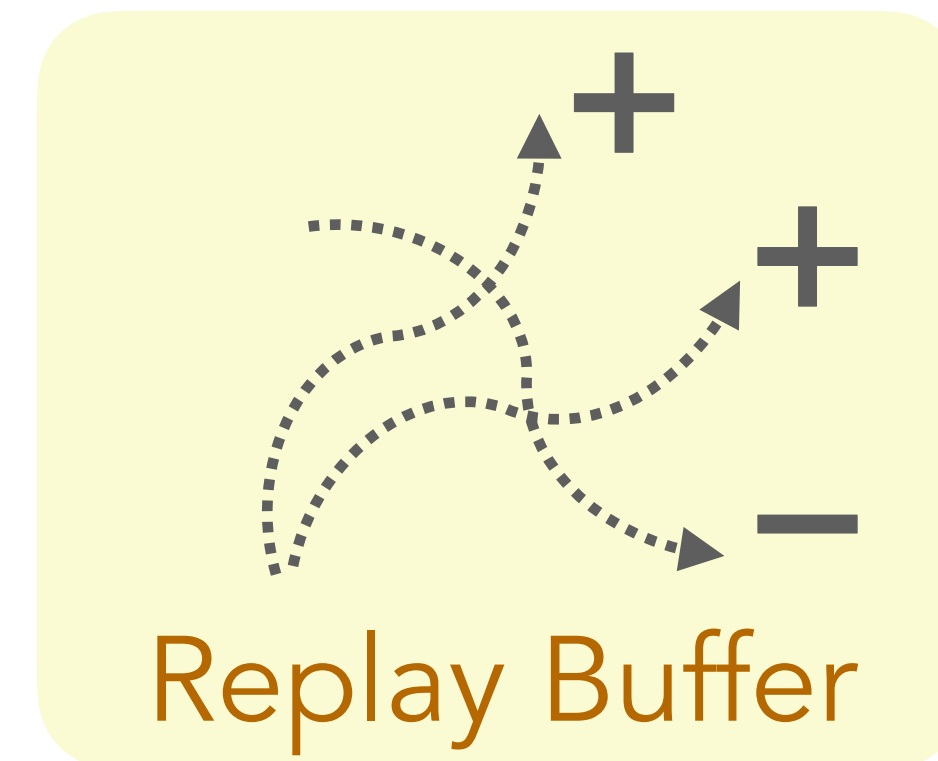
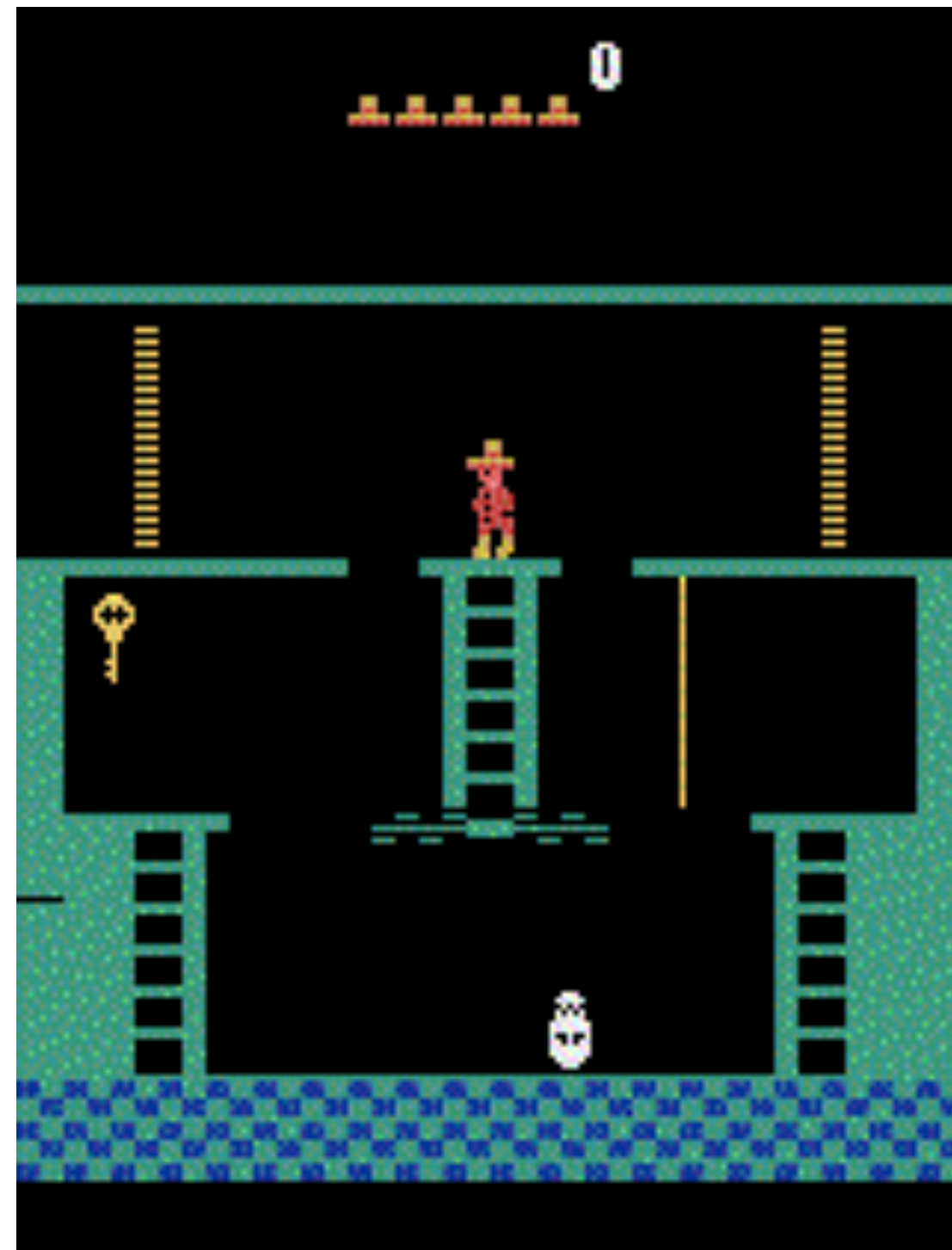


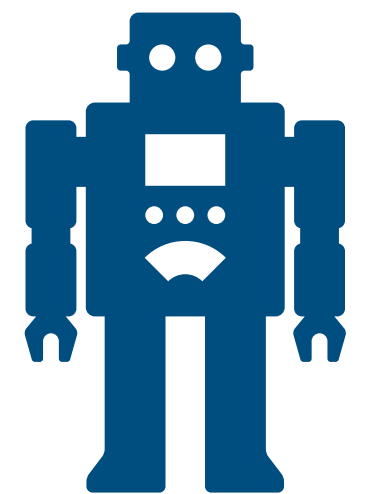
Leveraging Skills from Unlabeled Prior Data for Efficient Online Exploration

Max Wilcoxson*, Qiyang (Colin) Li*, Kevin Frans, Sergey Levine
UC Berkeley

Exploration is Hard

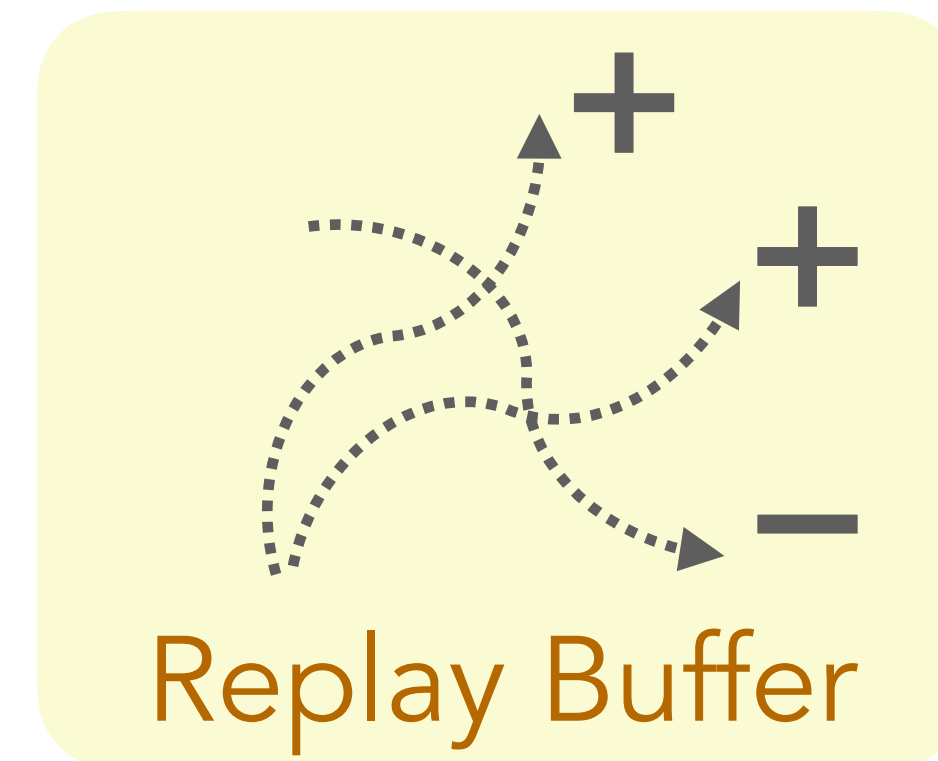
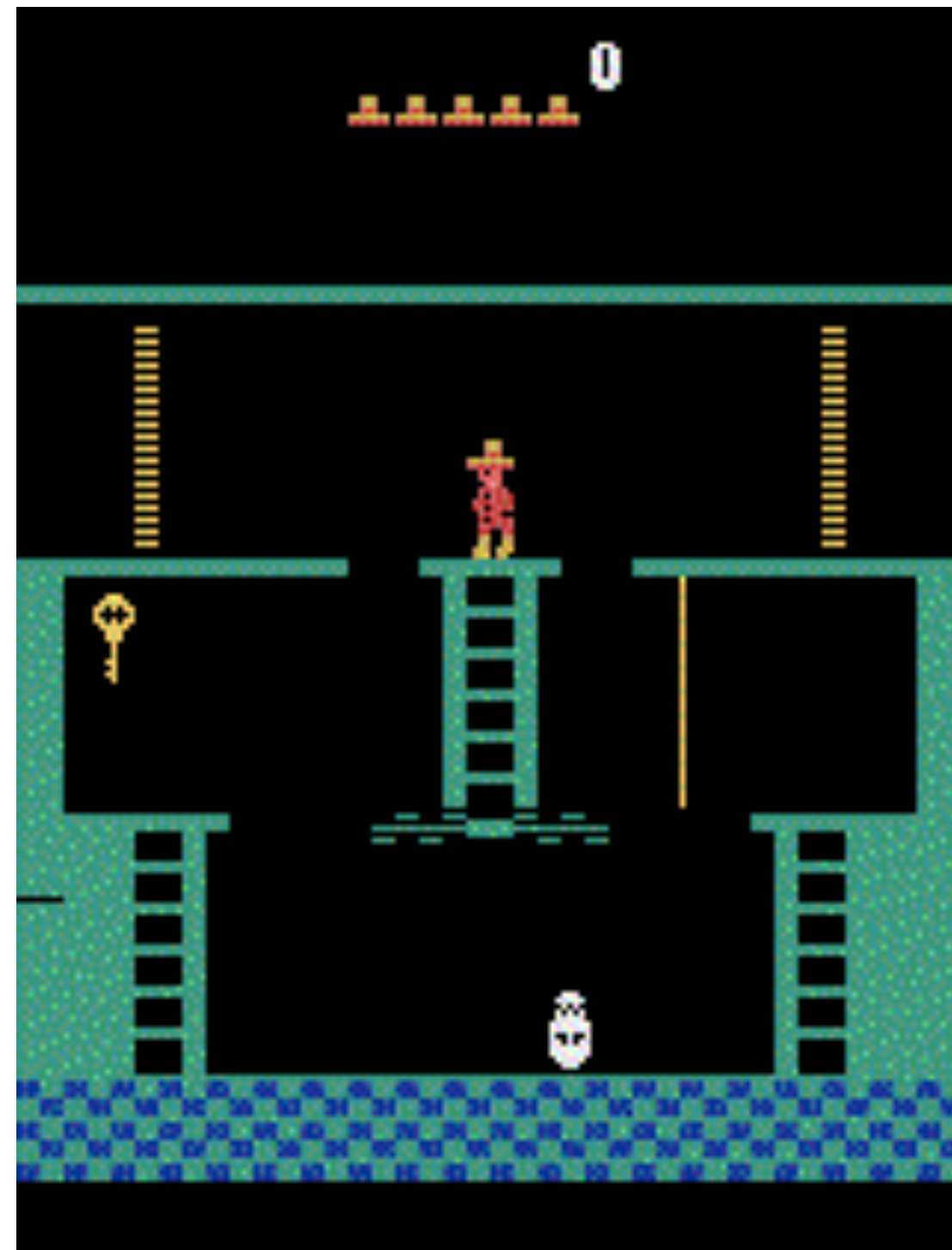


Online

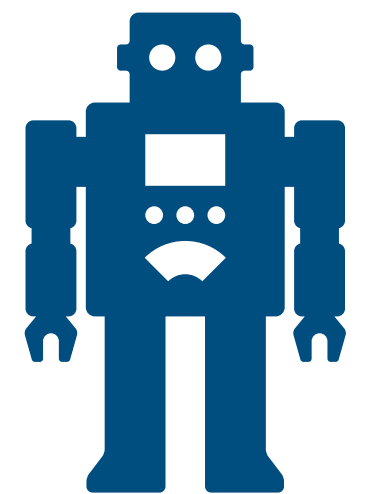


Online Agent

Exploration is Hard

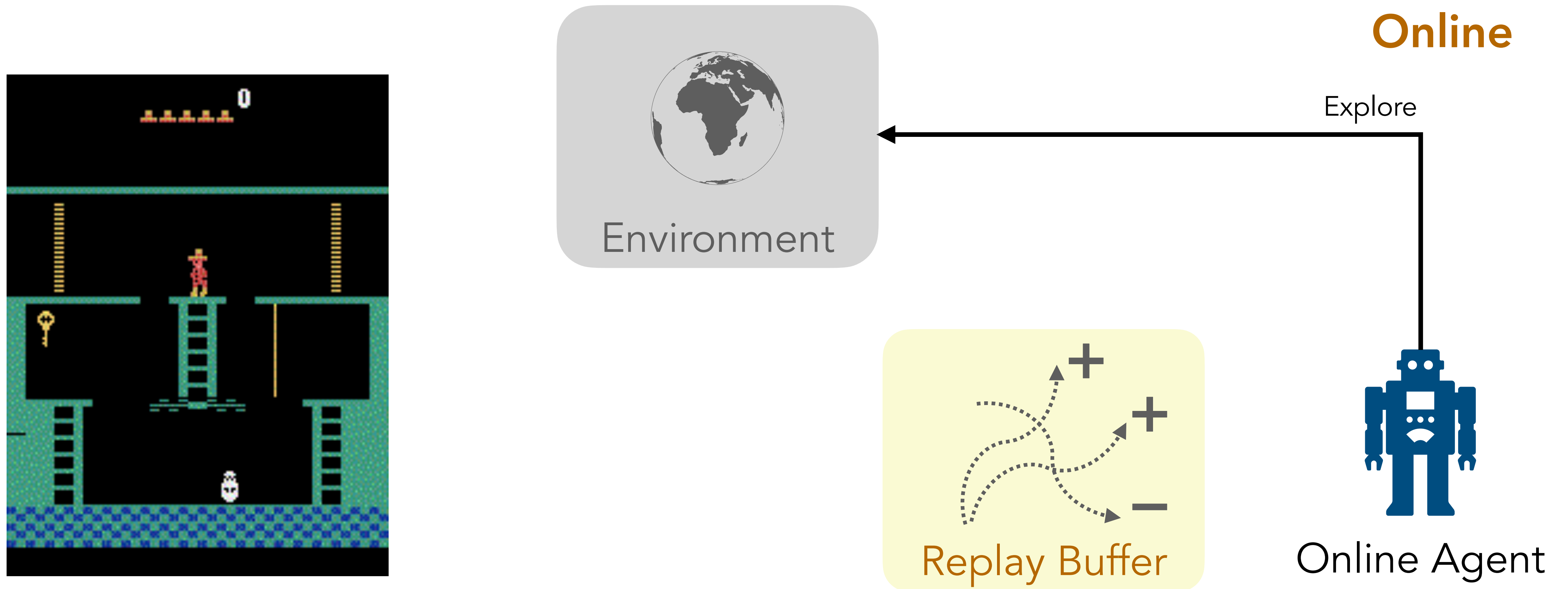


Online

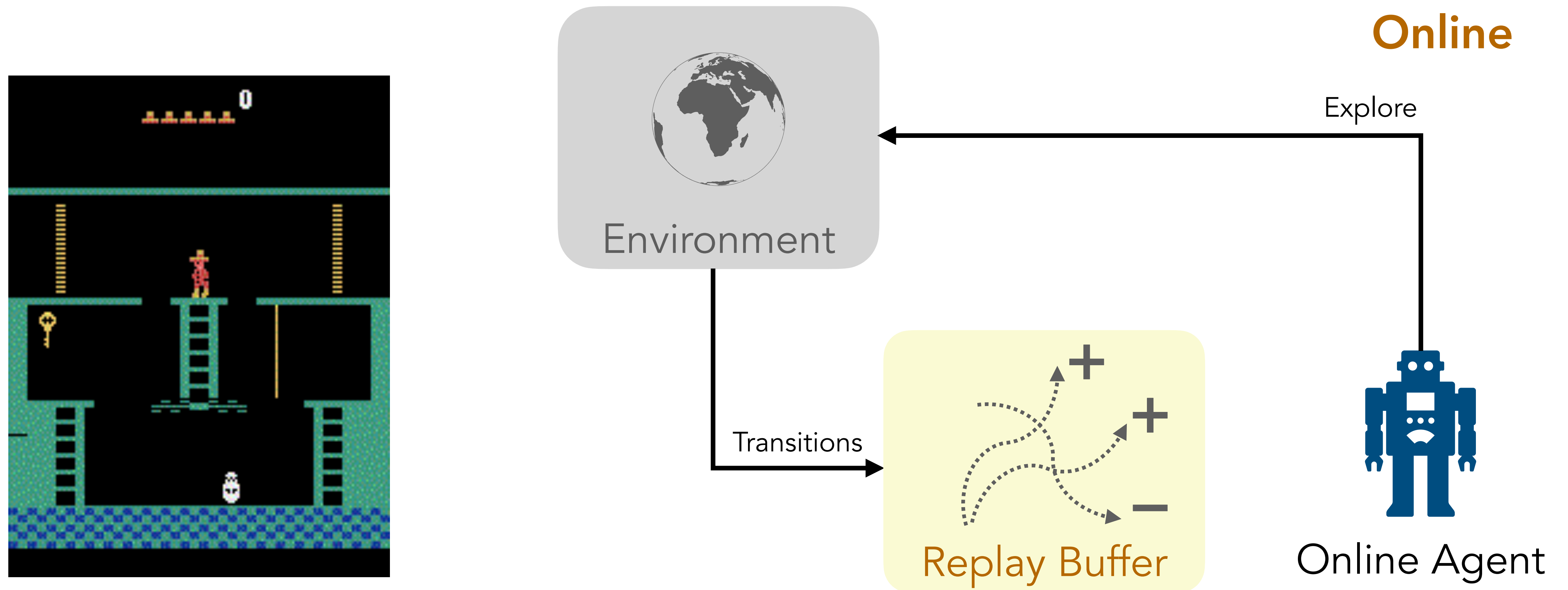


Online Agent

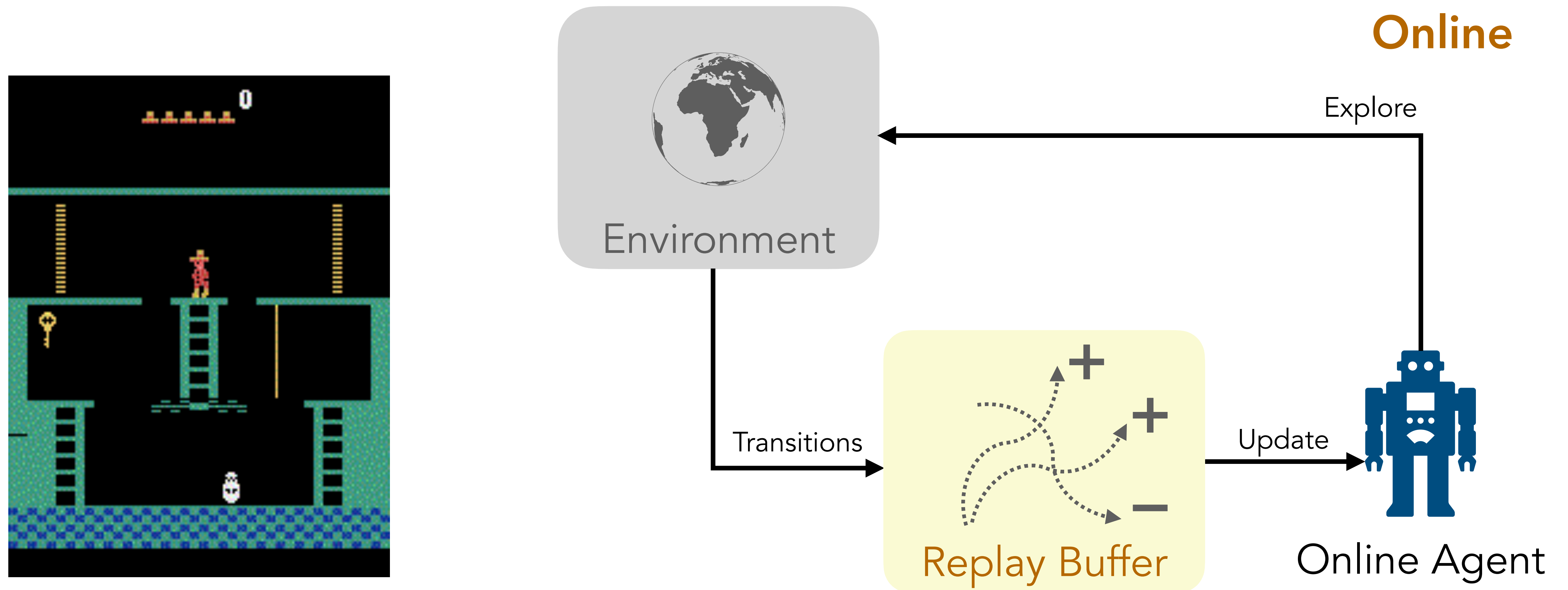
Exploration is Hard



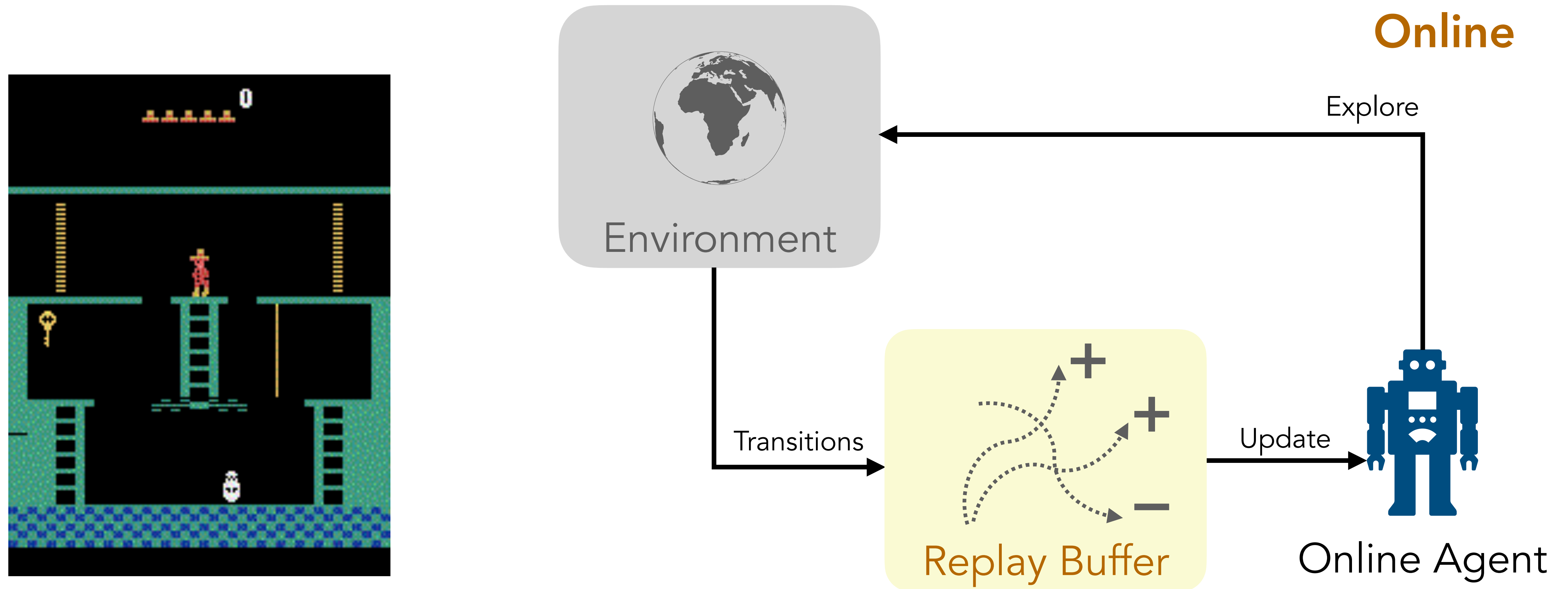
Exploration is Hard



Exploration is Hard



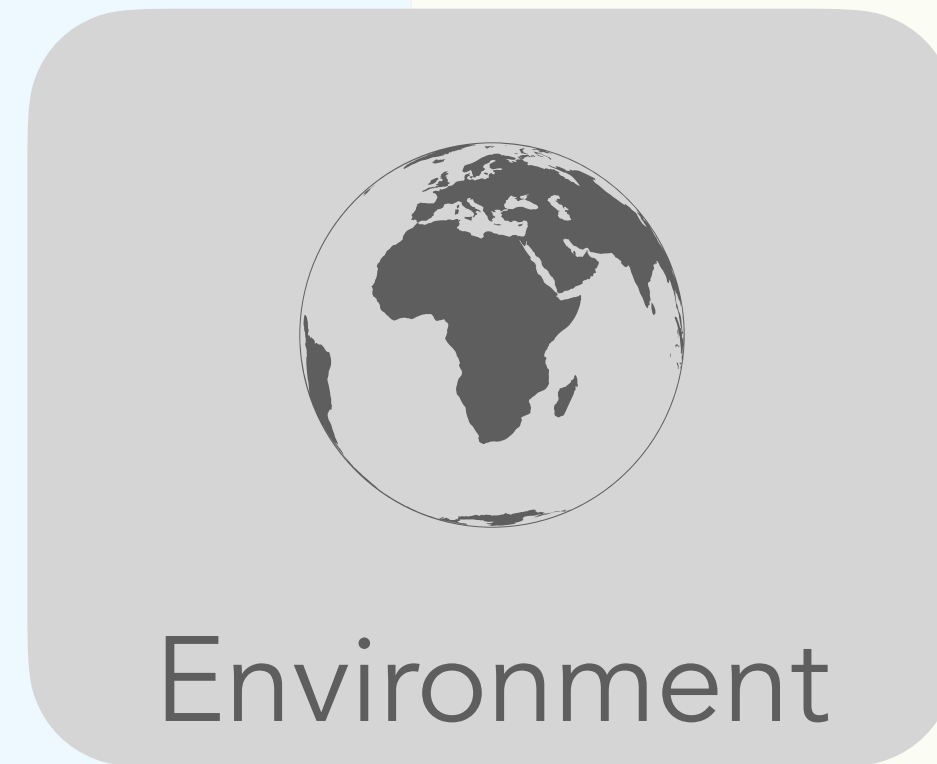
Exploration is Hard



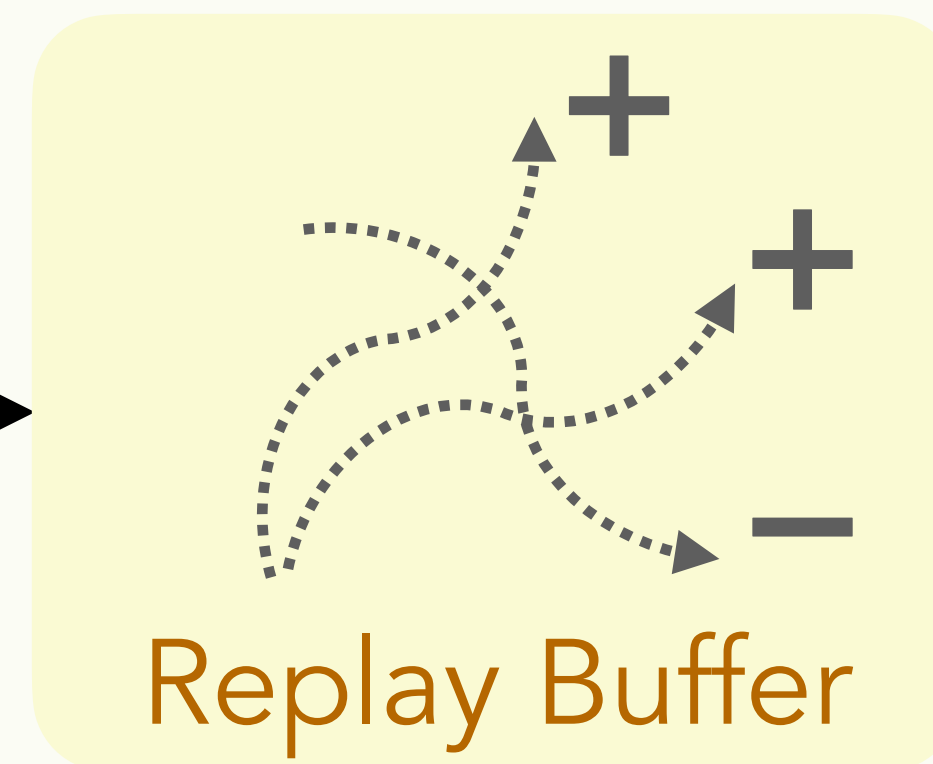
In the worst case, we must visit **every possible state** in the environment

Data-driven Exploration

Offline

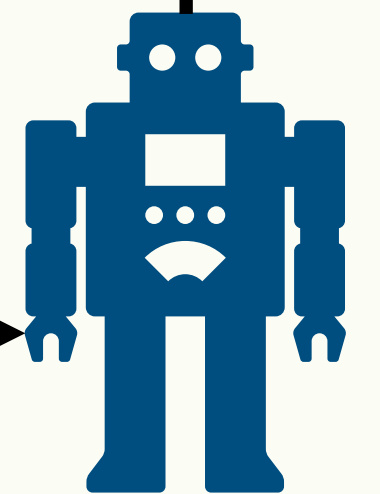


Labeled,
Task-Specific



Online

Explore

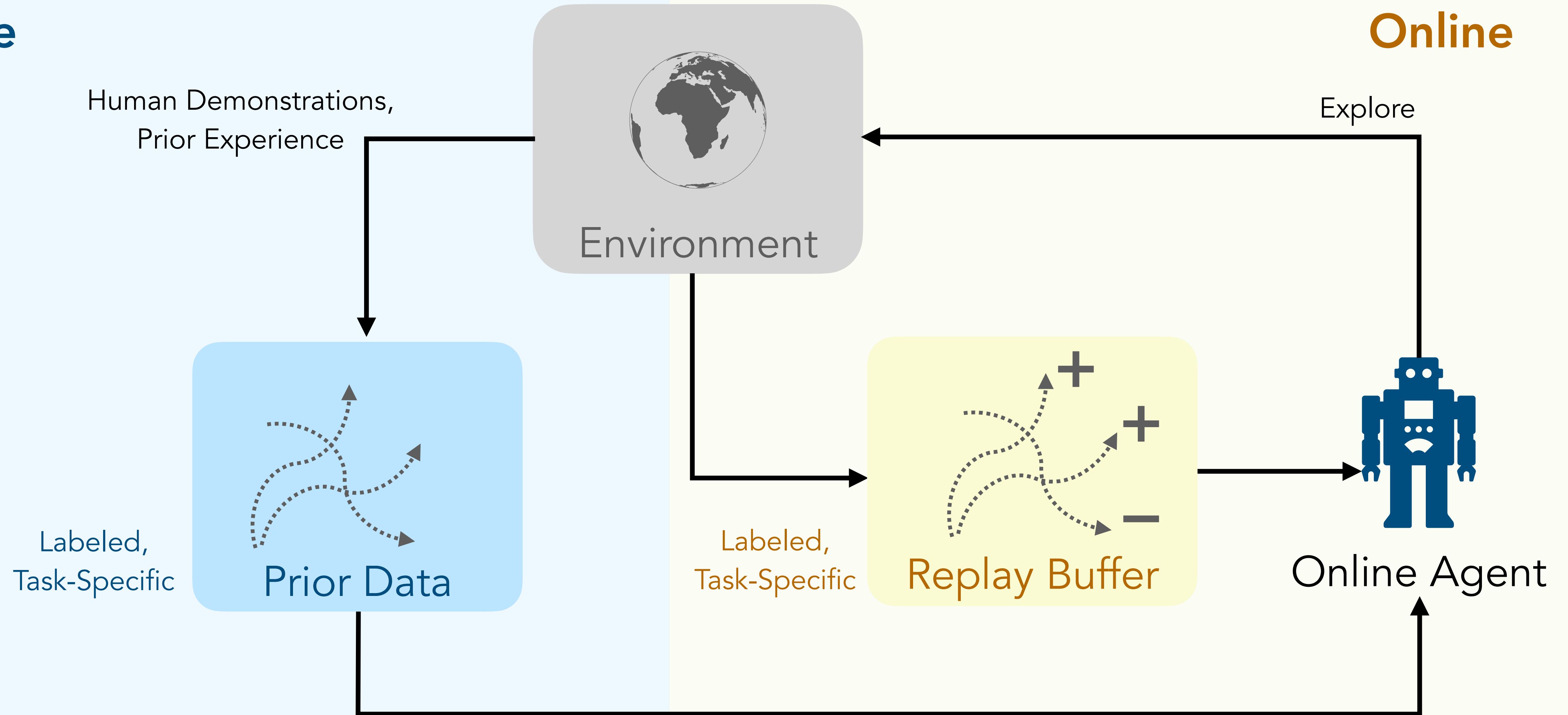


Online Agent

Data-driven Exploration

Offline

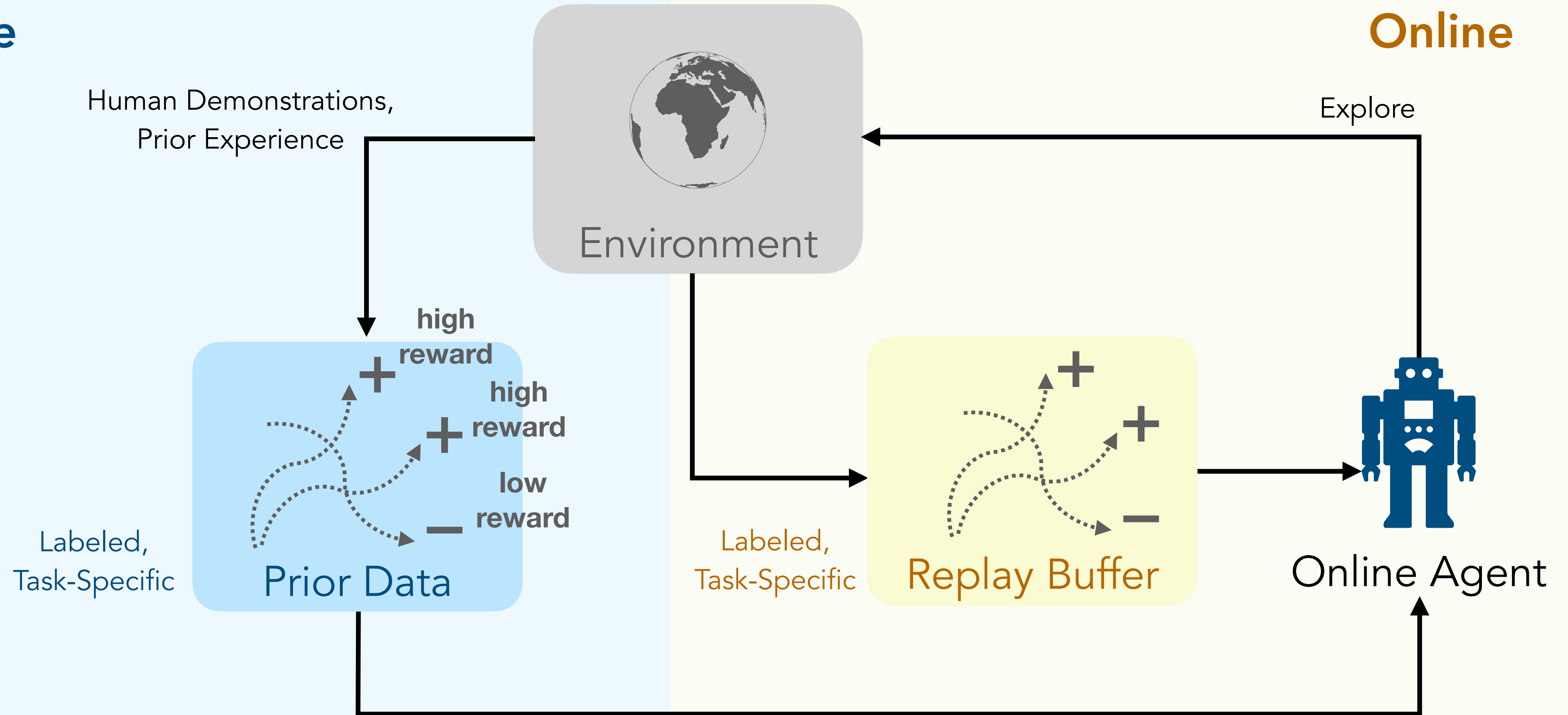
Online



Data-driven Exploration

Offline

Online

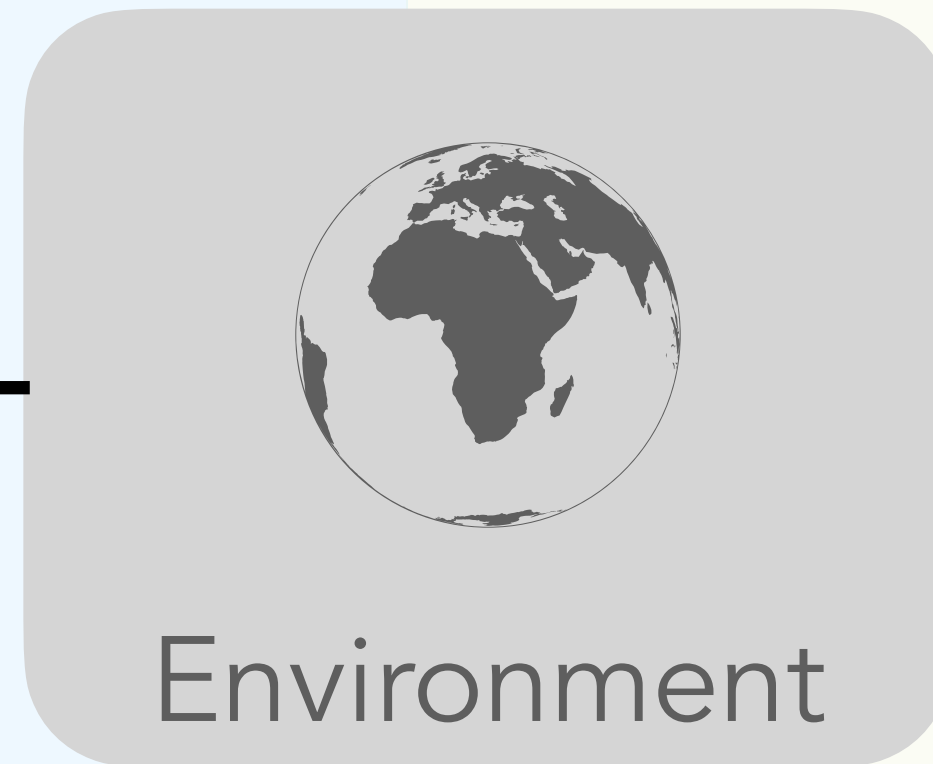


Data-driven Exploration

Offline

Online

Human Demonstrations,
Prior Experience



Explore

Expensive to collect, and doesn't
scale to large, multitask datasets!

Labeled,
Task-Specific

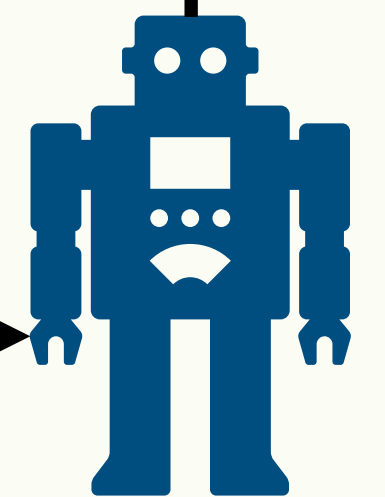
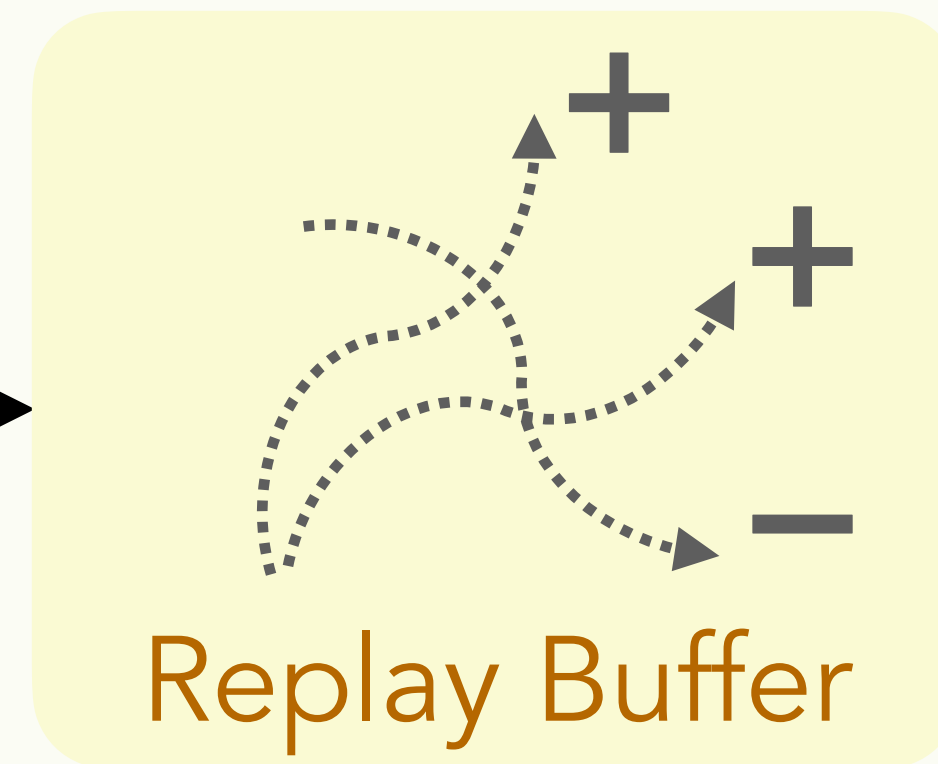
Prior Data

Labeled,
Task-Specific

Replay Buffer

Online Agent

limited to data collected
for the specific task at hand

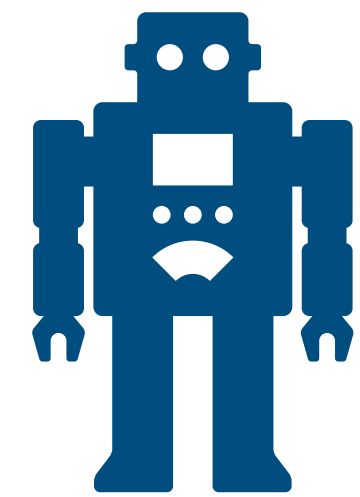


Data-driven Exploration

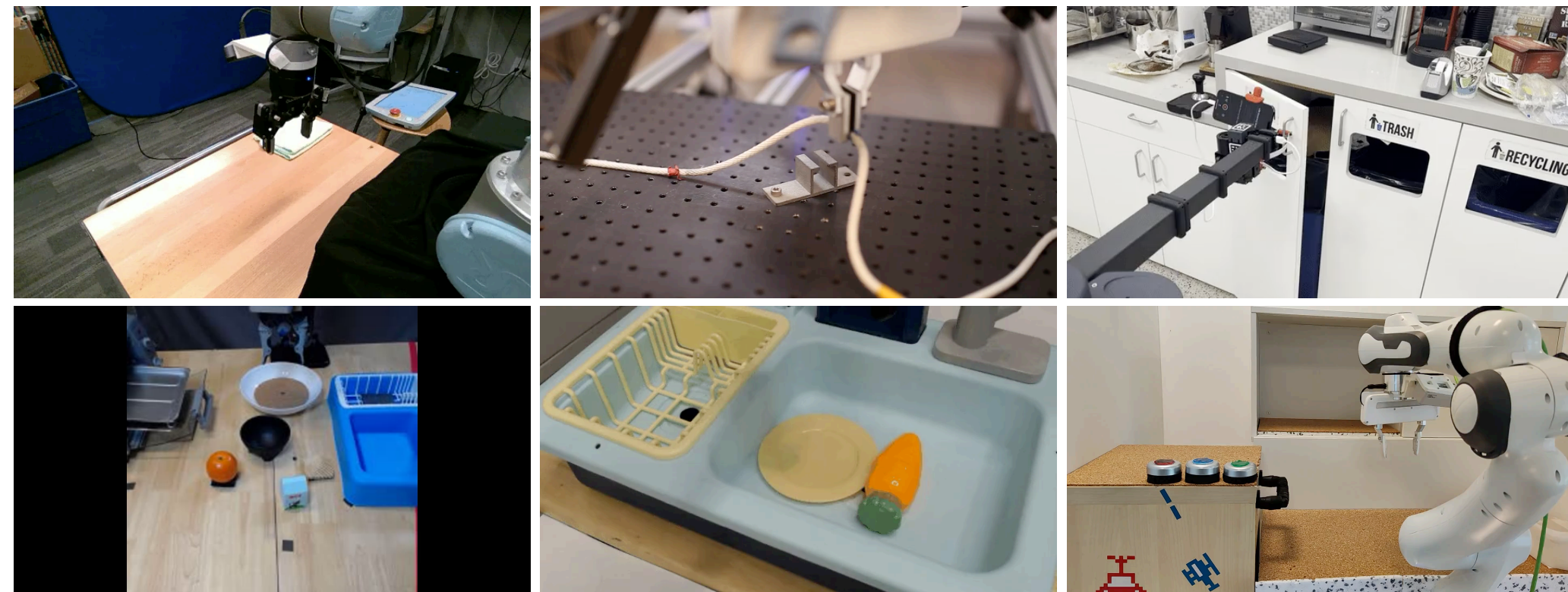
Unlabeled,
Task-Agnostic



???



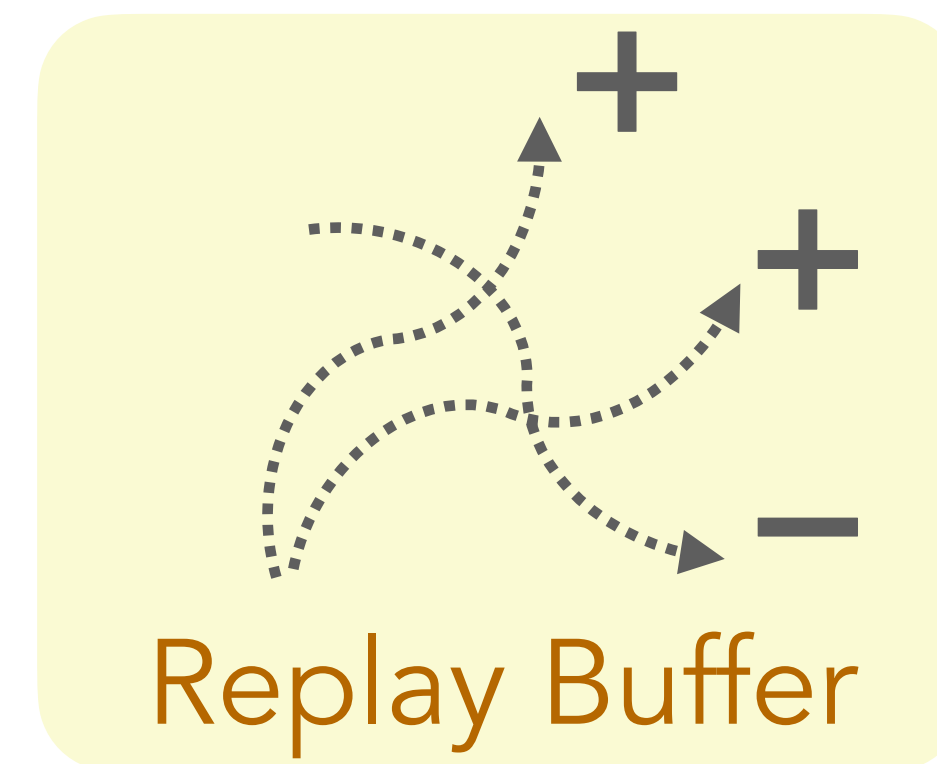
Online Policy



can leverage data collected for a range of tasks,
scales without expensive supervision

Labeled,
Task-Specific

Off-policy RL

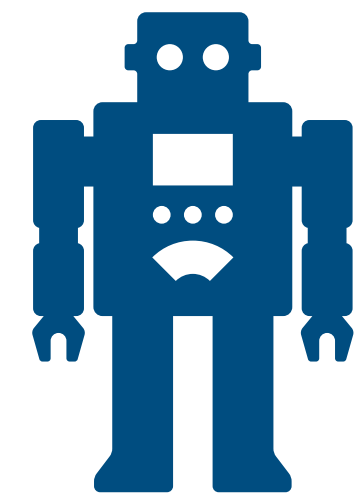


Data-driven Exploration

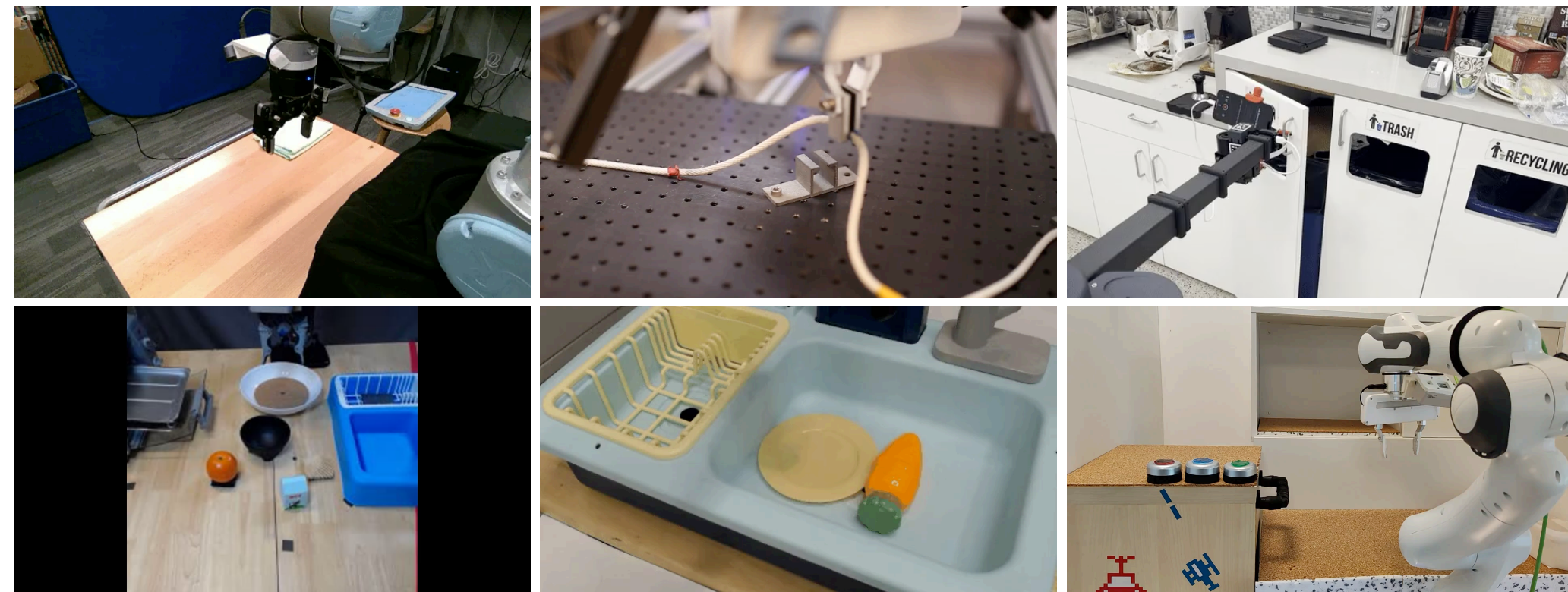
Unlabeled,
Task-Agnostic



???



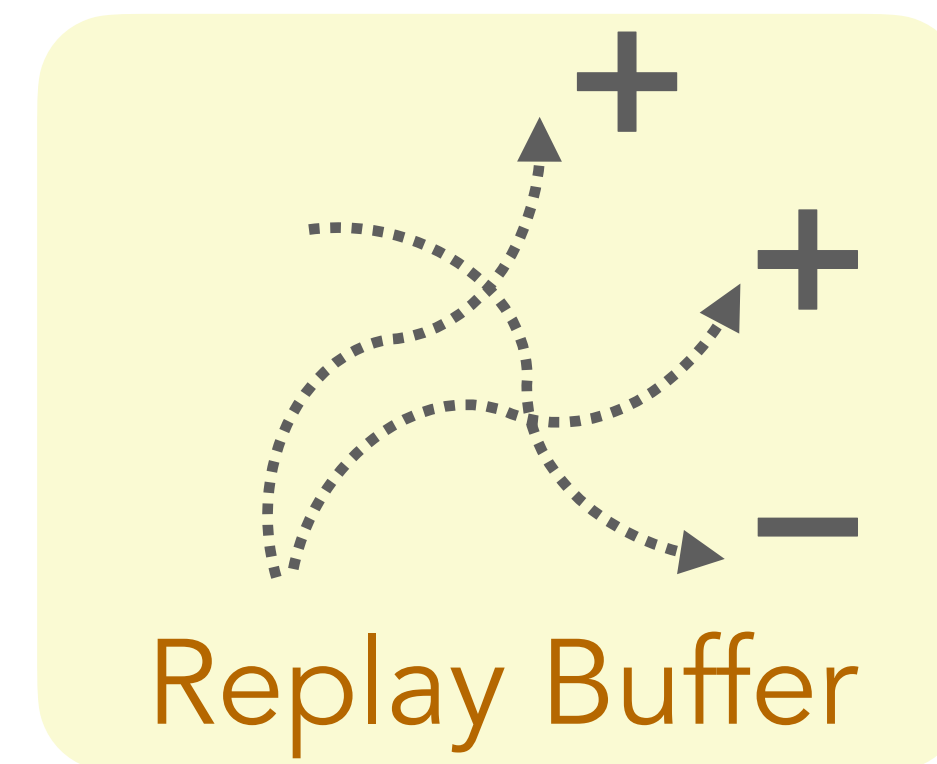
Online Policy



can leverage data collected for a range of tasks,
scales without expensive supervision

Labeled,
Task-Specific

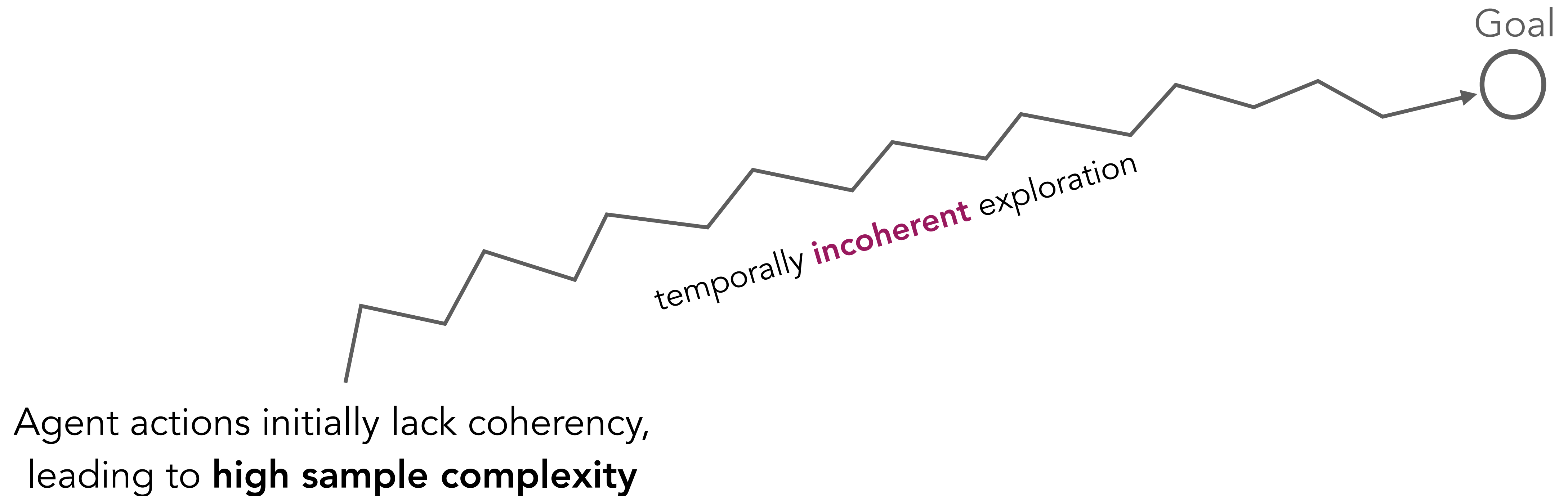
Off-policy RL



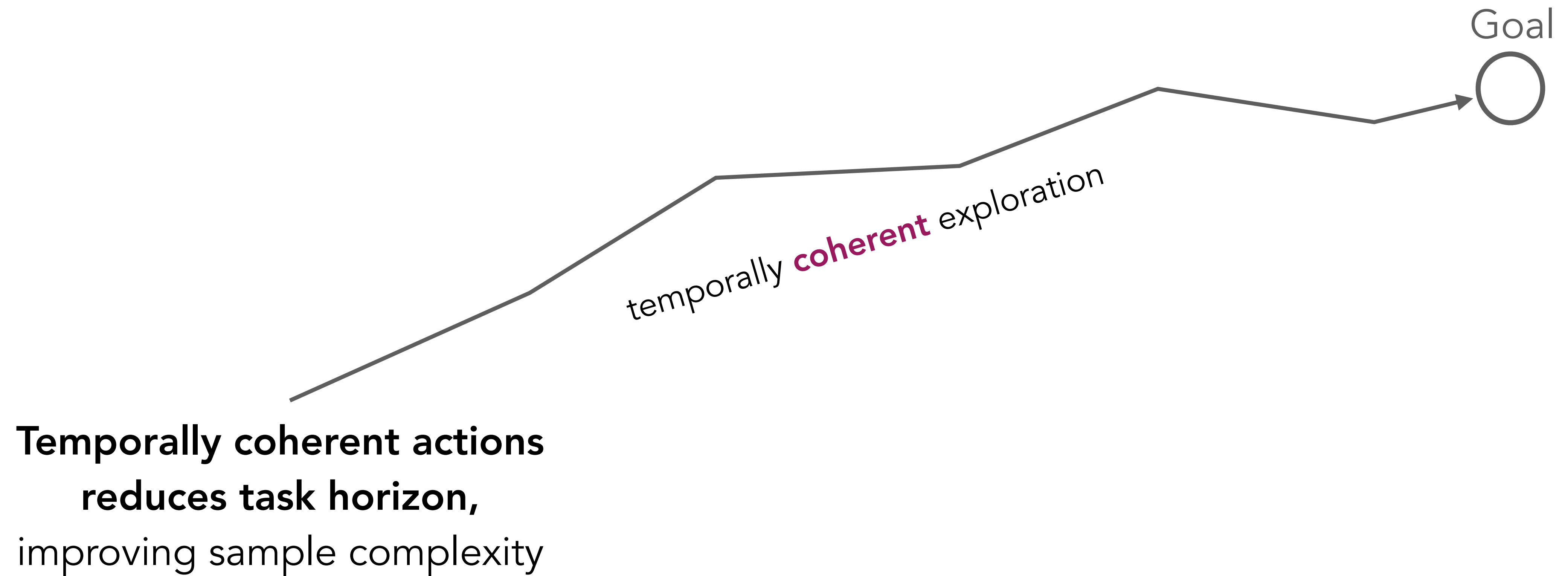
Data-driven Exploration



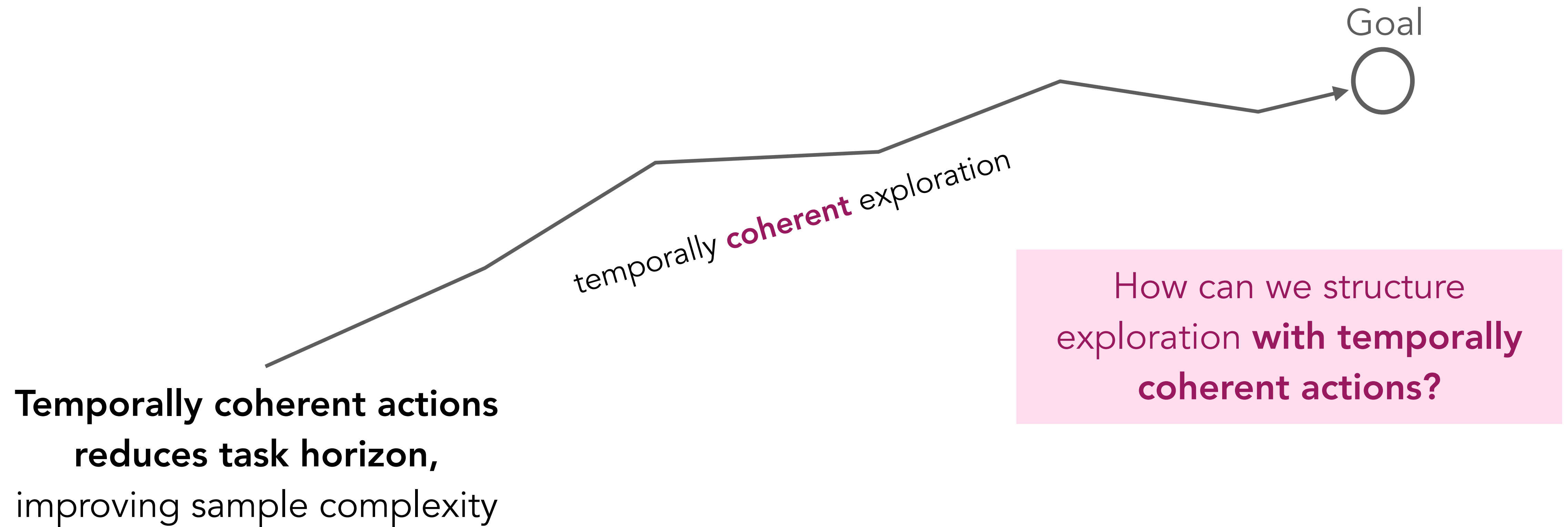
Challenge #1: Unstructured Exploration



Challenge #1: Unstructured Exploration



Challenge #1: Unstructured Exploration

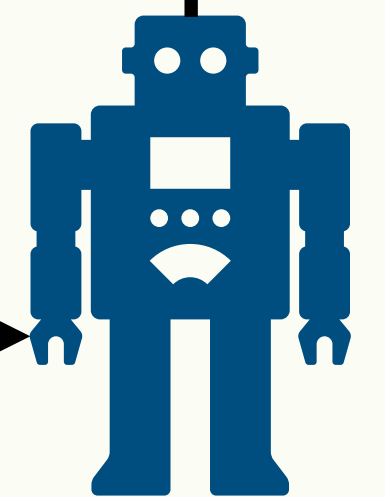
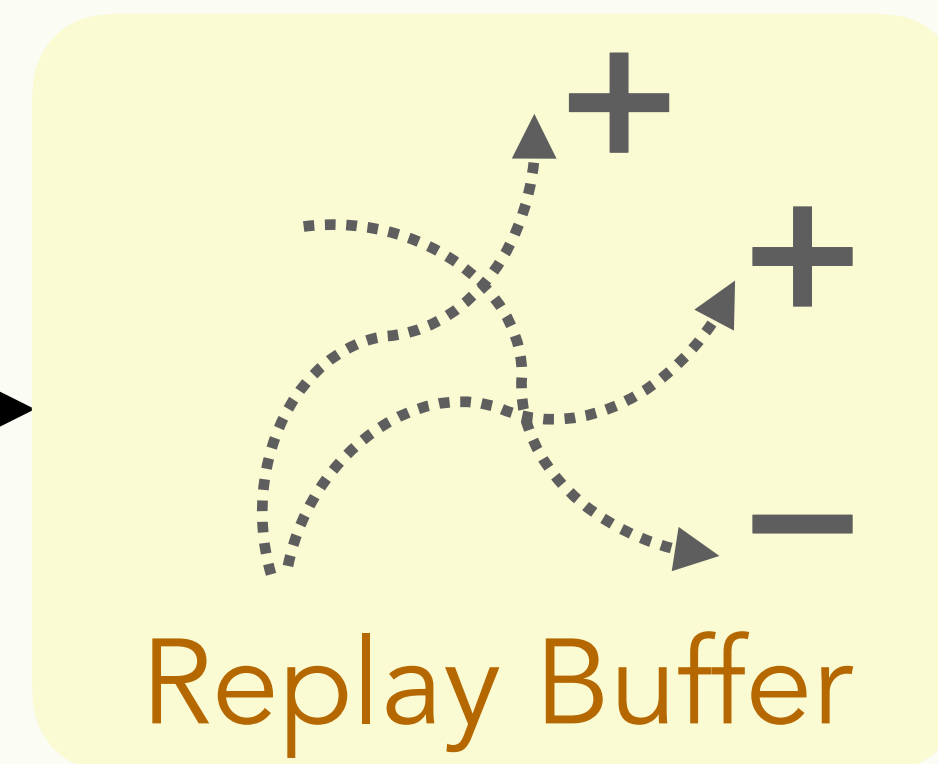
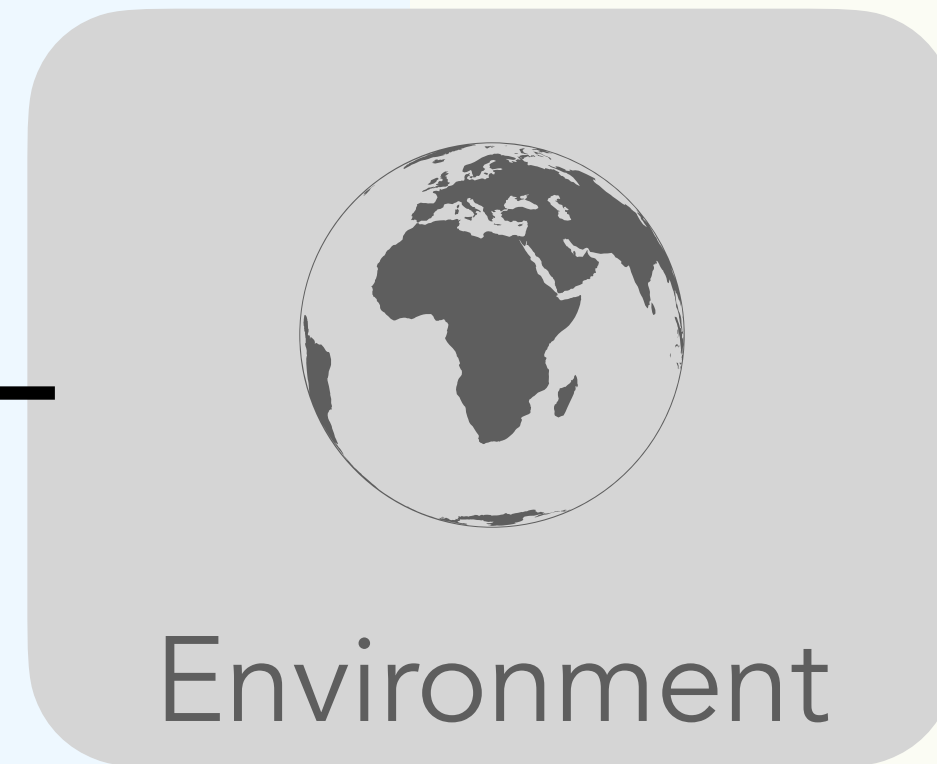
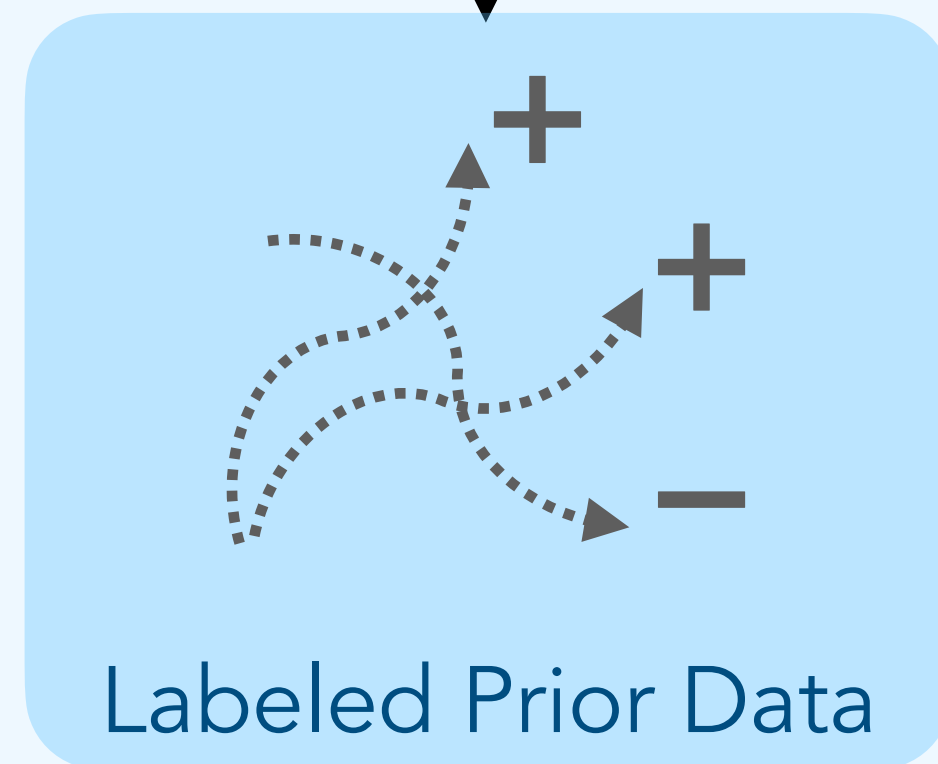


Challenge #2: No access to task-specific reward labels

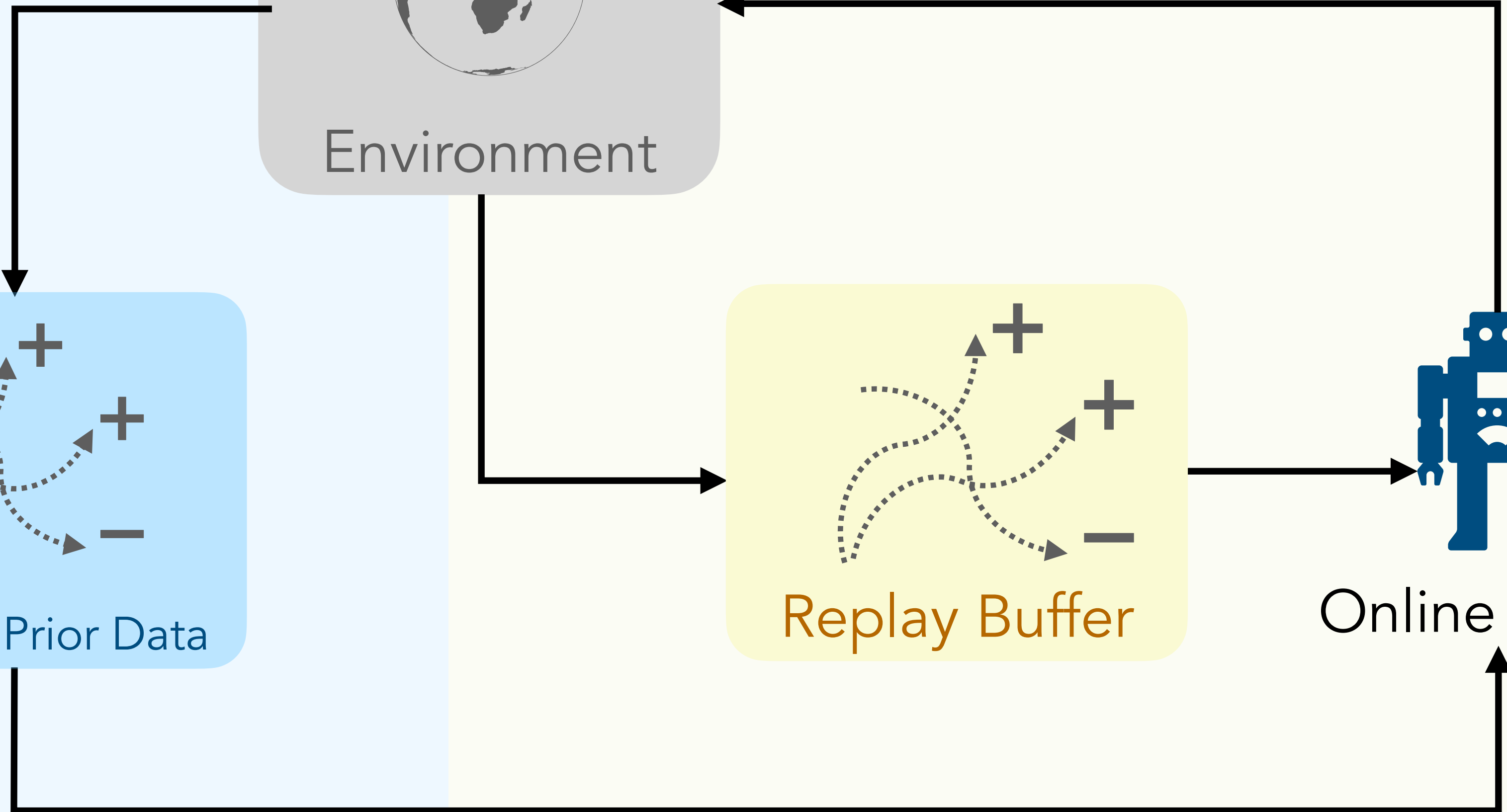
Offline

Online

Labeled prior data can
accelerate online learning



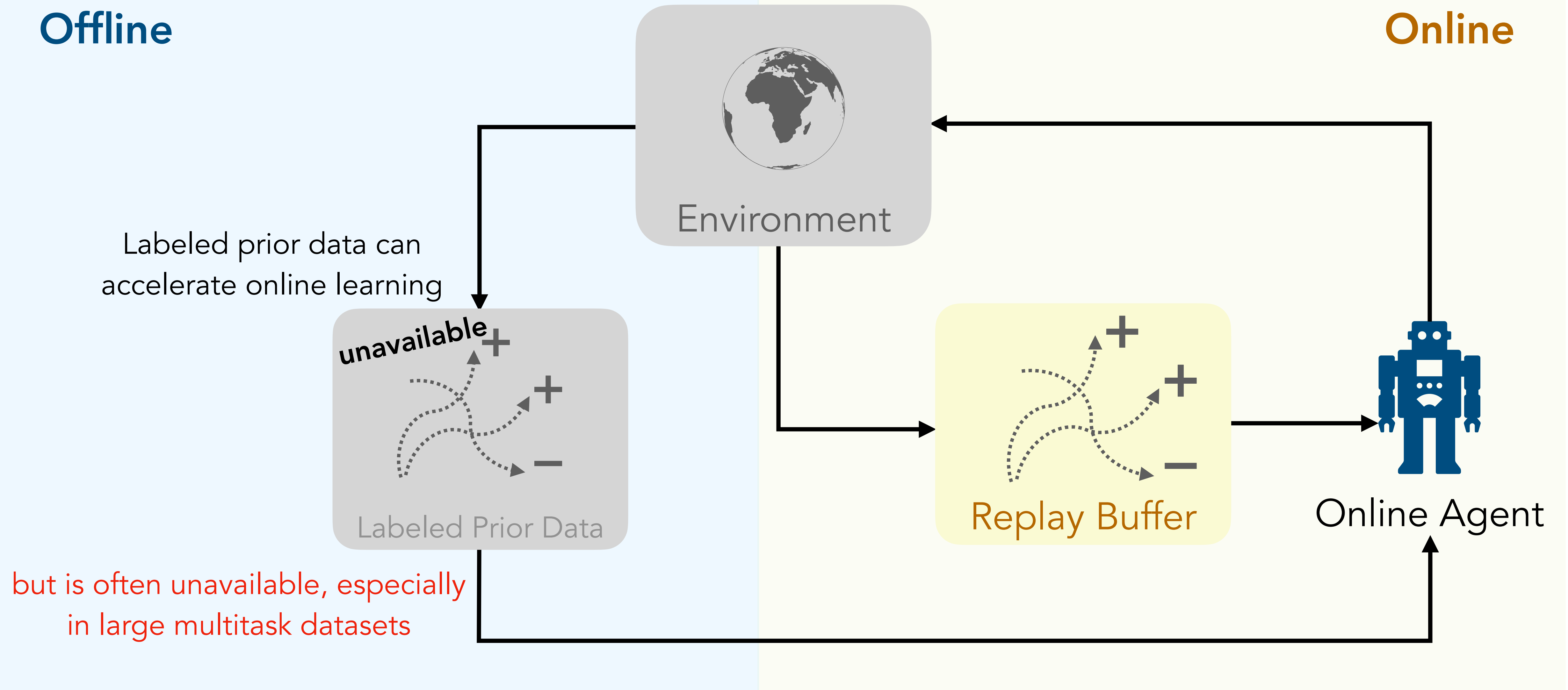
Online Agent



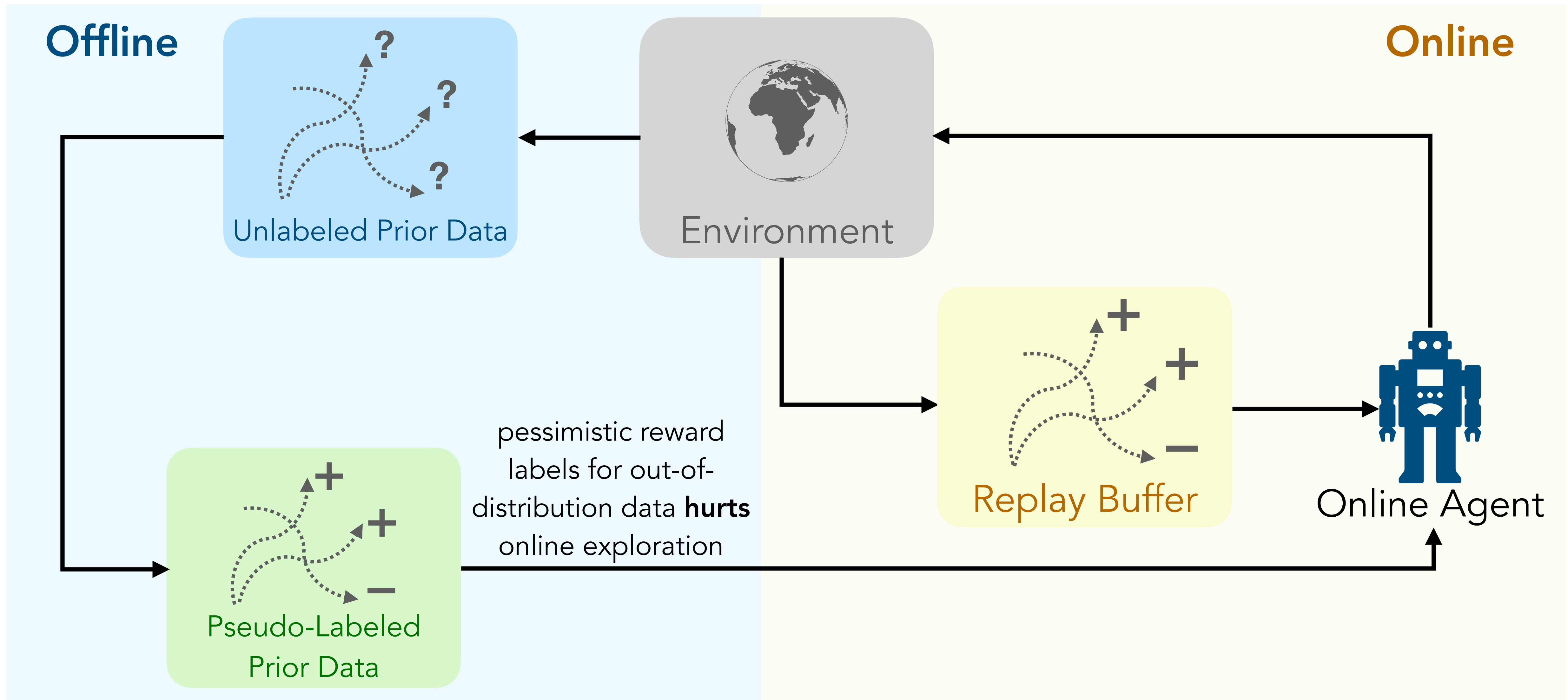
Challenge #2: No access to task-specific reward labels

Offline

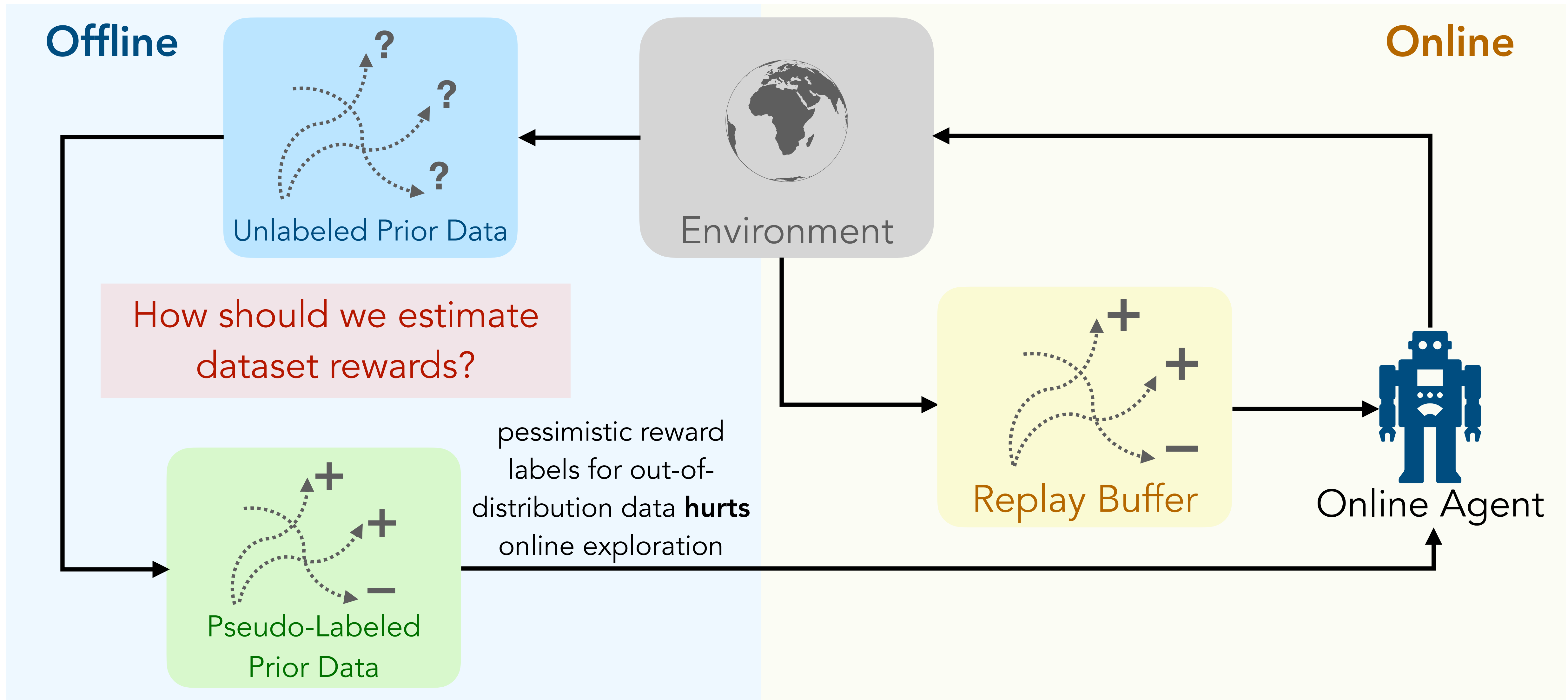
Online



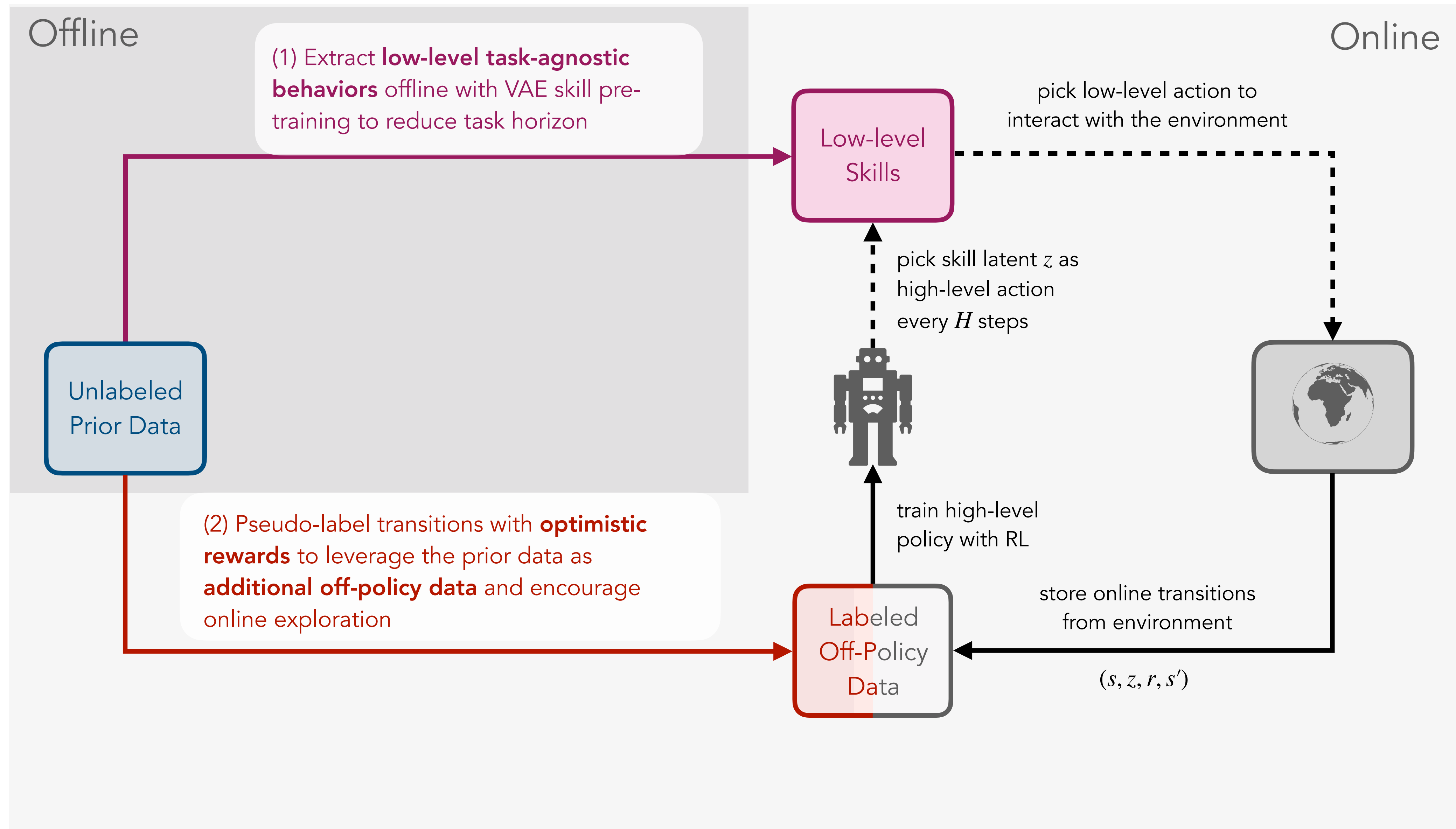
Challenge #2: No access to task-specific reward labels



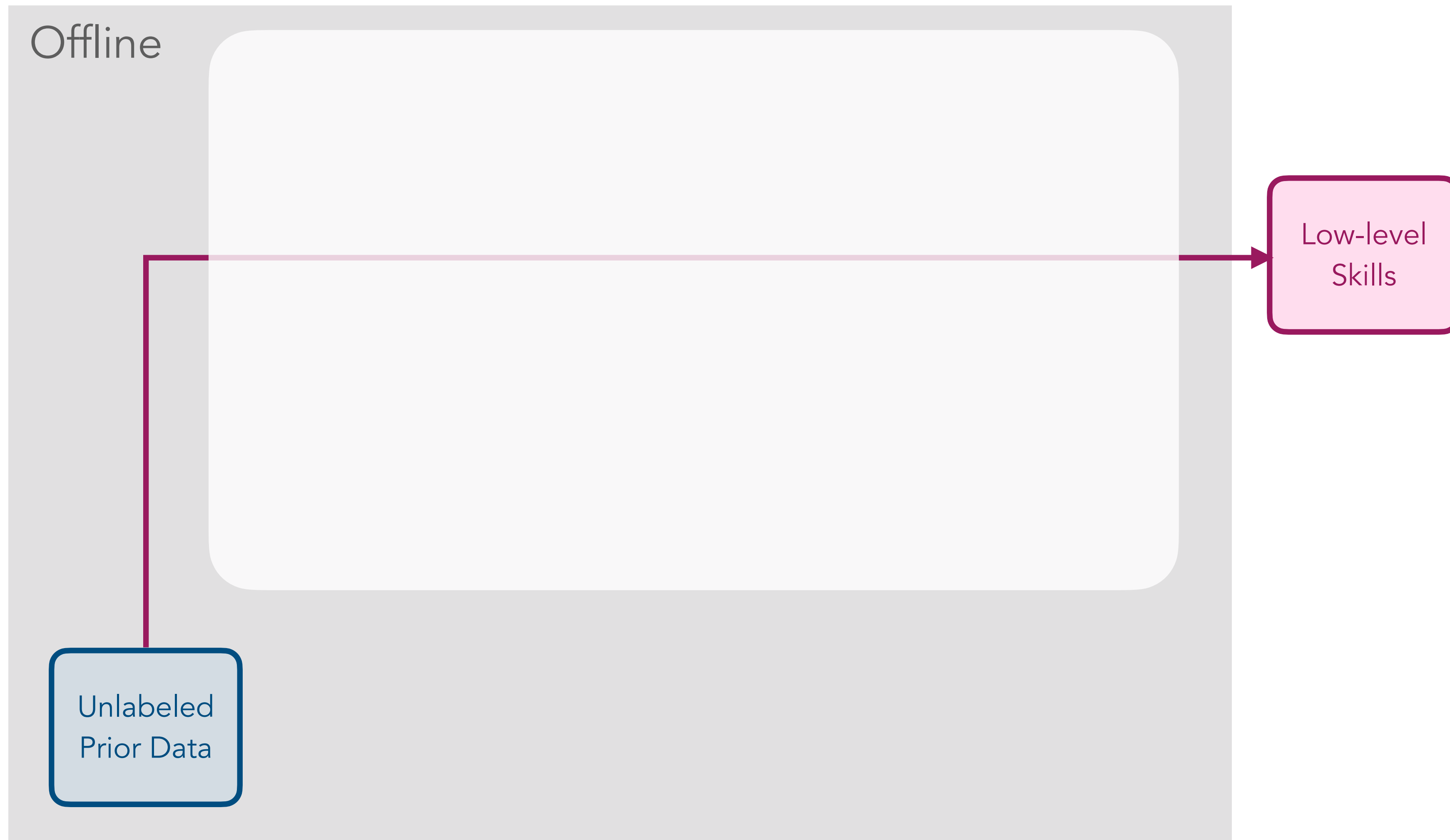
Challenge #2: No access to task-specific reward labels



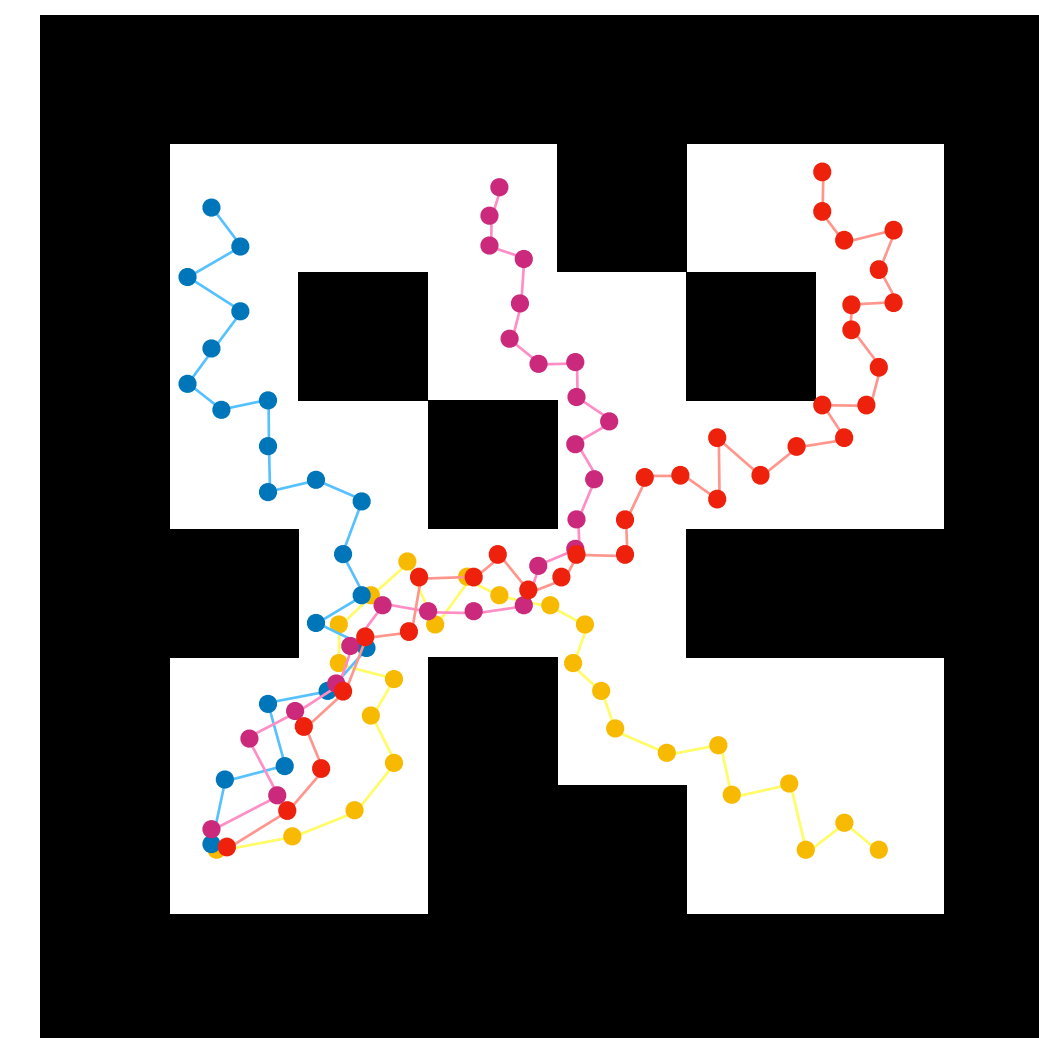
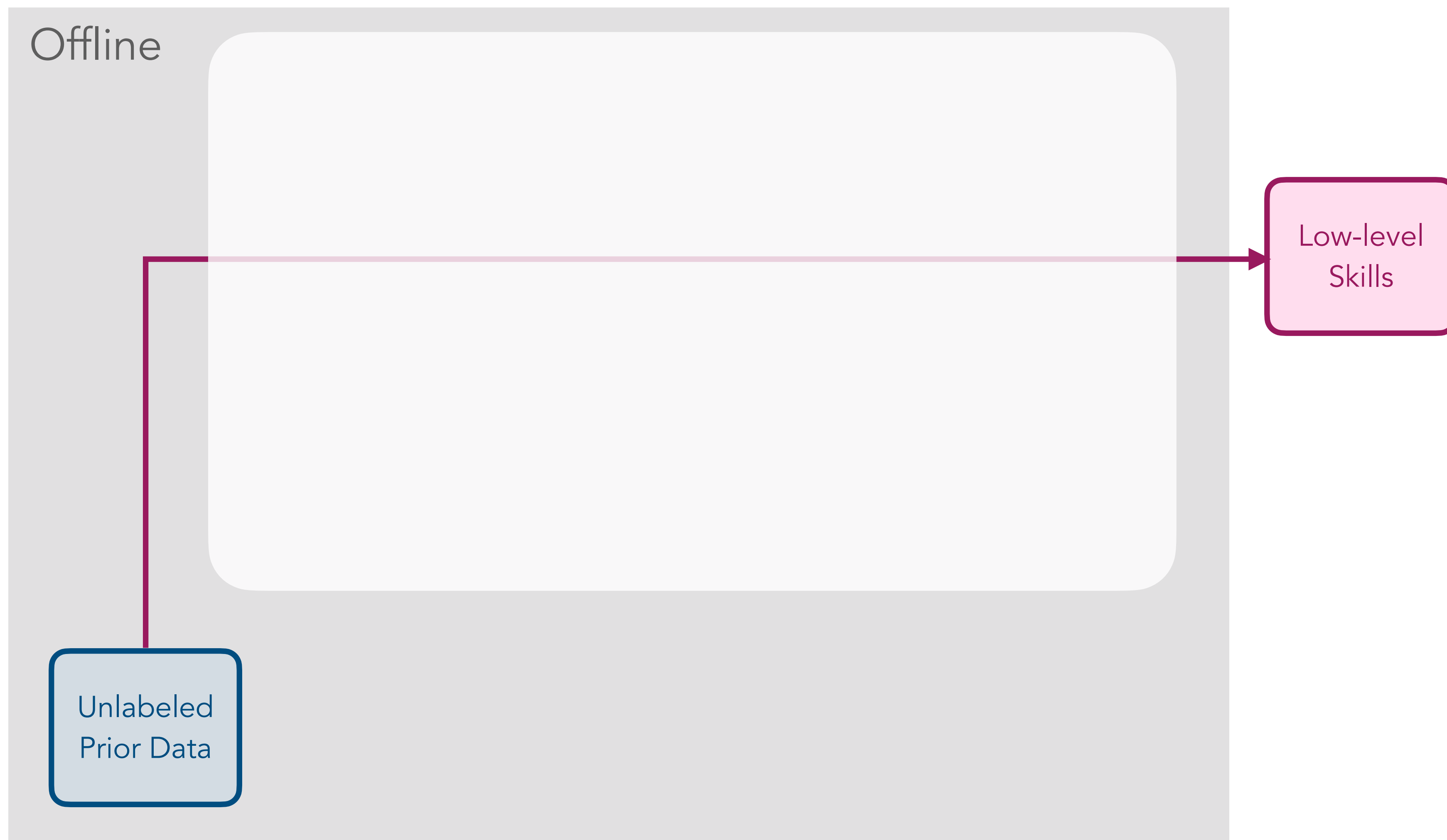
SUPE: Skills from Unlabeled Prior Data for Exploration



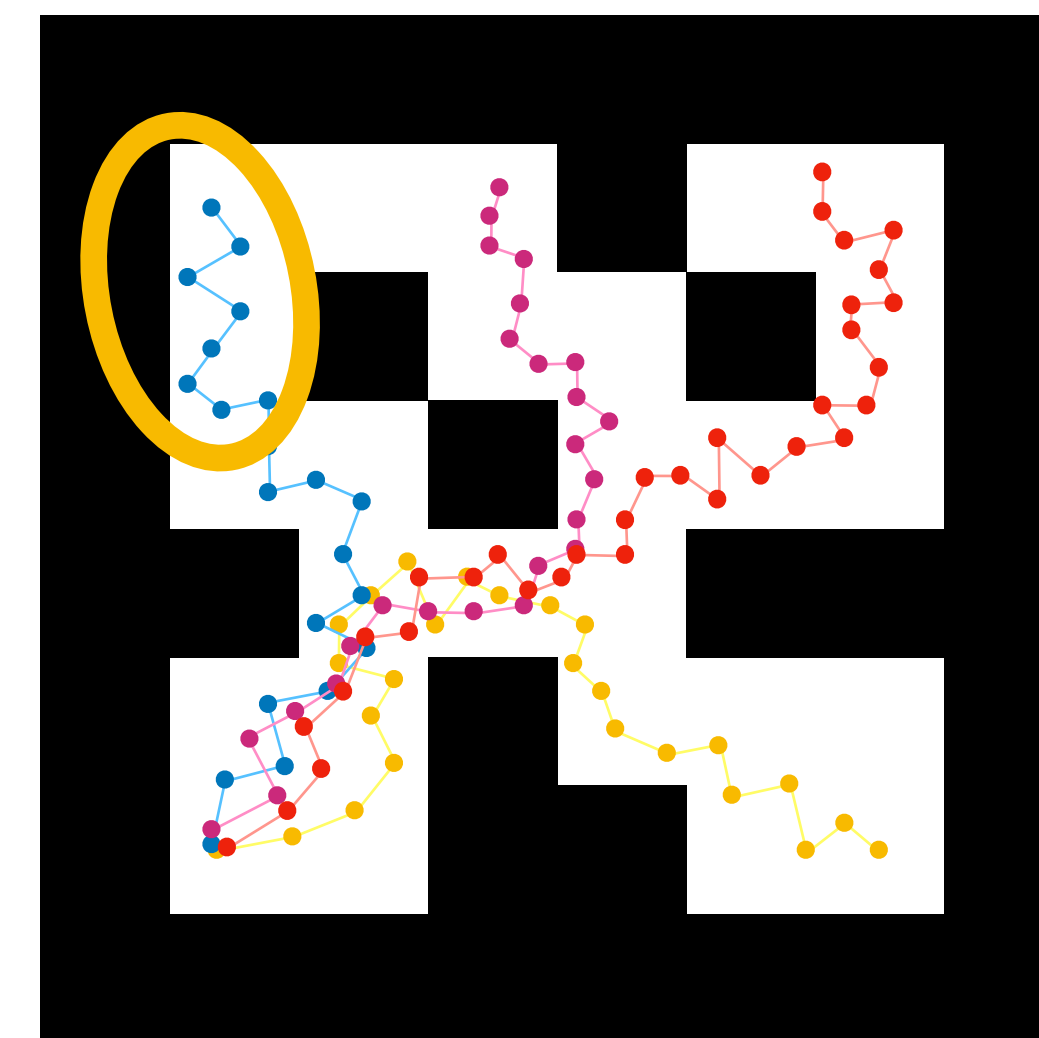
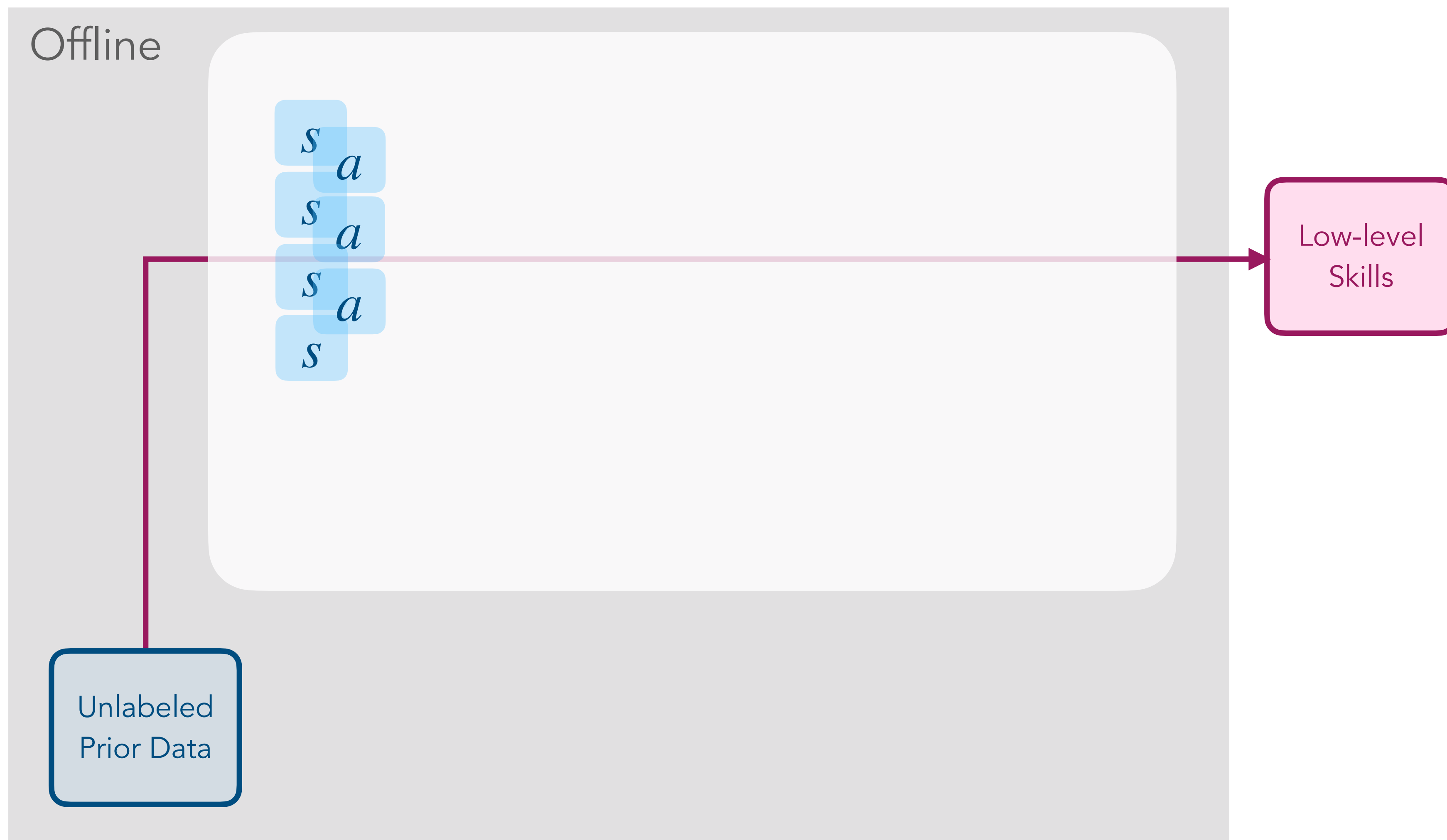
(1) Extract **low-level task-agnostic behaviors** offline
with VAE skill pre-training to reduce task horizon



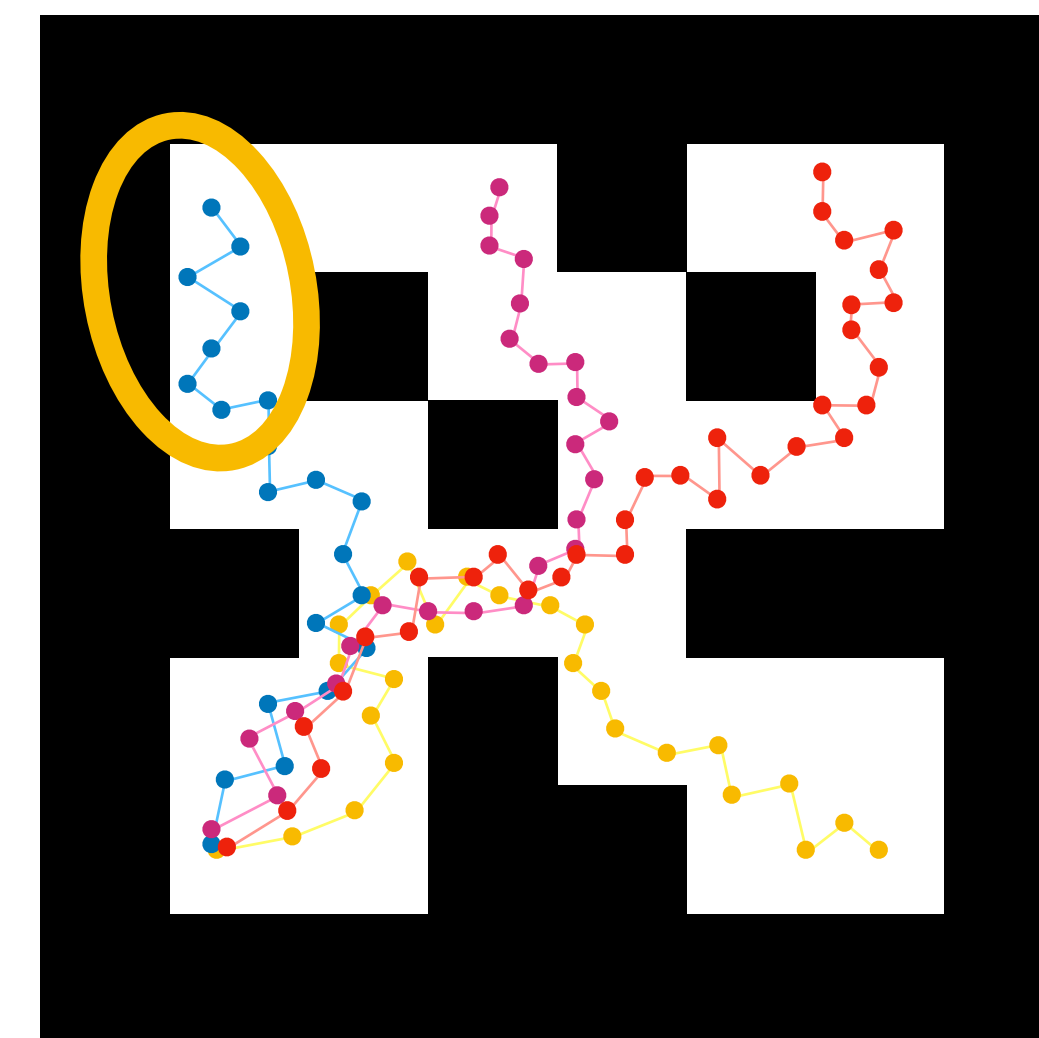
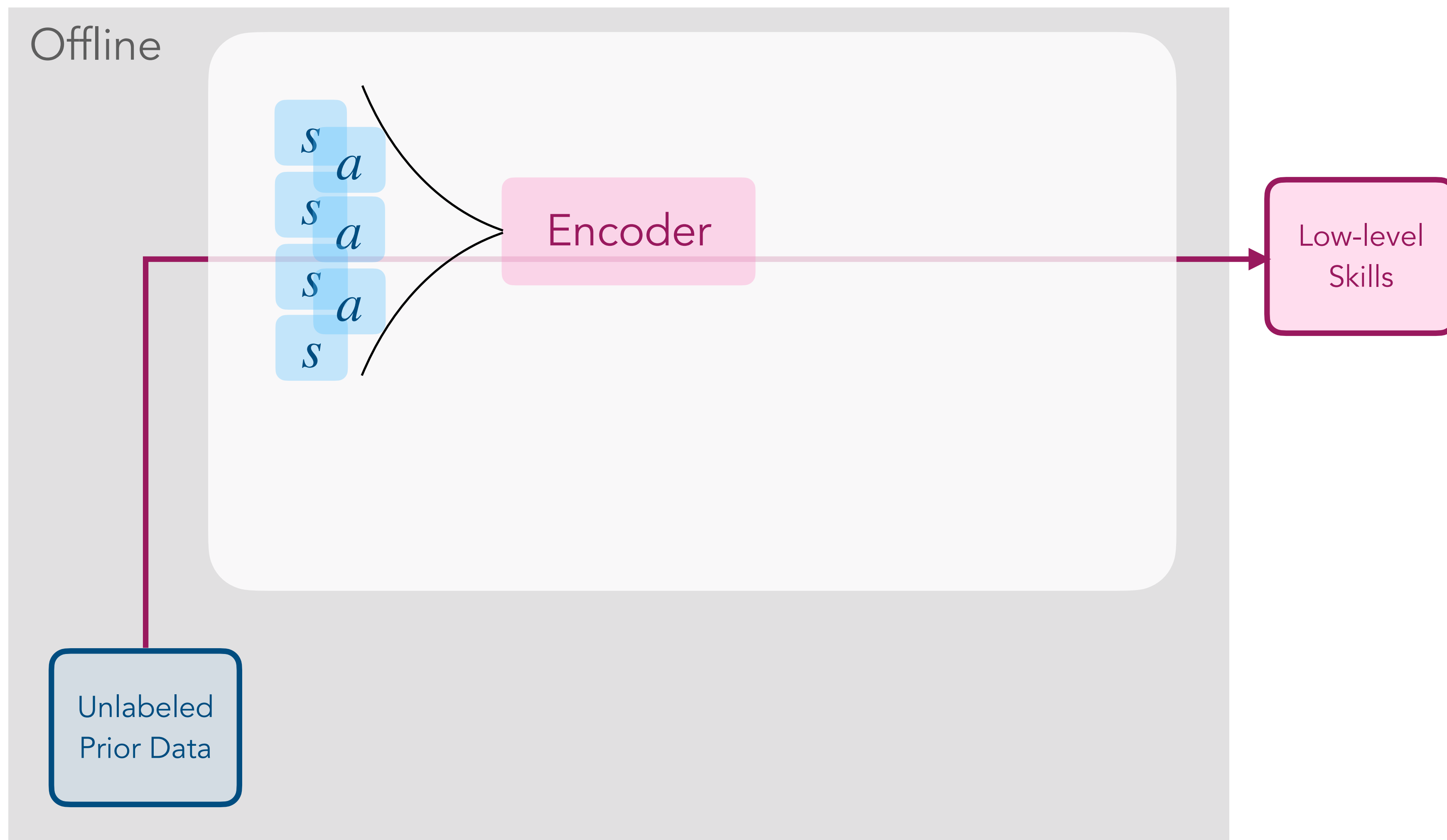
(1) Extract **low-level task-agnostic behaviors** offline with VAE skill pre-training to reduce task horizon



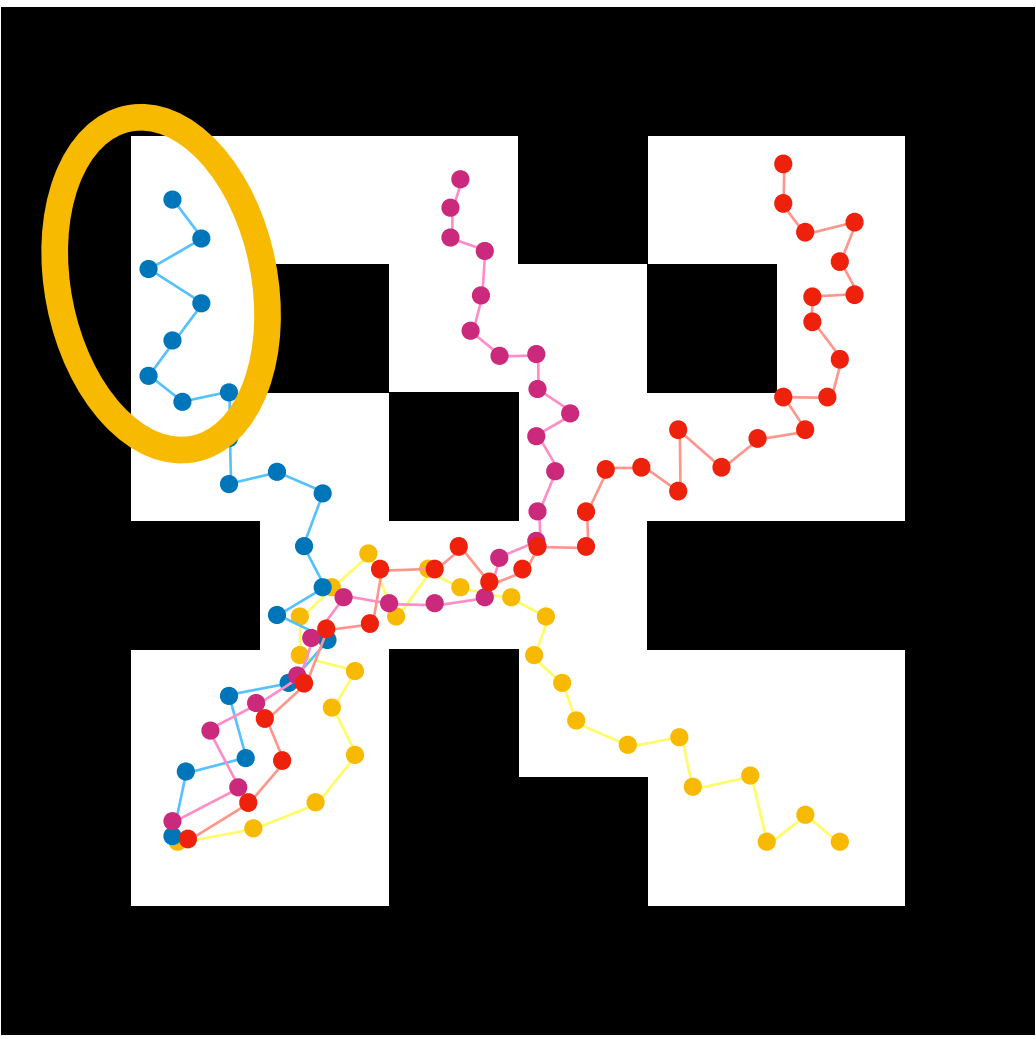
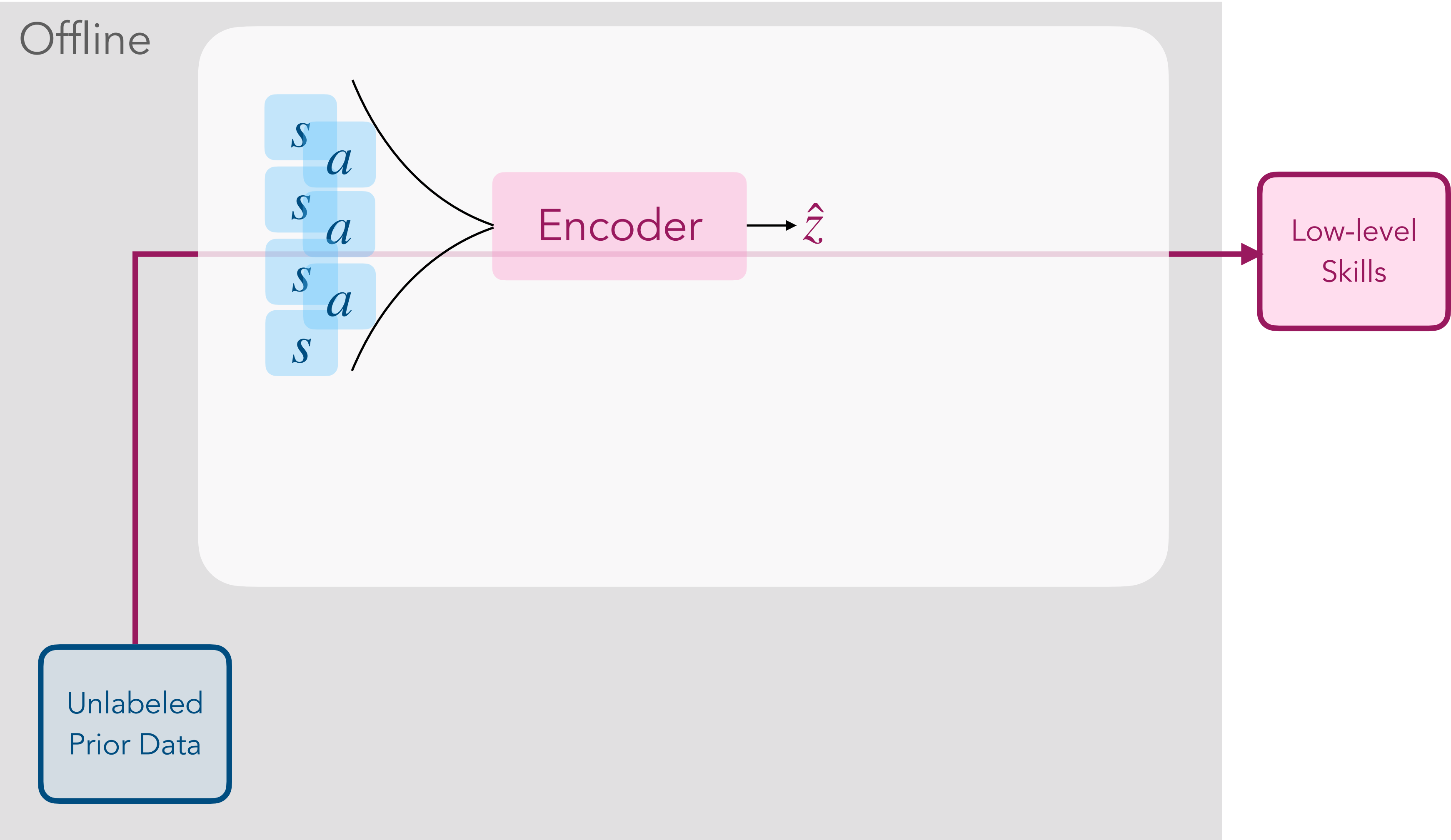
(1) Extract **low-level task-agnostic behaviors** offline with VAE skill pre-training to reduce task horizon



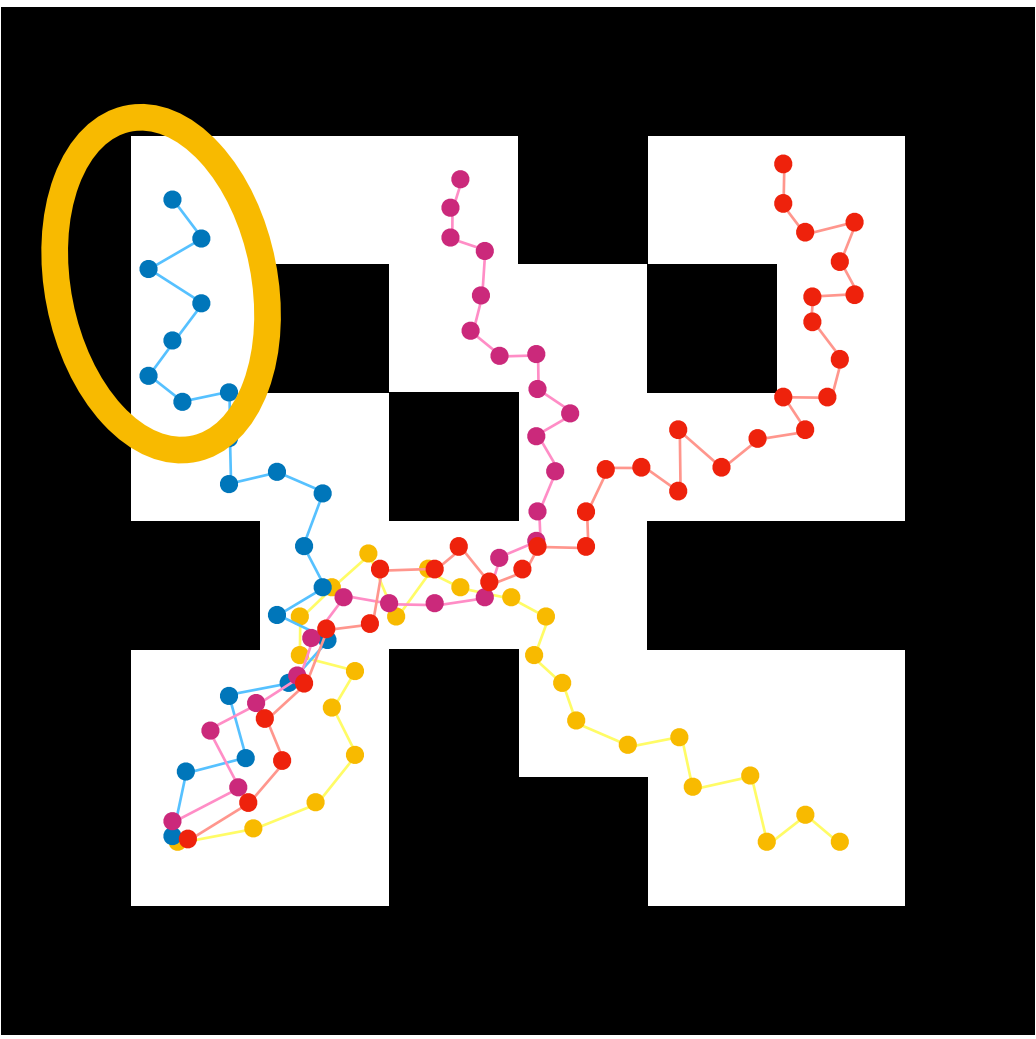
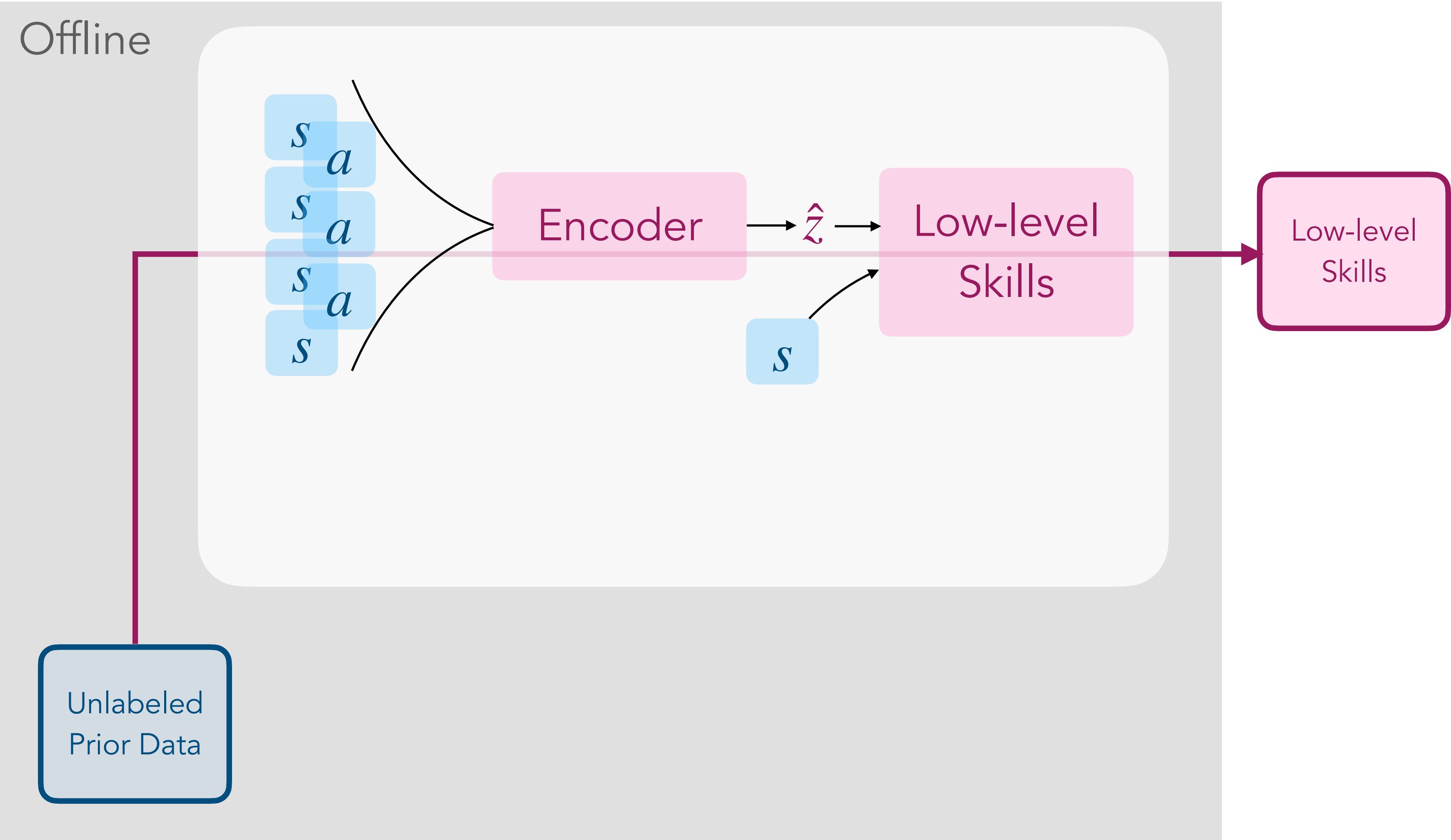
(1) Extract **low-level task-agnostic behaviors** offline with VAE skill pre-training to reduce task horizon



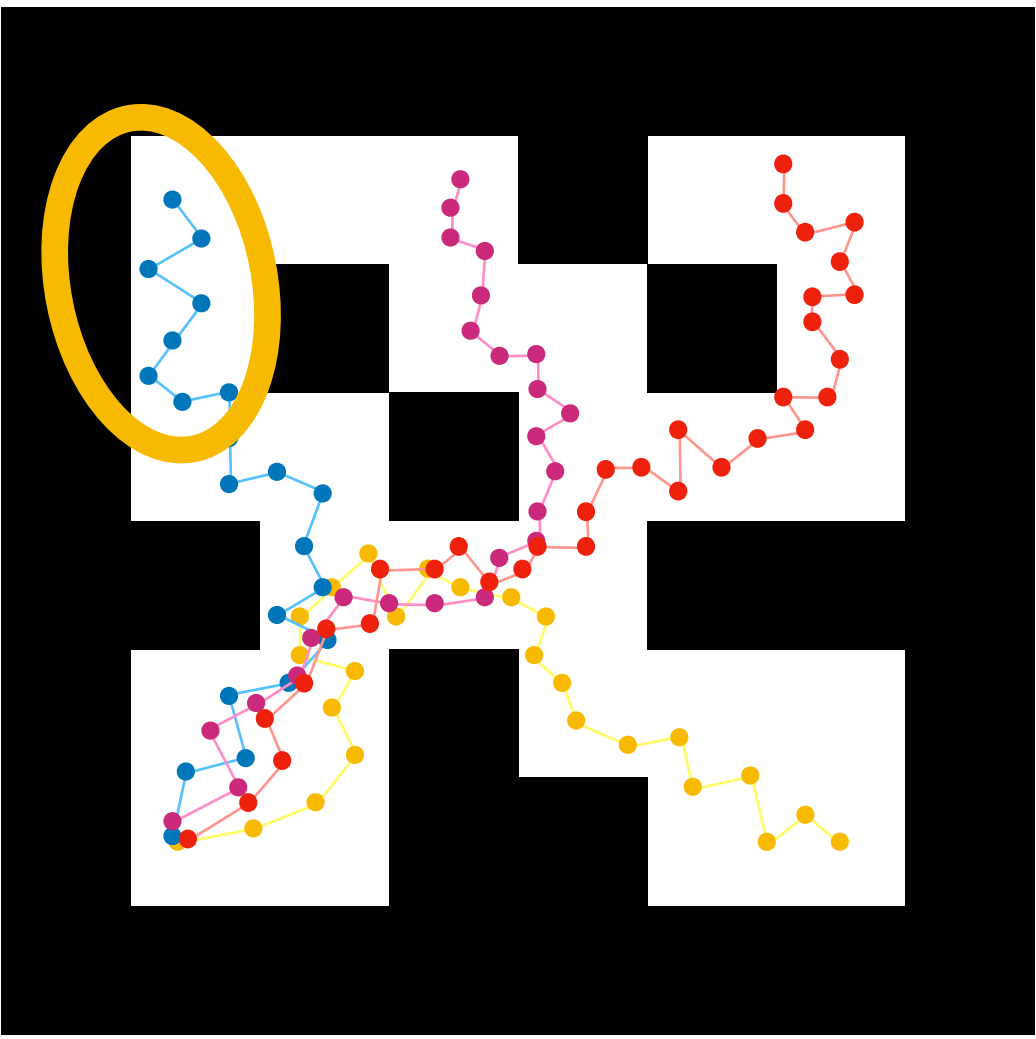
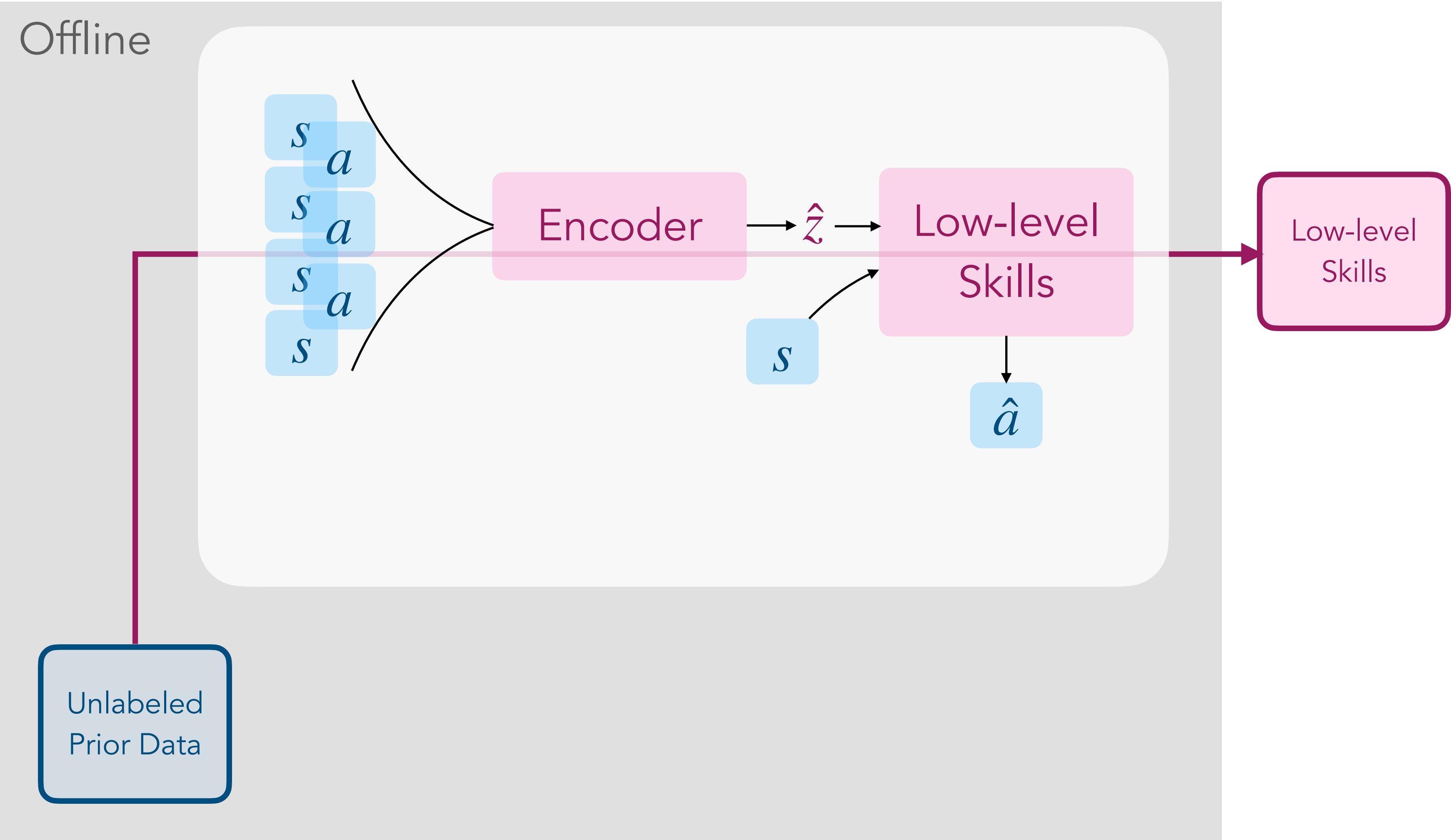
(1) Extract **low-level task-agnostic behaviors** offline with VAE skill pre-training to reduce task horizon



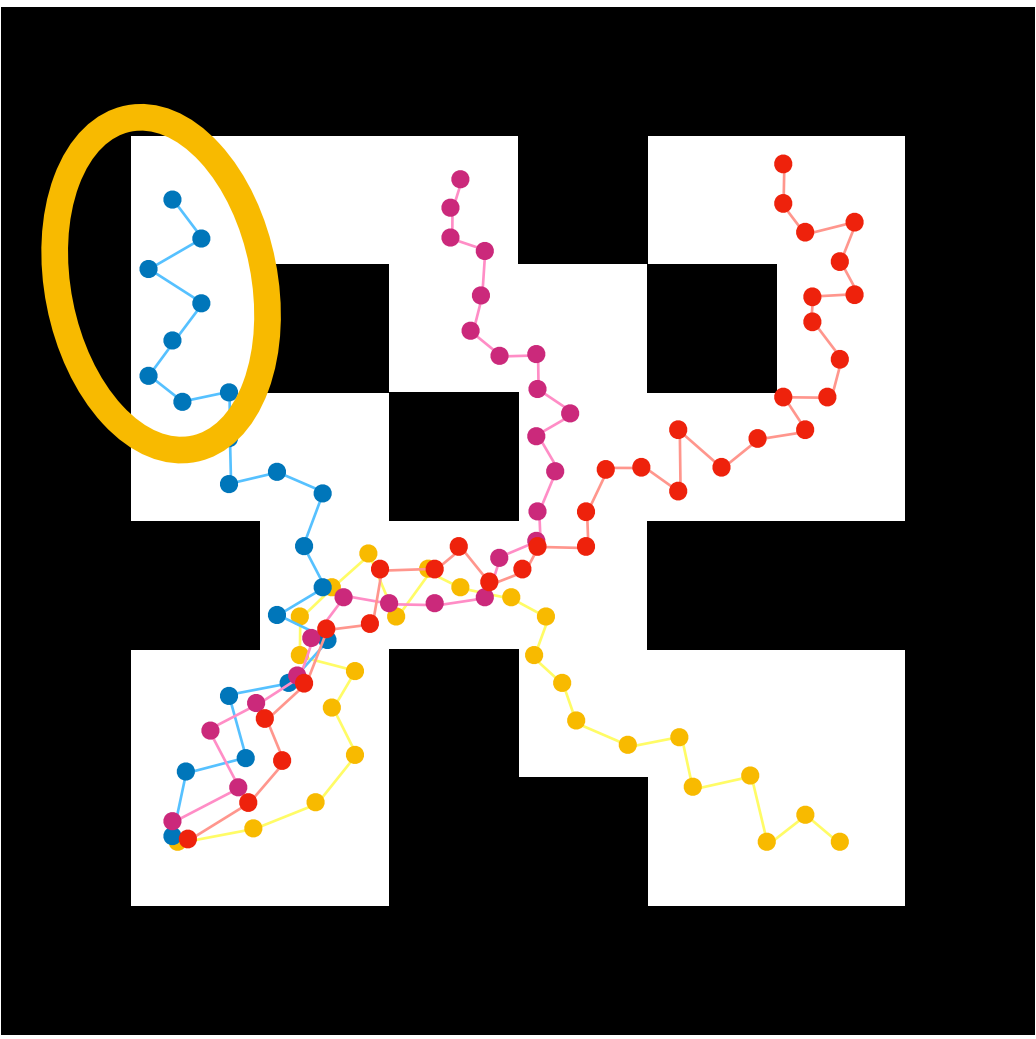
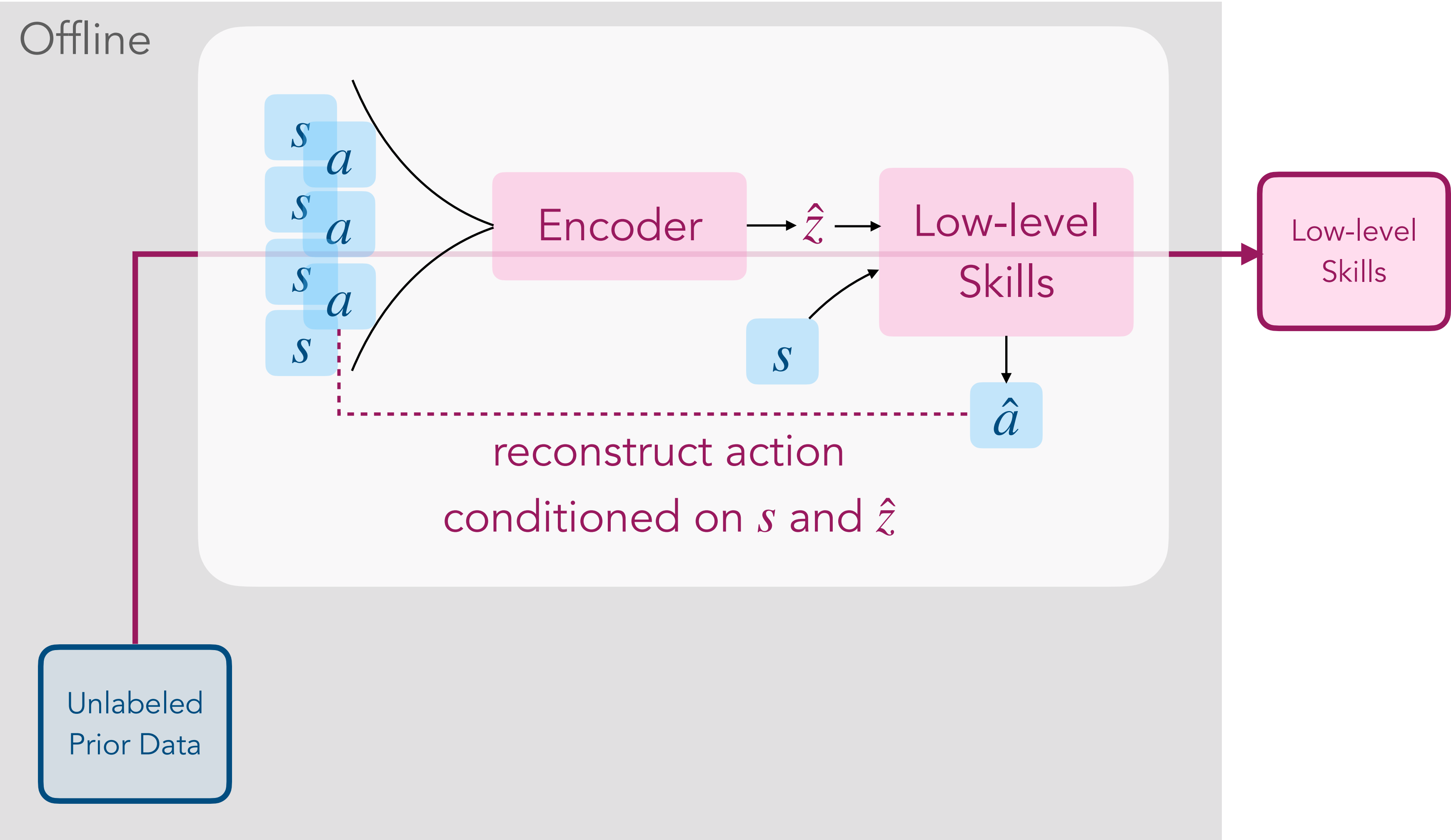
(1) Extract **low-level task-agnostic behaviors** offline with VAE skill pre-training to reduce task horizon



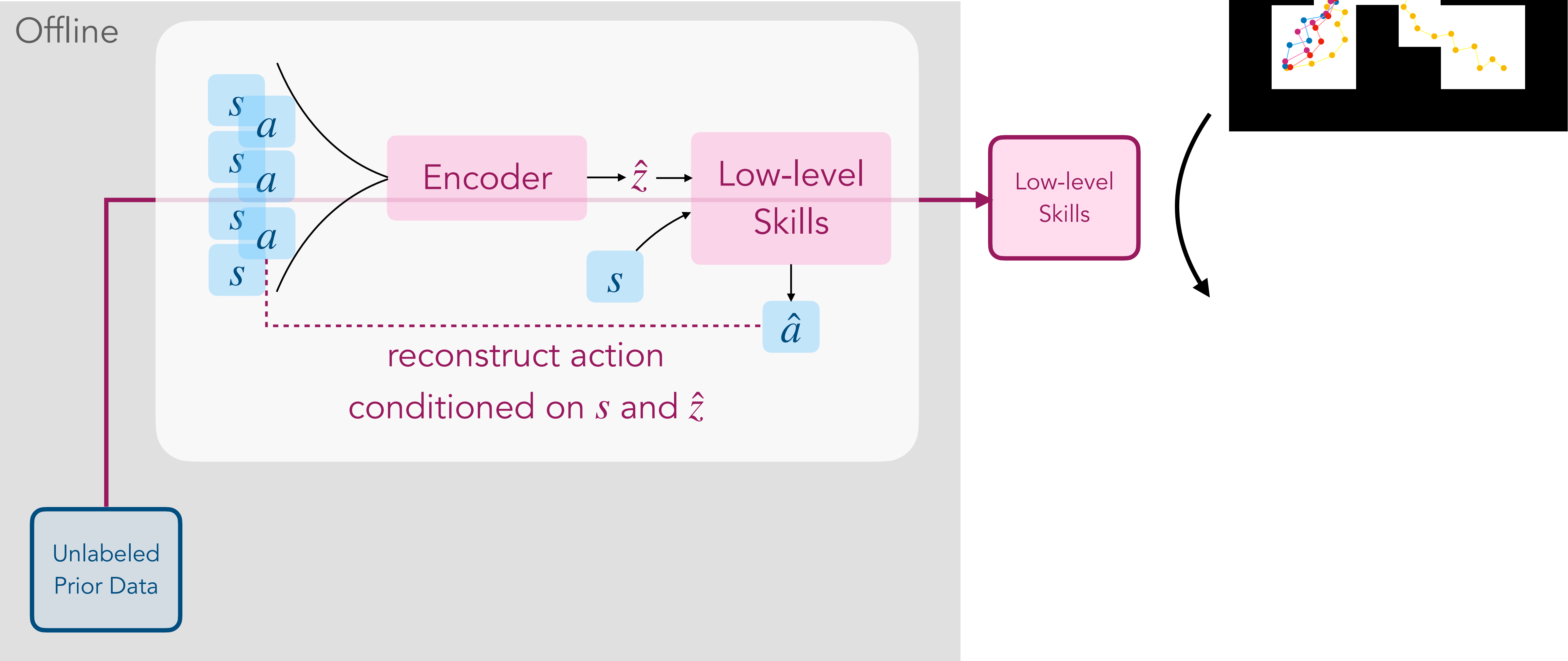
(1) Extract **low-level task-agnostic behaviors** offline with VAE skill pre-training to reduce task horizon



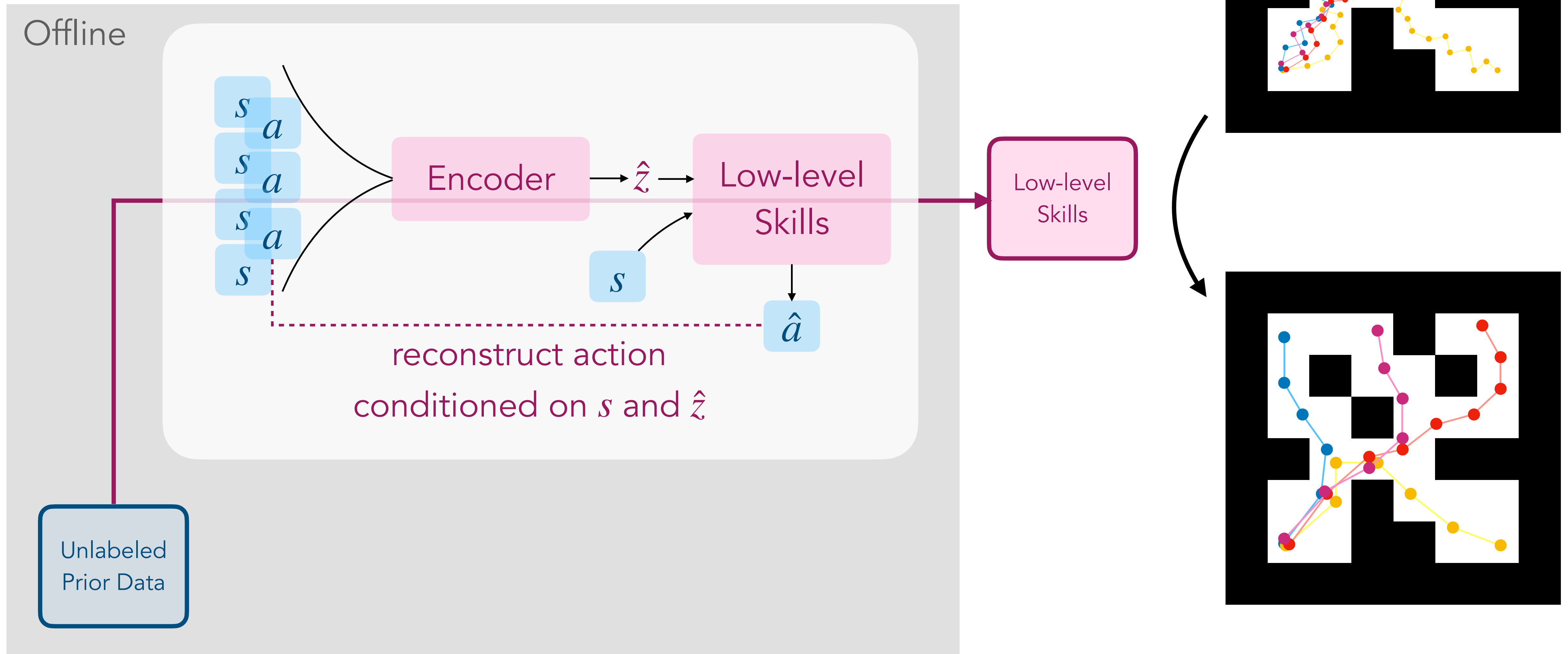
(1) Extract **low-level task-agnostic behaviors** offline with VAE skill pre-training to reduce task horizon



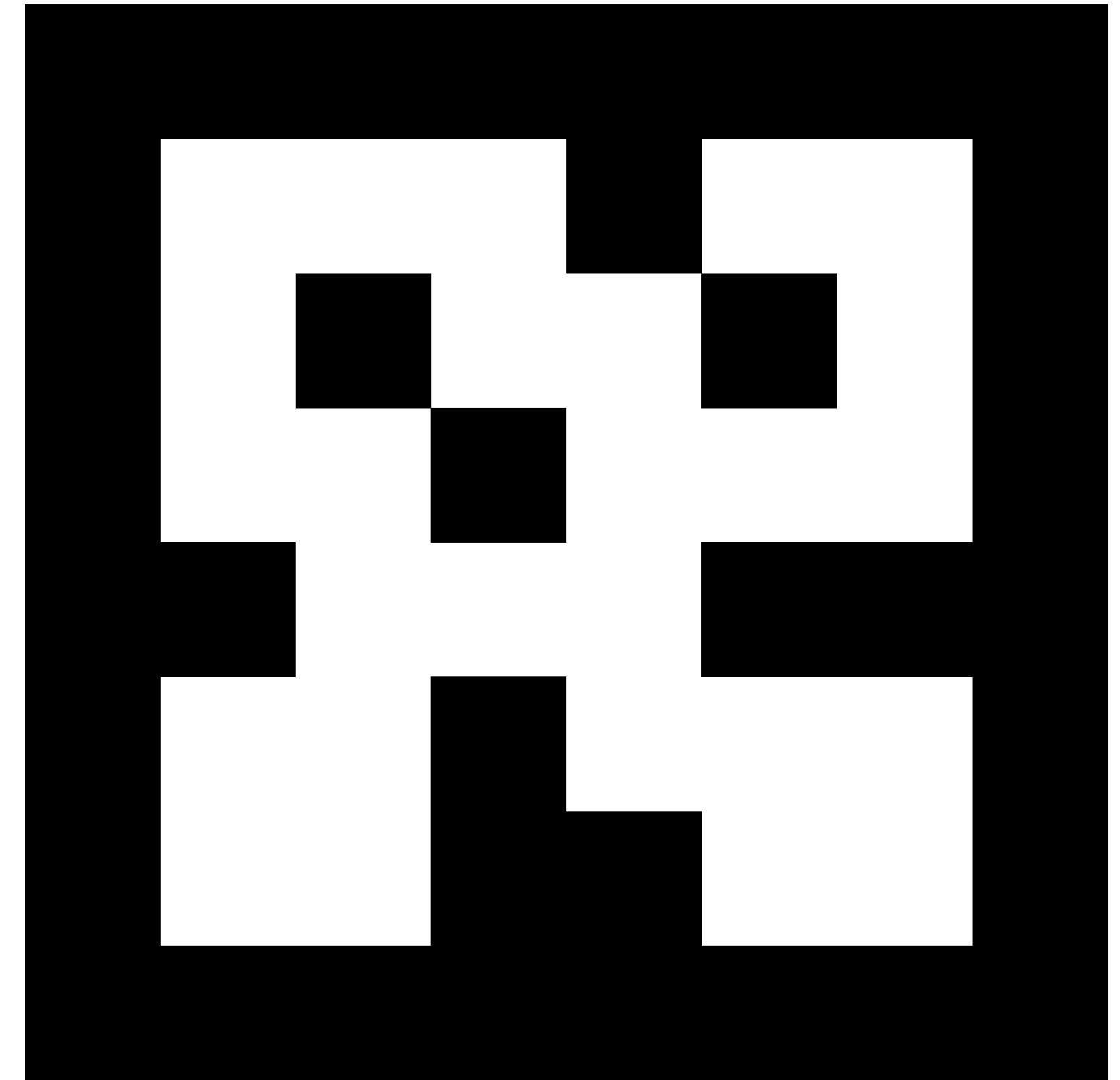
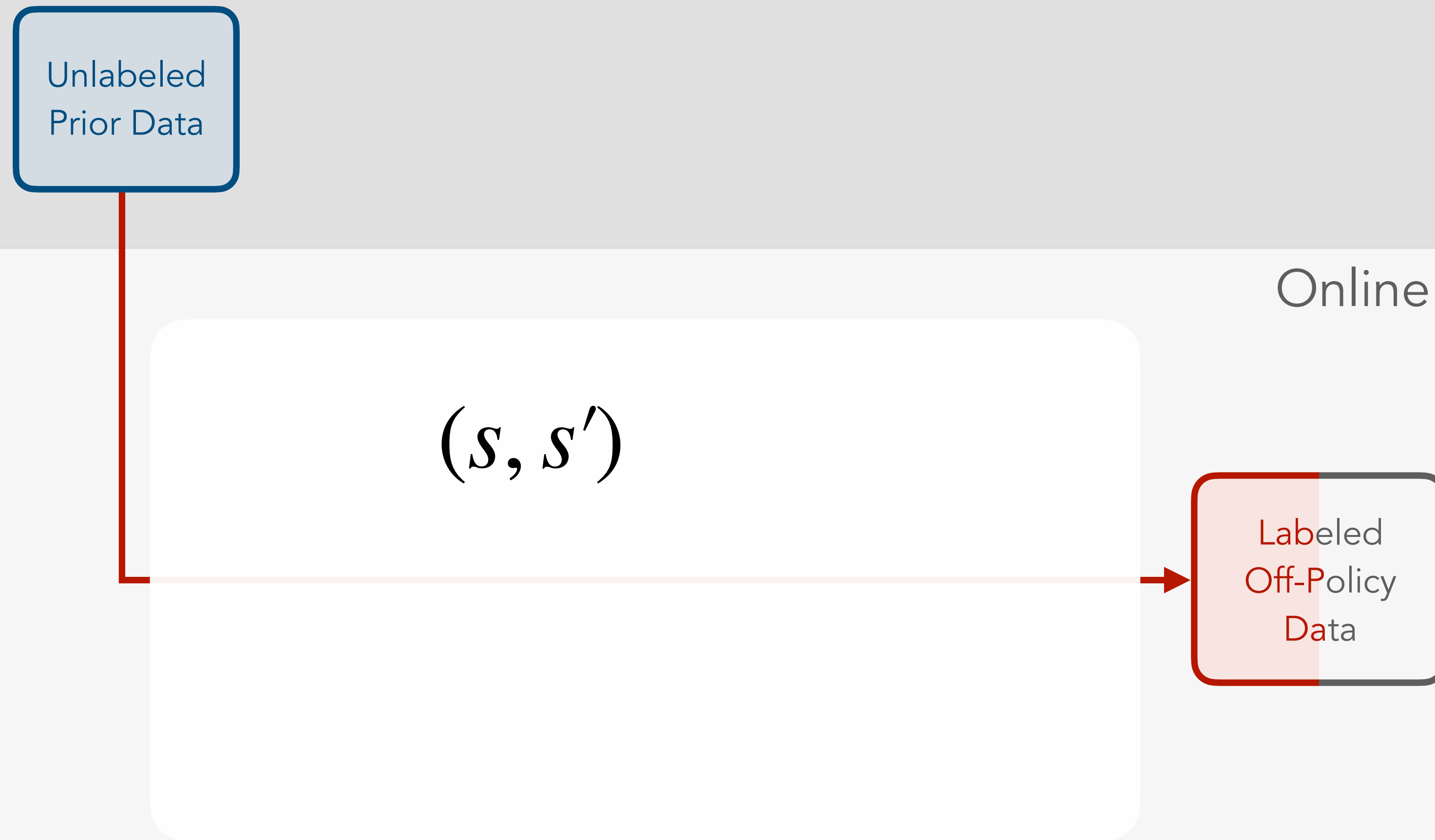
(1) Extract **low-level task-agnostic behaviors** offline with VAE skill pre-training to reduce task horizon



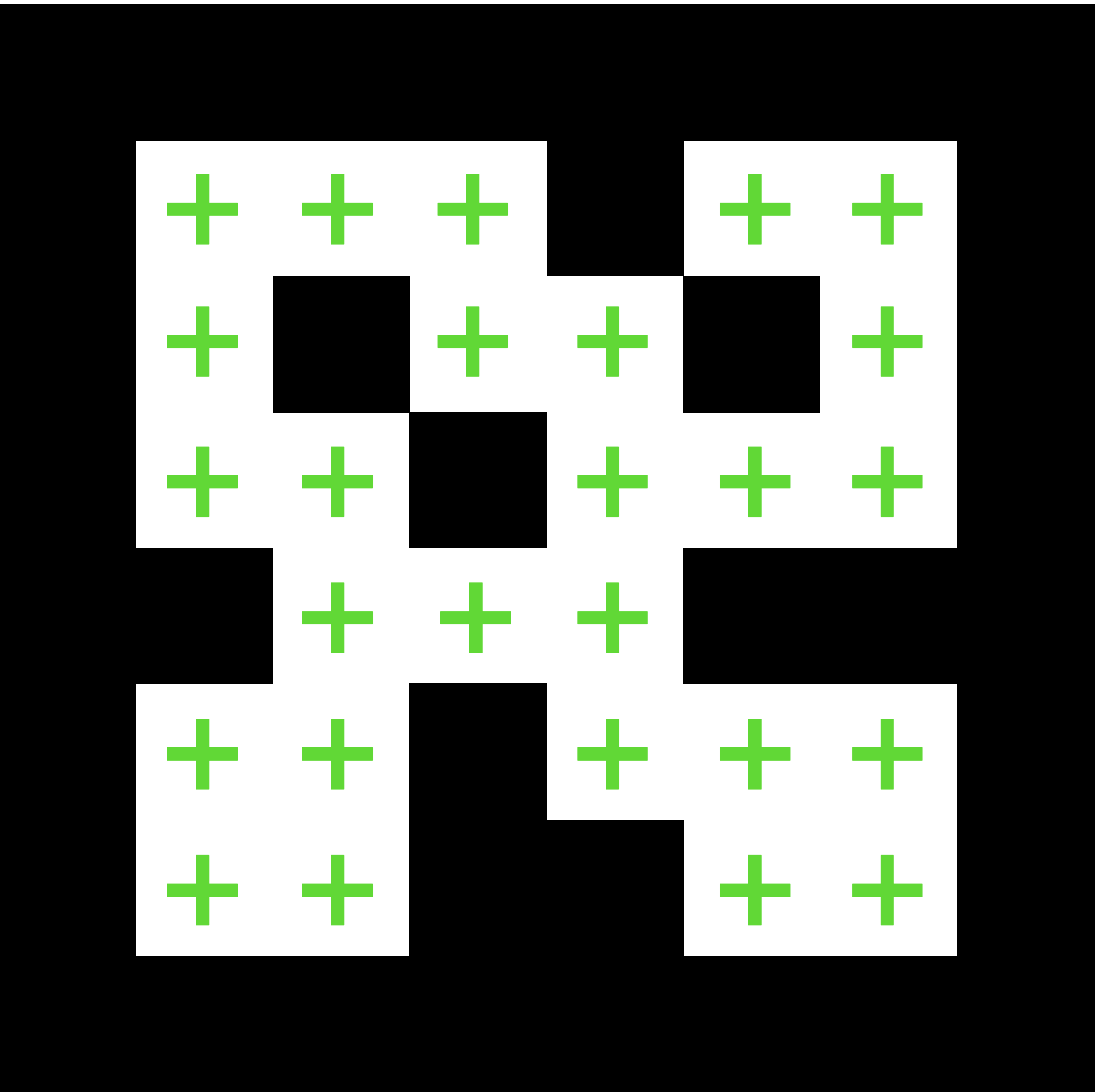
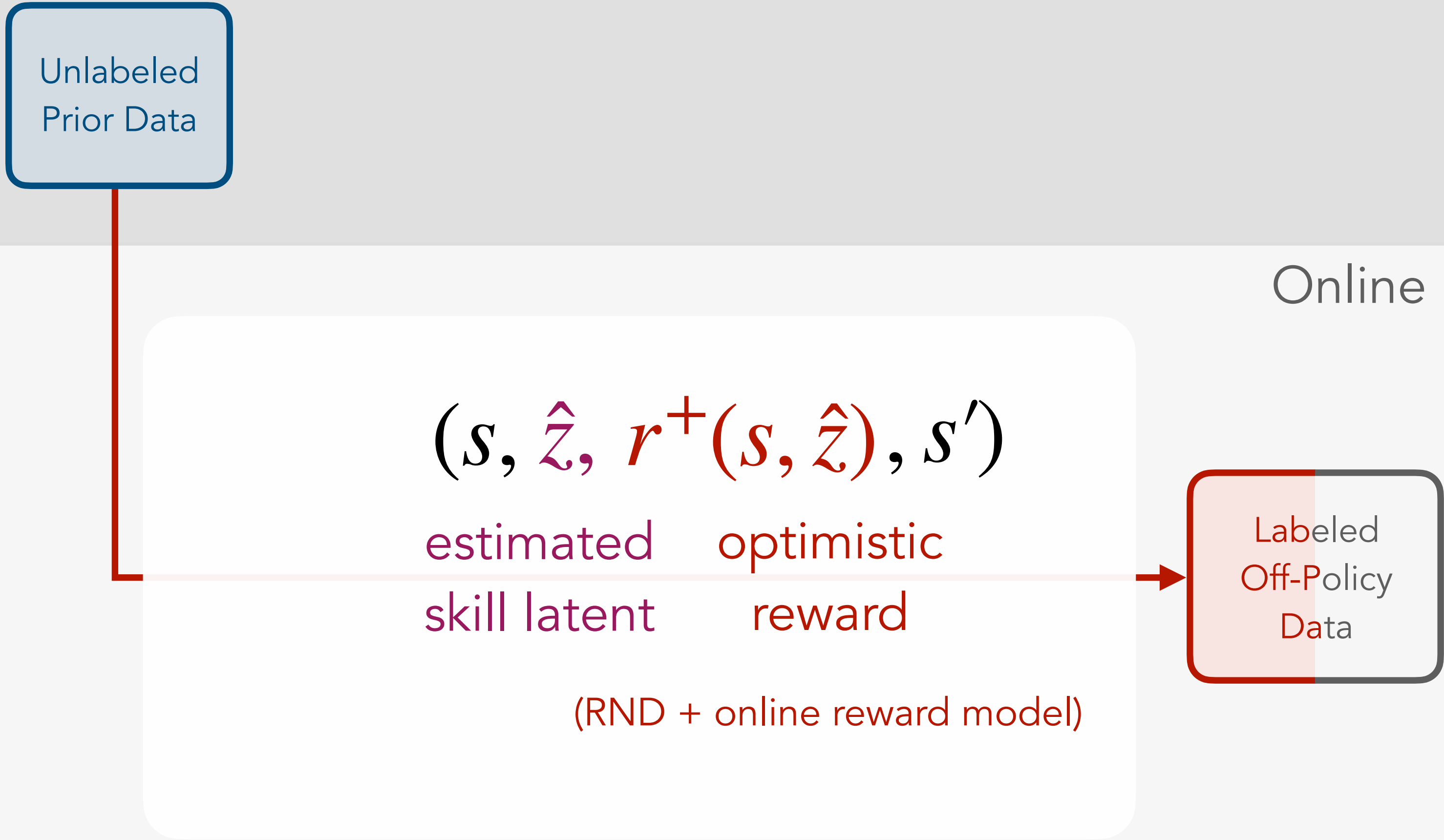
(1) Extract **low-level task-agnostic behaviors** offline with VAE skill pre-training to reduce task horizon



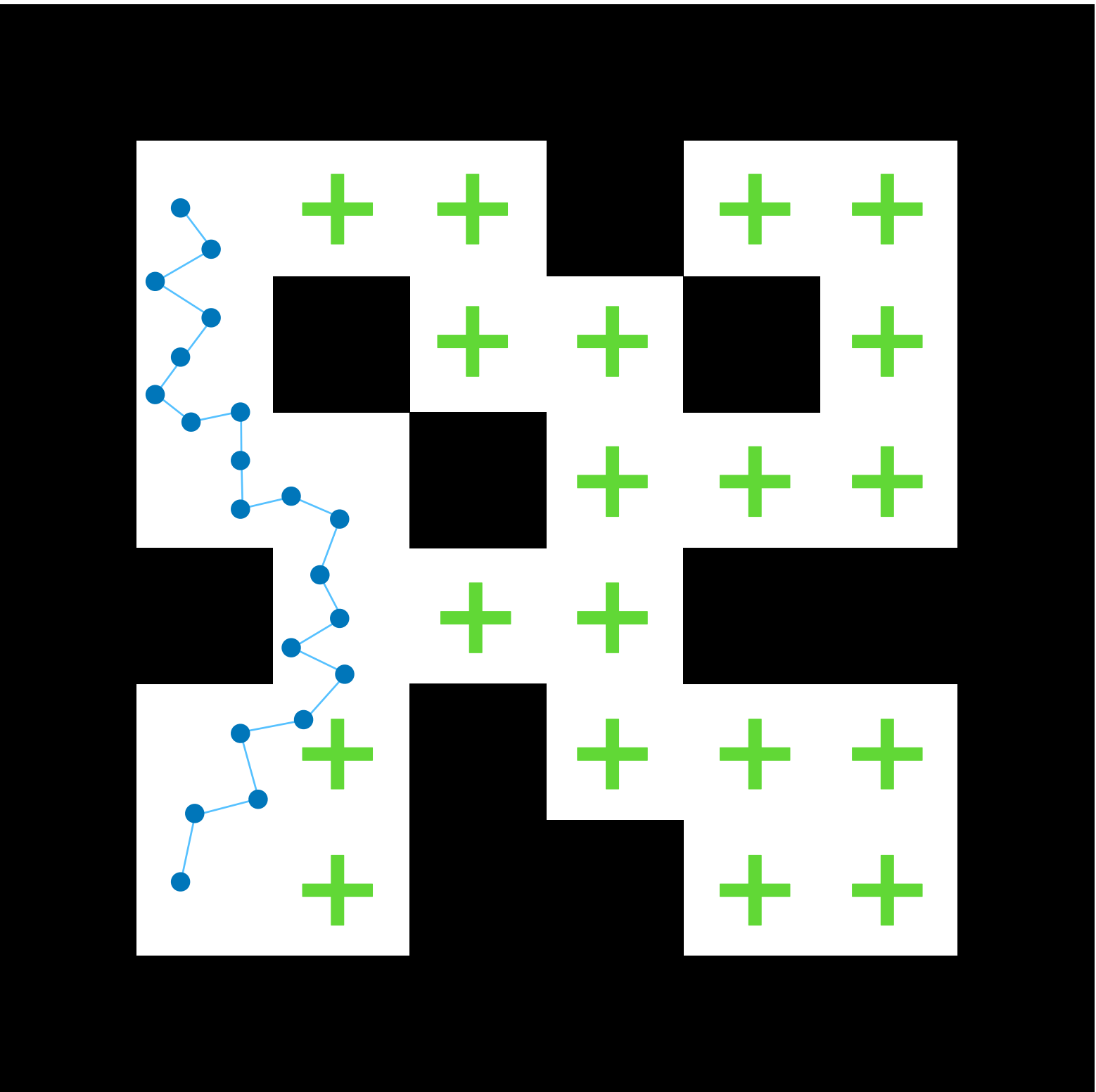
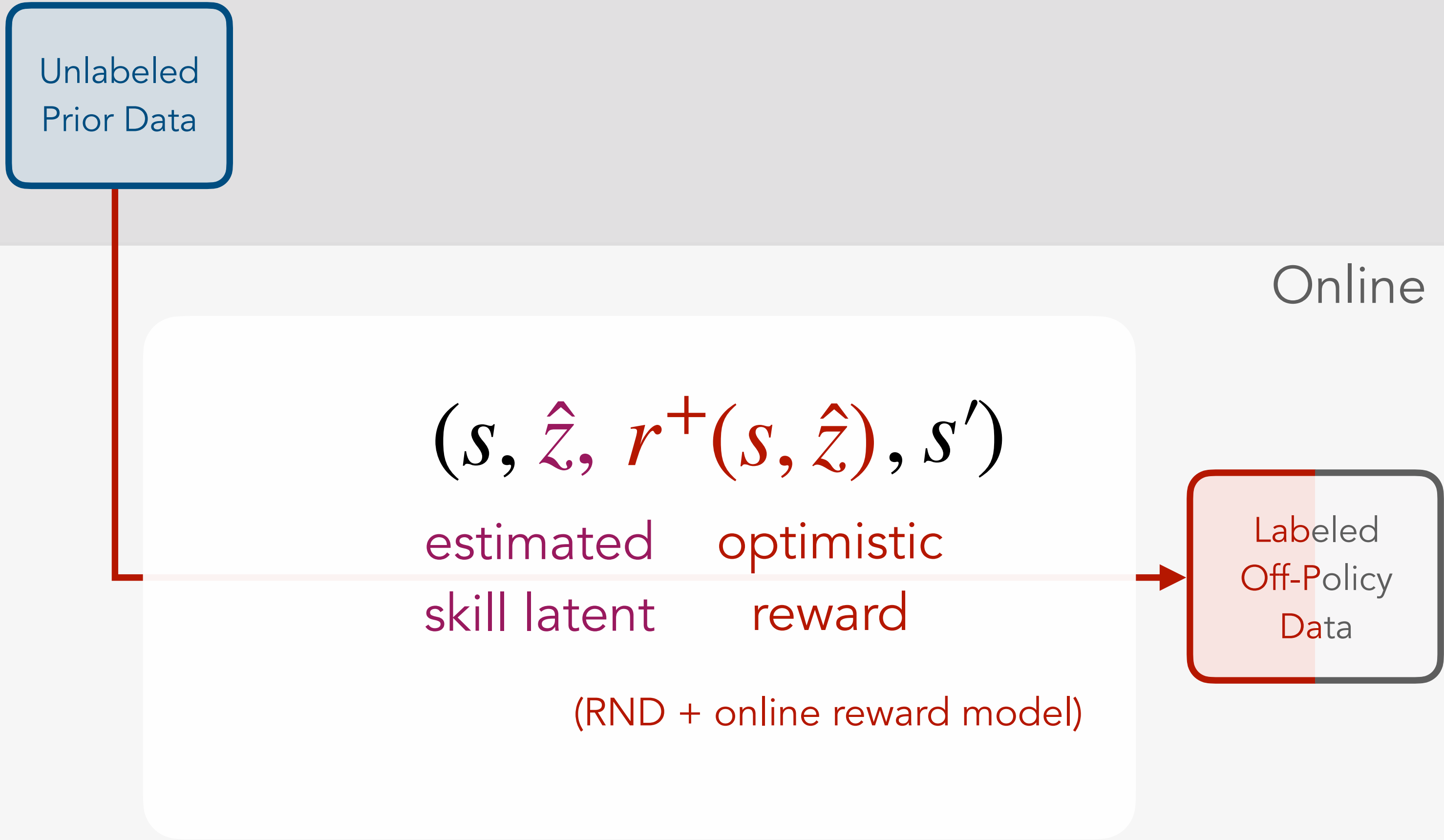
(2) Pseudo-label transitions with **optimistic rewards** to leverage the prior data as **additional off-policy data** and encourage online exploration



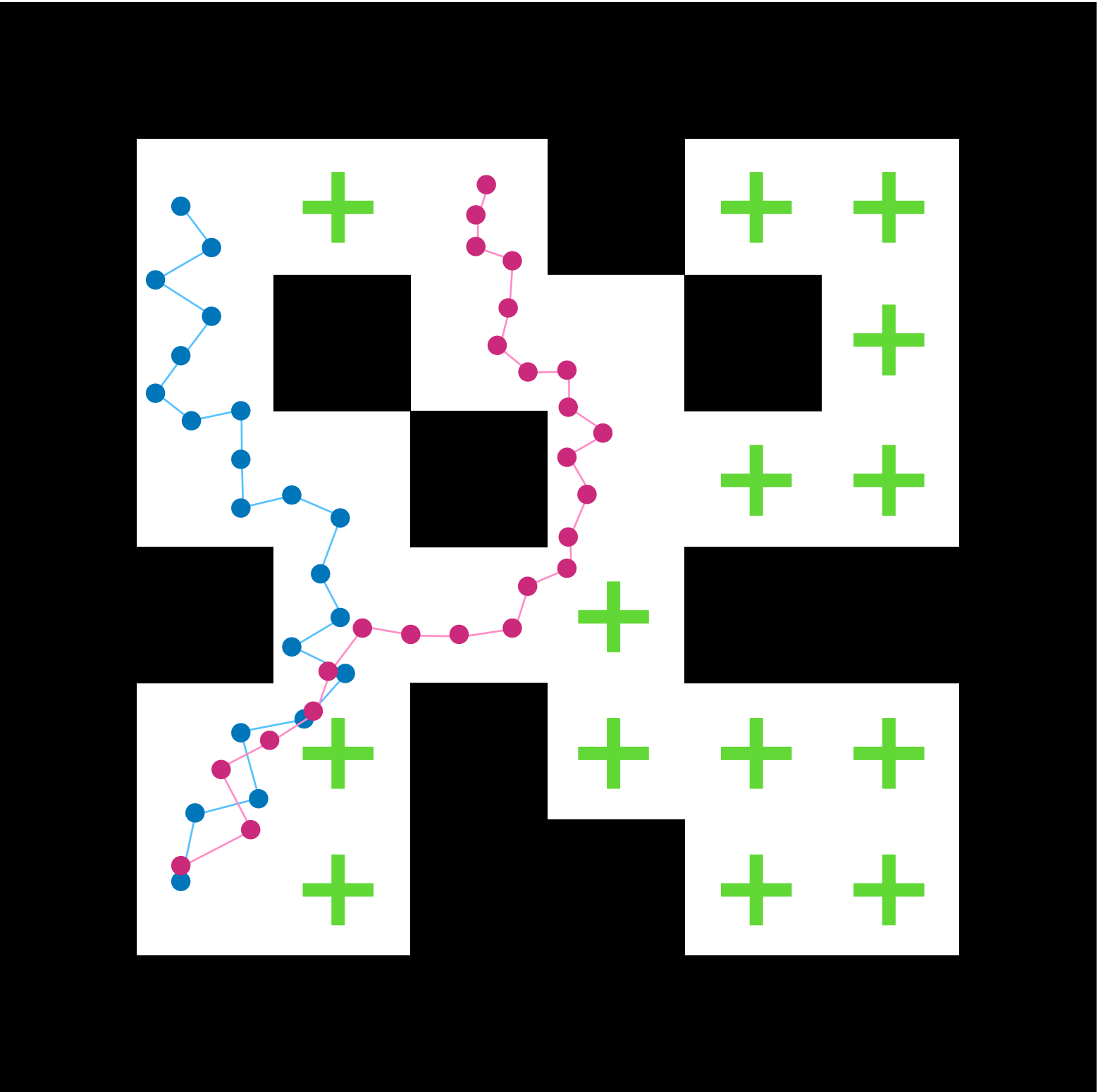
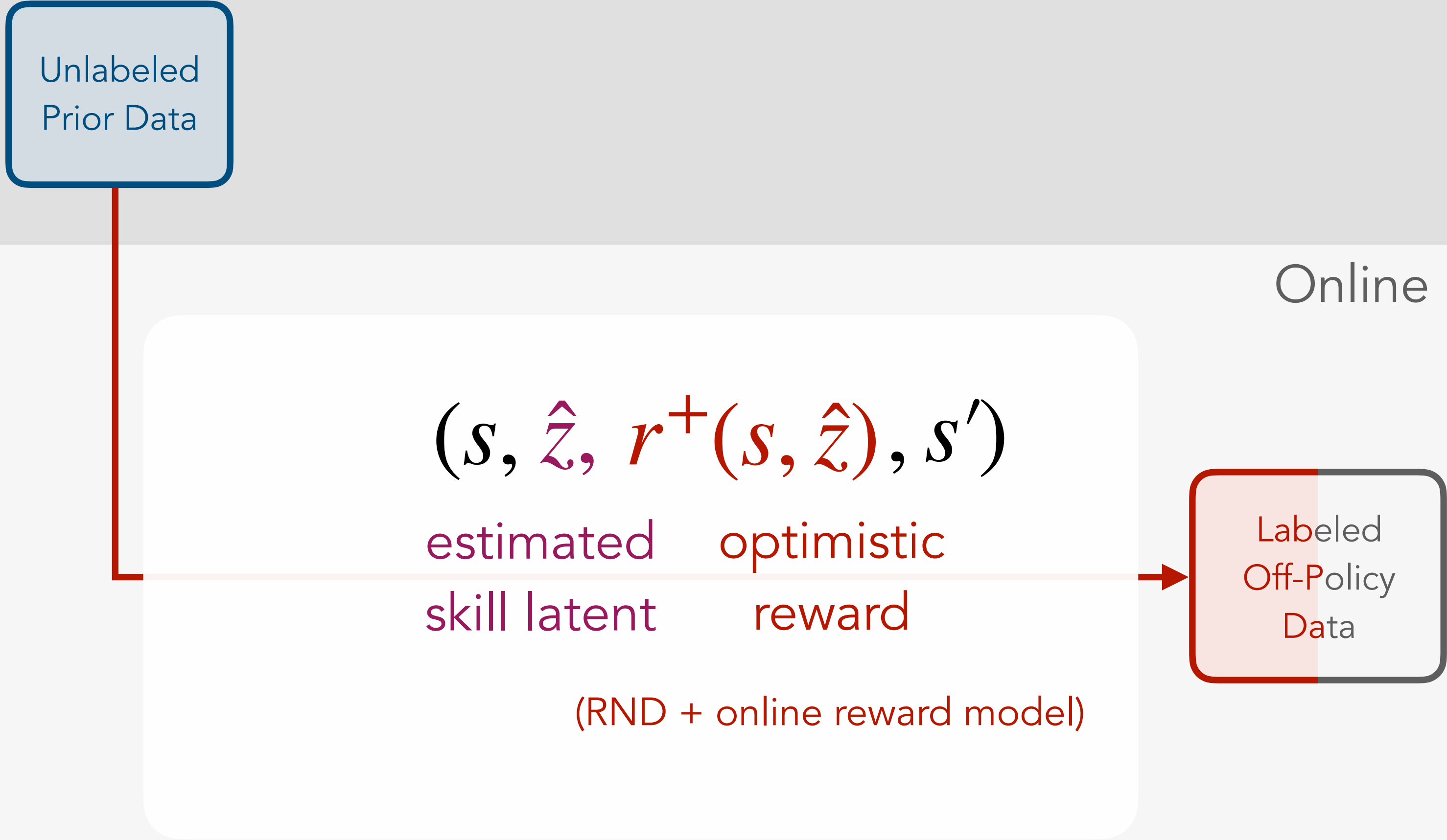
(2) Pseudo-label transitions with **optimistic rewards** to leverage the prior data as **additional off-policy data** and encourage online exploration



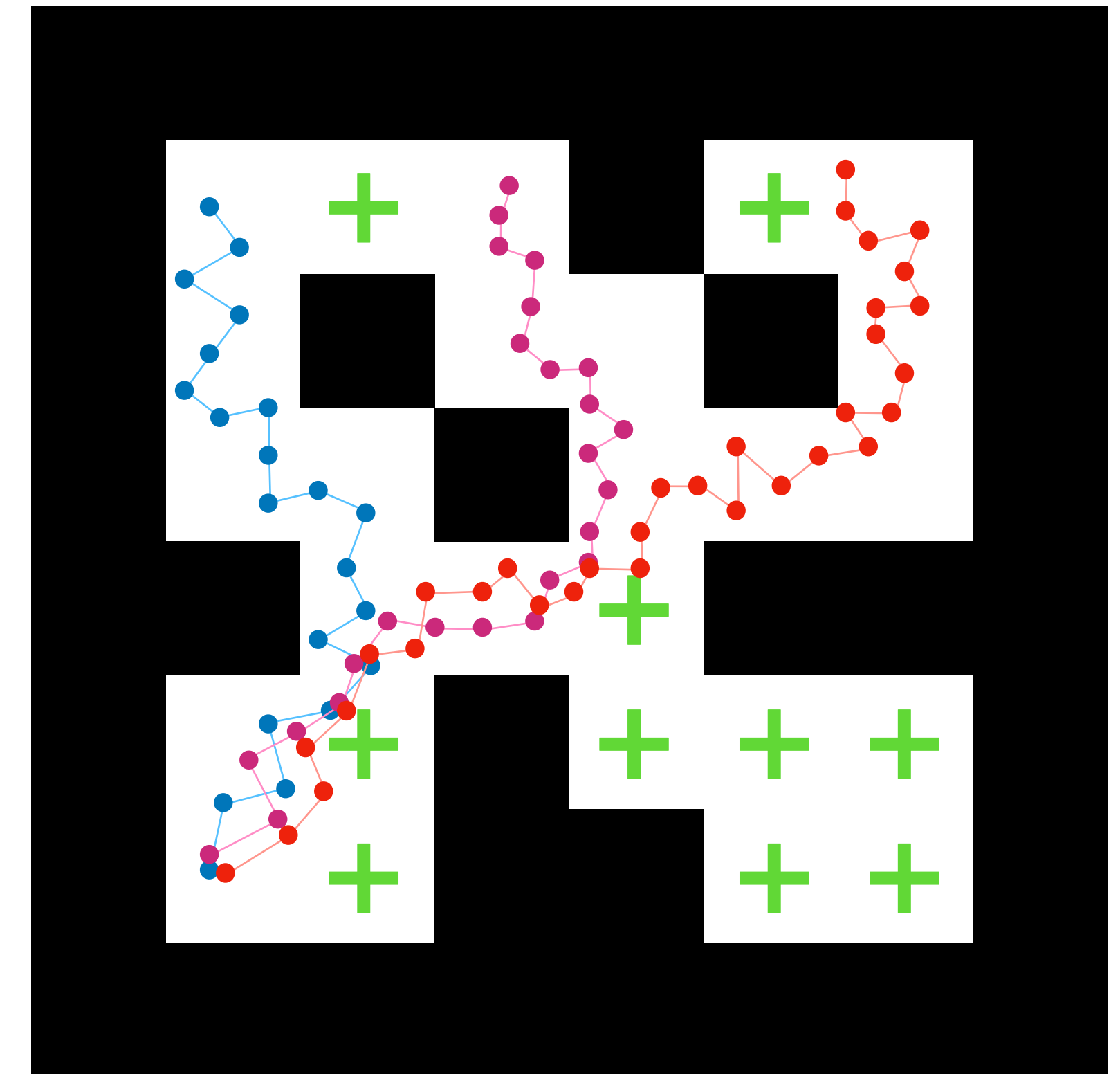
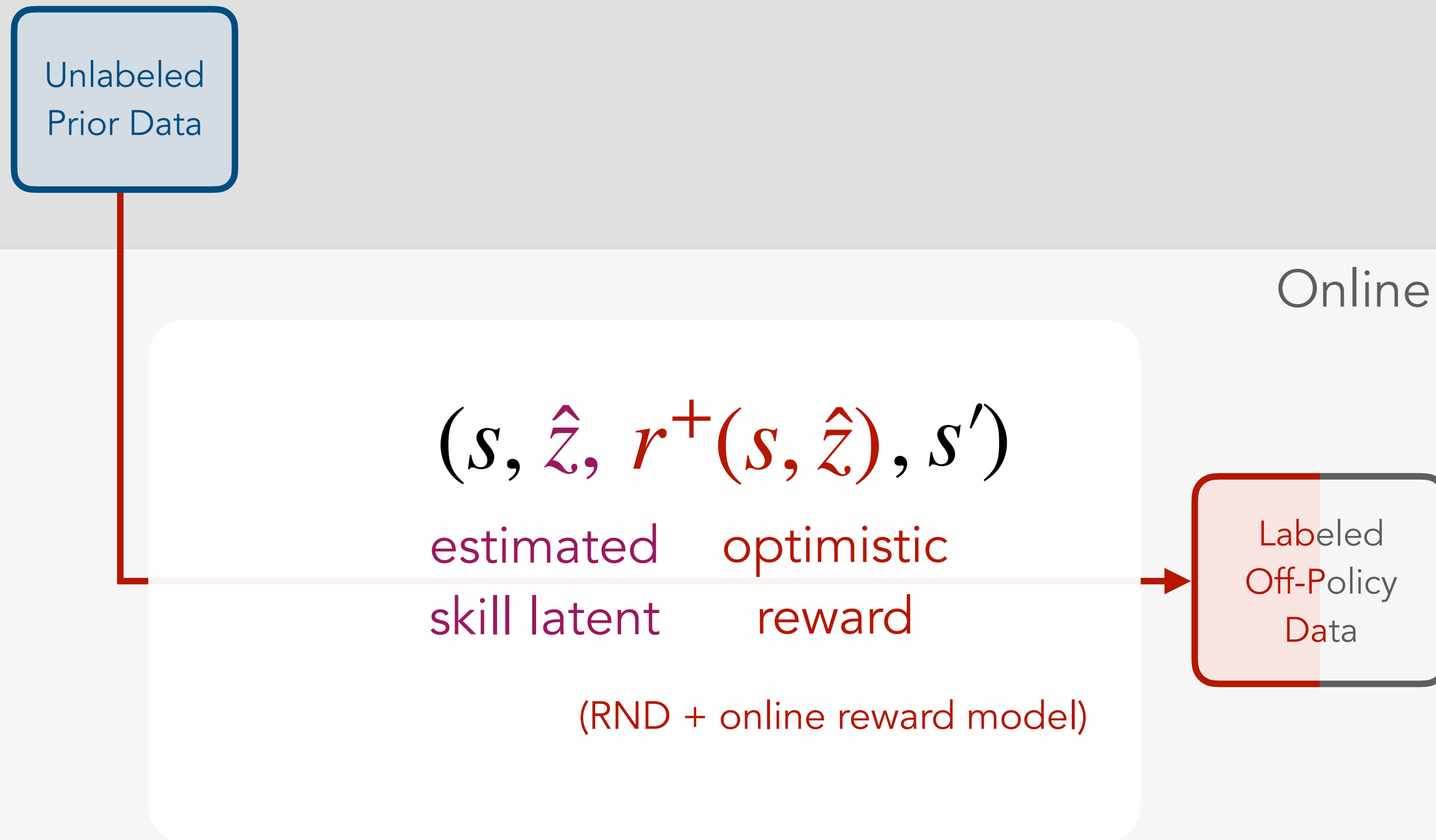
(2) Pseudo-label transitions with **optimistic rewards** to leverage the prior data as **additional off-policy data** and encourage online exploration



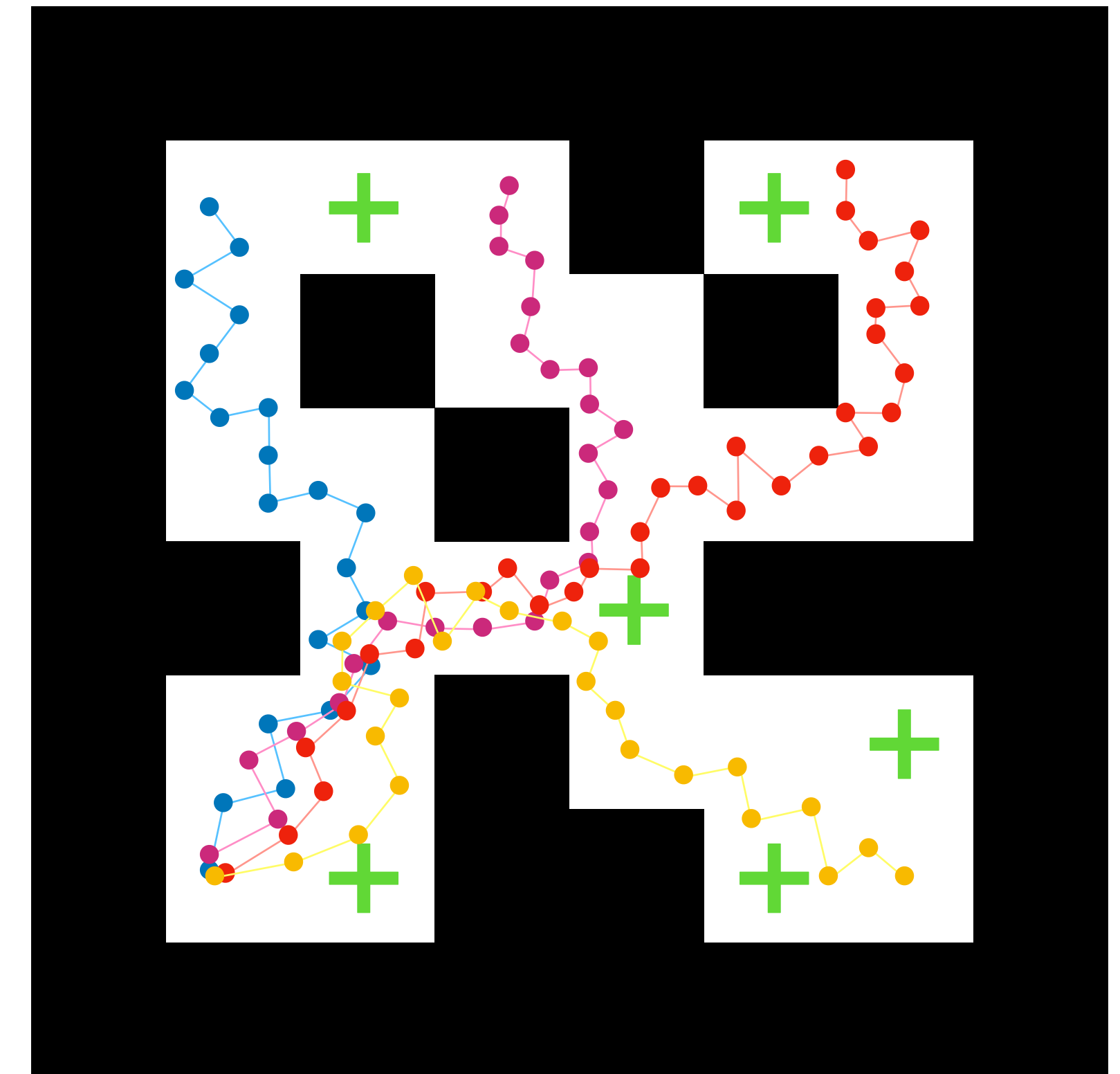
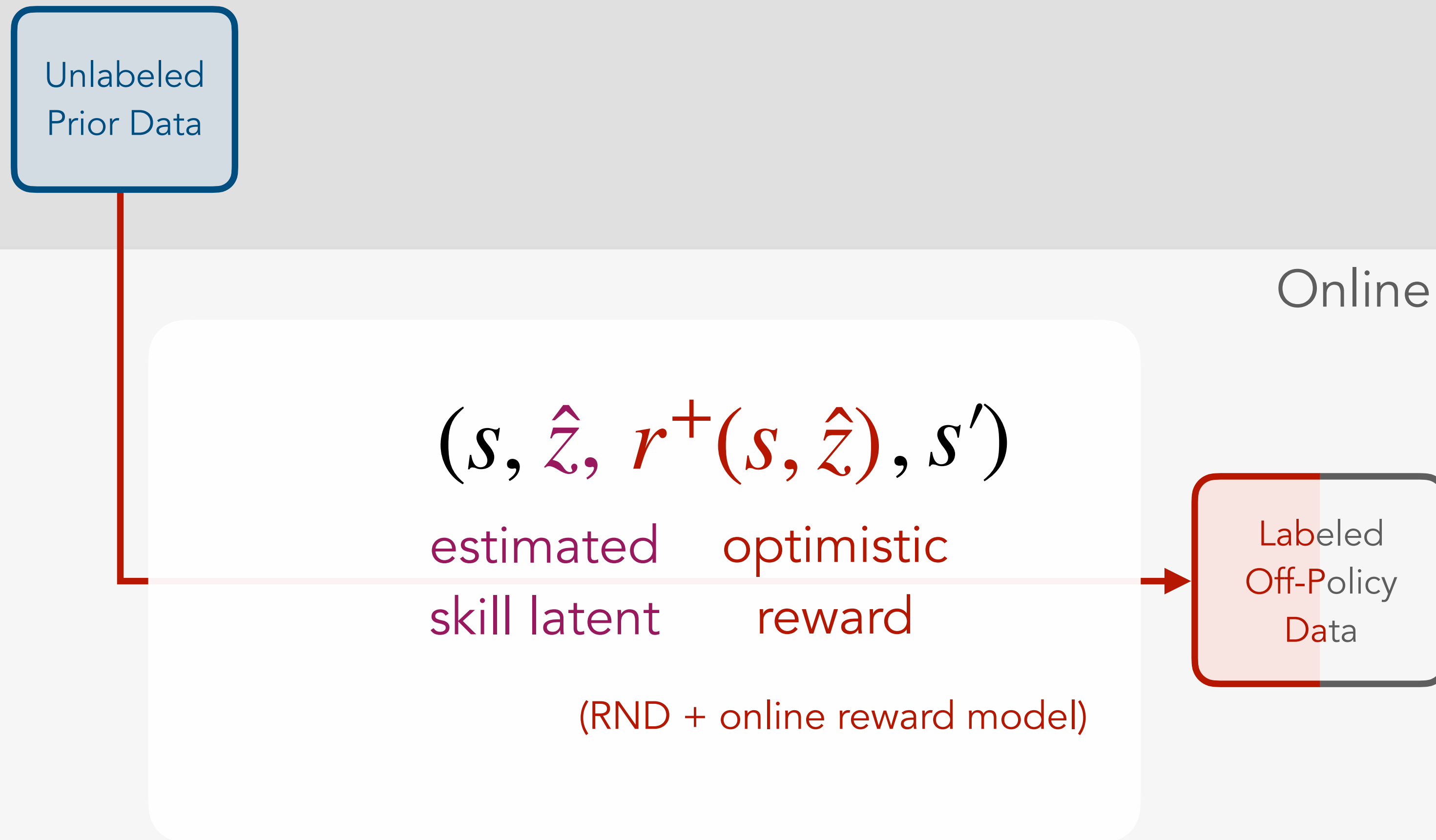
(2) Pseudo-label transitions with **optimistic rewards** to leverage the prior data as **additional off-policy data** and encourage online exploration



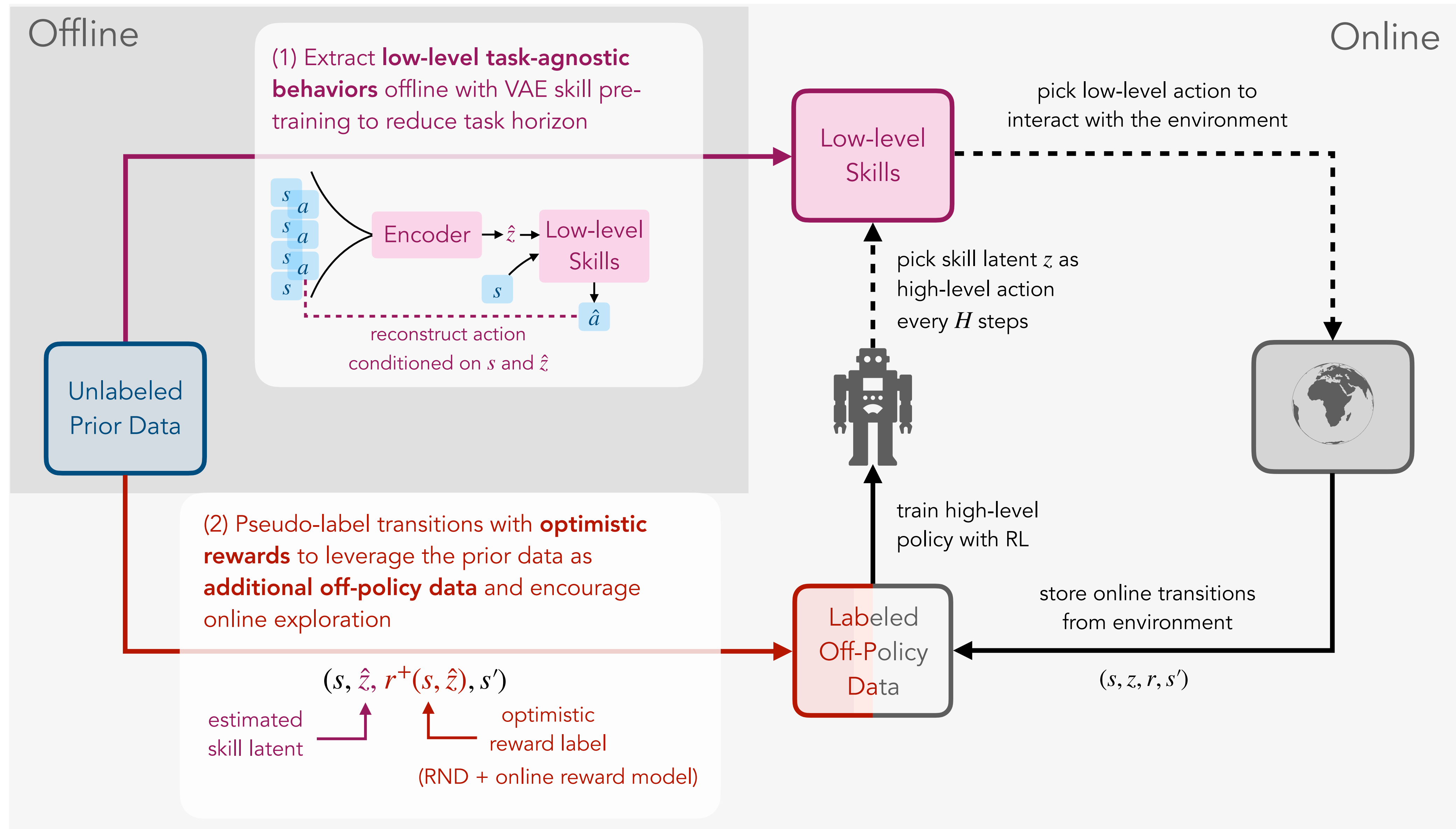
(2) Pseudo-label transitions with **optimistic rewards** to leverage the prior data as **additional off-policy data** and encourage online exploration



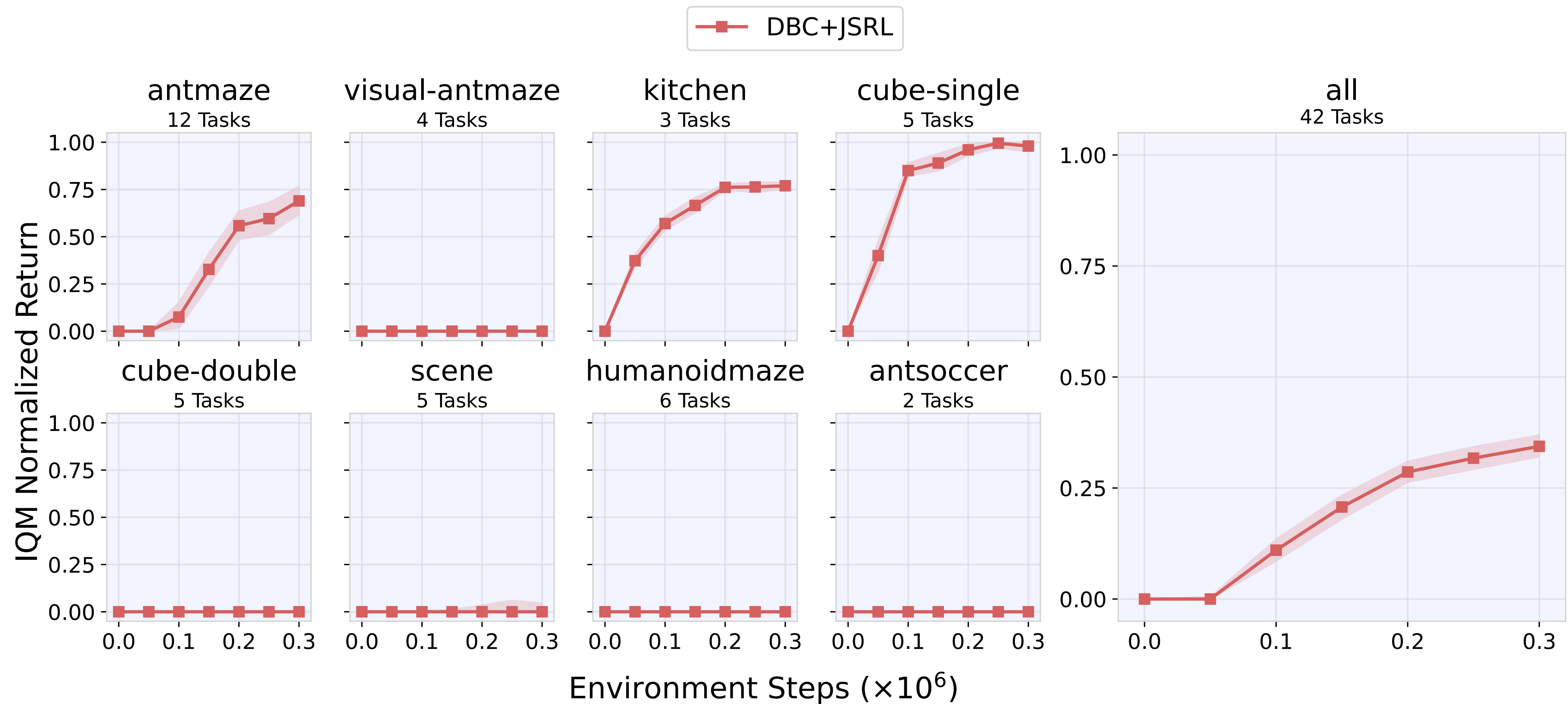
(2) Pseudo-label transitions with **optimistic rewards** to leverage the prior data as **additional off-policy data** and encourage online exploration



SUPE: Skills from Unlabeled Prior Data for Exploration

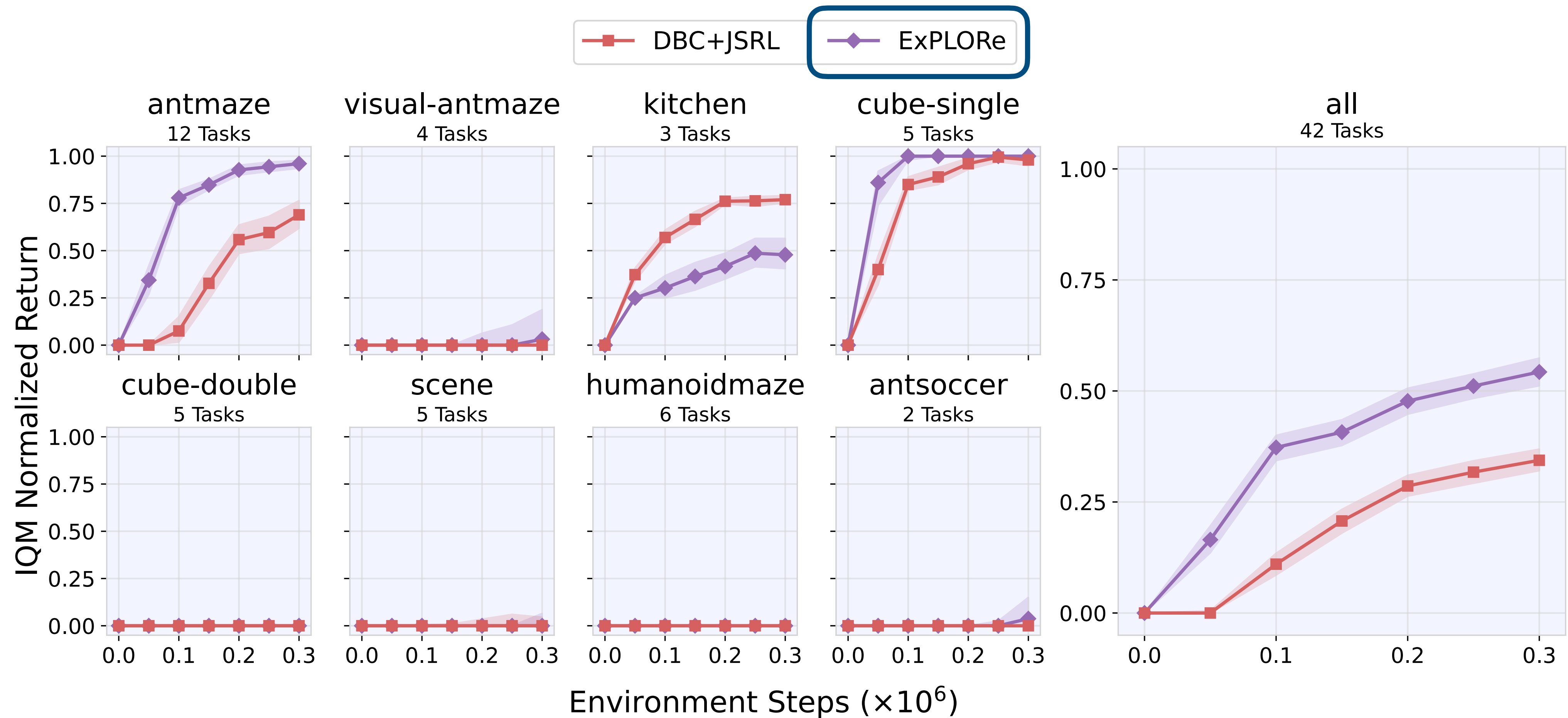


How well does *SUPE* do?

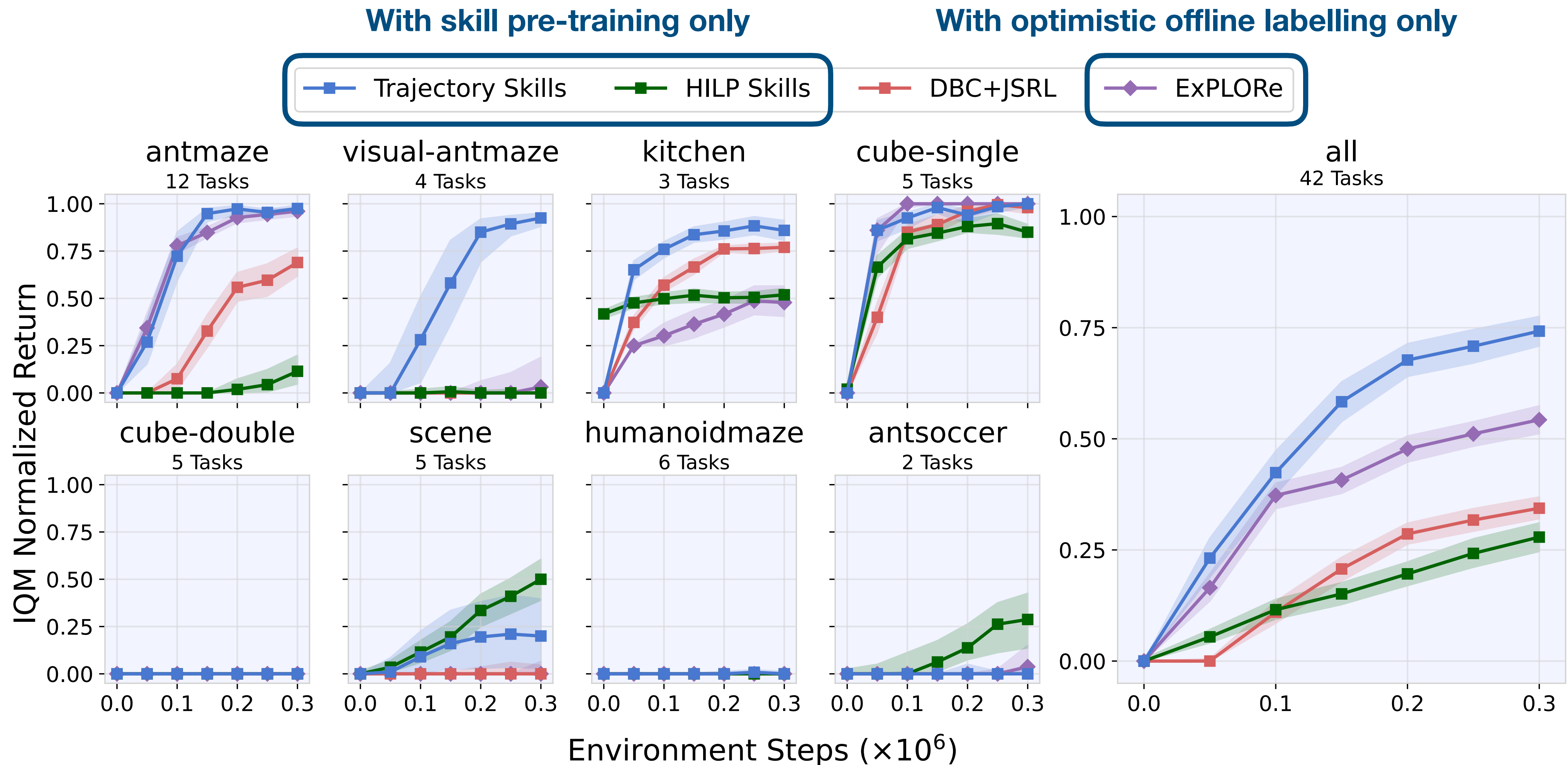


How well does *SUPE* do?

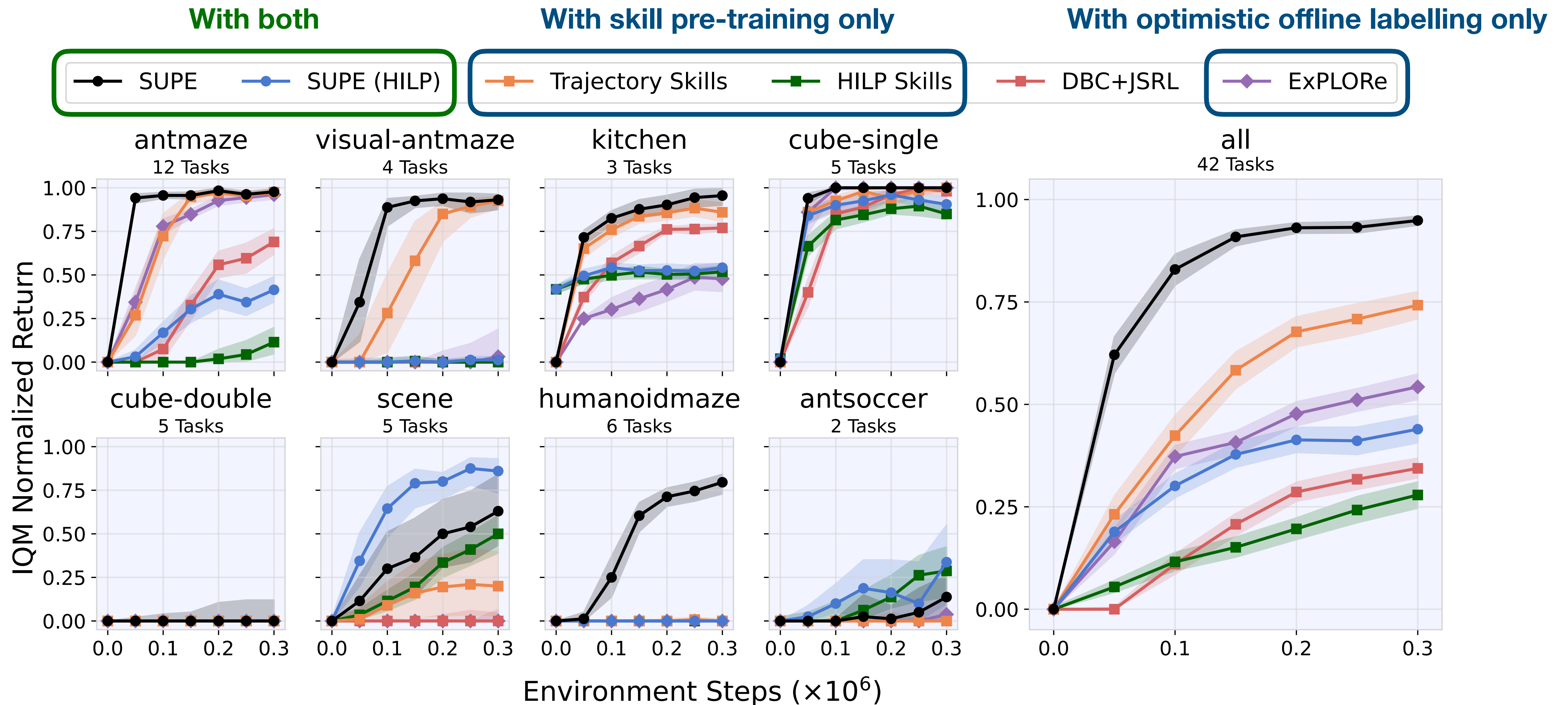
With optimistic offline labelling only



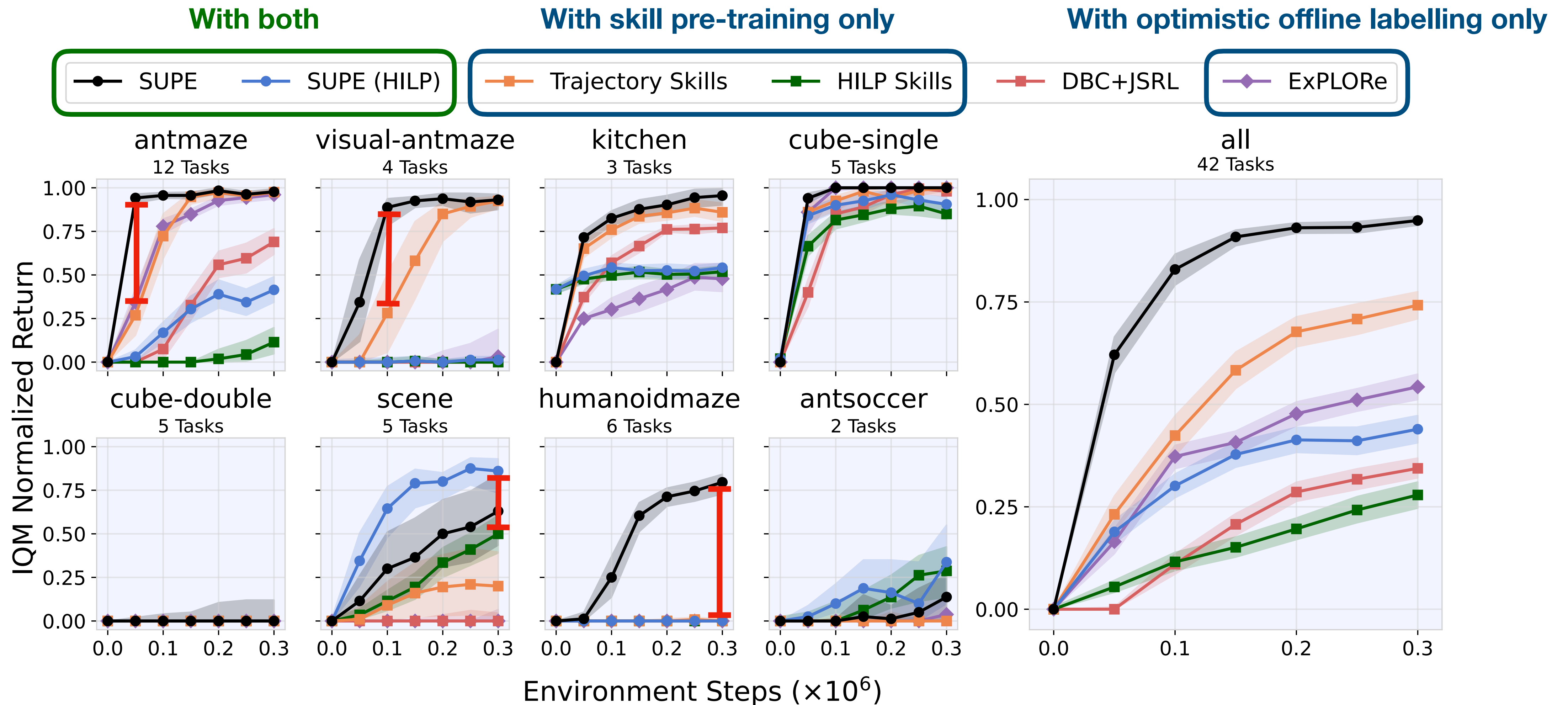
How well does *SUPE* do?



How well does *SUPE* do?

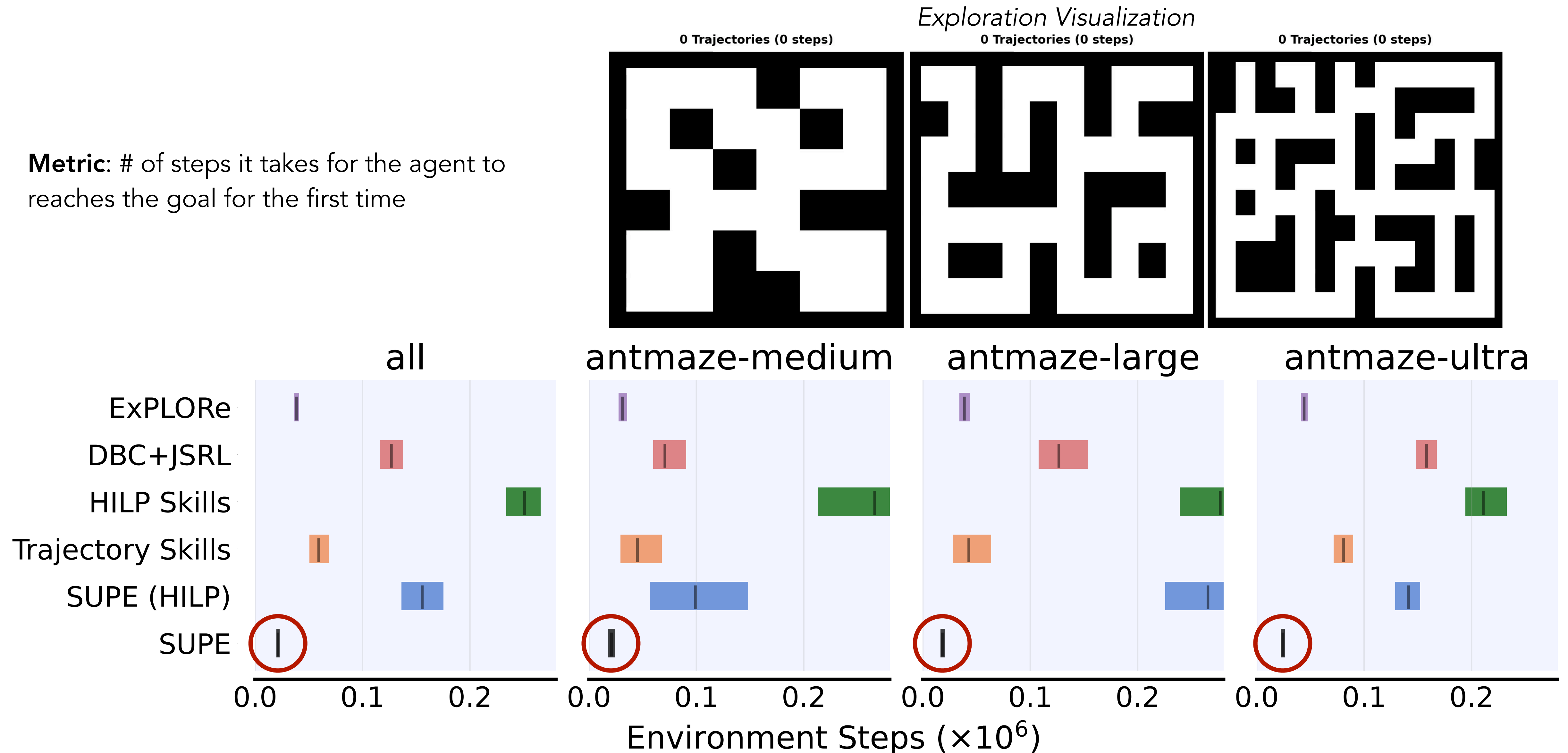


How well does *SUPE* do?



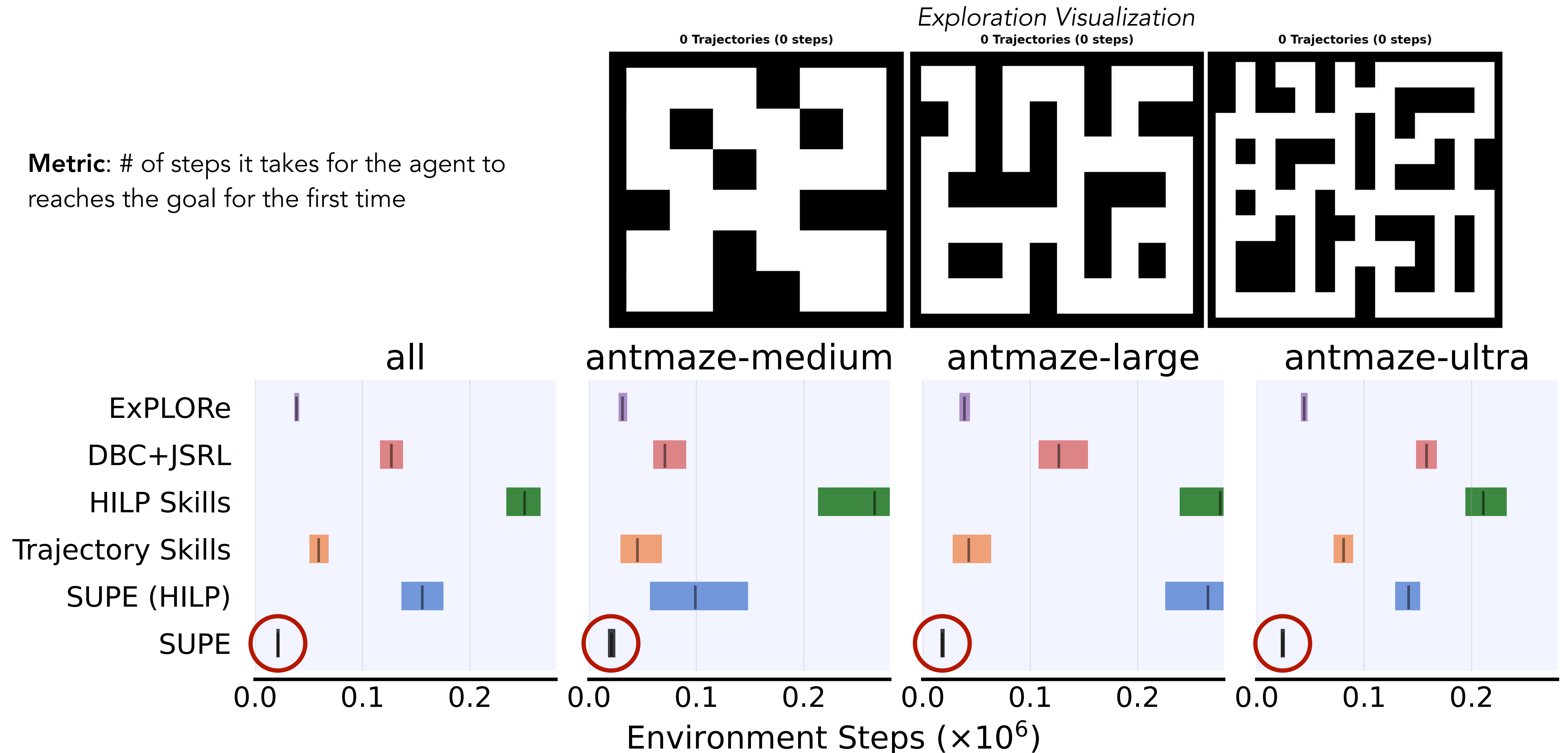
Does *SUPE* actually explore better?

Metric: # of steps it takes for the agent to reaches the goal for the first time



Does *SUPE* actually explore better?

Metric: # of steps it takes for the agent to reaches the goal for the first time



Thank you for your time!!

We have open-sourced our code, and it is available at <https://github.com/rail-berkeley/SUPE>.