③ $Q_\beta(s_t, a$

$t+1$



Finalist Action S

②

$$\cdots_t \mid \theta, \beta) \leftarrow r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1}\sim p(\cdot|s,a), a_{t+1}\sim\pi(\cdot|s)}[Q_{\hat{\beta}}(s_{t+1}, a_{t+1} \mid \cdots$$

**Policy Expansion Module**

$$\langle \pi_\beta | \pi_\theta \rangle \xrightarrow{\;③\;} \begin{bmatrix} Q^\theta_{t+1} \\ Q^\beta_{t+1} \end{bmatrix} \longrightarrow$$

Set

**Multi-action Evalu... Module**

$$\begin{bmatrix} Q^{mean}_{t+1} \\ Q^{std}_{t+1} \end{bmatrix} \xrightarrow{\;④\;}$$

$$a_{t+1} \sim \pi_{\theta,\beta}(\cdot \mid s_t)$$

$$④ \quad a_e = arg_{a\in D}\{Q_{mean}(s_{t+1}, a_t\cdots$$

Temporally - Align... Representation Lear...
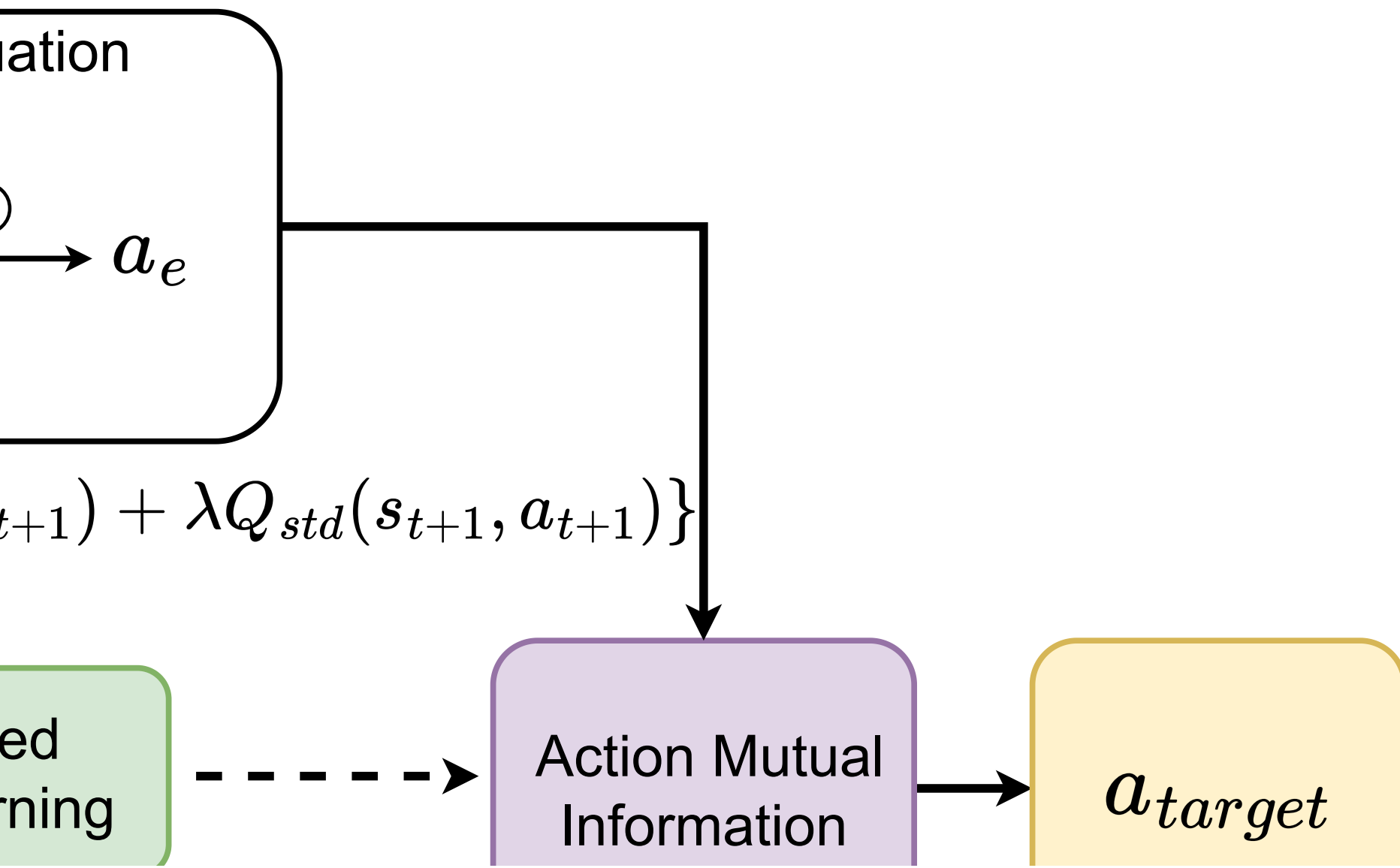
$\vartheta, \beta)]$

uation

$\rightarrow a_e$

$_{t+1}) + \lambda Q_{std}(s_{t+1}, a_{t+1})\}$

ed
rning

Action Mutual
Information

$a_{target}$

Candidate Action S

○ → Offline da

● → Online da

$t$

①  Time step

## Action-Oriented Module

Finalist Action Set

$$a_t \sim P(a_t)$$

$(s_t, a_t, r, s_{t+1})$

Offline    Online
Dataset

En

$(s_t, a_t, r, s_{t+1})$

$(s_{t+1})$

$a_o$

vironment

Agent