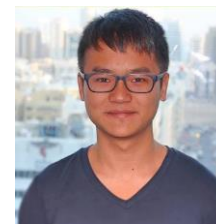


Greed is Bad or Better Exploration with Optimistic Actor Critic

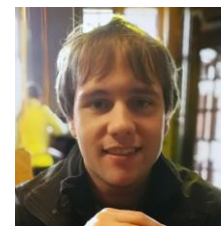
Kamil Ciosek
Microsoft Research Cambridge



Kamil Ciosek
MSR Cambridge



Quan Vuong
PhD student, UCSD



Robert Loftin
MSR Cambridge



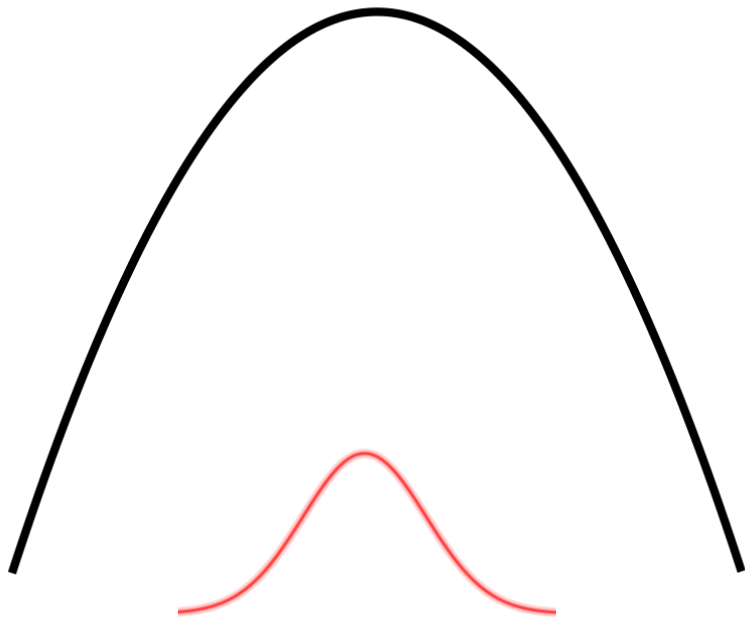
Katja Hofmann
MSR Cambridge

Policy Gradients are greedy

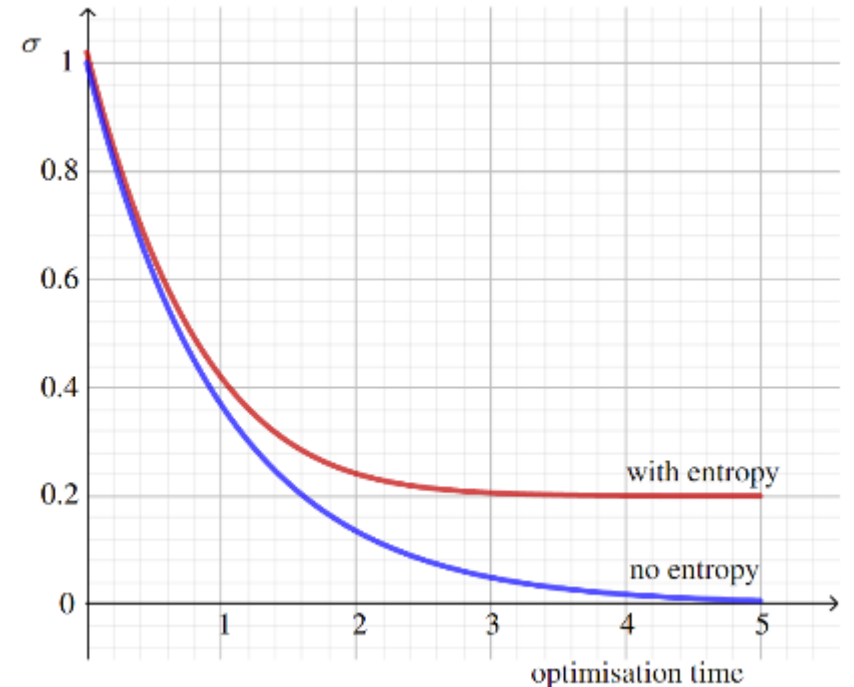
Policy gradients are *greedy*

Maximise $\nabla_{\theta} J$

What happens to the policy standard deviation?



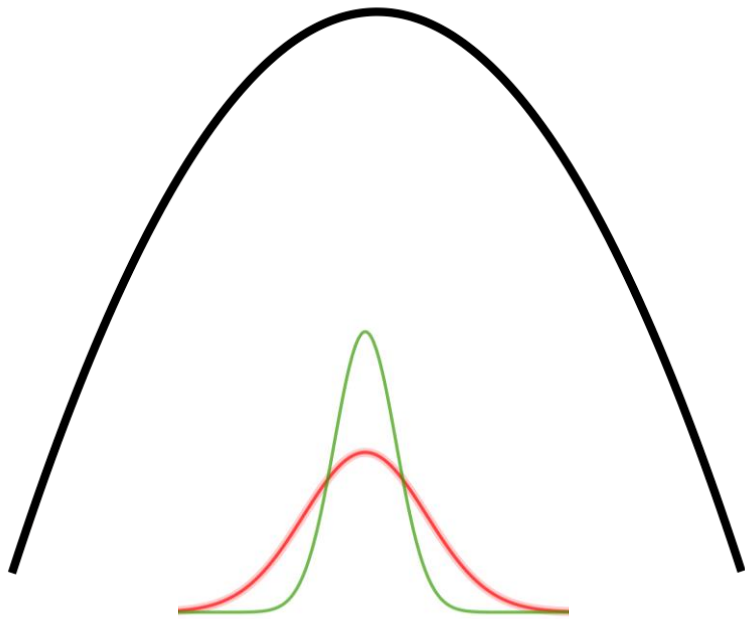
Consider a bandit with quadratic reward.



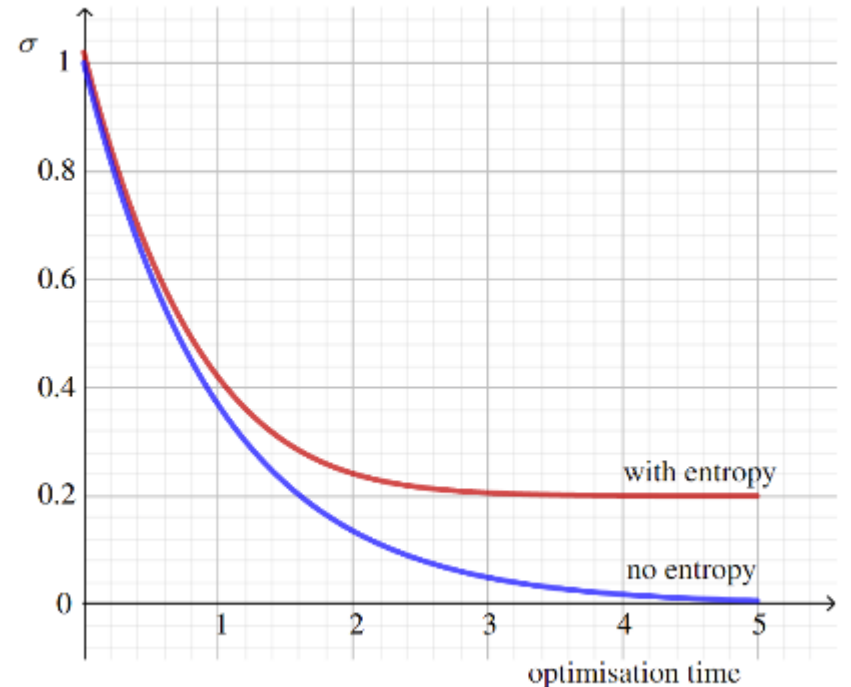
Policy gradients are *greedy*

Maximise $\nabla_{\theta} J$

What happens to the policy standard deviation?



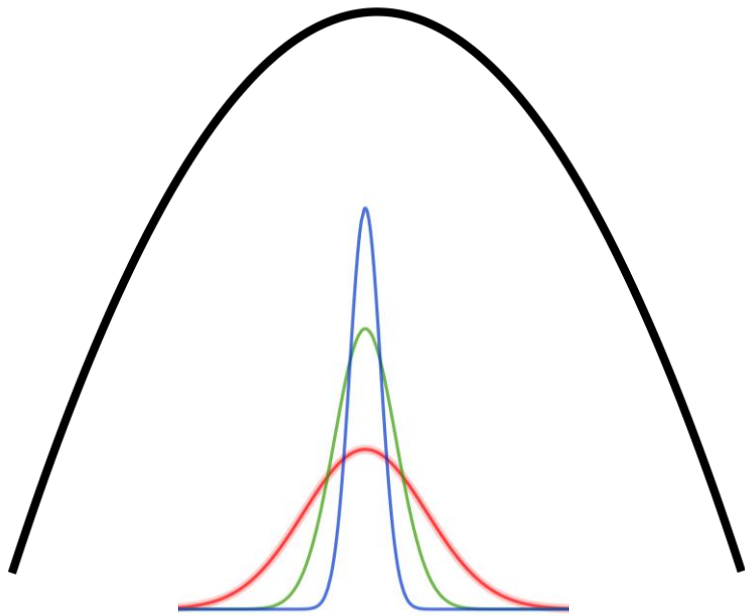
Consider a bandit with quadratic reward.



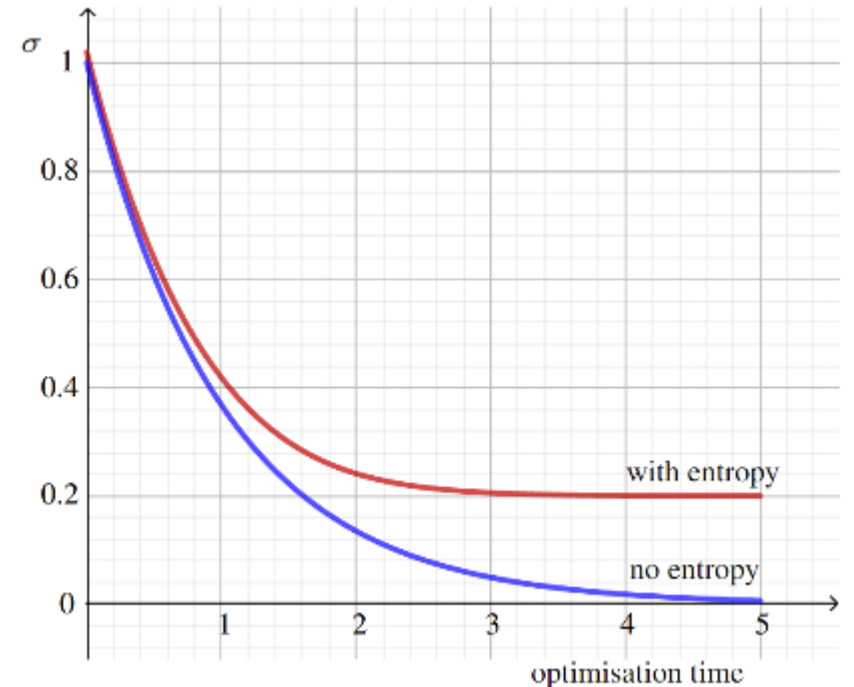
Policy gradients are *greedy*

Maximise $\nabla_{\theta} J$

What happens to the policy standard deviation?



Consider a bandit with quadratic reward.



Modern Policy Gradient
Methods use a Lower bound

Lower bound on critic

if this is too large...

$$\hat{Q}(s_t, a_t) \leftarrow R(s_t, a_t) + \gamma \check{Q}(s_{t+1}, a) \quad a \sim \pi_T(\cdot | s_{t+1})$$

...this becomes too large

+ effect amplified by policy optimisation

Lower bound on critic

if this is too large...

$$\hat{Q}(s_t, a_t) \leftarrow R(s_t, a_t) + \gamma \check{Q}(s_{t+1}, a) \quad a \sim \pi_T(\cdot | s_{t+1})$$

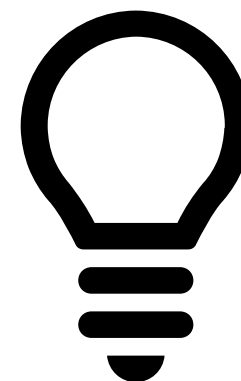
...this becomes too large

+ effect amplified by policy optimisation



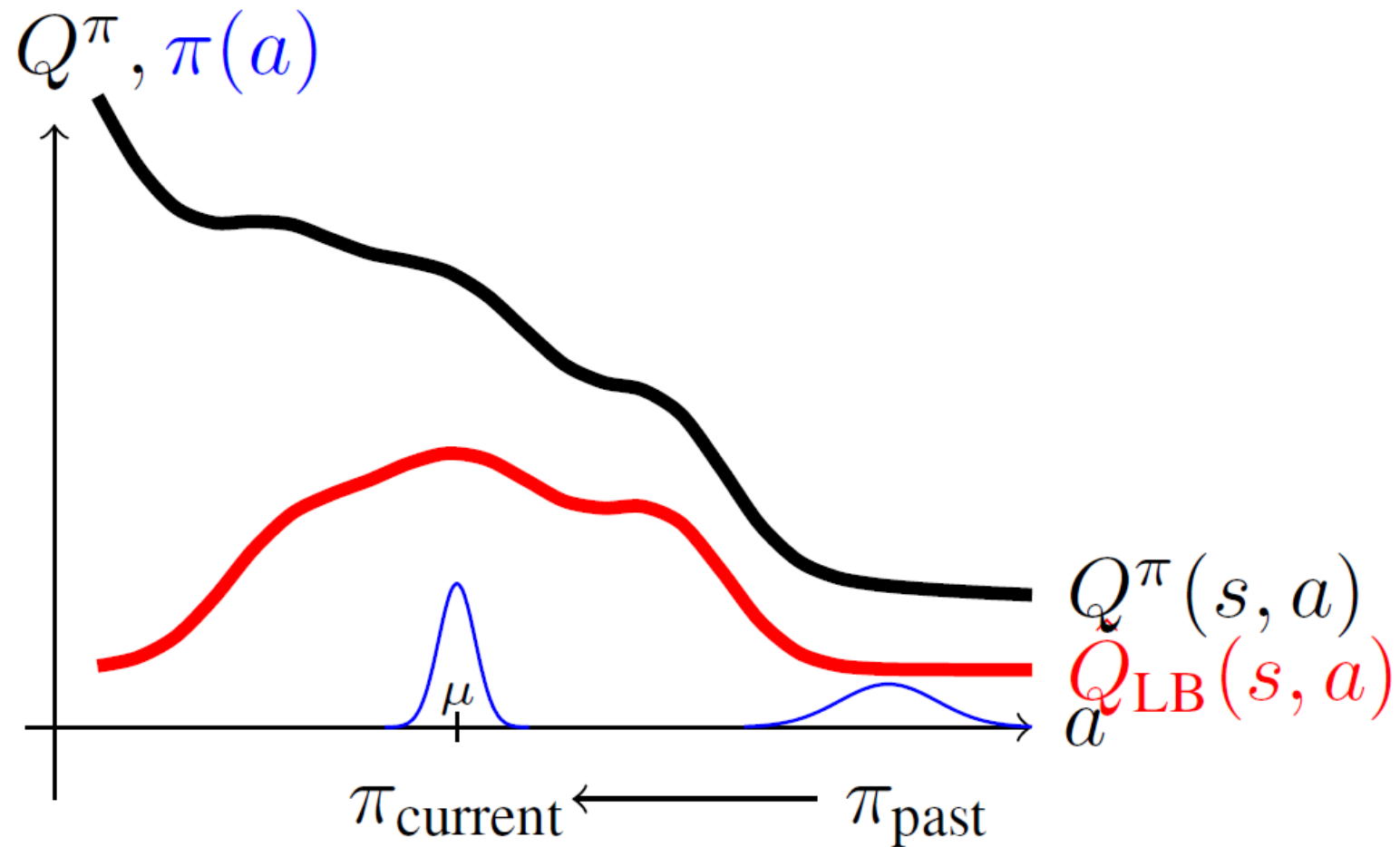
$$\hat{Q}_{\text{LB}}^{\{1,2\}}(s_t, a_t) \leftarrow R(s_t, a_t) + \gamma \min(\check{Q}_{\text{LB}}^1(s_{t+1}, a), \check{Q}_{\text{LB}}^2(s_{t+1}, a))$$

conservative update reduces overestimation

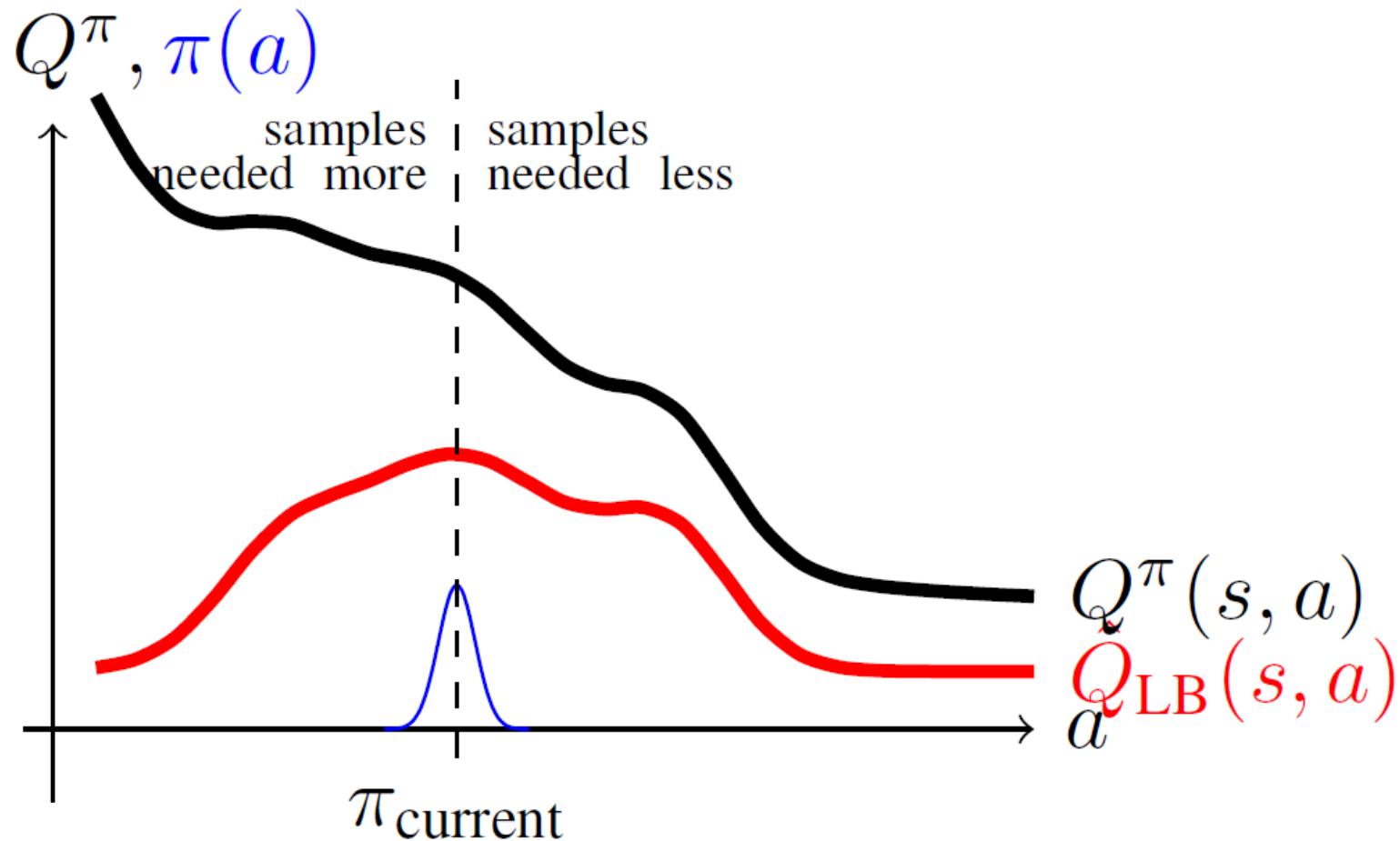


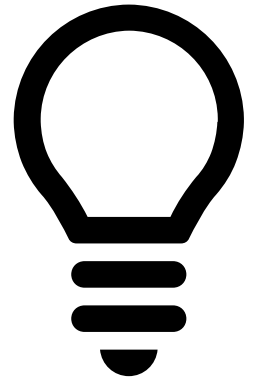
Greediness + Lower Bound
Lead To Problems

First problem: pessimistic underexploration



Second problem: directional uninformedness



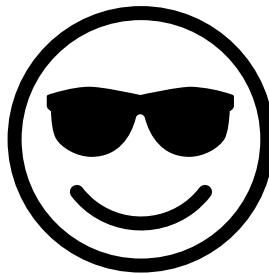


Solve these Problems by
Exploring with Upper Bound

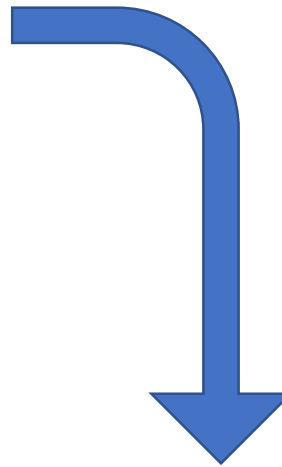
Use the bootstrap to make an upper bound.

$$\mu_Q(s, a) = \frac{1}{2} \left(\hat{Q}_{\text{LB}}^1(s, a) + \hat{Q}_{\text{LB}}^2(s, a) \right)$$

$$\sigma_Q(s, a) = \sqrt{\sum_{i \in \{1, 2\}} \frac{1}{2} \left(\hat{Q}_{\text{LB}}^i(s, a) - \mu_Q(s, a) \right)^2}$$



level of optimism

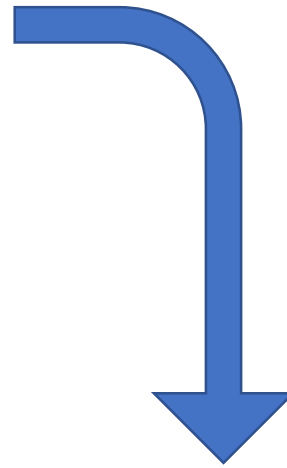


$$\hat{Q}_{\text{UB}}(s, a) = \mu_Q(s, a) + \beta_{\text{UB}} \sigma_Q(s, a)$$

How to choose the exploration policy

We want a policy that:

- Is close to target policy.
- Maximises the critic upper bound.

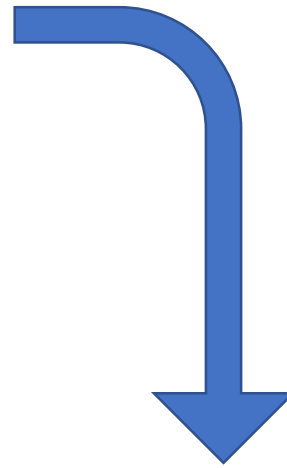


$$\mu_e, \Sigma_E = \underset{\substack{\mu, \Sigma: \\ \text{KL}(\mathcal{N}(\mu, \Sigma), \mathcal{N}(\mu_T, \Sigma_T)) \leq \delta}}{\arg \max} E_{a \sim \mathcal{N}(\mu, \Sigma)} [\hat{Q}_{\text{UB}}(s, a)]$$

How to choose the exploration policy

We want a policy that:

- Is close to target policy.
- Maximises the critic upper bound.



$$\mu_e, \Sigma_E = \underset{\substack{\mu, \Sigma: \\ \text{KL}(\mathcal{N}(\mu, \Sigma), \mathcal{N}(\mu_T, \Sigma_T)) \leq \delta}}{\arg \max} E_{a \sim \mathcal{N}(\mu, \Sigma)} [\bar{Q}_{\text{UB}}(s, a)]$$

Linearize!

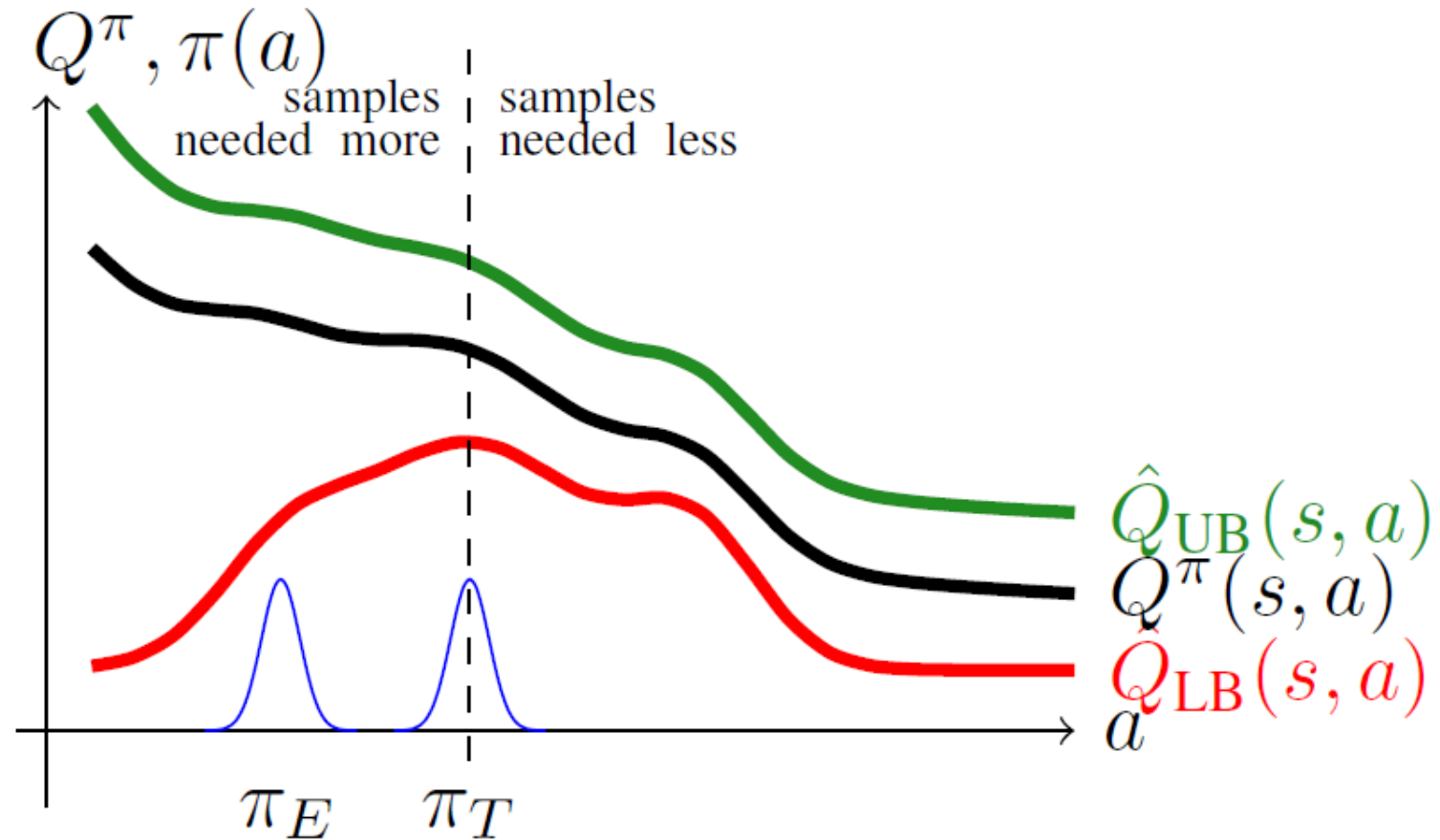
The OAC exploration policy (interpretation)

$$\pi_E = \mathcal{N}(\mu_E, \Sigma_E), \quad \mu_E = \mu_T + \underbrace{\frac{\sqrt{2\delta}}{\left\| \left[\nabla_a \hat{Q}_{UB}(s, a) \right]_{a=\mu_T} \right\|_{\Sigma_T}}}_{\text{shift}} \Sigma_T \left[\nabla_a \hat{Q}_{UB}(s, a) \right]_{a=\mu_T} \quad \text{and} \quad \Sigma_E = \Sigma_T.$$

OAC explores with a shifted policy!

shift in the direction given by upper bound.

OAC explores efficiently

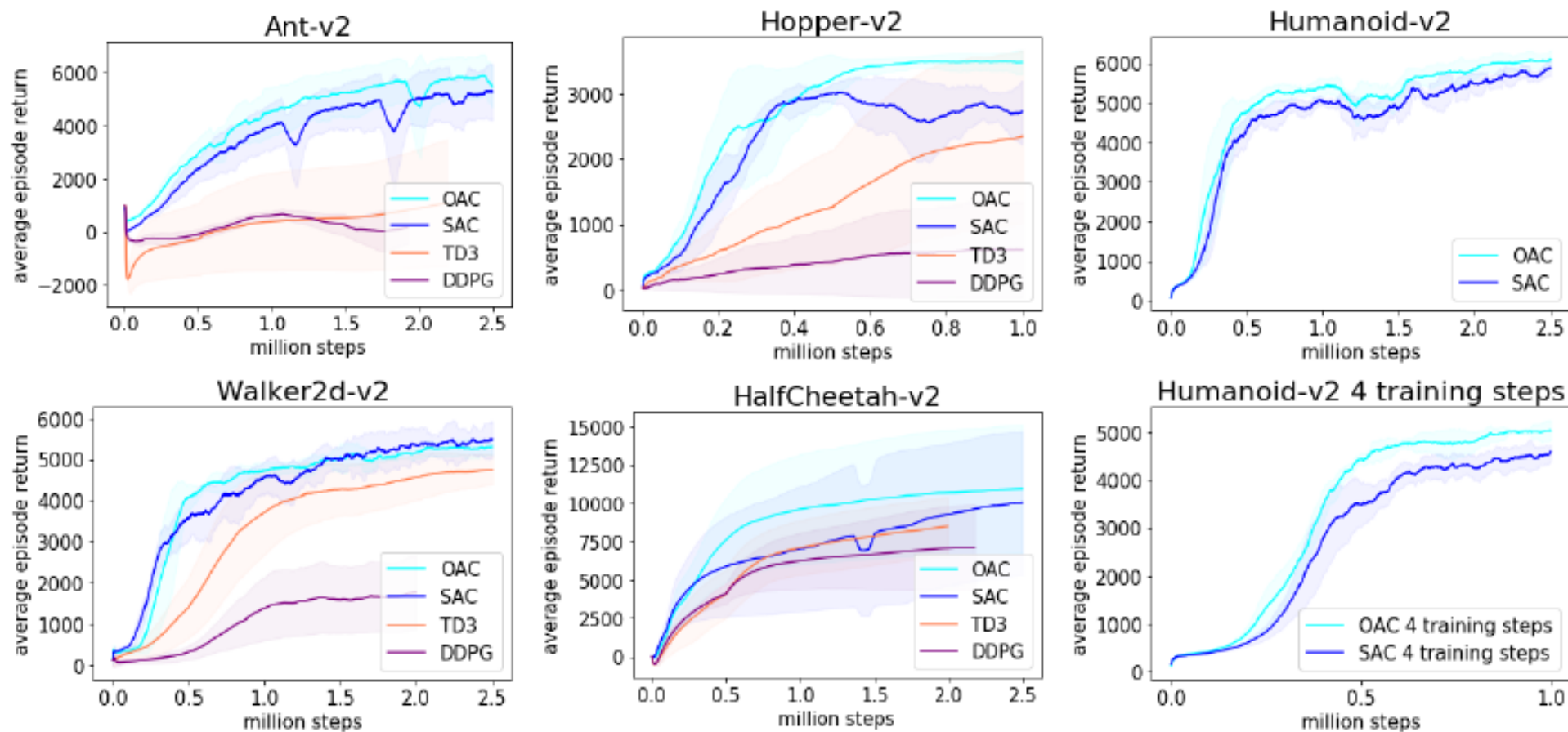


OAC avoids spurious maximum

OAC is directionally informed

It works!

It works!



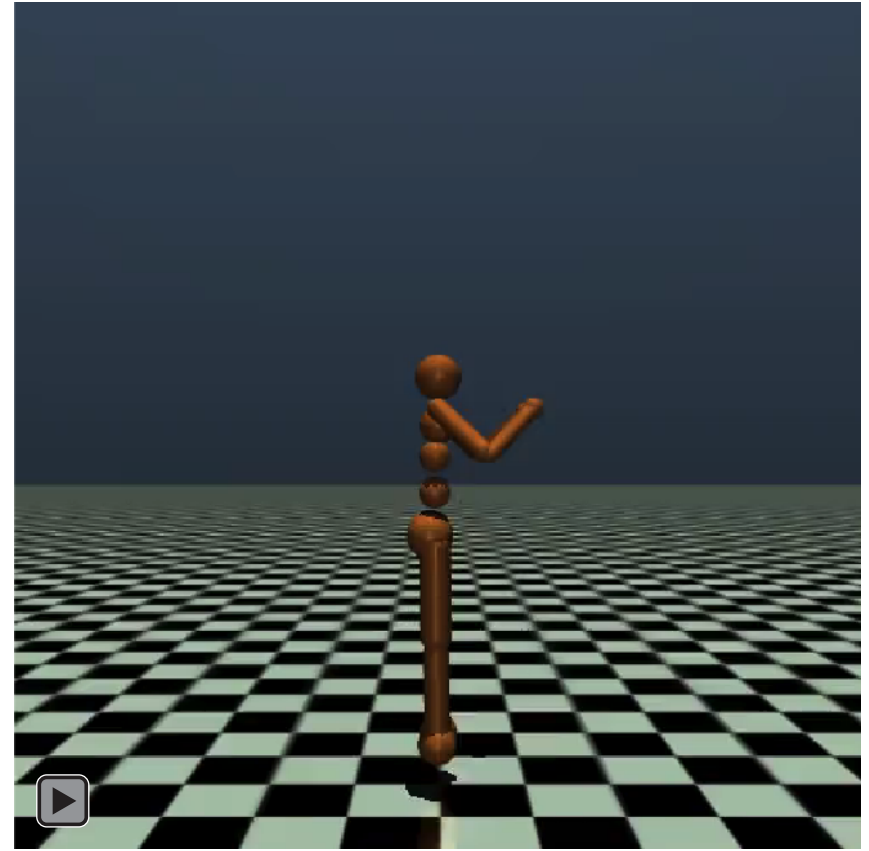
No hyperparameters were tuned on Humanoid!

Visual Comparison

SAC (previous state of the art)



OAC (our approach)





We have openings for interns, post-docs, researchers.

Kamil.Ciosek@Microsoft.com

Talk to me at the poster session!

Poster #179, starting at 5:30PM, East Exhibition Hall B+C