

① Robot skill learning: from real-world to simulation and back.

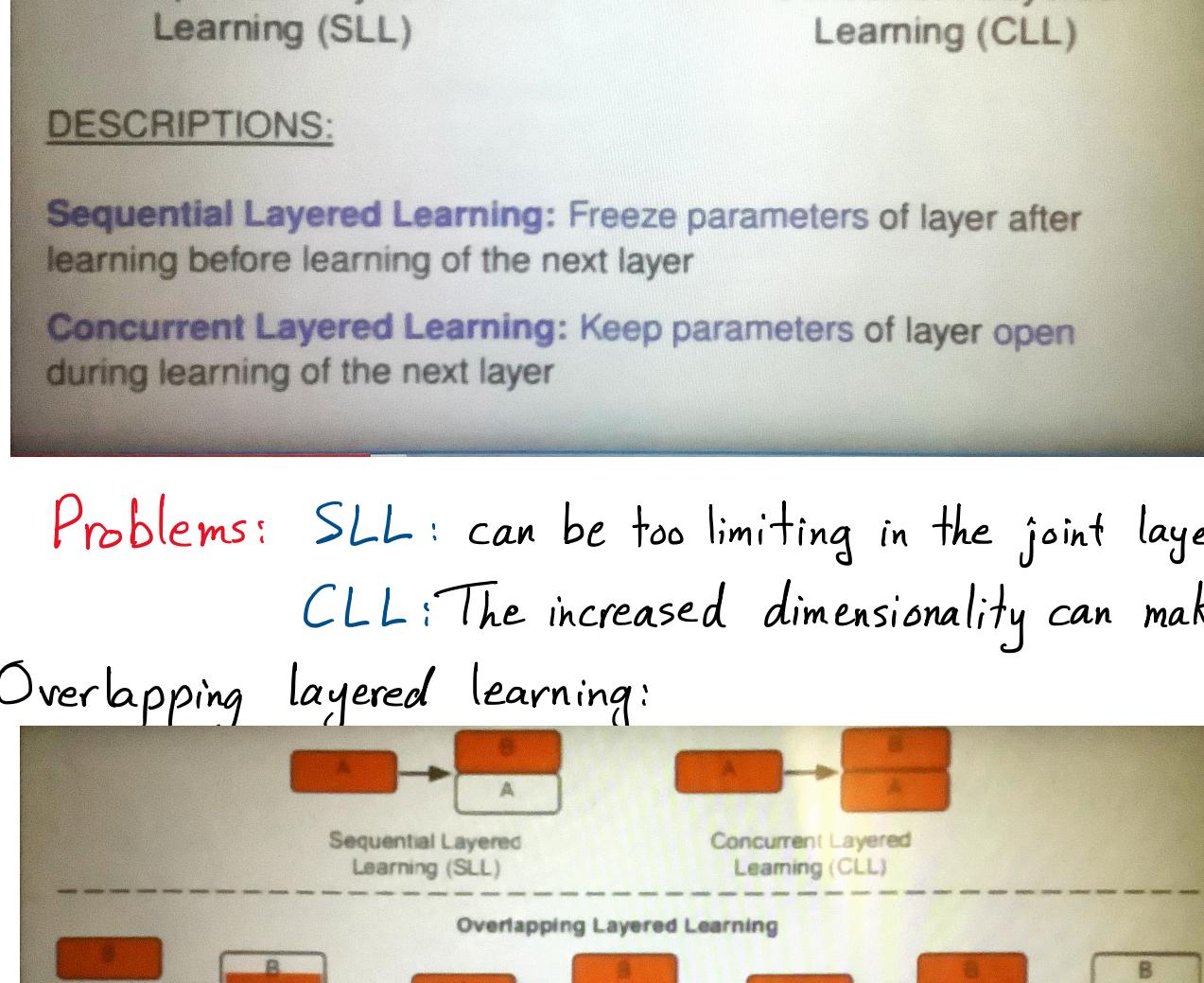
- Research question:

To what degree can autonomous intelligent agents learn in the presence of teammates and/or adversaries in realtime dynamic domain?

- Layered learning:

- For domains too complex for tractably mapping state features S to outputs O .
- Hierarchical subtask decomposition given $\{L_1, \dots, L_n\}$
- Machine learning: exploit data to train, adapt
Learning in one layer feeds into next layer.

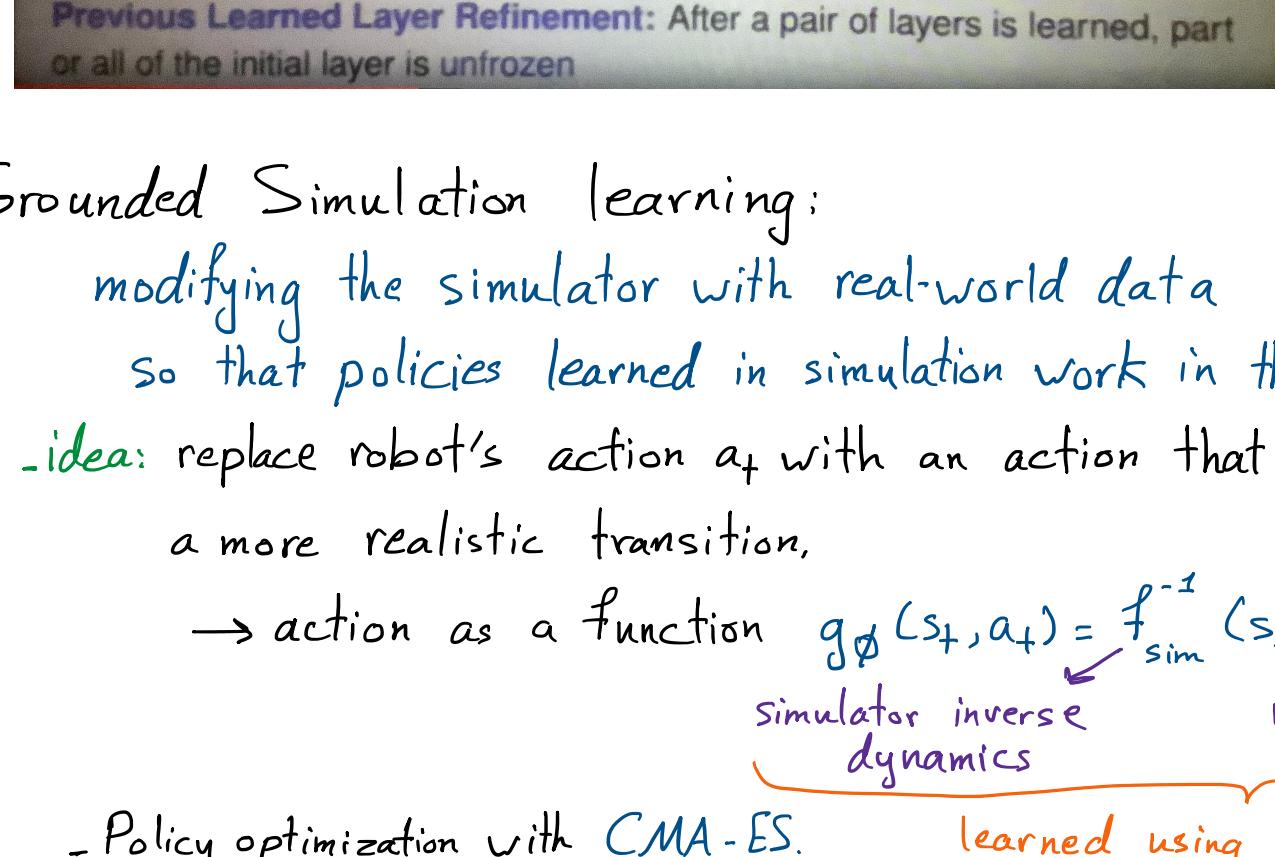
- Layered learning paradigms:



Problems: SLL: can be too limiting in the joint layer policy search space.

CLL: The increased dimensionality can make learning harder.

- Overlapping layered learning:



- Grounded Simulation learning:

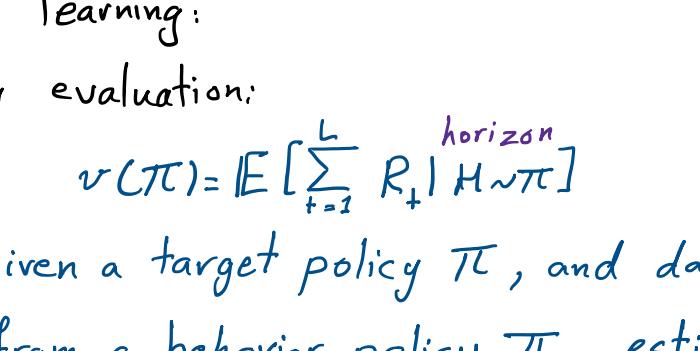
modifying the simulator with real-world data
so that policies learned in simulation work in the real-world.

idea: replace robot's action a_t with an action that produces a more realistic transition,

$$\rightarrow \text{action as a function } g_\theta(s_t, a_t) = f_{\text{sim}}^{-1}(s_t, f(s_t, a_t))$$

simulator inverse dynamics ↓ robot dynamics

- Policy optimization with CMA-ES. learned using supervised learning



- TEXPLORER for Robot RL [Hester & Stone '13]

- Ad Hoc Teams: Ad hoc team player is an individual

- Unknown team mates.

- Teammates likely sub-optimal: no control.

[Genter et al. '15]

② Efficient Robot skill learning:

- Off-policy policy evaluation:

$$\text{performance: } v(\pi) = \mathbb{E} \left[\sum_{t=1}^L R_t \mid H \right]$$

evaluation: given a target policy π , and data generated

from a behavior policy π_b , estimate $v(\pi) \rightarrow \hat{v}$

$$\text{metric: } \text{MSE}[\hat{v}] = (v(\pi) - \hat{v})^2$$

- Importance-Sampling policy evaluation:

Re-weight reward totals for each trajectory:

$$IS(\pi, H, \pi_b) := \left(\prod_{t=1}^L \frac{\pi(A_t | S_t)}{\pi_b(A_t | S_t)} \right) \times \sum_{t=1}^L R_t$$

↓
Relative likelihood

$$\text{- Behavior policy gradient: } \theta_{i+1} = \theta_i - \alpha \frac{\partial}{\partial \theta} \text{MSE}[IS(\pi, H, \pi_\theta)]$$

$$\frac{\partial}{\partial \theta} \text{MSE}[IS(\pi, H, \pi_\theta)] = \mathbb{E}_{\pi_\theta} \left[-IS(\pi, H, \pi_\theta)^2 \sum_{t=0}^L \frac{\partial}{\partial \theta} \log \pi_\theta(A_t | S_t) \right]$$

- Regression importance sampling: IS with maximum likelihood policy

1. learn π_b with supervised learning

$$\pi_\theta = \arg \max_{\theta} \sum_{(S, A) \in D} \log \pi_\theta(A | S)$$

2. Importance weight for trajectory H :

$$IS(\pi, H, \pi_\theta) = \prod_{t=1}^L \frac{\pi_\theta(A_t | S_t)}{\pi_\theta(A_t | S_t)} \sum_{t=0}^L \gamma^t R_t$$