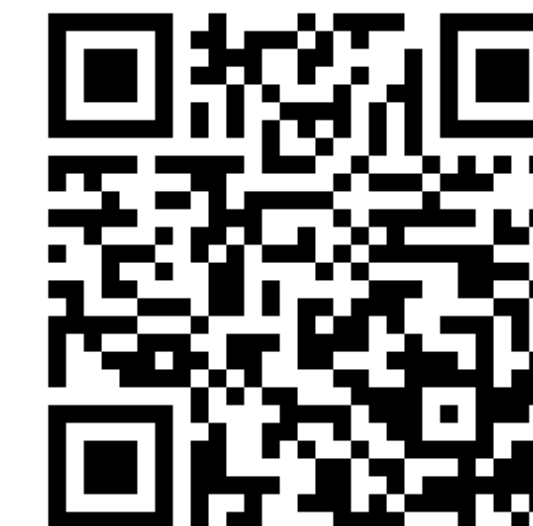


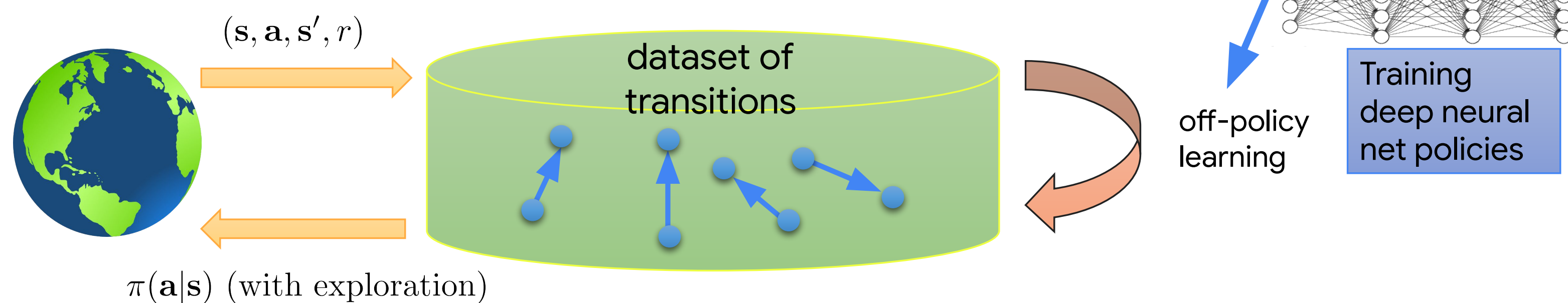
Implicit Under-Parameterization Inhibits Data Efficient Deep Reinforcement Learning

Aviral Kumar*, Rishabh Agarwal*, Dibya Ghosh, Sergey Levine



Paper: <https://openreview.net/pdf?id=O9bnihsFfXU>
 Video: <https://www.youtube.com/watch?v=dgnpGI2iNw8>

Modern Deep RL Algorithms



Q-Learning

1. Train Q-functions by minimizing TD Error:

$$E_{(s,a) \sim \pi_{\beta}(s,a)} [(Q_{\phi}(s,a) - (r(s,a) + \gamma E[Q_{\phi}(s',a')]))^2]$$
2. Collect new data in the environment by rolling out the learned policy.

Typically solved "approximately" using gradient descent for a fixed number of steps

How does this approximate optimization procedure affect learning?

Implicit Under-Parameterization

TD Error (Regressing to itself)

$$E_{(s,a) \sim \pi_{\beta}(s,a)} [(Q_{\phi}(s,a) - (r(s,a) + \gamma E[Q_{\phi}(s',a')]))^2]$$

Formalizing Implicit under-parameterization

$$Q_{\phi}(s,a) = \mathbf{w}^T \Phi_{\phi}(s,a) \quad \Phi_{\phi}(s,a) \in \mathbb{R}^{|S| \times |A| \times d}$$

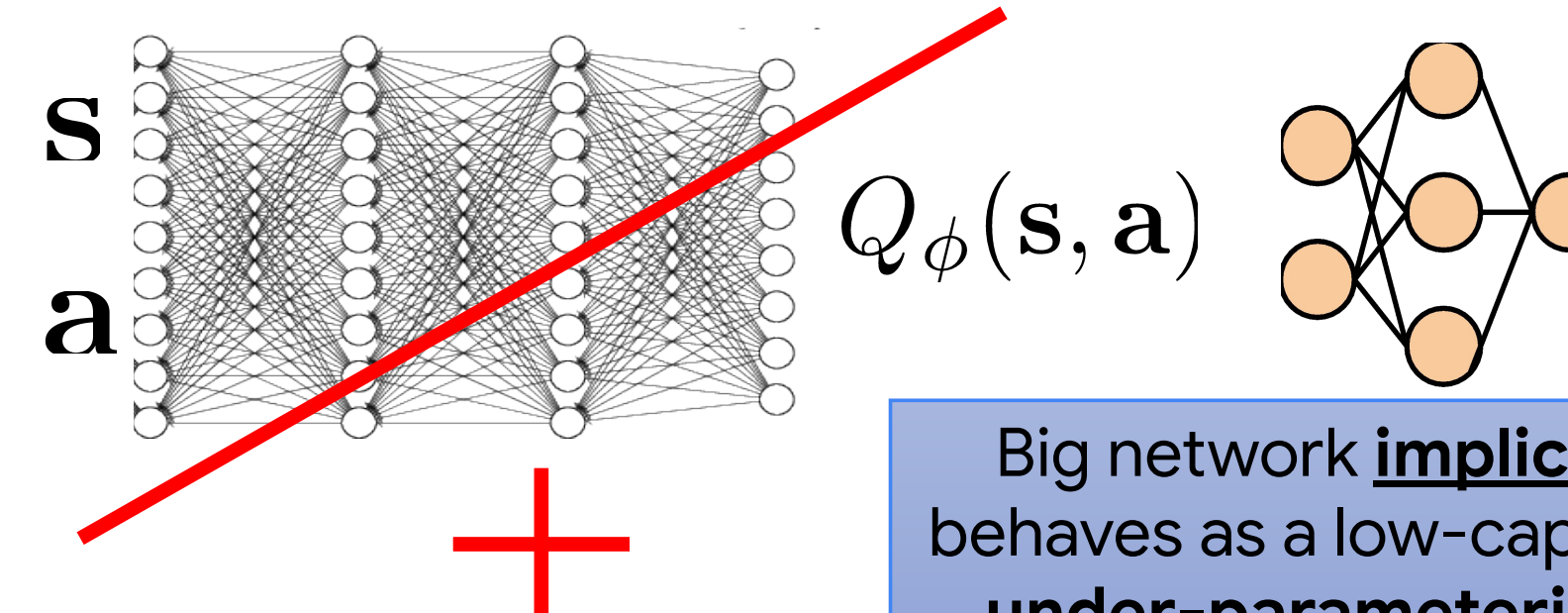
Learned features

Effective rank

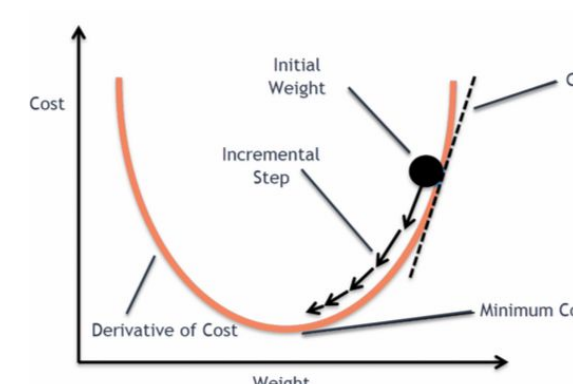
$$\text{srnk}_{\delta}(\Phi) = \min \left\{ k : \frac{\sum_{i=1}^k \sigma_i(\Phi)}{\sum_{i=1}^N \sigma_i(\Phi)} \geq 1 - \delta \right\}$$

Soft notion of rank of the matrix

Low rank => more aliasing => poor performance



Big network **implicitly** behaves as a low-capacity, **under-parameterized** network



Gradient descent optimizer

Data-Efficient Deep Reinforcement Learning

Data-Efficient Deep RL: Want to learn the most per unit amount of experience/data

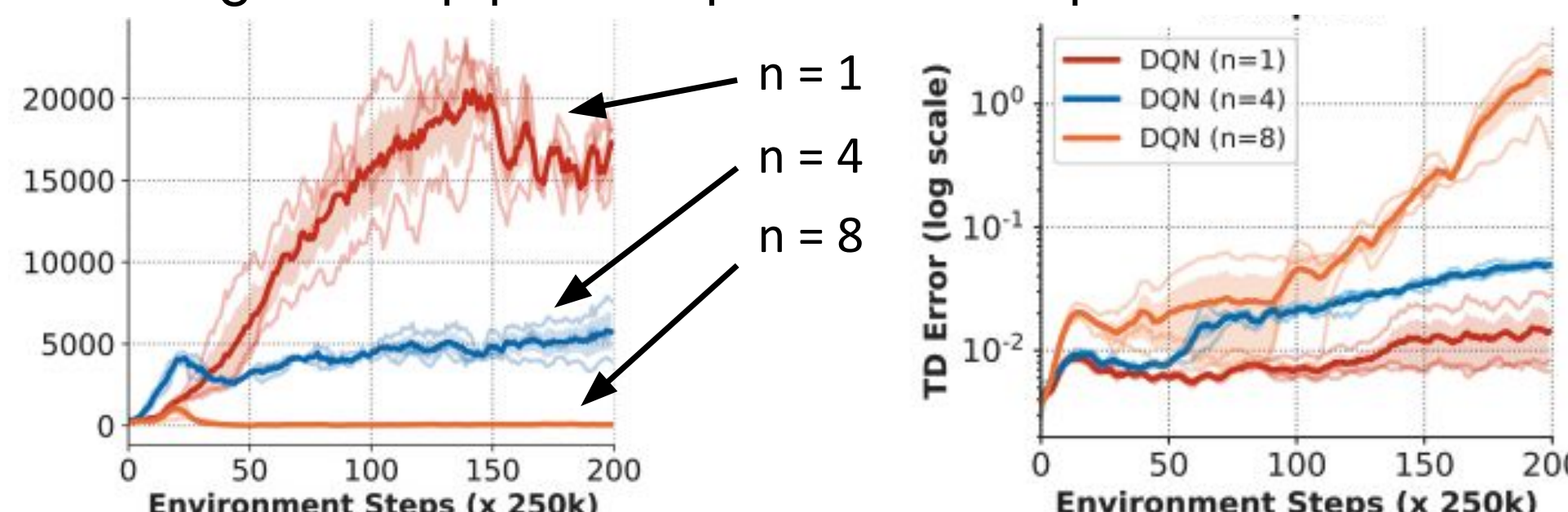
How can we obtain data-efficiency?

Reinforcement Learning:

Training ≥ 1 step per datapoint leads to poor



Typically 100-200 updates per datapoint



OK, maybe I need to prevent overfitting?

But training error is high with larger n

Supervised Learning:

Train more, control for statistical overfitting
 Train error = 0, validation error = high

Why do we see "underfitting" with more training?

What Causes Implicit Under-Parameterization?

Gradient descent leads to low rank solutions

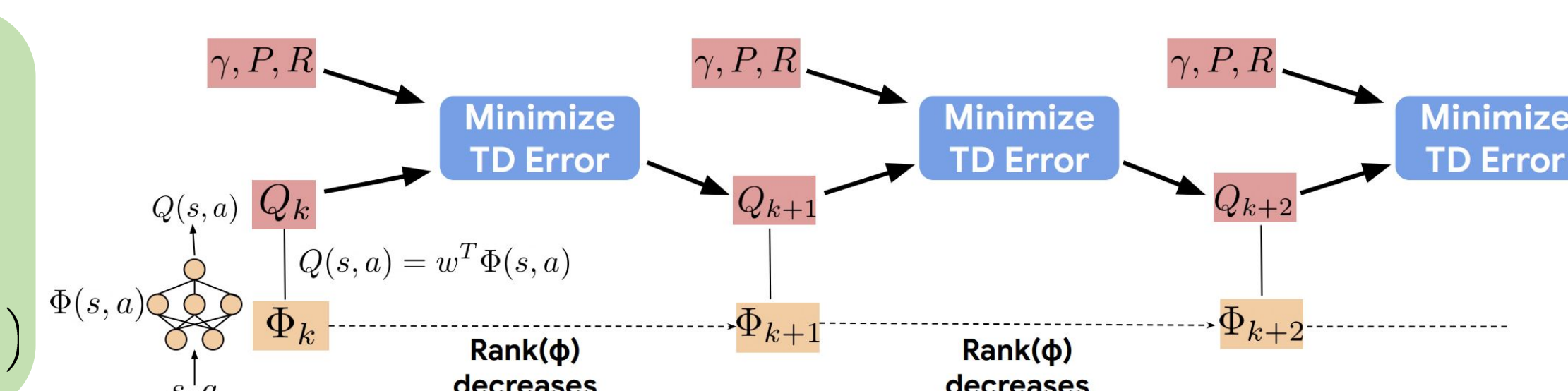
$$\min_X \|AX - y\|_2^2$$

$$\min_X \|AX - y\|_2^2 + \lambda f(X)$$

$f(X)$ makes large singular values larger.

$$\frac{\sigma_{\text{large}}}{\sigma_{\text{small}}} \uparrow \text{ over training}$$

Rank decrease effect due to supervised learning gets compounded due to bootstrapping



Bounded increase, that goes down with more training

$$\forall l \in \mathbb{N} \text{ and } t \geq k_l, \frac{\sigma_i(\mathbf{M}_t)}{\sigma_j(\mathbf{M}_t)} < \frac{\sigma_i(\mathbf{M}_{k_l})}{\sigma_j(\mathbf{M}_{k_l})} + \mathcal{O}\left(\left(\frac{\sigma_i(\mathbf{S})}{\sigma_j(\mathbf{S})}\right)^{k_l}\right)$$

Singular value ratio at time t

Singular value ratio some time before t

Fu*, Kumar* et al. Diagnosing Bottlenecks in Deep Q-Learning Algorithms. ICML 2019.
 Fedus*, Ramachandran*, Agarwal et al. Revisiting Fundamentals of Experience Replay. ICML 2020.

..analysis with kernel regression and deep linear nets in the paper