# MetaCURE: Meta Reinforcement Learning with Empowerment-Driven Exploration

**Jin Zhang**\*, **Jianhao Wang**\*, **Hao Hu, Tong Chen,**

**Yingfeng Chen, Changjie Fan, Chongjie Zhang**

IIIS, Tsinghua University

jin-zhan20@mails.tsinghua.edu.cn

**Machine Intelligence Group**

清华大学交叉信息研究院
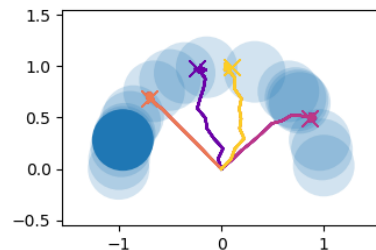Tsinghua University    Institute for Interdisciplinary Information Sciences

# Motivation



- Humans are experts in transferring knowledge

- Meta learning (Schmidhuber, J. 1987):
  - Meta-training: gain useful knowledge
  from previous tasks
  - Adaptation: adapt to new tasks with few-shot data

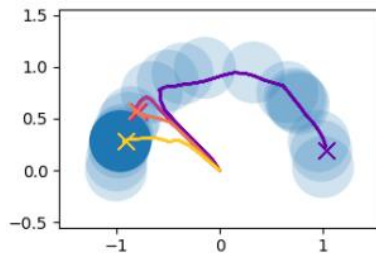- Meta-RL: how should we collect data in both phases?

# Exploration in Meta-RL

- How to explore in a new task?
  - Curiosity-driven methods?
    - Task-irrelevant distractors

  - Posterior sampling (Rakelly, Kate, et al. 2019)?
    - Exploitation policies may not explore effectively, as they are not optimized for exploration
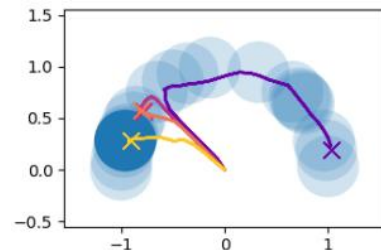


Posterior sampling



Efficient exploration

# Empowerment-Driven Exploration

- Meta-RL as task inference (Humplik, Jan, et al. 2019):
  - $\pi(a|s,z)$, $z$ is a latent variable containing task information.
- Exploration should support task inference.
  - Gain empowerment over the current task.
  - Objective: $\max I(C; \mathcal{K})$
    - $C$: exploration experience
    - $\mathcal{K}$: task identification
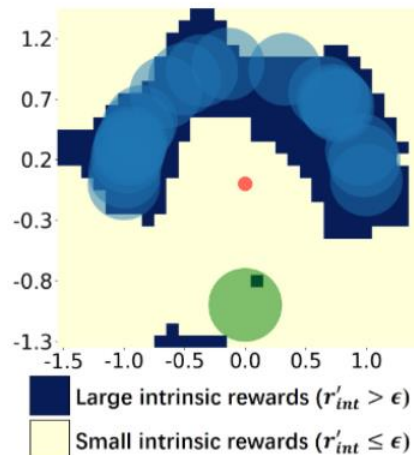
Efficient exploration

# Empowerment-Driven Exploration

- **Deriving intrinsic rewards**

  - $\max I(\mathcal{C}; \mathcal{K})$

  - $r'_{int}(c_{:t+1}, \kappa) = \underbrace{-\log p(r_t, s_{t+1} | c_{:t}, a_t)}_{L_{pred}(c_{:t+1})} + \underbrace{\log p(r_t, s_{t+1} | c_{:t}, a_t, \kappa)}_{-L_{pred}^{task}(\kappa, c_t)}$
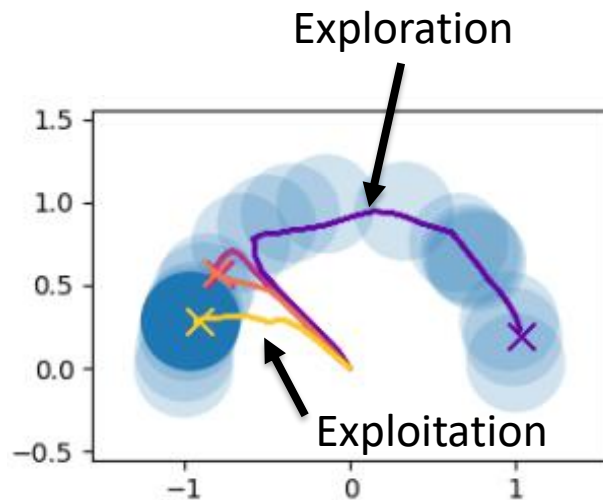
  - Subtraction of two model prediction errors!

    - $L_{pred}$: uncertainty given current experiences

    - $L_{pred}^{task}$: uncertainty given task identification

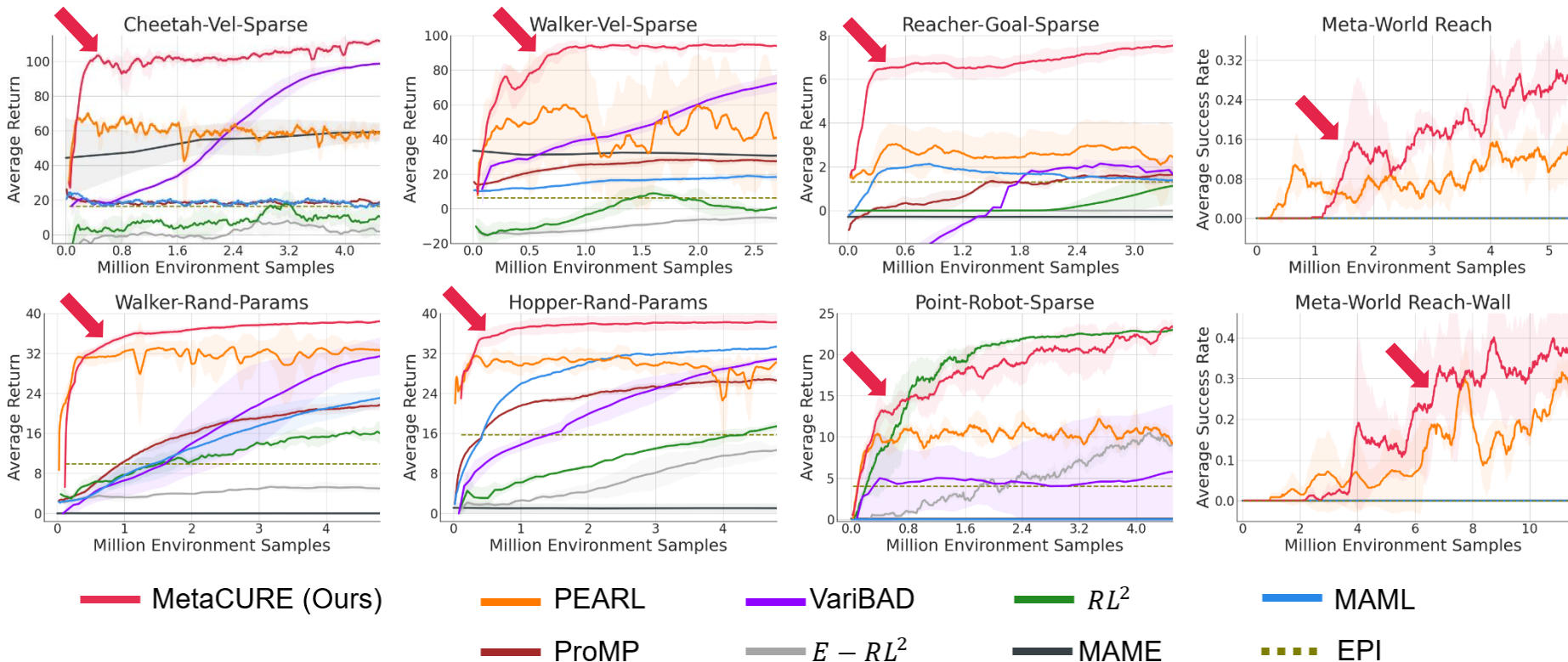    - Implication: only focus on uncertainty that helps task inference



Large intrinsic rewards ($r'_{int} > \epsilon$)

Small intrinsic rewards ($r'_{int} \leq \epsilon$)

# Separating Exploration and Exploitation

- **Exploration and exploitation naturally obtain different objectives!**
  - Exploration: obtain task information
  - Exploitation: maximize expected return



- **They should be two separate policies.**

# Results

# Take-Aways

- MetaCURE addresses the exploration problem in Meta-RL.

- Empowerment-driven exploration:
  - Maximize MI between exploration experiences and the task identification
- Separation of exploration and exploitation policies

- These ideas lead to superior performance on various hard sparse-reward Meta-RL benchmarks.

Thanks for your listening

Machine Intelligence Group

清华大学交叉信息研究院
Tsinghua University　Institute for Interdisciplinary Information Sciences