_\omega$ is used to estimate the value of final search states.

# Methodology



The DTS steps are as follows: (1) Do multiple $L$-step rollouts … (2) Estimate the *DTS value* $\hat{Q}$ for all legal actions, and the *search policy* $\pi_{\text{search}}$ selects the legal action with the maximum $\hat{Q}$. (3) Use *mixing technique* to trade-off between the *search policy* $\pi_{\text{search}}$ and the *actor* $\pi_{\theta}$ to choose the real action to be executed.
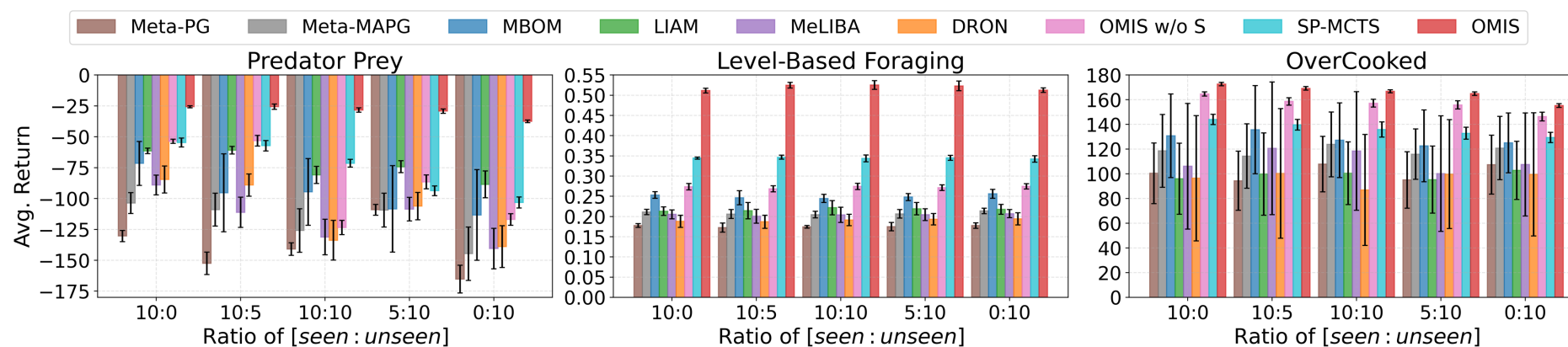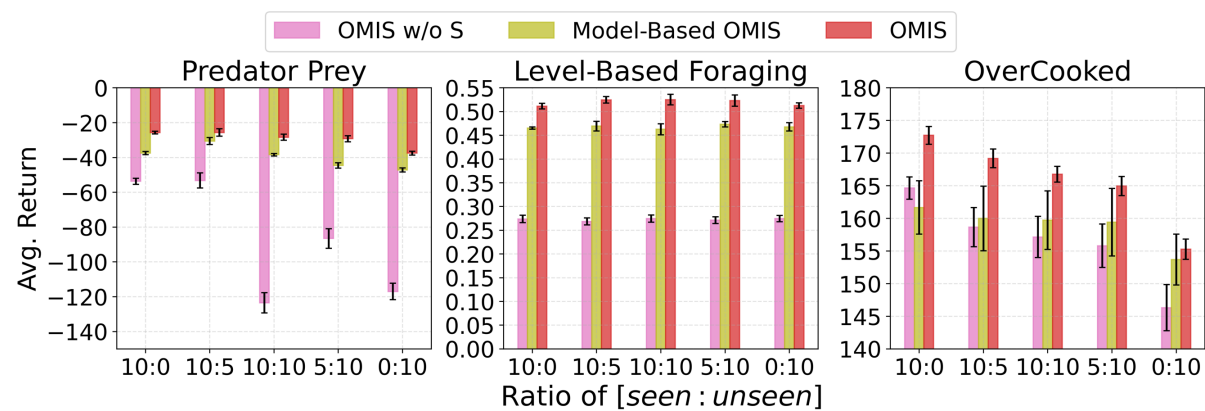
# Theoretical Results

- **Generalization Guarantee of OMIS w/o DTS** (Supported by **Lemma 4.1** and **Theorem 4.2**): OMIS w/o DTS can accurately recognize *seen opponent policies (policies in the training set of opponent policies* $\Pi^{\text{train}}$*)* and recognize *unseen opponent policies (policies not in* $\Pi^{\text{train}}$*)* as the most familiar seen ones with the measurement of Kullback-Leibler Divergence.

# Theoretical Results

- **Generalization Guarantee of OMIS w/o DTS** (Supported by **Lemma 4.1** and **Theorem 4.2**): OMIS w/o DTS can accurately recognize *seen opponent policies (policies in the training set of opponent policies* $\Pi^{\text{train}}$*)* and recognize *unseen opponent policies (policies not in* $\Pi^{\text{train}}$*)* as the most familiar seen ones with the measurement of Kullback-Leibler Divergence.

- **Policy Improvement of OMIS w/ DTS** (Supported by **Theorem 4.3**): OMIS w/ DTS avoids any gradient updates and theoretically provides improvement guarantees.
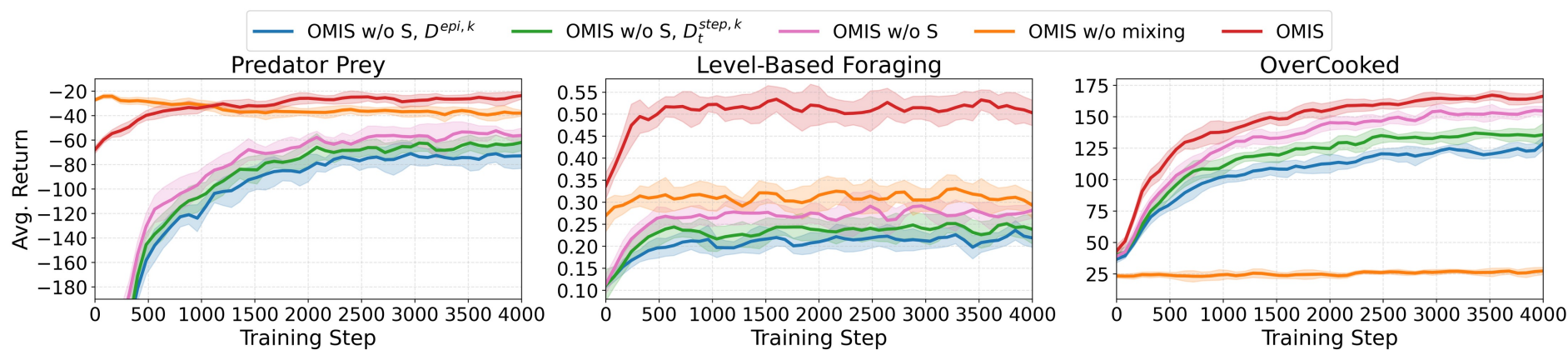
# Experimental Results



**Main Results**: OMIS effectively outperforms *PFAs*, *TFAs*, and *DTS-based* OM baselines when testing against *unknown non-stationary opponent agents* in competitive, cooperative and mixed environments.
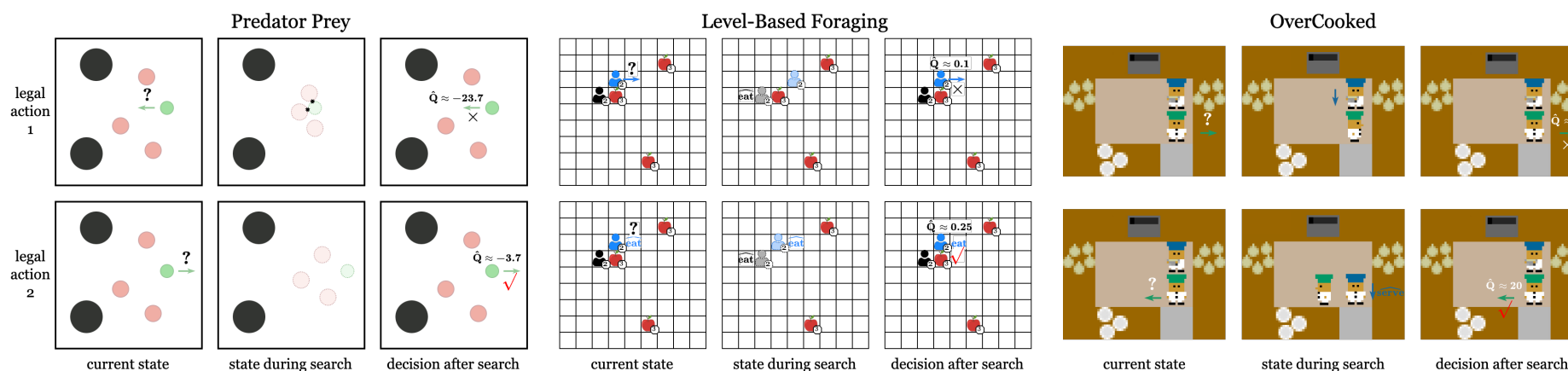
# Experimental Results



**Ablations**: OMIS can work effectively when the *transition dynamics* are *unknown* and *learned* instead.

# Experimental Results



**Ablations**: Each designed component of OMIS plays a positive role in boosting the performance.

# Experimental Results



**Visualizations**: When playing against opponents with a previously unseen policy, the DTS of OMIS promptly evaluates each legal action, predicts the opponent's actions and the state values during the DTS process, and ultimately selects the most advantageous action for the self-agent.
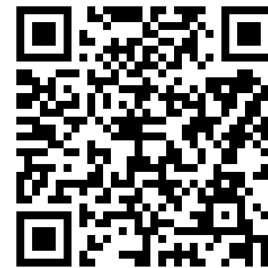
# Thank You!

Project
Website

Paper

Code