

TEXPLORE and APPLD

Learning Presentation

Nice Wang

wangxiaonannice@gmail.com

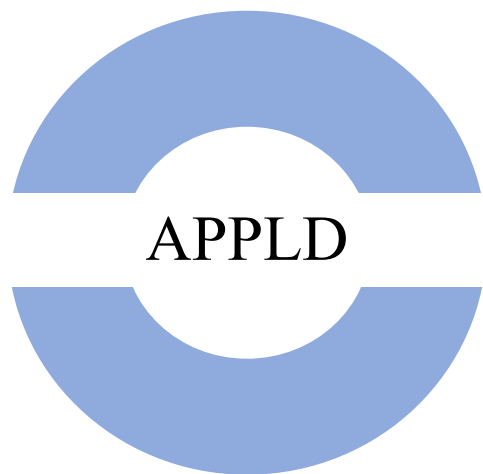
2023.03



Learning Summarization – Page 3

RTMBA Parallel Architecture with UCT(λ) Planning Method – Page 4

Thinking of Future Work – Page 5



Learning Summarization and Thinking of Future Work – Page 6

TEXPLORE: real-time sample-efficient reinforcement learning for robots

Challenges among RL to make it generally applicable to Robot control tasks:

1. Limited exploration: learn from few samples
2. Learn from continuous state representations
3. In face of sensor/actuator delays
4. Learn while taking actions continually in real-time (computationally efficient)

MDP formalism for the task:

$\mathcal{M}(S, A, R, T)$, S : states, A : actions, $R(s, a)$: reward, T : transition

a. For transition, that is posterior: $T(s, a, s') = P(s'|s, a)$

b. Policies come from value function:

$$\text{Bellman equation: } Q^*(s, a) = R(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q^*(s', a')$$

$$\text{The Optimal policy: } \pi(s) = \operatorname{argmax}_a Q^*(s, a)$$

Sample-efficient

Real-time

Learn delay information

Solutions in TEXPLORE:

1. Multi-threaded parallel architecture and MCTS -> RTMBA Parallel Architecture with UCT(λ) Planning Method

2. Model learning and predicting:

a. supervised learners:

n feature models: $s_i^{rel} = featModel_i(s, a)$ corresponds to i th feature of $s' - s$ learning relative transition

one reward model: $r = rewardModel(s, a)$

Better select inputs with correct delay

b. decision tree algorithm

c. using supervised model's ability of generalization: make predictions for unseen or infrequently visited states

3. k-Markov approach: additionally use k previous actions for search

4. Linear regression trees to model continuous domains

5. Random forests to provide targeted, limited exploration:

a. m trees for each model

b. each tree update each with prob. ω

randomness when choosing splits

c. when planning: $Q(s, a) = \frac{1}{m} \sum_{i=1}^m R_i(s, a) + \gamma \frac{1}{m} \sum_{i=1}^m \sum_{s'} P_i(s'|s, a) \max_{a'} Q(s', a')$

suboptimal

toward $s - a$ with higher values

avoiding $s - a$ has low values

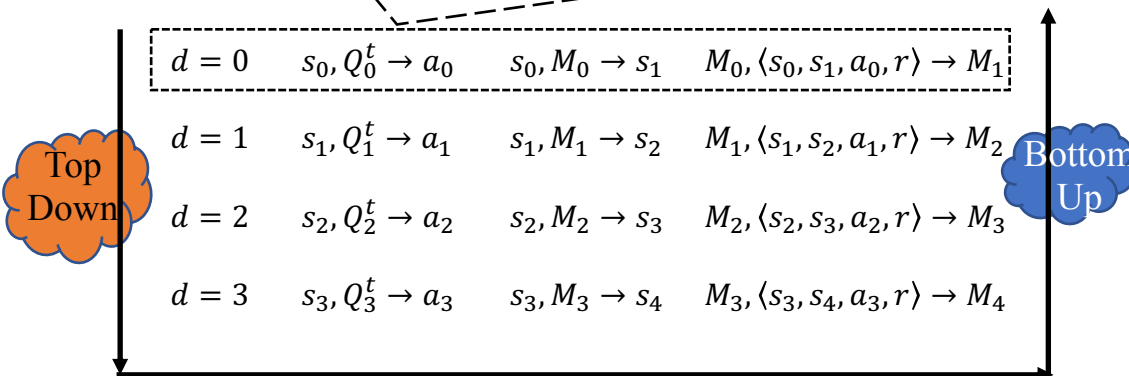
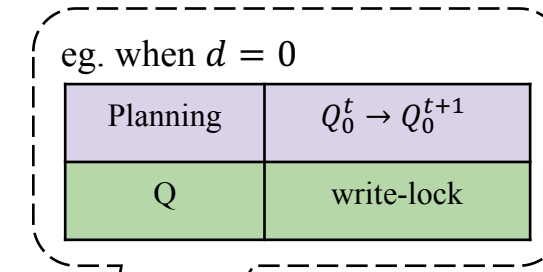
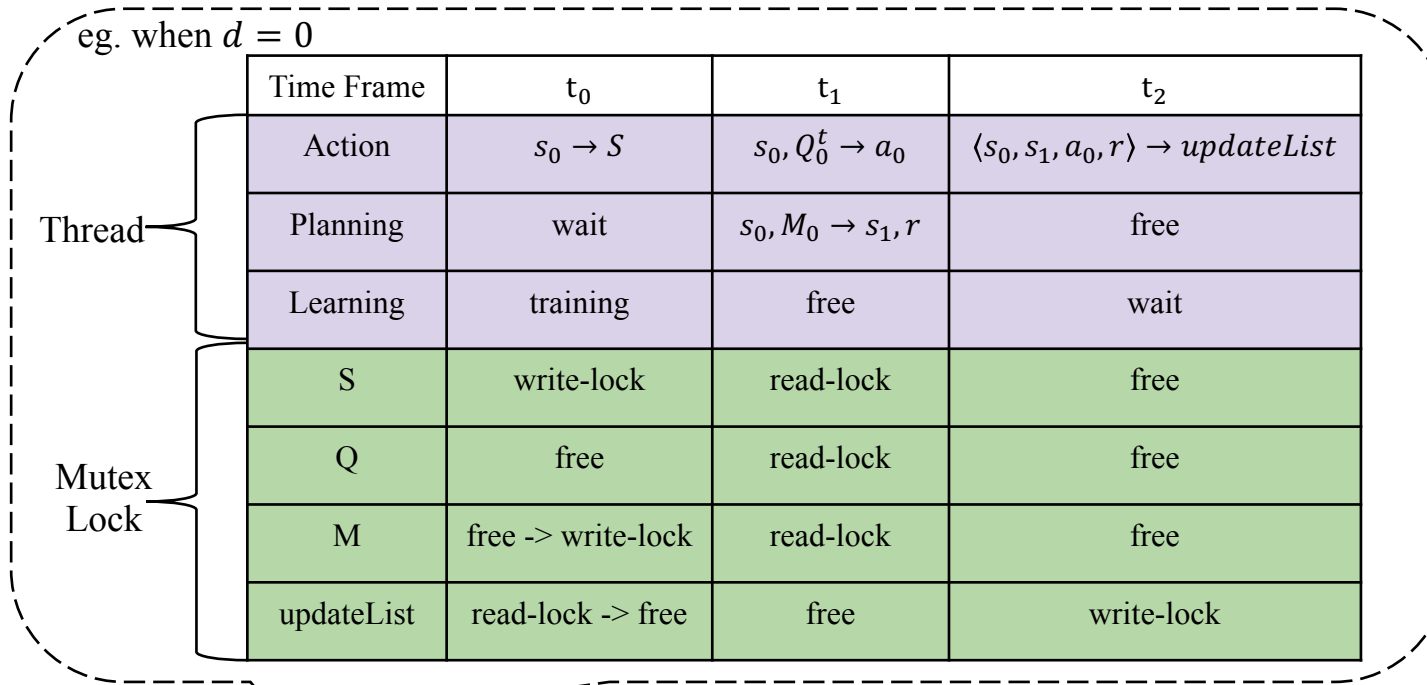
explore optimistic

avoid pessimistic

d. decision tree's ability of generalization

RTMBA Parallel Architecture with UCT(λ) Planning Method \Rightarrow Provide actions continually in real-time at whatever frequency is required.

Suppose that $maxDepth = 4$, rollout from state s_0 at depth 0 as following (Top-down and then Bottom-up):



$$Q_0^{t+1} = \alpha r + \alpha \gamma \{ \lambda r + \lambda \gamma \{ \lambda [r + \gamma (\lambda r + (1 - \lambda) a_{d=3})] + (1 - \lambda) a_{d=2} \} + (1 - \lambda) a_{d=1} \} + (1 - \alpha) Q_0^t$$

$$Q_1^{t+1} = \alpha r + \alpha \gamma \{ \lambda [r + \gamma (\lambda r + (1 - \lambda) a_{d=3})] + (1 - \lambda) a_{d=2} \} + (1 - \alpha) Q_1^t$$

$$Q_2^{t+1} = \alpha [r + \gamma (\lambda r + (1 - \lambda) a_{d=3})] + (1 - \alpha) Q_2^t$$

$$Q_3^{t+1} = \alpha r + (1 - \alpha) Q_3^t$$

Note:

1. Q_3^t : $Q(s_3, -)$ after updating t times

2. $a_{d=3} = \max_a Q_3^{t+1}$

Re-planning from history Q
 \rightarrow Speed up Q learning for new M

Thinking of Future Work of TEXPLORE

Supplementary Advantages of TEXPLORE:

1. Plan actions accordingly to a state-action visit-count to favor less-visited states

Shortages of TEXPLORE:

1. Random forests not quickly when effects of actions not generalization across states
2. Suboptimal:
 - a. s - a pair which can't be predicted from neighbors will not be learned
 - b. sometimes without high-rewarding s - a pairs
3. Partially observable

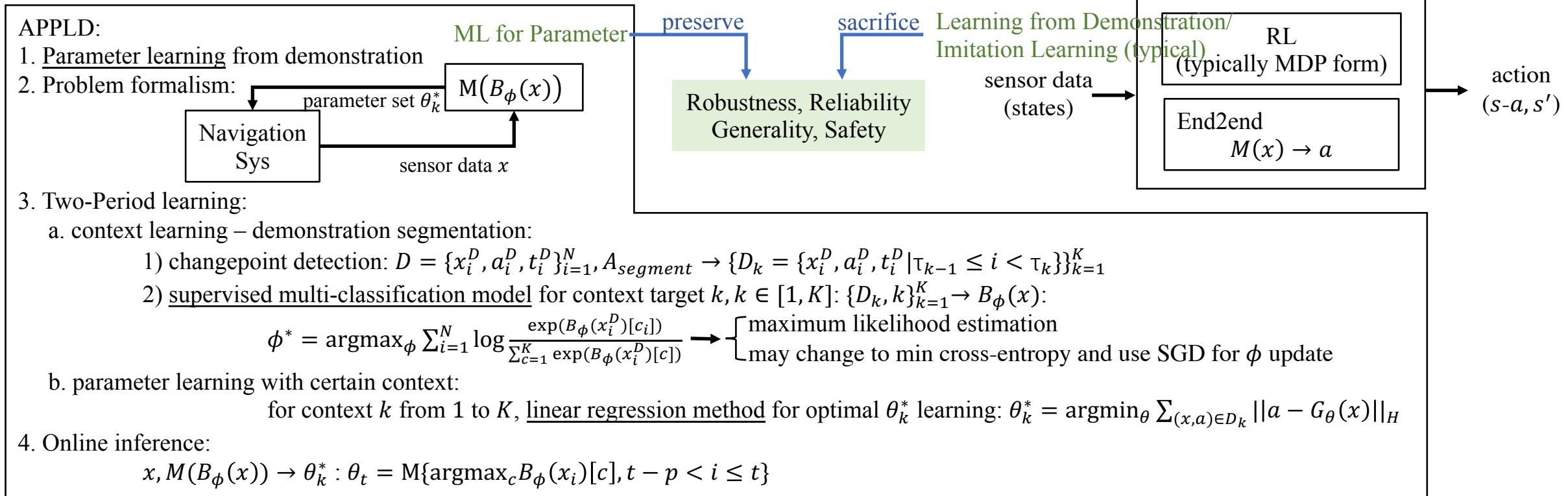
Some Thinking of Future Works:

1. For partially observable:
 - a. stochastic transition for unobservable
 - b. POMDPs
 2. Dealing with continuous actions
 3. Developmental and lifelong learning, external rewards
 4. Resetting visit counts according to model change frequency
 5. Further parallelization with multi-cores
 6. Modeling delays themselves as part of the state-space (reference [1]):
 - a. MDP \rightarrow RDMDP (random delay MDP):
 $RDMDP(MDP, p_\omega, p_\alpha)$, p_ω : observation delay, p_α : action delay
 - b. Delay-Correcting Actor-Critic (DCRC): use off-policy multi-step value estimation
 7. Gated decision trees applied in model-based reinforcement learning (eg. MoET, reference [2])
 8. Dealing with task that has high dimensionality of the action space: refer to RTE (reference [3])
 9. More suitable models (may such as time sequential prediction model) used in model-based reinforcement learning
 10. Online and few-shot learning with few samples within the framework of robotic reinforcement learning
- summarized by the paper
- ideas from other literature
- my own ideas

Reference:

- [1] Bouteiller, Yann, et al. "Reinforcement learning with random delays." International conference on learning representations. 2021.
- [2] Vasić, Marko, et al. "MoET: Mixture of Expert Trees and its application to verifiable reinforcement learning." Neural Networks 151 (2022): 34-47.
- [3] Chazizlygeroudis, Konstantinos, Vassilis Vassiliades, and Jean-Baptiste Mouret. "Reset-free trial-and-error learning for robot damage recovery." Robotics and Autonomous Systems 100 (2018): 236-250.

APPLD: Adaptive Planner Parameter Learning from Demonstration



Some Thinking of Future Works:

1. Clustering similar contexts together
2. Perform parameter learning and changepoint detection jointly
3. Param-learning from interventions/evaluative feedback (ref [1-2])
4. Reinforcement learning approach (for param or policy, ref [3-4])

Reference:

- [1] Wang, Zizhao, et al. "Appli: Adaptive planner parameter learning from interventions." 2021 IEEE international conference on robotics and automation (ICRA). IEEE, 2021.
- [2] Wang, Zizhao, et al. "Apple: Adaptive planner parameter learning from evaluative feedback." IEEE Robotics and Automation Letters 6.4 (2021): 7744-7749.
- [3] Xu, Zifan, et al. "Applr: Adaptive planner parameter learning from reinforcement." 2021 IEEE international conference on robotics and automation (ICRA). IEEE, 2021.
- [4] Spencer, Jonathan, et al. "Expert Intervention Learning: An online framework for robot learning from explicit and implicit human feedback." Autonomous Robots (2022): 1-15.