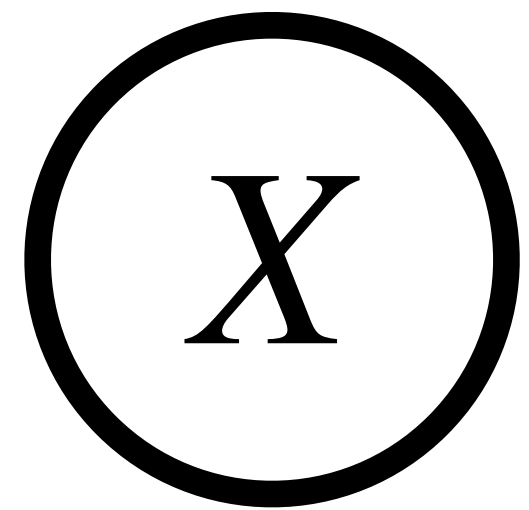


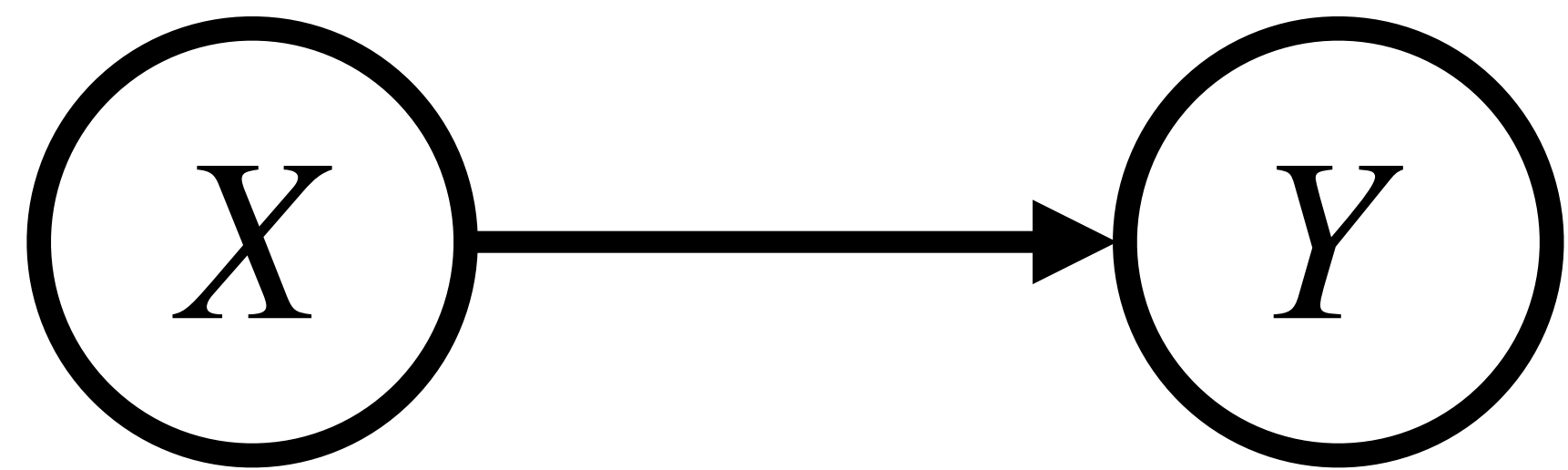
Causal Confounds in Sequential Decision Making

Gokul Swamy

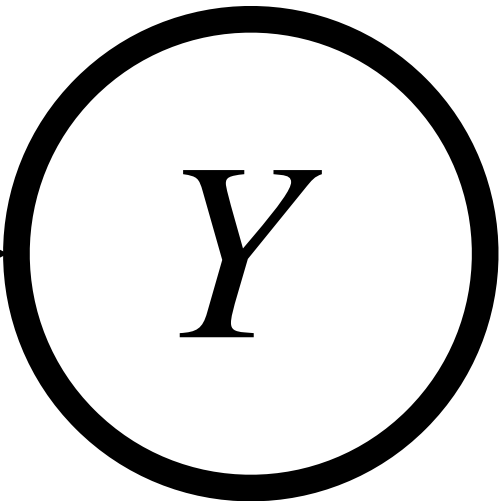
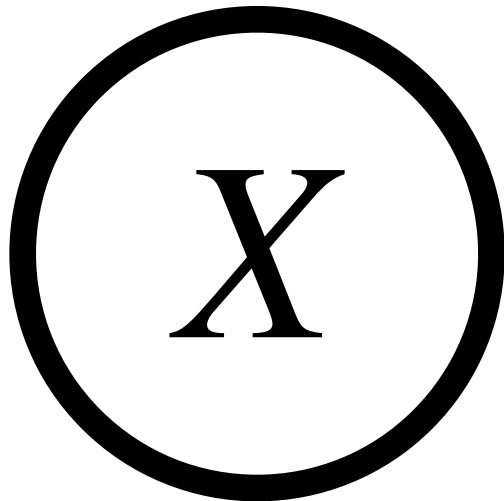


(joint work with Sanjiban Choudhury, Drew Bagnell, Steven Wu)





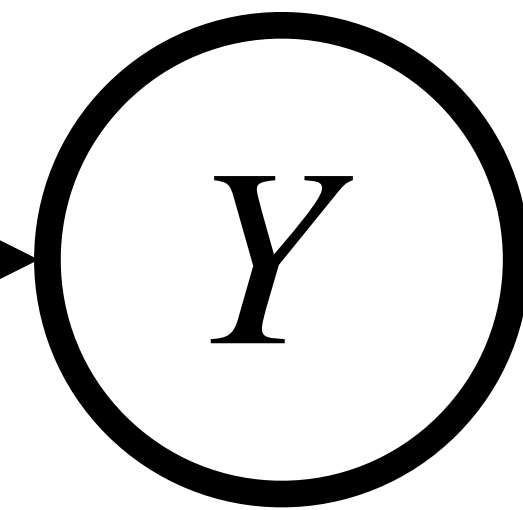
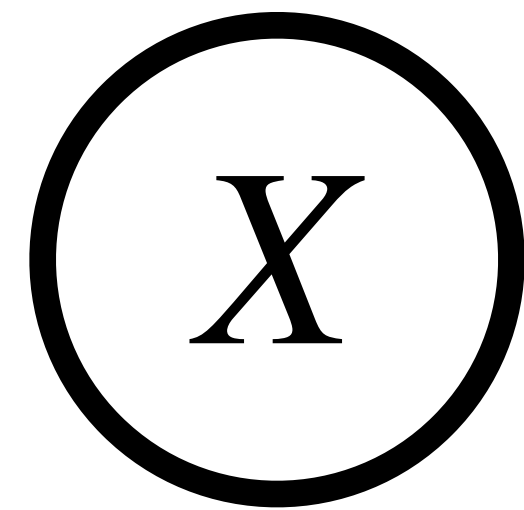
*Swimsuit
Sales*



*Ice-Cream
Sales*

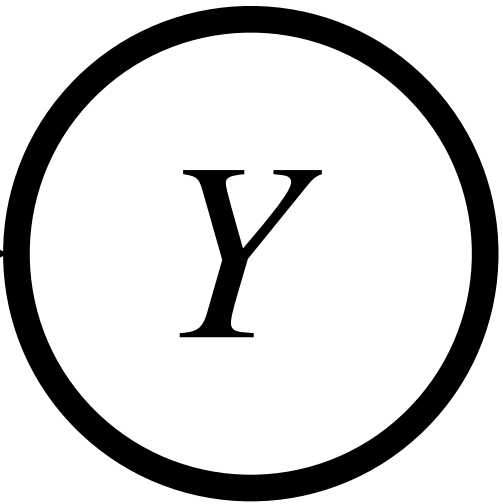
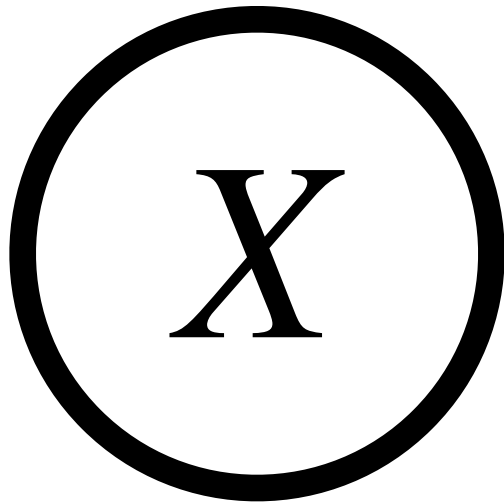


*Swimsuit
Sales*



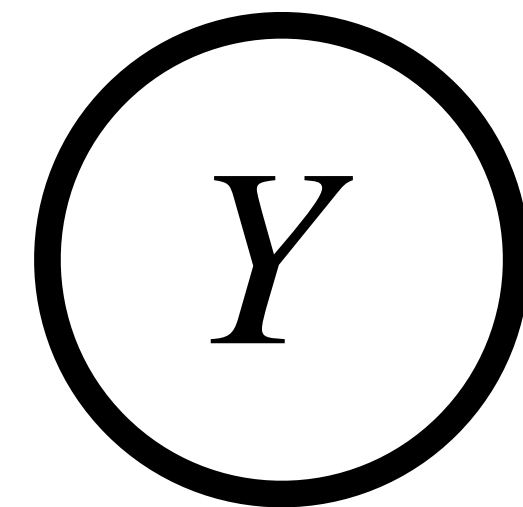
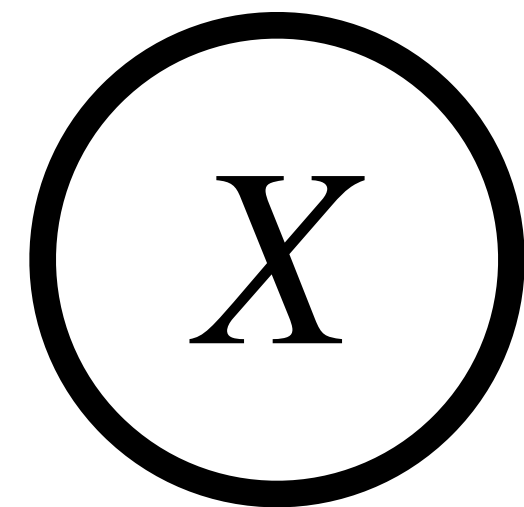
*Ice-Cream
Sales*

*Swimsuit
Sales*

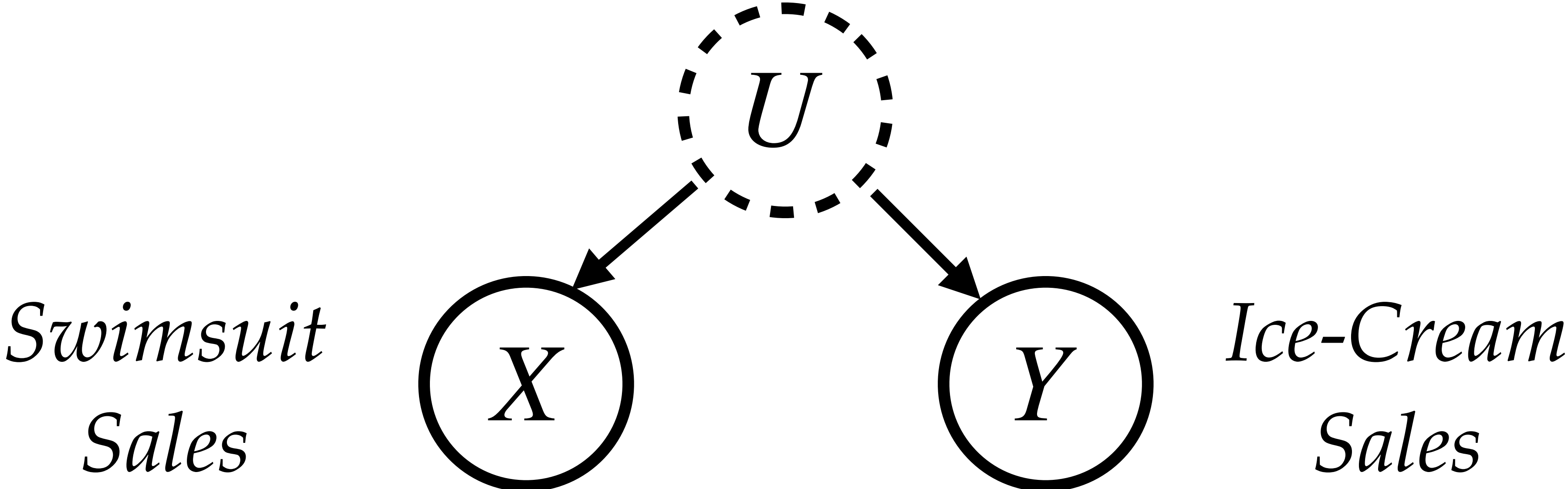


*Ice-Cream
Sales*

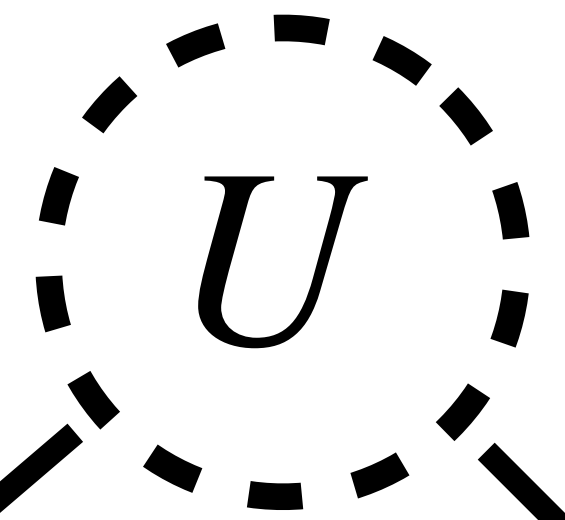
*Swimsuit
Sales*



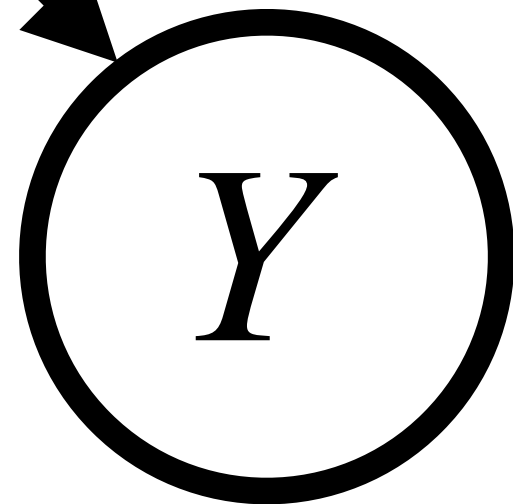
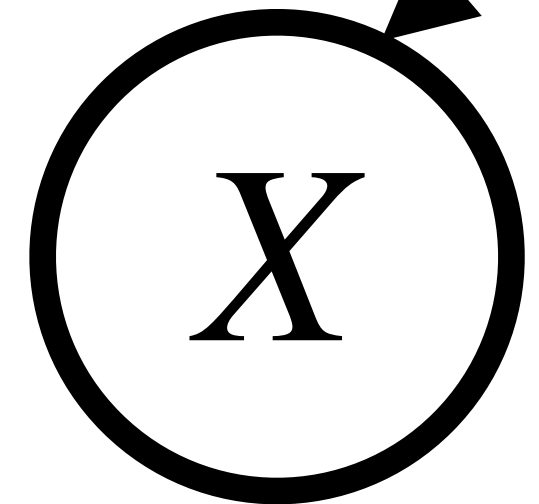
*Ice-Cream
Sales*



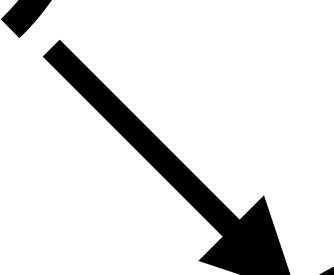
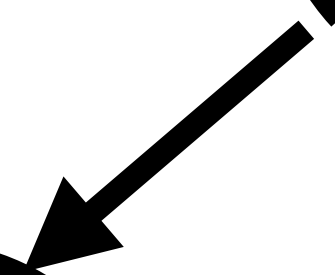
Temperature



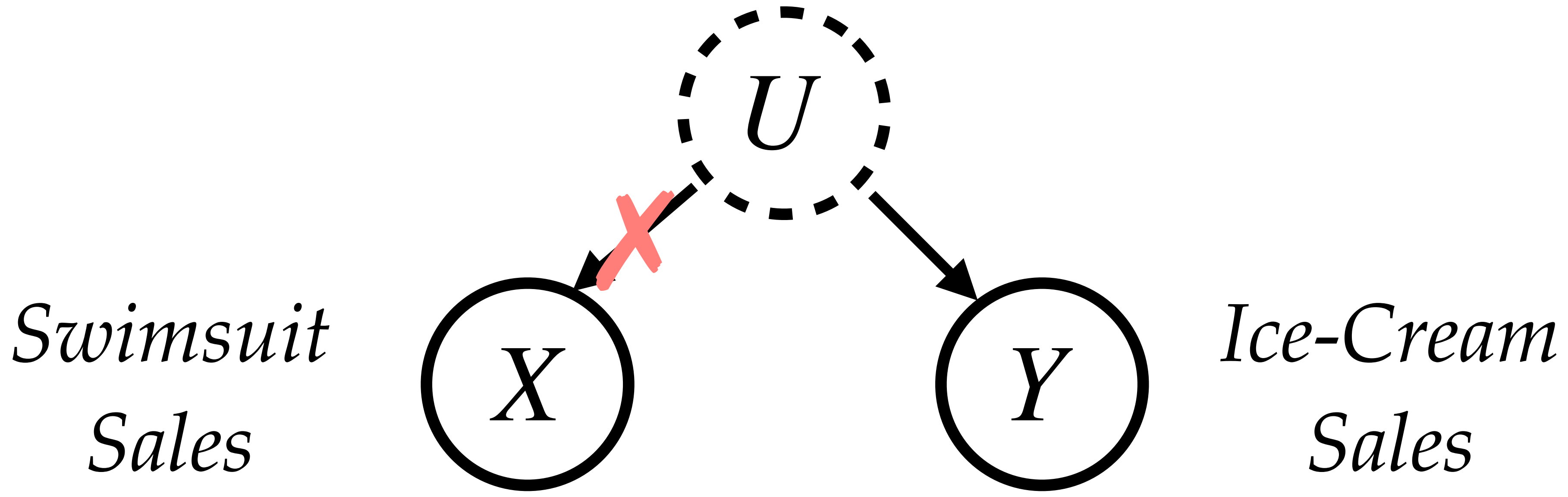
*Swimsuit
Sales*

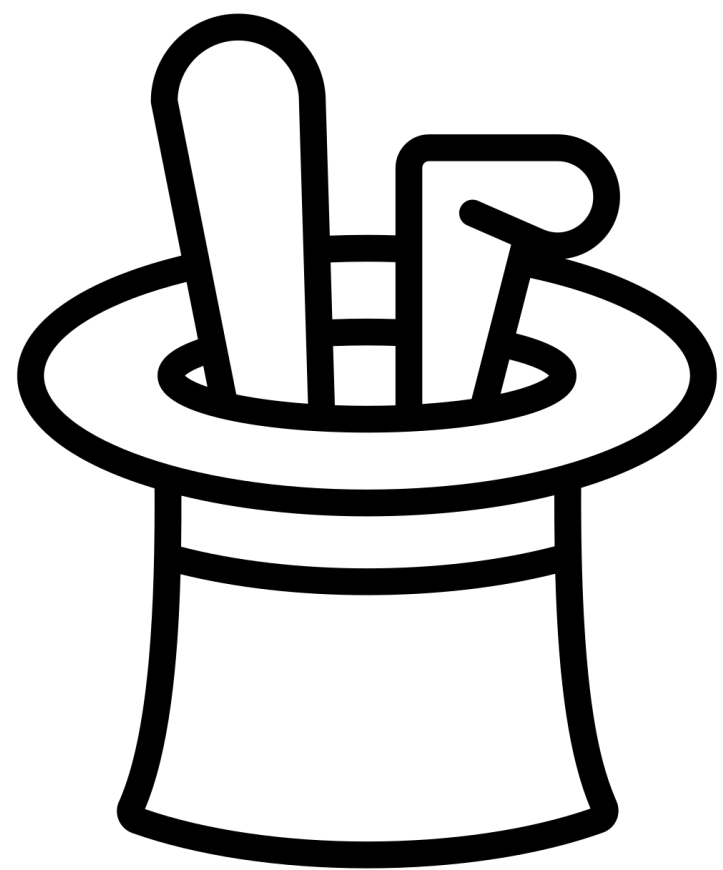


*Ice-Cream
Sales*

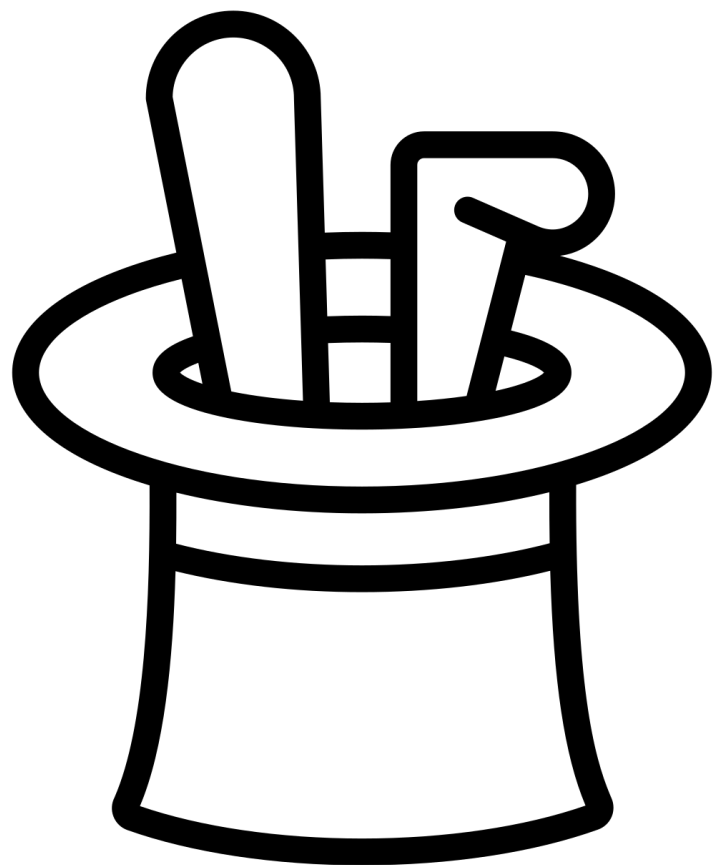


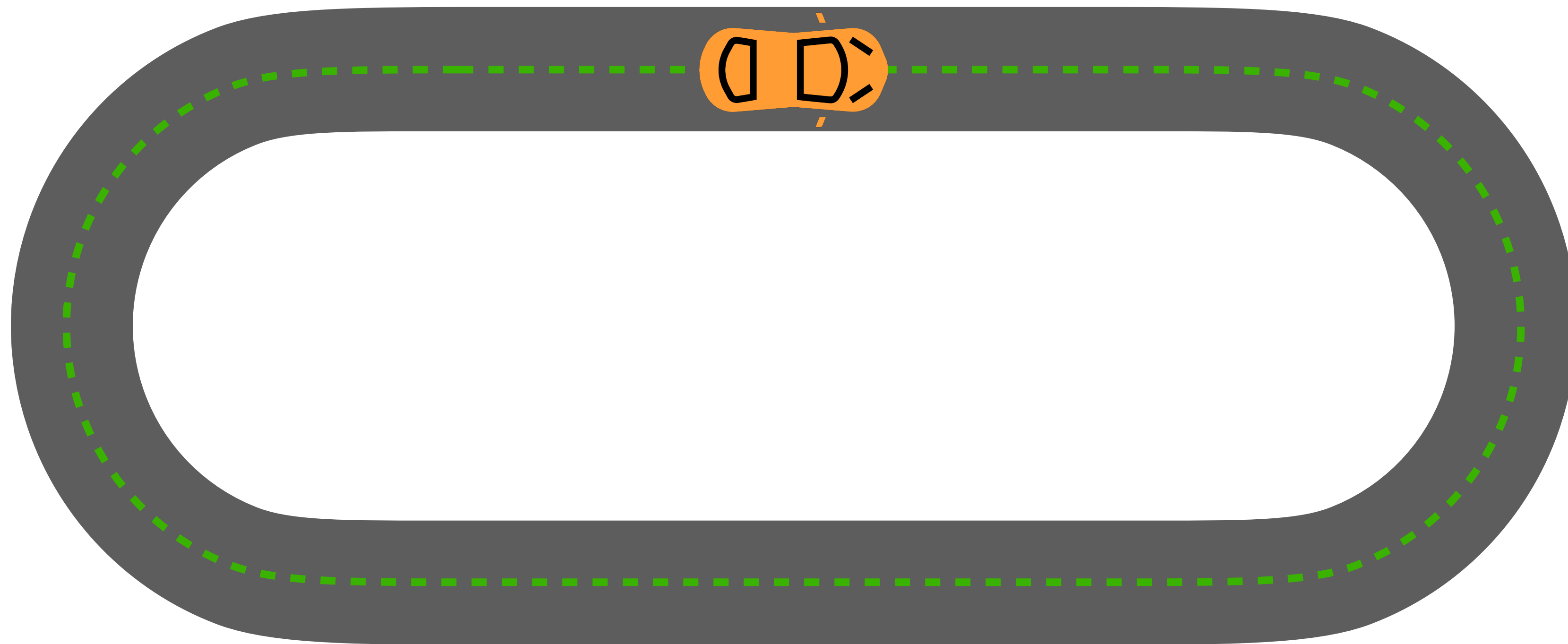
Temperature

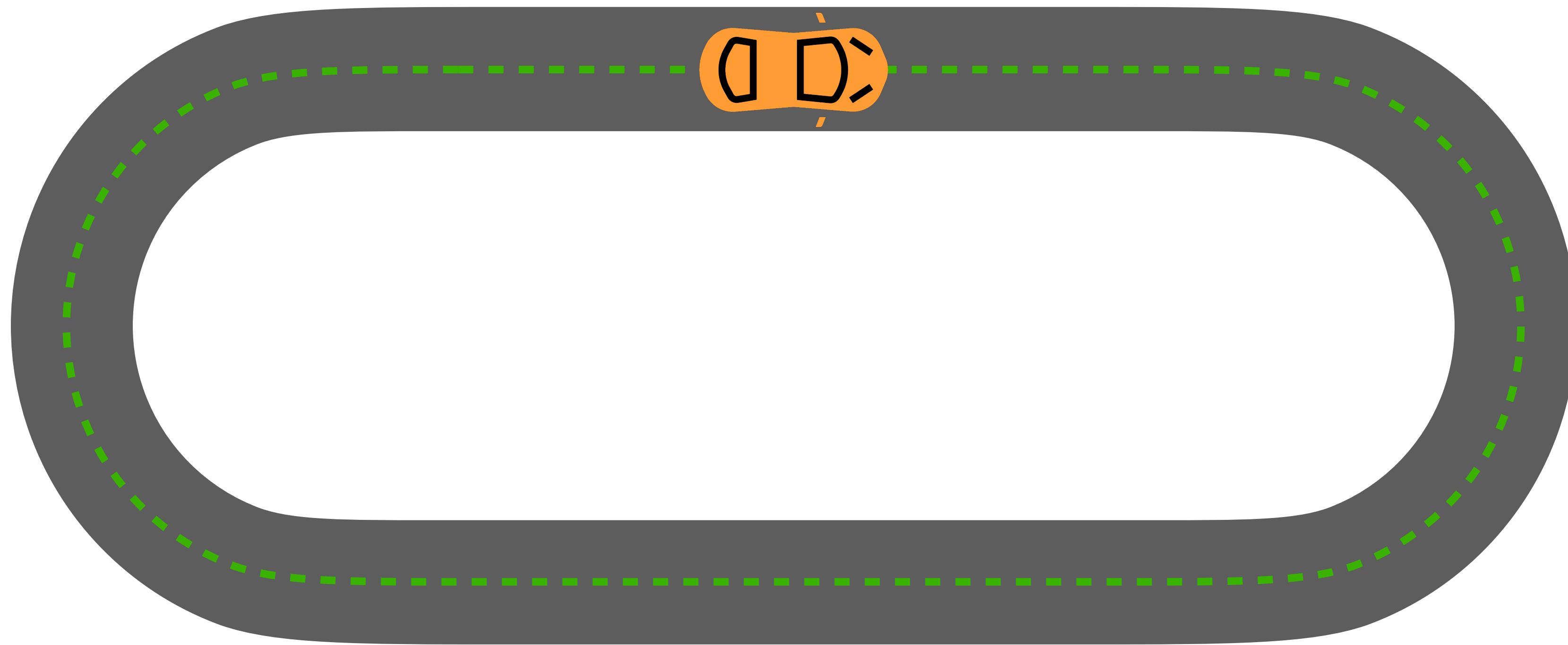




*Interventions happen via
interaction with
the environment in sequential
decision making.*

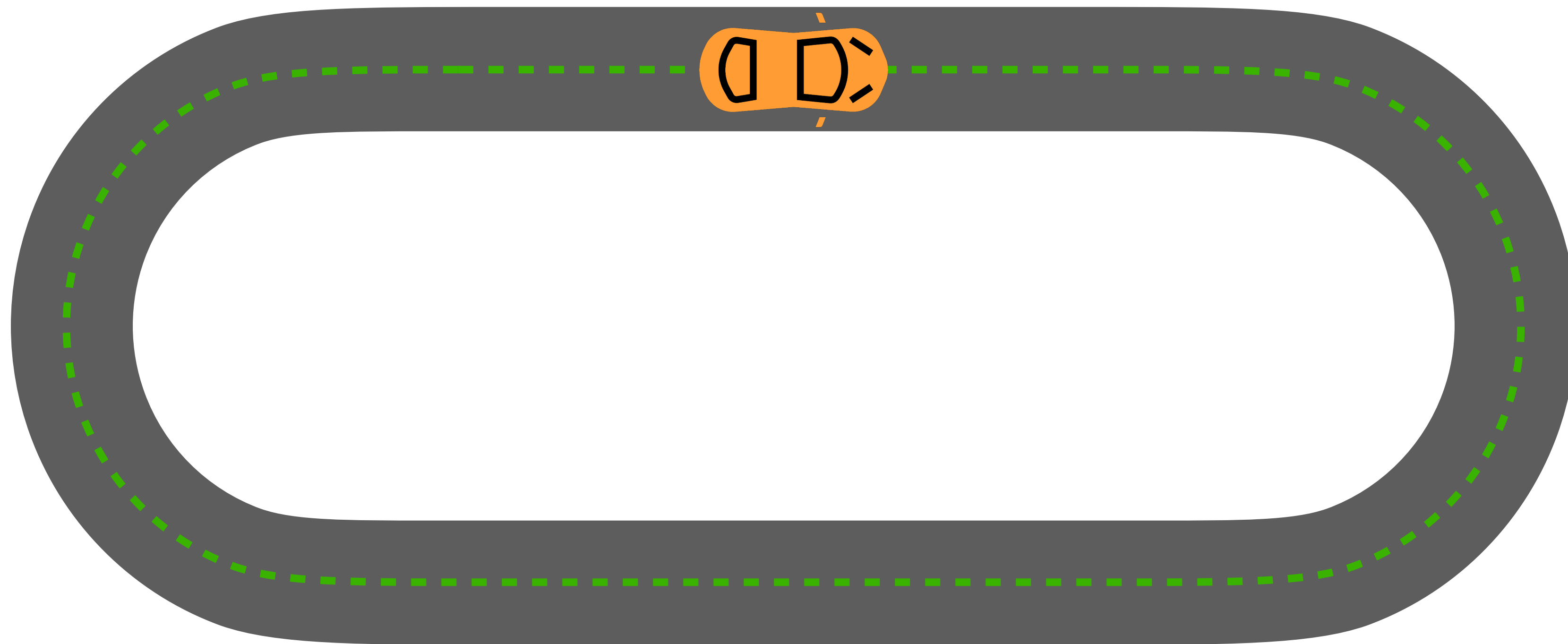




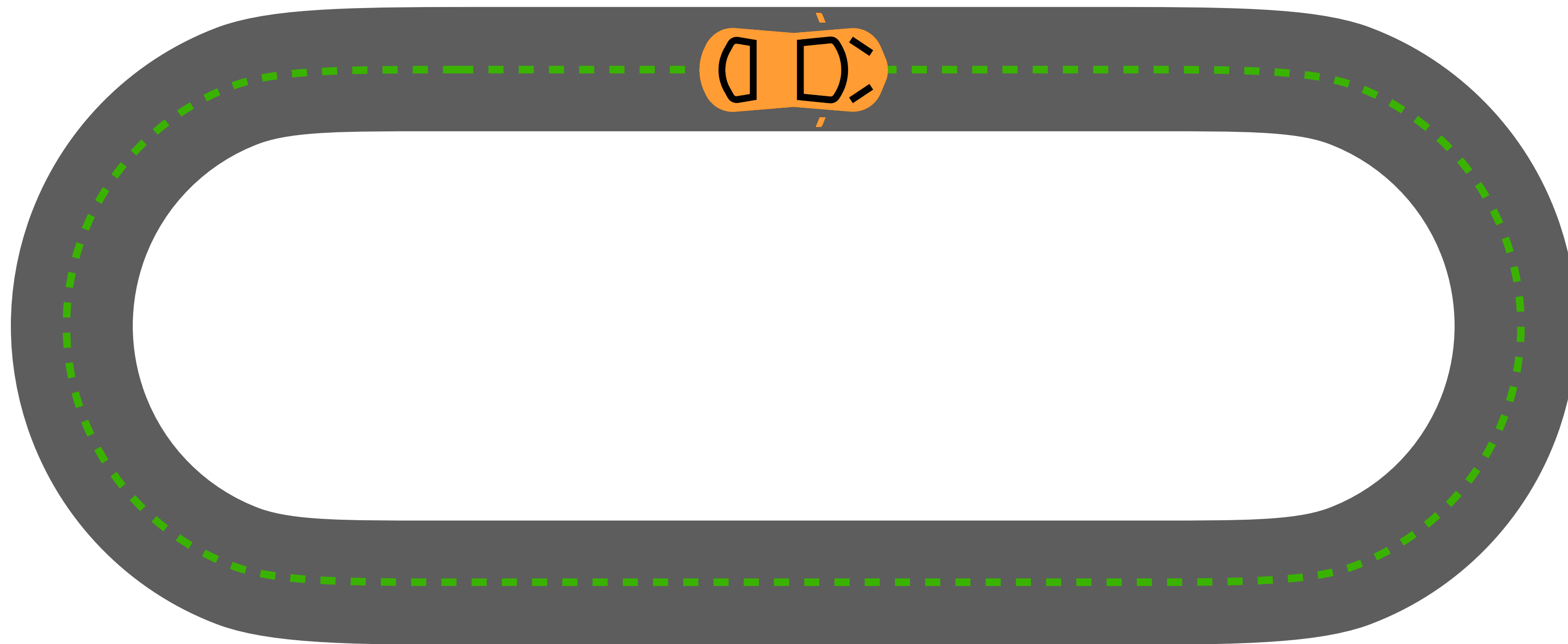


$\{s_1 \dots s_n\}$

$\{a_1 \dots a_n\}$

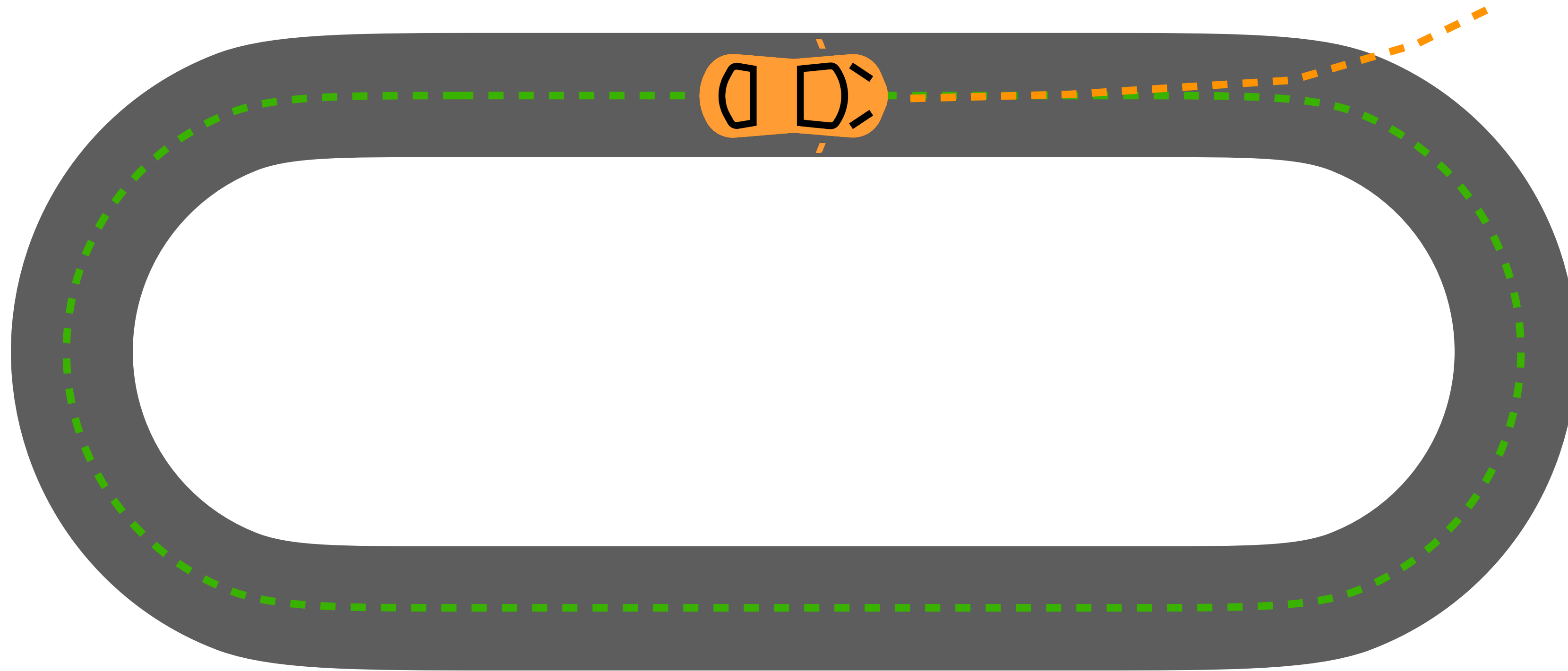


$$\{s_1 \dots s_n\} \mapsto \{a_1 \dots a_n\}$$



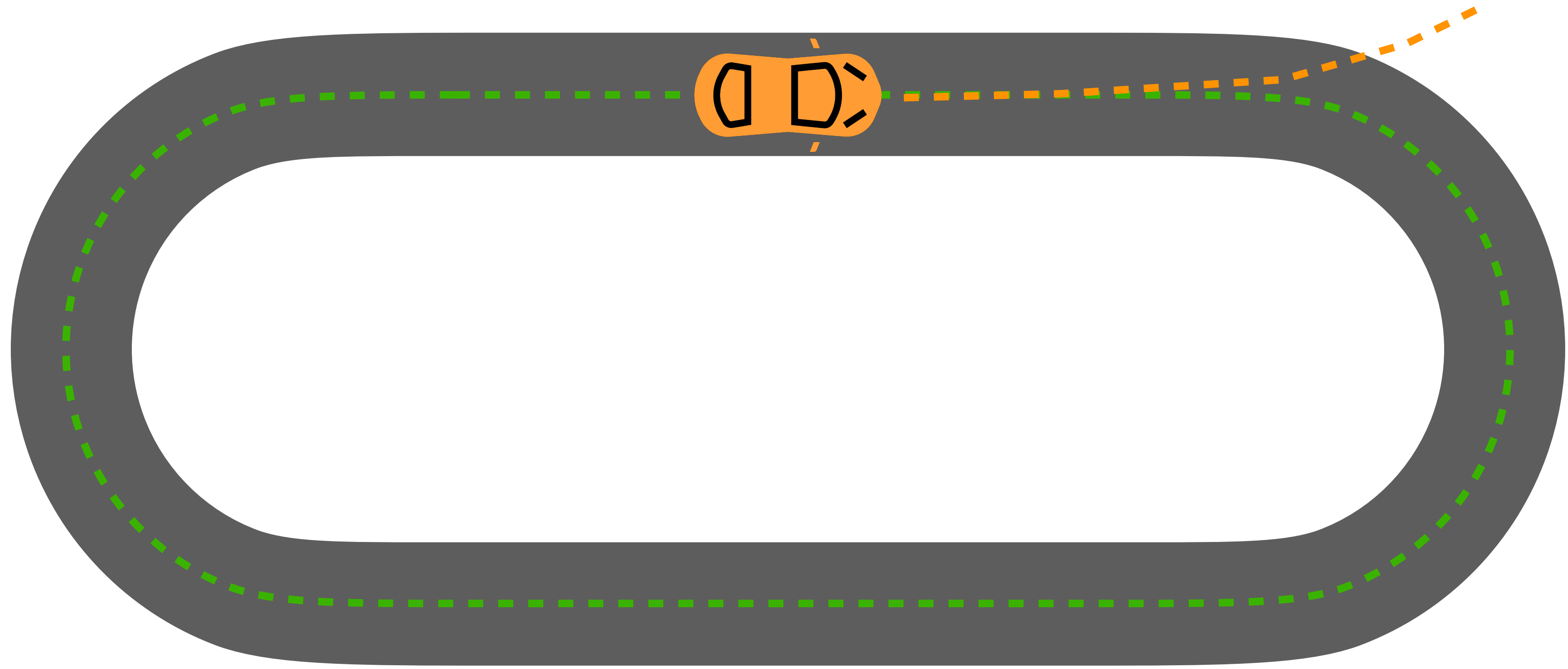
$$\{s_1 \dots s_n\} \mapsto \{a_1 \dots a_n\}$$

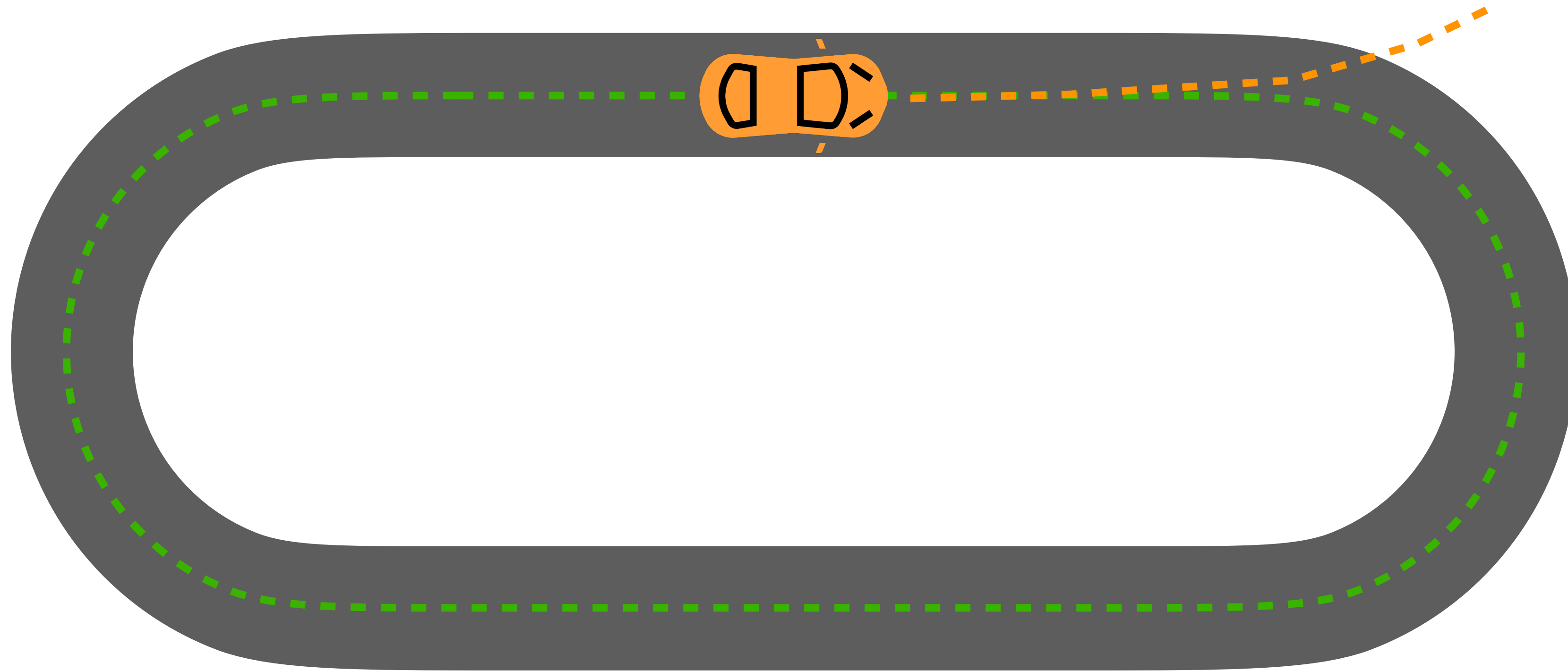
Behavioral Cloning



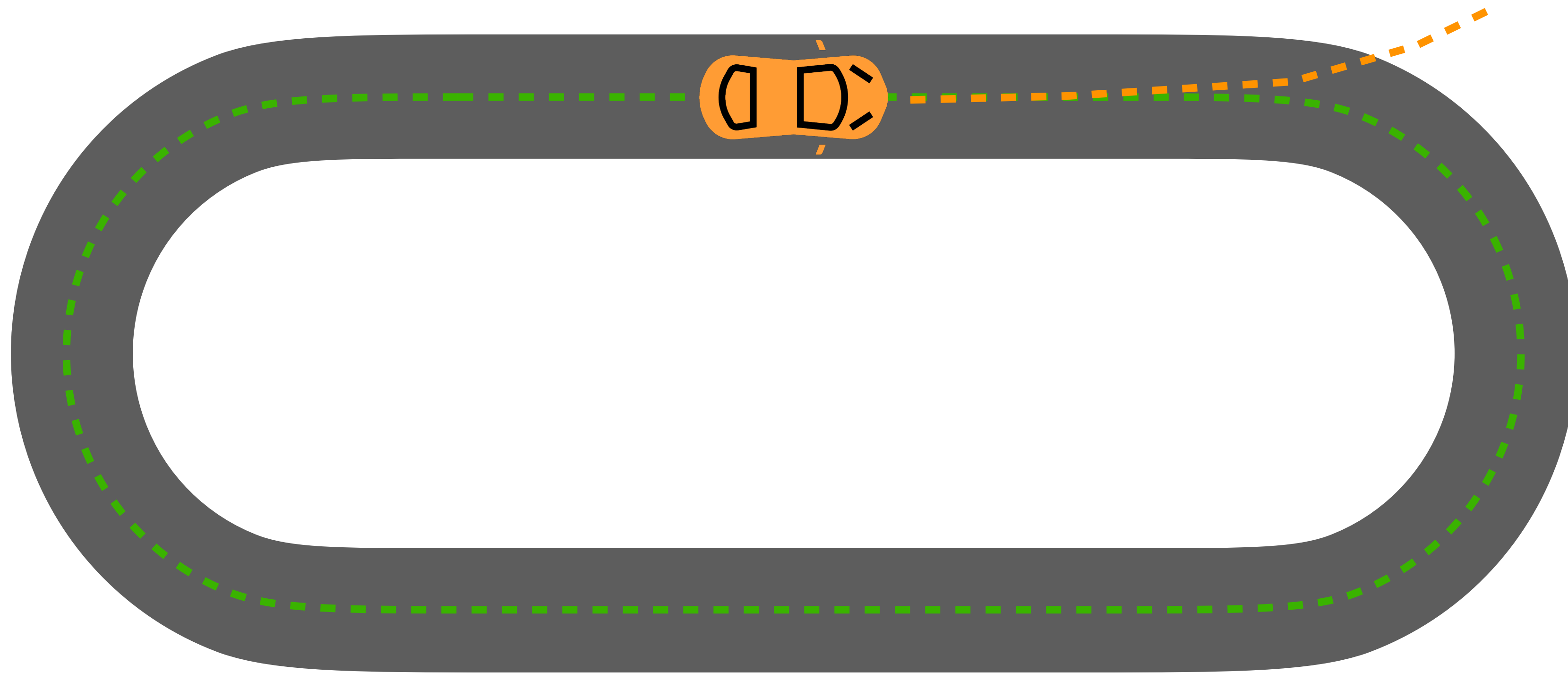
$$\{s_1 \dots s_n\} \mapsto \{a_1 \dots a_n\}$$

Behavioral Cloning



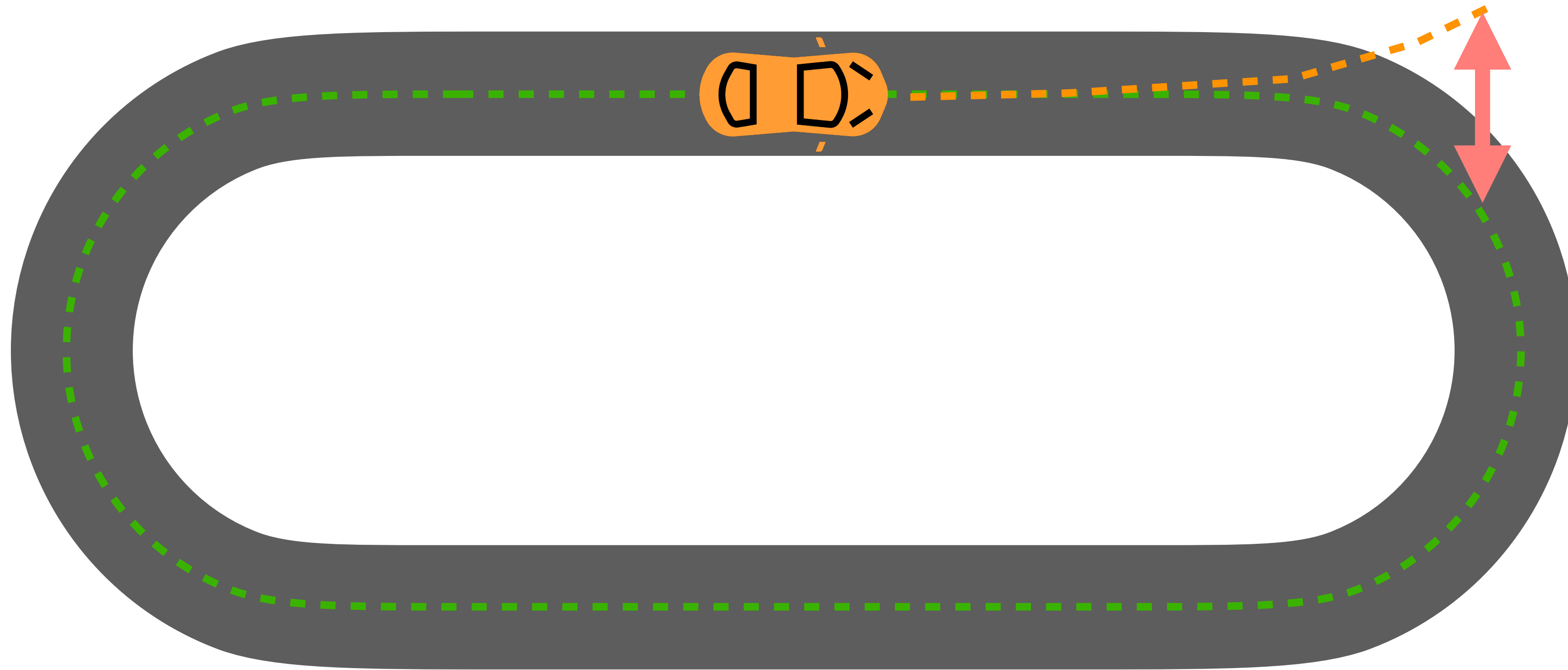


MaxEnt IRL / GAIL



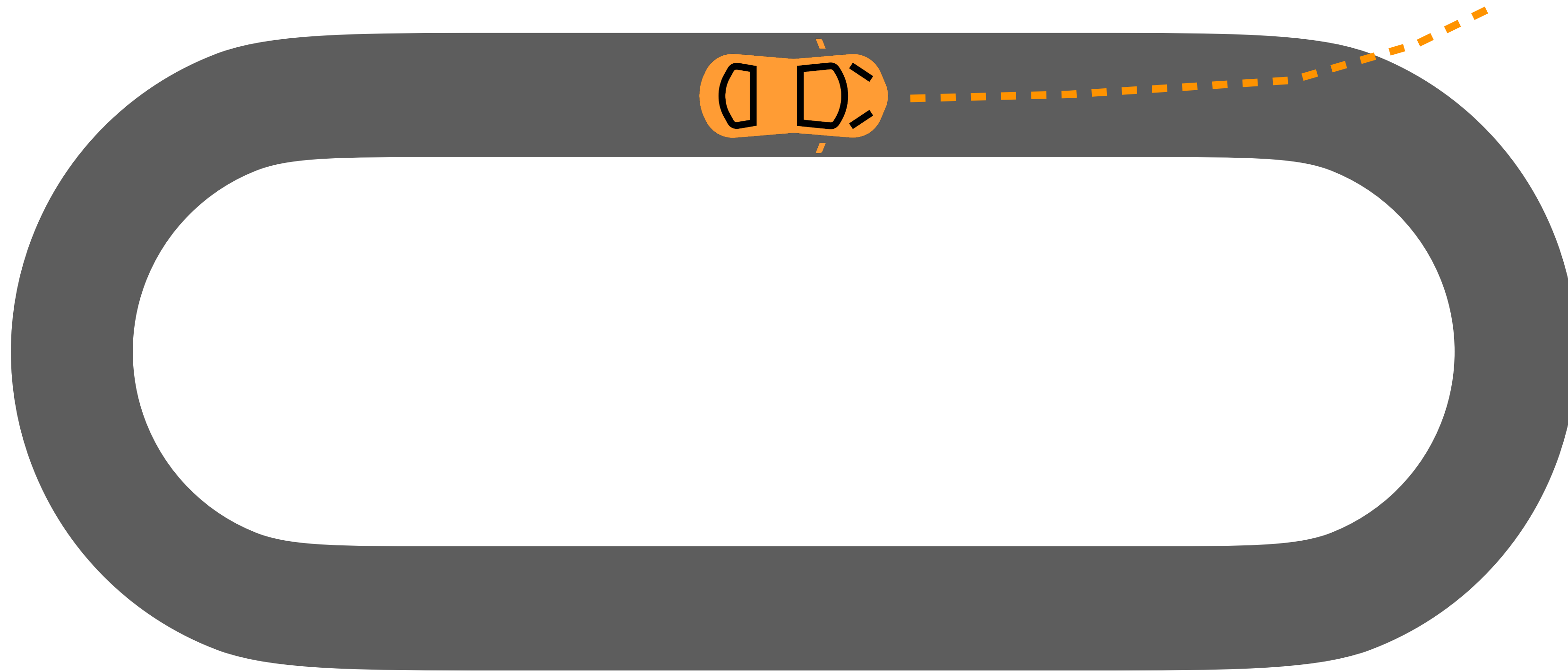
$\{s_1 \dots s_n\}$
 $\{a_1 \dots a_n\}$

MaxEnt IRL / GAIL



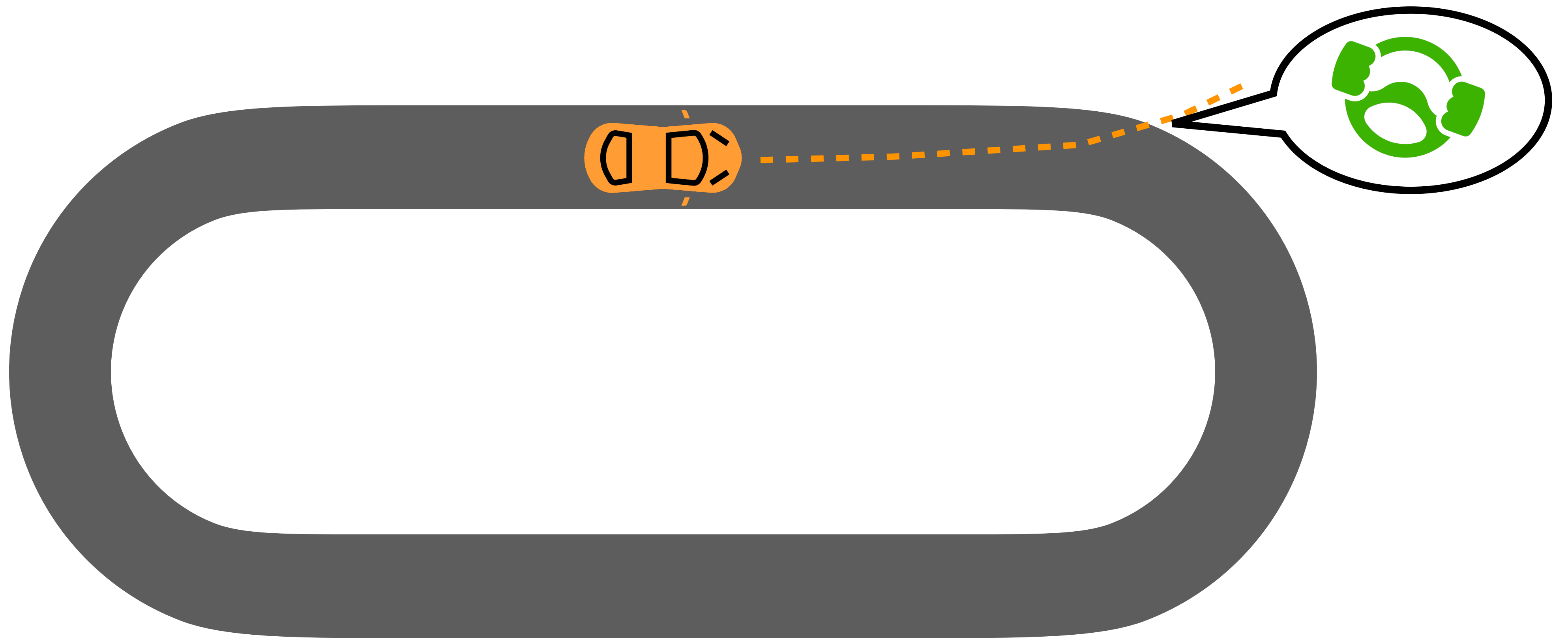
$$\begin{array}{ccc} \{s_1 \dots s_n\} & \longleftrightarrow & \{s_1 \dots s_n\} \\ \{a_1 \dots a_n\} & & \{a_1 \dots a_n\} \end{array}$$

MaxEnt IRL / GAIL



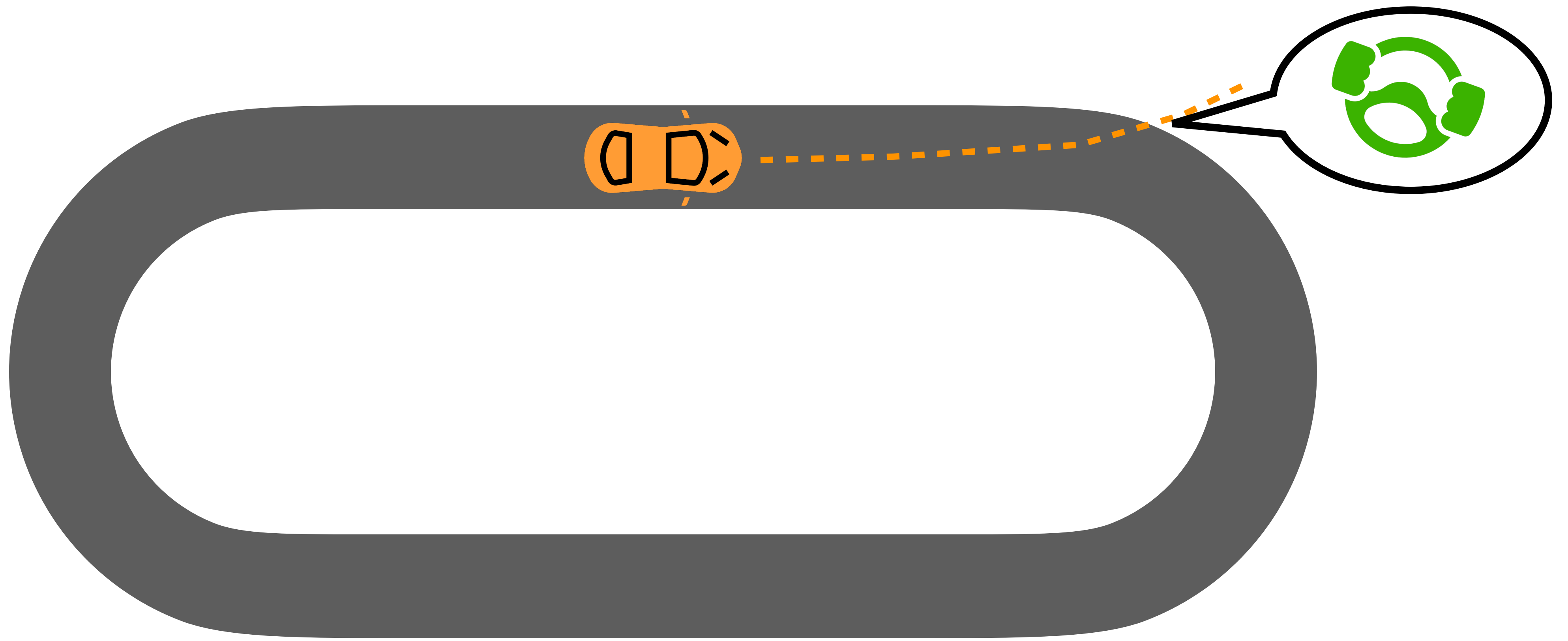
$\{s_1 \dots s_n\}$

Dagger



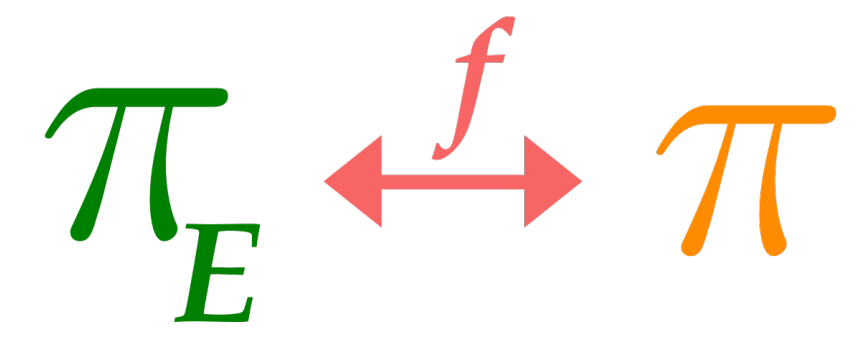
$\{s_1 \dots s_n\}$

Dagger

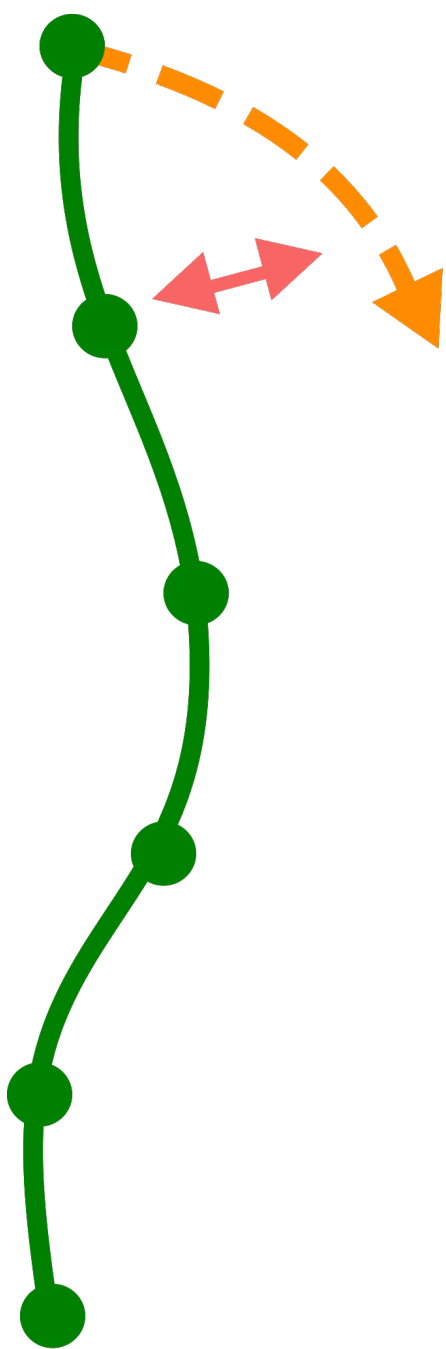


$$\{s_1 \dots s_n\} \mapsto \{a_1 \dots a_n\}$$

Dagger

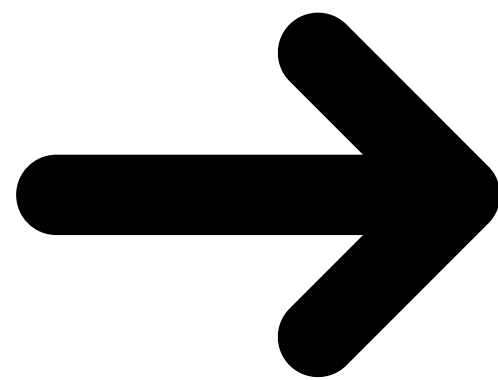
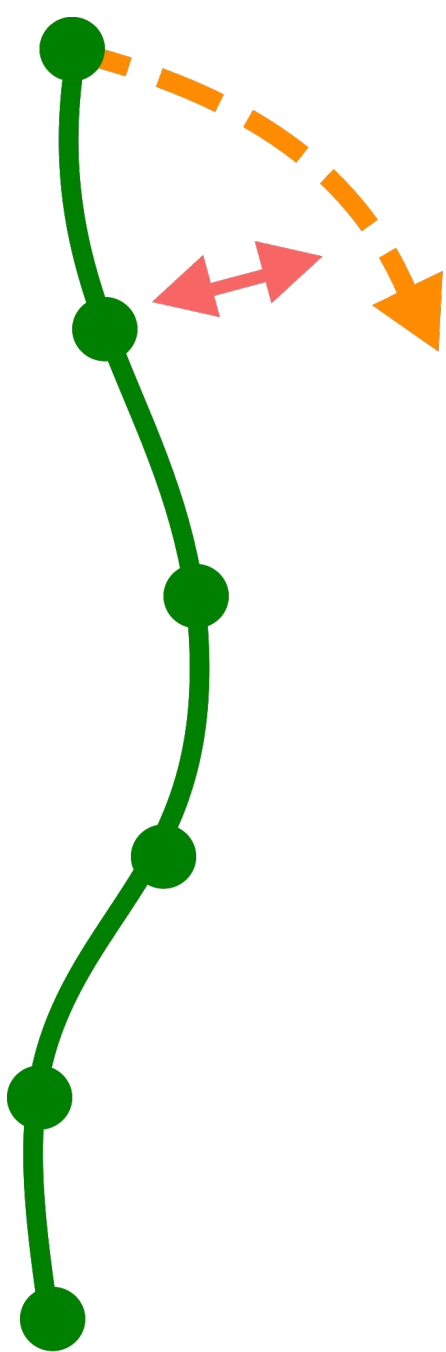


$$\pi_E \xleftrightarrow{f} \pi$$



Behavioral Cloning

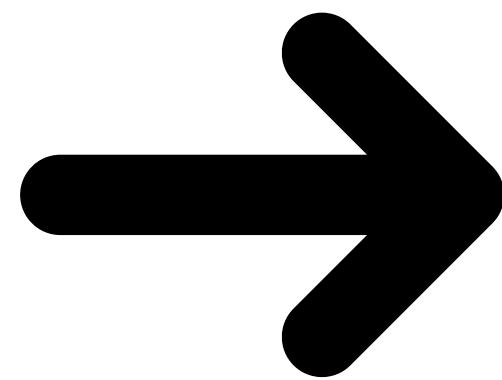
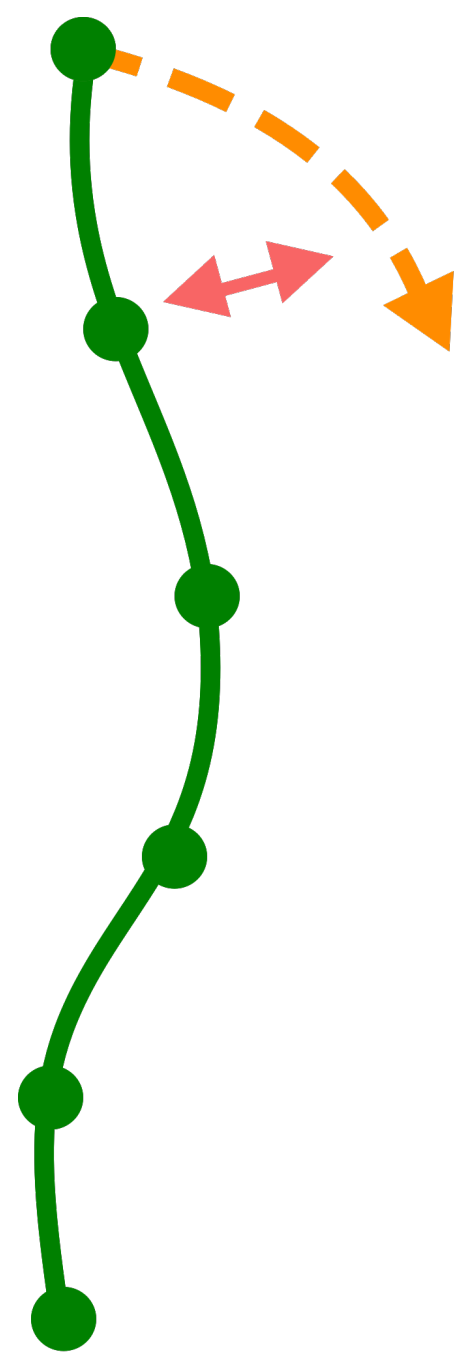
$$\pi_E \xleftrightarrow{f} \pi$$



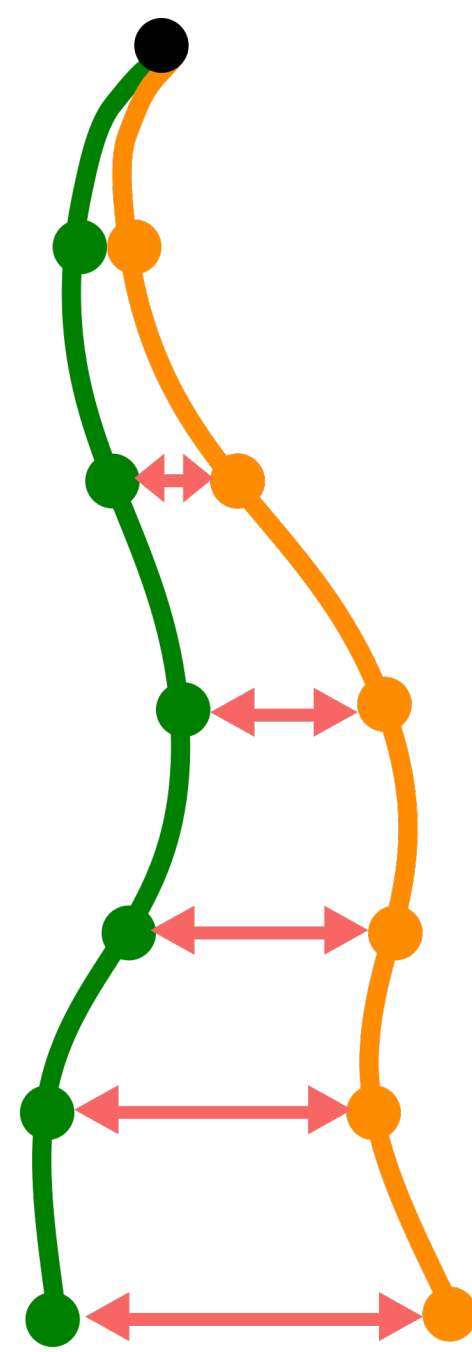
Environment

Behavioral Cloning

$$\pi_E \xleftrightarrow{f} \pi$$



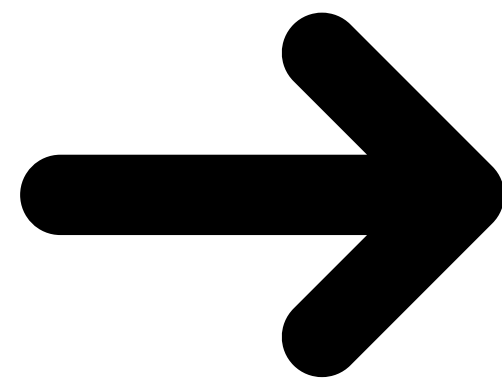
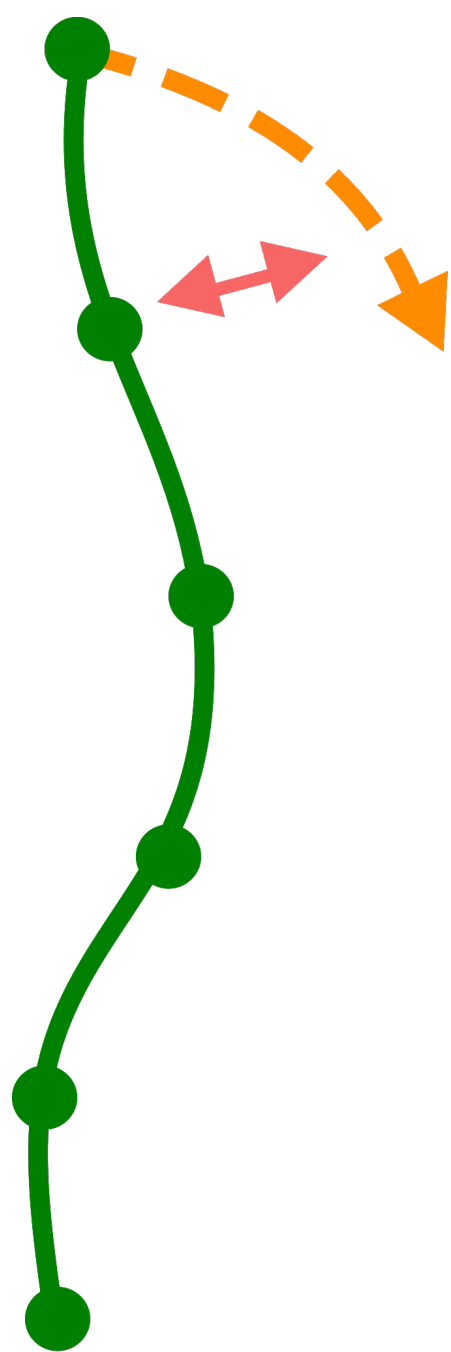
Environment



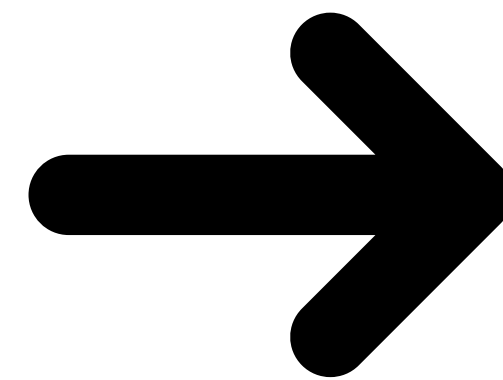
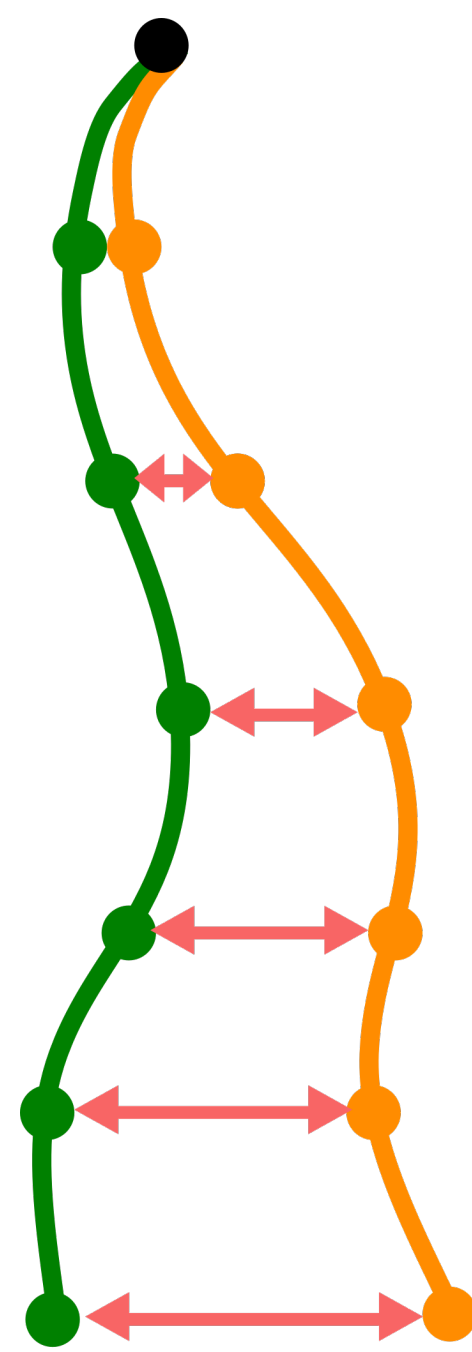
Behavioral Cloning

GAIL

$$\pi_E \xleftrightarrow{f} \pi$$



Environment

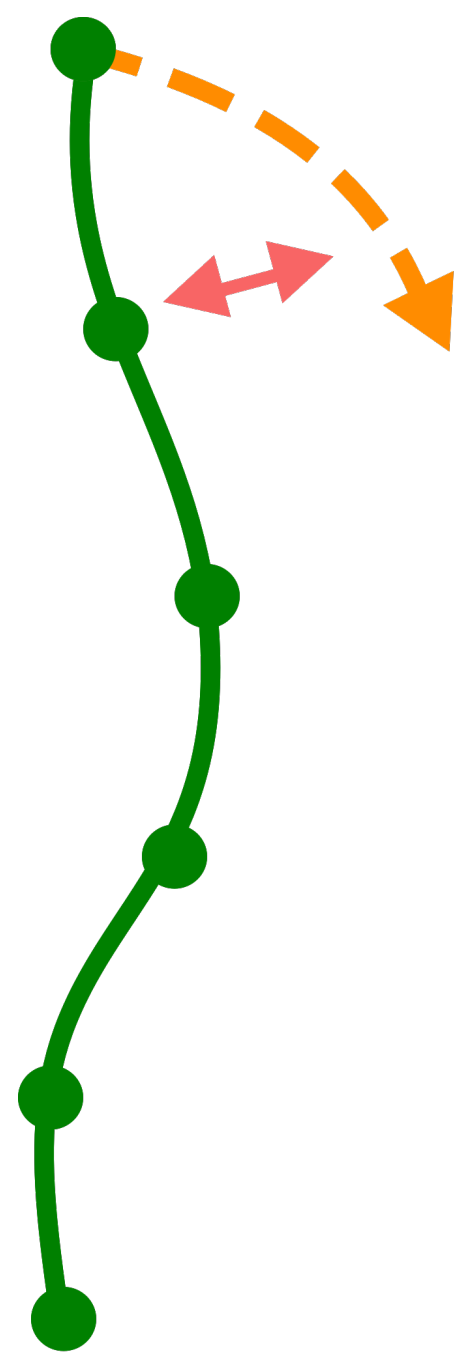


Query Expert

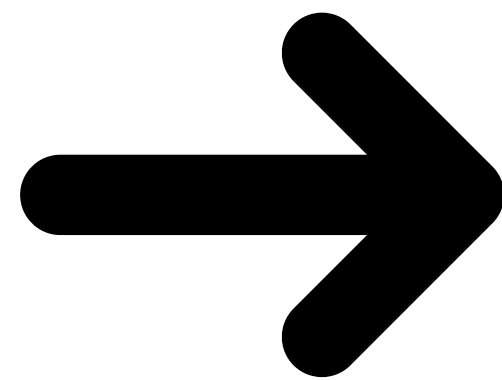
Behavioral Cloning

GAIL

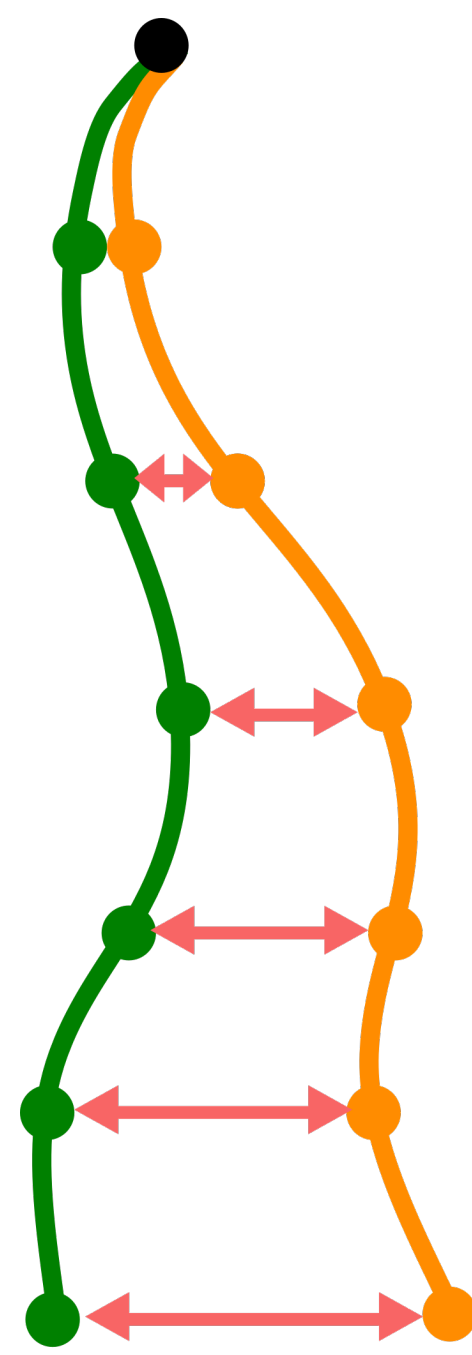
$$\pi_E \xleftrightarrow{f} \pi$$



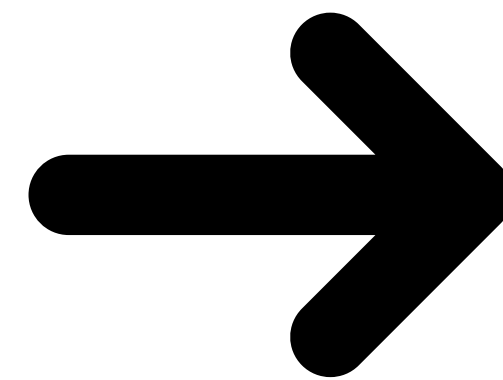
Behavioral Cloning



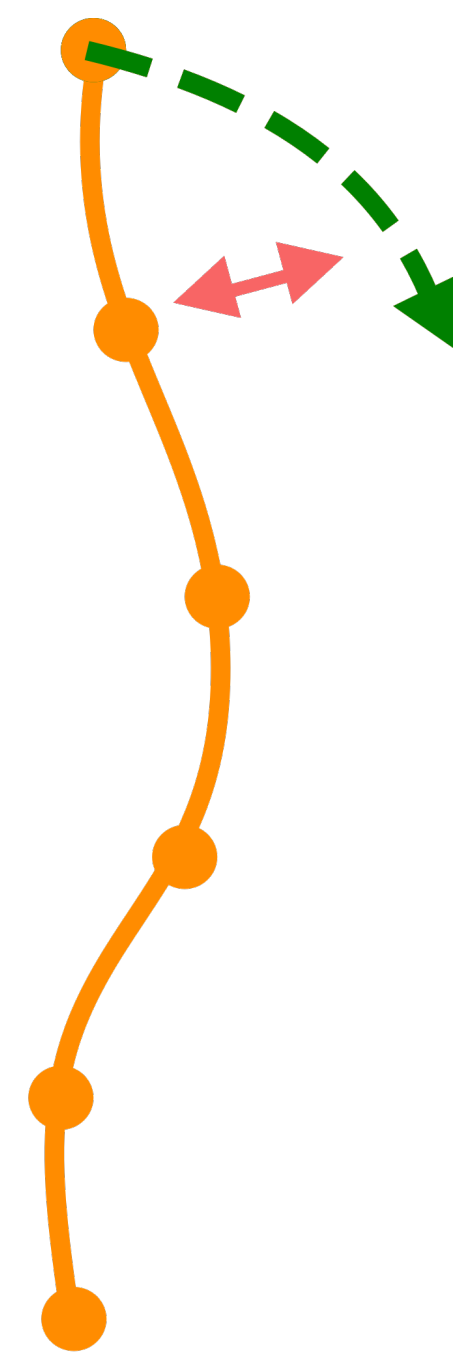
Environment



GAIL



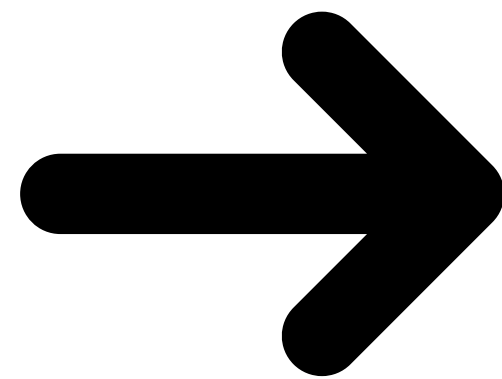
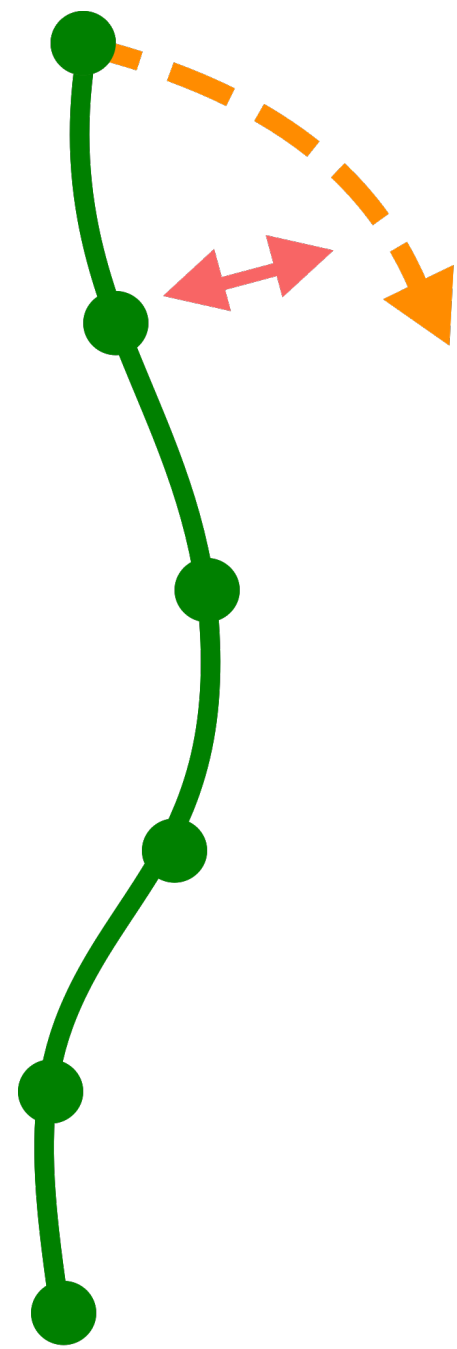
Query Expert



DAgger

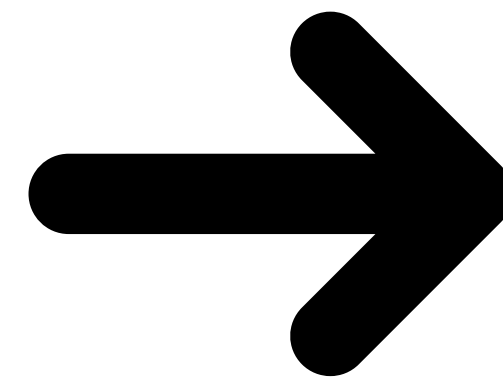
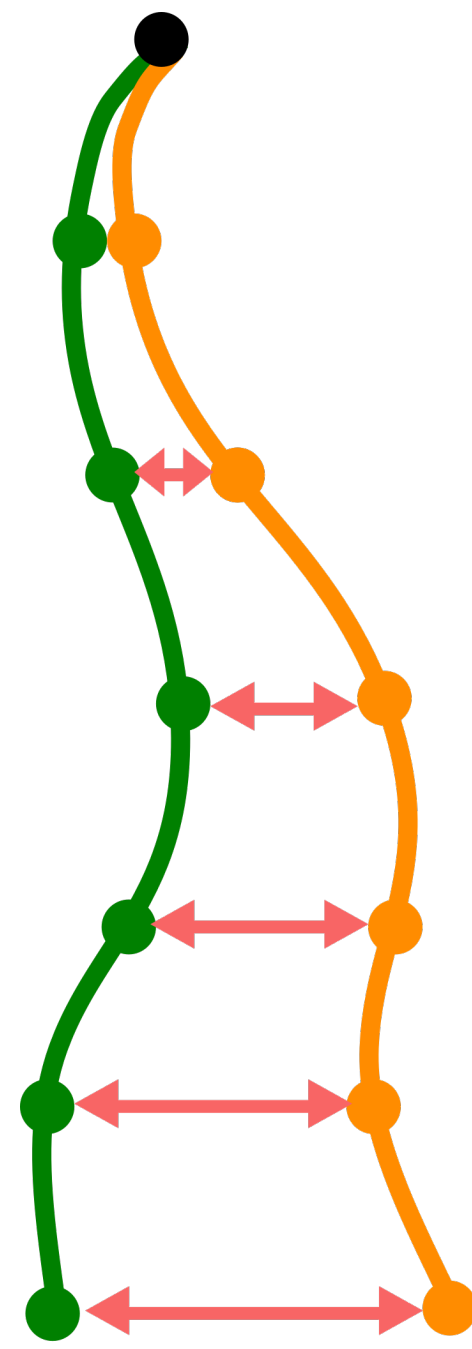
$$\pi_E \xleftrightarrow{f} \pi$$

Offline



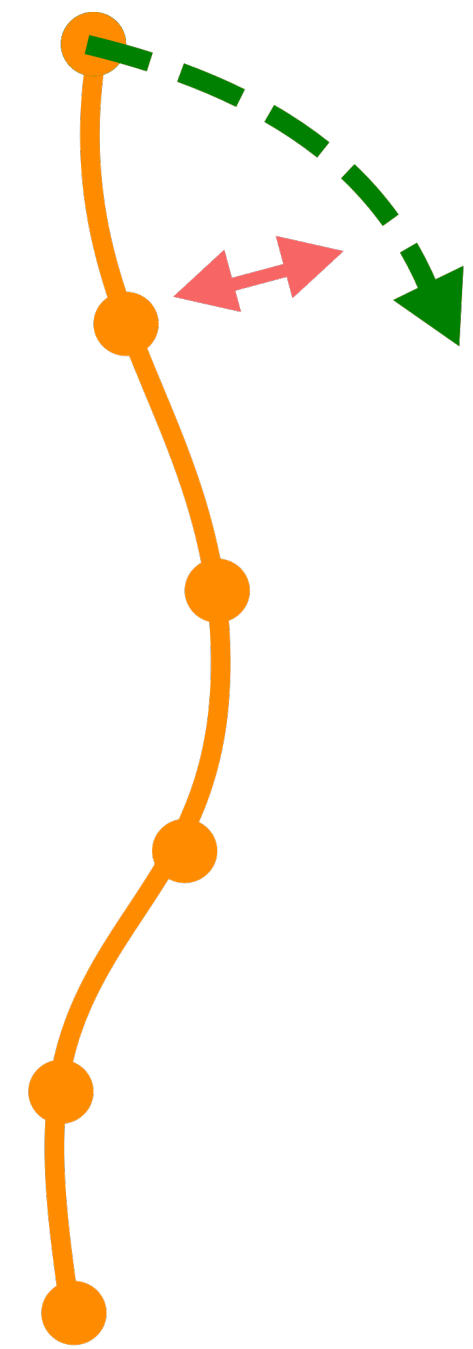
Environment

Online



Query Expert

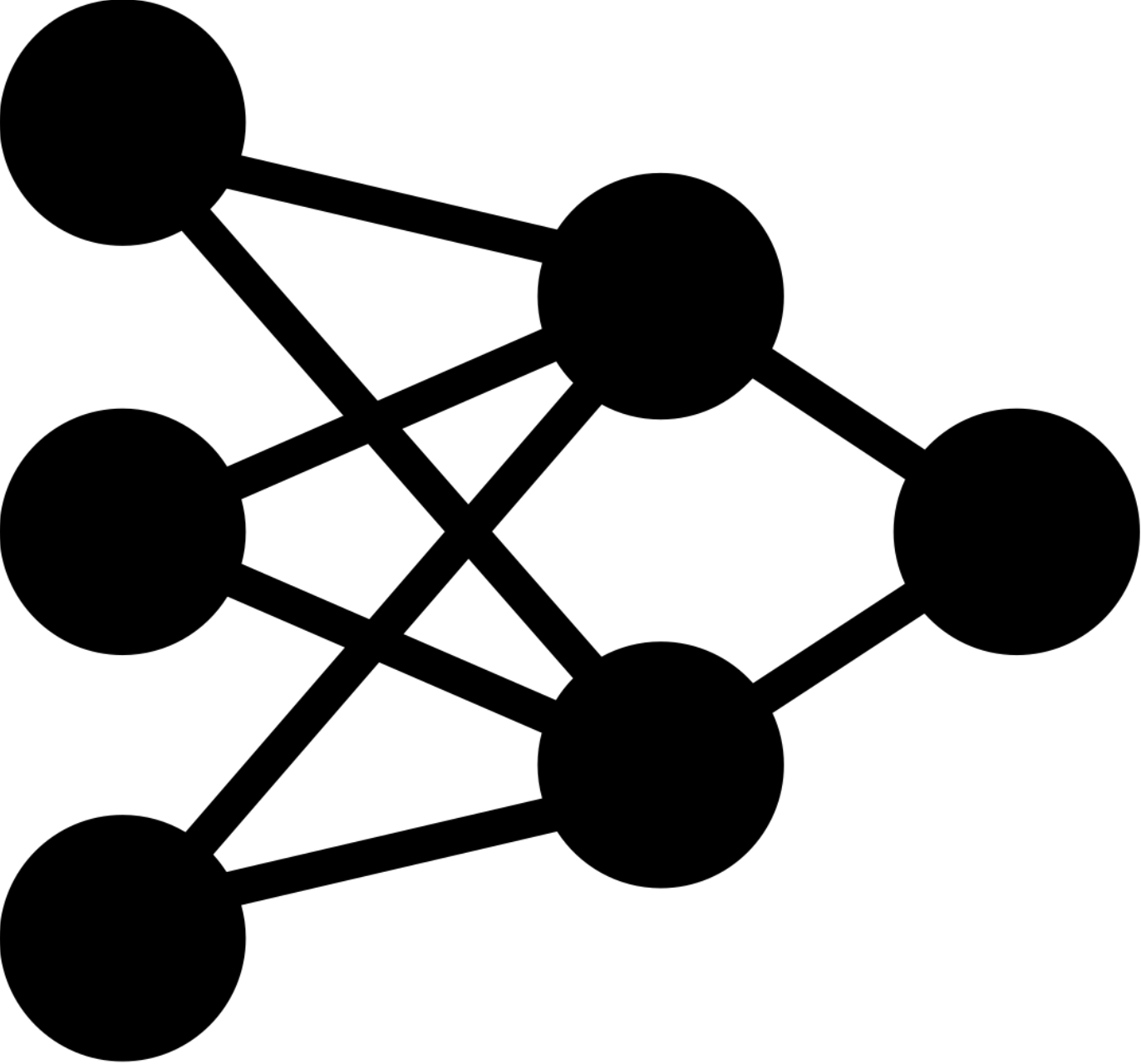
Interactive

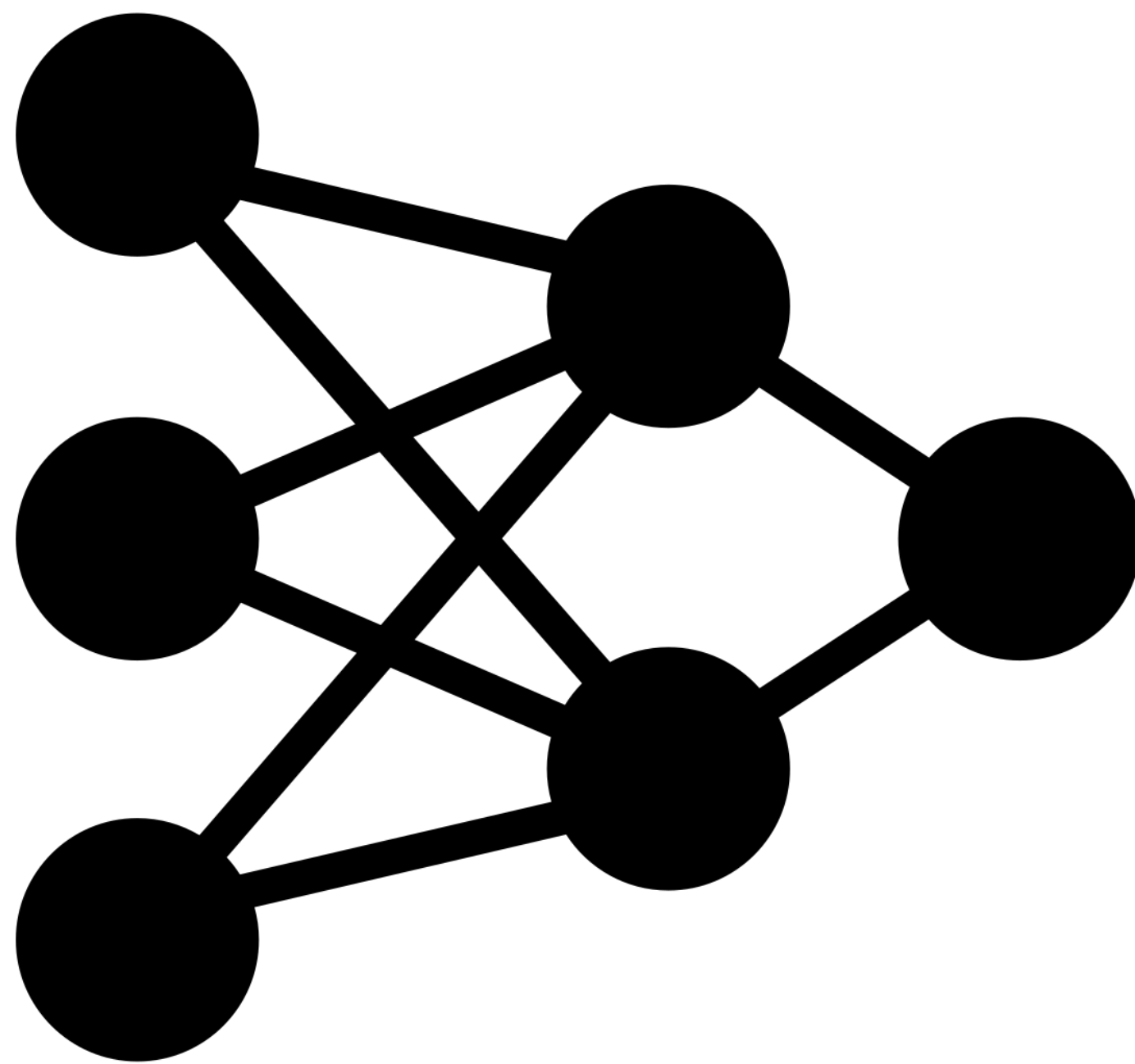


Behavioral Cloning

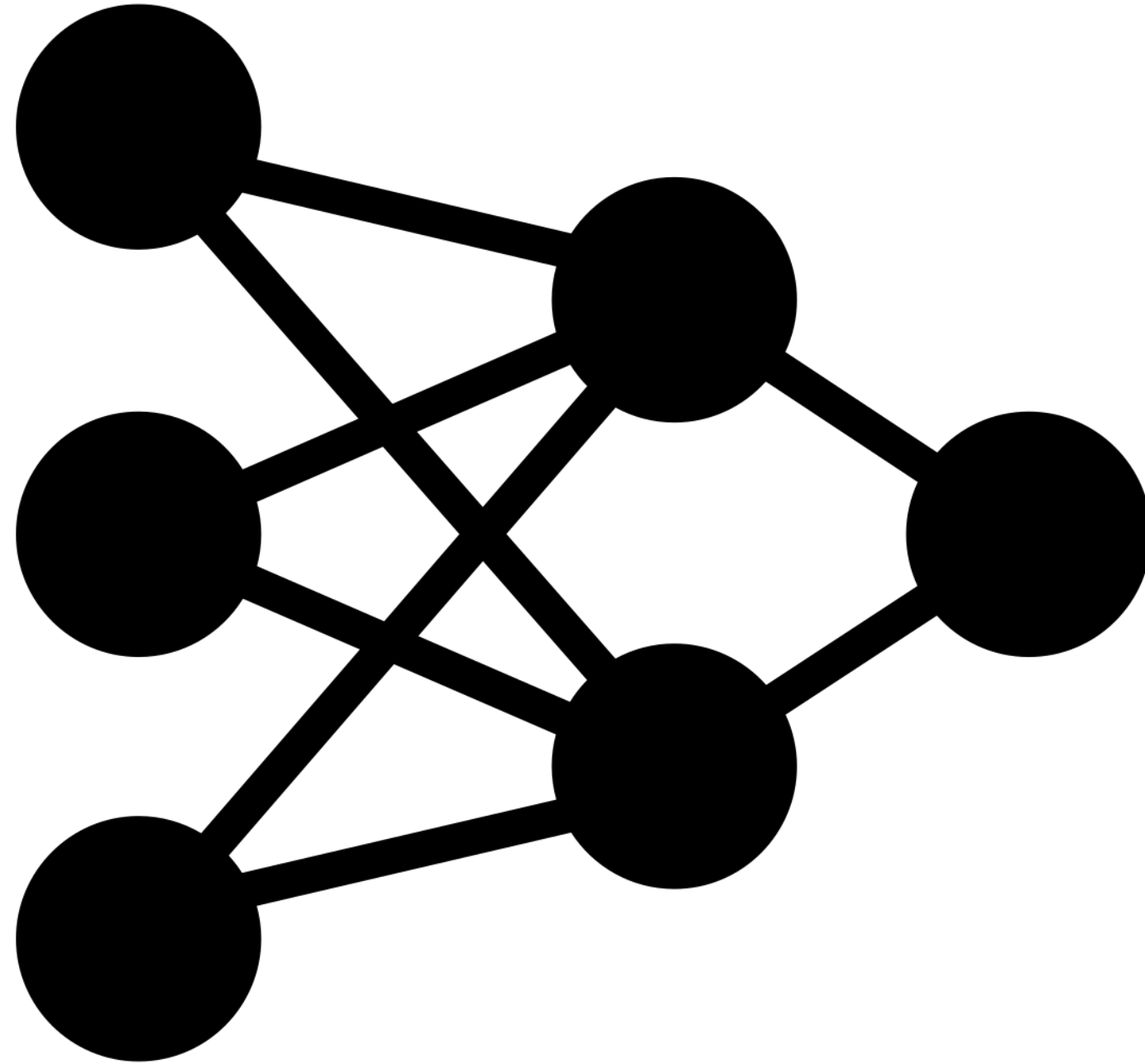
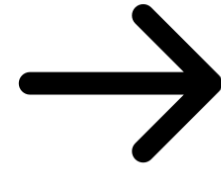
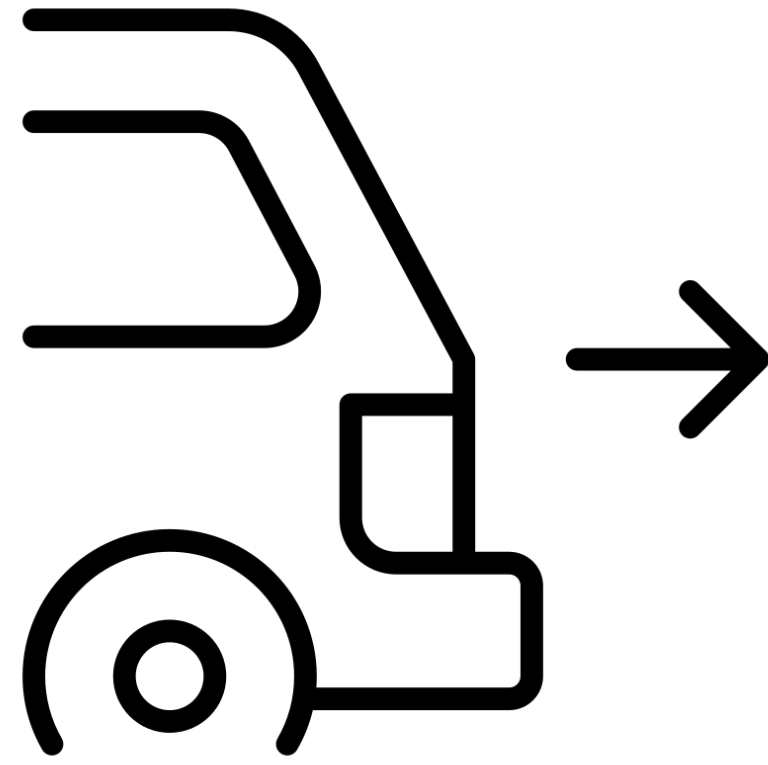
GAIL

DAgger

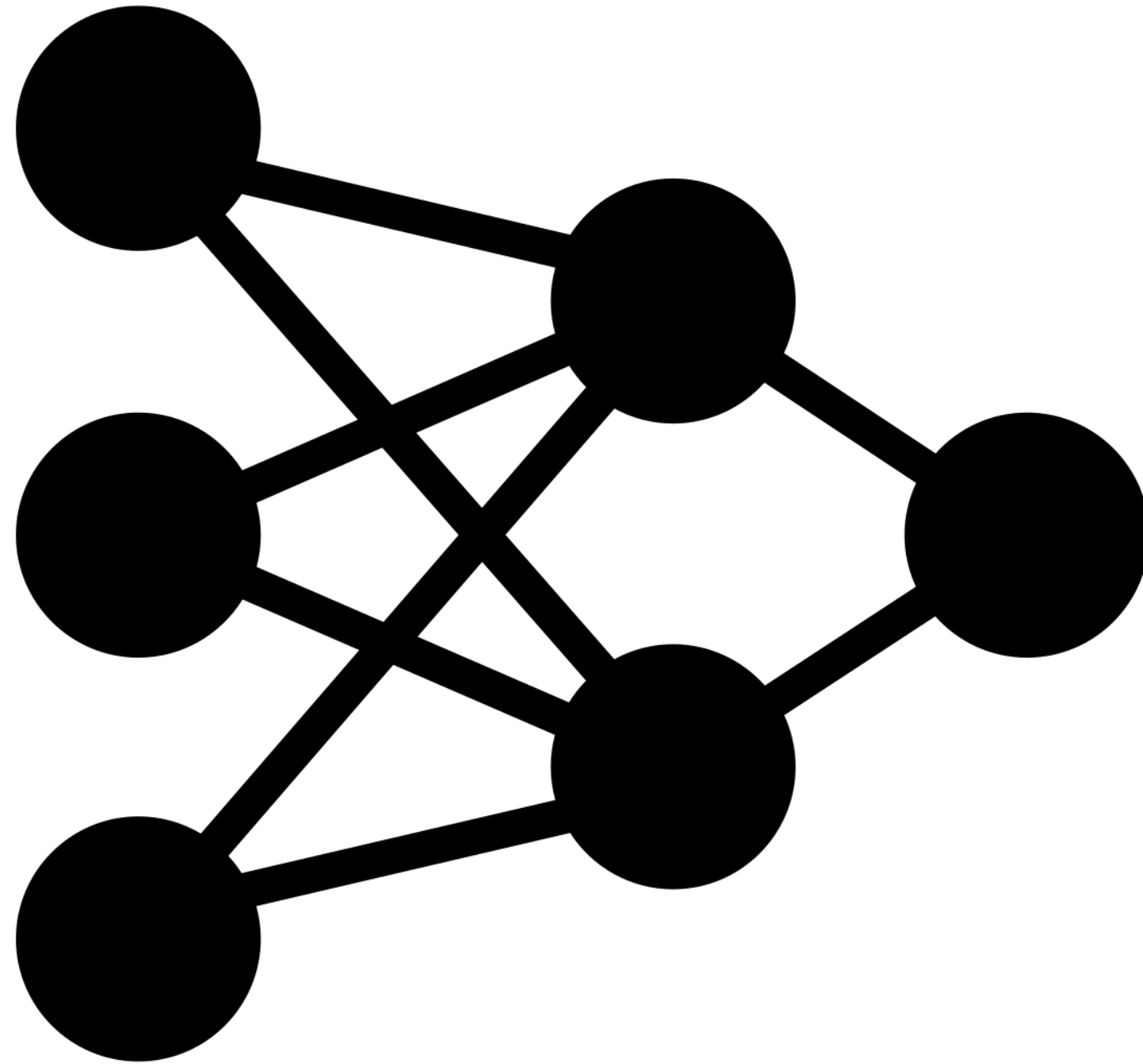
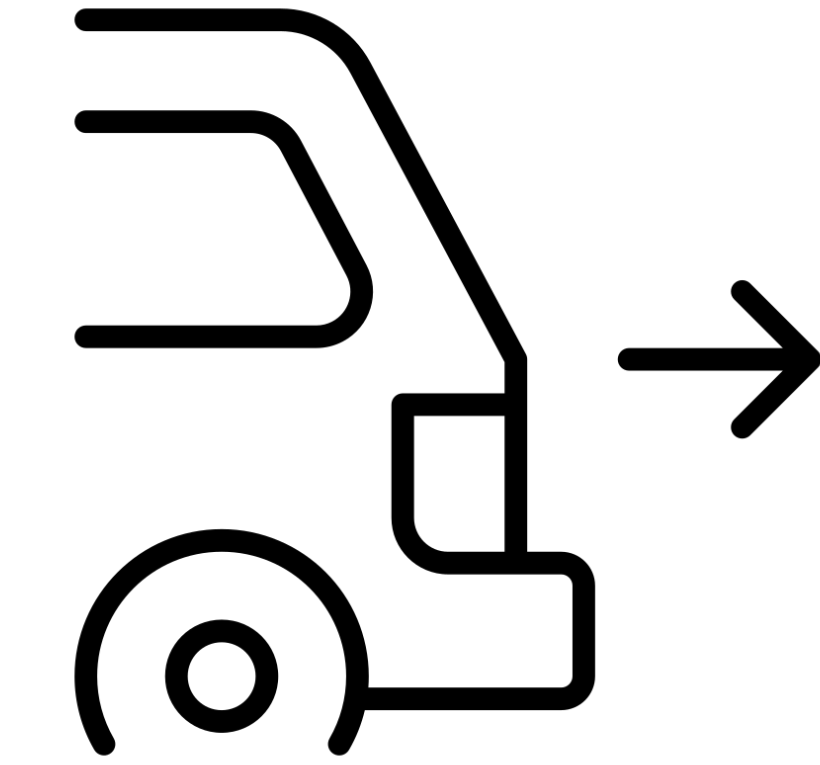




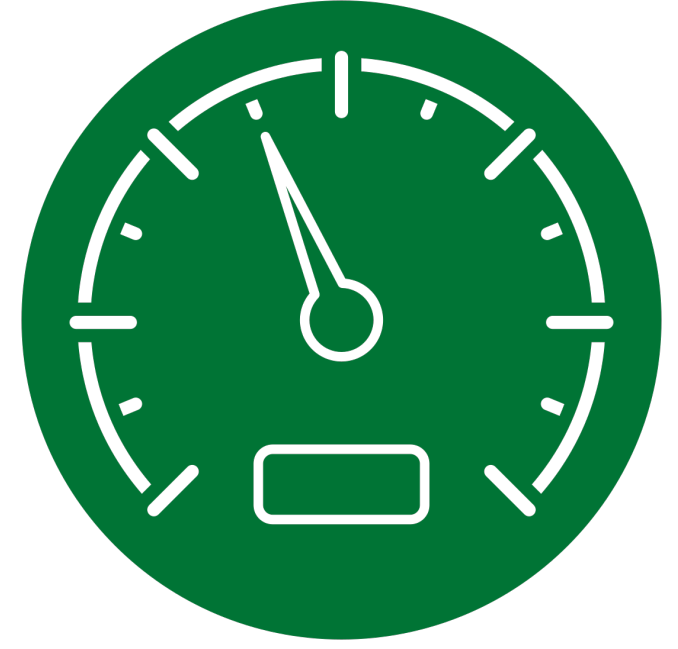
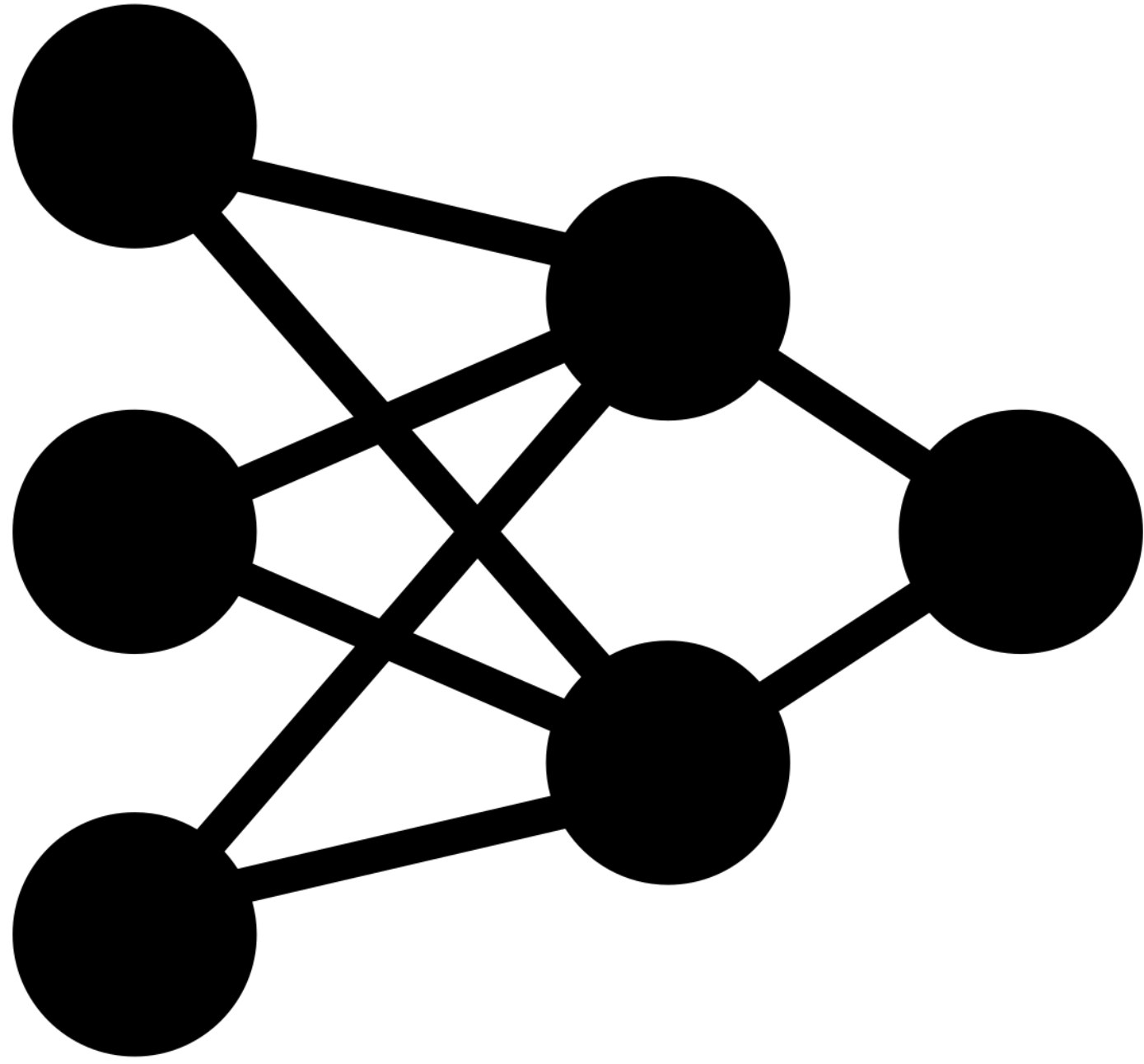
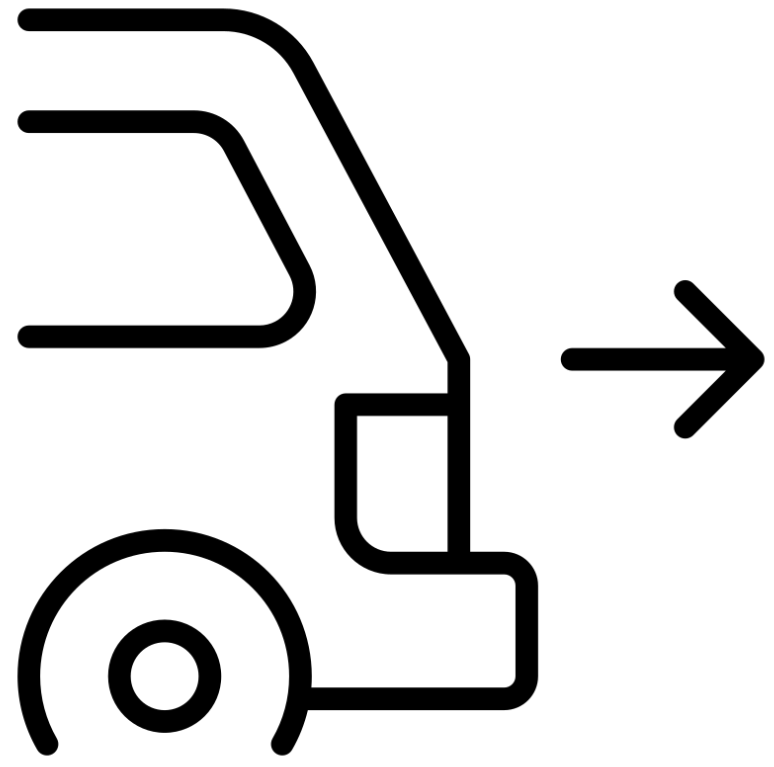
Brake?

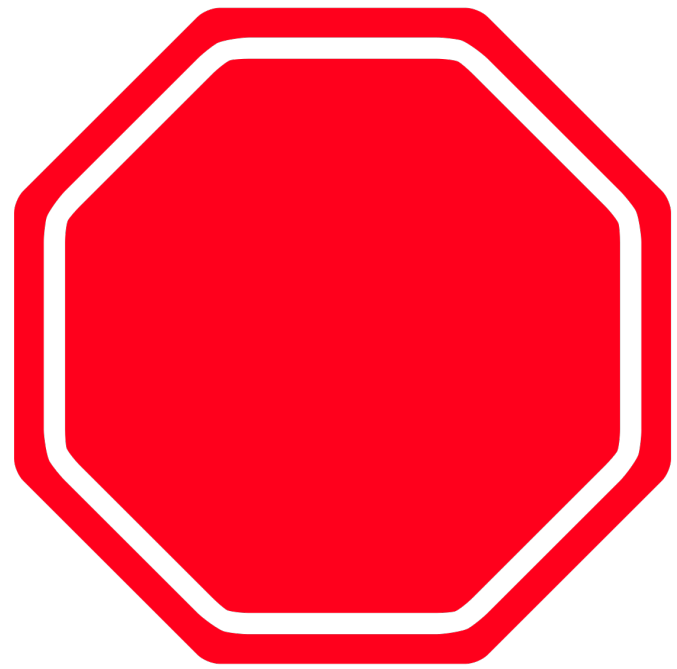
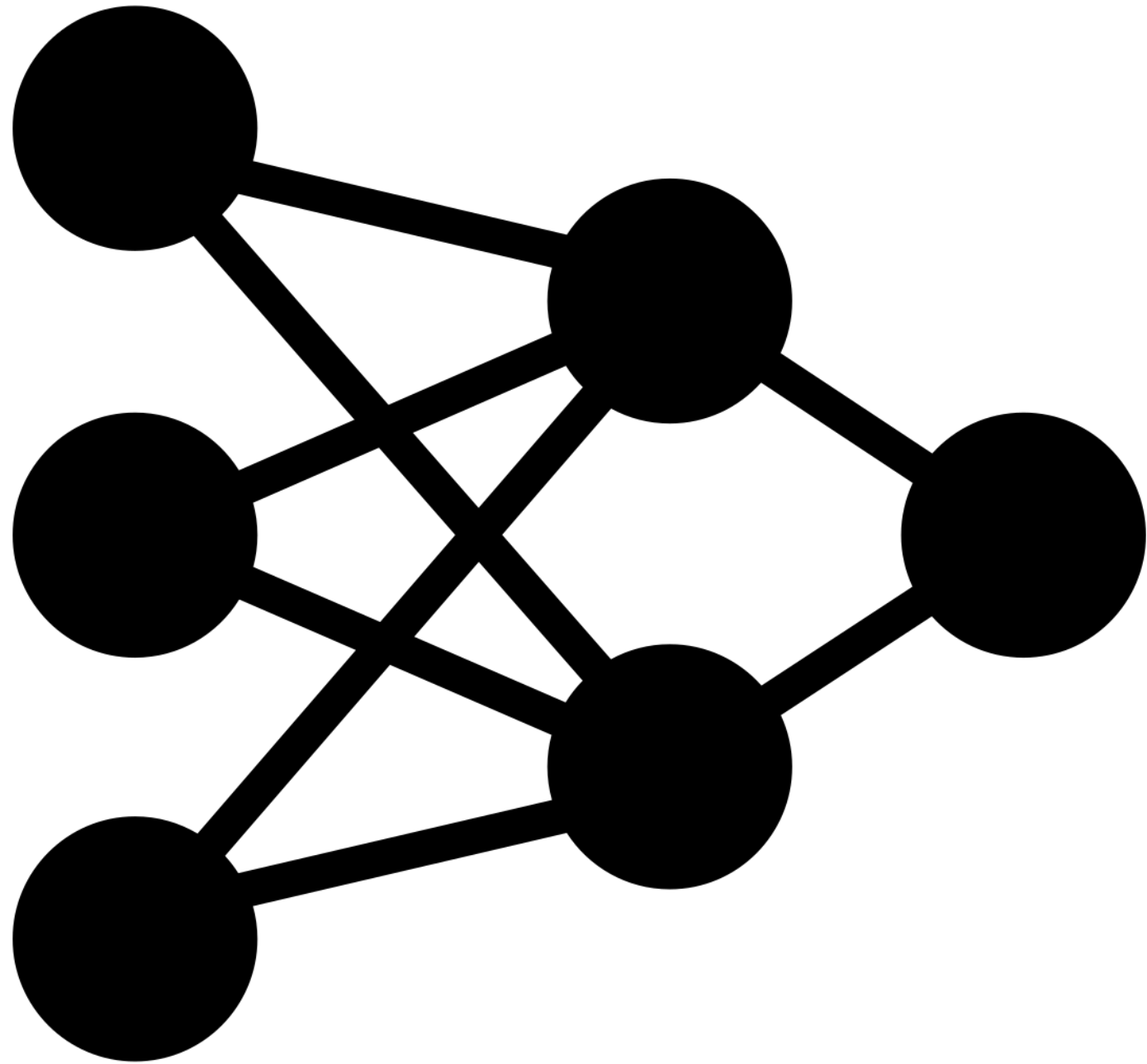


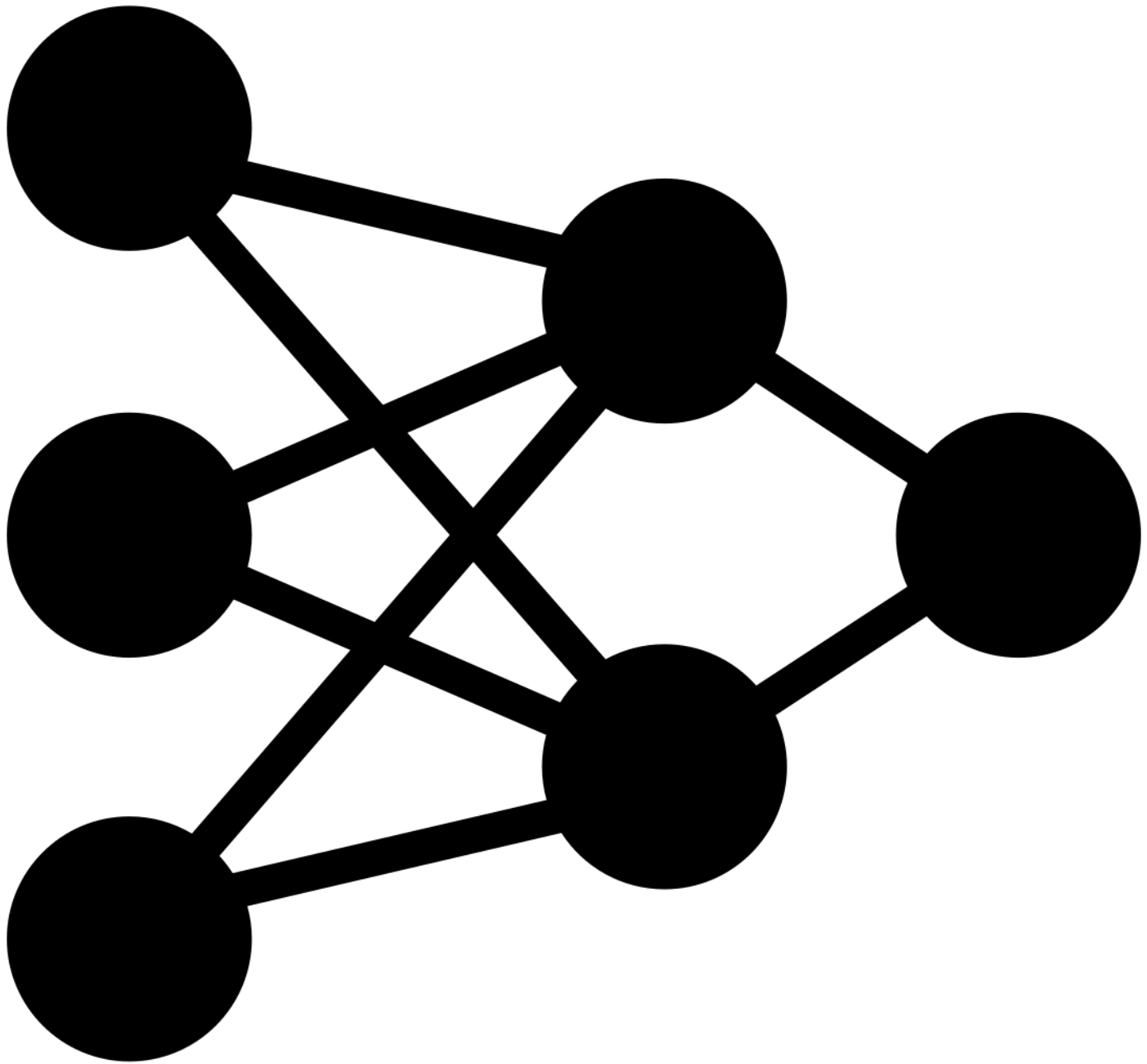
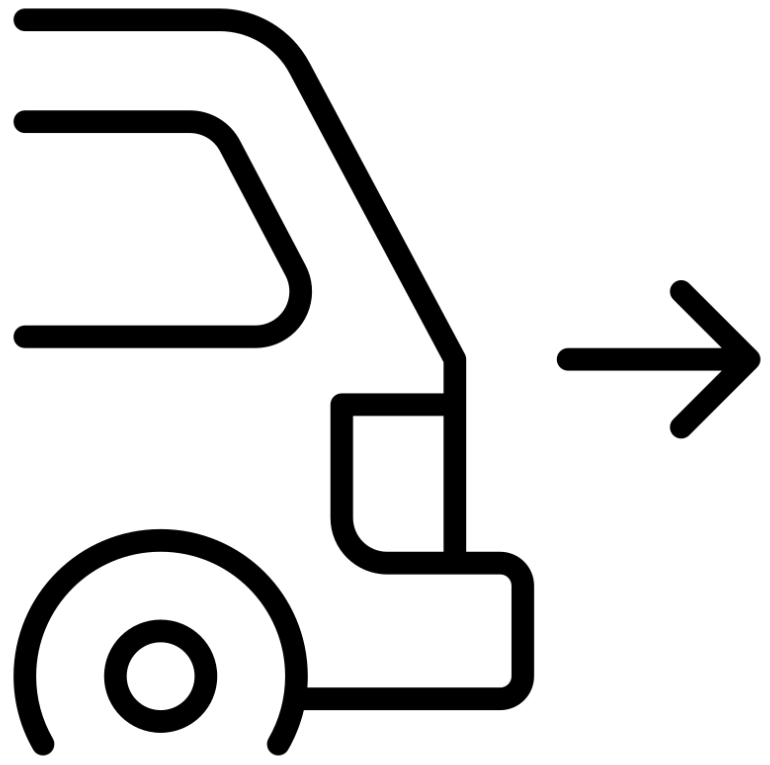
Brake?

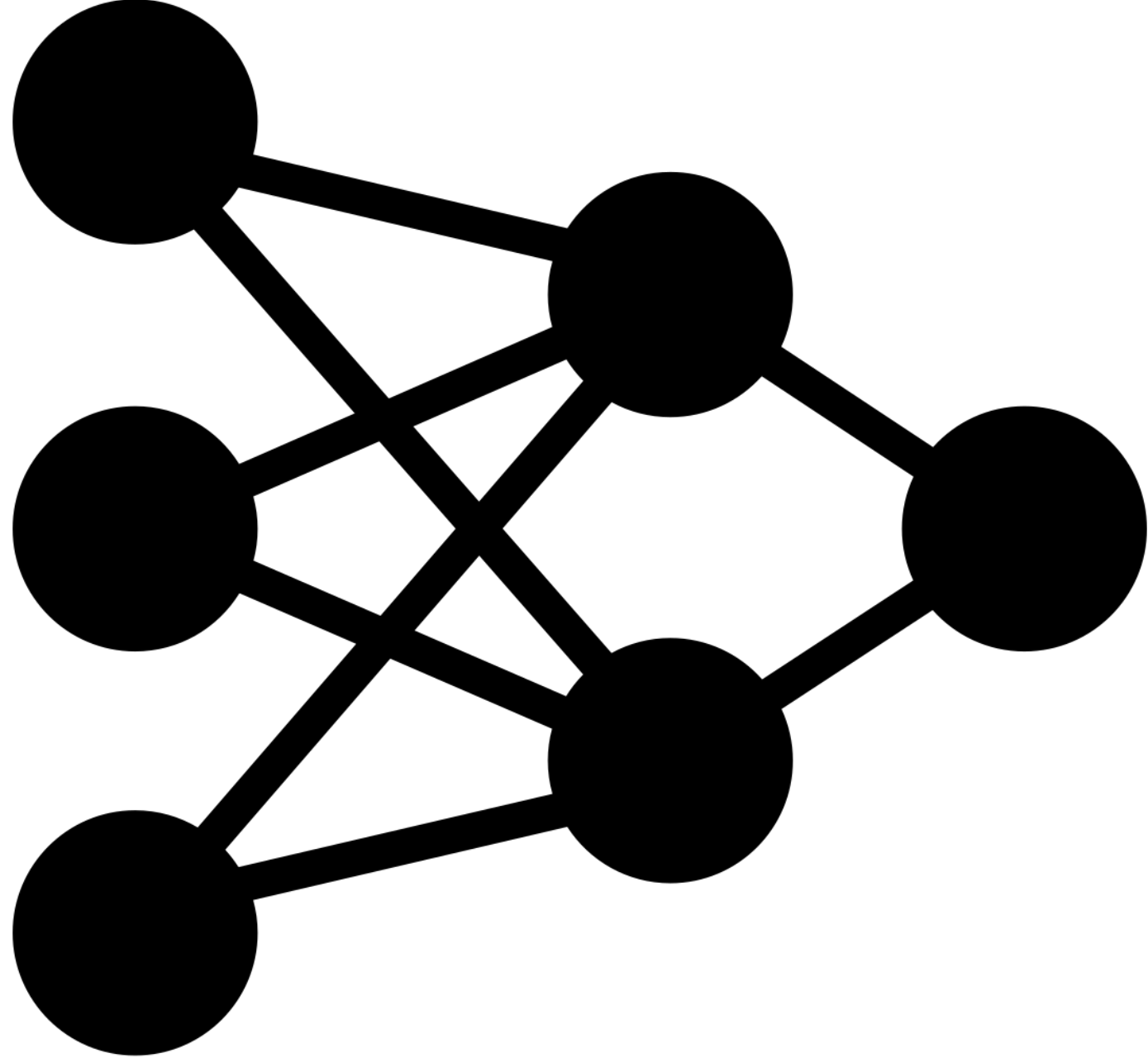
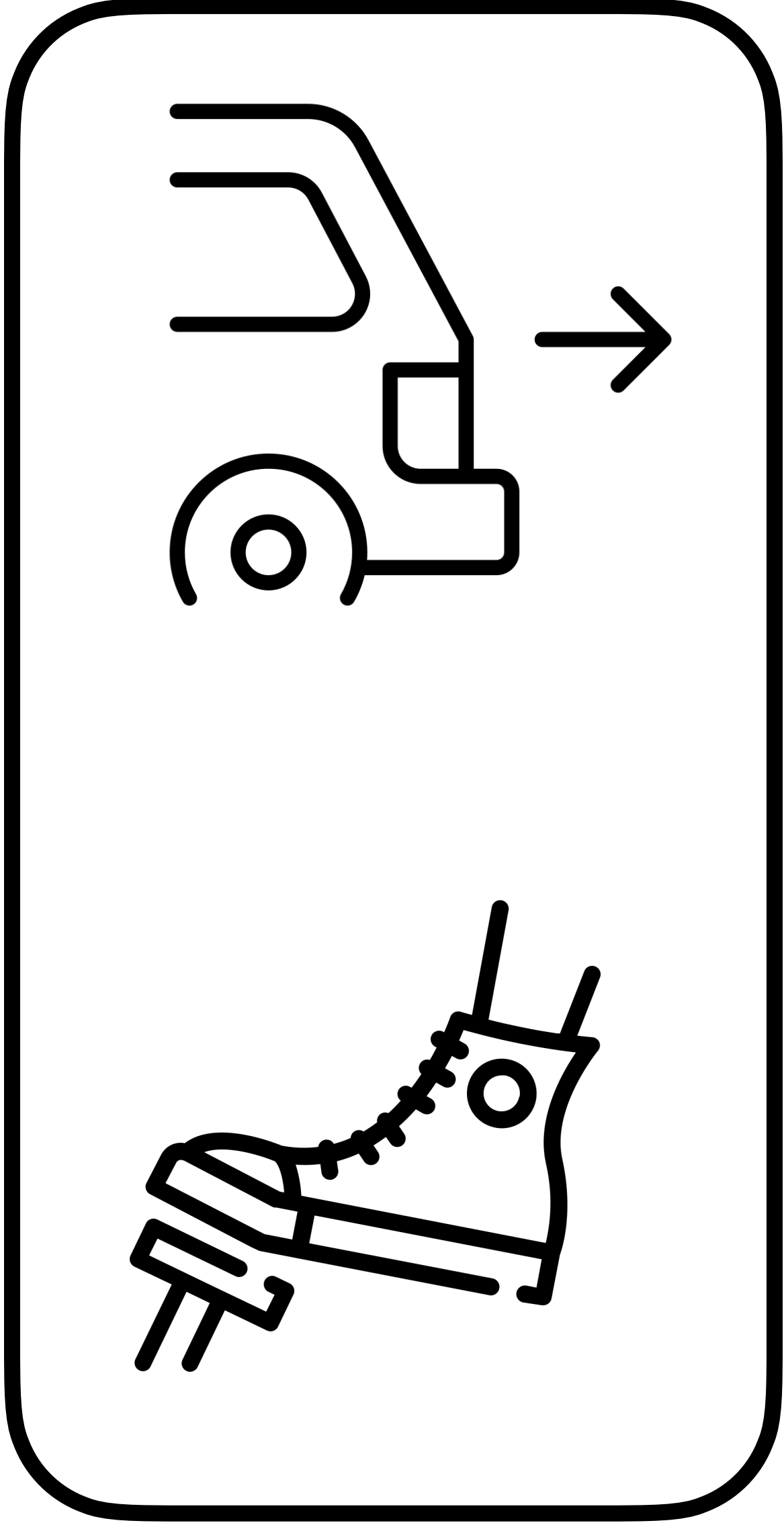


Brake?









Q: Would DAgger fix this problem?

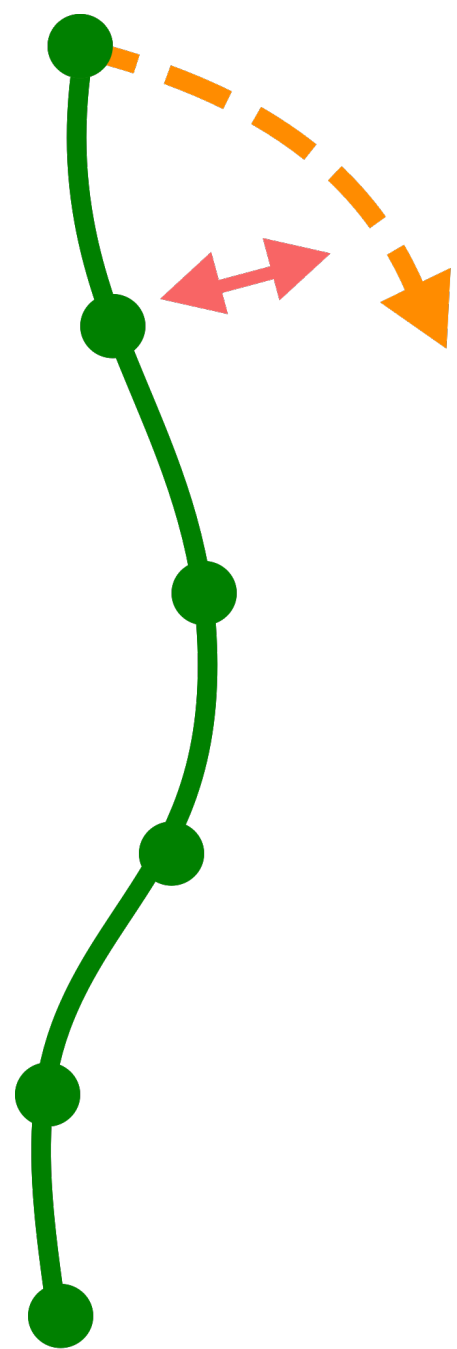
Q: Would DAgger fix this problem?

A: Yes, it's just covariate shift?

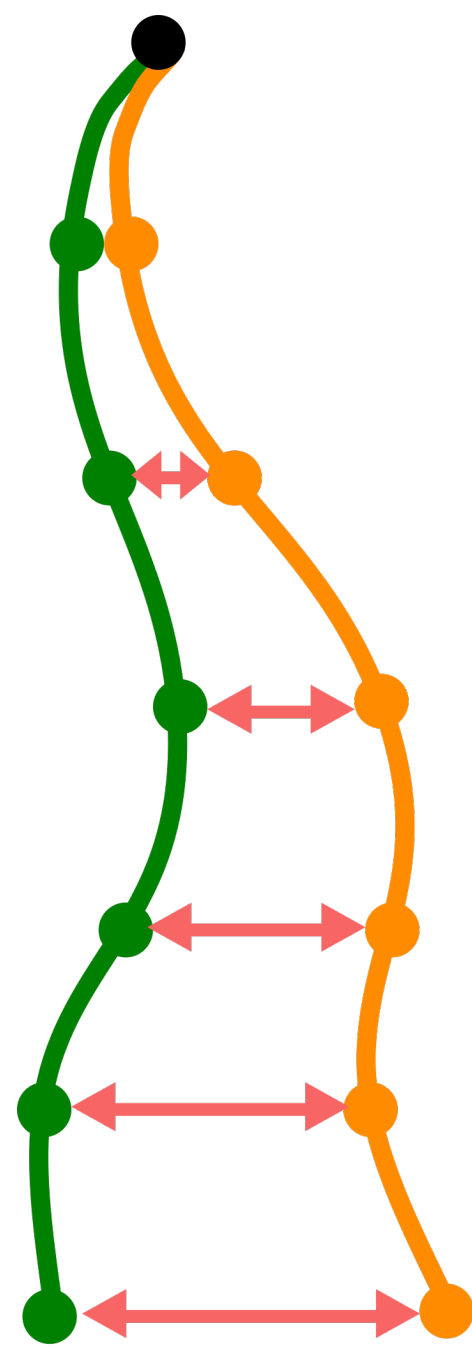
	Offline	Online	Interactive
Covariate Shift	✗	✓	✓
Hidden Context			
TCN			

$$\pi_E \xleftrightarrow{f} \pi$$

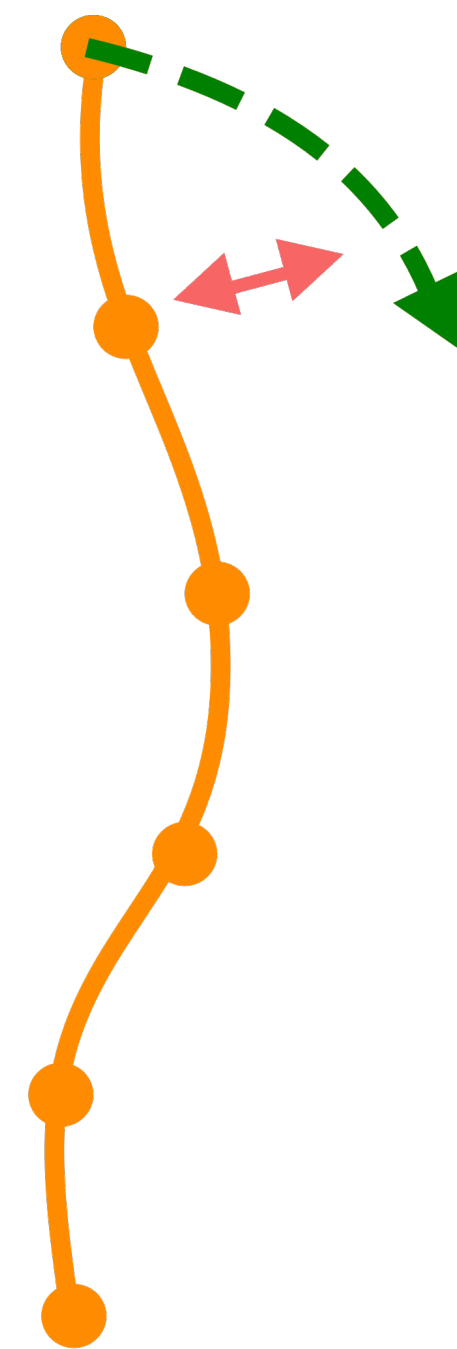
Offline



Online



Interactive



$$J(\pi_E) - J(\pi) \leq O(\epsilon T^2)$$

Behavioral Cloning ...

$$J(\pi_E) - J(\pi) \leq O(\epsilon T)$$

GAIL, MaxEnt IRL ...

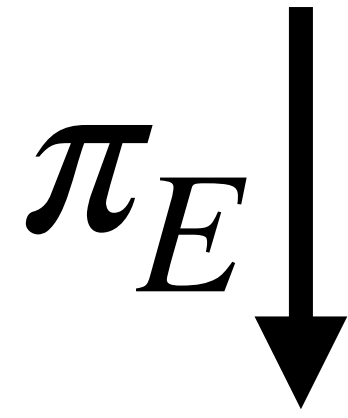
$$J(\pi_E) - J(\pi) \leq O(\epsilon HT)$$

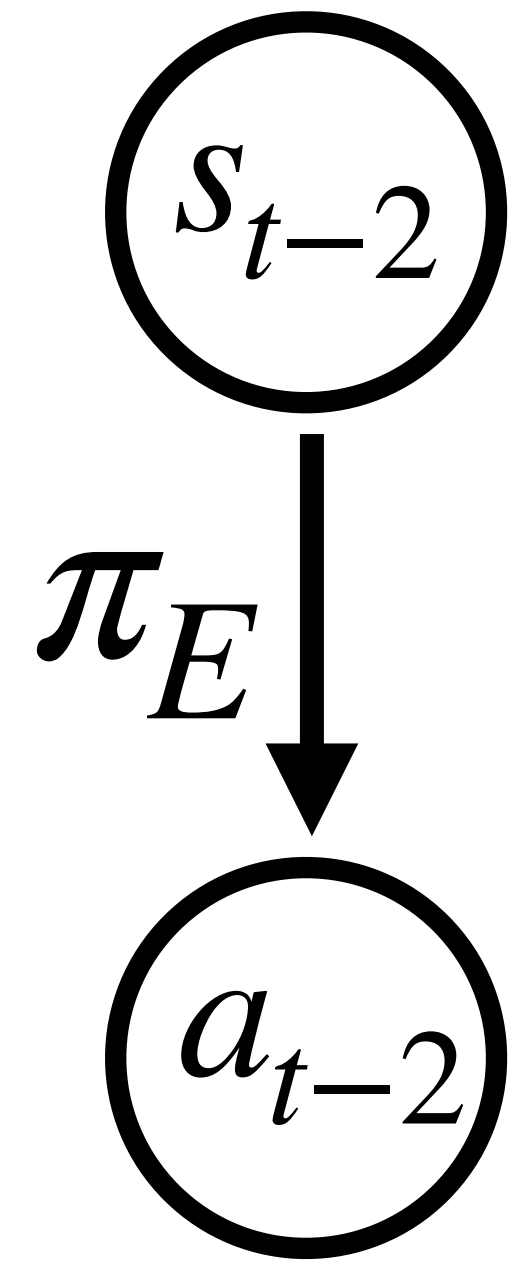
Dagger ...

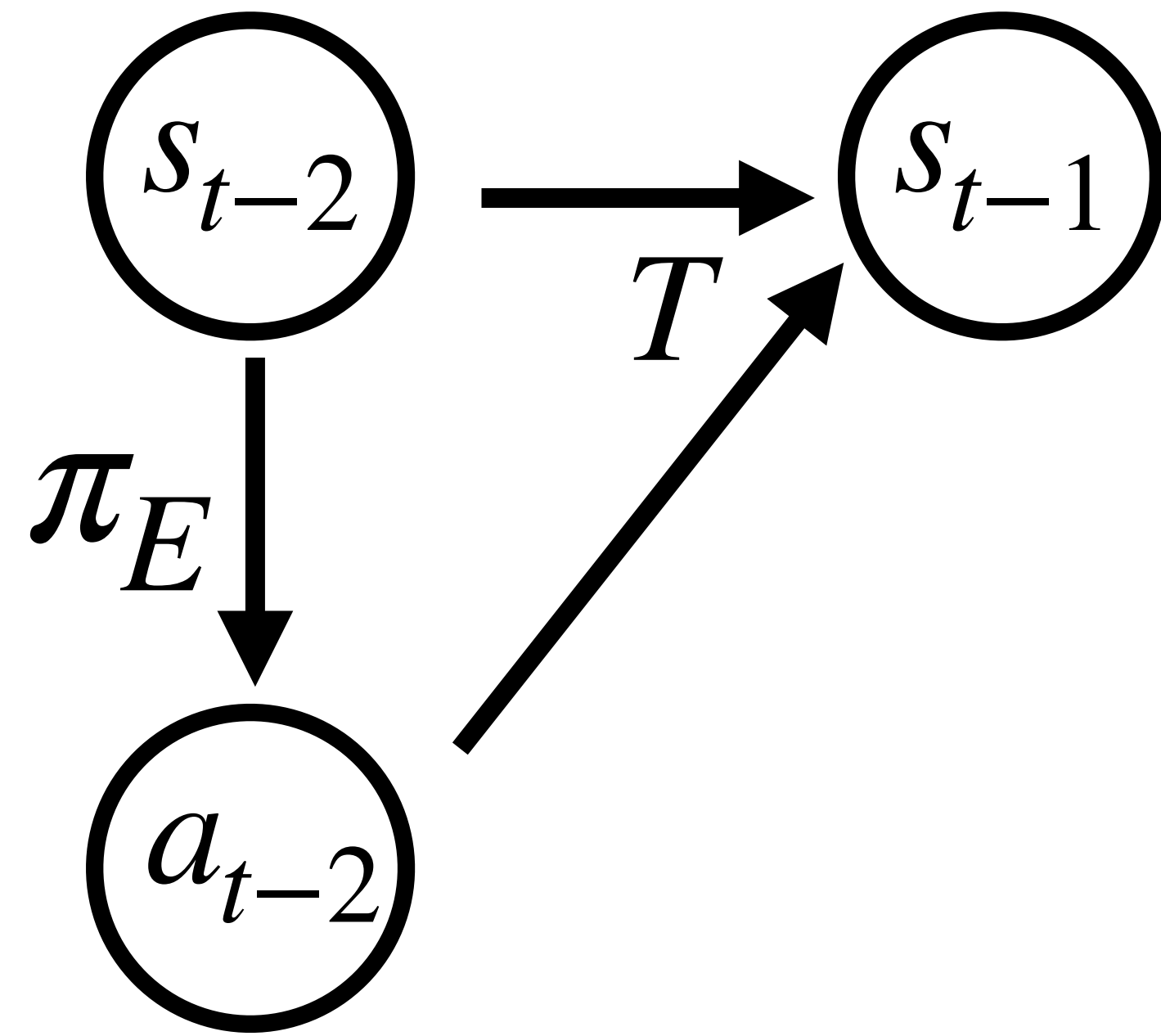
“Hence, a system trained with multiple frames would merely predict a steering angle equal to the current rate of turn as observed through the camera. This would lead to catastrophic behavior in test mode. The robot would simply turn in circles.”
— Muller et al., 2006

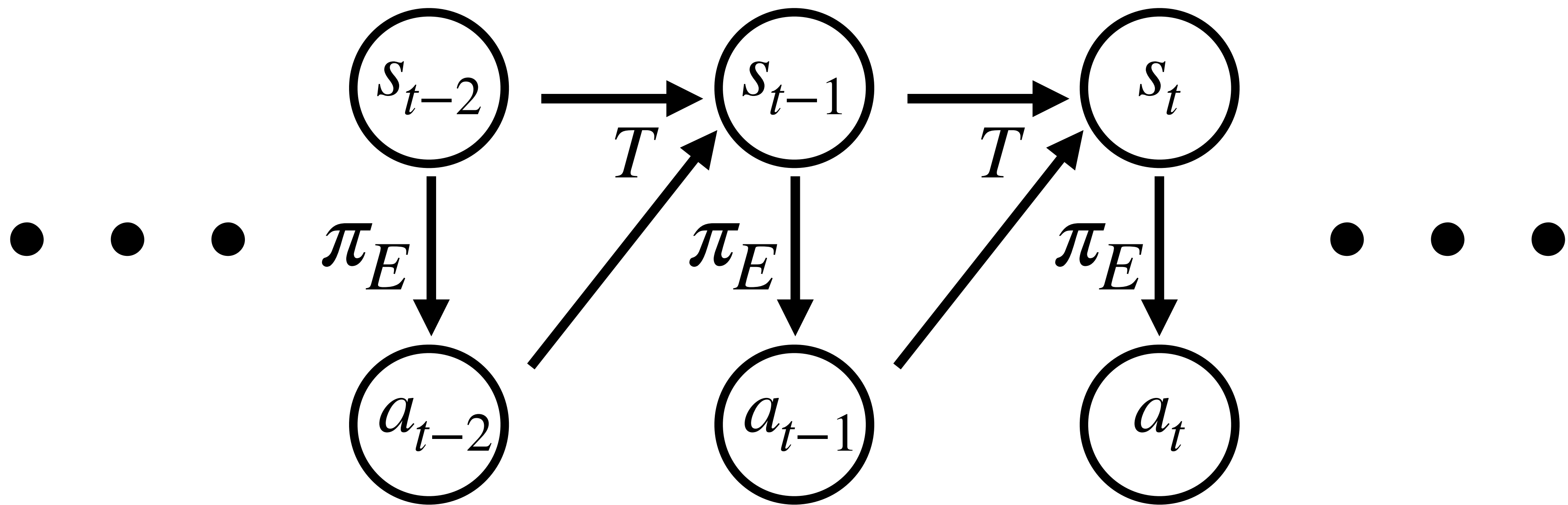
$$S_{t-2}$$

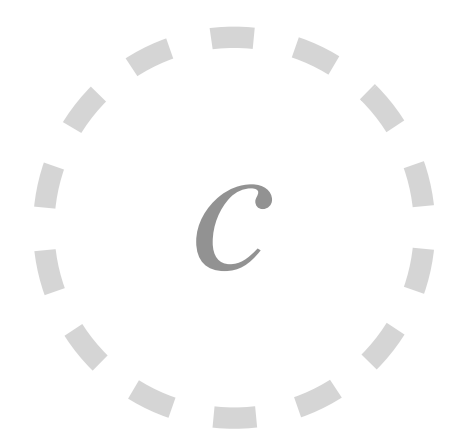
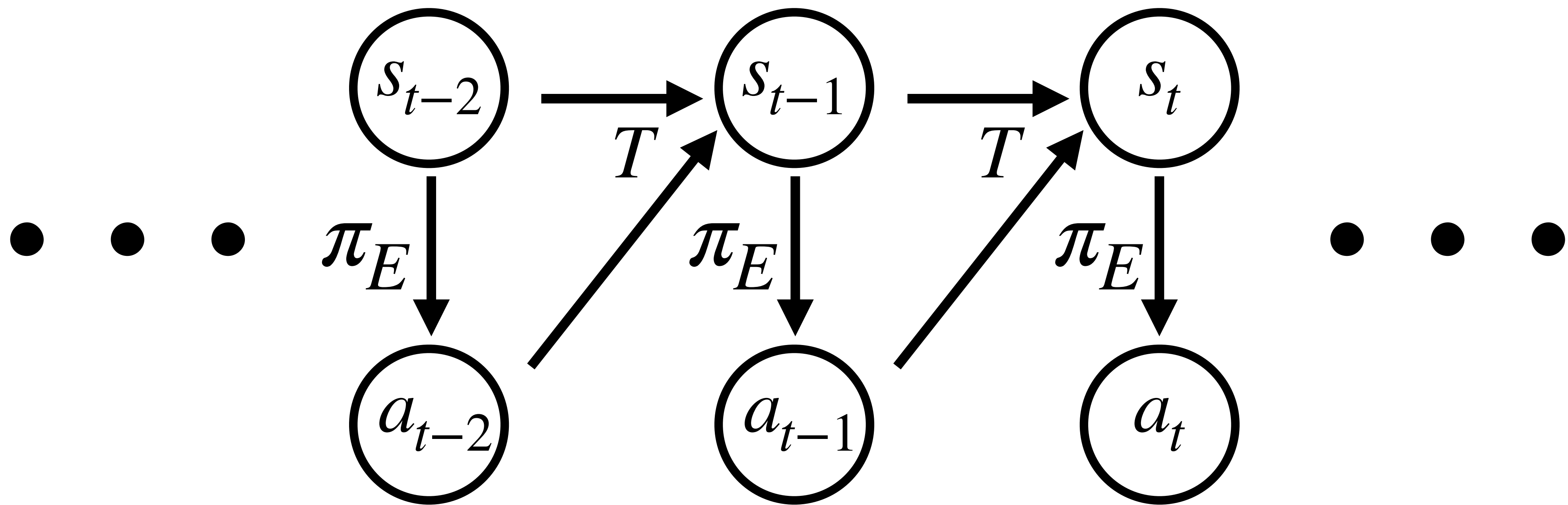
S_{t-2}

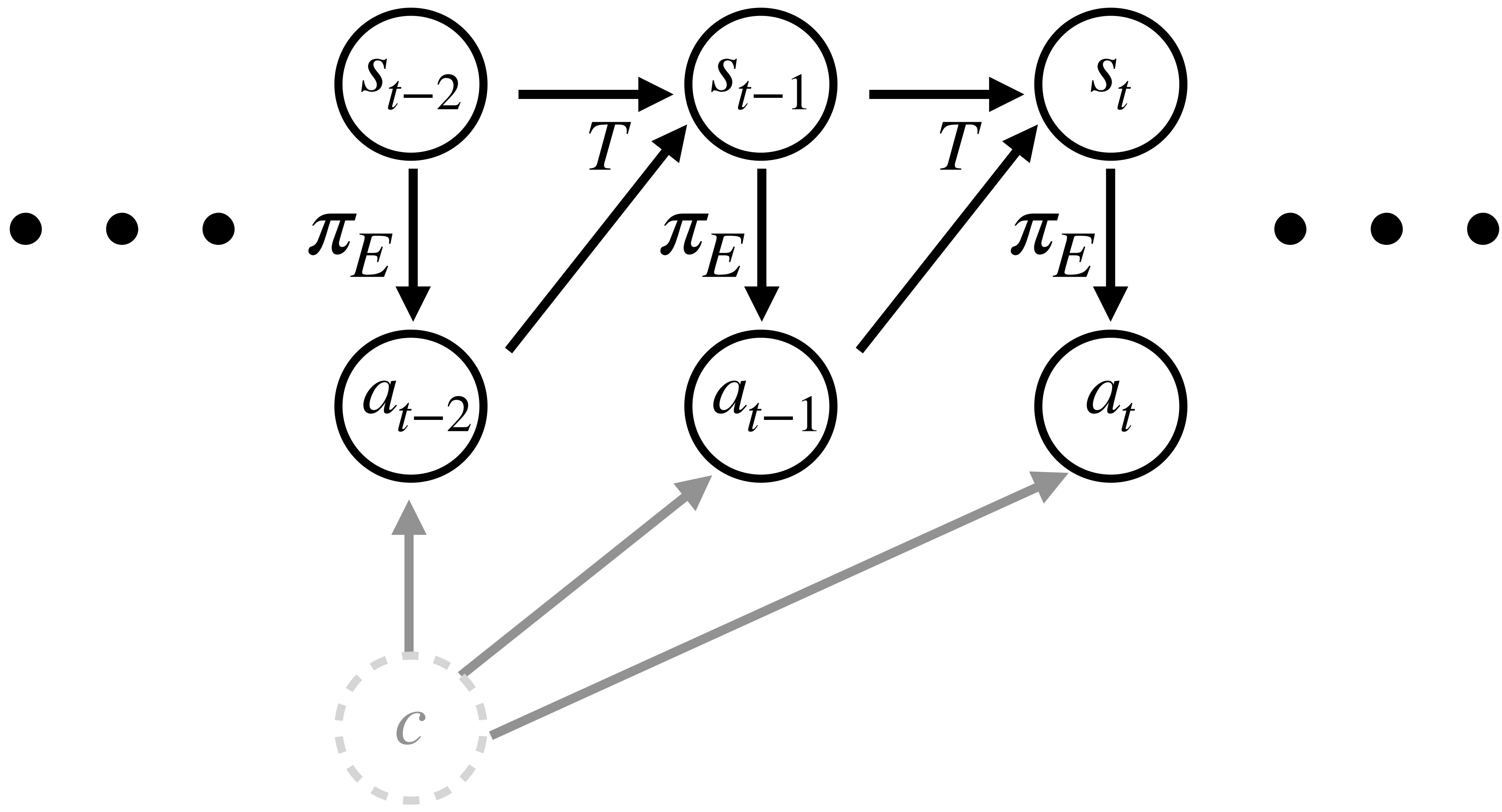


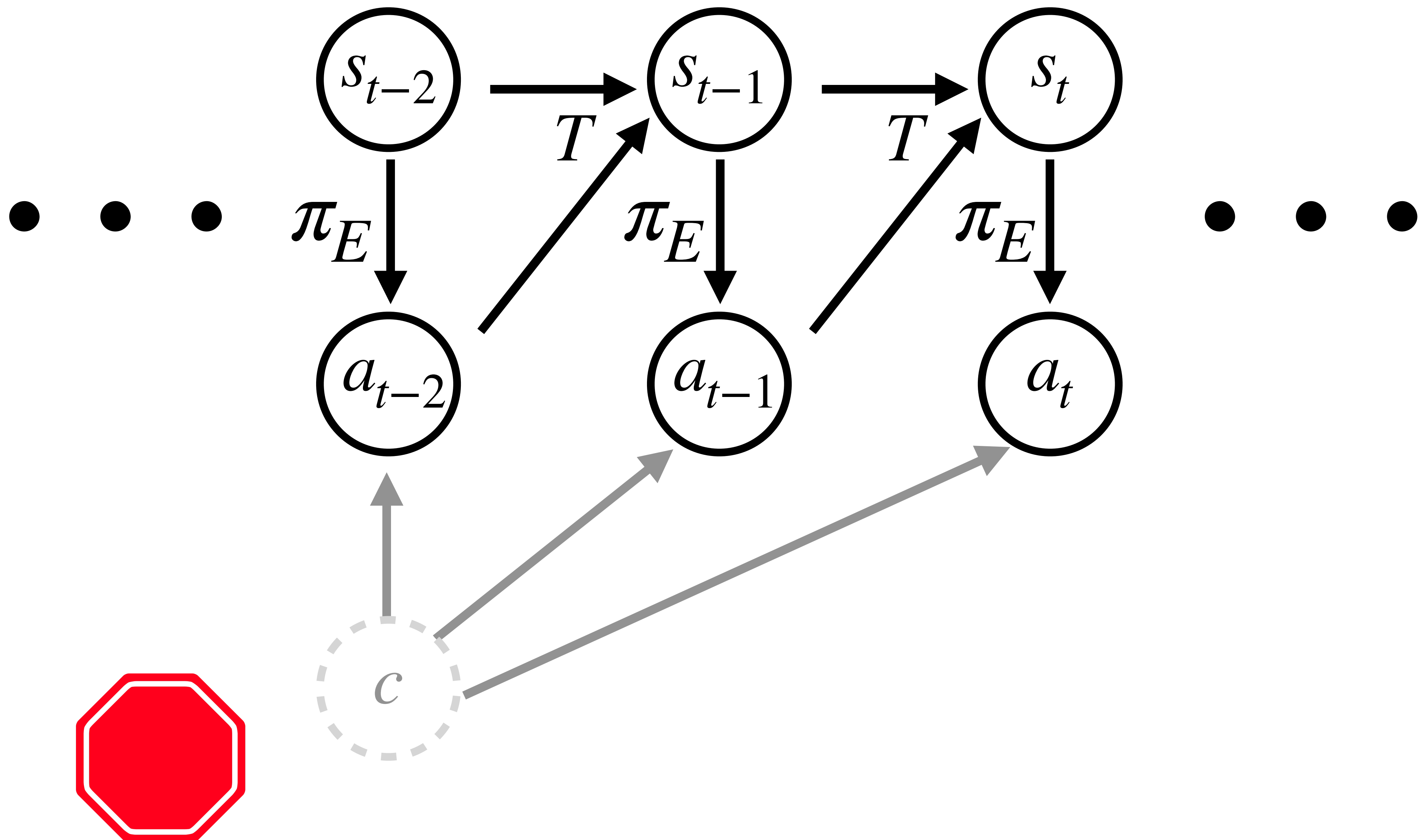


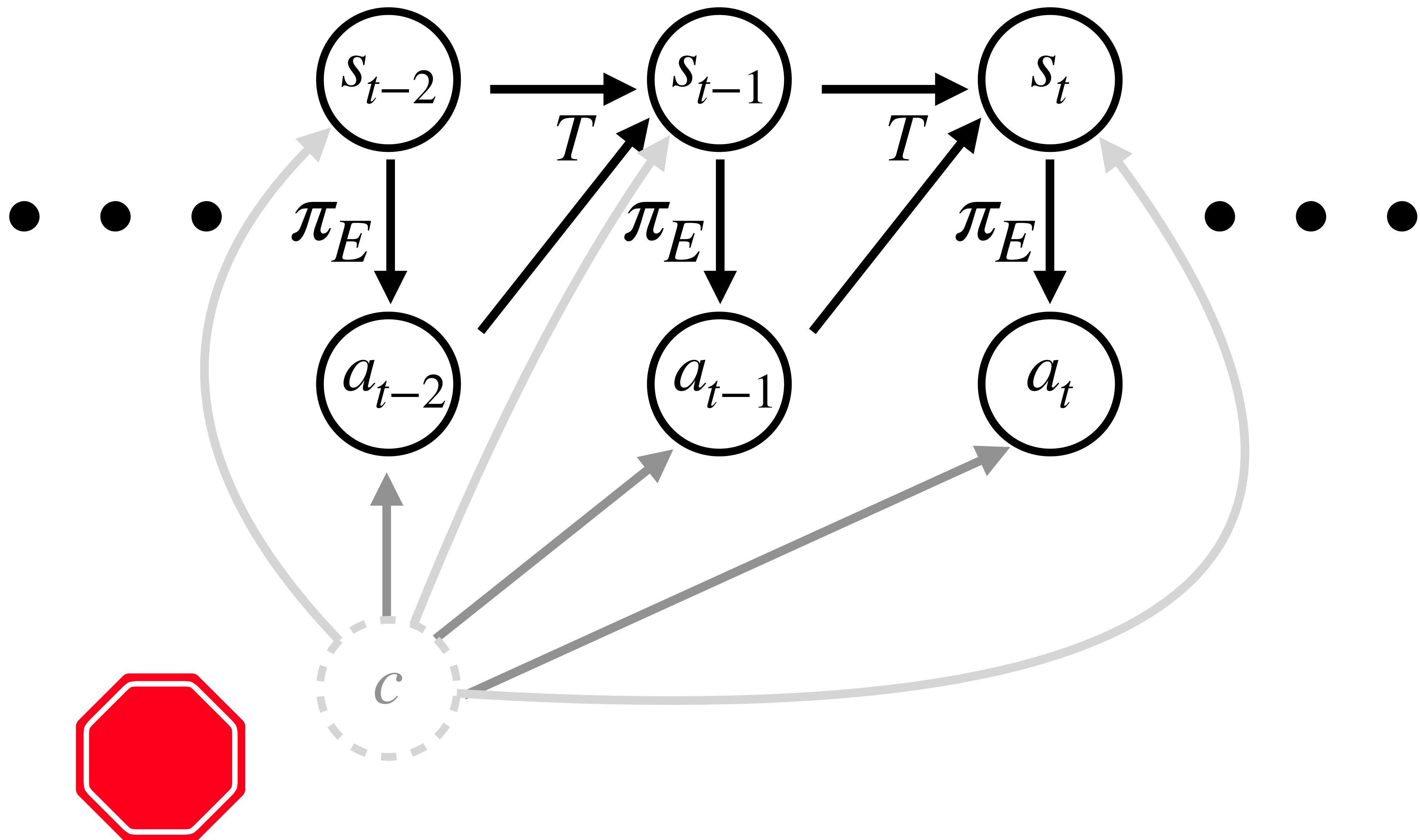


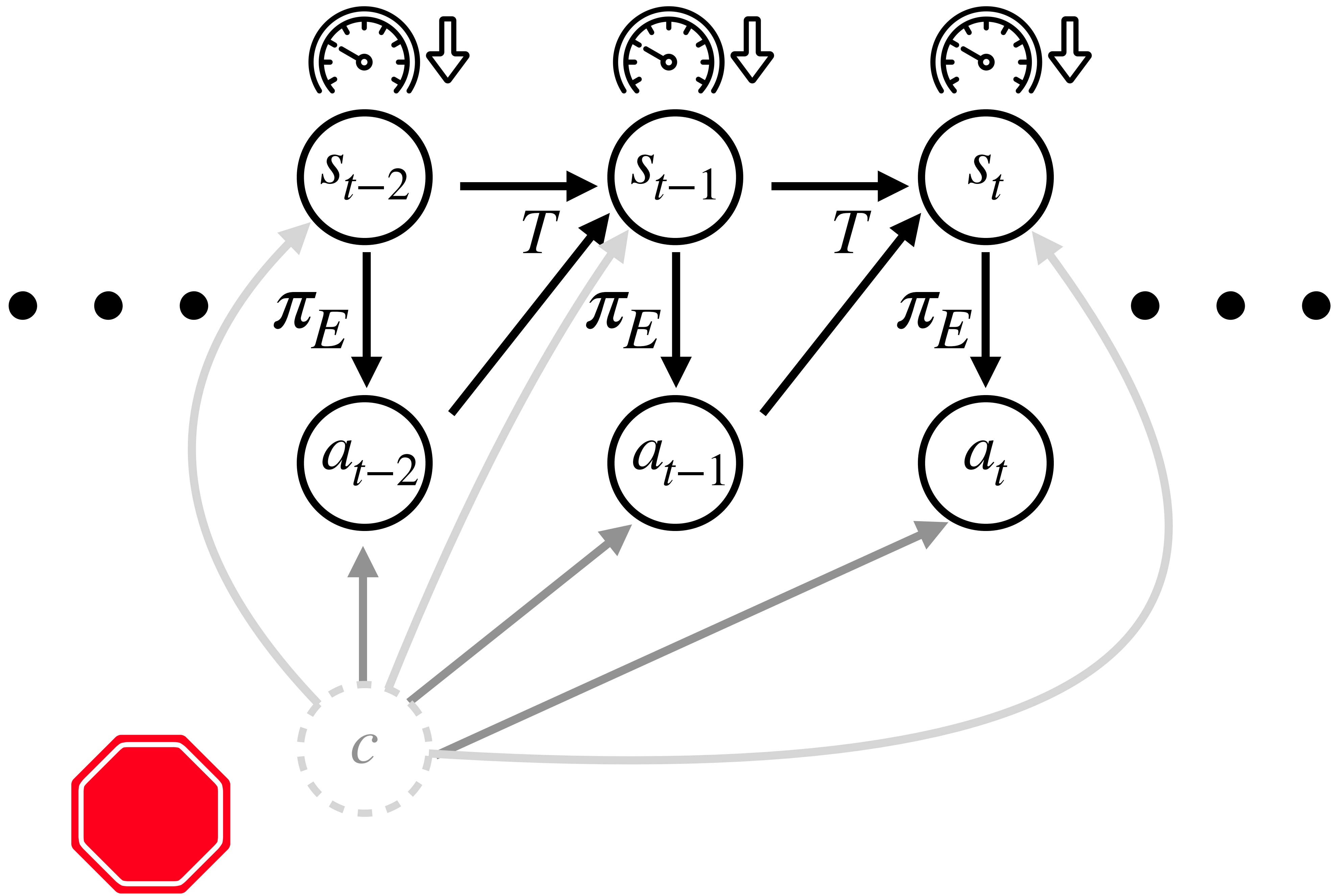












MDP

POMDP

State

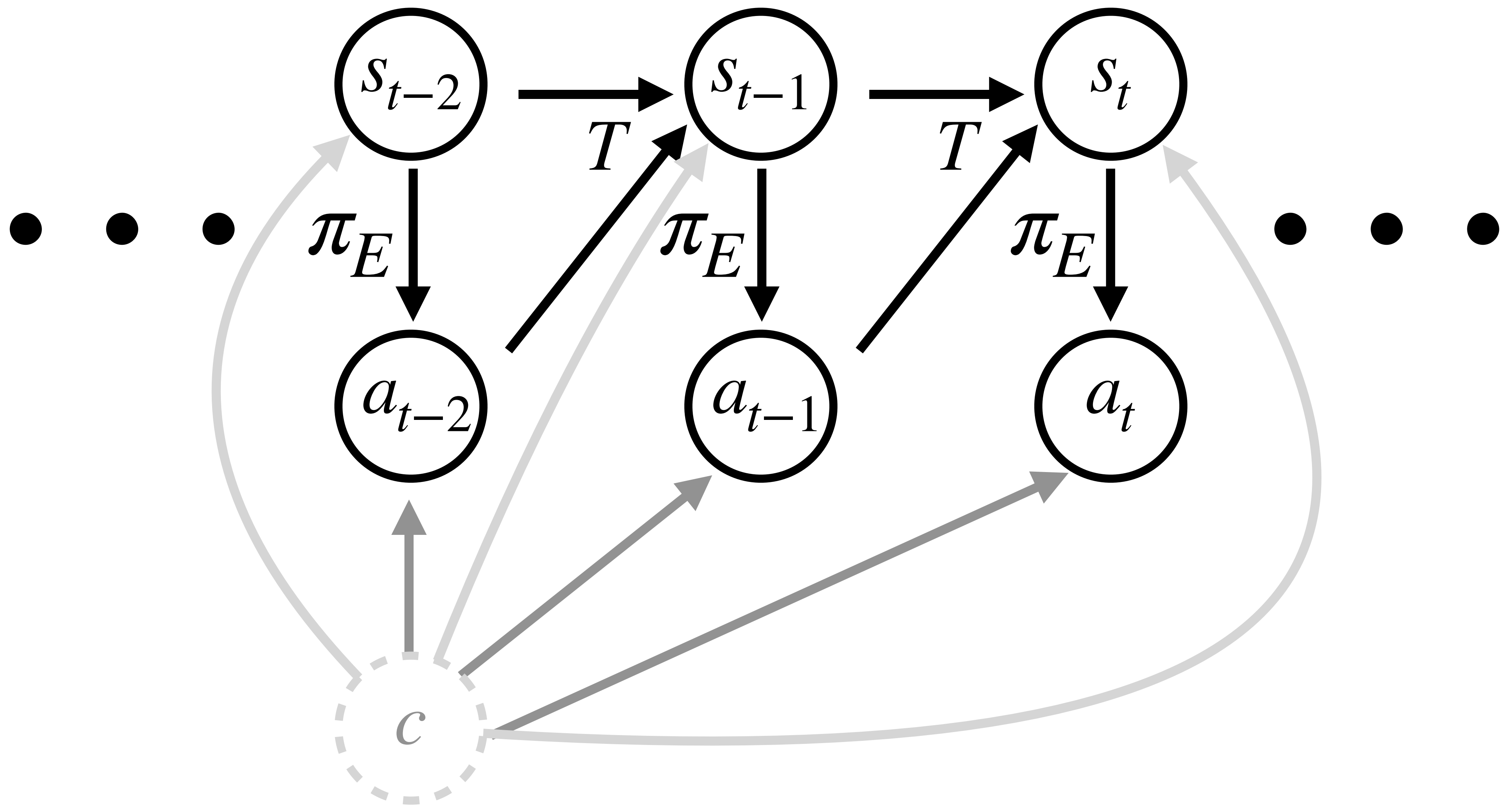
s_t

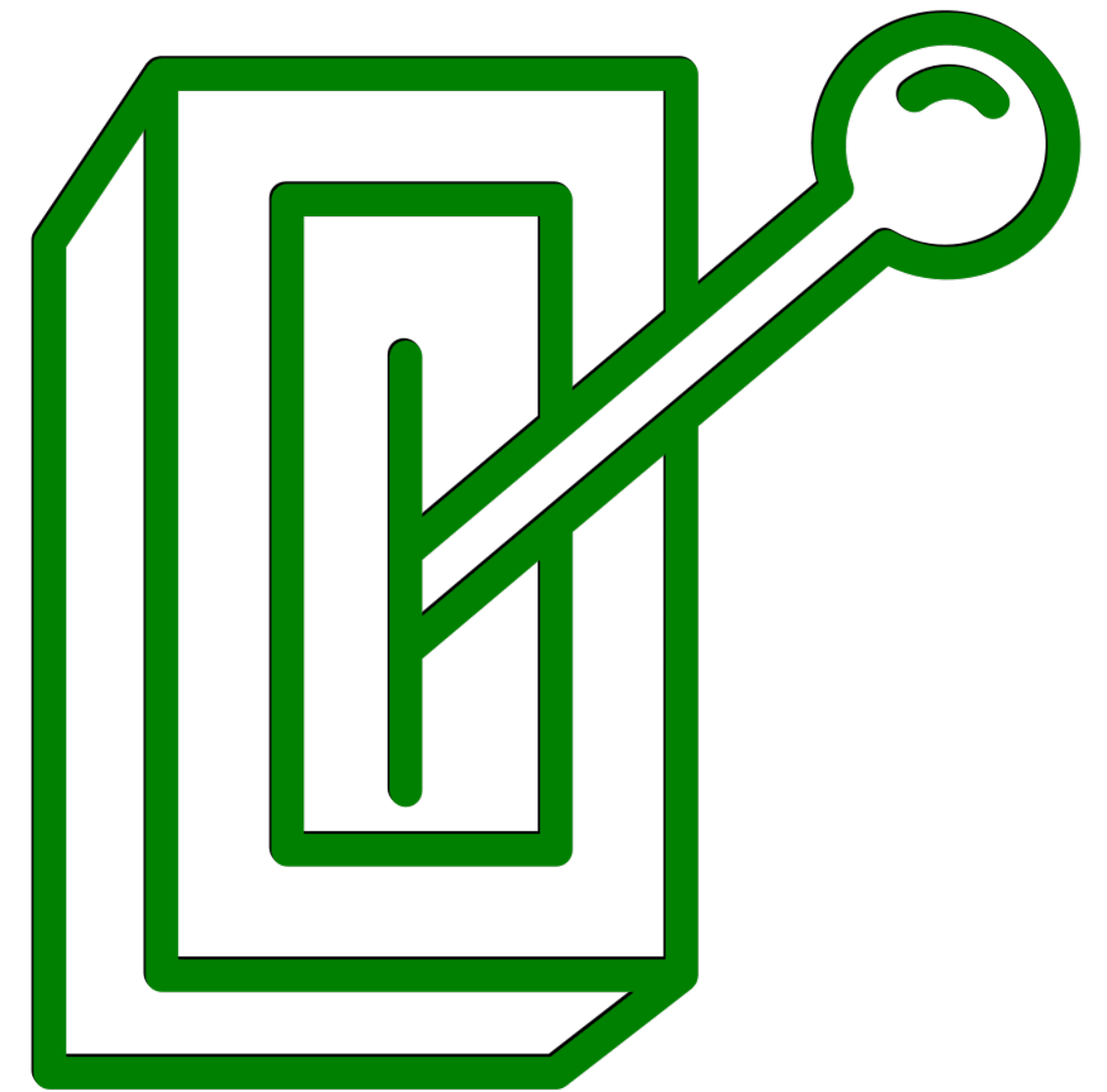
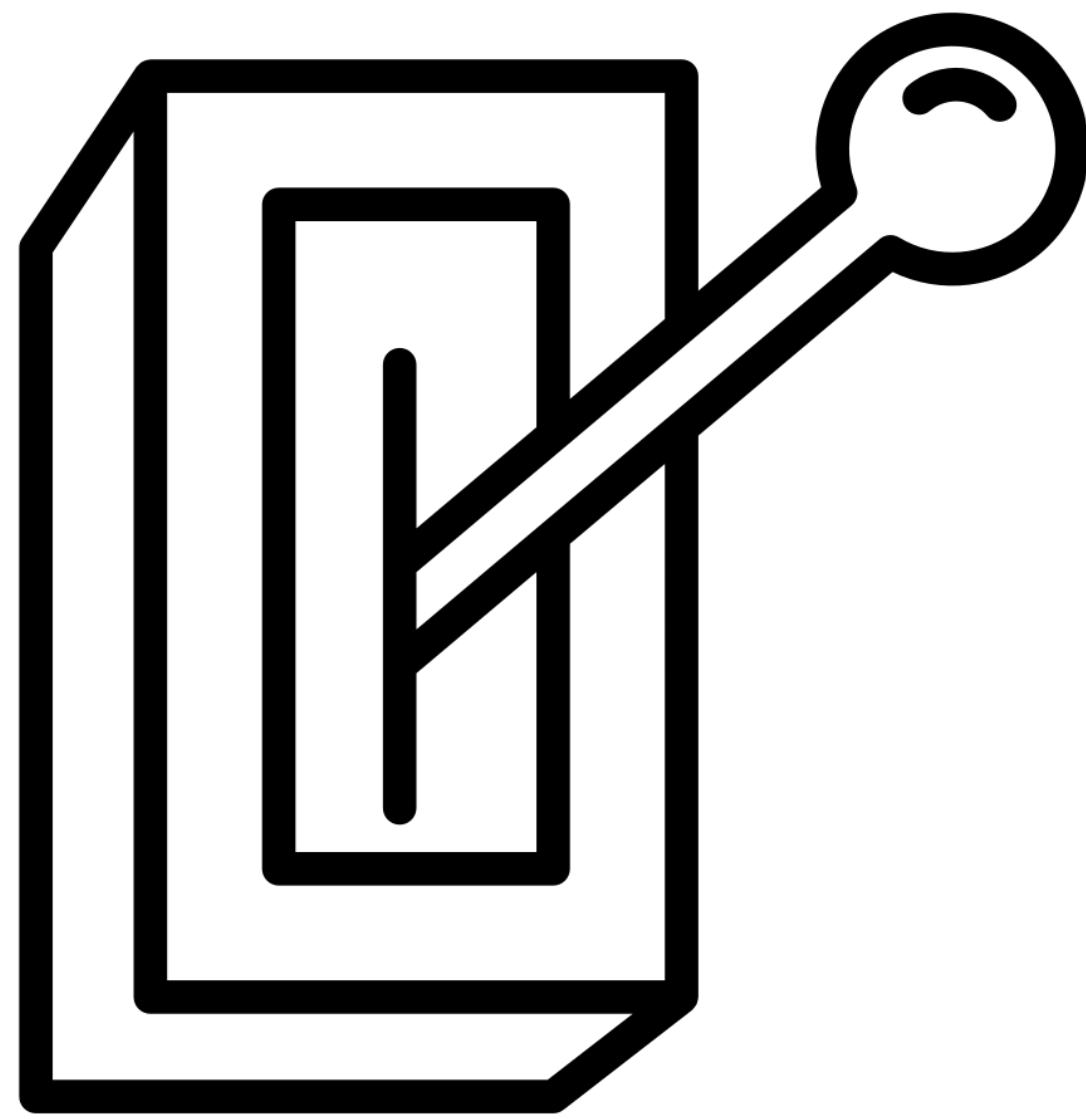
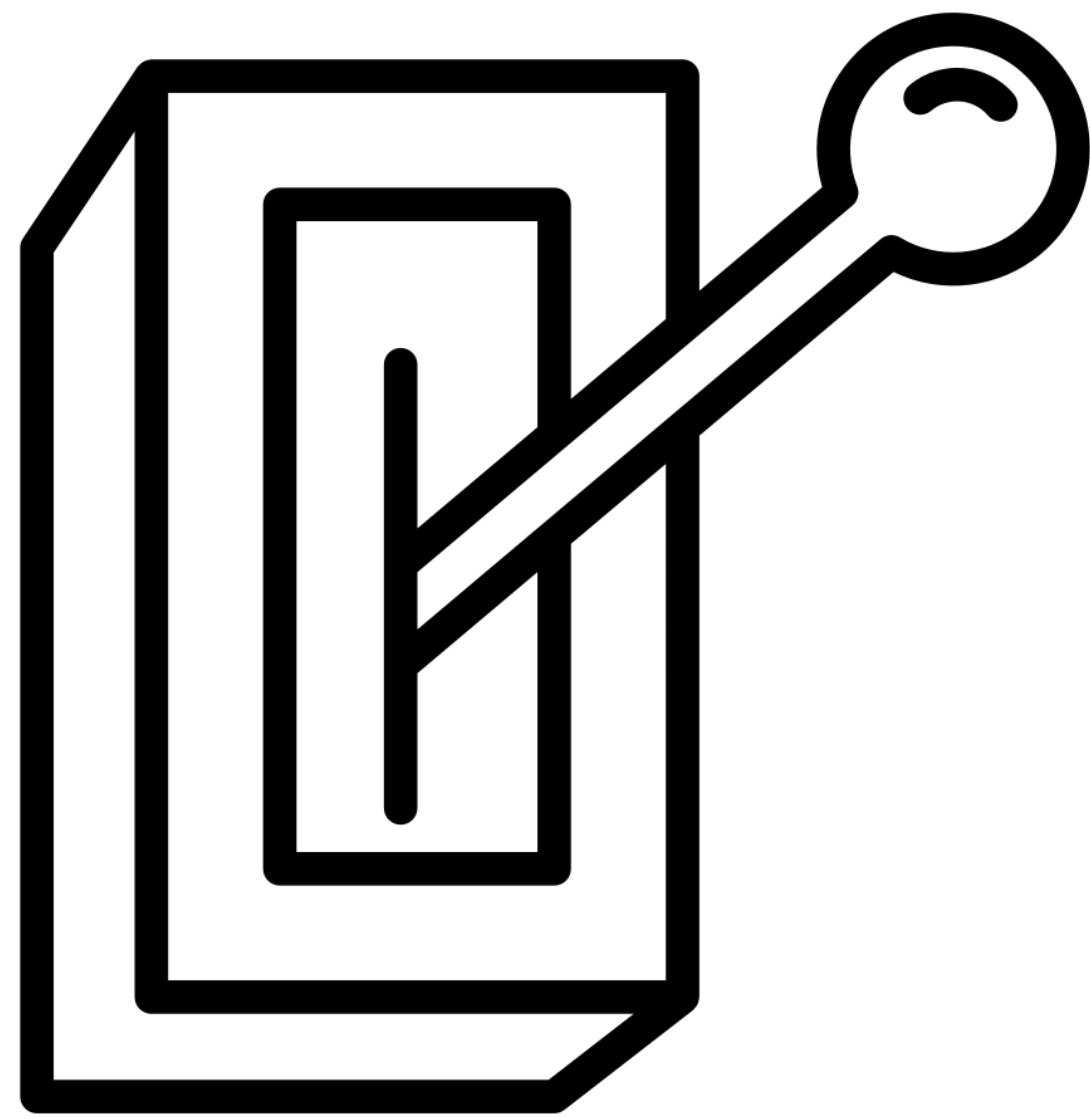
Policy

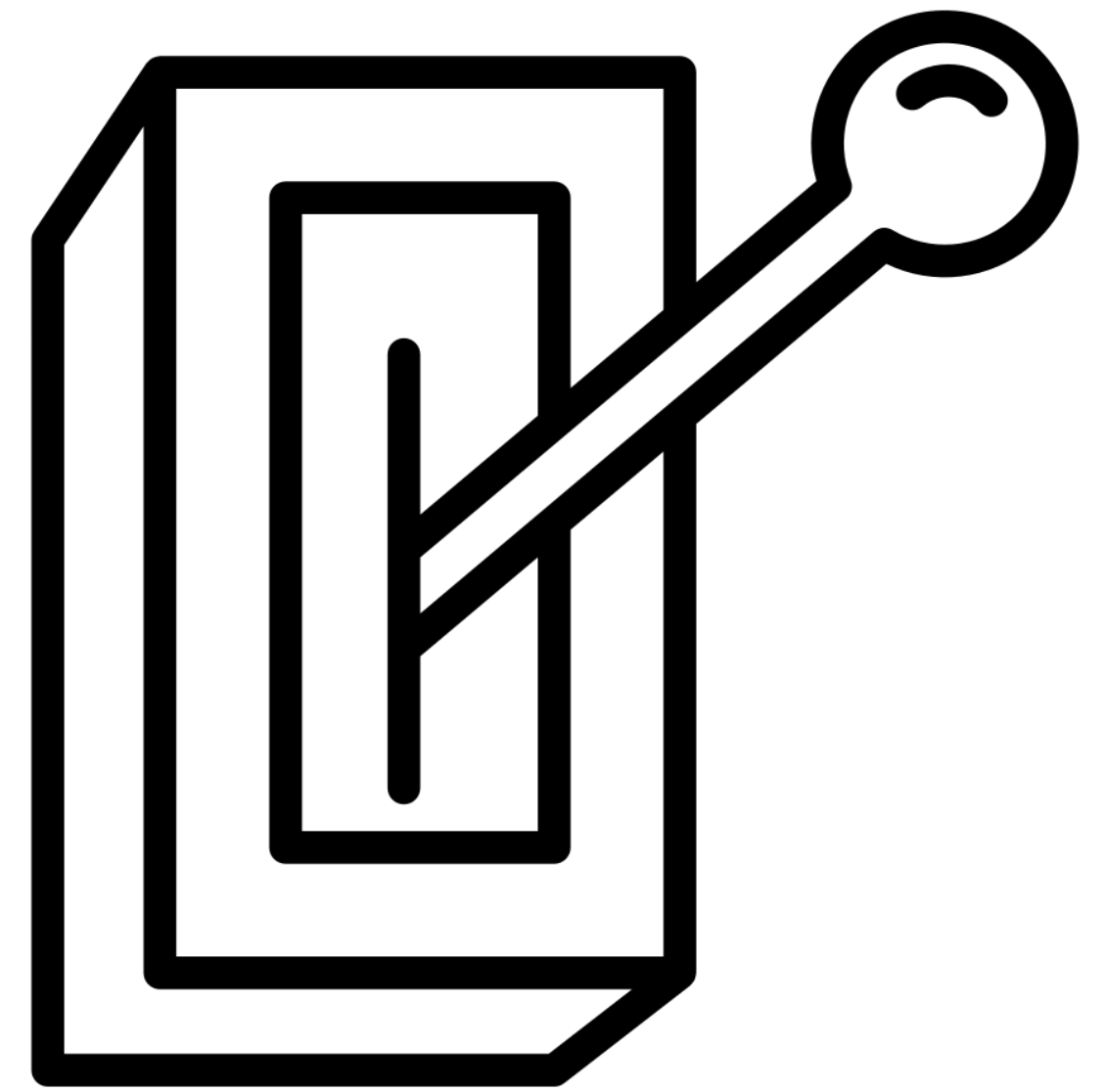
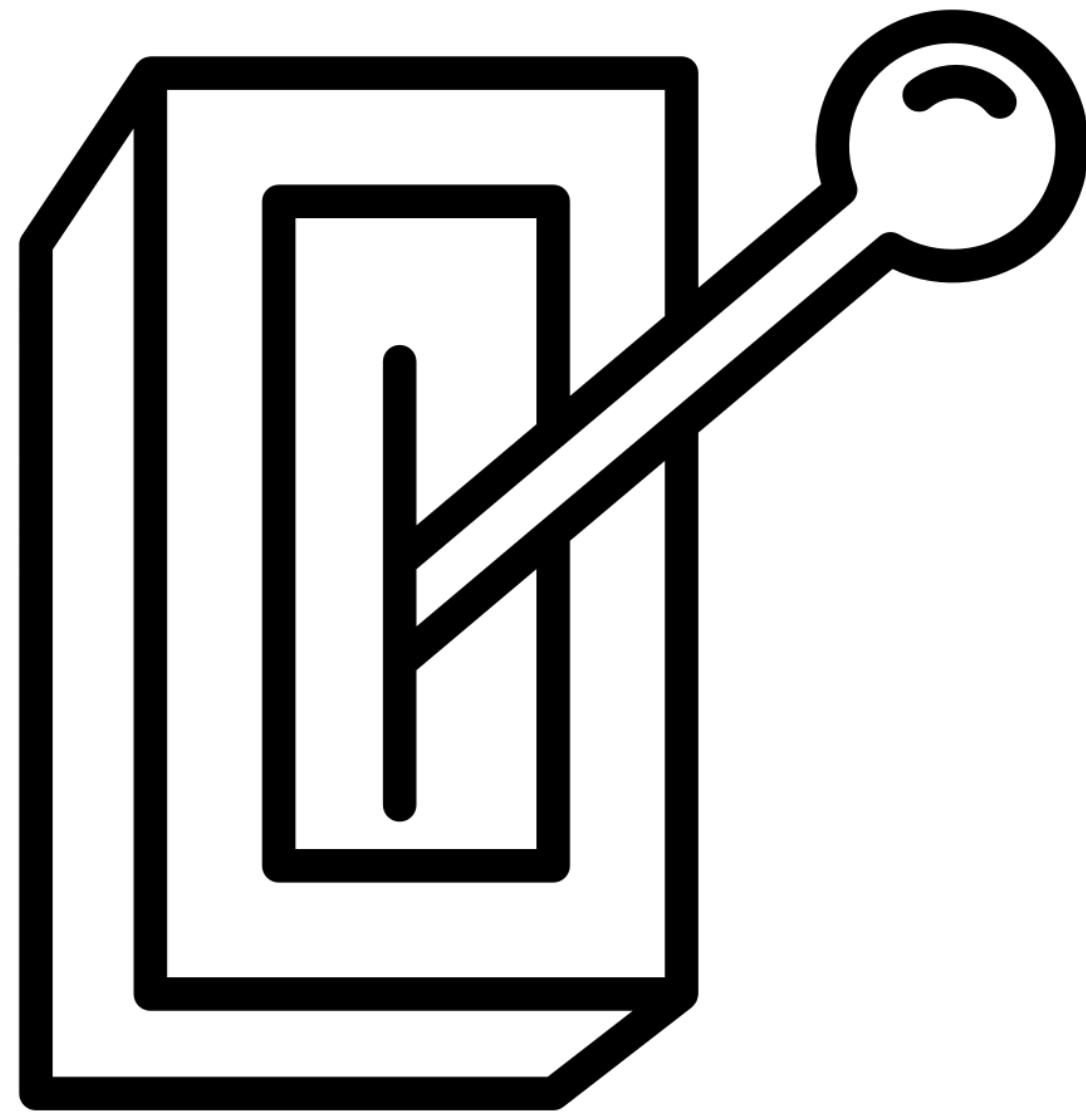
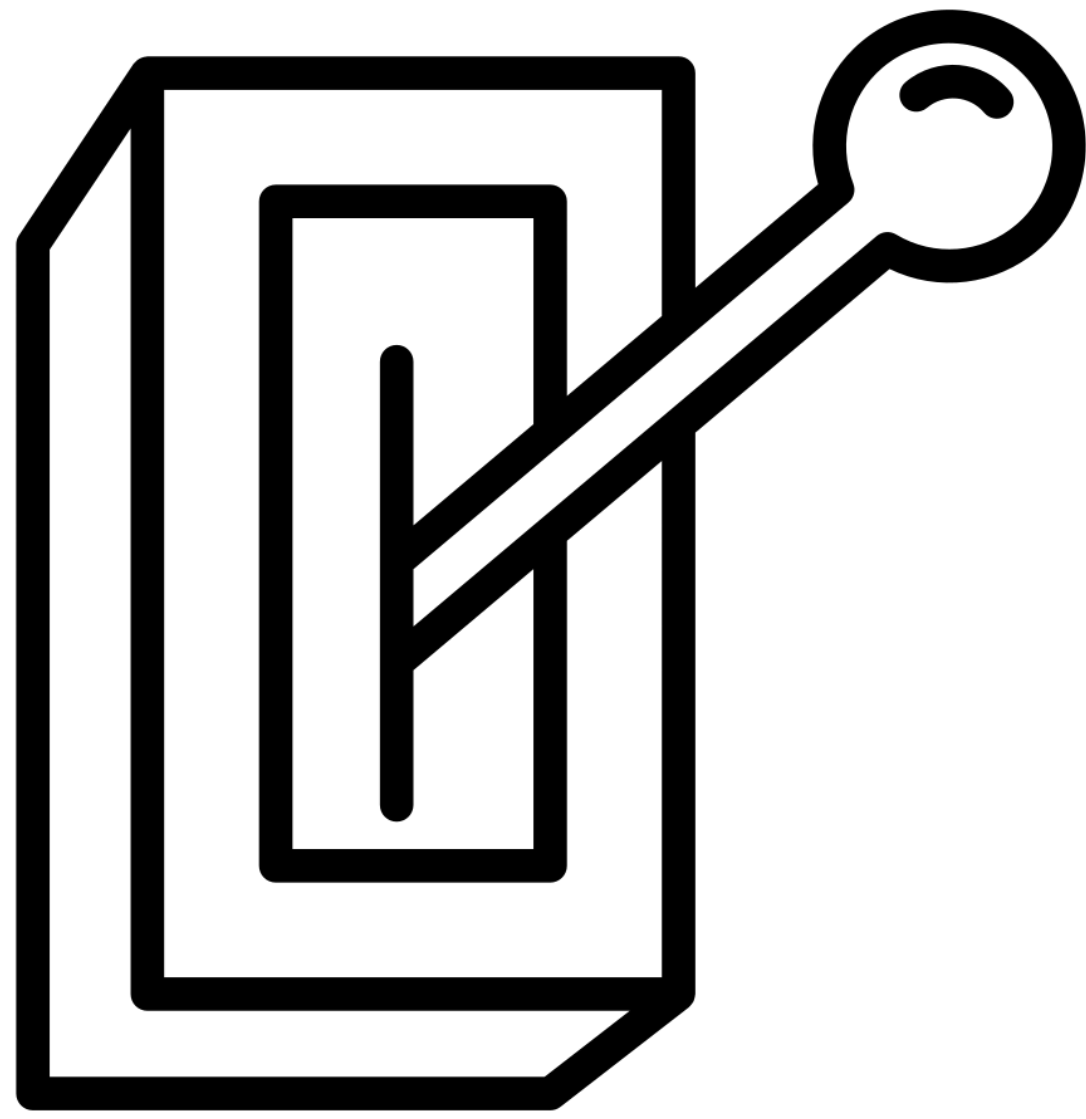
$\pi(\cdot | s_t)$

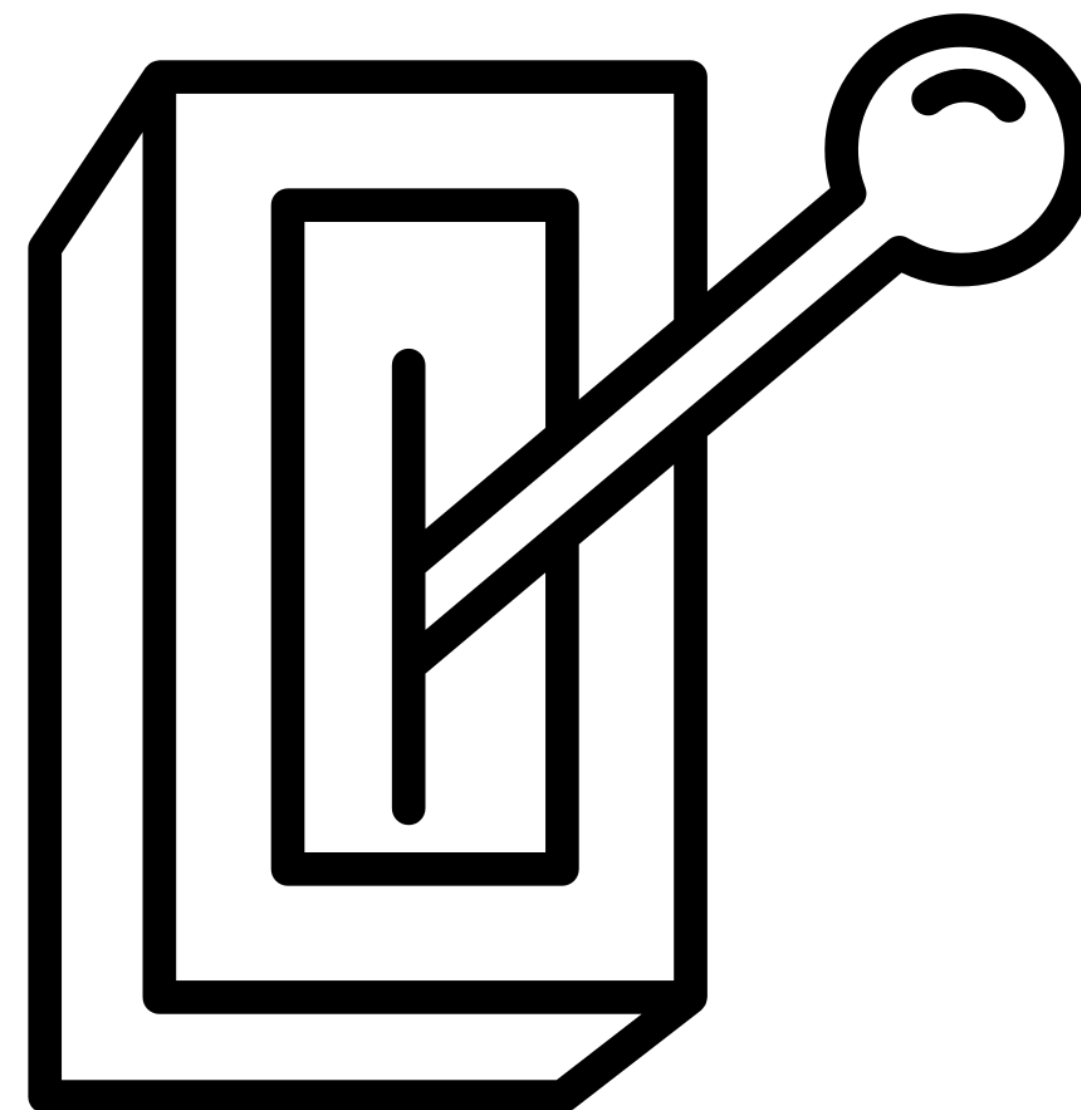
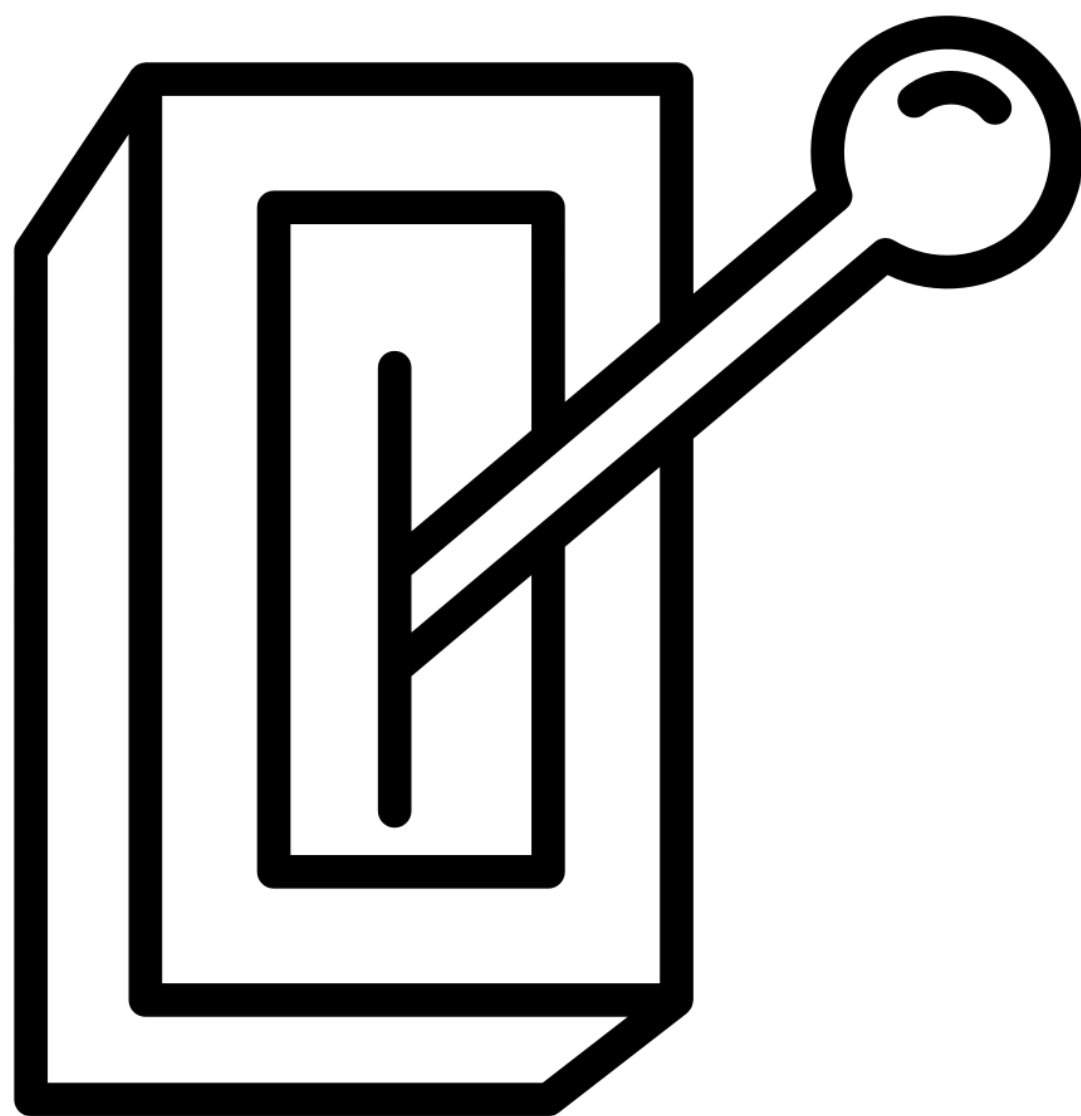
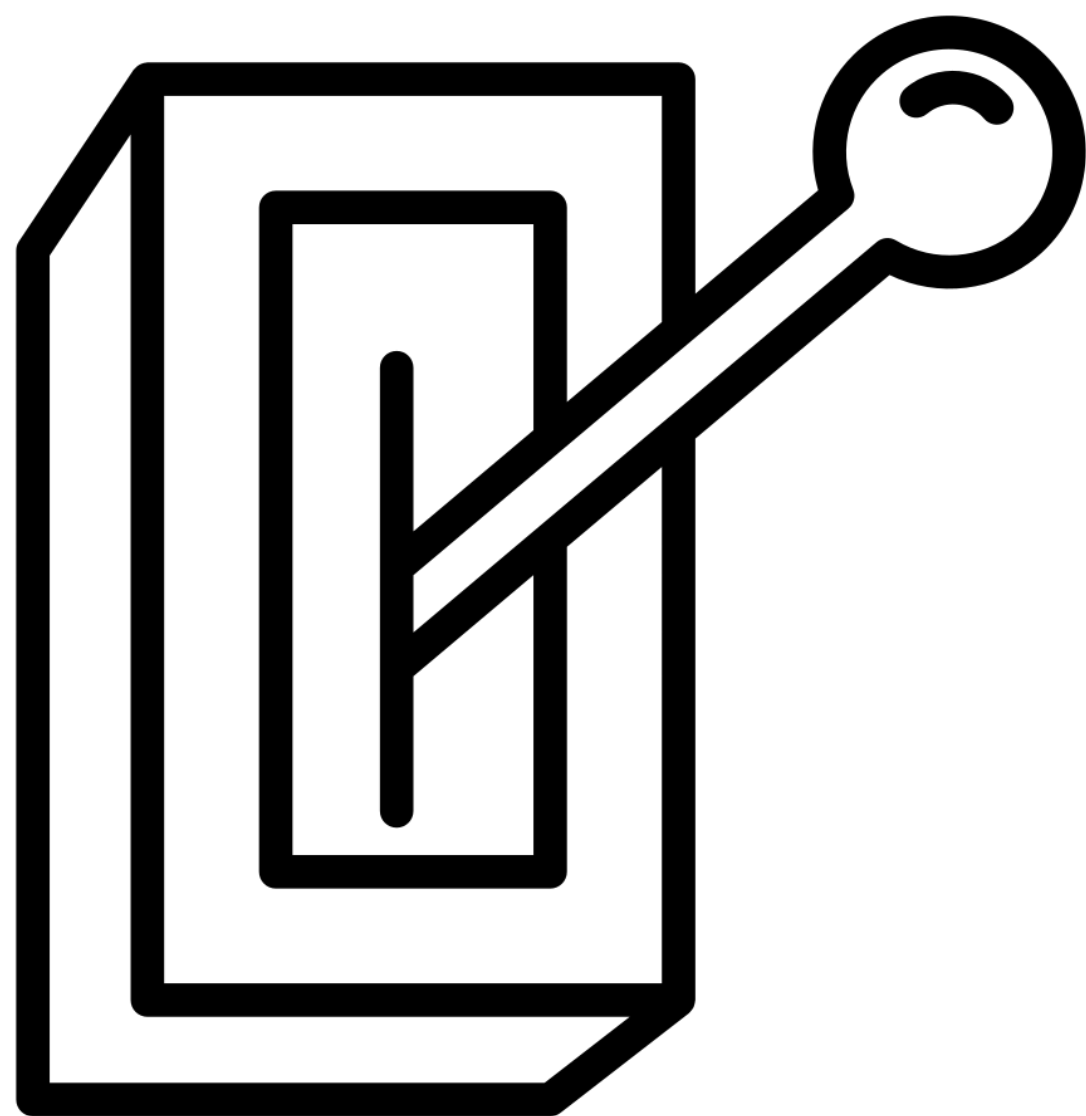
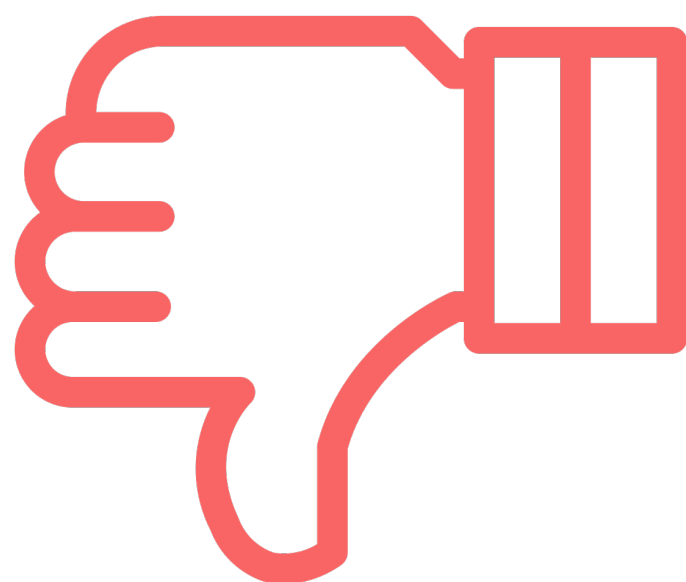
	MDP	POMDP
State	s_t	$p(s_t, c \mid s_1, a_1 \dots s_{t-1}, a_{t-1})$
Policy	$\pi(\cdot \mid s_t)$	

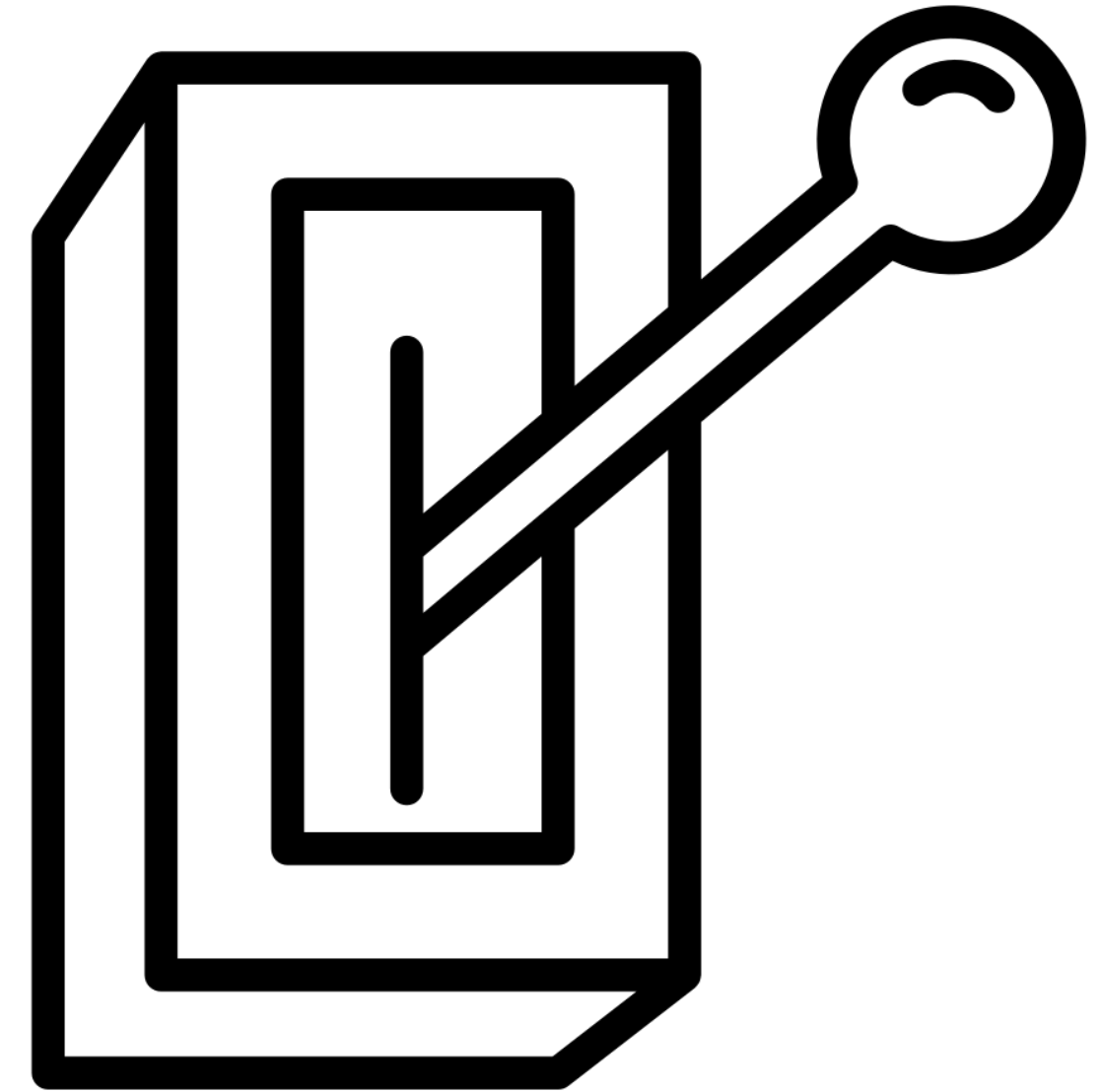
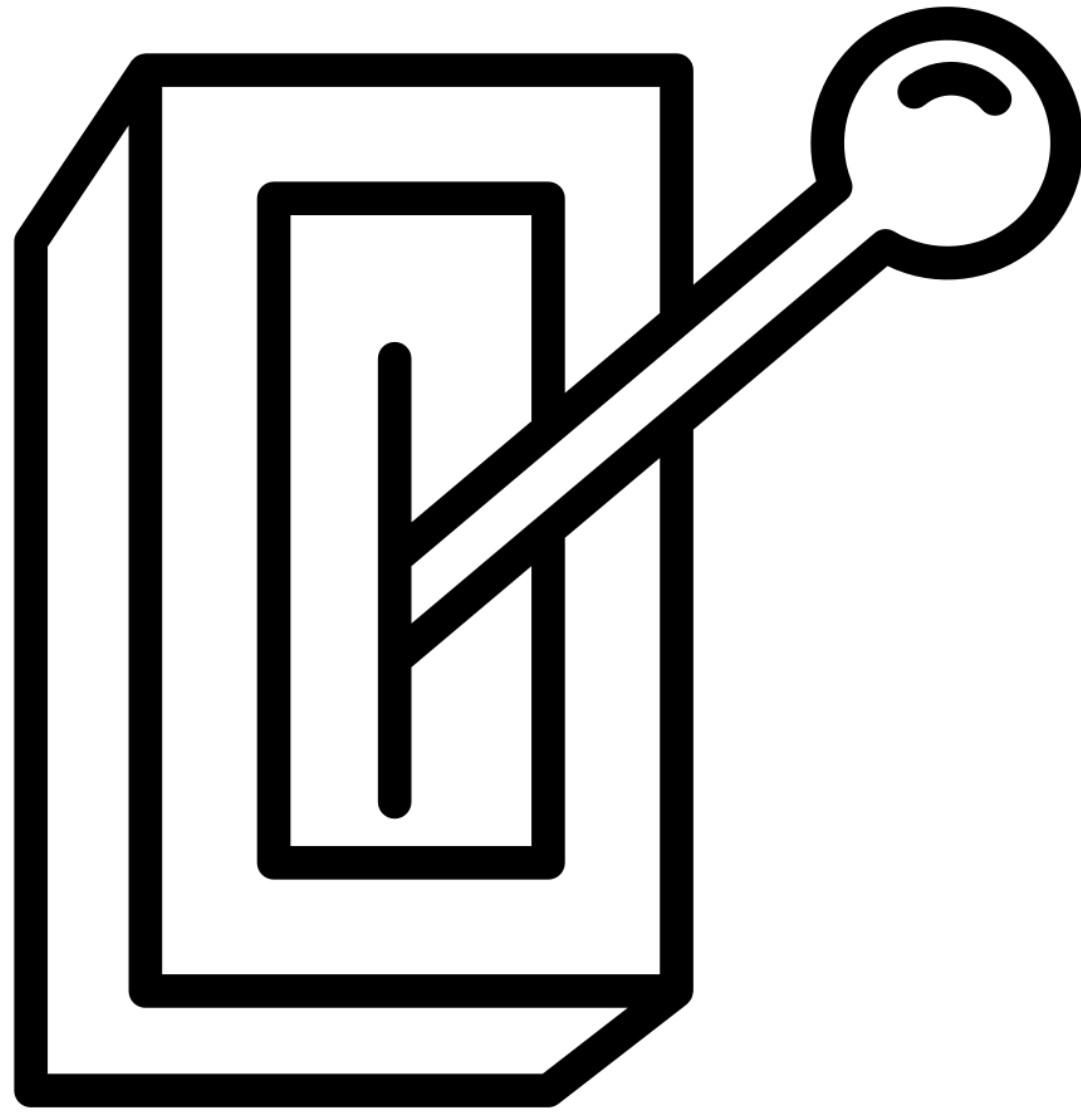
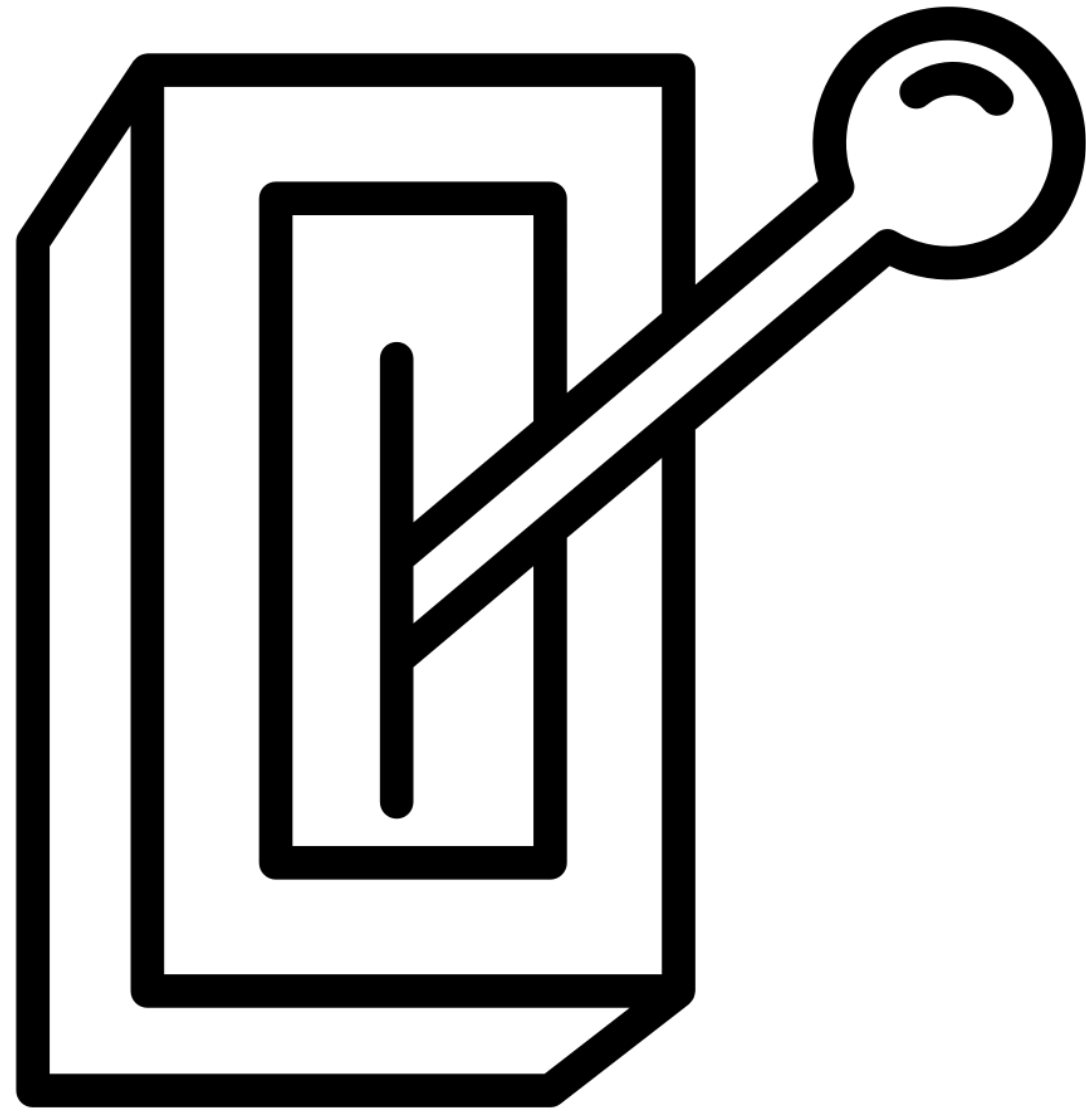
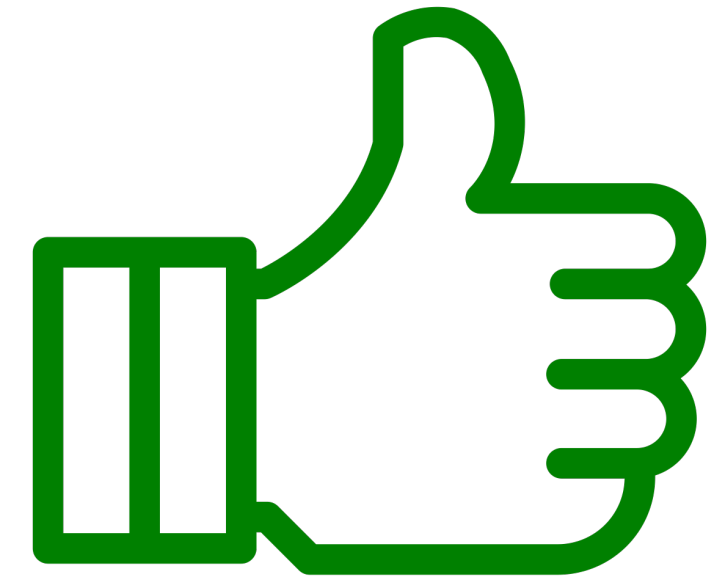
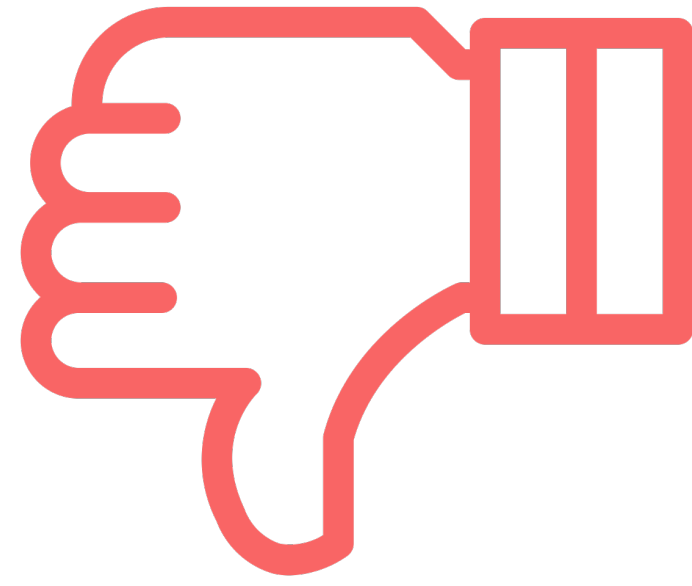
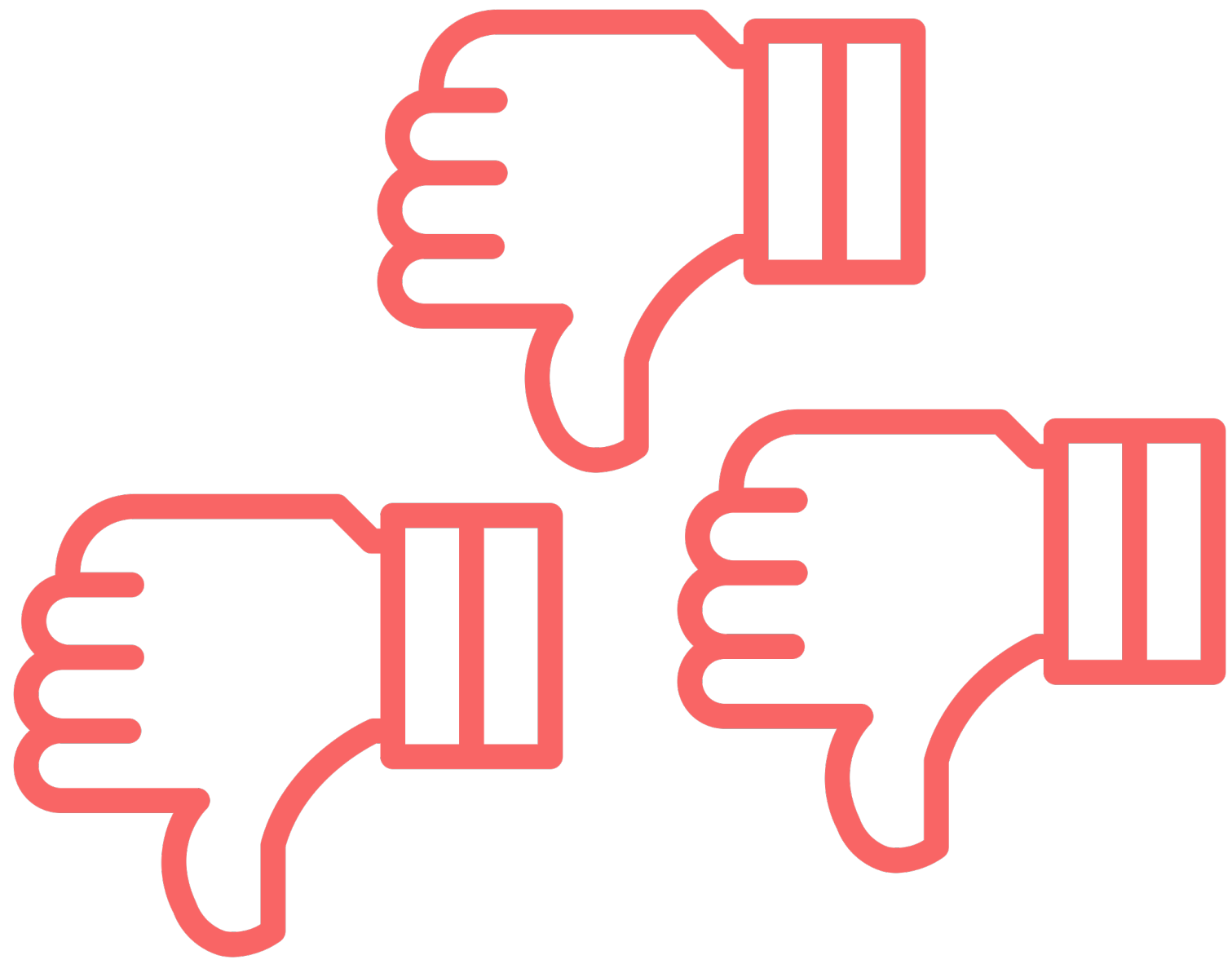
	MDP	POMDP
State	s_t	$p(s_t, c \mid s_1, a_1 \dots s_{t-1}, a_{t-1})$
Policy	$\pi(\cdot \mid s_t)$	$\pi(\cdot \mid s_1, a_1 \dots s_{t-1}, a_{t-1})$

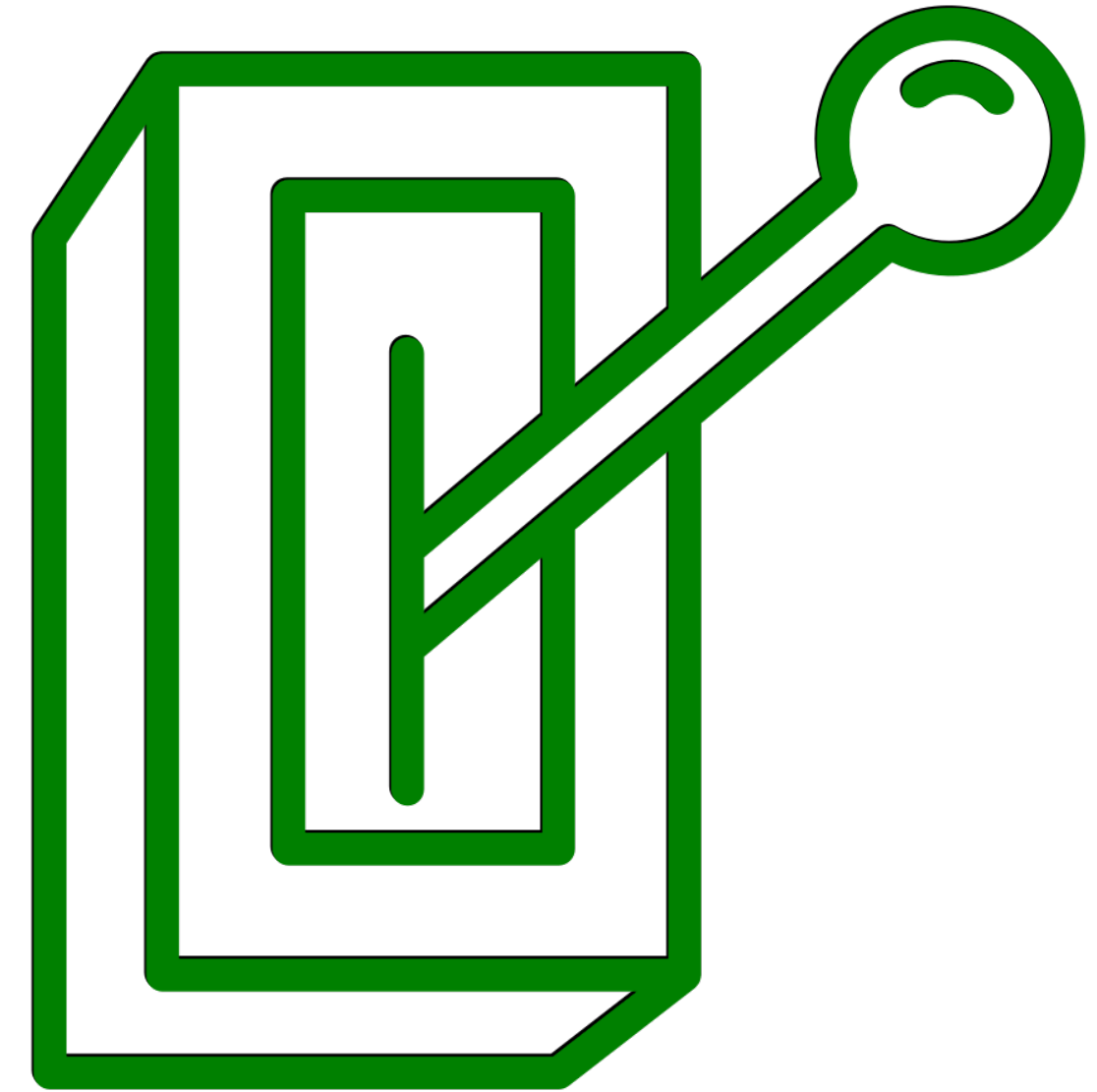
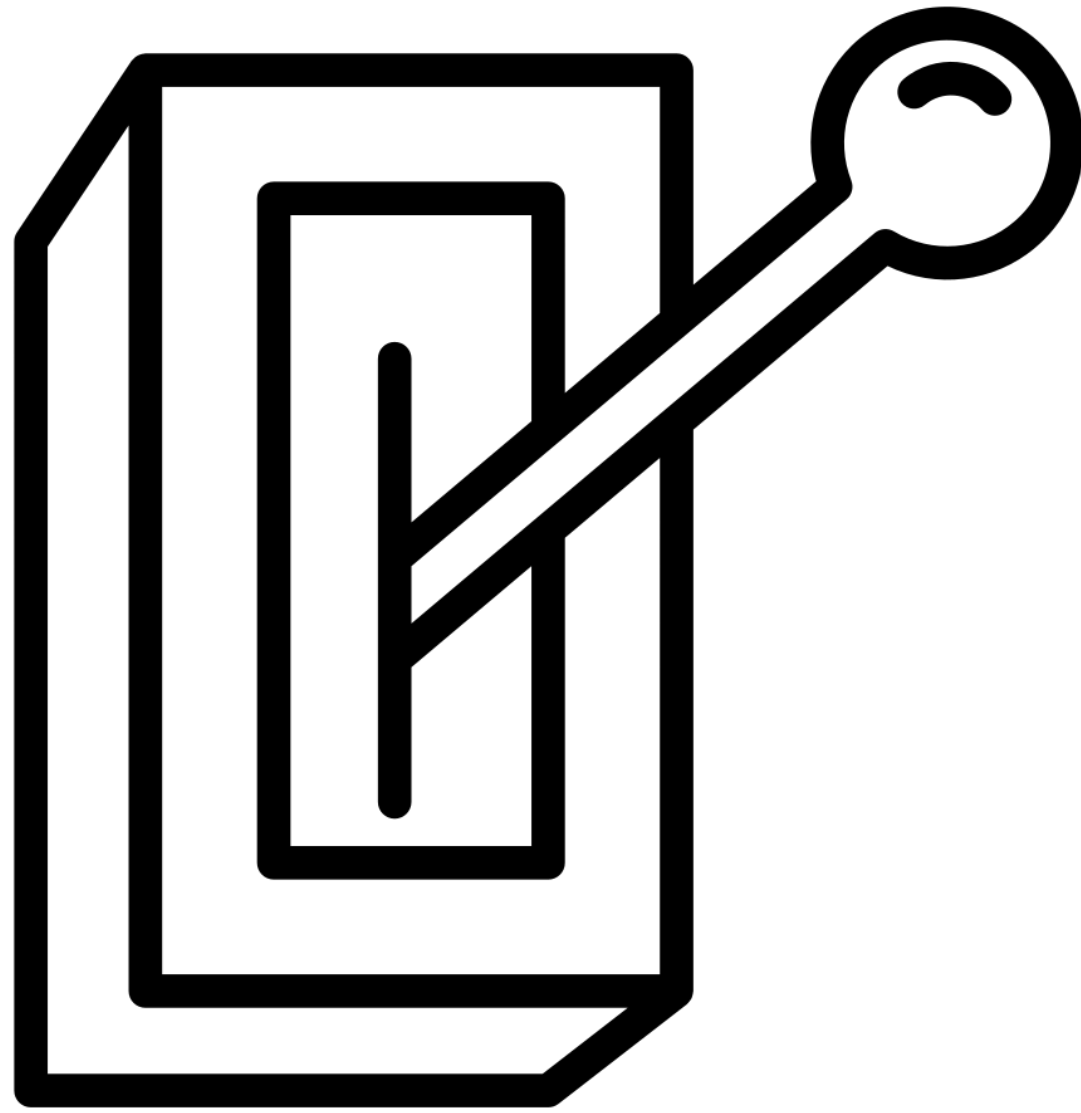
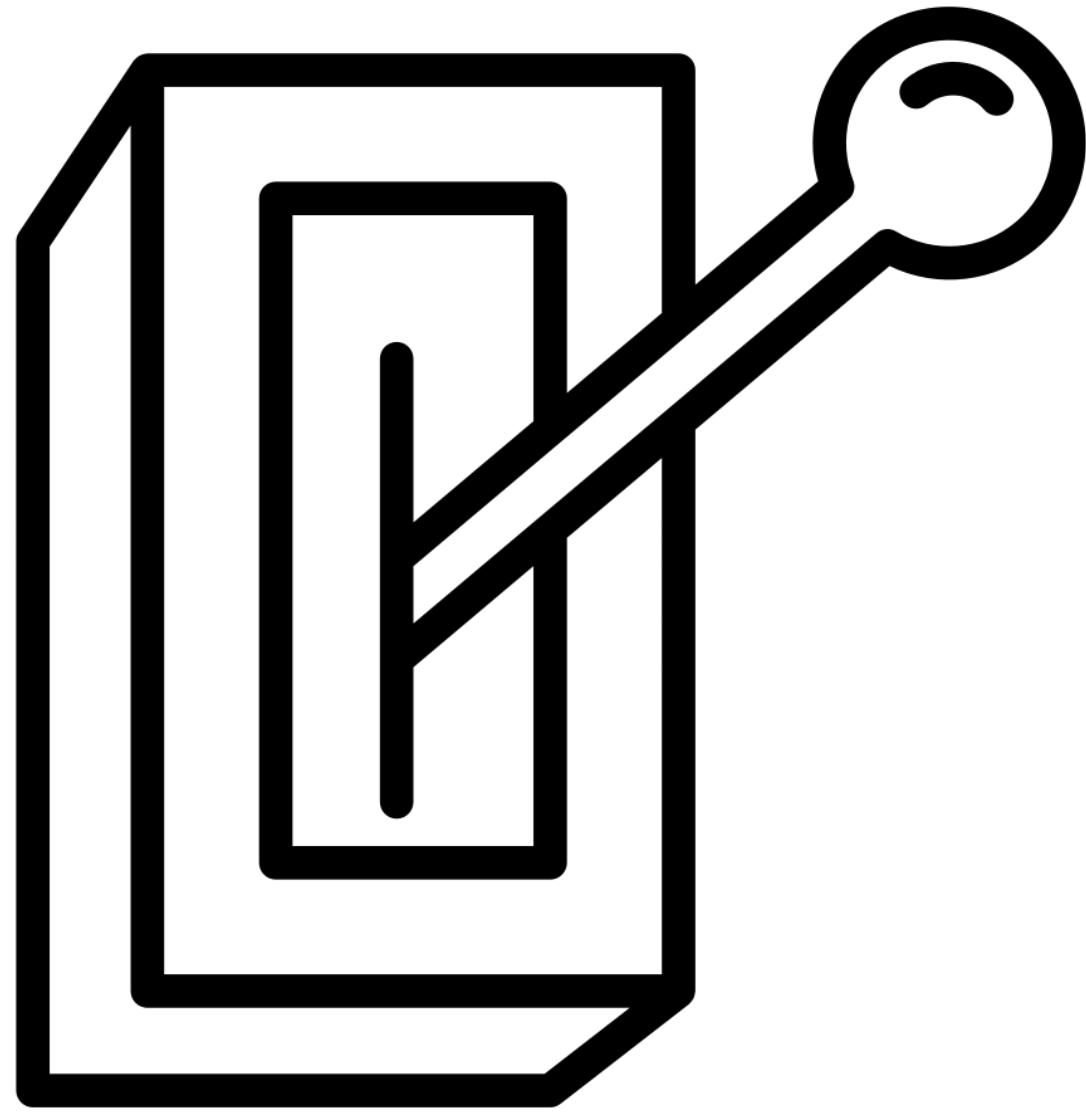
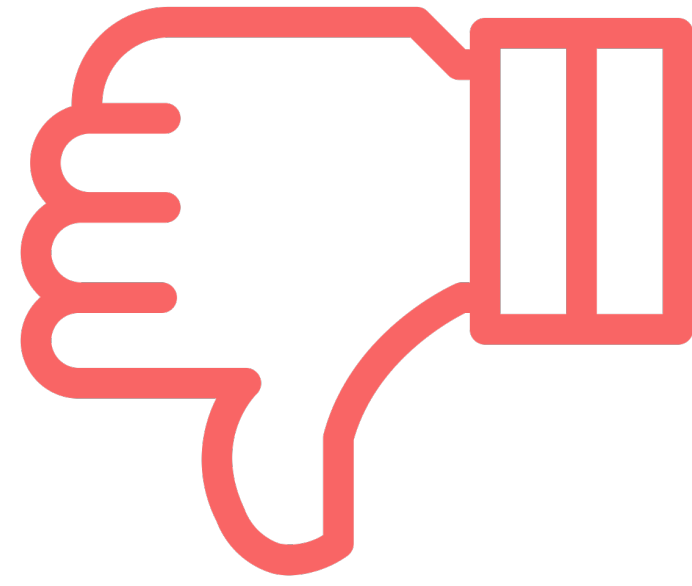
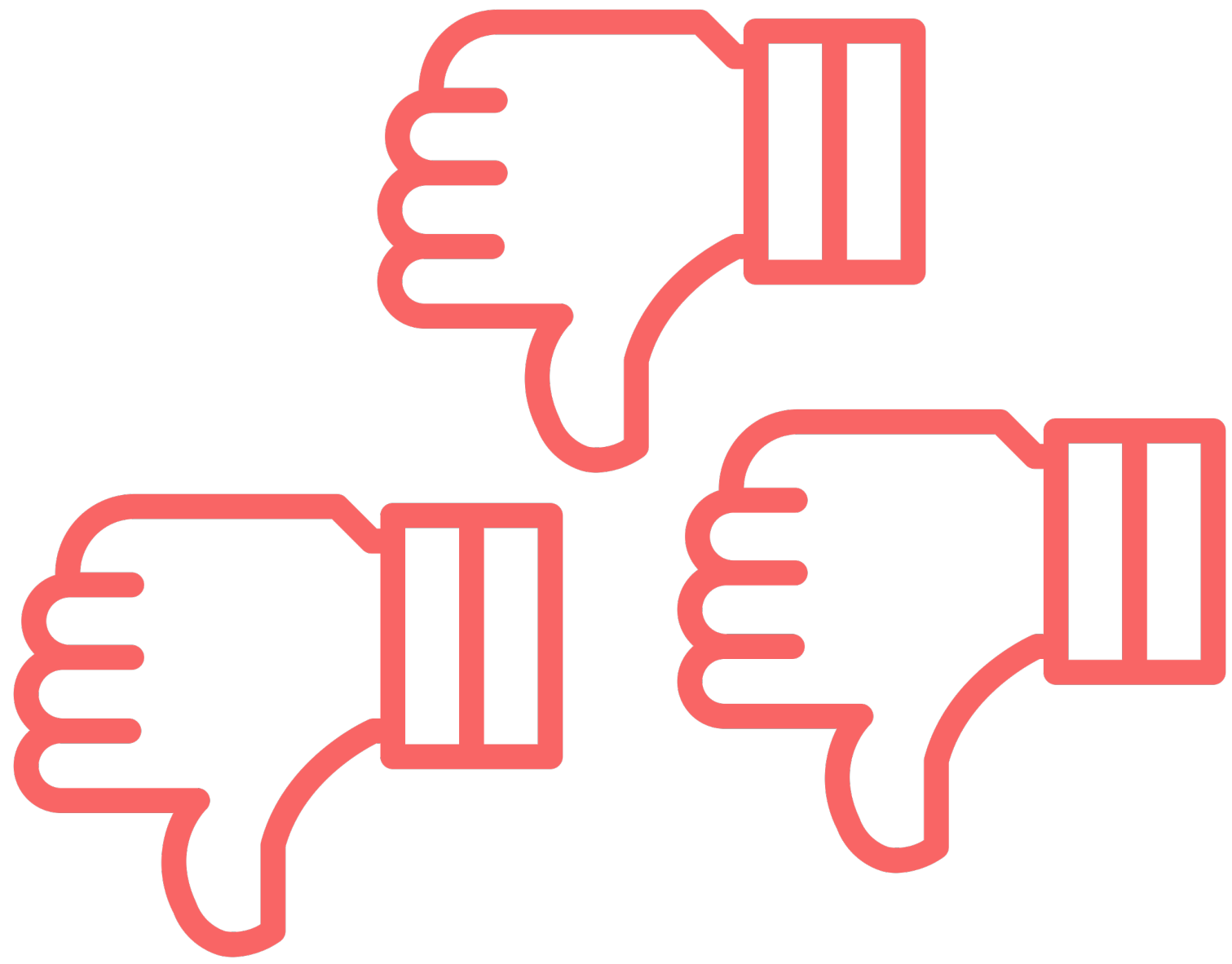












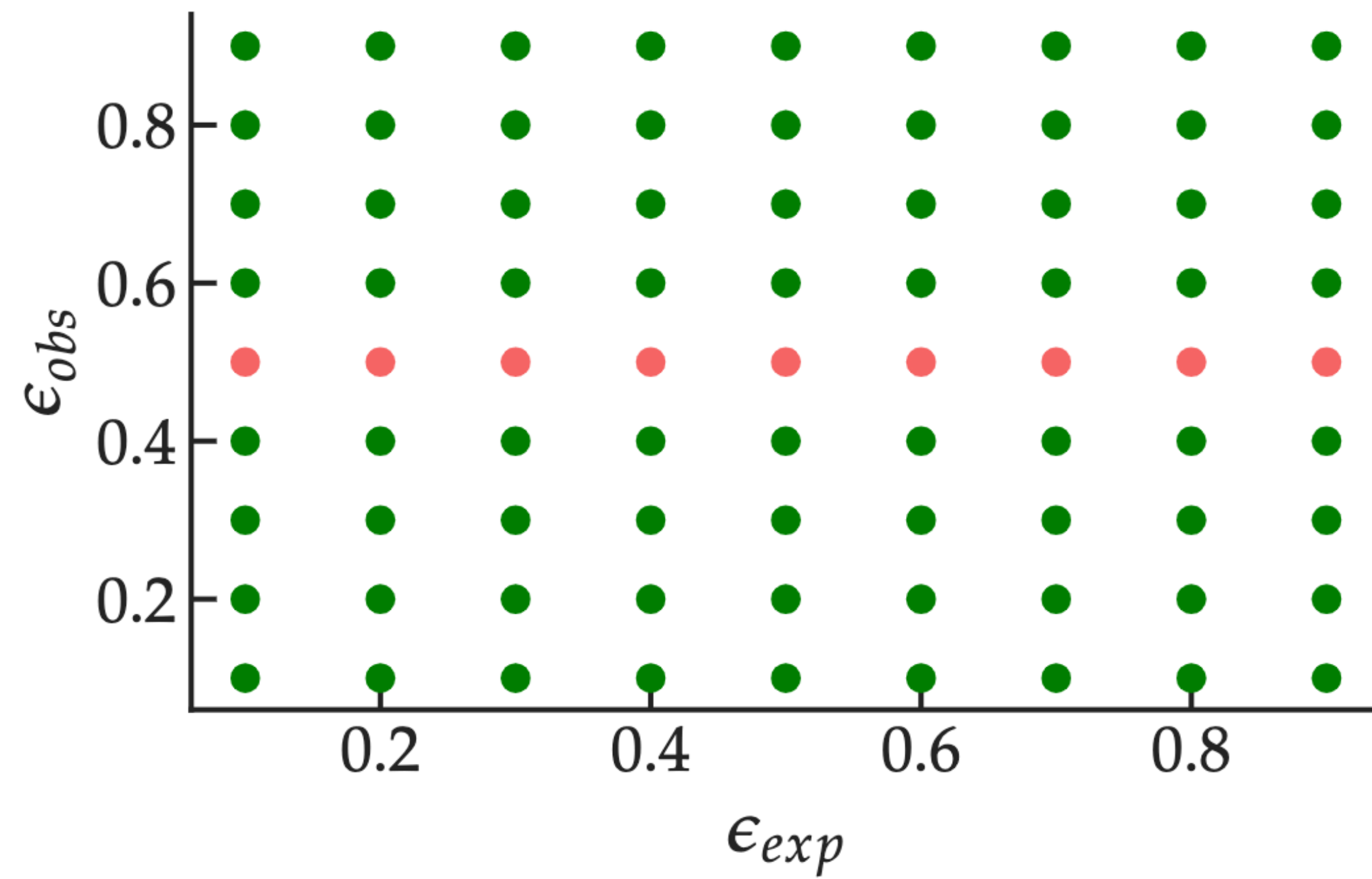
On-Policy (e.g. DAgger):

On-Policy (e.g. DAgger):

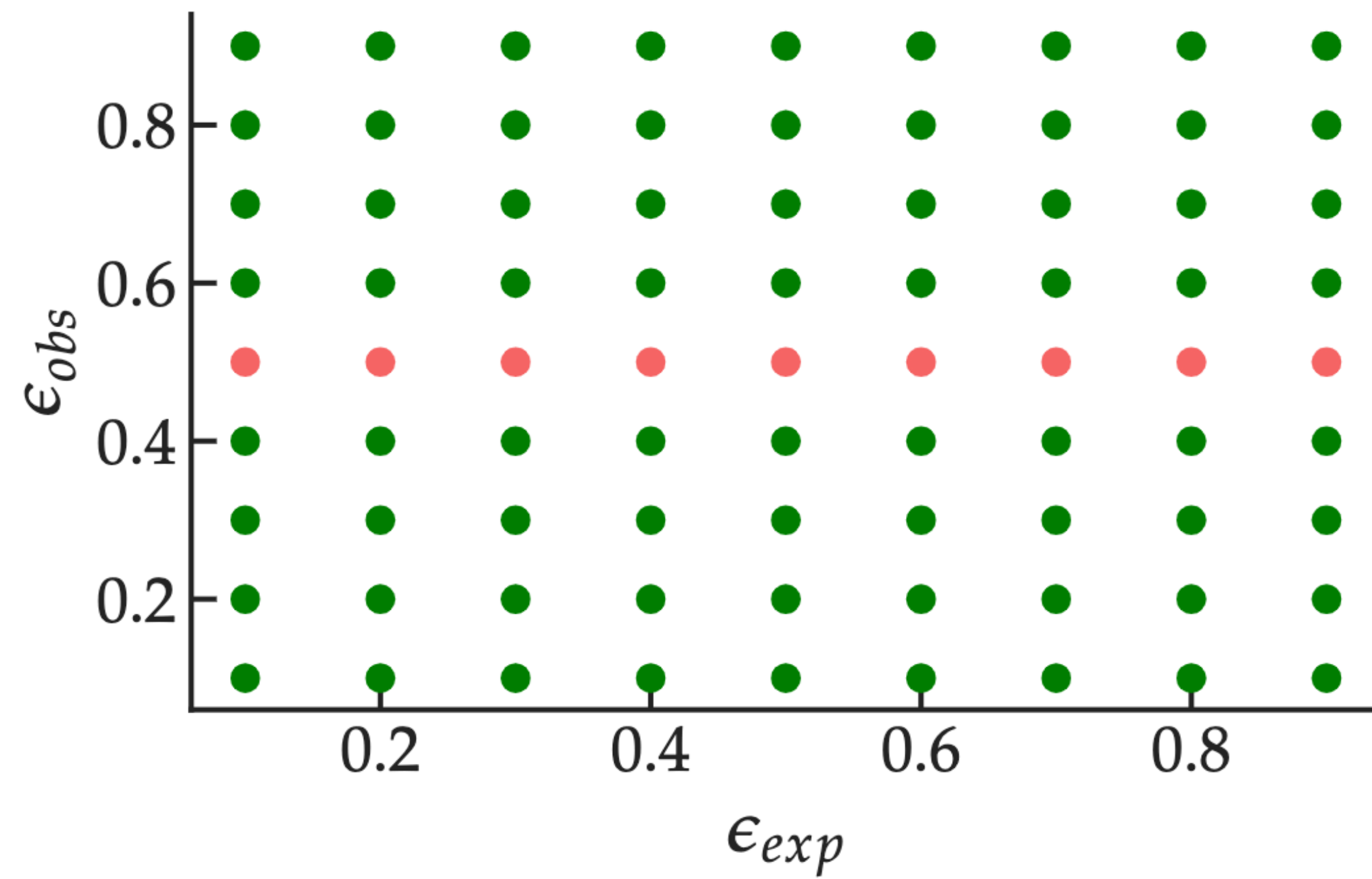
Off-Policy (e.g. BC):

On-Policy (e.g. DAgger):

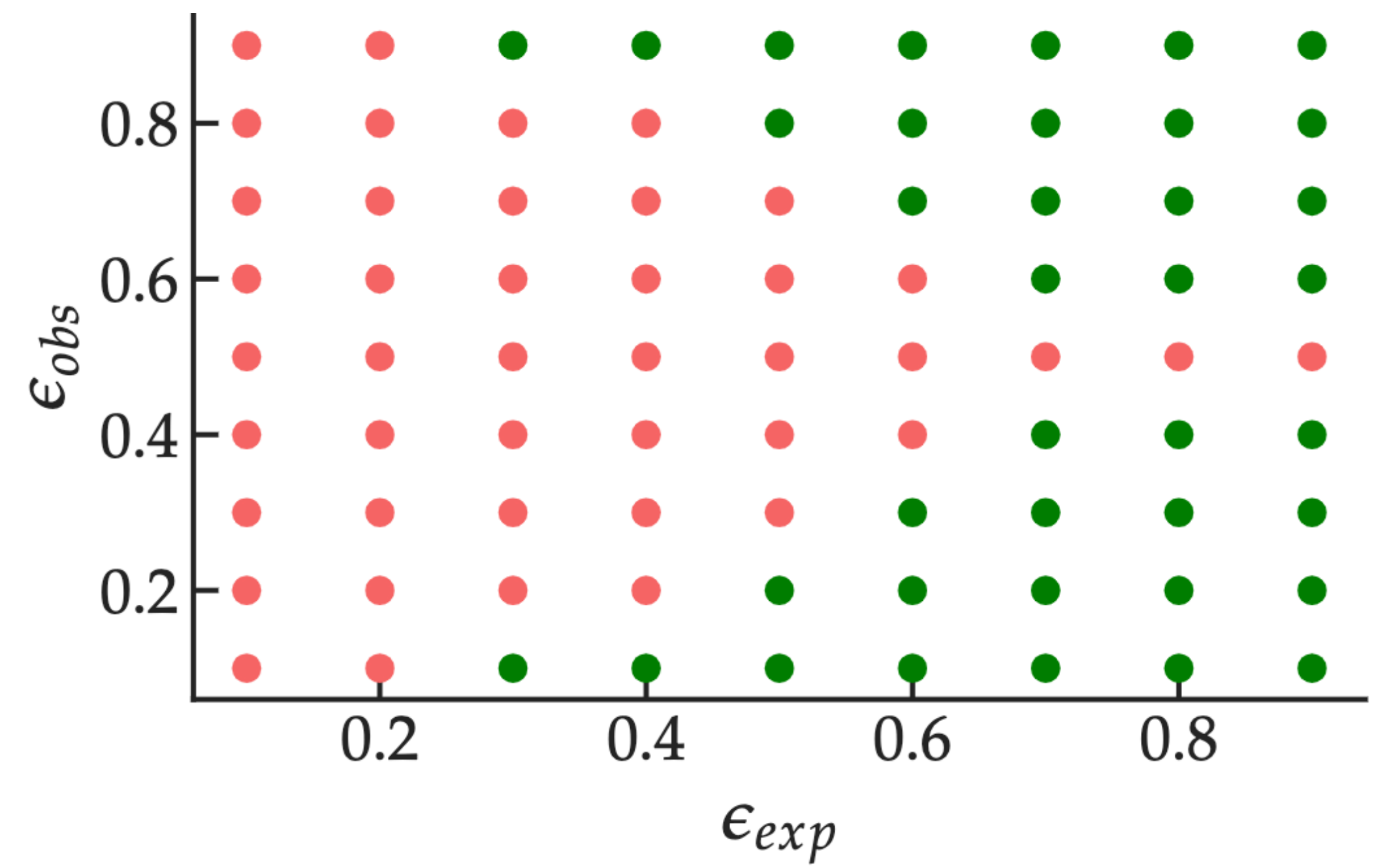
Off-Policy (e.g. BC):

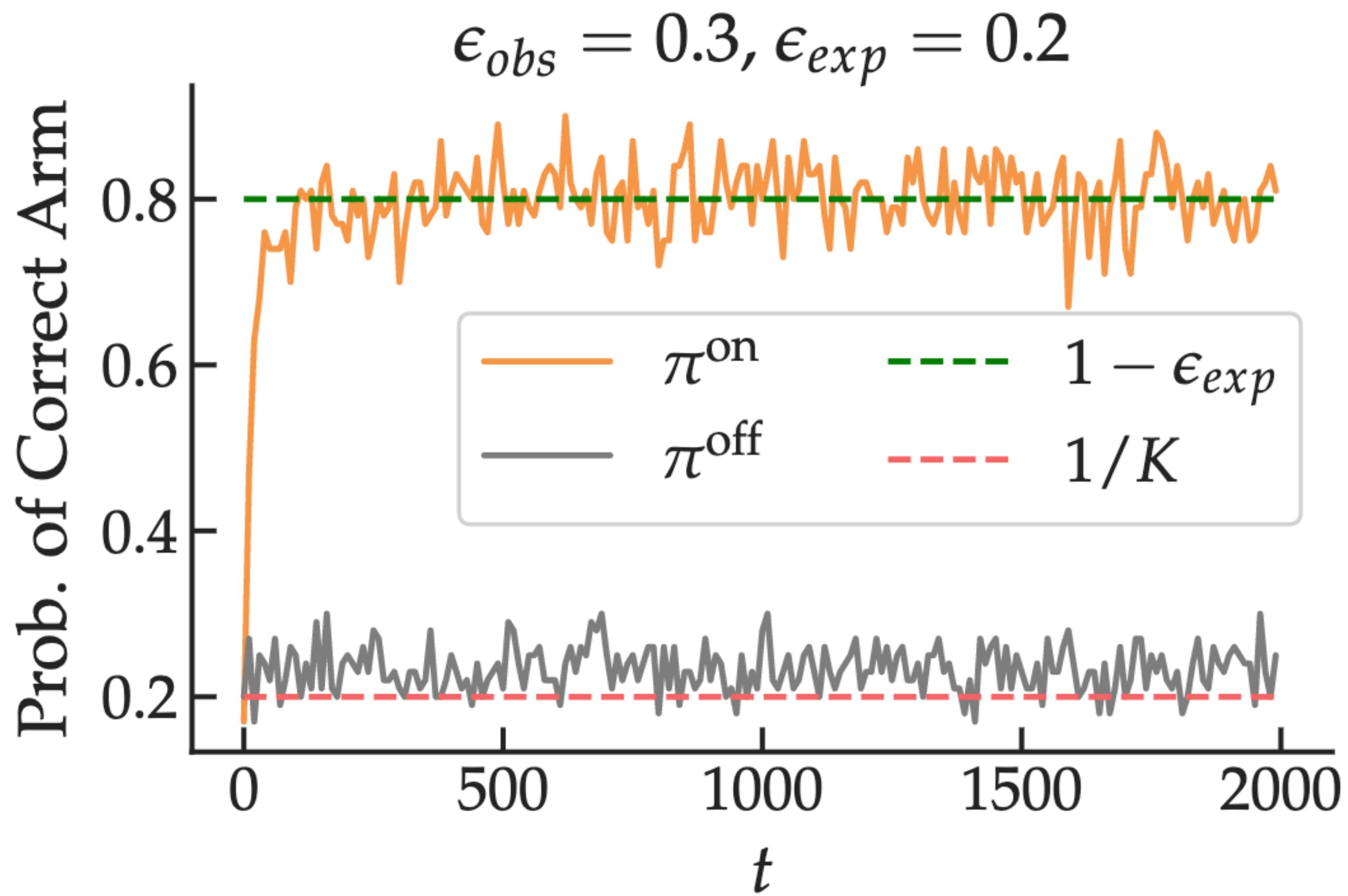


On-Policy (e.g. DAgger):



Off-Policy (e.g. BC):







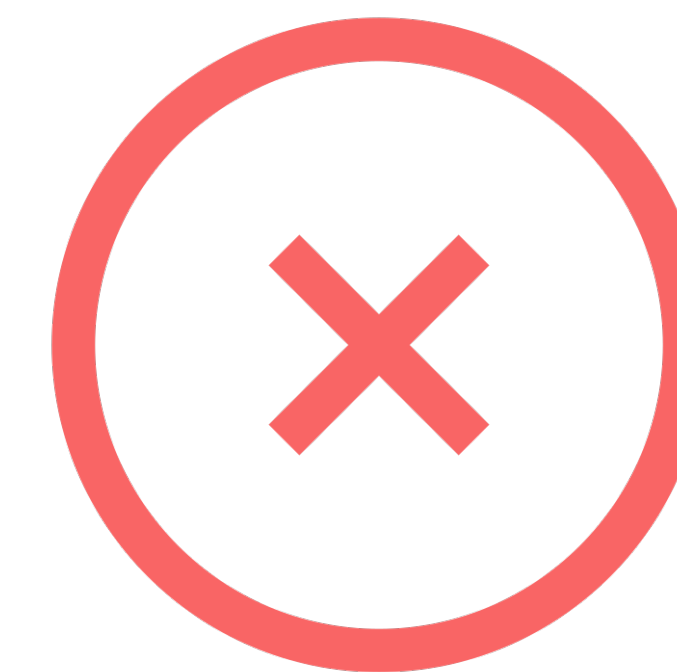
On-Policy:



On-Policy:



Off-Policy:



Train-time:

Train-time: $\pi(a_t | h_t) \approx p(a_t^E | s_1^E, a_1^E, \dots, s_t^E)$

Train-time: $\pi(a_t | h_t) \approx p(a_t^E | s_1^E, a_1^E, \dots, s_t^E)$

Test-time:

$$\begin{array}{ll} \textit{Train-time:} & \pi(a_t | h_t) \approx p(a_t^E | s_1^E, a_1^E, \dots, s_t^E) \\ \textit{Test-time:} & p(a_t^E | s_1, a_1, \dots, s_t) \end{array}$$

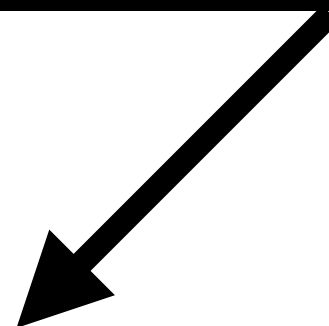
Train-time: $\pi(a_t | h_t) \approx p(a_t^E | s_1^E, a_1^E, \dots, s_t^E)$

Test-time: $p(a_t^E | s_1, a_1, \dots, s_t)$

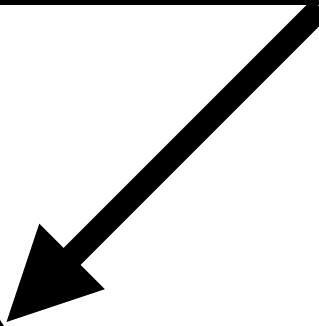
It's just covariate shift in the space of histories!

Hidden Context

Hidden Context



Hidden Context



Asymp. Realizability

Hidden Context

```
graph TD; A[Hidden Context] --> B[Asymp. Realizability]; B --> C[ ];
```

Asymp. Realizability

Hidden Context

```
graph TD; A[Hidden Context] --> B[Asymp. Realizability]; B --> C[Nonzero Early Error];
```

Asymp. Realizability

Nonzero Early Error

Hidden Context

```
graph TD; A[Hidden Context] --> B[Asymp. Realizability]; A --> C[ ]; B --> D[Nonzero Early Error];
```

The diagram consists of three rounded rectangular boxes with black outlines. The top box contains the text "Hidden Context". Two arrows originate from the bottom center of this box, pointing downwards and outwards to the left and right. The left arrow points to the top center of the middle box, which contains the text "Asymp. Realizability". The right arrow points to the right side of the page. A single arrow originates from the bottom center of the middle box and points downwards to the top center of the bottom box, which contains the text "Nonzero Early Error".

Asymp. Realizability

Nonzero Early Error

Hidden Context

```
graph TD; A[Hidden Context] --> B[Asymp. Realizability]; A --> C[Sequence Models]; B --> D[Nonzero Early Error]
```

Asymp. Realizability

Sequence Models

Nonzero Early Error

Hidden Context

```
graph TD; A[Hidden Context] --> B[Asymp. Realizability]; A --> C[Sequence Models]; B --> D[Nonzero Early Error]; C --> E[ ];
```

The diagram is a flowchart with four nodes. The top node is 'Hidden Context'. Two arrows point from it to 'Asymp. Realizability' on the left and 'Sequence Models' on the right. From 'Asymp. Realizability', an arrow points down to 'Nonzero Early Error'. From 'Sequence Models', an arrow points down to an empty space.

Asymp. Realizability

Sequence Models

Nonzero Early Error

Hidden Context

```
graph TD; A[Hidden Context] --> B[Asymp. Realizability]; A --> C[Sequence Models]; B --> D[Nonzero Early Error]; C --> E[H-space Cov. Shift]
```

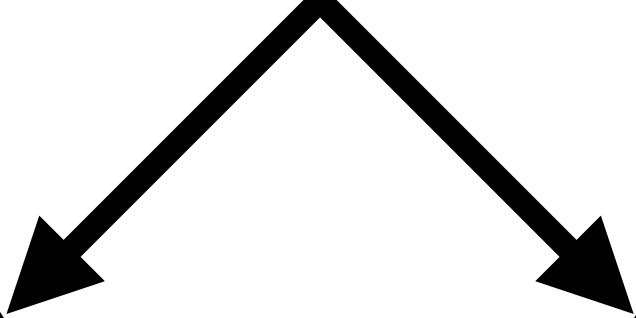
Asymp. Realizability

Sequence Models

Nonzero Early Error

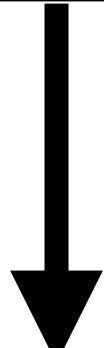
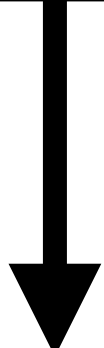
H-space Cov. Shift

Hidden Context



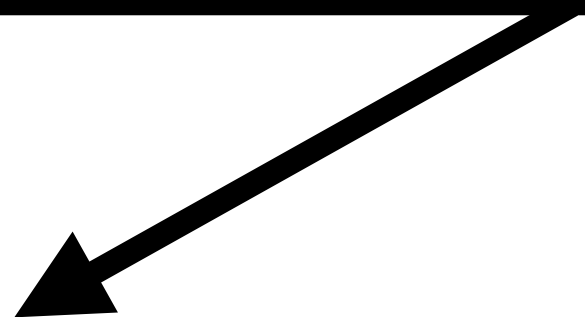
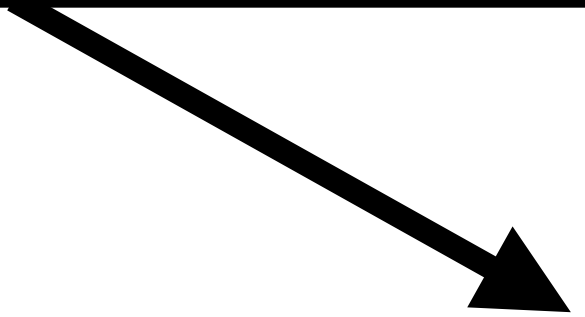
Asymp. Realizability

Sequence Models

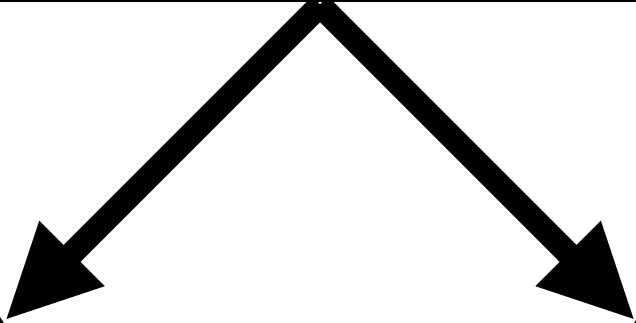


Nonzero Early Error

H-space Cov. Shift

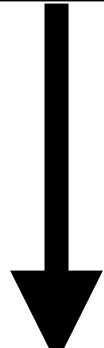
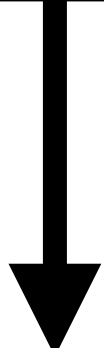


Hidden Context



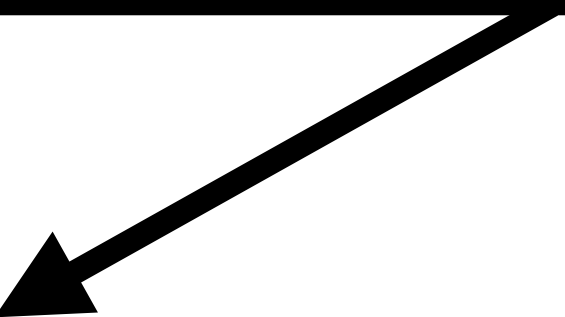
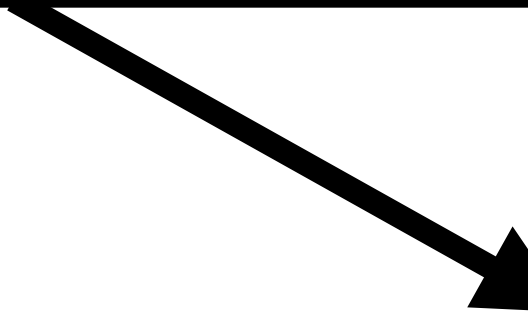
Asymp. Realizability

Sequence Models

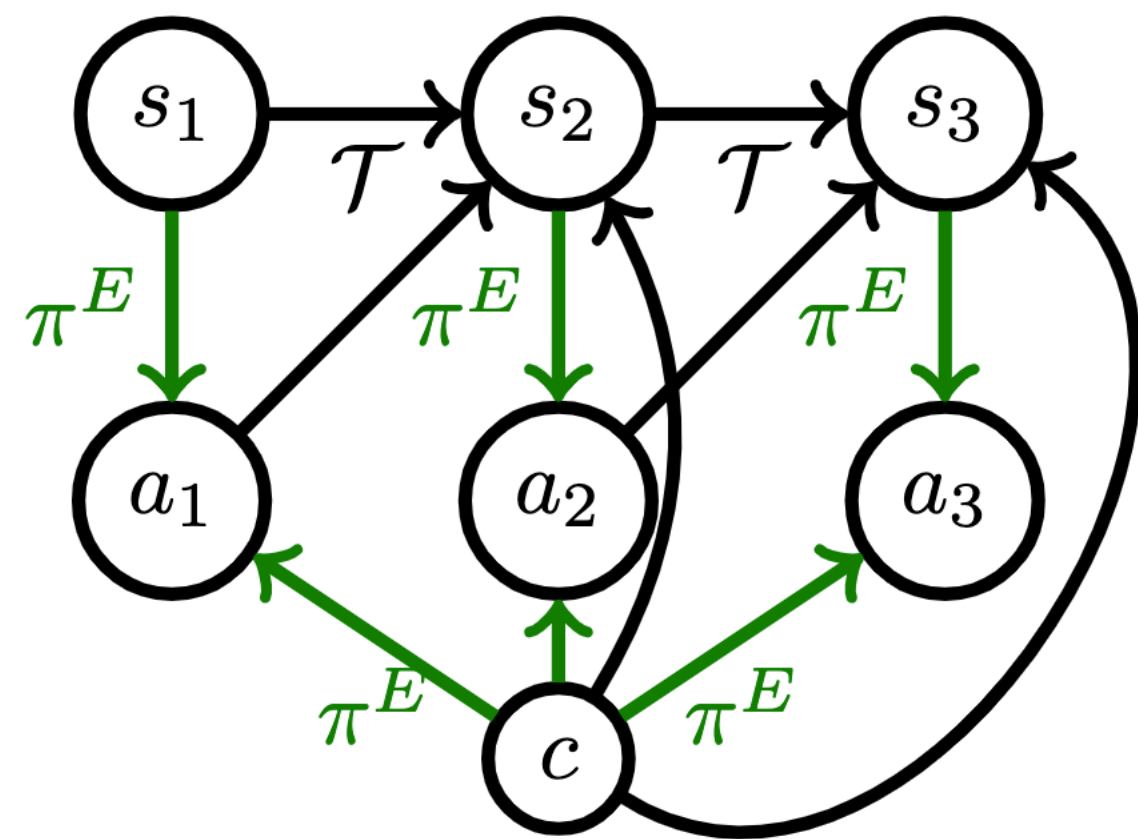


Nonzero Early Error

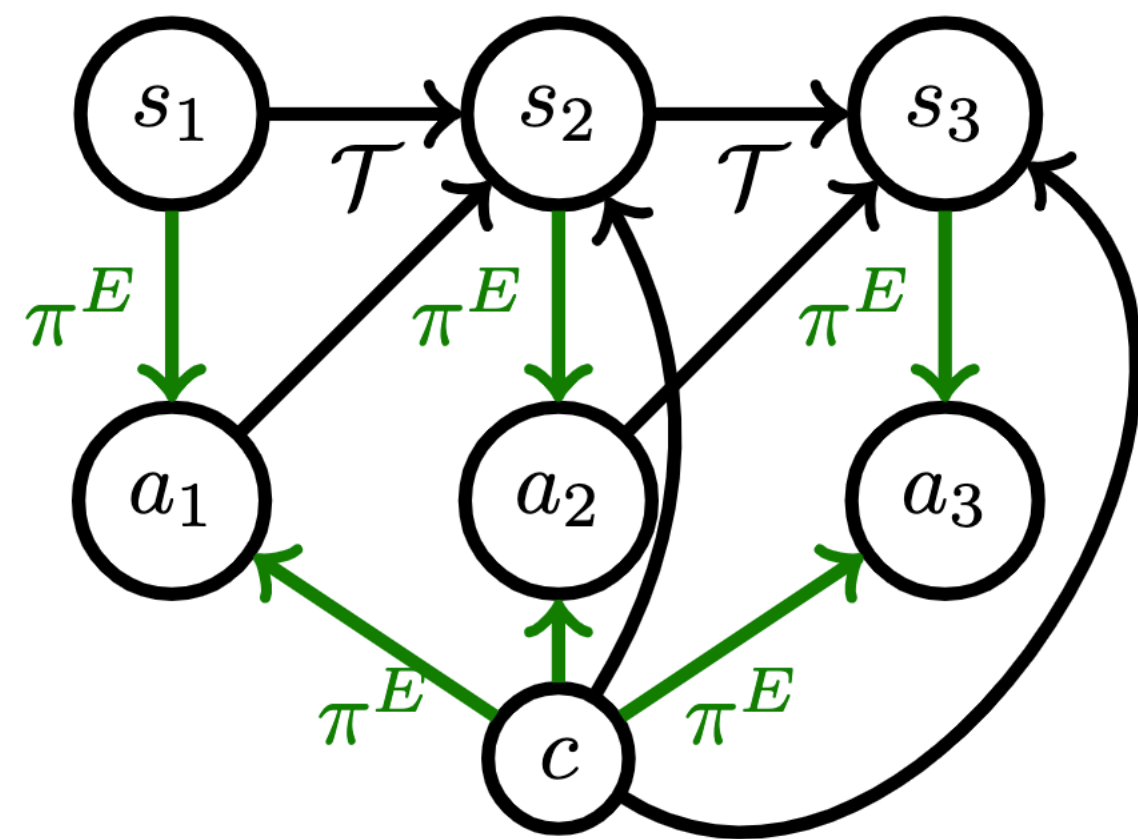
H-space Cov. Shift



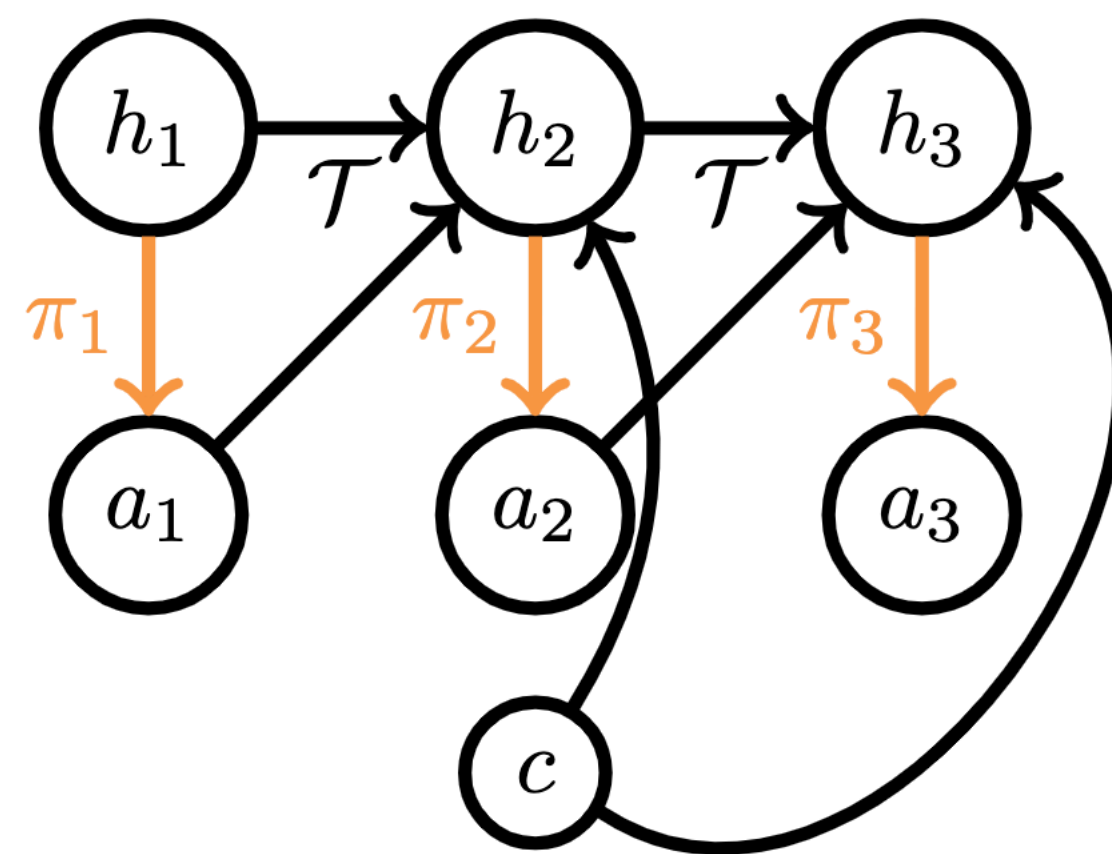
Offline IL Fails



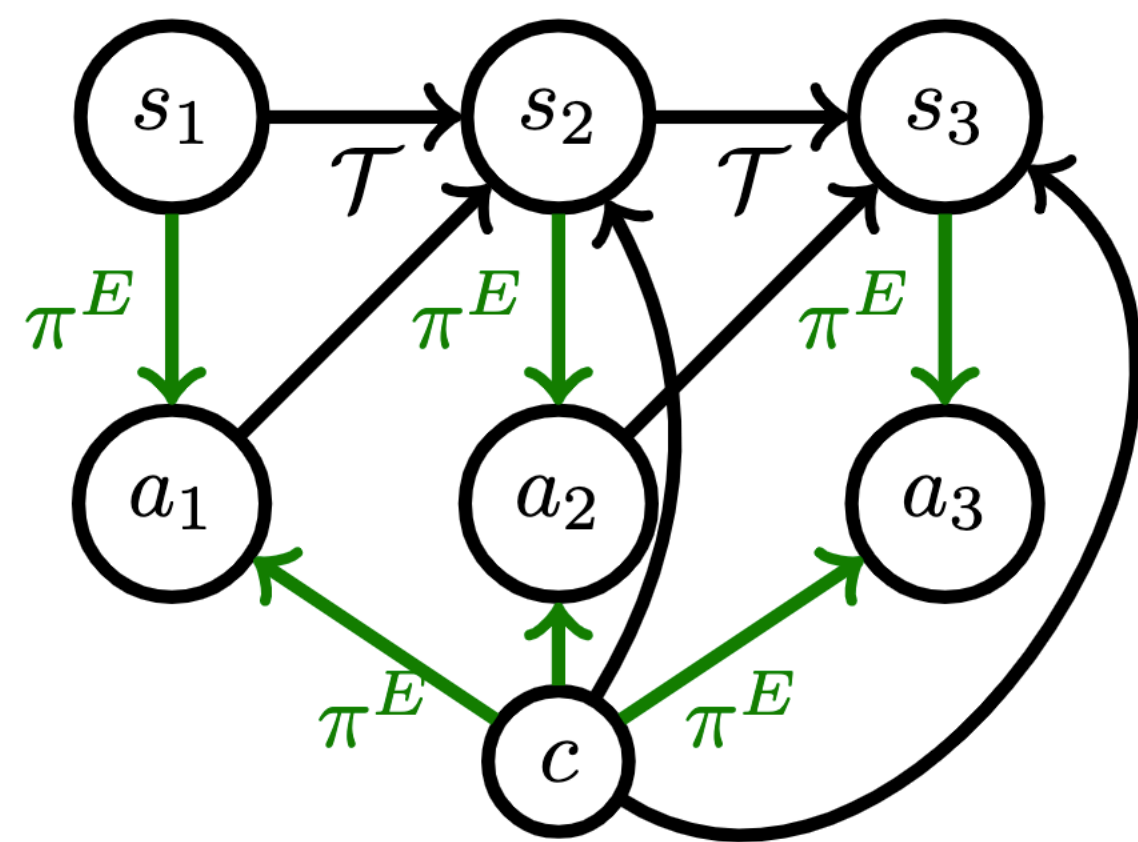
(a) $\tau \sim \pi^E$



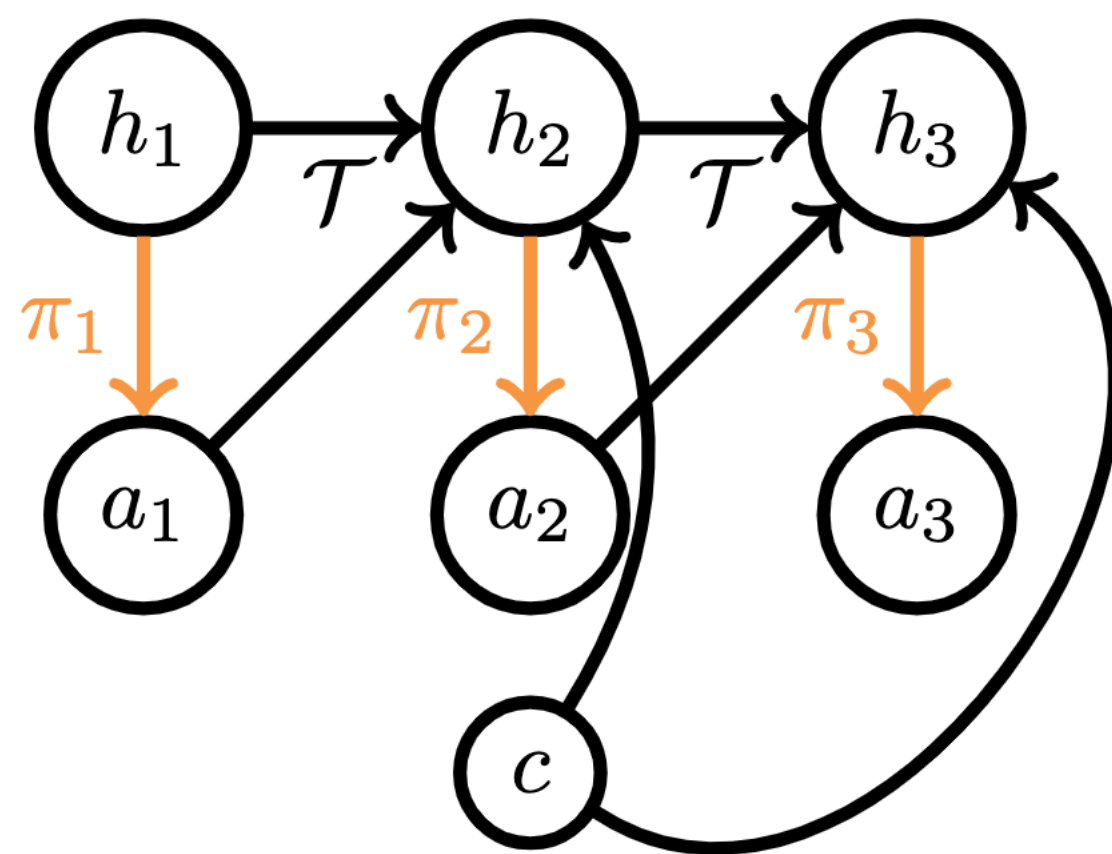
(a) $\tau \sim \pi^E$



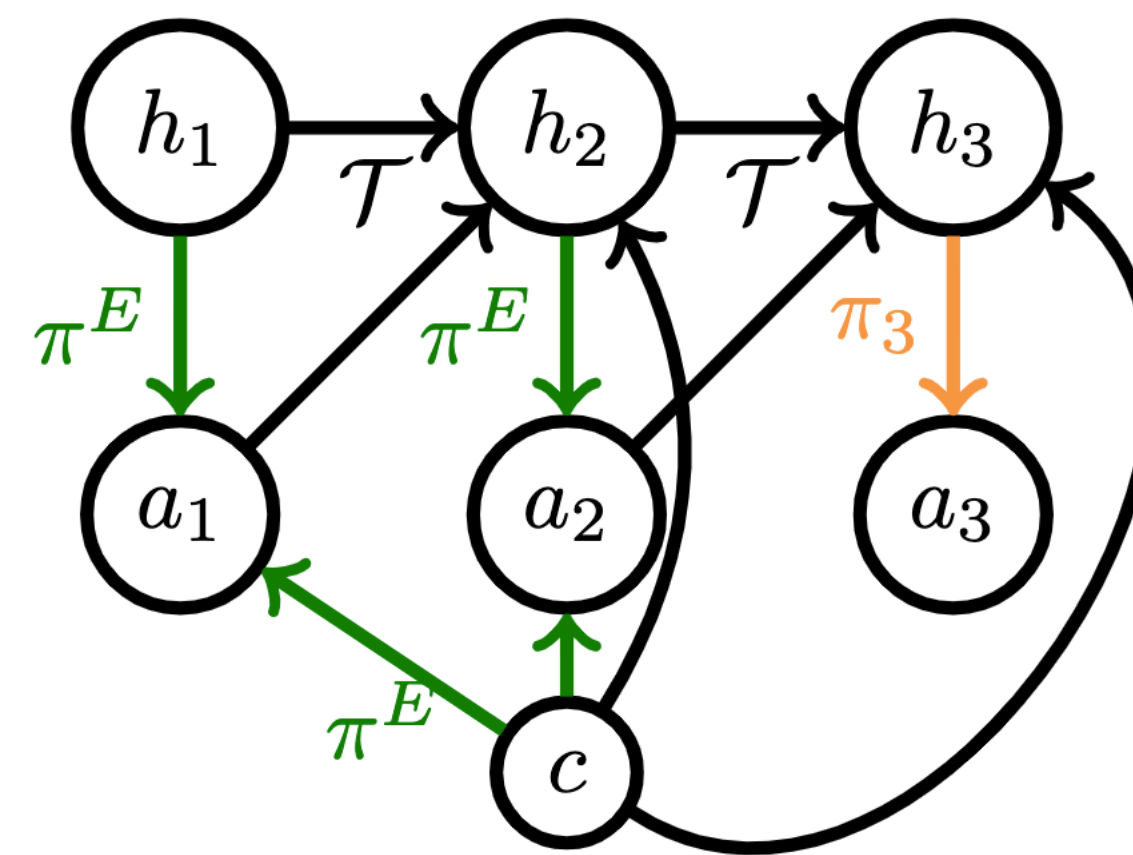
(b) On-Policy



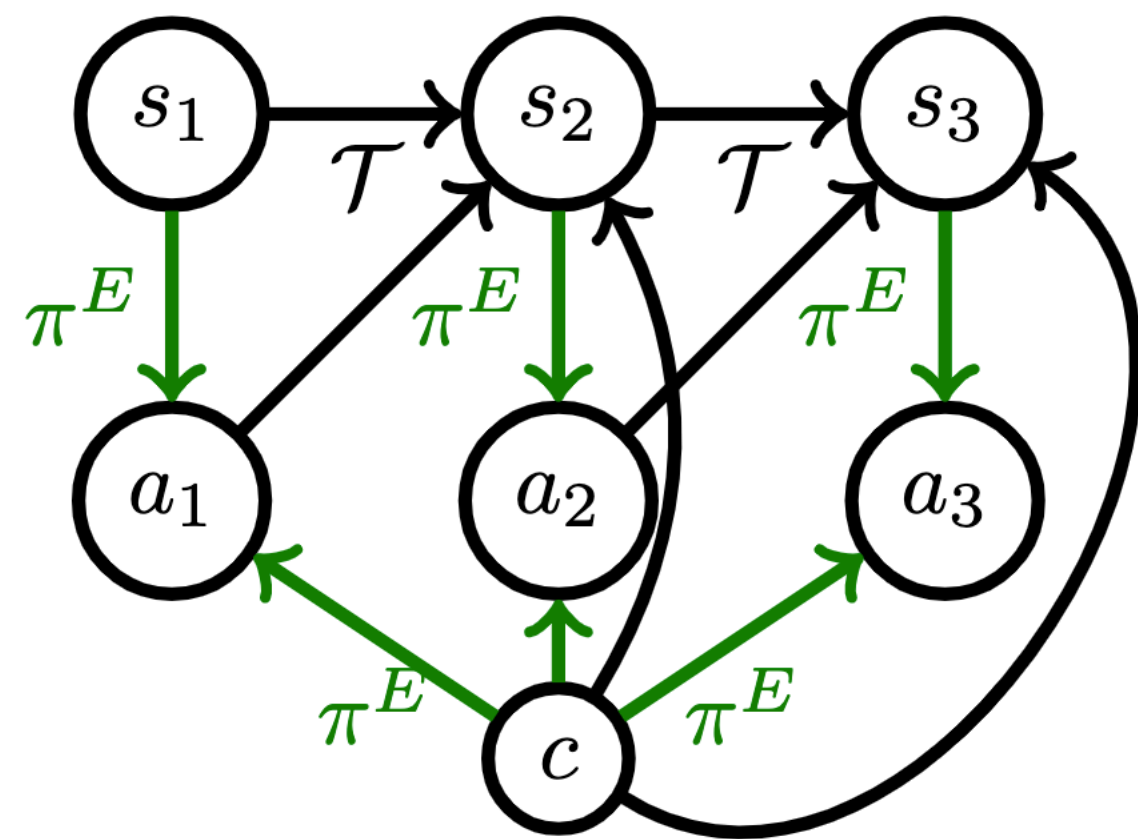
(a) $\tau \sim \pi^E$



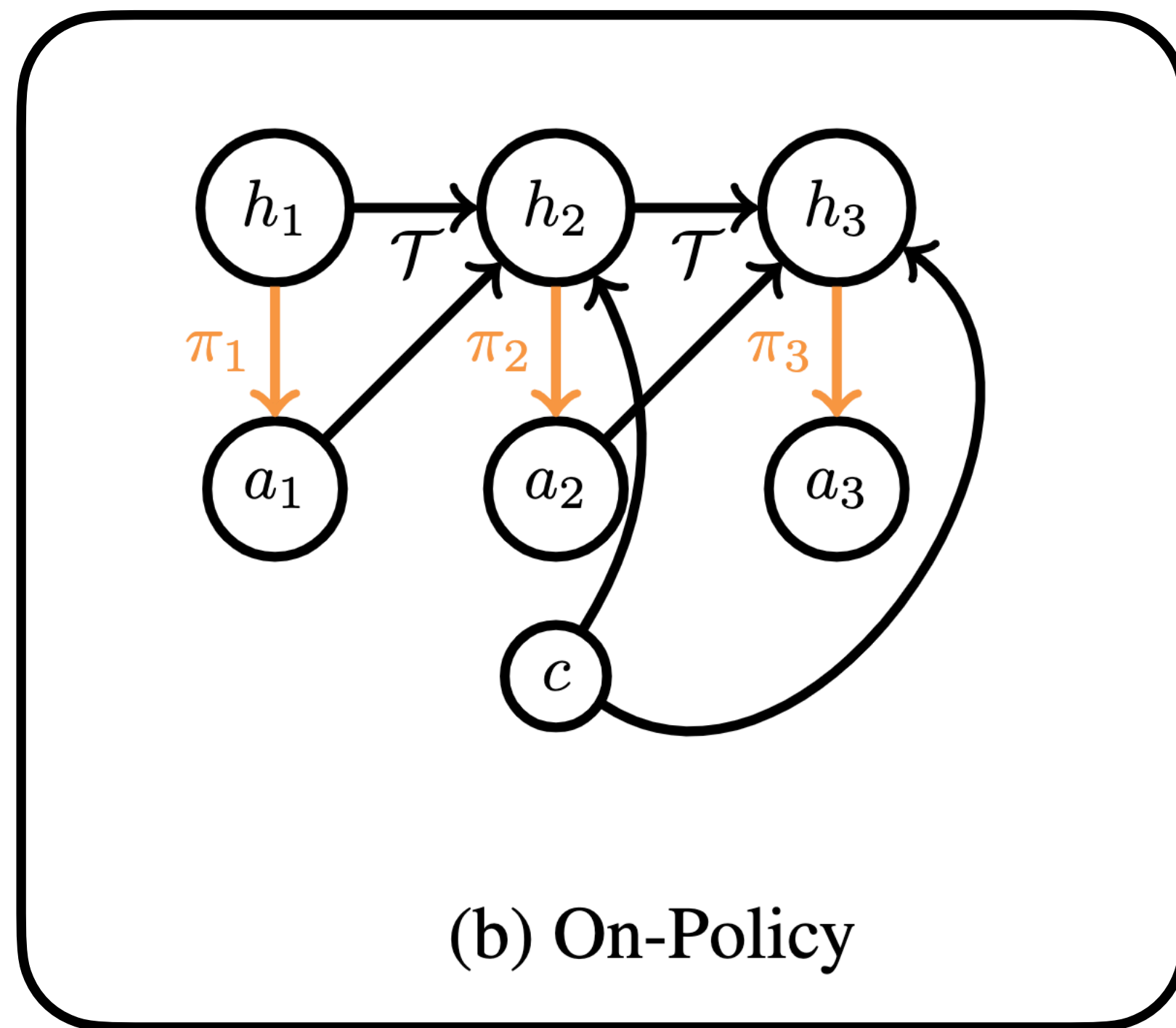
(b) On-Policy



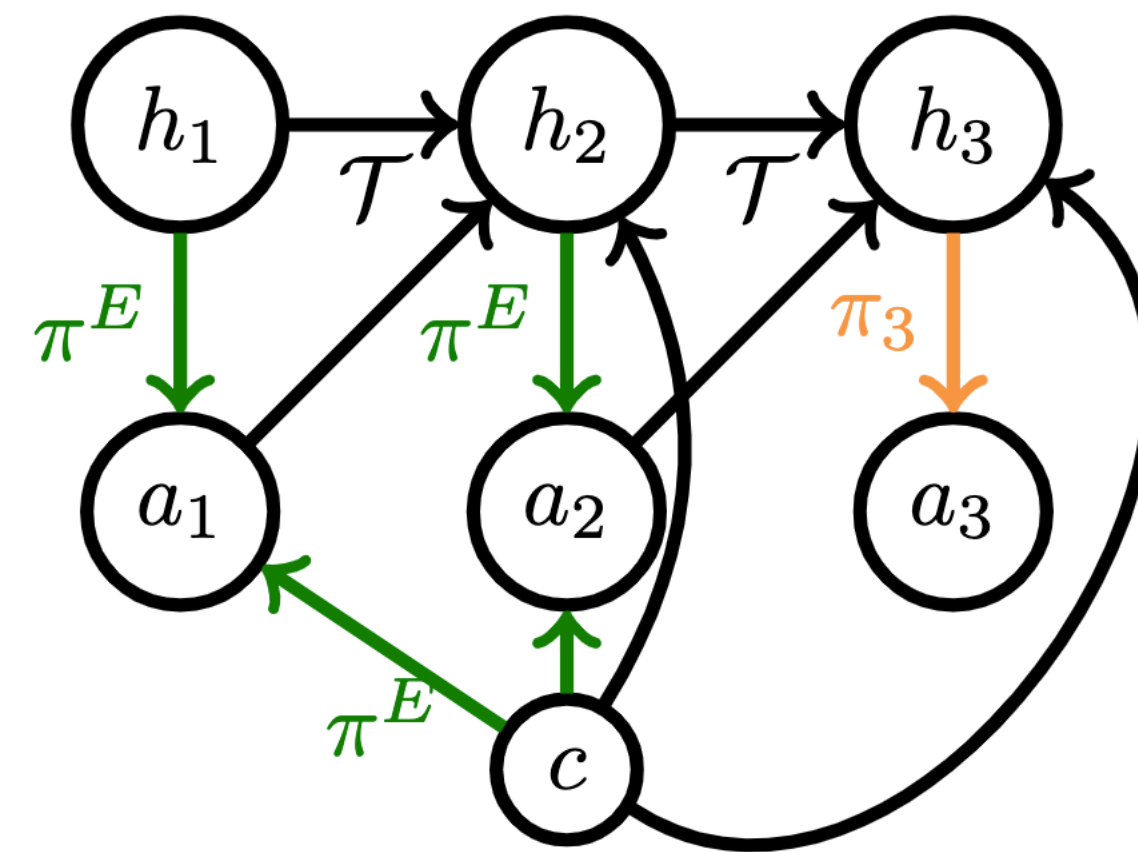
(c) Off-Policy



(a) $\tau \sim \pi^E$



(b) On-Policy



(c) Off-Policy

$$p_{\text{on}}(c, h_t) \propto p(\tau; \pi) \propto p(c)p(s_1) \prod_{i=1}^{t-1} \mathcal{T}(s_{i+1} | s_i, a_i, c)$$

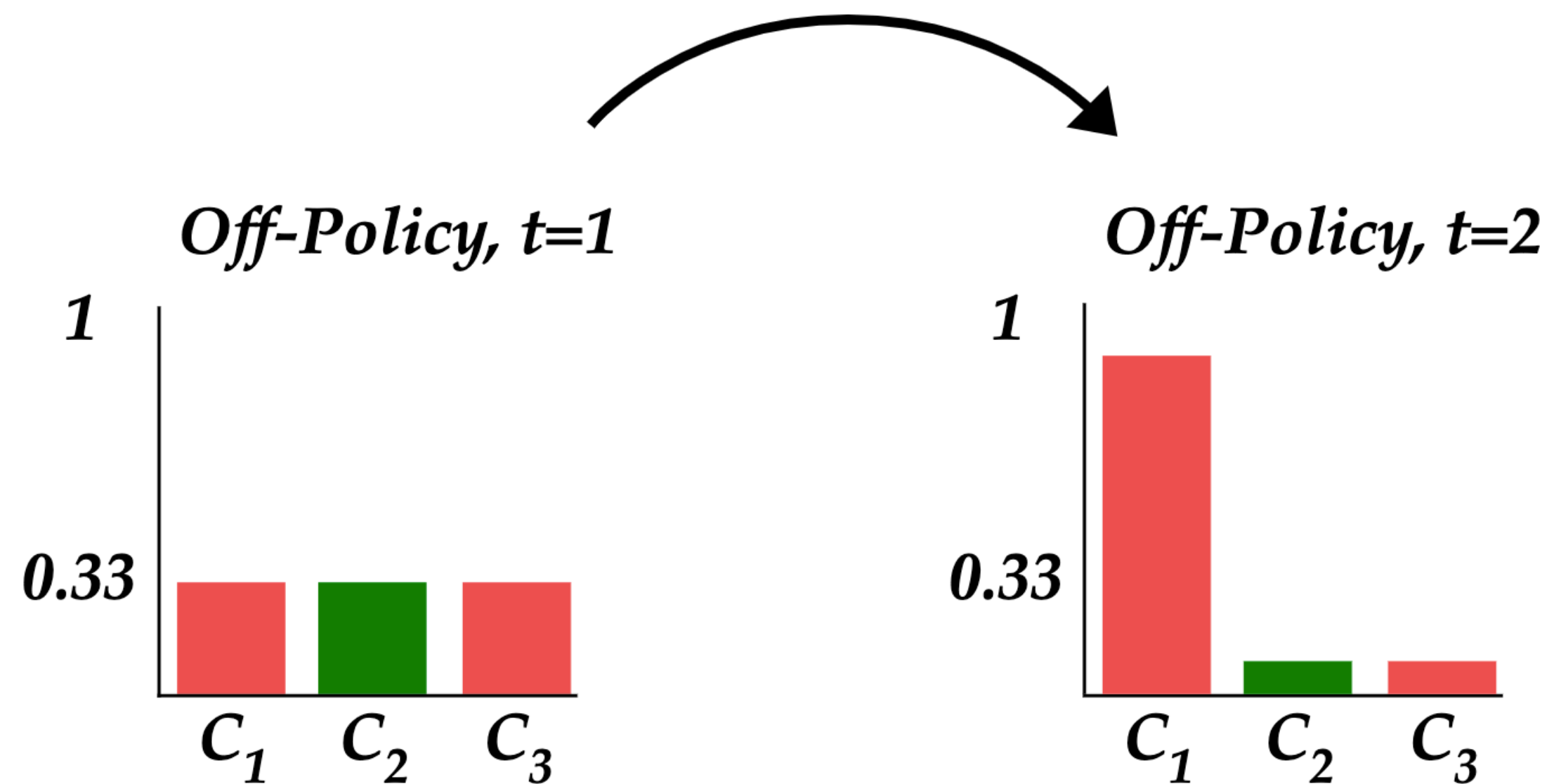
$$p_{\text{on}}(c, h_t) \propto p(\tau; \pi) \propto p(c)p(s_1) \prod_{i=1}^{t-1} \mathcal{T}(s_{i+1} | s_i, a_i, c)$$

$$p_{\text{off}}(c, h_t) \propto p(\tau; \pi^E) \propto p(c)p(s_1) \prod_{i=1}^{t-1} \pi^E(a_i | c, s_i) \mathcal{T}(s_{i+1} | s_i, a_i, c)$$

$$p_{\text{on}}(c, h_t) \propto p(\tau; \pi) \propto p(c)p(s_1) \prod_{i=1}^{t-1} \mathcal{T}(s_{i+1} | s_i, a_i, c)$$

$$p_{\text{off}}(c, h_t) \propto p(\tau; \pi^E) \propto p(c)p(s_1) \prod_{i=1}^{t-1} \pi^E(a_i | c, s_i) \mathcal{T}(s_{i+1} | s_i, a_i, c)$$

Learner picks arm 1 randomly



Theorem (informal): Off-policy learners have a value difference to the expert bounded by the sum of their errors (tight) while on-policy learners have one dependent on their asymptotic error.

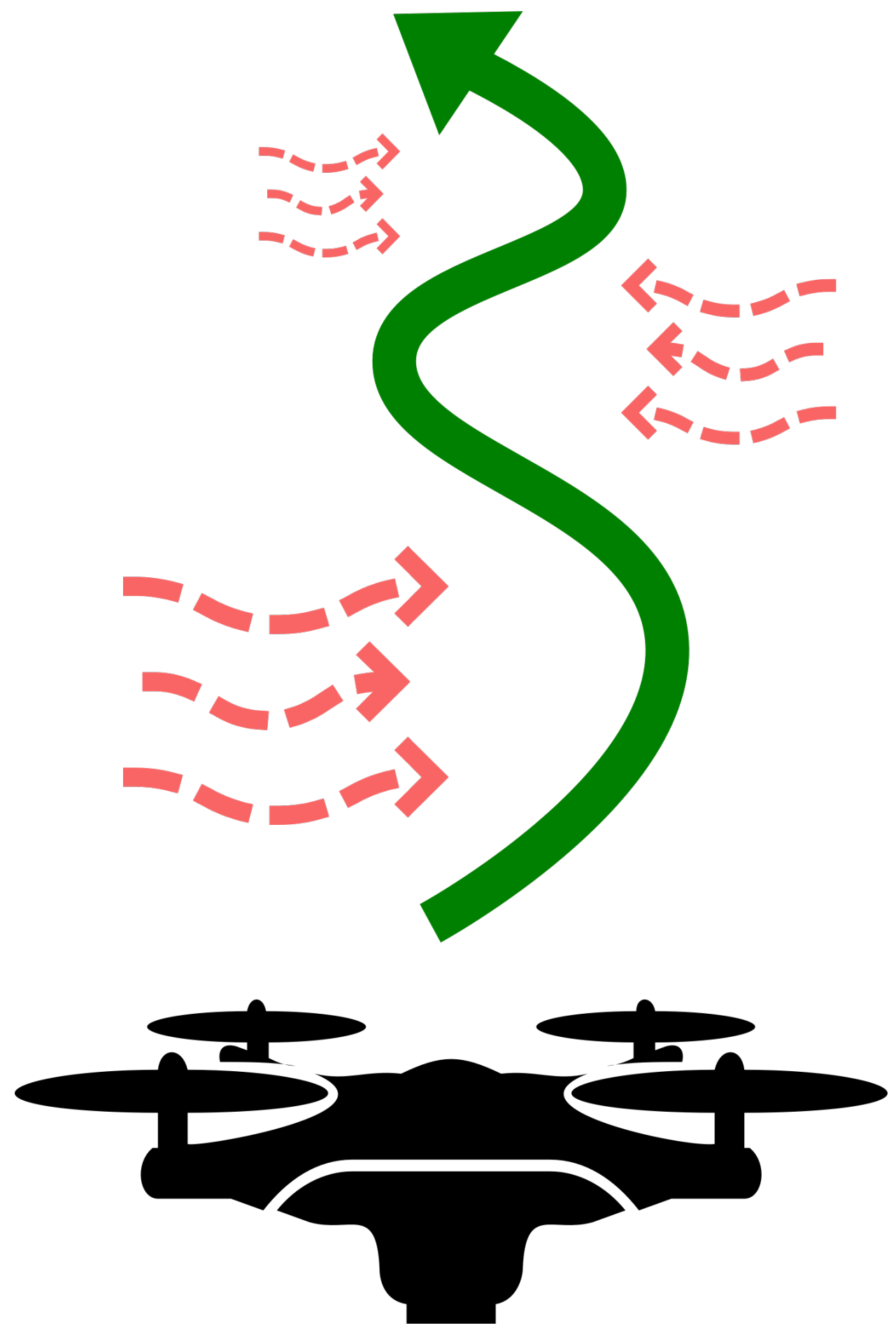
	Offline	Online	Interactive
Covariate Shift	✗	✓	✓
Hidden Context	✗	✓ w/ History	✓ w/ History
TCN			

*“Actually, since we were fitting a model to a time-series, **samples tend to be correlated in time** [...] Thus, when leaving out a sample in cross validation, **we actually left out a large window (16 seconds) of data** around that sample, to diminish this bias.”*

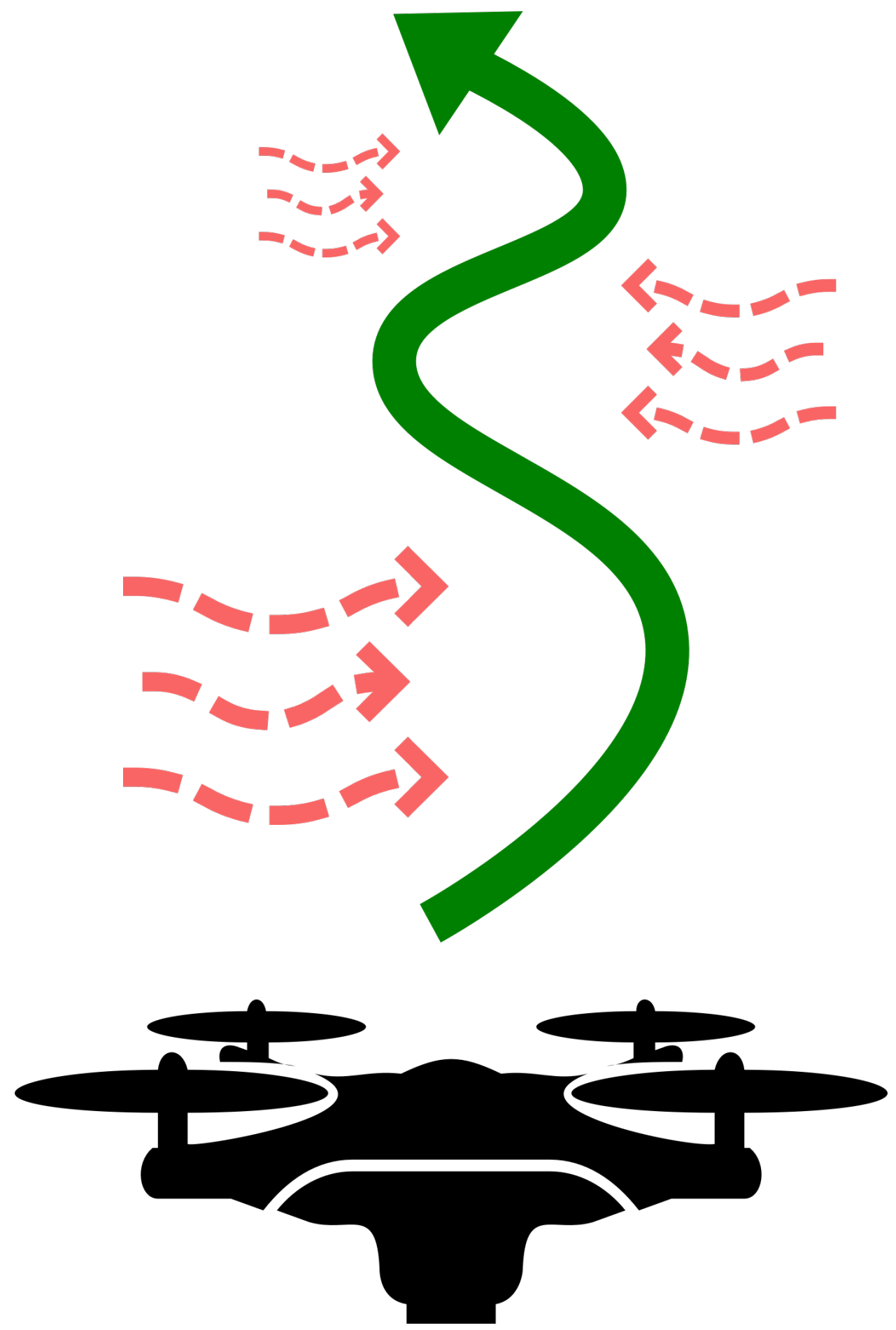
— Ng et al., 2003



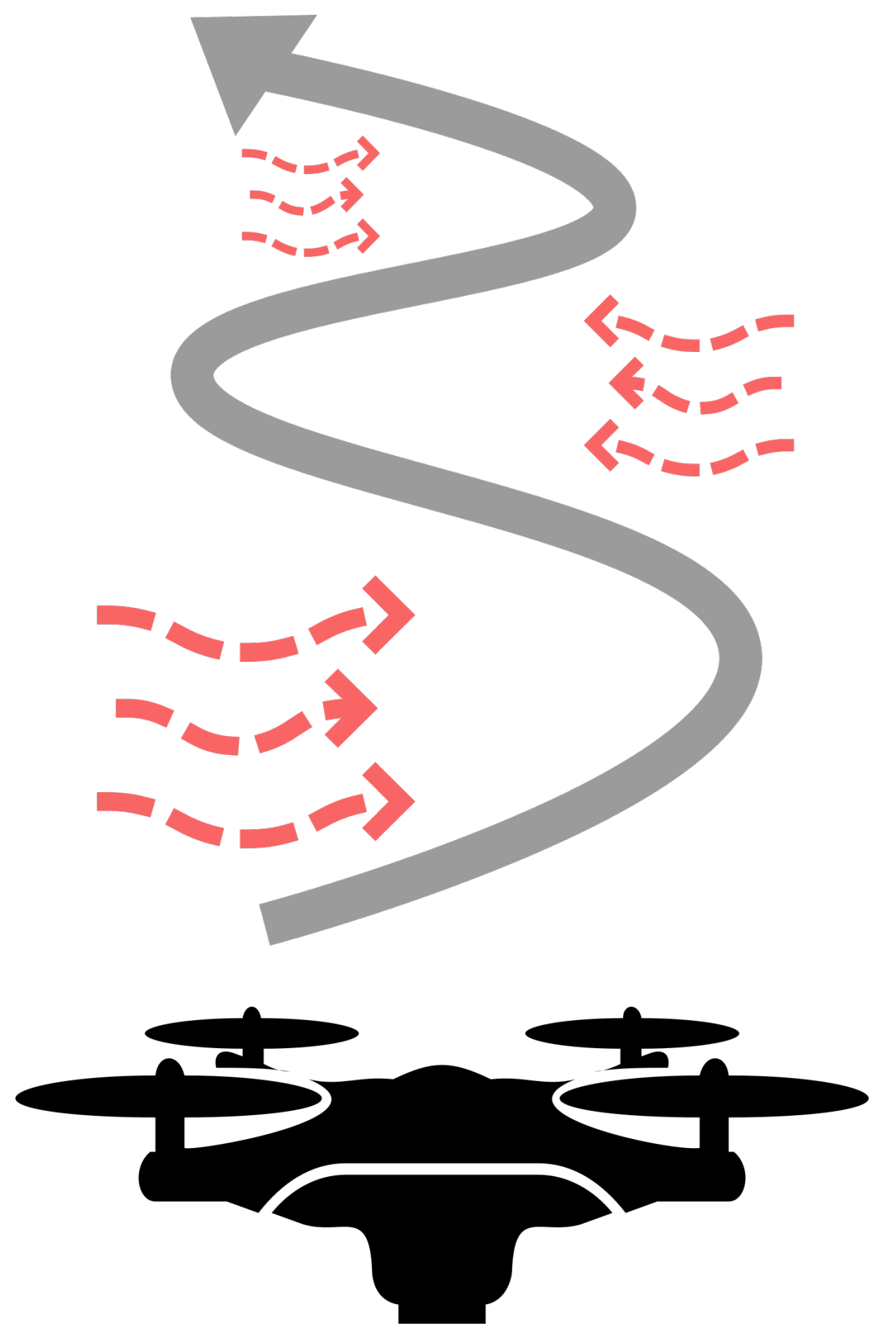
πE



π_E

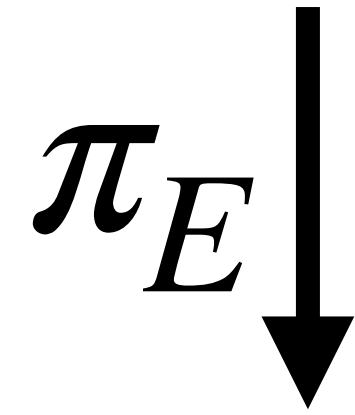


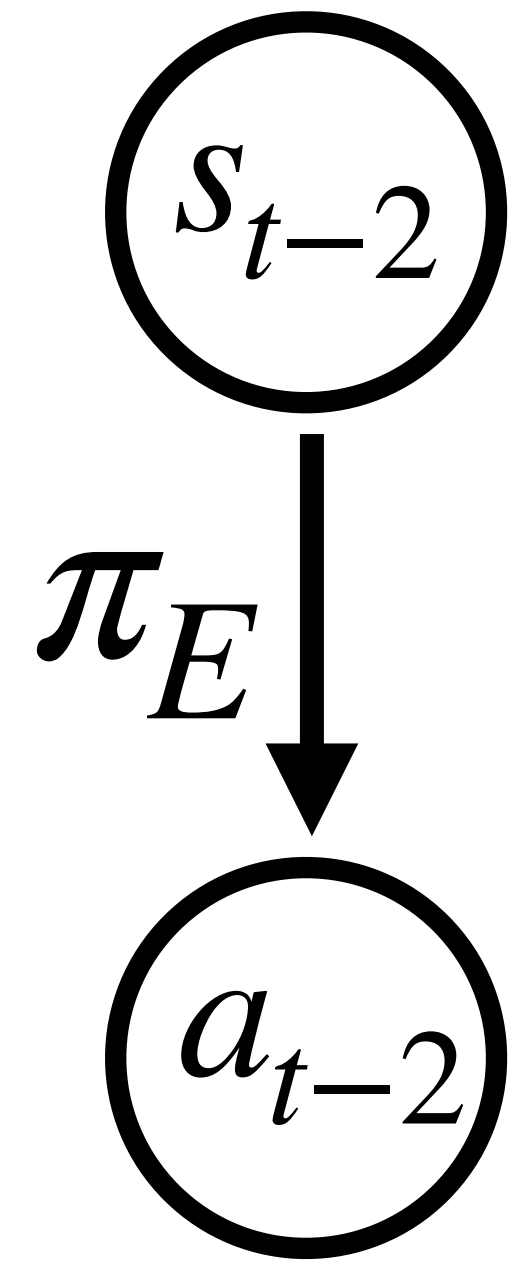
π_{BC}

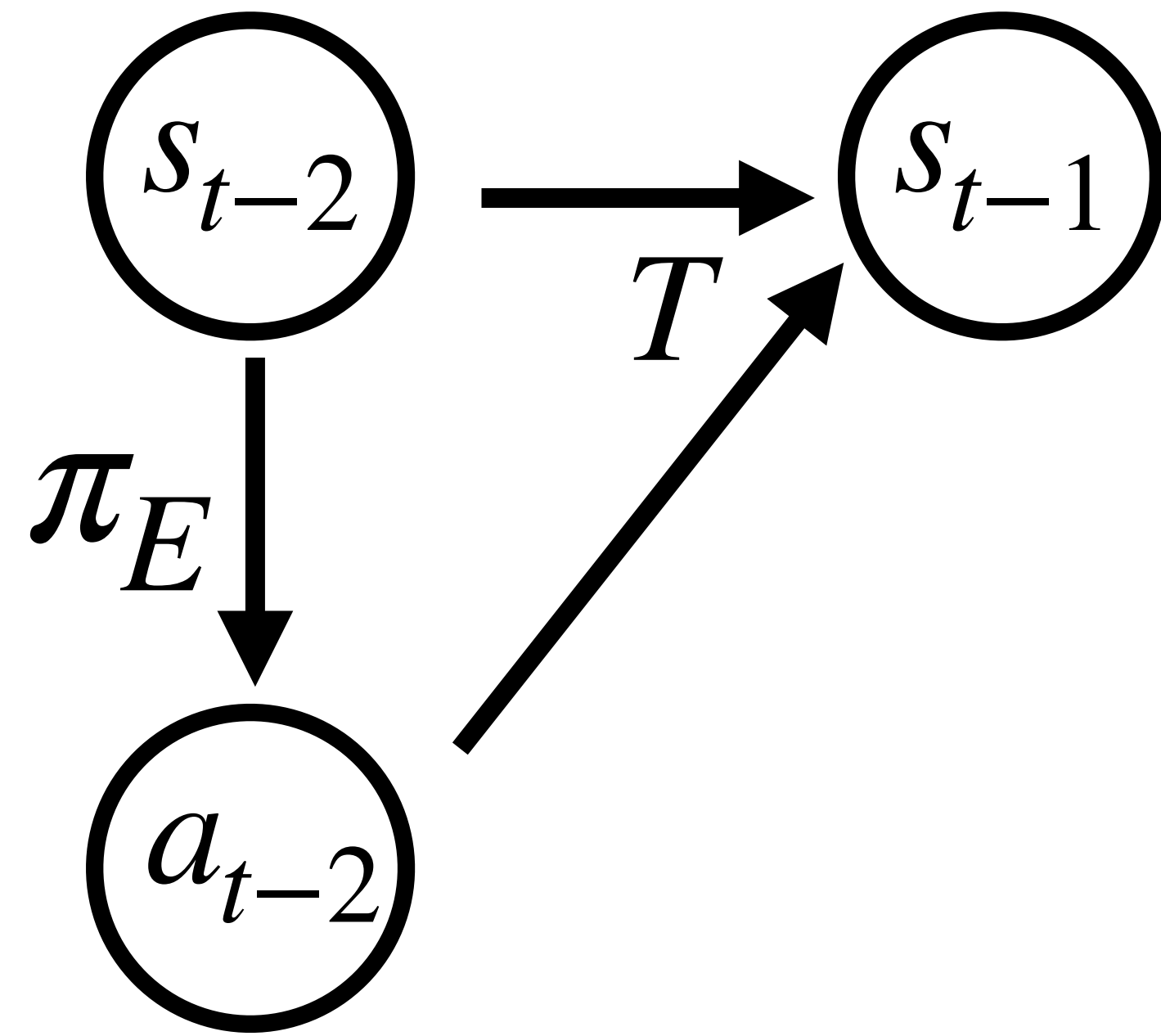


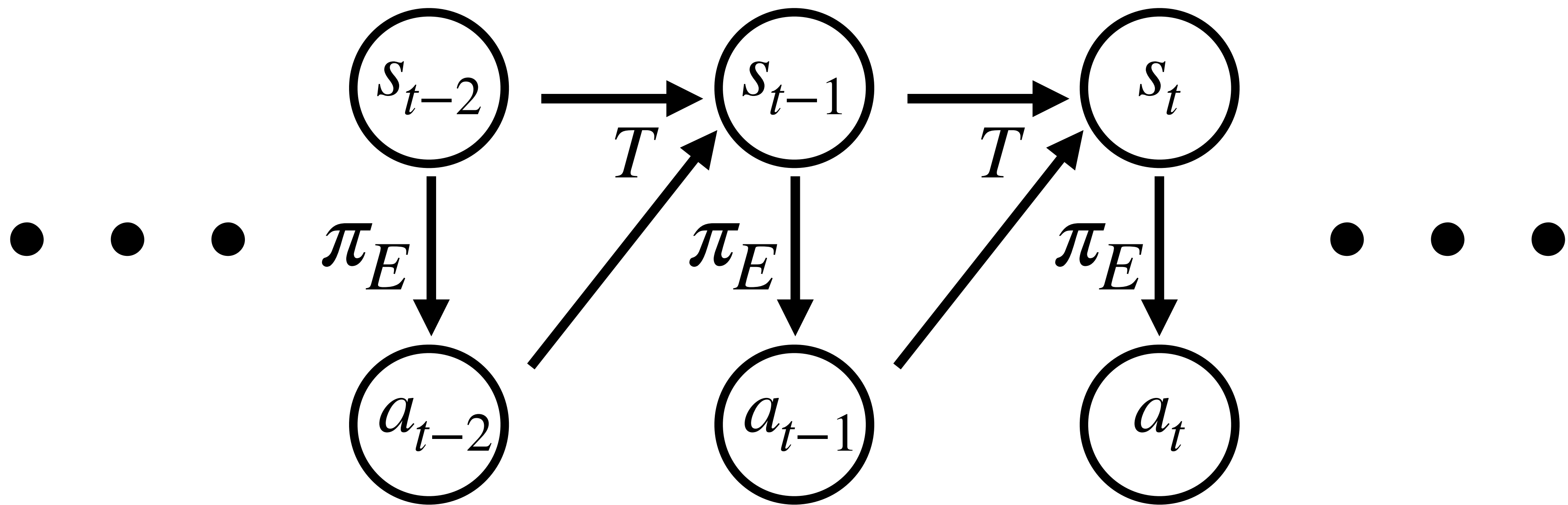
$$S_{t-2}$$

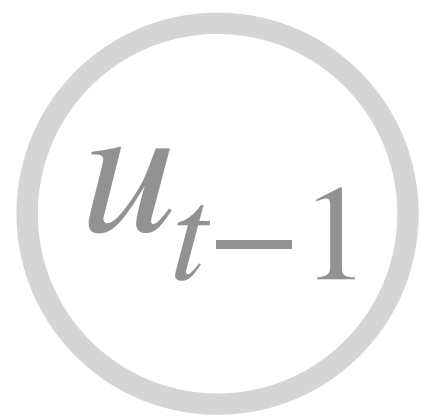
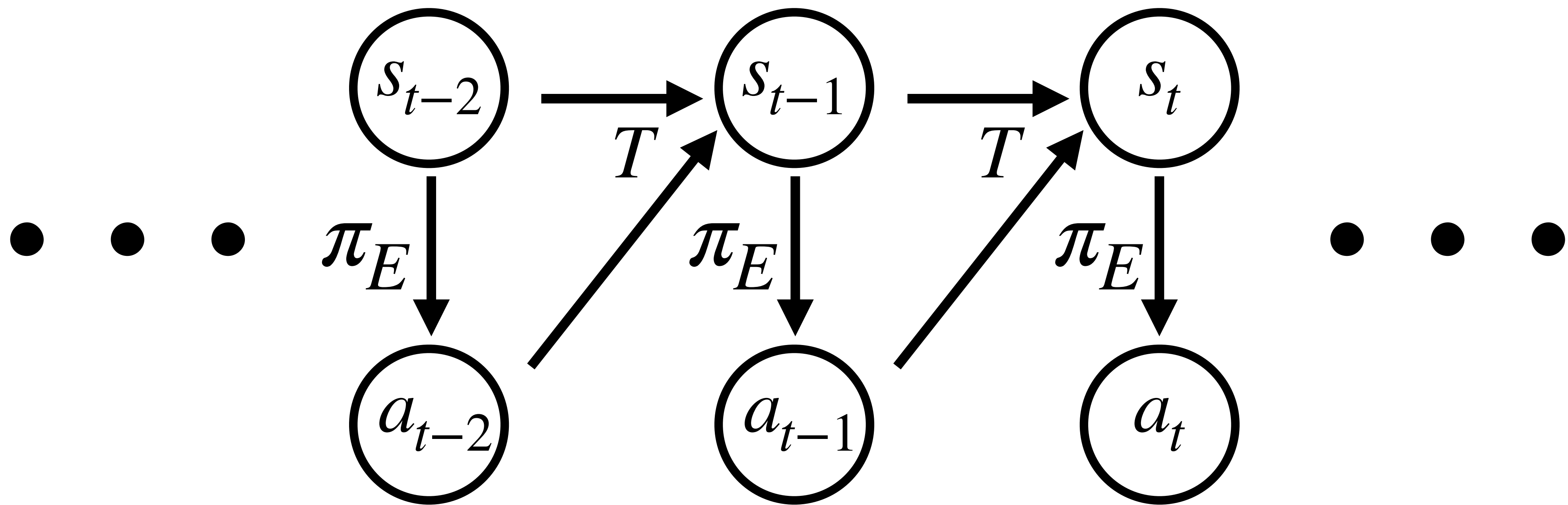
S_{t-2}

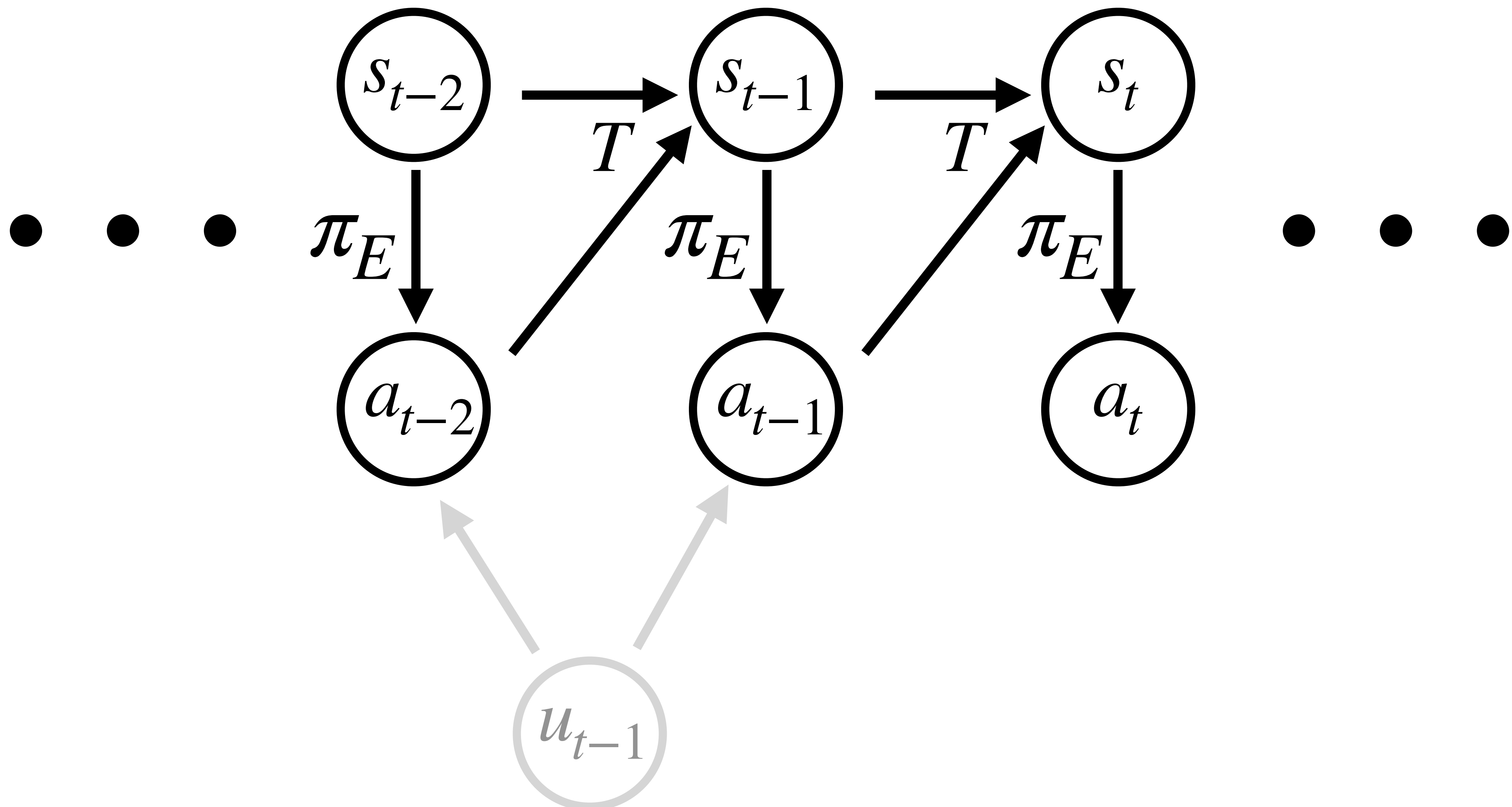


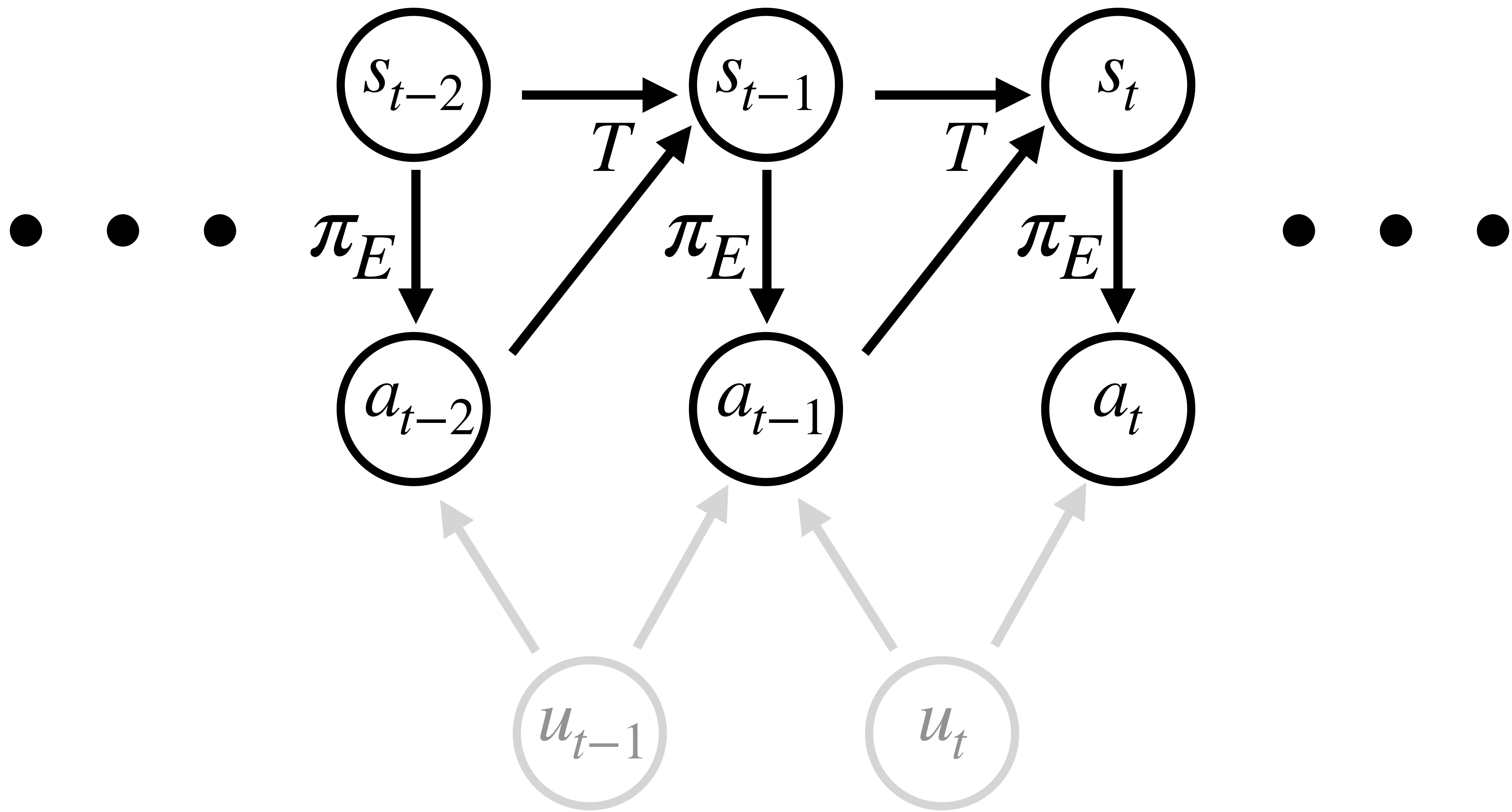


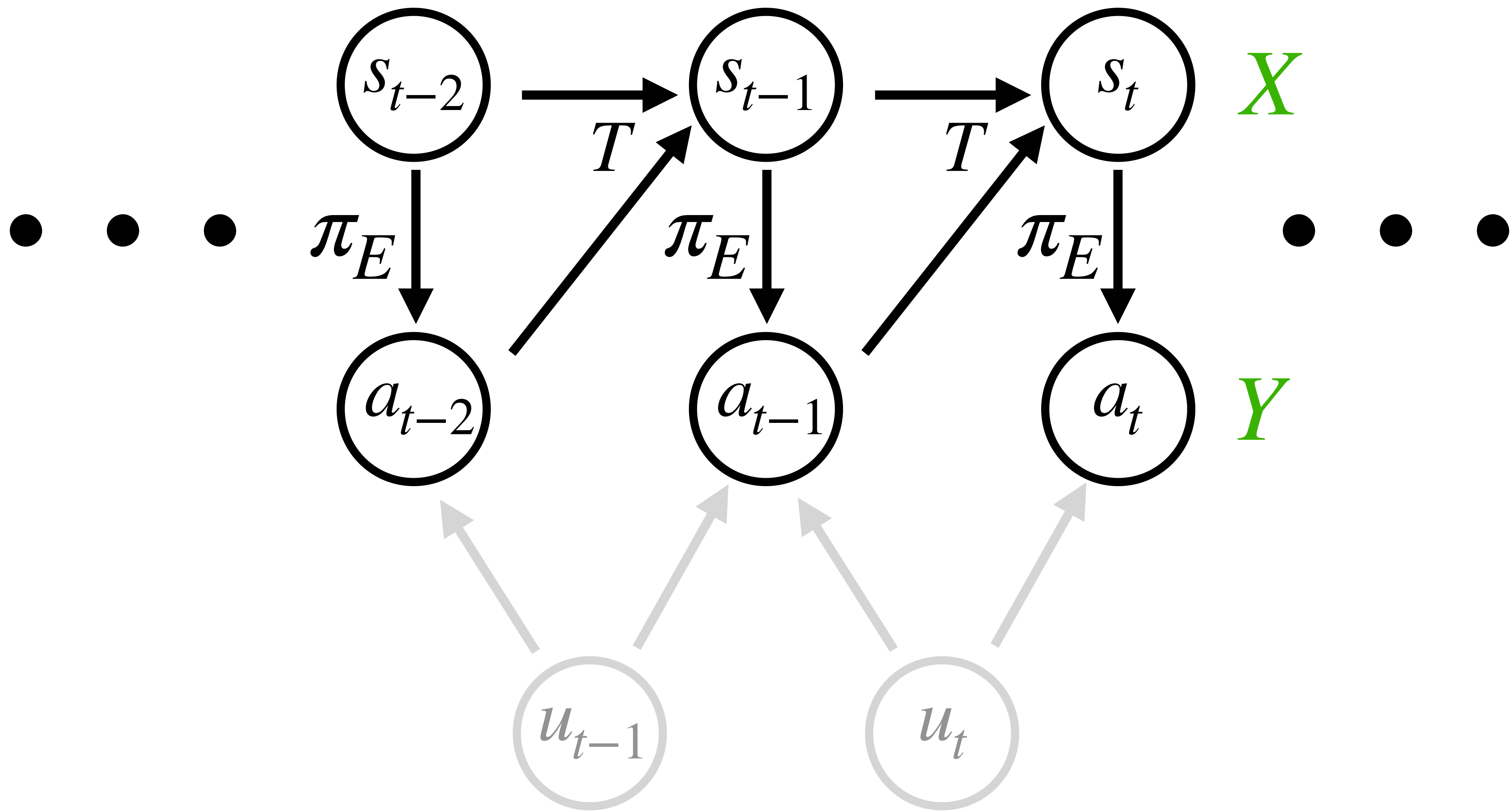


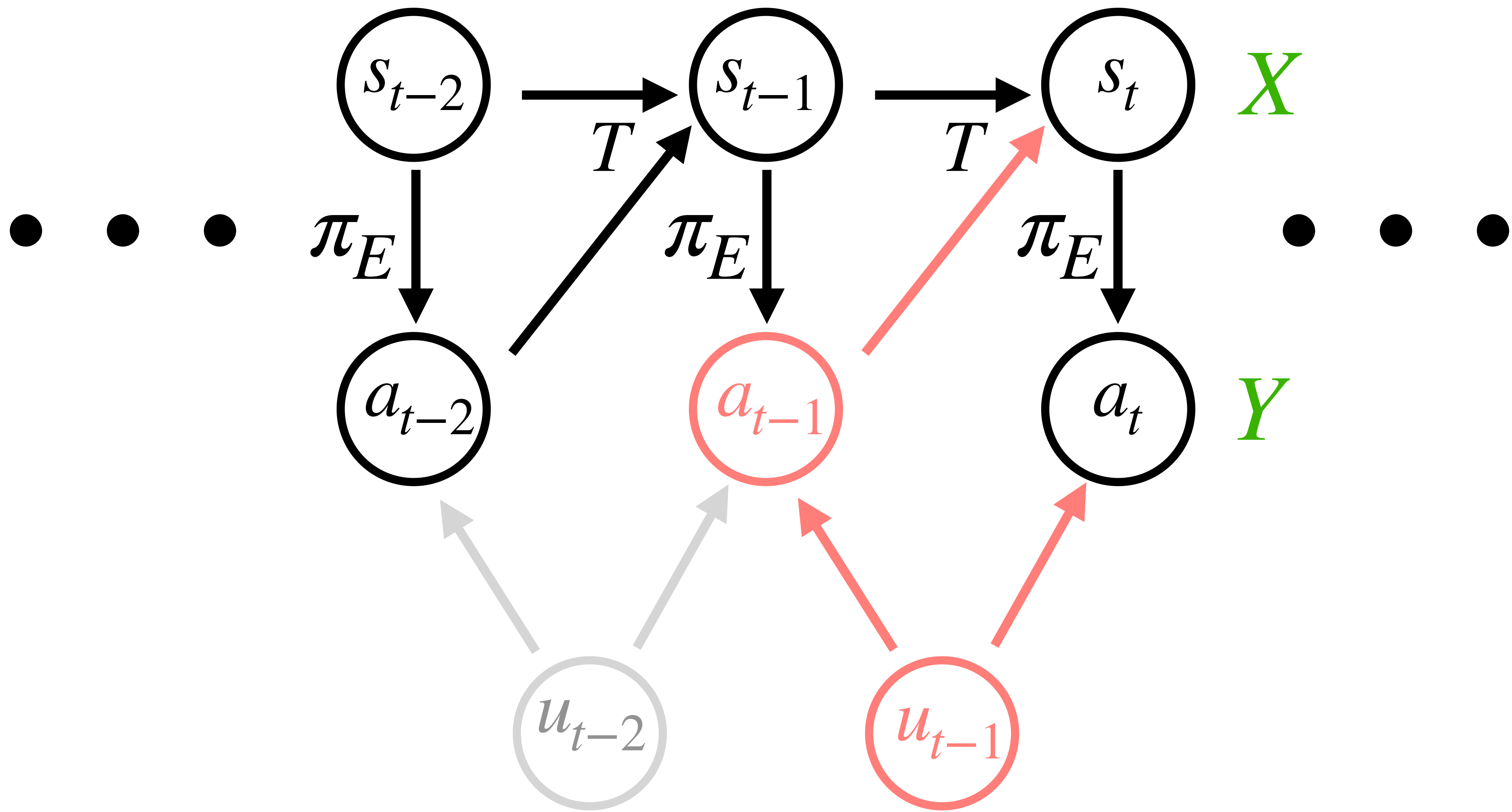




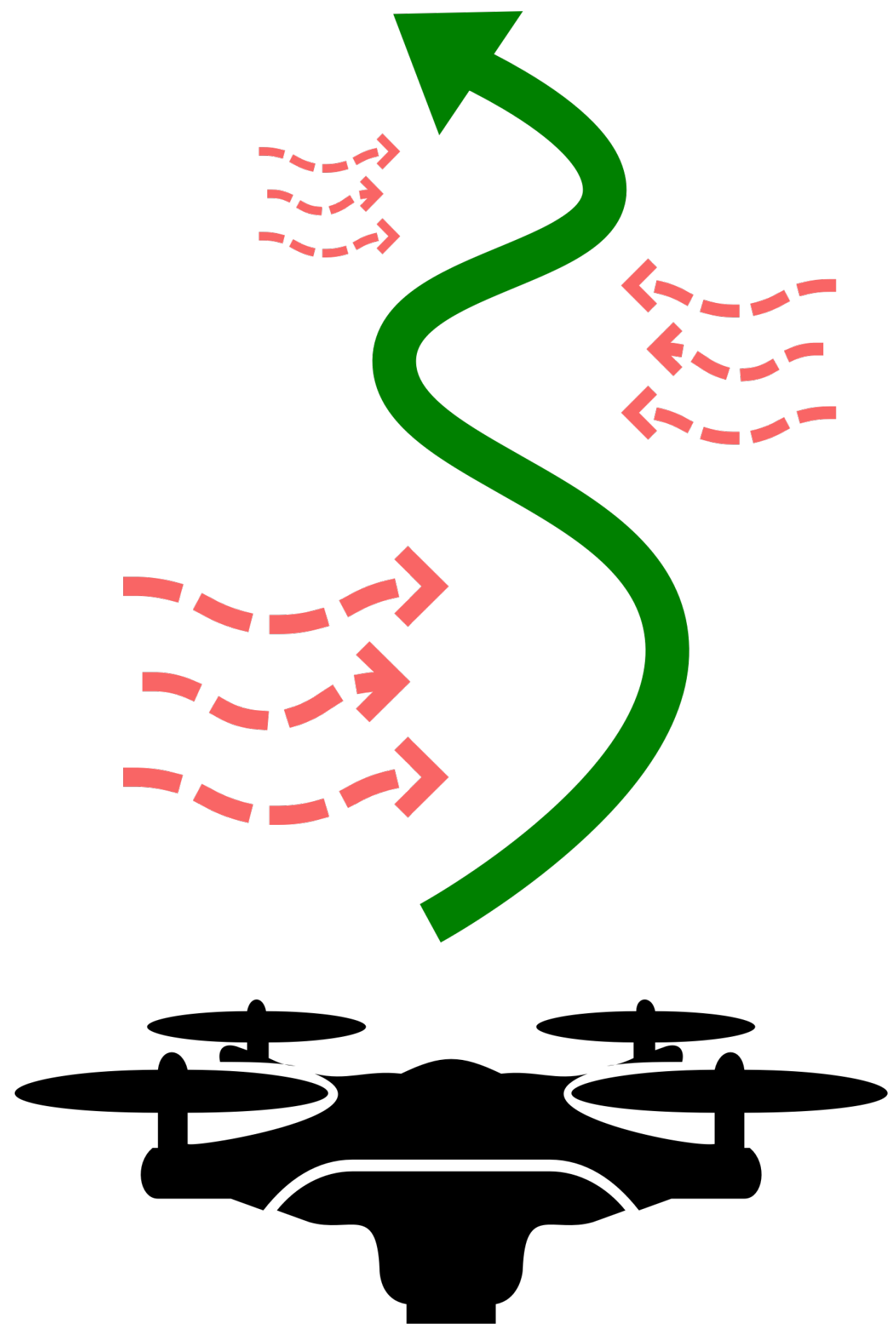




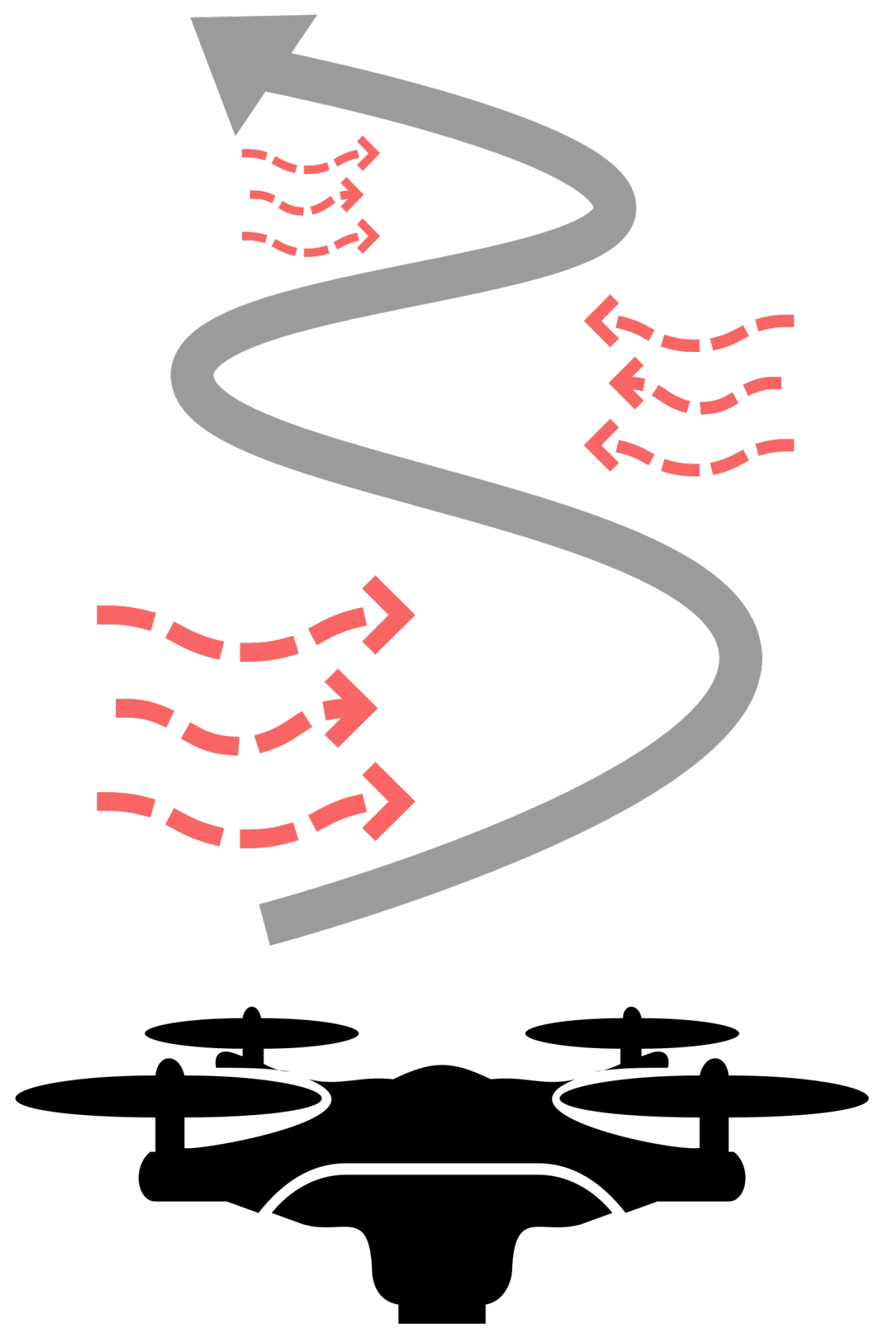




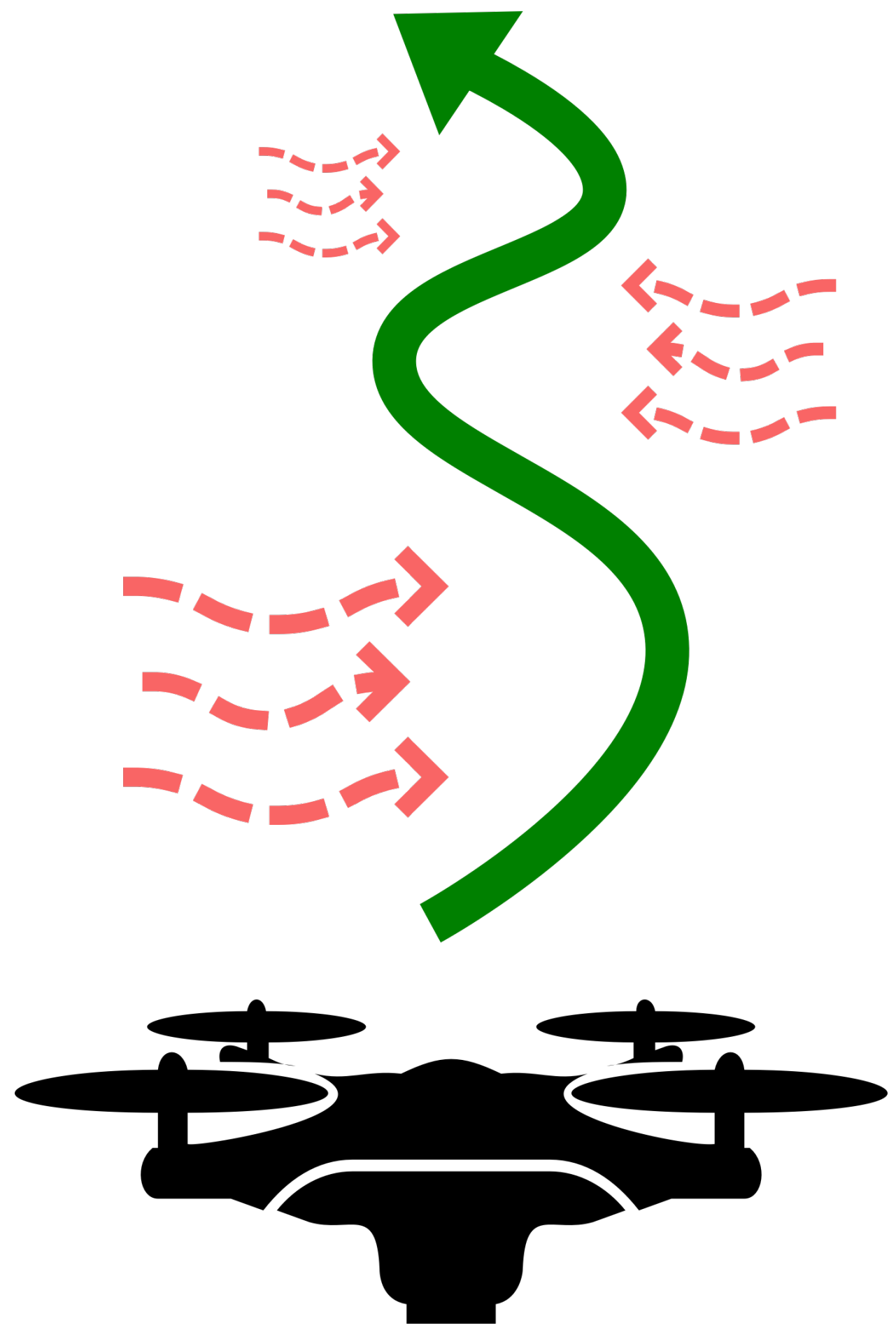
π_E



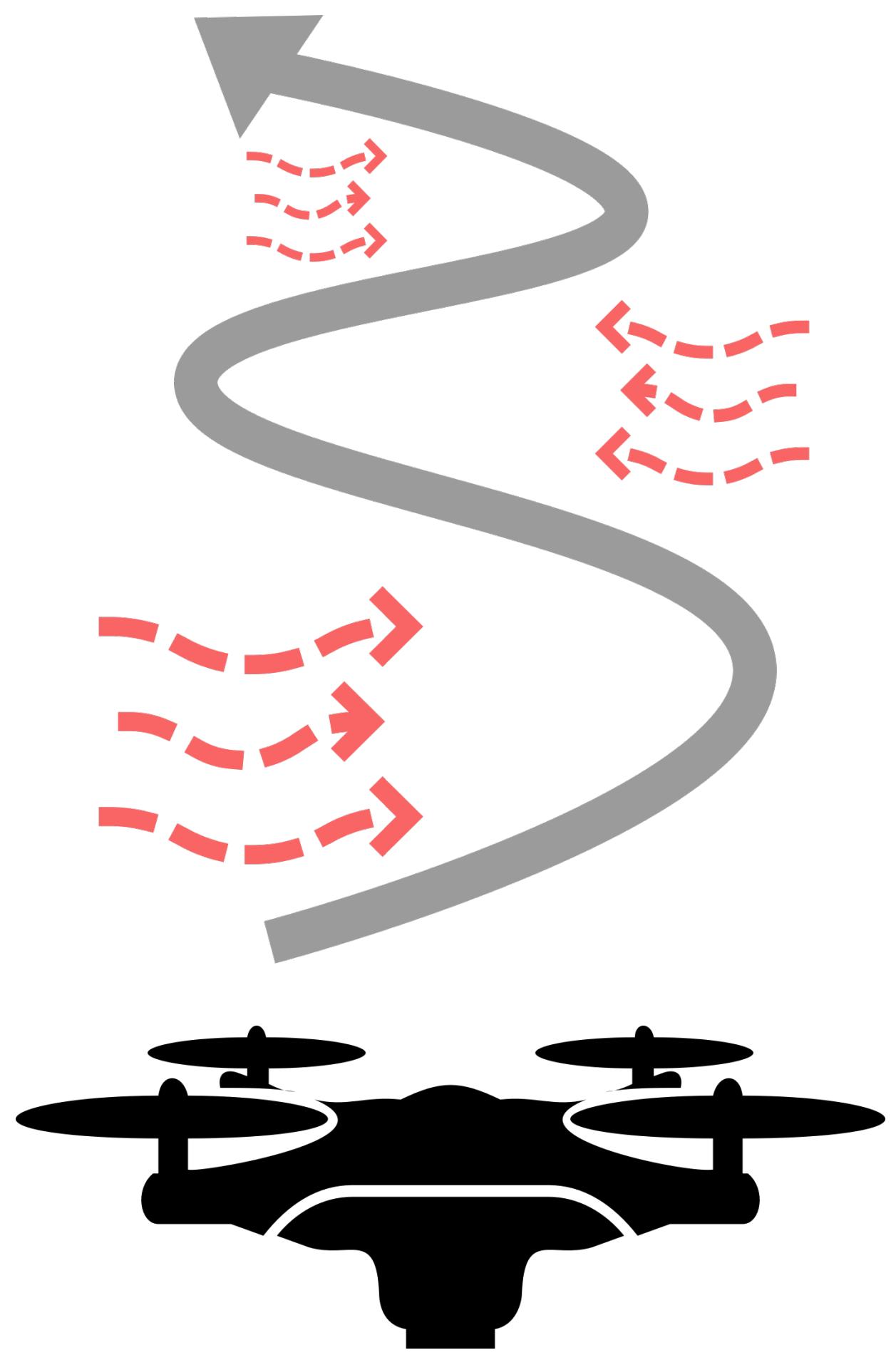
π_{BC}



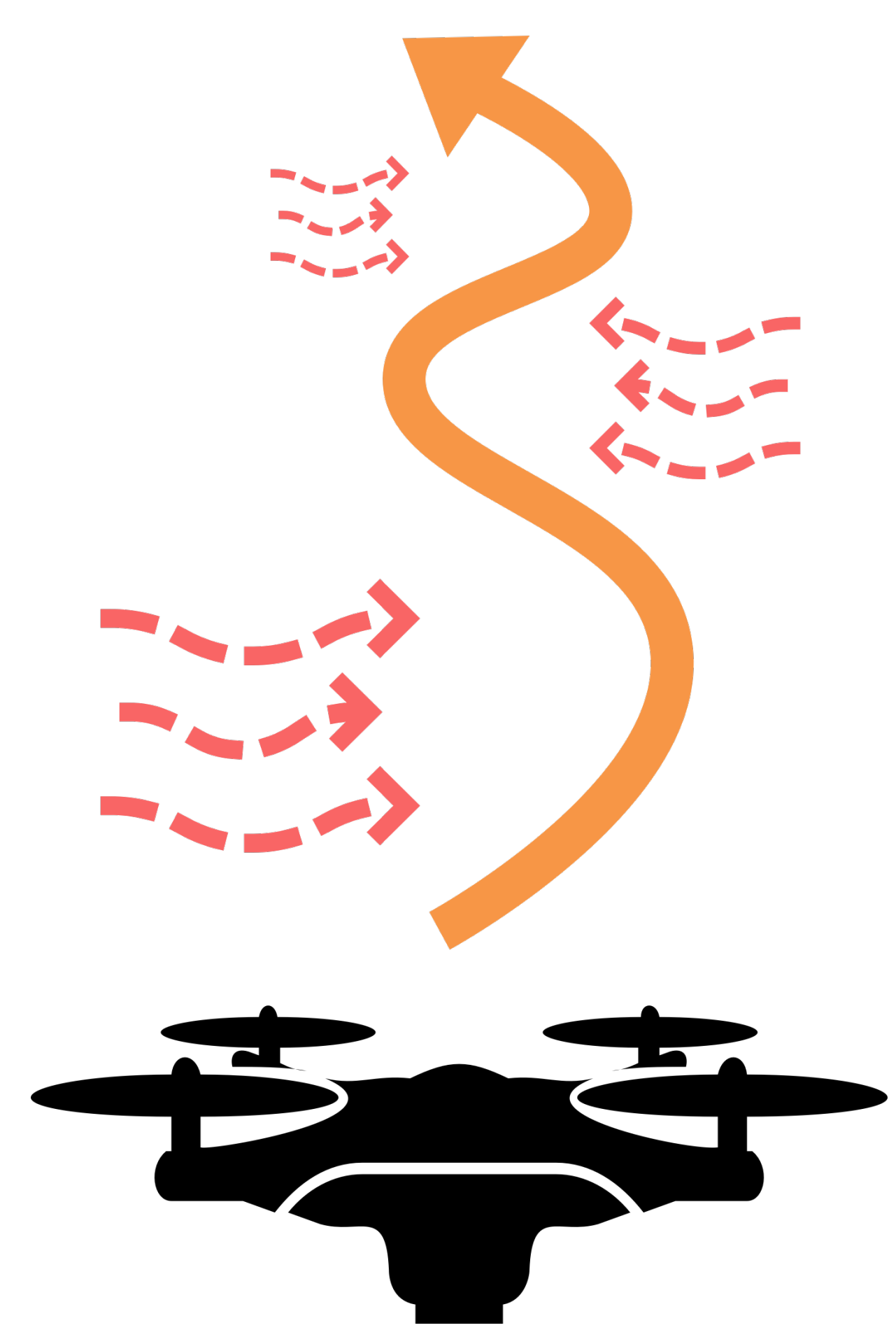
π_E

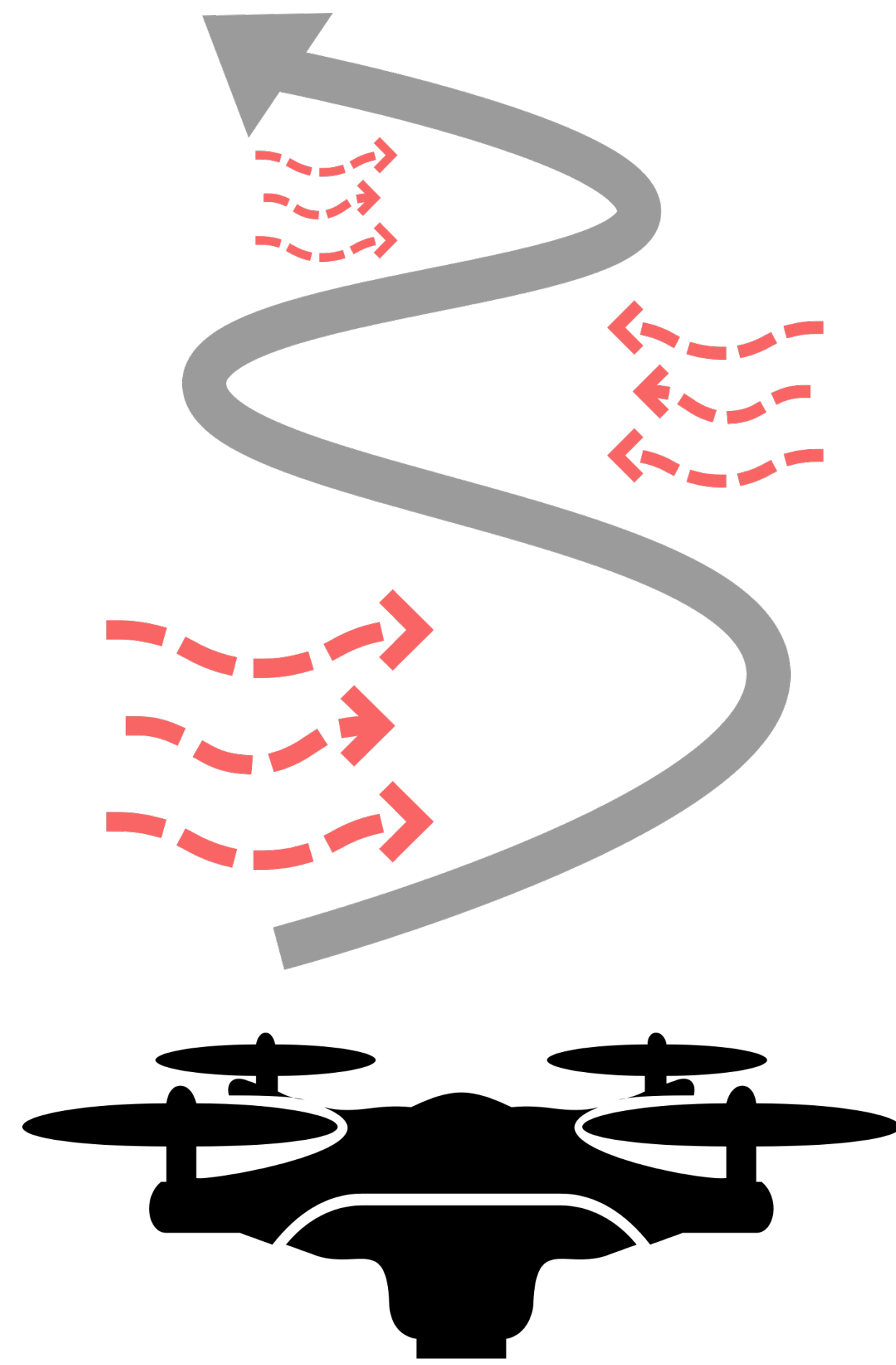


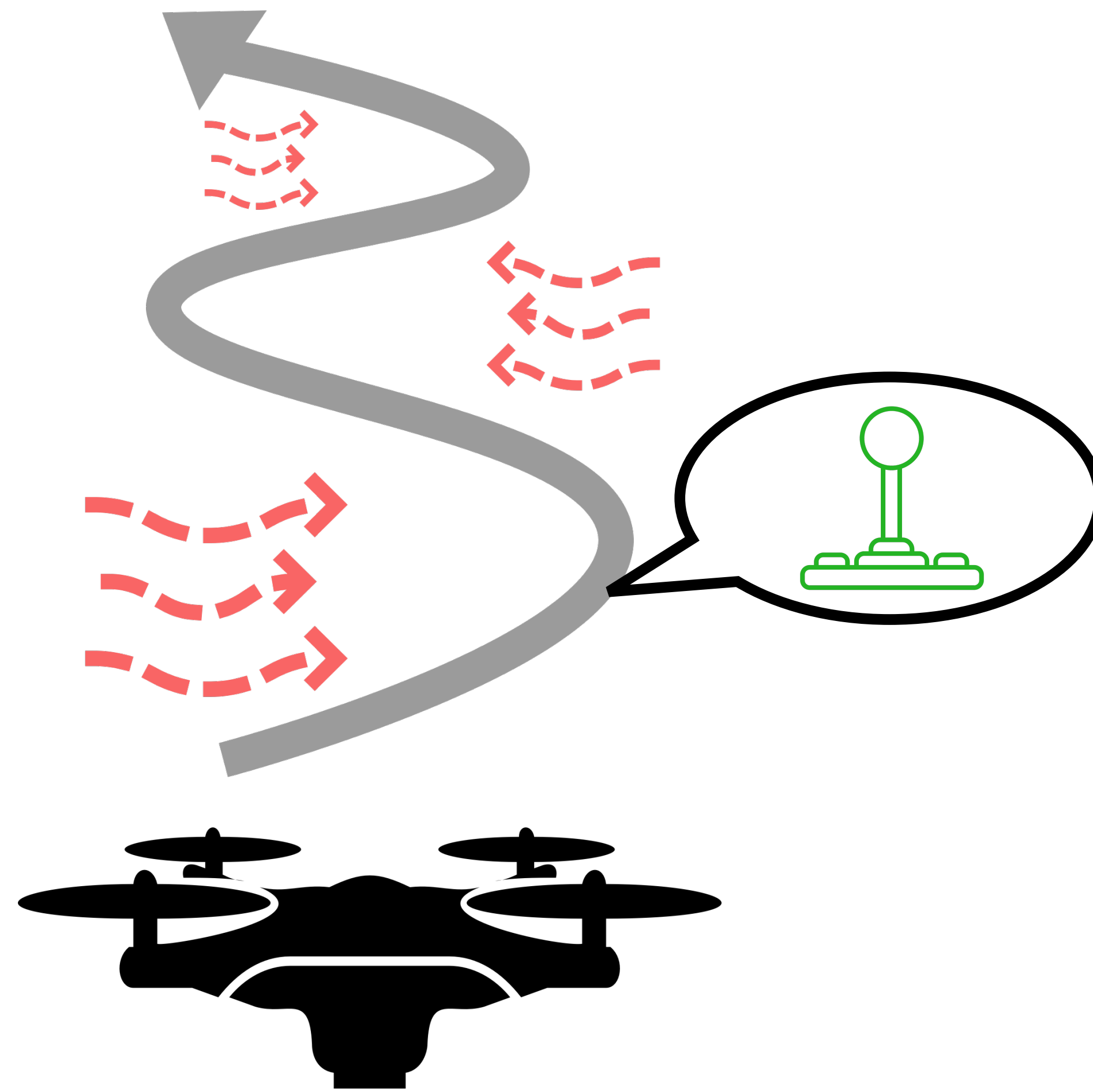
π_{BC}

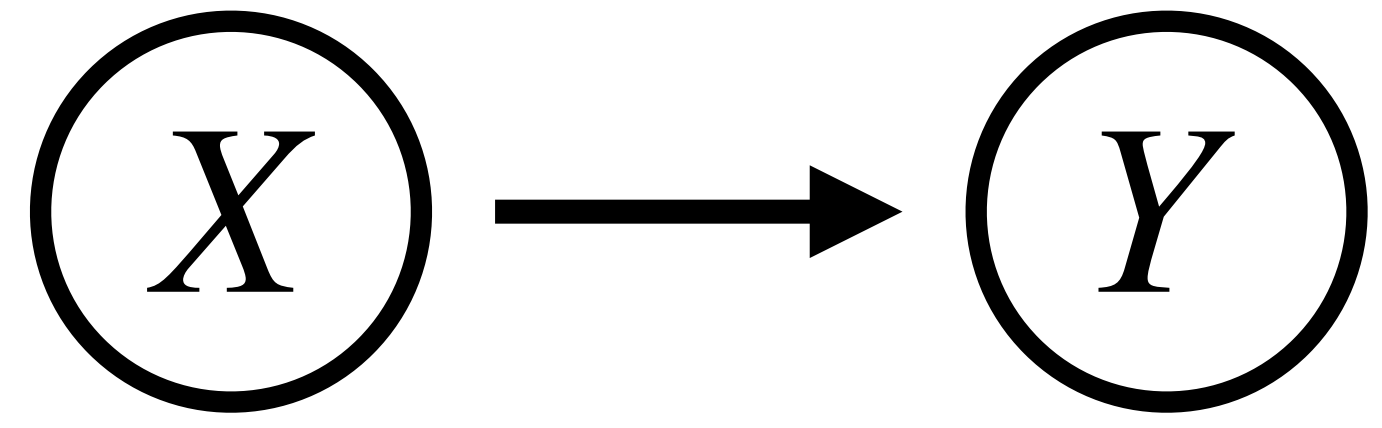


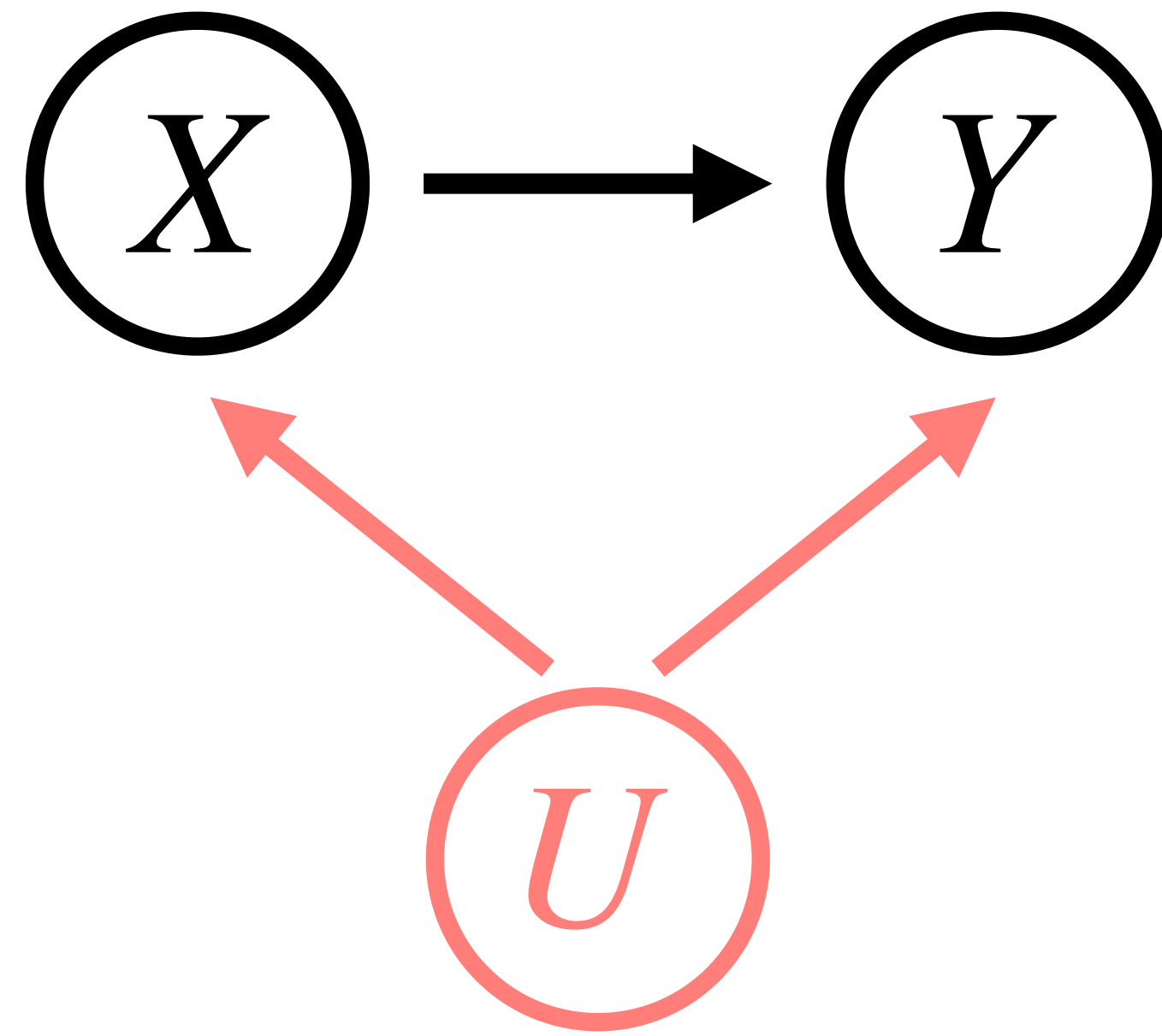
π

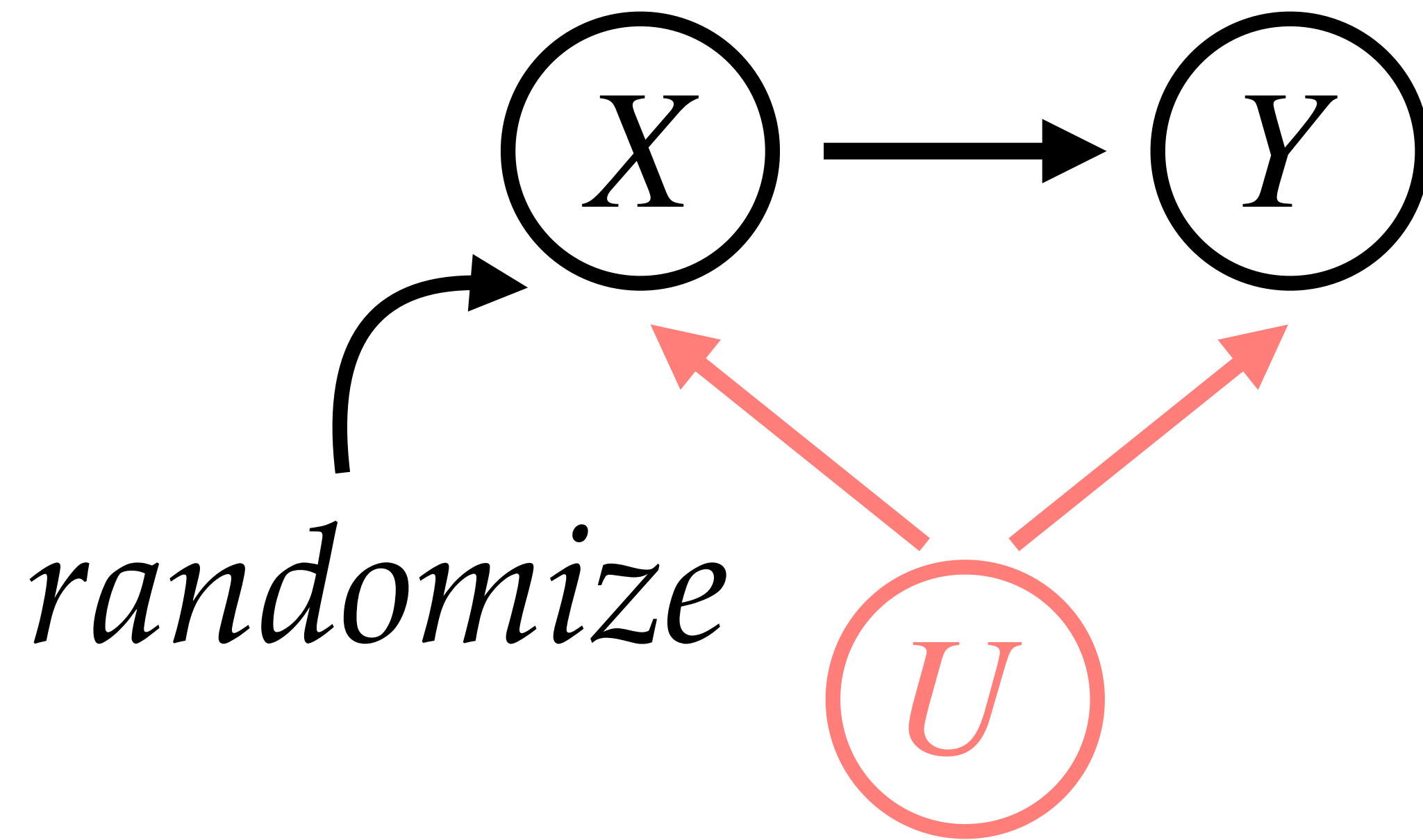


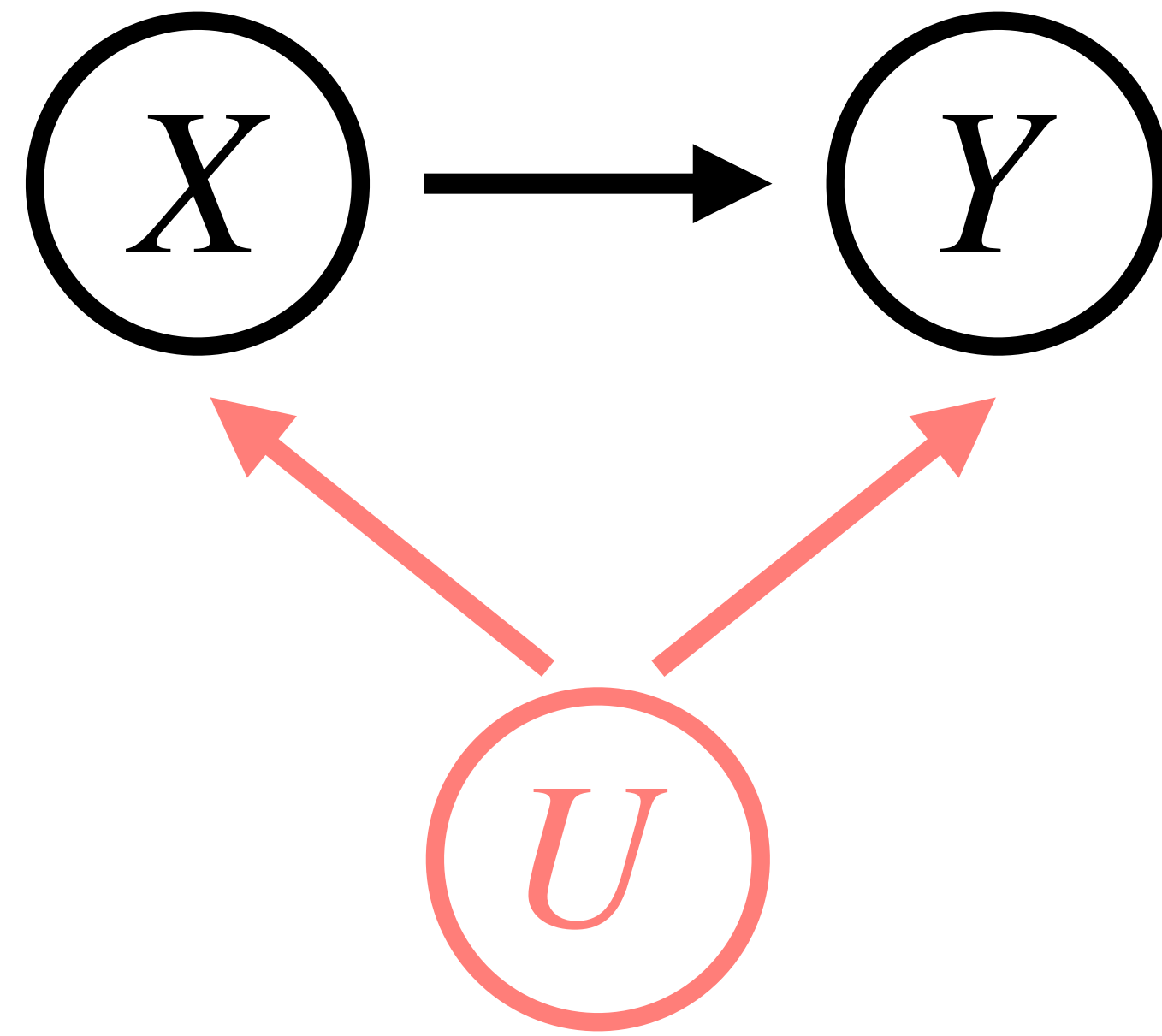


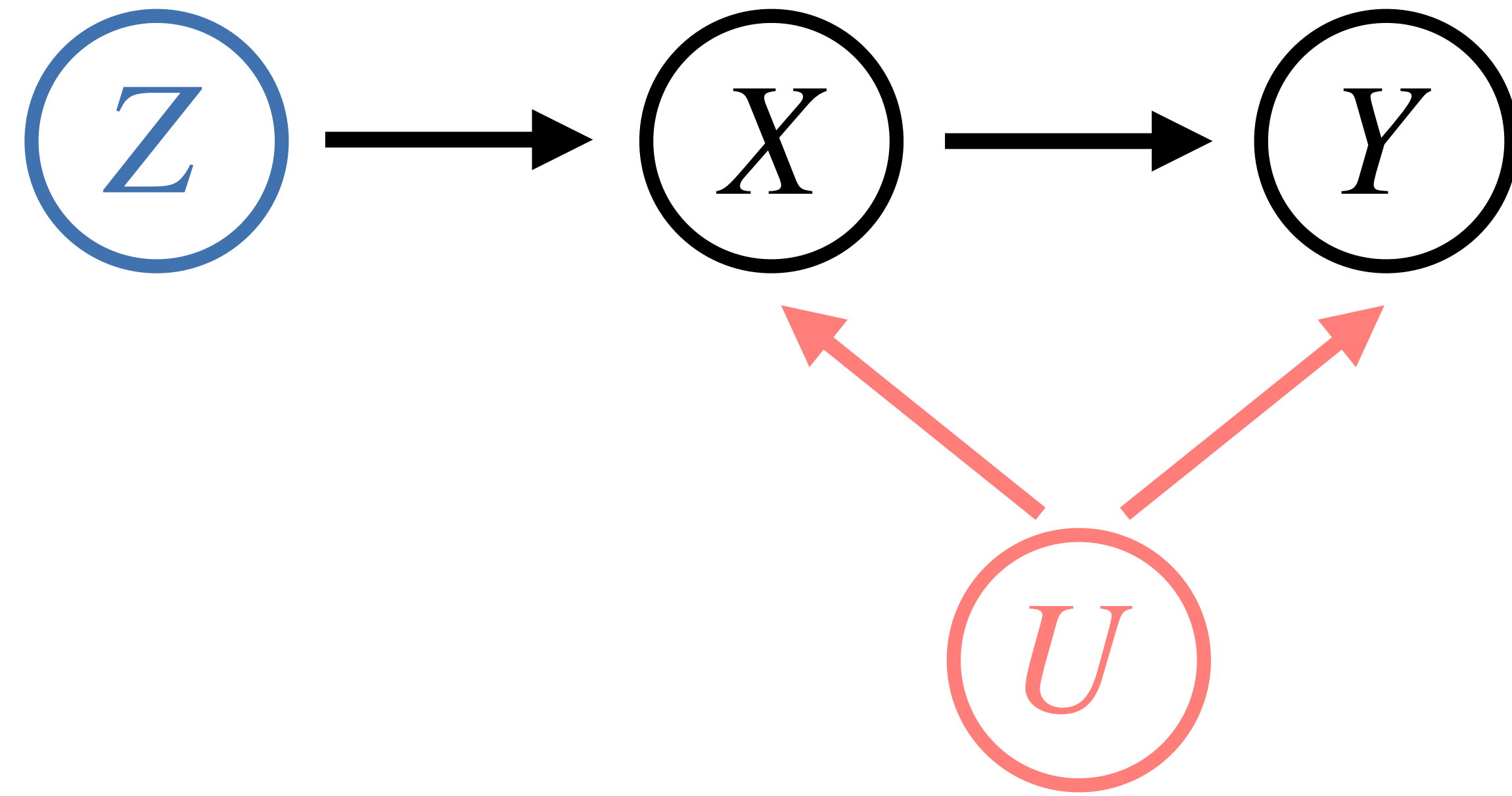


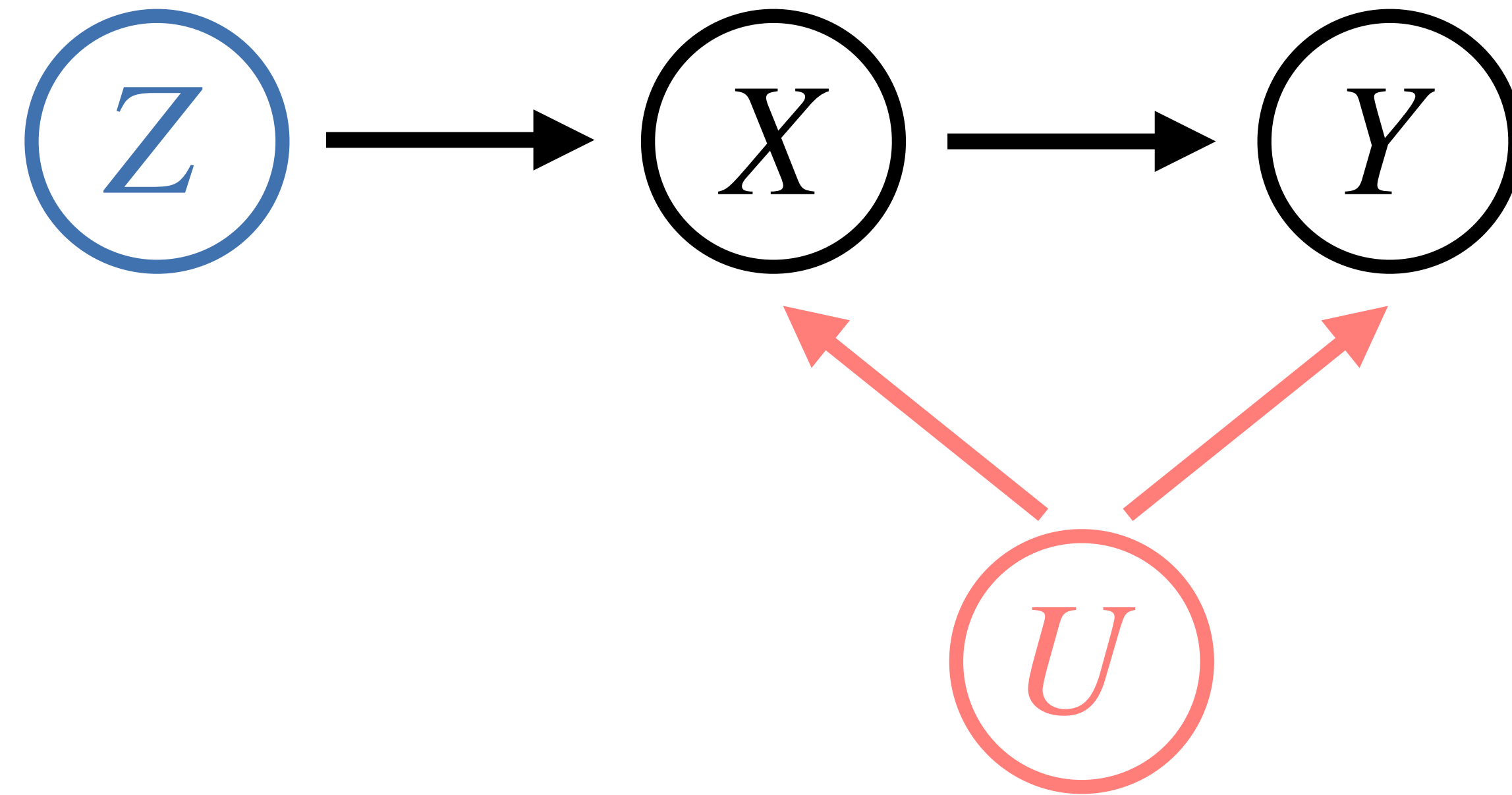




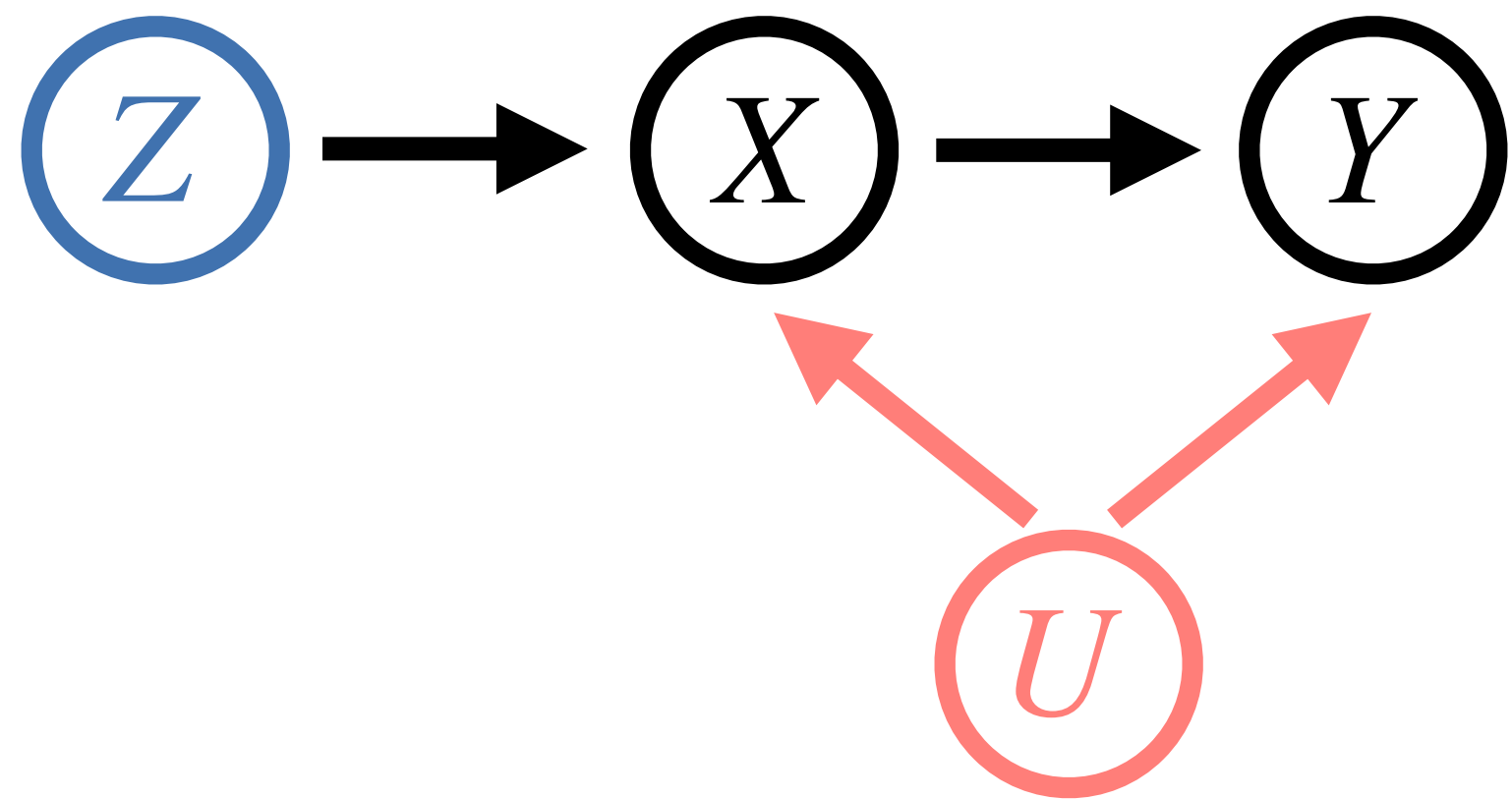


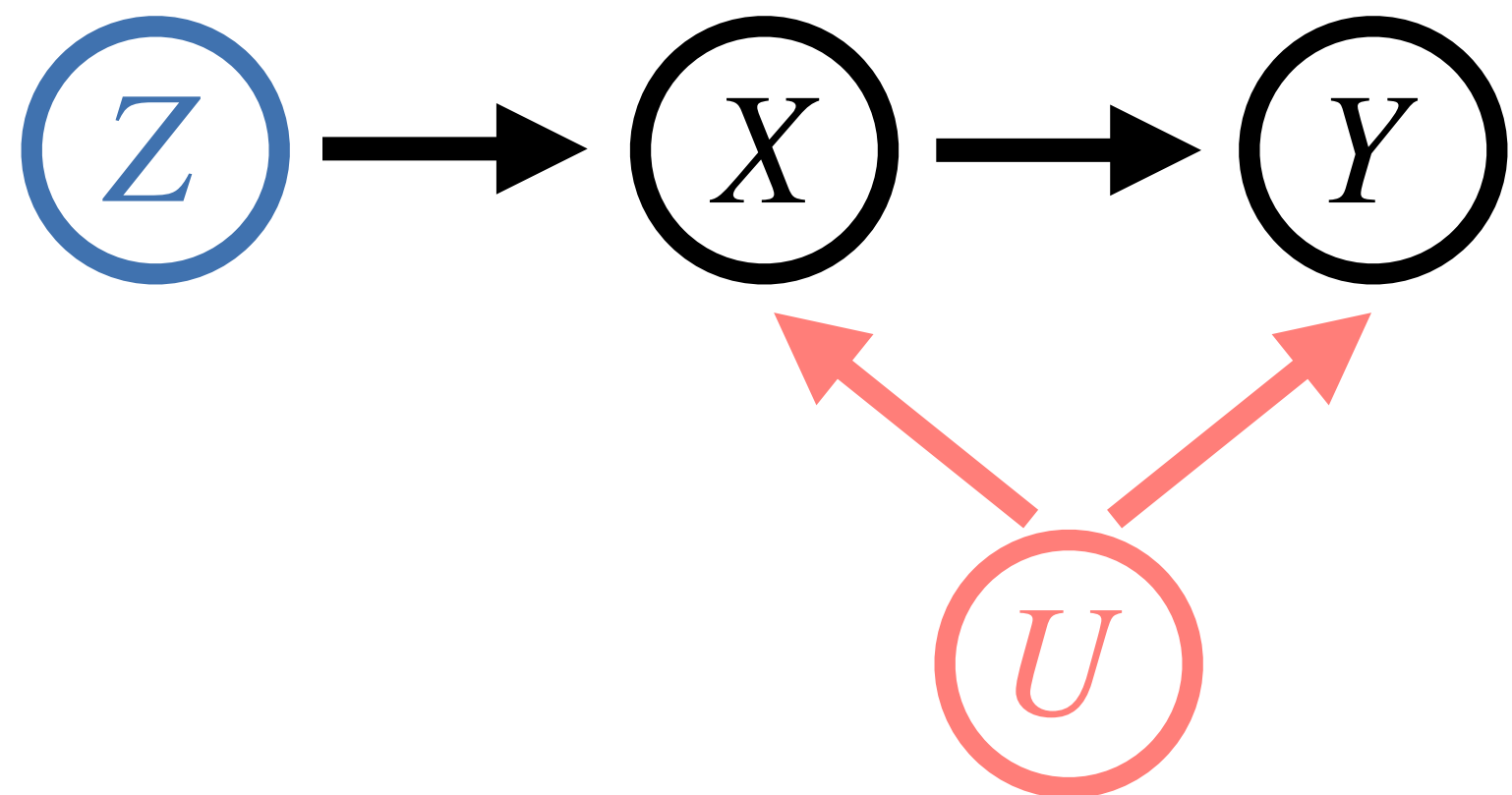




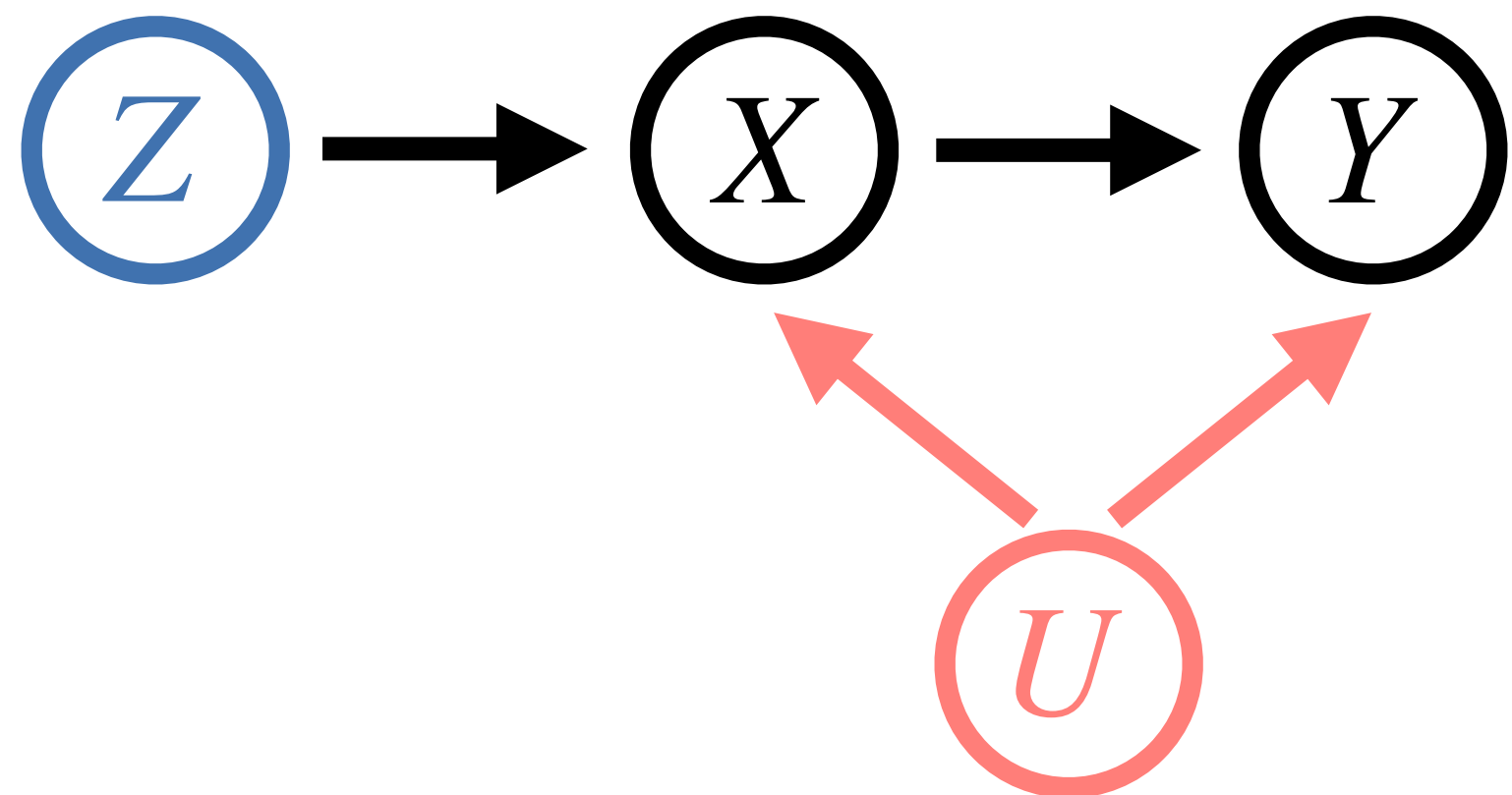


Key Idea: We can condition on instrument Z to counter the effect of confounder U on X .



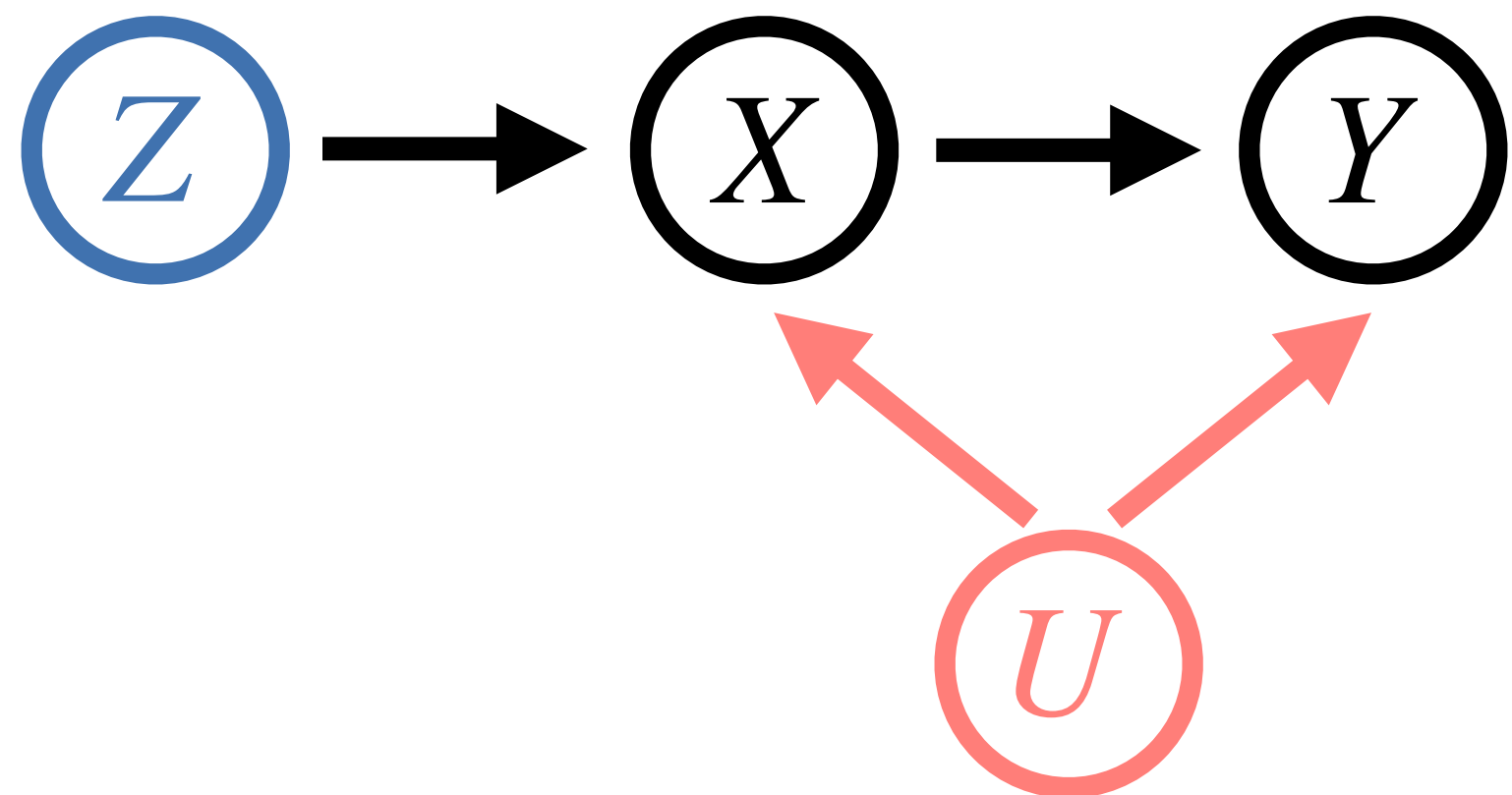


$$X = g(Z, U)$$



$$X = g(Z, U)$$

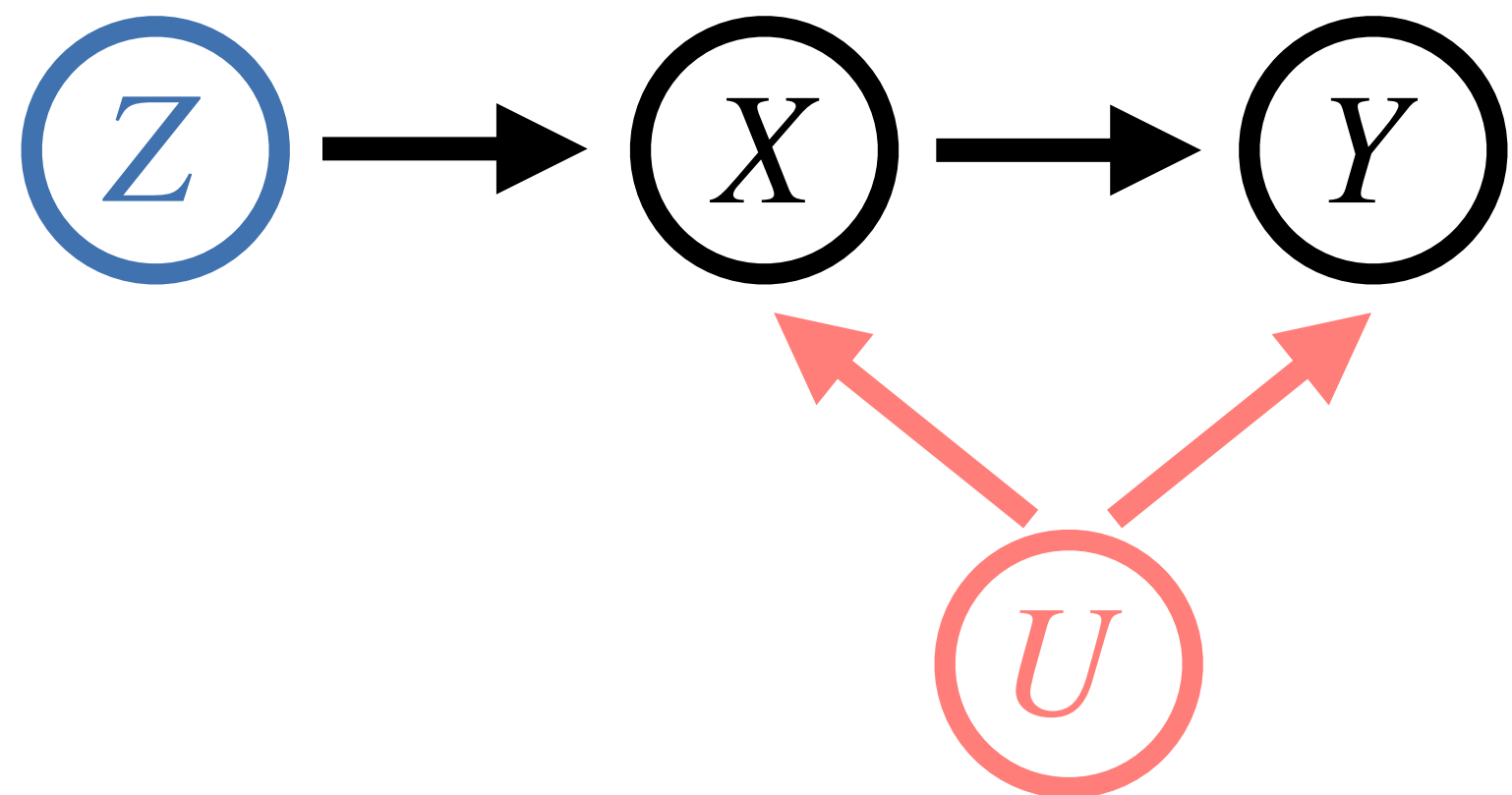
$$Y = h(X) + U$$



$$X = g(Z, U)$$

$$Y = h(X) + U$$

$$\mathbb{E}[U] = 0$$

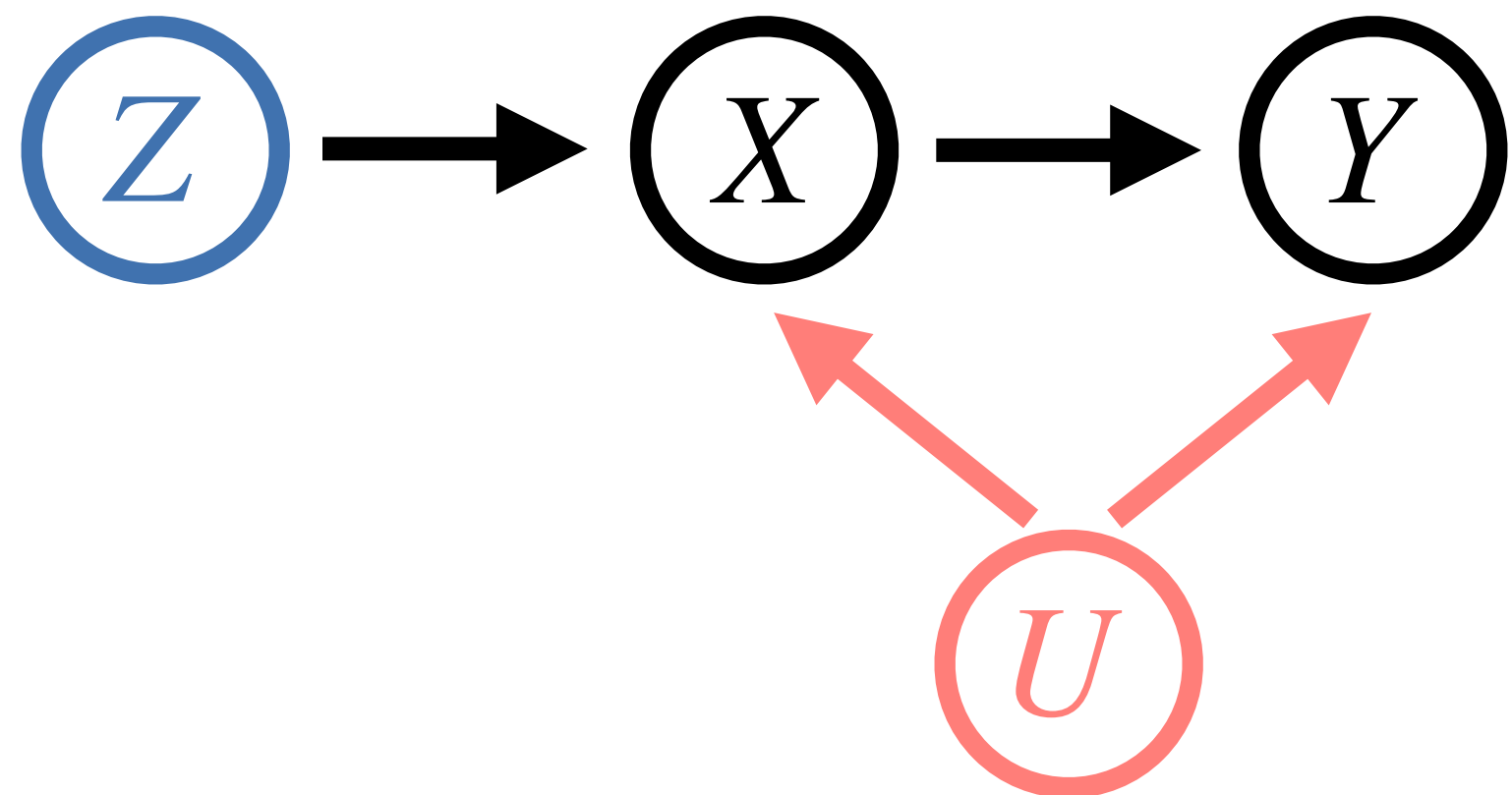


$$X = g(Z, U)$$

$$Y = h(X) + U$$

$$\mathbb{E}[U] = 0$$

$$0 = \mathbb{E}[U] = \mathbb{E}[U | z] = \mathbb{E}[Y - h(X) | z]$$



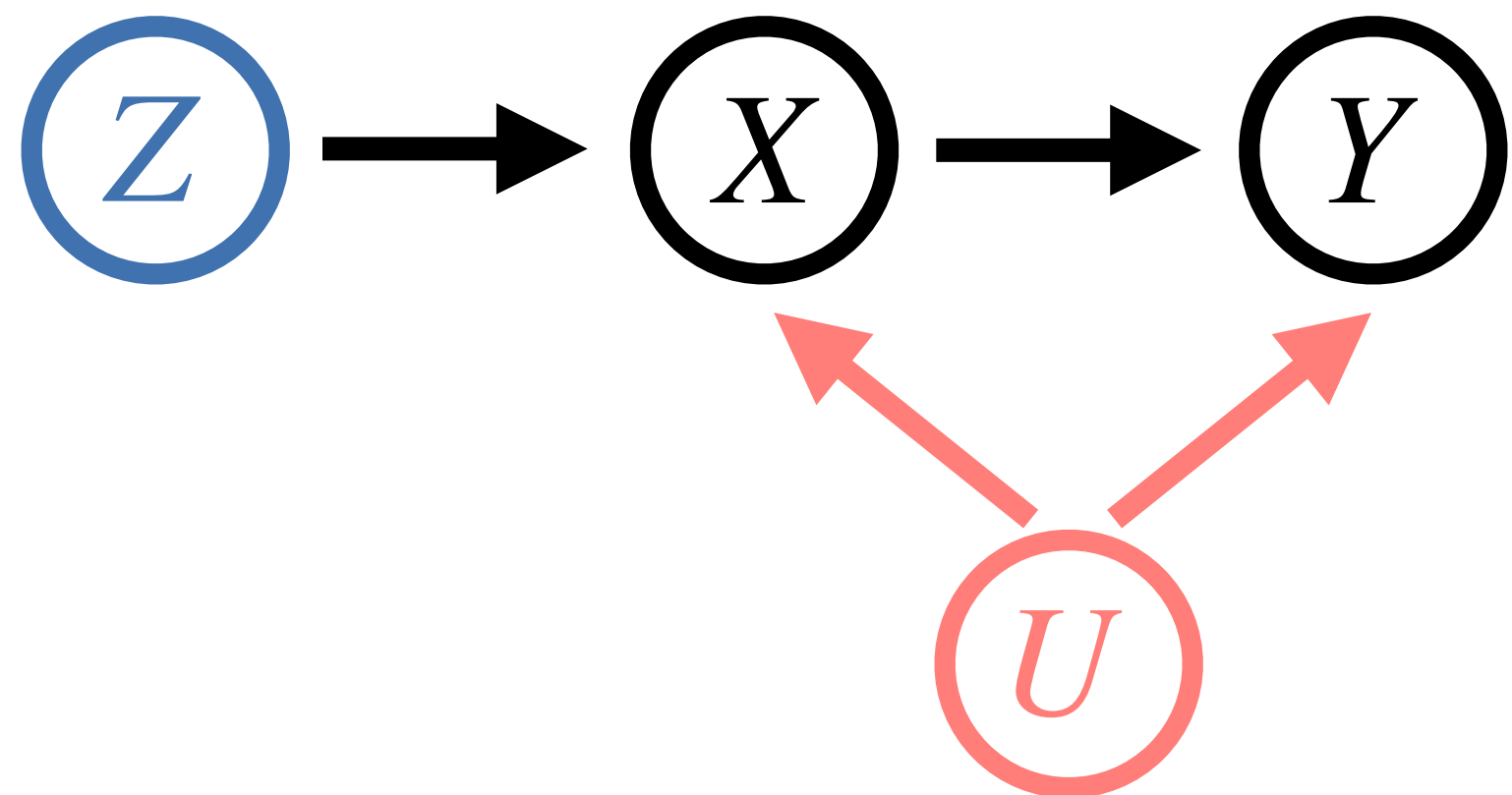
$$X = g(Z, U)$$

$$Y = h(X) + U$$

$$\mathbb{E}[U] = 0$$

$$0 = \mathbb{E}[U] = \mathbb{E}[U | z] = \mathbb{E}[Y - h(X) | z]$$

$$\Rightarrow \mathbb{E}[Y | z] = \mathbb{E}[h(X) | z], \forall z$$



$$X = g(Z, U)$$

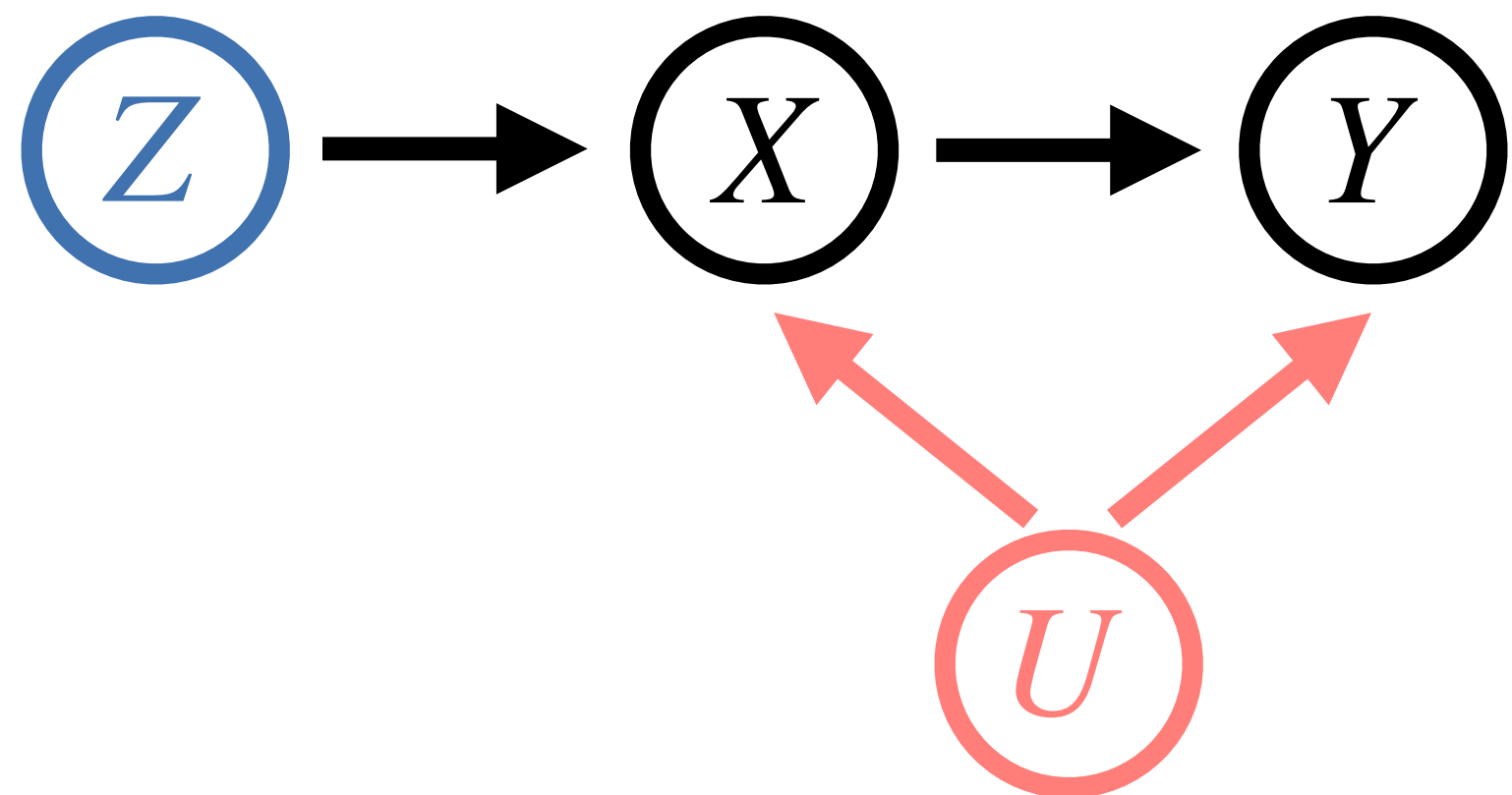
$$Y = h(X) + U$$

$$\mathbb{E}[U] = 0$$

$$0 = \mathbb{E}[U] = \mathbb{E}[U | z] = \mathbb{E}[Y - h(X) | z]$$

$$\Rightarrow \mathbb{E}[Y | z] = \mathbb{E}[h(X) | z], \forall z$$

$$\Rightarrow \min_h \mathbb{E}_z[(\mathbb{E}[Y | z] - \mathbb{E}[h(X) | z])^2]$$



$$X = g(Z, U)$$

$$Y = h(X) + U$$

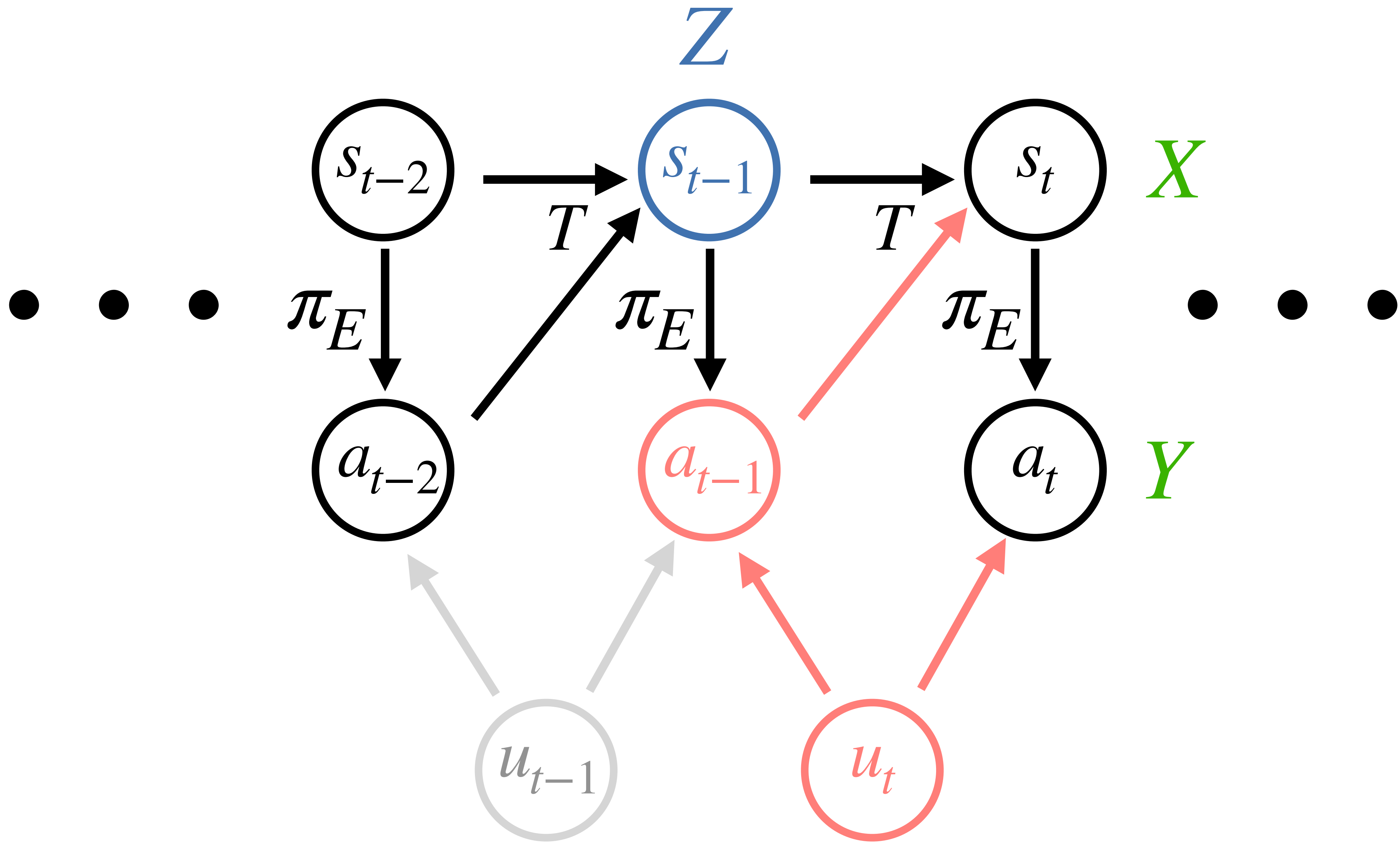
$$\mathbb{E}[U] = 0$$

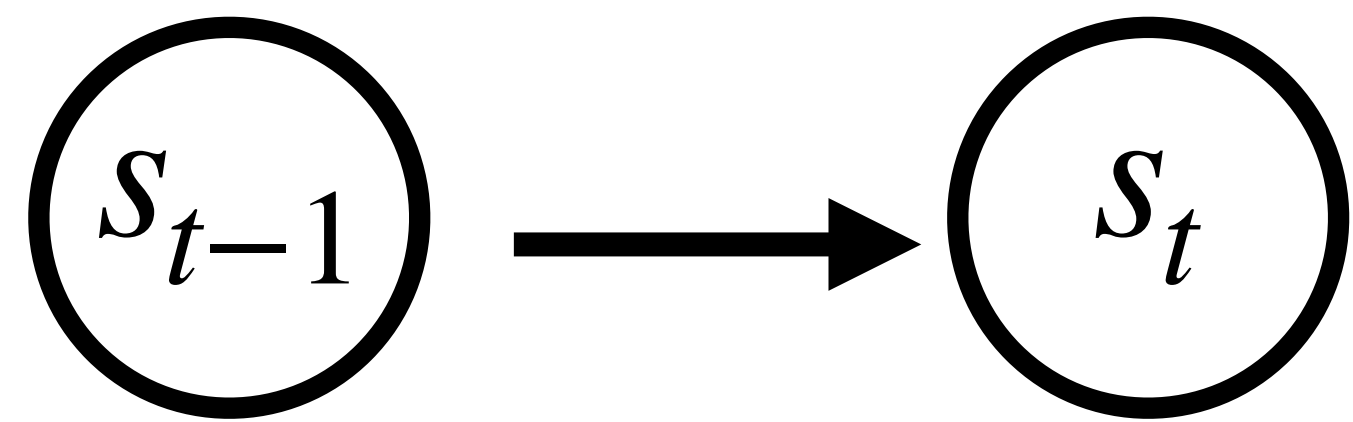
$$0 = \mathbb{E}[U] = \mathbb{E}[U | z] = \mathbb{E}[Y - h(X) | z]$$

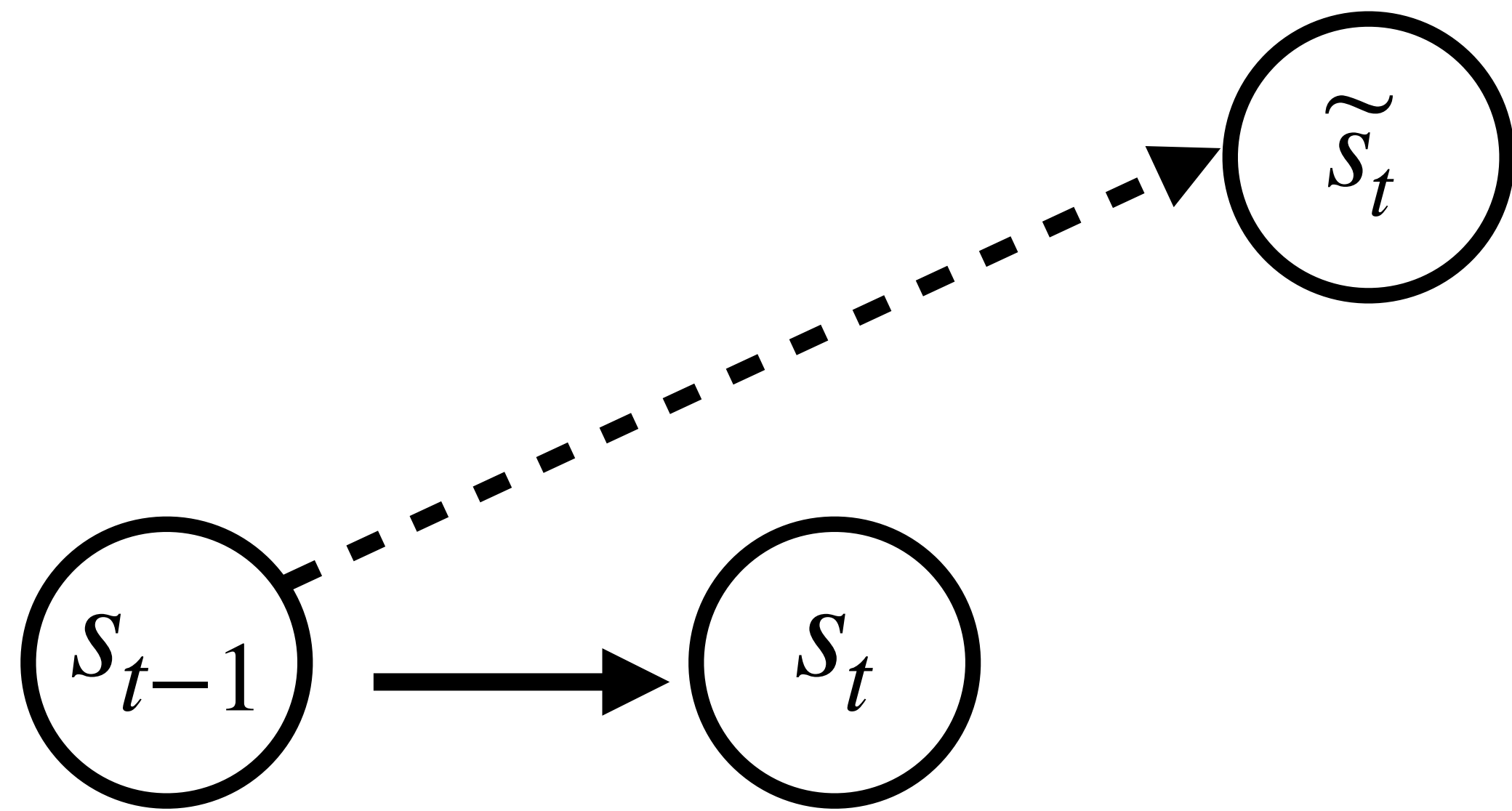
$$\Rightarrow \mathbb{E}[Y | z] = \mathbb{E}[h(X) | z], \forall z$$

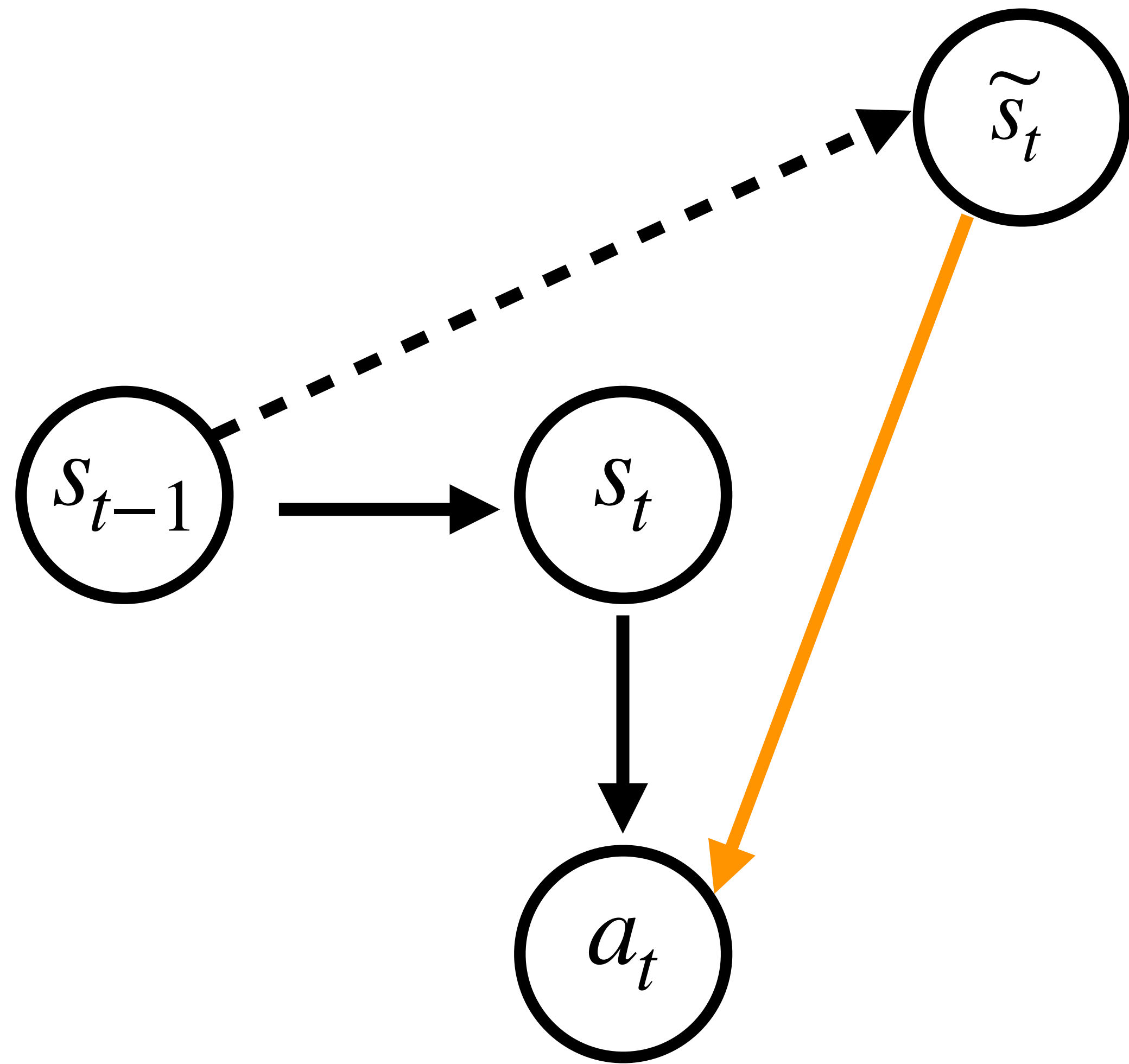
$$\Rightarrow \min_h \mathbb{E}_z[(\mathbb{E}[Y | z] - \mathbb{E}[h(X) | z])^2]$$

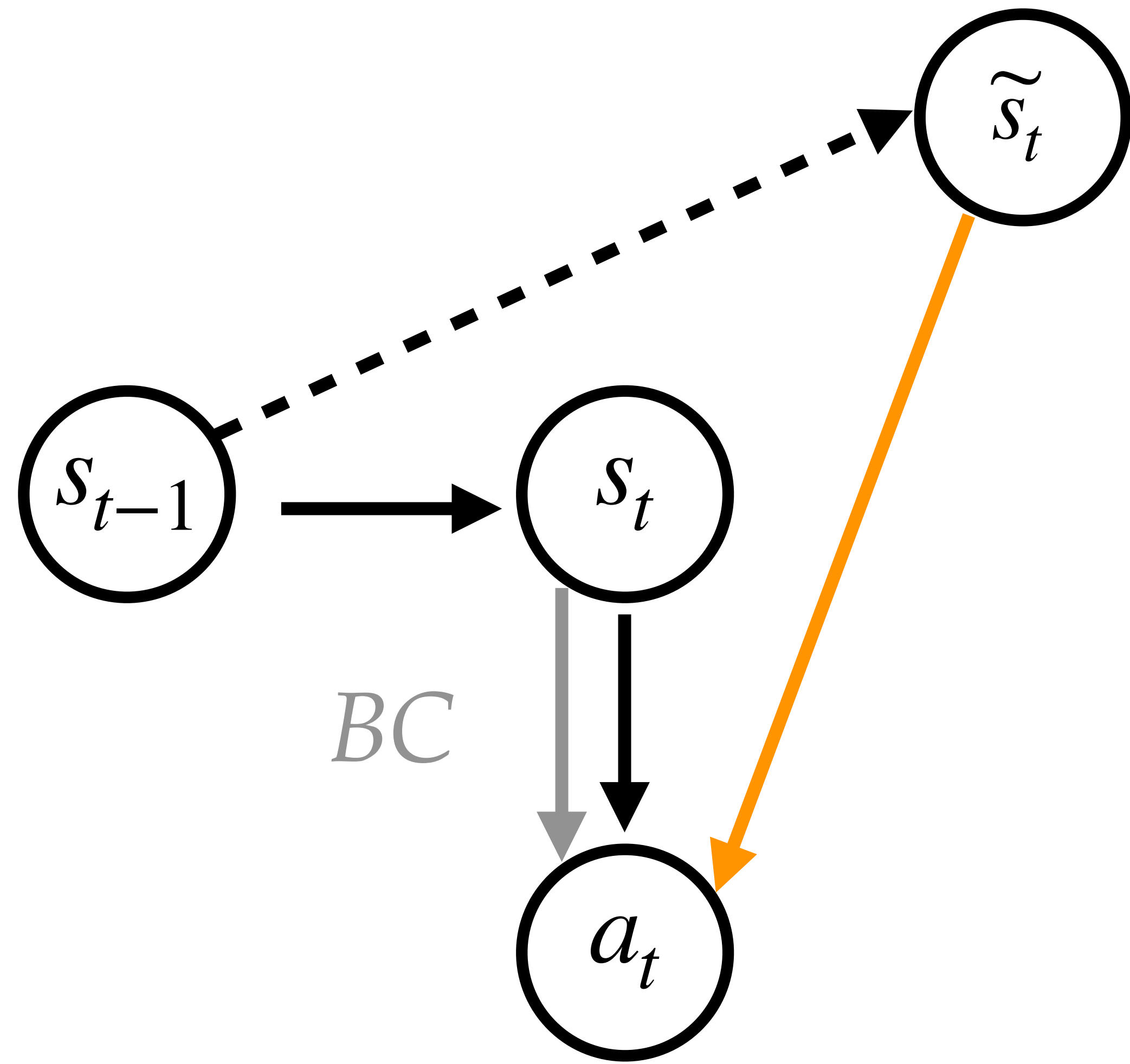
$$\Leftrightarrow \min_h \max_f \mathbb{E}_z[2(Y - h(X))f(Z) - f^2(Z)]$$

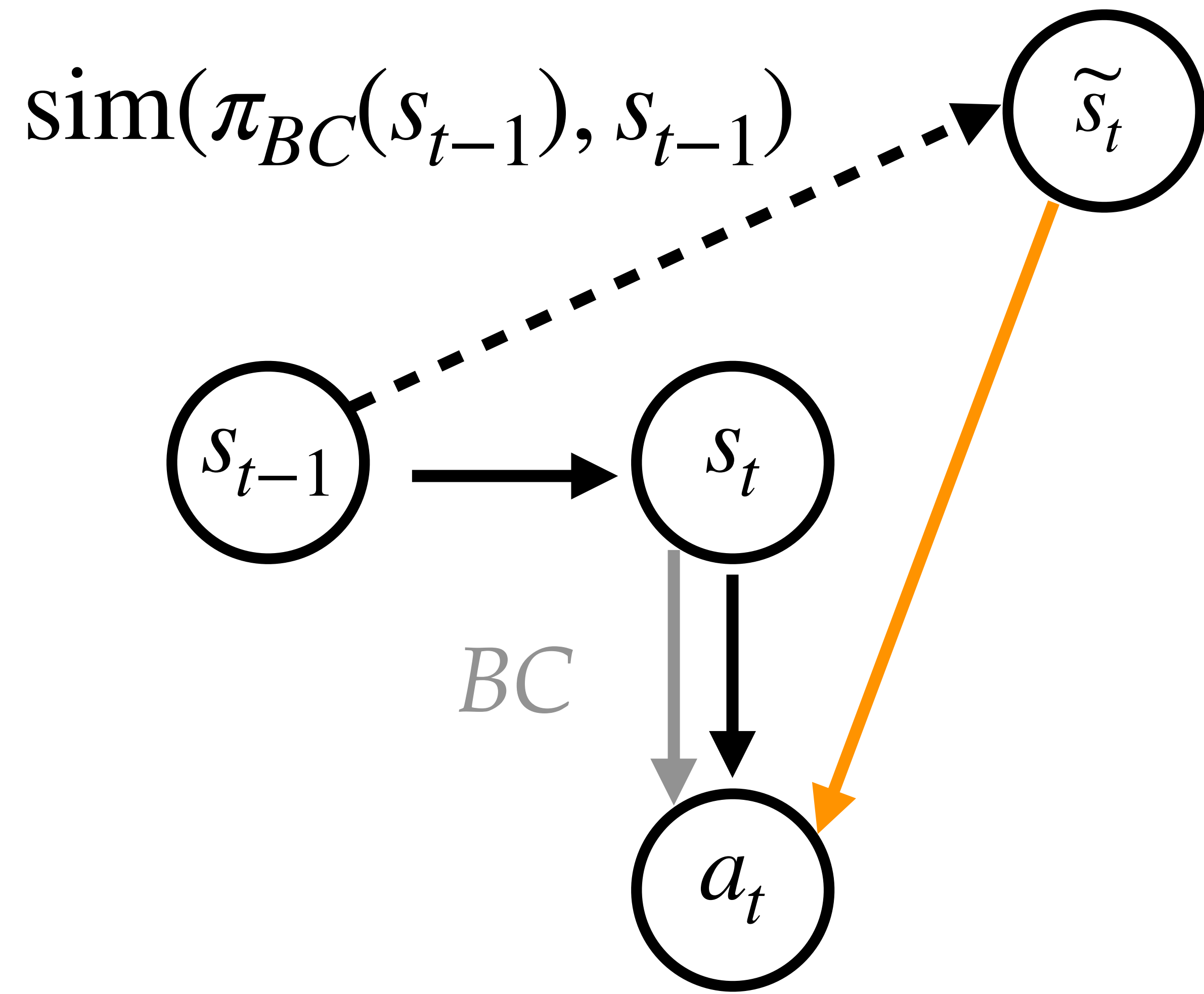


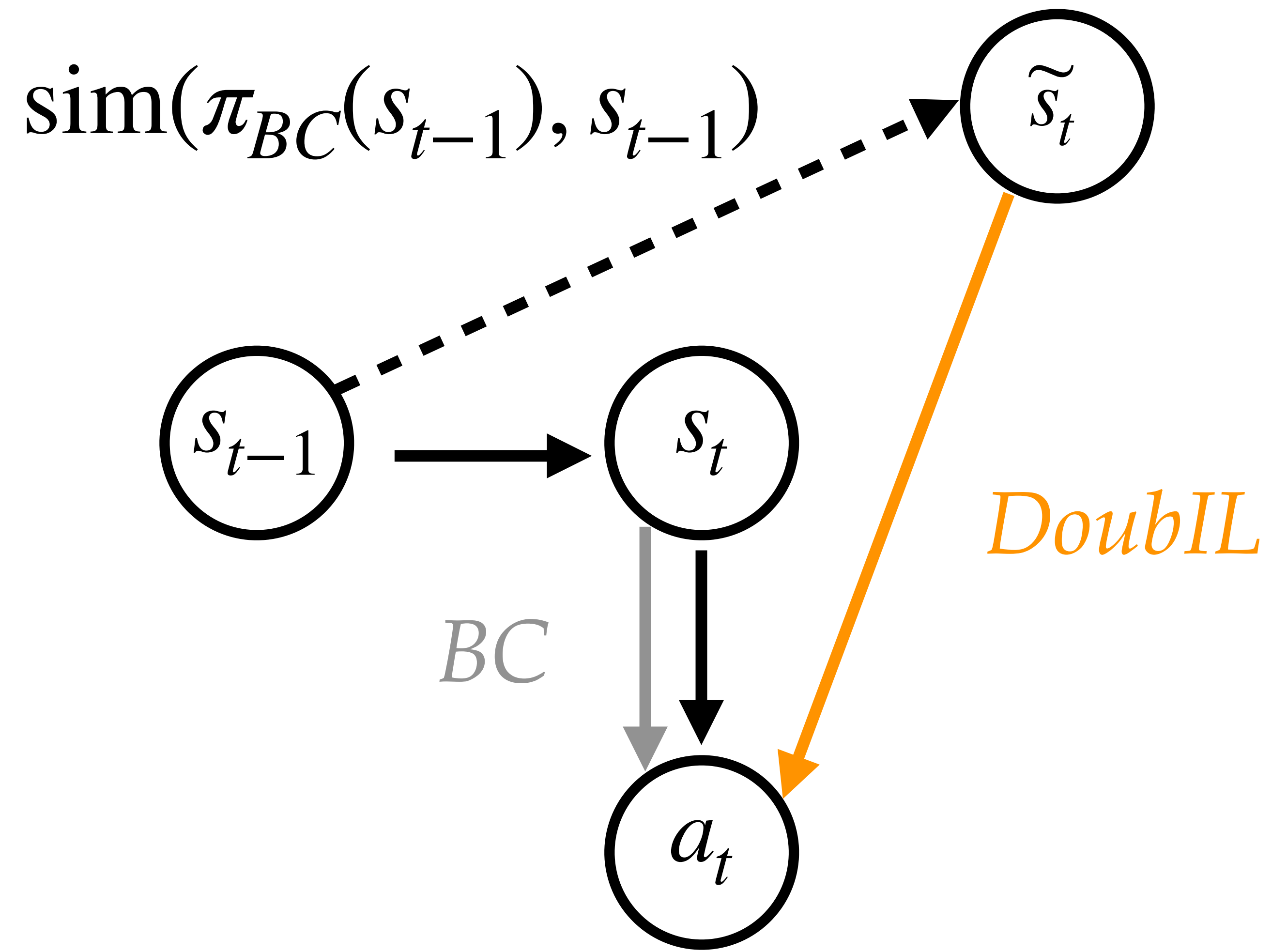


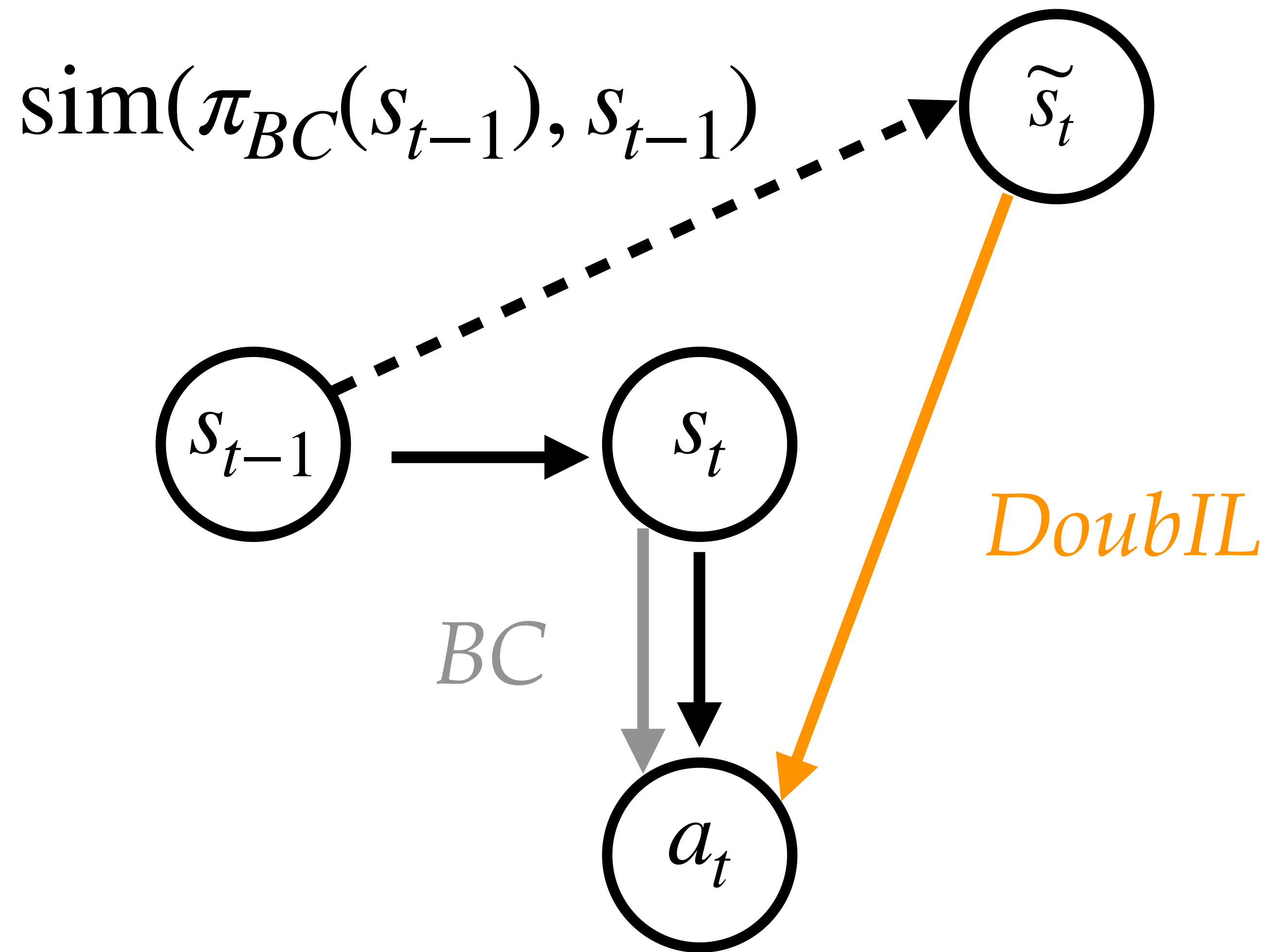












$$J(\pi_E) - J(\pi) \leq c(\sqrt{\epsilon} + \sqrt{\delta})\kappa(\Pi)T^2$$

$$\min_{\pi} \max_f \mathbb{E}[2(a_t - \pi(s_t))f(s_{t-1}) - f(s_{t-1})^2]$$

$$\min_{\pi} \max_f \mathbb{E}[2(a_t - \pi(s_t))f(s_{t-1}) - f(s_{t-1})^2]$$

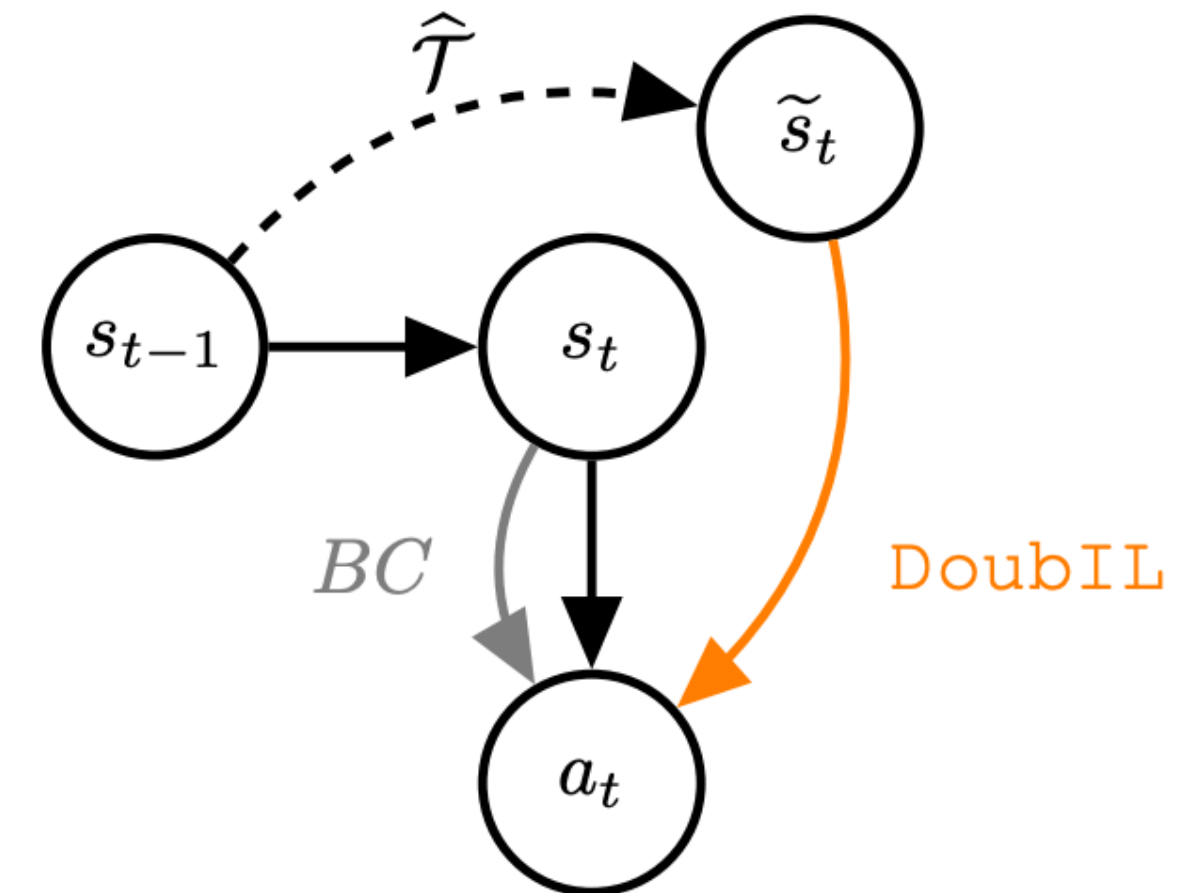
$$J(\pi_E) - J(\pi) \leq c\sqrt{\epsilon}\kappa(\Pi)T^2$$

Instrumental Variable Imitation Learning

generative modeling

game-theoretic

DoubIL



$$J(\pi_E) - J(\pi) \leq c(\sqrt{\epsilon} + \sqrt{\delta})\kappa(\Pi)T^2$$

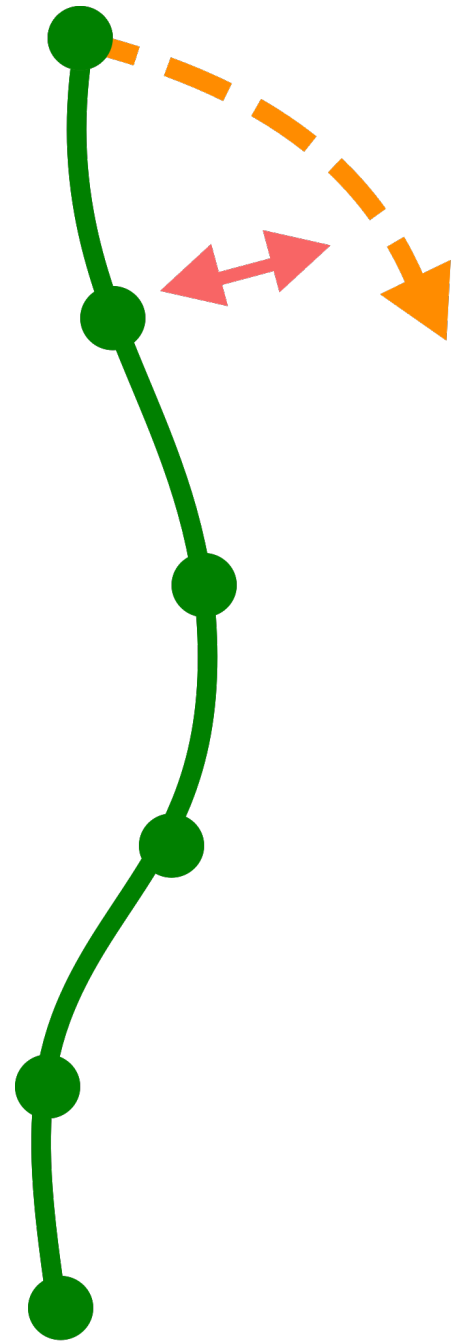
ResiduIL

$$\min_{\pi} \max_f \mathbb{E}[2(a_t - \pi(s_t))f(s_{t-1}) - f(s_{t-1})^2]$$

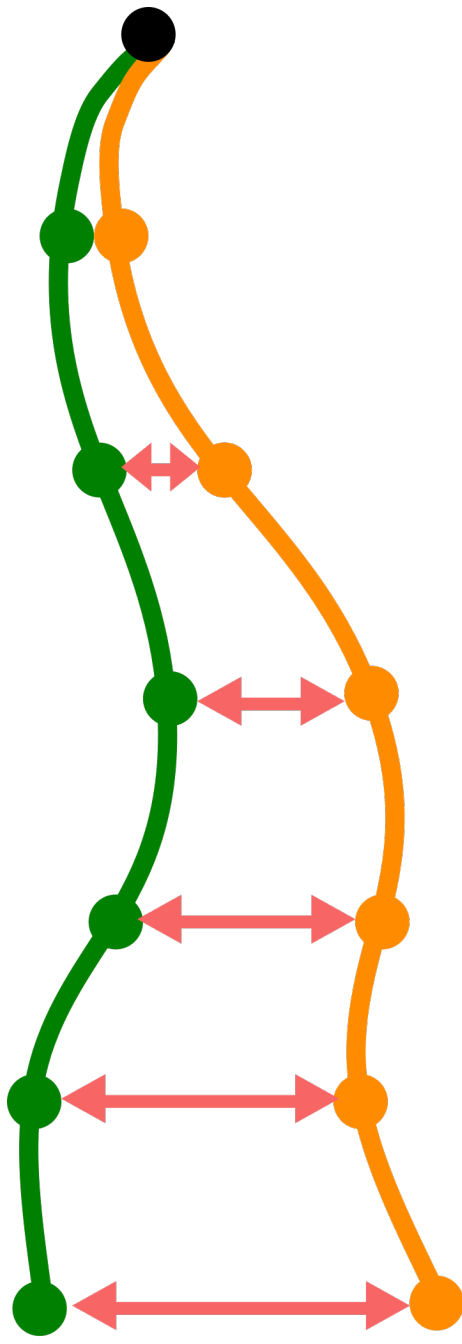
$$J(\pi_E) - J(\pi) \leq c\sqrt{\epsilon}\kappa(\Pi)T^2$$

$$\pi_E \xleftrightarrow{f} \pi$$

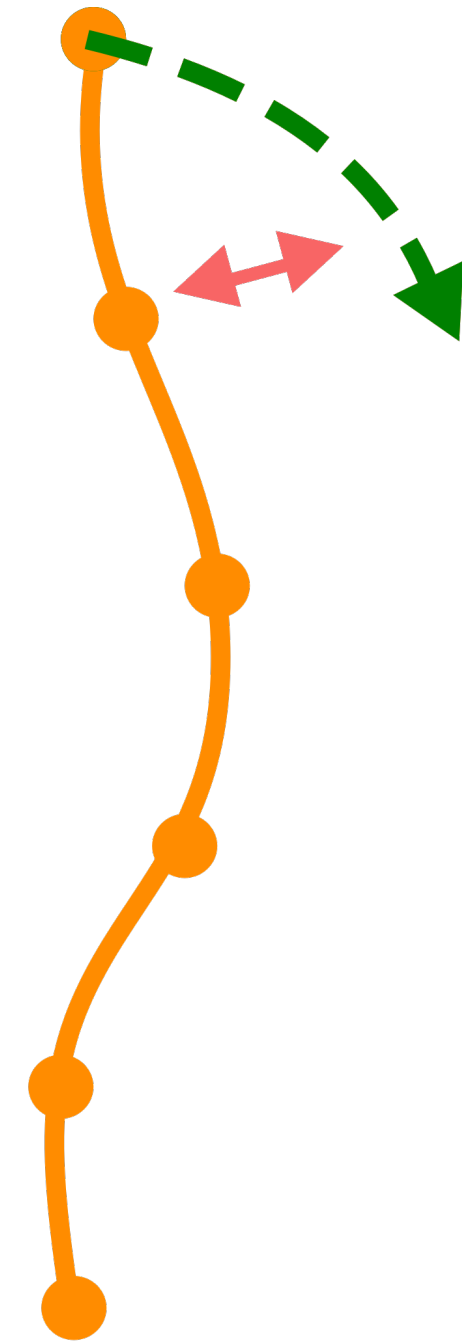
Offline



Online



Interactive

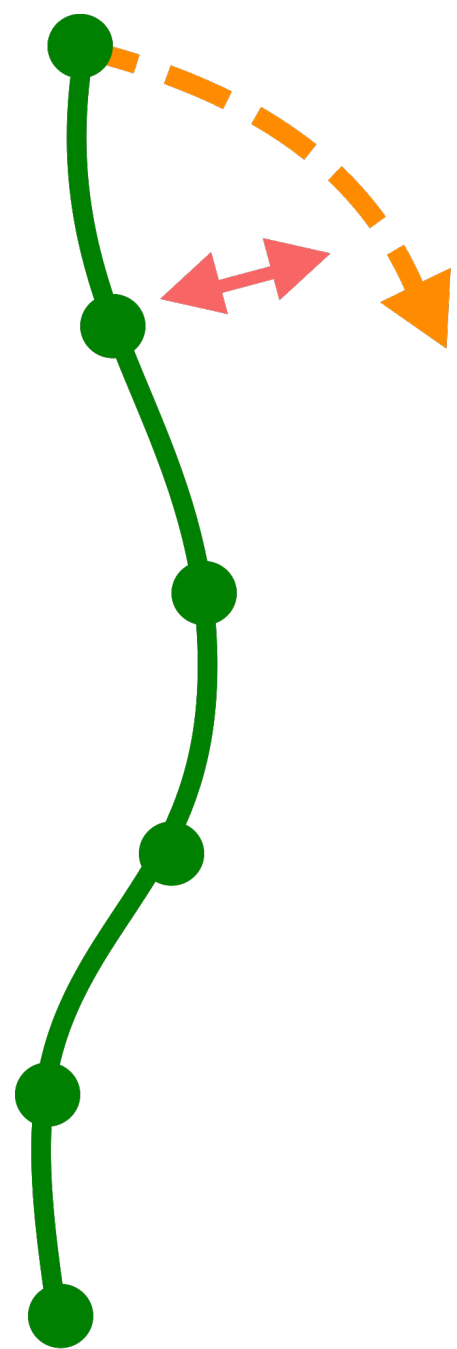


*Inconsistent,
IVR Consistent*

Consistent

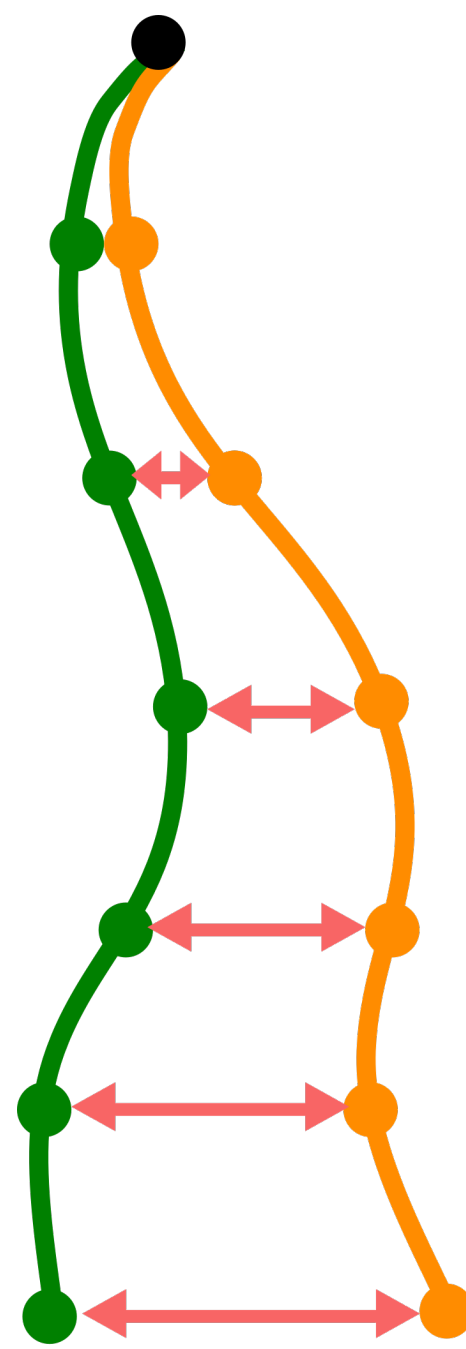
$$\pi_E \xleftrightarrow{f} \pi$$

Offline



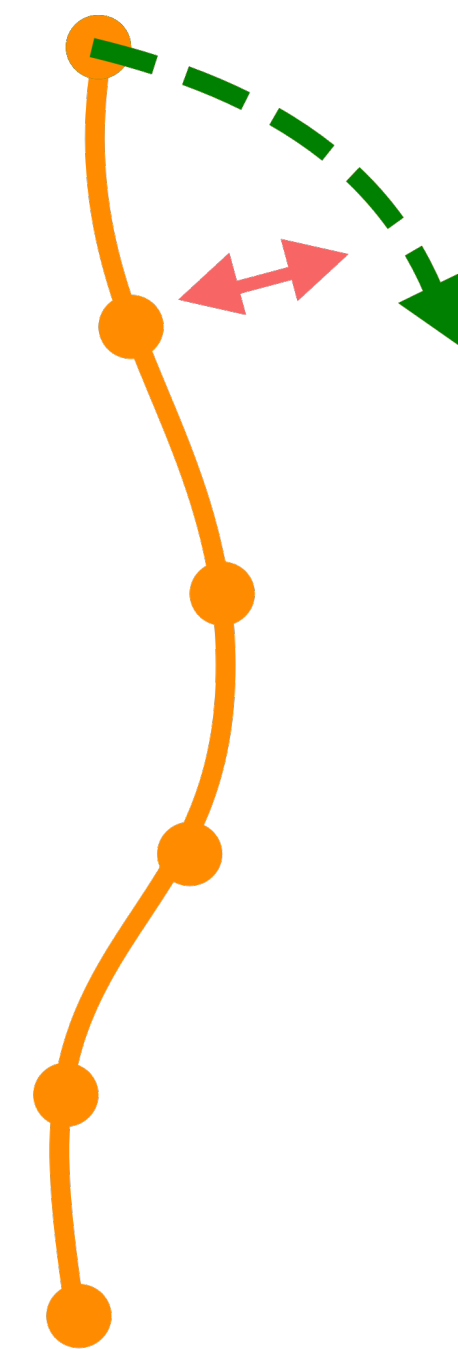
*Inconsistent,
IVR Consistent*

Online

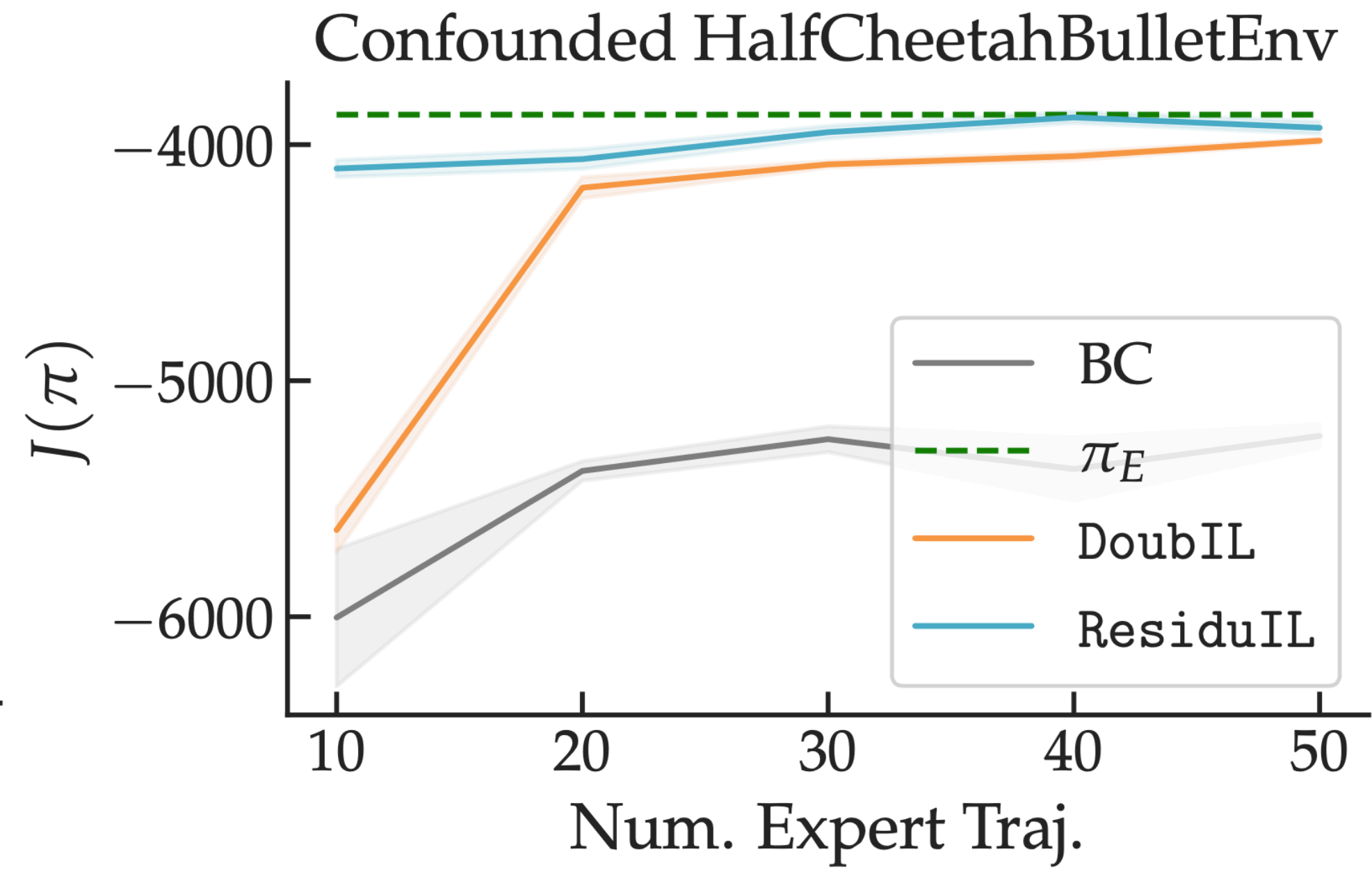
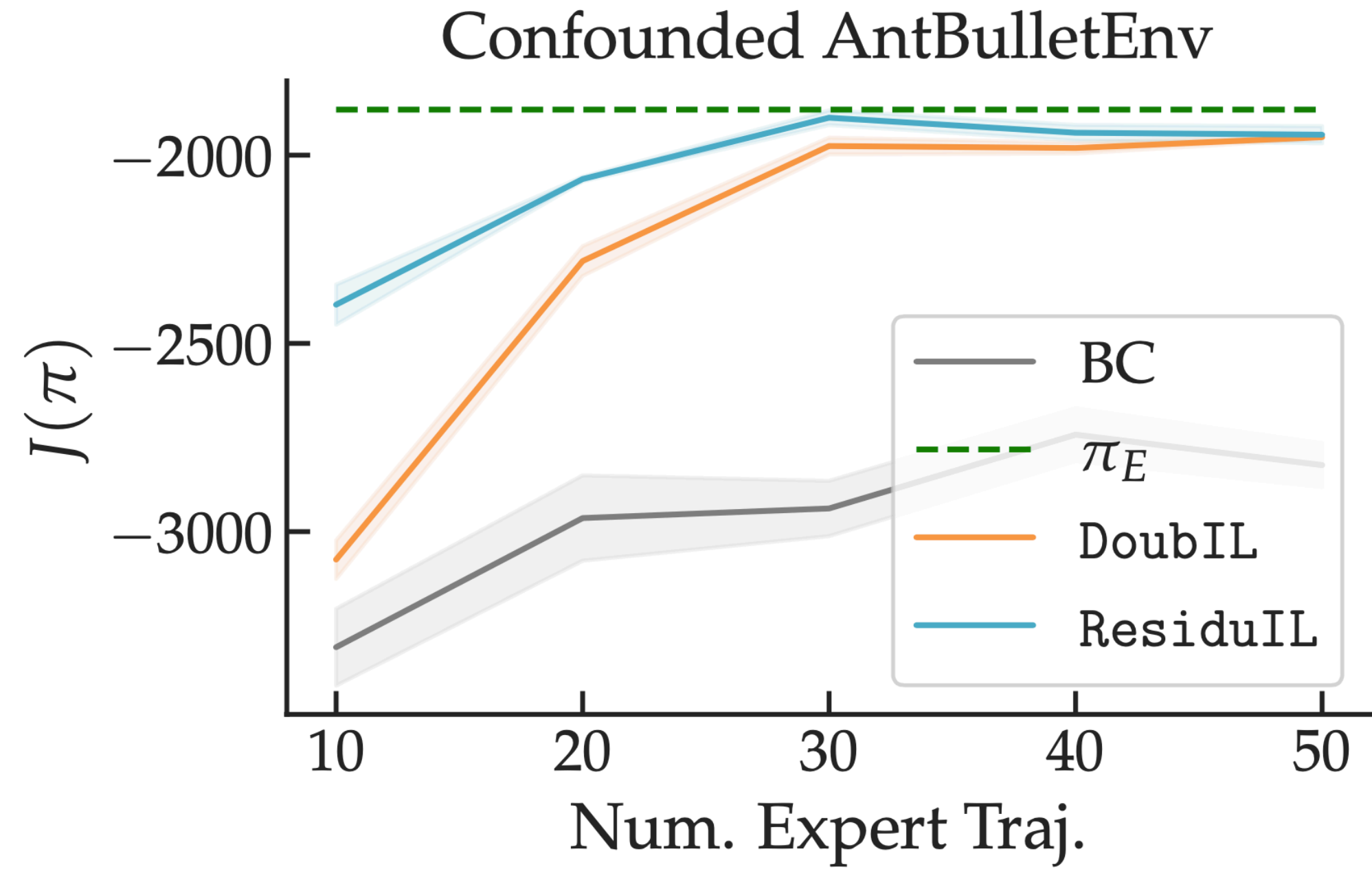
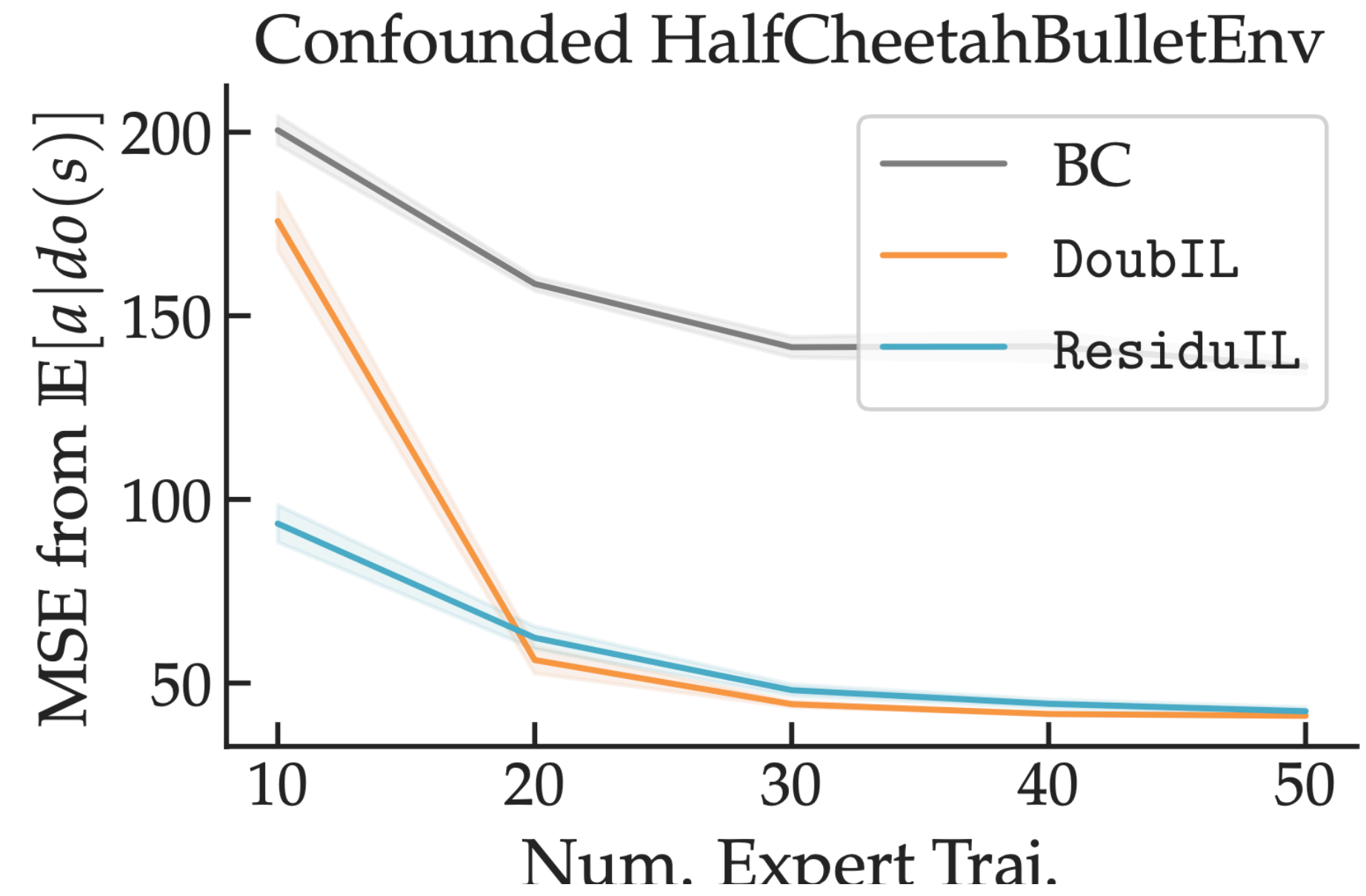
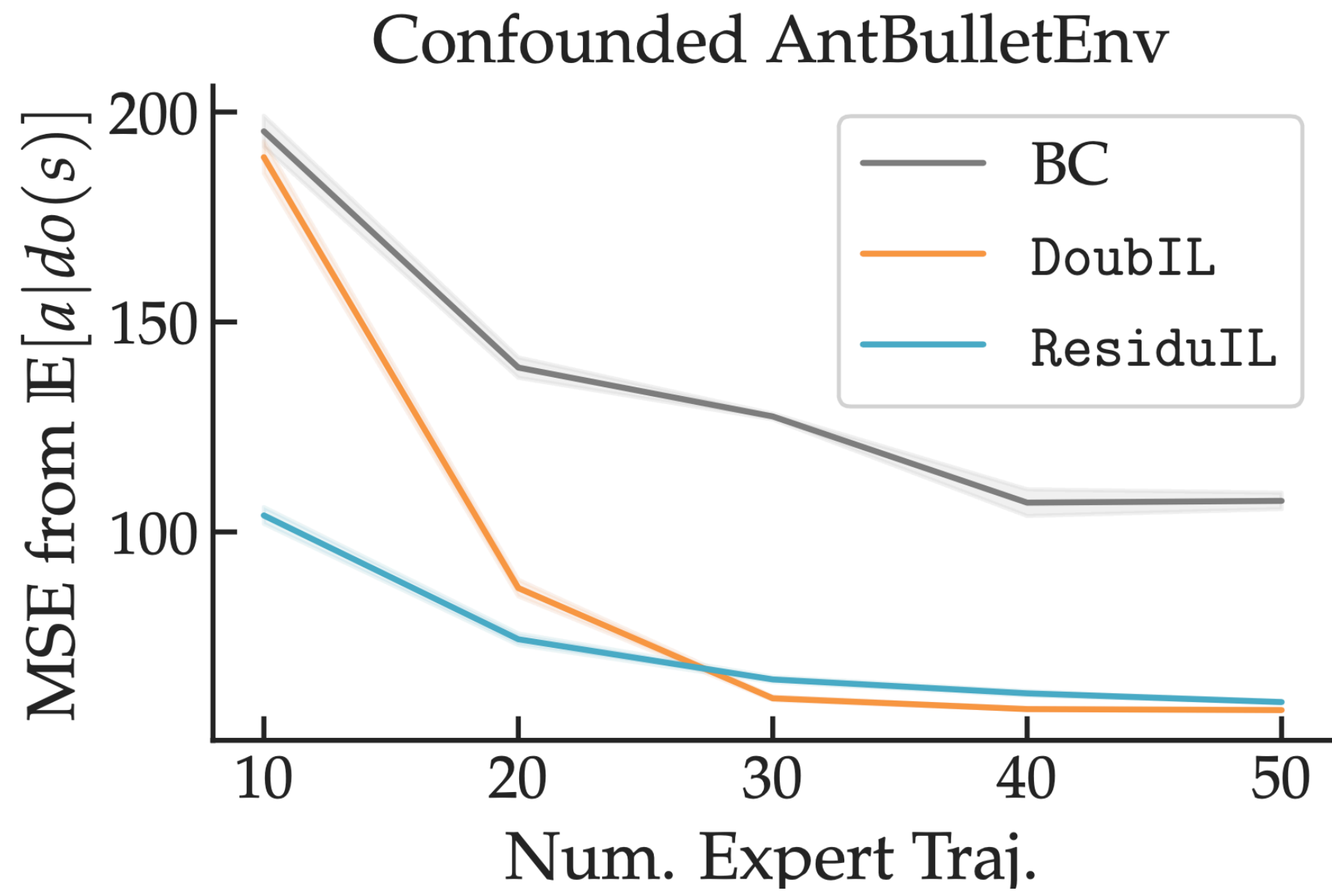


*Inconsistent,
Hybrid?*

Interactive

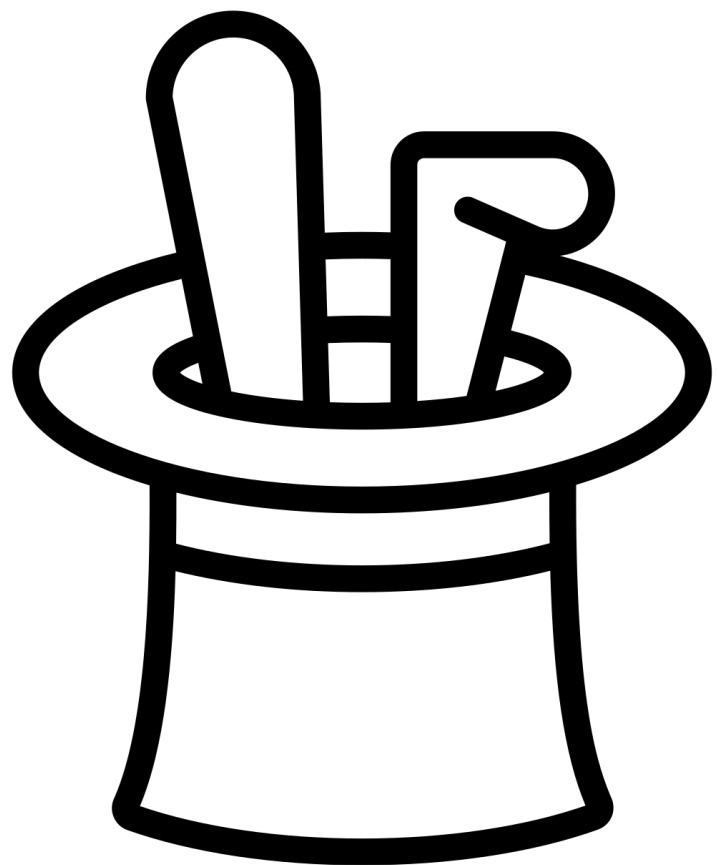


Consistent



	Offline	Online	Interactive
Covariate Shift	✗	✓	✓
Hidden Context	✗	✓ w/ History	✓ w/ History
TCN	✓ w/ IVR	✓ w/ IVR	✓

*Interventions happen via
interaction with
the environment in sequential
decision making.*



Thanks!

[https://goku1.dev/
gswamy@cmu.edu](https://goku1.dev/gswamy@cmu.edu)

