**(a) Conservative exploration framework**
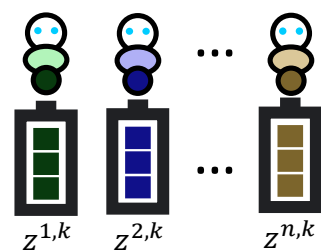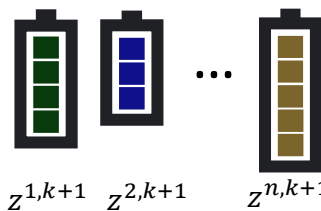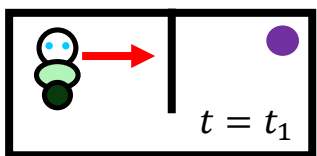
Unified Epigraph Optimization (Sec. 3-1)

$$\max_z z := \sum_i z^i \quad \text{s.t.} \quad \max_{\boldsymbol{\pi}} \min\{ J_{\text{ext}}(\boldsymbol{\pi}) - J_{\text{ext}}(\hat{\boldsymbol{\pi}}) , J_{\text{int,i}}(\pi^i) - z^i \}$$

Outer optimization (Alg. 2)

$z^{1,k} \quad z^{2,k} \quad \cdots \quad z^{n,k}$

Iterate $k \to k+1$

$z^{1,k+1} \quad z^{2,k+1} \quad \cdots \quad z^{n,k+1}$

Inner optimization (Alg. 1)

$t = t_1$: $\min\{ A^{\text{ext}}(s_{t_1}, \boldsymbol{a}_{t_1}) , V^{\text{int}}(s^i_{t_1}) - z^i_{t_1} \}$

$t = t_2$: $\min\{ A^{\text{ext}}(s_{t_2}, \boldsymbol{a}_{t_2}) , V^{\text{int}}(s^i_{t_2}) - z^i_{t_2} \}$

Budget dynamics
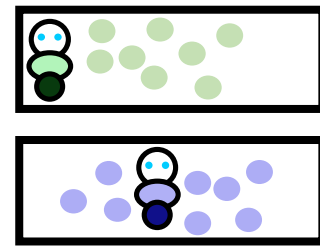$\gamma z^i_{t_2+1} = z^i_{t_2} - r^{\text{int}}_{i,t_2}$

$t = t_3$: $\min\{ A^{\text{ext}}(s_{t_3}, \boldsymbol{a}_{t_3}) , V^{\text{int}}(s^i_{t_3}) - z^i_{t_3} \}$
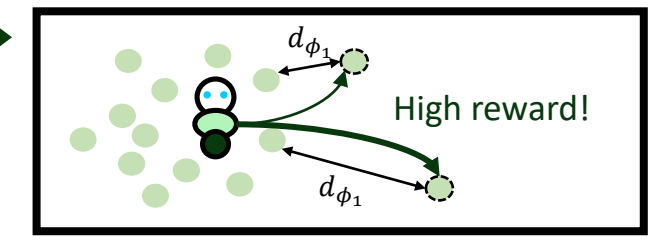
: exploration budget ($z^i_t$)  : task progress  : exploration

**(b) Successor distance-based intrinsic reward**

$r^{\text{int}}_{i,t}$

Factorized Per-Agent Episodic Novelty : $r^{\text{int}}_{i,t}$ (Sec. 3-2)

**Agent-centric data**

**Per-Agent Episodic novelty**

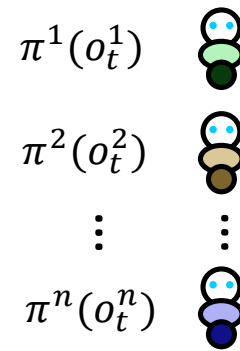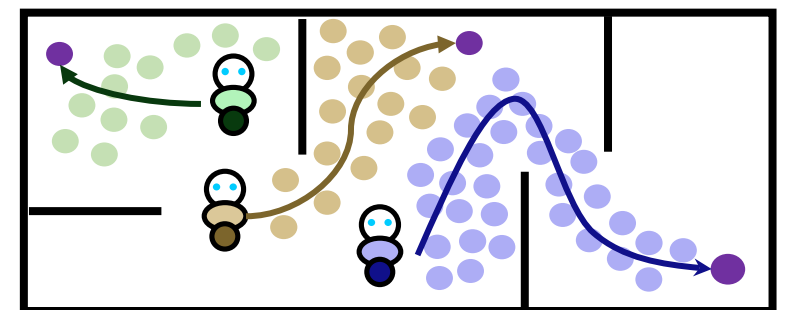$d_{\phi_1}$
$d_{\phi_2}$
$d_{\phi_n}$

$d_{\phi_1}$

High reward!

$d_{\phi_1}$

: visited state of ego agent  : new state  , : visited state of other agent  $d_{\phi_i}$: SD network

**(c) Distributed execution**

$\{\pi^i\}$

**Agents**

$\pi^1(o^1_t)$
$\pi^2(o^2_t)$
$\vdots$
$\pi^n(o^n_t)$

$s_t, r^{\text{ext}}_t$

$\boldsymbol{a}_t = (a^1_t, ..., a^n_t)$

**Environment**

, , : visited state of agent 1, 2, $n$  : goal