

# Bellman rank and Exploration with Function Approximation

# 3 core challenges of RL

Long-term planning

Generalization

Exploration

# 3 core challenges of RL

✓ Long-term planning

Approximate DP



✓ Generalization

Exploration X

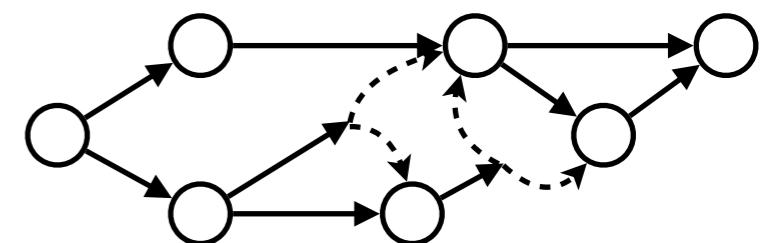
# 3 core challenges of RL

✓ Long-term planning

Approximate DP



PAC-MDP



✗ Generalization

Exploration ✓

# 3 core challenges of RL

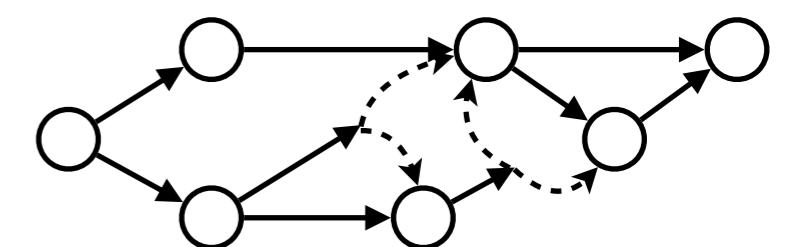
Long-term planning

Approximate DP



?

PAC-MDP



Generalization

Exploration

# 3 core challenges of RL

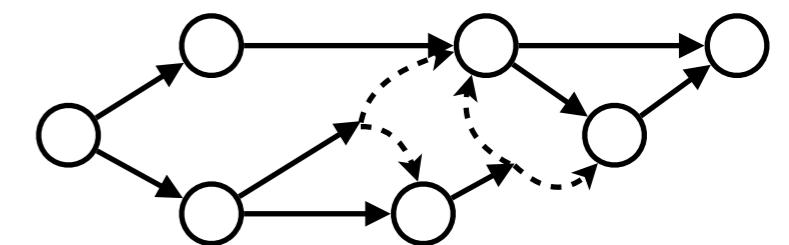
(Dynamic Programming)

Long-term planning

Approximate DP



PAC-MDP



?

Generalization  
(Supervised Learning)

Exploration  
(Multi-Armed Bandit)

# 3 core challenges of RL

*Bellman equation*

(Dynamic Programming)

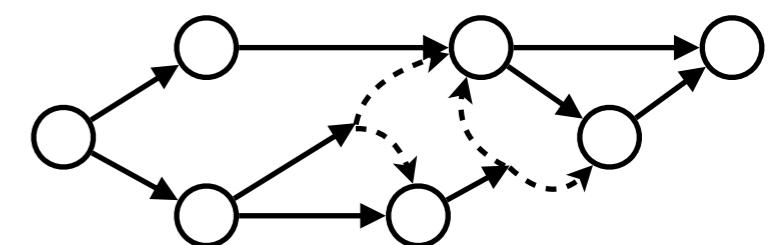
Long-term planning

Approximate DP



?

PAC-MDP

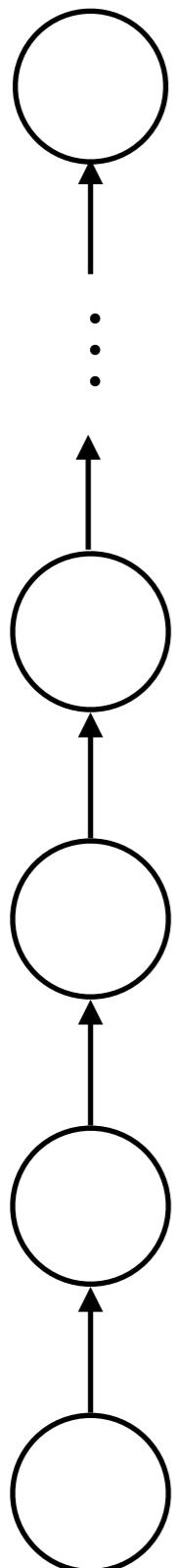


Generalization  
(Supervised Learning)

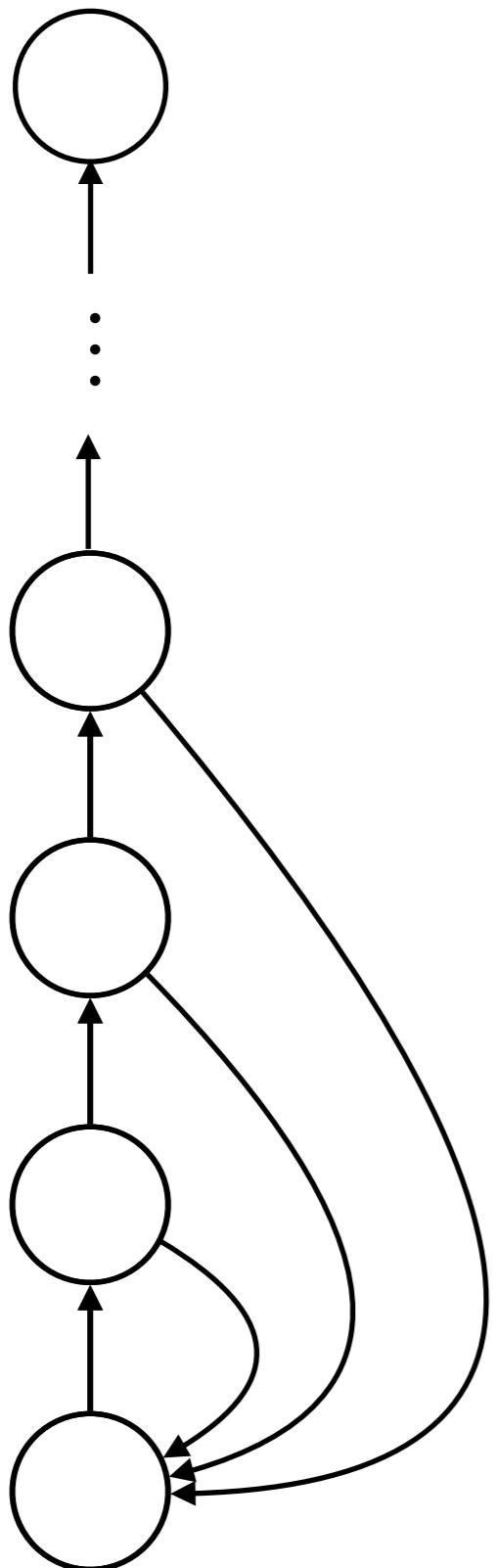
*Statistical complexity*  
(e.g., *VC-dimension*)

Exploration  
(Multi-Armed Bandit)  
*Optimism in face  
of uncertainty*

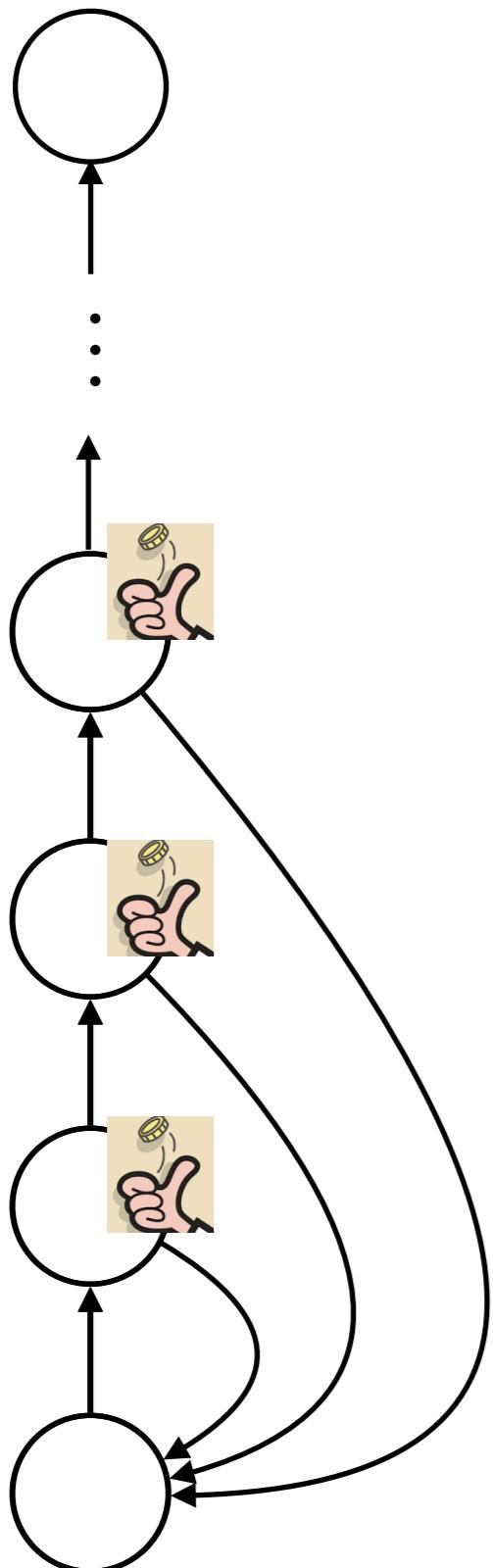
Random exploration can be inefficient



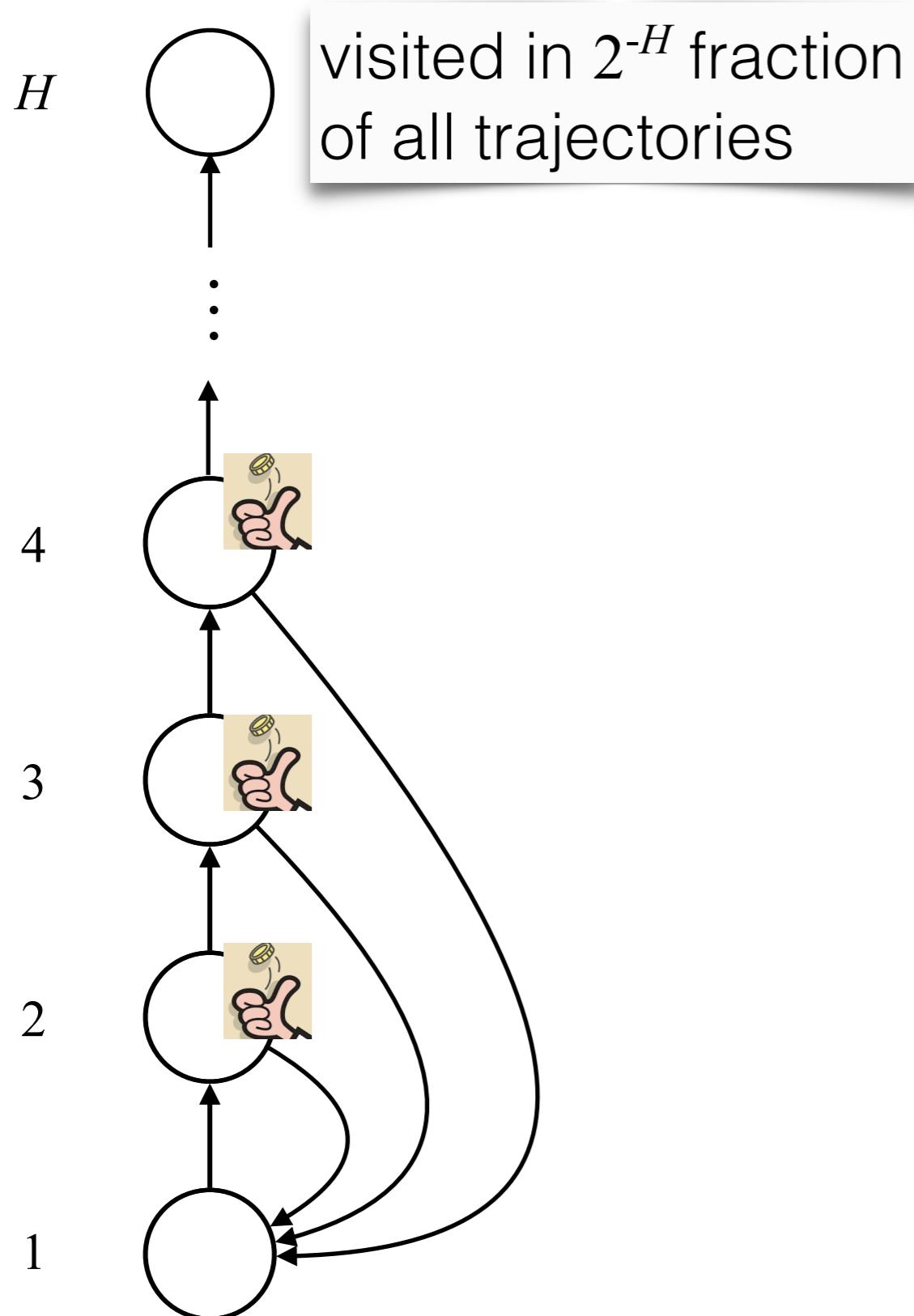
Random exploration can be inefficient



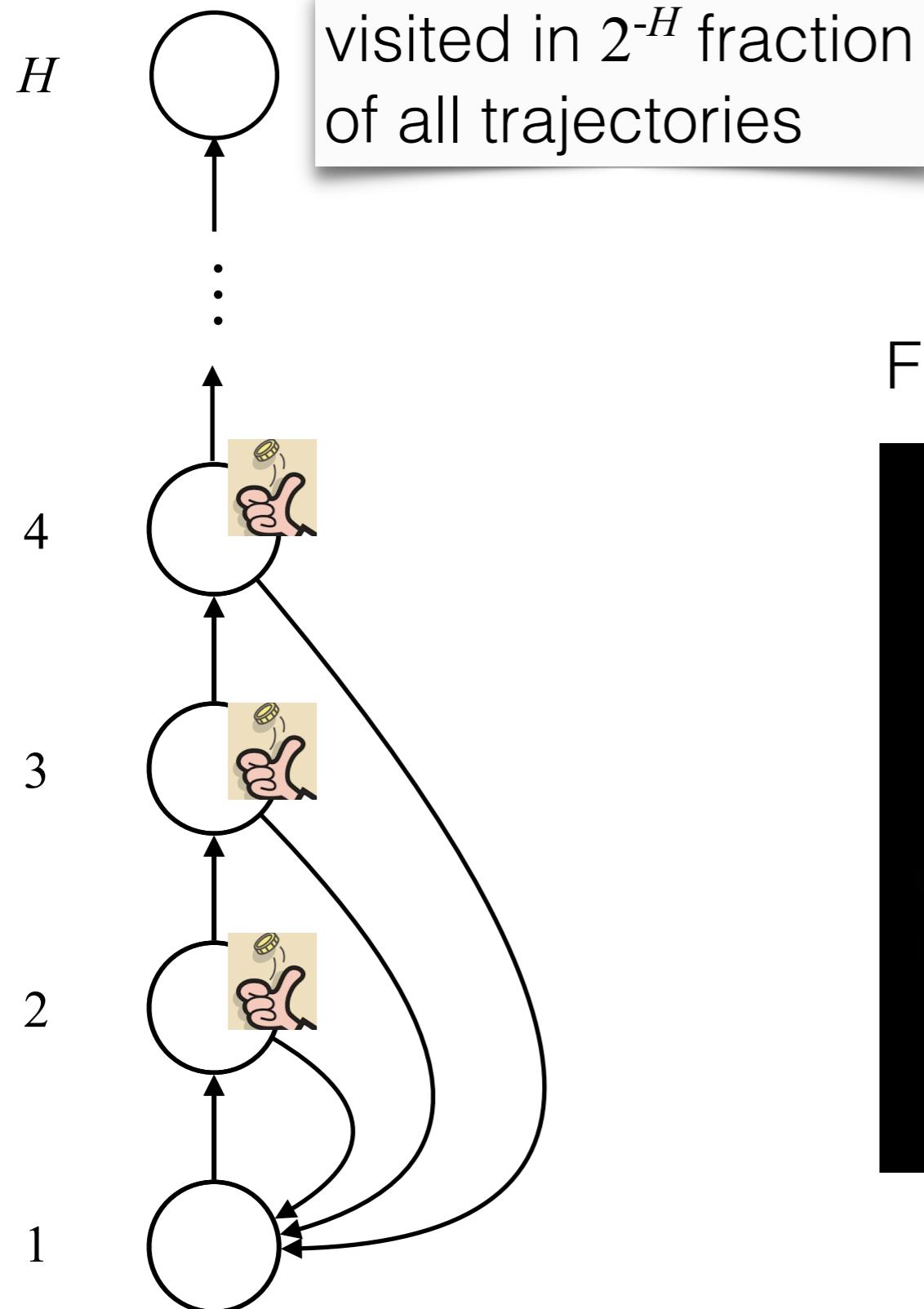
# Random exploration can be inefficient



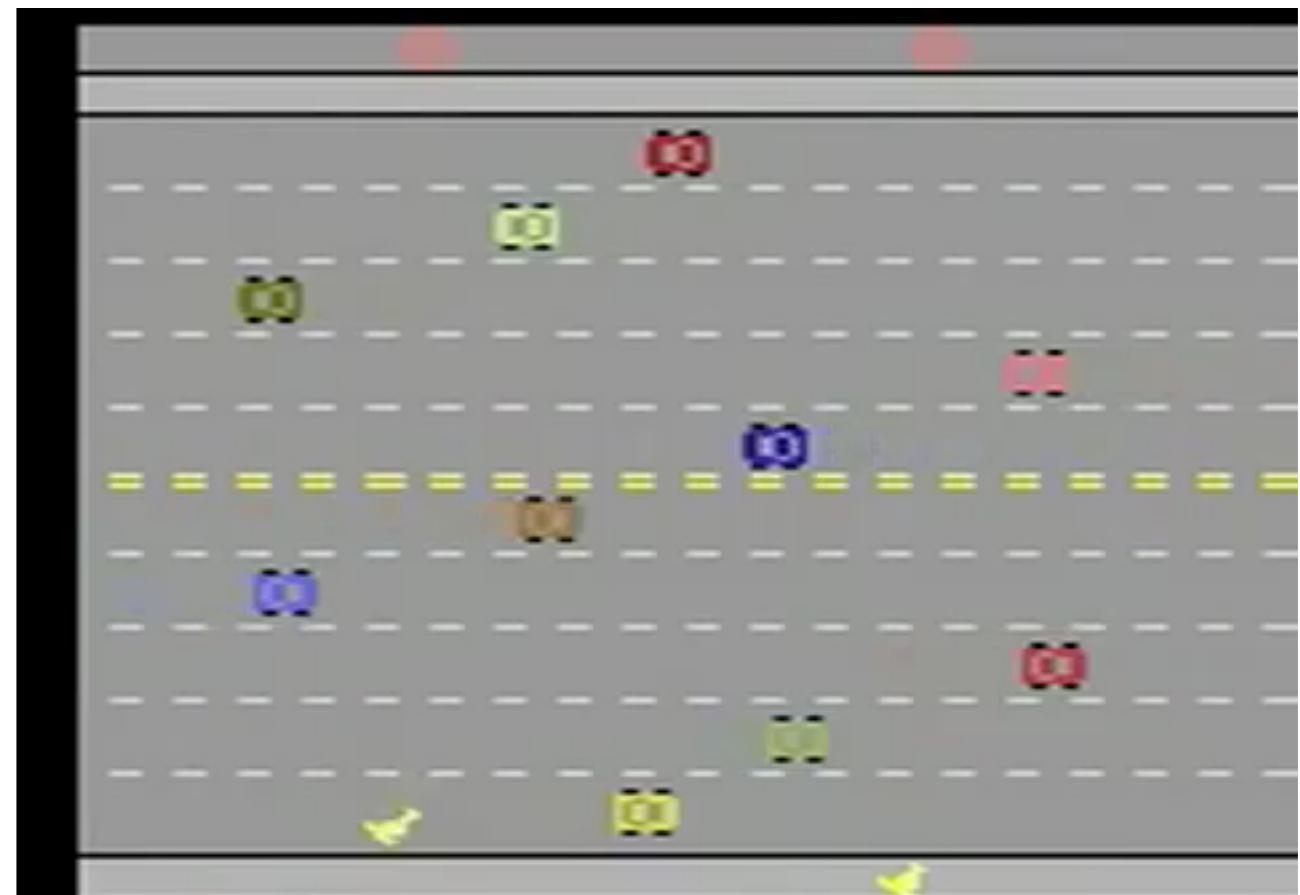
# Random exploration can be inefficient



# Random exploration can be inefficient

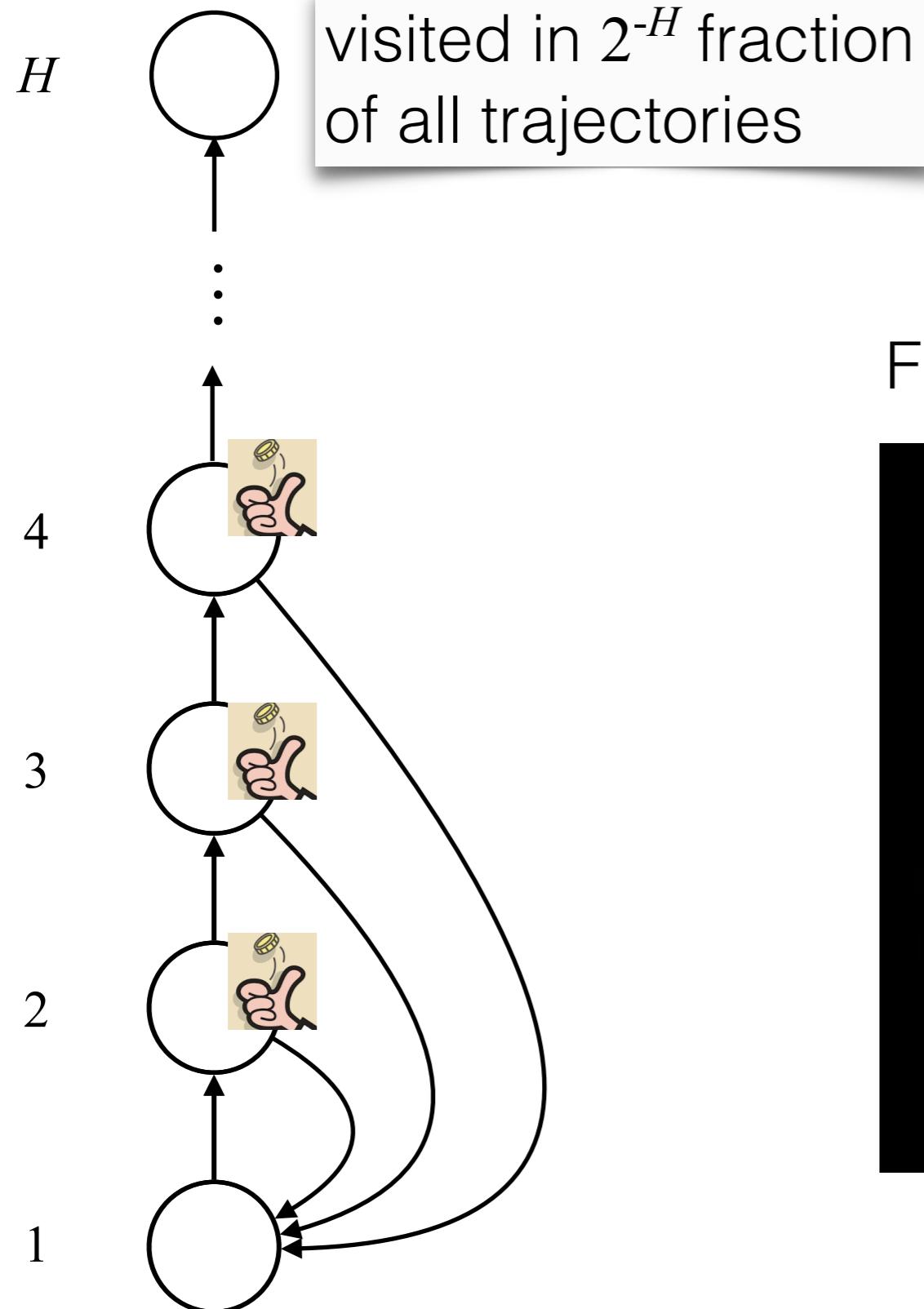


Freeway (one of the Atari games)

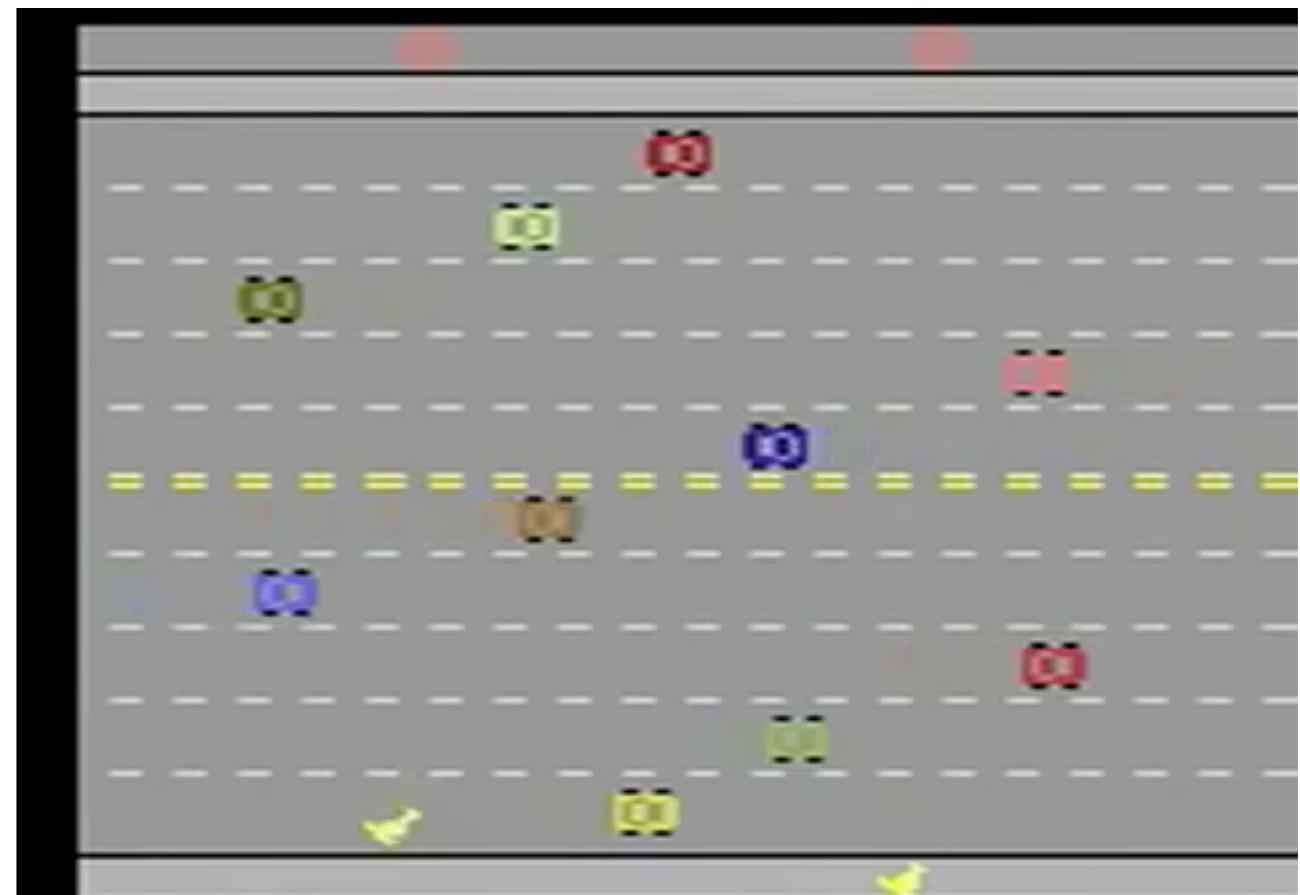


“Freeway + RL”: <https://youtu.be/44CiIPmlimQ>

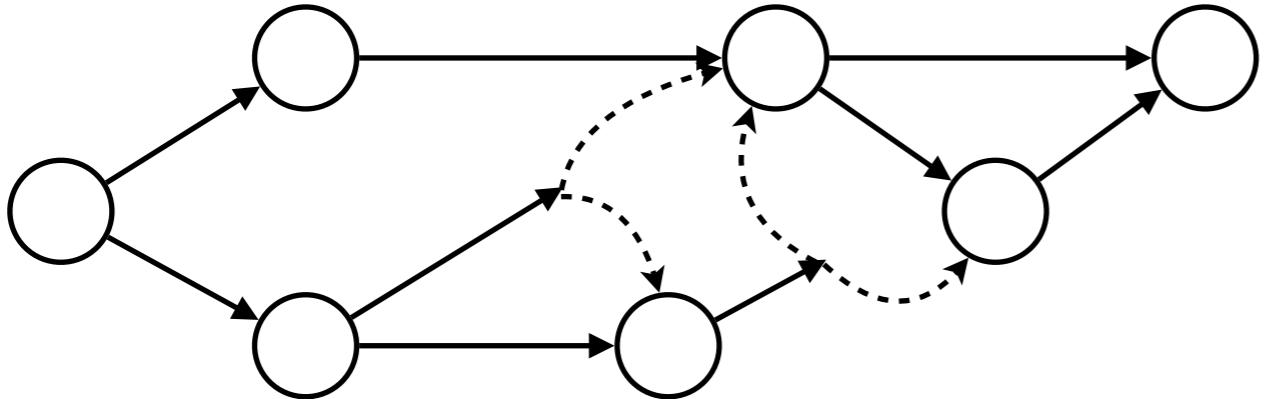
# Random exploration can be inefficient



Freeway (one of the Atari games)

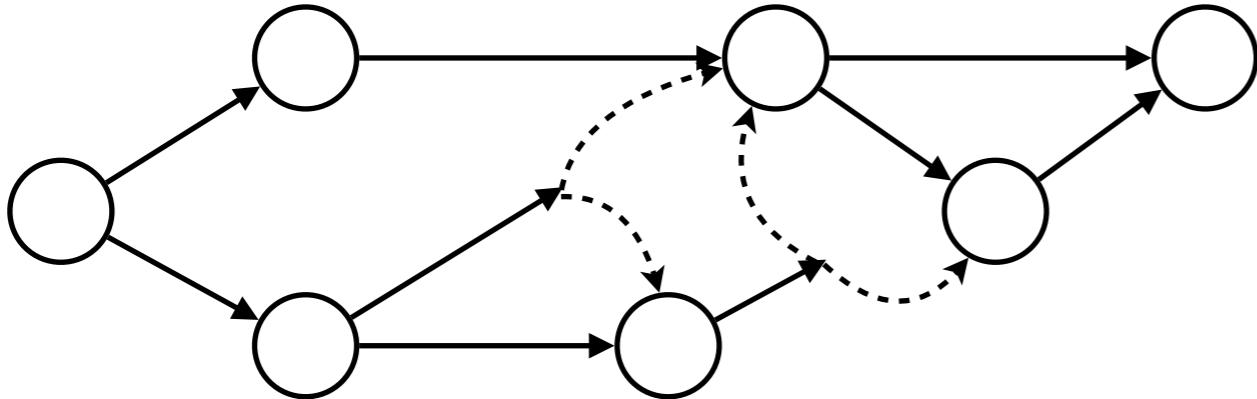


“Freeway + RL”: <https://youtu.be/44CiIPmlimQ>



“tabular RL”

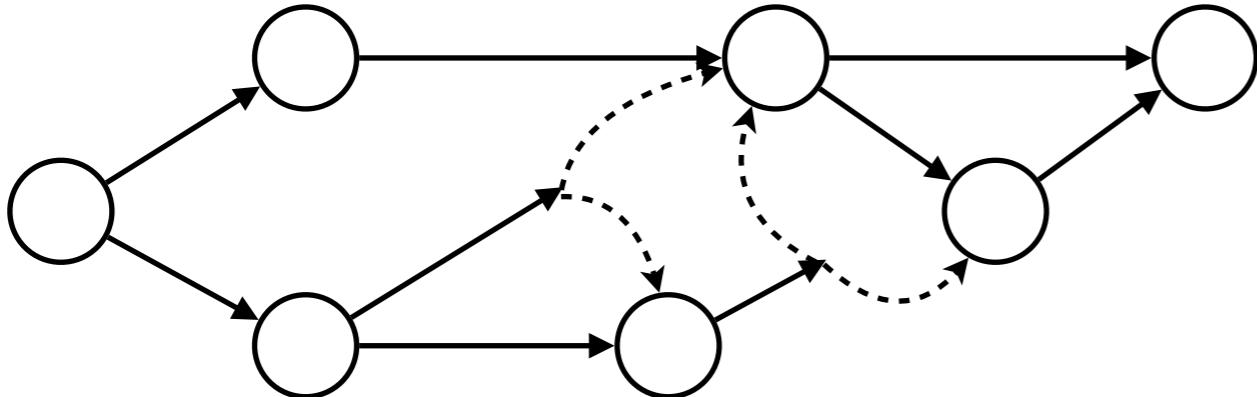
Exploration in small state space is tractable



“tabular RL”

Exploration in small state space is tractable

- Optimize chances for reaching under-visited states

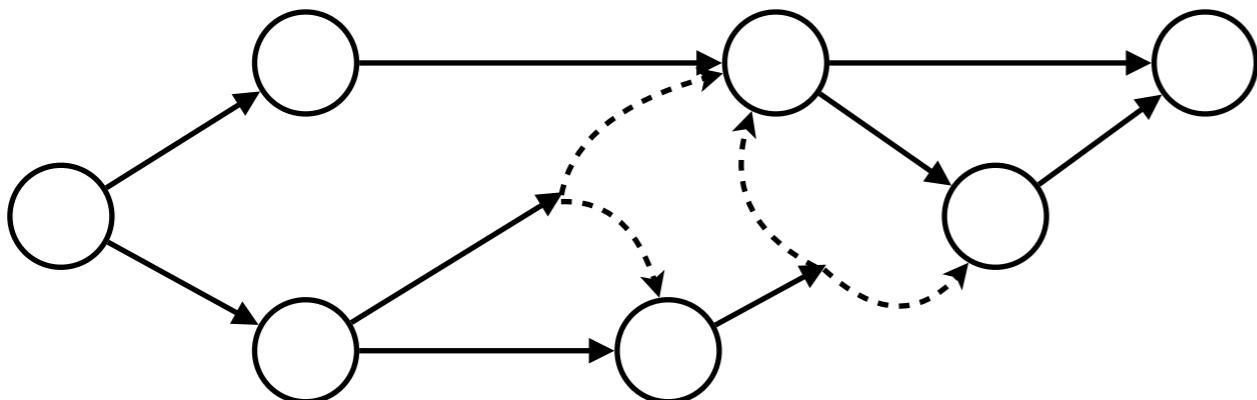


“tabular RL”

Exploration in small state space is tractable

- Optimize chances for reaching under-visited states
- Sample complexity =  $\text{poly}(|S|)$  (and  $|A|, H, 1/\varepsilon, 1/\delta$ )

“PAC-MDP” [Kearns & Singh’98] [Brafman & Tennenholz’02] ...

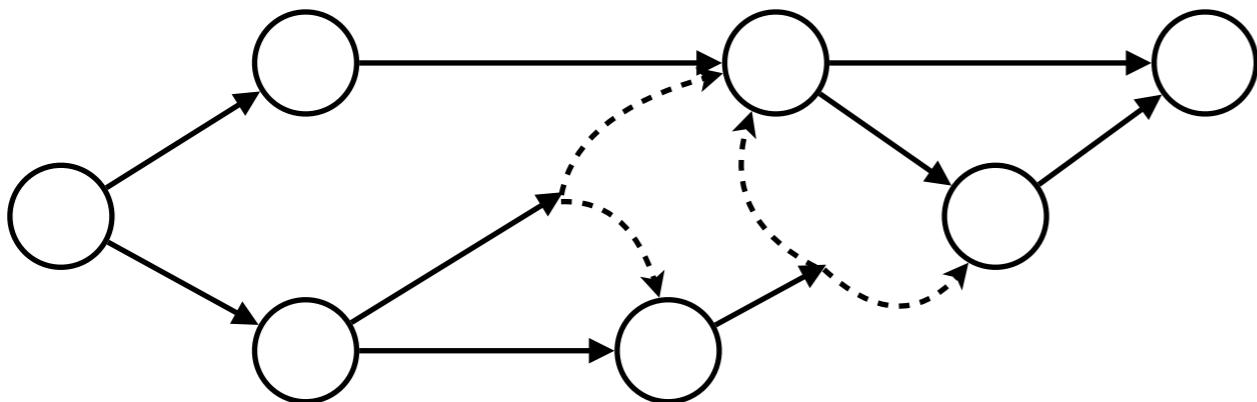


Generalization ?  
 • Large state space

“tabular RL”

Exploration in small state space is tractable

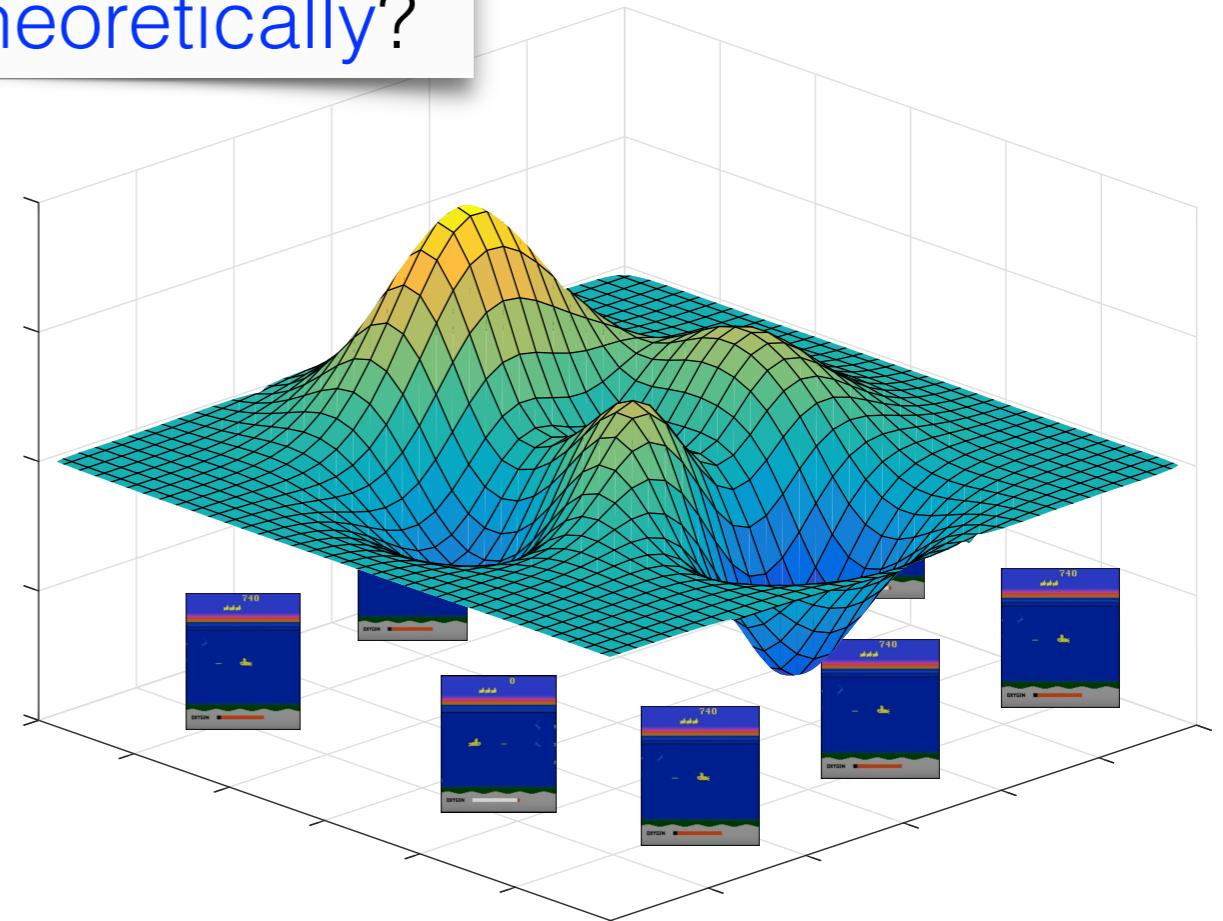
- Optimize chances for reaching under-visited states
  - Sample complexity =  $\text{poly}(|S|)$  (and  $|A|, H, 1/\varepsilon, 1/\delta$ )
- “PAC-MDP” [Kearns & Singh’98] [Brafman & Tennenholtz’02] ...



Generalization  
• Large state space

Systematic exploration in  
large state spaces, at least  
information-theoretically?

Exploration  
• Learner gathers own data



# Formal Model

- Episodic MDP with horizon  $H$
- In each **episode**: for  $h = 1, \dots, H$ , learner

# Formal Model

- Episodic MDP with horizon  $H$
- In each **episode**: for  $h = 1, \dots, H$ , learner
  - observes **state feature**  $x_h \in X$  (possibly infinite) (w.l.o.g.  $x_1 = x^0$ )

# Formal Model

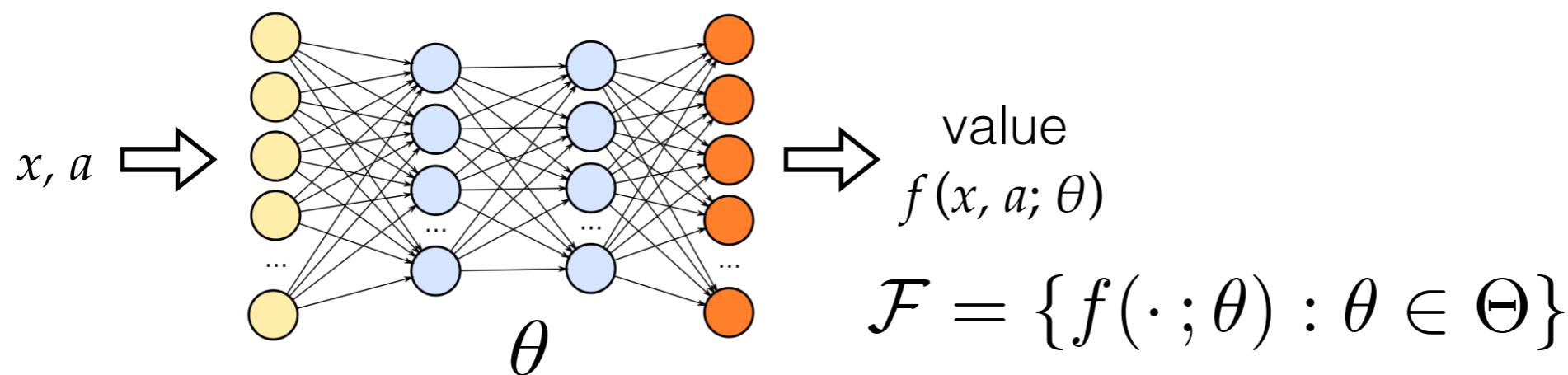
- Episodic MDP with horizon  $H$
- In each **episode**: for  $h = 1, \dots, H$ , learner
  - observes **state feature**  $x_h \in X$  (possibly infinite) (w.l.o.g.  $x_1 = x^0$ )
  - chooses **action**  $a_h \in A$  (finite & manageable)

# Formal Model

- Episodic MDP with horizon  $H$
- In each **episode**: for  $h = 1, \dots, H$ , learner
  - observes **state feature**  $x_h \in X$  (possibly infinite) (w.l.o.g.  $x_1 = x^0$ )
  - chooses **action**  $a_h \in A$  (finite & manageable)
  - receives **reward**  $r_h \in \mathbb{R}$  (bounded)

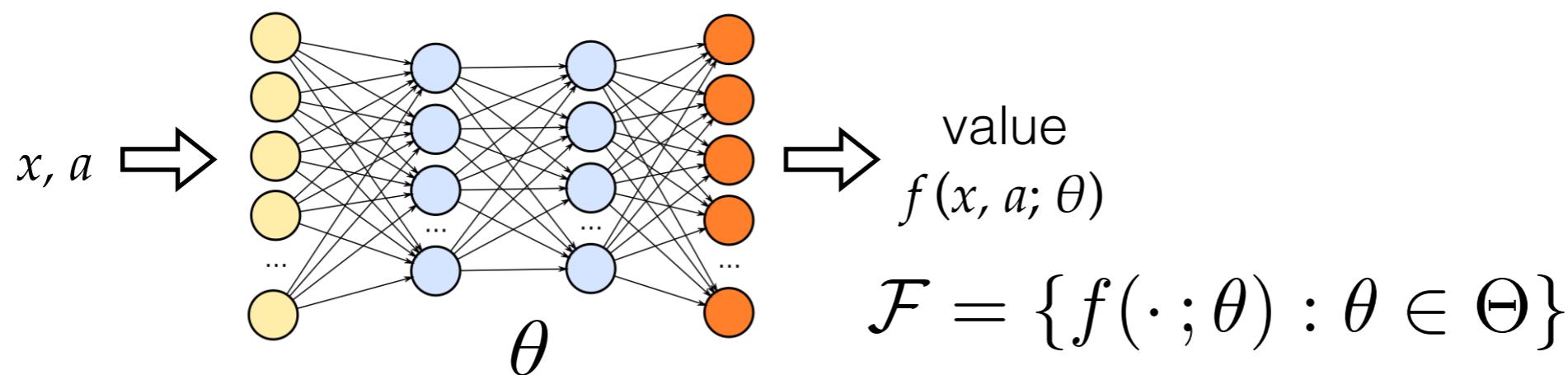
# Formal Model

- Episodic MDP with horizon  $H$
- In each **episode**: for  $h = 1, \dots, H$ , learner
  - observes **state feature**  $x_h \in X$  (possibly infinite) (w.l.o.g.  $x_1 = x^0$ )
  - chooses **action**  $a_h \in A$  (finite & manageable)
  - receives **reward**  $r_h \in \mathbb{R}$  (bounded)
- Learning goal: given  $F$  such that  $Q^* \in F$ , (will relax)



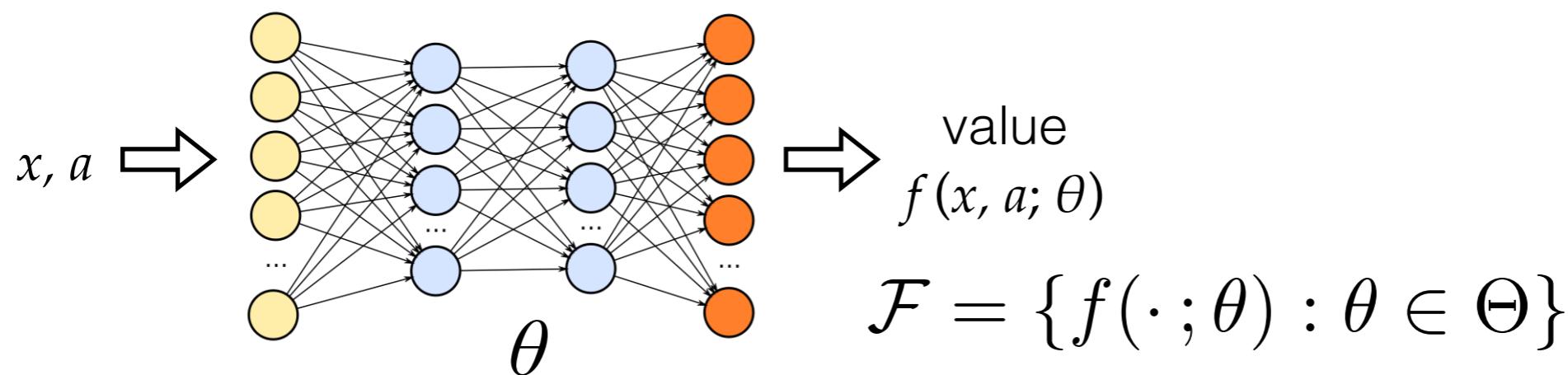
# Formal Model

- Episodic MDP with horizon  $H$
- In each **episode**: for  $h = 1, \dots, H$ , learner
  - observes **state feature**  $x_h \in X$  (possibly infinite) (w.l.o.g.  $x_1 = x^0$ )
  - chooses **action**  $a_h \in A$  (finite & manageable)
  - receives **reward**  $r_h \in \mathbb{R}$  (bounded)
- Learning goal: given  $F$  such that  $Q^* \in F$ , (will relax)  
w.p.  $1 - \delta$ , find policy  $\pi$  s.t.  $J(\pi^*) - J(\pi) \leq \varepsilon$



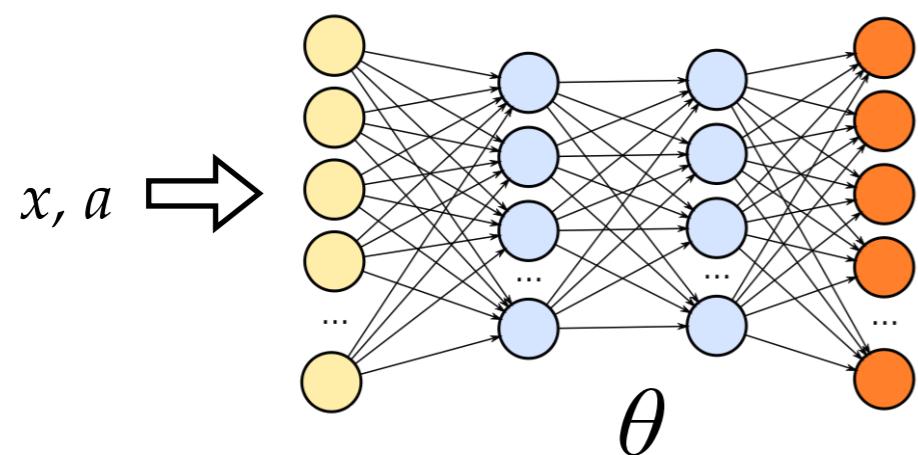
# Formal Model

- Episodic MDP with horizon  $H$
- In each **episode**: for  $h = 1, \dots, H$ , learner
  - observes **state feature**  $x_h \in X$  (possibly infinite) (w.l.o.g.  $x_1 = x^0$ )
  - chooses **action**  $a_h \in A$  (finite & manageable)
  - receives **reward**  $r_h \in \mathbb{R}$  (bounded)
- Learning goal: given  $F$  such that  $Q^* \in F$ , (will relax)  
w.p.  $1 - \delta$ , find policy  $\pi$  s.t.  $J(\pi^*) - J(\pi) \leq \varepsilon$   
using  $\text{poly}(|A|, H, \log|F|, 1/\varepsilon, 1/\delta)$  episodes. (can extend to VC-dim)



# Formal Model

- Episodic MDP with horizon  $H$
- In each **episode**: for  $h = 1, \dots, H$ , learner
  - observes **state feature**  $x_h \in X$  (possibly infinite) (w.l.o.g.  $x_1 = x^0$ )
  - chooses **action**  $a_h \in A$  (finite & manageable)
  - receives **reward**  $r_h \in \mathbb{R}$  (bounded)
- Learning goal: given  $F$  such that  $Q^* \in F$ , (will relax)  
w.p.  $1 - \delta$ , find policy  $\pi$  s.t.  $J(\pi^*) - J(\pi) \leq \varepsilon$   
using  $\text{poly}(|A|, H, \log|F|, 1/\varepsilon, 1/\delta)$  episodes. (can extend to VC-dim)



$$\mathcal{F} = \{f(\cdot; \theta) : \theta \in \Theta\}$$

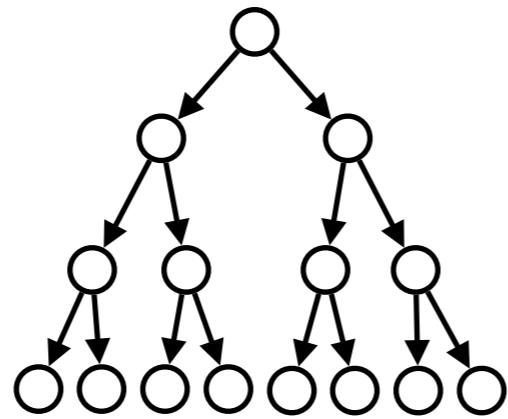
exponential (in  $H$ )  
lower bound exists!  
[Krishnamurthy et al'16]

## Proof of lower bound

- Idea: we are allowed unbounded # of states — use a depth- $H$  complete tree to essentially emulate MAB w/  $|A|^H$  arms

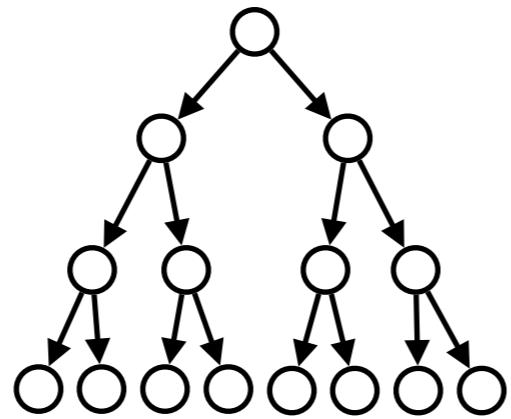
## Proof of lower bound

- Idea: we are allowed unbounded # of states — use a depth- $H$  complete tree to essentially emulate MAB w/  $|A|^H$  arms
- Recall that sample complexity lower bound for MAB is  $\# \text{arms}/\varepsilon^2$



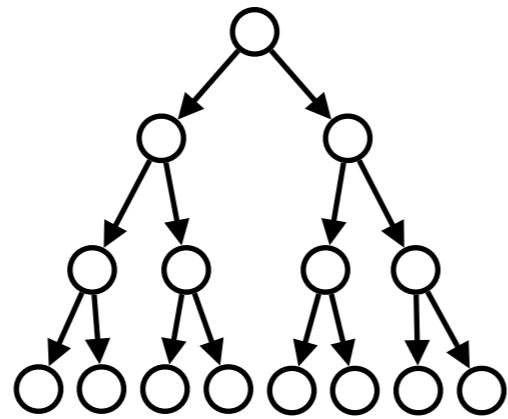
# Proof of lower bound

- Idea: we are allowed unbounded # of states — use a depth- $H$  complete tree to essentially emulate MAB w/  $|A|^H$  arms
  - Recall that sample complexity lower bound for MAB is  $\# \text{arms} / \varepsilon^2$
  - Without function approximation: exponential sample complexity for exploration algorithms



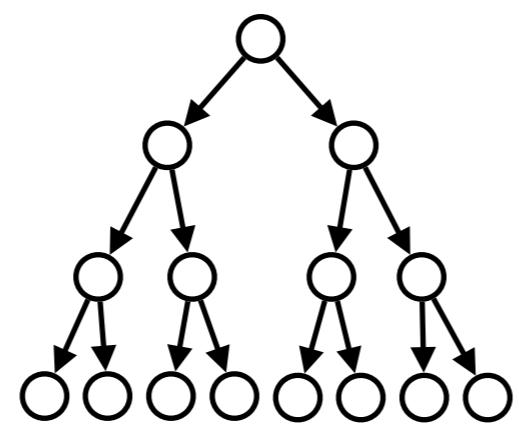
## Proof of lower bound

- Idea: we are allowed unbounded # of states — use a depth- $H$  complete tree to essentially emulate MAB w/  $|A|^H$  arms
- Recall that sample complexity lower bound for MAB is  $\# \text{arms}/\varepsilon^2$
- Without function approximation: exponential sample complexity for exploration algorithms
- Remain to show: function approx. does not help



# Proof of lower bound

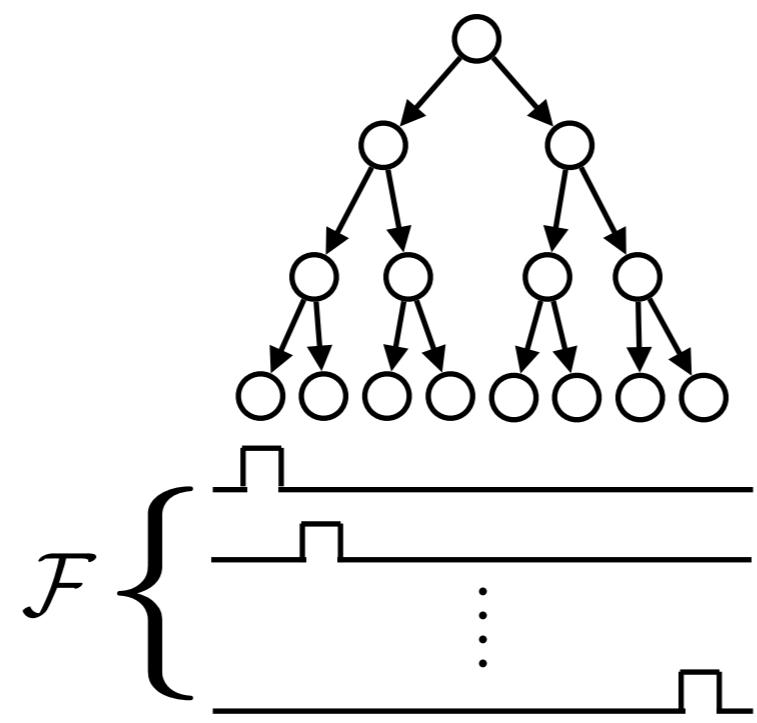
# Show: func. approx. does not help:



# Proof of lower bound

Show: func. approx. does not help:

- Let  $F$  be the collection of  $Q^*$  from all MDPs in family

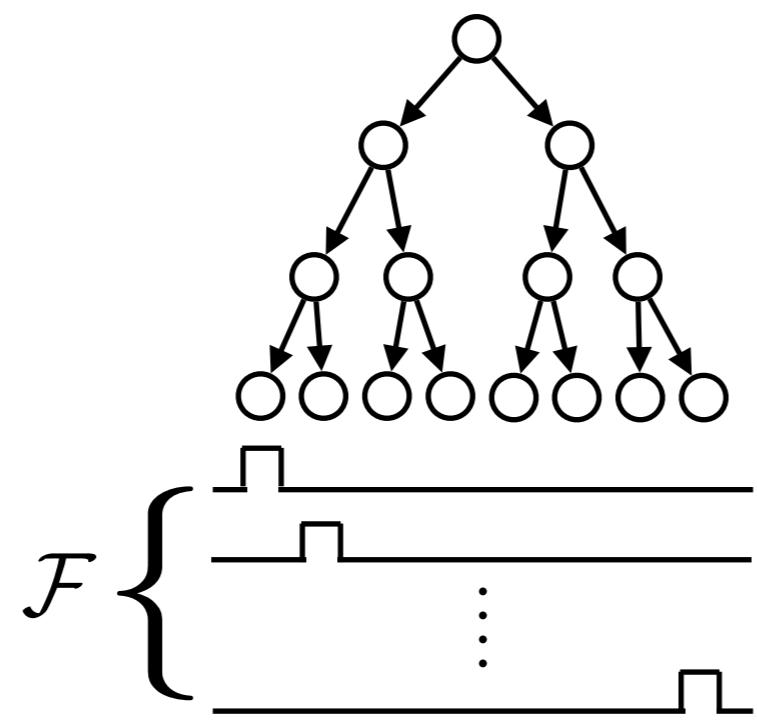


Construction from [Krishnamurthy et al'16]

# Proof of lower bound

Show: func. approx. does not help:

- Let  $F$  be the collection of  $Q^*$  from all MDPs in family
- $\log|F| = H \log|A|$ , always realizable

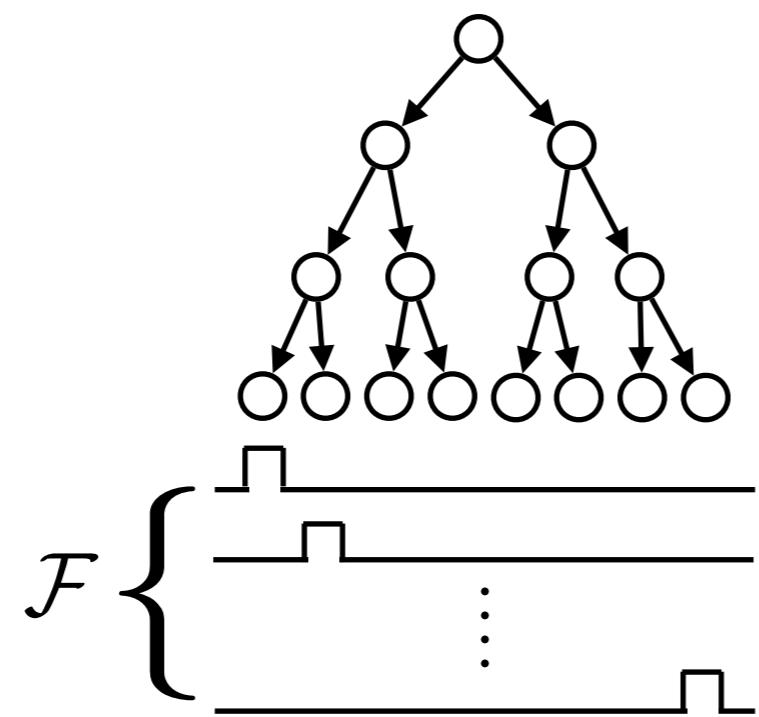


Construction from [Krishnamurthy et al'16]

# Proof of lower bound

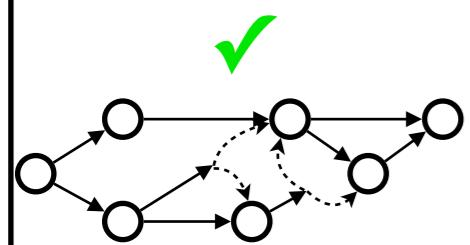
Show: func. approx. does not help:

- Let  $F$  be the collection of  $Q^*$  from all MDPs in family
- $\log|F| = H \log|A|$ , always realizable
- In lower bound proof, alg is allowed to specialize to the problem family — giving  $F$  and  $G$  does not help

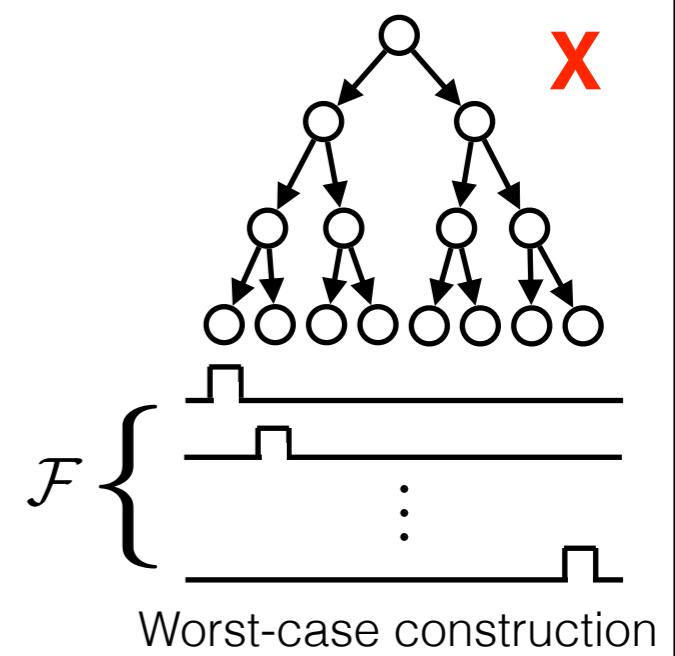


Construction from [Krishnamurthy et al'16]

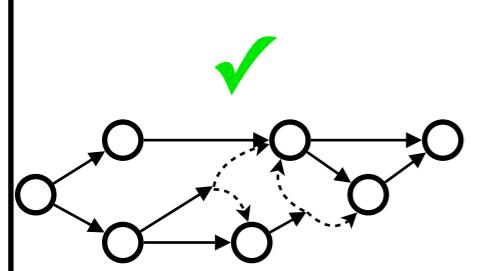
# Zoo of RL Exploration



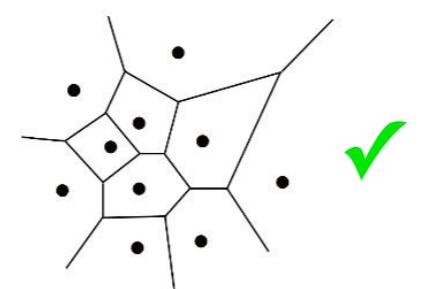
Finite MDPs  
[Kearns & Singh'98]  
(small #states)



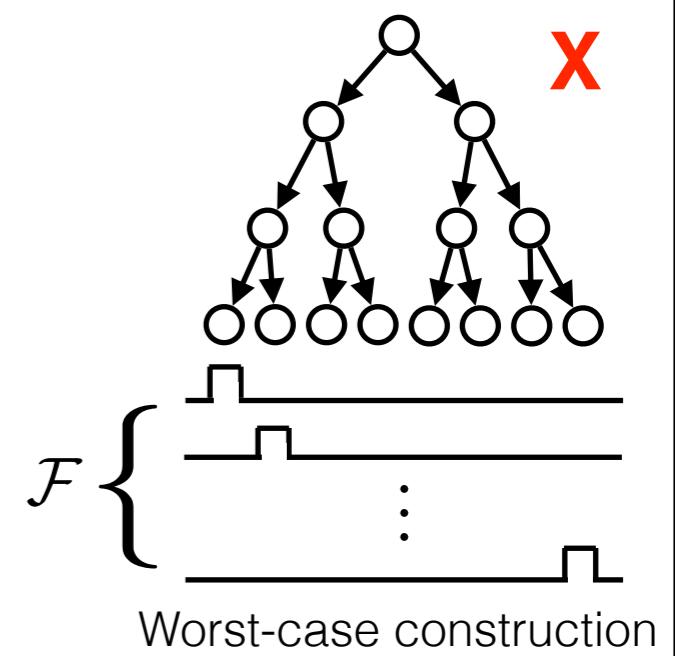
# Zoo of RL Exploration



Finite MDPs  
[Kearns & Singh'98]  
(small #states)

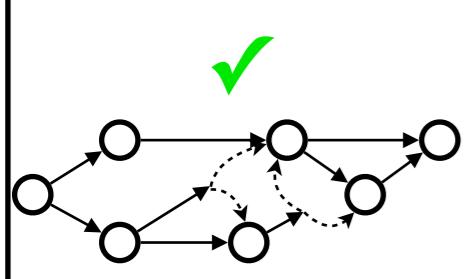


Metric space [Kakade et al'03]  
Abstraction [Li'09]  
(small #abstract states)

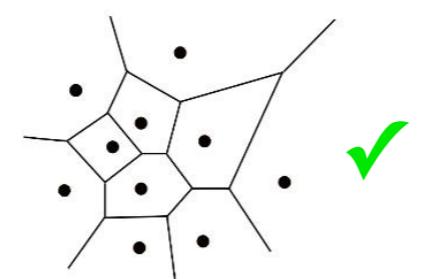


Worst-case construction

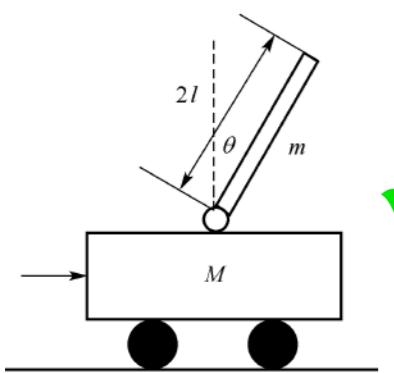
# Zoo of RL Exploration



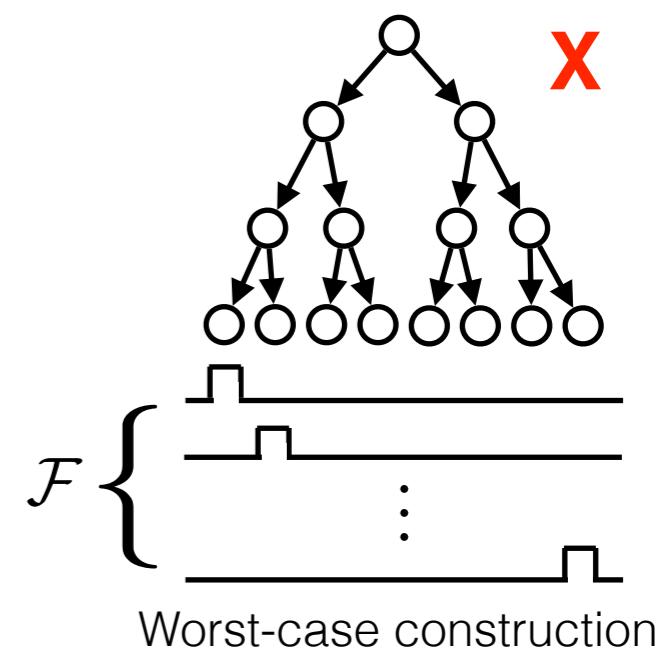
Finite MDPs  
[Kearns & Singh'98]  
(small #states)



Metric space [Kakade et al'03]  
Abstraction [Li'09]  
(small #abstract states)

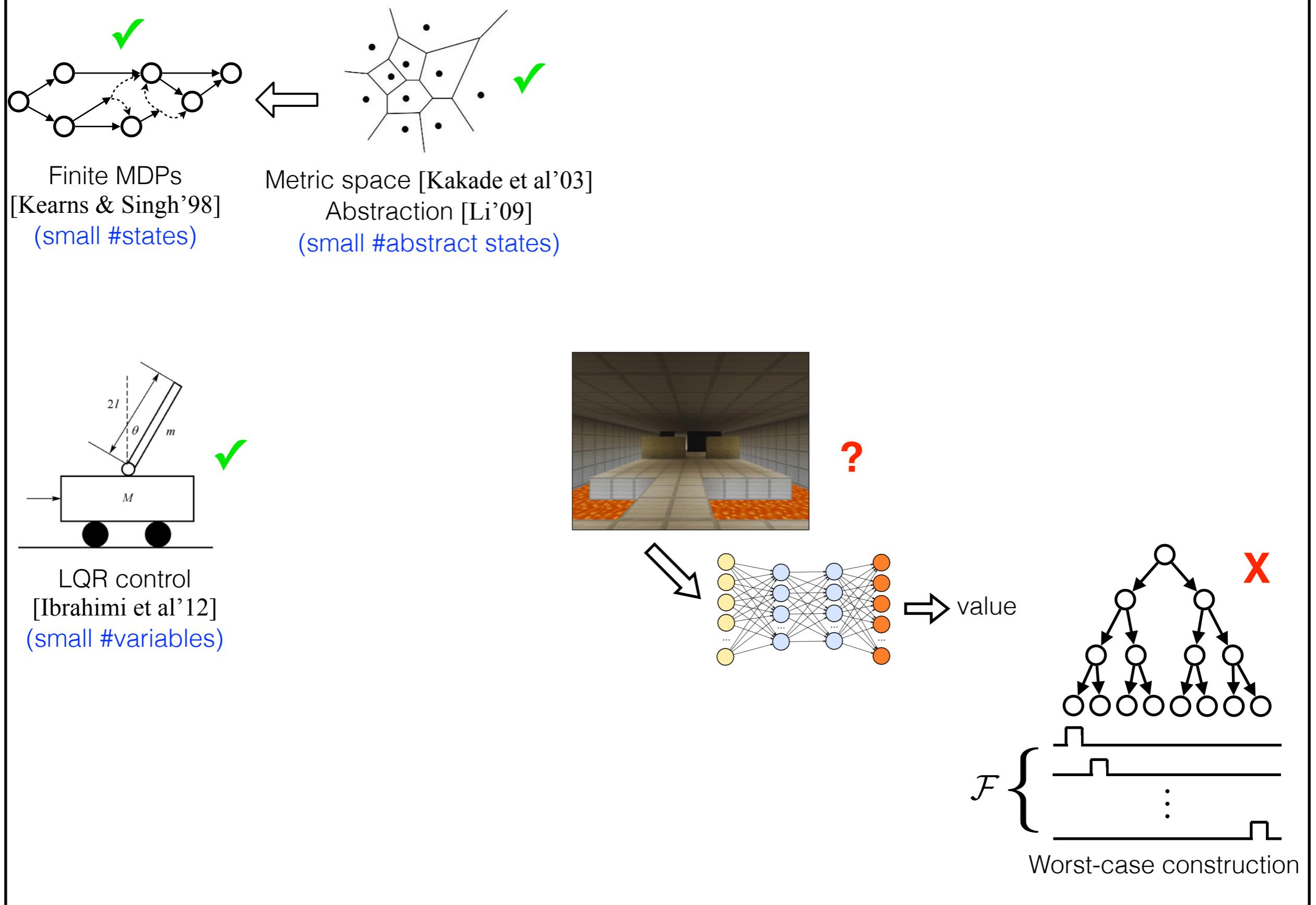


LQR control  
[Ibrahim et al'12]  
(small #variables)

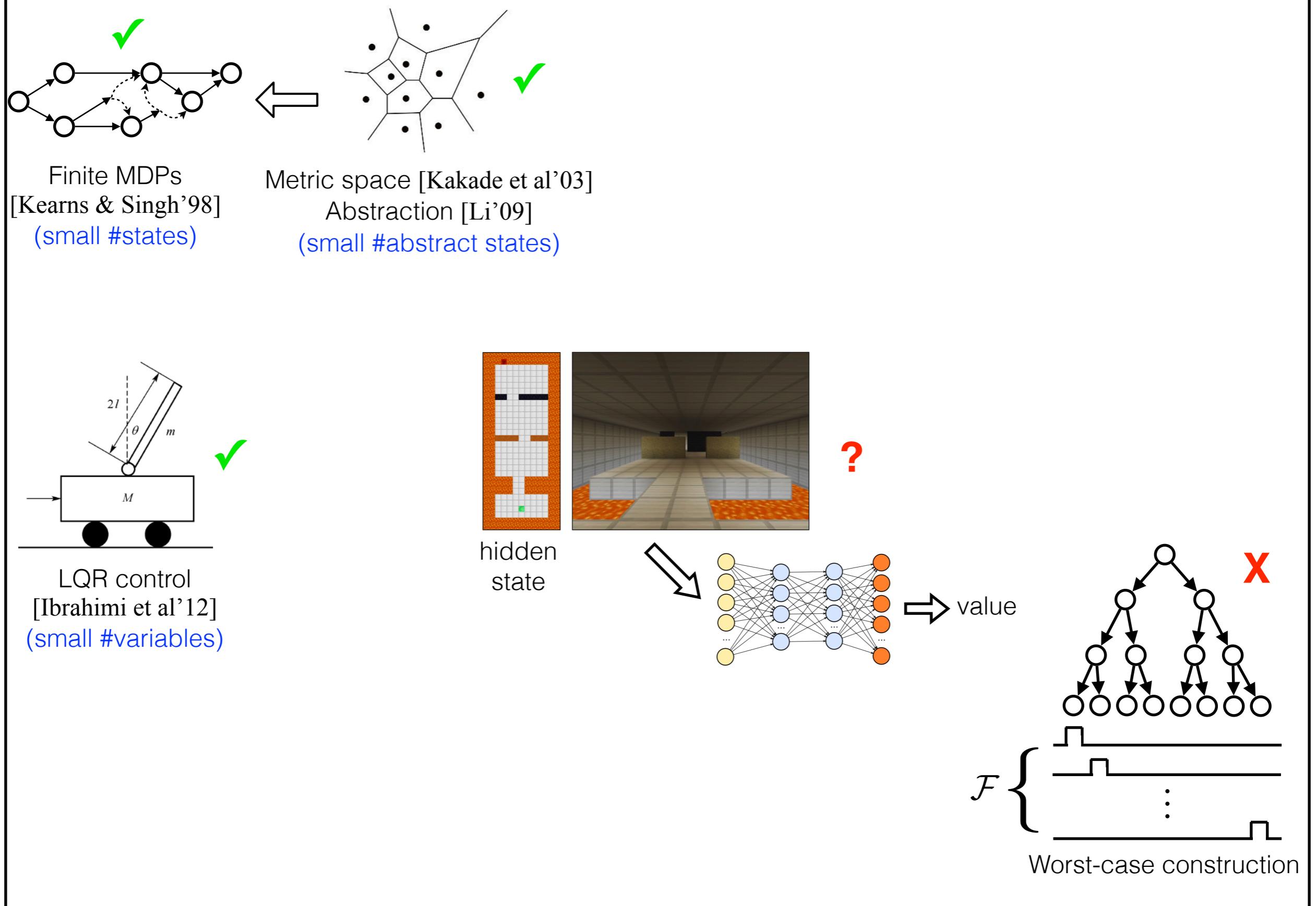


Worst-case construction

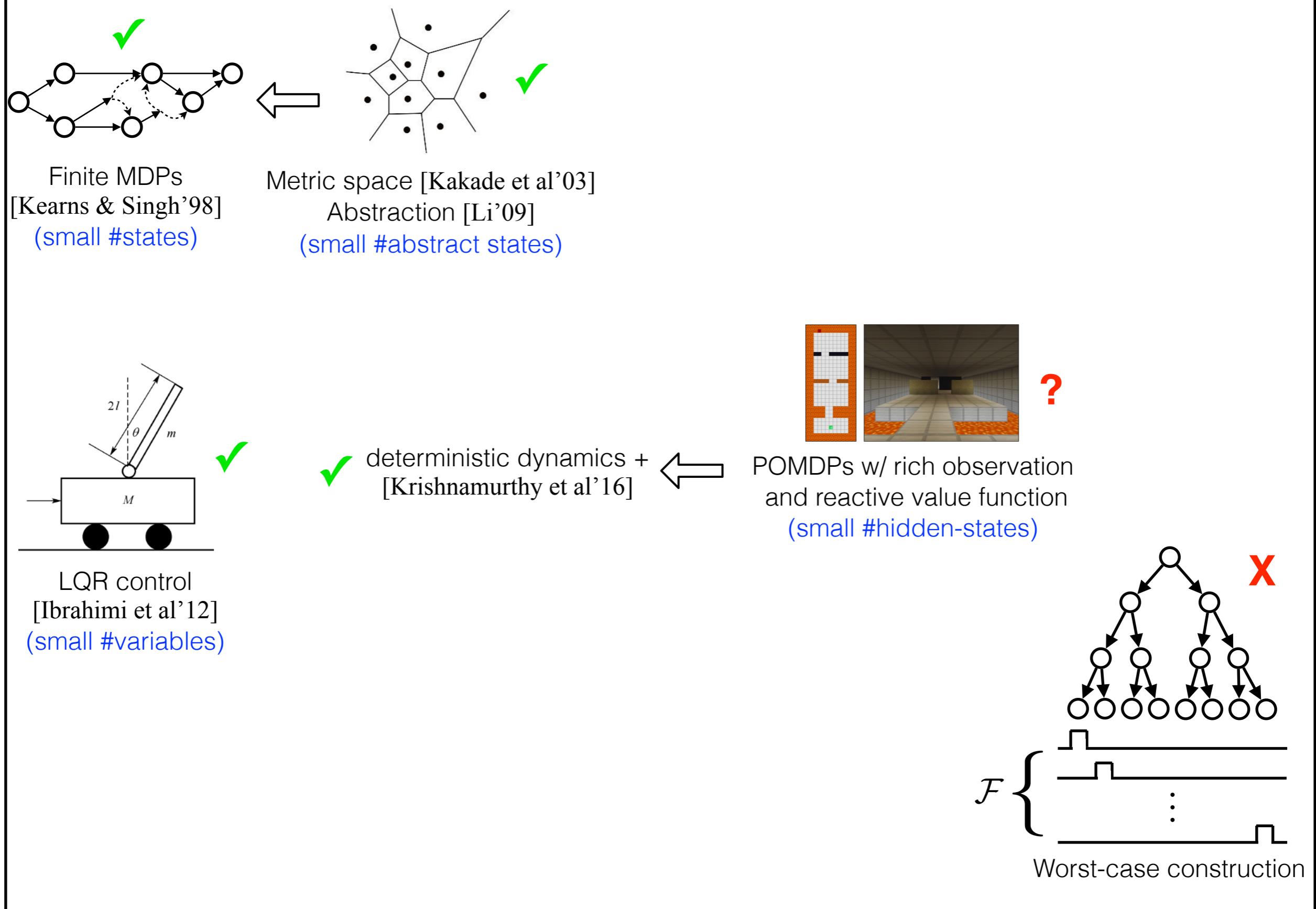
# Zoo of RL Exploration



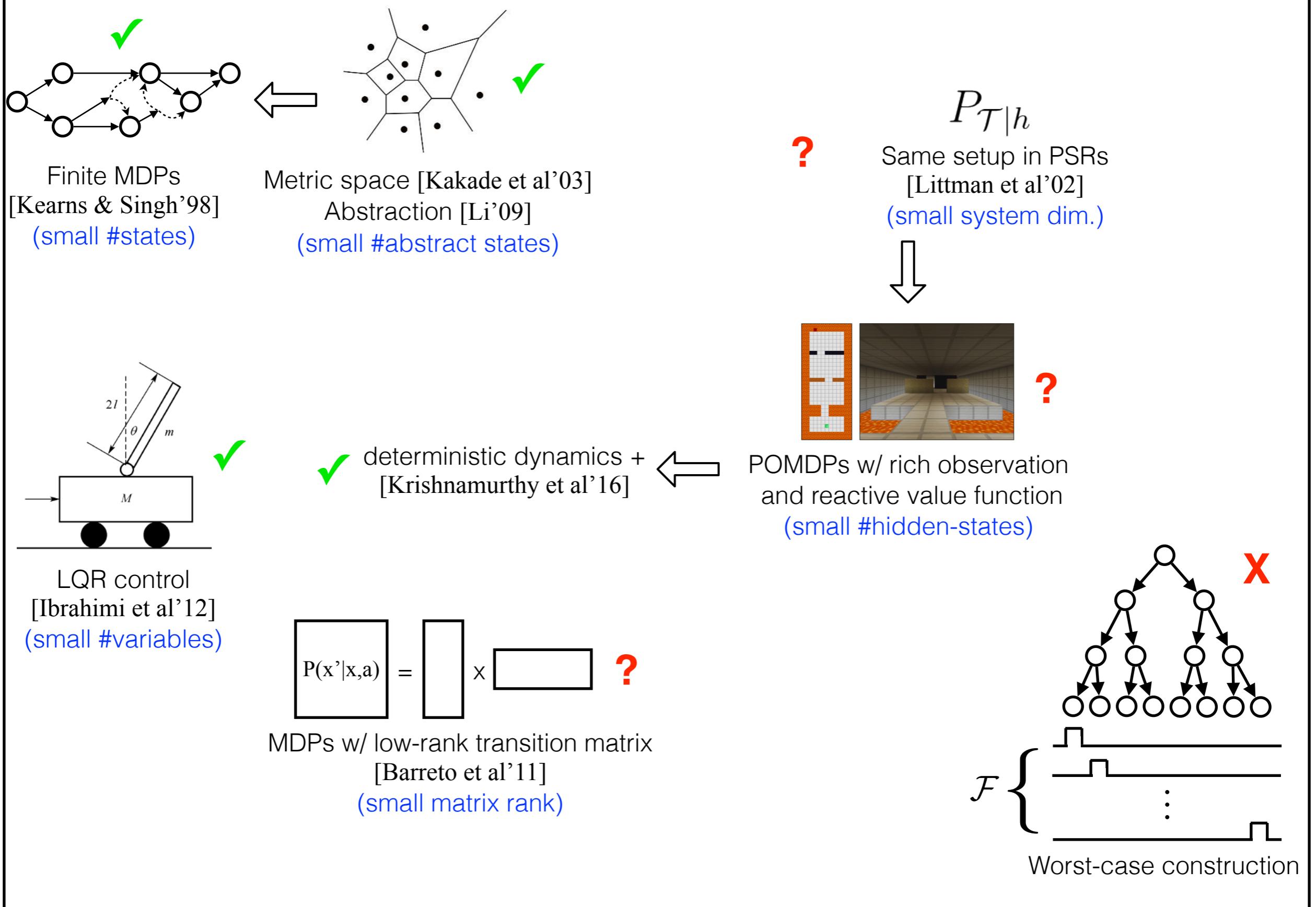
# Zoo of RL Exploration



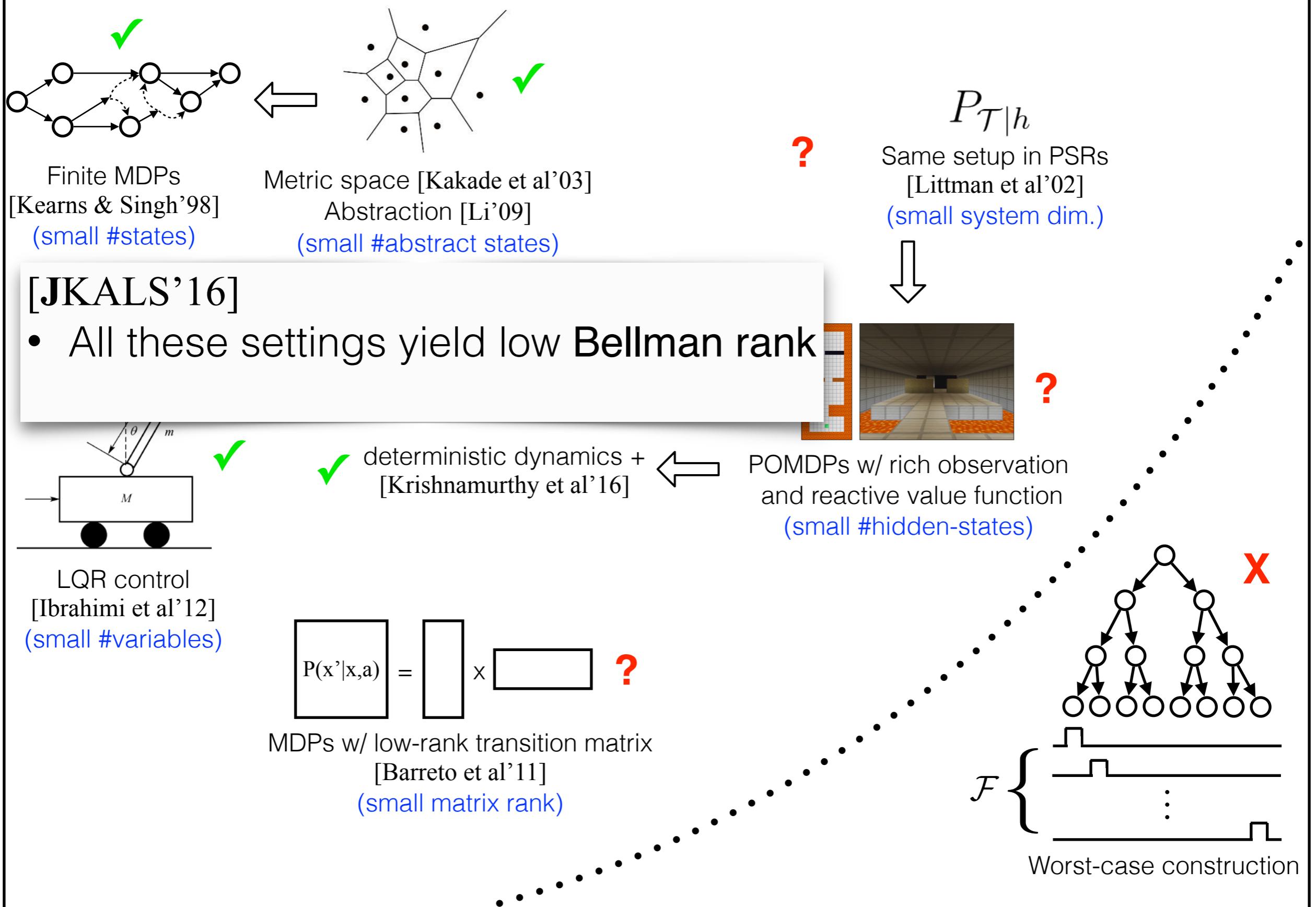
# Zoo of RL Exploration



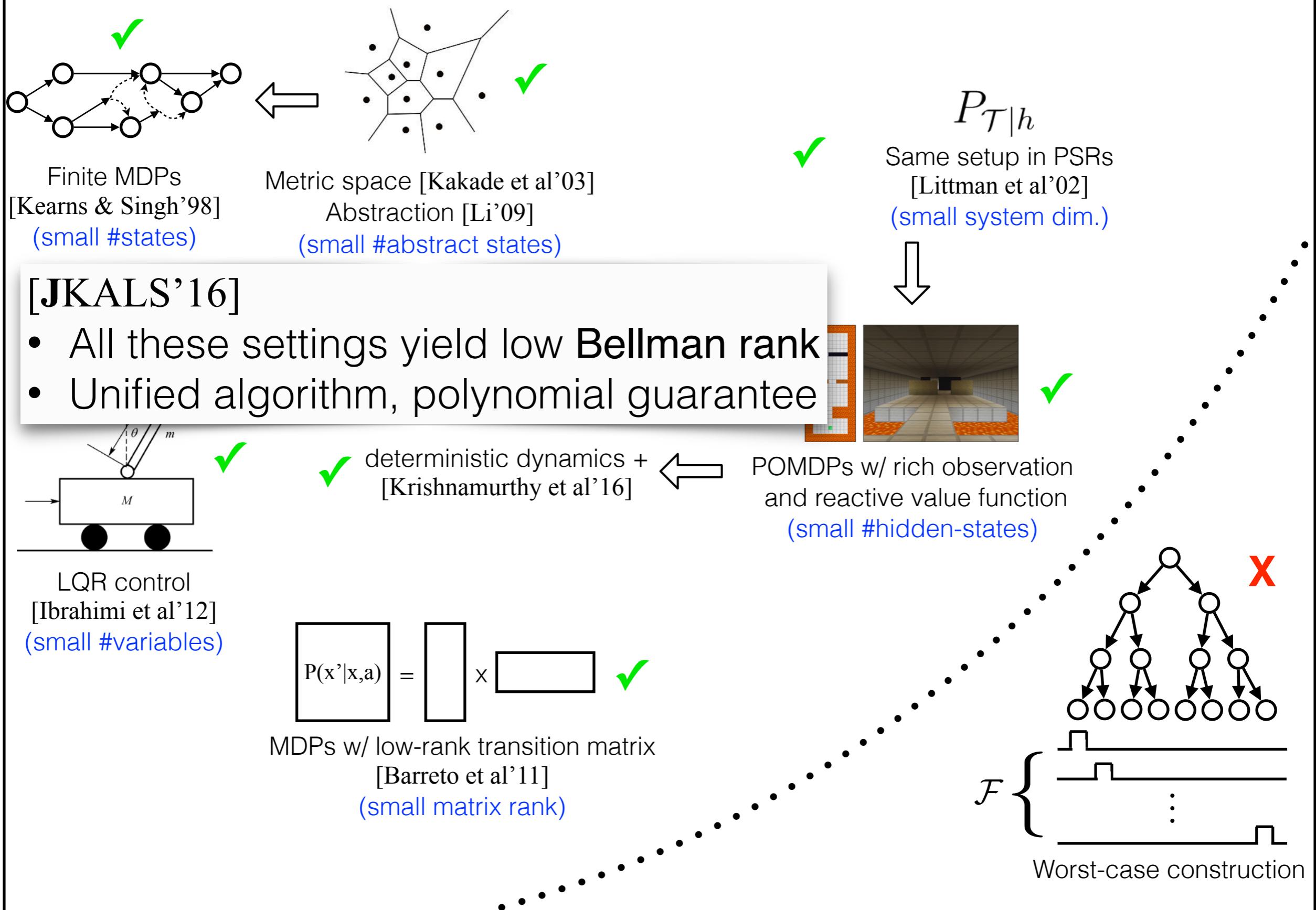
# Zoo of RL Exploration



# Zoo of RL Exploration



# Zoo of RL Exploration



# Defining Bellman rank

## Step 1: Average Bellman Error

- Bellman error of  $f$  at  $(x_h, a_h)$

$$f(x_h, a_h) - \mathbb{E}_{r_h, x_{h+1} | x_h, a_h} \left[ r_h + \max_{a \in \mathcal{A}} f(x_{h+1}, a) \right]$$

# Defining Bellman rank

## Step 1: Average Bellman Error

- Bellman error of  $f$  at  $(x_h, a_h)$

$$f(x_h, a_h) - \mathbb{E}_{r_h, x_{h+1} | x_h, a_h} \left[ r_h + \max_{a \in \mathcal{A}} f(x_{h+1}, a) \right]$$

- $Q^*$  has 0 Bellman error for all  $(x_h, a_h)$ .

# Defining Bellman rank

## Step 1: Average Bellman Error

- Bellman error of  $f$  at  $(x_h, a_h)$

$$f(x_h, a_h) - \mathbb{E}_{r_h, x_{h+1} | x_h, a_h} \left[ r_h + \max_{a \in \mathcal{A}} f(x_{h+1}, a) \right]$$

- $Q^*$  has 0 Bellman error for all  $(x_h, a_h)$ .
- Average Bellman error of  $f$  is the linear combination of its Bellman errors over  $(x_h, a_h)$

# Defining Bellman rank

## Step 1: Average Bellman Error

- Bellman error of  $f$  at  $(x_h, a_h)$

$$f(x_h, a_h) - \mathbb{E}_{r_h, x_{h+1} | x_h, a_h} \left[ r_h + \max_{a \in \mathcal{A}} f(x_{h+1}, a) \right]$$

- $Q^*$  has 0 Bellman error for all  $(x_h, a_h)$ .
- Average Bellman error of  $f$  is the linear combination of its Bellman errors over  $(x_h, a_h)$ 
  - Weights: distribution over  $x_h$  induced by policy  $\pi$ .

$$\mathcal{E}^h(f, \pi) := \mathbb{E}_{\substack{a_{1:h-1} \sim \pi \\ a_h \sim f}} [f(x_h, a_h) - r_h - \max_{a \in \mathcal{A}} f(x_{h+1}, a)]$$

$a_h = \arg \max f(x_h, \cdot)$

# Defining Bellman rank

## Step 1: Average Bellman Error

- Bellman error of  $f$  at  $(x_h, a_h)$

$$f(x_h, a_h) - \mathbb{E}_{r_h, x_{h+1} | x_h, a_h} \left[ r_h + \max_{a \in \mathcal{A}} f(x_{h+1}, a) \right]$$

- $Q^*$  has 0 Bellman error for all  $(x_h, a_h)$ .
- Average Bellman error of  $f$  is the linear combination of its Bellman errors over  $(x_h, a_h)$ 
  - Weights: distribution over  $x_h$  induced by policy  $\pi$ .

$$\mathcal{E}^h(f, \pi) := \mathbb{E}_{\substack{a_{1:h-1} \sim \pi \\ a_h \sim f}} [f(x_h, a_h) - r_h - \max_{a \in \mathcal{A}} f(x_{h+1}, a)]$$

$a_h = \arg \max f(x_h, \cdot)$

- $\mathcal{E}^h(Q^*, \pi) = 0$  for all  $\pi$  and  $h$ .

# Defining Bellman rank

## Step 2: Bellman error matrices

$$\textcolor{red}{f} \in \mathcal{F}$$

$$\pi \in \Pi_{\mathcal{F}}$$

$$\mathcal{E}^h(\textcolor{red}{f}, \pi) := \mathbb{E}_{\substack{a_{1:h-1} \sim \pi \\ a_h \sim \textcolor{red}{f}}} [\textcolor{red}{f}(x_h, a_h) - r_h - \max_{a \in \mathcal{A}} \textcolor{red}{f}(x_{h+1}, a)]$$

# Defining Bellman rank

## Step 2: Bellman error matrices

$$\textcolor{red}{f} \in \mathcal{F}$$

$\pi \in \Pi_{\mathcal{F}}$ $\downarrow$ class of greedy policies induced from $F$ : $\Pi_{\mathcal{F}} := \{x \mapsto \arg \max f(x, \cdot) : f \in \mathcal{F}\}$	$\cdots \cdots \cdots \quad \mathcal{E}^h(\textcolor{red}{f}, \pi) :=$ $\mathbb{E}_{\substack{a_{1:h-1} \sim \pi \\ a_h \sim \textcolor{red}{f}}} [\textcolor{red}{f}(x_h, a_h) - r_h - \max_{a \in \mathcal{A}} \textcolor{red}{f}(x_{h+1}, a)]$
---	--

# Defining Bellman rank

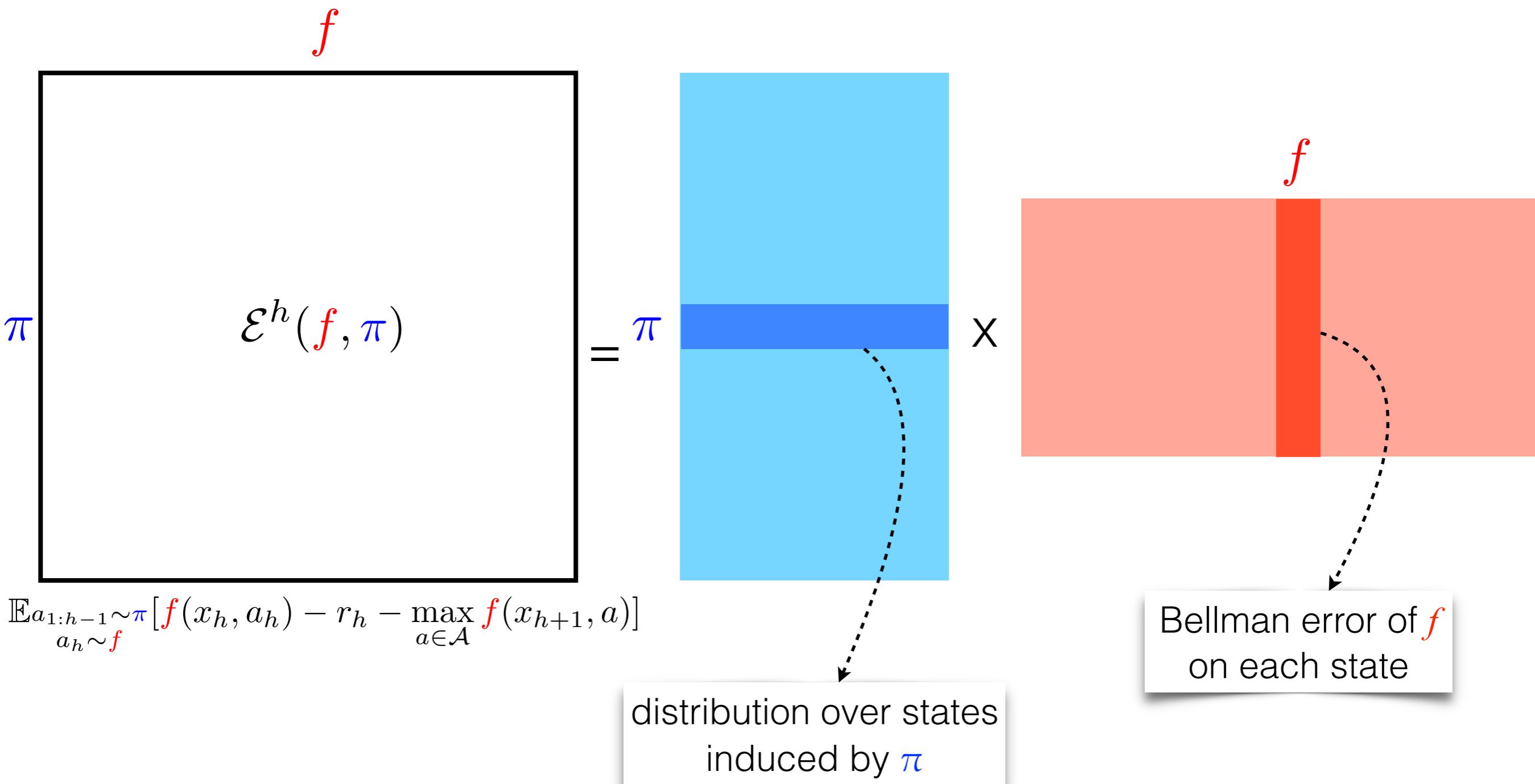
## Step 2: Bellman error matrices

$$f \in \mathcal{F}$$

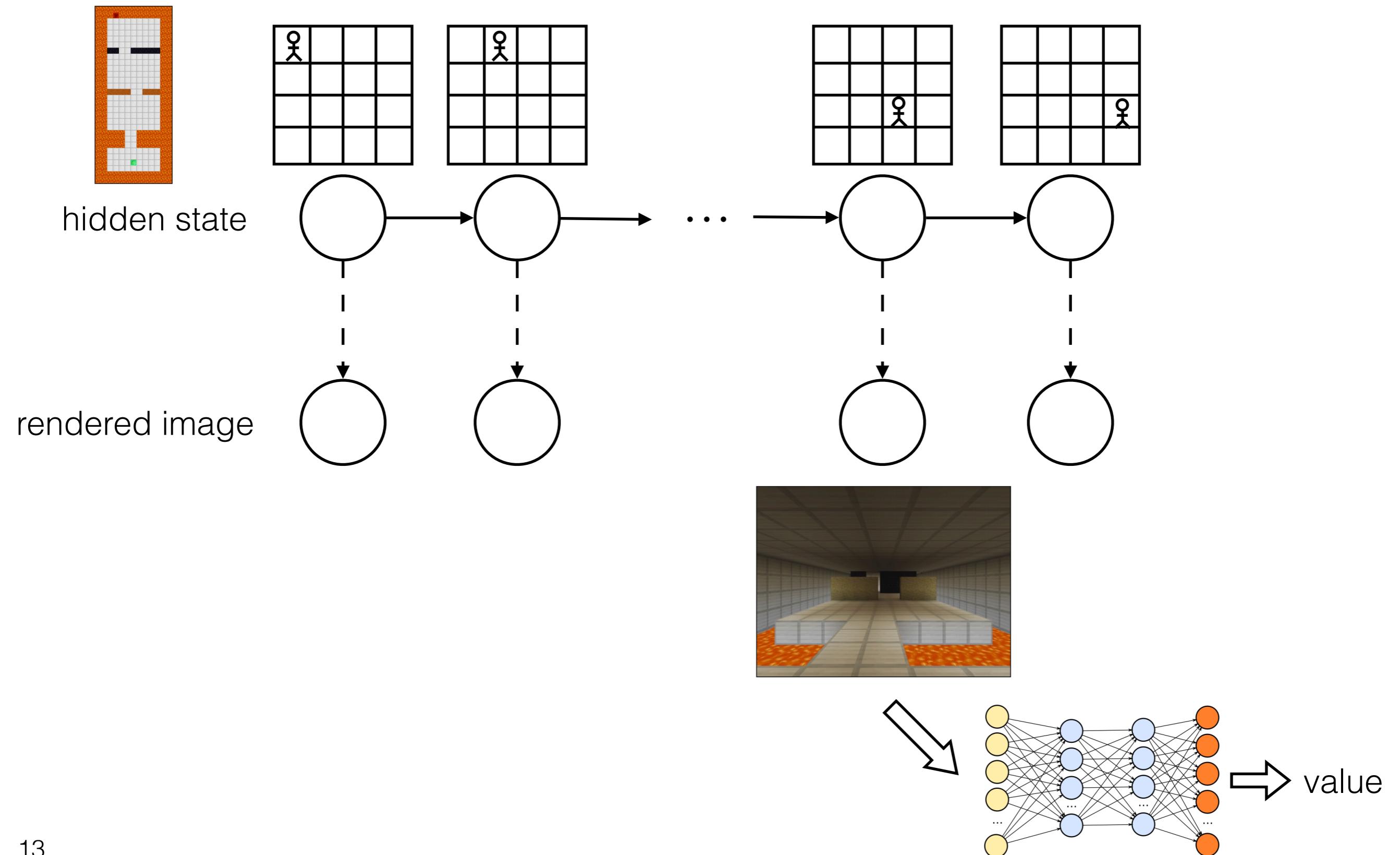
$$\pi \in \Pi_{\mathcal{F}} \quad \dots \dots \dots \quad \mathcal{E}^h(f, \pi) := \\ \mathbb{E}_{\substack{a_{1:h-1} \sim \pi \\ a_h \sim f}} [f(x_h, a_h) - r_h - \max_{a \in \mathcal{A}} f(x_{h+1}, a)]$$

**Definition:** *Bellman rank* is an uniform upper bound on the rank of matrices  $[\mathcal{E}^h(\textcolor{red}{f}, \pi)]_{\pi, f}$  over  $h = 1, 2, \dots, H$ .

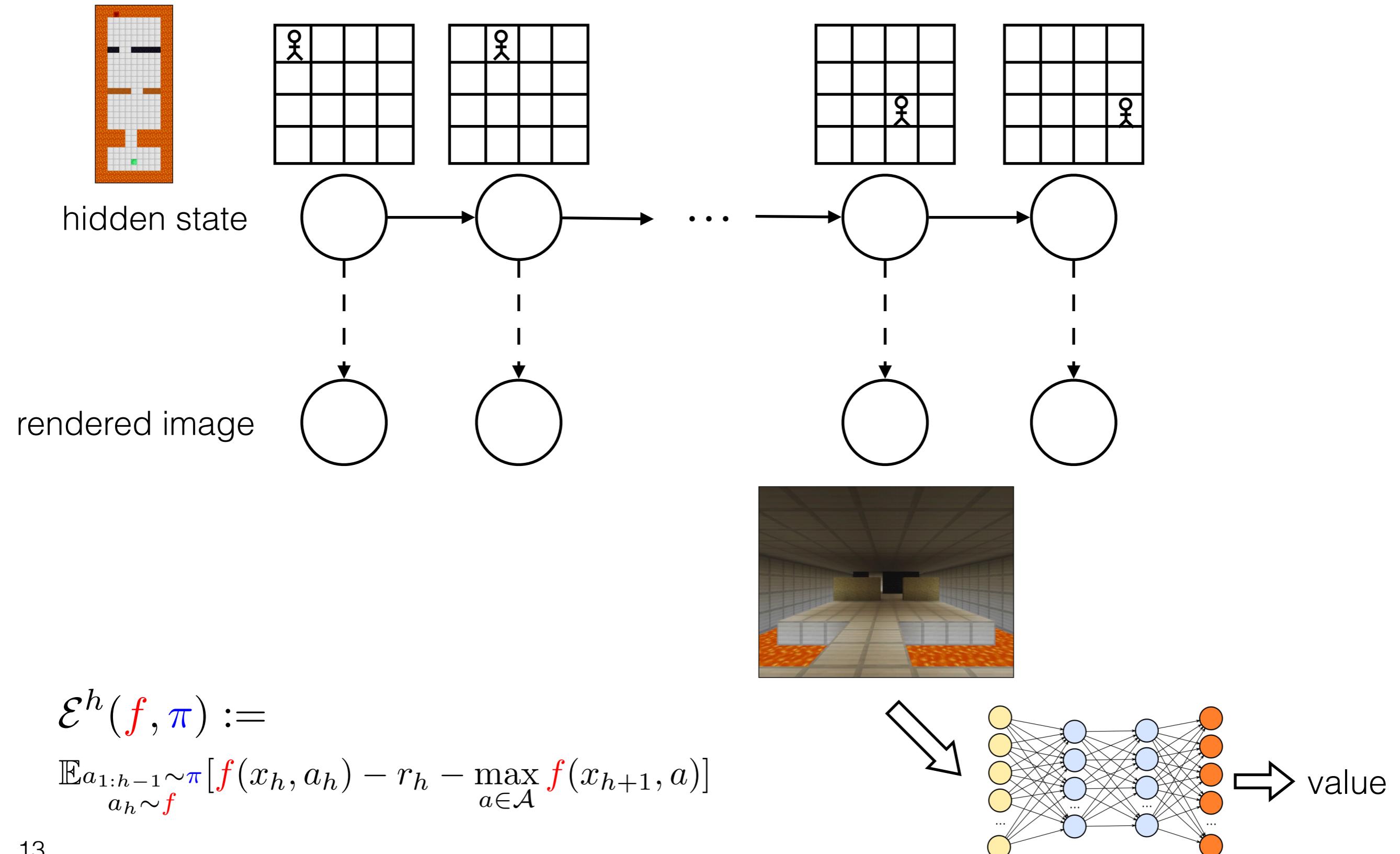
# Tabular MDP: Bellman rank $\leq \#\text{states}$



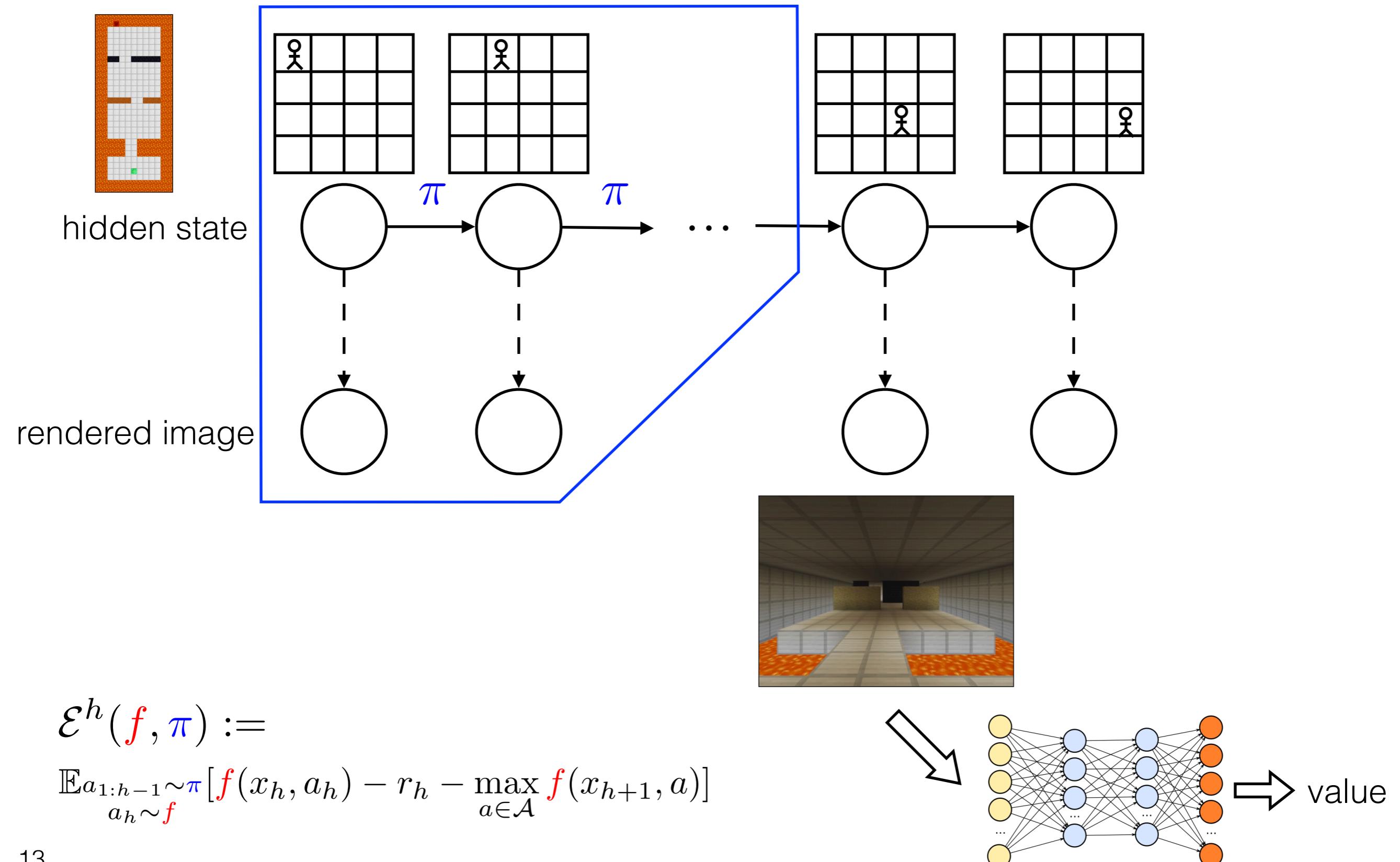
“Visual grid-world”: Bellman rank  $\leq$  # hidden states



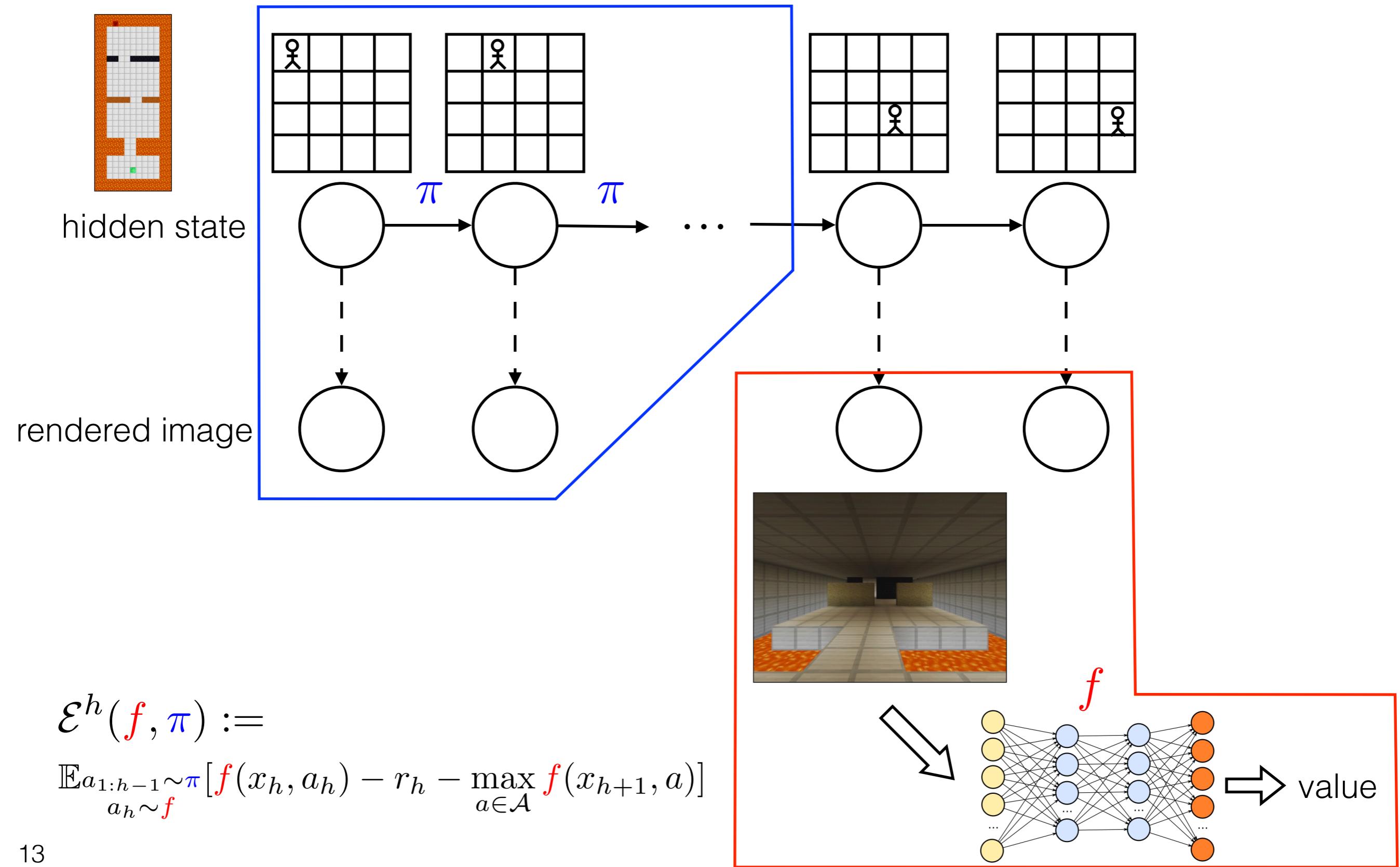
# “Visual grid-world”: Bellman rank $\leq$ # hidden states



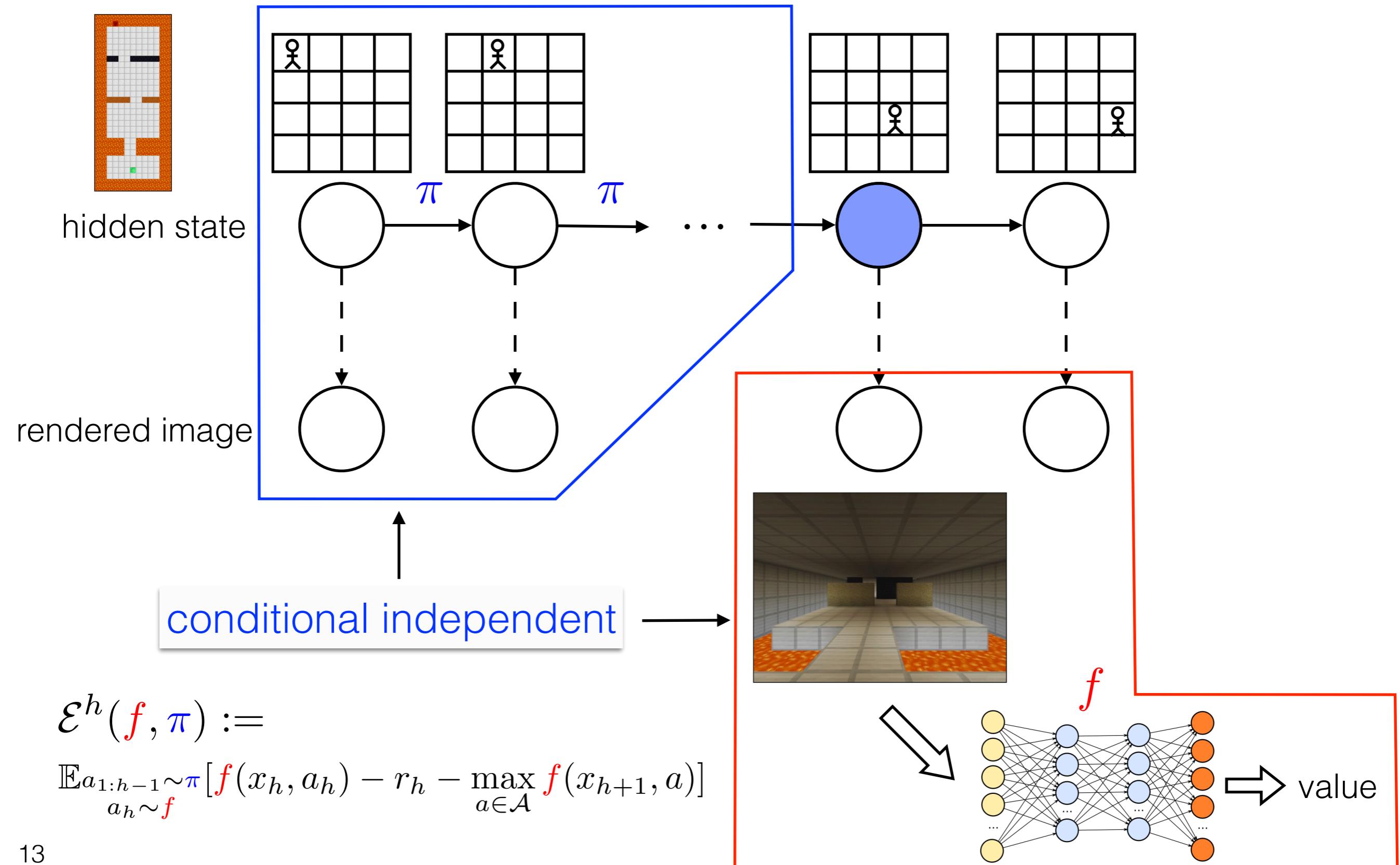
# “Visual grid-world”: Bellman rank $\leq$ # hidden states



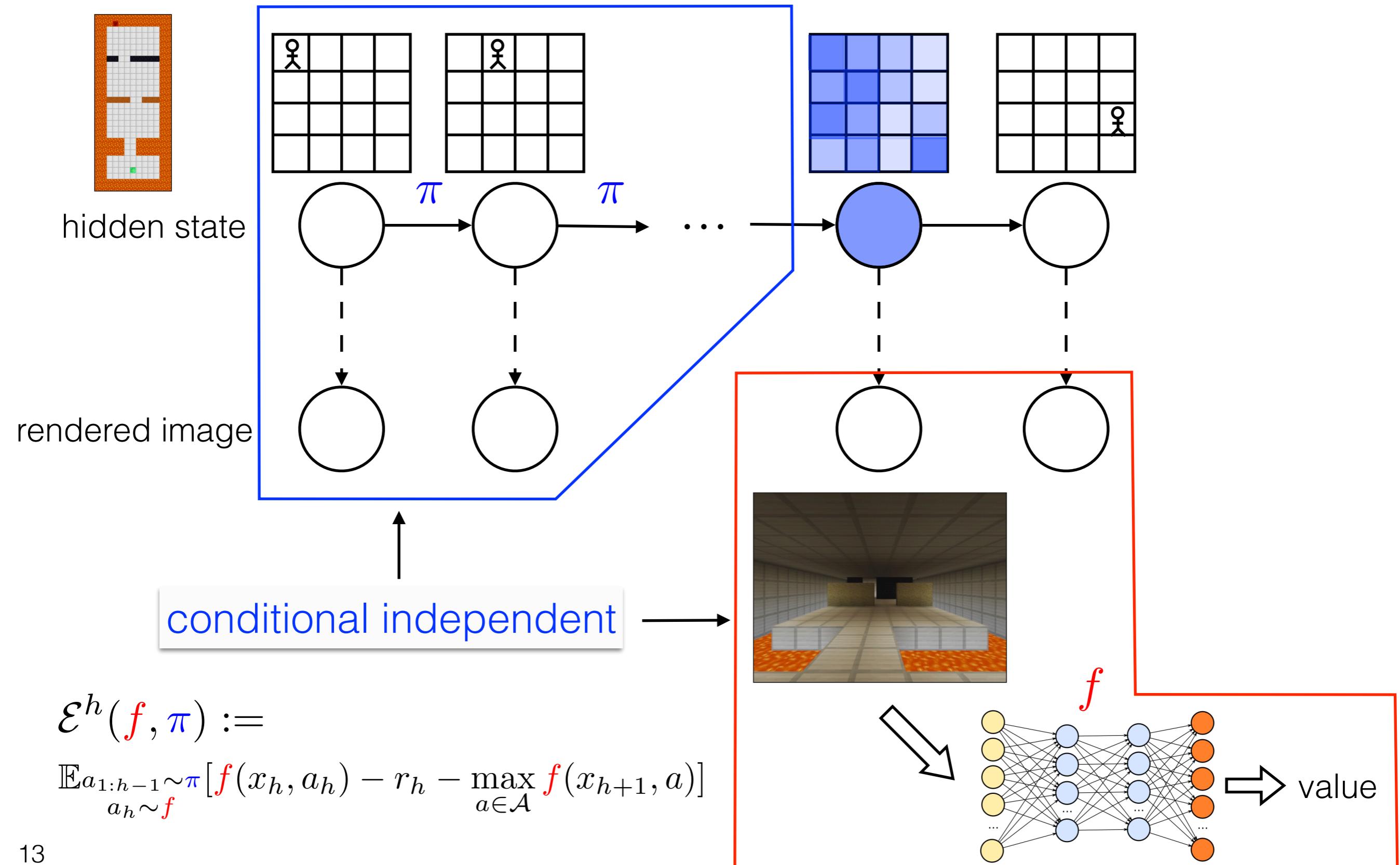
# “Visual grid-world”: Bellman rank $\leq$ # hidden states



# “Visual grid-world”: Bellman rank $\leq$ # hidden states



# “Visual grid-world”: Bellman rank $\leq$ # hidden states



# $Q^*$ -irrelevant abstractions

## $Q^*$ -irrelevant abstractions

- Number of abstract states is small

## $Q^*$ -irrelevant abstractions

- Number of abstract states is small
- Challenge: abstract state does not “block” influence from past

## $Q^*$ -irrelevant abstractions

- Number of abstract states is small
- Challenge: abstract state does not “block” influence from past
- Witness statistics: for each possible  $(x, a, r, x')$

$$\Pr_{a_{1:h-1} \sim \pi} [x_h = x, r_h = r, x_{h+1} = x' \mid \text{do } a_h = a]$$

## $Q^*$ -irrelevant abstractions

- Number of abstract states is small
- Challenge: abstract state does not “block” influence from past
- Witness statistics: for each possible  $(x, a, r, x')$

$$\Pr_{a_{1:h-1} \sim \pi} [x_h = x, r_h = r, x_{h+1} = x' \mid \text{do } a_h = a]$$

## $Q^*$ -irrelevant abstractions

- Number of abstract states is small
- Challenge: abstract state does not “block” influence from past
- Witness statistics: for each possible  $(x, a, r, x')$

$$\Pr_{a_{1:h-1} \sim \pi} [x_h = x, r_h = r, x_{h+1} = x' \mid \text{do } a_h = a]$$

- Dimension: (#abstract states)<sup>2</sup> \* (# actions) \* (# possible values for reward)

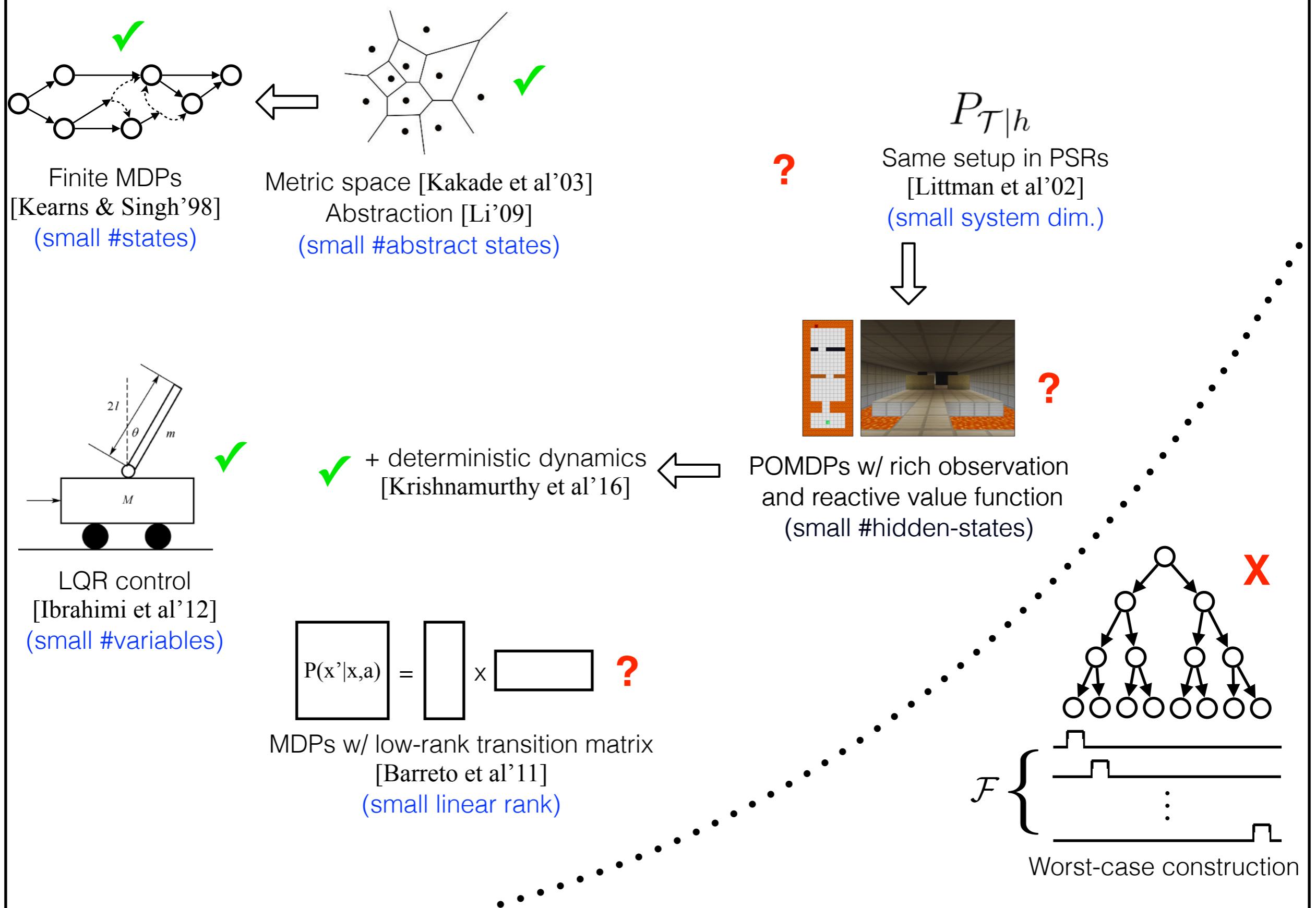
## $Q^*$ -irrelevant abstractions

- Number of abstract states is small
- Challenge: abstract state does not “block” influence from past
- Witness statistics: for each possible  $(x, a, r, x')$

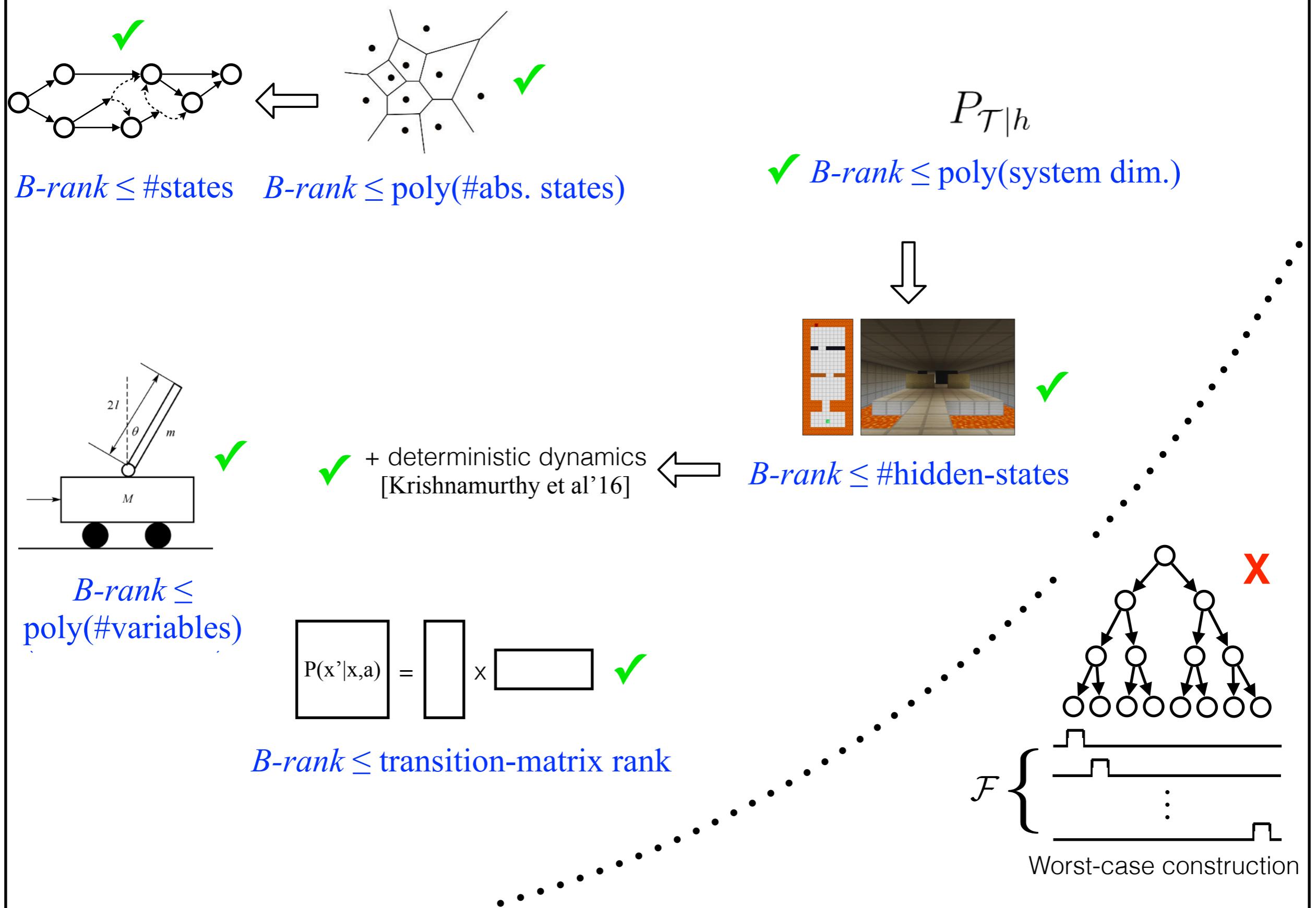
$$\Pr_{a_{1:h-1} \sim \pi} [x_h = x, r_h = r, x_{h+1} = x' \mid \text{do } a_h = a]$$

- Dimension: (#abstract states)<sup>2</sup> \* (# actions) \* (# possible values for reward)
  - Reward can always be discretized (and incur a small error)

# Zoo of RL Exploration



# Zoo of RL Exploration



# New algorithm: OLIVE

**(Optimism-Led Iterative Value-function Elimination)**

$F_1 := F.$  // version space

(Ignoring statistical slackness parameters)

For iteration  $t=1, 2, \dots$

# New algorithm: OLIVE

(**O**ptimism-**L**ed **I**terative **V**alue-function **E**limination)

$F_1 := F$ . // version space

(Ignoring statistical slackness parameters)

For iteration  $t=1, 2, \dots$

- Choose  $f_t$  as the  $f \in F_t$  that maximizes  $v_f := \max_{a \in \mathcal{A}} f(x^0, a)$

# New algorithm: OLIVE

(Optimism-Led Iterative Value-function Elimination)

$F_1 := F$ . // version space

(Ignoring statistical slackness parameters)

For iteration  $t=1, 2, \dots$

- Choose  $f_t$  as the  $f \in F_t$  that maximizes  $v_f := \max_{a \in \mathcal{A}} f(x^0, a)$
- Estimate the value of  $\pi_t$  — the greedy policy of  $f_t$ .

# New algorithm: OLIVE

(Optimism-Led Iterative Value-function Elimination)

$F_1 := F$ . // version space

(Ignoring statistical slackness parameters)

For iteration  $t=1, 2, \dots$

- Choose  $f_t$  as the  $f \in F_t$  that maximizes  $v_f := \max_{a \in \mathcal{A}} f(x^0, a)$
- Estimate the value of  $\pi_t$  — the greedy policy of  $f_t$ .
  - If  $J(\pi^t) \geq v_{f_t}$  return  $\pi_t$ .

# New algorithm: OLIVE

(Optimism-Led Iterative Value-function Elimination)

$F_1 := F$ . // version space

(Ignoring statistical slackness parameters)

For iteration  $t=1, 2, \dots$

- Choose  $f_t$  as the  $f \in F_t$  that maximizes  $v_f := \max_{a \in \mathcal{A}} f(x^0, a)$
- Estimate the value of  $\pi_t$  — the greedy policy of  $f_t$ .
  - If  $J(\pi_t) \geq v_{f_t}$  return  $\pi_t$ .

Estimate by MC evaluation

# New algorithm: OLIVE

(Optimism-Led Iterative Value-function Elimination)

$F_1 := F$ . // version space

(Ignoring statistical slackness parameters)

For iteration  $t=1, 2, \dots$

- Choose  $f_t$  as the  $f \in F_t$  that maximizes  $v_f := \max_{a \in \mathcal{A}} f(x^0, a)$
- Estimate the value of  $\pi_t$  — the greedy policy of  $f_t$ .
  - If  $J(\pi_t) \geq v_{f_t}$  ( $\geq v_{Q^*} = J(\pi^*)$ ), return  $\pi_t$ .

# New algorithm: OLIVE

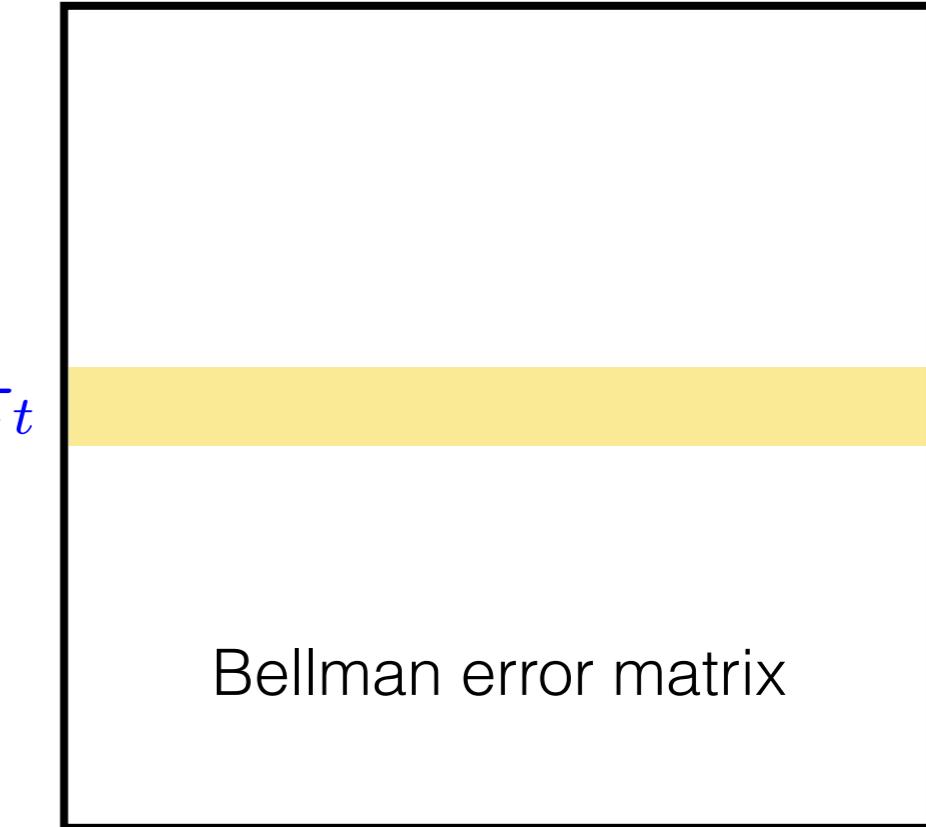
(Optimism-Led Iterative Value-function Elimination)

$F_1 := F$ . // version space

(Ignoring statistical slackness parameters)

For iteration  $t=1, 2, \dots$

- Choose  $f_t$  as the  $f \in F_t$  that maximizes  $v_f := \max_{a \in \mathcal{A}} f(x^0, a)$
- Estimate the value of  $\pi_t$  — the greedy policy of  $f_t$ .
  - If  $J(\pi_t) \geq v_{f_t}$  ( $\geq v_{Q^*} = J(\pi^*)$ ), return  $\pi_t$ .
- Estimate  $\mathcal{E}^h(f, \pi_t)$  for all  $f, h$ .



# New algorithm: OLIVE

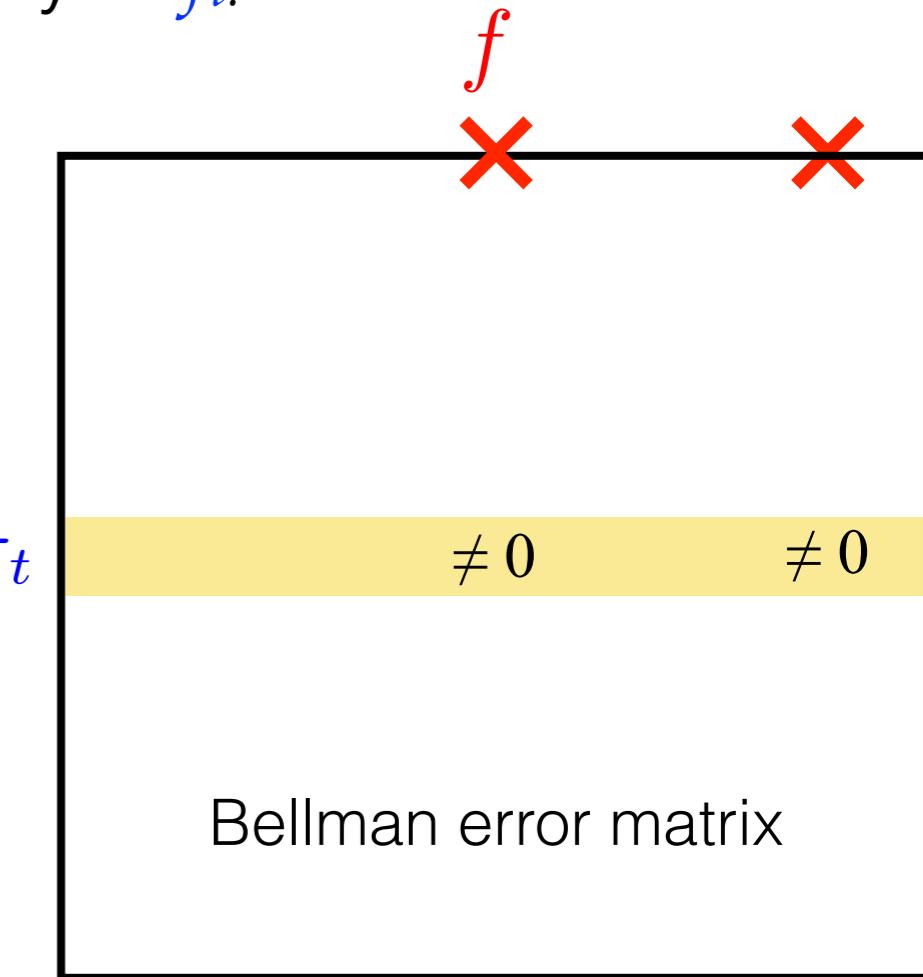
(Optimism-Led Iterative Value-function Elimination)

$F_1 := F$ . // version space

(Ignoring statistical slackness parameters)

For iteration  $t=1, 2, \dots$

- Choose  $f_t$  as the  $f \in F_t$  that maximizes  $v_f := \max_{a \in \mathcal{A}} f(x^0, a)$
- Estimate the value of  $\pi_t$  — the greedy policy of  $f_t$ .
  - If  $J(\pi_t) \geq v_{f_t}$  ( $\geq v_{Q^*} = J(\pi^*)$ ), return  $\pi_t$ .
- Estimate  $\mathcal{E}^h(f, \pi_t)$  for all  $f, h$ .
- Eliminate  $f$  s.t.  $\mathcal{E}^h(f, \pi_t) \neq 0, \forall h$   
 $\Rightarrow F_{t+1}$ .



# New algorithm: OLIVE

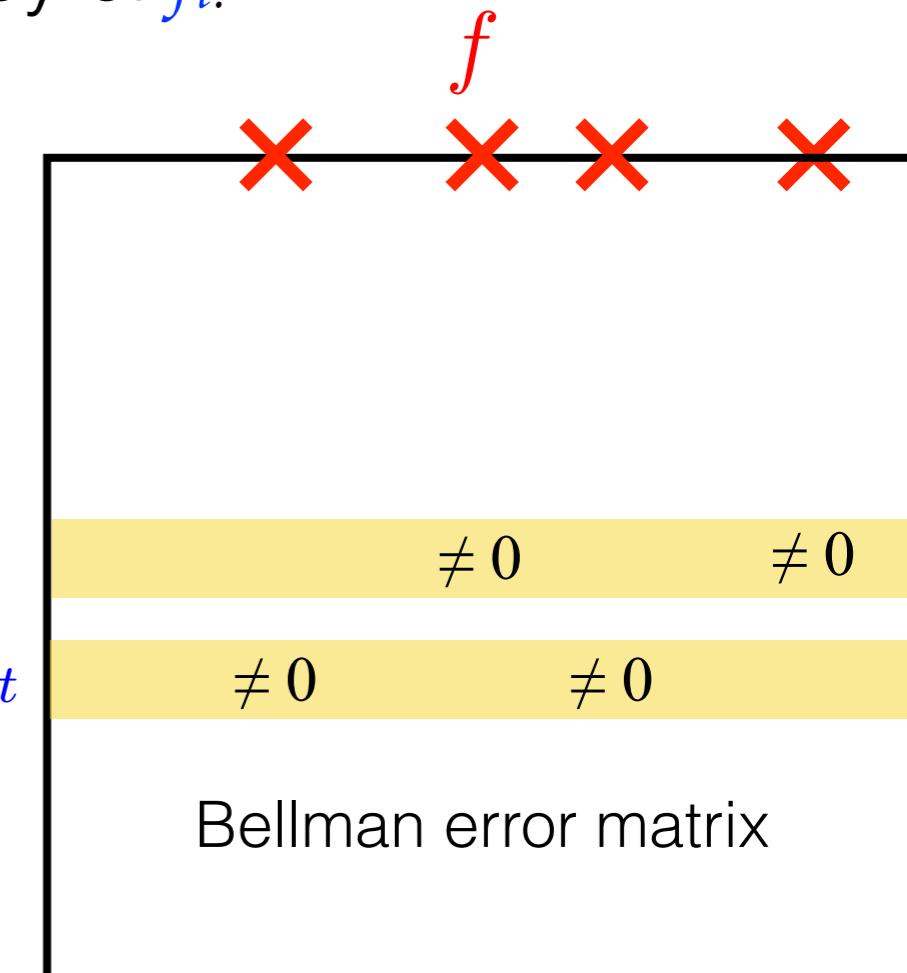
(Optimism-Led Iterative Value-function Elimination)

$F_1 := F$ . // version space

(Ignoring statistical slackness parameters)

For iteration  $t=1, 2, \dots$

- Choose  $f_t$  as the  $f \in F_t$  that maximizes  $v_f := \max_{a \in \mathcal{A}} f(x^0, a)$
- Estimate the value of  $\pi_t$  — the greedy policy of  $f_t$ .
  - If  $J(\pi_t) \geq v_{f_t}$  ( $\geq v_{Q^*} = J(\pi^*)$ ), return  $\pi_t$ .
- Estimate  $\mathcal{E}^h(f, \pi_t)$  for all  $f, h$ .
- Eliminate  $f$  s.t.  $\mathcal{E}^h(f, \pi_t) \neq 0, \forall h$   
 $\Rightarrow F_{t+1}.$



Bellman error matrix

# Sample complexity analysis

For iteration  $t=1, 2, \dots$

- Estimate the value of  $\pi_t$  — the greedy policy of  $f_t$ .

How many sample trajectories  
needed?

- Estimate  $\mathcal{E}^h(f, \pi_t)$  for all  $f, h$ .

# Sample complexity analysis

For iteration  $t=1, 2, \dots$

Run  $\pi_t$  for  $O(1/\varepsilon^2)$  episodes — Done.

- Estimate the value of  $\pi_t$  — the greedy policy of  $f_t$ .

How many sample trajectories  
needed?

- Estimate  $\mathcal{E}^h(f, \pi_t)$  for all  $f, h$ .

# Sample complexity analysis

For iteration  $t=1, 2, \dots$

Run  $\pi_t$  for  $O(1/\varepsilon^2)$  episodes — Done.

- Estimate the value of  $\pi_t$  — the greedy policy of  $f_t$ .  
How many sample trajectories  
needed?
- Estimate  $\mathcal{E}^h(f, \pi_t)$  for all  $f, h$ .  $\mathbb{E}_{a_{1:h-1} \sim \pi_t, a_h \sim f} [f \cdots]$

# Sample complexity analysis

For iteration  $t=1, 2, \dots$

Run  $\pi_t$  for  $O(1/\varepsilon^2)$  episodes — Done.

- Estimate the value of  $\pi_t$  — the greedy policy of  $f_t$ .

How many sample trajectories  
needed?

- Estimate  $\mathcal{E}^h(f, \pi_t)$  for all  $f, h$ .  $\mathbb{E}_{a_{1:h-1} \sim \pi_t, a_h \sim f} [f \cdots]$

- Naive: collect data with  $a_{1:h-1} \sim \pi_t, a_h \sim f$  for each  $f$

# Sample complexity analysis

For iteration  $t=1, 2, \dots$

Run  $\pi_t$  for  $O(1/\varepsilon^2)$  episodes — Done.

- Estimate the value of  $\pi_t$  — the greedy policy of  $f_t$ .

How many sample trajectories  
needed?

- Estimate  $\mathcal{E}^h(f, \pi_t)$  for all  $f, h$ .  $\mathbb{E}_{a_{1:h-1} \sim \pi_t, a_h \sim f} [f \cdots]$

- Naive: collect data with  $a_{1:h-1} \sim \pi_t, a_h \sim f$  for each  $f$
- $|F|$  samples — too many

# Sample complexity analysis

For iteration  $t=1, 2, \dots$

Run  $\pi_t$  for  $O(1/\varepsilon^2)$  episodes — Done.

- Estimate the value of  $\pi_t$  — the greedy policy of  $f_t$ .

How many sample trajectories  
needed?

- Estimate  $\mathcal{E}^h(f, \pi_t)$  for all  $f, h$ .  $\mathbb{E}_{a_{1:h-1} \sim \pi_t, a_h \sim f} [f \cdots]$

- Naive: collect data with  $a_{1:h-1} \sim \pi_t, a_h \sim f$  for each  $f$
- $|F|$  samples — too many
- Instead:  $a_{1:h-1} \sim \pi_t, a_h \sim \text{Unif}(A)$  & Importance Sampling

# Sample complexity analysis

For iteration  $t=1, 2, \dots$

Run  $\pi_t$  for  $O(1/\varepsilon^2)$  episodes — Done.

- Estimate the value of  $\pi_t$  — the greedy policy of  $f_t$ .

How many sample trajectories  
needed?

- Estimate  $\mathcal{E}^h(f, \pi_t)$  for all  $f, h$ .  $\mathbb{E}_{a_{1:h-1} \sim \pi_t, a_h \sim f} [f \cdots]$

- Naive: collect data with  $a_{1:h-1} \sim \pi_t, a_h \sim f$  for each  $f$
- $|F|$  samples — too many
- Instead:  $a_{1:h-1} \sim \pi_t, a_h \sim \text{Unif}(A)$  & Importance Sampling
- 1 sample of size  $O(|A| \log |F| / \varepsilon^2)$  — works for all  $f$  simultaneously

# Sample complexity analysis

For iteration  $t=1, 2, \dots$

How many iterations???

Run  $\pi_t$  for  $O(1/\varepsilon^2)$  episodes — Done.

- Estimate the value of  $\pi_t$  — the greedy policy of  $f_t$ .

How many sample trajectories  
needed?

- Estimate  $\mathcal{E}^h(f, \pi_t)$  for all  $f, h$ .  $\mathbb{E}_{a_{1:h-1} \sim \pi_t, a_h \sim f} [f \dots]$

- Naive: collect data with  $a_{1:h-1} \sim \pi_t, a_h \sim f$  for each  $f$
- $|F|$  samples — too many
- Instead:  $a_{1:h-1} \sim \pi_t, a_h \sim \text{Unif}(A)$  & Importance Sampling
- 1 sample of size  $O(|A| \log |F| / \varepsilon^2)$  — works for all  $f$  simultaneously

# Sample complexity analysis

Claim: If no statistical errors, # **iterations**  $\leq$  **Bellman rank**.

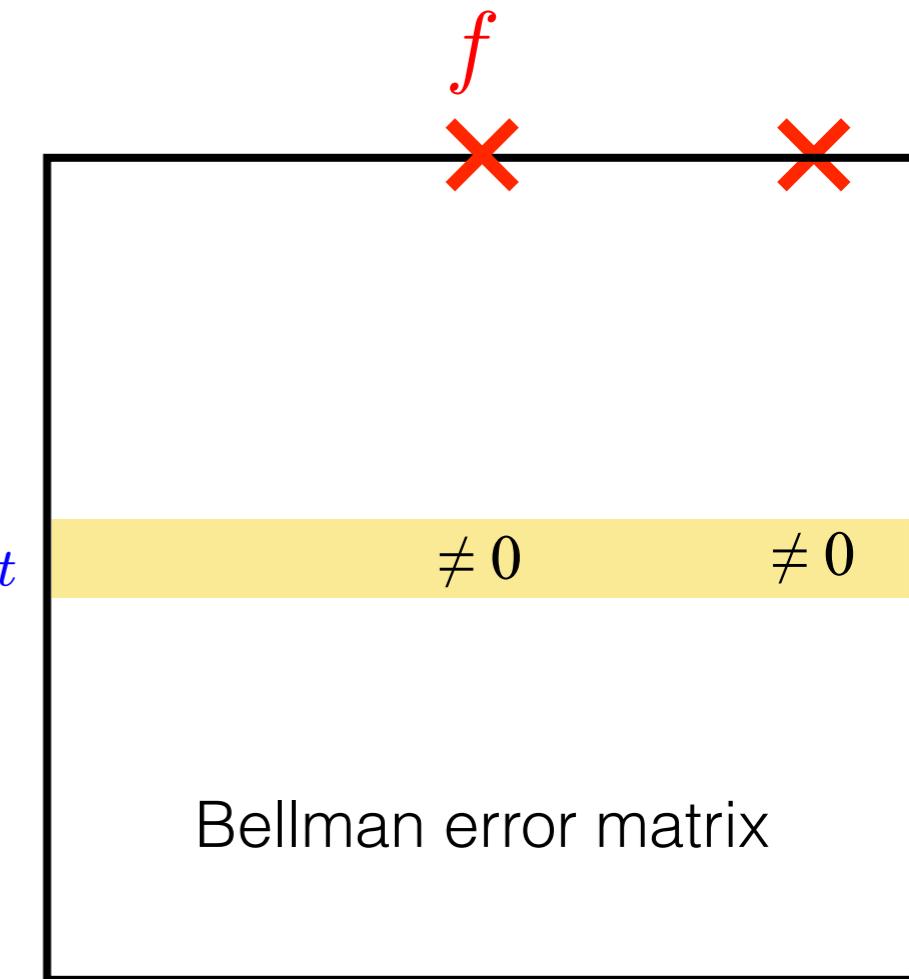
$f$

$\pi_t$

Bellman error matrix

# Sample complexity analysis

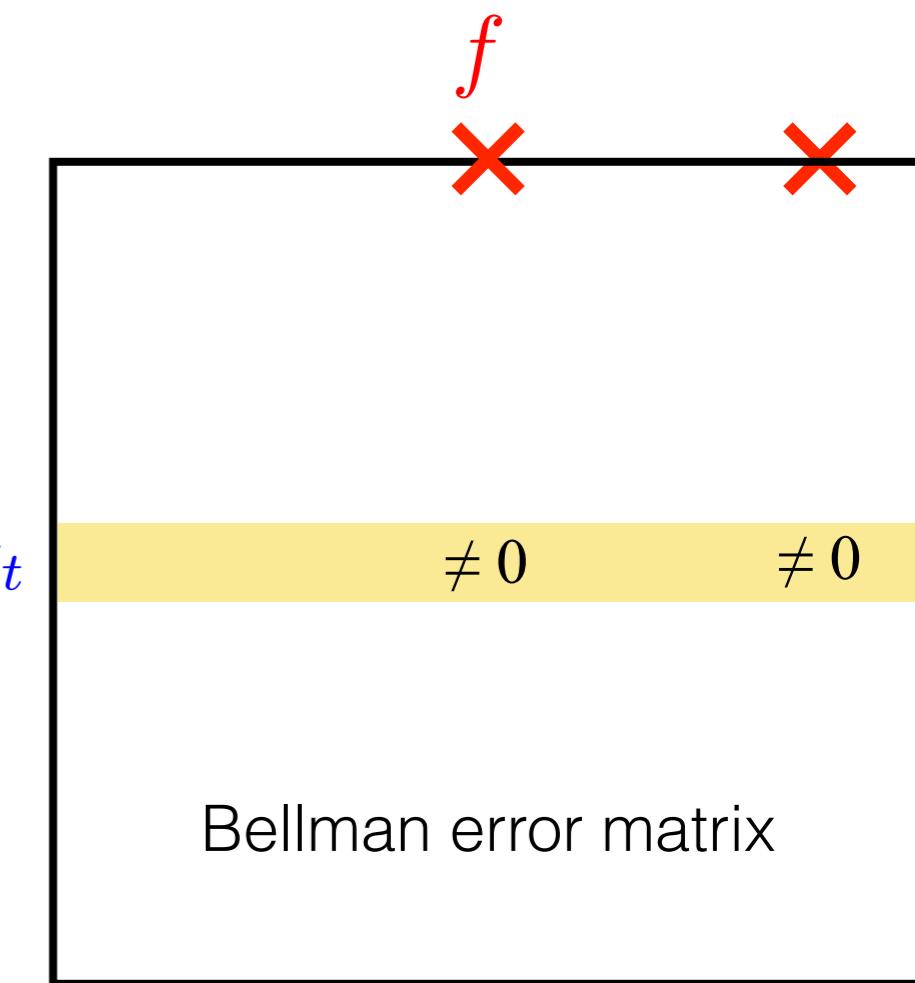
Claim: If no statistical errors, # **iterations**  $\leq$  **Bellman rank**.



# Sample complexity analysis

Claim: If no statistical errors, # **iterations**  $\leq$  **Bellman rank**.

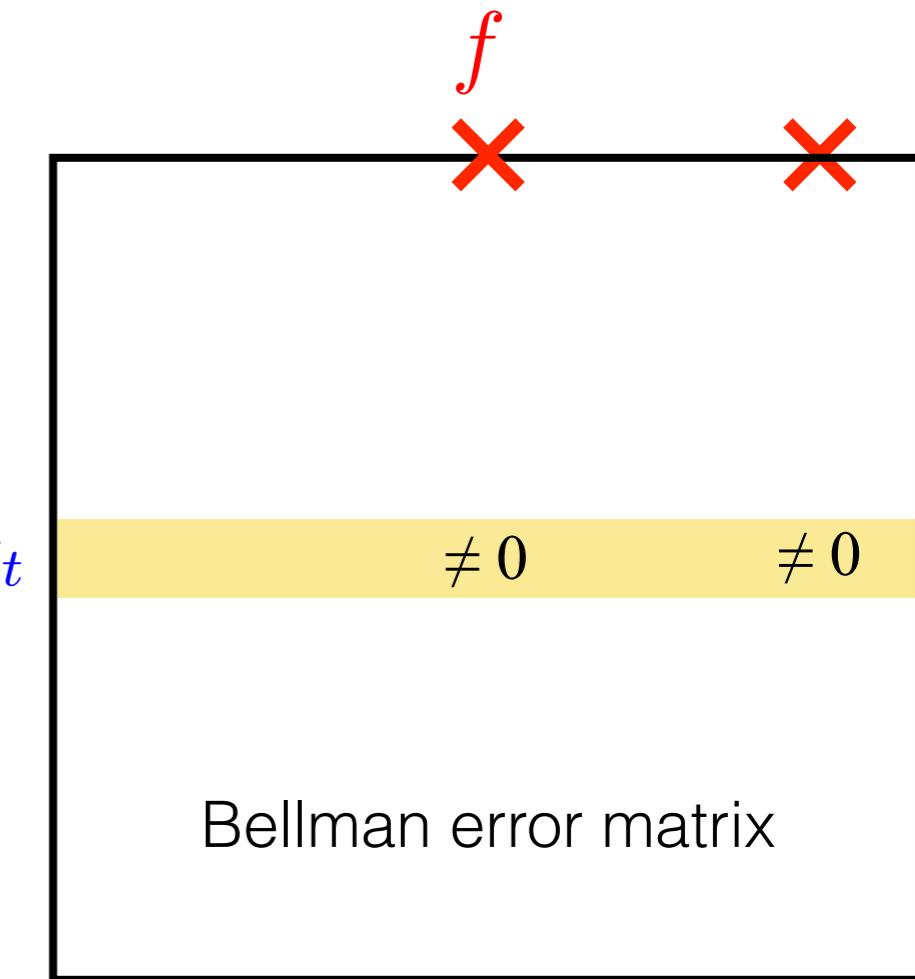
- All surviving  $f$  have all-0 columns so far



# Sample complexity analysis

Claim: If no statistical errors, # **iterations**  $\leq$  **Bellman rank**.

- All surviving  $f$  have all-0 columns so far
- Will show: some  $f$  has “ $\neq 0$ ” in the next iteration

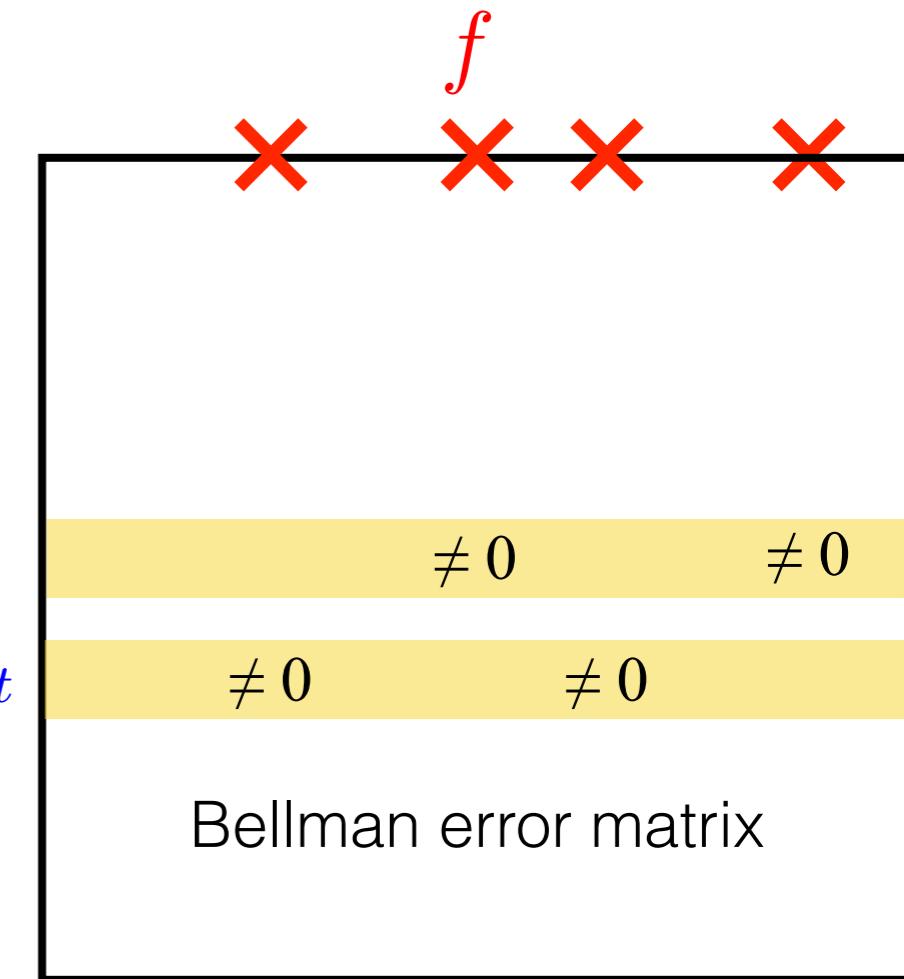


Bellman error matrix

# Sample complexity analysis

Claim: If no statistical errors, # **iterations**  $\leq$  **Bellman rank**.

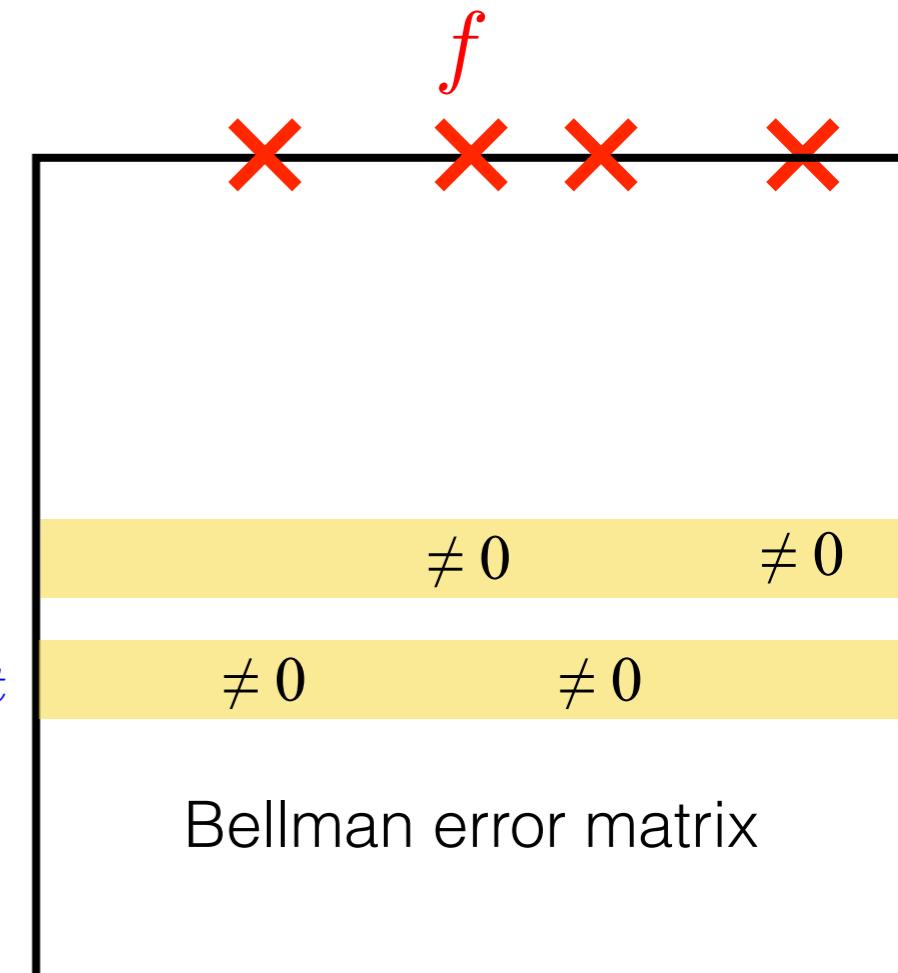
- All surviving  $f$  have all-0 columns so far
- Will show: some  $f$  has “ $\neq 0$ ” in the next iteration



# Sample complexity analysis

Claim: If no statistical errors, # **iterations**  $\leq$  **Bellman rank**.

- All surviving  $f$  have all-0 columns so far
- Will show: some  $f$  has “ $\neq 0$ ” in the next iteration
- Then: linearly independent rows  $\Rightarrow$  #iterations  $\leq$  matrix rank



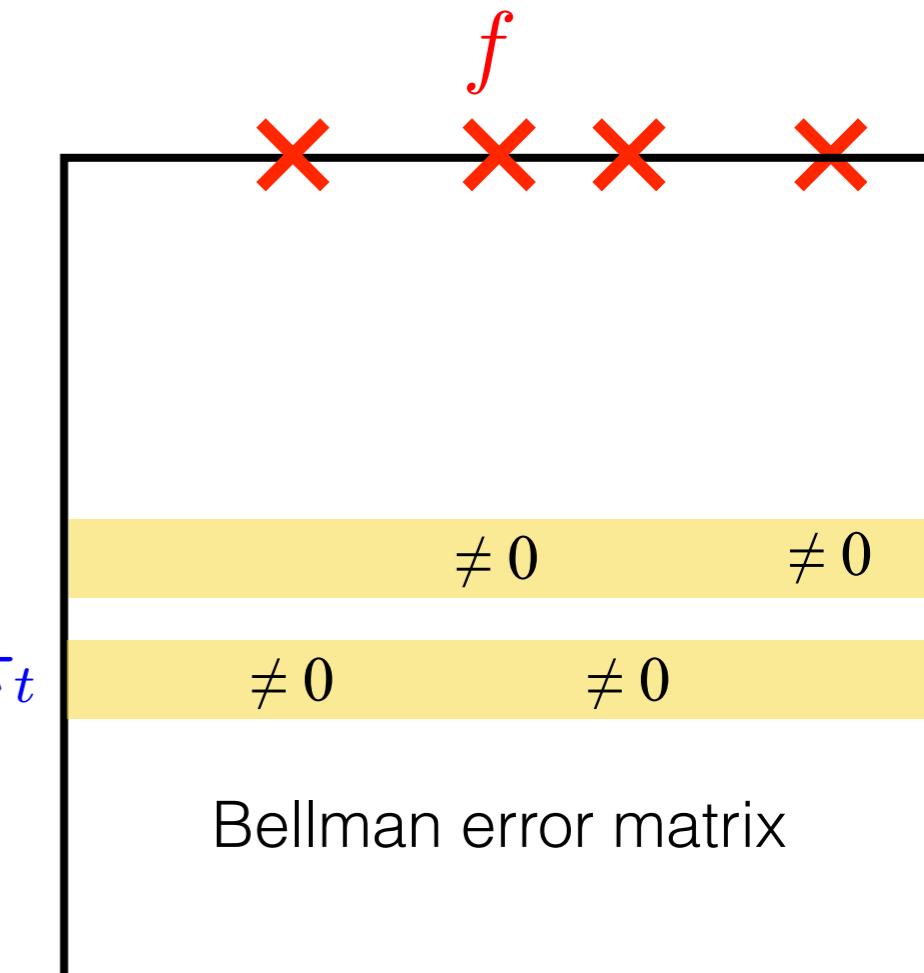
# Sample complexity analysis

Claim: If no statistical errors, # **iterations**  $\leq$  **Bellman rank**.

- All surviving  $f$  have **all-0 columns** so far
  - Will show: some  $f$  has “ $\neq 0$ ” in the next iteration
- Then: linearly independent rows  $\Rightarrow$  #iterations  $\leq$  matrix rank

$f_t$  has “ $\neq 0$ ” unless terminate:  
(recall  $\pi_t$  is greedy wrt  $f_t$ )

$$v_{f_t} - J(\pi_t) = \sum_{h=1}^H \mathcal{E}^h(f_t, \pi_t)$$



# Sample complexity analysis

Claim: If no statistical errors, # **iterations**  $\leq$  **Bellman rank**.

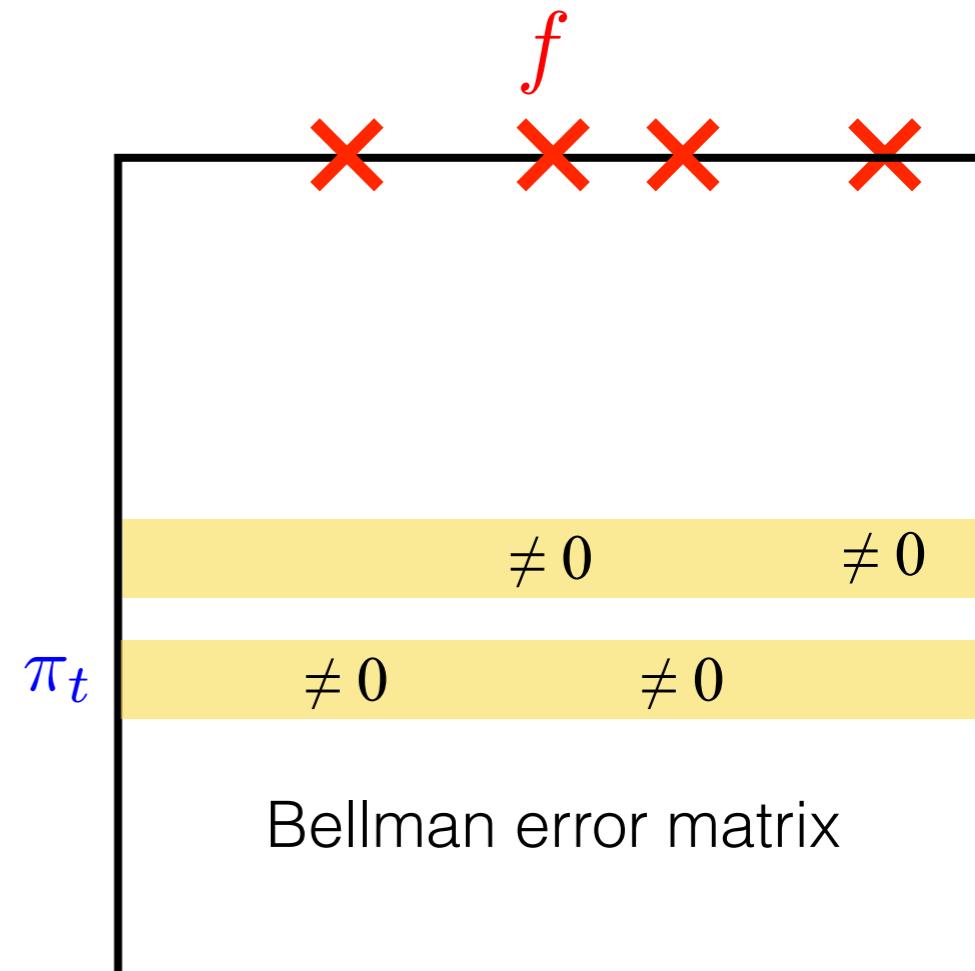
- All surviving  $f$  have **all-0 columns** so far
  - Will show: some  $f$  has “ $\neq 0$ ” in the next iteration
- Then: linearly independent rows  $\Rightarrow$  #iterations  $\leq$  matrix rank

$f_t$  has “ $\neq 0$ ” unless terminate:  
(recall  $\pi_t$  is greedy wrt  $f_t$ )

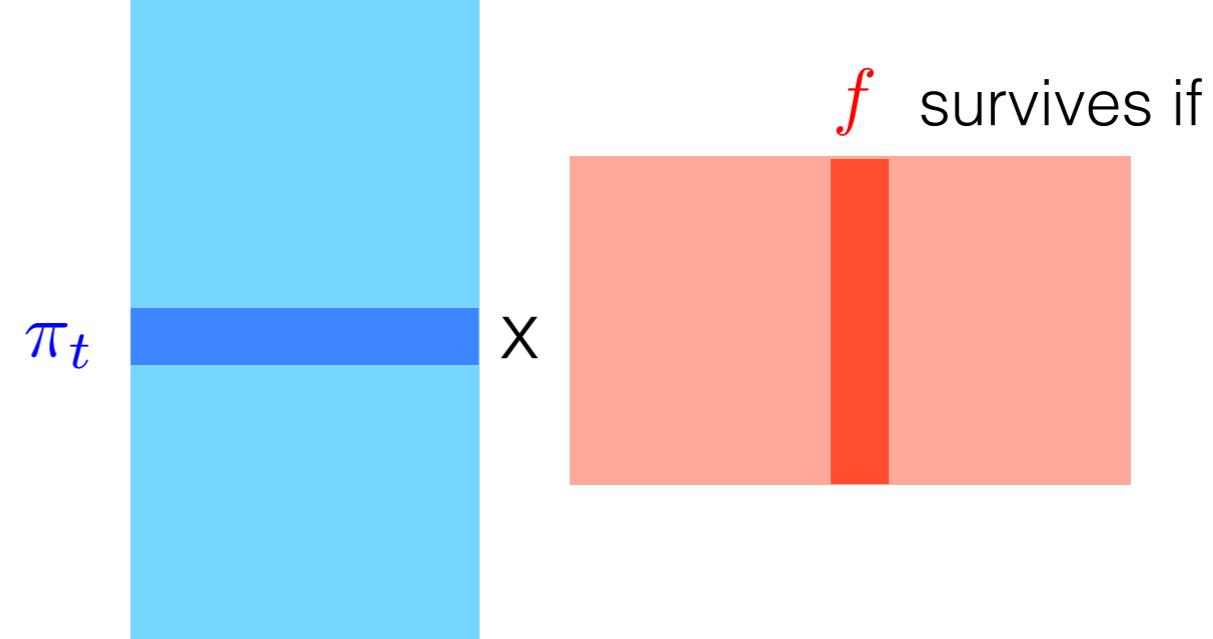
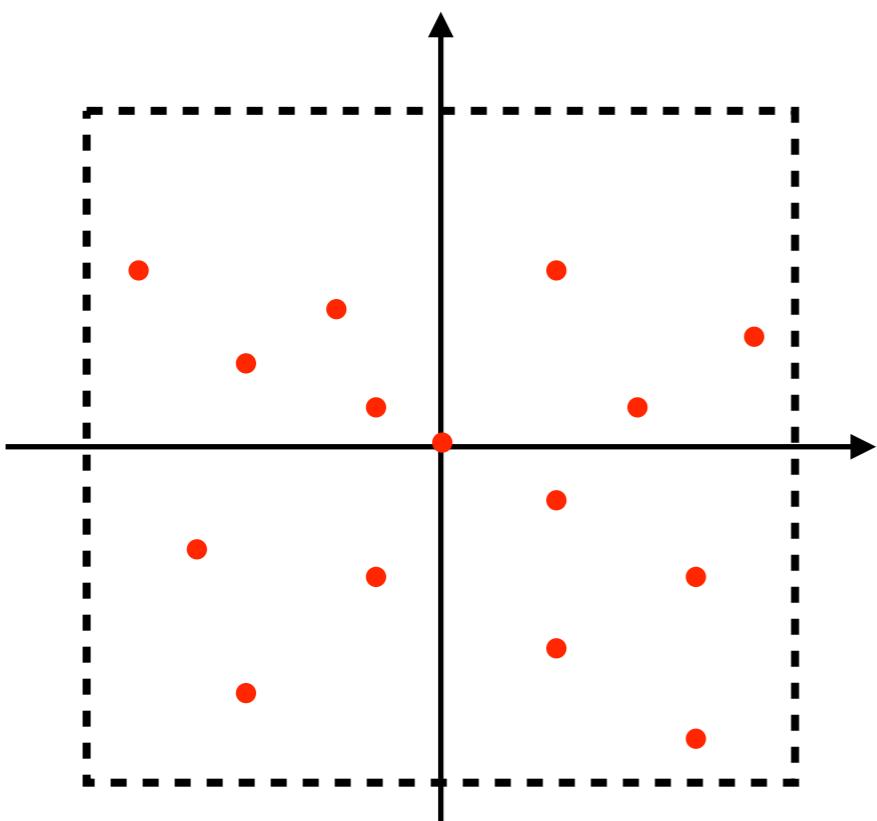
$$0 < v_{f_t} - J(\pi_t) = \sum_{h=1}^H \mathcal{E}^h(f_t, \pi_t)$$



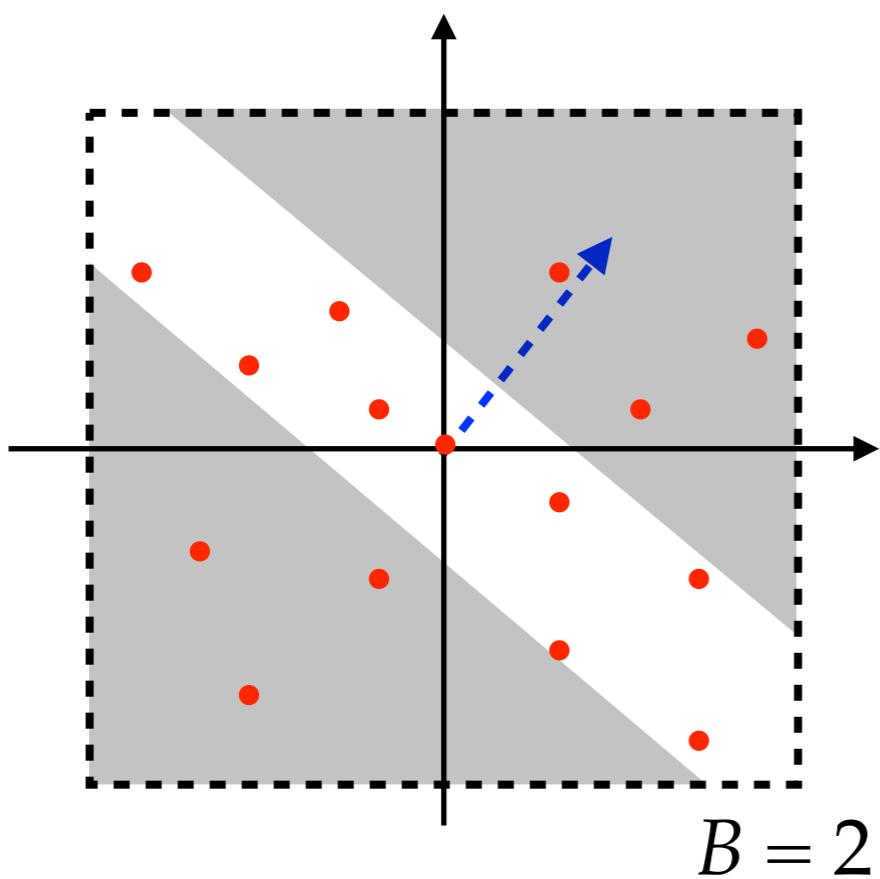
Optimized:  $v_{f_t} \geq v_{Q^*} = J(\pi^*)$



# Sample complexity of OLIVE

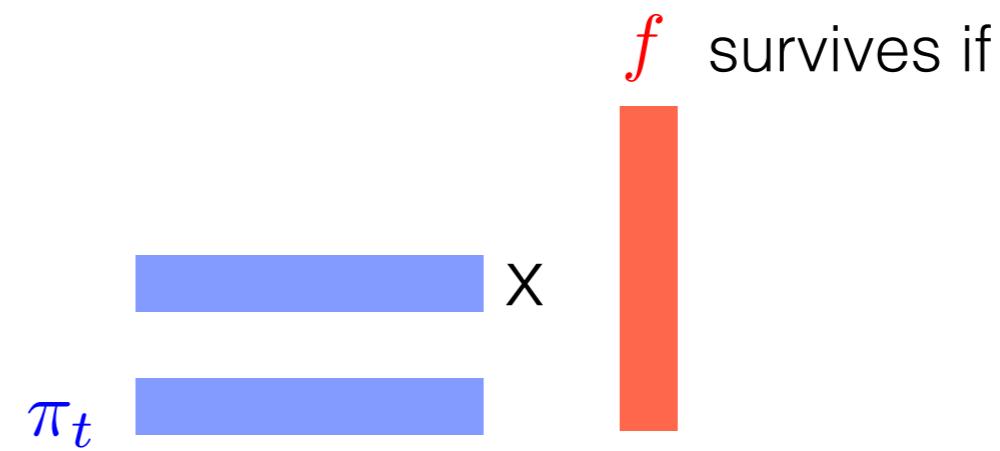
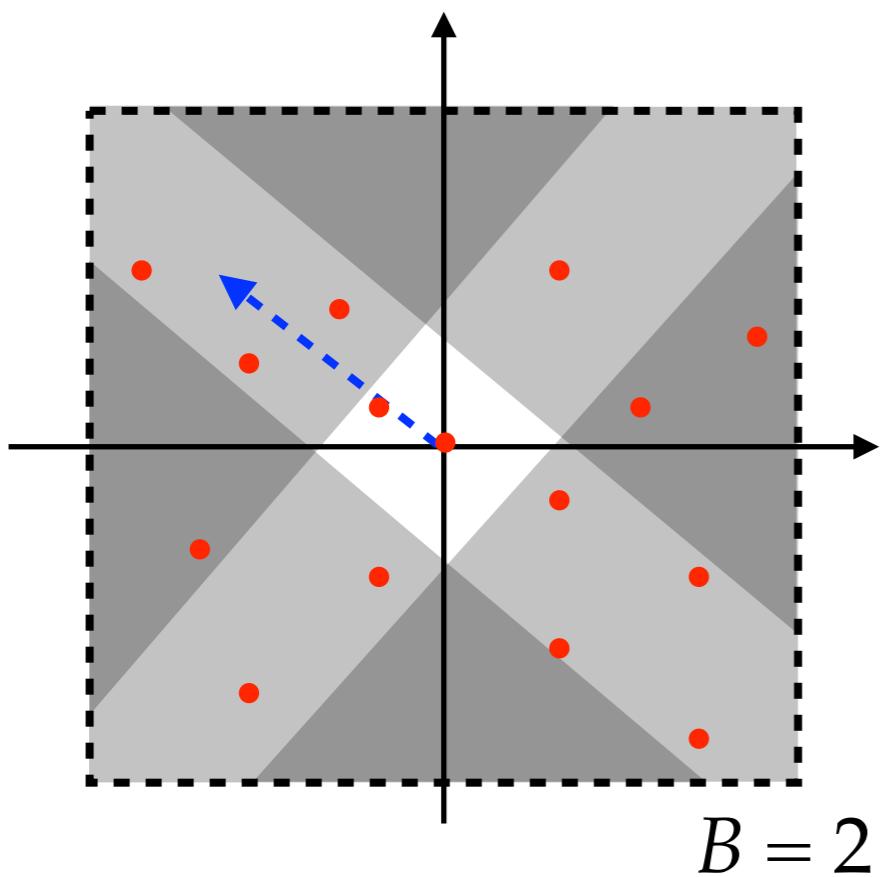


# Sample complexity of OLIVE



$\pi_t$  X  $f$  survives if

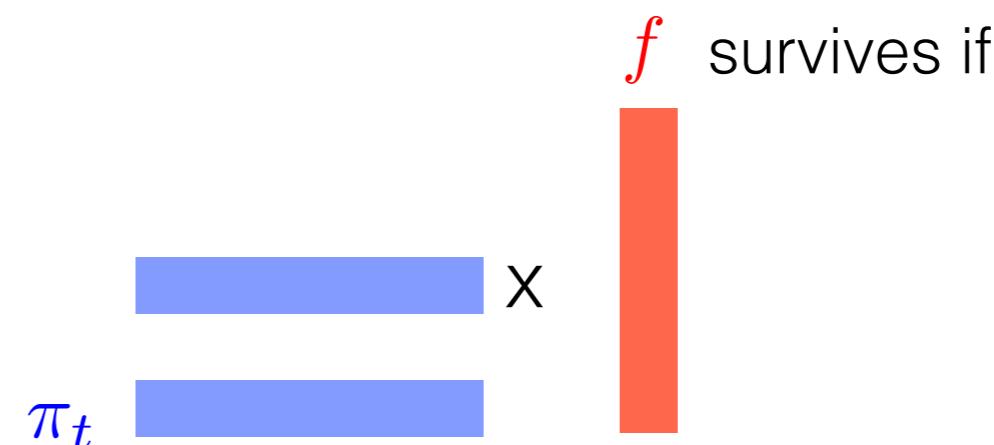
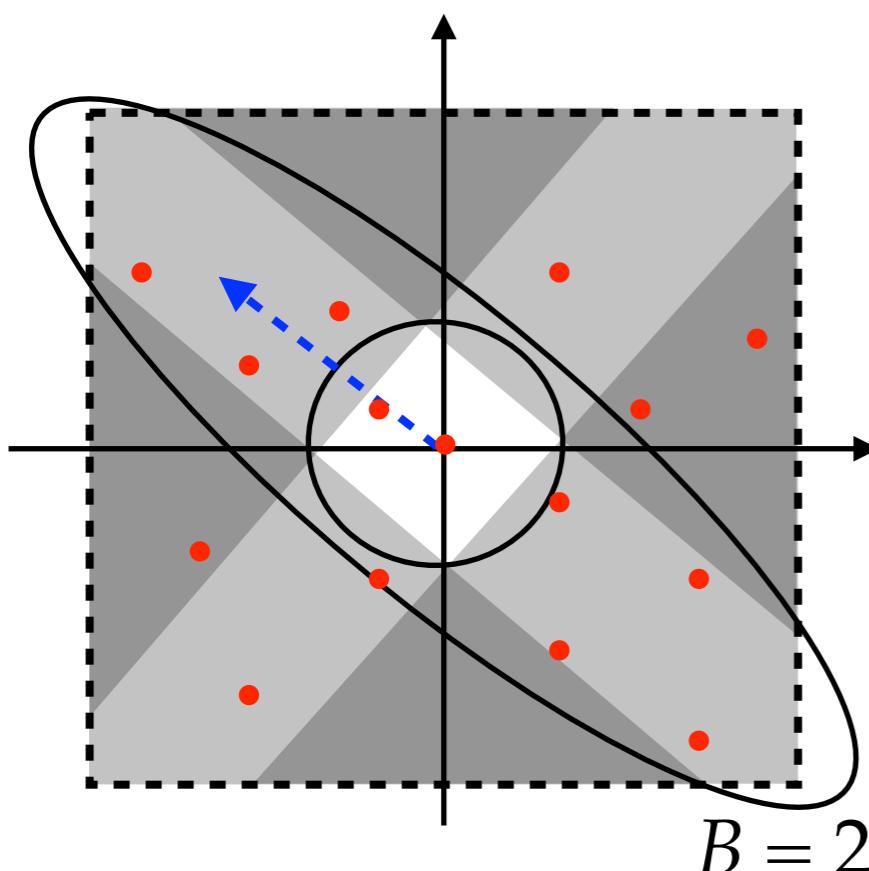
# Sample complexity of OLIVE



# Sample complexity of OLIVE

**Theorem:** If  $Q^* \in \mathcal{F}$ , w.p.  $\geq 1-\delta$ , OLIVE returns a  $\varepsilon$ -optimal policy after acquiring the following number of trajectories

Bellman rank  $\tilde{O} \left( \frac{\mathbf{B}^2 H^3 |\mathcal{A}|}{\epsilon^2} \log(|\mathcal{F}|/\delta) \right)$



# Bellman Equations revisited

$$\mathbb{E}_{\substack{a_{1:h-1} \sim \pi' \\ a_h \sim f}} [f(x_h, a_h) - r_h - \max_{a \in \mathcal{A}} f(x_{h+1}, a)] = 0$$

# Bellman Equations revisited

$$\mathbb{E}_{\substack{a_1:h-1 \sim \pi' \\ a_h \sim f}} [f(x_h, a_h) - r_h - \max_{a \in \mathcal{A}} f(x_{h+1}, a)] = 0$$

- $f$  on non-greedy actions never used!

# Bellman Equations revisited

$$\mathbb{E}_{\substack{a_{1:h-1} \sim \pi' \\ a_h \sim \pi}} [g(x_h) - r_h - g(x_{h+1})] = 0$$

- $f$  on non-greedy actions never used!
- Reparametrize:  $f \Rightarrow (g, \pi); F \Rightarrow G, \Pi$ .

# Bellman Equations revisited

$$\mathbb{E}_{\substack{a_{1:h-1} \sim \pi' \\ a_h \sim \pi}} [g(x_h) - r_h - g(x_{h+1})] = 0$$

- $f$  on non-greedy actions never used!
- Reparametrize:  $f \Rightarrow (g, \pi); F \Rightarrow G, \Pi$ .
- Bellman equations for policy evaluation

# Bellman Equations revisited

$$\mathbb{E}_{\substack{a_{1:h-1} \sim \pi' \\ a_h \sim \pi}} [g(x_h) - r_h - g(x_{h+1})] = 0$$

- $f$  on non-greedy actions never used!
- Reparametrize:  $f \Rightarrow (g, \pi); F \Rightarrow G, \Pi$ .
- Bellman equations for policy evaluation
  - Even if  $\pi^* \notin \Pi$ , can still compete with *any*  $\pi \in \Pi$  whose policy-specific value function is (approx.) in  $G$

# Bellman Equations revisited

$$\mathbb{E}_{\substack{a_{1:h-1} \sim \pi' \\ a_h \sim \pi}} [g(x_h) - r_h - g(x_{h+1})] = 0$$

- $f$  on non-greedy actions never used!
- Reparametrize:  $f \Rightarrow (g, \pi); F \Rightarrow G, \Pi$ .
- Bellman equations for policy evaluation
  - Even if  $\pi^* \notin \Pi$ , can still compete with *any*  $\pi \in \Pi$  whose policy-specific value function is (approx.) in  $G$
  - Allow infinite classes with VC-type dimensions

# Computational Efficiency

[Dann+JKALS, arXiv'18]

- OLIVE requires solving a constrained optimization problem
  - $\textcolor{red}{f} \in \mathcal{F}_t \Leftrightarrow \textcolor{red}{f} \in \mathcal{F}, \mathcal{E}^h(\textcolor{red}{f}, \pi_{t'}) \neq 0, \forall h \in [H], t' \in [t - 1]$
  - $f_t = \max v_{\textcolor{red}{f}}$ , subject to the constraints.

# Computational Efficiency

[Dann+JKALS, arXiv'18]

- OLIVE requires solving a constrained optimization problem
  - $\textcolor{red}{f} \in \mathcal{F}_t \Leftrightarrow \textcolor{red}{f} \in \mathcal{F}, \mathcal{E}^h(\textcolor{red}{f}, \pi_{t'}) \neq 0, \forall h \in [H], t' \in [t - 1]$
  - $f_t = \max v_{\textcolor{red}{f}}$ , subject to the constraints.
- How to access  $F$  (or  $G, \Pi$ )?

# Computational Efficiency

[Dann+JKALS, arXiv'18]

- OLIVE requires solving a constrained optimization problem
  - $\textcolor{red}{f} \in \mathcal{F}_t \Leftrightarrow \textcolor{red}{f} \in \mathcal{F}, \mathcal{E}^h(\textcolor{red}{f}, \pi_{t'}) \neq 0, \forall h \in [H], t' \in [t - 1]$
  - $f_t = \max v_{\textcolor{red}{f}}$ , subject to the constraints.
- How to access  $F$  (or  $G, \Pi$ )?
  - **Oracles.** E.g.,

# Computational Efficiency

[Dann+JKALS, arXiv'18]

- OLIVE requires solving a constrained optimization problem
  - $\textcolor{red}{f} \in \mathcal{F}_t \Leftrightarrow \textcolor{red}{f} \in \mathcal{F}, \mathcal{E}^h(\textcolor{red}{f}, \pi_{t'}) \neq 0, \forall h \in [H], t' \in [t - 1]$
  - $f_t = \max v_{\textcolor{red}{f}}$ , subject to the constraints.
- How to access  $F$  (or  $G, \Pi$ )?
  - **Oracles.** E.g.,
    - **Cost-sensitive Classification** for  $\Pi \subset (X \rightarrow A)$   
Given  $\{(x^i \in X, c^i \in R^A)\}_{i \in [n]}$ , oracle minimizes  $\sum_{i=1}^n c^i(\pi(x^i))$

# Computational Efficiency

[Dann+JKALS, arXiv'18]

- OLIVE requires solving a constrained optimization problem
  - $\textcolor{red}{f} \in \mathcal{F}_t \Leftrightarrow \textcolor{red}{f} \in \mathcal{F}, \mathcal{E}^h(\textcolor{red}{f}, \pi_{t'}) \neq 0, \forall h \in [H], t' \in [t - 1]$
  - $f_t = \max v_{\textcolor{red}{f}}$ , subject to the constraints.
- How to access  $F$  (or  $G, \Pi$ )?
  - **Oracles.** E.g.,
    - **Cost-sensitive Classification** for  $\Pi \subset (X \rightarrow A)$   
Given  $\{(x^i \in X, c^i \in R^A)\}_{i \in [n]}$ , oracle minimizes  $\sum_{i=1}^n c^i(\pi(x^i))$
    - Linear optimization, squared-loss regression for  $G \subset (X \rightarrow R)$

# Computational Efficiency

[Dann+JKALS, arXiv'18]

- OLIVE requires solving a constrained optimization problem
  - $\textcolor{red}{f} \in \mathcal{F}_t \Leftrightarrow \textcolor{red}{f} \in \mathcal{F}, \mathcal{E}^h(\textcolor{red}{f}, \pi_{t'}) \neq 0, \forall h \in [H], t' \in [t - 1]$
  - $f_t = \max v_{\textcolor{red}{f}}$ , subject to the constraints.
- How to access  $F$  (or  $G, \Pi$ )?
  - **Oracles.** E.g.,
    - **Cost-sensitive Classification** for  $\Pi \subset (X \rightarrow A)$   
Given  $\{(x^i \in X, c^i \in R^A)\}_{i \in [n]}$ , oracle minimizes  $\sum_{i=1}^n c^i(\pi(x^i))$
    - Linear optimization, squared-loss regression for  $G \subset (X \rightarrow R)$
    - Can we **reduce** the computation of OLIVE to **oracles**?

# Computational Efficiency

[Dann+JKALS, arXiv'18]

- No polynomial reduction exists

# Computational Efficiency

[Dann+JKALS, arXiv'18]

- No polynomial reduction exists
  - NP-hard even in **tabular** MDPs

# Computational Efficiency

[Dann+JKALS, arXiv'18]

- No polynomial reduction exists
  - NP-hard even in **tabular** MDPs
  - ERM also NP-hard — “absorbs” hardness?

# Computational Efficiency

[Dann+JKALS, arXiv'18]

- No polynomial reduction exists
  - NP-hard even in **tabular** MDPs
  - ERM also NP-hard — “absorbs” hardness?
  - Common oracles are **efficient** in the **tabular** case  
i.e.,  $|X|$  has finite cardinality,  $\Pi = X \rightarrow A$

# Computational Efficiency

[Dann+JKALS, arXiv'18]

- No polynomial reduction exists
  - NP-hard even in **tabular** MDPs
  - ERM also NP-hard — “absorbs” hardness?
  - Common oracles are **efficient** in the **tabular** case  
i.e.,  $|X|$  has finite cardinality,  $\Pi = X \rightarrow A$
- More recent advances: sample & computationally efficient alg for:

# Computational Efficiency

[Dann+JKALS, arXiv'18]

- No polynomial reduction exists
  - NP-hard even in **tabular** MDPs
  - ERM also NP-hard — “absorbs” hardness?
  - Common oracles are **efficient** in the **tabular** case  
i.e.,  $|X|$  has finite cardinality,  $\Pi = X \rightarrow A$
- More recent advances: sample & computationally efficient alg for:
  - linear MDPs (see previous lectures)

# Computational Efficiency

[Dann+JKALS, arXiv'18]

- No polynomial reduction exists
  - NP-hard even in **tabular** MDPs
  - ERM also NP-hard — “absorbs” hardness?
  - Common oracles are **efficient** in the **tabular** case  
i.e.,  $|X|$  has finite cardinality,  $\Pi = X \rightarrow A$
- More recent advances: sample & computationally efficient alg for:
  - linear MDPs (see previous lectures)
  - “block MDPs” (see previous “visual gridworld” example): latent-state decoding

# Computational Efficiency

[Dann+JKALS, arXiv'18]

- No polynomial reduction exists
  - NP-hard even in **tabular** MDPs
  - ERM also NP-hard — “absorbs” hardness?
  - Common oracles are **efficient** in the **tabular** case  
i.e.,  $|X|$  has finite cardinality,  $\Pi = X \rightarrow A$
- More recent advances: sample & computationally efficient alg for:
  - linear MDPs (see previous lectures)
  - “block MDPs” (see previous “visual gridworld” example): latent-state decoding
  - Check out COLT’21 tutorial: <https://rltheorybook.github.io/colt21tutorial>

# Detailed Analysis (with Statistical Errors)

$B$  (Bellman rank)

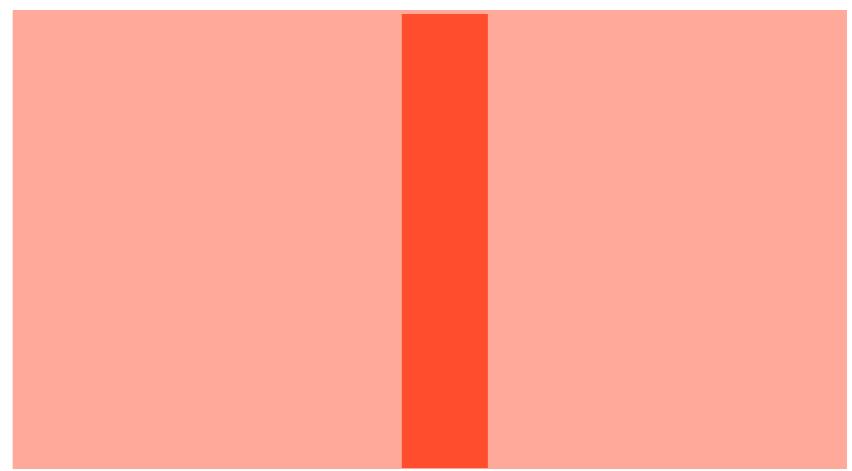
$f$

$\pi_{t-1}$

$=$   
 $\pi_{t-1}$

$\times$

$f$



$B$  (Bellman rank)

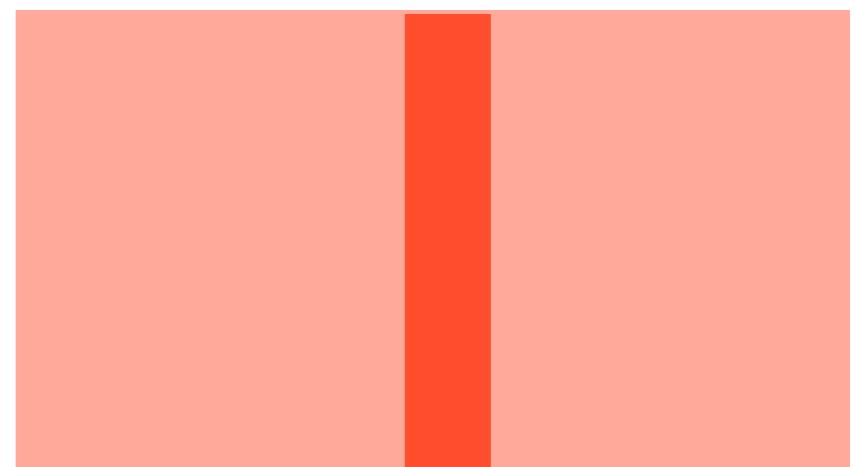
$f$

$\pi_{t-1}$

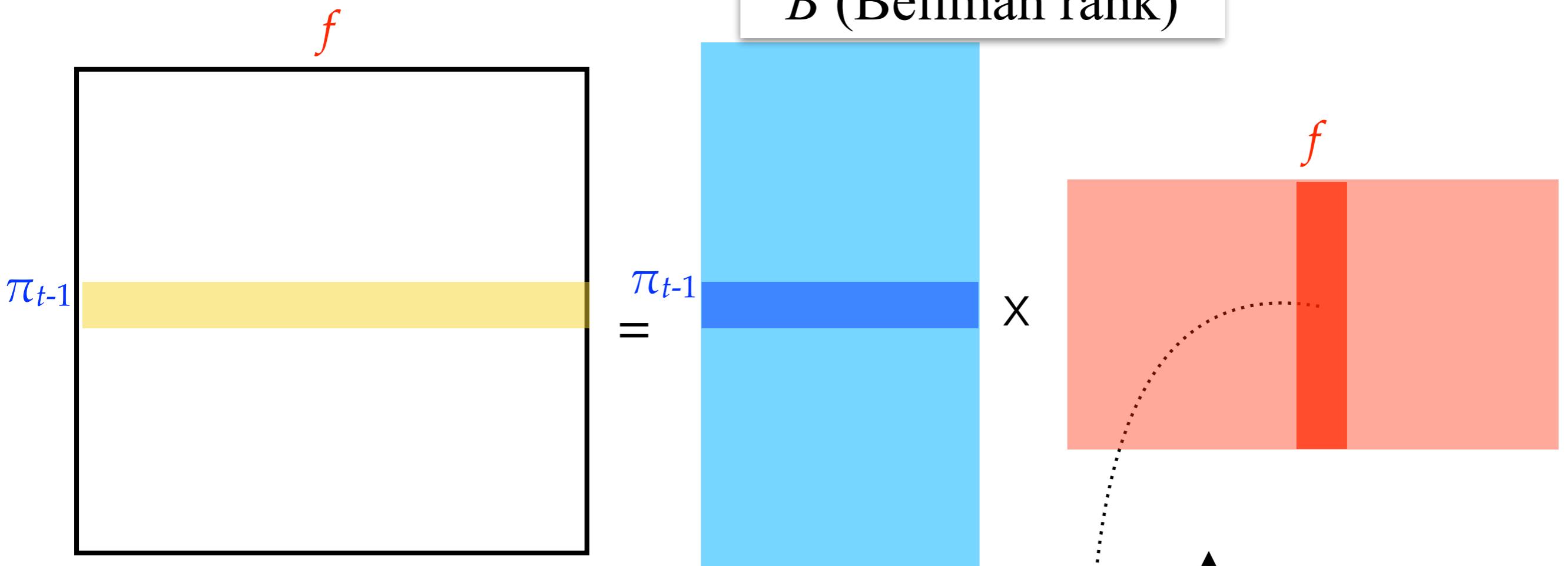
$=$   
 $\pi_{t-1}$

$\times$

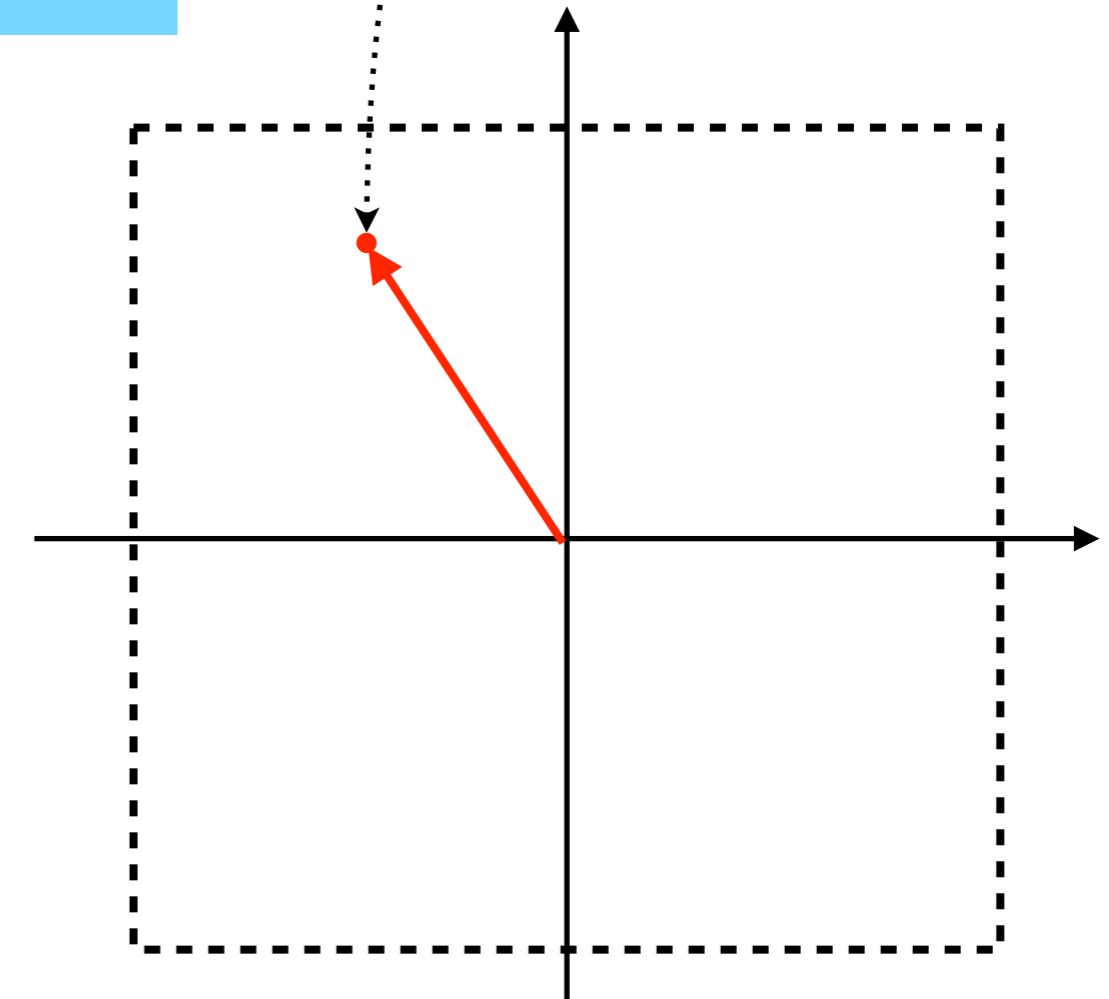
$f$



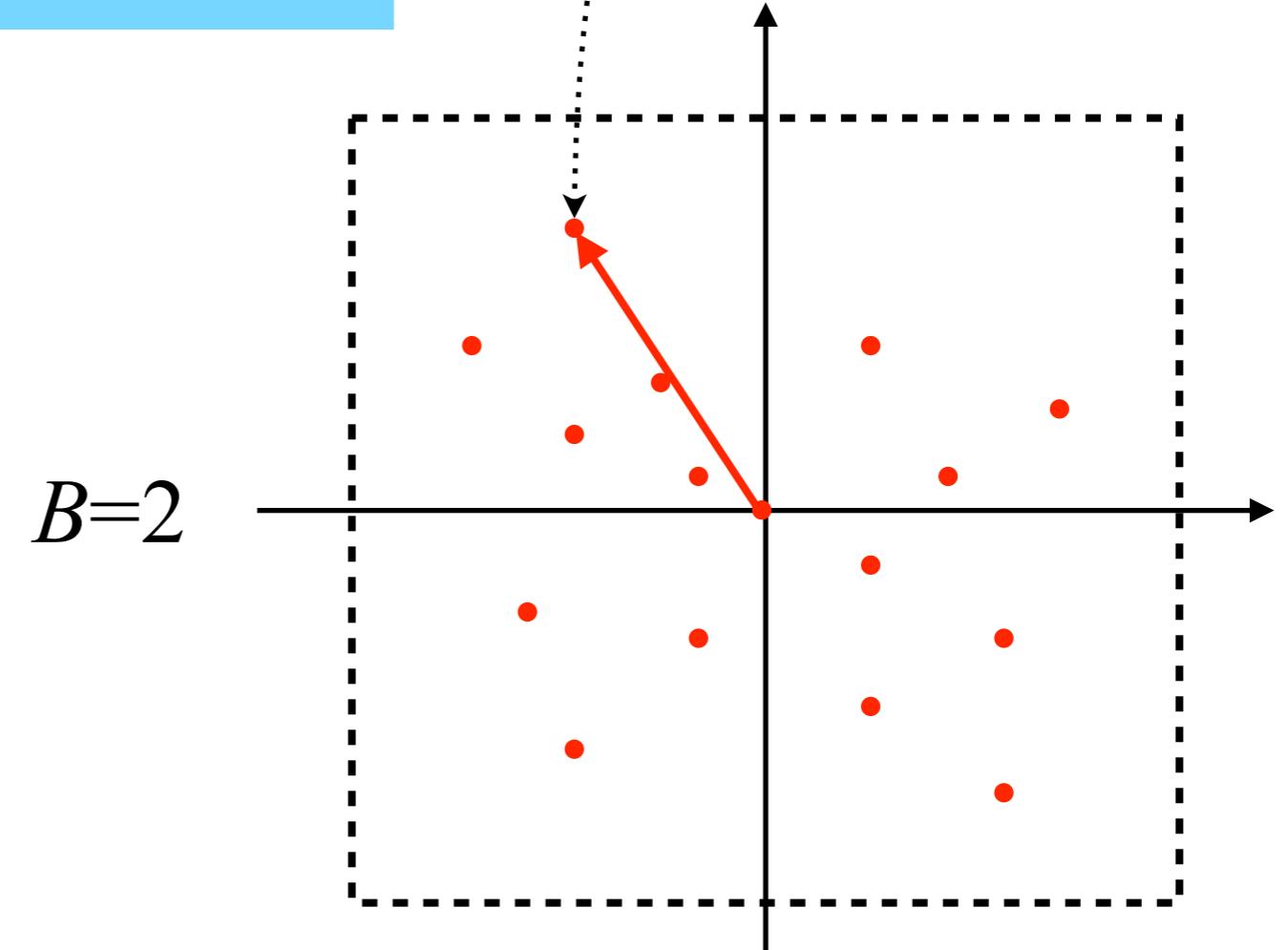
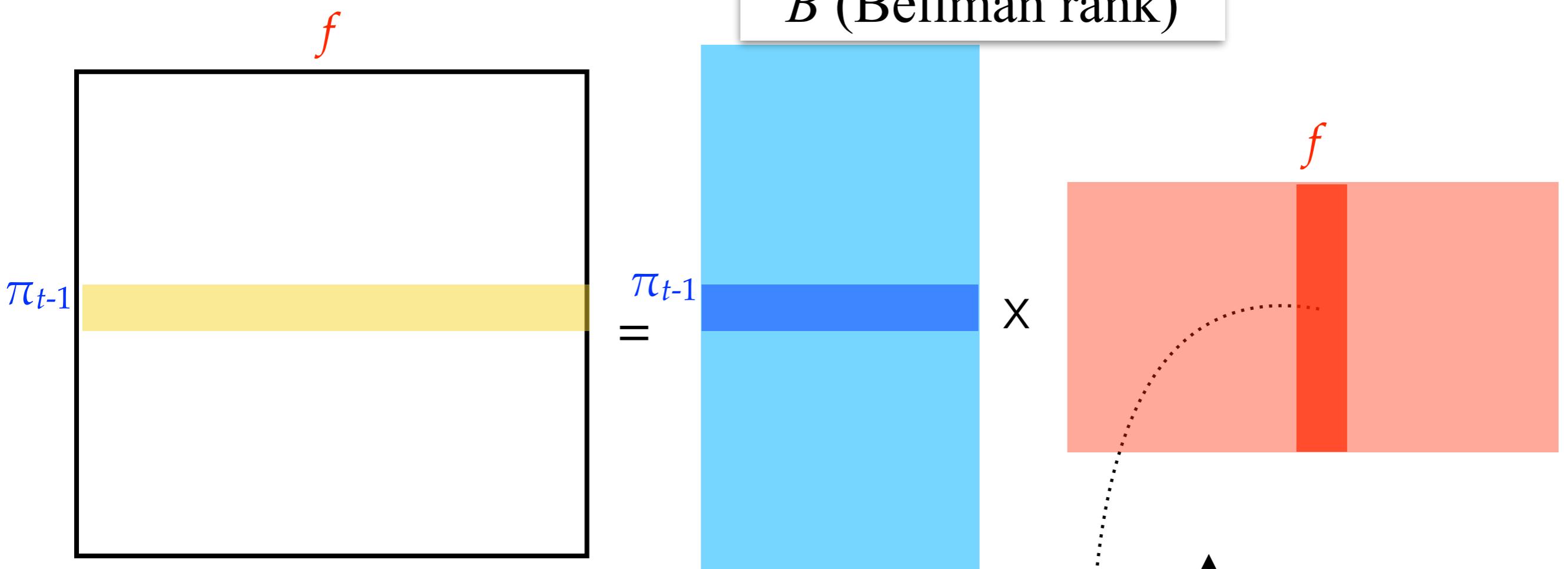
$B$  (Bellman rank)



$B=2$

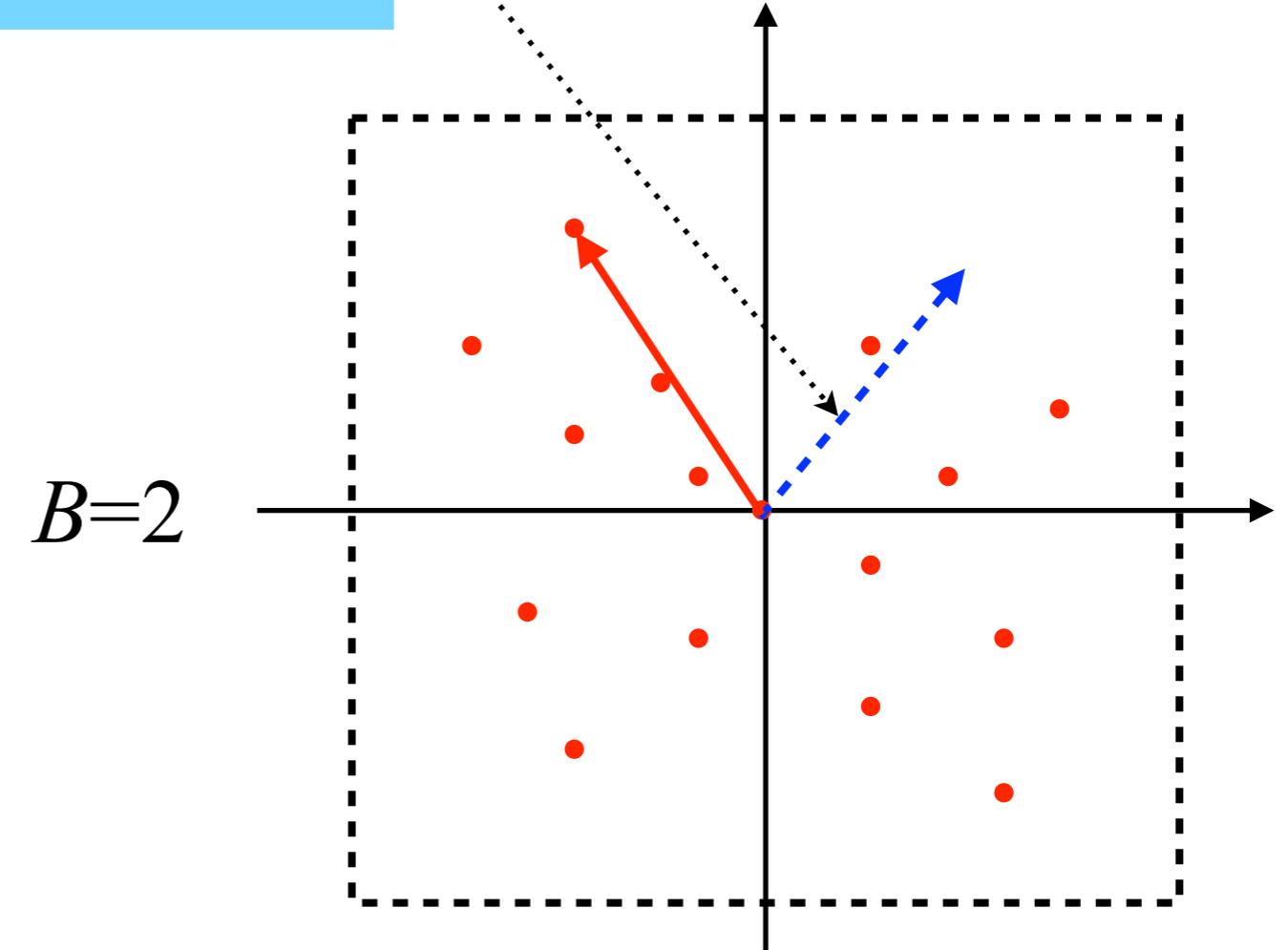


$B$  (Bellman rank)

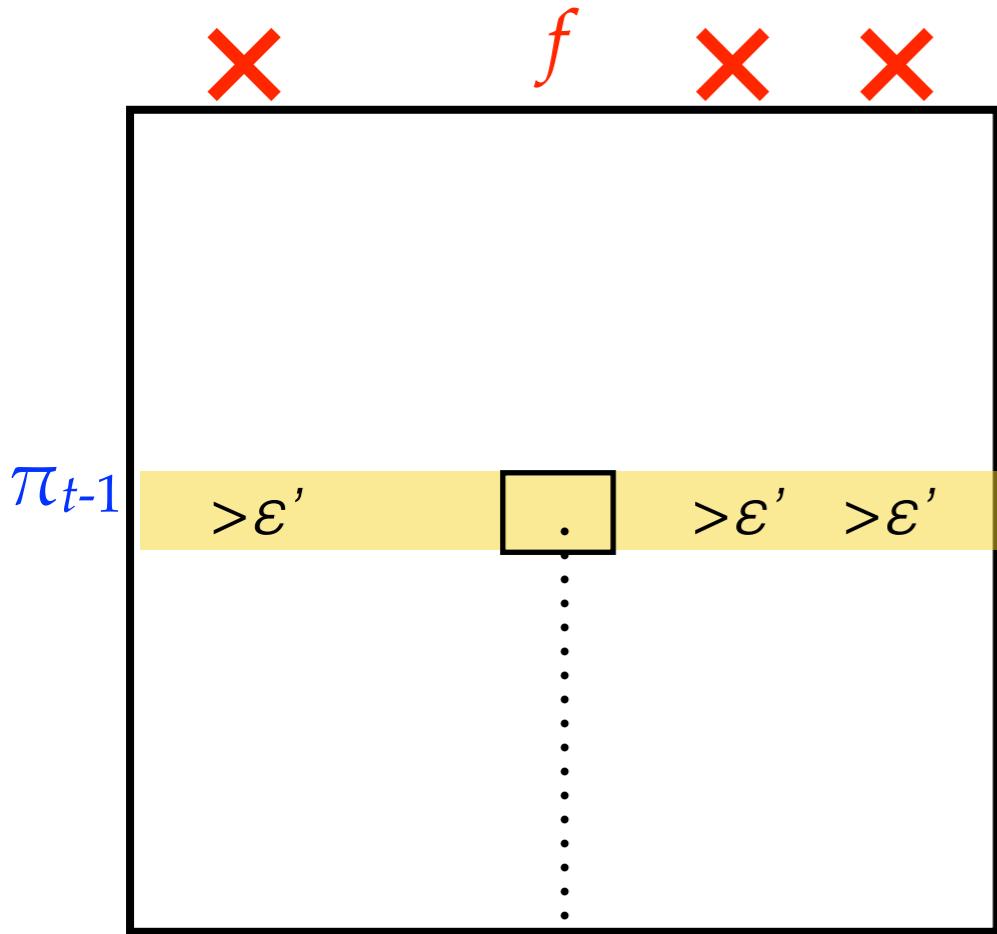


$B$  (Bellman rank)

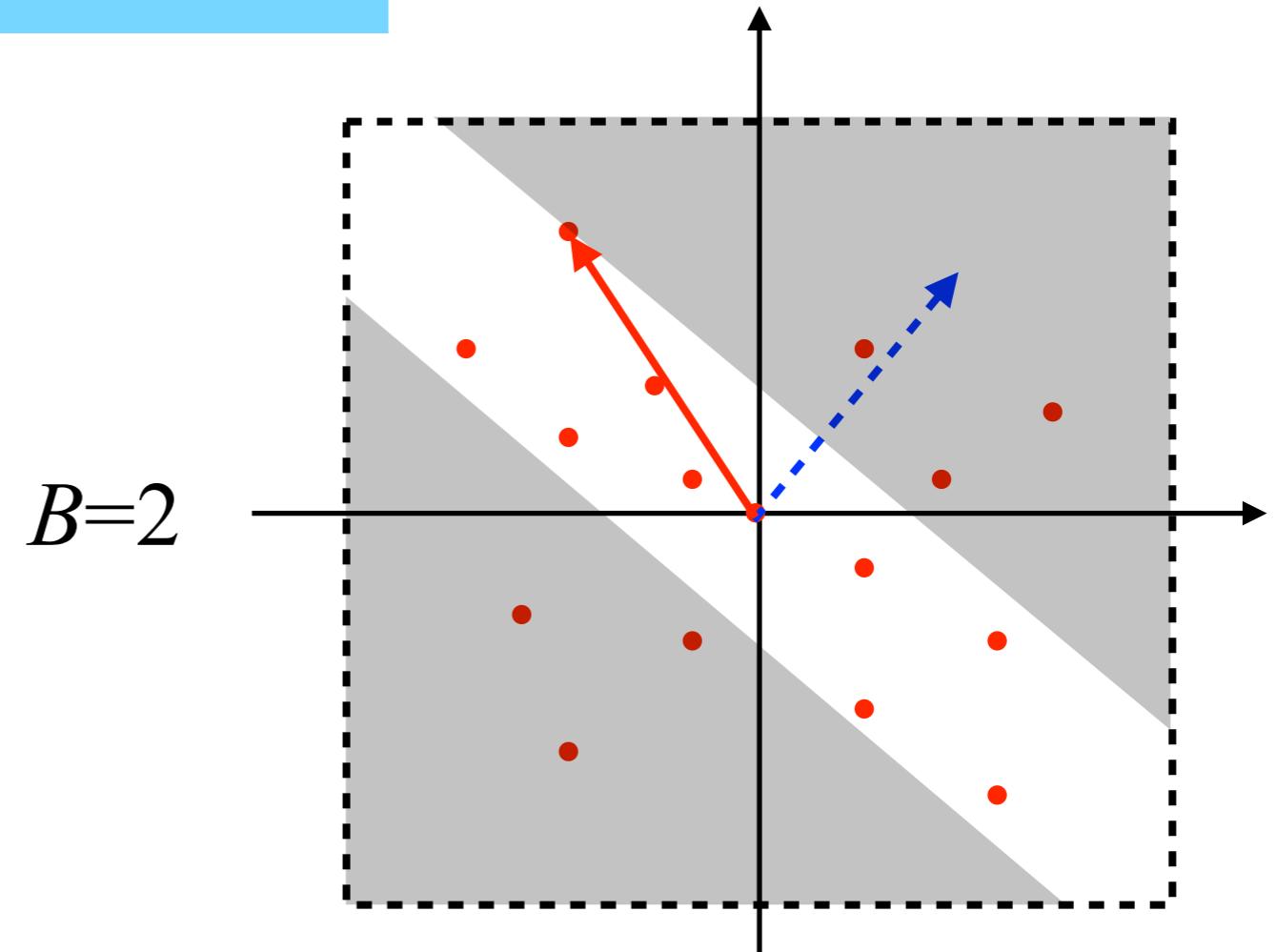
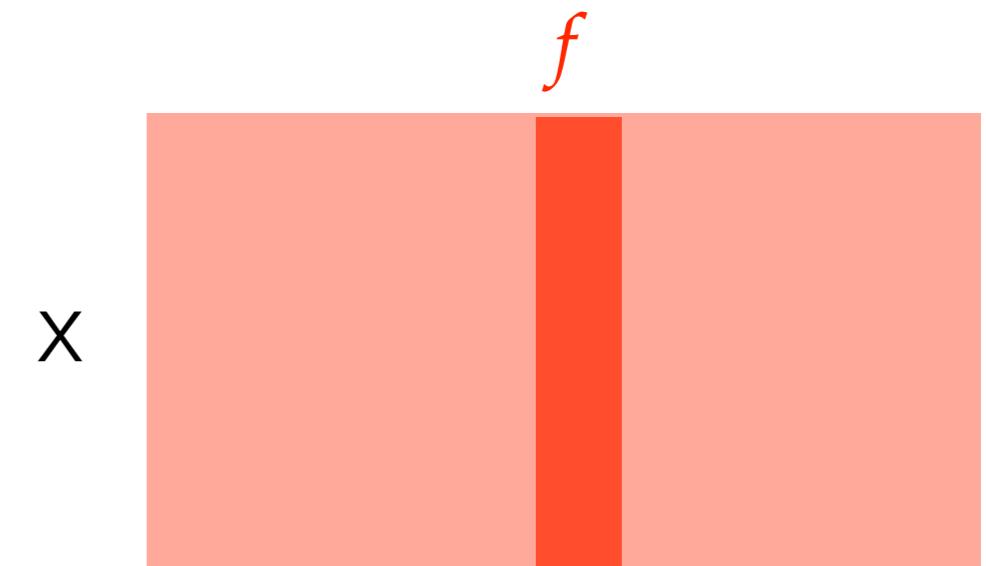
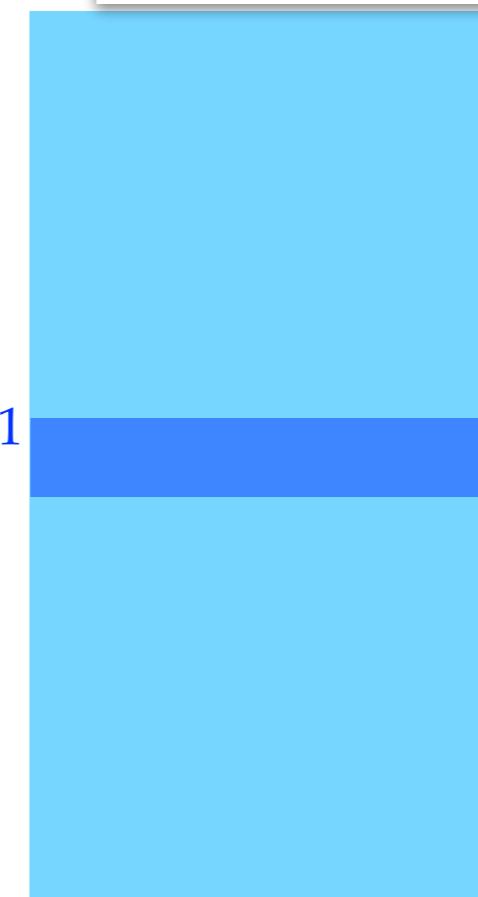
$$\pi_{t-1} \xrightarrow{f} \pi_{t-1} = \text{blue rectangle} \times \text{orange rectangle}$$



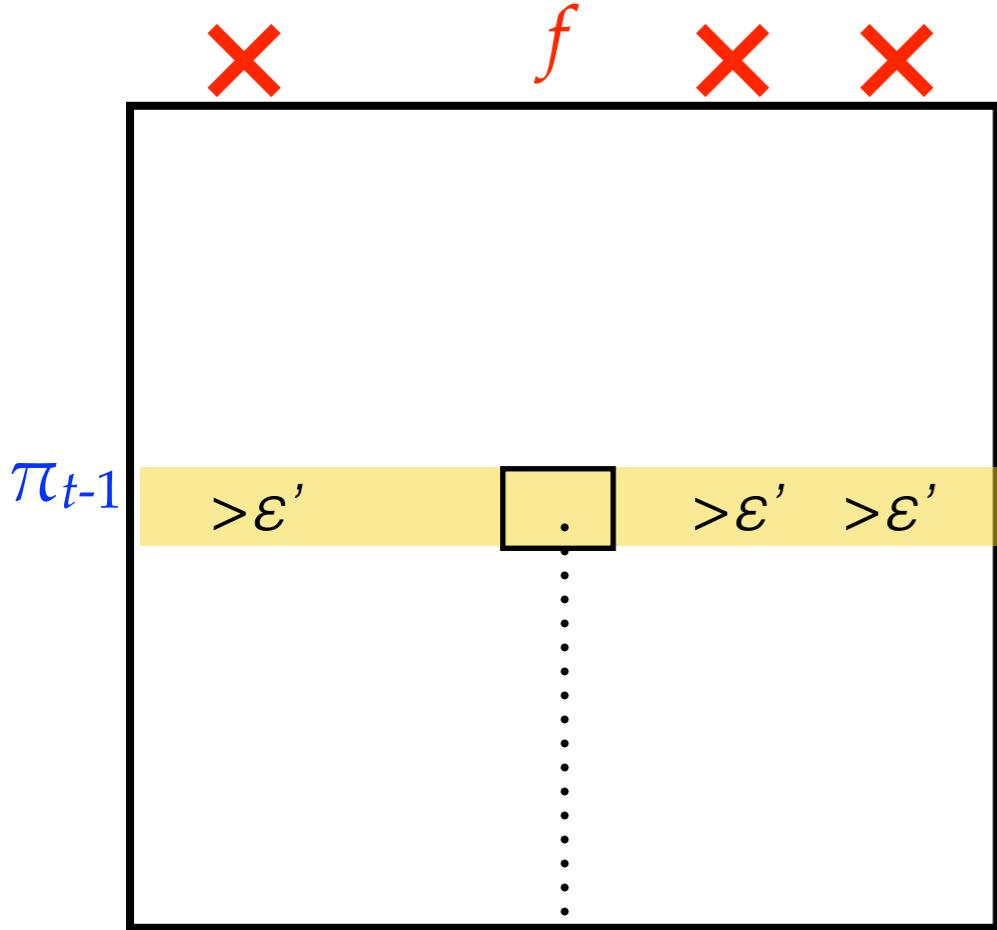
$B$  (Bellman rank)



$=$



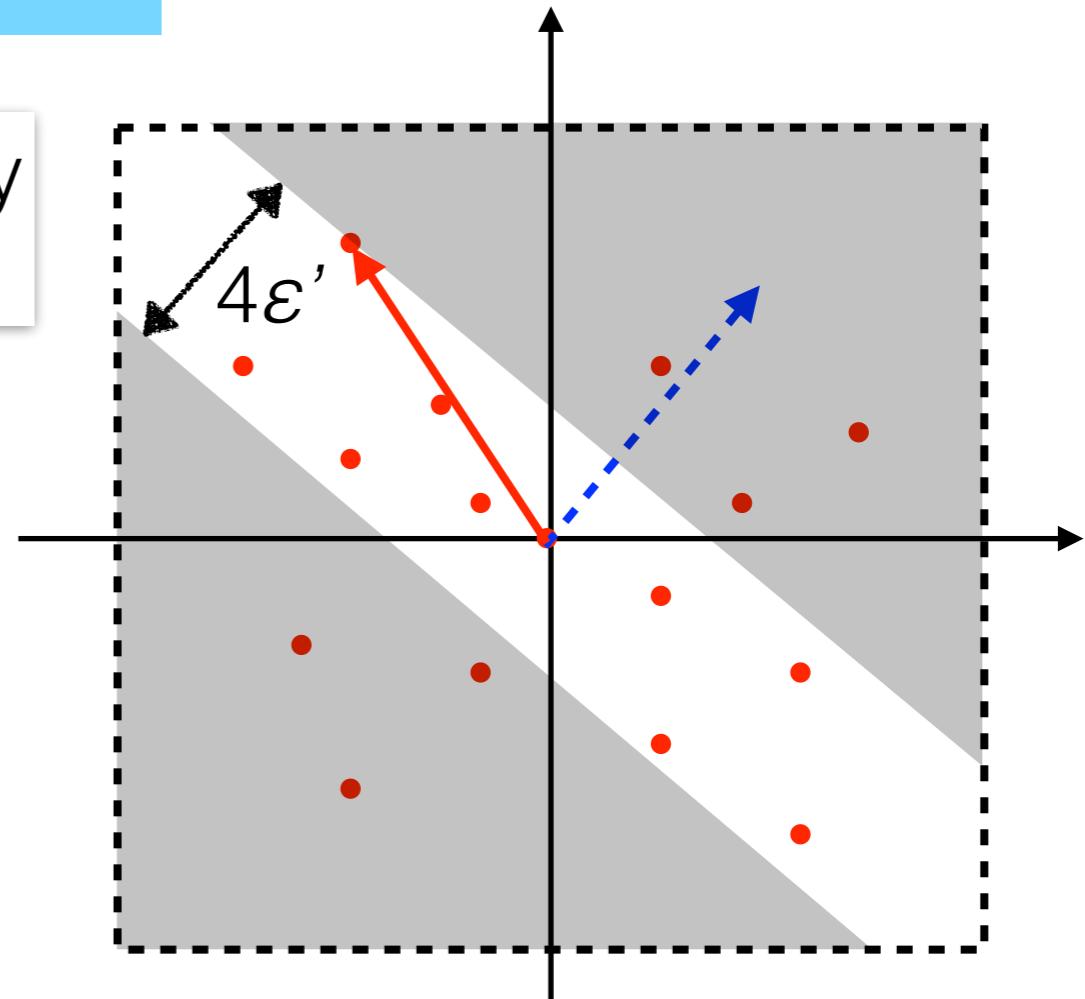
$B$  (Bellman rank)

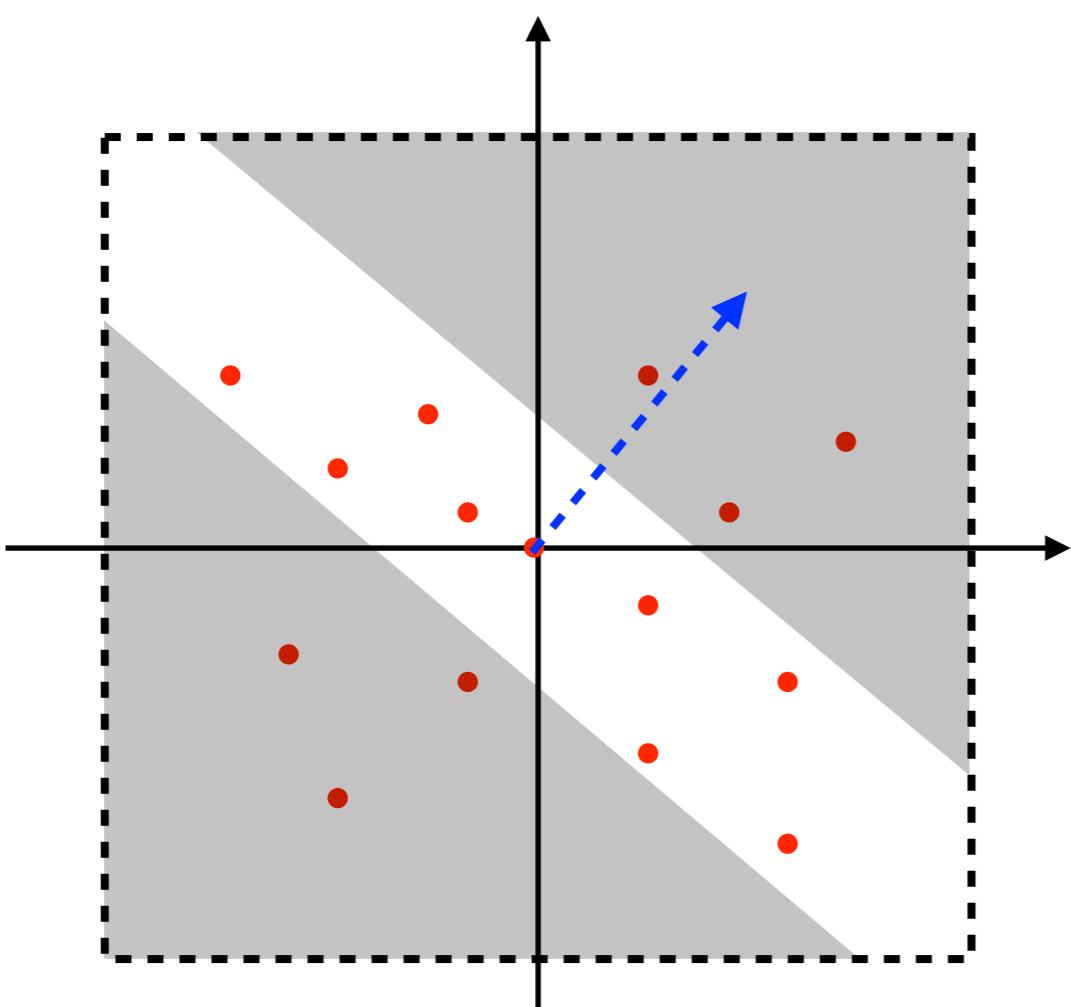


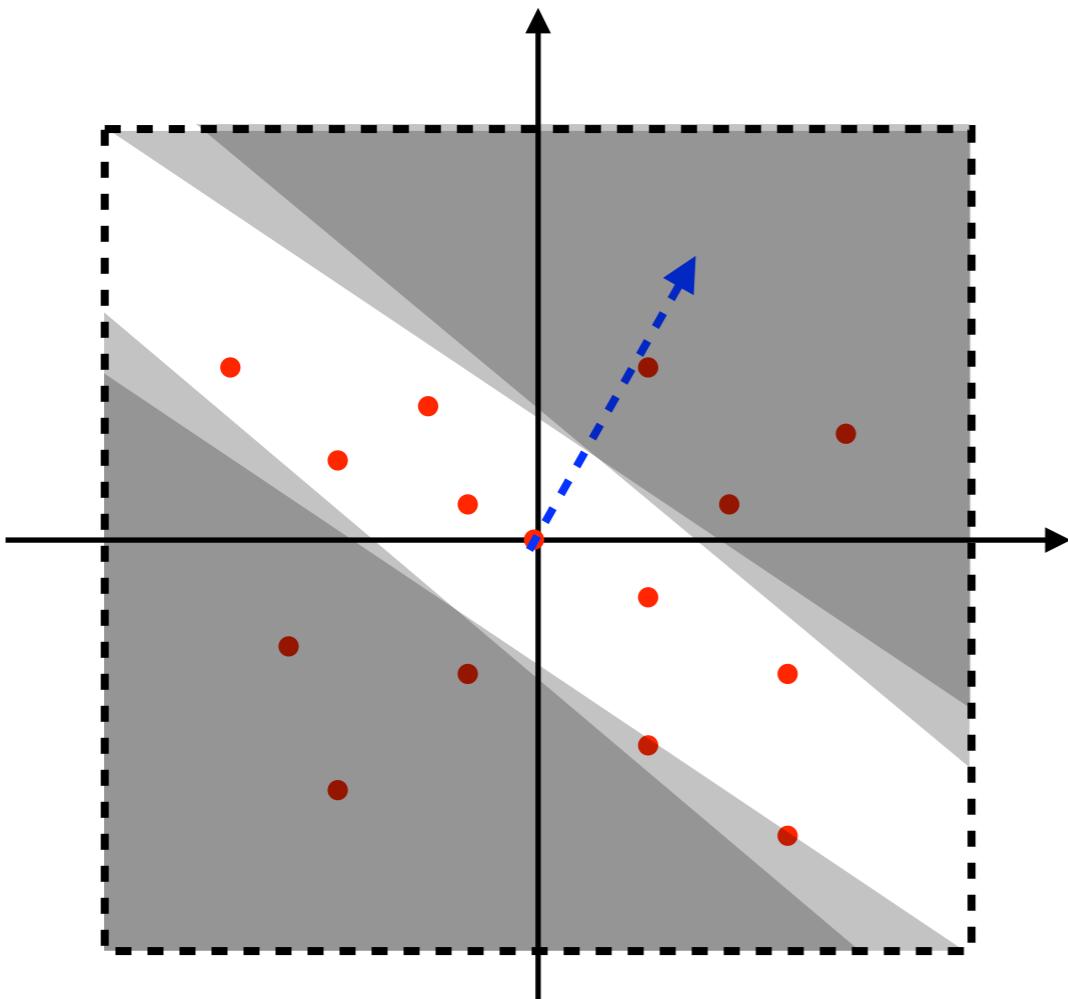
$\langle \rightarrow, \dashrightarrow \rangle$

$\varepsilon'$  controlled by  
sample size

$B=2$

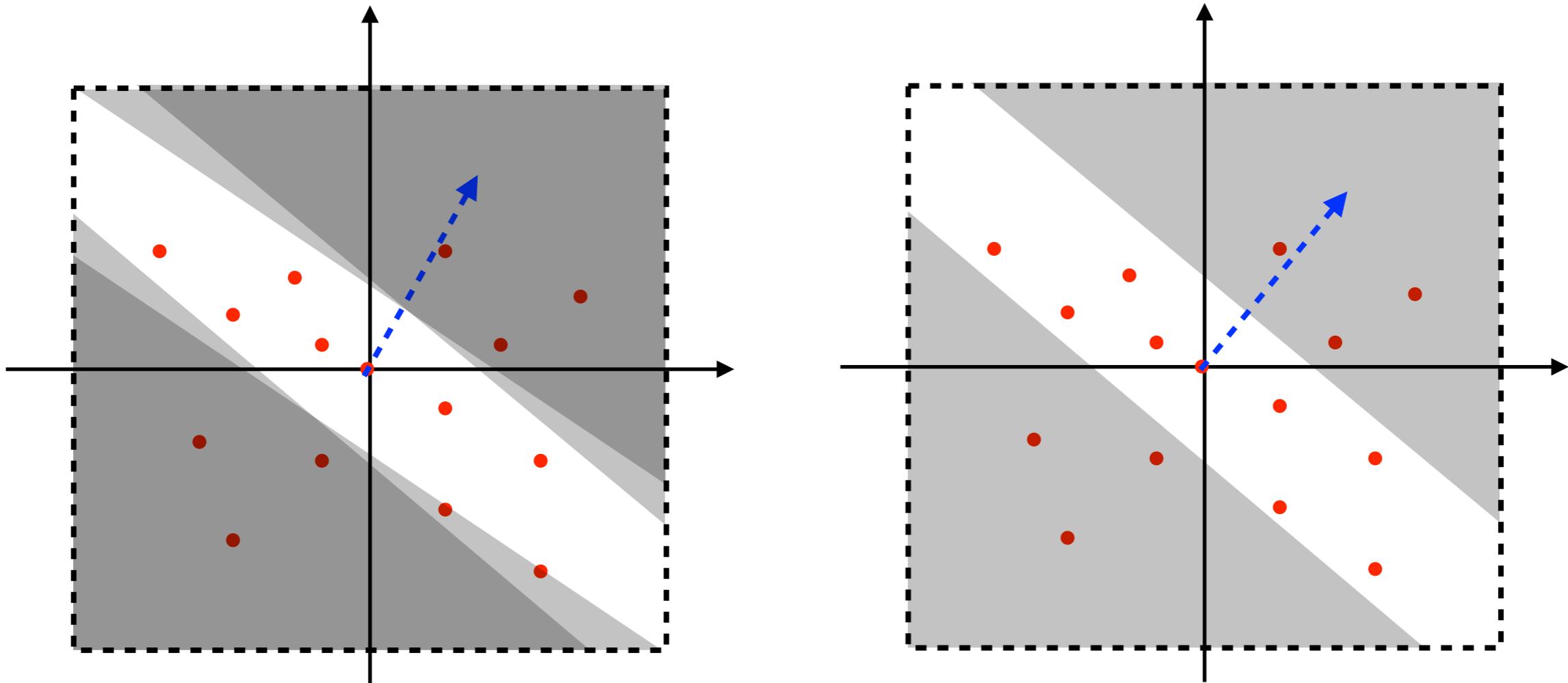






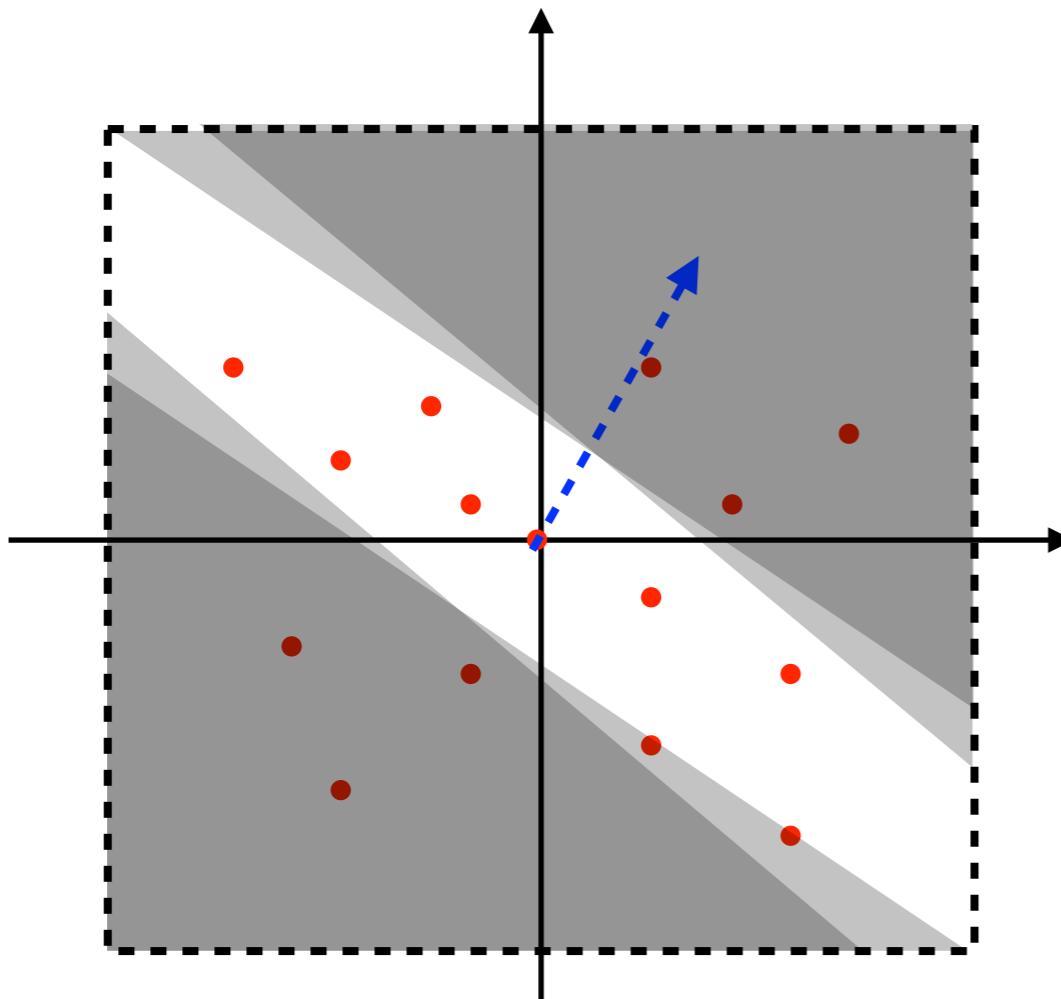
## inefficient exploration

- new distribution is similar to previous ones
- area of white space shrinks slowly



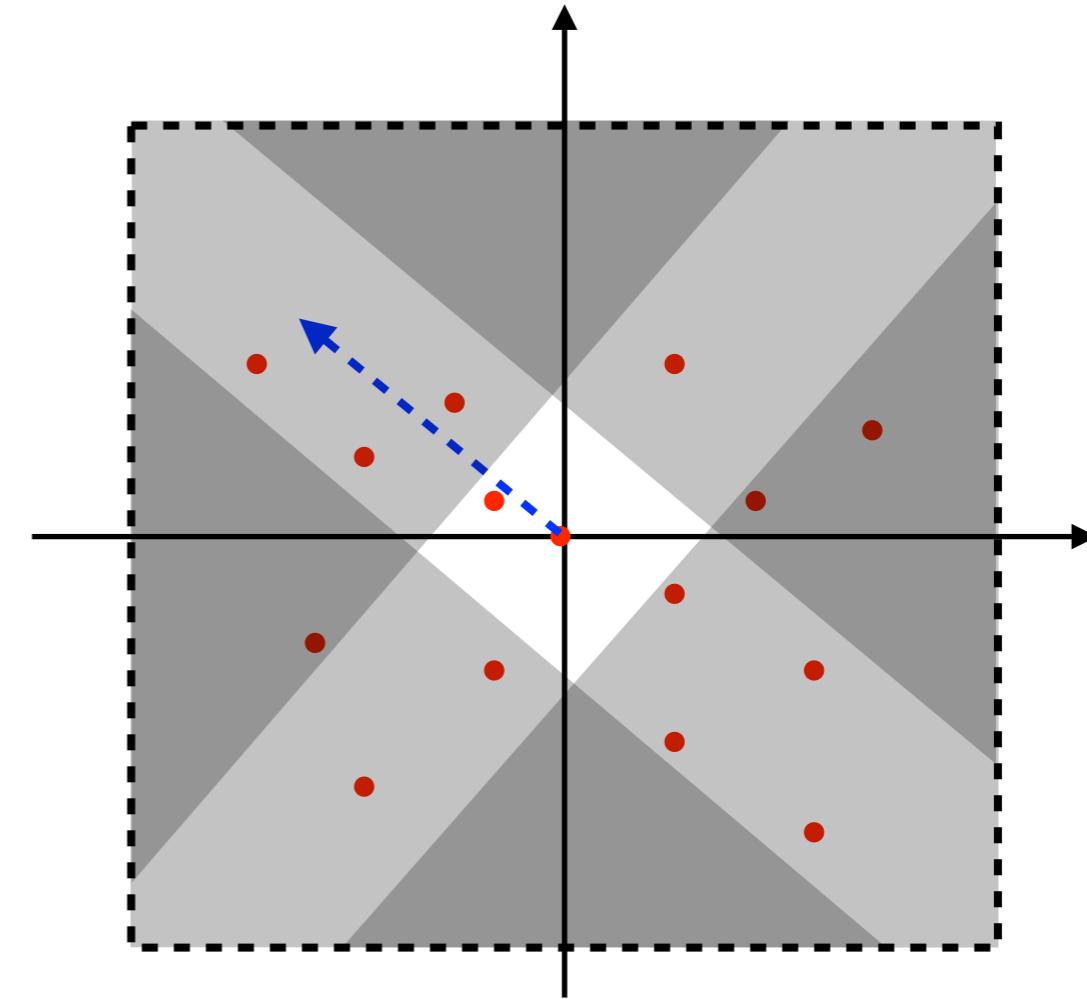
## inefficient exploration

- new distribution is similar to previous ones
- area of white space shrinks slowly



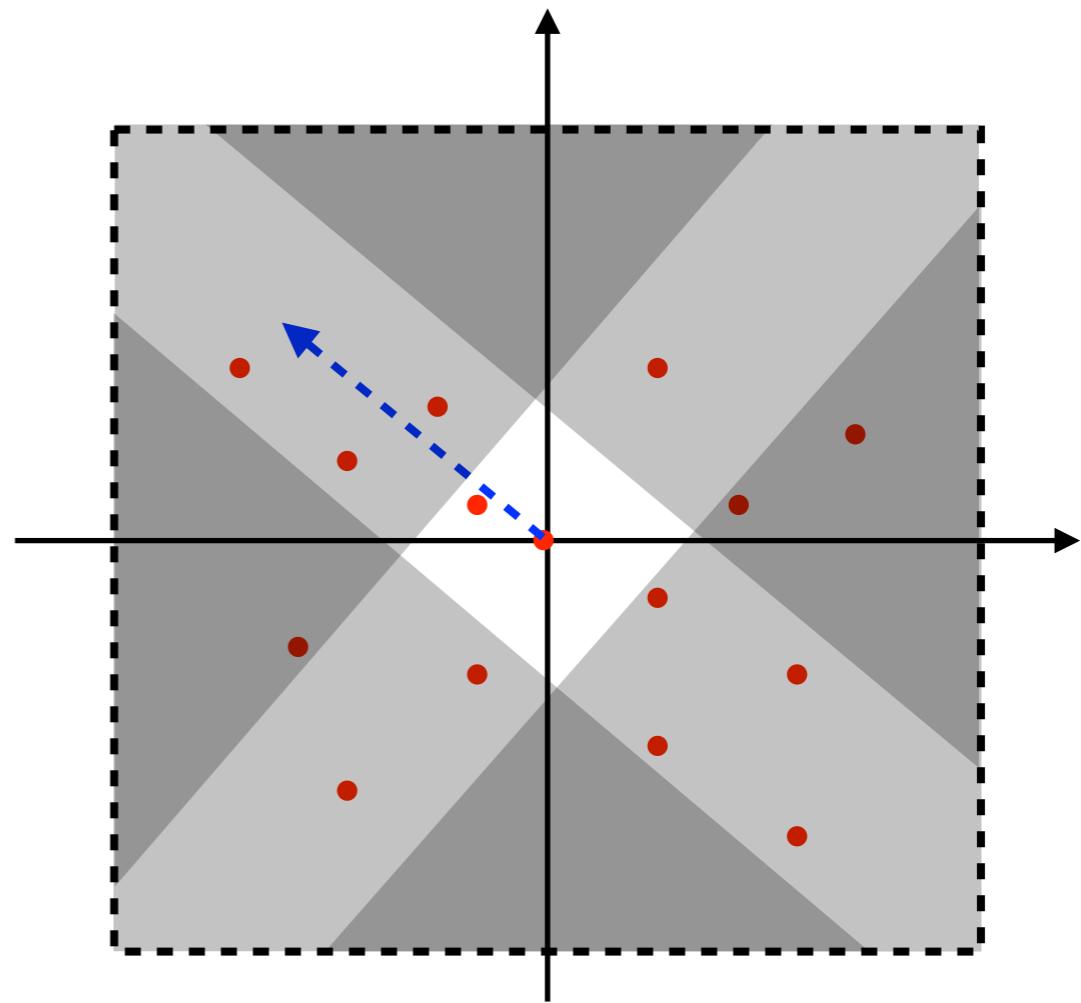
## inefficient exploration

- new distribution is similar to previous ones
- area of white space shrinks slowly



## efficient exploration

- new distribution is different from previous ones
- area of white space shrinks quickly

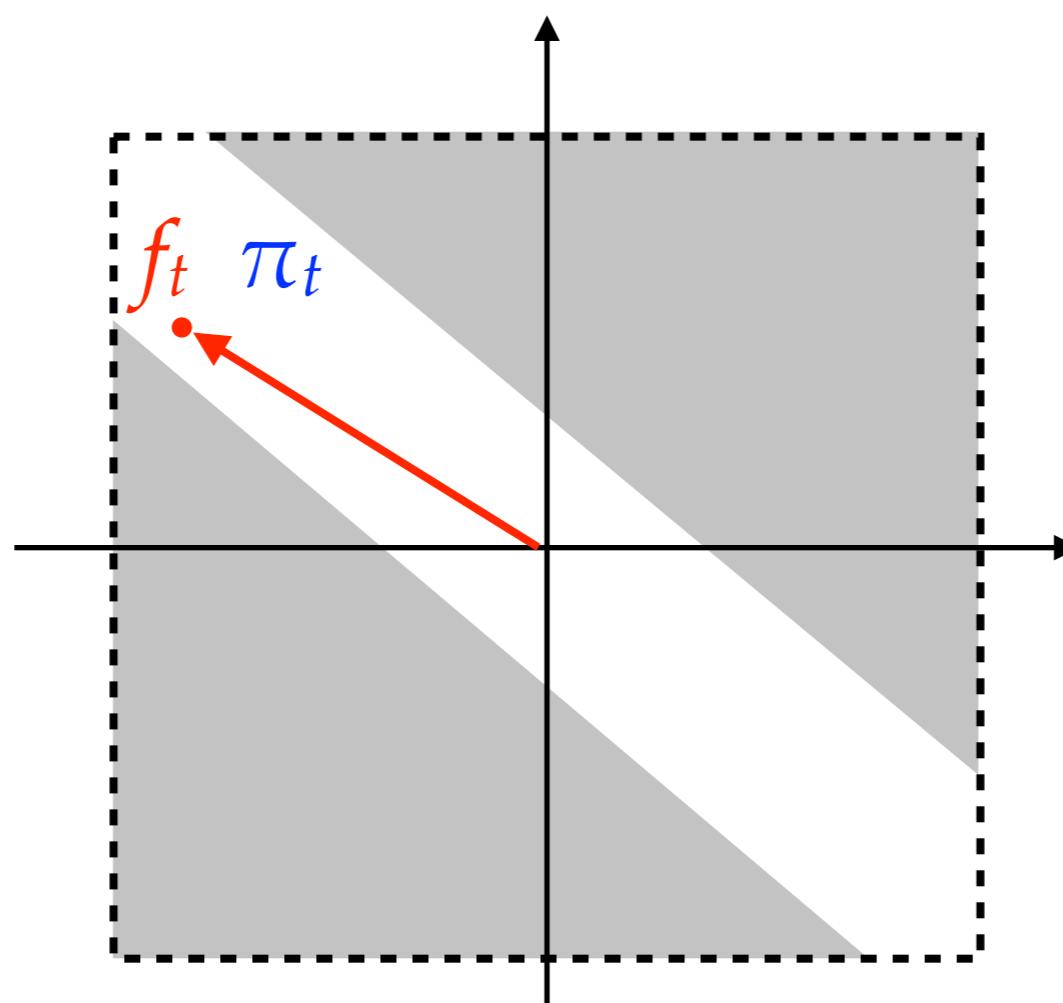


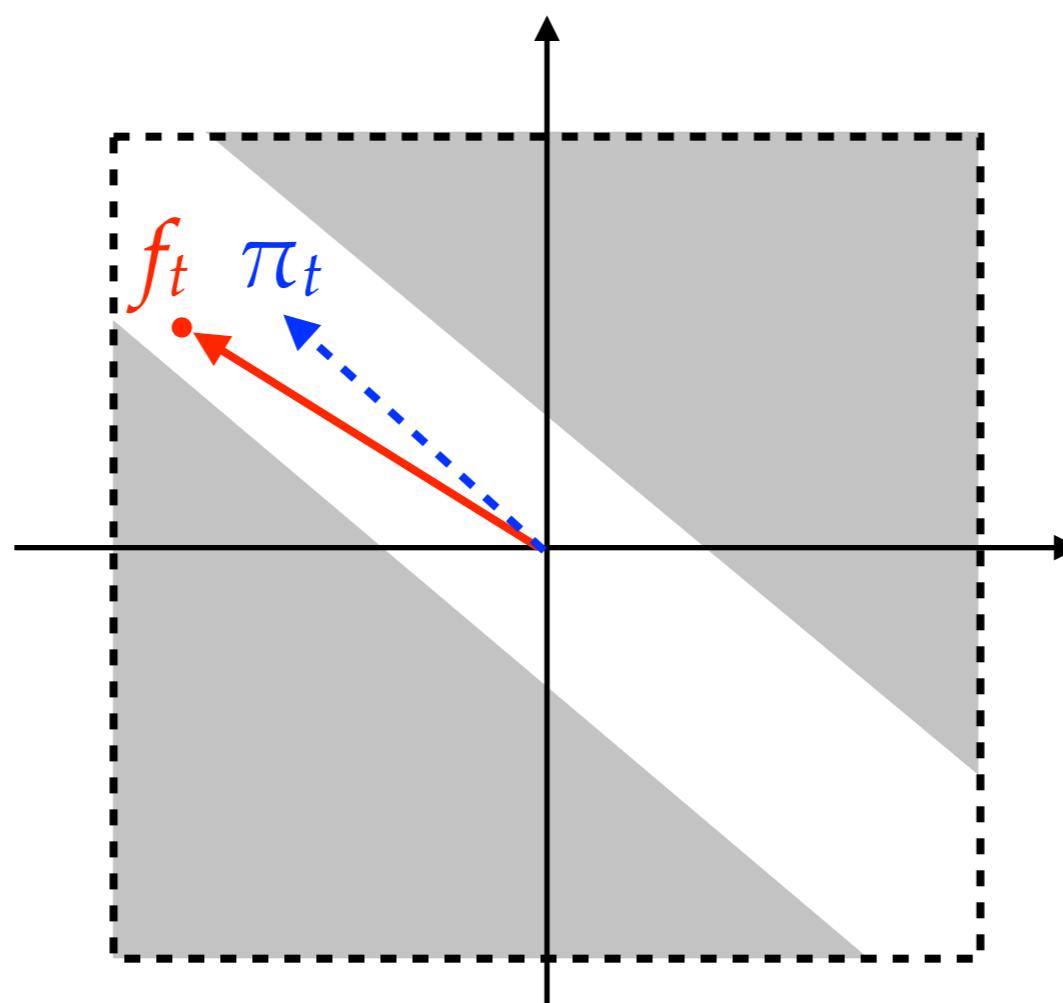
efficient exploration

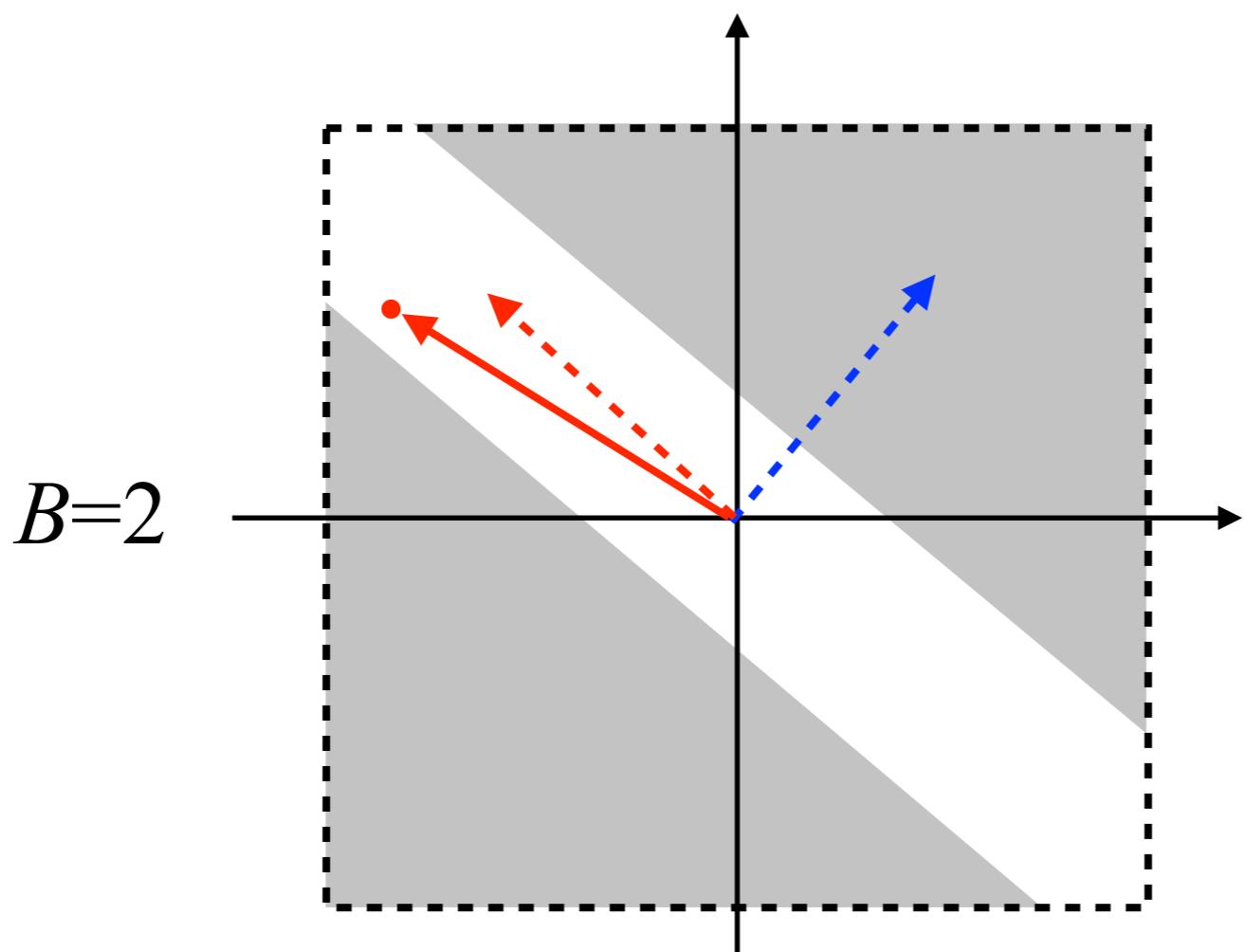
algorithm

analysis

- new distribution is different from previous ones
- area of white space shrinks quickly



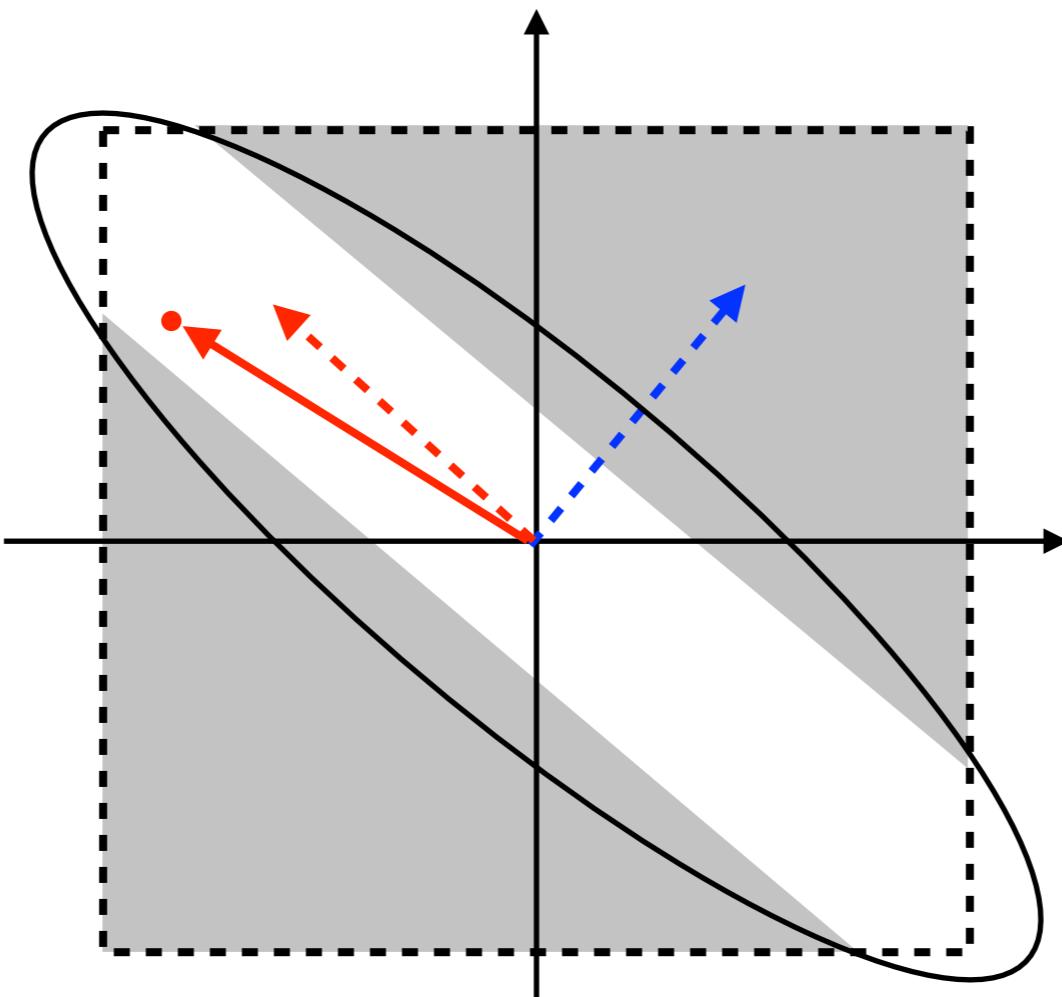




$B=2$

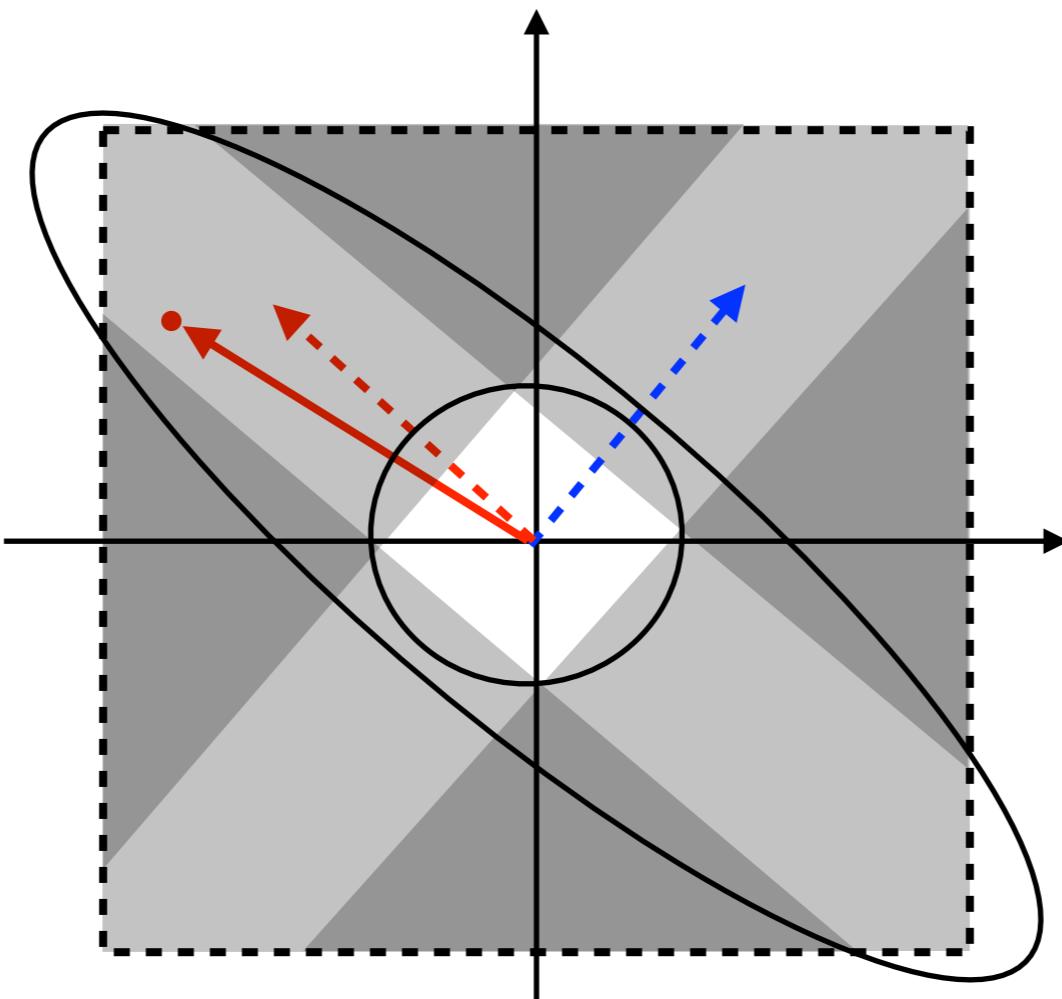
$B$

$B=2$

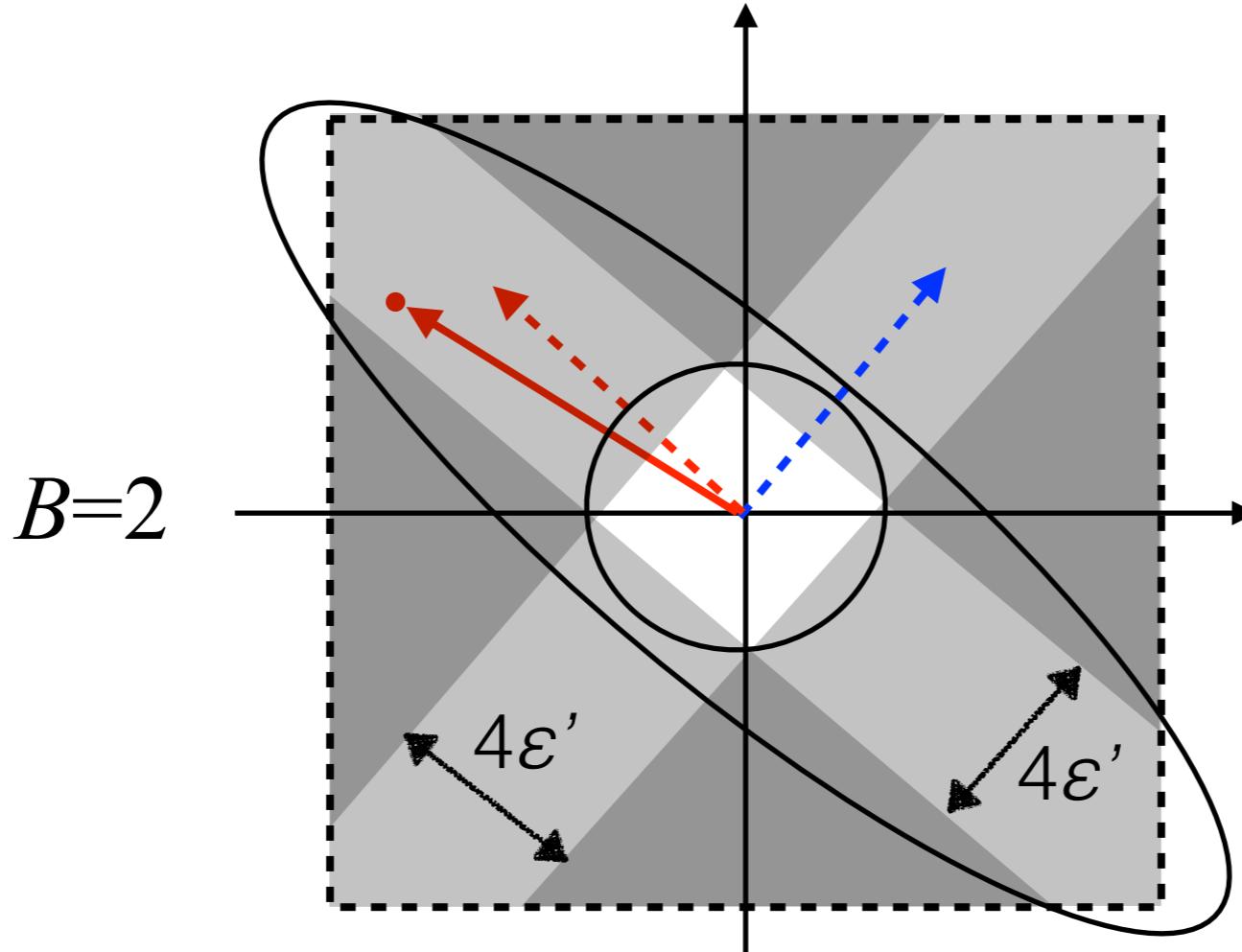


$B$

$B=2$

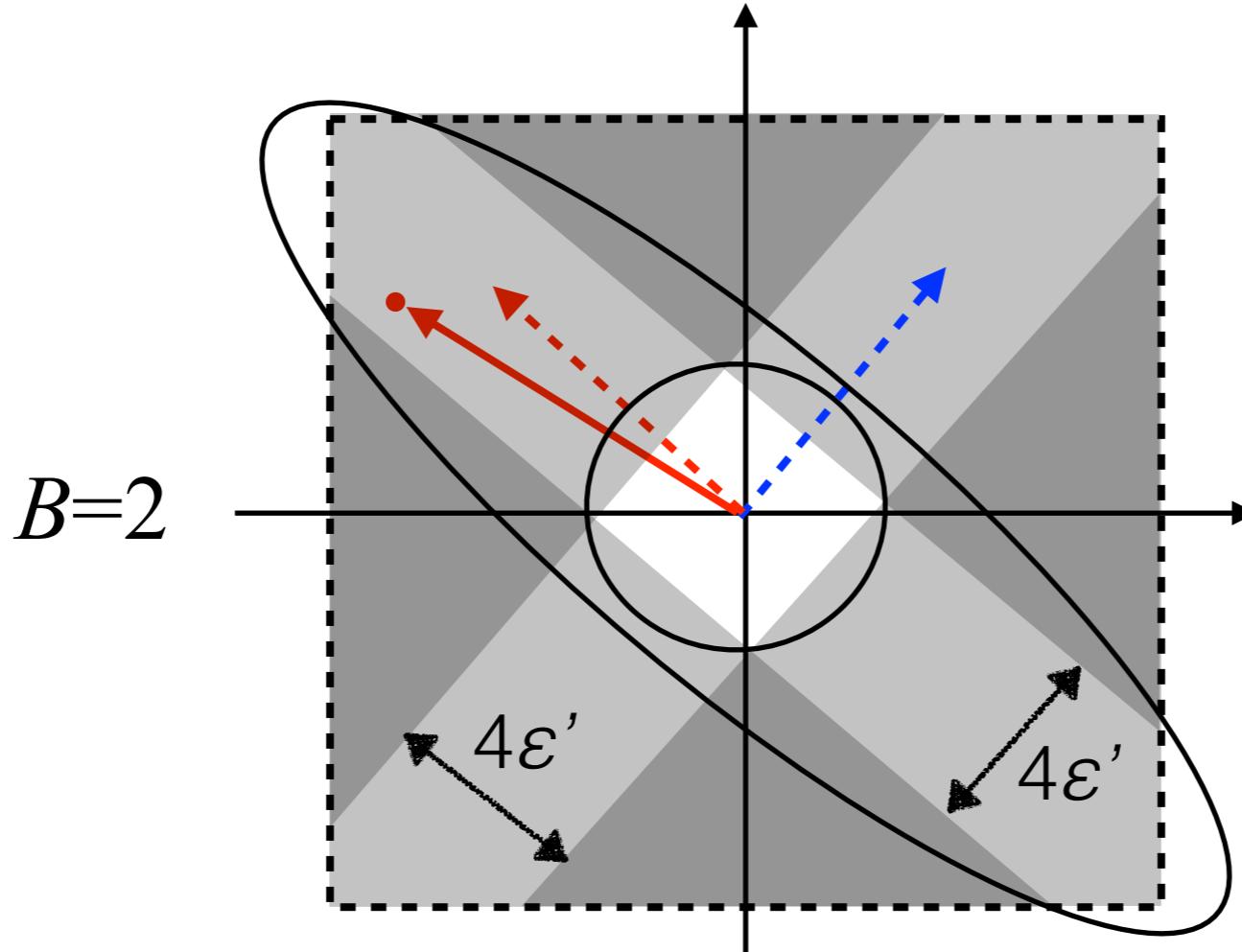


$B$



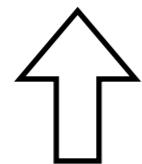
Adaptation of [Todd, 1982]:  
 Ellipsoid volume shrinks exponentially if

$$|\langle \rightarrow, \dashrightarrow \rangle| \geq 3\sqrt{B} \times 2\varepsilon'$$



Adaptation of [Todd, 1982]:  
 Ellipsoid volume shrinks exponentially if

$$|\langle \rightarrow, \dashrightarrow \rangle| \geq 3\sqrt{B} \times 2\epsilon'$$



controlled by sub-optimality



controlled by sample size