# VAST: Value Function Factorization with Variable Agent Sub-Teams

## NeurIPS 2021

**Thomy Phan[1]**, Fabian Ritz[1], Lenz Belzner[2],

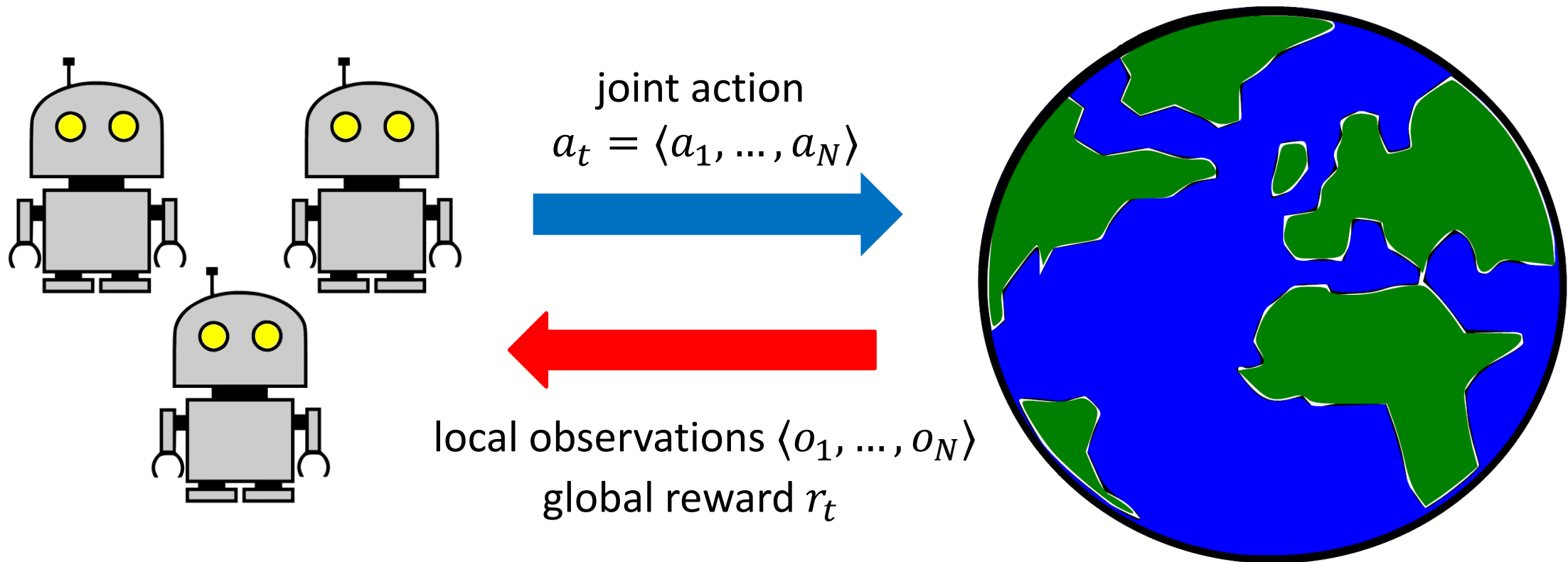Philipp Altmann[1], Thomas Gabor[1], Claudia Linnhoff-Popien[1]

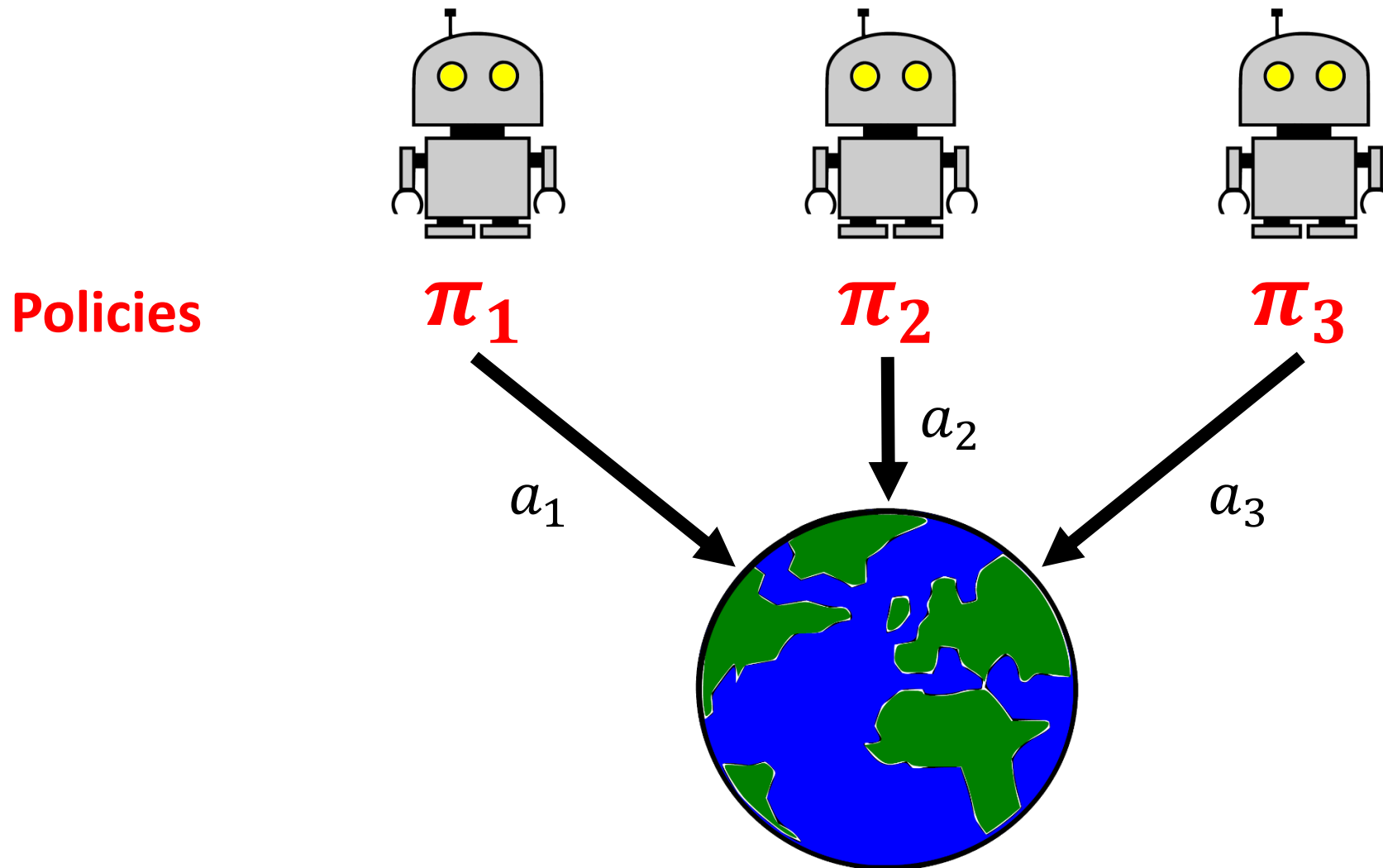[1]LMU Munich, [2]Technische Hochschule Ingolstadt

# Preliminaries

Thomy Phan, thomy.phan@ifi.lmu.de

# Cooperative Multi-Agent Systems (MAS)



joint action
$$a_t = \langle a_1, \ldots, a_N \rangle$$

local observations $\langle o_1, \ldots, o_N \rangle$
global reward $r_t$
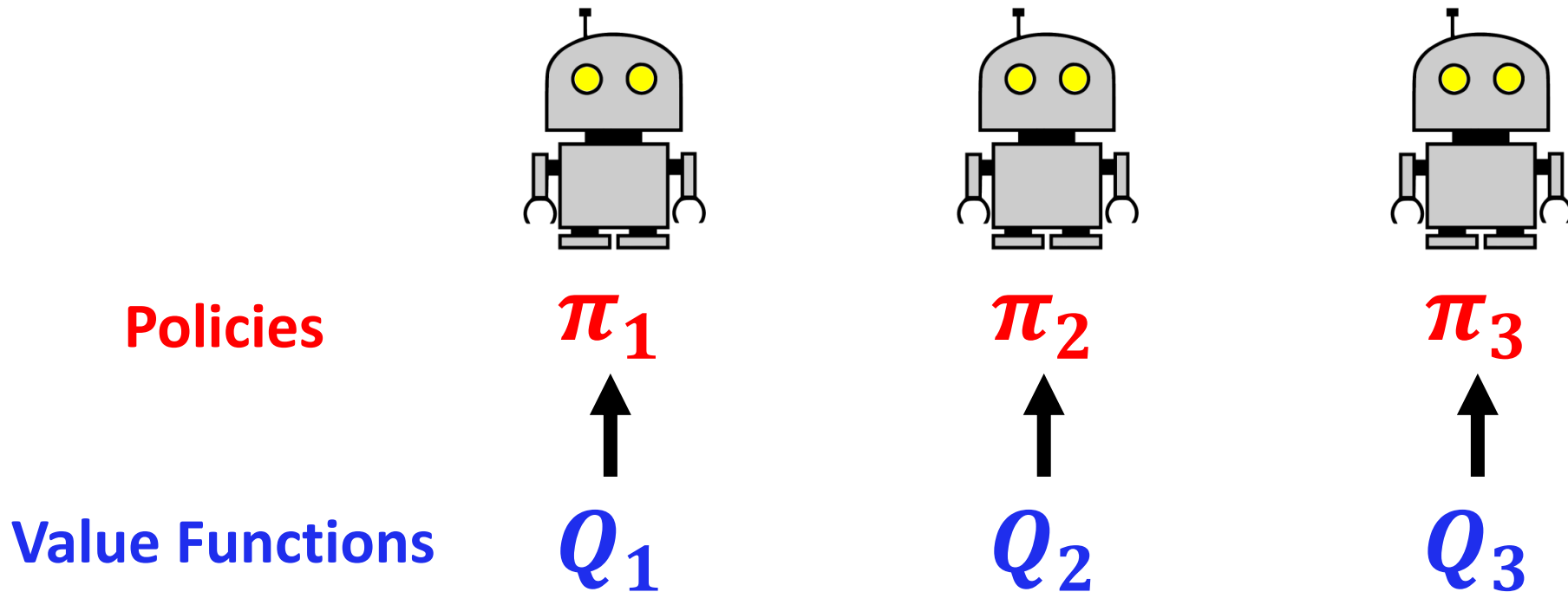
**Goal:** Maximize expectation of the return $\sum_{k=1}^{\infty} \gamma^k r_{t+k}$

# Multi-Agent Reinforcement Learning

**Policies**

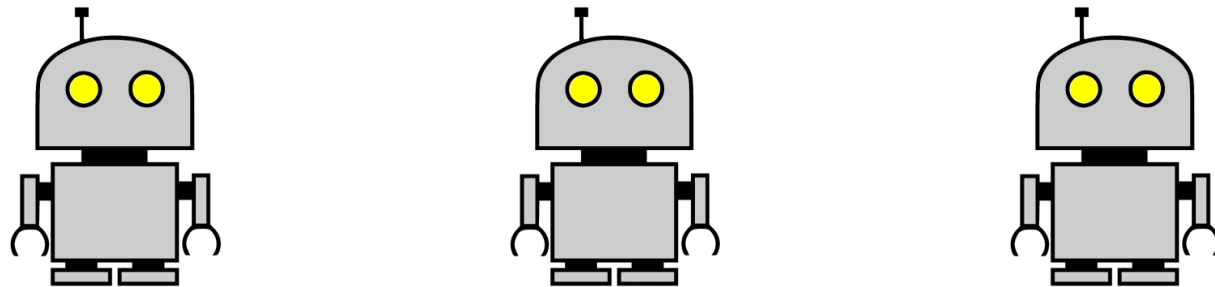$\boldsymbol{\pi_1}$ $\qquad\qquad$ $\boldsymbol{\pi_2}$ $\qquad\qquad$ $\boldsymbol{\pi_3}$

$a_1$

$a_2$

$a_3$

# **Value-based Multi-Agent Reinforcement Learning**

**Policies**     $\boldsymbol{\pi_1}$          $\boldsymbol{\pi_2}$          $\boldsymbol{\pi_3}$

**Value Functions**   $\boldsymbol{Q_1}$          $\boldsymbol{Q_2}$          $\boldsymbol{Q_3}$

$$Q_i(\tau_i, a_i) = \mathbb{E}[\sum_{k=1}^{\infty} \gamma^k r_{t+k}]$$

# Value Function Factorization



**Local Value Functions**

$$Q_1 \qquad Q_2 \qquad Q_3$$

$$\Psi$$

**Factorization Operator**

VDN
QMIX
QTRAN
W-QMIX
Qatten
QPLEX
...

**Centralized Value Function**

$$Q_{tot} = \mathbb{E}\left[\sum_{k=1}^{\infty} \gamma^k r_{t+k}\right]$$

# Individual-Global-Max (IGM) Consistency

$$argmax_{\boldsymbol{a}_t}\boldsymbol{Q}_{tot}(\boldsymbol{\tau}_t, \boldsymbol{a}_t) = \begin{pmatrix} argmax_{a_{t,1}}Q_1(\tau_{t,1}, a_{t,1}) \\ ... \\ argmax_{a_{t,N}}Q_N(\tau_{t,N}, a_{t,N}) \end{pmatrix}$$

Factorization operators must ensure IGM consistency
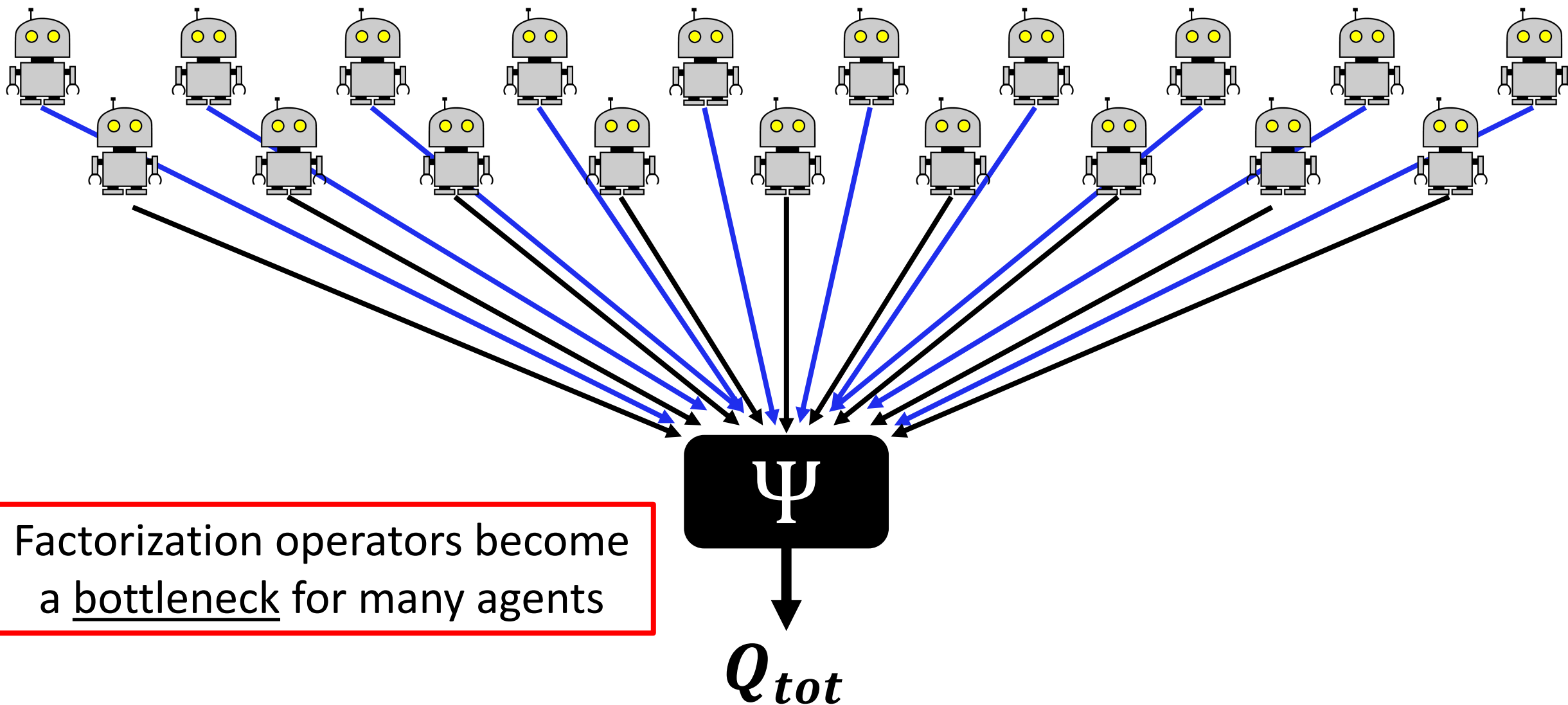
$$\Psi$$

$$Q_{tot}$$

Factorization operators become a <u>bottleneck</u> for many agents

$$\Psi$$

$$Q_{tot}$$

# Value Function Factorization with Variable Agent Sub-Teams (VAST)

$$Q_{tot}$$

# Idea of VAST: Sub-Team Assignment



$$\Psi$$

$$Q_{tot}$$

# Idea of VAST: Factorization on Sub-Team Values



$Q_{t,1}^G$   $Q_{t,2}^G$   $Q_{t,3}^G$

$\Psi$

VAST preserves IGM consistency
(given <u>any</u> sub-team assignment)

$Q_{tot}$

# Meta-Gradient Learning for Sub-Team Assignment

- **Idea:** Optimize sub-team assignments using a <u>high-level objective</u> $J$
    - Decide at each state $s_t$ which sub-team $k$ agent $i$ should be assigned to
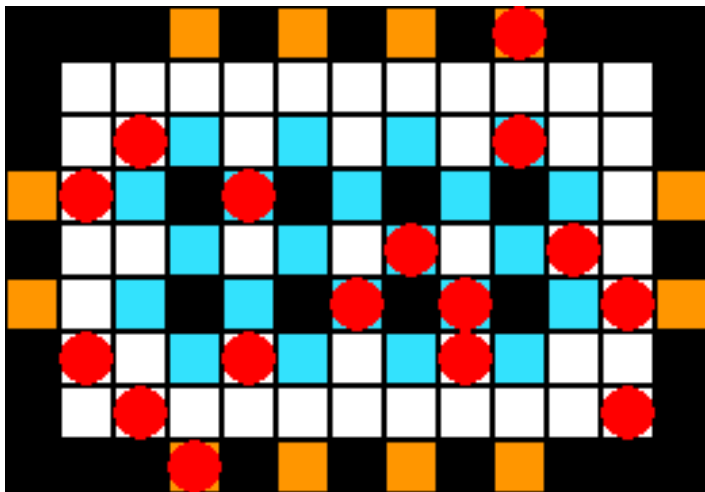    - Learn meta-policy $\mathcal{X}(k|s_t, i, \tau_{t,i})$ via gradient ascent w.r.t. objective $J$

$$\hat{A}(k, i, s_t, \tau_{t,i}) \nabla log \mathcal{X}(k|s_t, i, \tau_{t,i})$$

    - Advantage function $\hat{A}$ can be defined using domain knowledge or reward-based metrics (e.g., return, TD-error, …)
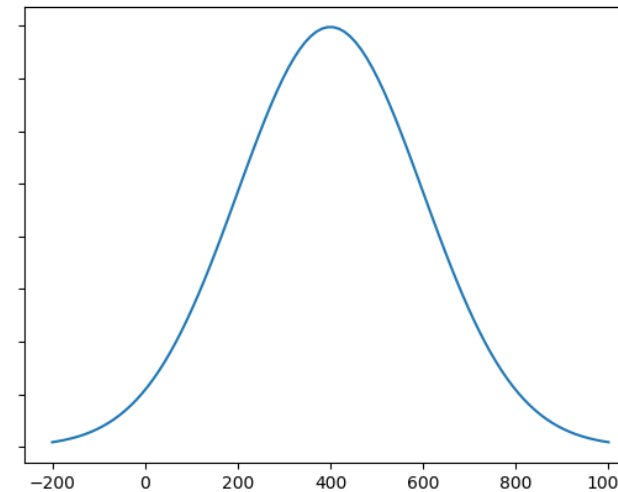
# Results

# Evaluation Domains



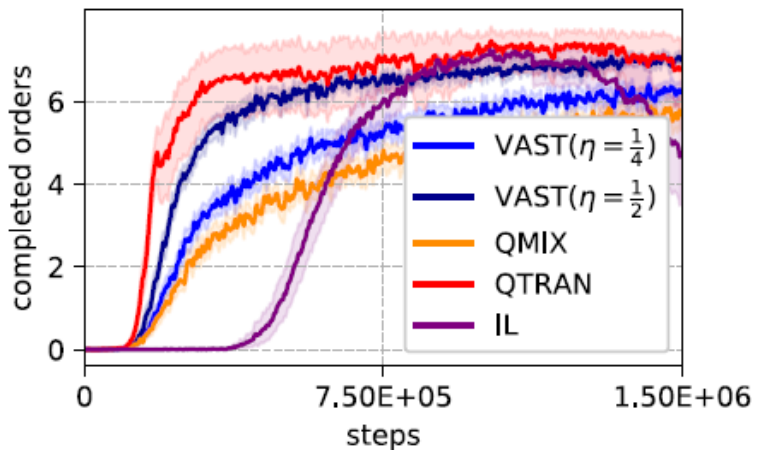**Warehouse**
(4 – 16 agents)

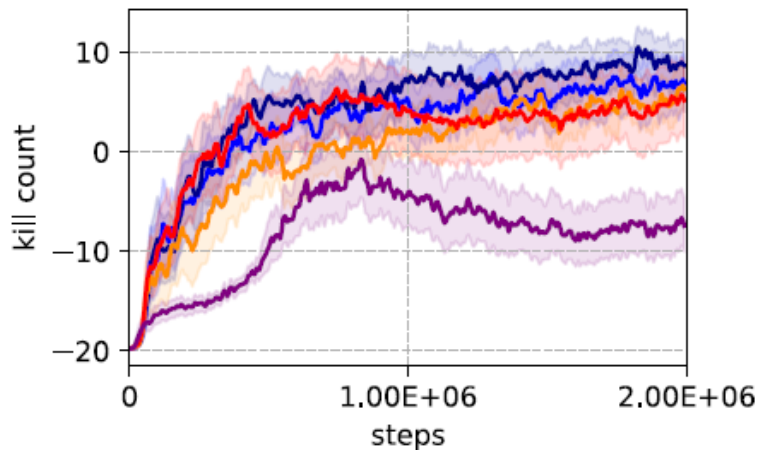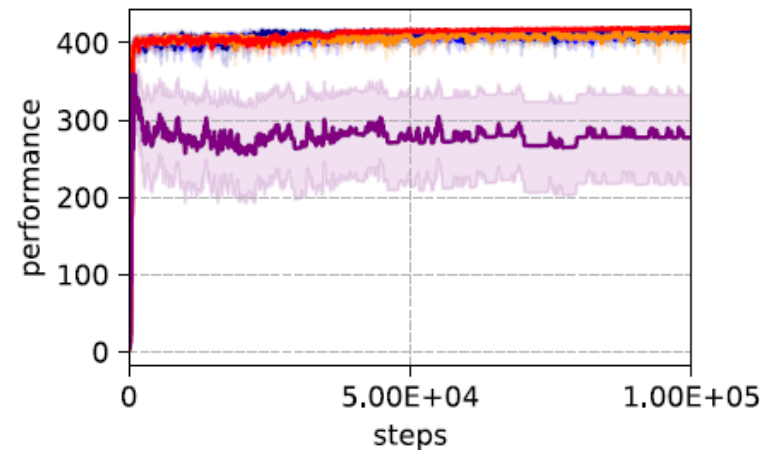**Battle**
(20 – 80 agents)

**Gaussian Squeeze**
(200 – 800 agents)
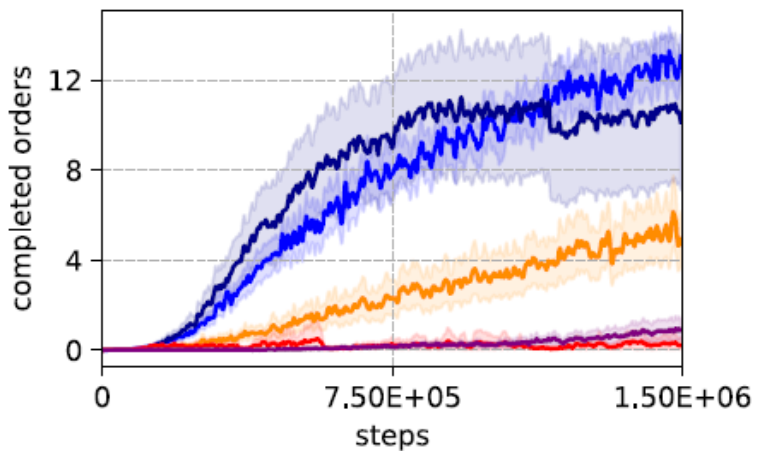
# State-of-the-Art Comparison
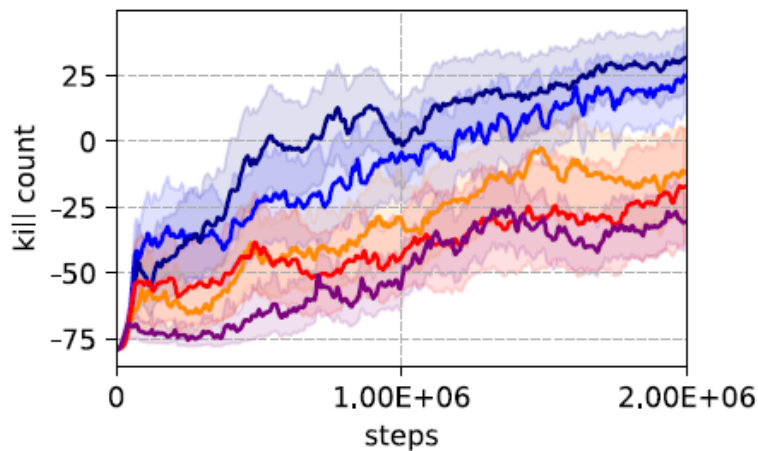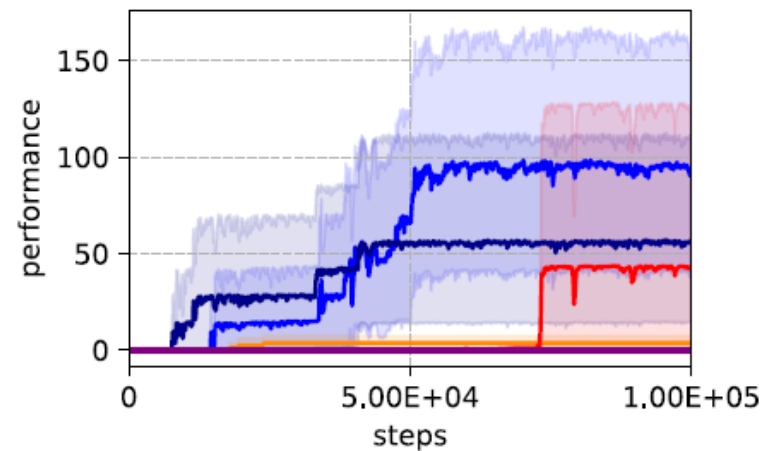


(a) *Warehouse[4]*

(b) *Battle[20]*

(c) *GaussianSqueeze[200]*

(d) *Warehouse[16]*

(e) *Battle[80]*

(f) *GaussianSqueeze[800]*

# Meta-Gradient Generated Sub-Teams in Battle[80]

Thomy Phan, thomy.phan@ifi.lmu.de

# Conclusion
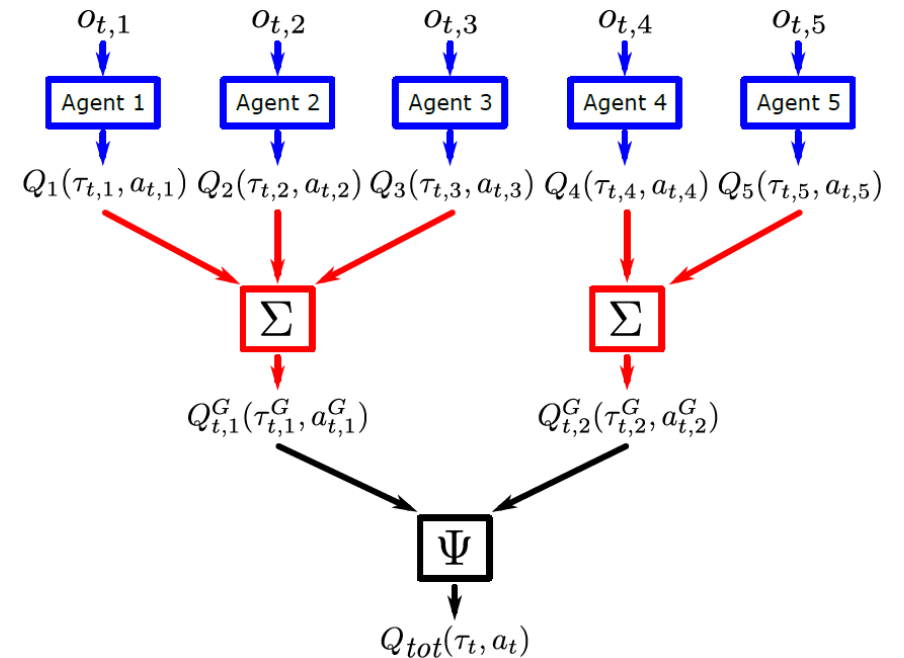
Thomy Phan, thomy.phan@ifi.lmu.de
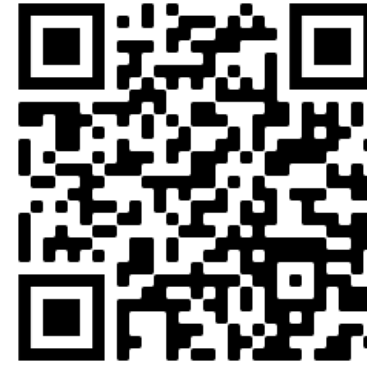
# Conclusion and Future Work

- VAST can improve scalability of value factorization w.r.t. many agents

- IGM consistency is preserved by VAST

- Meta-gradient based sub-teams can improve performance of VAST

**Future Work**

- Deeper hierarchies of sub-teams

- Non-linear factorization of sub-team values

**Code available at** ➡️

# VAST: Value Function Factorization with Variable Agent Sub-Teams
## NeurIPS 2021

**Thomy Phan**[1], Fabian Ritz[1], Lenz Belzner[2],

Philipp Altmann[1], Thomas Gabor[1], Claudia Linnhoff-Popien[1]

[1]LMU Munich, [2]Technische Hochschule Ingolstadt