# LAGMA: LAtent Goal-guided Multi-Agent Reinforcement Learning

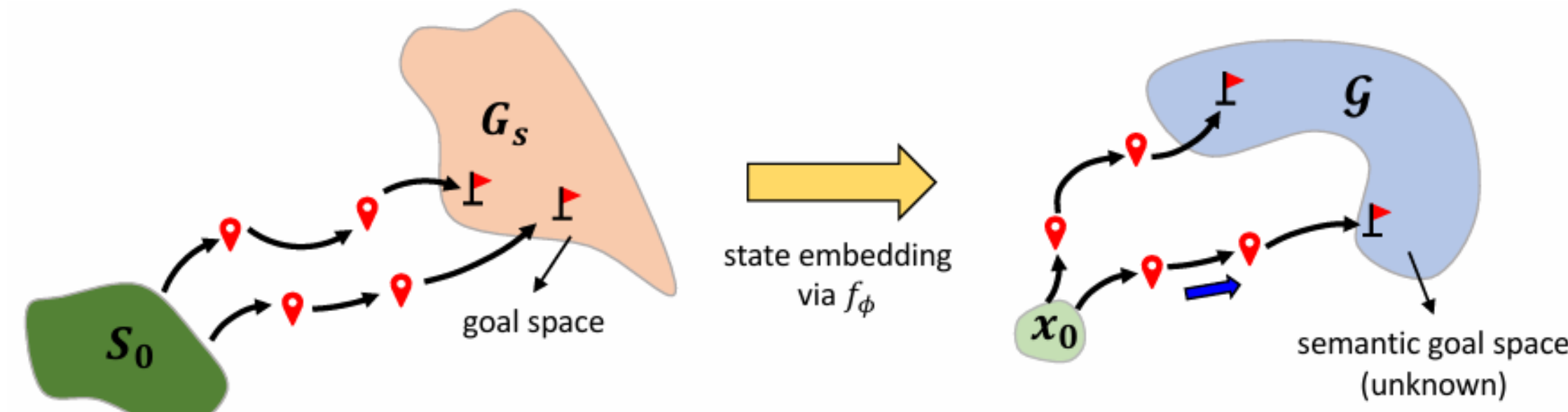Hyungho Na[1] and Il-Chul Moon[1,2]

[1]KAIST, [2]summary.ai

Paper   Code

## Motivation: Efficient Training for MARL

- Goal-conditioned RL (GCRL) has shown good performance in a single agent task, such as complex pathfinding with sparse reward
- GCRL concept has been limitedly applied to MARL task
  - Goal is not explicitly known
  - Partial observation and decentralized execution
  - complex coordination rather than the shortest path finding



- Cooperative MARL problem
  → finding trajectories toward semantic goals in latent space
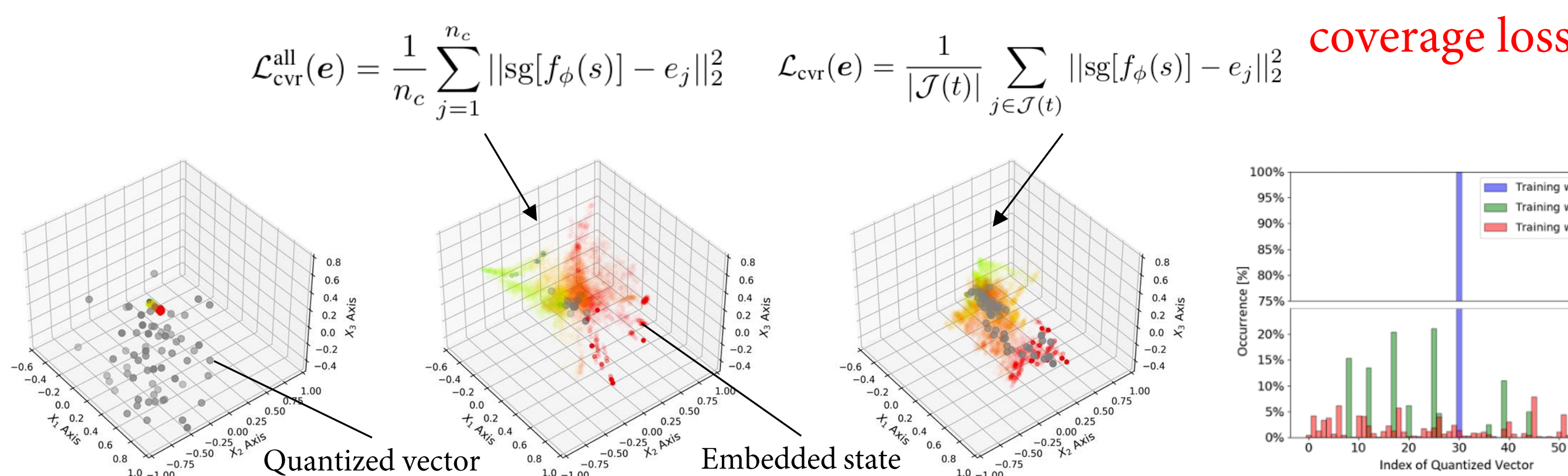
## Key Concept: Goal-reaching Trajectory

- Achieving a common goal in cooperative tasks → **Goal-reaching**
  - $\mathcal{T}$ satisfying the following is considered **Goal-reaching denoted as $\mathcal{T}^*$**

$$\mathcal{T} := \{s_0, \boldsymbol{a_0}, r_0, s_1, \boldsymbol{a_1}, r_1, ..., s_T\} \text{ such that } \Sigma_{t=0}^{T-1} r_t = R_{\max}$$

  - For $\forall s \in \mathcal{T}^*$, $\tau_{s_t}^* := \{s_t, s_{t+1}, ... s_T\}$ is a goal-reaching trajectory

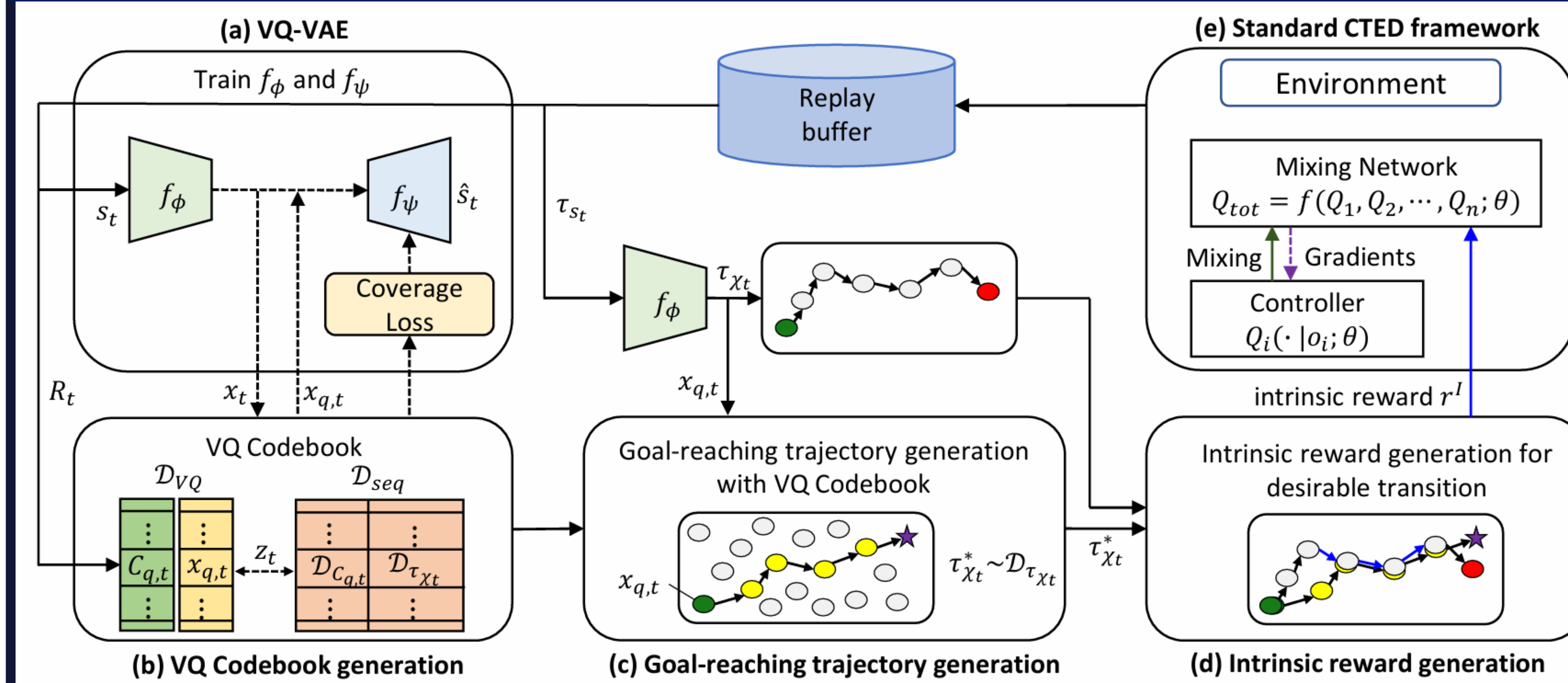## Embedding via modified VQ-VAE

- For quantized embedding space construction
  - VQ-VAE: $\mathcal{L}_{VQ}(\phi, \psi, \boldsymbol{e}) =$

$$||f_\psi([x = f_\phi(s)]_q) - s||_2^2 + \lambda_{vq}||sg[f_\phi(s)] - x_q||_2^2 + \lambda_{commit}||f_\phi(s) - sg[x_q]||_2^2$$

  reconstruction loss    VQ loss    commitment loss

  - Modified VQ-VAE: $\mathcal{L}_{VQ}^{tot}(\phi, \psi, \boldsymbol{e})f = \mathcal{L}_{VQ}(\phi, \psi, \boldsymbol{e}) + \lambda_{cvr}\mathcal{L}_{cvr}(\boldsymbol{e})$

  coverage loss

$$\mathcal{L}_{cvr}^{all}(\boldsymbol{e}) = \frac{1}{n_c}\sum_{j=1}^{n_c}||sg[f_\phi(s)] - e_j||_2^2 \quad \mathcal{L}_{cvr}(\boldsymbol{e}) = \frac{1}{|\mathcal{J}(t)|}\sum_{j\in\mathcal{J}(t)}||sg[f_\phi(s)] - e_j||_2^2$$



Quantized vector     Embedded state

## LAGMA generates a goal-reaching trajectory in latent space and incentivizes transitions towards this reference trajectory



(a) VQ-VAE
(b) VQ Codebook generation
(c) Goal-reaching trajectory generation
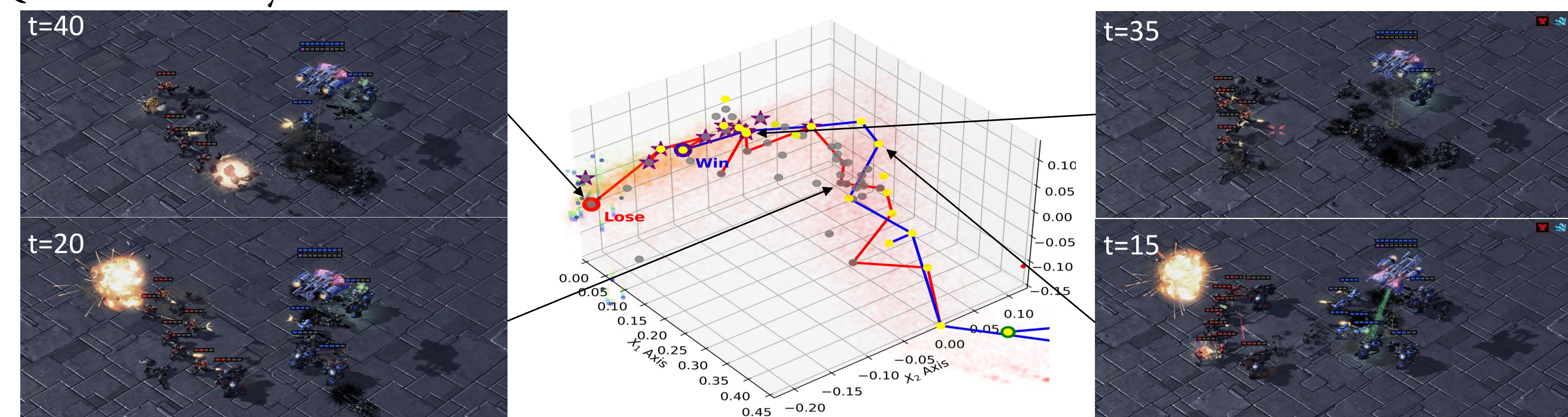(d) Intrinsic reward generation
(e) Standard CTED framework

- **(1) Modified VQ-VAE** is developed for quantized embedding space construction
- **(2) Goal-reaching trajectory** is generated via **extended codebook**
- **(3) Latent Goal-guided intrinsic reward** guarantees a better convergence on optimal policy
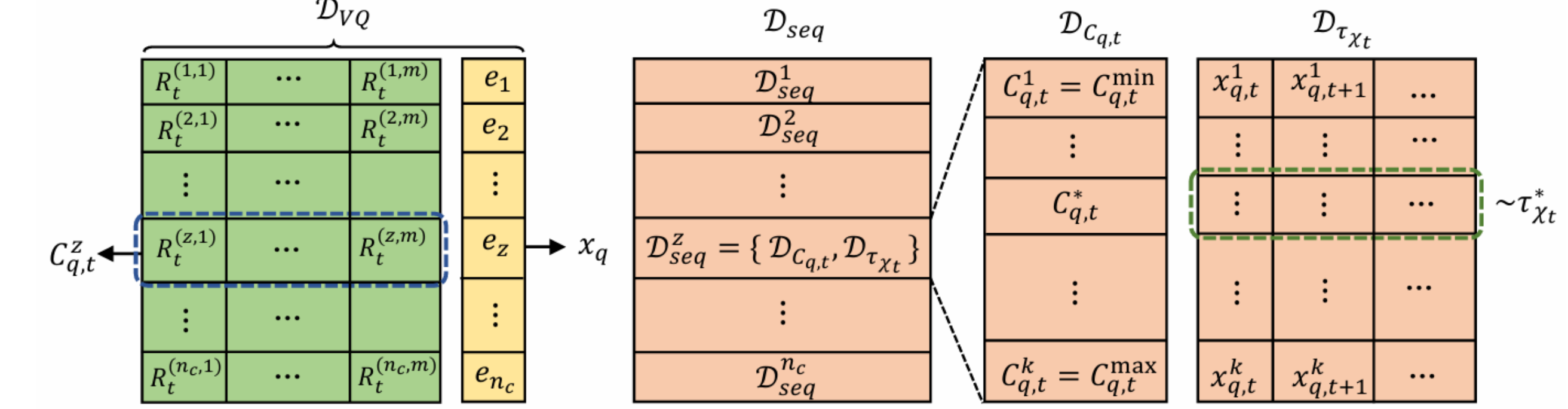
## Visualization in Latent Space

- Overall Objective:

$$\mathcal{L}(\theta) = \left(r^{ext} + r^I + \gamma\max_{\boldsymbol{a'}}Q_{\theta^-}^{tot}(s', \boldsymbol{a'}) - Q_\theta^{tot}(s, \boldsymbol{a})\right)^2$$

- Qualitative Analysis:



t=40    t=35    t=20    t=15

## Goal-reaching Trajectory Generation

- The value of quantized vector is computed as $C_{q,t}(x_{q,t}) = \frac{1}{N_{x_{q,t}}}\sum_{j=1}^{N_{x_{q,t}}}R_t^j(x_{q,t})$
- Note that: $R_t = \Sigma_{i=t}^{T-1}\gamma^{i-t}r_i$ and $\tau_{\chi_t} = [f_\phi(\tau_{x_t})]_q = \{x_{q,t}, x_{q,t+1}, x_{q,t+2}, ..., x_{q,T}\}$
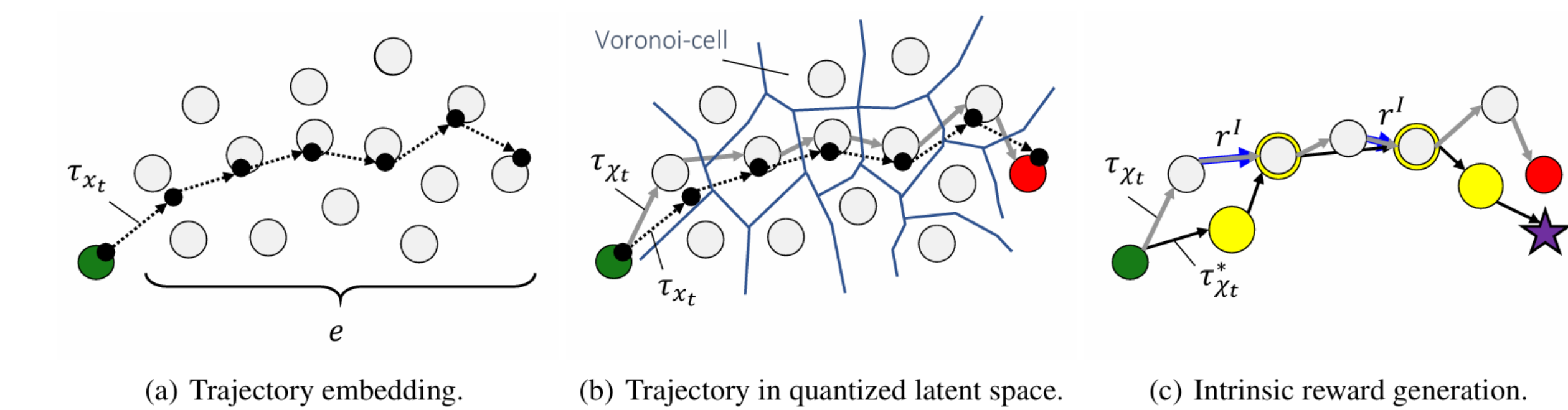


## Training with Latent Goal-guided Incentive

For $s' \in \tau_{\chi_t}^*$, $r^I(s') := \gamma(C_{q,t}(s') - \max_{\boldsymbol{a'}}Q_{\theta^-}(s', \boldsymbol{a'}))$ guarantees

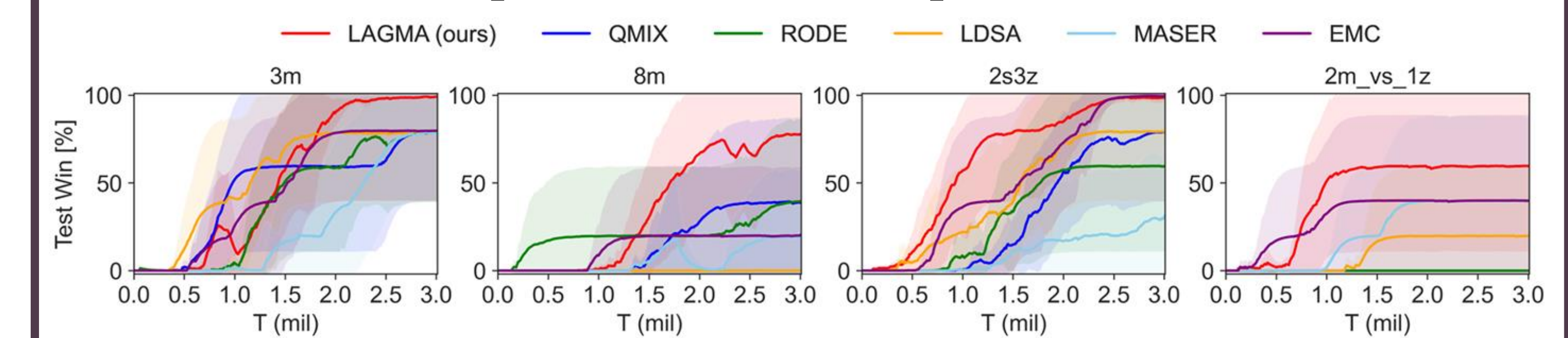$$y = r(s, \boldsymbol{a}) + r^I(s') + \gamma V_{\theta^-}(s') \to y^*$$

where $y^* = r(s, \boldsymbol{a}) + \gamma V^*(s')$    if $x_{q,t+1} \in \tau_{\chi_t}^*$ and $x_{q,t+1} \neq x_{q,t}$



(a) Trajectory embedding.    (b) Trajectory in quantized latent space.    (c) Intrinsic reward generation.

## Experiments

- Performance comparison in SMAC (sparse reward)



LAGMA (ours)   QMIX   RODE   LDSA   MASER   EMC

3m    8m    2s3z    2m_vs_1z

- Ablation on key components
  - Evaluated with controlled parameters



LAGMA (ours)   LAGMA (CL-All)   LAGMA (No-CL)   LAGMA (Cqt-No-Upd)   LAGMA (Cq0)

3m    2s3z    3m    2s3z