# Reinforcement Learning with Robot Foundation Models: Data-Driven RL
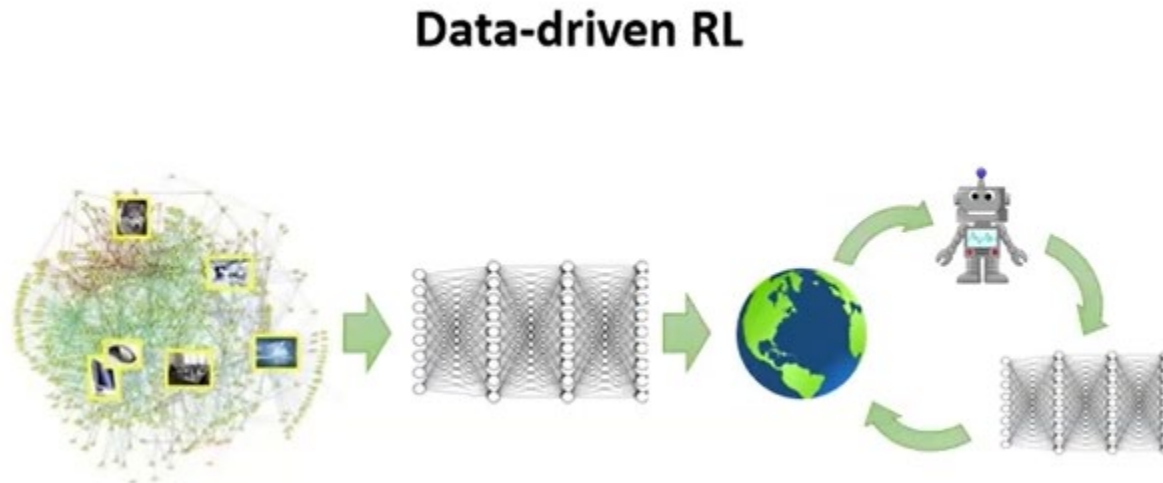
Sung-eui Yoon

Scalable Graphics, Vision and Robotics

SGVR Lab
KAIST

# Class Objectives

- **Data-driven RL**
  - Offline RL algorithms
  - Online finetuning from offline initializations
  - Making all this work with big scalable models (unclear yet)
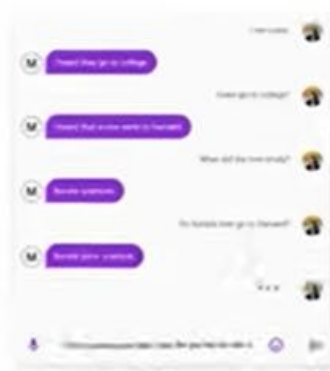  - Covered in draft of my book



Data-driven RL
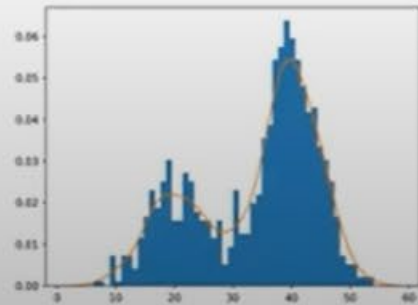
Ack.: Sergey Levine's talk slides

SGVR Lab
KAIST

# From Modern data-driven AI (Estimation) to Goals: Rethinking Foundation Models
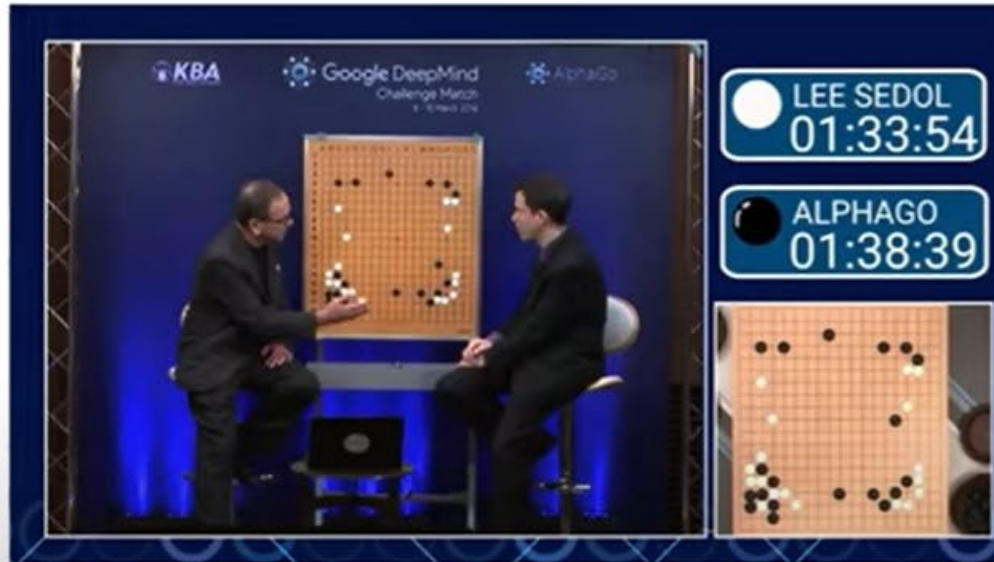


$$p_\theta(\mathbf{x}) \qquad p_\theta(\mathbf{y}|\mathbf{x})$$

SGVR Lab
KAIST

# Emergent Behavior in RL vs. Human-Like Imitation in Data-Driven Models



Impressive because no person had thought of it!

"Move 37" in Lee Sedol AlphaGo match: reinforcement learning "discovers" a move that surprises everyone

Impressive because it looks like something a person might draw!

SGVR Lab
KAIST

# So, where are we now?



Data-Driven AI

Explaining a joke

Prompt

Explain this joke:

Joke: Did you see that Google just hired an eloquent whale for their TPU team? It showed them how to communicate between two different pods!
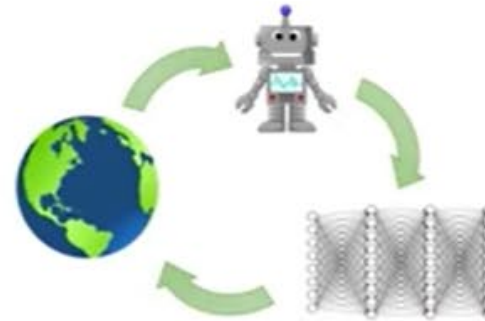
Model Response

Prediction: TPUs are a type of computer chip that Google uses for deep learning. A "pod" is a group of TPUs. A "pod" is also a group of whales. The joke is that the whale is able to communicate between two groups of whales, but the speaker is pretending that the whale is able to communicate between two groups of TPUs.

+ learns about the real world from data

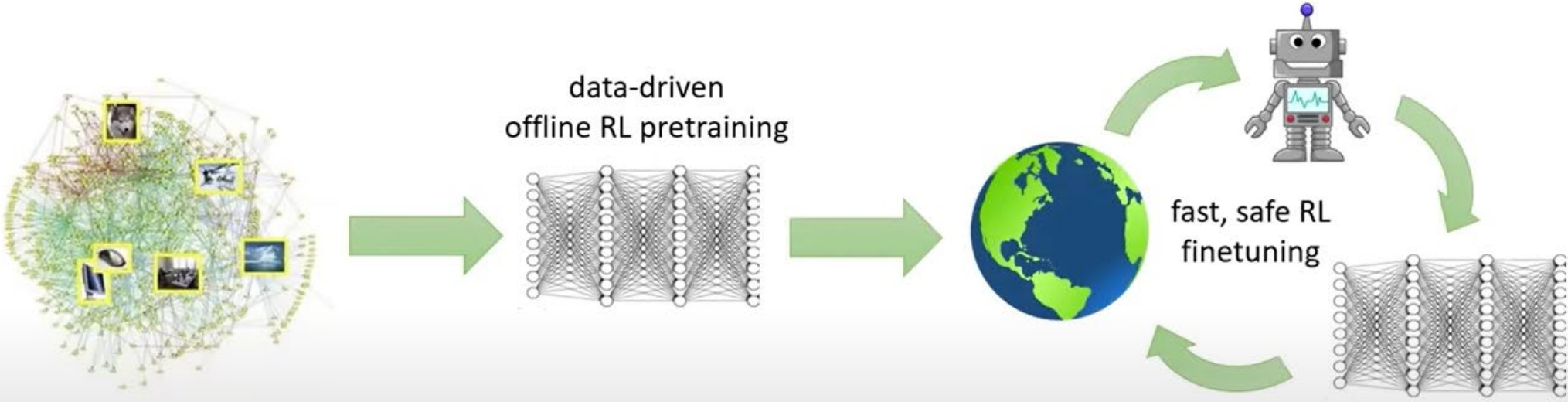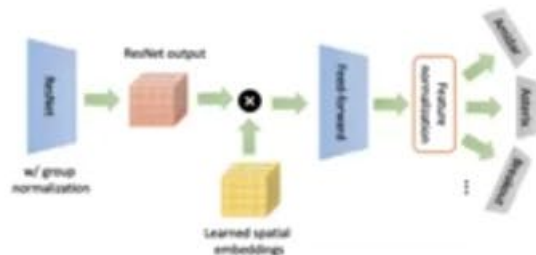- doesn't try to do **better** than the data

Reinforcement Learning

+ optimizes a goal with emergent behavior

- doesn't make use of real-world data

SGVR Lab
KAIST

# The Recipe?



data-driven
offline RL pretraining

fast, safe RL
finetuning

SGVR Lab
KAIST

# What can we accomplish when combining data and optimization?



Data-driven RL algorithms



Robotic foundation models and RL



RL with generative models

SGVR Lab
KAIST

# What do we need to figure out?
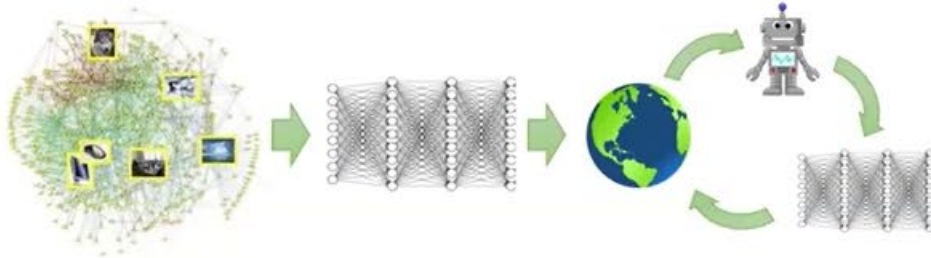


**Classic RL**

- Online, lifelong learning process
- Starts from scratch
- Largely **trial and error** driven
- **Central problems:**
  - Exploration
  - Sample efficiency
  - Optimization performance

**Data-driven RL**

- Offline pretraining + online finetuning
- Always start from data
- Largely **representation learning** driven
- **Central problems:**
  - Distributional shift
  - Scalability and stability
  - Representation learning with big models

SGVR Lab
KAIST

# To break this down..



Data-driven RL

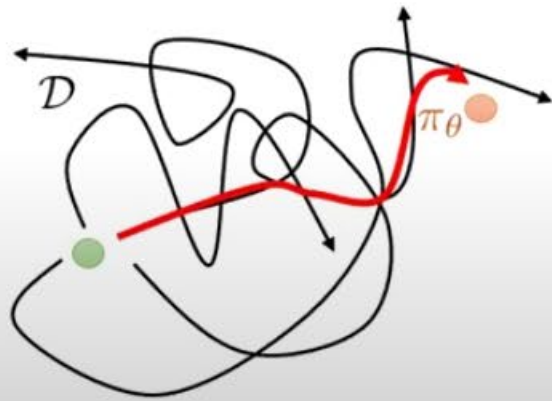1. Offline RL algorithms — We understand this pretty well

2. Online finetuning from offline initializations — We understand this a little

3. Making all this work with big, scalable models — We hardly understand this at all

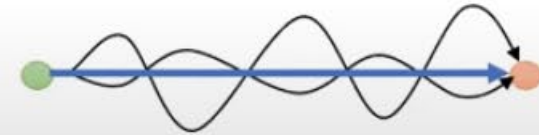- Offline pretraining + online finetuning
- Always start from data
- Largely **representation learning** driven
- **Central problems:**
  - Distributional shift
  - Scalability and stability
  - Representation learning with big models

SGVR Lab
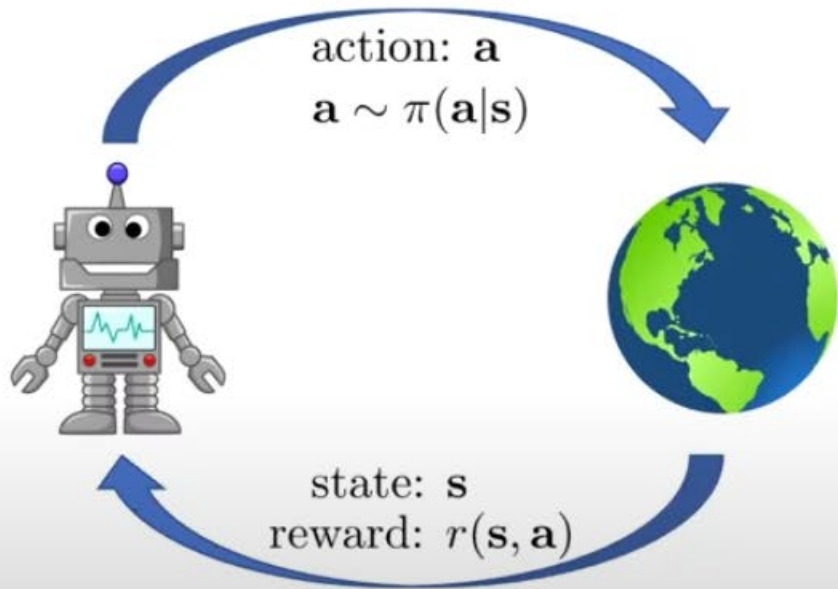KAIST

# What do we expect offline RL to do?



"Macro-scale" stitching

But this is just the clearest example!

"Micro-scale" stitching:

SGVR Lab
KAIST

# Off-policy RL: a quick primer



RL objective: $\max_{\pi} \sum_{t=1}^{T} E_{\mathbf{s}_t, \mathbf{a}_t \sim \pi}[r(\mathbf{s}_t, \mathbf{a}_t)]$

Q-function: $Q^{\pi}(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^{T} E_{\mathbf{s}_{t'}, \mathbf{a}_{t'} \sim \pi}[r(\mathbf{s}_{t'}, \mathbf{a}_{t'})|\mathbf{s}_t, \mathbf{a}_t]$

$\pi(\mathbf{a}|\mathbf{s}) = 1$ if $\mathbf{a} = \arg\max_{\mathbf{a}} Q^{\pi}(\mathbf{s}, \mathbf{a})$

$Q^{\star}(\mathbf{s}, \mathbf{a}) = r(\mathbf{s}, \mathbf{a}) + \max_{\mathbf{a}'} Q^{\star}(\mathbf{s}', \mathbf{a}')$

enforce this equation at all states!

$\text{minimize} \sum_i (Q(\mathbf{s}_i, \mathbf{a}_i) - [r(\mathbf{s}_i, \mathbf{a}_i) + \max_{\mathbf{a}_i'} Q(\mathbf{s}_i', \mathbf{a}_i')])^2$

action: $\mathbf{a}$
$\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})$

state: $\mathbf{s}$
reward: $r(\mathbf{s}, \mathbf{a})$

SGVR Lab
KAIST

# Some principles for offline RL

$$Q(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + \max_{\mathbf{a}'} Q(\mathbf{s}', \mathbf{a}')$$ (crossed out)

$$Q(\mathbf{s}, \mathbf{a}) \leftarrow \underbrace{r(\mathbf{s}, \mathbf{a}) + E_{\mathbf{a}' \sim \pi_{\text{new}}}[Q(\mathbf{s}', \mathbf{a}')]}_{y(\mathbf{s}, \mathbf{a})}$$

expect good accuracy when $\pi_\beta(\mathbf{a}|\mathbf{s}) = \pi_{\text{new}}(\mathbf{a}|\mathbf{s})$

even *worse*: $\pi_{\text{new}} = \arg\max_\pi E_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})}[Q(\mathbf{s}, \mathbf{a})]$

what is the objective?

$$\min_Q E_{(\mathbf{s}, \mathbf{a}) \sim \pi_\beta(\mathbf{s}, \mathbf{a})}\left[(Q(\mathbf{s}, \mathbf{a}) - y(\mathbf{s}, \mathbf{a}))^2\right]$$

behavior policy          target value

how often does *that* happen?

how well it does          how well it *thinks* it does (Q-values)
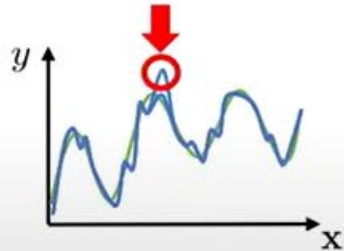


**"Stabilizing Off-Policy Q-Learning via Bootstrapping Error Reduction"**
**NeurIPS 2019**

12

SGVR Lab
KAIST

# Some principles for offline RL

➢ Many different methods, similar principles seem to be effective:

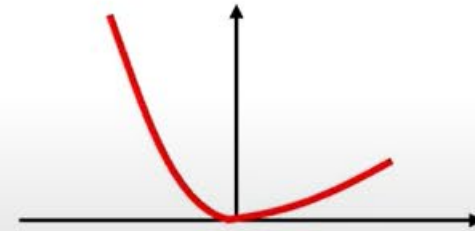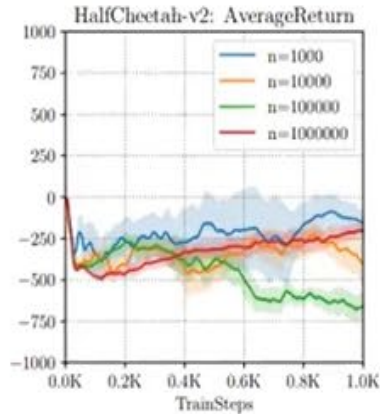| use value-based methods (i.e., Q-learning or Q-function actor-critic) | somehow fix the distributional shift problem |
|---|---|

$$D_{KL}(\pi(\mathbf{a}|\mathbf{s})\|\pi_\beta(\mathbf{a}|\mathbf{s})) \leq \epsilon$$

pessimism (e.g., CQL)         policy constraints (e.g., BRAC)         avoid OOD actions in updates
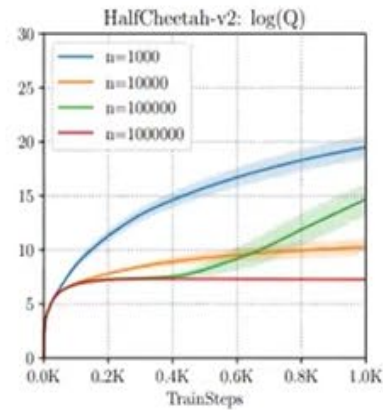                                                                        (e.g., AWAC, IQL)

SGVR Lab
KAIST

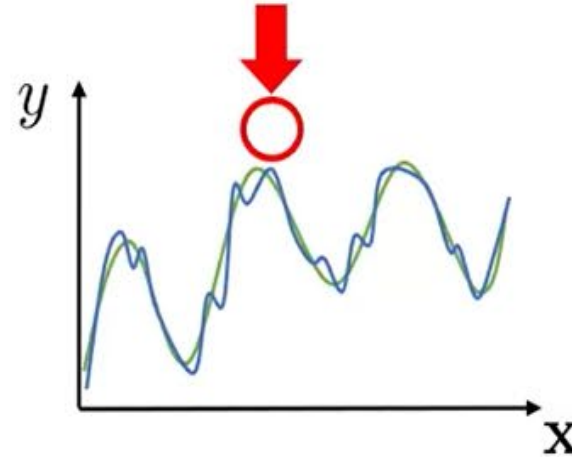# Conservative Q-learning



how well it does

how well it *thinks*
it does (Q-values)

$$\hat{Q}^\pi = \arg\min_Q \max_\pi \alpha E_{\mathbf{s}\sim D, \mathbf{a}\sim\pi(\mathbf{a}|\mathbf{s})}[Q(\mathbf{s},\mathbf{a})] - \alpha E_{\mathbf{s},\mathbf{a}\sim D}[Q(\mathbf{s},\mathbf{a})] \Big\}$$ term to push down big Q-values

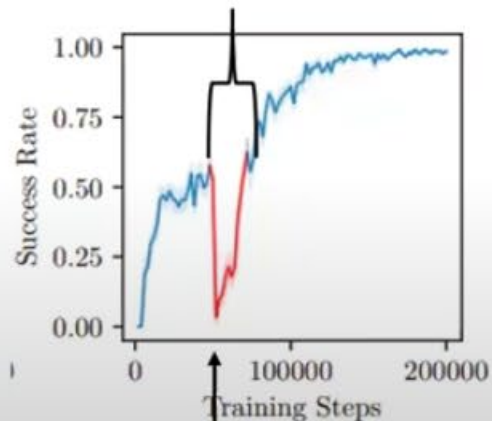regular objective $\Big\{ +E_{(\mathbf{s},\mathbf{a},\mathbf{s}')\sim D}\Big[(Q(\mathbf{s},\mathbf{a}) - (r(\mathbf{s},\mathbf{a}) + E_\pi[Q(\mathbf{s}',\mathbf{a}')]))^2\Big]$

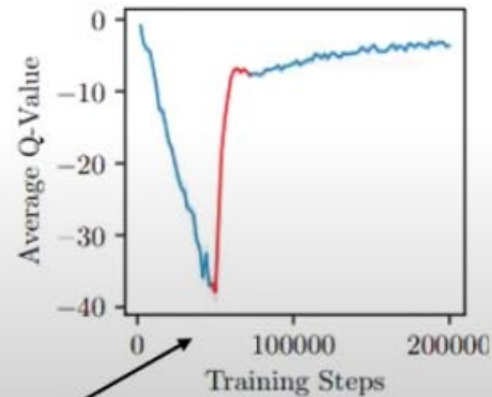**"Conservative Q-Learning for Offline Reinforcement Learning."**
**NeurIPS 2020**

# What about online finetuning?

$$\hat{Q}^\pi = \arg\min_Q \max_\pi \alpha E_{\mathbf{s}\sim D, \mathbf{a}\sim\pi(\mathbf{a}|\mathbf{s})}[Q(\mathbf{s},\mathbf{a})] - \alpha E_{\mathbf{s},\mathbf{a}\sim D}[Q(\mathbf{s},\mathbf{a})]$$

$$+ E_{(\mathbf{s},\mathbf{a},\mathbf{s}')\sim D}\left[(Q(\mathbf{s},\mathbf{a}) - (r(\mathbf{s},\mathbf{a}) + E_\pi[Q(\mathbf{s}',\mathbf{a}')]))^2\right]$$
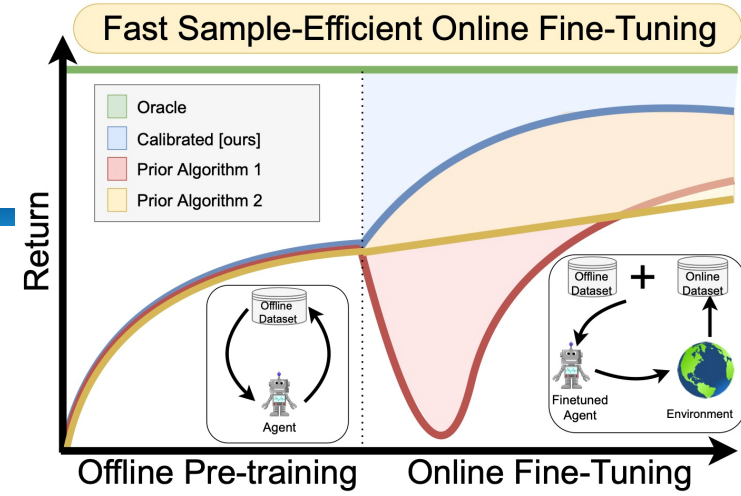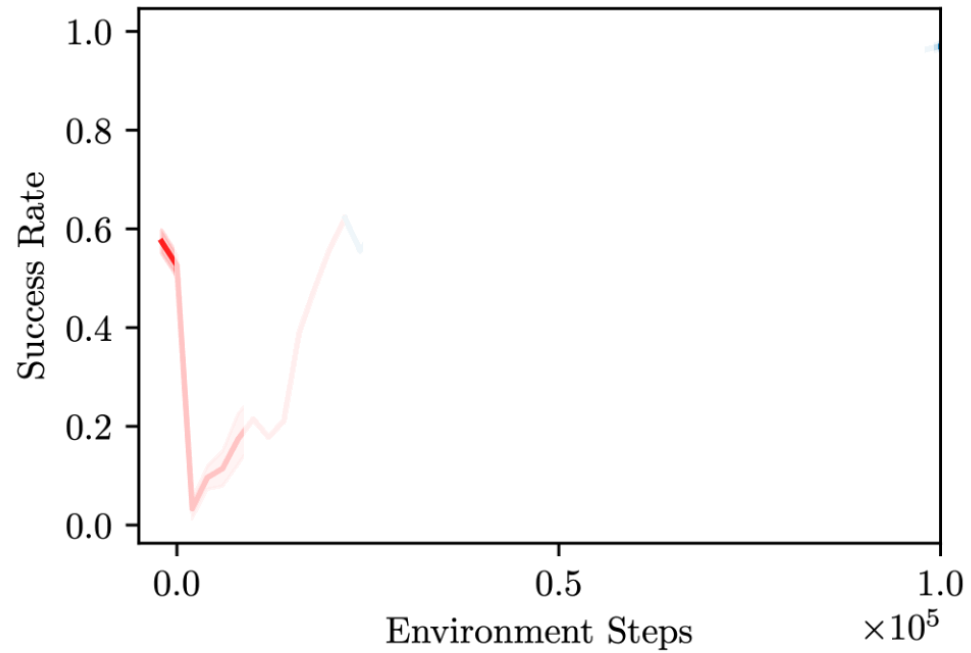
this period wasted recovering
offline performance

online training starts here (at 50k steps)

SGVR Lab
KAIST

# Example

Return

Oracle
Calibrated [ours]
Prior Algorithm 1
Prior Algorithm 2

Offline Pre-training    Online Fine-Tuning



**Policy Unlearning in Online Fine-Tuning**

SGVR Lab
KAIST

# Simple solution: Calibration

**Key idea:** need to put the offline-trained value function on the **right scale**

(we refer to this as *calibration*)

**Calibration:** learned values *upper bound* the values of *some* real policy

**Conservatism:** learned values *lower bound* the values of the *learned* policy

**How do we ensure that our Q-function is calibrated?**

before:

$$\min_Q \max_\pi \alpha E_{\mathbf{s}\sim D, \mathbf{a}\sim\pi(\mathbf{a}|\mathbf{s})}[Q(\mathbf{s}, \mathbf{a})]$$

now:

$$\min_Q \max_\mu \alpha E_{\mathbf{s}\sim D, \mathbf{a}\sim\pi(\mathbf{a}|\mathbf{s})}[\max\{Q(\mathbf{s}, \mathbf{a}), V^\mu(\mathbf{s})\}]$$

$$V^{\pi_\beta}(\mathbf{s}_t) \approx \sum_{t'=t}^{\infty} \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{a}_{t'})$$

how to get this?

SGVR Lab
KAIST

# The method: Cal-QL

**Calibration:** learned values *upper bound* the values of *some* real policy
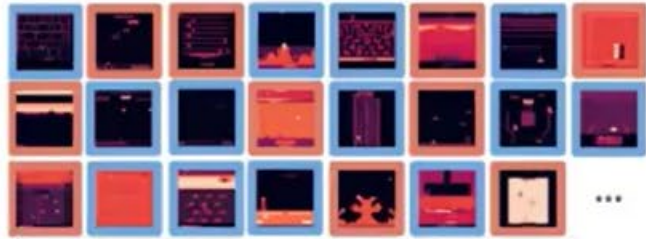
**Conservatism:** learned values *lower bound* the values of the *learned* policy
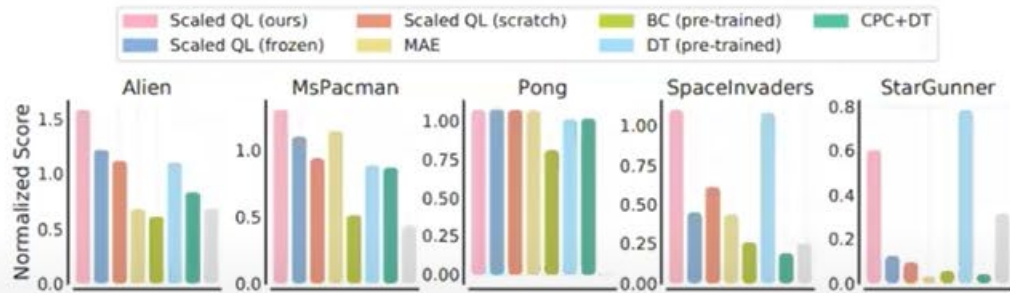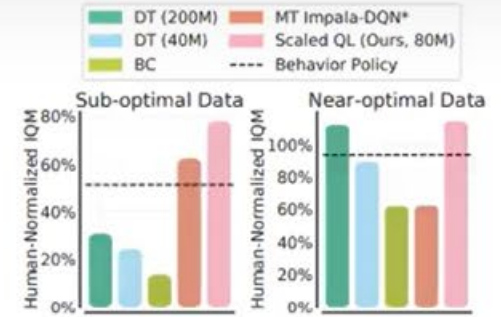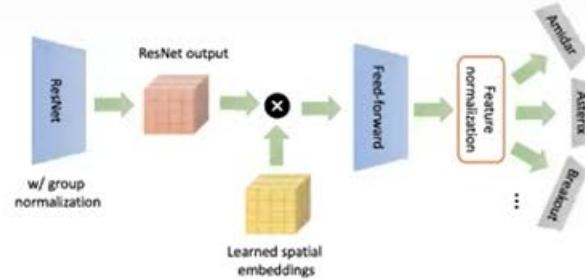
ensures "conservatism"          ensures "calibration"

$$\hat{Q}^{\pi} = \arg\min_{Q}\max_{\pi} \alpha E_{\mathbf{s}\sim D, \mathbf{a}\sim\pi(\mathbf{a}|\mathbf{s})}[\max\{Q(\mathbf{s},\mathbf{a}), V^{\pi_{\beta}}(\mathbf{s})\}] - \alpha E_{\mathbf{s},\mathbf{a}\sim D}[Q(\mathbf{s},\mathbf{a})]$$

$$+ E_{(\mathbf{s},\mathbf{a},\mathbf{s}')\sim D}\left[(Q(\mathbf{s},\mathbf{a}) - r(\mathbf{s},\mathbf{a}) + E_{\pi}[Q(\mathbf{s}',\mathbf{a}')])^2\right]$$

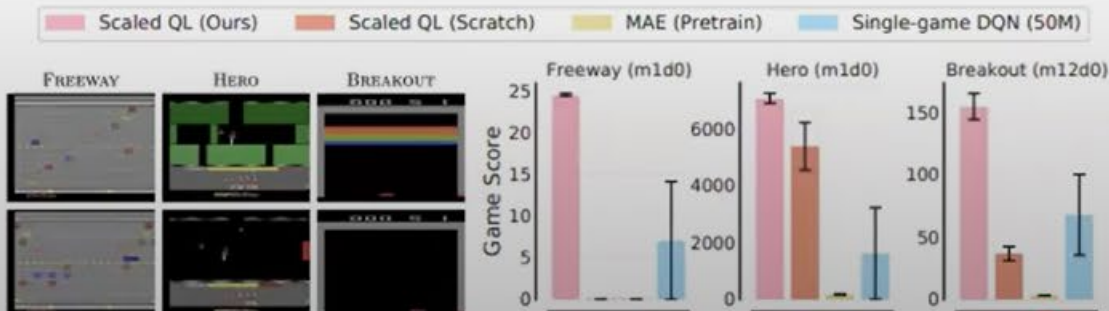Cal-QL: Calibrated Offline RL Pre-Training for Efficient Online Fine-Tuning, NeurIPS 23

SGVR Lab
KAIST

**18**

# How about working with big models?



Offline Q-Learning on Diverse Multi-Task
Data Both Scales and Generalizes, ICLR 23

# The representation learning mystery

- Seems to be annoyingly hard to make this work with large transformer models
- Seems to require larger models with more capacity than we might expect (from, e.g., imitation learning)

**Something about RL (i.e., TD learning) seems "harder" than supervised learning**

instability

action **resampled** from data distribution

$$Q(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + \gamma E_{\mathbf{a}' \sim \pi(\mathbf{a}|\mathbf{s})}[Q(\mathbf{s}', \mathbf{a}')]$$

$$Q(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + \gamma Q(\mathbf{s}', \mathbf{a}')$$

action is always **exactly** the action in the dataset

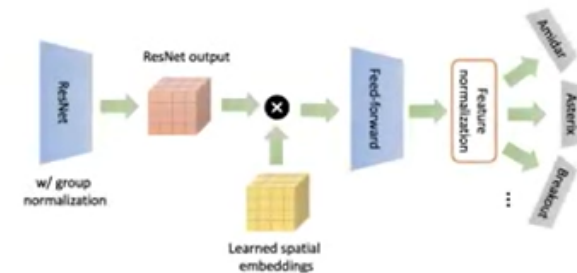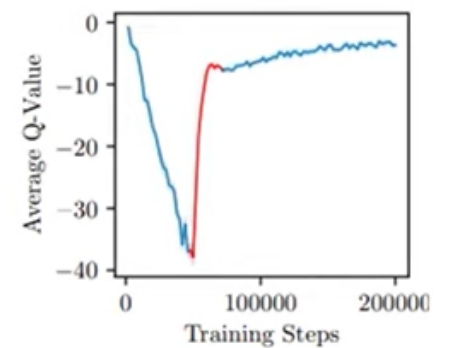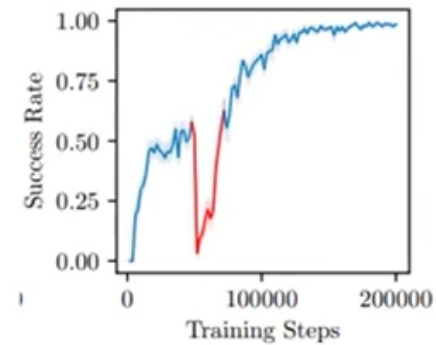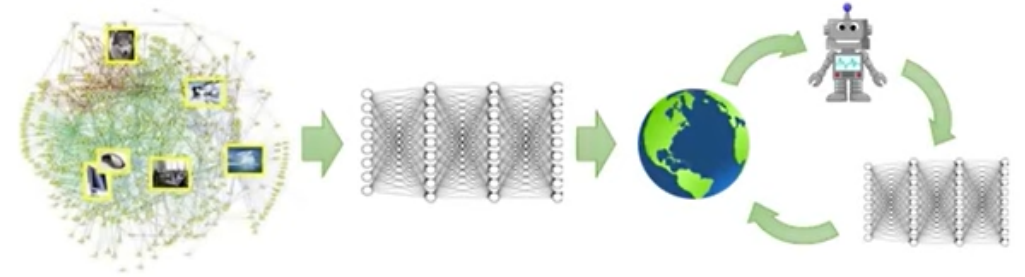$\phi(\mathbf{s}, \mathbf{a})^T \phi(\mathbf{s}', \mathbf{a}')$

TD learning updates value functions based on the difference between the predicted value and a new, partially observed value—hence, "temporal difference."

**"DR3: Value-Based Deep** Reinforcement Learning Requires Explicit Regularization" NeurIPS 2021

20

# Summary and takeaways

- Offline RL is an essential component of data-driven RL
  - We must handle the distributional shift between the offline data distribution and the new policy



- Online RL finetuning from offline initializations presents new challenges
  - We must be able to finetune via online RL without losing the benefits of the offline initialization



- Doing this with large models presents yet more challenges
  - Harder to make RL algorithms as scalable as supervised learning algorithms

SGVR Lab
KAIST

# Next Time

- **Robotics foundation models**

SGVR Lab
KAIST

# Homework

- **Come up with one question on what we have discussed today**
  - Write a question two times before the mid-term exam

- **Browse two papers**
  - Submit their summaries online before the Mon. Class

SGVR Lab
KAIST