

# **Restless Bandits with Average Reward: Breaking the Uniform Global Attractor Assumption**

**Presenter: Yige Hong (CMU)**

**Joint work with Weina Wang (CMU), Qiaomin Xie (UW-Madison), and Yudong Chen (UW-Madison)**

# Setting: restless bandits



# Setting: restless bandits



$N = 3$  arms

# Setting: restless bandits

$S_1$

$S_2$

$S_3$

$N = 3$  arms

# Setting: restless bandits

$$S_1 \xrightarrow{A_1 \in \{0,1\}}$$

$$S_2 \xrightarrow{A_2 \in \{0,1\}}$$

$$S_3 \xrightarrow{A_3 \in \{0,1\}}$$

$N = 3$  arms

# Setting: restless bandits

$$\$ = r(S_1, A_1)$$

$$S_1 \xrightarrow{A_1 \in \{0,1\}}$$

$$\$ = r(S_2, A_2)$$

$$S_2 \xrightarrow{A_2 \in \{0,1\}}$$

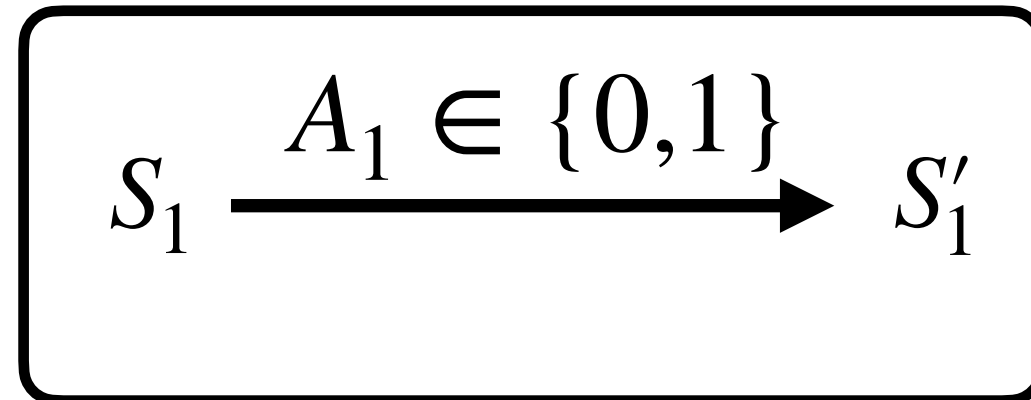
$$\$ = r(S_3, A_3)$$

$$S_3 \xrightarrow{A_3 \in \{0,1\}}$$

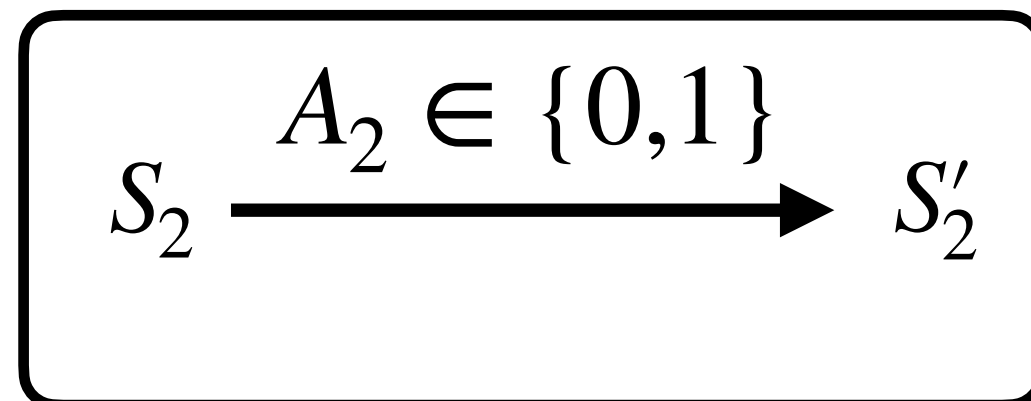
$N = 3$  arms

# Setting: restless bandits

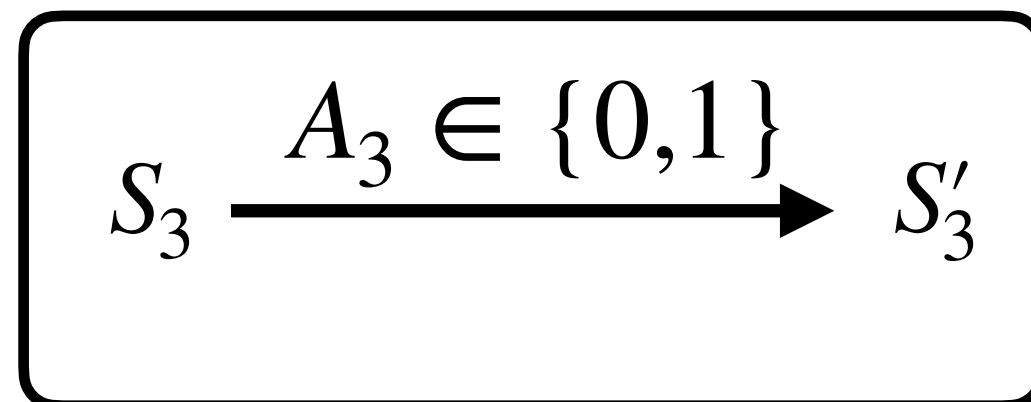
$$\$ = r(S_1, A_1)$$



$$\$ = r(S_2, A_2)$$



$$\$ = r(S_3, A_3)$$



$N = 3$  arms

# Setting: restless bandits

$$\$ = r(S_1, A_1)$$

$$S_1 \xrightarrow[A_1 \in \{0,1\}]{P(\cdot | S_1, A_1)} S'_1$$

$$\$ = r(S_2, A_2)$$

$$S_2 \xrightarrow[A_2 \in \{0,1\}]{P(\cdot | S_2, A_2)} S'_2$$

$$\$ = r(S_3, A_3)$$

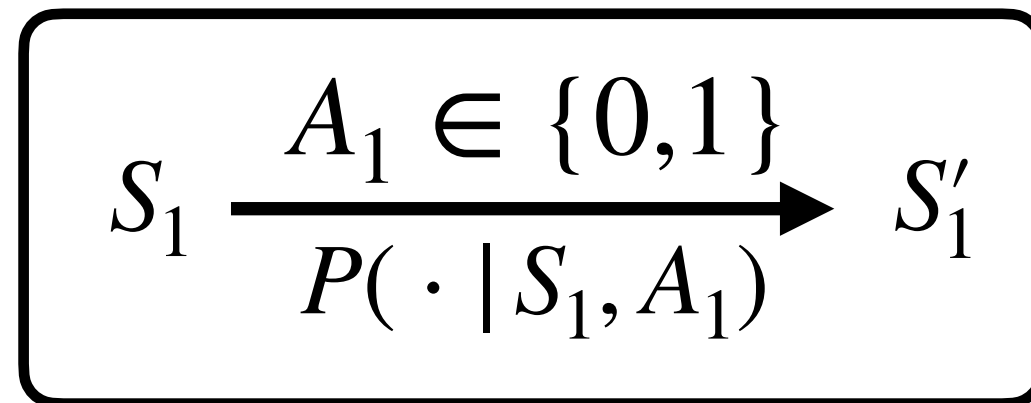
$$S_3 \xrightarrow[A_3 \in \{0,1\}]{P(\cdot | S_3, A_3)} S'_3$$

$N = 3$  arms

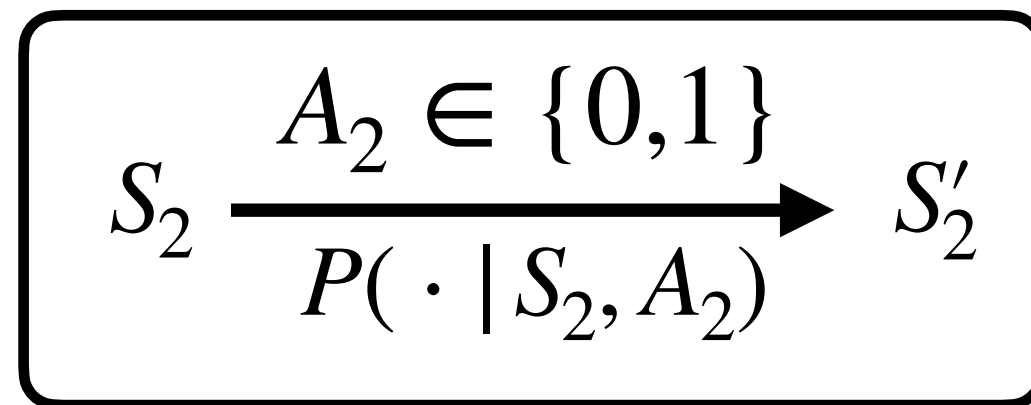


# Setting: restless bandits

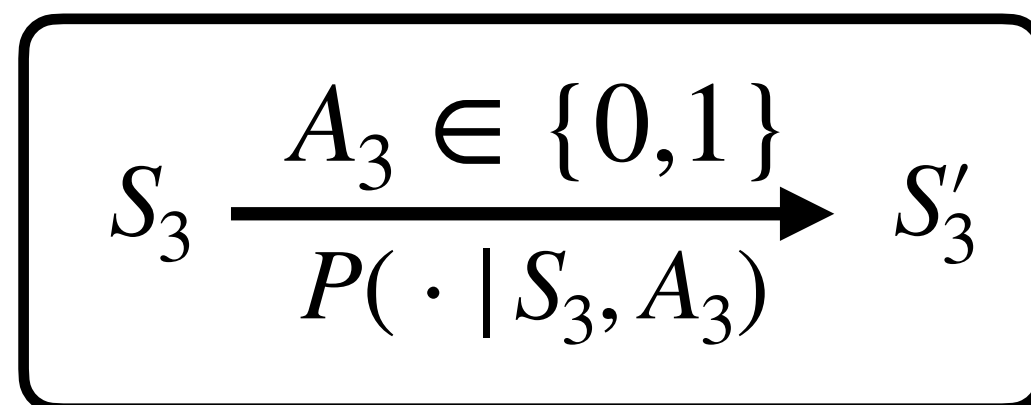
$$\$ = r(S_1, A_1)$$



$$\$ = r(S_2, A_2)$$



$$\$ = r(S_3, A_3)$$

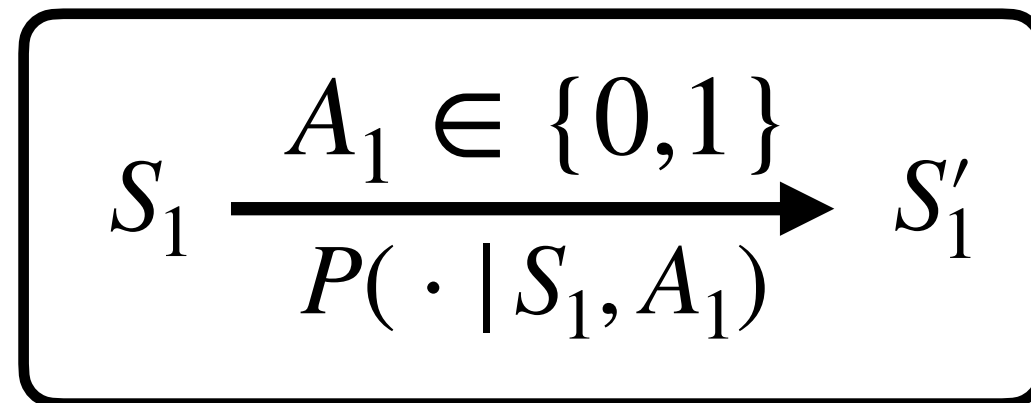


$N = 3$  arms

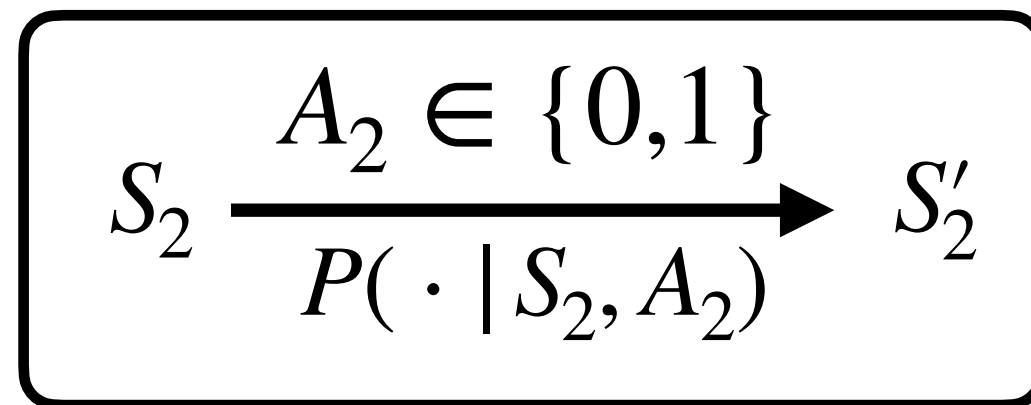
$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

# Setting: restless bandits

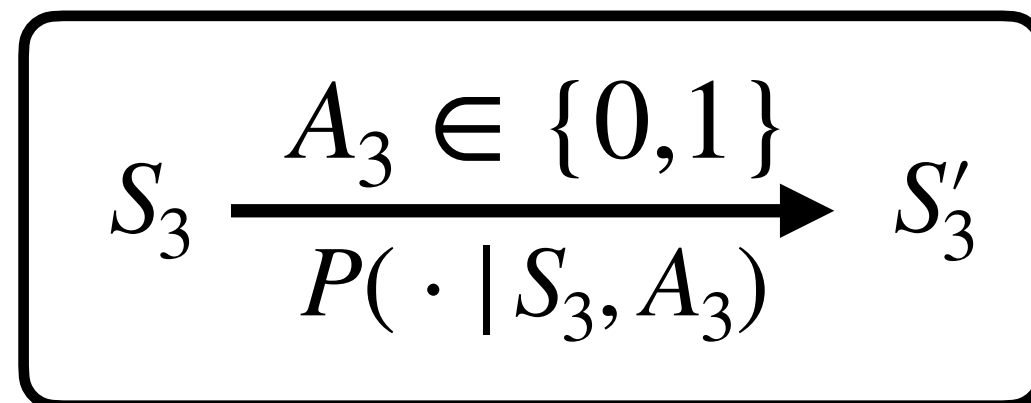
$$\$ = r(S_1, A_1)$$



$$\$ = r(S_2, A_2)$$



$$\$ = r(S_3, A_3)$$



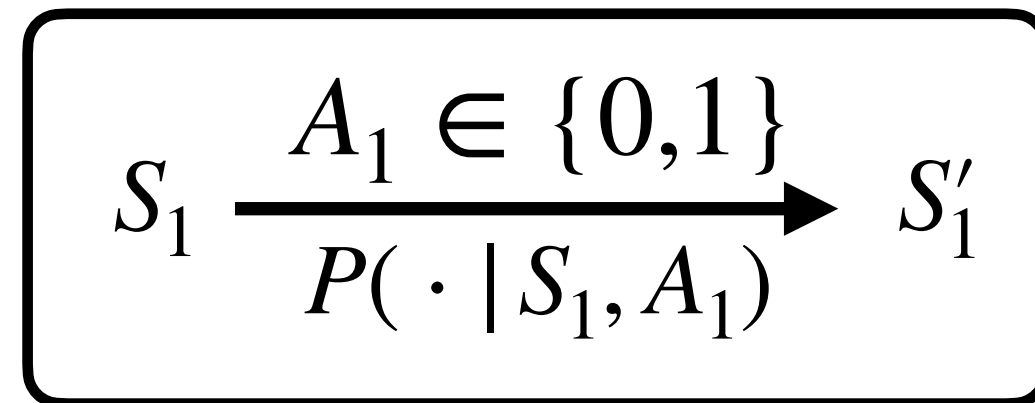
$N = 3$  arms

$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

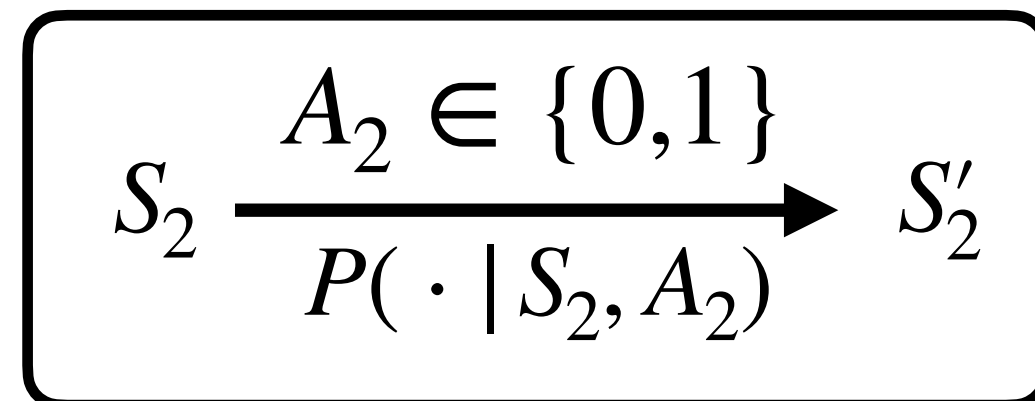
$$\pi: (S_1, S_2, \dots, S_N) \mapsto (A_1, A_2, \dots, A_N)$$

# Setting: restless bandits

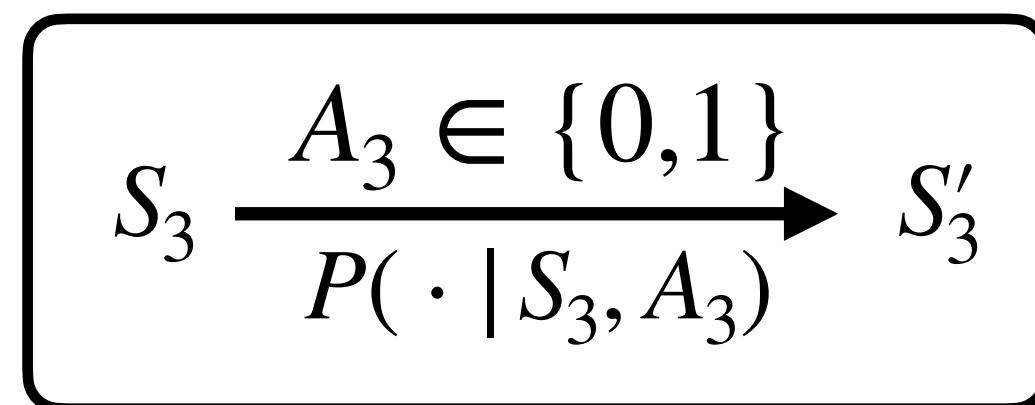
$$\text{\$} = r(S_1, A_1)$$



$$\text{\$} = r(S_2, A_2)$$



$$\text{\$} = r(S_3, A_3)$$



$N = 3$  arms

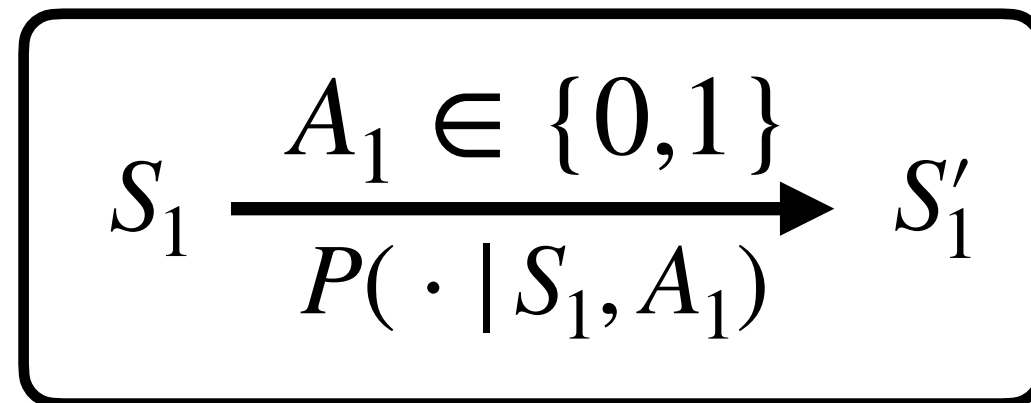
$\max_{\pi} V_N^{\pi} \triangleq$  long run average reward under policy  $\pi$

$$\text{s.t. } \sum_{i=1}^N A_i = \alpha N, \text{ any time slot}$$

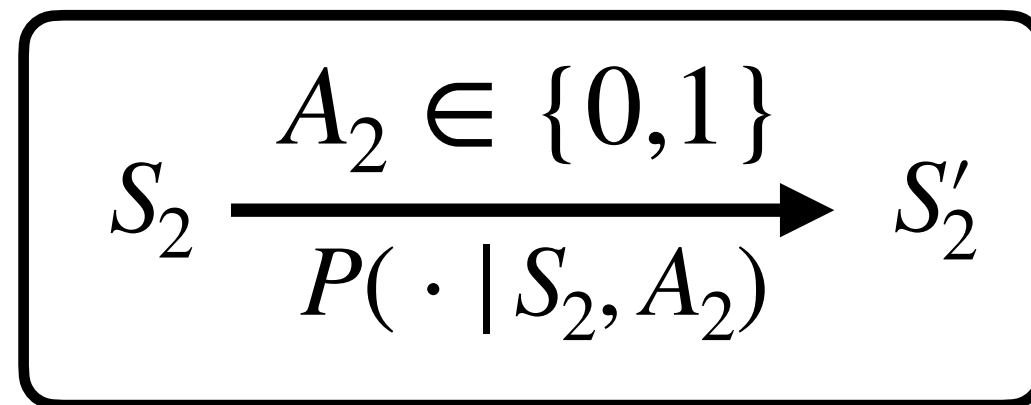
$$\pi: (S_1, S_2, \dots, S_N) \mapsto (A_1, A_2, \dots, A_N)$$

# Setting: restless bandits

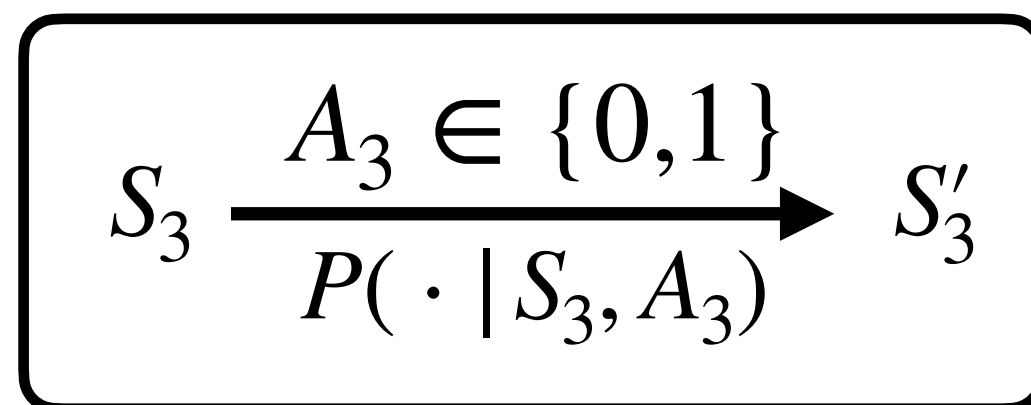
$$\$ = r(S_1, A_1)$$



$$\$ = r(S_2, A_2)$$



$$\$ = r(S_3, A_3)$$



$N = 3$  arms

$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

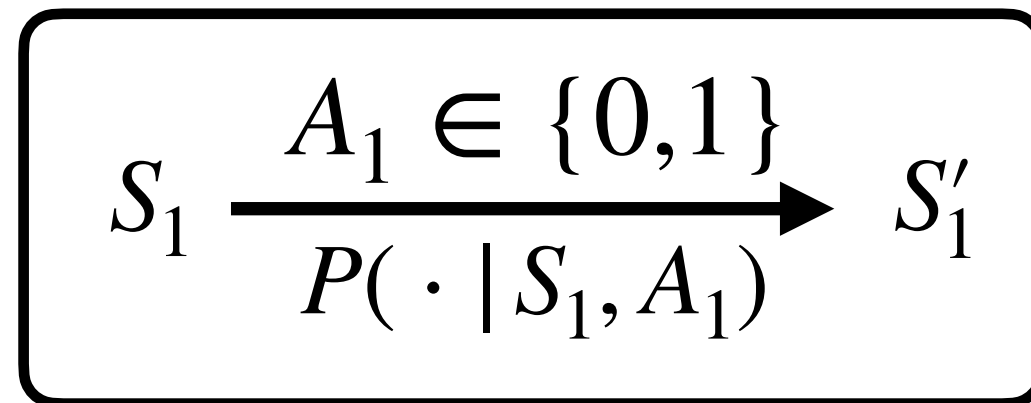
$$\text{s.t. } \sum_{i=1}^N A_i = \alpha N, \text{ any time slot}$$

$$0 < \alpha < 1$$

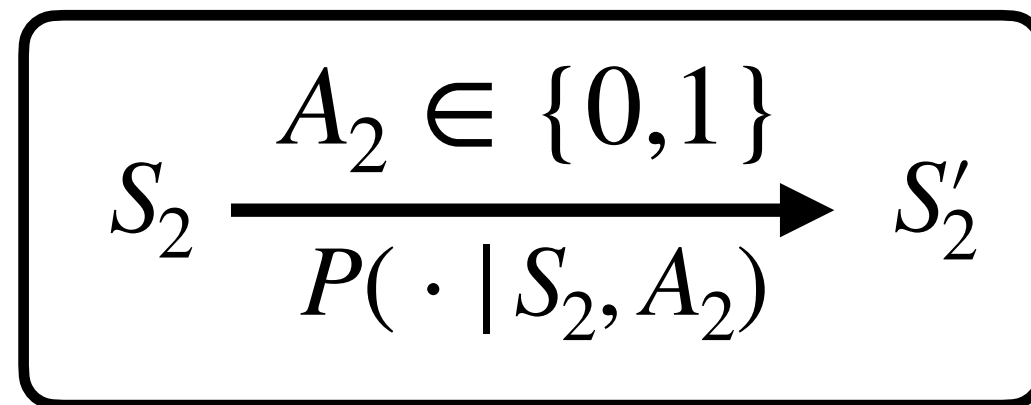
$$\pi: (S_1, S_2, \dots, S_N) \mapsto (A_1, A_2, \dots, A_N)$$

# Setting: restless bandits

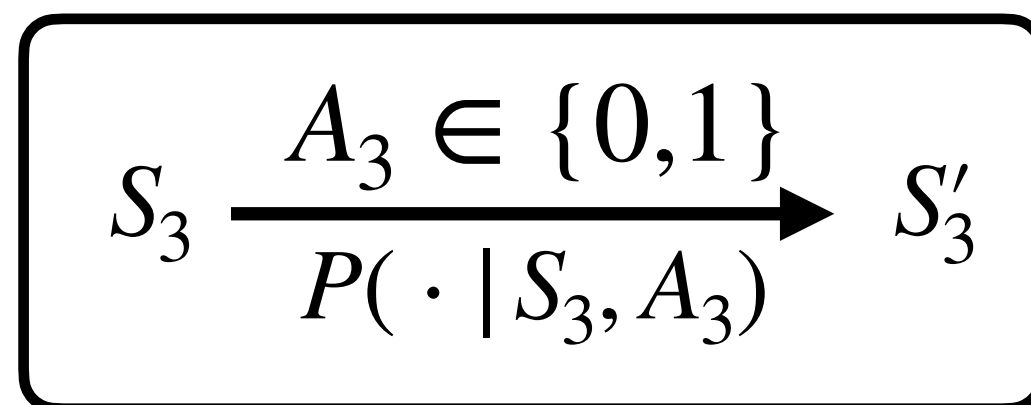
$$\$ = r(S_1, A_1)$$



$$\$ = r(S_2, A_2)$$



$$\$ = r(S_3, A_3)$$



Full information

$N = 3$  arms

$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

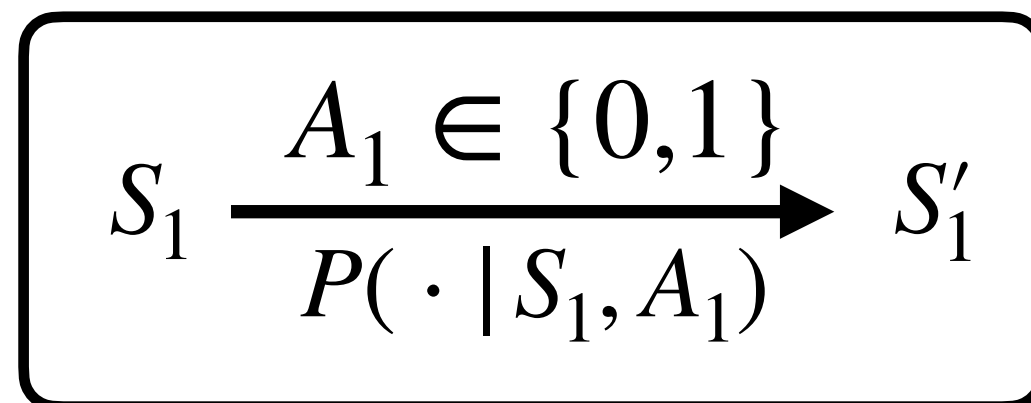
$$\text{s.t. } \sum_{i=1}^N A_i = \alpha N, \text{ any time slot}$$

$$0 < \alpha < 1$$

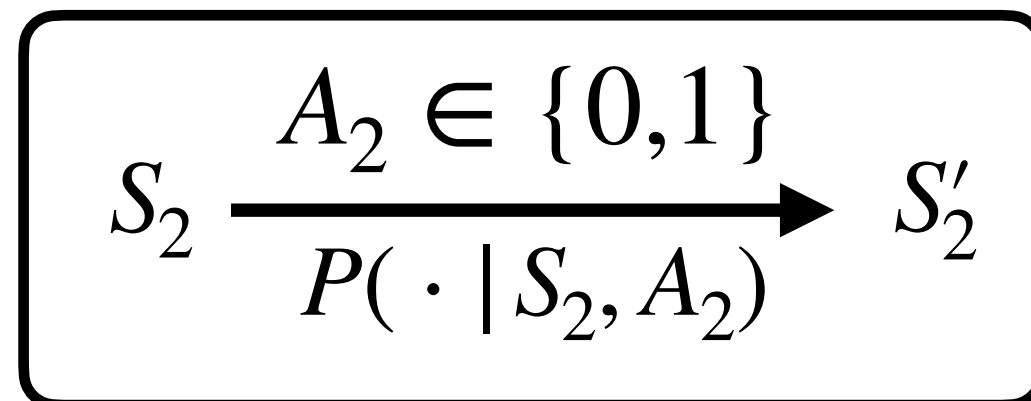
$$\pi: (S_1, S_2, \dots, S_N) \mapsto (A_1, A_2, \dots, A_N)$$

# Setting: restless bandits

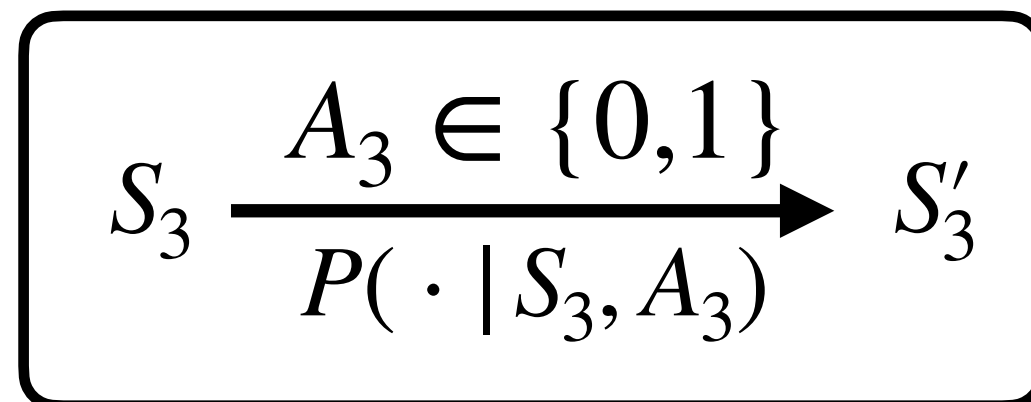
$$\$ = r(S_1, A_1)$$



$$\$ = r(S_2, A_2)$$



$$\$ = r(S_3, A_3)$$



$N = 3$  arms

$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

$$\text{s.t. } \sum_{i=1}^N A_i = \alpha N, \text{ any time slot}$$

$$0 < \alpha < 1$$

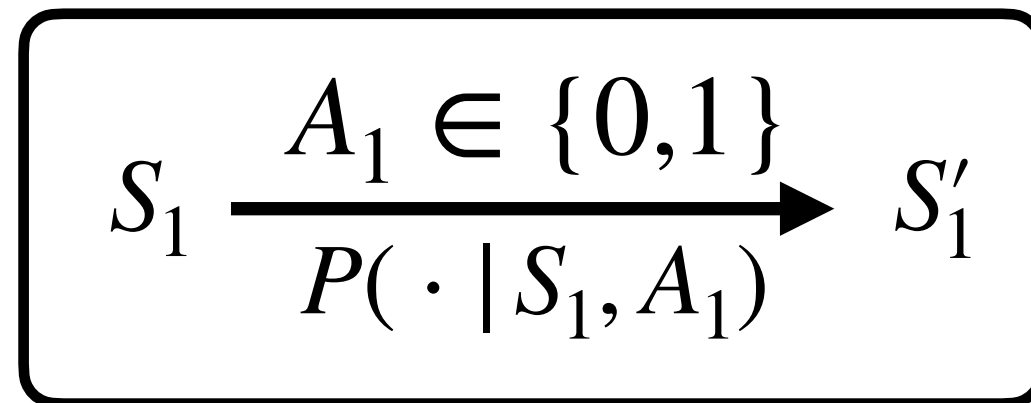
$$\pi: (S_1, S_2, \dots, S_N) \mapsto (A_1, A_2, \dots, A_N)$$

Full information

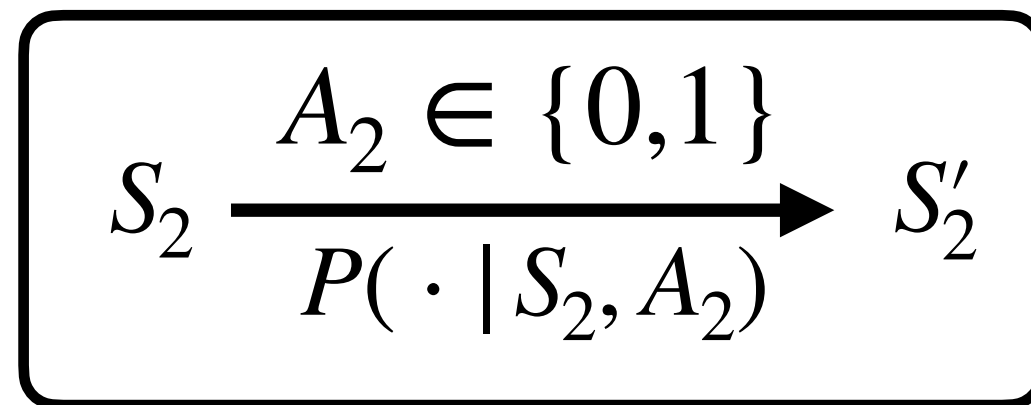
If  $N$  large, high dimensional;  
PSPACE hard to solve

# Goal: asymptotic optimality for large N

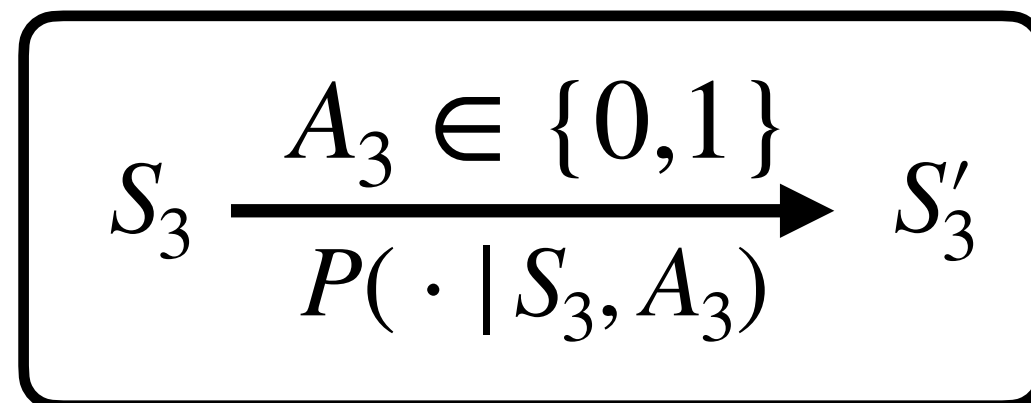
$$\$ = r(S_1, A_1)$$



$$\$ = r(S_2, A_2)$$



$$\$ = r(S_3, A_3)$$



$N = 3$  arms

$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

$$\text{s.t. } \sum_{i=1}^N A_i = \alpha N, \text{ any time slot}$$

$$0 < \alpha < 1$$

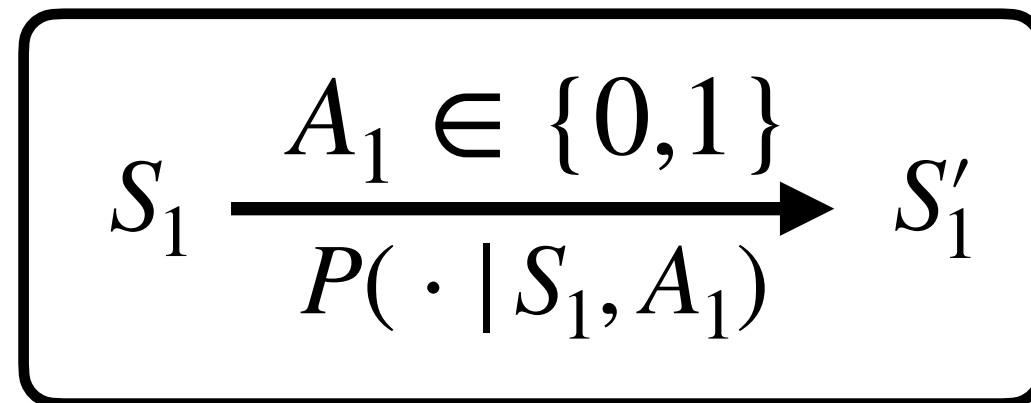
$$\pi: (S_1, S_2, \dots, S_N) \mapsto (A_1, A_2, \dots, A_N)$$

Full information

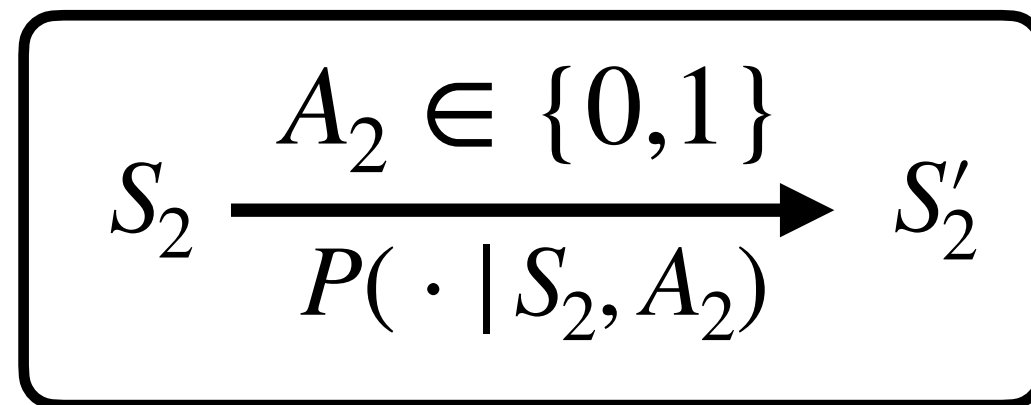
If  $N$  large, high dimensional;  
PSPACE hard to solve

# Goal: asymptotic optimality for large N

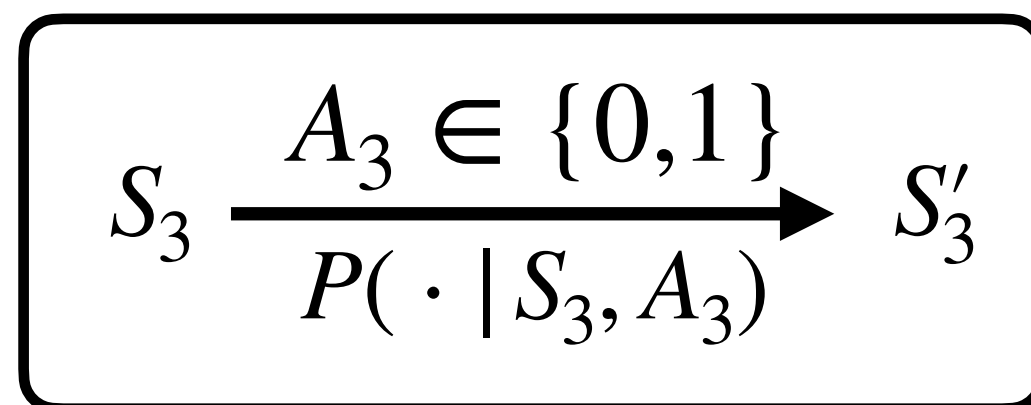
$$\$ = r(S_1, A_1)$$



$$\$ = r(S_2, A_2)$$



$$\$ = r(S_3, A_3)$$



$N = 3$  arms

Full information

$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

$$\pi: (S_1, S_2, \dots, S_N) \mapsto (A_1, A_2, \dots, A_N)$$

$$\text{s.t. } \sum_{i=1}^N A_i = \alpha N, \text{ any time slot}$$

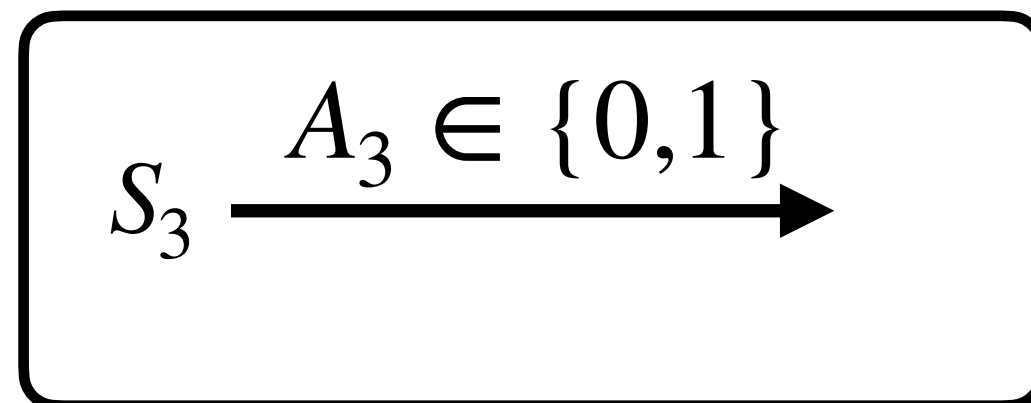
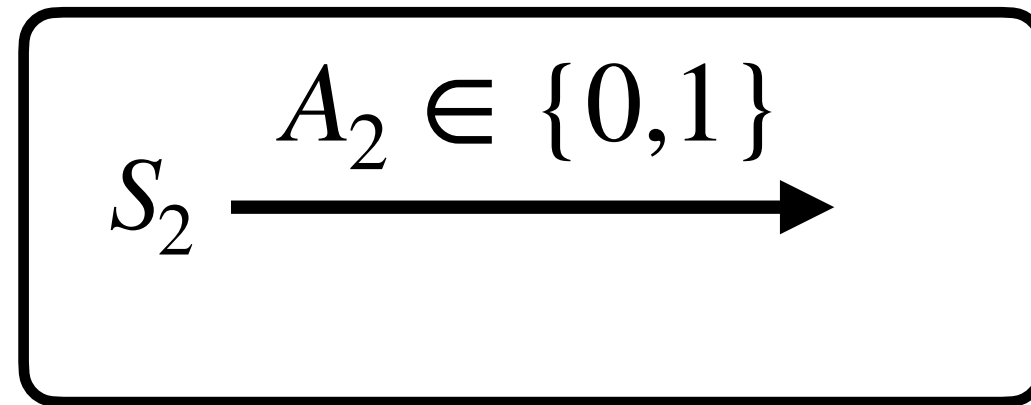
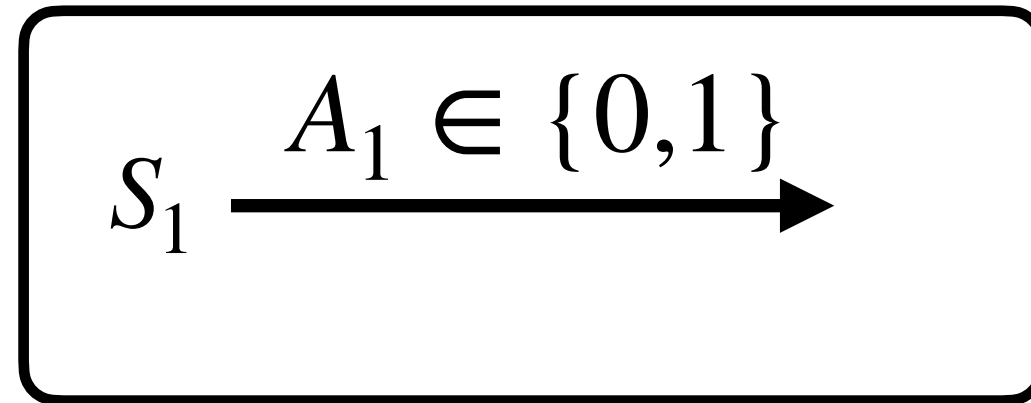
$$0 < \alpha < 1$$

If  $N$  large, high dimensional;  
PSPACE hard to solve

Goal: computationally efficient algorithm for finding  $\pi$   
such that  $V_N^* - V_N^{\pi} \rightarrow 0$  as  $N \rightarrow \infty$



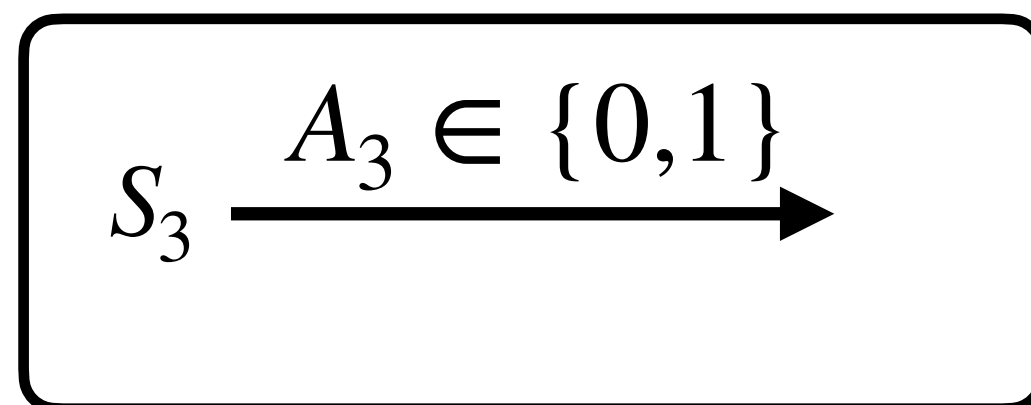
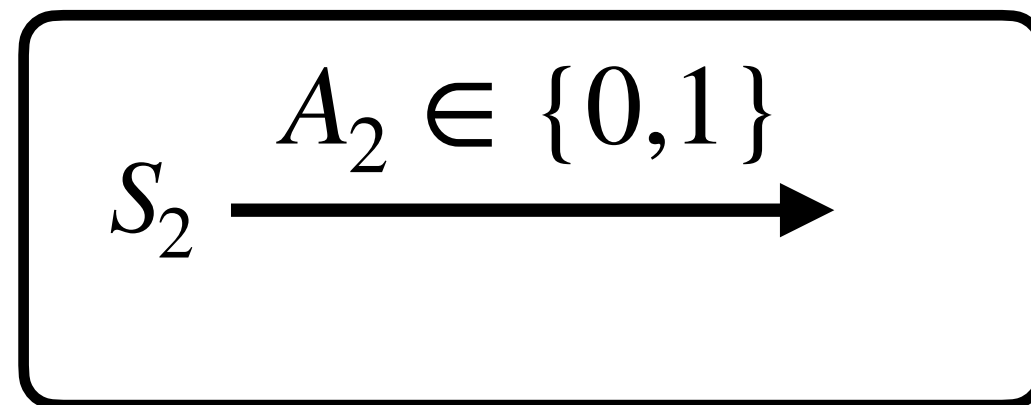
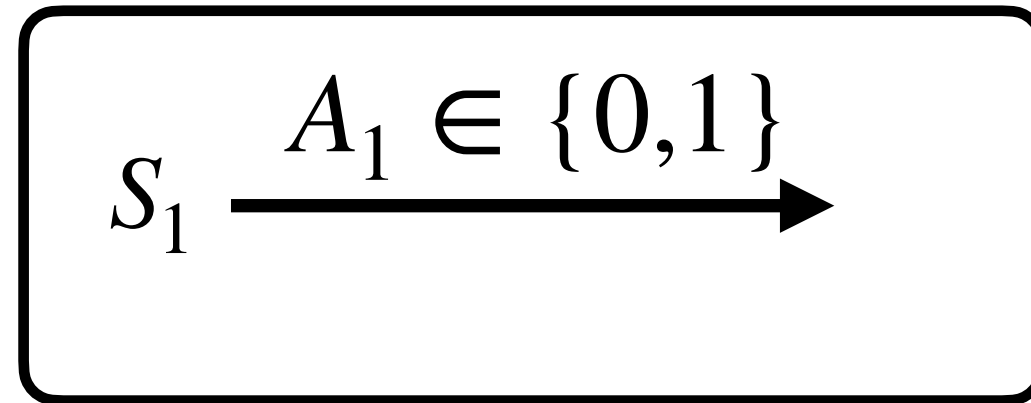
# Prior work: priority policies based on relaxation



$\max_{\pi} V_N^{\pi} \triangleq$  long run average reward under policy  $\pi$

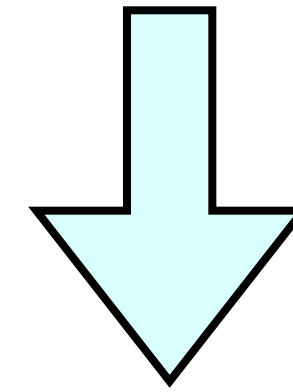
s.t.  $\sum_{i=1}^N A_i = \alpha N$ , any time slot

# Prior work: priority policies based on relaxation

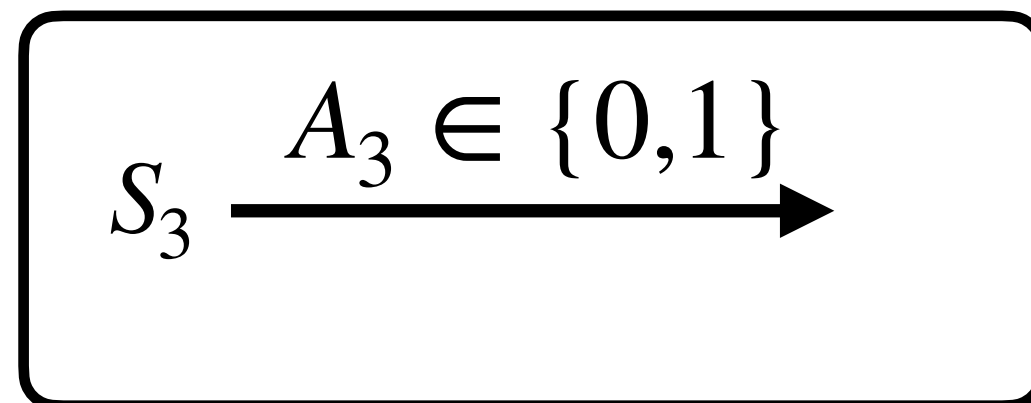
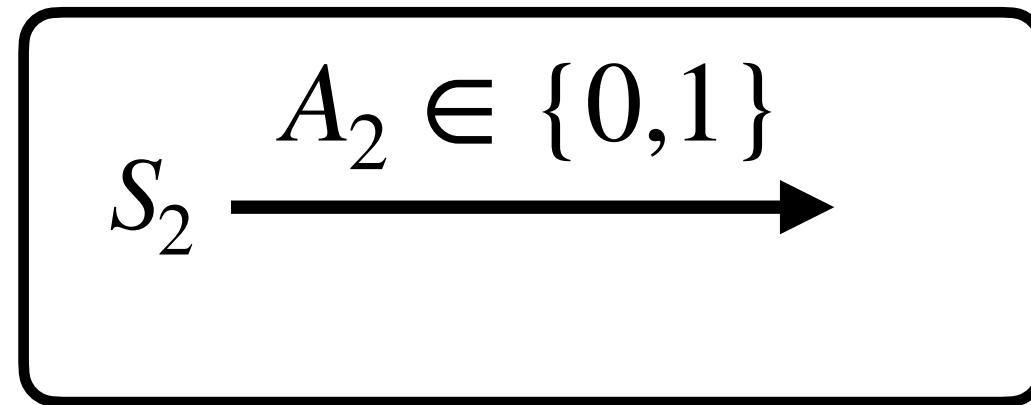
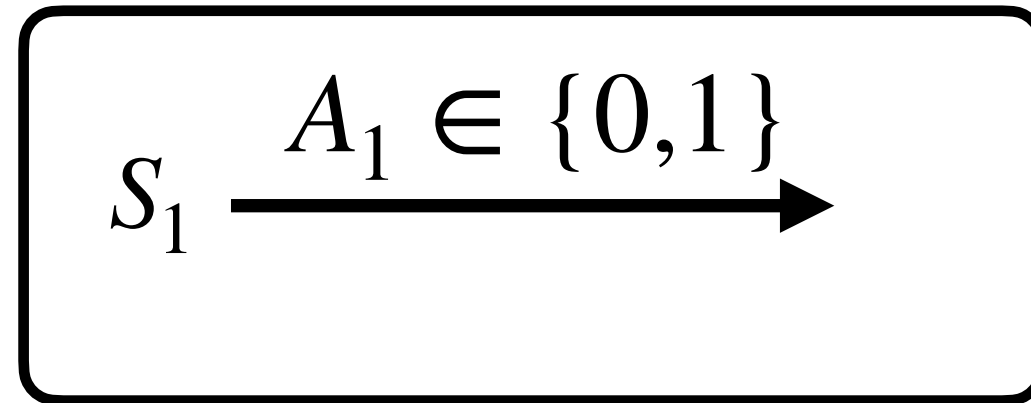


$\max_{\pi} V_N^{\pi} \triangleq$  long run average reward under policy  $\pi$

s.t.  $\sum_{i=1}^N A_i = \alpha N$ , any time slot

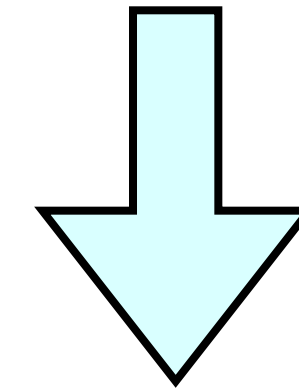


# Prior work: priority policies based on relaxation



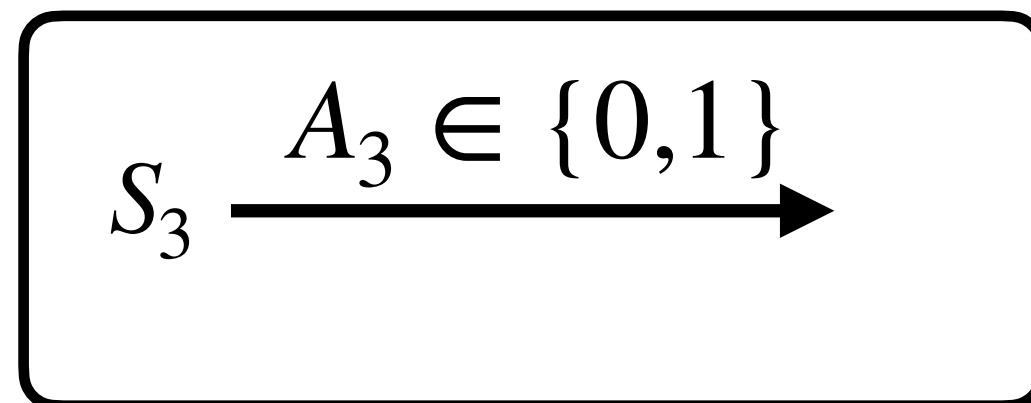
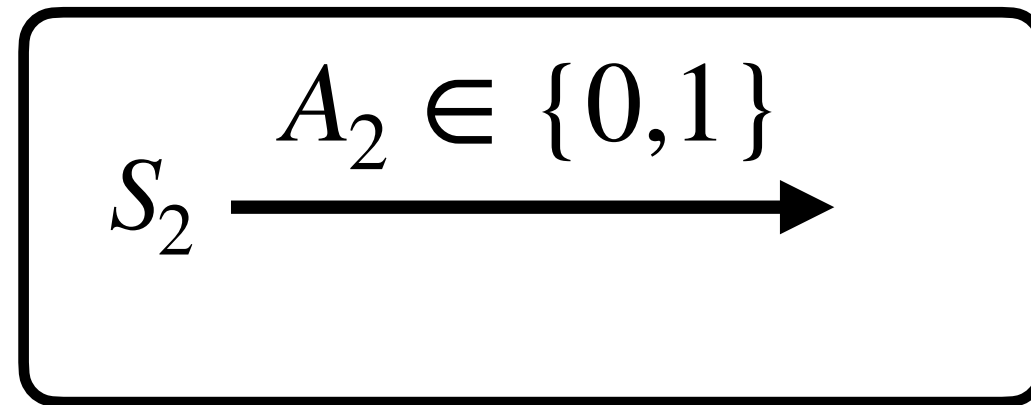
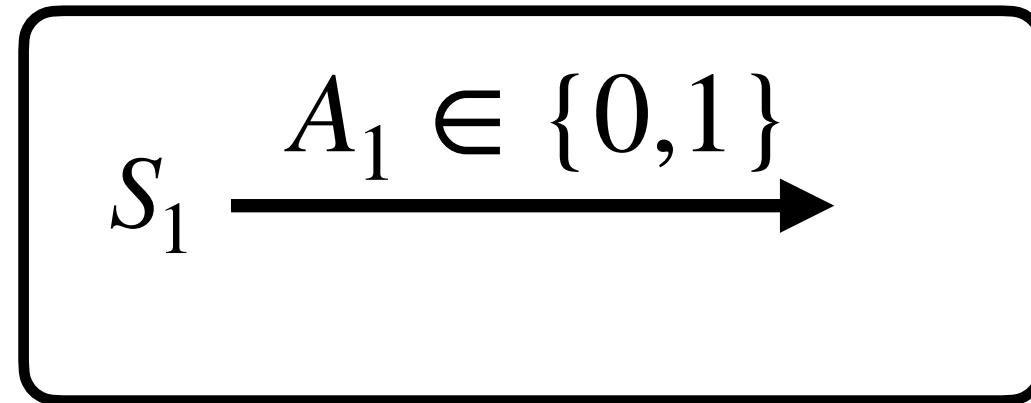
$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

$$\text{s.t. } \sum_{i=1}^N A_i = \alpha N, \text{ any time slot}$$



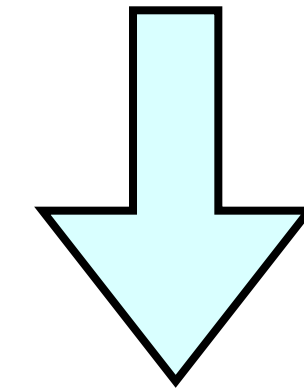
relax

# Prior work: priority policies based on relaxation



$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

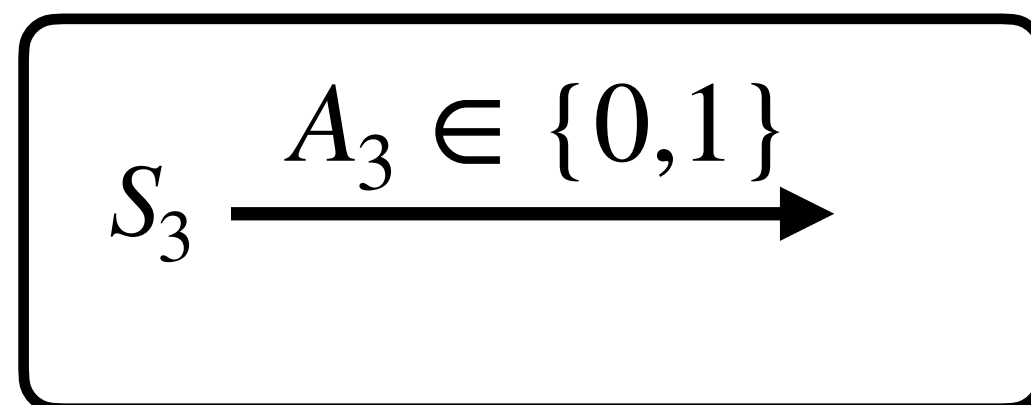
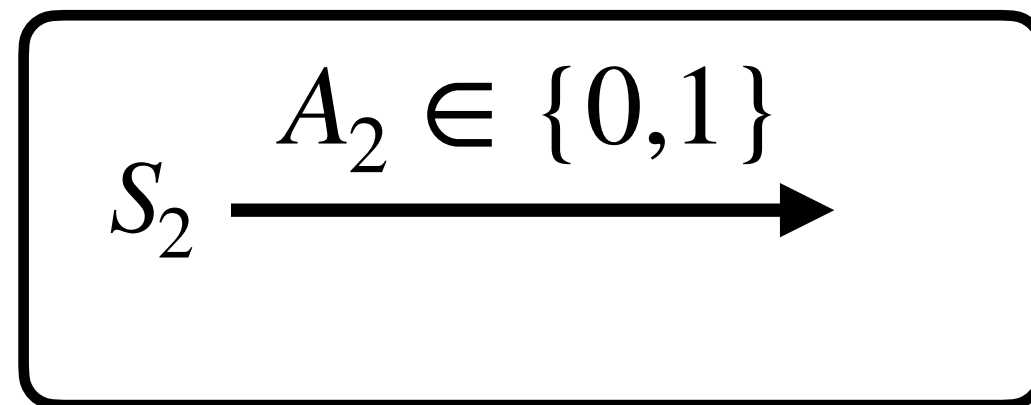
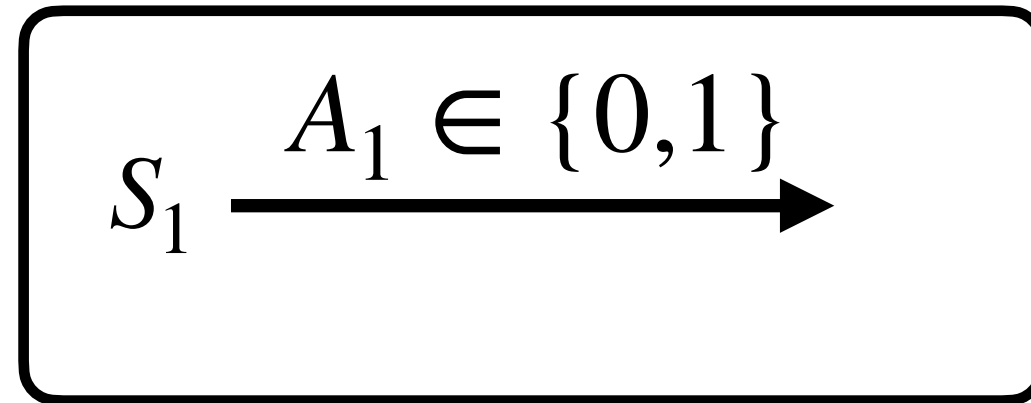
$$\text{s.t. } \sum_{i=1}^N A_i = \alpha N, \text{ any time slot}$$



relax

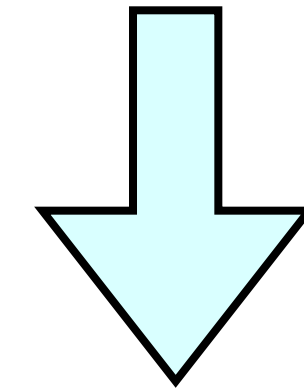
single-armed problem

# Prior work: priority policies based on relaxation



$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

$$\text{s.t. } \sum_{i=1}^N A_i = \alpha N, \text{ any time slot}$$

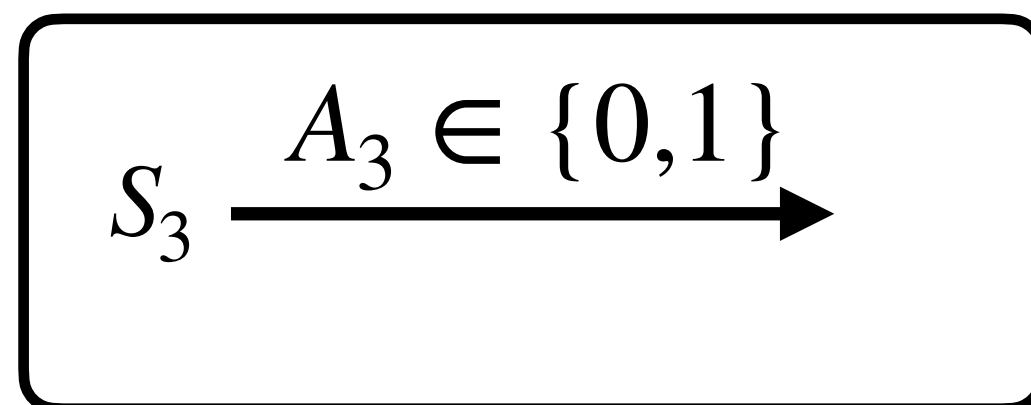
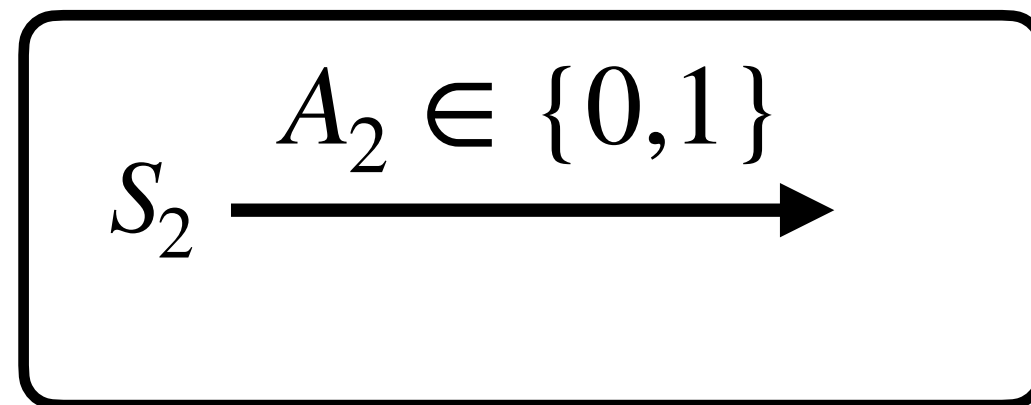
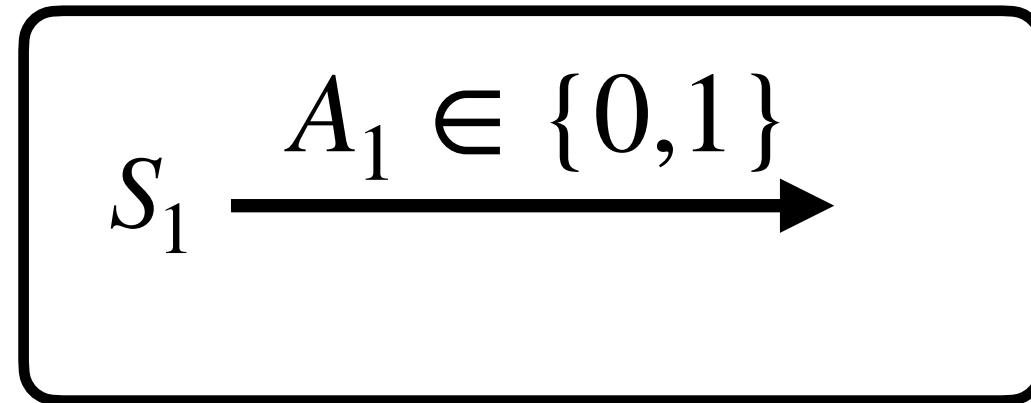


relax

single-armed problem

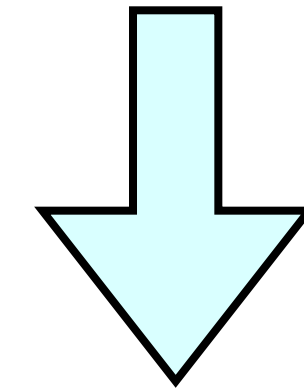
independent  
of N

# Prior work: priority policies based on relaxation



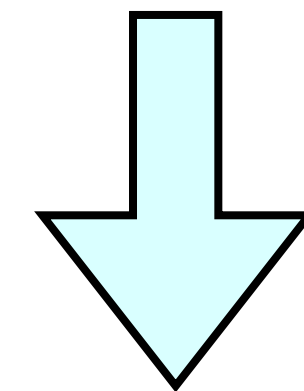
$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

$$\text{s.t. } \sum_{i=1}^N A_i = \alpha N, \text{ any time slot}$$



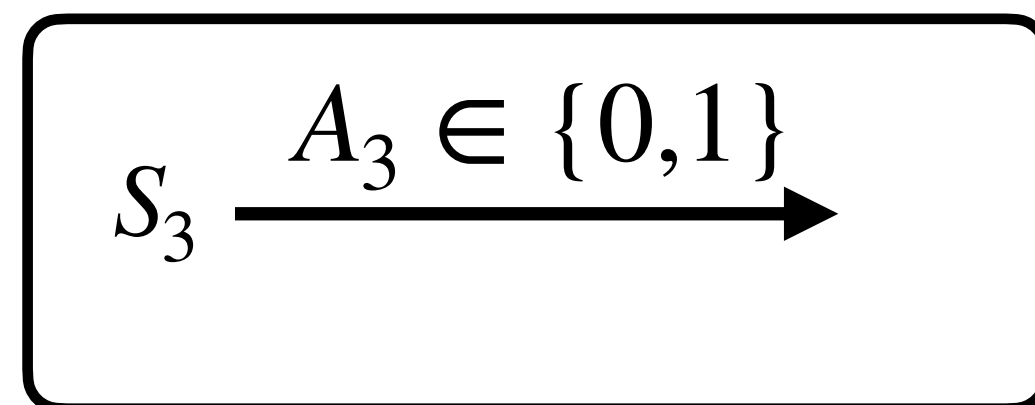
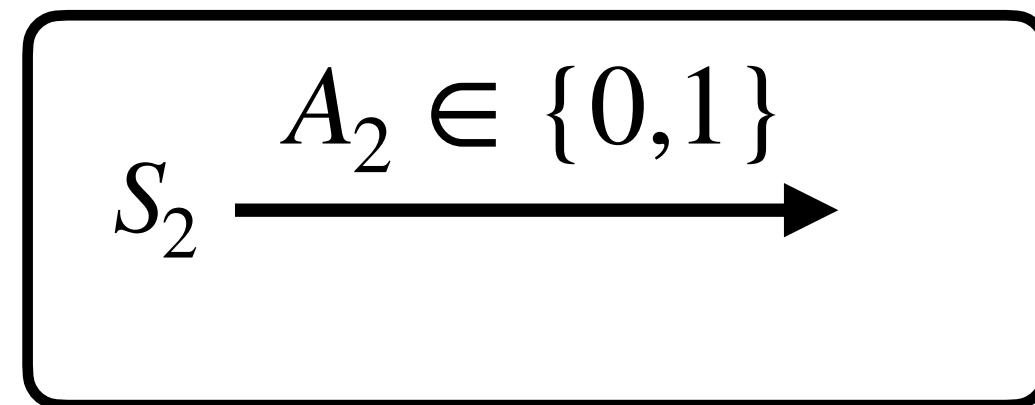
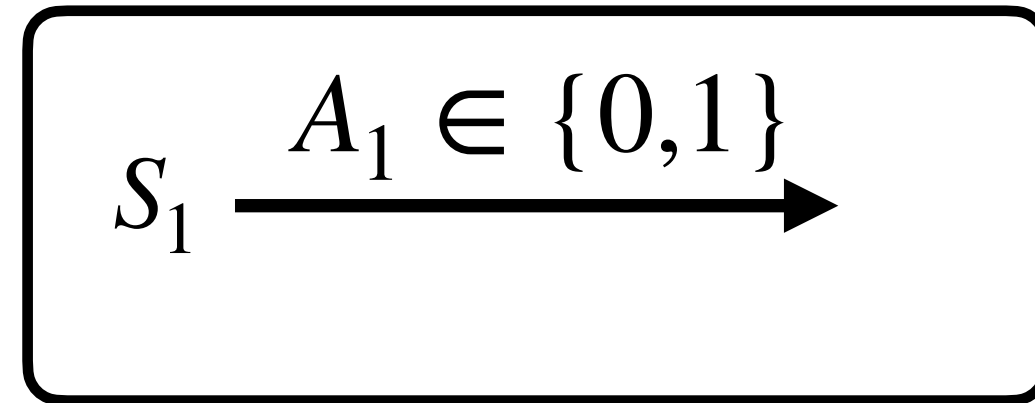
relax

single-armed problem



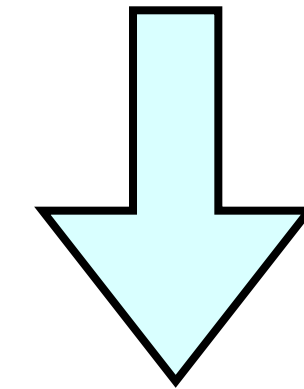
independent  
of N

# Prior work: priority policies based on relaxation



$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

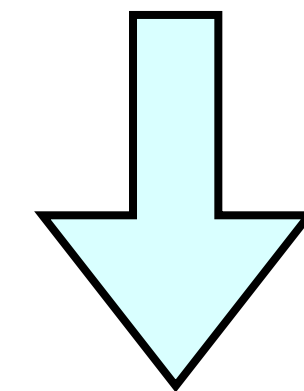
$$\text{s.t. } \sum_{i=1}^N A_i = \alpha N, \text{ any time slot}$$



relax

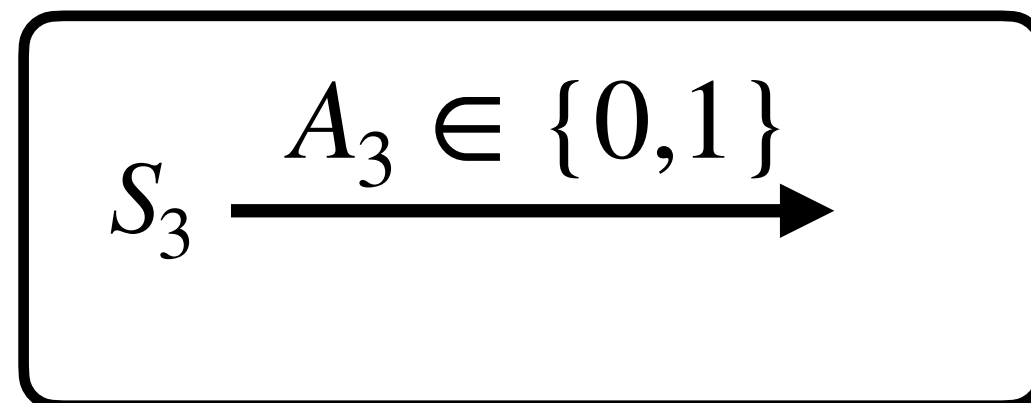
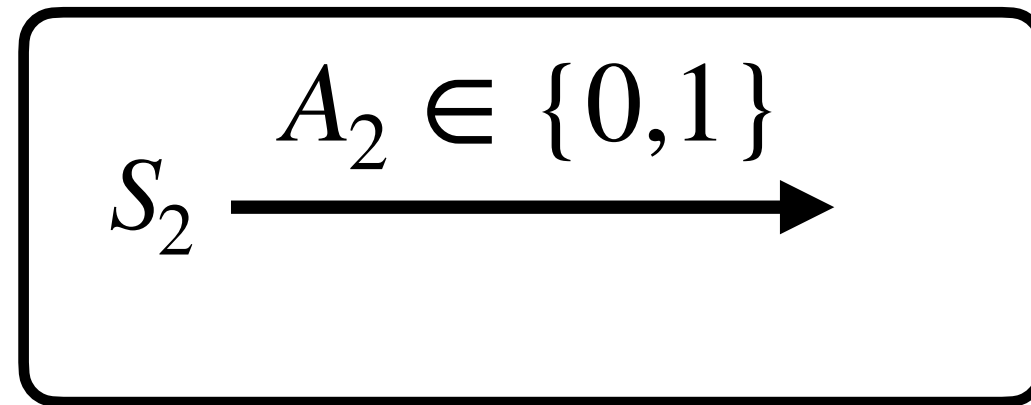
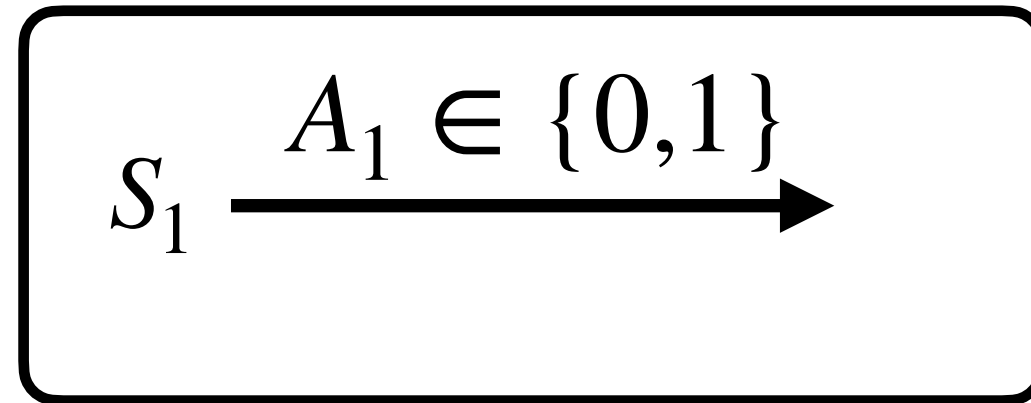
single-armed problem

independent  
of N



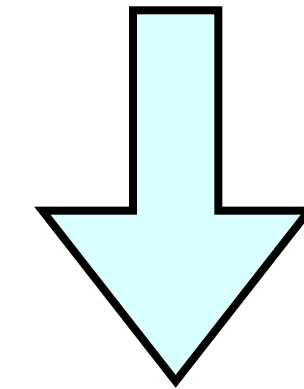
priorities for activating arms

# Prior work: priority policies based on relaxation



$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

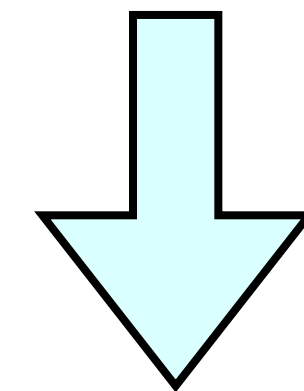
$$\text{s.t. } \sum_{i=1}^N A_i = \alpha N, \text{ any time slot}$$



relax

single-armed problem

independent  
of N



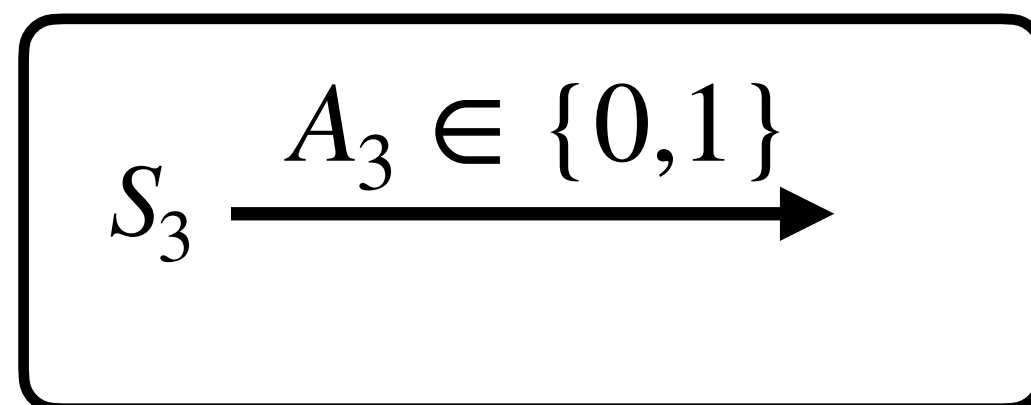
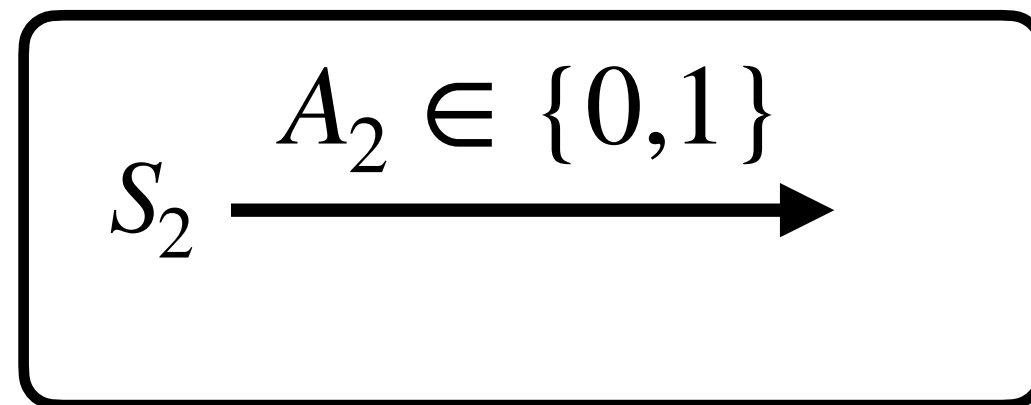
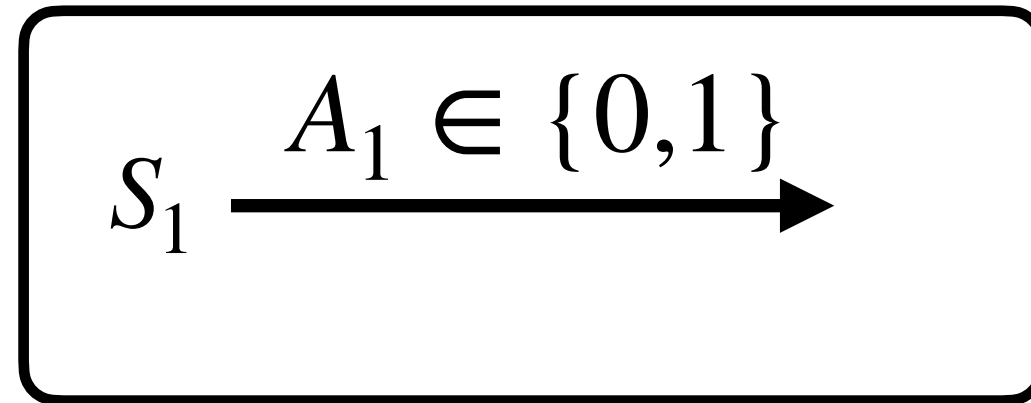
priorities for activating arms

Whittle index [Whi88, GGY20]

LP Priority [Ver16, GGY22]

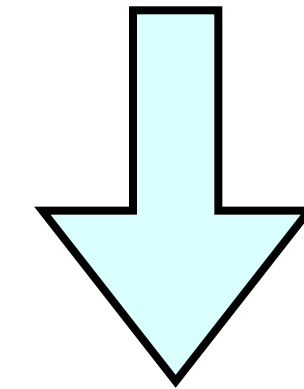


# Prior work: priority policies based on relaxation



$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

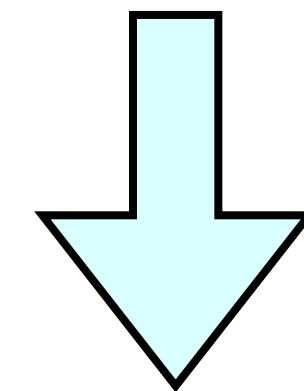
$$\text{s.t. } \sum_{i=1}^N A_i = \alpha N, \text{ any time slot}$$



relax

single-armed problem

independent  
of N



priorities for activating arms

Whittle index [Whi88, GGY20]

LP Priority [Ver16, GGY22]

**Asymptotically  
optimal?**

**Key limitation: require UGAP assumption**

# Key limitation: require UGAP assumption

All prior policies [WW90][Ver16][GGY20][GGY22]  
need *Uniform Global Attractor Property (UGAP)*  
to be asymptotic optimal

# Key limitation: require UGAP assumption

All prior policies [WW90][Ver16][GGY20][GGY22]  
need *Uniform Global Attractor Property (UGAP)*  
to be asymptotic optimal

UGAP is tricky and hard to verify

# Key limitation: require UGAP assumption

All prior policies [WW90][Ver16][GGY20][GGY22]  
need *Uniform Global Attractor Property (UGAP)*  
to be asymptotic optimal

**UGAP is tricky and hard to verify**

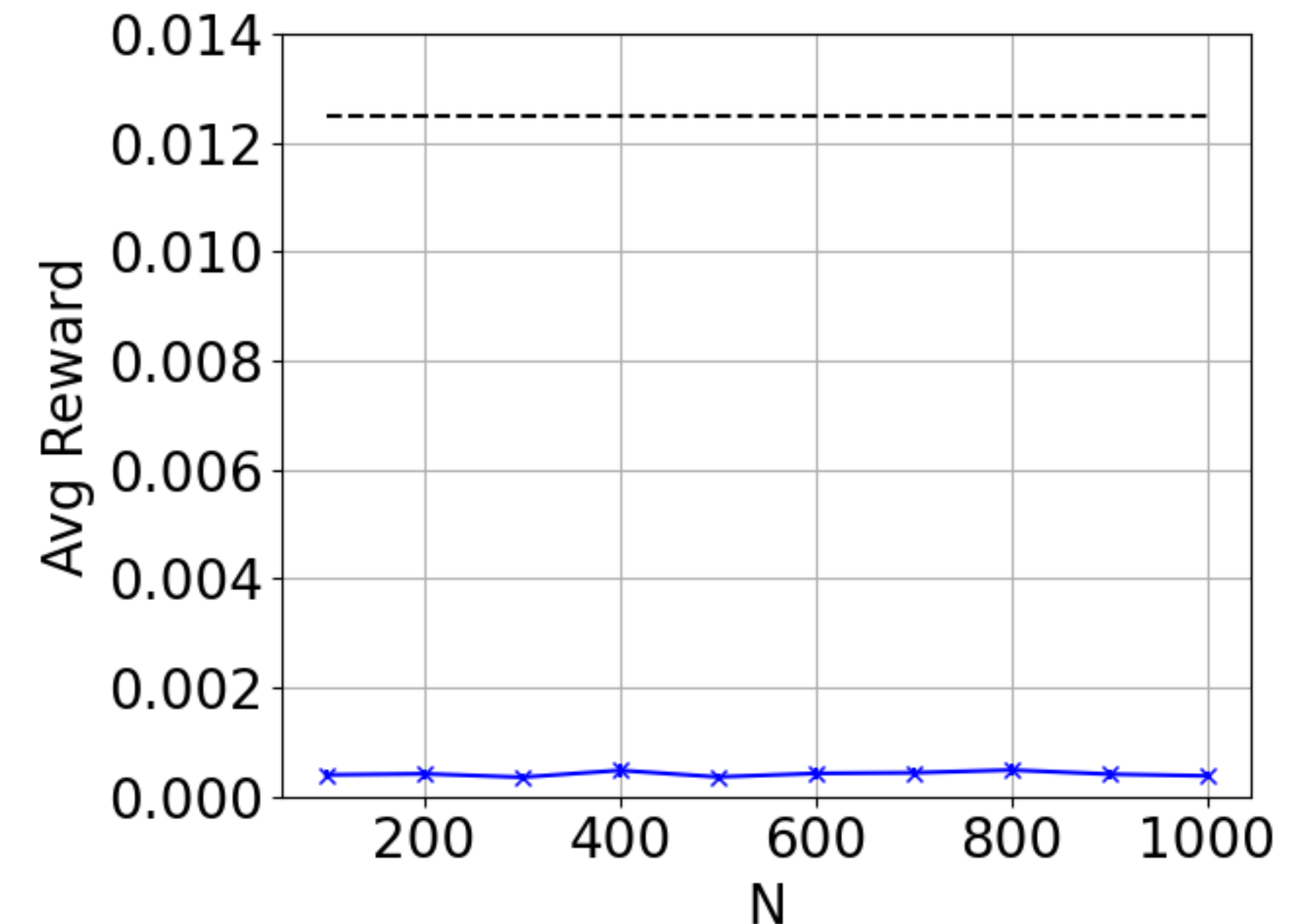
UGAP: no bad “local optimum”

# Key limitation: require UGAP assumption

All prior policies [WW90][Ver16][GGY20][GGY22] need *Uniform Global Attractor Property (UGAP)* to be asymptotic optimal

**UGAP is tricky and hard to verify**

UGAP: no bad “local optimum”

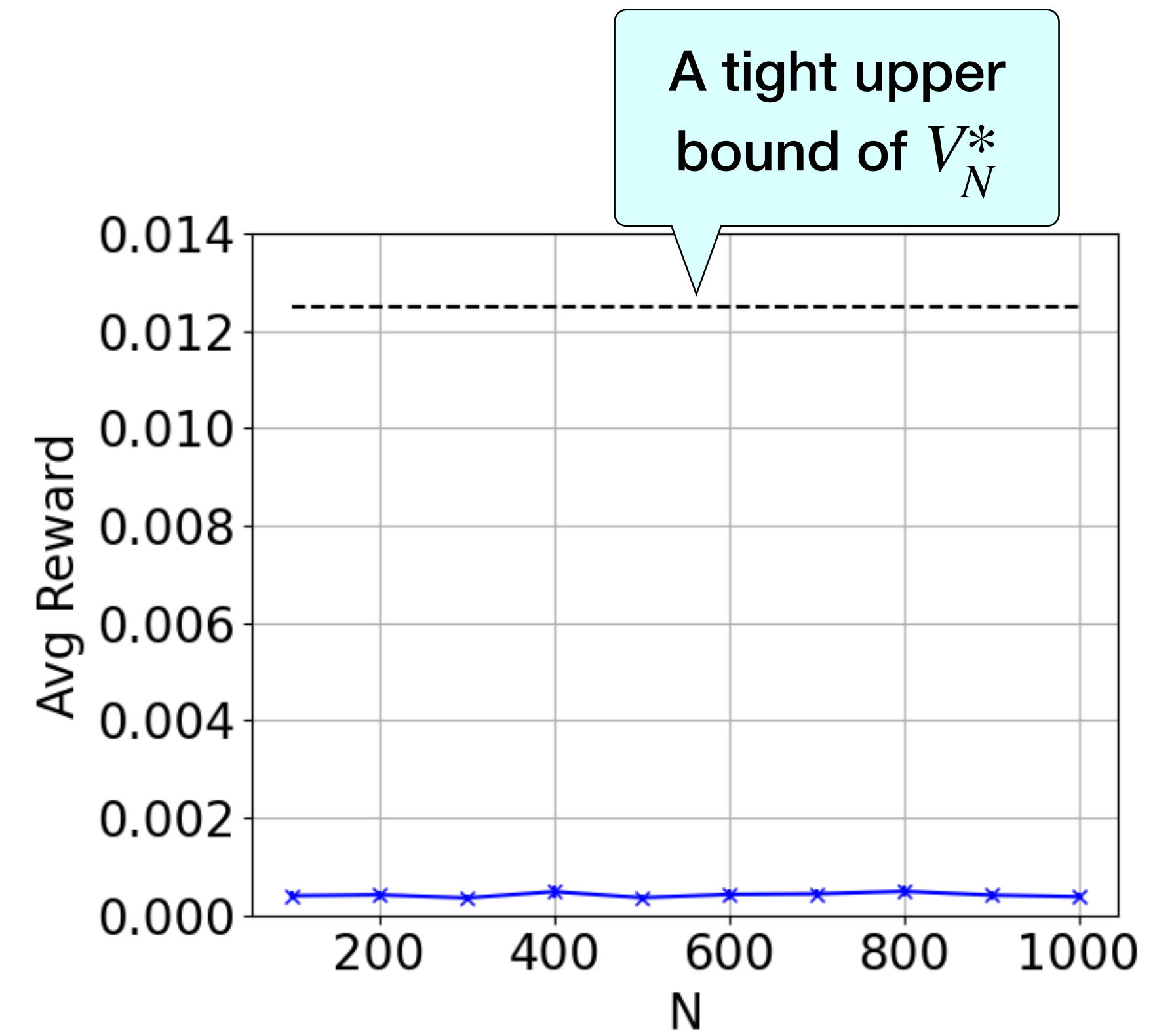


# Key limitation: require UGAP assumption

All prior policies [WW90][Ver16][GGY20][GGY22] need *Uniform Global Attractor Property (UGAP)* to be asymptotic optimal

**UGAP is tricky and hard to verify**

UGAP: no bad “local optimum”

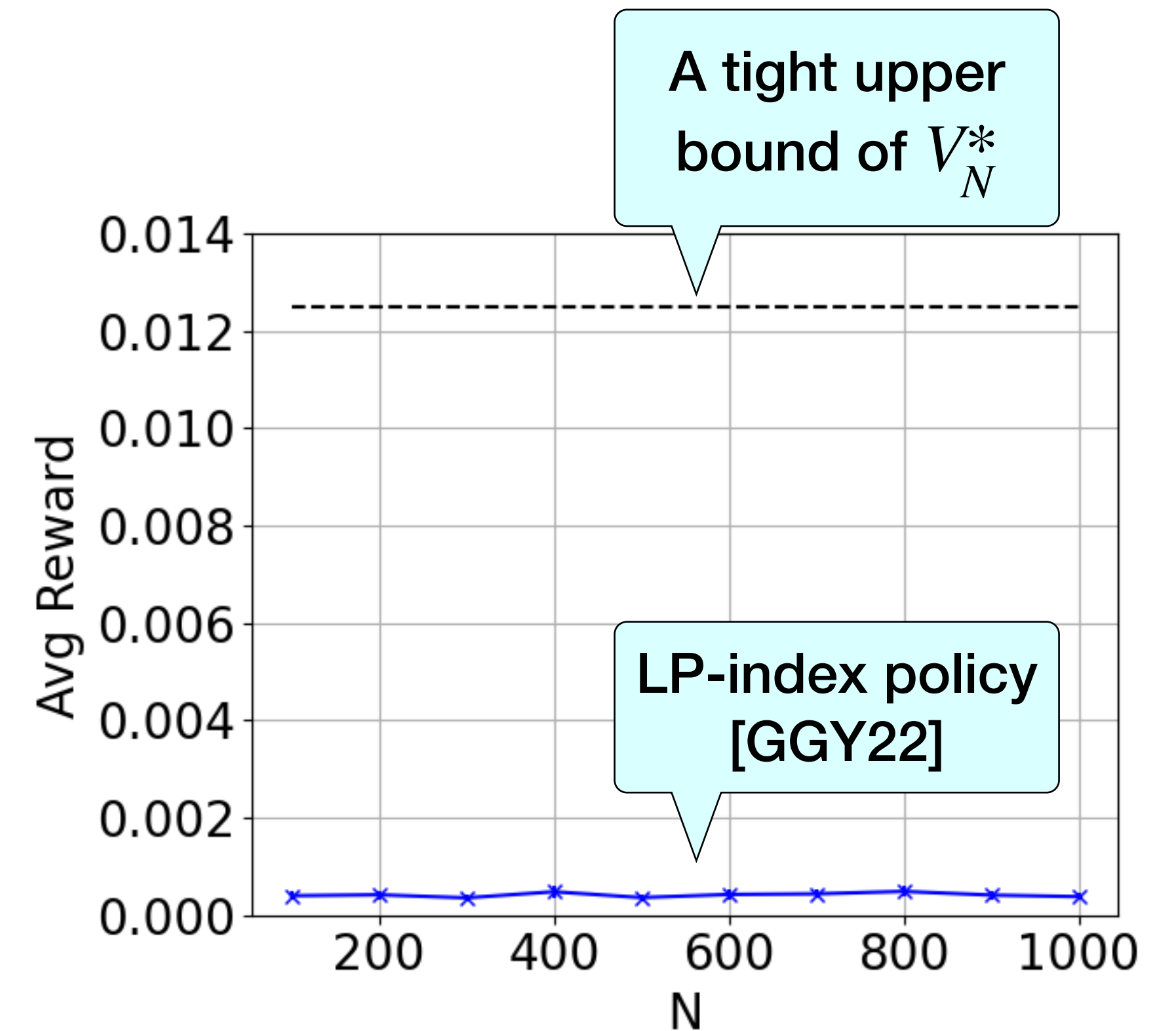


# Key limitation: require UGAP assumption

All prior policies [WW90][Ver16][GGY20][GGY22] need *Uniform Global Attractor Property (UGAP)* to be asymptotic optimal

**UGAP is tricky and hard to verify**

UGAP: no bad “local optimum”







**Can we remove UGAP?**

**Can we remove UGAP?**

**YES!**

Can we remove UGAP?

**YES!**

**We propose the first policy that is asymptotically optimal without UGAP.**

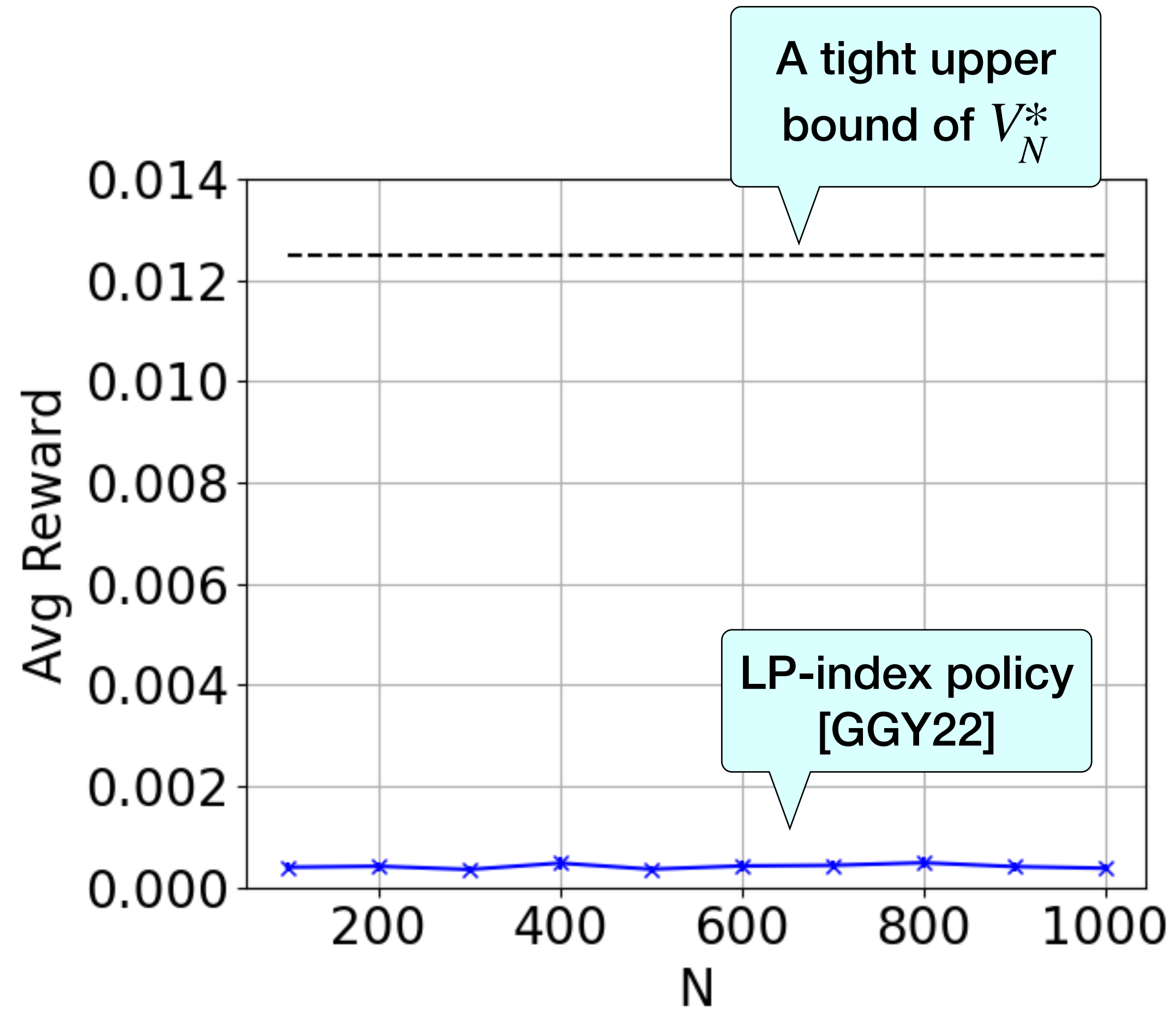
**Can we remove UGAP?**

**YES!**

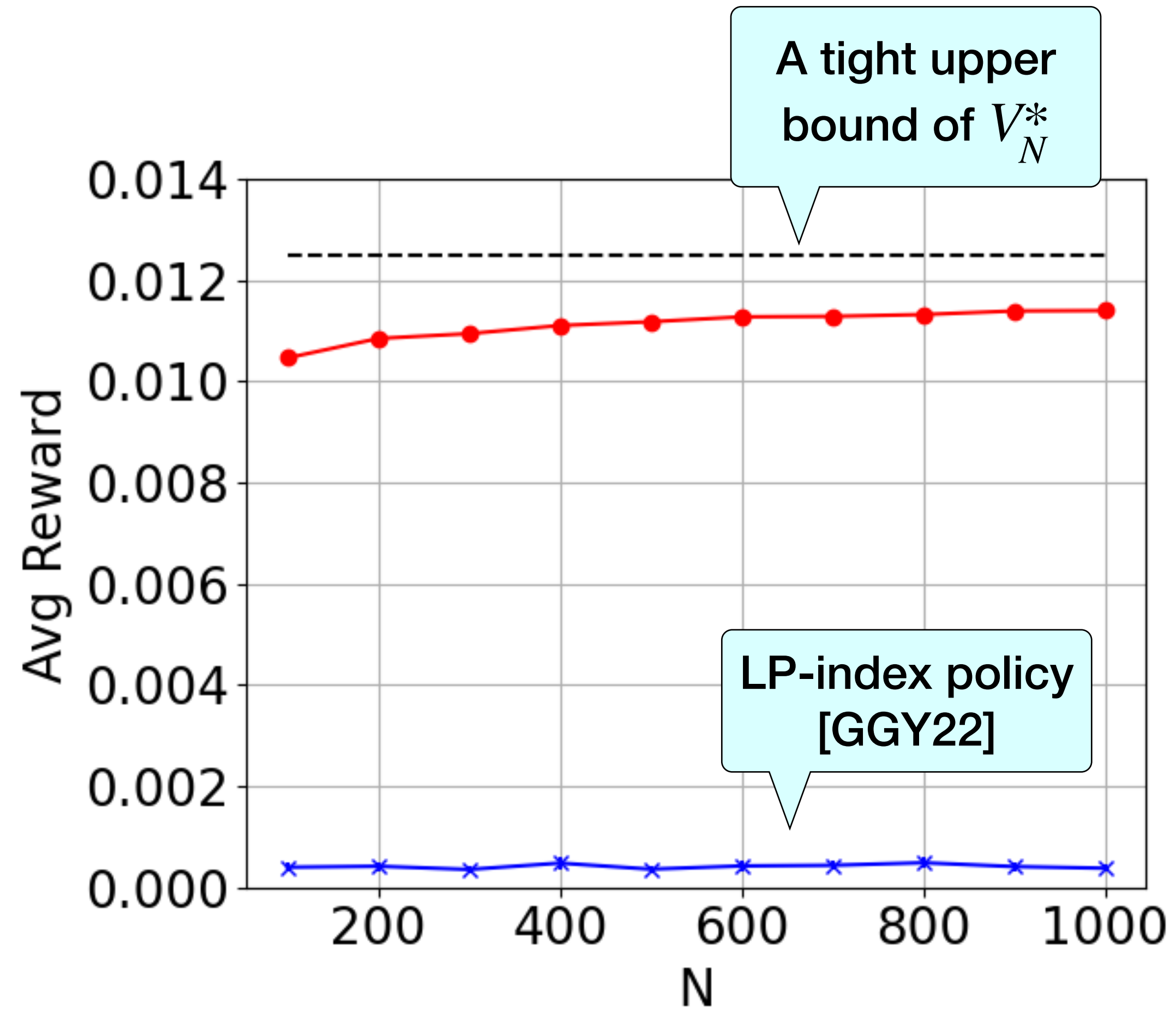
**We propose the first policy that is asymptotically optimal without UGAP.**

(See our paper for specific conditions.)

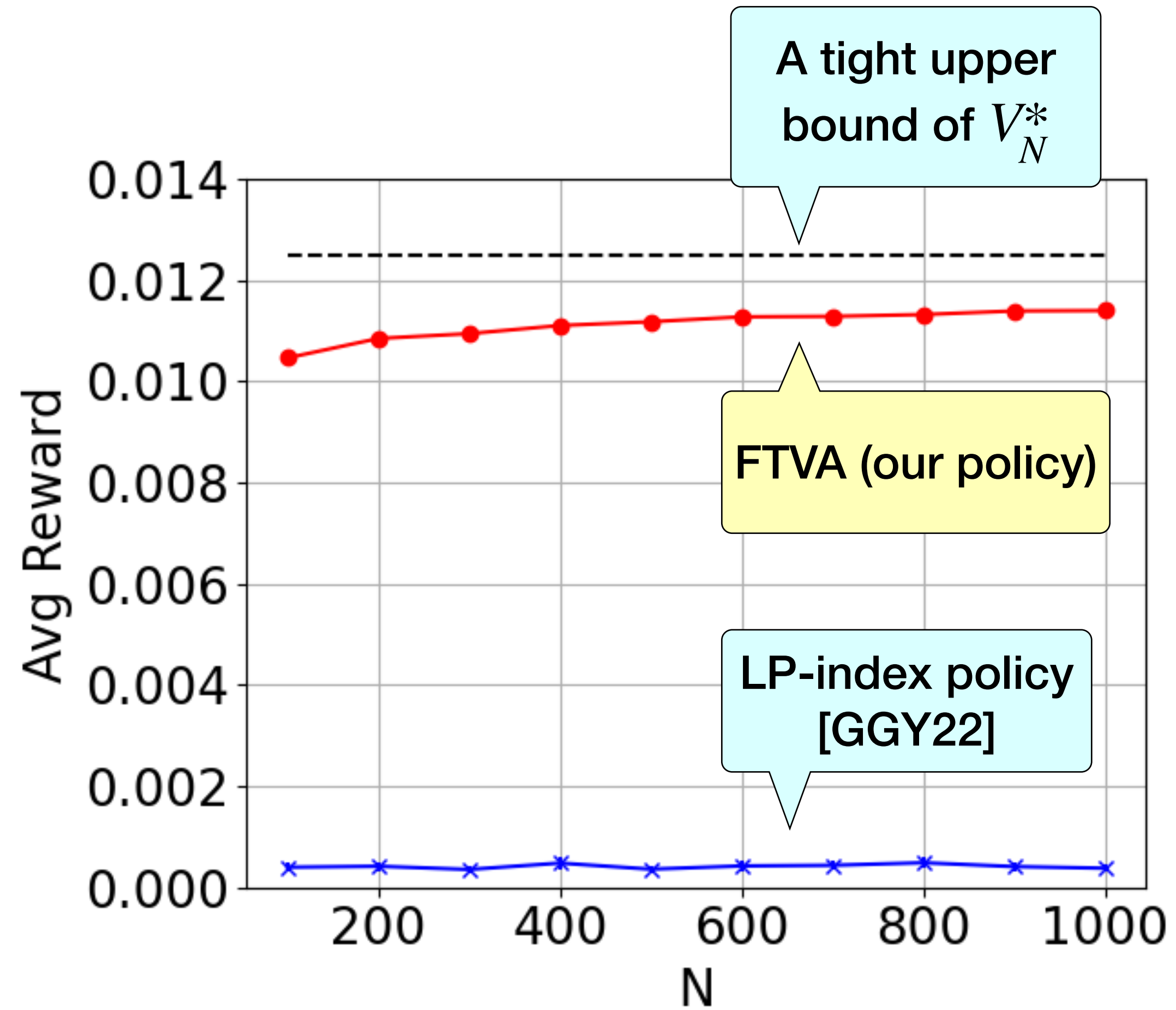
# Our policy: Follow-the-Virtual-Advice (FTVA)



# Our policy: Follow-the-Virtual-Advice (FTVA)

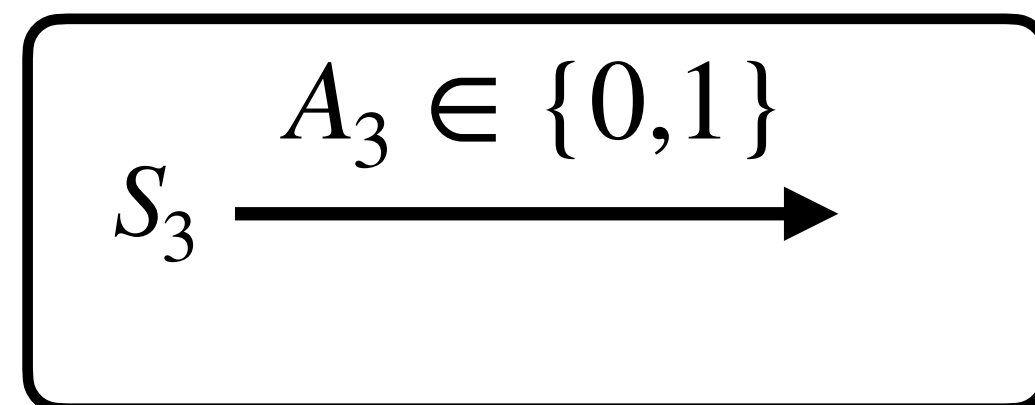
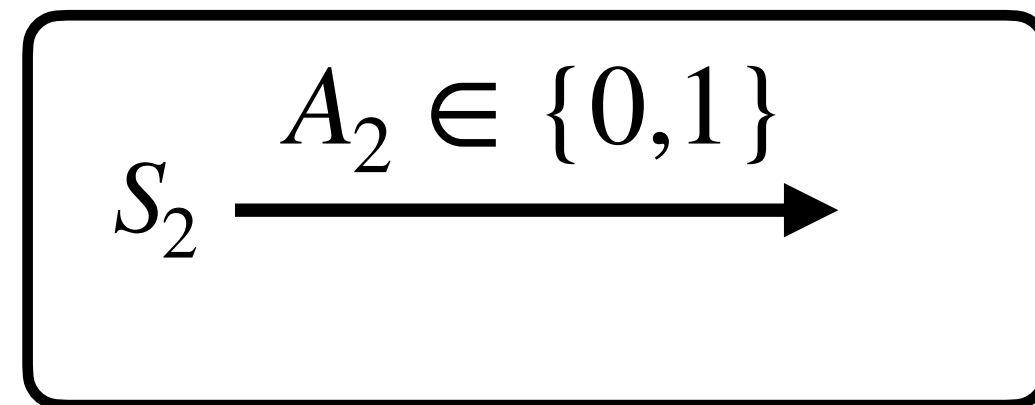
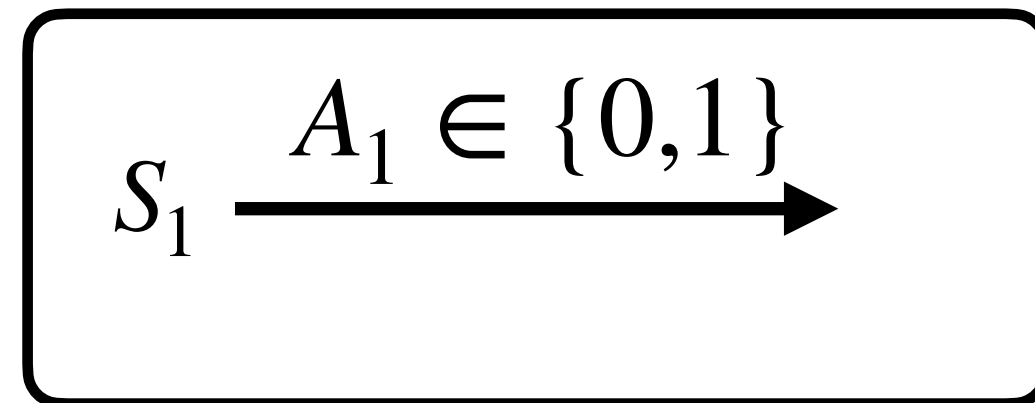


# Our policy: Follow-the-Virtual-Advice (FTVA)





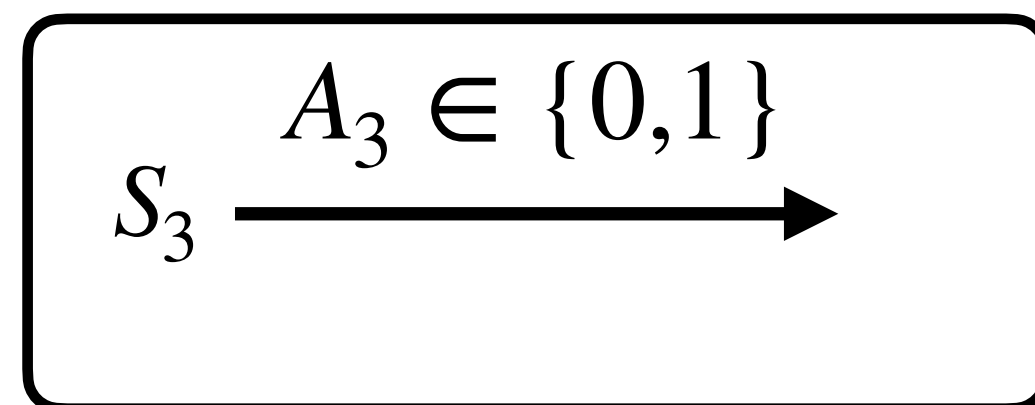
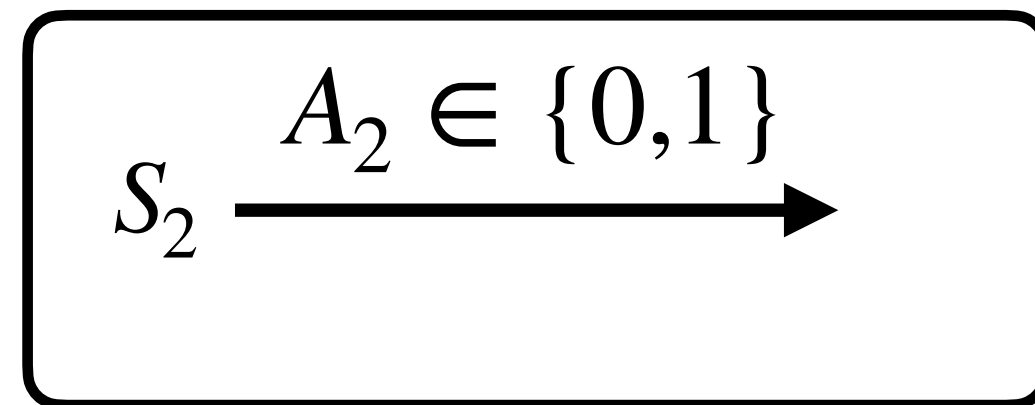
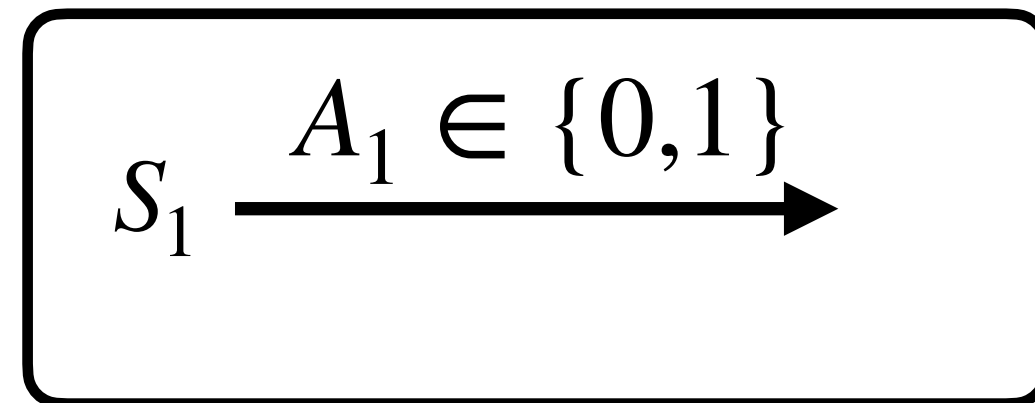
# Our policy: Follow-the-Virtual-Advice (FTVA)



$\max_{\pi} V_N^{\pi} \triangleq$  long run average reward under policy  $\pi$

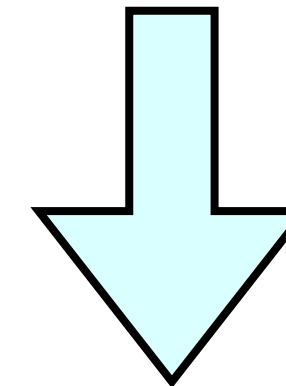
s.t.  $\sum_{i=1}^N A_i = \alpha N$ , any time slot

# Our policy: Follow-the-Virtual-Advice (FTVA)



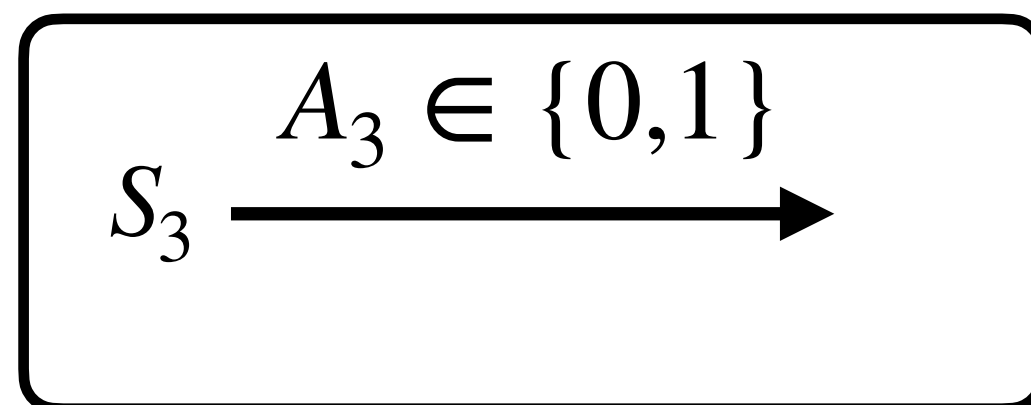
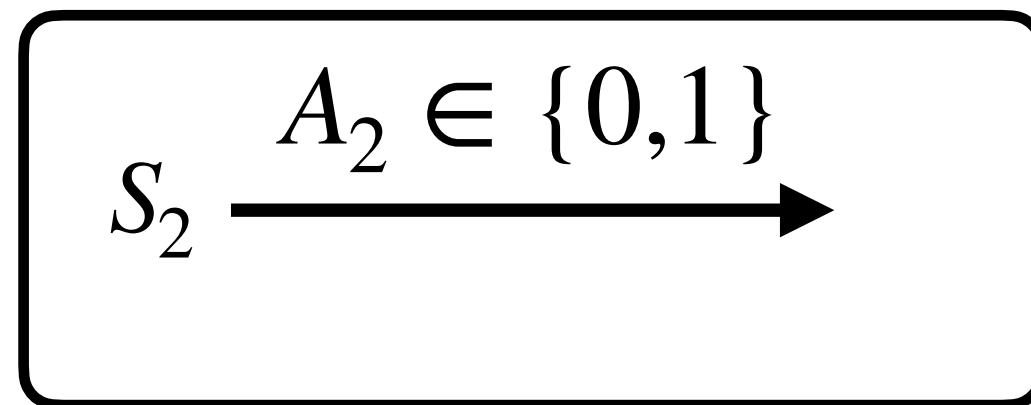
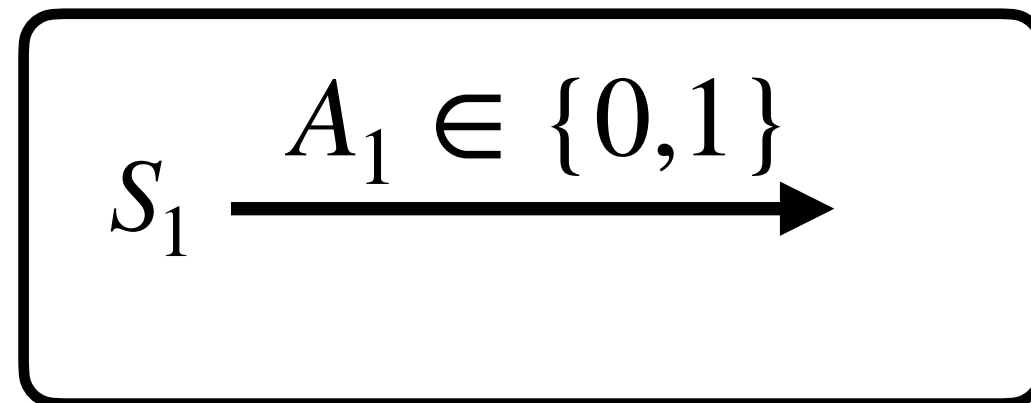
$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

$$\text{s.t. } \sum_{i=1}^N A_i = \alpha N, \text{ any time slot}$$



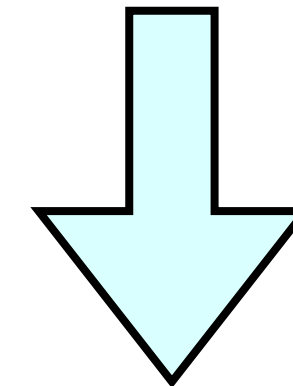
relax

# Our policy: Follow-the-Virtual-Advice (FTVA)



$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

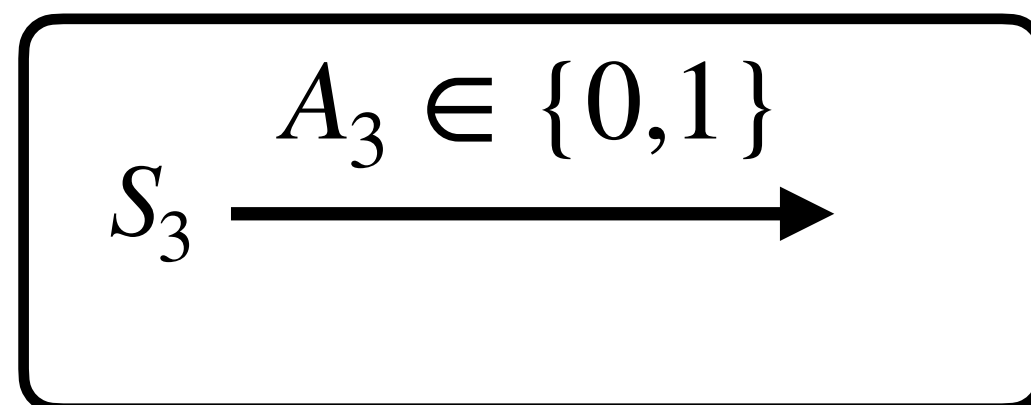
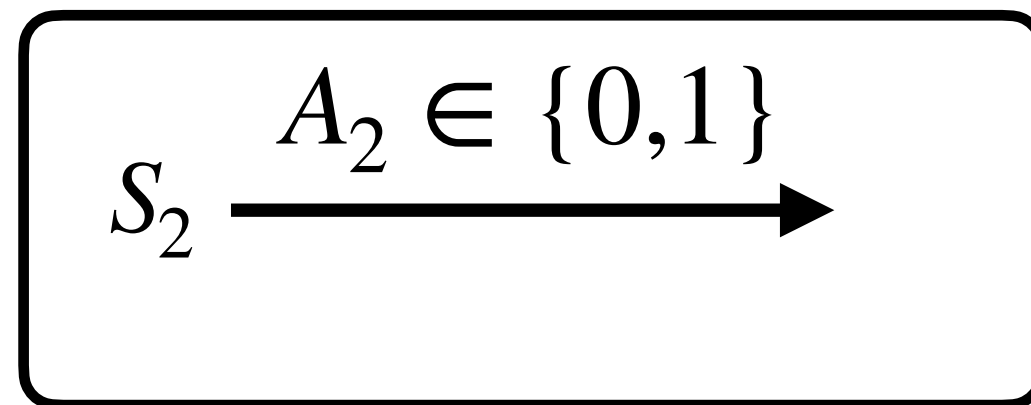
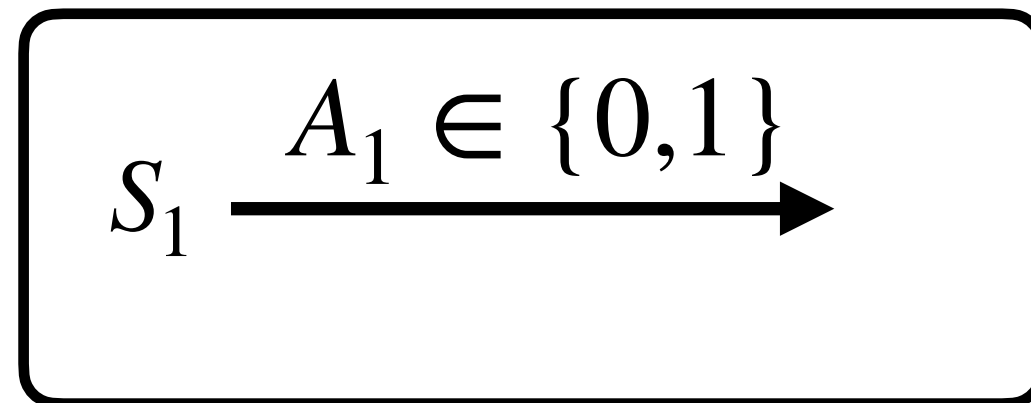
$$\text{s.t. } \sum_{i=1}^N A_i = \alpha N, \text{ any time slot}$$



relax

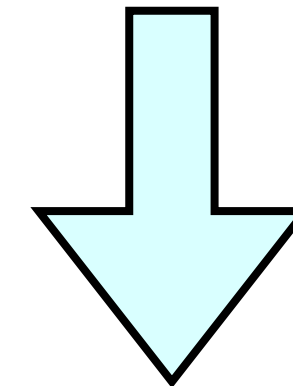
single-armed problem

# Our policy: Follow-the-Virtual-Advice (FTVA)



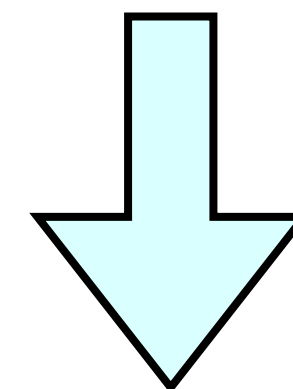
$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

$$\text{s.t. } \sum_{i=1}^N A_i = \alpha N, \text{ any time slot}$$

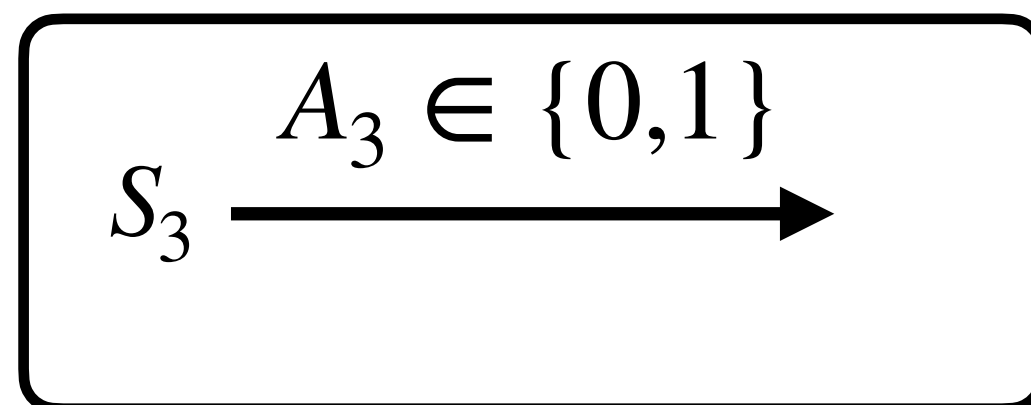
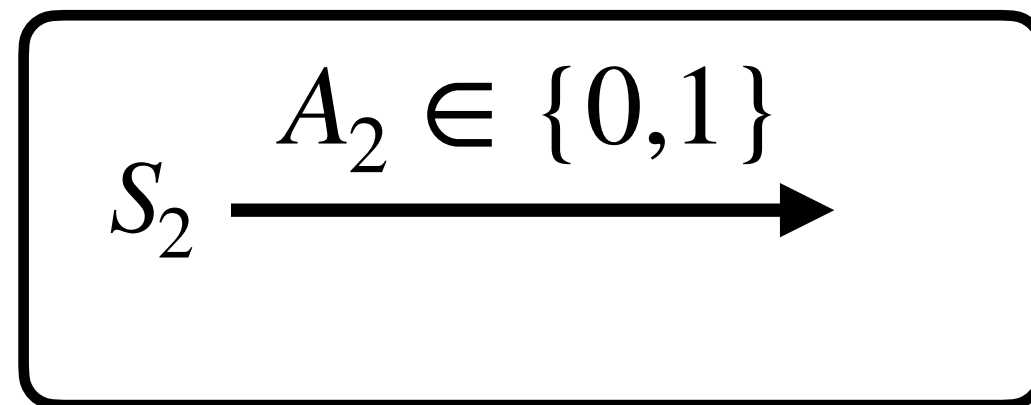
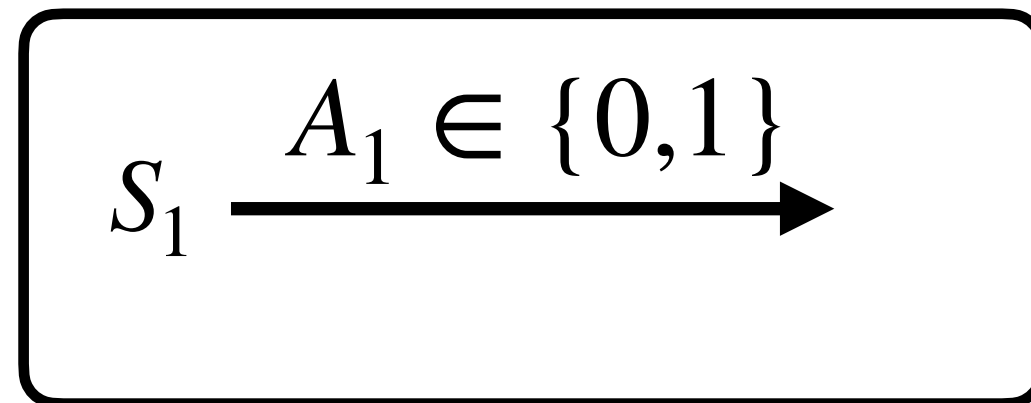


relax

single-armed problem

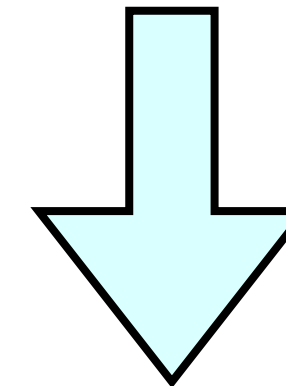


# Our policy: Follow-the-Virtual-Advice (FTVA)



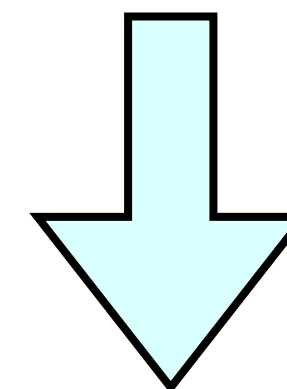
$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

$$\text{s.t. } \sum_{i=1}^N A_i = \alpha N, \text{ any time slot}$$



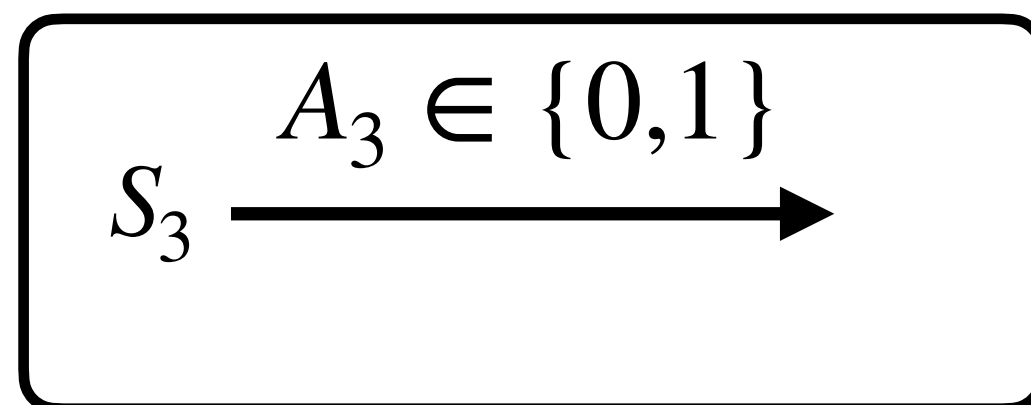
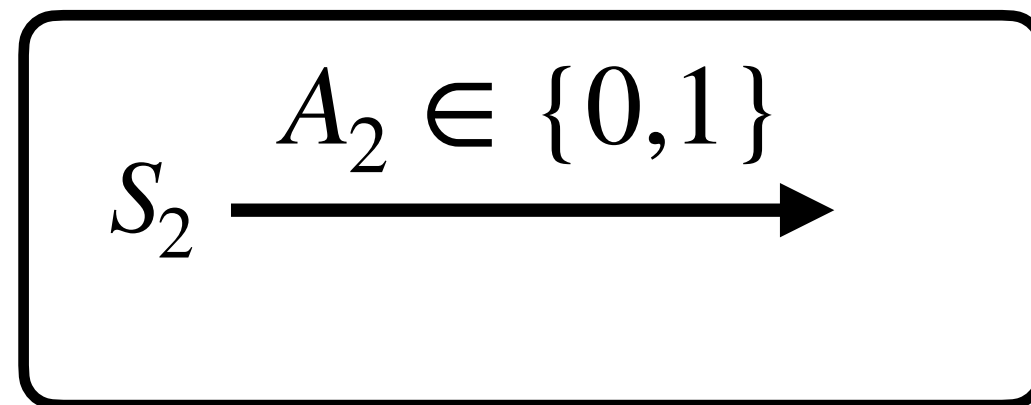
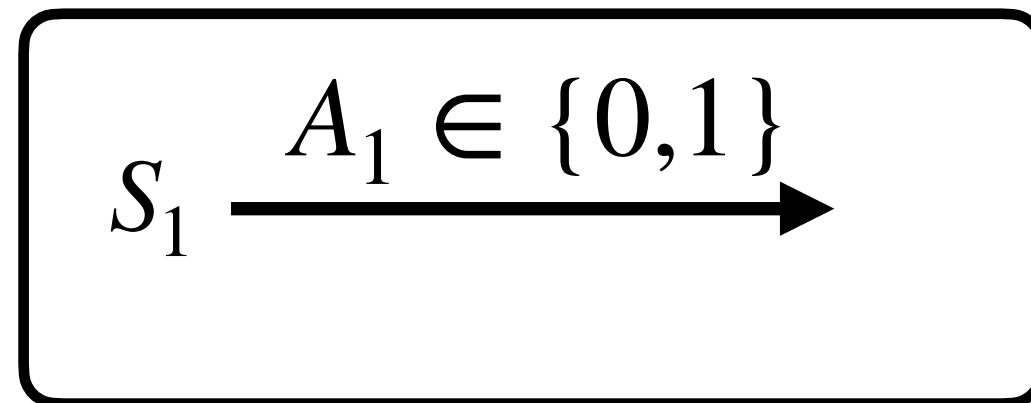
relax

single-armed problem



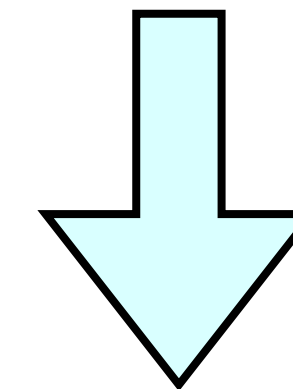
single-armed policy  $\bar{\pi}(a | s)$

# Our policy: Follow-the-Virtual-Advice (FTVA)



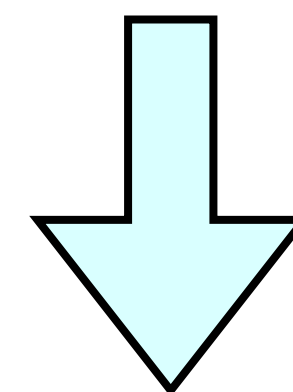
$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

$$\text{s.t. } \sum_{i=1}^N A_i = \alpha N, \text{ any time slot}$$



relax

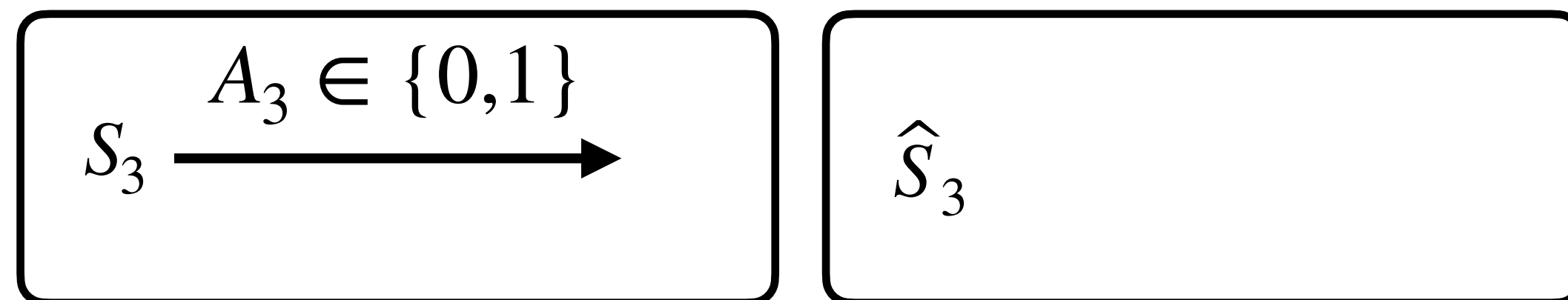
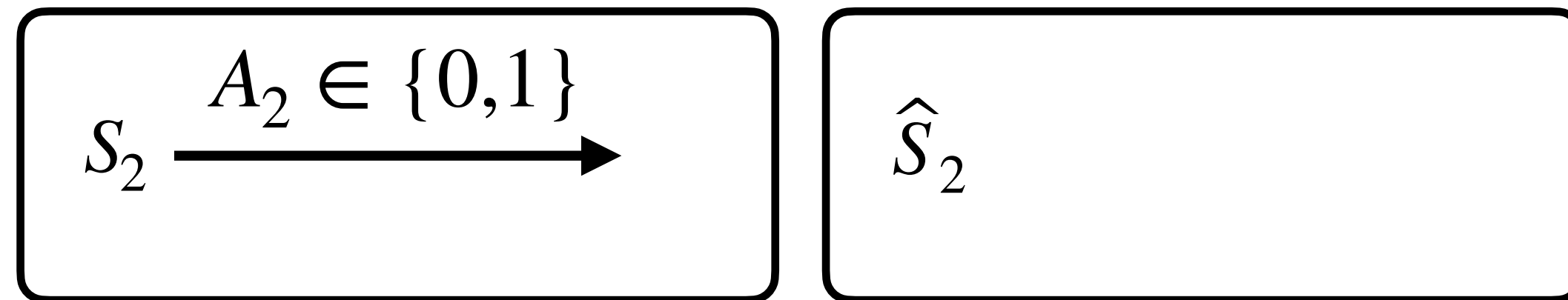
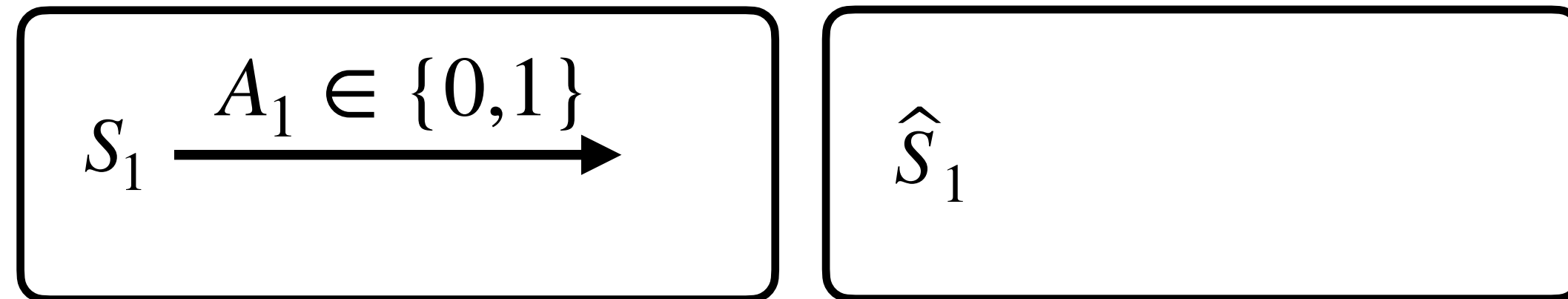
single-armed problem



single-armed policy  $\bar{\pi}(a | s)$

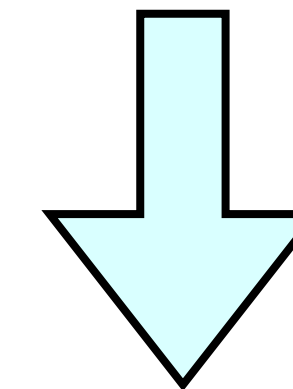
generate ideal actions

# Our policy: Follow-the-Virtual-Advice (FTVA)



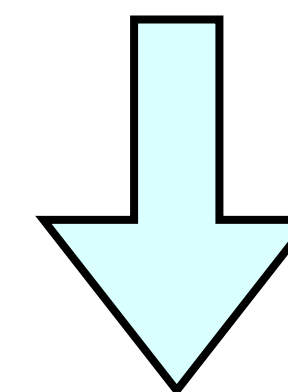
$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

$$\text{s.t. } \sum_{i=1}^N A_i = \alpha N, \text{ any time slot}$$



relax

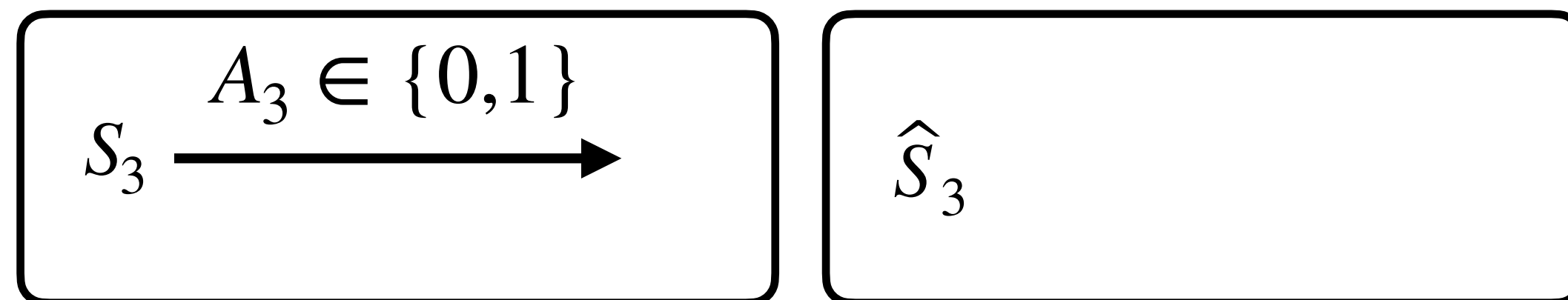
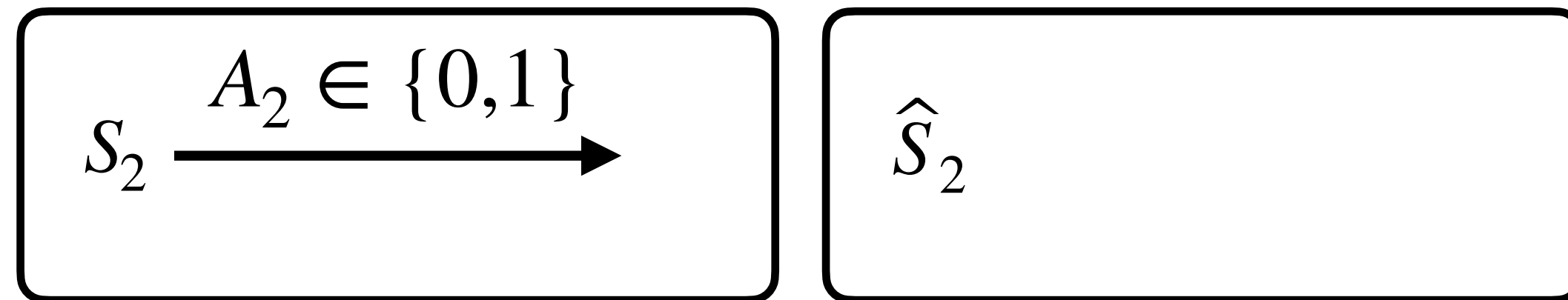
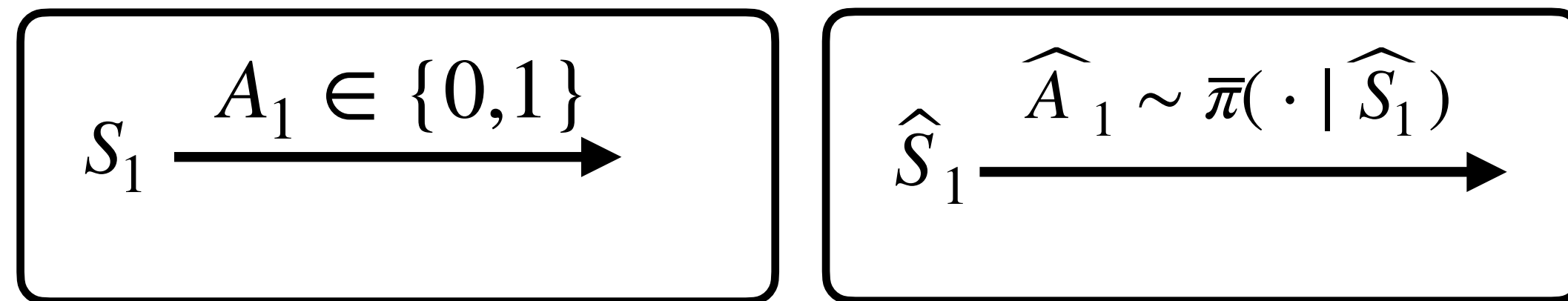
single-armed problem



single-armed policy  $\bar{\pi}(a | s)$

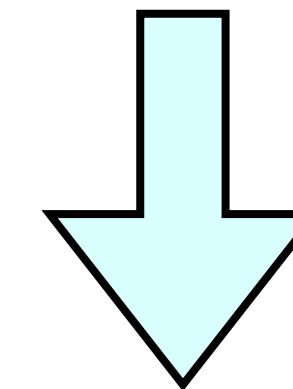
generate ideal actions

# Our policy: Follow-the-Virtual-Advice (FTVA)



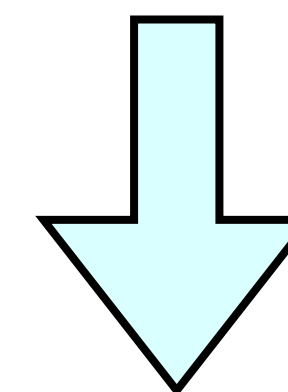
$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

$$\text{s.t. } \sum_{i=1}^N A_i = \alpha N, \text{ any time slot}$$



relax

single-armed problem

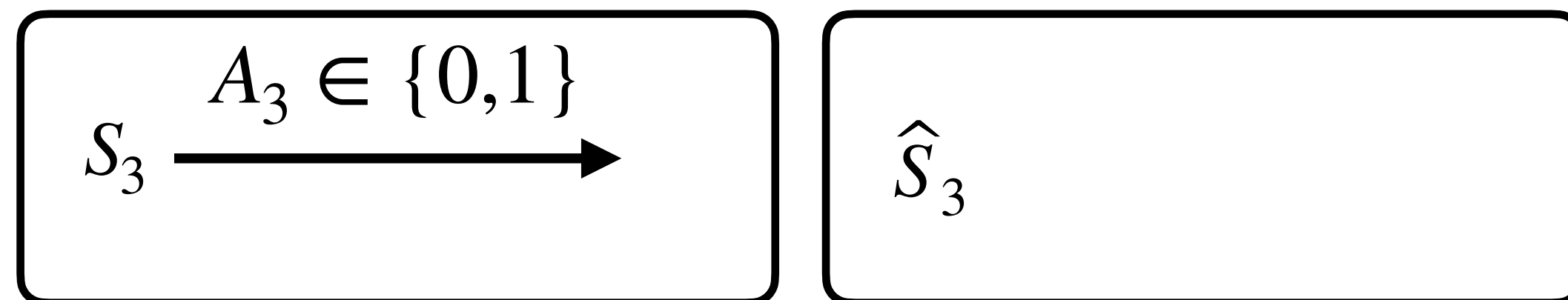
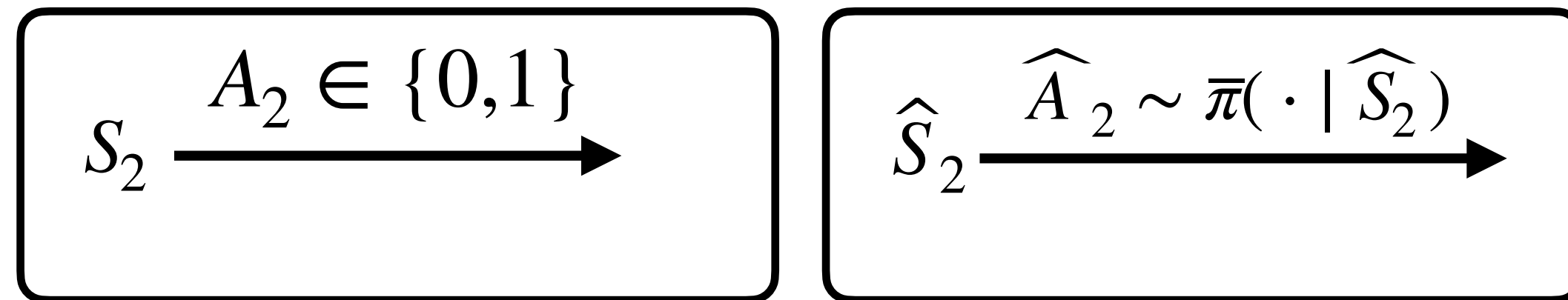
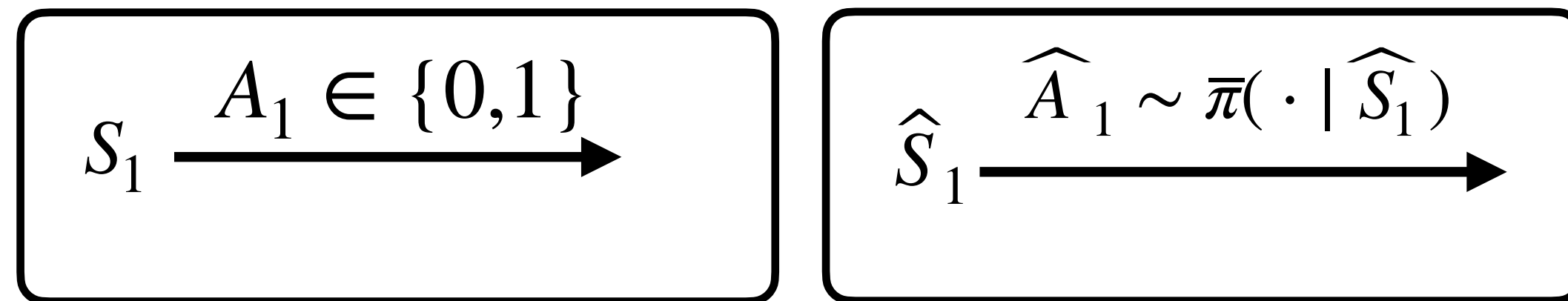


single-armed policy  $\bar{\pi}(a | s)$

generate ideal actions

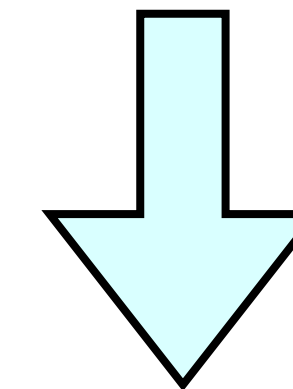


# Our policy: Follow-the-Virtual-Advice (FTVA)



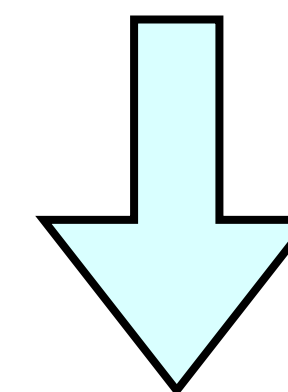
$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

$$\text{s.t. } \sum_{i=1}^N A_i = \alpha N, \text{ any time slot}$$



relax

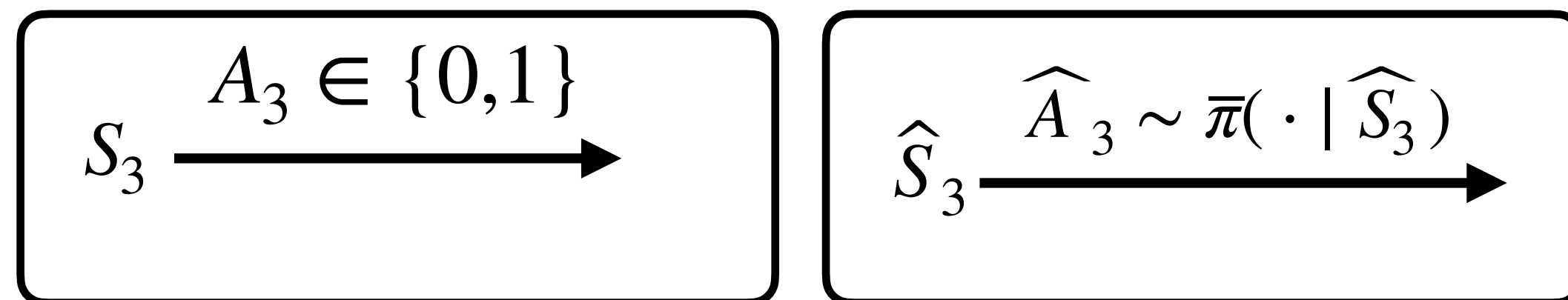
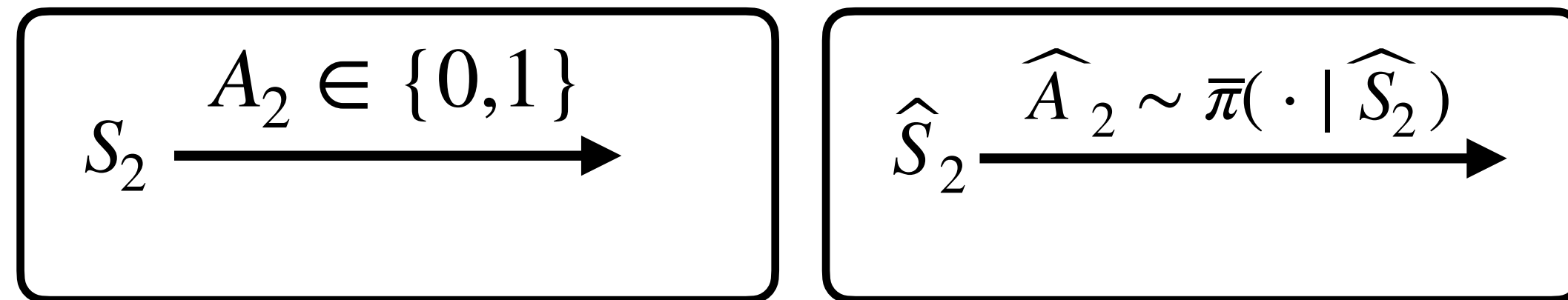
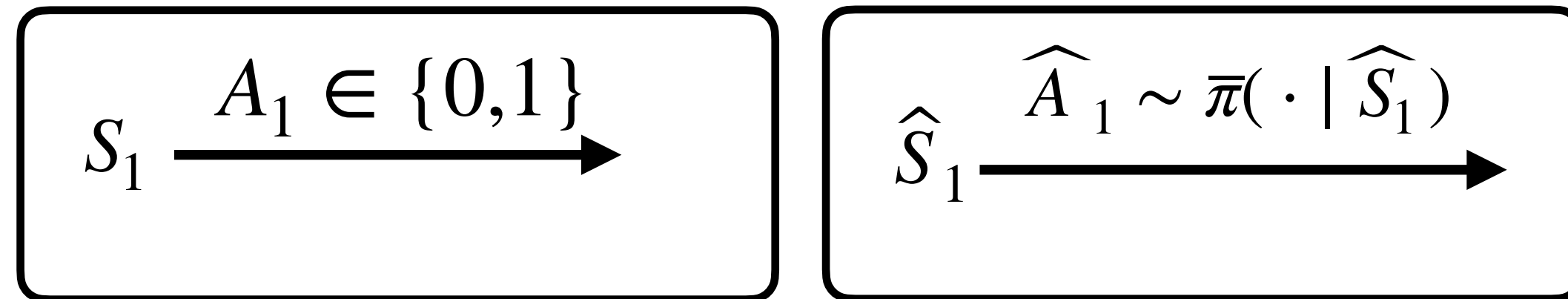
single-armed problem



single-armed policy  $\bar{\pi}(a | s)$

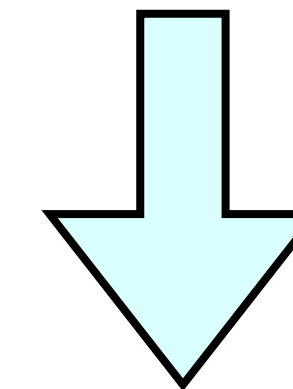
generate ideal actions

# Our policy: Follow-the-Virtual-Advice (FTVA)



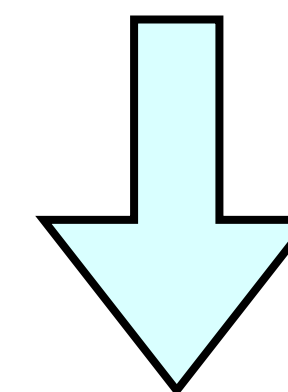
$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

$$\text{s.t. } \sum_{i=1}^N A_i = \alpha N, \text{ any time slot}$$



relax

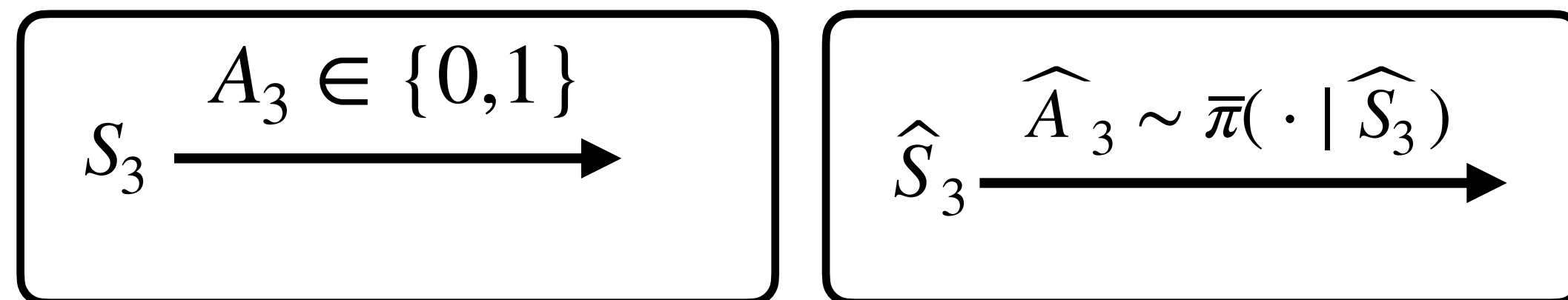
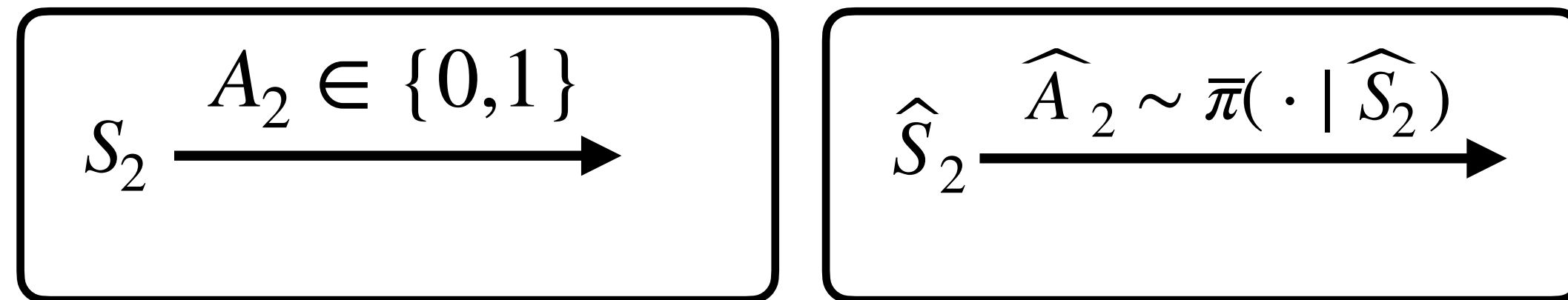
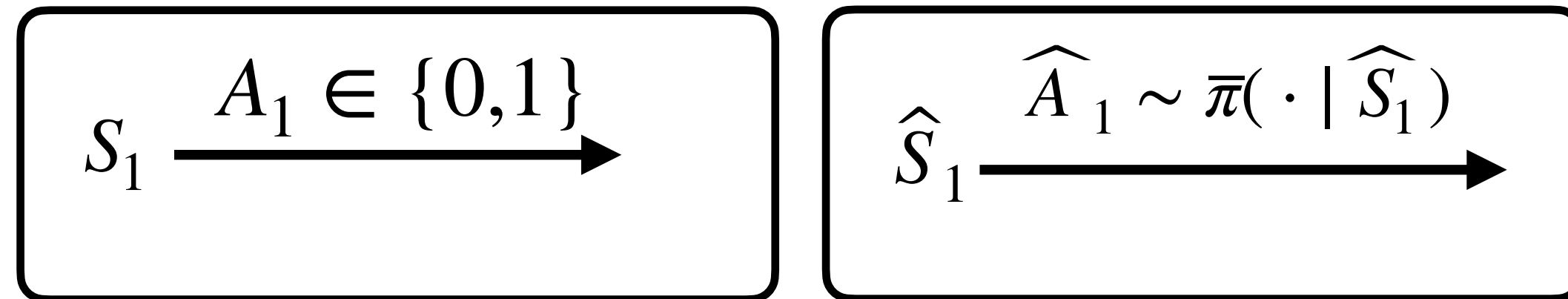
single-armed problem



single-armed policy  $\bar{\pi}(a | s)$

generate ideal actions

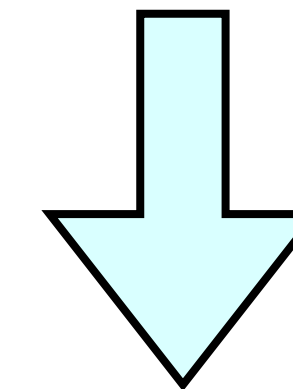
# Our policy: Follow-the-Virtual-Advice (FTVA)



no constraint

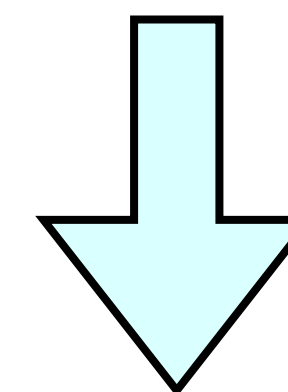
$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

$$\text{s.t. } \sum_{i=1}^N A_i = \alpha N, \text{ any time slot}$$



relax

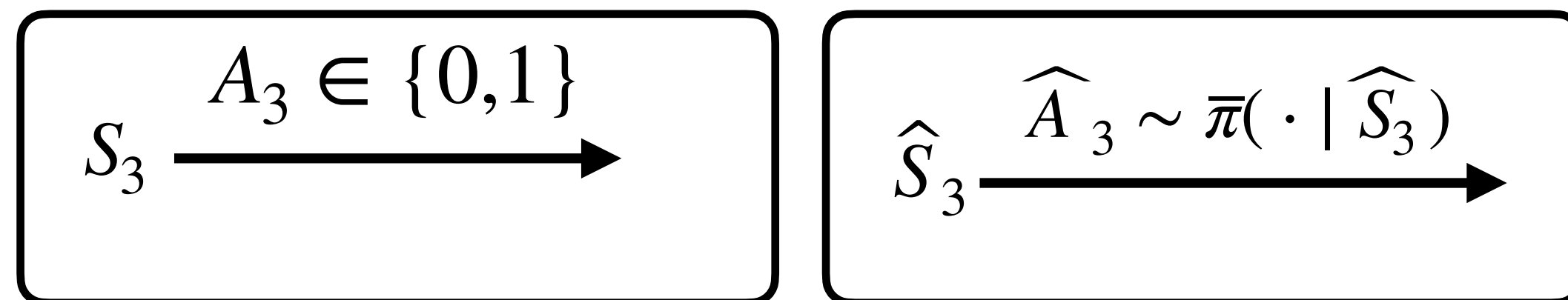
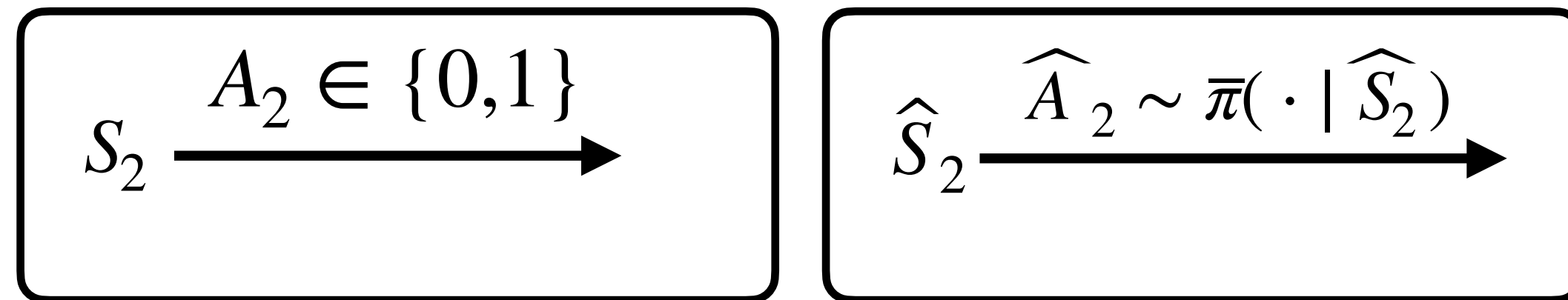
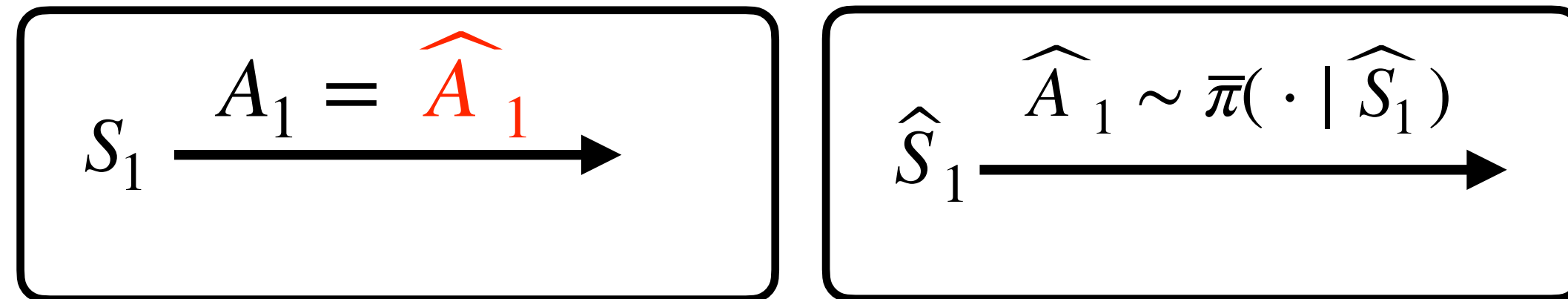
single-armed problem



single-armed policy  $\bar{\pi}(a | s)$

generate ideal actions

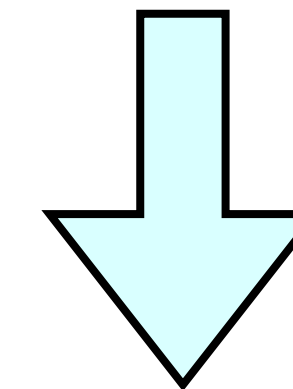
# Our policy: Follow-the-Virtual-Advice (FTVA)



no constraint

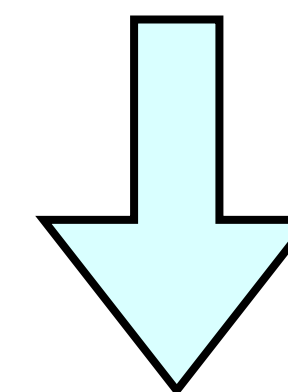
$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

$$\text{s.t. } \sum_{i=1}^N A_i = \alpha N, \text{ any time slot}$$



relax

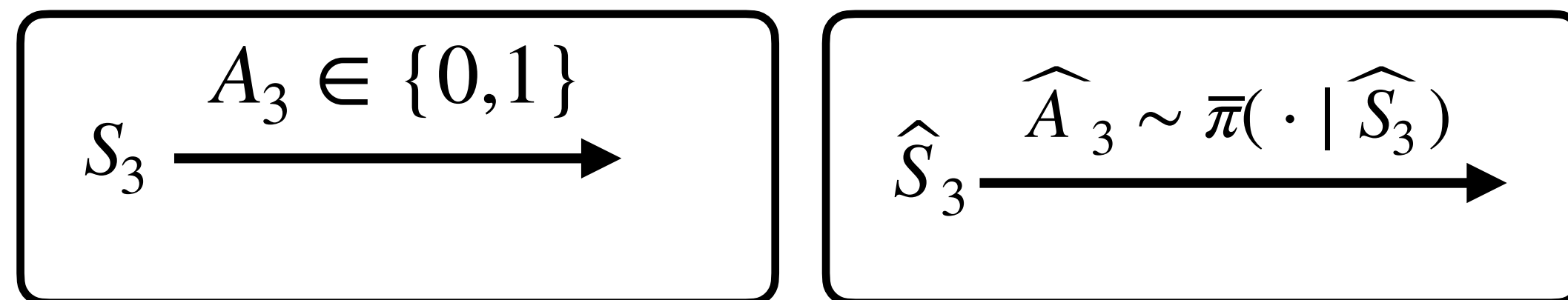
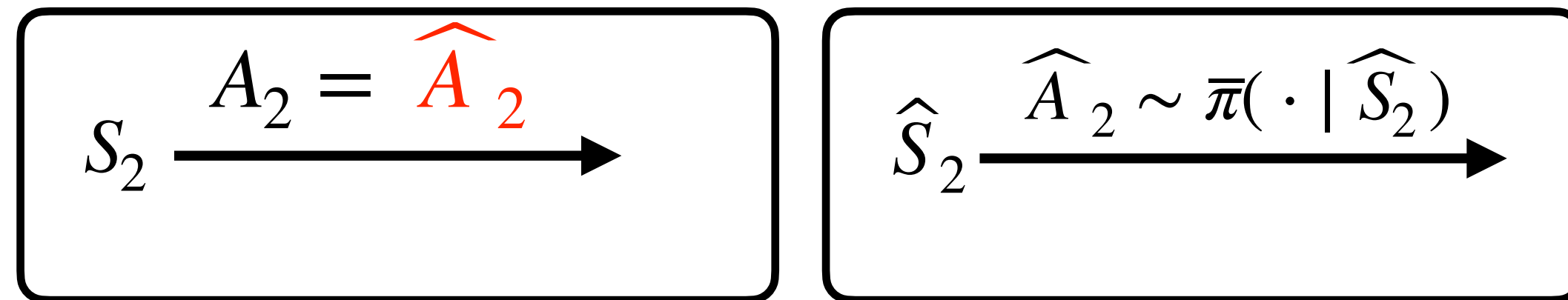
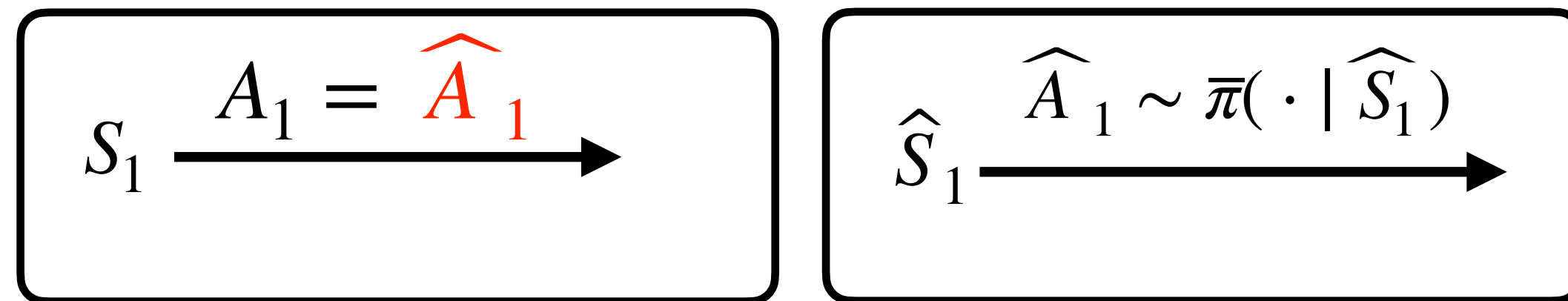
single-armed problem



single-armed policy  $\bar{\pi}(a | s)$

generate ideal actions

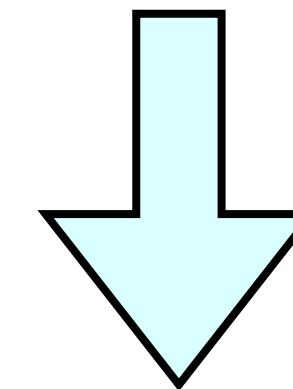
# Our policy: Follow-the-Virtual-Advice (FTVA)



no constraint

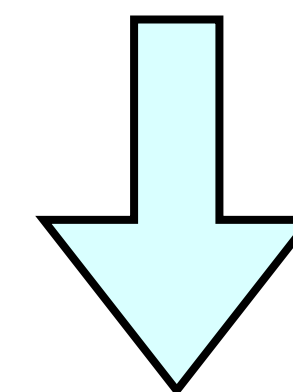
$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

$$\text{s.t. } \sum_{i=1}^N A_i = \alpha N, \text{ any time slot}$$



relax

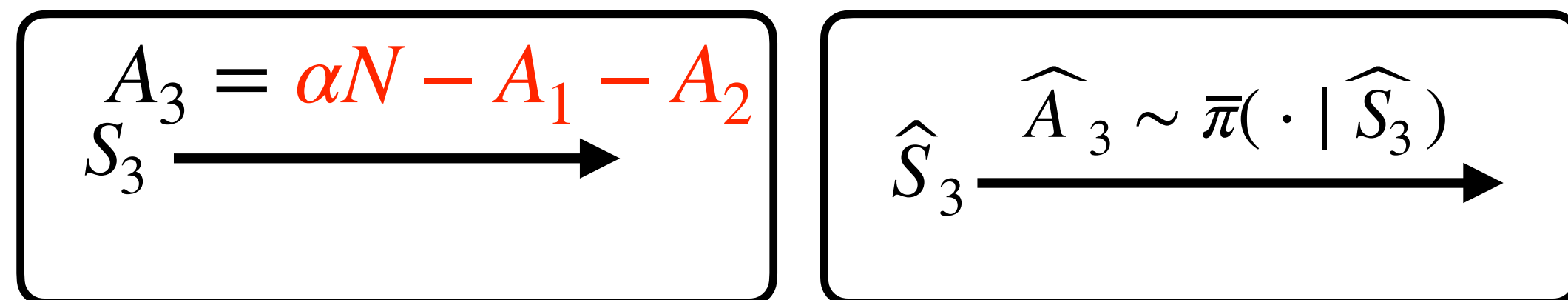
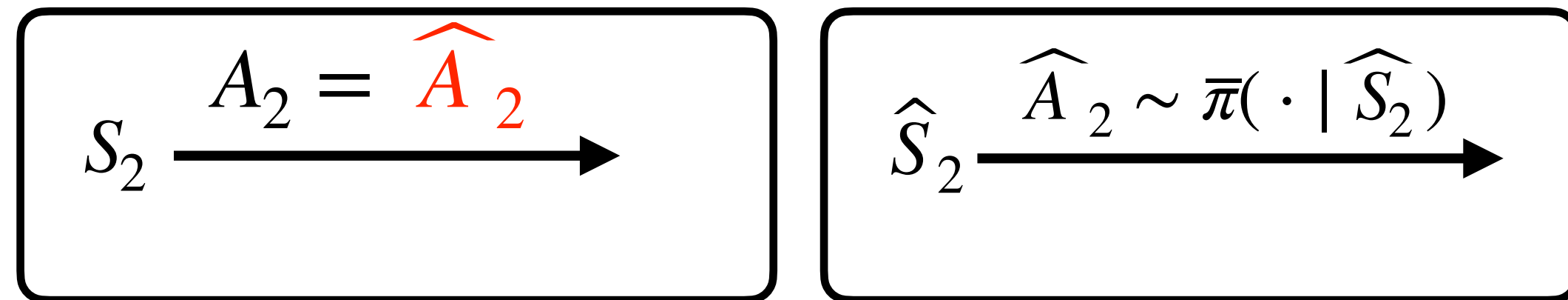
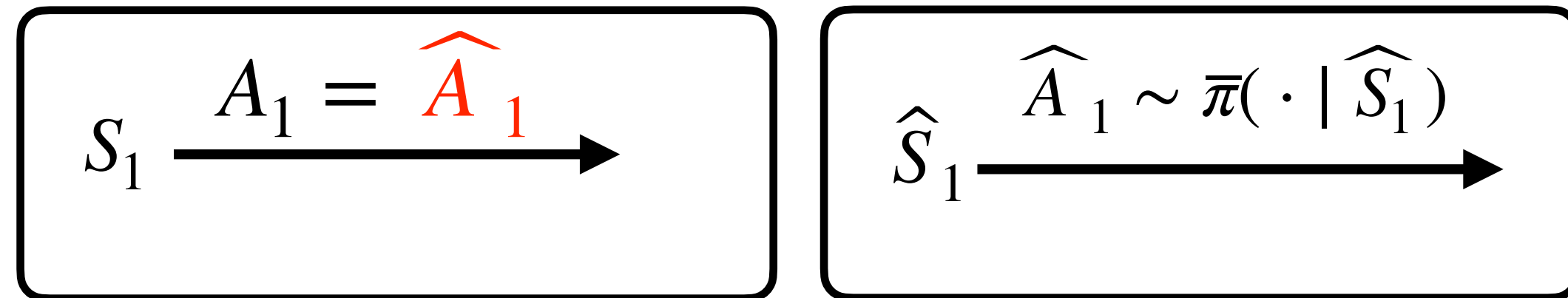
single-armed problem



single-armed policy  $\bar{\pi}(a | s)$

generate ideal actions

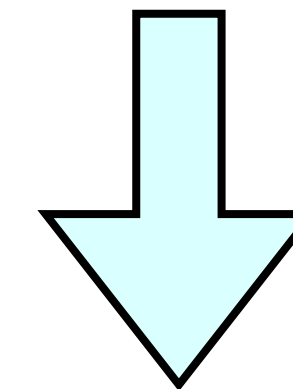
# Our policy: Follow-the-Virtual-Advice (FTVA)



no constraint

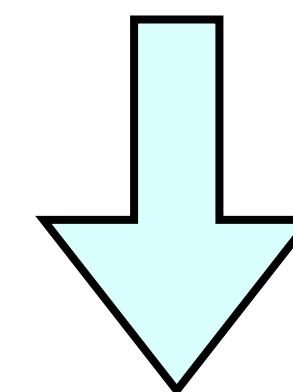
$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

$$\text{s.t. } \sum_{i=1}^N A_i = \alpha N, \text{ any time slot}$$



relax

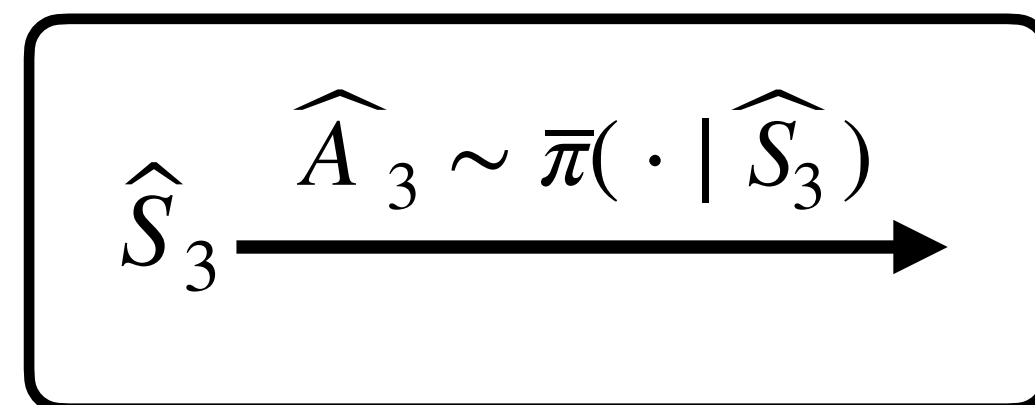
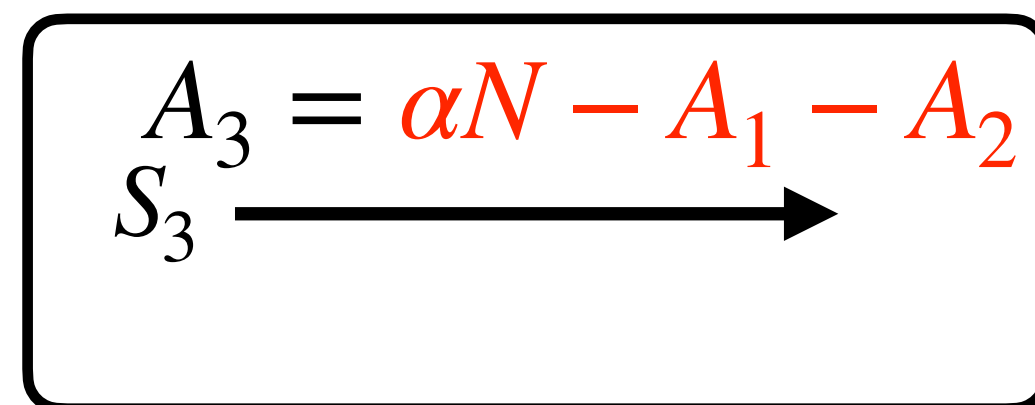
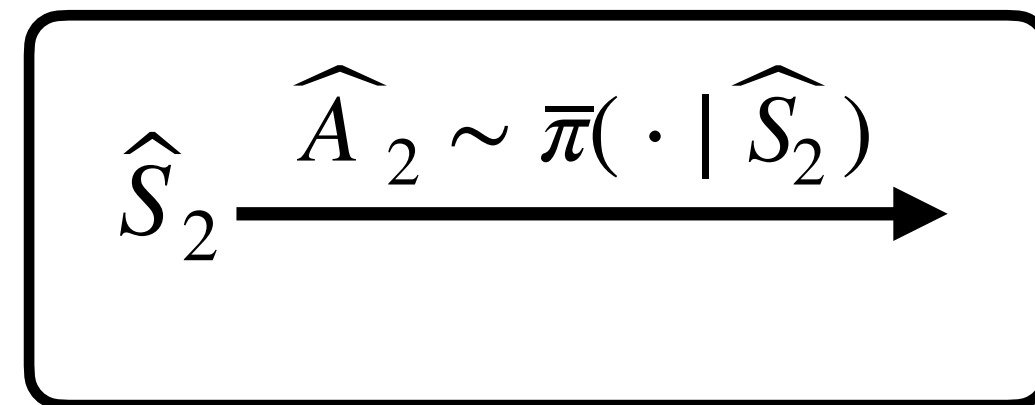
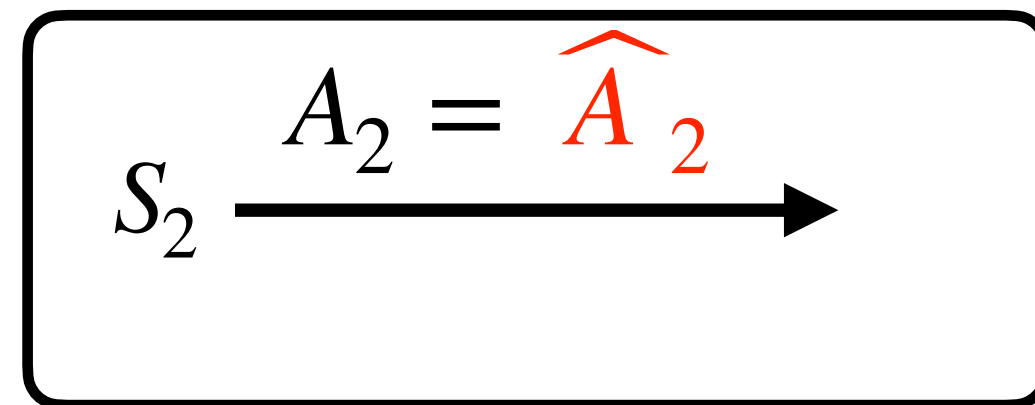
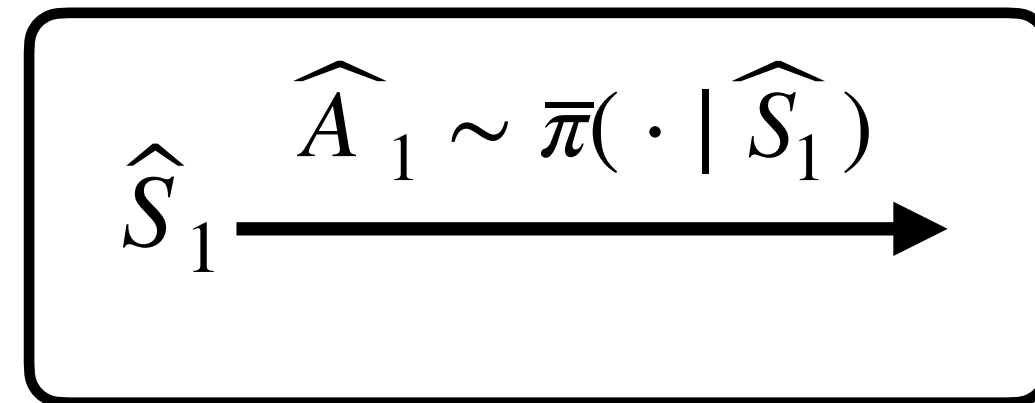
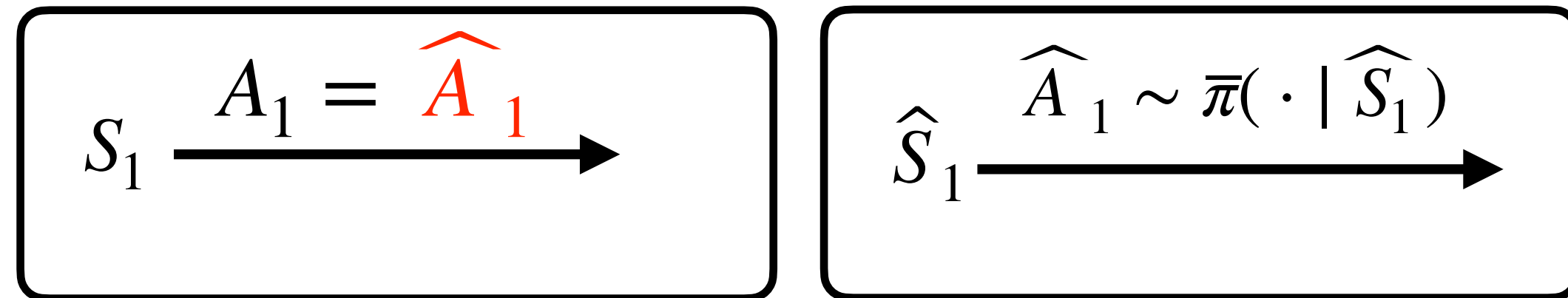
single-armed problem



single-armed policy  $\bar{\pi}(a | s)$

generate ideal actions

# Our policy: Follow-the-Virtual-Advice (FTVA)

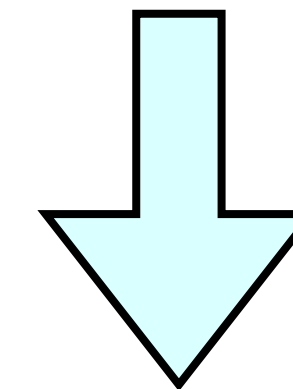


with constraint

no constraint

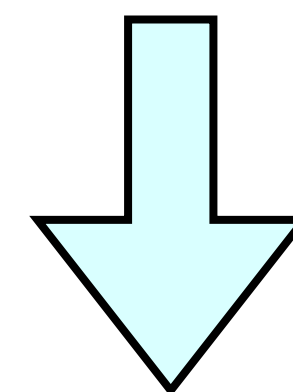
$$\max_{\pi} V_N^{\pi} \triangleq \text{long run average reward under policy } \pi$$

$$\text{s.t. } \sum_{i=1}^N A_i = \alpha N, \text{ any time slot}$$



relax

single-armed problem



single-armed policy  $\bar{\pi}(a | s)$

generate ideal actions