

Pretraining for Intelligent Agent

Offline Data and Simulation Perspective

Hojoon Lee, Byungkun Lee

KAIST AI

Short Bio



Hojoon Lee

- Ph.D student at KAIST AI, advised by Jaegul Choo.
- During M.S, his papers have been rejected 7 times in a row.
- Research Interest: Representation learning for decision making.



Byungkun Lee

- Ph.D student at KAIST AI, advised by Jaegul Choo.
- During M.S, his papers have been rejected 6 times in a row.
- Research Interest: I'm exploring my interest based on eps-greedy.

What is Intelligent Robot?



TSLA LIVE

What is Intelligence?

The ability to perceive, plan, act, and adapt.

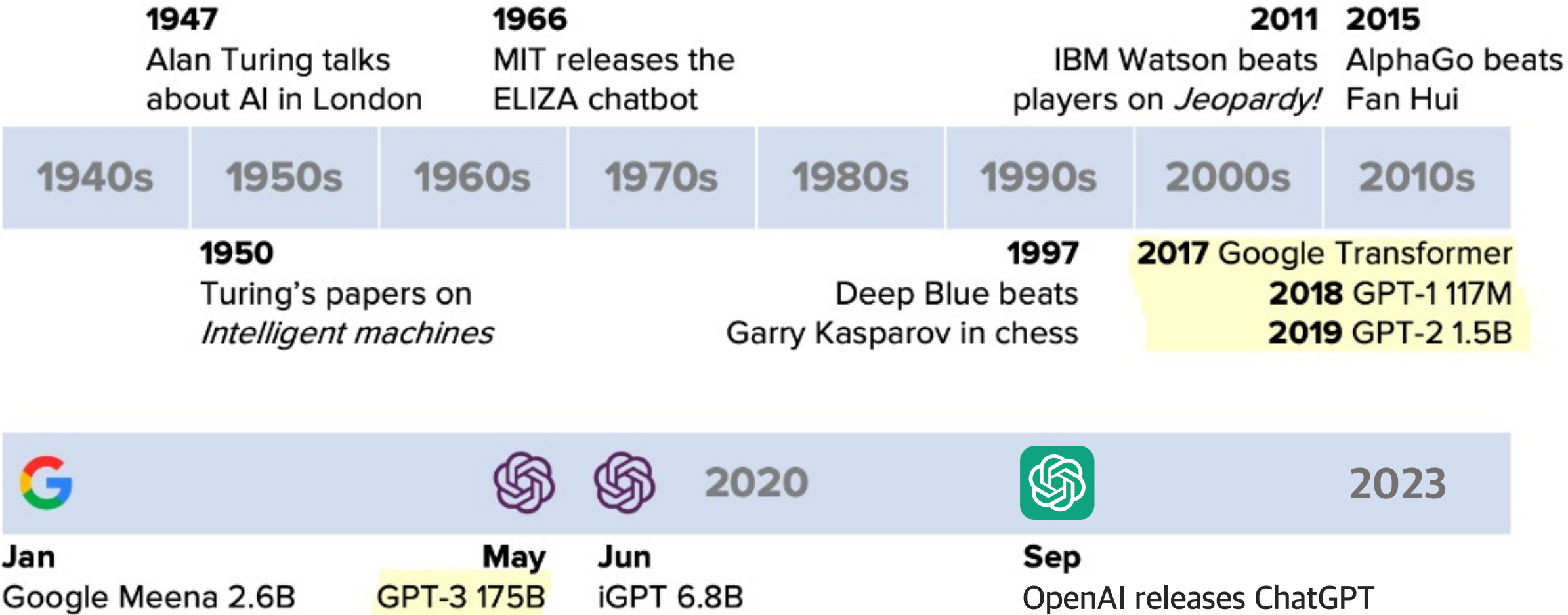
Patrick Winston

What is Artificial Intelligence?

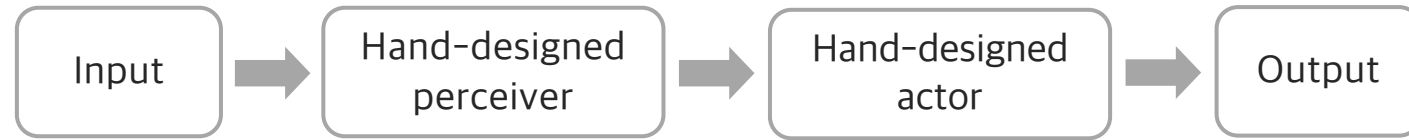
The study of computations that can learn to perceive, plan, act, and adapt.

Patrick Winston

History of Artificial Intelligence 1947-2023

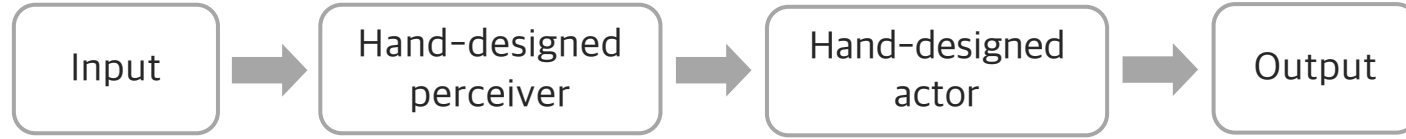


Rule-based system



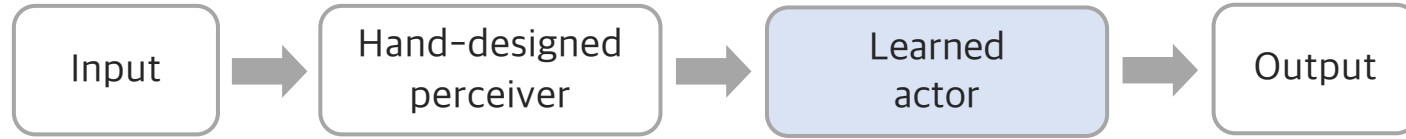
IBM DeepBlue

Rule-based system



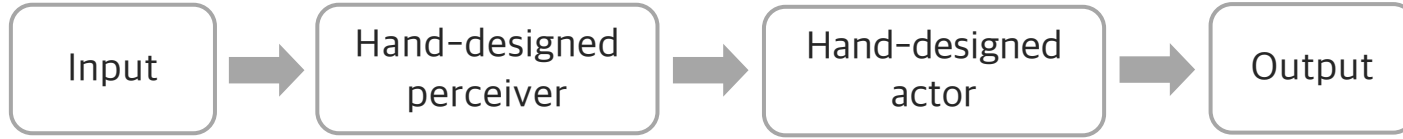
IBM DeepBlue

Classical machine-learning system (act)



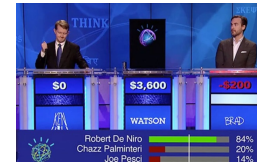
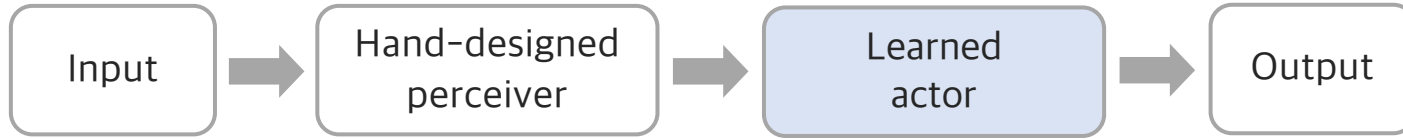
IBM Watson

Rule-based system



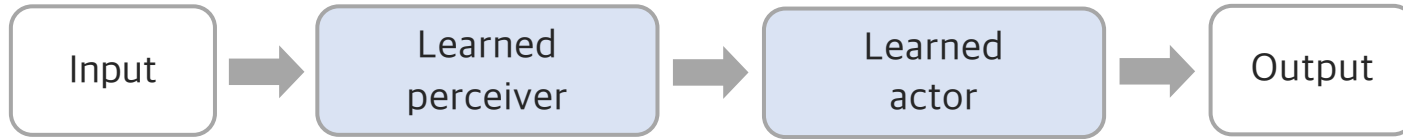
IBM DeepBlue

Classical machine-learning system (act)



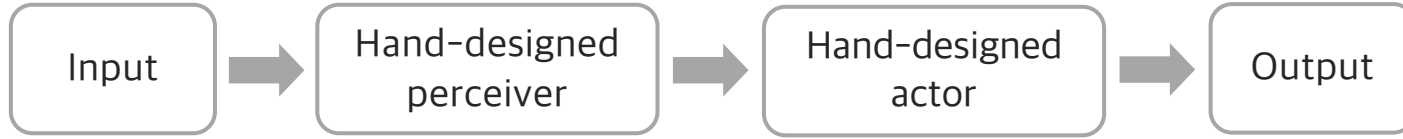
IBM Watson

Self-supervised learning system (perceive, act)



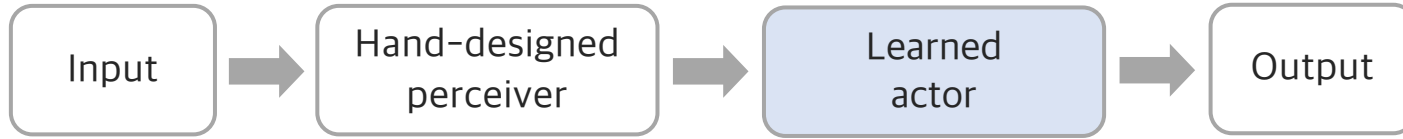
OpenAI ChatGPT

Rule-based system



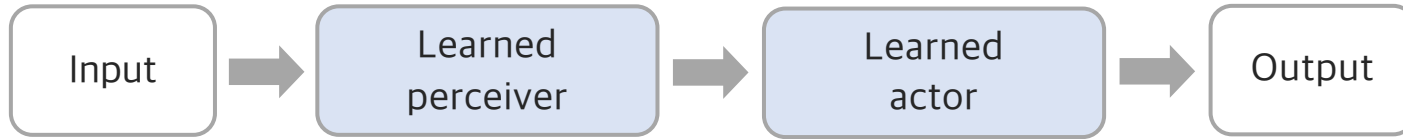
IBM DeepBlue

Classical machine-learning system (act)



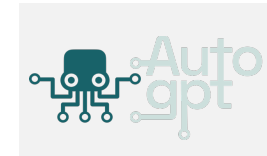
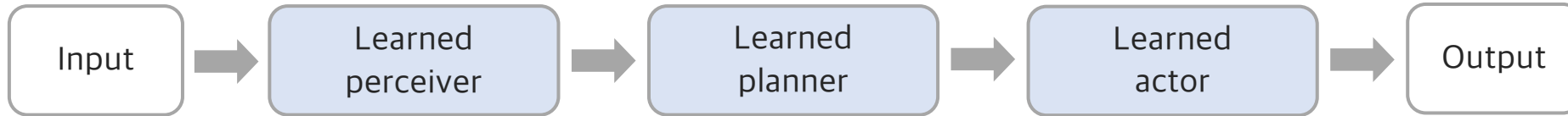
IBM Watson

Self-supervised learning system (perceive, act)



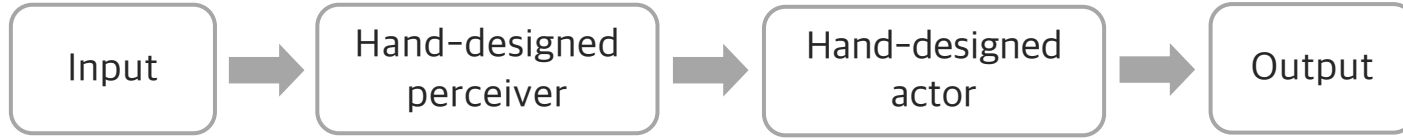
OpenAI ChatGPT

Self-supervised planning system (perceive, plan, act)



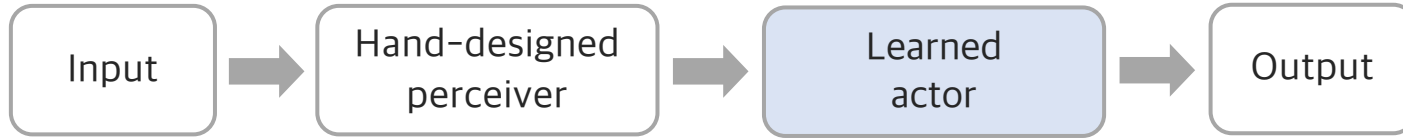
AutoGPT

Rule-based system



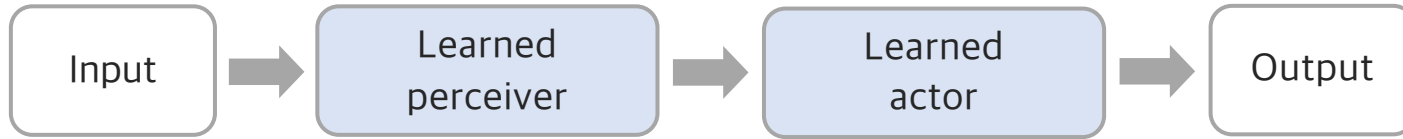
IBM DeepBlue

Classical machine-learning system (act)



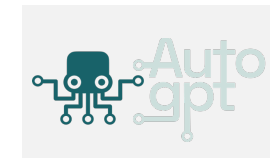
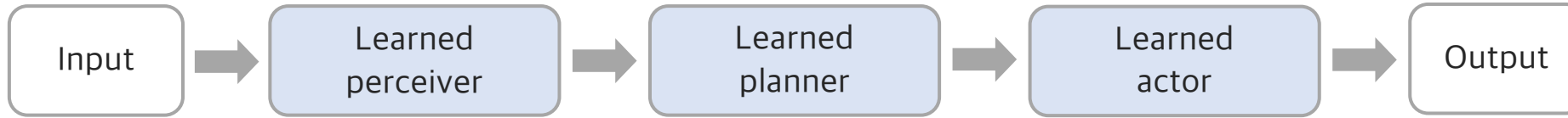
IBM Watson

Self-supervised learning system (perceive, act)



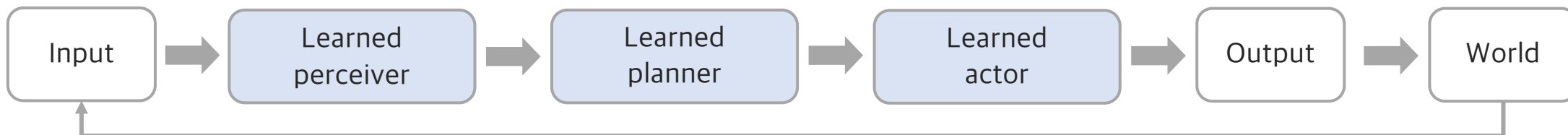
OpenAI ChatGPT

Self-supervised planning system (perceive, plan, act)



AutoGPT

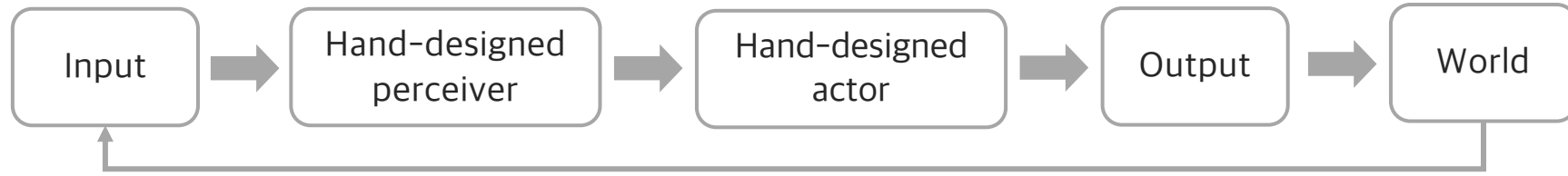
Life-long learning system (perceive, plan, act, adapt)



None

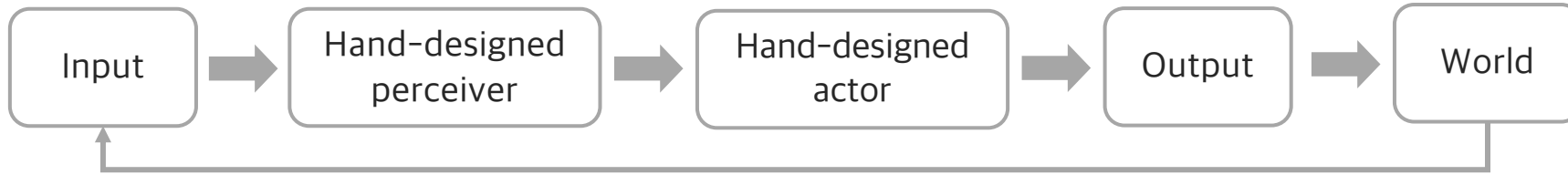
Where are we now for Intelligent Robot?

Rule-based system



B-Dynamics Spot

Rule-based system



B-Dynamics Spot

Somewhat mixed system

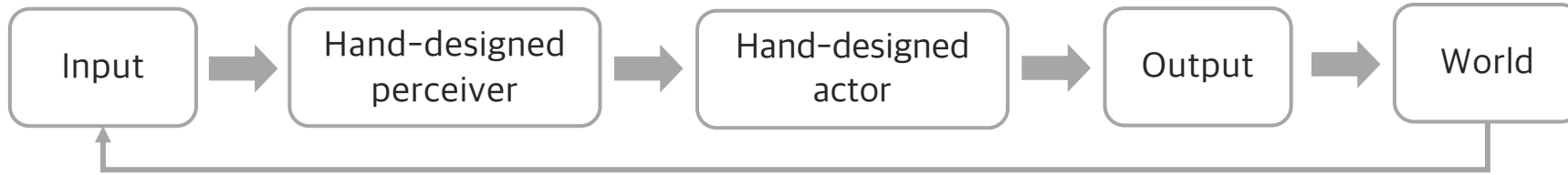


Tesla Bot



Tesla AutoPilot

Rule-based system



B-Dynamics Spot

Somewhat mixed system

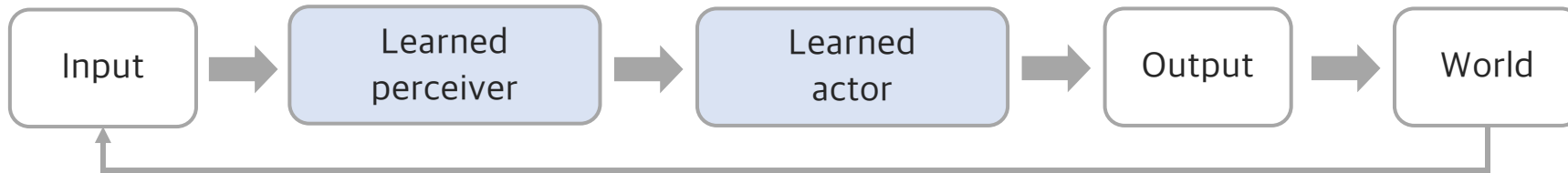


Tesla Bot



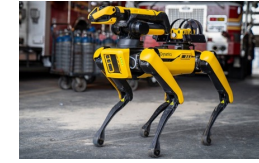
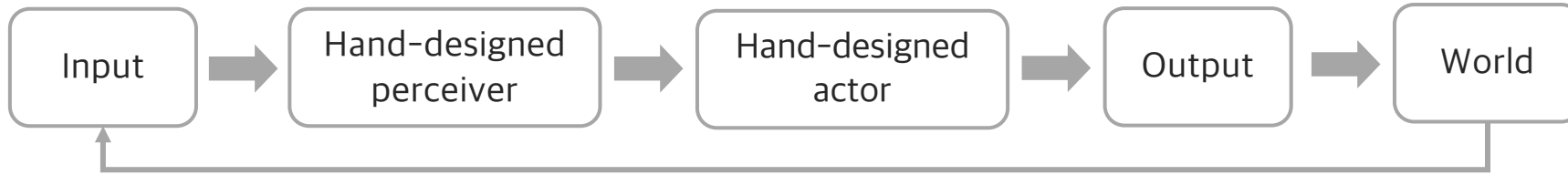
Tesla AutoPilot

Self-supervised learning system (perceive, act)



None

Rule-based system



B-Dynamics Spot

Somewhat mixed system

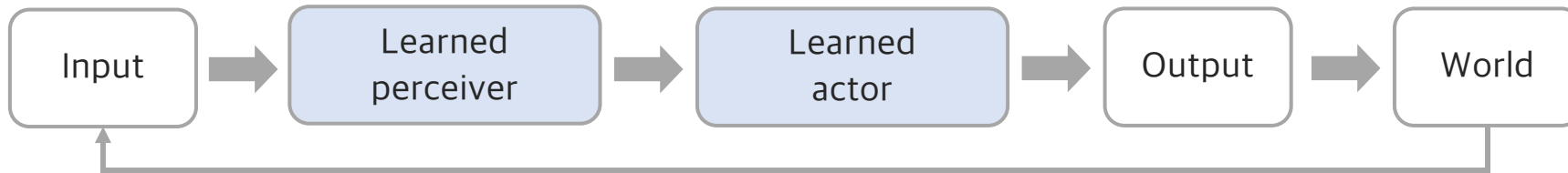


Tesla Bot



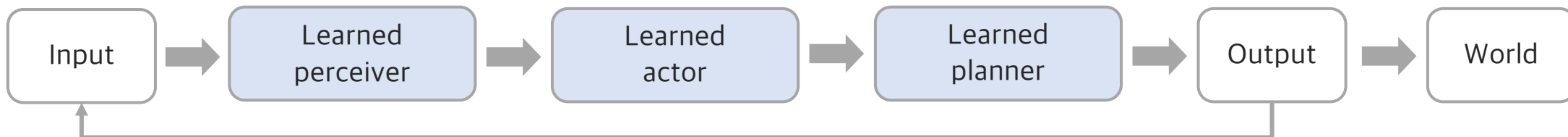
Tesla AutoPilot

Self-supervised learning system (perceive, act)



None

Self-supervised planning system (perceive, act, plan)



None

How can we make Intelligent Robot?

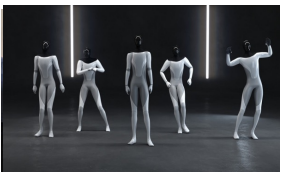
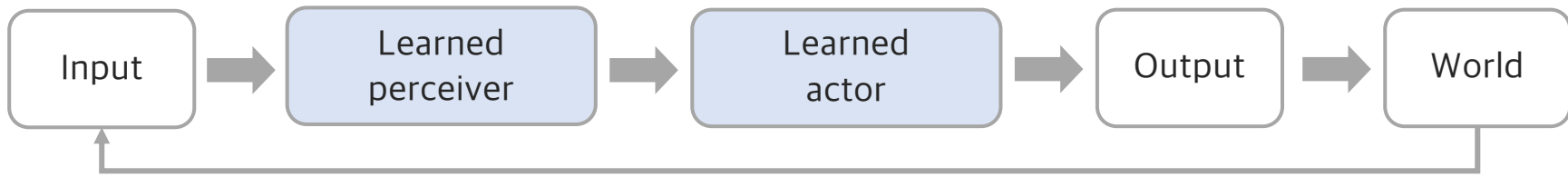


Tesla Bot



Tesla AutoPilot

How?



Tesla SuperBot



Tesla SuperPilot

Source of Training Robot

- Online Dataset
 - (+) explore the uncertain region.
 - (-) time-consuming, **expensive to collect.**

Source of Training Robot

- Online Dataset
 - (+) explore the uncertain region.
 - (-) time-consuming, **expensive to collect.**
- Offline Dataset
 - (+) cheaper than online dataset.
 - (-) distribution shift, cannot explore the uncertain region.



RT-X
1M robot trajectory



EGO-4D
ego-centric video

Source of Training Robot

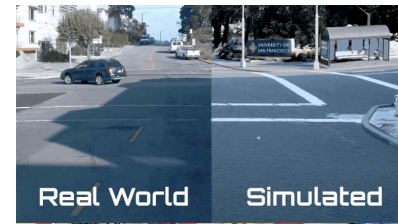
- Online Dataset
 - (+) explore the uncertain region.
 - (-) time-consuming, **expensive to collect.**
- Offline Dataset
 - (+) cheaper than online dataset.
 - (-) distribution shift, cannot explore the uncertain region.
- Simulator
 - (+) cheaper than online & offline dataset.
 - (-) larger distribution shift, hard to construct.



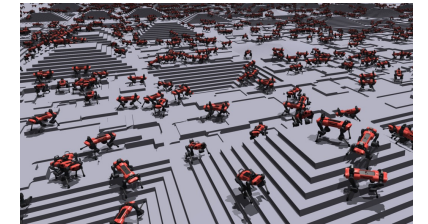
RT-X
1M robot trajectory



EGO-4D
ego-centric video

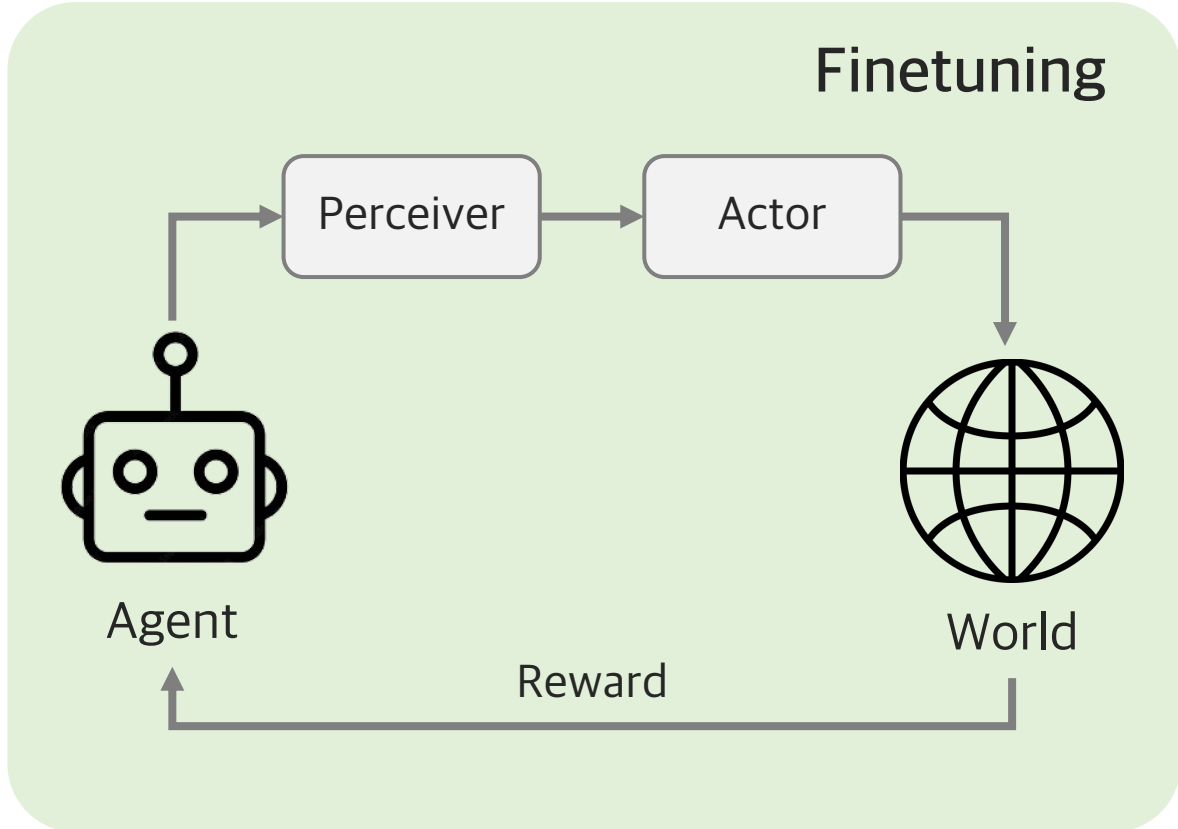
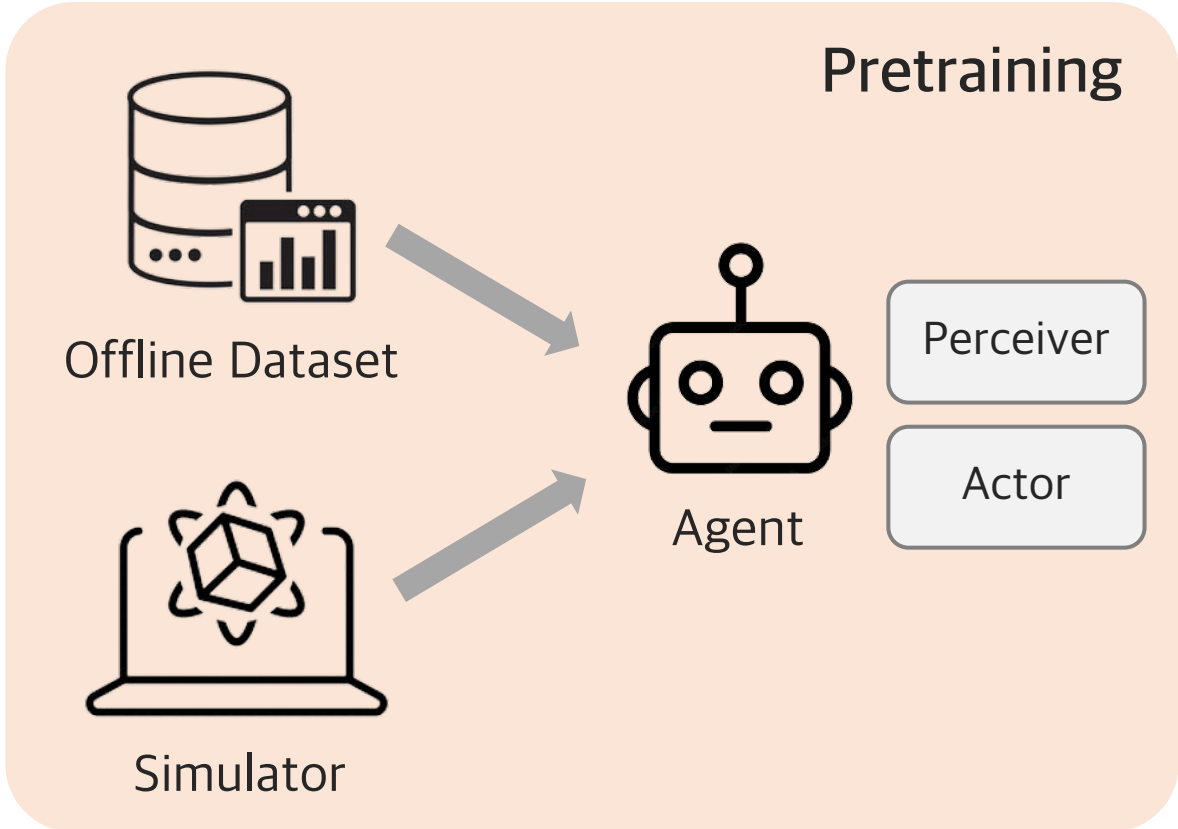


Simulation City
Waymo

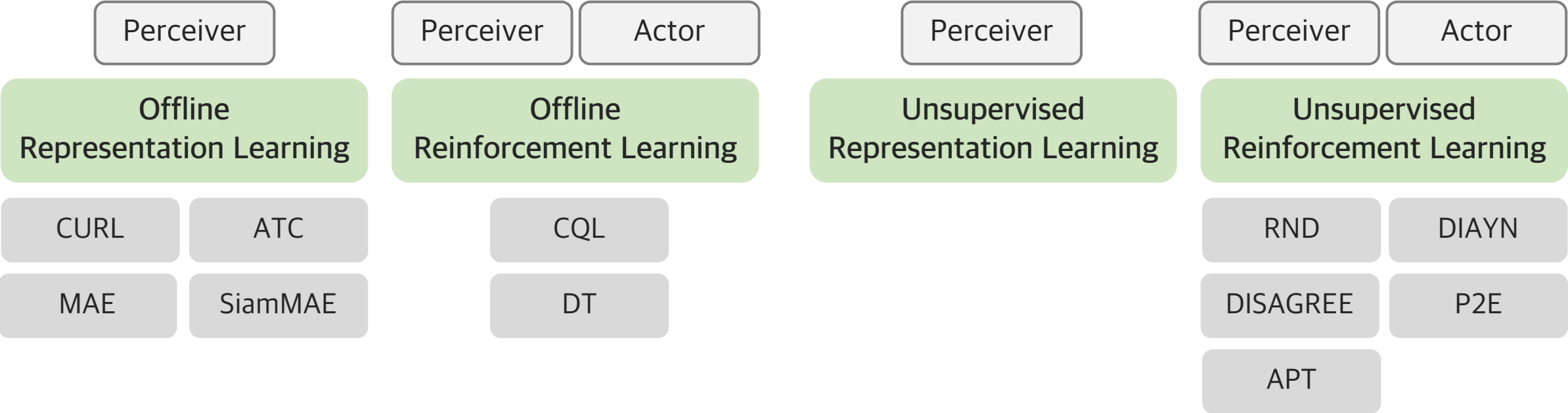


Issac Gym
NVIDIA

Pretrain-then-Finetune paradigm



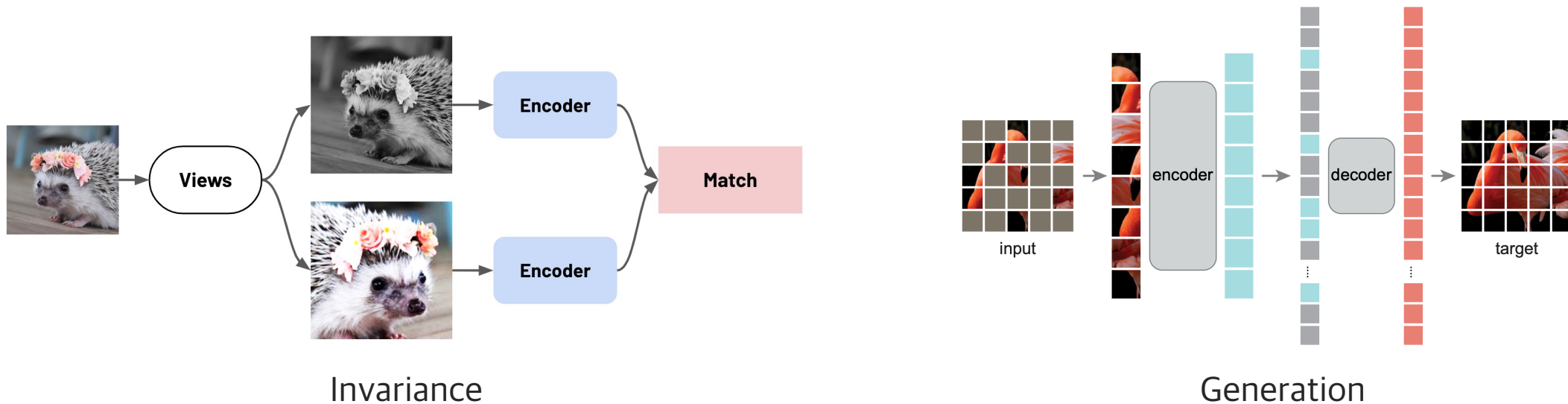
Taxonomy of Pretraining Algorithms



Pretraining from Offline Dataset

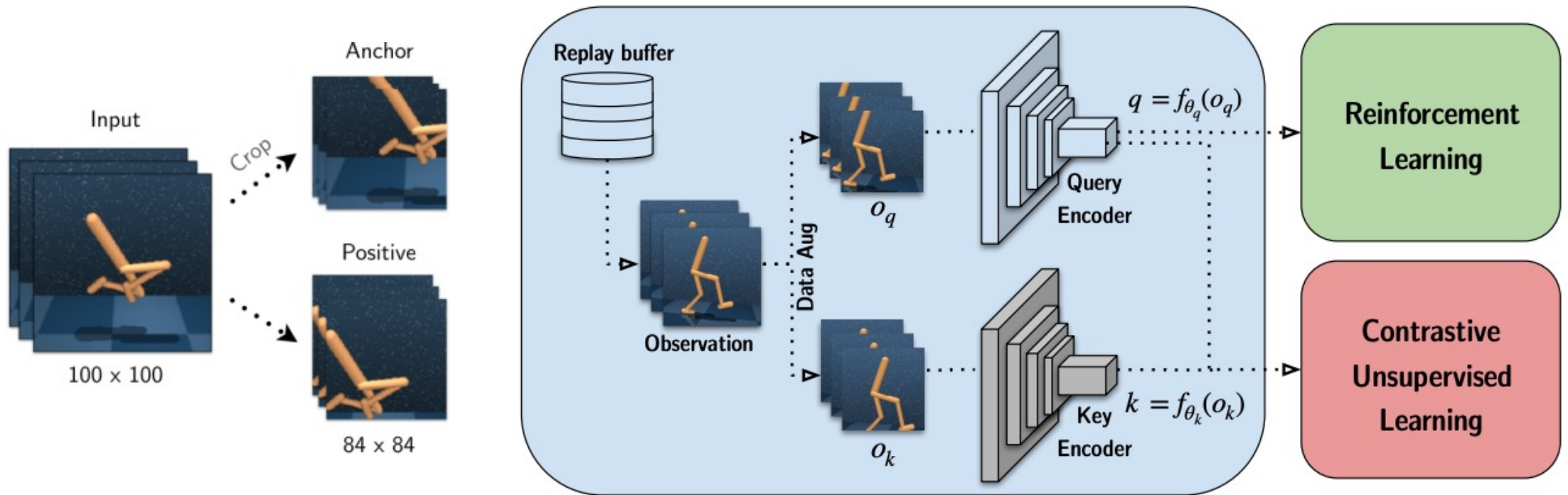
Offline Representation Learning

- Train effective **Perceiver** from offline dataset.
- Invariance vs Generation
 - Let X be data, $Z(X)$ be representation.
 - Let I denotes mutual information of two random variables.
 - **Invariance**: Maximize $I(Z(X_1); Z(X_2))$ where X_1, X_2 are invariant data.
 - **Generation**: Maximize $I(Z(\bar{X}); X)$ where \bar{X} is perturbed image of X .



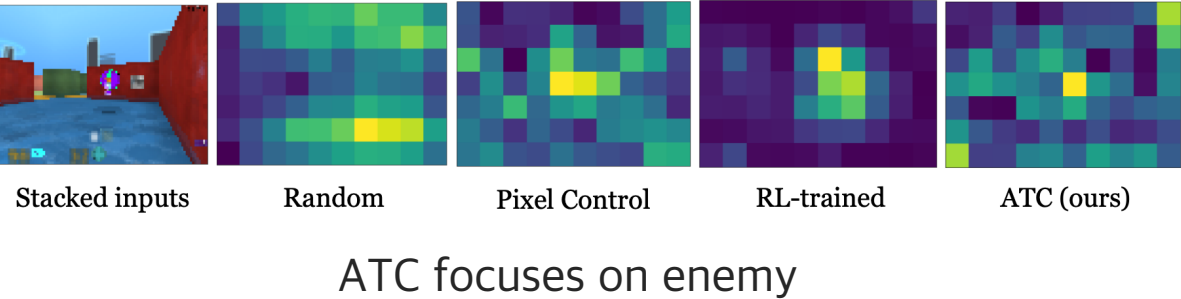
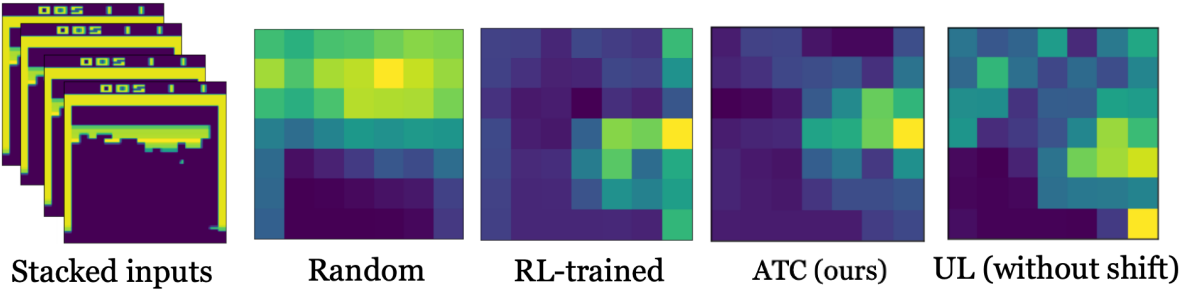
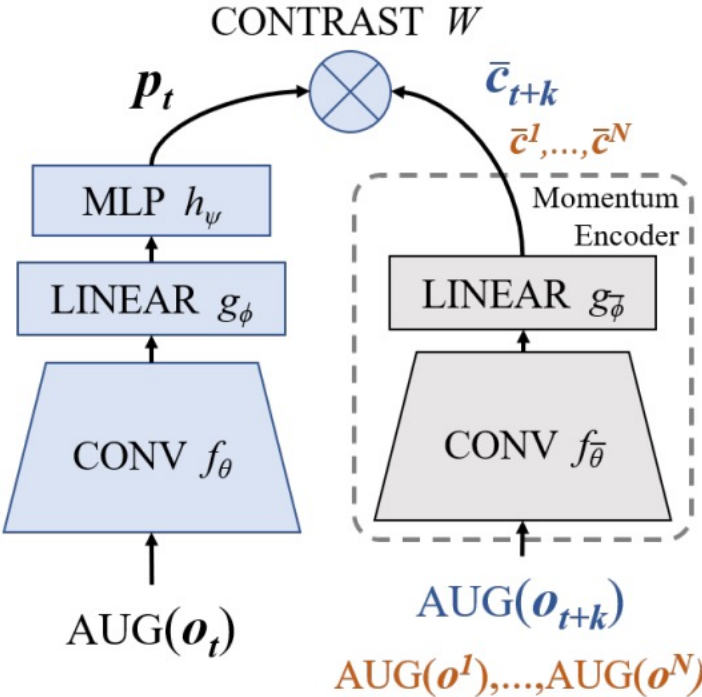
Offline Representation Learning

- CURL (Learning Spatial Invariance)
 - Maximize $I(Z(X_1); Z(X_2))$ where X_1, X_2 are same data with different augmentation.



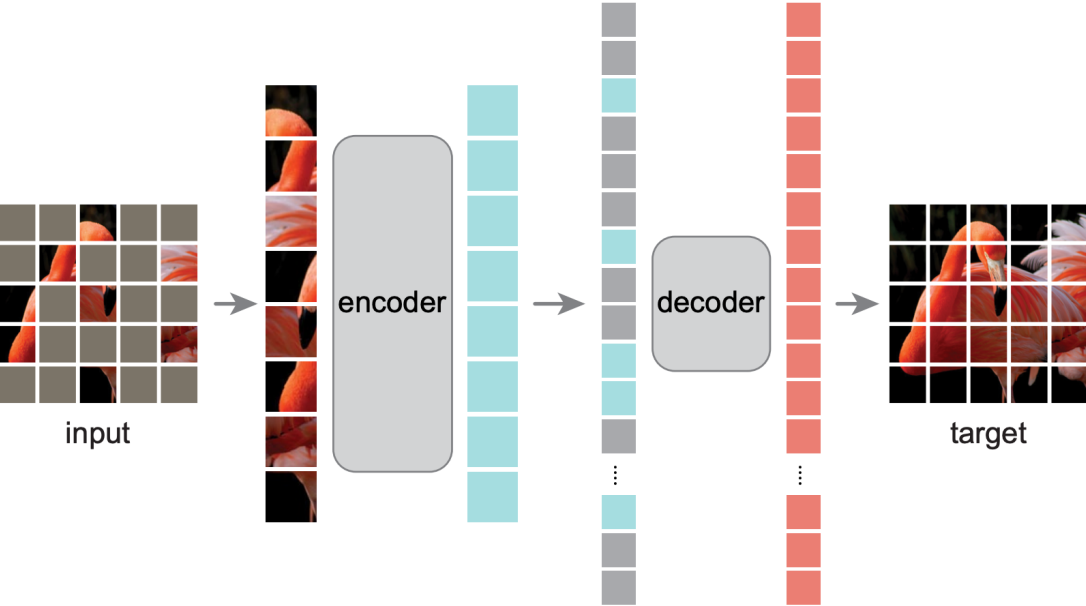
Offline Representation Learning

- ATC (Learning Spatiotemporal Invariance)
 - Maximize $I(Z(X_{t+k}); Z(X_t))$ where X_{t+k}, X_t are data from same trajectory.

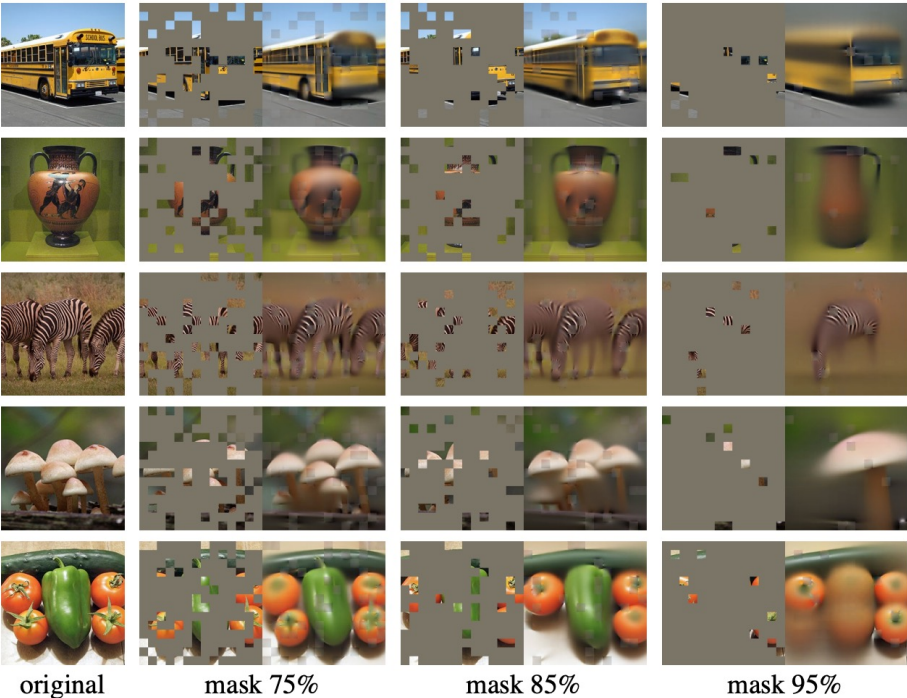


Offline Representation Learning

- MAE (Spatial Generation)
 - Maximize $I(Z(\bar{X}); X)$ where \bar{X} is perturbed image of X .



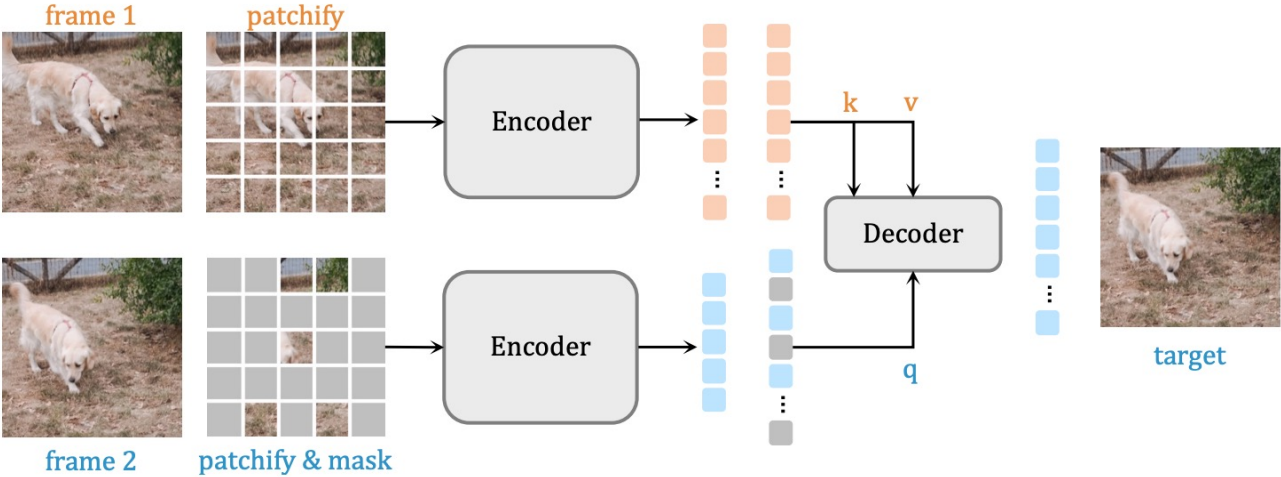
Masked autoencoder (like BERT)



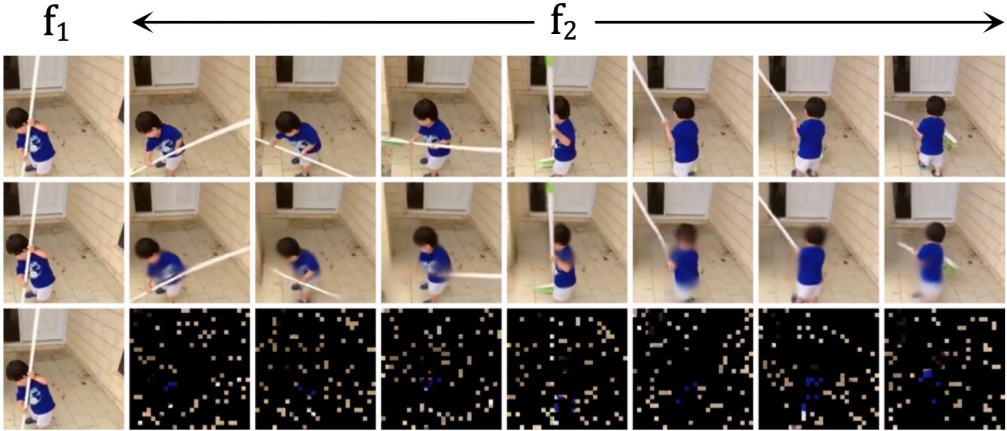
Reconstruction

Offline Representation Learning

- SiamMAE (Spatiotemporal Generation)
 - Maximize $I(Z(\bar{X}_{t+k}); X_t)$ where \bar{X}_{t+k} is perturbed image of X_{t+k} .



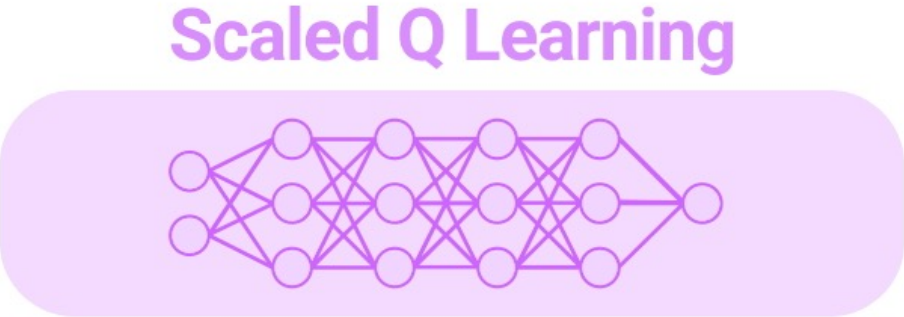
Masked temporal autoencoder



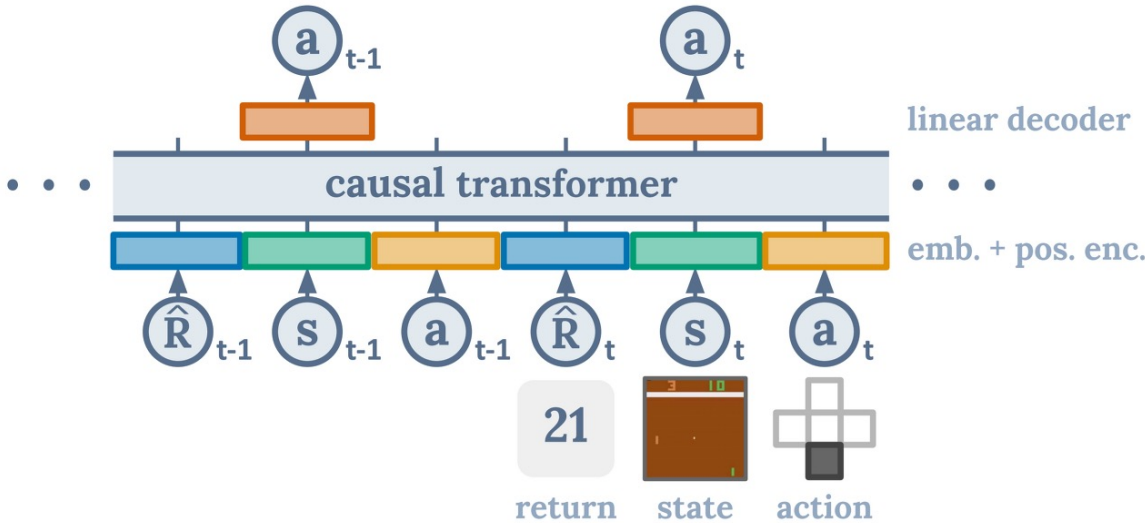
Reconstruction

Offline Reinforcement Learning

- Train effective Perceiver Actor from offline dataset.
- Q-learning vs Sequence Modeling



Q-Learning (Optimizing Bellman Equation)

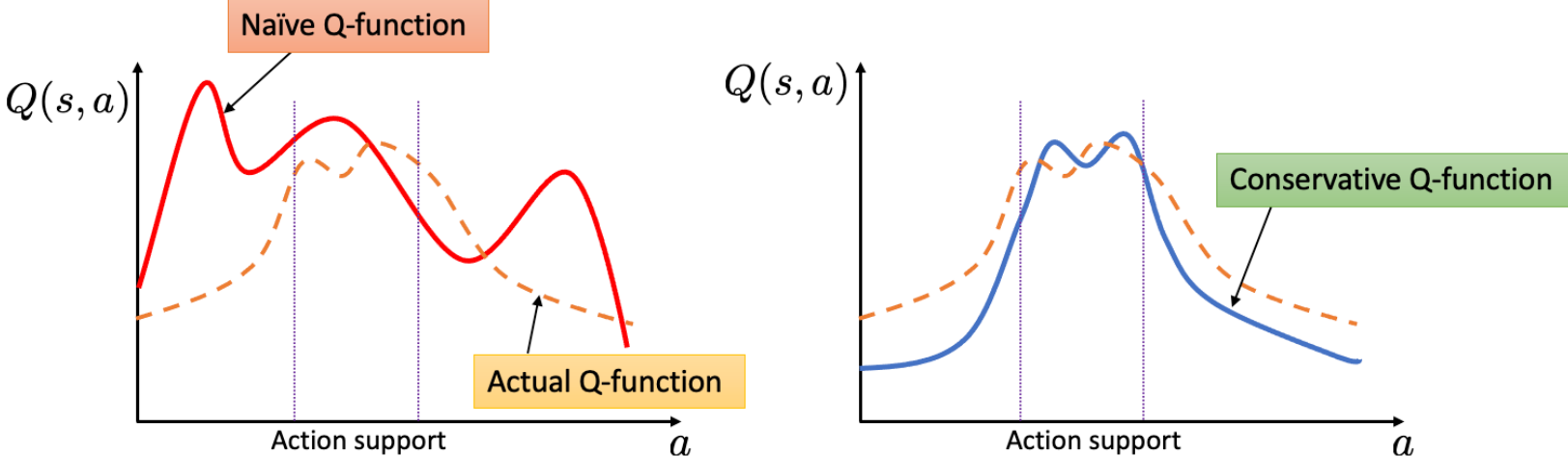


Sequence Modeling

Offline Reinforcement Learning

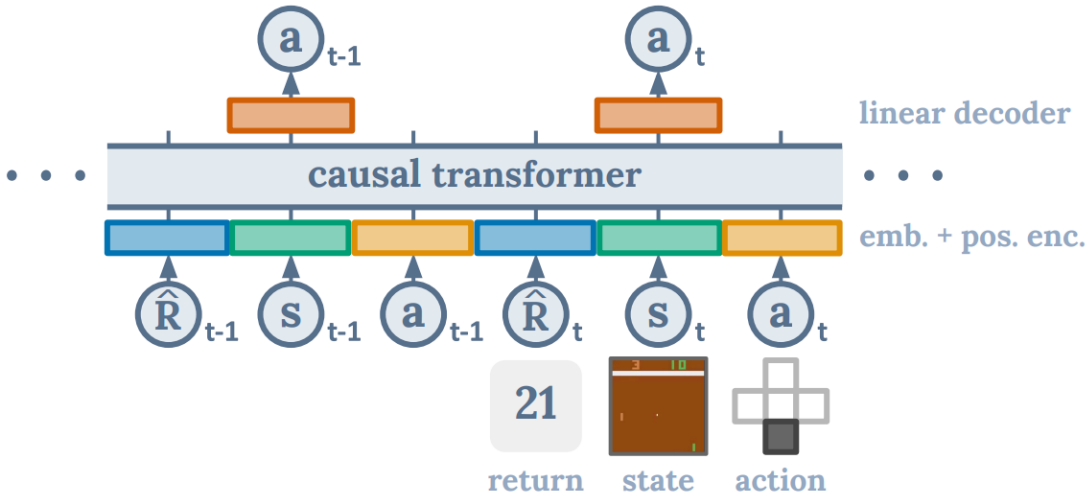
- CQL (Conservative Q-Learning)
 - Applying Q-learning to offline dataset will cause **extrapolation error**.
 - Extrapolation error brings overestimation biases.
 - Solve it by conservative Q-value estimation to unseen actions.

$$\hat{Q}_{\text{CQL}}^\pi := \arg \min_Q \alpha \cdot \left(\underbrace{\mathbb{E}_{\mathbf{s} \sim \mathcal{D}, \mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} [Q(\mathbf{s}, \mathbf{a})]}_{\text{minimize Q-values}} - \underbrace{\mathbb{E}_{\mathbf{s} \sim \mathcal{D}, \mathbf{a} \sim \hat{\pi}_\beta(\mathbf{a}|\mathbf{s})} [Q(\mathbf{s}, \mathbf{a})]}_{\text{maximize Q-values under data}} \right) + \frac{1}{2} \underbrace{\mathbb{E}_{\mathbf{s}, \mathbf{a}, \mathbf{s}' \sim \mathcal{D}} \left[\left(Q - \hat{B}^\pi Q \right)^2 \right]}_{\text{standard Bellman error}}$$

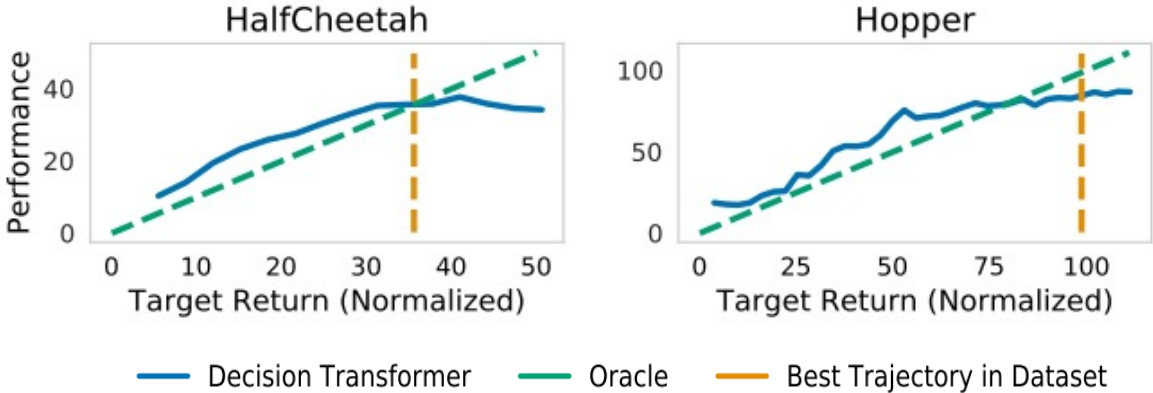


Offline Reinforcement Learning

- DT (Decision Transformer)
 - Formulate RL as a big sequence modeling problem.



Causal Transformer as RL agent

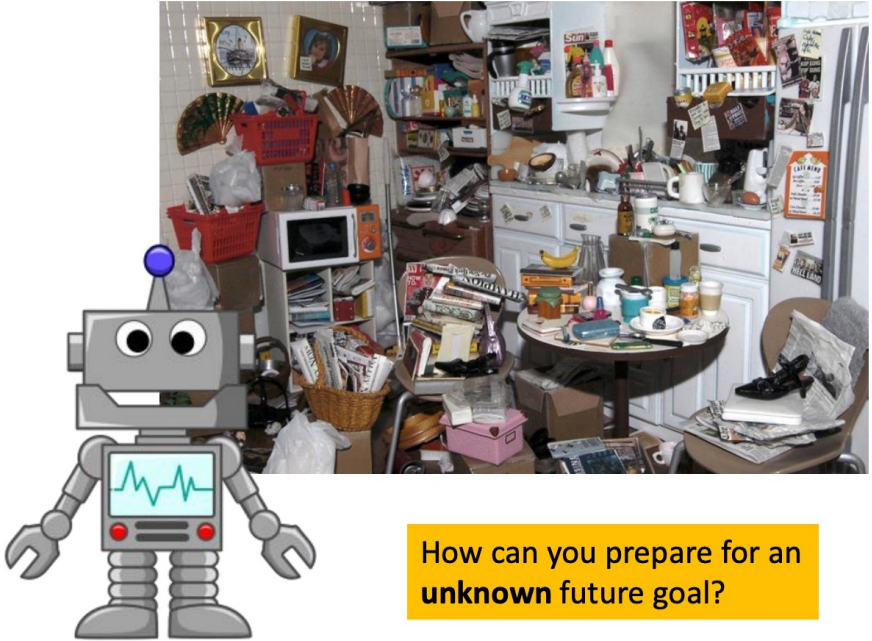


Reward-conditioned policy

Pretraining from Simulator

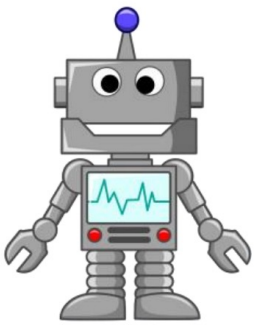
Unsupervised Reinforcement Learning

- Train effective Perceiver Actor from simulator.
- Assumes we do not have an access to a pre-defined reward function.



How can you prepare for an **unknown** future goal?

training time: unsupervised

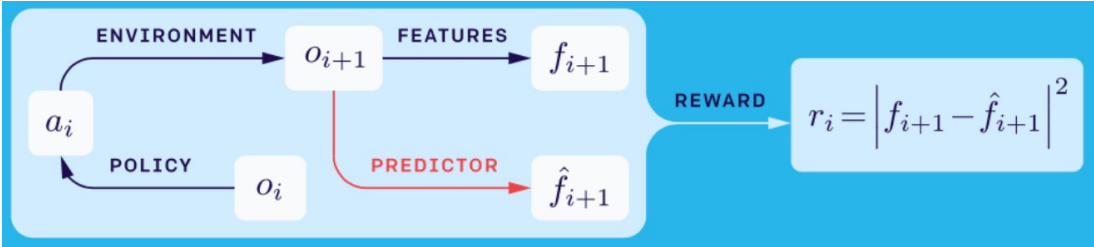


Unsupervised Reinforcement Learning

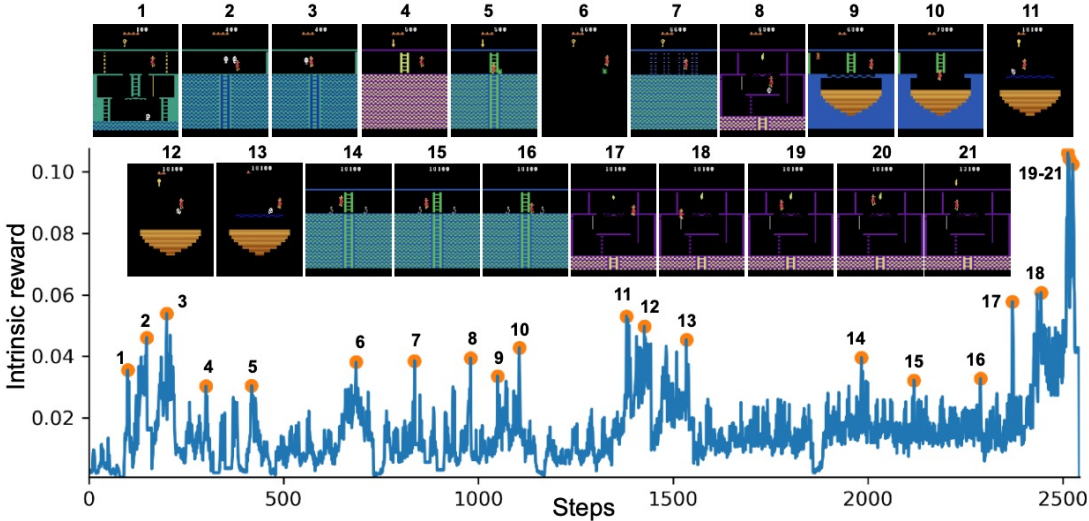
- Curiosity-driven exploration
 - Explore 'curious' states.
- Data coverage maximization
 - Maximize the 'coverage of data' collected through pretraining.
- Skill discovery
 - Learn task-agnostic skills.
- World model
 - Learn dynamics of the environment.

Unsupervised Reinforcement Learning

- RND (Curiosity Driven Exploration)
 - Low prediction error: high reward



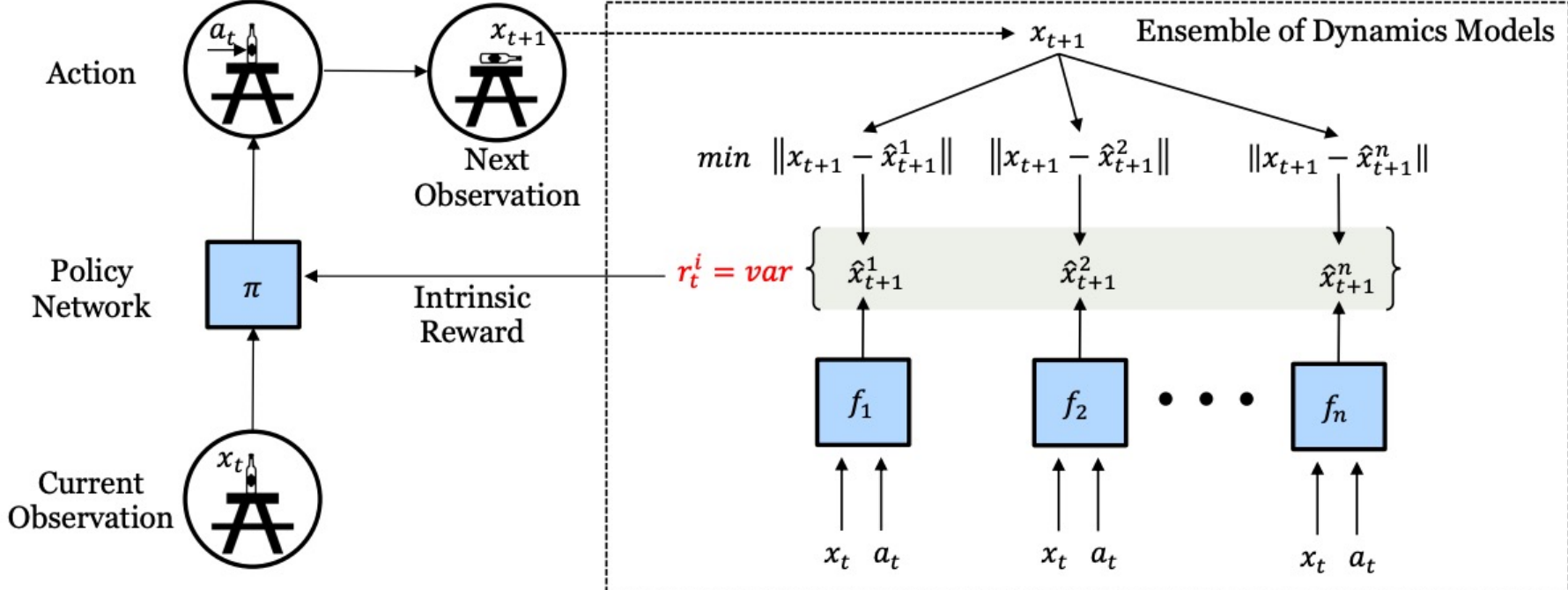
Prediction error as reward.



Spikes in reward correspond to meaningful events.

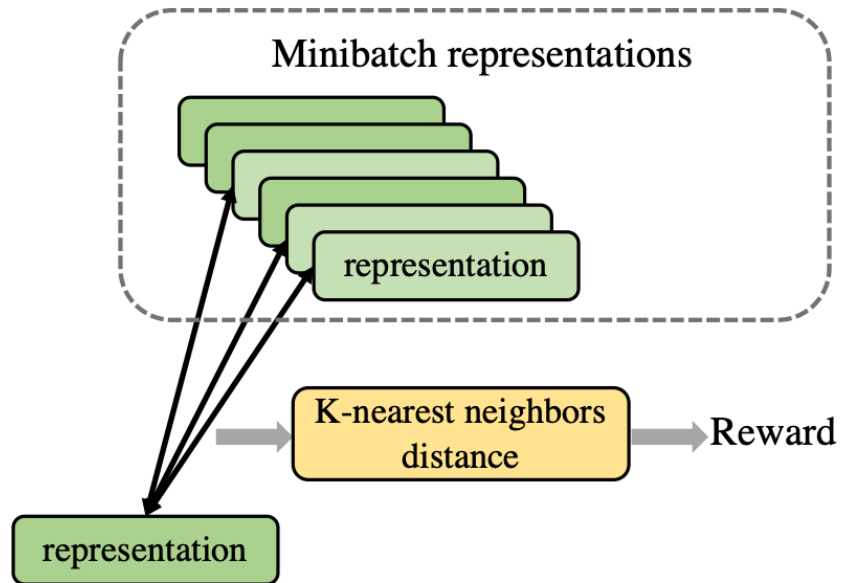
Unsupervised Reinforcement Learning

- DISAGREE (Curiosity Driven Exploration)
 - Disagreement among ensembles: high reward



Unsupervised Reinforcement Learning

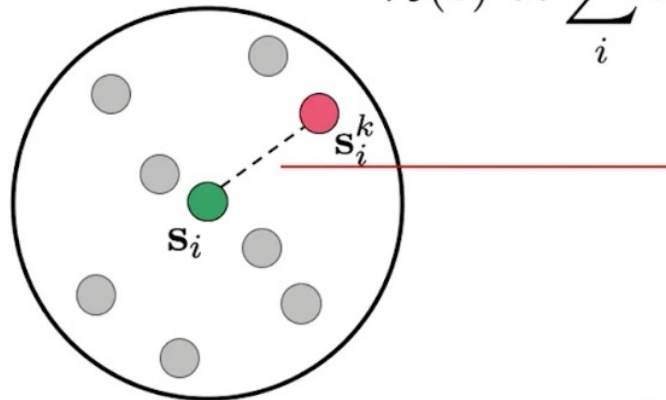
- APT (Data-coverage maximization)
 - Maximize the state entropy ($H(d_\pi)$) in replay buffer



State entropy as reward

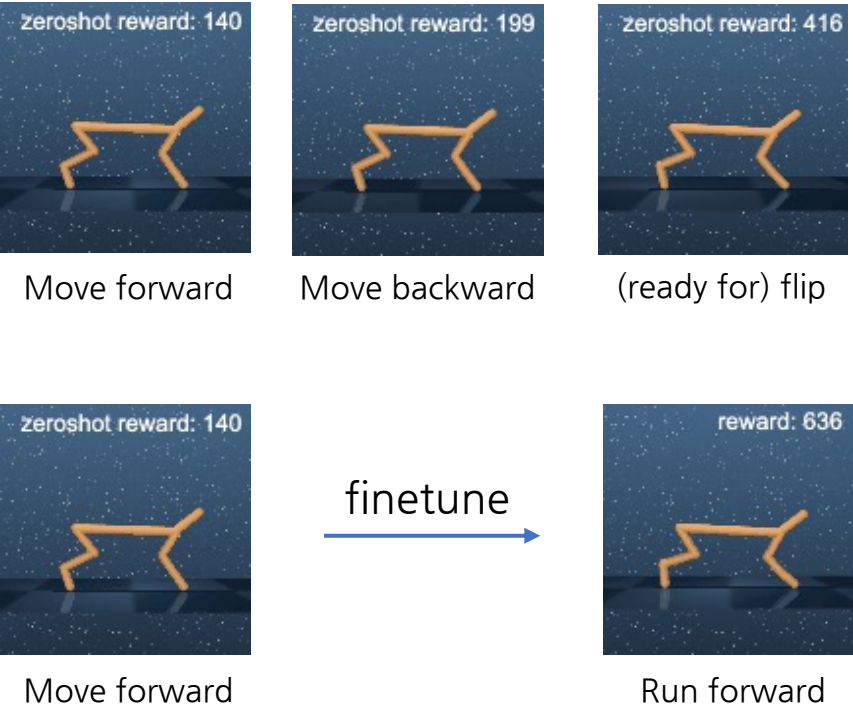
$$\mathcal{H}(s) = -\mathbb{E}_{s \sim p(s)} [\log p(s)]$$

$$\hat{\mathcal{H}}(s) \propto \sum_i \log(\|s_i - s_i^k\|)$$



Unsupervised Reinforcement Learning

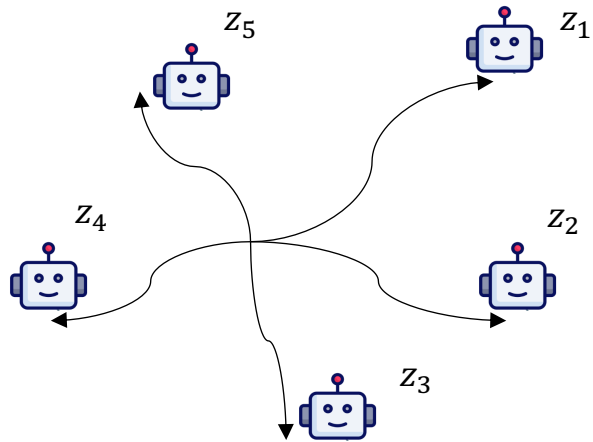
- DIAYN (Skill-Discovery)
 - Maximize the mutual information between state and skill $I(s; z)$



Diversity is all you need: learning skills without a reward function, Eysenbach et al. ICLR 2019.

Unsupervised Reinforcement Learning

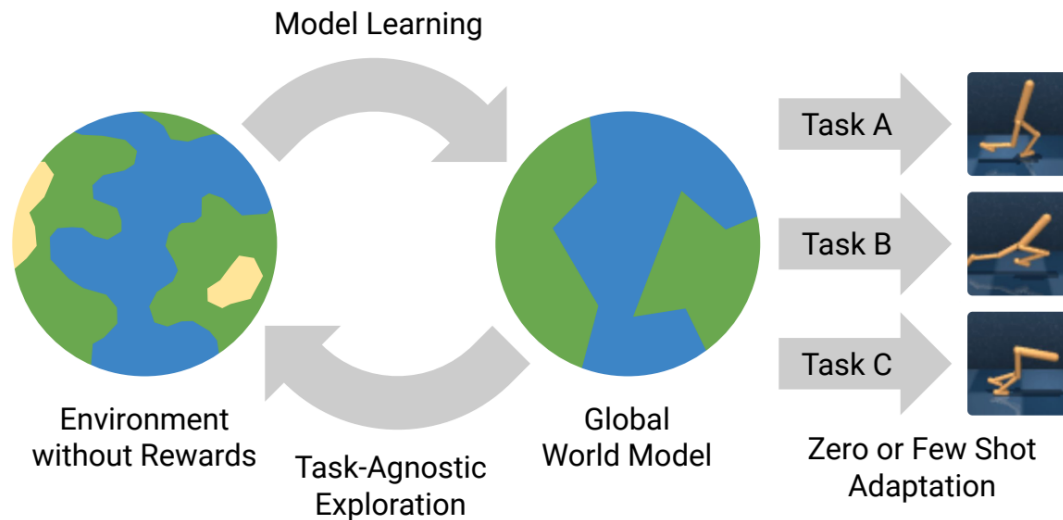
- DIAYN (Skill-Discovery)
 - Maximize the mutual information between state and skill $I(s; z)$
 1. Sample skill $z \sim p(z)$
 2. Rollout trajectory $\tau \sim \pi(a|s, z)$
 3. Maximize mutual information between skill z and state s



Mutual information as reward

Unsupervised Reinforcement Learning

- Plan2Explore (World-Model)
 - Learn various components to represent the environment.
 1. Dynamics model $f(s_{t+1}|s_t, a_t)$
 2. (optional) reward predictor, image encoder/decoder, ...

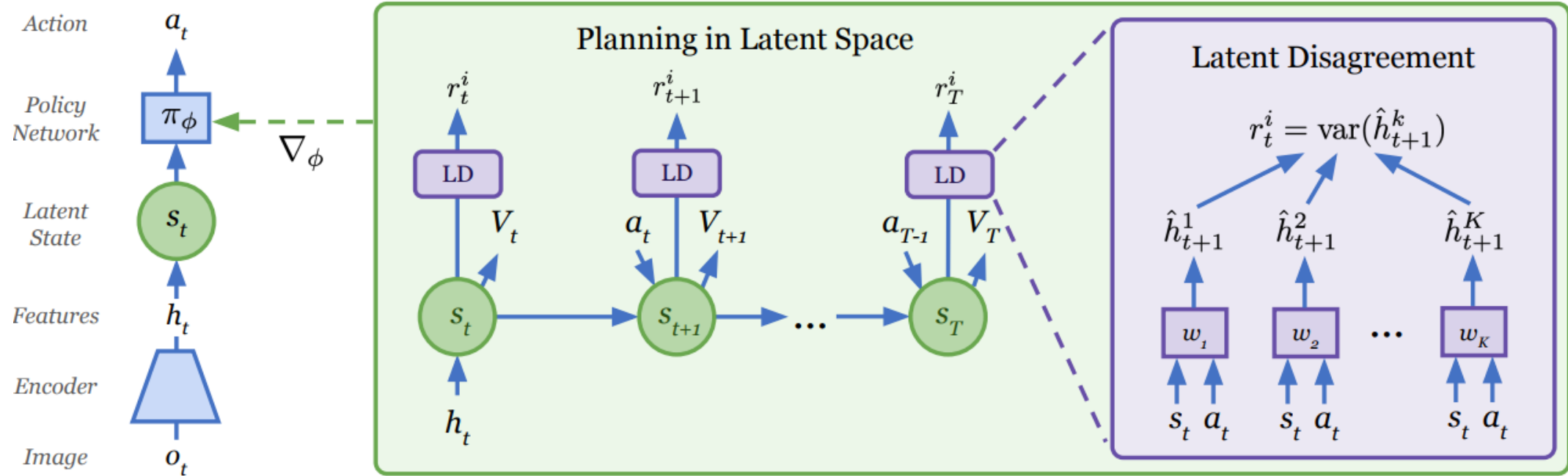


Pretraining: learn world model

Downstream task: now we can **plan** with world model

Unsupervised Reinforcement Learning

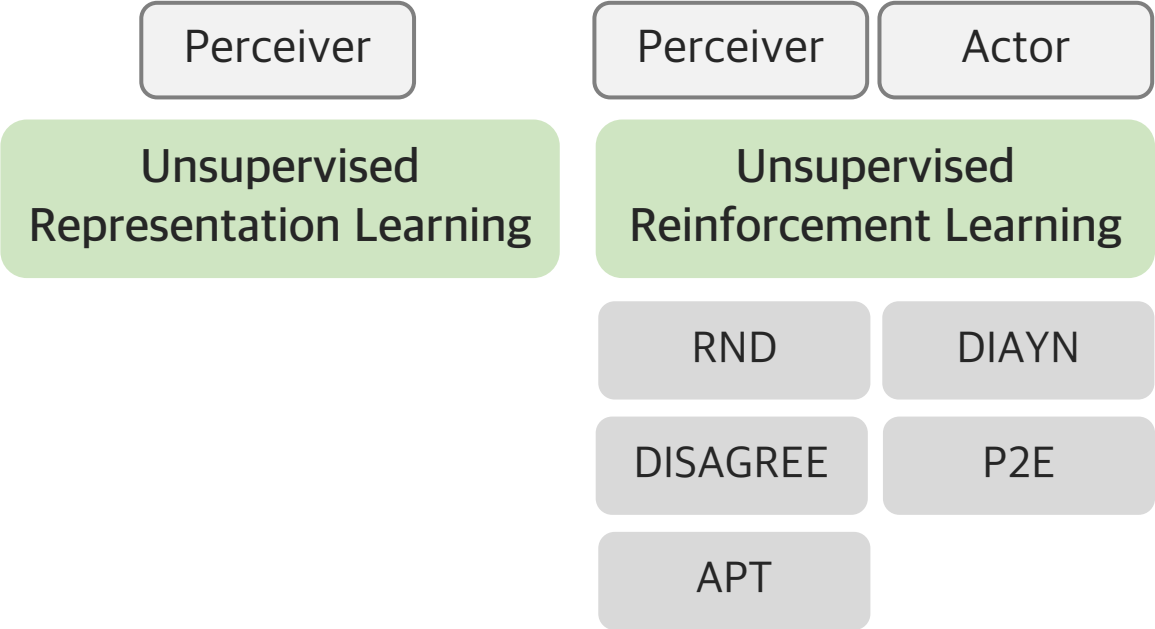
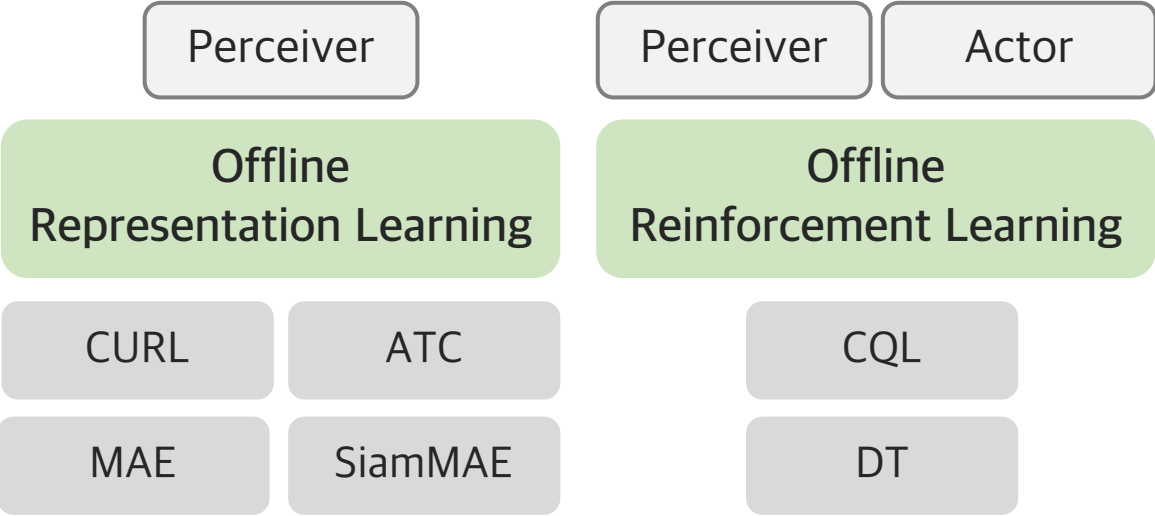
- Plan2Explore (World-Model)
 - Learn world model with disagreement



Planning to explore via self-supervised world models, Sekar et al. ICML 2020.

Summary: Let's Train Robot

Summary



Q & A