# Discrete Codebook World Models for Continuous Control

**Aidan Scannell, Mohammadreza Nakhaei, Kalle Kujanpää, Yi Zhao, Kevin Luck, Arno Solin, Joni Pajarinen**

**University of Edinburgh**
**Finnish Center for Artificial Intelligence (FCAI)**
**Aalto University**

FCAI

fcai.fi

# World Models

$$p(s_{t+1}, r_t \mid s_t, a_t)$$

# World Models

$$p(\boxed{s_{t+1}}, r_t \mid s_t, a_t)$$

# World Models

$$p(\boxed{s_{t+1}}, \boxed{r_t} \mid s_t, a_t)$$
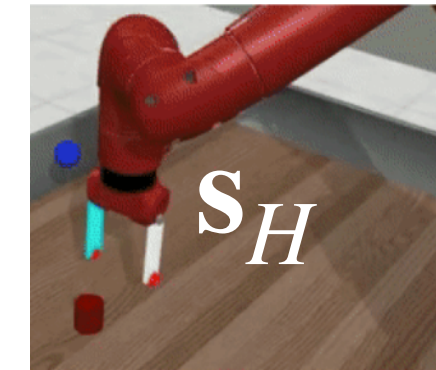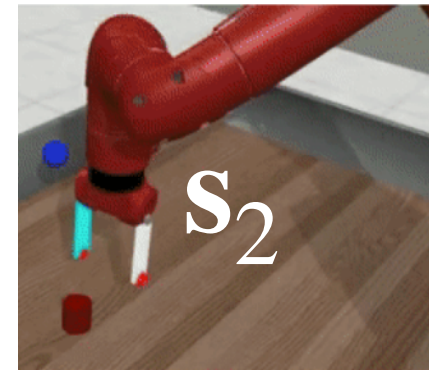
# World Models

$$p(s_{t+1}, r_t \mid s_t, \boxed{a_t})$$

# World Models

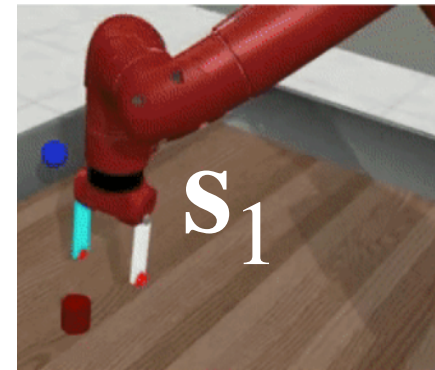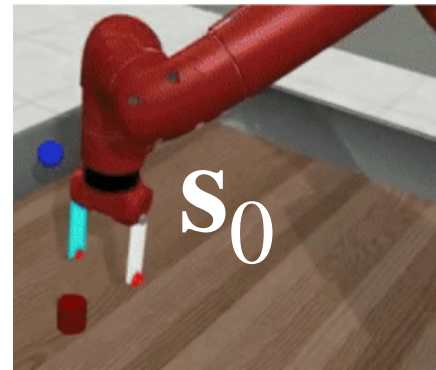$$p(s_{t+1}, r_t \mid s_t, a_t)$$

# World Models

$$p(s_{t+1}, r_t \mid s_t, a_t)$$

# World Models

# World Models



$\mathbf{s}_0$

$\mathbf{s}_1$

$\mathbf{s}_2$

$\mathbf{s}_H$

FCAI

fcai.fi

# World Models



$\mathbf{s}_0$ Encoder

$\mathbf{s}_1$ Encoder

$\mathbf{s}_2$ Encoder

$\mathbf{s}_H$ Encoder

FCAI

fcai.fi

# World Models

# World Models

# World Models



**FCAI**

**fcai.fi**

# World Models

# World Models

# World Models



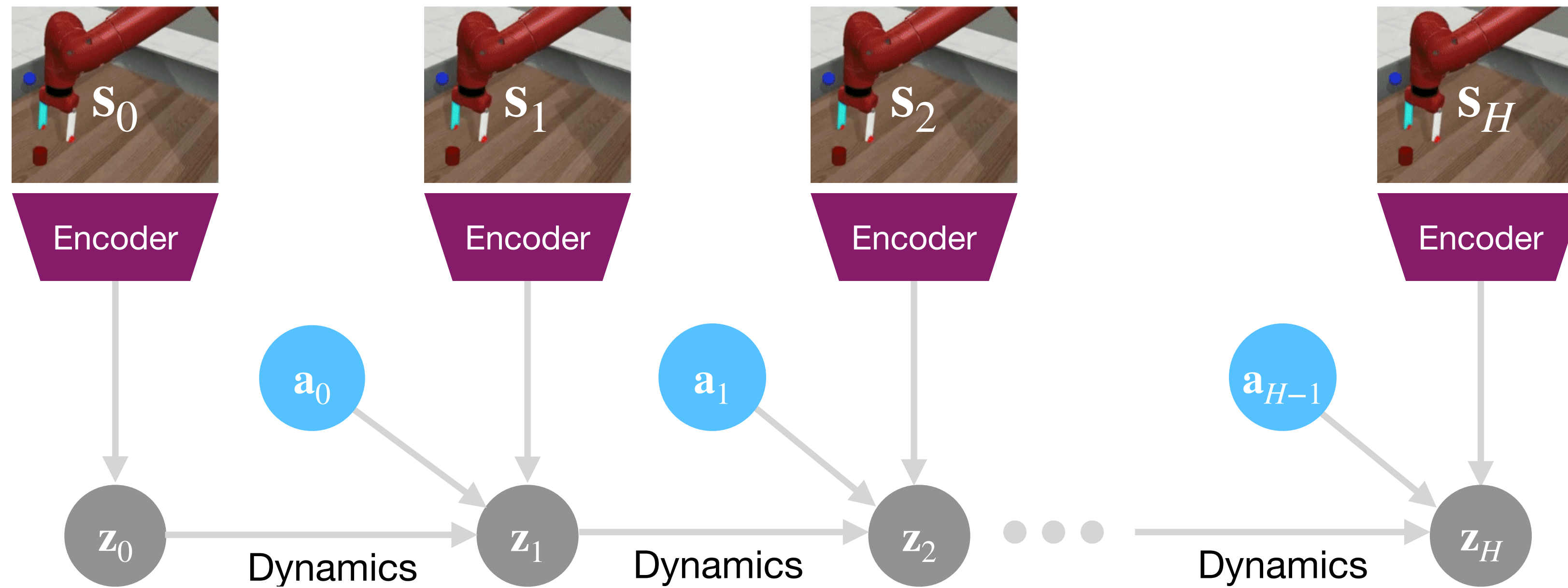**FCAI**                                                    **fcai.fi**

# Latent Space Design Choices

# Latent Space Design Choices

| | Discrete Latent States? | Discrete Encoding Type | Stochastic Dynamics? | Reconstruction? |
|---|---|---|---|---|
| **DreamerV3** | ✅ | One-hot | ✅ | ✅ |
| **TD-MPC2** | ❌ | N/A | ❌ | ❌ |

Danijar Hafner, et al. Mastering diverse domains through world models. arXiv preprint arXiv:2301.04104, 2023.

Nicklas Hansen, et al. TD-MPC2: Scalable, Robust World Models for Continuous Control. In The Twelfth International Conference on Learning Representations, October 2023.

**FCAI**

**fcai.fi**

# Latent Space Design Choices

**1. Do discrete latent spaces offer benefits over continuous ones?**

FCAI

fcai.fi

# Latent Space Design Choices

1. Do discrete latent spaces offer benefits over continuous ones?
2. How does the choice of discrete encoding (e.g., one-hot, label, or codebook encodings) affect performance?

FCAI

fcai.fi

# Latent Space Design Choices

1. Do discrete latent spaces offer benefits over continuous ones?
2. How does the choice of discrete encoding (e.g., one-hot, label, or codebook encodings) affect performance?
3. Are there advantages to modelling the latent dynamics stochastically rather than deterministically?

# Latent Space Design Choices

| | **Discrete Latent States?** | **Discrete Encoding Type** | **Stochastic Dynamics?** | **Reconstruction?** |
|---|---|---|---|---|
| **DreamerV3** | ✅ | One-hot | ✅ | ✅ |
| **TD-MPC2** | ❌ | N/A | ❌ | ❌ |

Danijar Hafner, et al. Mastering diverse domains through world models. arXiv preprint arXiv:2301.04104, 2023.

Nicklas Hansen, et al. TD-MPC2: Scalable, Robust World Models for Continuous Control. In The Twelfth International Conference on Learning Representations, October 2023.

**FCAI**

**fcai.fi**

# Latent Space Design Choices

| | Discrete Latent States? | Discrete Encoding Type | Stochastic Dynamics? | Reconstruction? |
|---|---|---|---|---|
| **DreamerV3** | ✅ | One-hot | ✅ | ✅ |
| **TD-MPC2** | ❌ | N/A | ❌ | ❌ |
| **DC-MPC (ours)** | ✅ | Codebook | ✅ | ❌ |

Danijar Hafner, et al. Mastering diverse domains through world models. arXiv preprint arXiv:2301.04104, 2023.

Nicklas Hansen, et al. TD-MPC2: Scalable, Robust World Models for Continuous Control. In The Twelfth International Conference on Learning Representations, October 2023.

**FCAI**

**fcai.fi**

# DC-MPC: World Model Training

# DC-MPC: World Model Training

# DC-MPC: World Model Training



$\mathbf{s}_0$

Encoder

# DC-MPC: World Model Training



$s_0$

Encoder

$x_0$

# DC-MPC: World Model Training



$s_0$

Encoder

$x_0$

# DC-MPC: World Model Training



$s_0$

Encoder

$x_0$

Codebook $\mathscr{C}$

$c^{(1)}$

$c^{(2)}$

$c^{(3)}$

$c^{(|\mathscr{C}|)}$

Number of codes in codebook $|\mathscr{C}| = 15$

Number of channels $b = 2$

# DC-MPC: World Model Training



$\mathbf{s}_0$

Encoder

$\mathbf{x}_0$

Codebook $\mathscr{C}$

$\mathbf{c}^{(1)}$
$\mathbf{c}^{(2)}$
$\mathbf{c}^{(3)}$

$\mathbf{c}^{(|\mathscr{C}|)}$

Number of codes in codebook $|\mathscr{C}| = 15$

Number of channels $b = 2$

FCAI

fcai.fi

# DC-MPC: World Model Training



$\mathbf{s}_0$

Encoder

$\mathbf{x}_0$

$\mathbf{c}_0$

Codebook $\mathscr{C}$

$\mathbf{c}^{(1)}$

$\mathbf{c}^{(2)}$

$\mathbf{c}^{(3)}$

$\mathbf{c}^{(|\mathscr{C}|)}$

Number of codes in codebook $|\mathscr{C}| = 15$

Number of channels $b = 2$

**FCAI**

**fcai.fi**

# DC-MPC: World Model Training



$\mathbf{s}_0$

Encoder

$\mathbf{x}_0$

$\mathbf{c}_0$

Codebook $\mathscr{C}$

$\mathbf{c}^{(1)}$

$\mathbf{c}^{(2)}$

$\mathbf{c}^{(3)}$

$\mathbf{c}^{(|\mathscr{C}|)}$

Number of codes in codebook $|\mathscr{C}| = 15$

Number of channels $b = 2$

FCAI

fcai.fi

# DC-MPC: World Model Training



$\mathbf{s}_0$

Encoder

$\mathbf{x}_0$

$\mathbf{a}_0$

$\mathbf{c}_0$

Dynamics

$p(\mathbf{c}_1 \mid \mathbf{c}_0, \mathbf{a}_0)$

Codebook $\mathscr{C}$

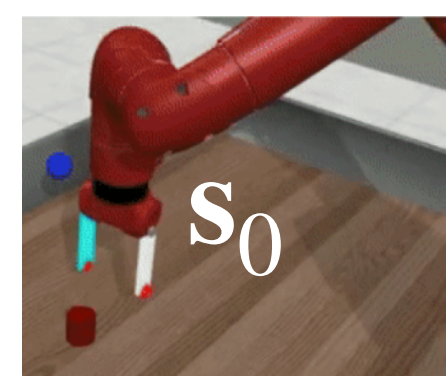$\mathbf{c}^{(1)}$

$\mathbf{c}^{(2)}$

$\mathbf{c}^{(3)}$

$\mathbf{c}^{(|\mathscr{C}|)}$

Number of codes in codebook $|\mathscr{C}| = 15$

Number of channels $b = 2$

FCAI

fcai.fi

# DC-MPC: World Model Training



$\mathbf{s}_0$

Encoder

$\mathbf{x}_0$

$\mathbf{a}_0$

Dynamics

$\mathbf{c}_0$

$p(\mathbf{c}_1 \mid \mathbf{c}_0, \mathbf{a}_0)$

$\hat{r}_1$

Codebook $\mathscr{C}$

$\mathbf{c}^{(1)}$

$\mathbf{c}^{(2)}$

$\mathbf{c}^{(3)}$

$\mathbf{c}^{(|\mathscr{C}|)}$

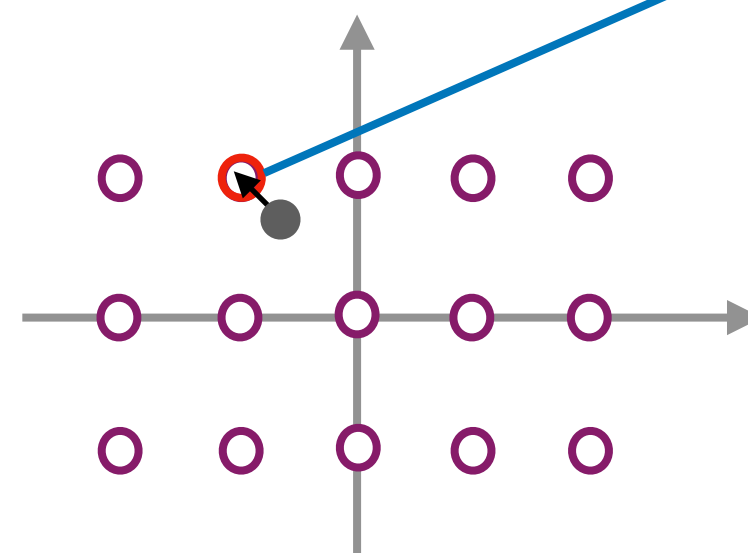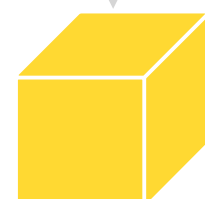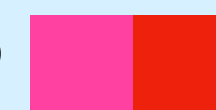Number of codes in codebook $|\mathscr{C}| = 15$

Number of channels $b = 2$

FCAI

fcai.fi

# DC-MPC: World Model Training



$\mathbf{s}_0$

Encoder

$\mathbf{x}_0$

$\mathbf{a}_0$

$\mathbf{c}_0$

Dynamics

$p(\mathbf{c}_1 \mid \mathbf{c}_0, \mathbf{a}_0)$

$\sim$

$\hat{\mathbf{c}}_1$

$\hat{r}_1$

Codebook $\mathscr{C}$

$\mathbf{c}^{(1)}$

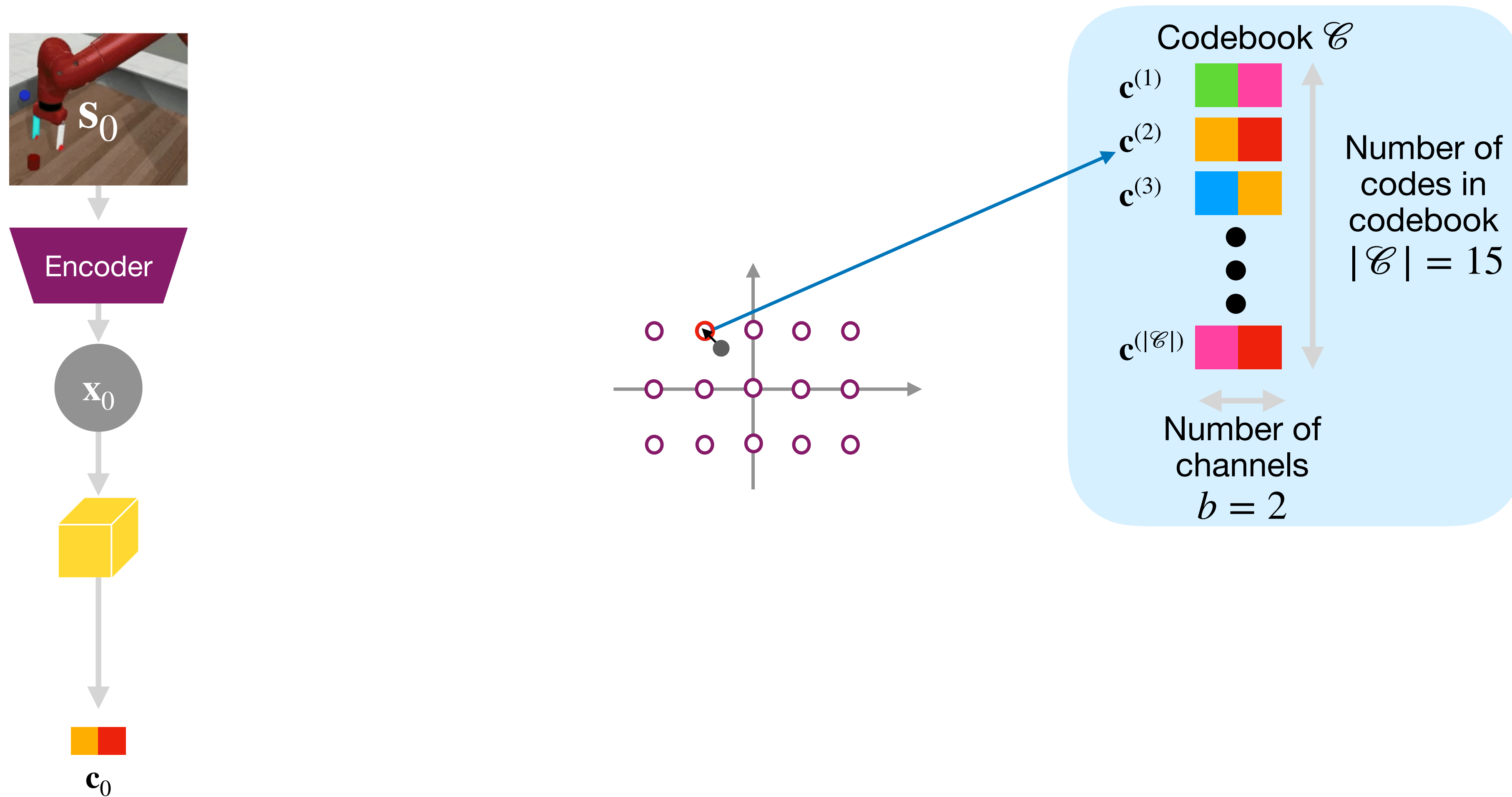$\mathbf{c}^{(2)}$

$\mathbf{c}^{(3)}$

$\mathbf{c}^{(|\mathscr{C}|)}$

Number of codes in codebook $|\mathscr{C}| = 15$

Number of channels $b = 2$

FCAI

fcai.fi

# DC-MPC: World Model Training



$\mathbf{s}_0$

Encoder

$\mathbf{x}_0$

$\mathbf{c}_0$

Dynamics

$p(\mathbf{c}_1 \mid \mathbf{c}_0, \mathbf{a}_0)$

$\mathbf{a}_0$

$\sim$

$\hat{\mathbf{c}}_1$

Dynamics

$p(\hat{\mathbf{c}}_2 \mid \hat{\mathbf{c}}_1, \mathbf{a}_1)$

$\mathbf{a}_1$

$\hat{r}_1$

Codebook $\mathscr{C}$

$\mathbf{c}^{(1)}$

$\mathbf{c}^{(2)}$

$\mathbf{c}^{(3)}$

$\mathbf{c}^{(|\mathscr{C}|)}$

Number of codes in codebook $|\mathscr{C}| = 15$

Number of channels $b = 2$

# DC-MPC: World Model Training



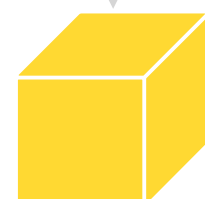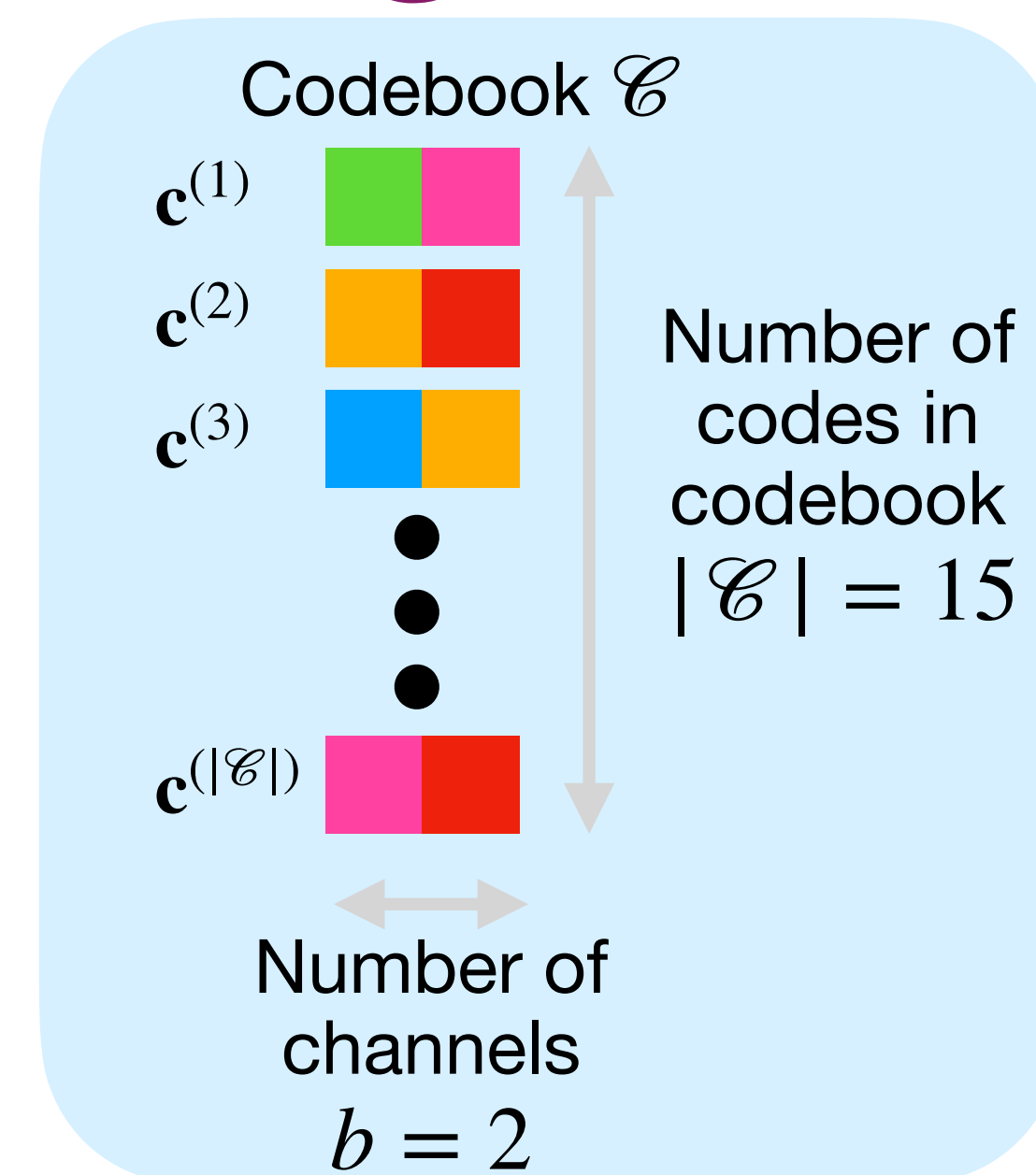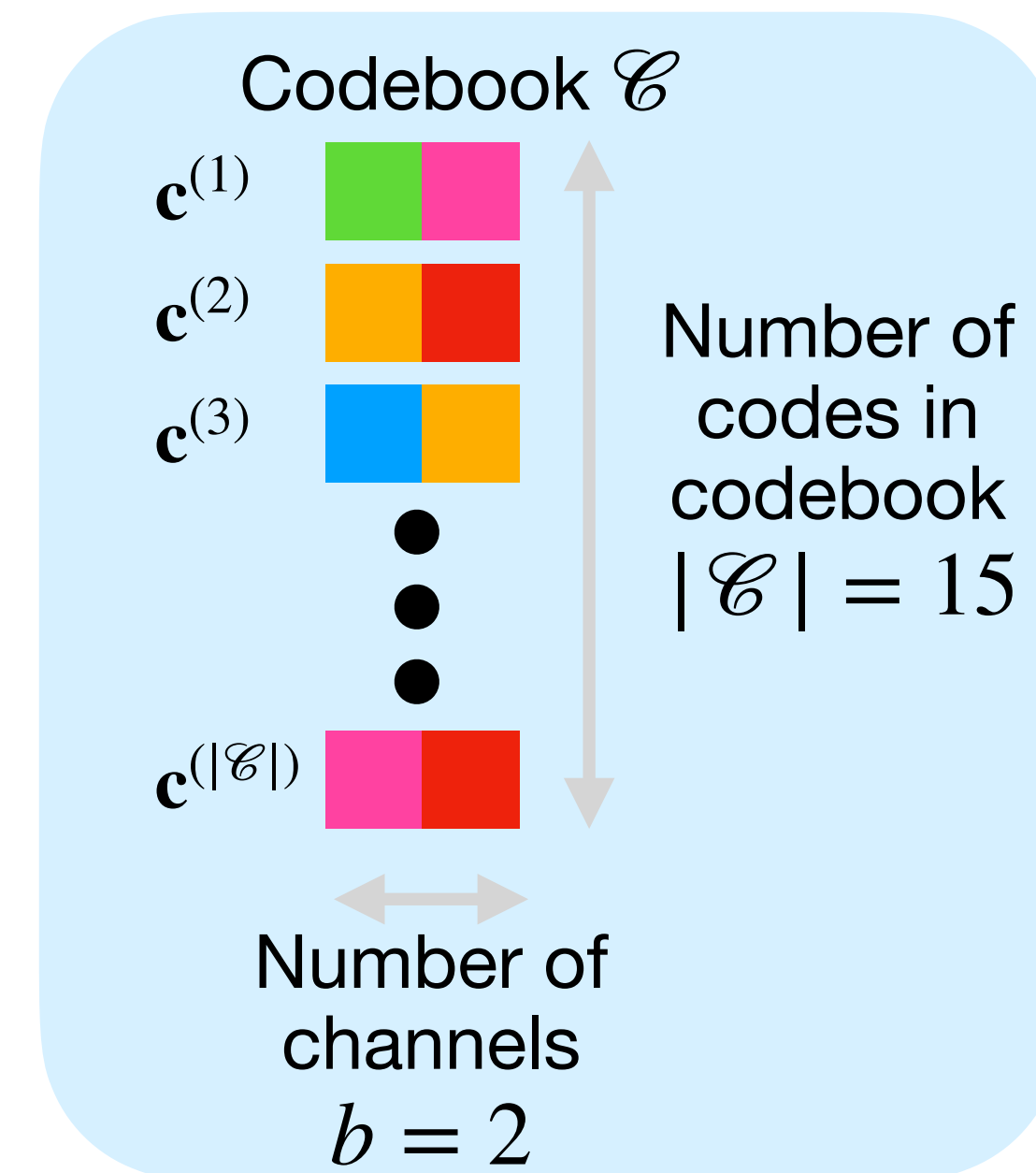Codebook $\mathscr{C}$

$\mathbf{c}^{(1)}$
$\mathbf{c}^{(2)}$
$\mathbf{c}^{(3)}$

$\mathbf{c}^{(|\mathscr{C}|)}$

Number of codes in codebook $|\mathscr{C}| = 15$

Number of channels $b = 2$

$\mathbf{s}_0$

Encoder

$\mathbf{x}_0$

$\mathbf{c}_0$

$\mathbf{a}_0$

Dynamics

$p(\mathbf{c}_1 \mid \mathbf{c}_0, \mathbf{a}_0)$

$\sim$

$\hat{\mathbf{c}}_1$

$\mathbf{a}_1$

Dynamics

$p(\hat{\mathbf{c}}_2 \mid \hat{\mathbf{c}}_1, \mathbf{a}_1)$

$\hat{r}_1$

$\hat{r}_2$

FCAI

fcai.fi

# DC-MPC: World Model Training



$\mathbf{s}_0$

Encoder

$\mathbf{x}_0$

$\mathbf{c}_0$

$\mathbf{a}_0$

Dynamics

$p(\mathbf{c}_1 \mid \mathbf{c}_0, \mathbf{a}_0)$

$\hat{r}_1$

$\sim$

$\hat{\mathbf{c}}_1$

$\mathbf{a}_1$

Dynamics

$p(\hat{\mathbf{c}}_2 \mid \hat{\mathbf{c}}_1, \mathbf{a}_1)$

$\hat{r}_2$

Codebook $\mathscr{C}$

$\mathbf{c}^{(1)}$

$\mathbf{c}^{(2)}$

$\mathbf{c}^{(3)}$

$\mathbf{c}^{(|\mathscr{C}|)}$

Number of codes in codebook $|\mathscr{C}| = 15$

Number of channels $b = 2$

FCAI

fcai.fi

# DC-MPC: World Model Training



$\mathbf{s}_0$

Encoder

$\mathbf{x}_0$

Codebook $\mathscr{C}$

$\mathbf{c}^{(1)}$

$\mathbf{c}^{(2)}$

$\mathbf{c}^{(3)}$

$\mathbf{c}^{(|\mathscr{C}|)}$

Number of codes in codebook $|\mathscr{C}| = 15$

Number of channels $b = 2$

$\mathbf{c}_0$

$\mathbf{a}_0$

Dynamics

$p(\mathbf{c}_1 \mid \mathbf{c}_0, \mathbf{a}_0)$

$\hat{r}_1$

$\sim$

$\hat{\mathbf{c}}_1$

$\mathbf{a}_1$

Dynamics

$p(\hat{\mathbf{c}}_2 \mid \hat{\mathbf{c}}_1, \mathbf{a}_1)$

$\hat{r}_2$

$\hat{\mathbf{c}}_{H-1}$

$\mathbf{a}_{H-1}$

Dynamics
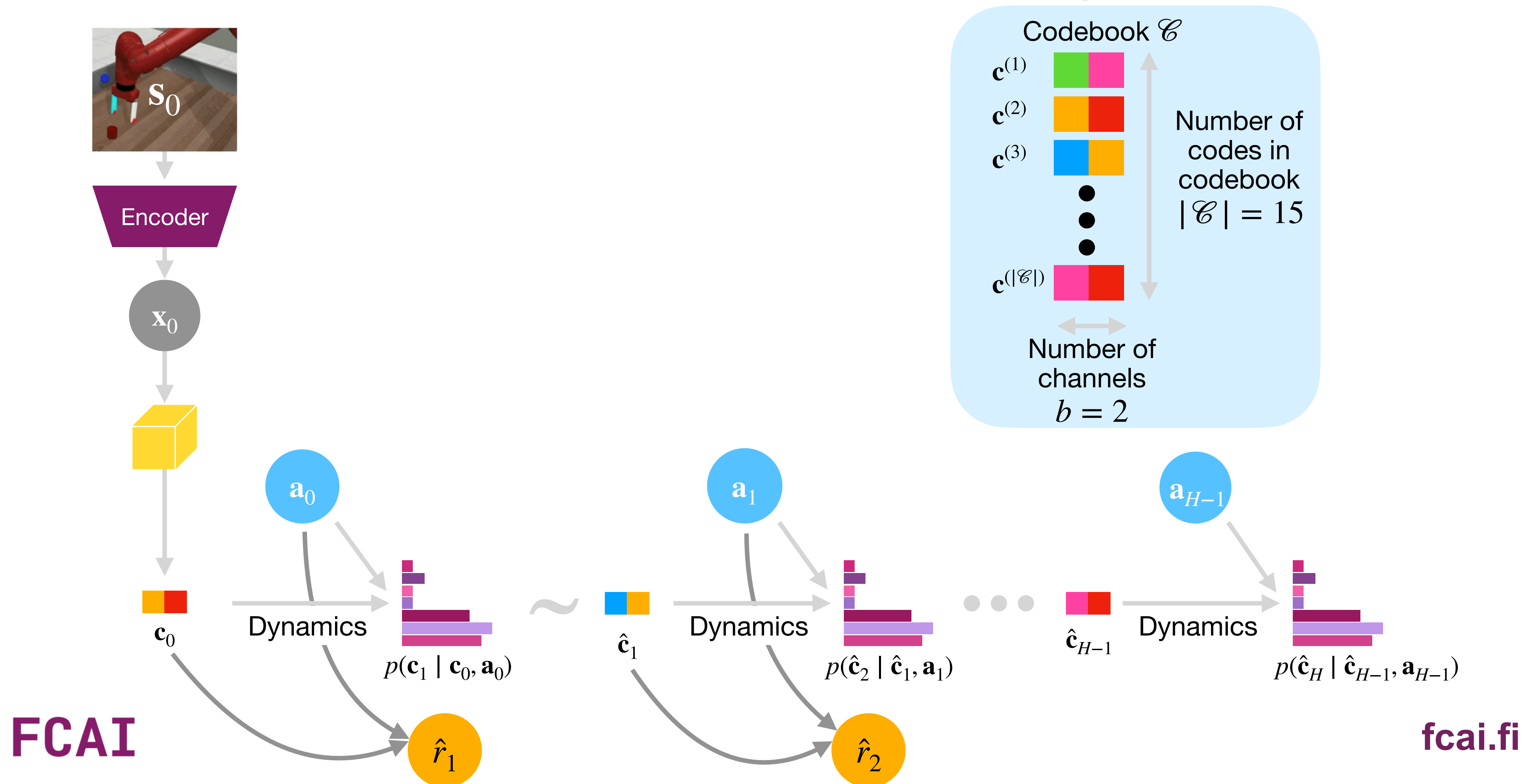
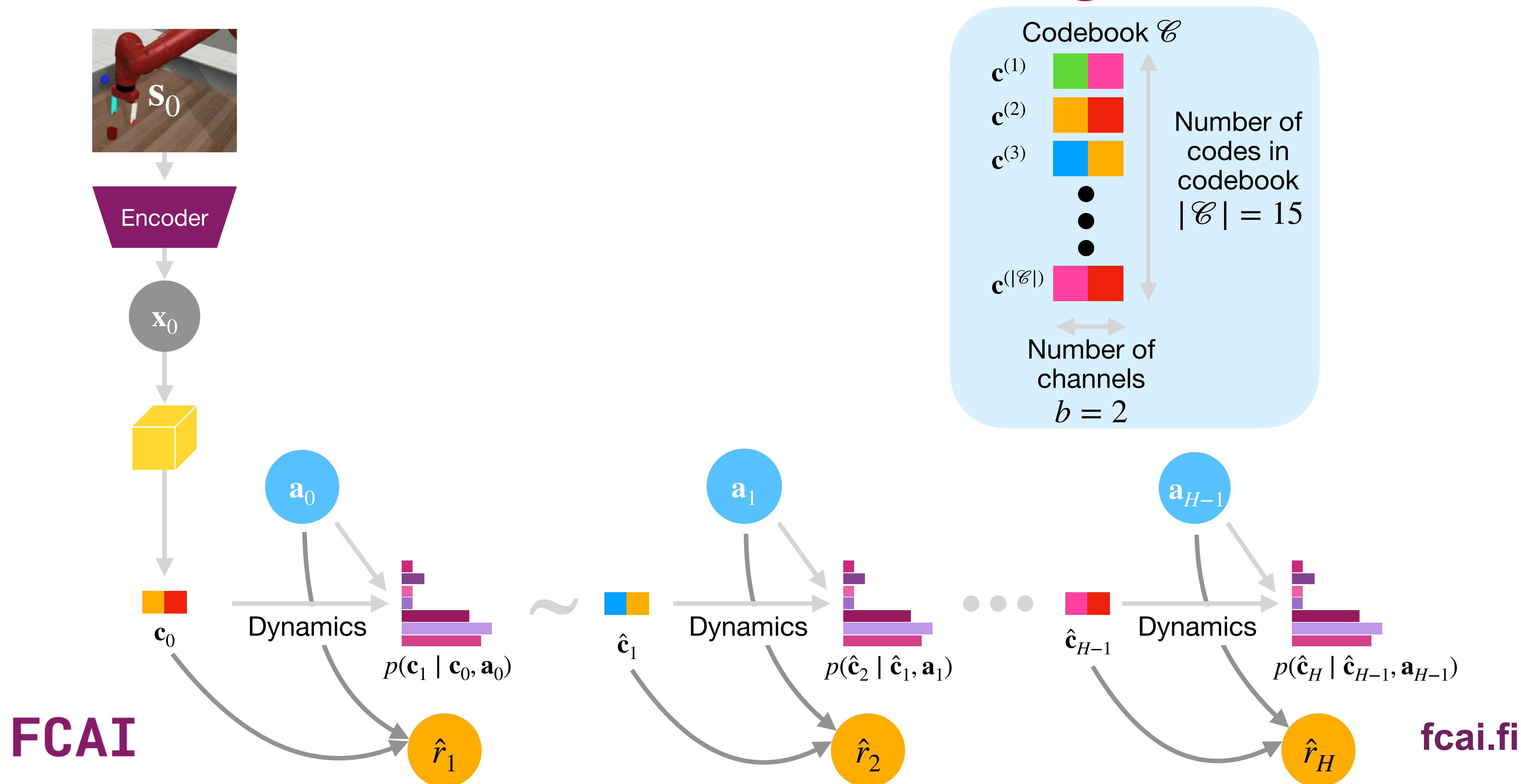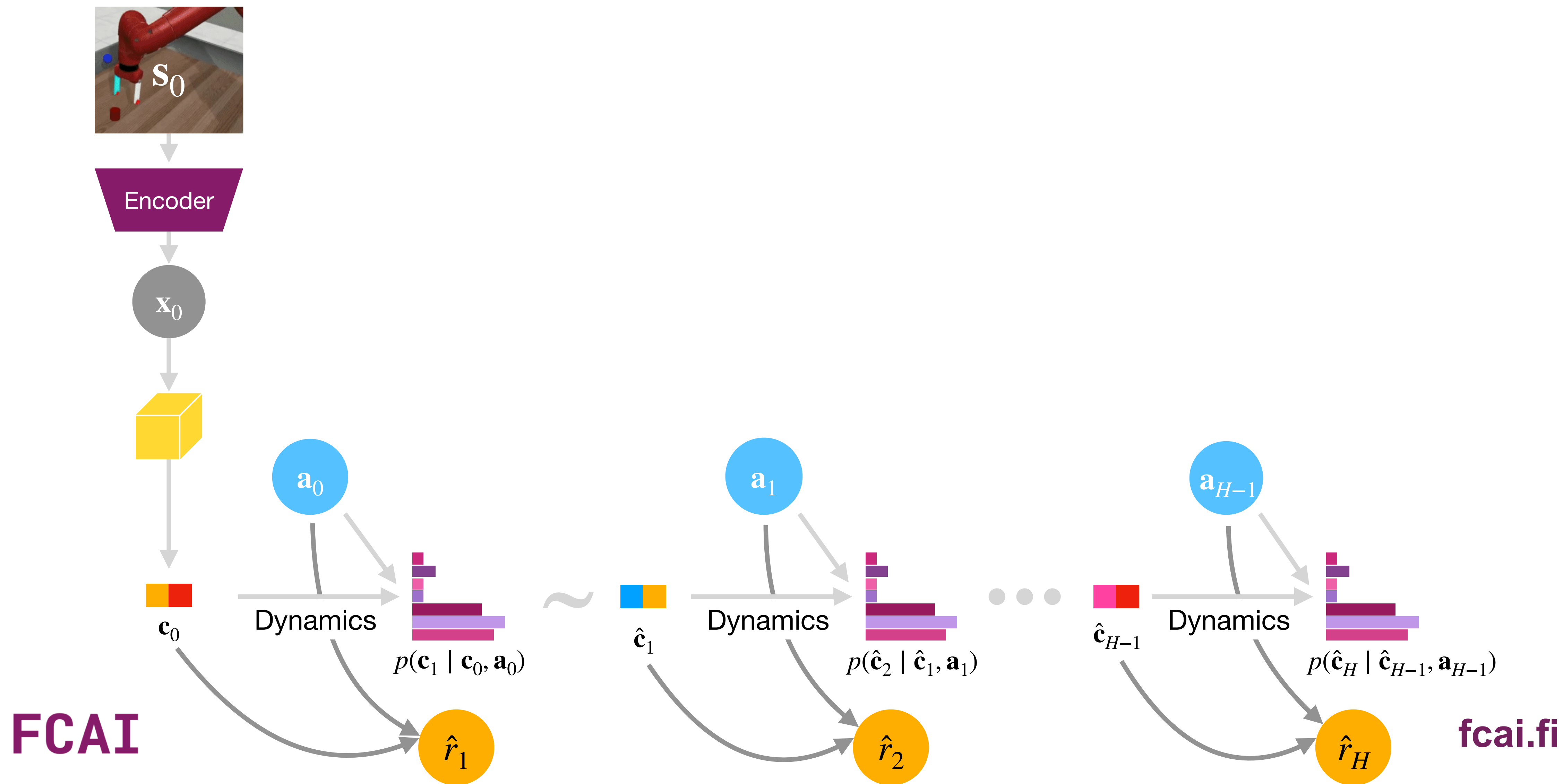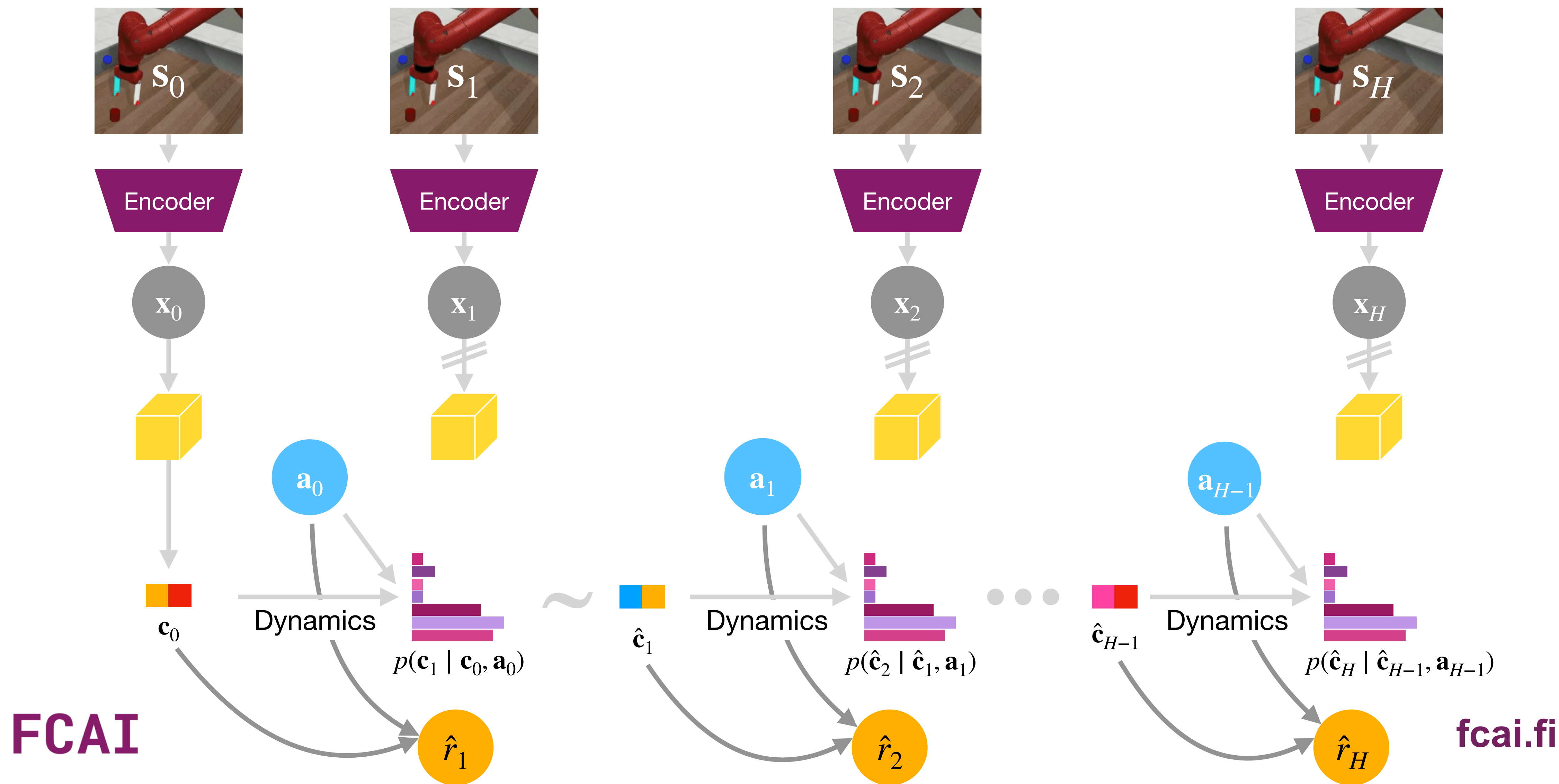$p(\hat{\mathbf{c}}_H \mid \hat{\mathbf{c}}_{H-1}, \mathbf{a}_{H-1})$

FCAI

fcai.fi

# DC-MPC: World Model Training

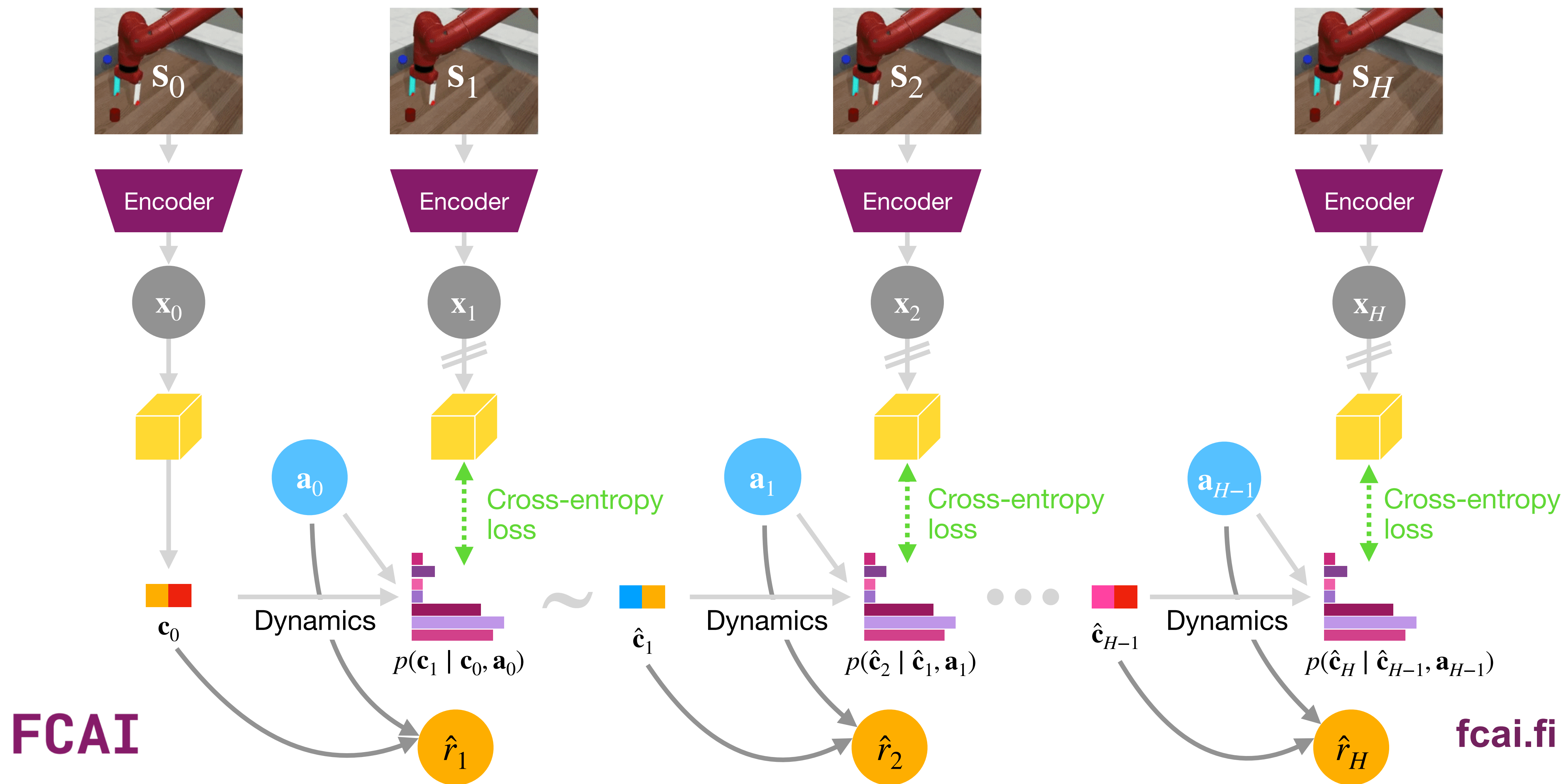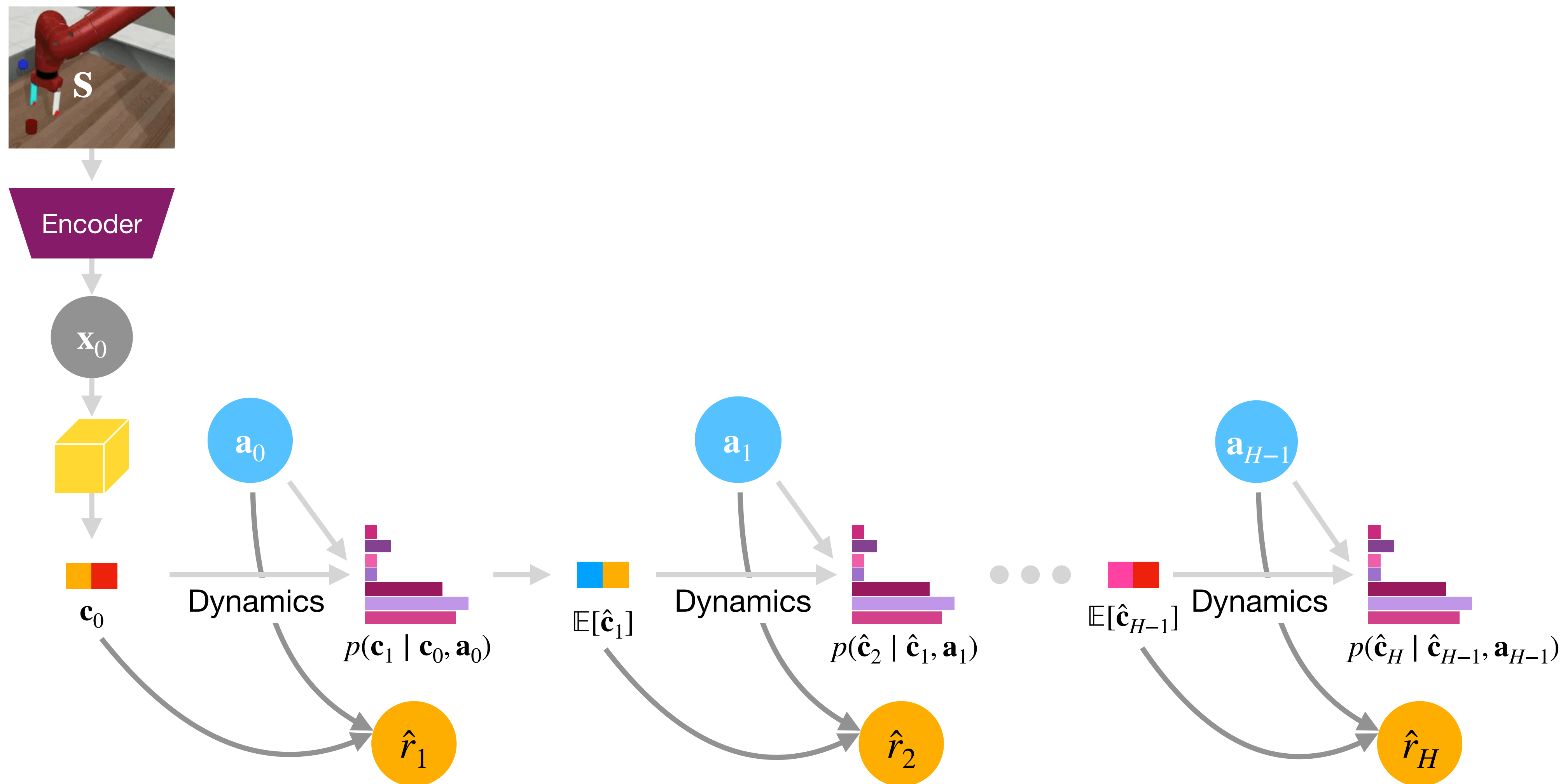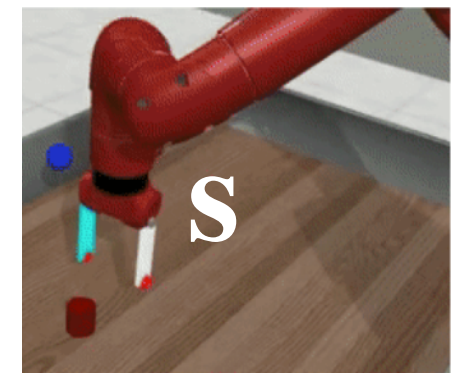# DC-MPC: World Model Training

# DC-MPC: World Model Training



$\mathbf{s}_0$ $\mathbf{s}_1$ $\mathbf{s}_2$ $\mathbf{s}_H$

Encoder Encoder Encoder Encoder

$\mathbf{x}_0$ $\mathbf{x}_1$ $\mathbf{x}_2$ $\mathbf{x}_H$

$\mathbf{a}_0$ $\mathbf{a}_1$ $\mathbf{a}_{H-1}$

$\mathbf{c}_0$ Dynamics $p(\mathbf{c}_1 \mid \mathbf{c}_0, \mathbf{a}_0)$ $\sim$ $\hat{\mathbf{c}}_1$ Dynamics $p(\hat{\mathbf{c}}_2 \mid \hat{\mathbf{c}}_1, \mathbf{a}_1)$ $\hat{\mathbf{c}}_{H-1}$ Dynamics $p(\hat{\mathbf{c}}_H \mid \hat{\mathbf{c}}_{H-1}, \mathbf{a}_{H-1})$

$\hat{r}_1$ $\hat{r}_2$ $\hat{r}_H$

FCAI fcai.fi

# DC-MPC: World Model Training

# DC-MPC: Decision-time Planning



$\mathbf{S}$

Encoder

$\mathbf{x}_0$

$\mathbf{a}_0$     $\mathbf{a}_1$     $\mathbf{a}_{H-1}$

$\mathbf{c}_0$    Dynamics    $p(\mathbf{c}_1 \mid \mathbf{c}_0, \mathbf{a}_0)$    $\mathbb{E}[\hat{\mathbf{c}}_1]$    Dynamics    $p(\hat{\mathbf{c}}_2 \mid \hat{\mathbf{c}}_1, \mathbf{a}_1)$    $\mathbb{E}[\hat{\mathbf{c}}_{H-1}]$    Dynamics    $p(\hat{\mathbf{c}}_H \mid \hat{\mathbf{c}}_{H-1}, \mathbf{a}_{H-1})$

$\hat{r}_1$     $\hat{r}_2$     $\hat{r}_H$

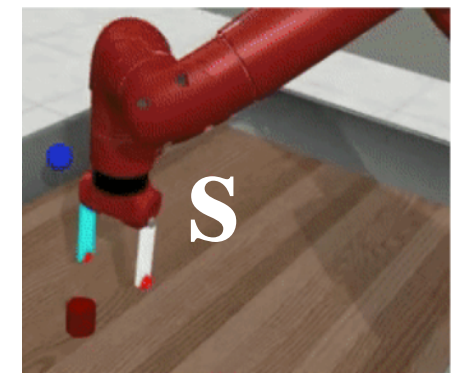FCAI          fcai.fi

# DC-MPC: Decision-time Planning



$$J(\mathbf{a}_{0:H}, \mathbf{s}) = \gamma^H Q_\psi(\hat{\mathbf{c}}_H, \mathbf{a}_H) + \sum_{h=0}^{H-1} \gamma^h R_\xi(\hat{\mathbf{c}}_h, \mathbf{a}_h)$$
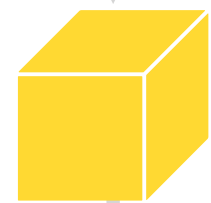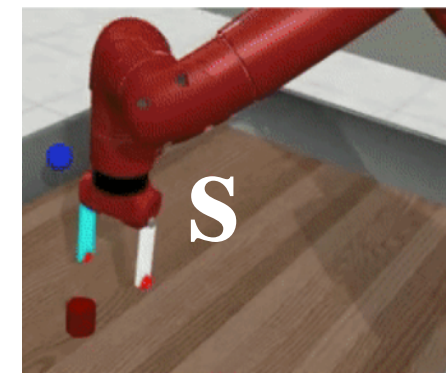
Encoder

$\mathbf{x}_0$

$\mathbf{a}_0$

$\mathbf{a}_1$

$\mathbf{a}_{H-1}$

$\mathbf{c}_0$

Dynamics

$p(\mathbf{c}_1 \mid \mathbf{c}_0, \mathbf{a}_0)$

$\mathbb{E}[\hat{\mathbf{c}}_1]$

Dynamics

$p(\hat{\mathbf{c}}_2 \mid \hat{\mathbf{c}}_1, \mathbf{a}_1)$

$\mathbb{E}[\hat{\mathbf{c}}_{H-1}]$

Dynamics

$p(\hat{\mathbf{c}}_H \mid \hat{\mathbf{c}}_{H-1}, \mathbf{a}_{H-1})$

$\hat{r}_1$

$\hat{r}_2$

$\hat{r}_H$

FCAI

fcai.fi

# DC-MPC: Decision-time Planning



**Reward func.**

$$J(\mathbf{a}_{0:H}, \mathbf{s}) = \gamma^H Q_\psi(\hat{\mathbf{c}}_H, \mathbf{a}_H) + \boxed{\sum_{h=0}^{H-1} \gamma^h R_\xi(\hat{\mathbf{c}}_h, \mathbf{a}_h)}$$

FCAI

fcai.fi

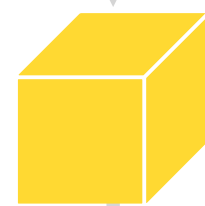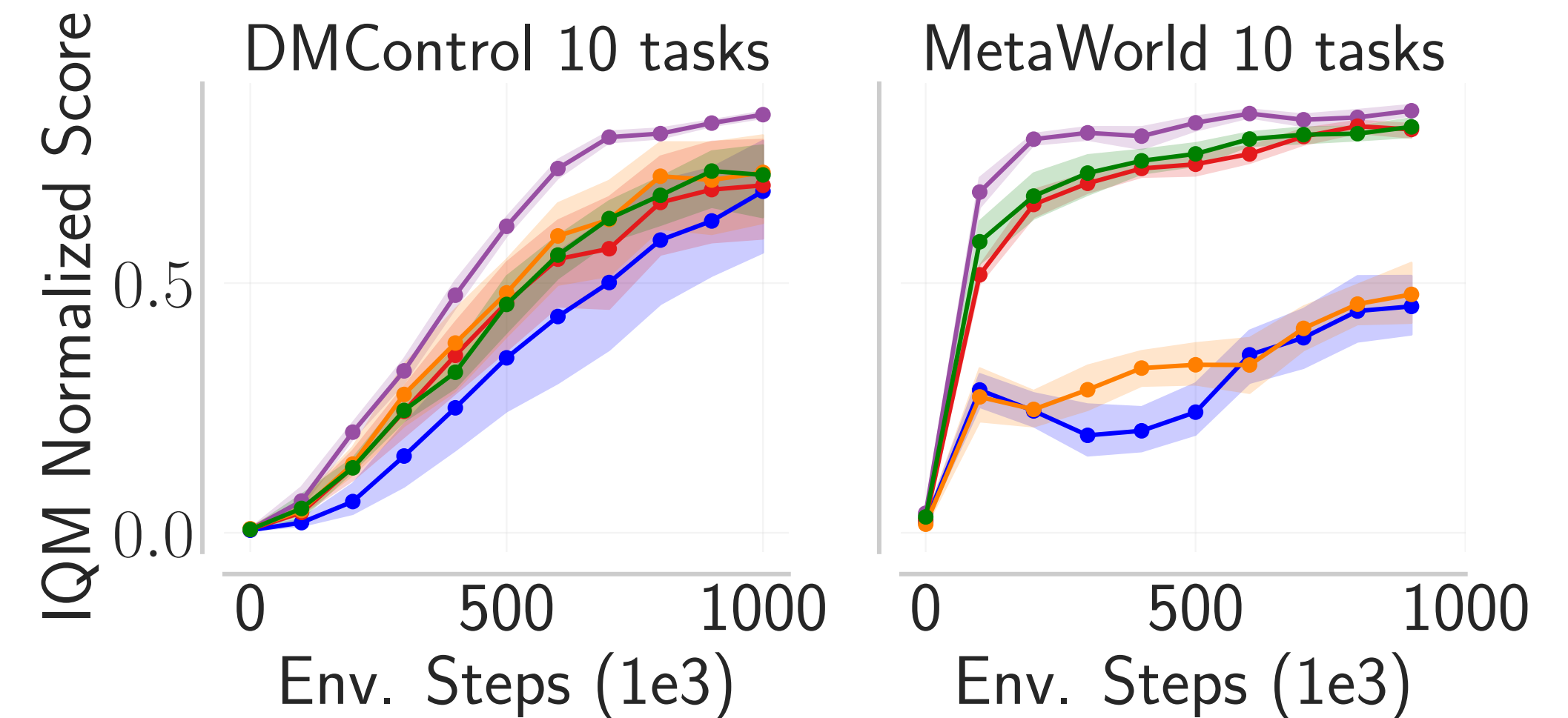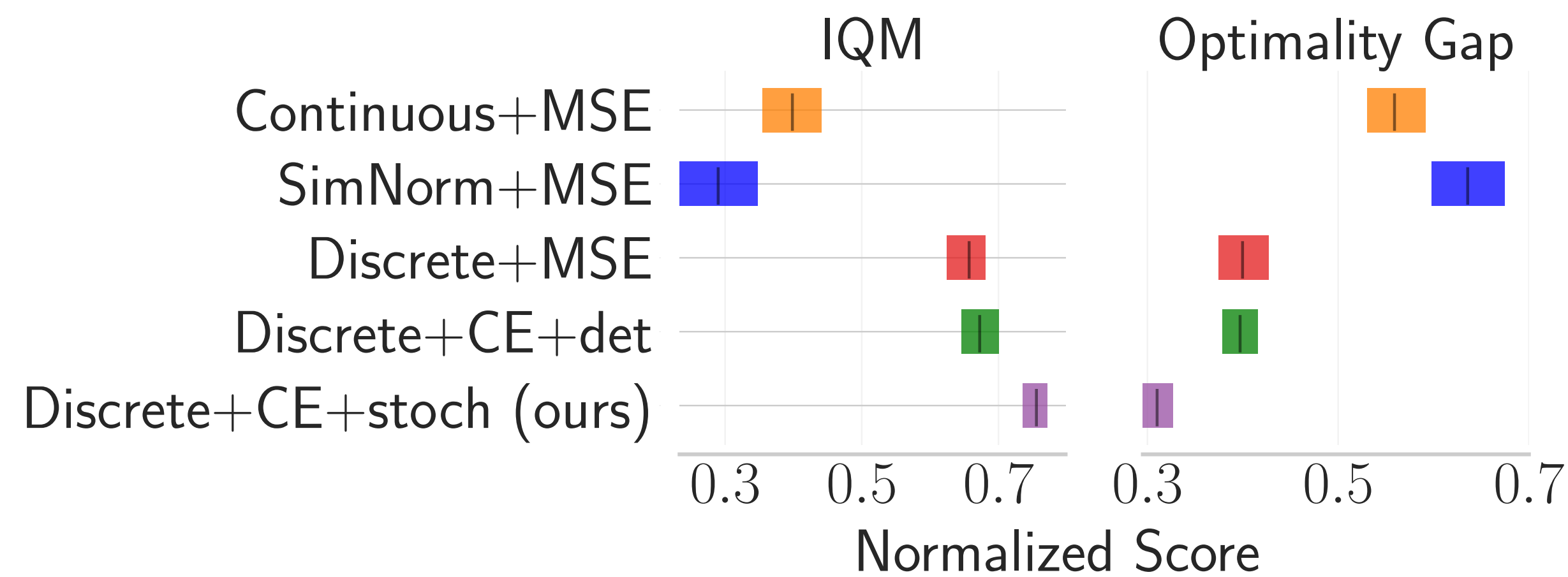# DC-MPC: Decision-time Planning



$$J(\mathbf{a}_{0:H}, \mathbf{s}) = \underbrace{\gamma^H Q_\psi(\hat{\mathbf{c}}_H, \mathbf{a}_H)}_{\text{Bootstrap with action-value}} + \underbrace{\sum_{h=0}^{H-1} \gamma^h R_\xi(\hat{\mathbf{c}}_h, \mathbf{a}_h)}_{\text{Reward func.}}$$

**Bootstrap with action-value**

**Reward func.**

Encoder

$\mathbf{x}_0$

$\mathbf{c}_0$

$\mathbf{a}_0$

$\mathbf{a}_1$

$\mathbf{a}_{H-1}$

Dynamics

$p(\mathbf{c}_1 \mid \mathbf{c}_0, \mathbf{a}_0)$

$\mathbb{E}[\hat{\mathbf{c}}_1]$

$p(\hat{\mathbf{c}}_2 \mid \hat{\mathbf{c}}_1, \mathbf{a}_1)$

$\mathbb{E}[\hat{\mathbf{c}}_{H-1}]$

$p(\hat{\mathbf{c}}_H \mid \hat{\mathbf{c}}_{H-1}, \mathbf{a}_{H-1})$

$\hat{r}_1$

$\hat{r}_2$
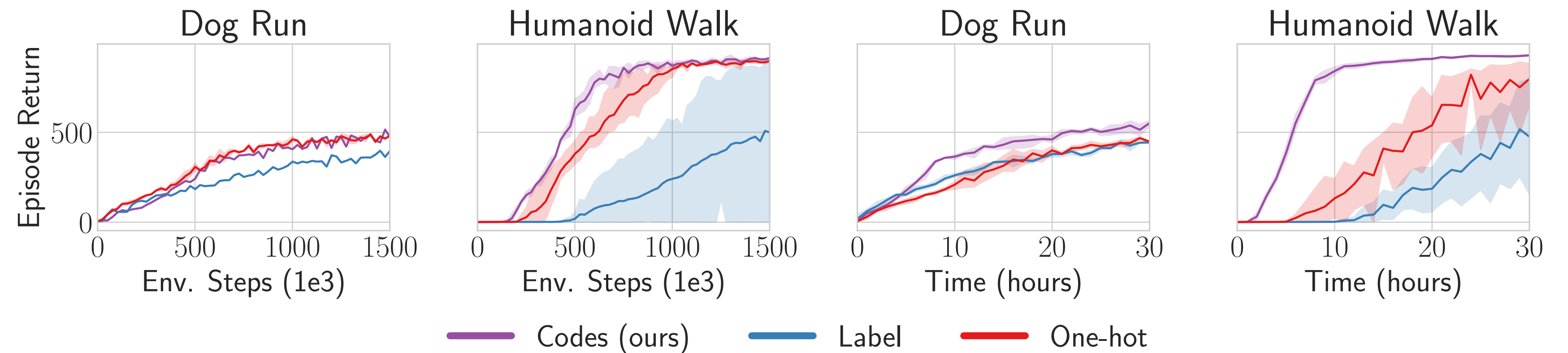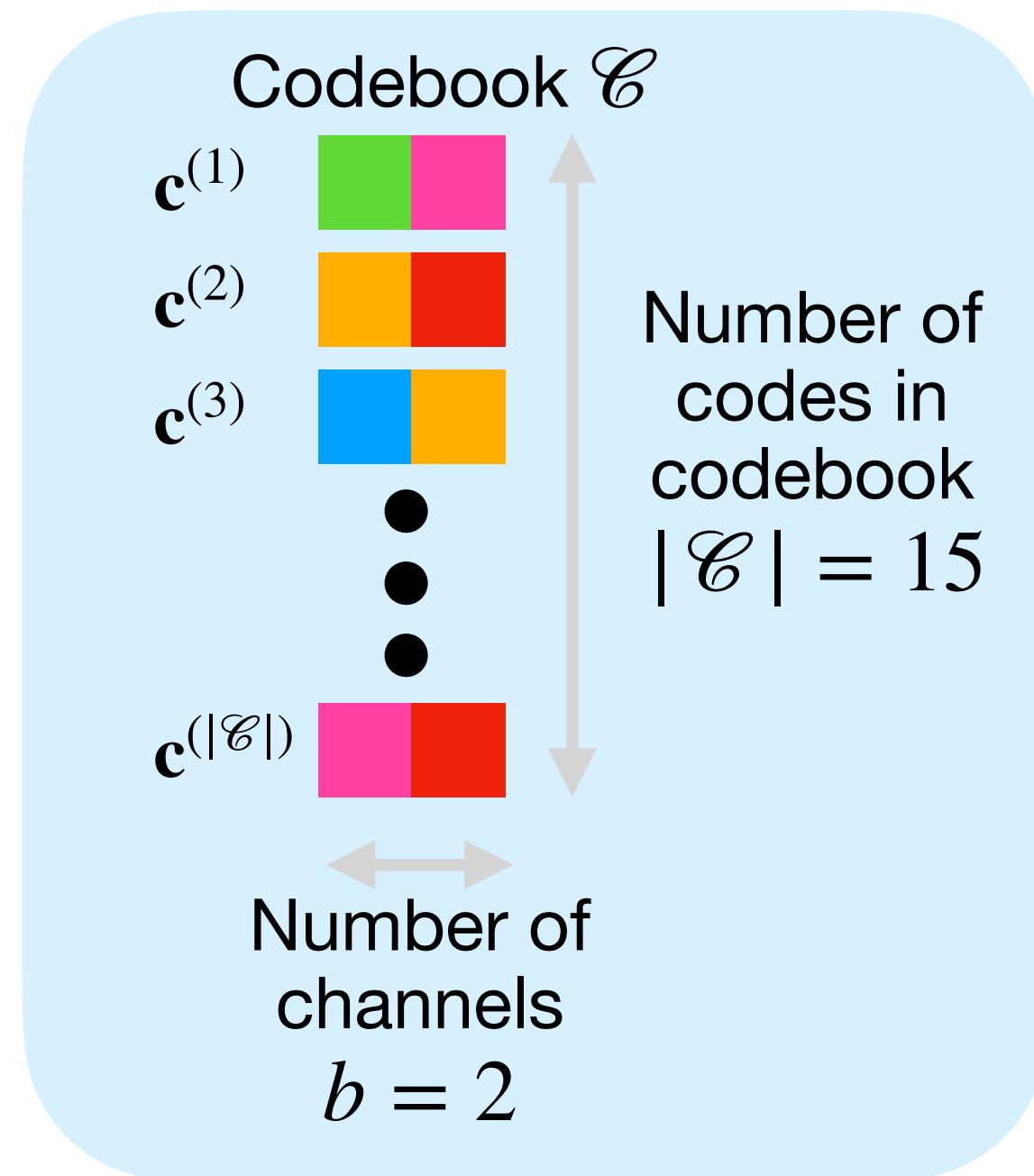
$\hat{r}_H$

FCAI

fcai.fi

# Why Does DC-MPC Work So Well?
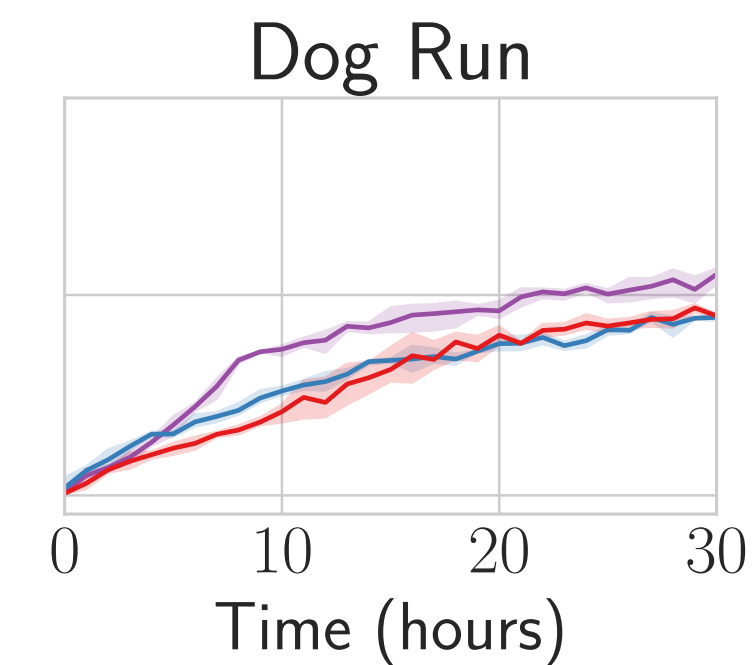## Combination of Discrete Codebook and Stochastic Dynamics
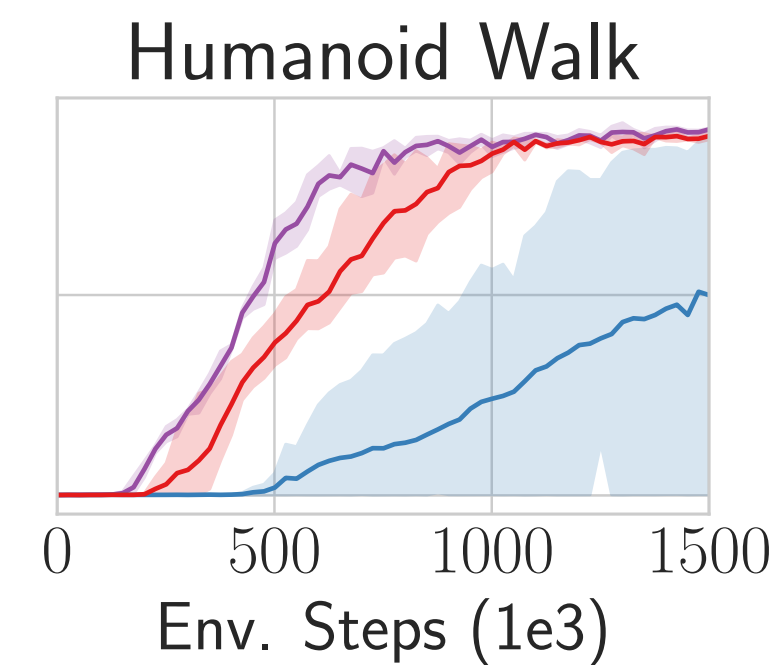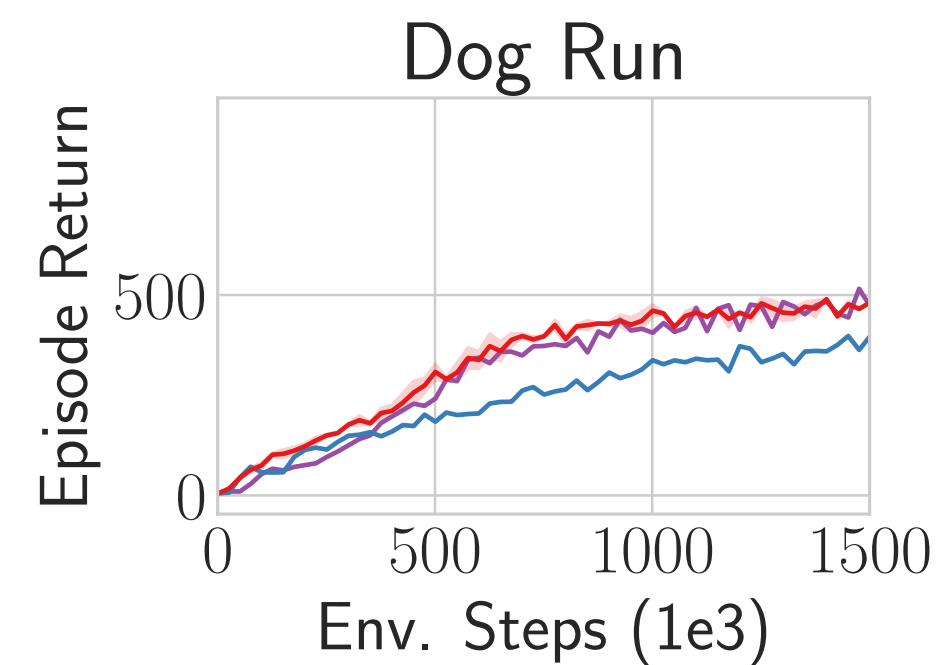
**FCAI**

**fcai.fi**

# Comparison of Different Discrete Encodings
## Codebook > One-hot > Label

# Comparison of Different Discrete Encodings
## Codebook > One-hot > Label



Codebook $\mathscr{C}$

$\mathbf{c}^{(1)}$

$\mathbf{c}^{(2)}$

$\mathbf{c}^{(3)}$

$\mathbf{c}^{(|\mathscr{C}|)}$

Number of codes in codebook $|\mathscr{C}| = 15$

Number of channels $b = 2$



Dog Run — Humanoid Walk — Dog Run — Humanoid Walk

Episode Return — 500 — 0

Env. Steps (1e3) — Env. Steps (1e3) — Time (hours) — Time (hours)
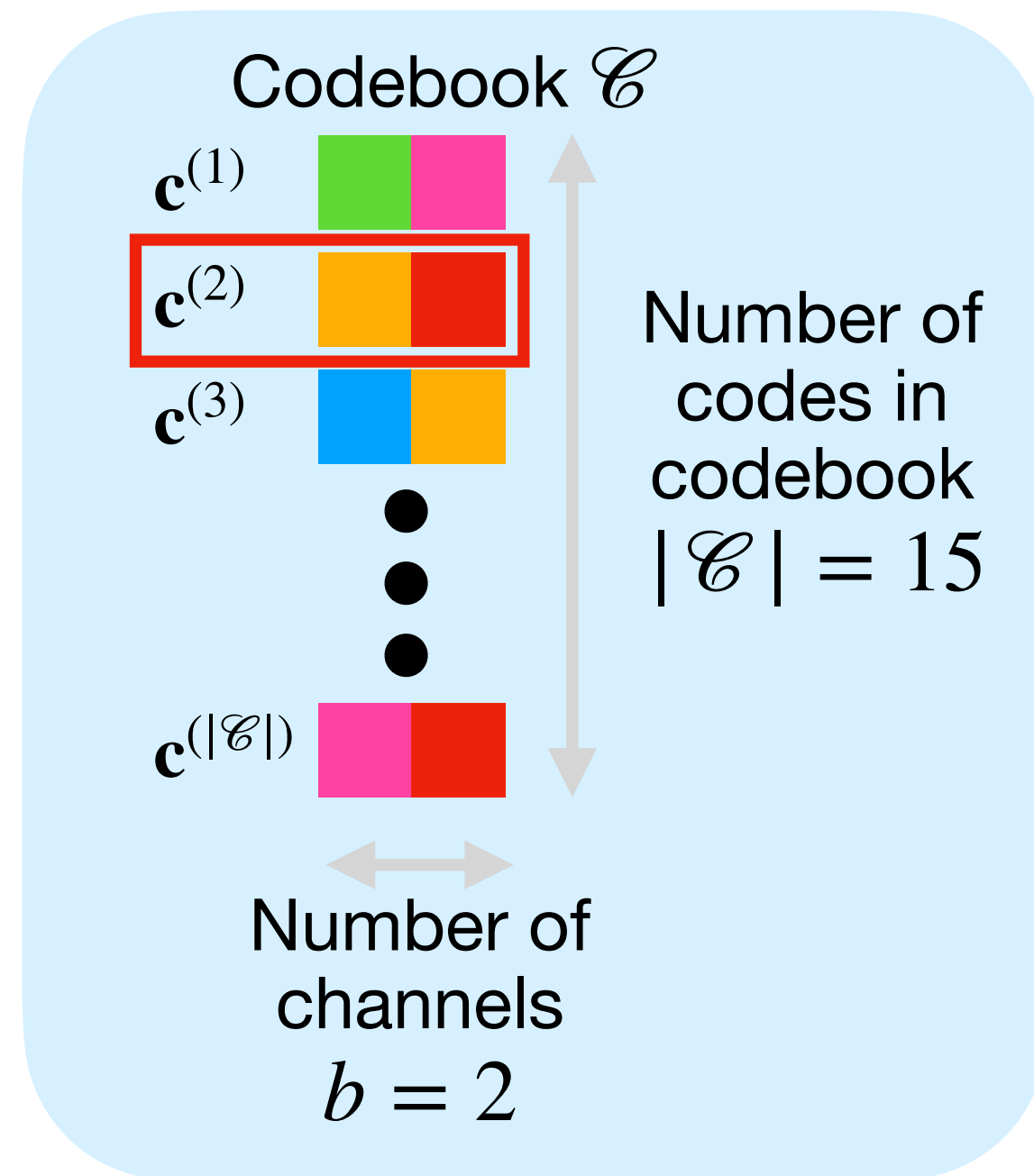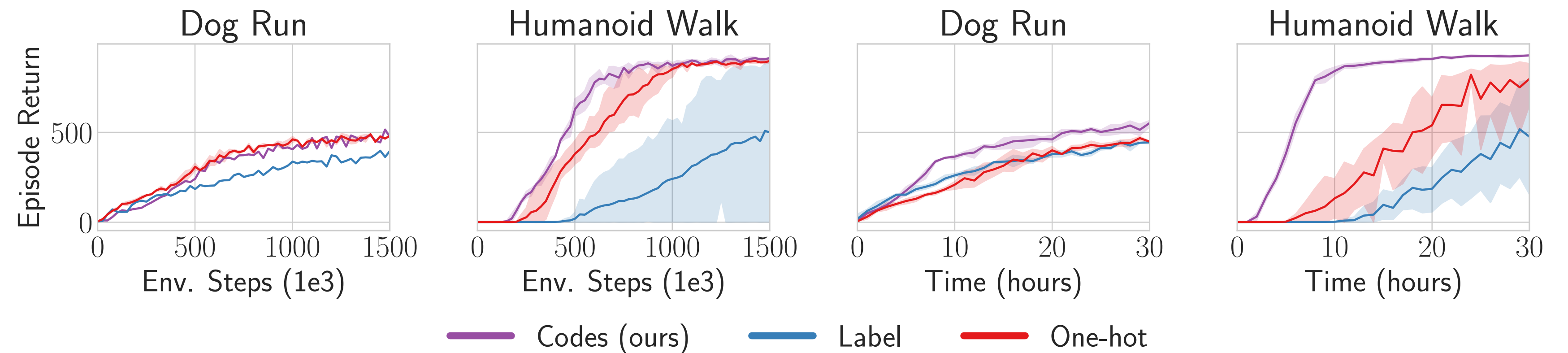
Codes (ours) — Label — One-hot

# Comparison of Different Discrete Encodings
## Codebook > One-hot > Label



$$\mathbf{e}_{\text{code}} = \mathbf{c}^{(2)} = \{-0.5, 1\}$$

Codebook $\mathscr{C}$

$\mathbf{c}^{(1)}$

$\mathbf{c}^{(2)}$

$\mathbf{c}^{(3)}$

$\mathbf{c}^{(|\mathscr{C}|)}$

Number of codes in codebook $|\mathscr{C}| = 15$

Number of channels $b = 2$

Dog Run

Humanoid Walk

Dog Run

Humanoid Walk

Episode Return

Env. Steps (1e3)

Env. Steps (1e3)

Time (hours)

Time (hours)

— Codes (ours)    — Label    — One-hot

FCAI

**fcai.fi**

# Comparison of Different Discrete Encodings
## Codebook > One-hot > Label



Codebook $\mathscr{C}$

$\mathbf{c}^{(1)}$

$\mathbf{c}^{(2)}$

$\mathbf{c}^{(3)}$

$\mathbf{c}^{(|\mathscr{C}|)}$

Number of codes in codebook $|\mathscr{C}| = 15$

Number of channels $b = 2$

$\mathbf{e}_{code} = \mathbf{c}^{(2)} = \{-0.5, 1\}$

$\mathbf{e}_{label} = 2$



Dog Run — Humanoid Walk — Dog Run — Humanoid Walk

Episode Return

Env. Steps (1e3) — Env. Steps (1e3) — Time (hours) — Time (hours)

Codes (ours) — Label — One-hot

# Comparison of Different Discrete Encodings
## Codebook > One-hot > Label



Codebook $\mathscr{C}$

$\mathbf{c}^{(1)}$

$\mathbf{c}^{(2)}$

$\mathbf{c}^{(3)}$

Number of codes in codebook $|\mathscr{C}| = 15$

$\mathbf{c}^{(|\mathscr{C}|)}$

Number of channels $b = 2$
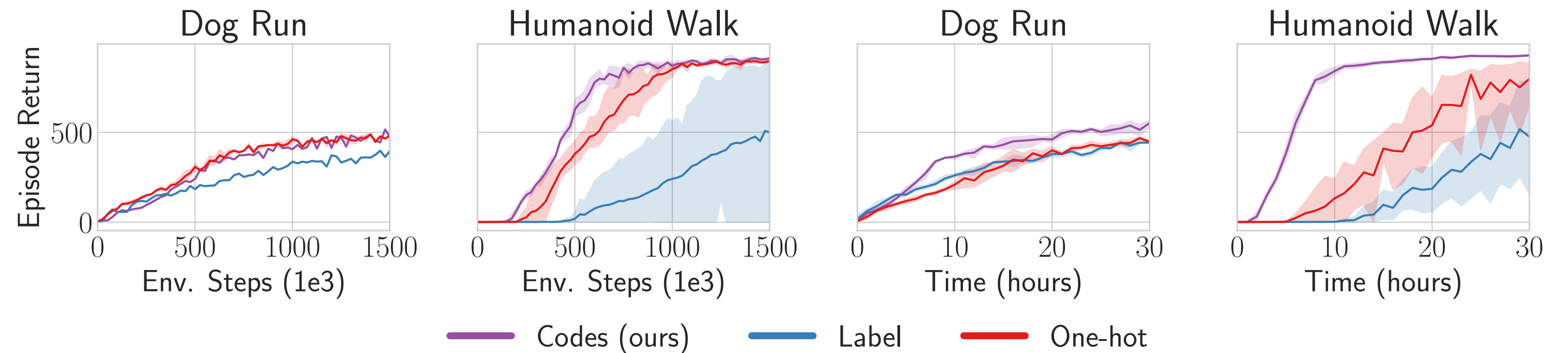
$$\mathbf{e}_{\text{code}} = \mathbf{c}^{(2)} = \{-0.5, 1\}$$

$$\mathbf{e}_{\text{label}} = 2$$

$$\mathbf{e}_{\text{one-hot}} = \{0,1,0,0,0,0,0,0,0,0,0,0,0,0,0\}$$



Codes (ours)   Label   One-hot

FCAI                                                              fcai.fi

**Email:**     ascannel@ed.ac.uk

**Website:**  www.aidanscannell.com/dcmpc

**Poster 28506:**

Wednesday 23rd April 10 am - 12.30 pm (GMT+8)

**FCAI**                                        **fcai.fi**