# Exact POMDP Solutions:

# $\alpha$-vectors

# Recap

# Recap

- POMDP

# Recap

- POMDP

$$(S, A, O, R, T, Z, \gamma)$$

# Recap

- POMDP
- Belief Updates

$$(S, A, O, R, T, Z, \gamma)$$

# Recap

- POMDP
- Belief Updates

$$(S, A, O, R, T, Z, \gamma)$$

$$b_t(s) = P(s_t = s \mid h_t)$$

# Recap

- POMDP
- Belief Updates

$$(S, A, O, R, T, Z, \gamma)$$

$$b_t(s) = P(s_t = s \mid h_t)$$

$$b' = \tau(b, a, o)$$

# Recap

- POMDP
- Belief Updates

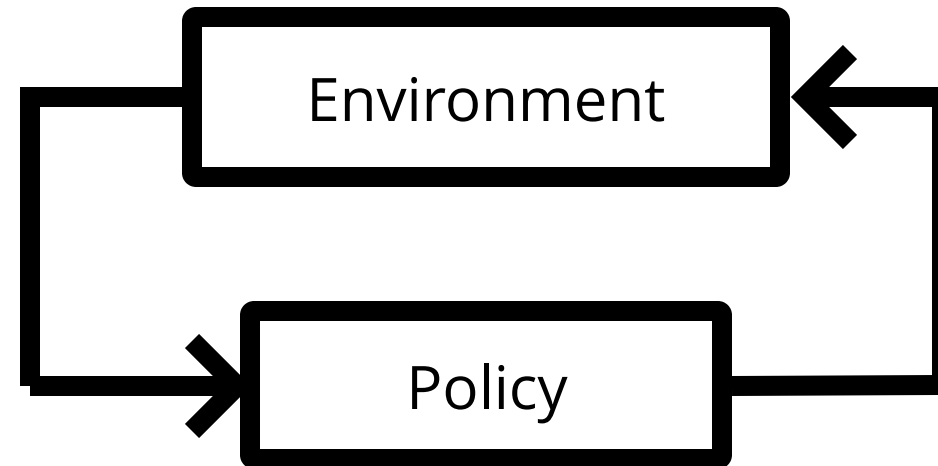$$(S, A, O, R, T, Z, \gamma)$$
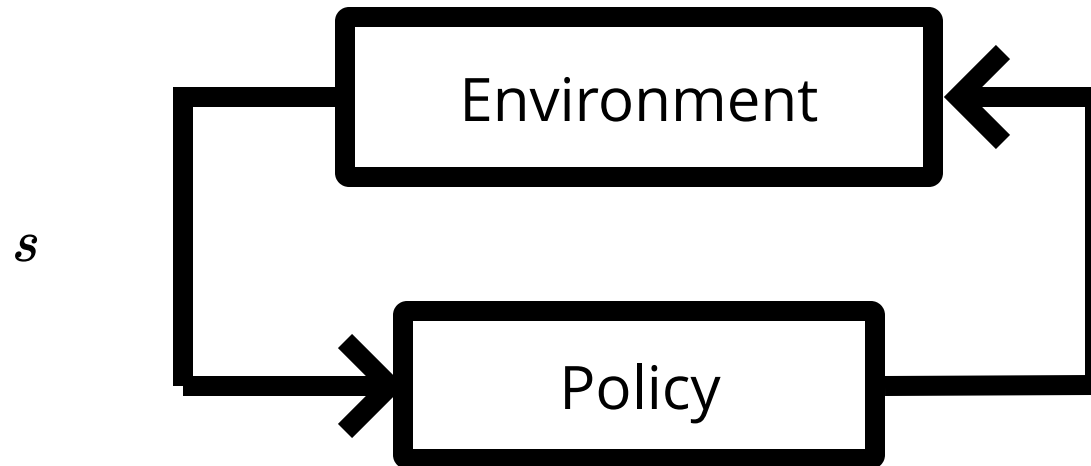
$$b_t(s) = P(s_t = s \mid h_t)$$

$$b' = \tau(b, a, o)$$

$$b'(s') \propto Z(o \mid a, s') \sum_s T(s' \mid s, a) \, b(s)$$

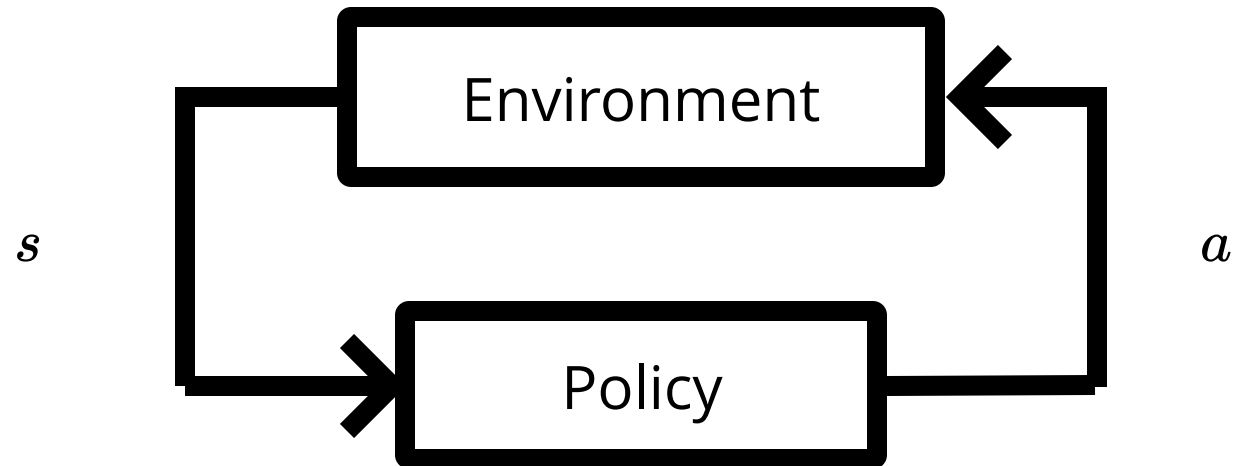# MDP Sense-Plan-Act Loop
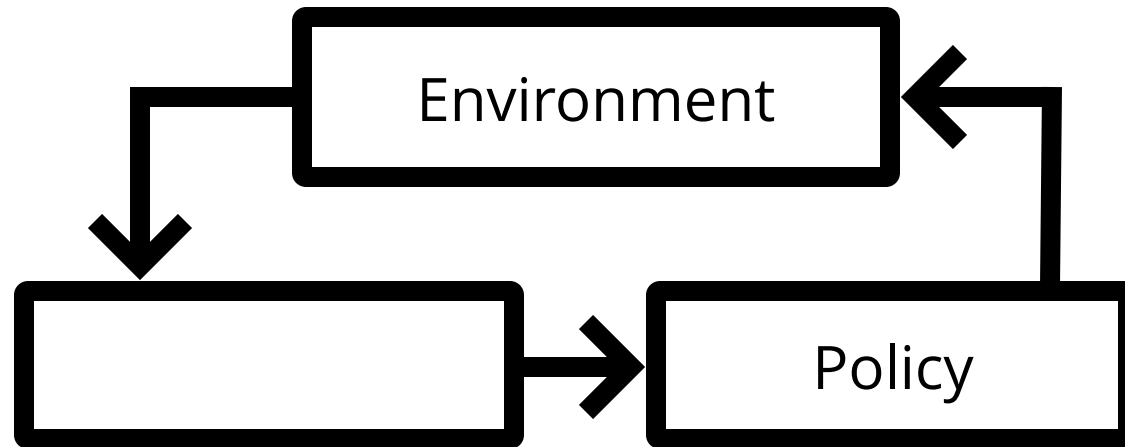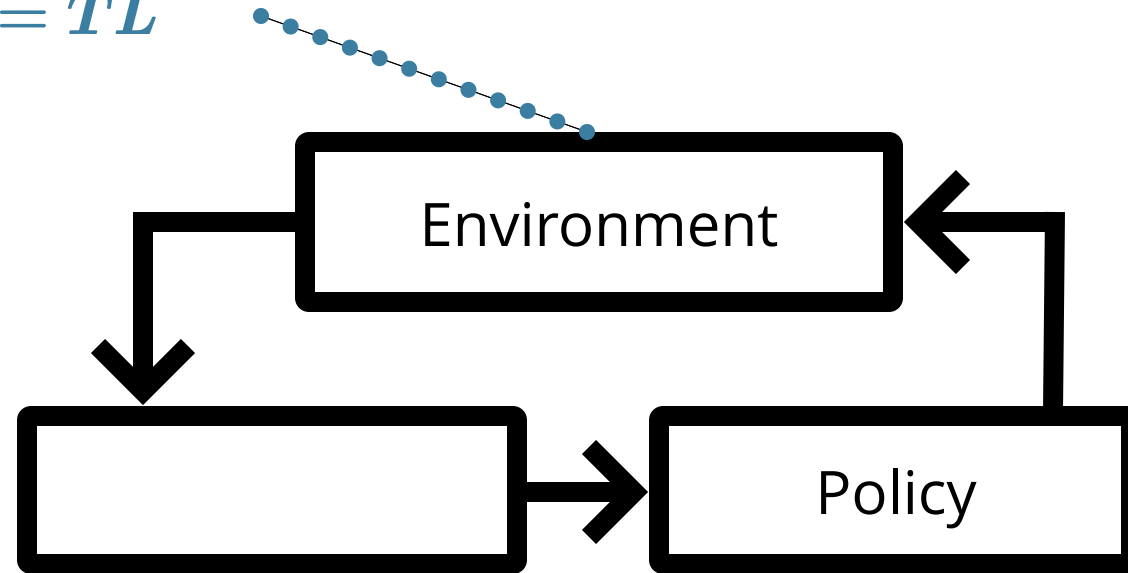
# MDP Sense-Plan-Act Loop

# MDP Sense-Plan-Act Loop



$s$

$a$

# POMDP Sense-Plan-Act Loop

# POMDP Sense-Plan-Act Loop

**True State**

$$s = TL$$

Environment

Policy

# POMDP Sense-Plan-Act Loop

**True State**

$$s = TL$$

**Observation**

$$o = TL$$

Environment

Policy

# POMDP Sense-Plan-Act Loop

**True State**

$$s = TL$$

**Observation**

$$o = TL$$

Environment

(Options below)

Policy

# POMDP Sense-Plan-Act Loop

**True State**

$$s = TL$$

**Environment**

**Observation**

$$o = TL$$

**(Options below)** → **Policy**

**Option 1: History**

# POMDP Sense-Plan-Act Loop

**True State**

$$s = TL$$

**Environment**

$a$

**Observation**

$$o = TL$$

(Options below) → Policy

**Option 1: History**   $h$

**History:** $h_t =$

$$(b_0, a_0, o_1, a_1, \ldots a_{t-1}, o_t)$$

# POMDP Sense-Plan-Act Loop

**True State**

$$s = TL$$

Environment

$a$

**Observation**

$$o = TL$$

(Options below)

Policy

**Option 1: History**

$h$

**Option 2: Belief Updater**

**History:** $h_t =$

$$(b_0, a_0, o_1, a_1, \ldots a_{t-1}, o_t)$$

# POMDP Sense-Plan-Act Loop

**True State**

$$s = TL$$

Environment

$a$

**Observation**

$$o = TL$$

(Options below) → Policy

**Option 1: History**

$h$  $b$

**Option 2: Belief Updater**

History: $h_t =$

**Belief:** $b_t = P(s_t \mid h_t)$

$$(b_0, a_0, o_1, a_1, \ldots a_{t-1}, o_t)$$

$TL$  $TR$

# Exercise 1: Crying Baby Belief Update

# Exercise 1: Crying Baby Belief Update

$S = \{h, \neg h\}$

$A = \{f, \neg f\}$

$O = \{c, \neg c\}$ .

# Exercise 1: Crying Baby Belief Update

$S = \{h, \neg h\}$

$A = \{f, \neg f\}$

$O = \{c, \neg c\}$

$R(s, a) = R(s) + R(a)$

# Exercise 1: Crying Baby Belief Update

$S = \{h, \neg h\}$

$A = \{f, \neg f\}$

$O = \{c, \neg c\}$

$R(s, a) = R(s) + R(a)$

$R(s) = \begin{cases} -10 \text{ if } s = h \\ 0 \text{ otherwise} \end{cases}$

# Exercise 1: Crying Baby Belief Update

$S = \{h, \neg h\}$

$A = \{f, \neg f\}$

$O = \{c, \neg c\}$

$R(s, a) = R(s) + R(a)$

$R(s) = \begin{cases} -10 \text{ if } s = h \\ 0 \text{ otherwise} \end{cases}$

$R(a) = \begin{cases} -5 \text{ if } a = f \\ 0 \text{ otherwise} \end{cases}$

# Exercise 1: Crying Baby Belief Update

$S = \{h, \neg h\}$        $T(h \mid h, \neg f) = 1.0$

$A = \{f, \neg f\}$        $T(h \mid \neg h, \neg f) = 0.1$

$O = \{c, \neg c\}$        $T(\neg h \mid \cdot, f) = 1.0$

$R(s, a) = R(s) + R(a)$

$$R(s) = \begin{cases} -10 \text{ if } s = h \\ 0 \text{ otherwise} \end{cases}$$

$$R(a) = \begin{cases} -5 \text{ if } a = f \\ 0 \text{ otherwise} \end{cases}$$

# Exercise 1: Crying Baby Belief Update

$S = \{h, \neg h\}$    $T(h \mid h, \neg f) = 1.0$

$A = \{f, \neg f\}$    $T(h \mid \neg h, \neg f) = 0.1$

$O = \{c, \neg c\}$    $T(\neg h \mid \cdot, f) = 1.0$

$$b'(s') \propto Z(o \mid a, s') \sum_s T(s' \mid s, a)\, b(s)$$

$b'(h) \propto \; {}^{z(}$

$R(s, a) = R(s) + R(a)$

$R(s) = \begin{cases} -10 \text{ if } s = h \\ 0 \text{ otherwise} \end{cases}$

$R(a) = \begin{cases} -5 \text{ if } a = f \\ 0 \text{ otherwise} \end{cases}$

$Z(c \mid \cdot, h) = 0.8)$

$Z(c \mid \cdot, \neg h) = 0.1$

# Exercise 1: Crying Baby Belief Update

$$S = \{h, \neg h\} \qquad T(h \mid h, \neg f) = 1.0$$

$$A = \{f, \neg f\} \qquad T(h \mid \neg h, \neg f) = 0.1$$

$$O = \{c, \neg c\} \qquad T(\neg h \mid \cdot, f) = 1.0$$

$$b'(s') \propto Z(o \mid a, s') \sum_s T(s' \mid s, a)\, b(s)$$

$$R(s, a) = R(s) + R(a)$$

$$R(s) = \begin{cases} -10 \text{ if } s = h \\ 0 \text{ otherwise} \end{cases}$$

$$R(a) = \begin{cases} -5 \text{ if } a = f \\ 0 \text{ otherwise} \end{cases}$$

$$Z(c \mid \cdot, h) = 0.8)$$

$$Z(c \mid \cdot, \neg h) = 0.1$$

$$\gamma = 0.9$$

# Exercise 1: Crying Baby Belief Update

$S = \{h, \neg h\}$

$A = \{f, \neg f\}$

$O = \{c, \neg c\}$

$T(h \mid h, \neg f) = 1.0$

$T(h \mid \neg h, \neg f) = 0.1$

$T(\neg h \mid \cdot, f) = 1.0$

$R(s, a) = R(s) + R(a)$

$$R(s) = \begin{cases} -10 \text{ if } s = h \\ 0 \text{ otherwise} \end{cases}$$

$$R(a) = \begin{cases} -5 \text{ if } a = f \\ 0 \text{ otherwise} \end{cases}$$

$Z(c \mid \cdot, h) = 0.8)$

$Z(c \mid \cdot, \neg h) = 0.1$

$\gamma = 0.9$

$$b'(s') \propto Z(o \mid a, s') \sum_s T(s' \mid s, a)\, b(s)$$

Starting at a $b(h) = 0$, calculate

$b'$ with $a = \neg f$ and $o = c$.

$b'(h) \propto \overset{Z(c \mid \neg f, h)}{0.8} \quad + \quad (\overset{T(h \mid \neg h, \neg f)}{0.1} \overset{b(\neg h)}{(1)} + \overset{T(h \mid h, \neg f)}{1.0} (0))$

$b'(h) \propto 0.08$

$b'(\neg h) \propto 0.09$

$b'(h) = 0.08 / (0.08 + 0.09) \simeq 47\%$
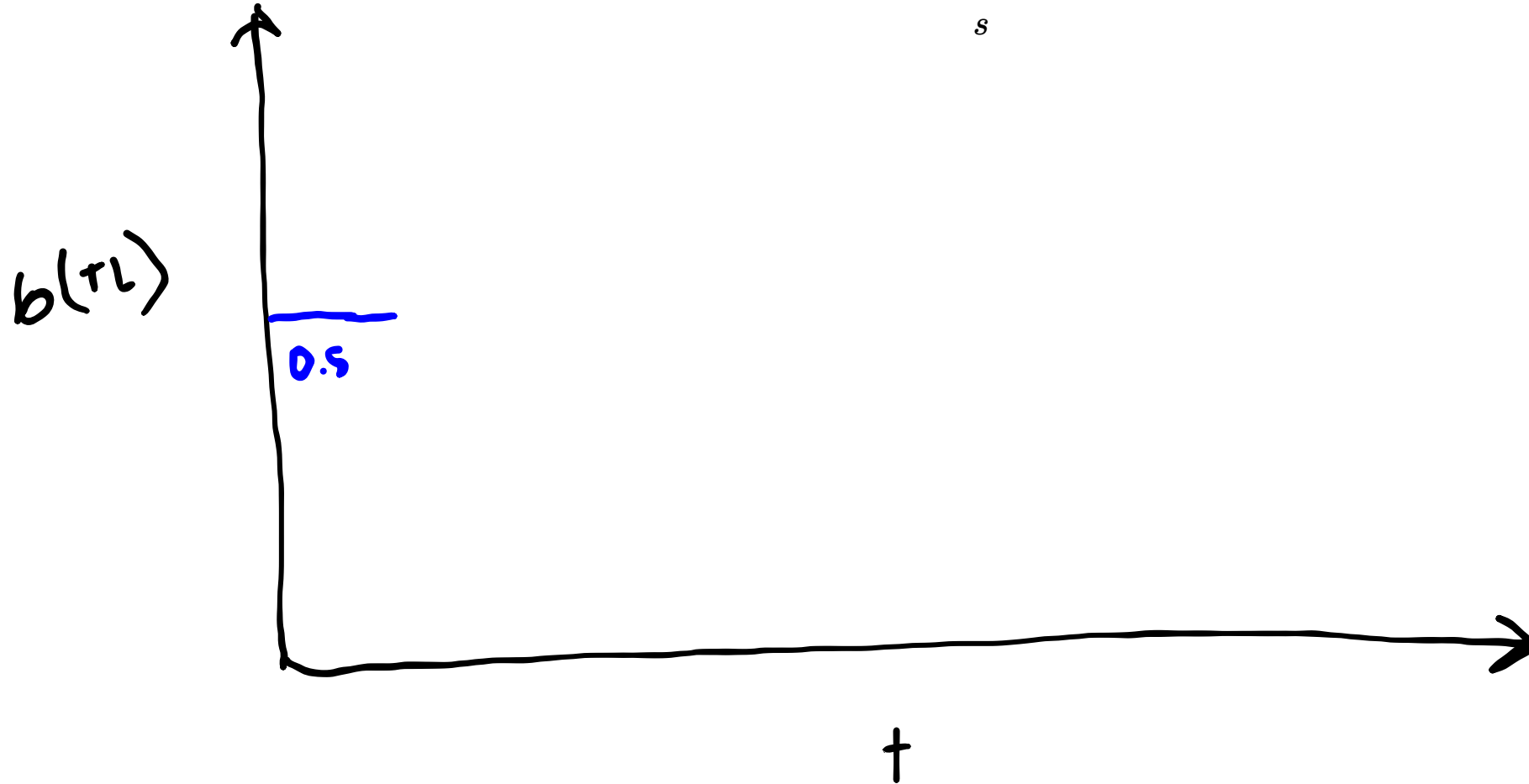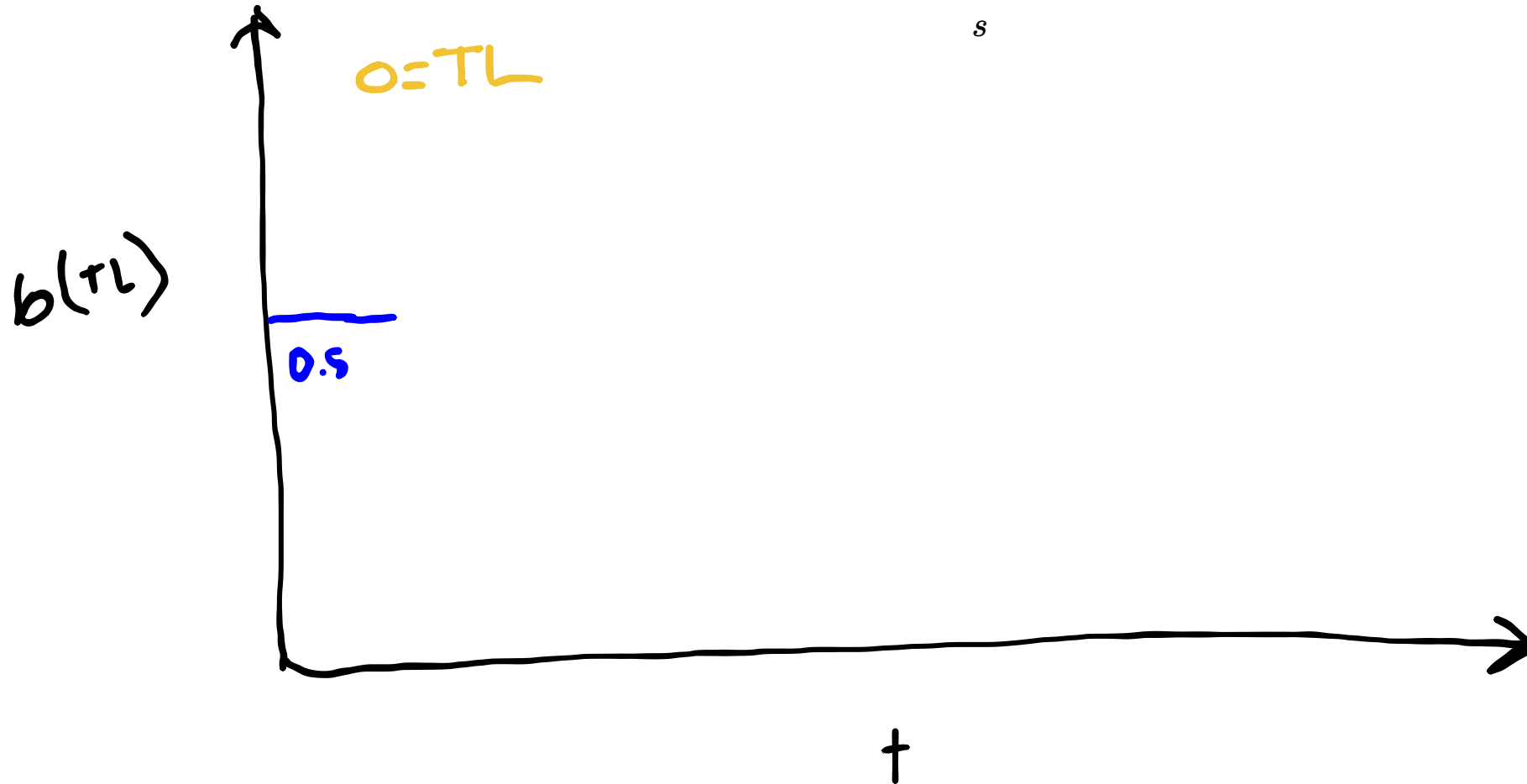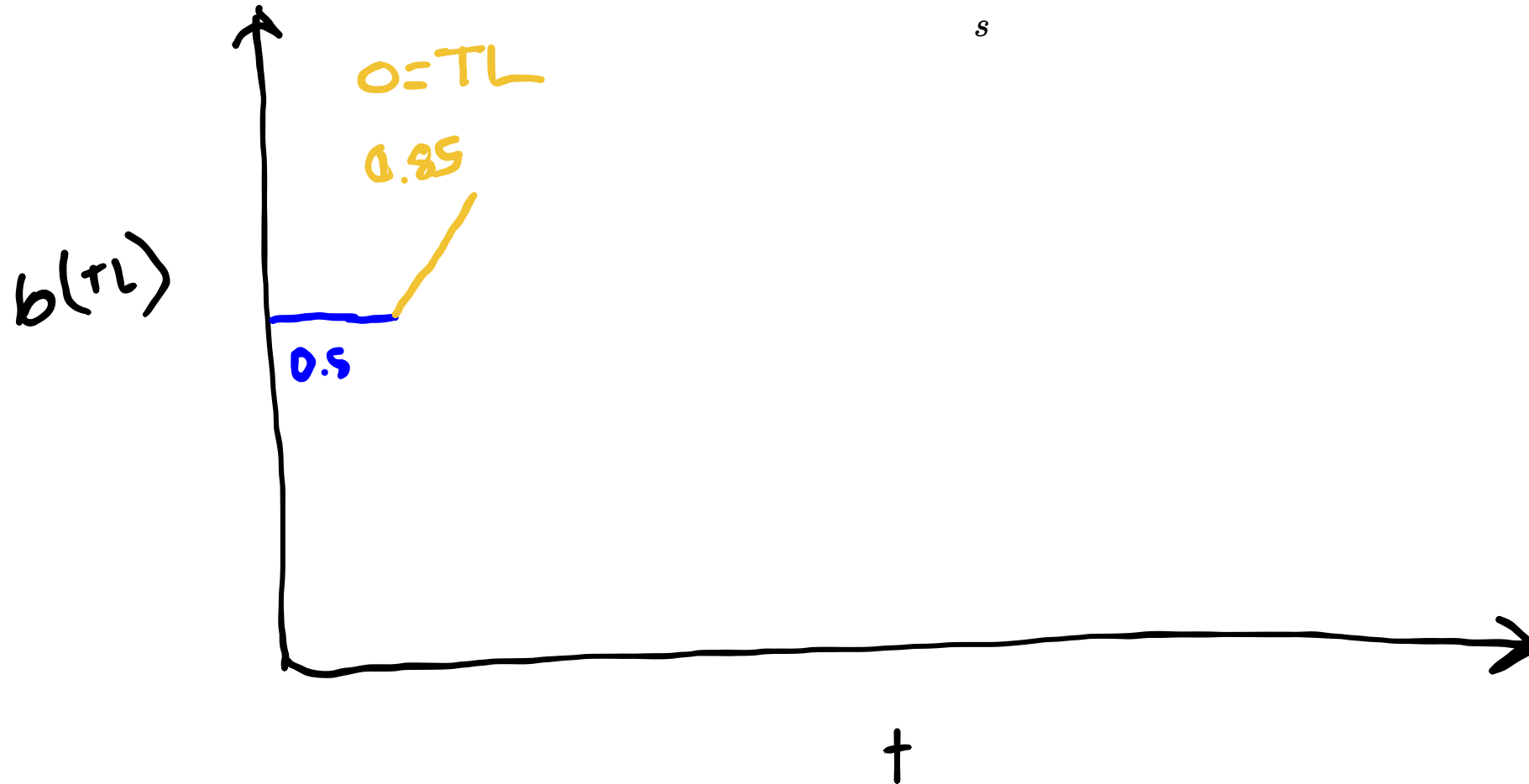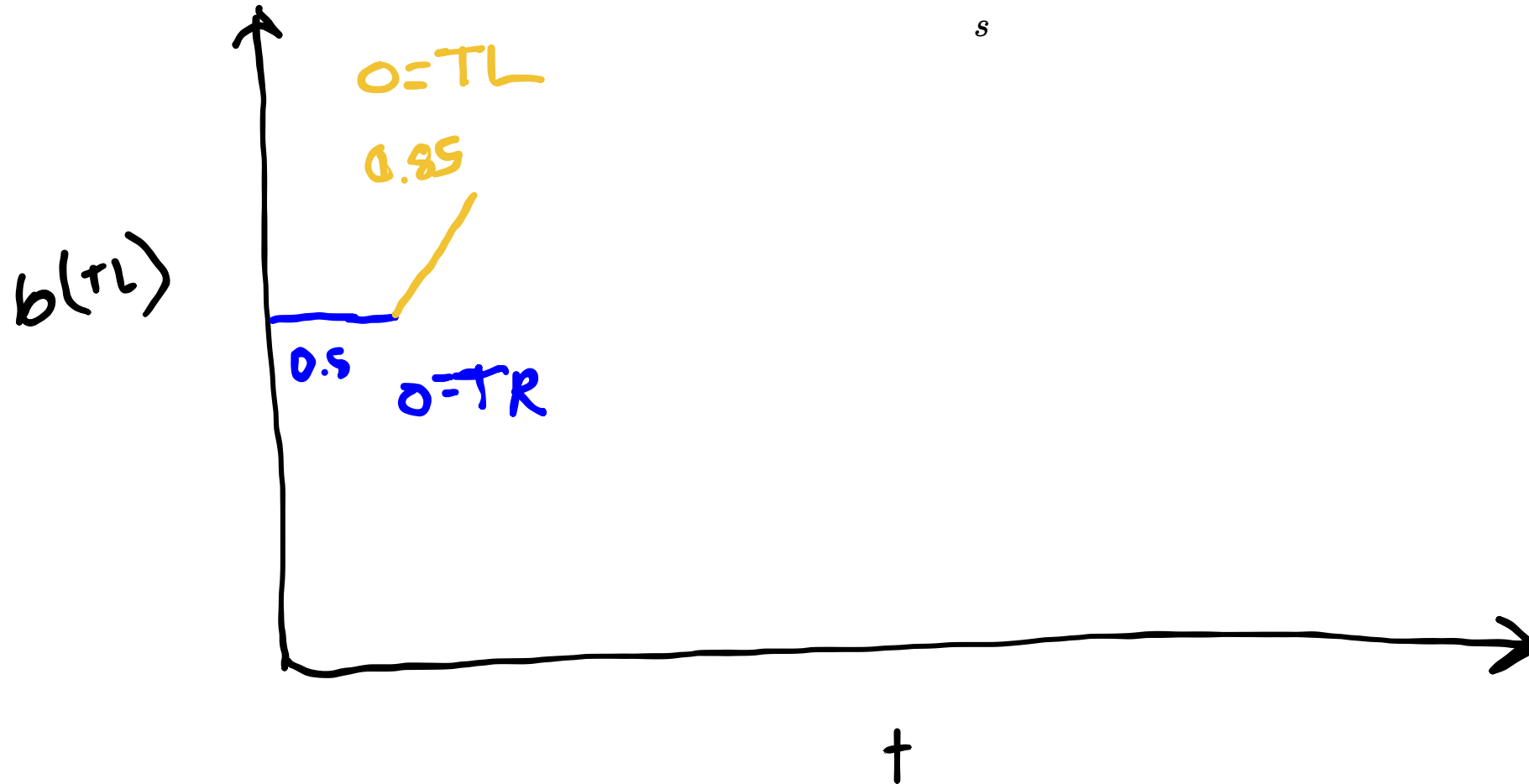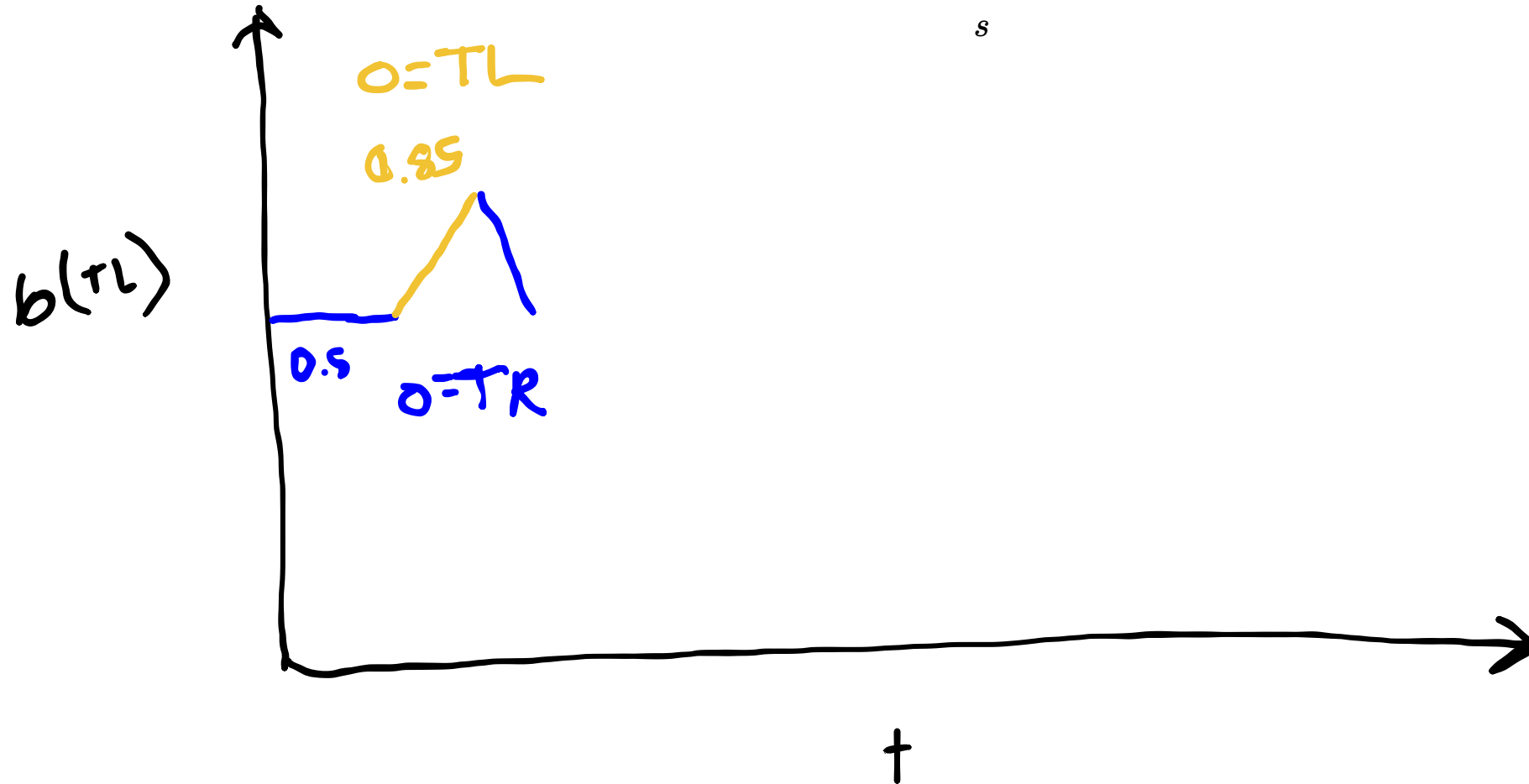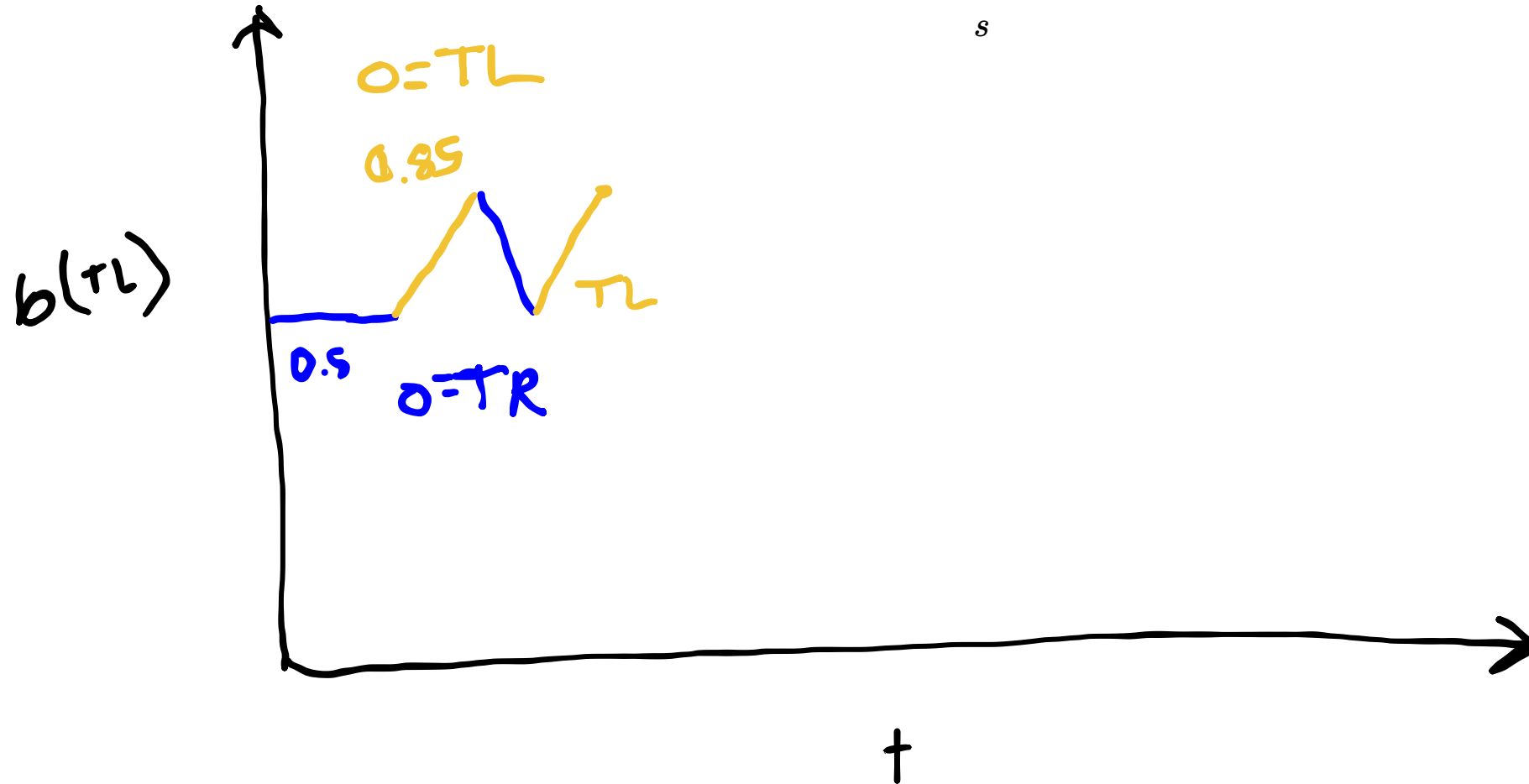
$b'(\neg h) = 53\%$

# Belief Dynamics

# Belief Dynamics

$$b'(s') \propto Z(o \mid a, s') \sum_s T(s' \mid s, a)\, b(s)$$

# Belief Dynamics

$$b'(s') \propto Z(o \mid a, s') \sum_s T(s' \mid s, a)\, b(s)$$
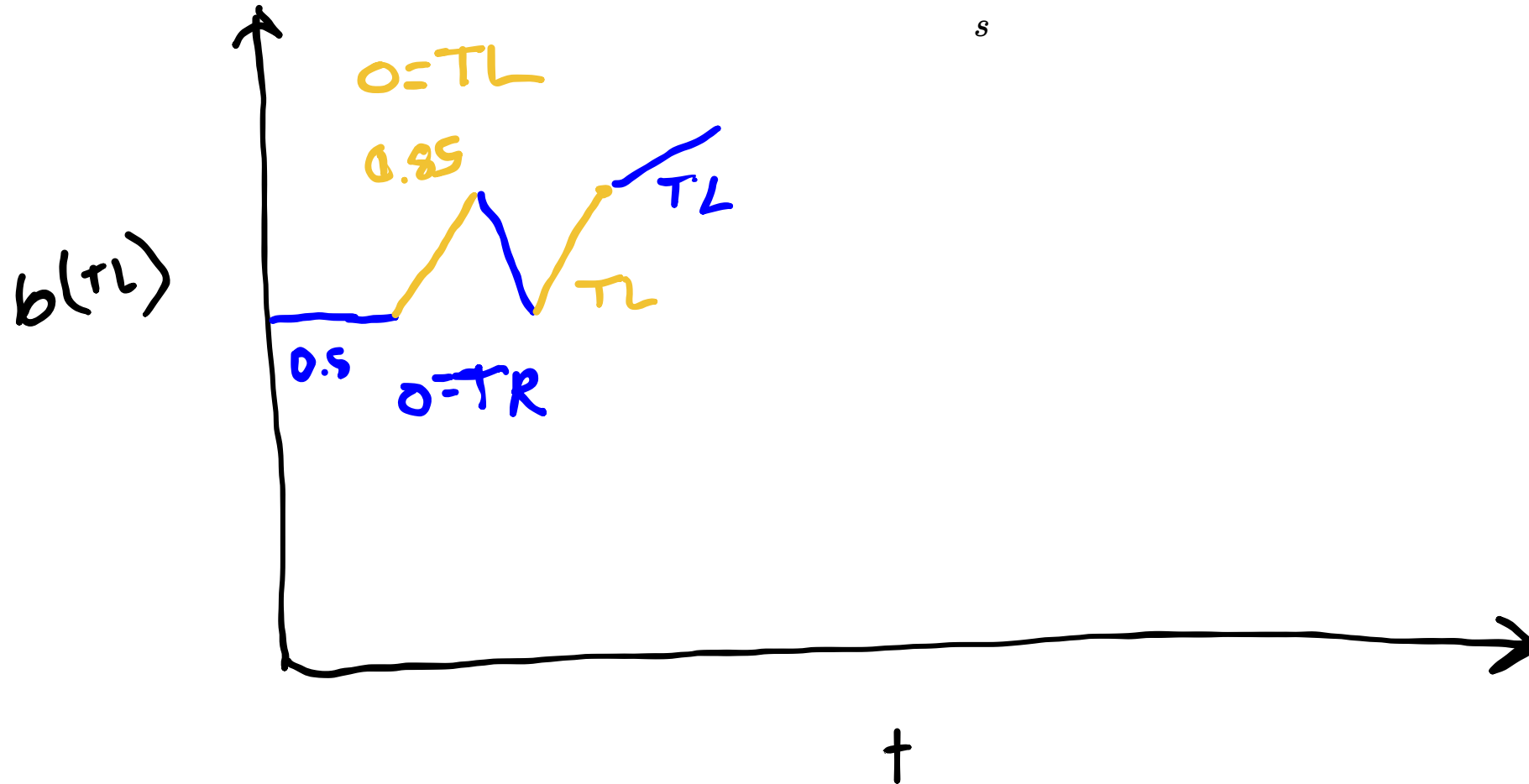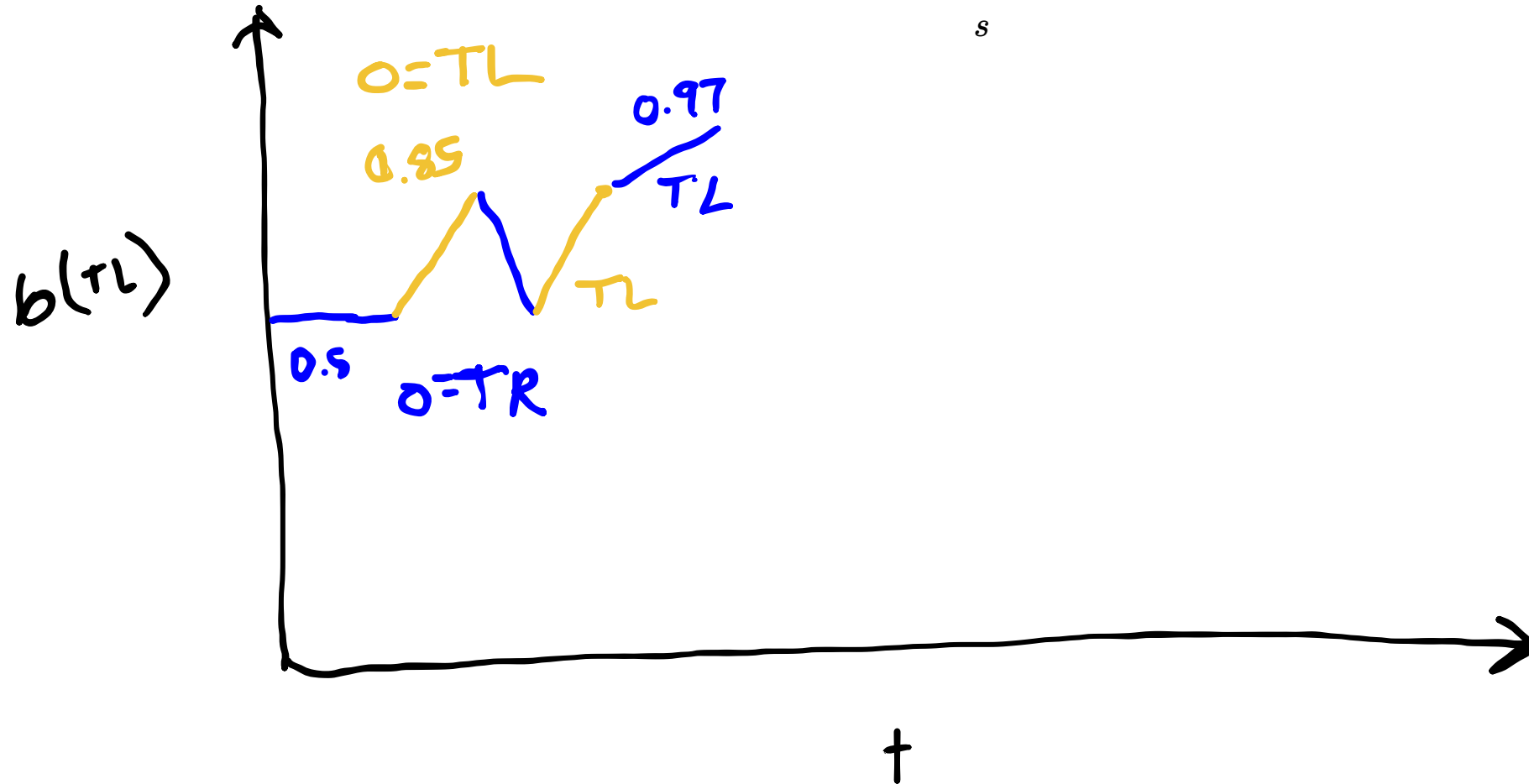
$b(TL)$

†

# Belief Dynamics

$$b'(s') \propto Z(o \mid a, s') \sum_s T(s' \mid s, a) \, b(s)$$

# Belief Dynamics

$$b'(s') \propto Z(o \mid a, s') \sum_s T(s' \mid s, a) \, b(s)$$

$b(TL)$

0.5

t

# Belief Dynamics

$$b'(s') \propto Z(o \mid a, s') \sum_s T(s' \mid s, a)\, b(s)$$



o=TL

$b(\tau_L)$

0.5

t

# Belief Dynamics

$$b'(s') \propto Z(o \mid a, s') \sum_s T(s' \mid s, a) \, b(s)$$

# Belief Dynamics

$$b'(s') \propto Z(o \mid a, s') \sum_s T(s' \mid s, a) \, b(s)$$
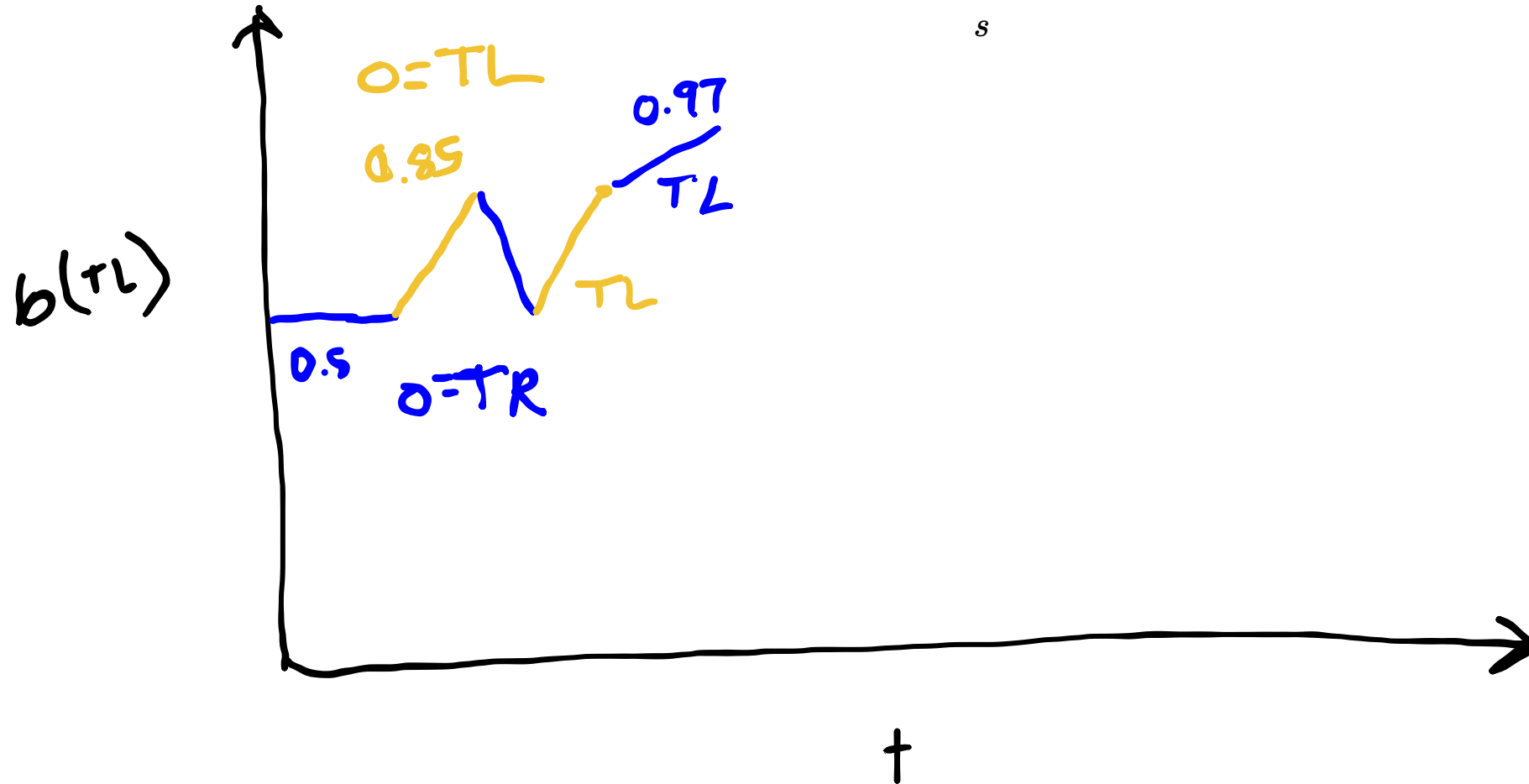
# Belief Dynamics

$$b'(s') \propto Z(o \mid a, s') \sum_s T(s' \mid s, a)\, b(s)$$

# Belief Dynamics

$$b'(s') \propto Z(o \mid a, s') \sum_s T(s' \mid s, a) \, b(s)$$
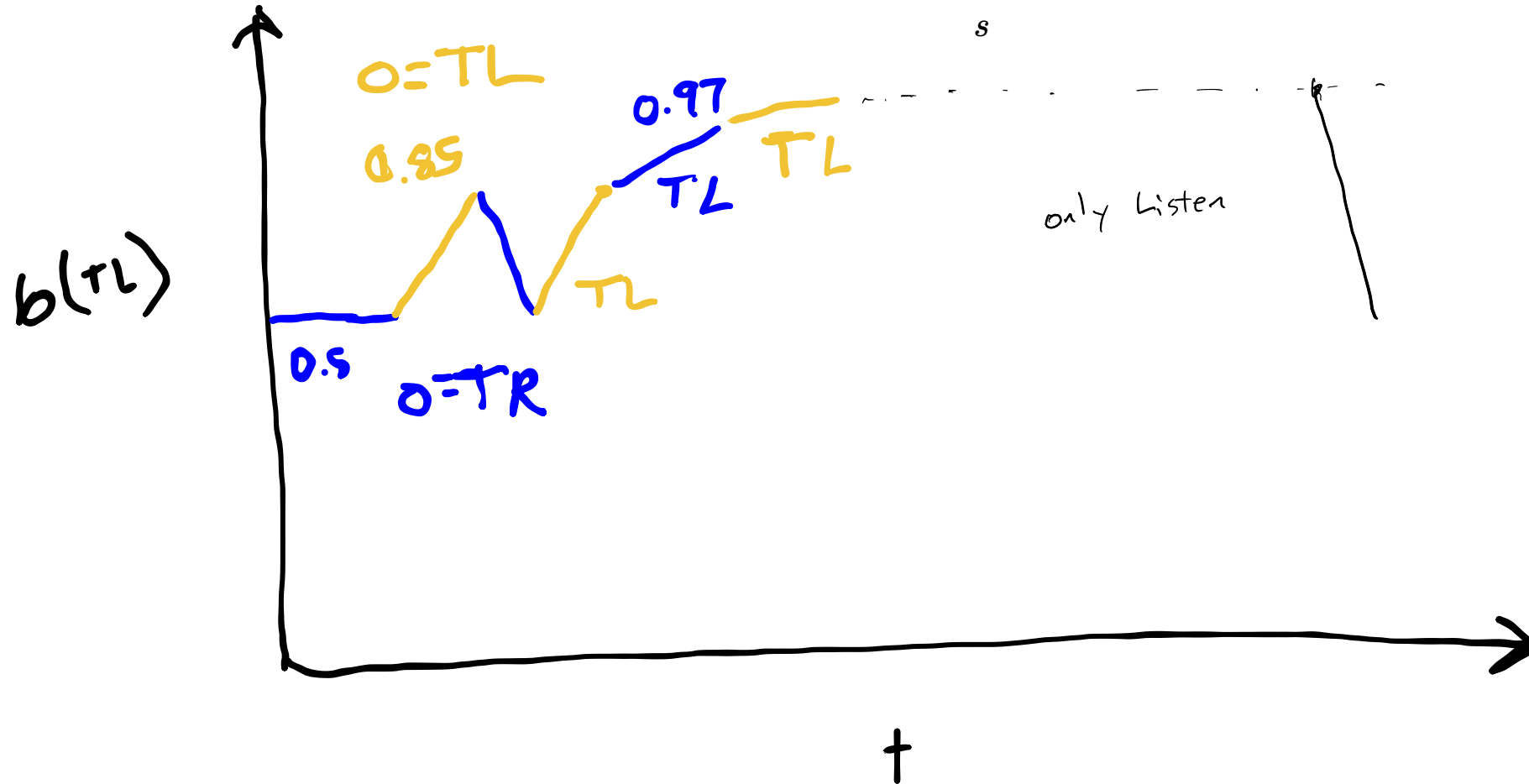
# Belief Dynamics

$$b'(s') \propto Z(o \mid a, s') \sum_s T(s' \mid s, a)\, b(s)$$

# Belief Dynamics
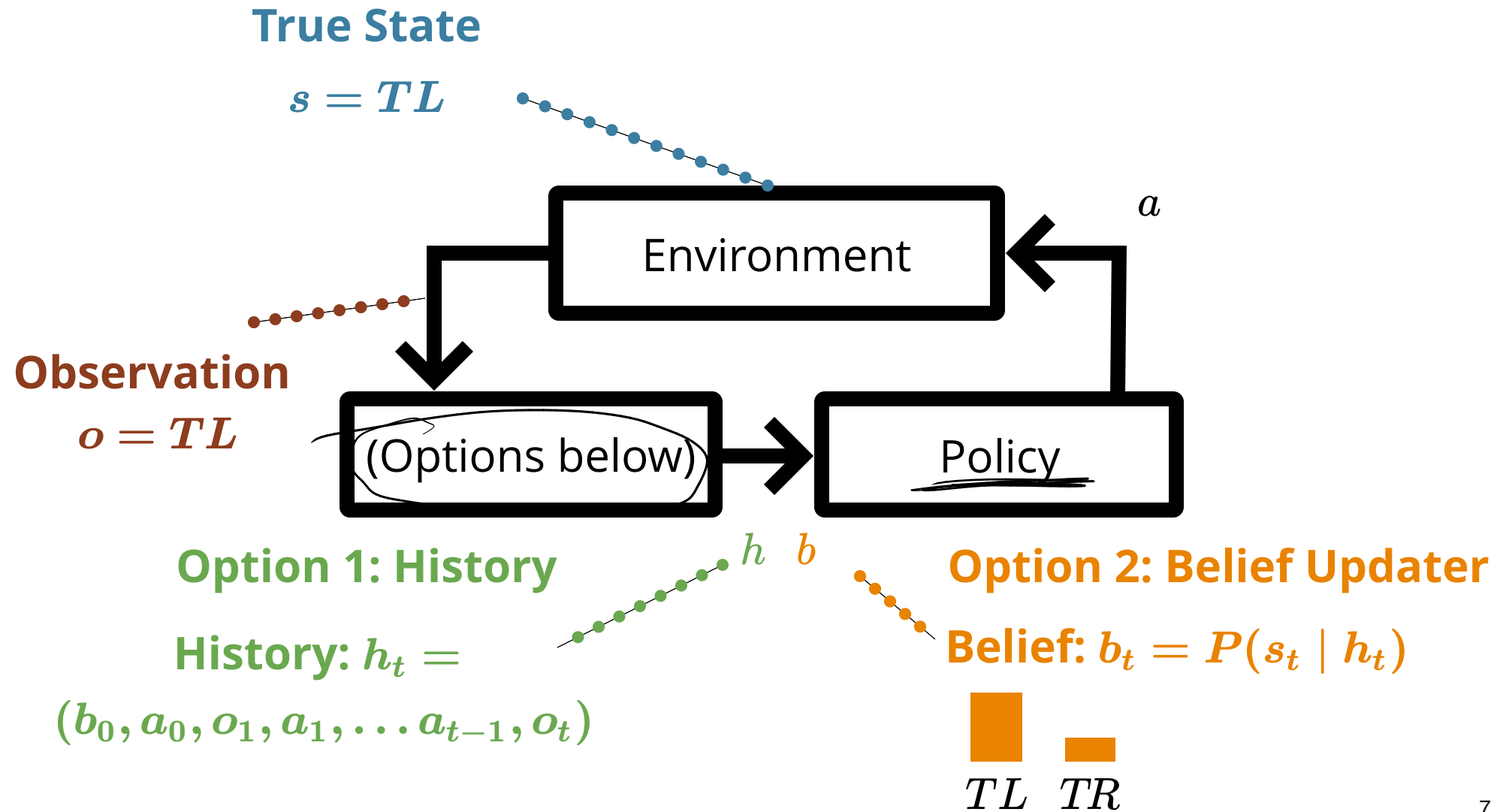
$$b'(s') \propto Z(o \mid a, s') \sum_s T(s' \mid s, a)\, b(s)$$

# Belief Dynamics

$$b'(s') \propto Z(o \mid a, s') \sum_s T(s' \mid s, a) \, b(s)$$

# Belief Dynamics

$$b'(s') \propto Z(o \mid a, s') \sum_s T(s' \mid s, a)\, b(s)$$
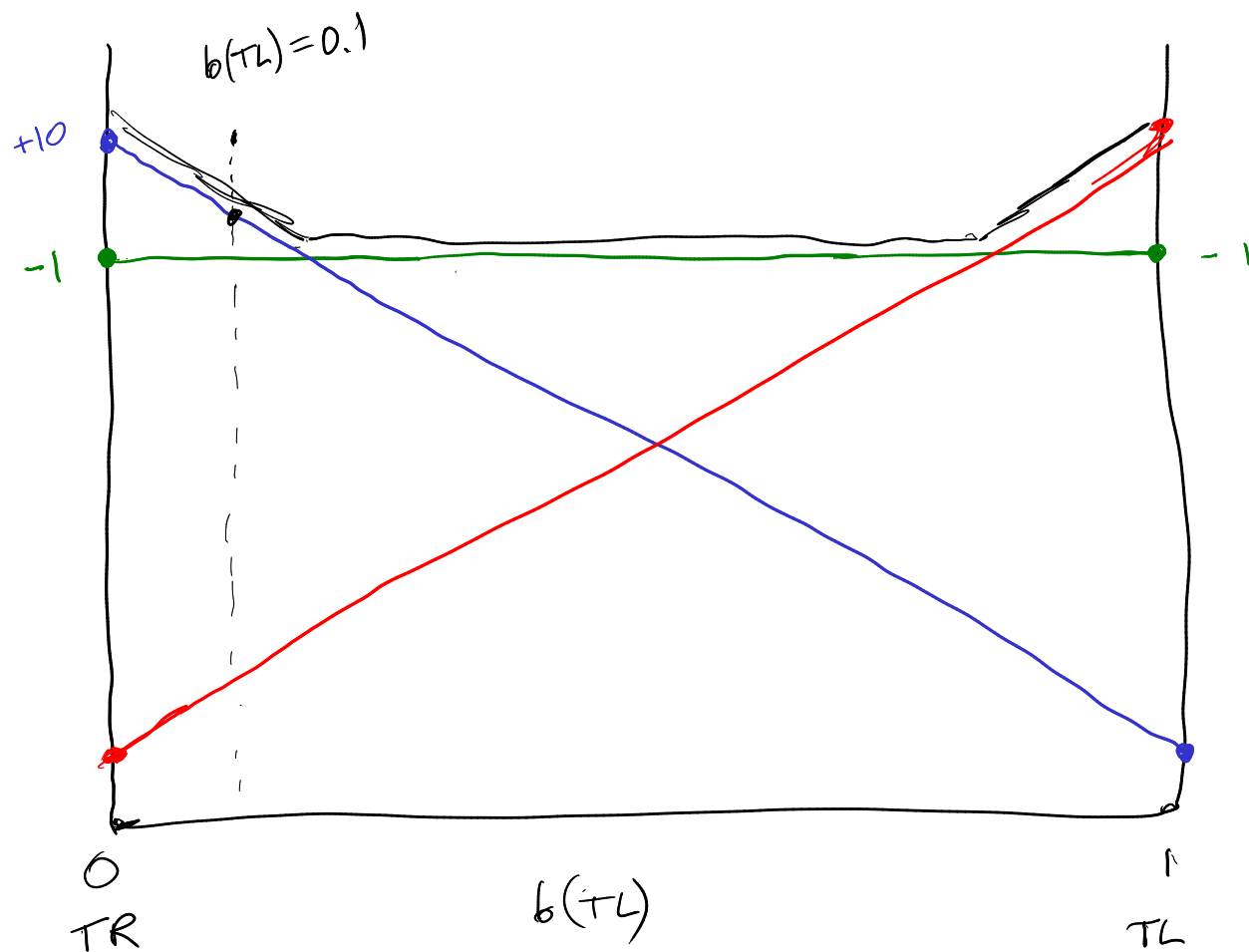
# POMDP Sense-Plan-Act Loop

**True State**

$$s = TL$$

Environment

$a$

**Observation**

$$o = TL$$

(Options below) → Policy

$h$  $b$

**Option 1: History**

**Option 2: Belief Updater**

**History:** $h_t =$

**Belief:** $b_t = P(s_t \mid h_t)$

$$(b_0, a_0, o_1, a_1, \ldots a_{t-1}, o_t)$$

$TL$  $TR$

# Guiding Quesiton

How do we calculate the optimal action in a POMDP?

# One-step utility

Reward
+10 empty
-1 Listen
-100 tiger

$[R(s',a), R(s^2,a)...]$
$[b(s'), b(s^2)...b(s^3)]$

$$R(b,a) = \bar{r}^a \cdot b$$

one-step
$\alpha$- vector

$b(TL) = 0.1$

+10

$a = L$

$R(TR, L) = -1$
$R(TL, L) = -1$

$a = OL$
$R(TR, OL) = +10$
$R(TL, OL) = -100$

$a = OR$
$R(TR, OR) = -100$
$R(TL, OR) = 10$

-1

-1

$O$
$TR$

$b(TL)$

$1$
$TL$

$$R(b,a) = \sum b(s) R(s,a) = b(TL) R(TL,a) + (1 - b(TL)) R(TR,a)$$

# One-step utility

Reward: +10 empty door
         -1  Listen
         -100 Tiger



$R(b,a)$
$= \underset{s \sim b}{E} [R(s,a)]$

$a = OL$      $a = L$      $a = OR$

$10$
$-1$

$-100$

$0$   $0.1$                    $b(TL)$                    $1$

$R(b,a) = \bar{r}_a \cdot b$

$\alpha$-vector

$R(b,a) = 10 \cdot (1 - b(TL)) - 100 \, b(TL) = -110 \, b + 10$

10

# Exercise 2: Crying Baby 1-Step Utility

$S = \{h, \neg h\}$

$A = \{f, \neg f\}$

$O = \{c, \neg c\}$

$T(h \mid h, \neg f) = 1.0$

$T(h \mid \neg h, \neg f) = 0.1$

$T(\neg h \mid \cdot, f) = 1.0$

Draw the 1-step utility $\alpha$-vectors for the Crying Baby problem.

$R(s, a) = R(s) + R(a)$

$R(s) = \begin{cases} -10 \text{ if } s = h \\ 0 \text{ otherwise} \end{cases}$

$R(a) = \begin{cases} -5 \text{ if } a = f \\ 0 \text{ otherwise} \end{cases}$

$Z(c \mid \cdot, h) = 0.8)$

$Z(c \mid \cdot, \neg h) = 0.1$

$\gamma = 0.9$



f    ¬f

0

-5

-10

R(h,f) - 15

b(h)

s = ¬h

s = h

# Alpha Vectors for Conditional Plans

# Alpha Vectors for Conditional Plans

Conditional Plans: fixed-depth history-based policies

# Alpha Vectors for Conditional Plans
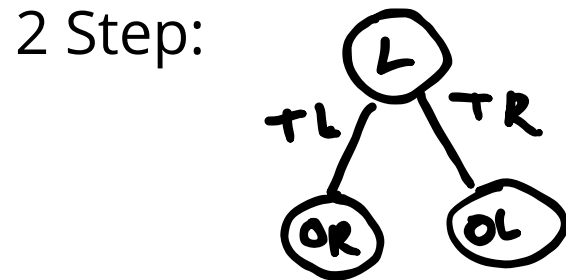
**Conditional Plans: fixed-depth history-based policies**

1 Step:

# Alpha Vectors for Conditional Plans

## Conditional Plans: fixed-depth history-based policies

1 Step:

# Alpha Vectors for Conditional Plans
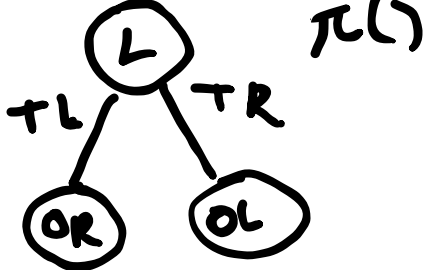
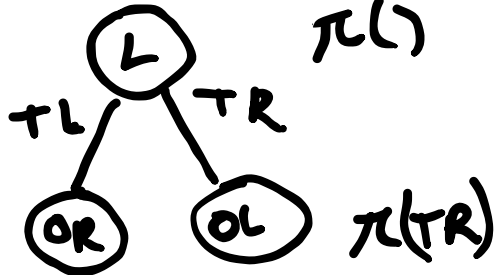**Conditional Plans: fixed-depth history-based policies**

1 Step:

2 Step:

# Alpha Vectors for Conditional Plans

## Conditional Plans: fixed-depth history-based policies

1 Step:

(L)  (OL)  (OR)

2 Step:

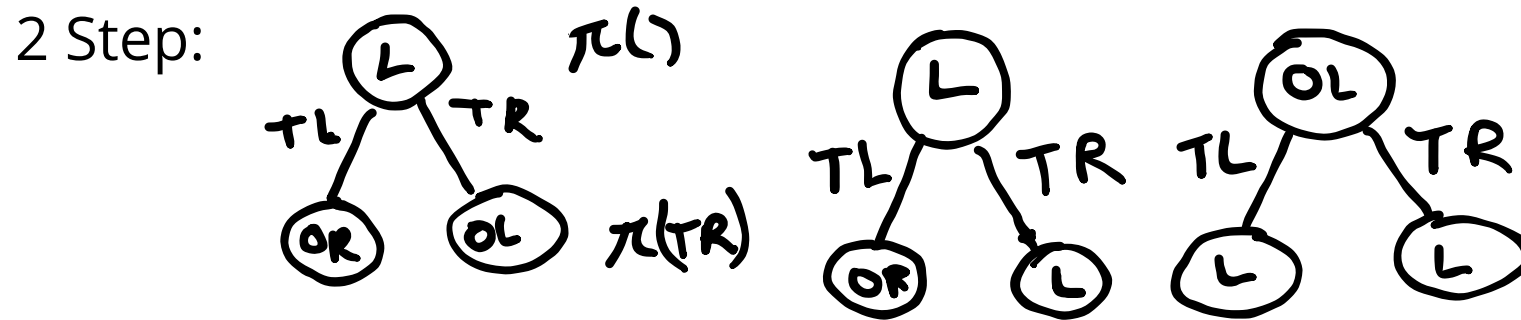# Alpha Vectors for Conditional Plans

Conditional Plans: fixed-depth history-based policies

1 Step:  (L)  (OL)  (OR)

2 Step:



$\pi()$

# Alpha Vectors for Conditional Plans

Conditional Plans: fixed-depth history-based policies

1 Step:

L    OL    OR

2 Step:

L    $\pi()$

TL    TR

OR    OL    $\pi(TR)$

# Alpha Vectors for Conditional Plans

**Conditional Plans: fixed-depth history-based policies**

1 Step: (L) (OL) (OR)

2 Step:

# Alpha Vectors for Conditional Plans

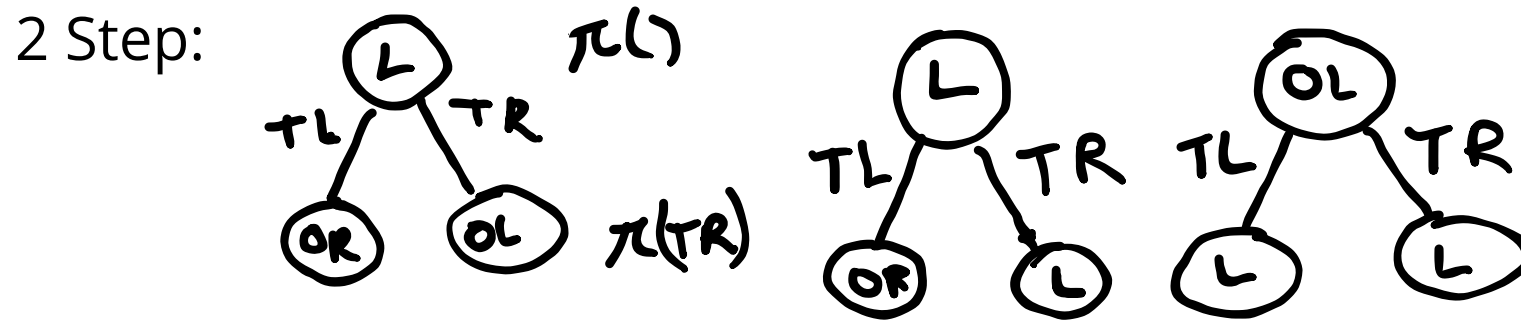Conditional Plans: fixed-depth history-based policies

1 Step:



2 Step:

# Alpha Vectors for Conditional Plans

## Conditional Plans: fixed-depth history-based policies

1 Step:



2 Step:



$$|A|^{\frac{(|O|^{h}-1)}{(|O|-1)}}$$

# Alpha Vectors for Conditional Plans

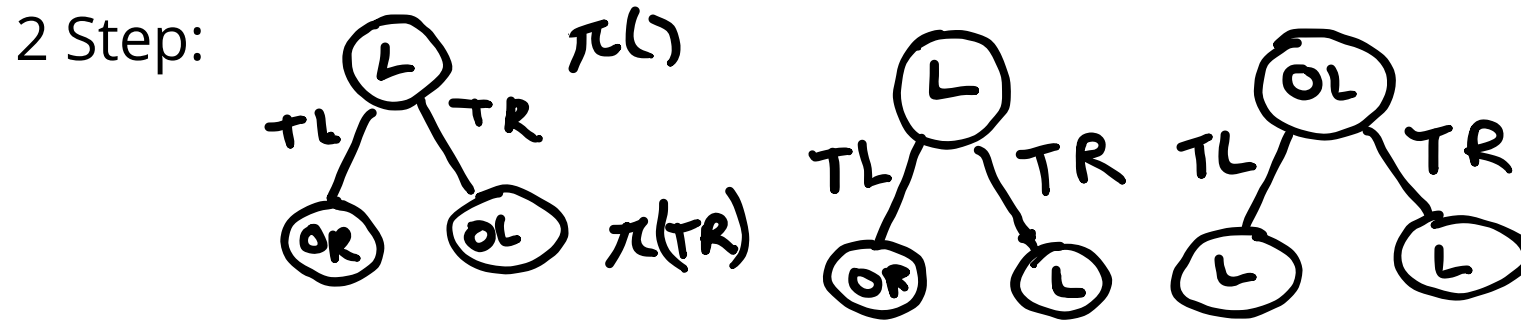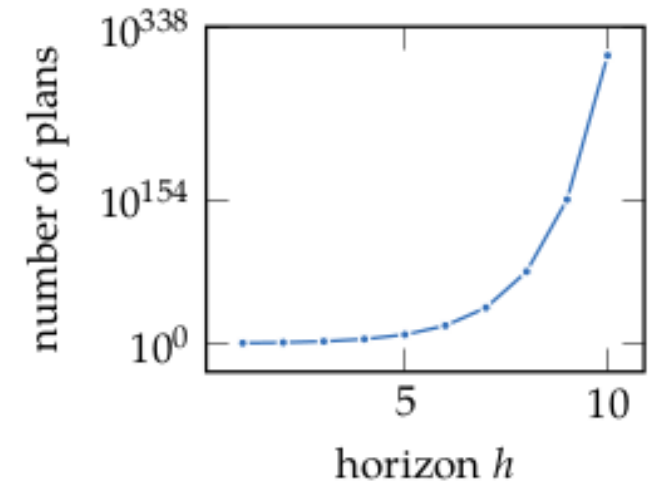## Conditional Plans: fixed-depth history-based policies

1 Step:



2 Step:



$$|A|^{\frac{(|O|^h - 1)}{(|O| - 1)}}$$
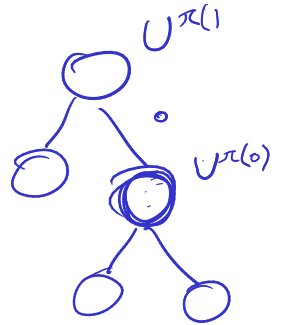
27 two step plans!

# Alpha Vectors for Conditional Plans

## Conditional Plans: fixed-depth history-based policies

1 Step:



2 Step:



$$|A|^{\frac{(|O|^h - 1)}{(|O| - 1)}}$$

27 two step plans!

# Alpha Vectors for Conditional Plans

# Alpha Vectors for Conditional Plans

For 1-step: $U^\pi(s) = R(s, \pi())$

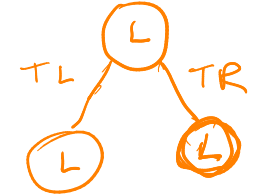# Alpha Vectors for Conditional Plans

For 1-step: $U^\pi(s) = R(s, \pi())$

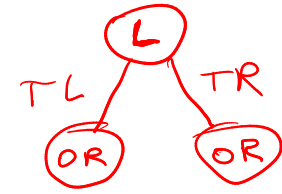$$U^\pi(s) = R(s, \pi()) + \gamma \left[ \sum_{s'} T(s' \mid s, \pi()) \sum_o \cancel{Z}(o \mid \pi(), s') U^{\pi(o)}(s') \right]$$



$U^{\pi(1)}$

$U^{\pi(0)}$

$\boxed{L}$ $\pi_L$

$U^{\pi_L()}(\cdot) = -1$

$\boxed{OL}$ $\pi_{OL}$

$U^{\pi_{OL}()}(TR) = 10$
$U^{\pi_{OL}()}(TL) = -100$

$\boxed{OR}$ $\pi_{OR}$

$U^{\pi_{OR}()}(TR) = -100$
$U^{\pi_{OR}()}(TL) = 10$

8.5

−1.95

−7.175

$b(TL)$

−96

$L$
TL — TR
$L$ $L$

$U^{\pi()}(\cdot) = -1 + \gamma(-1)$

$L$
TL — TR
$OR$ $OR$

$U^\pi(TL) = -1 + \gamma \, 10$
$= 8.5$

$U^\pi(TR) = -1 + \gamma(-100)$
$= -96$

$L$
TL — TR
$OR$ $OL$

$(T(TL|TL,L)=1)$
$(T(TR|TL,L)=0)$

$U^{\pi()}(TL) = -1 + \gamma \left( 0.85 \cdot 10 + 0.15 \cdot (-100) \right)$
$\uparrow Z(TL|L,TL) \uparrow U^{\pi(TL)}(TL)$
$= -7.175$

$U^{\pi()}(TR) = -7.175$
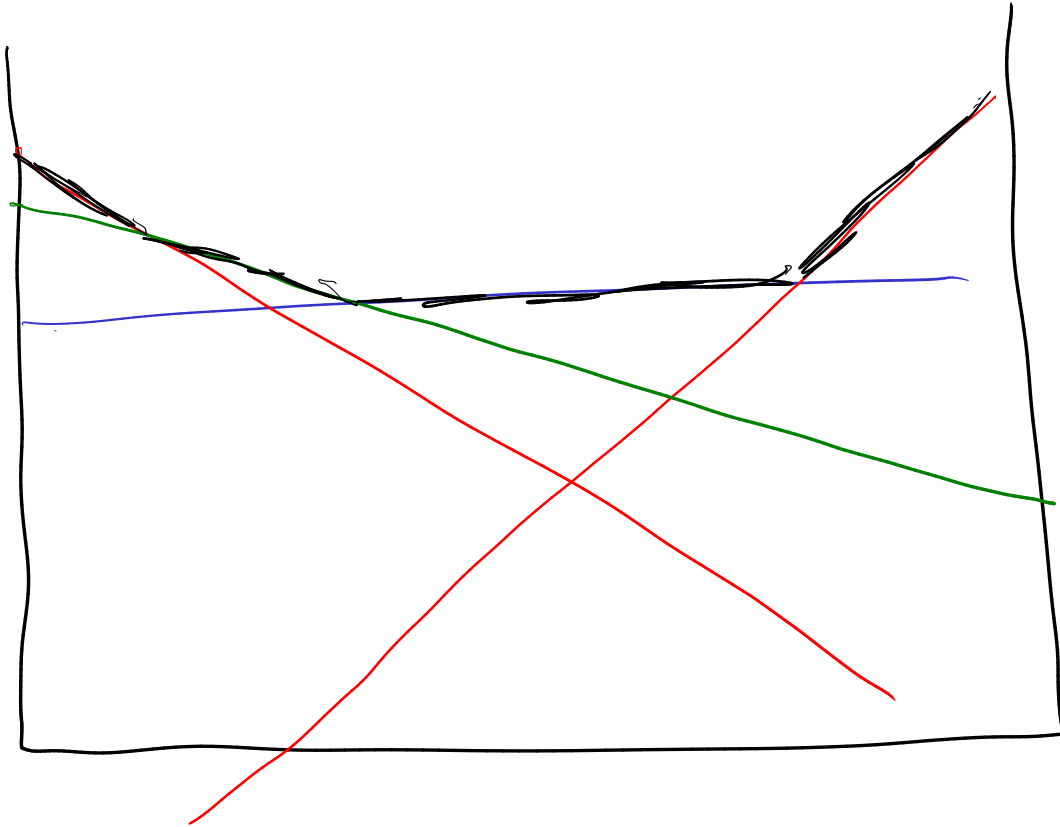
$OR$
TL — TR
$L$ $L$

$U^\pi(TL) = 9.5$
$U^\pi(TR) = -100.95$

13.2

# POMDP Value Functions

# POMDP Value Functions



$$V^*(b) = \max_{\alpha \in \Gamma} \alpha^\top b$$

# Exercise: 2-Step Crying Baby $\alpha$ Vectors

$S = \{h, \neg h\}$

$A = \{f, \neg f\}$

$O = \{c, \neg c\}$

$T(h \mid h, \neg f) = 1.0$

$T(h \mid \neg h, \neg f) = 0.1$

$T(\neg h \mid \cdot, f) = 1.0$

$R(s, a) = R(s) + R(a)$
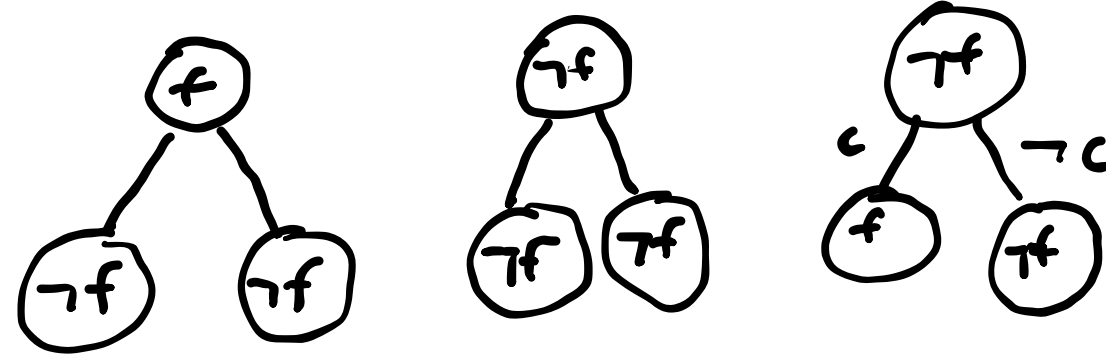
$R(s) = \begin{cases} -10 \text{ if } s = h \\ 0 \text{ otherwise} \end{cases}$

$R(a) = \begin{cases} -5 \text{ if } a = f \\ 0 \text{ otherwise} \end{cases}$
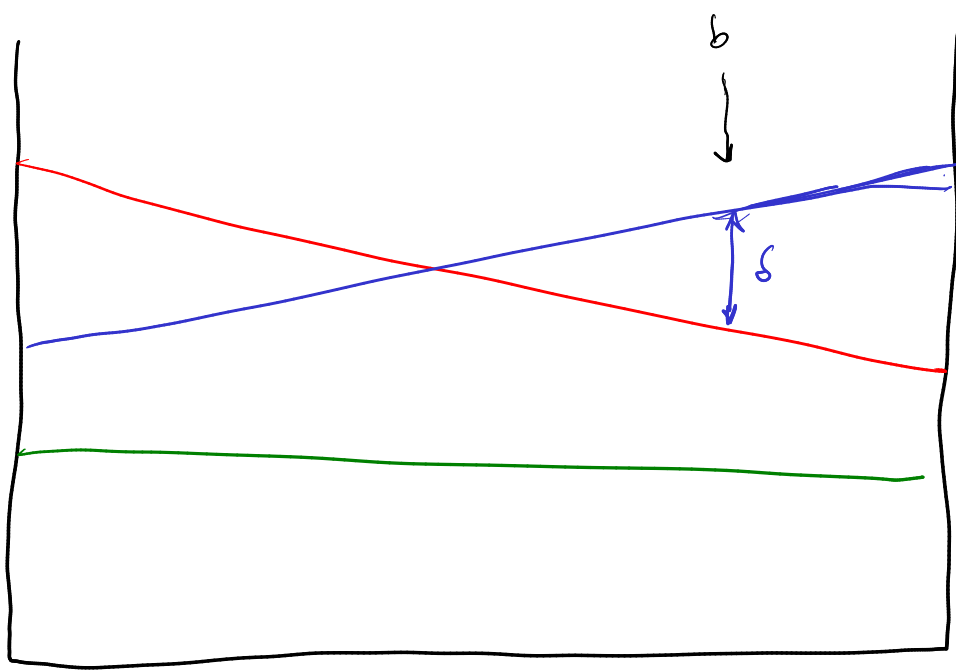
$Z(c \mid \cdot, h) = 0.8)$

$Z(c \mid \cdot, \neg h) = 0.1$

$\gamma = 0.9$

$U^\pi(s) = R(s, \pi()) + \gamma \left[ \sum_{s'} T(s' \mid s, \pi()) \sum_{o} O(o \mid \pi(), s') U^{\pi(o)}(s') \right]$

# $\alpha$-Vector Pruning



maximize $\delta$
$\delta$ $b$
subject to $b \geq 0$
$1^\top b = 1$ $\Big\}$ enforce $b$ is probability

$\alpha^\top b > \alpha'^\top b + \delta \quad \forall \alpha' \in \Gamma$

- If there is a positive $\delta$ solution then $\alpha$ is not dominated

- $b$ is sometimes called the "witress"

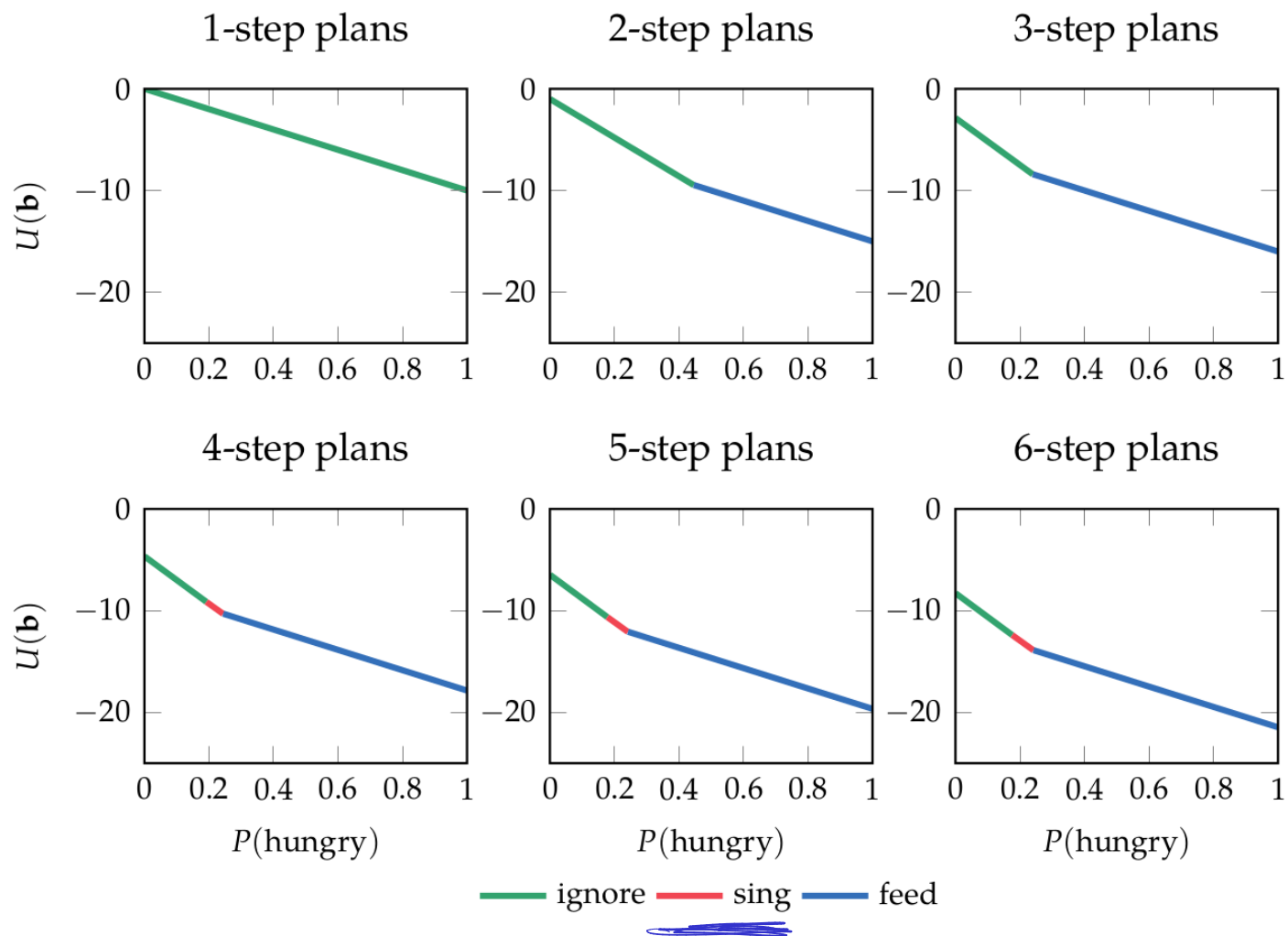# Alpha Vector Expansion

# POMDP Value Iteration (horizon $d$)

$\Gamma^0 \leftarrow \emptyset$

for $n \in 1 \dots d$

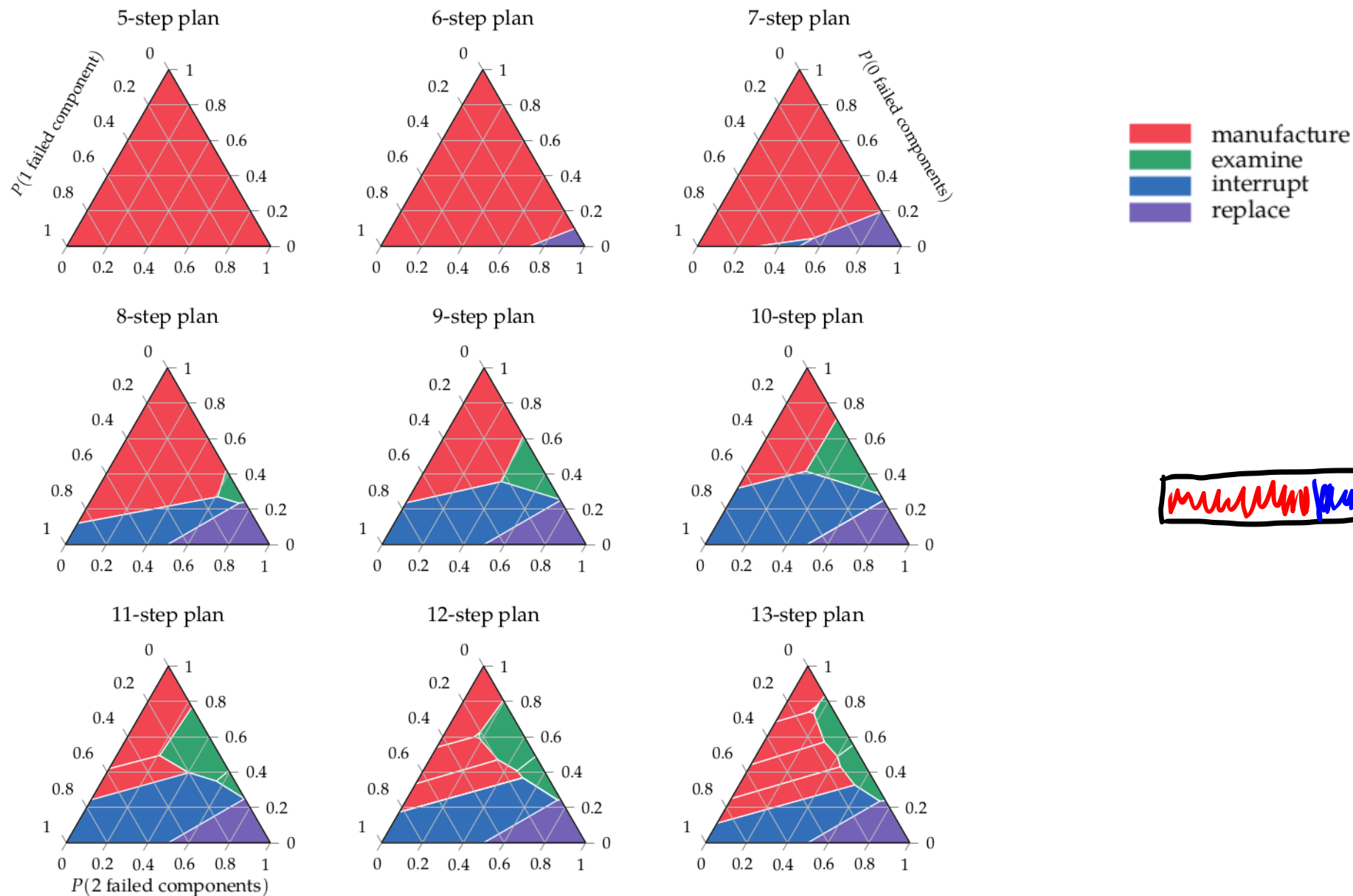    Construct $\Gamma^n$ by expanding with $\Gamma^{n-1}$

    Prune $\Gamma^n$

# Finite Horizon POMDP Value Iteration

# Finite Horizon POMDP Value Iteration

# Recap

# Recap

- A POMDP is an MDP on the _____

# Recap

- A POMDP is an MDP on the <u>belief space</u>

# Recap

- A POMDP is an MDP on the <u>belief space</u>
- The value function of a discrete POMDP can be represented by a set of _____

# Recap

- A POMDP is an MDP on the <u>belief space</u>
- The value function of a discrete POMDP can be represented by a set of <u>$\alpha$-vectors</u>

# Recap

- A POMDP is an MDP on the <u>belief spac</u>e
- The value function of a discrete POMDP can be represented by a set of <u>$\alpha$-vectors</u>
- Each $\alpha$ vector corresponds to a _____

# Recap

- A POMDP is an MDP on the <u>belief spac</u>e
- The value function of a discrete POMDP can be represented by a set of <u>$\alpha$-vectors</u>
- Each $\alpha$ vector corresponds to a <u>conditional plan</u>