# Skill-based Model-based Reinforcement Learning

Lucy Xiaoyang Shi, Joseph J. Lim, Youngwoon Lee.

2023.02.23

김나영

# Skill-based Model-based Reinforcement Learning

Lucy Xiaoyang Shi, Joseph J. Lim, Youngwoon Lee.

## Motivation

- Planning every action for long-horizon tasks is not practical
- Humans efficiently plan with high-level skills to solve complex tasks

## Approach

- Skill-based Model-based RL framework (SkiMo), which directly predicts the skill outcomes, rather than predicting all small details in the intermediate states, step by step

## Results & New finding

- SkiMo extends the temporal horizon of model-based approaches and improves the sample efficiency for model-based RL and skill-based RL

## Discussion & Comments

- SkiMo is a general framework that can be extended to RGB, depth, and tactile observations

# Motivation

Idea from Human Intelligence
- The ability to plan abstractly for solving complex tasks
- It can be used to scale the model to long-horizon tasks by reducing the search space of behaviors

Suggestion
- A novel skill-based and model-based reinforcement learning (RL) method,
  which learns a model and a policy in a high-level skill space,
  enabling accurate long-term prediction and efficient long-term planning

# Model-Based RL Model

# Skill-Based RL Model

Mechanism
- Learning <u>a flat single-step dynamics model</u>, which predicts the next state from the current state and action

Pros
- Can be used to simulate imaginary trajectories, which <u>improves sample efficiency</u> over model-free alternatives

Cons
- Only limited success in long-horizon tasks due to
  - (1) inaccurate long-term prediction
  - (2) computationally expensive search

Mechanism
- To solve long-horizon tasks by acting <u>with multi-action subroutines (skills)</u>

Pros
- (1) Enables systematic <u>long-range exploration</u> to plan farther into the future
- (2) Requires <u>a shorter horizon for policy optimization</u>, which makes long-horizon downstream tasks more tractable

Cons
- Requires a few million to billion environment interactions to learn

4

Open microwave | Move kettle

(a) Flat dynamics model without skills

(b) Flat dynamics model with skills

(c) Skill dynamics model with skills

## Skill-Based RL Model

## Model-Based RL Model

Directly predicts the resultant state <u>after skill execution</u>, without neeing to model every intermediate step and low-level action

Predicts the immediate next state <u>after one action execution</u>

# Wu et al.

- A temporally-extended dynamics model

Limit
- Conditions on <u>low-level actions</u> rather than skills
- Only used for low-level planning

# Shaha et al.

- Learns a skill dynamics model

Limit
- A limited set of <u>discrete, manually-defined skills</u>

# SkiMo

Mechanism
- Extract the skill space from data
- Devise a skill-level dynamics model

Meaning
- SkiMo is the first work that jointly <u>learns skills and a skill dynamics model from data</u> for model-based RL

# Preliminaries

- Formulate a problem as a Markov decision process

Unlabeled Offline Data
- Assume reward-free task-agnostic dataset, which is a set of <u>N state-action trajectories</u>
- Do not assume this dataset contains solutions for the downstream task,
  tackling the downstream task <u>requires recomposition</u> of skills learned from diverse trajectories

Skill-based RL
- Skills = A sequence of actions with <u>a fixed horizon H</u>
- Parameterize skills as a <u>skill latent z</u> and <u>skill policy π</u>, that maps <u>a skill latent</u> and <u>state</u>
- (Step 1) <u>The skill latent</u> and <u>skill policy</u> can be trained using variational auto-encoder (VAE)
  - <u>A skill encoder</u> embeds a sequence of transitions into a skill latent z
  - <u>A skill policy</u> decodes it back to the original actions sequence
- (Step 2) Learn <u>a skill prior p(z|s)*</u> to guide the downstream task policy to explore promising skills

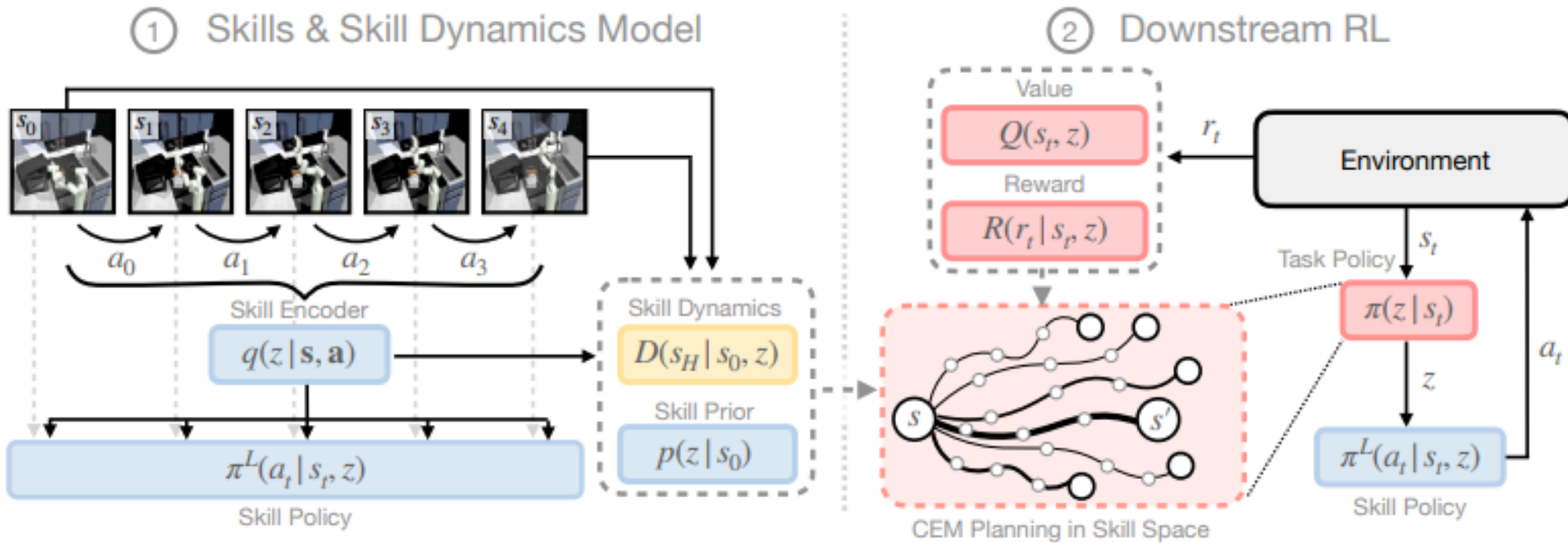* A skill prior : the skill distribution in the offline data                          7

# SkiMo Model Components

- (1) Skill policy
- (2) Skill dynamics model
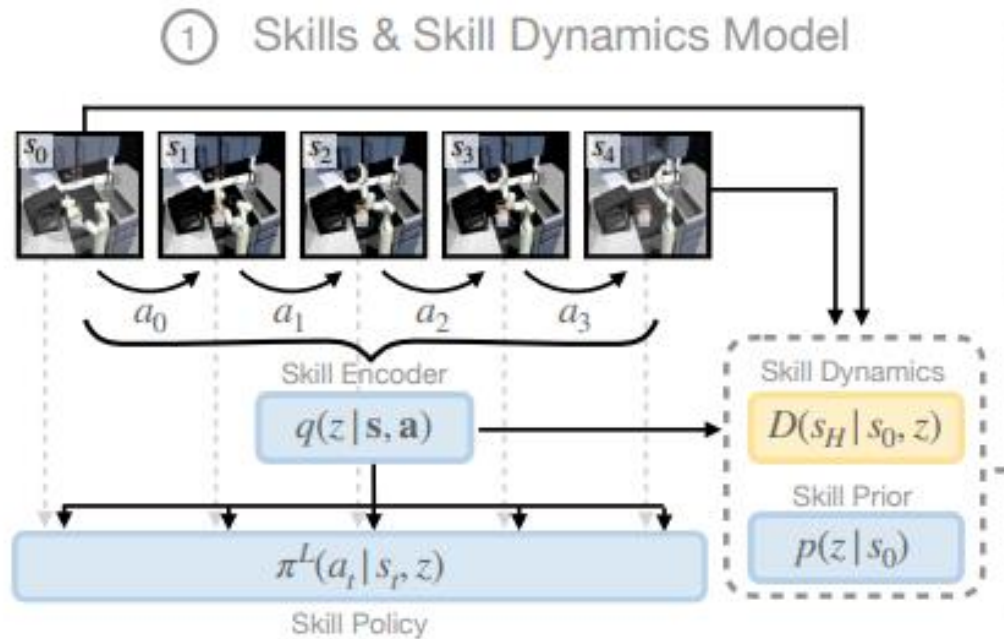- (3) Task policy

# SkiMo Mechanism

- (Step 1) <u>A state encoder</u> encodes an observation s into the latent state h
- (Step 2) Then, given a skill z, <u>the skill dynamics model</u> predicts the skill effect in the latent space
- (Step 3) <u>The task policy</u>, <u>reward function</u>, and <u>value function</u> predict a skill, reward, and value on the (imagined) latent state, respectively

① Skills & Skill Dynamics Model ② Downstream RL

- SkiMo consists of two phases
  - (1) Learning the skill dynamics model and skills from an offline dataset
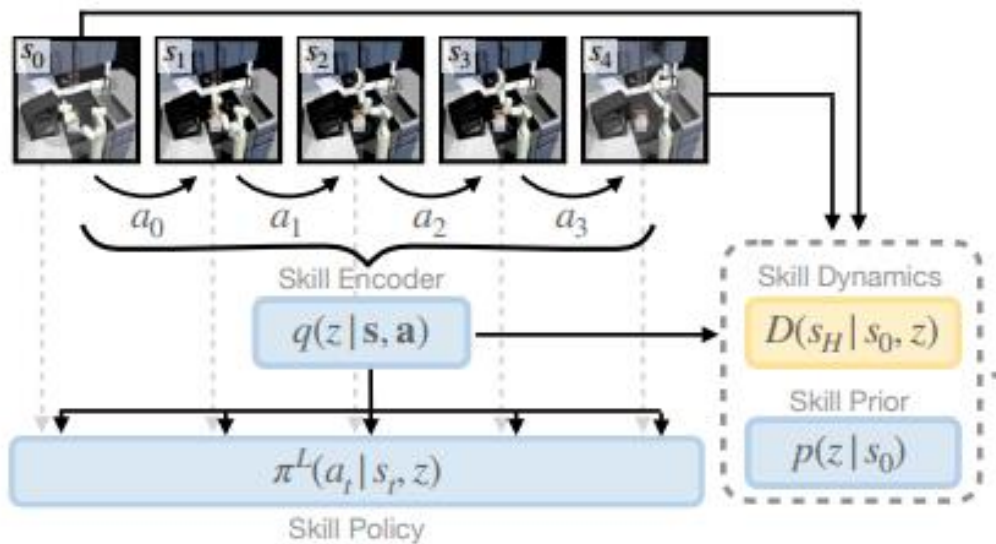  - (2) Downstream task learning with the skill dynamics mode

# Pre-Training Skill Dynamics Model and Skills from Task-agnostic Data



- SkiMo leverages offline data to extract
  - (1) <u>skills</u> for temporal abstraction of actions
  - (2) <u>skill dynamics</u> for skill-level planning
  - on a latent state space
  - (3) <u>a skill prior</u> to guide exploration

- <u>Jointly learn</u> a skill policy and skill dynamics model, in a self-supervised manner
  - Shape <u>the latent skill space Z</u> and <u>state embedding</u>

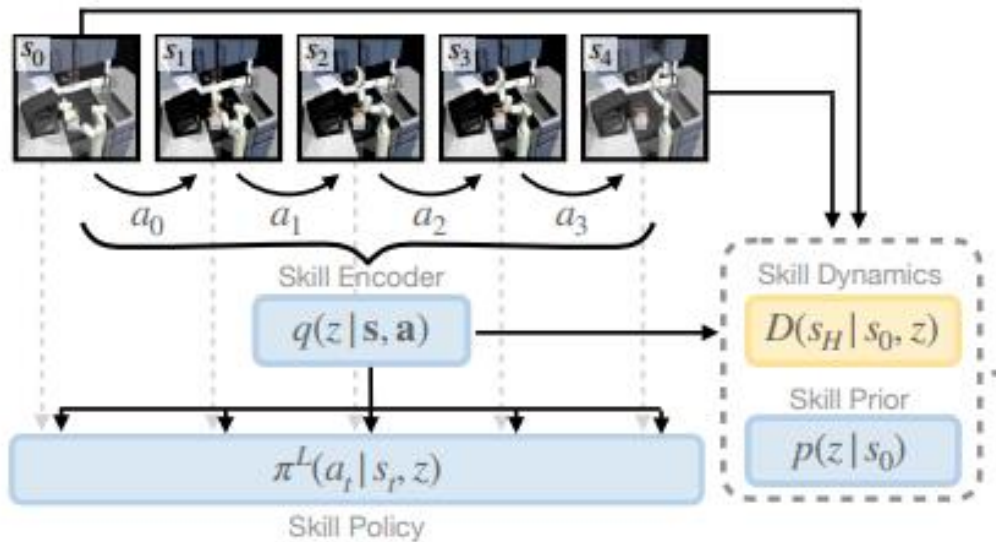# Pre-Training Skill Dynamics Model and Skills from Task-agnostic Data



(1) Skill policy
- To learn a low-dimensional <u>skill latent space Z*</u>, we train a conditional VAE on the offline dataset that reconstructs the action sequence through a skill embedding

- Given H consecutive states and actions,
  - (1) <u>a skill encoder</u> predicts a skill embedding z
  - (2) <u>a skill decoder</u> recontructs the original action sequence from z

$$\mathcal{L}_{\text{VAE}} = \mathbb{E}_{(\mathbf{s},\mathbf{a})_{0:H-1} \sim \mathcal{D}} \left[ \frac{\lambda_{\text{BC}}}{H} \sum_{i=0}^{H-1} \underbrace{(\pi_\theta^L(\mathbf{s}_i, \mathbf{z}) - \mathbf{a}_i)^2}_{\text{Behavioral cloning}} + \beta \cdot \underbrace{KL\big(q_\theta(\mathbf{z}|(\mathbf{s},\mathbf{a})_{0:H-1}) \parallel p(\mathbf{z})\big)}_{\text{Embedding regularization}} \right], \quad (2)$$

* A skill latent space : it encodes action sequences

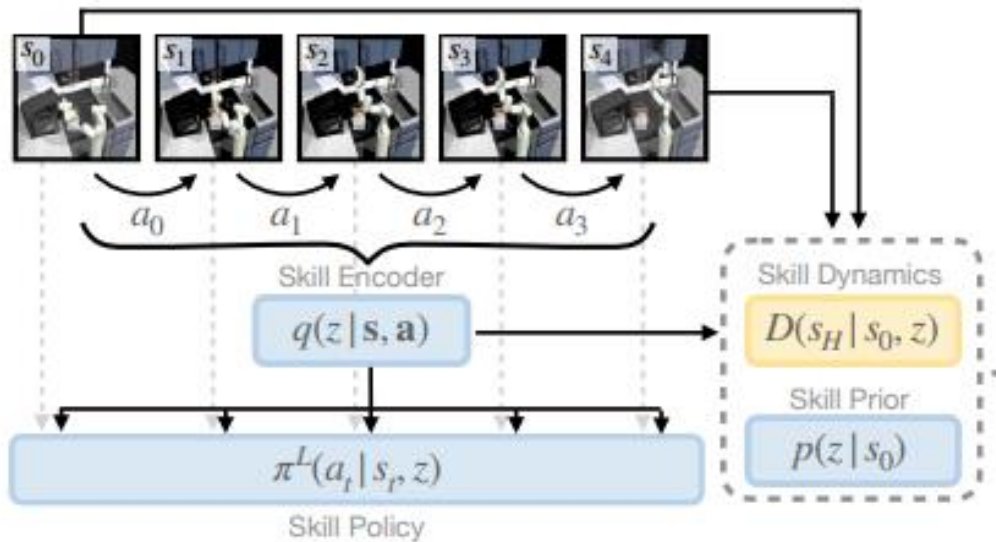# Pre-Training Skill Dynamics Model and Skills from Task-agnostic Data



(2) Skill dynamics model
- (1) Learns to predict the latent state H-steps ahead conditioned on a skill, for N sequential skill transitions using the latent state consistency loss

- (2) To prevent a trivial solution and encode rich information from observations, we additionally train an observation decoder using the observation reconstruction loss

$$\mathcal{L}_{\text{REC}} = \mathbb{E}_{(\mathbf{s},\mathbf{a})_{0:NH} \sim \mathcal{D}} \left[ \sum_{i=0}^{N-1} \left[ \underbrace{\lambda_{\text{O}} \| \mathbf{s}_{iH} - O_\theta(E_\psi(\mathbf{s}_{iH})) \|_2^2}_{\text{Observation reconstruction}} + \underbrace{\lambda_{\text{L}} \| D_\psi(\hat{\mathbf{h}}_{iH}, \mathbf{z}_{iH}) - E_{\psi^-}(\mathbf{s}_{(i+1)H}) \|_2^2}_{\text{Latent state consistency}} \right] \right] \quad (3)$$

# Pre-Training Skill Dynamics Model and Skills from Task-agnostic Data



(3) Skill prior
- is trained by minimizing <u>the KL divergence</u> between output distributions of <u>the skill encoder</u> and <u>the skill prior</u>

$$\mathcal{L}_{\text{SP}} = \mathbb{E}_{(\mathbf{s},\mathbf{a})_{0:H-1} \sim \mathcal{D}} \left[ \lambda_{\text{SP}} \cdot KL\Big( \mathbf{sg}(q_\theta(\mathbf{z}|\mathbf{s}_{0:H-1}, \mathbf{a}_{0:H-1})) \,\|\, p_\theta(\mathbf{z}|\mathbf{s}_0) \Big) \right], \qquad (4)$$
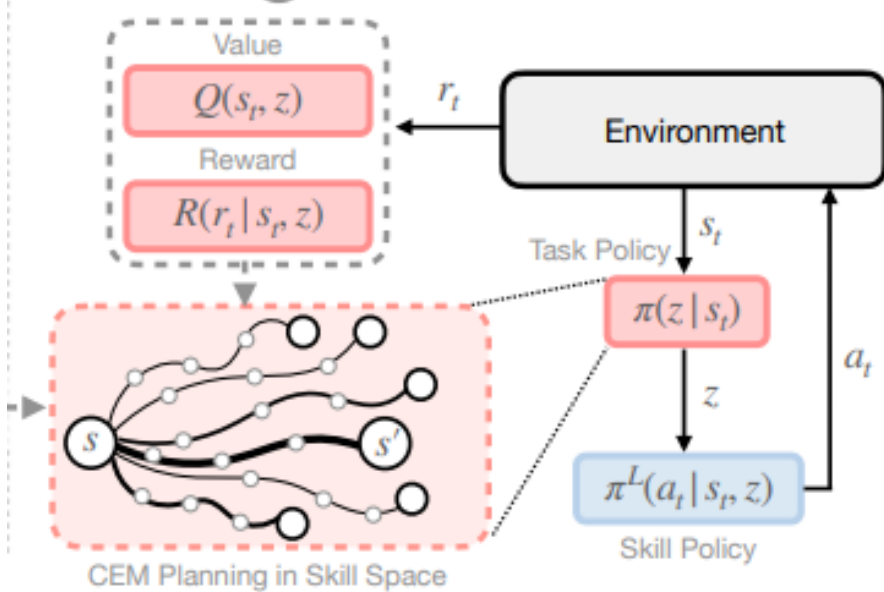
# Pre-Training Skill Dynamics Model and Skills from Task-agnostic Data

$$\mathcal{L} = \mathcal{L}_{\text{VAE}} + \mathcal{L}_{\text{REC}} + \mathcal{L}_{\text{SP}}$$

- Jointly train the policy, model, and prior, which leads to a well-shaped skill latent space that is optimized for both skill reconstruction and long-term prediction

14

# Downstream Task Learning with Learned Skill Dynamics Model



② Downstream RL

Value
$Q(s_t, z)$
Reward
$R(r_t \mid s_t, z)$

Environment

$r_t$

Task Policy
$\pi(z \mid s_t)$

$s_t$

$a_t$

$z$

$\pi^L(a_t \mid s_t, z)$

Skill Policy

CEM Planning in Skill Space

- (1) Learns a high-level <u>task policy</u> that outputs a latent skill embedding z
- (2) skill embedding z is then translated into a sequence of H actions using the pre-trained <u>skill policy</u>

# Downstream Task Learning with Learned Skill Dynamics Model

STEP 1
- The skill dynamics model and task policy can generate imaginary rollouts in the skill space by repeating
  - (1) sampling a skill
  - (2) predicting H-step future after executing the skill

STEP 2
- To evaluate imaginary rollouts, we train a reward function and Q-value function

$$\mathcal{L}'_{REC} = \mathbb{E}_{\mathbf{s}_t, \mathbf{z}_t, \mathbf{s}_{t+H}, r_t \sim \mathcal{D}} \left[ \underbrace{\lambda_{\mathrm{L}} \| D_\psi(\hat{\mathbf{h}}_t, \mathbf{z}_t) - E_{\psi^-}(\mathbf{s}_{t+H}) \|_2^2}_{\text{Latent state consistency}} + \underbrace{\lambda_{\mathrm{R}} \| r_t - R_\phi(\hat{\mathbf{h}}_t, \mathbf{z}_t) \|_2^2}_{\text{Reward prediction}} \right.$$

$$\left. + \underbrace{\lambda_{\mathrm{V}} \| r_t + \gamma Q_{\phi^-}(\hat{\mathbf{h}}_{t+H}, \pi_\phi(\hat{\mathbf{h}}_{t+H})) - Q_\phi(\hat{\mathbf{h}}_t, \mathbf{z}_t) \|_2^2}_{\text{Value prediction}} \right]. \quad (6)$$

STEP 3
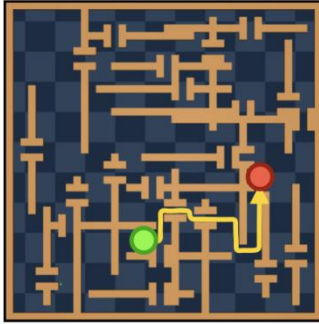- Finetune the skill dynamics model and state encoder on the downstream task

STEP 4
- Train a high-level task policy to maximize the estimated Q-value while regularizing it to the pre-trained skill prior
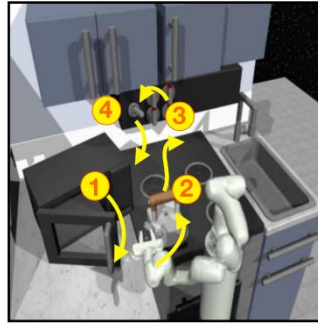
$$\mathcal{L}_{RL} = \mathbb{E}_{\mathbf{s}_t \sim \mathcal{D}} \left[ - Q_\phi(\hat{\mathbf{h}}_t, \pi_\phi(\mathrm{sg}(\hat{\mathbf{h}}_t))) + \alpha \cdot KL\big(\pi_\phi(\mathbf{z}_t | \mathrm{sg}(\hat{\mathbf{h}}_t)) \, \| \, p_\theta(\mathbf{z}_t | \mathbf{s}_t)\big) \right]. \quad (7)$$
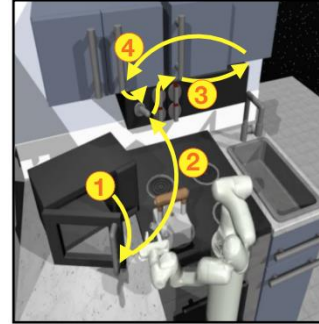
16

# Tasks
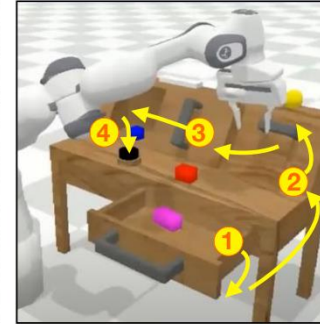


(a) Maze          (b) Kitchen          (c) Mis-aligned Kitchen          (d) CALVIN

- Compare SkiMo with prior model-based RL and skill-based RL methods on four long-horizon tasks with sparse rewards

Maze
- (goal) to reach the fixed goal region in red
- (reward) a sparse reward of 100 only when it reaches the goal

Kitchen
- (goal) to perform four sequential subtasks
- (reward) a reward of 1 for every sub-task completion in order

Mis-aligned Kitchen
- (goal) to perform four sequential subtasks, has a low sub-task transition probability
- (reward) a reward of 1 for every sub-task completion in order

CALVIN
- (goal) to perform four sequential subtasks
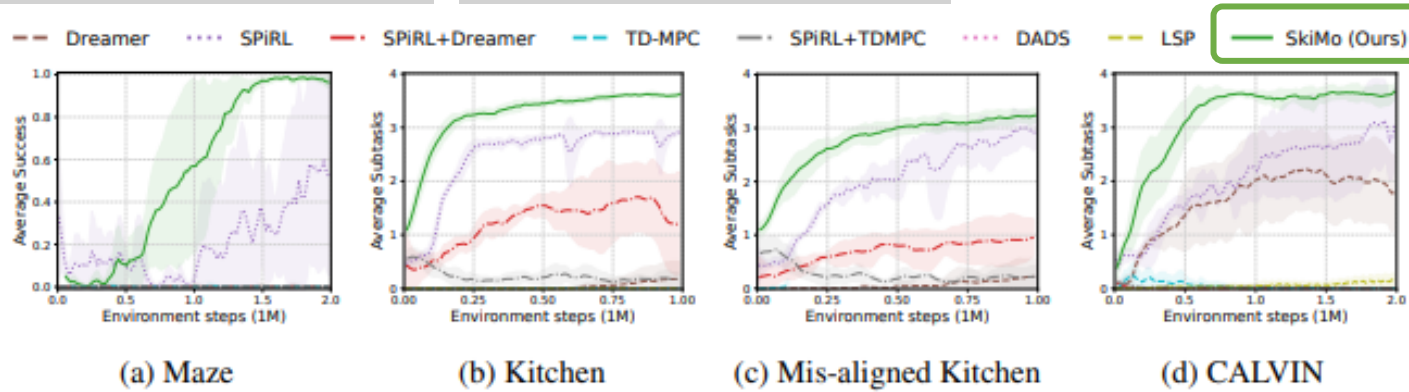- (reward) a reward of 1 for every sub-task completion in order

17

# Results



Figure 4: Learning curves of our method and baselines. All averaged over 5 random seeds.

Maze
- A hard exploration problem <u>due to the sparsity of the reward</u>: the agent only receives reward after taking 1,000+ steps to reach the goal

Kitchen
- SkiMo reaches the same performance with <u>5x less environment interactions</u> than SPiRL

Mis-aligned Kitchen
- makes the downstream learning harder because <u>the skill prior offers less meaningful regularization to the policy</u>

CALVIN
- Is very task-agnostic: any particular sub-task transition has probability lower than 0.1% on average
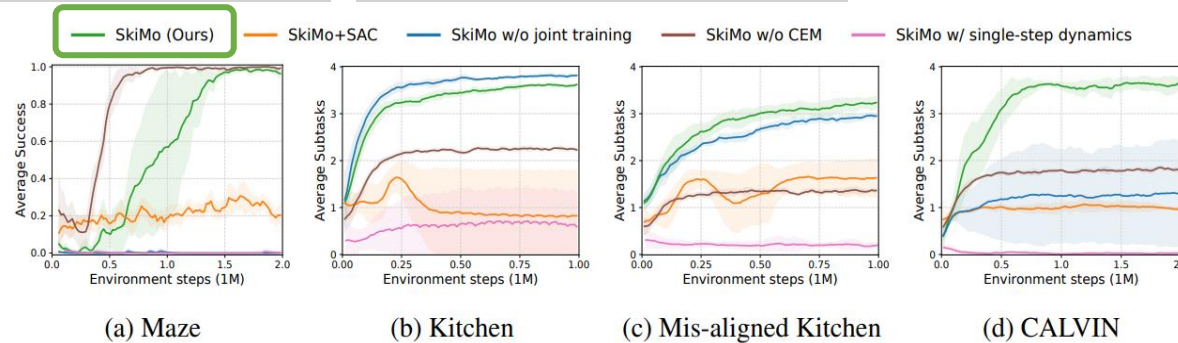
18

# Ablation Studies



Figure 5: Learning curves of our method and ablated methods. All averaged over 5 random seeds.

Model-based VS. Model-free
- SkiMo achieves better asymptotic performance and higher sample efficiency across all tasks than SkiMo+SAC, which uses model-free RL

Joint training of skills and skill dynamic model

CEM planning
- SkiMo learns significantly better and faster in Kitchen, Misaligned Kitchen, and CALVIN than SkiMo w/o CEM

Skill dynamics model
- The skill dynamic model can make accurate long-horizon predictions for planning due to significantly less compounding errors

- (1) A skill dynamics model <u>reduces the long-term future prediction error</u> via temporal abstraction
- (2) <u>Without needing to plan step-by-step</u>, downstream RL over the skill space allows for efficient and accurate temporally-extended reasoning
- (3) Joint training of the skill dynamics and skills further <u>improves the sample efficiency</u> by learning skills conducive to predict their consequences

Limitation and future work
- (1) While this method extracts <u>fixed-length skills</u> from offline data, the lengths of semantic skills may vary based on the contexts and goals
- (2) Further, although this experiment only focused <u>on state-based inputs</u>, SkiMo is a general framework that can be extended to RGB, depth, and tactile observations

# Skill-based Model-based Reinforcement Learning

Lucy Xiaoyang Shi, Joseph J. Lim, Youngwoon Lee.

2023.02.23

김나영