



TEMPLE: LEARNING TEMPLATE OF TRANSITIONS FOR SAMPLE EFFICIENT MULTI-TASK RL

Yanchao Sun*, Xiangyu Yin† and Furong Huang*
yys@umd.edu, yinxiangyu@bupt.edu.cn, furongh@umd.edu

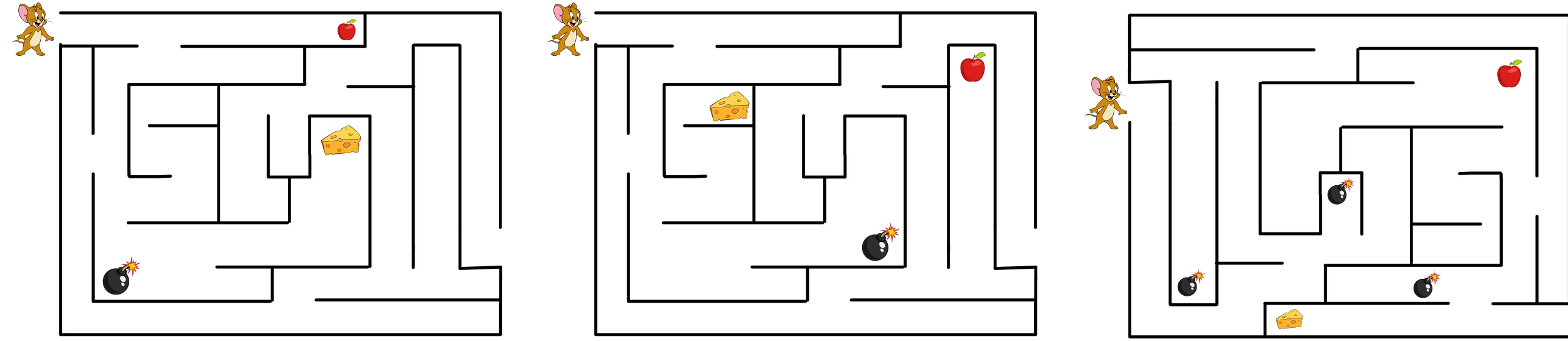
*University of Maryland, College Park, †Beijing University of Posts and Telecommunications, China



DEPARTMENT OF
COMPUTER SCIENCE

MULTI-TASK REINFORCEMENT LEARNING

Multi-task Reinforcement Learning (MTRL) studies the problem of efficiently learning a series of tasks by **knowledge transfer**.



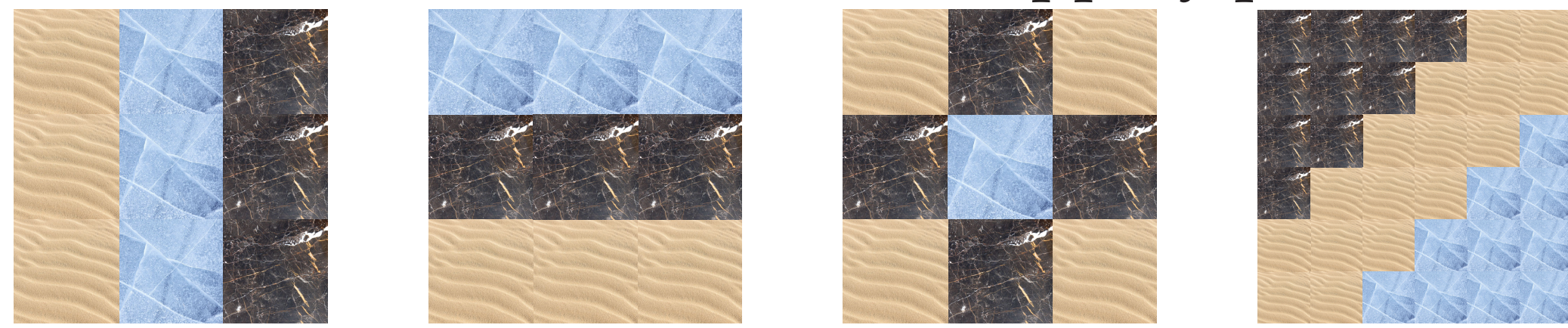
Challenges in MTRL:

- Guarantee the sample efficiency?
- Trade-off between correctness and efficiency?
- Tasks with different state/action spaces?

Our Contributions: our proposed algorithms achieve **SOTA sample complexity** and work for tasks with **varying state/action space**.

A MOTIVATING EXAMPLE

Various "landforms" → various slippery probabilities.



Sand: never slip; Marble: slip with prob 0.2; Ice: slip with prob 0.4.
G types of landforms → G^N different kinds of mazes (N: # of grids)

Our Goal: to extract and utilize **modular similarities** among and within tasks.

TRANSITION TEMPLATE

The Traditional Representation of Dynamics

The dynamics of a state-action pair can be represented as:

$$\theta(s, a) = [p(s_1|s, a), p(s_2|s, a), \dots, p(s_S|s, a), r(s, a)]$$

Transition Template (TT): A New Representation of Dynamics

By permuting the transition probability vector, we get:

$$\mathbf{g}(s, a) = [\text{desc}(p(\cdot|s, a)), r(s, a)]$$

1	2	3	4	5
6	7	8	9	10
11	12	13	14	15
16	17	18	19	20
21	22	23	24	25

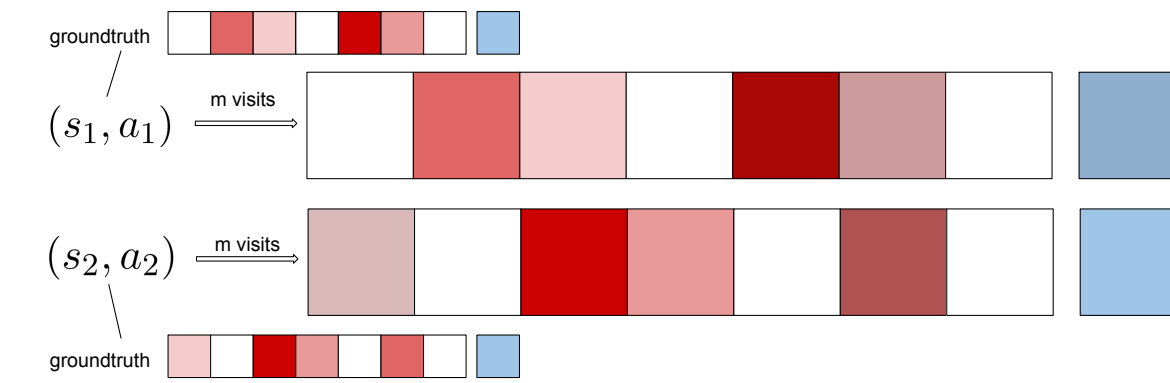
For all 100 s-a pairs in the 5×5 grid world, there are only 2 distinct TTs.

Transition Templates can capture modular similarities

LEARNING WITH TRANSITION TEMPLATES

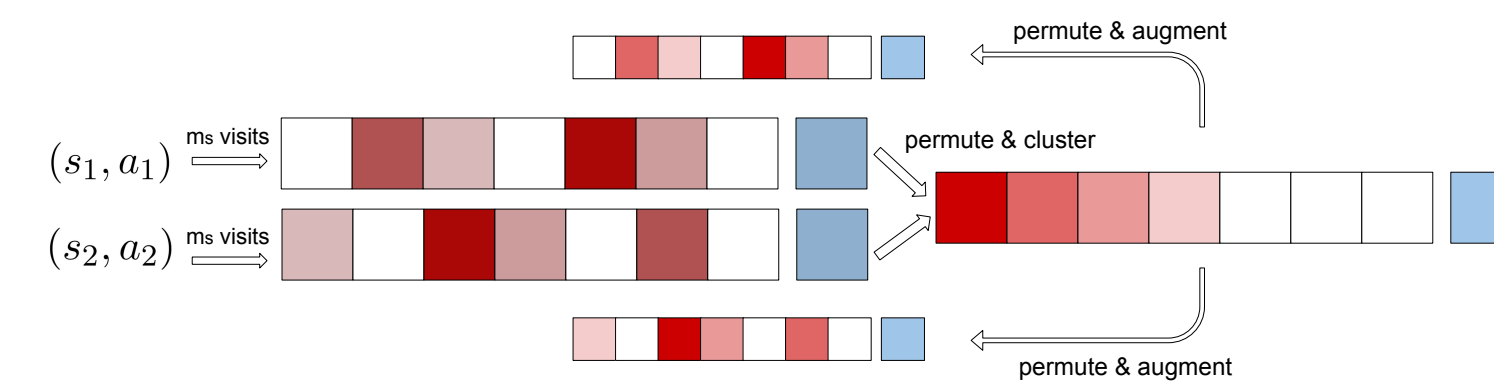
The Old Way: Direct Estimation

- Obtain $\hat{\theta}(s, a) = [\frac{\mathbf{n}(s, a, \cdot); R(s, a)}{n(s, a)}]$ once $n(s, a)$ reaches a **known threshold** m ;



Our Way: Augmented Estimation with Transition Templates

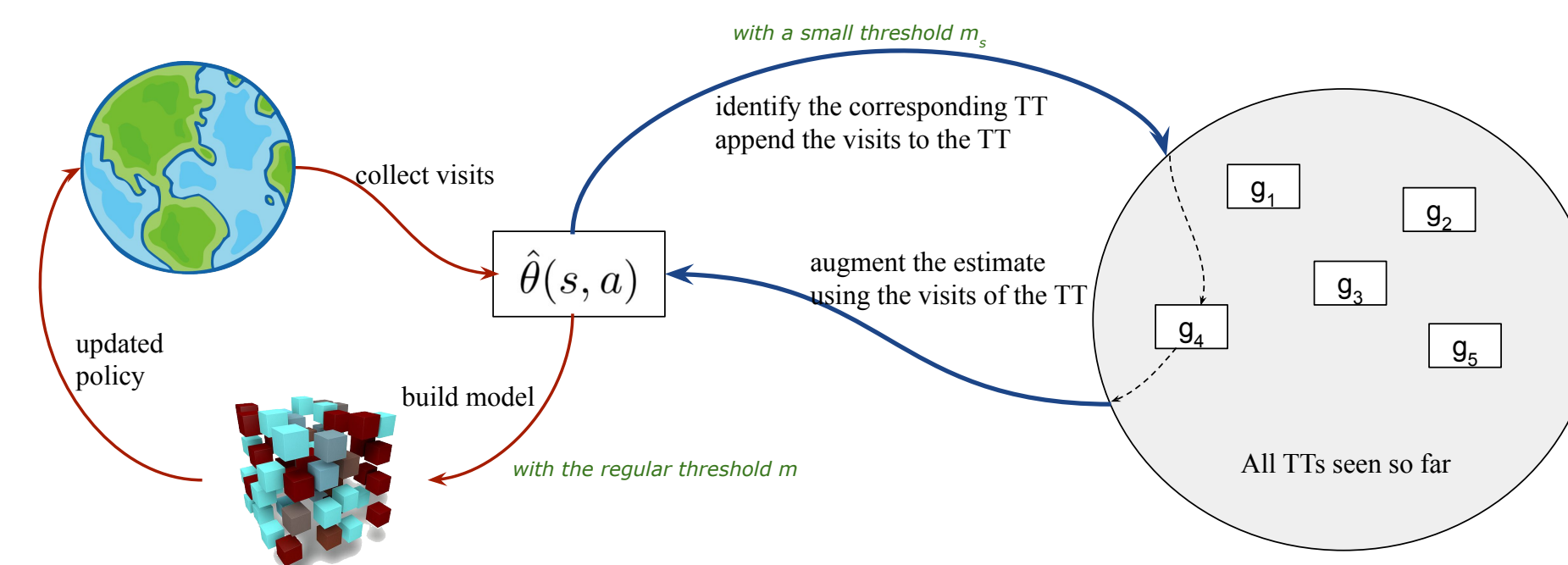
1. **Rough estimation:** obtain $\hat{\theta}(s, a) = [\frac{\mathbf{n}(s, a, \cdot); R(s, a)}{n(s, a)}]$ once $n(s, a)$ reaches a **small known threshold** m_s ;
2. **Permutation:** permute $\hat{\theta}(s, a)$ to its corresponding TT $\tilde{\mathbf{g}}(s, a)$;
3. **Template identification:** identify the permuted estimate $\tilde{\mathbf{g}}(s, a)$ as one of the existing groups of TTs \mathbf{g} .
4. **Augmentation:** obtain a more confident estimate $\hat{\theta}(s, a)$ by permuting back the accumulated knowledge of its corresponding \mathbf{g} .



ALGORITHMS FOR MTRL

O-TempLe: Online Template Learning

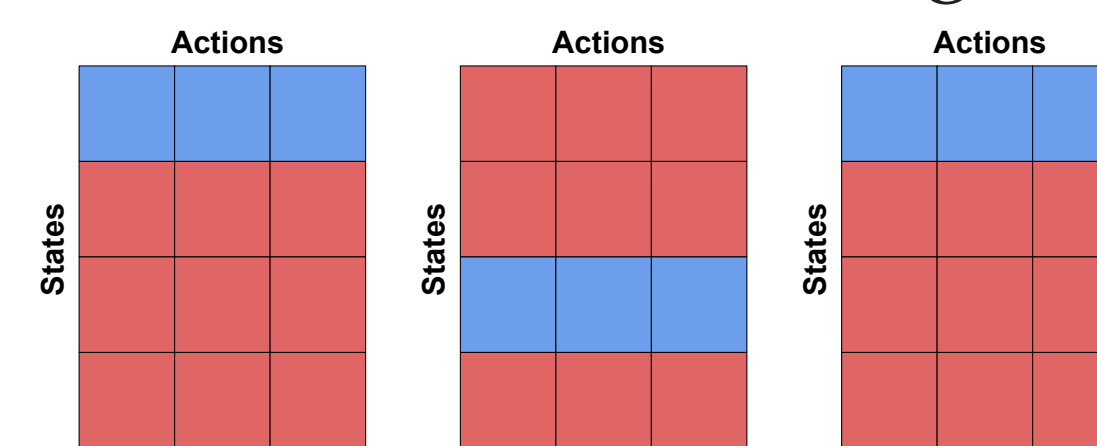
- Streaming-in tasks \mathcal{M} may have **different** state space S , action space A , dynamics $p(\cdot|\cdot, \cdot)$ and rewards $r(\cdot, \cdot)$.
- # of MDPs in \mathcal{M} can be arbitrarily large.



1. **Cluster** estimations to groups of TTs;
2. **Augment** the estimation of dynamics

FM-TempLe: Finite-Model Template Learning

- If # of MDPs is small, we can **further accelerate** the learning.
- Phase 1:** collect and cluster MDP models according to their TTs
Phase 2: identify the model of any new task by its TTs.



THEORETICAL RESULTS

Sample Complexity of O-TempLe

Suppose there are G underlying TTs in total. For any $\epsilon > 0, 1 > \delta > 0$, running O-TempLe on T tasks, each for at least $\mathcal{O}(\frac{D^2 SA}{\omega^2} \ln \frac{1}{\delta})$ steps, generates at most $\tilde{\mathcal{O}}(\frac{SGV_{\max}^3}{\epsilon^3(1-\gamma)^3} + \frac{TS AV_{\max}}{\omega^2 \epsilon(1-\gamma)})$ non- ϵ -optimal steps, with probability at least $1 - \delta$.

D : the diameter of the MDP;

ω : the error tolerance for TT identification, see the paper for more details.

Remarks. (1) The sample complexity of RMax (for T tasks) is $\tilde{\mathcal{O}}(\frac{TS^2 AV_{\max}^3}{\epsilon^3(1-\gamma)^3})$. (2) O-TempLe achieves **linear dependence** on S and A , the cardinality of state space and action space.

Sample Complexity of FM-TempLe

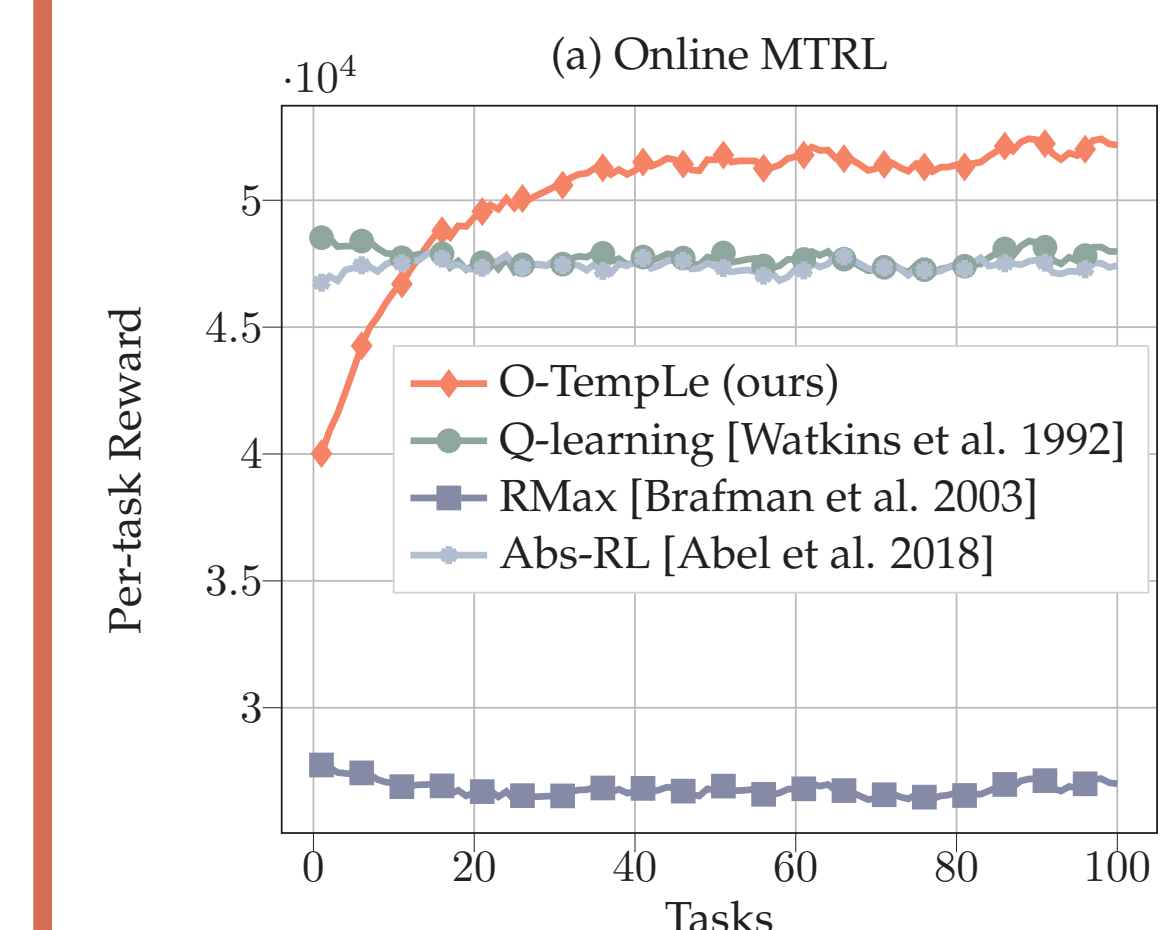
FM-TempLe on T tasks follows ϵ -optimal policies for all but $\tilde{\mathcal{O}}(\frac{SGV_{\max}^3}{\epsilon^3(1-\gamma)^3} + \frac{T_1 SA V_{\max}}{\omega^2 \epsilon(1-\gamma)} + \frac{(T-T_1) DC^2 V_{\max}}{\omega^2 \epsilon(1-\gamma)})$ steps with probability at least $1 - \delta$.

$T_1 = \Omega(\frac{1}{p_{\min}} \ln \frac{C}{\delta})$ is the number of tasks in the first phase

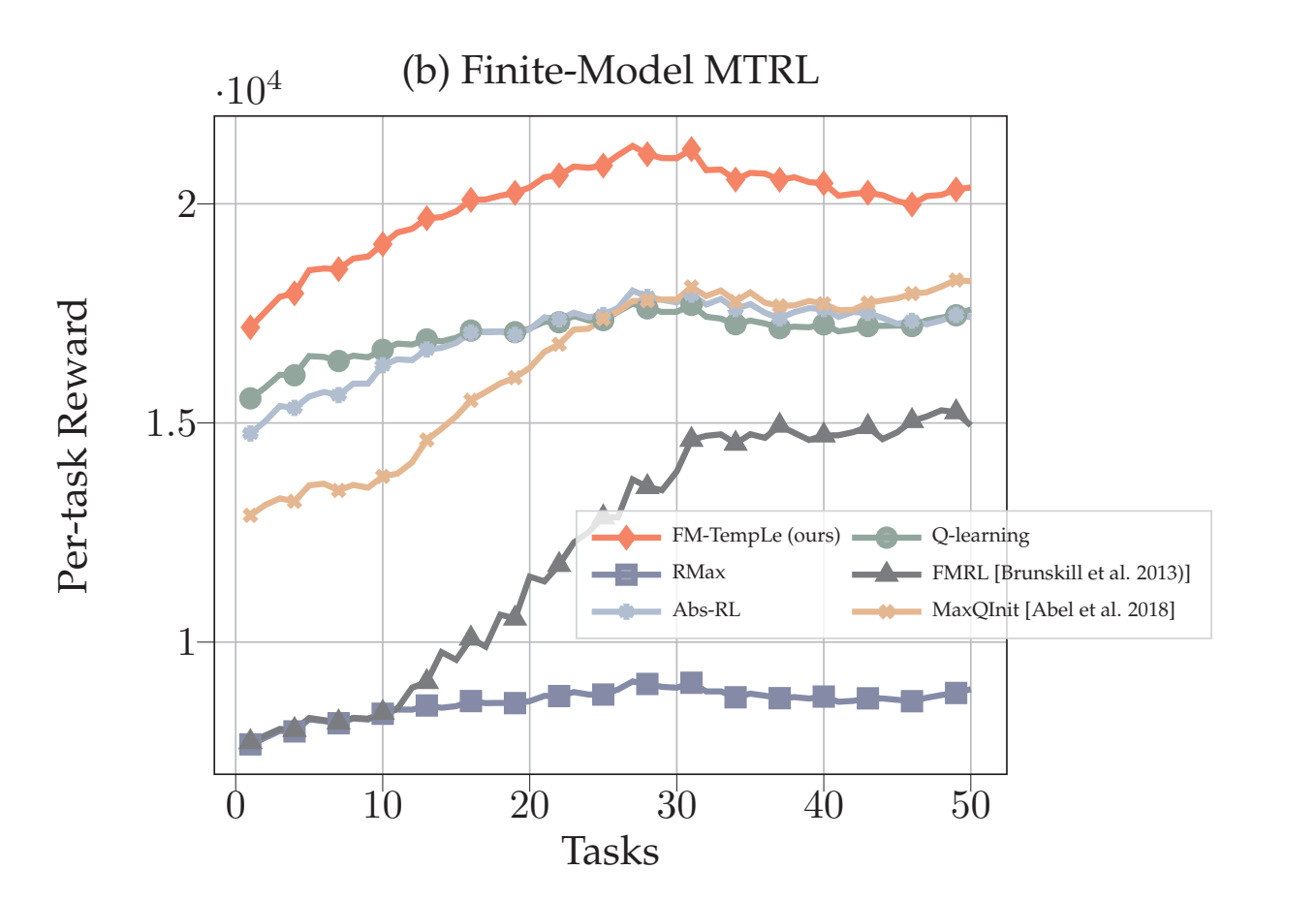
p_{\min} is the minimal probability for a task to be drawn from \mathcal{M} .

Remarks. (1) If $DC^2 < SA$ and $T \gg T_1$, FM-TempLe requires less samples than O-TempLe. (2) When T is large and T_1 is small, FM-TempLe can get rid of the dependence on S and A .

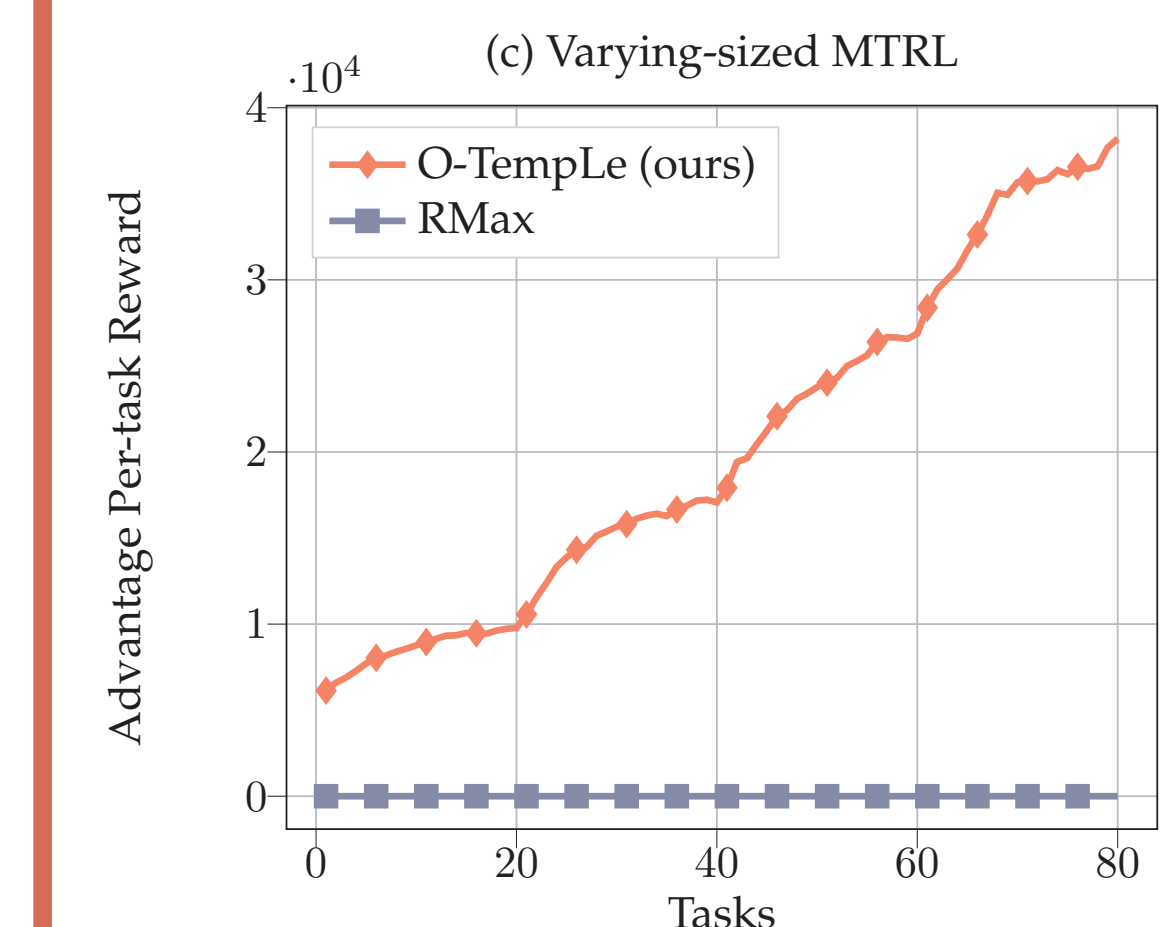
EXPERIMENT RESULTS



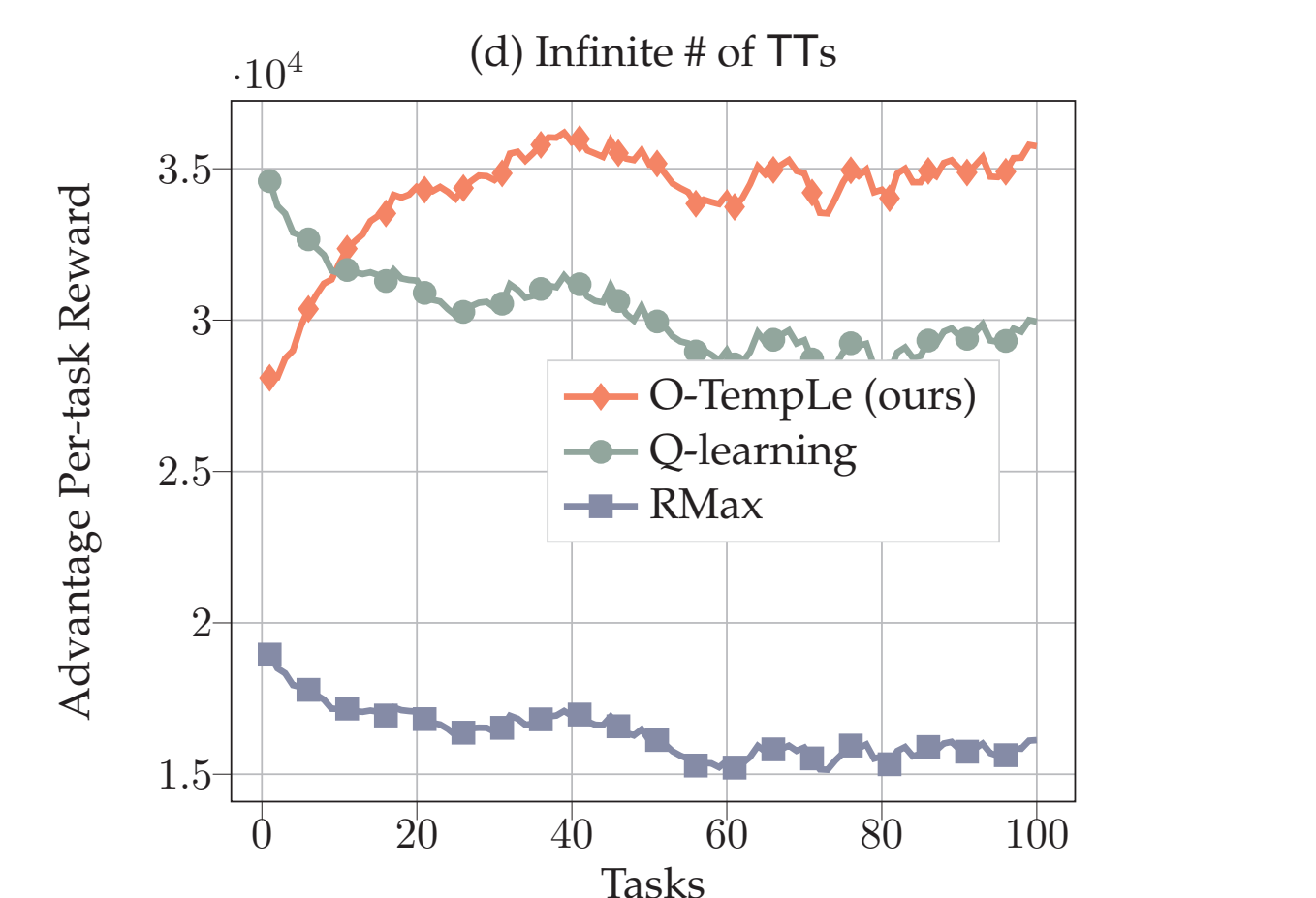
(a) 100 maze tasks with random combinations of 3 "landforms"



(b) 50 tasks sampled from 2 underlying maze models.



(c) mazes with different sizes.



(d) landform ~ mixture of Gaussian