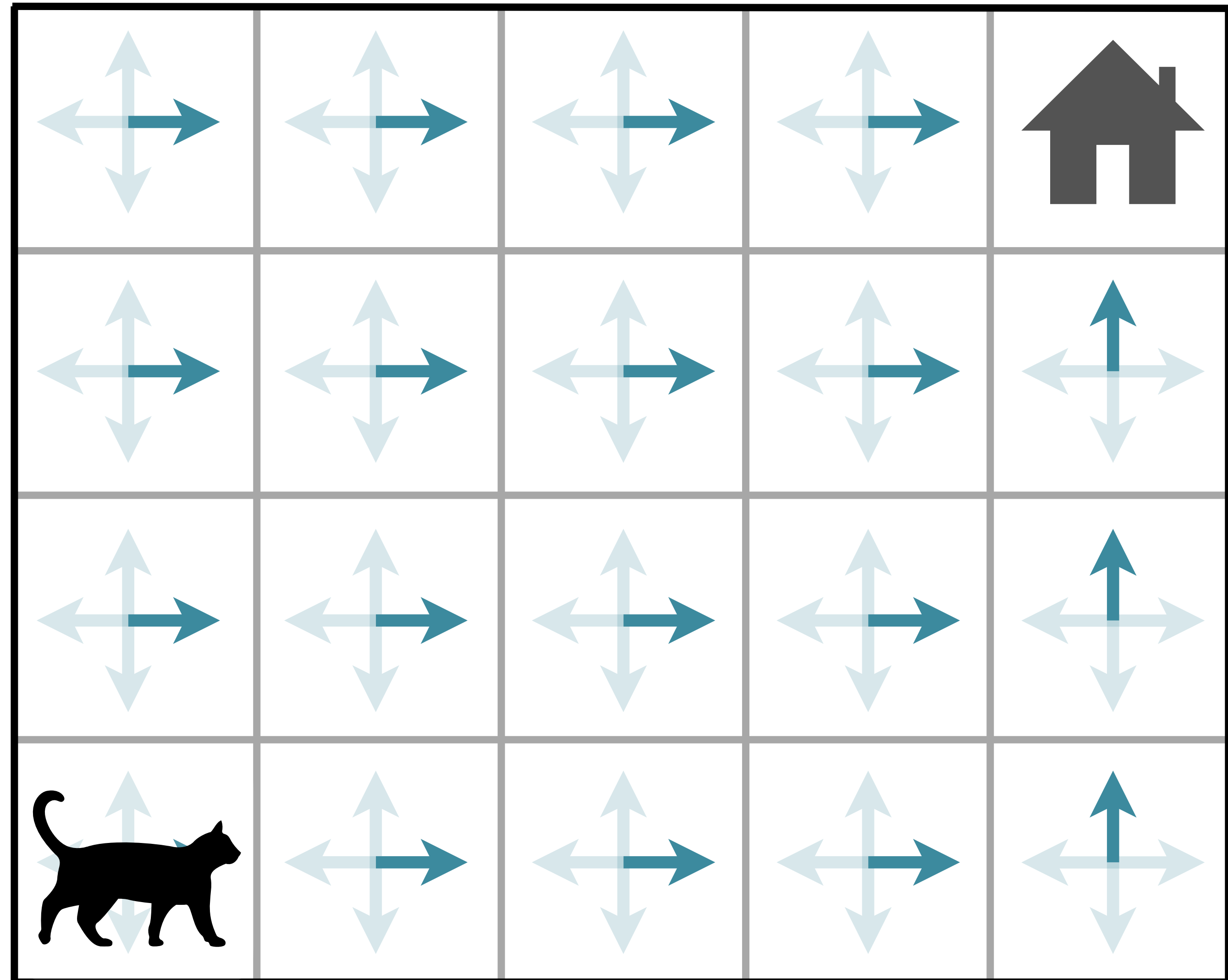


# Importance Sampling Ratio Placement for Gradient-TD Methods

Andy Patterson



# Roadmap

- **Importance Sampling (warm-up)**
- Off-policy TD(0) isr placement
- IS variance
- Gradient-TD placements

# Importance Sampling

**Sample:**  $x \sim b$

**Estimate:**  $\mathbb{E}_{\pi}[X]$

# Importance Sampling

**Sample:**  $x \sim b$

**Estimate:**  $\mathbb{E}_\pi[X]$

# Importance Sampling

**Sample:**  $x \sim b$

**Estimate:**  $\mathbb{E}_{\pi}[X]$

# Derivation of Importance Sampling

$$\mathbb{E}_{\pi}[X]$$

# Derivation of Importance Sampling

$$\mathbb{E}_{\pi}[X] \doteq \sum_{x \in X} x \pi(x)$$

# Derivation of Importance Sampling

$$\begin{aligned}\mathbb{E}_{\pi}[X] &\doteq \sum_{x \in X} x \pi(x) \\ &= \sum_{x \in X} x \pi(x) \frac{b(x)}{b(x)}\end{aligned}$$



# Derivation of Importance Sampling

$$\begin{aligned}\mathbb{E}_{\pi}[X] &\doteq \sum_{x \in X} x \pi(x) \\ &= \sum_{x \in X} x \pi(x) \frac{b(x)}{b(x)} \\ &= \sum_{x \in X} x \frac{\pi(x)}{b(x)} b(x)\end{aligned}$$

# Derivation of Importance Sampling

$$\mathbb{E}_{\pi}[X] \doteq \sum_{x \in X} x \pi(x)$$

$$= \sum_{x \in X} x \pi(x) \frac{b(x)}{b(x)}$$

$$= \sum_{x \in X} x \frac{\pi(x)}{b(x)} b(x)$$

**Importance  
sampling ratio**



# Derivation of Importance Sampling

$$\begin{aligned}\mathbb{E}_{\pi}[X] &\doteq \sum_{x \in X} x \pi(x) \\ &= \sum_{x \in X} x \pi(x) \frac{b(x)}{b(x)} \\ &= \sum_{x \in X} x \rho(x) b(x)\end{aligned}$$

# Derivation of Importance Sampling

$$\mathbb{E}_{\pi}[X] = \sum_{x \in X} x \rho(x) b(x)$$

# Derivation of Importance Sampling

$$\mathbb{E}_{\pi}[X] = \sum_{x \in X} x \rho(x) b(x)$$

# Derivation of Importance Sampling

$$\mathbb{E}_{\pi}[X] = \sum_{x \in X} x \rho(x) b(x)$$

# Derivation of Importance Sampling

$$\begin{aligned}\mathbb{E}_{\pi}[X] &= \sum_{x \in X} x \rho(x) b(x) \\ &= \mathbb{E}_b[X \rho(X)]\end{aligned}$$

# Derivation of Importance Sampling

$$\begin{aligned}\mathbb{E}_{\pi}[X] &= \sum_{x \in X} x \rho(x) b(x) \\ &= \mathbb{E}_b[X \rho(X)]\end{aligned}$$



# Roadmap

- **Importance Sampling**
- Off-policy TD(0) isr placement
- IS variance
- Gradient-TD placements

# Roadmap

- Importance Sampling
- **Off-policy TD(0) isr placement**
- IS variance
- Gradient-TD placements

# Off-Policy TD(0)

$$\delta = \rho(r + \gamma v') - v$$

$$w \leftarrow w + \alpha \delta x$$

# Off-Policy TD(0)

$$\delta = \rho(r + \gamma v') - v$$

$$\delta^+ = \rho(r + \gamma v' - v)$$

# Off-Policy TD(0)

$$\delta = \rho(r + \gamma v') - v$$

**Precup, Sutton, Singh (2000)**

**Precup, Sutton, Dasgupta (2001)**

$$\delta^+ = \rho(r + \gamma v' - v)$$

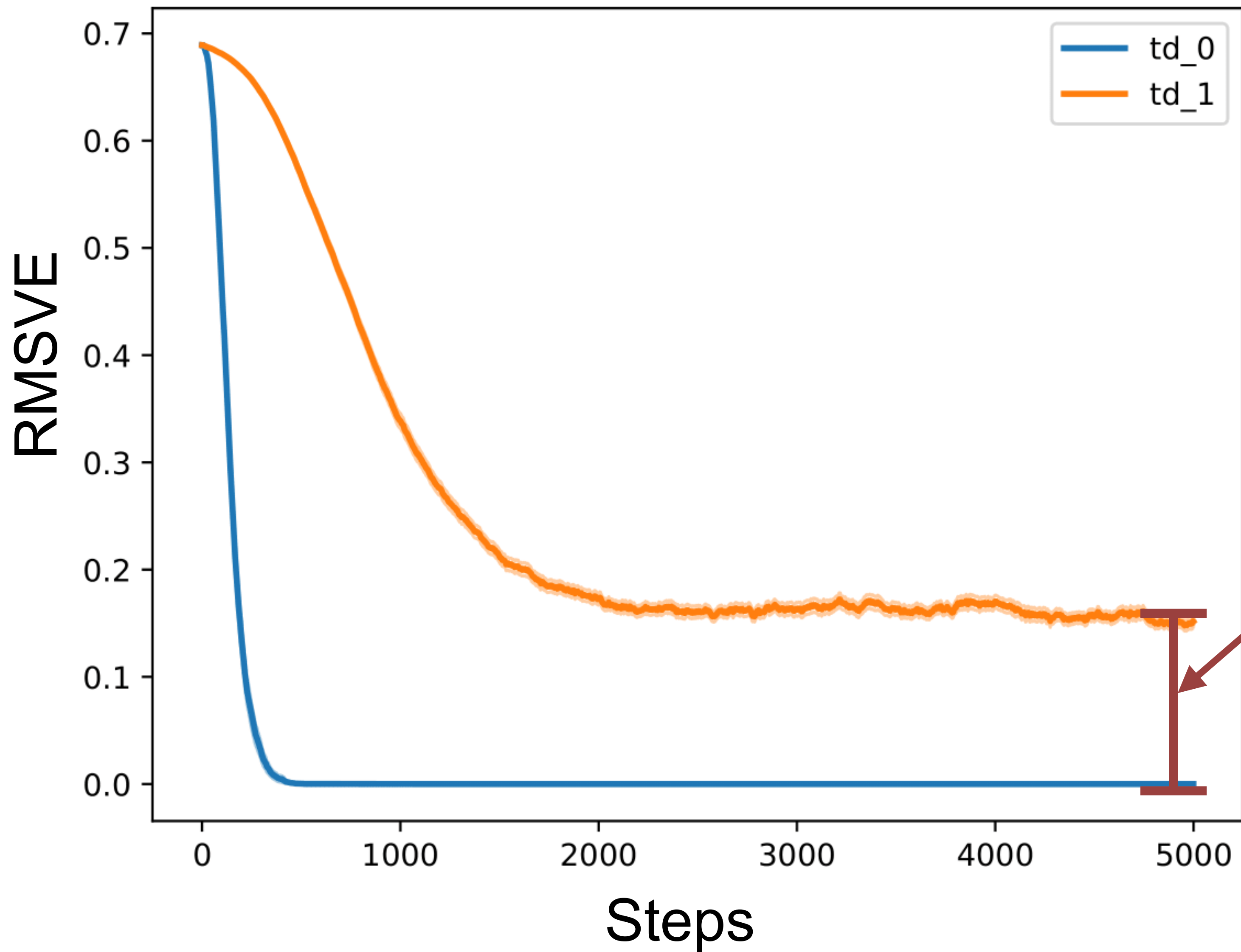
**Maei (2011)**

**van Hasselt, Mahmood, Sutton (2014)**

**Mahmood, van Hasselt, Sutton (2014)**

**td\_0:**  $\delta^+ = \rho(r + \gamma v' - v)$

**td\_1:**  $\delta = \rho(r + \gamma v') - v$



**Biased?**

$$\delta = \rho(r + \gamma v') - v$$

$$\delta^+ = \rho(r + \gamma v' - v)$$

---

$$\delta = \rho(r + \gamma v') - v$$

$$\delta^+ = \rho(r + \gamma v' - v)$$

---

$$\mathbb{E}_b[\delta] = \mathbb{E}_b[\rho(r + \gamma v') - v]$$

$$\mathbb{E}_b[\delta^+] = \mathbb{E}_b[\rho(r + \gamma v' - v)]$$



$$\delta = \rho(r + \gamma v') - v$$

$$\delta^+ = \rho(r + \gamma v' - v)$$

---

$$\begin{aligned}\mathbb{E}_b[\delta] &= \mathbb{E}_b[\rho(r + \gamma v') - v] \\ &= \mathbb{E}_b[\rho(r + \gamma v')] - \mathbb{E}_b[v]\end{aligned}$$

$$\begin{aligned}\mathbb{E}_b[\delta^+] &= \mathbb{E}_b[\rho(r + \gamma v' - v)] \\ &= \mathbb{E}_b[\rho(r + \gamma v')] - \mathbb{E}_b[\rho v]\end{aligned}$$

$$\delta = \rho(r + \gamma v') - v$$

$$\delta^+ = \rho(r + \gamma v' - v)$$

---

$$\begin{aligned}\mathbb{E}_b[\delta] &= \mathbb{E}_b[\rho(r + \gamma v') - v] \\ &= \mathbb{E}_b[\rho(r + \gamma v')] - \mathbb{E}_b[v]\end{aligned}$$

$$\begin{aligned}\mathbb{E}_b[\delta^+] &= \mathbb{E}_b[\rho(r + \gamma v' - v)] \\ &= \mathbb{E}_b[\rho(r + \gamma v')] - \mathbb{E}_b[\rho v]\end{aligned}$$

$$\delta = \rho(r + \gamma v') - v$$

$$\delta^+ = \rho(r + \gamma v' - v)$$

---

$$\begin{aligned}\mathbb{E}_b[\delta] &= \mathbb{E}_b[\rho(r + \gamma v') - v] \\ &= \mathbb{E}_b[\rho(r + \gamma v')] - \mathbb{E}_b[v]\end{aligned}$$

$$\begin{aligned}\mathbb{E}_b[\delta^+] &= \mathbb{E}_b[\rho(r + \gamma v' - v)] \\ &= \mathbb{E}_b[\rho(r + \gamma v')] - \mathbb{E}_b[\rho v]\end{aligned}$$

$$\mathbb{E}_b[v]$$

$$\mathbb{E}_b[\rho v]$$

$$\delta = \rho(r + \gamma v') - v$$

$$\delta^+ = \rho(r + \gamma v' - v)$$

---

$$\begin{aligned}\mathbb{E}_b[\delta] &= \mathbb{E}_b[\rho(r + \gamma v') - v] \\ &= \mathbb{E}_b[\rho(r + \gamma v')] - \mathbb{E}_b[v]\end{aligned}$$

$$\begin{aligned}\mathbb{E}_b[\delta^+] &= \mathbb{E}_b[\rho(r + \gamma v' - v)] \\ &= \mathbb{E}_b[\rho(r + \gamma v')] - \mathbb{E}_b[\rho v]\end{aligned}$$

$$\mathbb{E}_b[v] \stackrel{?}{=} \mathbb{E}_b[\rho v]$$

$$\delta = \rho(r + \gamma v') - v$$

$$\delta^+ = \rho(r + \gamma v' - v)$$

---

$$\begin{aligned}\mathbb{E}_b[\delta] &= \mathbb{E}_b[\rho(r + \gamma v') - v] \\ &= \mathbb{E}_b[\rho(r + \gamma v')] - \mathbb{E}_b[v]\end{aligned}$$

$$\begin{aligned}\mathbb{E}_b[\delta^+] &= \mathbb{E}_b[\rho(r + \gamma v' - v)] \\ &= \mathbb{E}_b[\rho(r + \gamma v')] - \mathbb{E}_b[\rho v]\end{aligned}$$

$$\begin{aligned}\mathbb{E}_b[v] &\stackrel{?}{=} \mathbb{E}_b[\rho v] \\ v &= v\mathbb{E}_b[\rho]\end{aligned}$$

$$\delta = \rho(r + \gamma v') - v$$

$$\delta^+ = \rho(r + \gamma v' - v)$$

---

$$\begin{aligned}\mathbb{E}_b[\delta] &= \mathbb{E}_b[\rho(r + \gamma v') - v] \\ &= \mathbb{E}_b[\rho(r + \gamma v')] - \mathbb{E}_b[v]\end{aligned}$$

$$\begin{aligned}\mathbb{E}_b[\delta^+] &= \mathbb{E}_b[\rho(r + \gamma v' - v)] \\ &= \mathbb{E}_b[\rho(r + \gamma v')] - \mathbb{E}_b[\rho v]\end{aligned}$$

$$\mathbb{E}_b[v] \stackrel{?}{=} \mathbb{E}_b[\rho v]$$

$$v = v \mathbb{E}_b[\rho]$$

$$\mathbb{E}_b[\rho] = \sum \frac{\pi(x)}{b(x)} b(x)$$

$$= \sum \pi(x) = 1$$

$$\delta = \rho(r + \gamma v') - v$$

$$\delta^+ = \rho(r + \gamma v' - v)$$

---

$$\begin{aligned}\mathbb{E}_b[\delta] &= \mathbb{E}_b[\rho(r + \gamma v') - v] \\ &= \mathbb{E}_b[\rho(r + \gamma v')] - \mathbb{E}_b[v]\end{aligned}$$

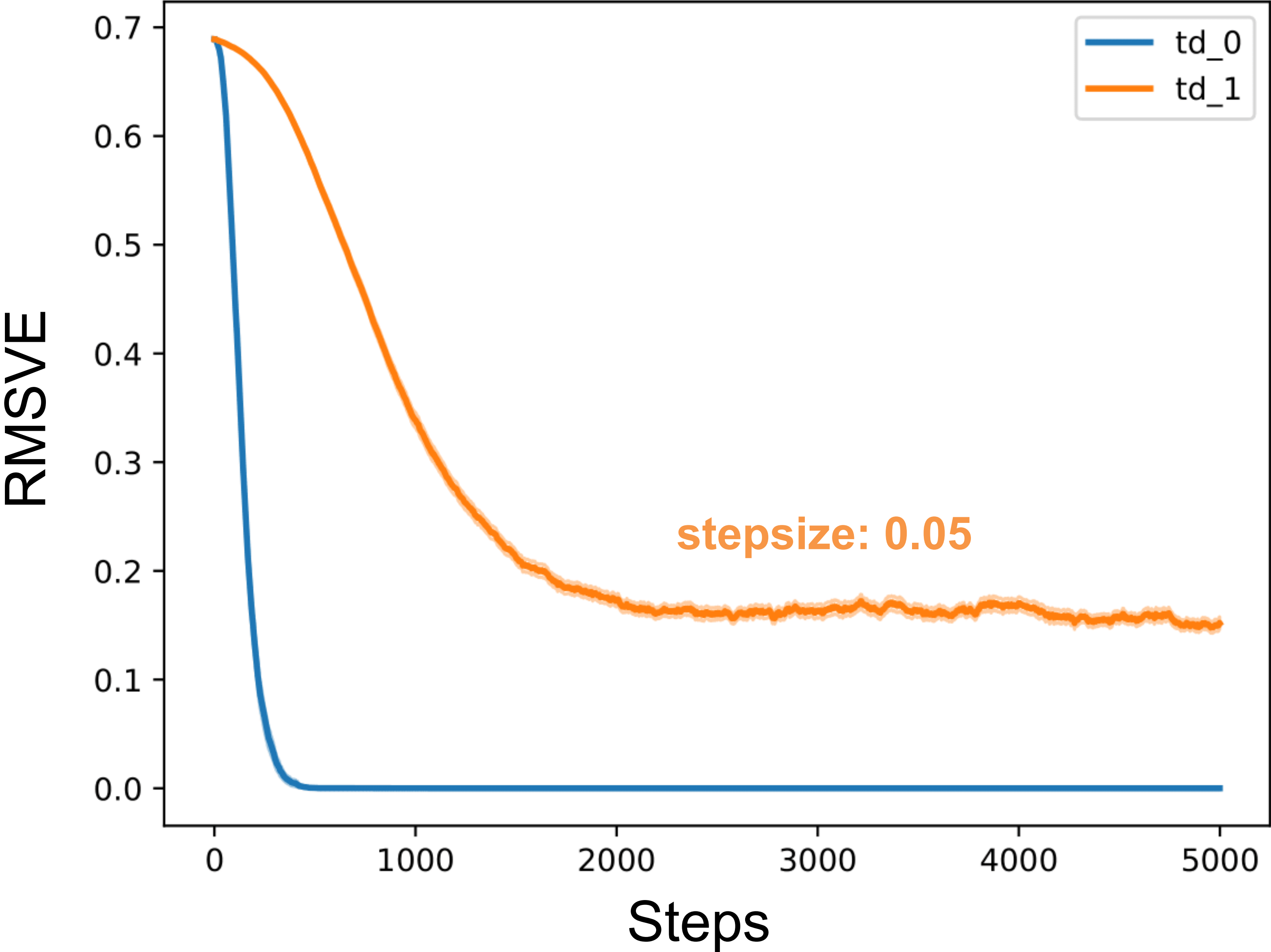
$$\begin{aligned}\mathbb{E}_b[\delta^+] &= \mathbb{E}_b[\rho(r + \gamma v' - v)] \\ &= \mathbb{E}_b[\rho(r + \gamma v')] - \mathbb{E}_b[\rho v]\end{aligned}$$

$$\mathbb{E}_b[v] \stackrel{?}{=} \mathbb{E}_b[\rho v]$$

$$v = v \mathbb{E}_b[\rho]$$

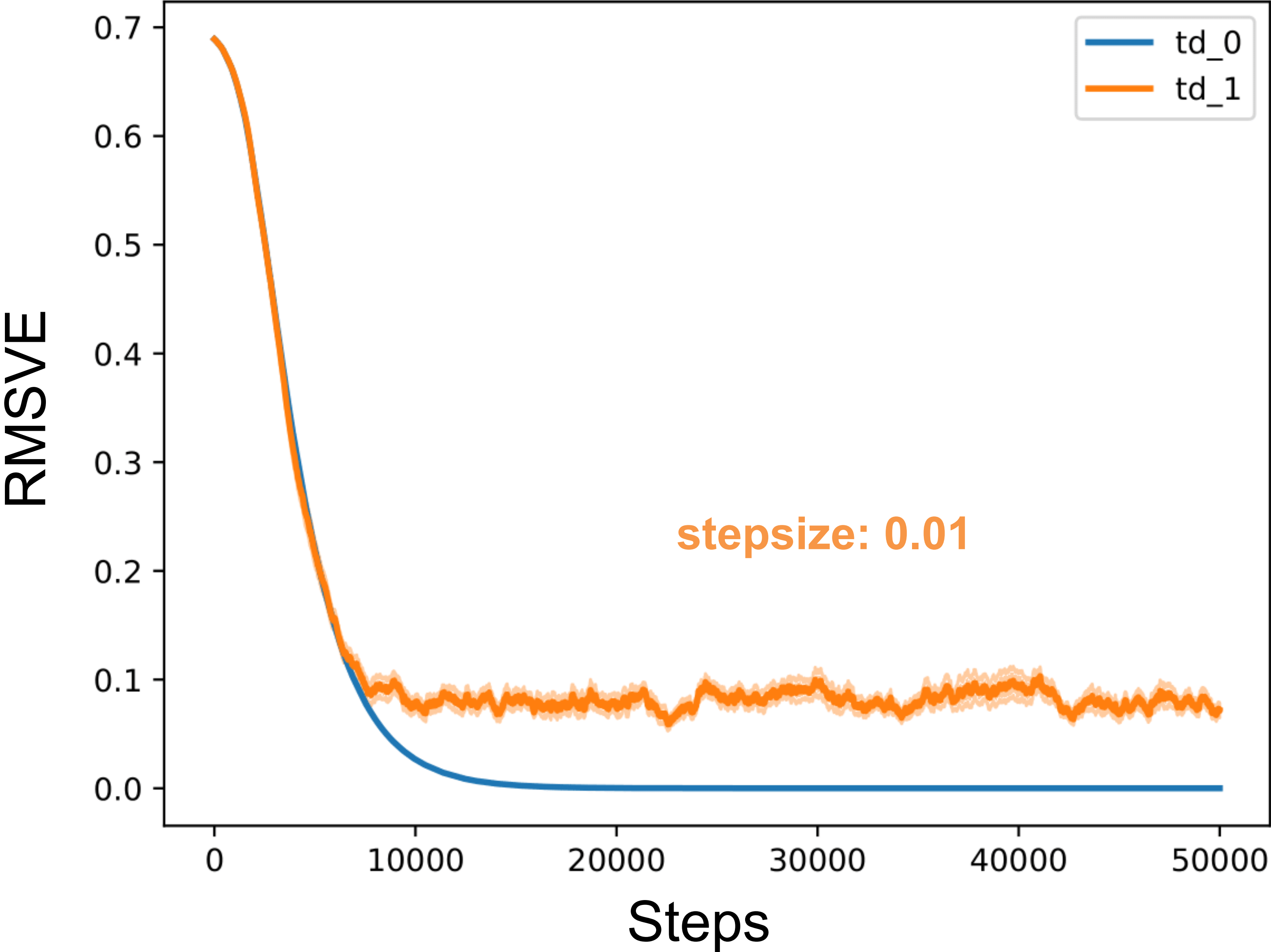
$$v = v$$

# 5k steps

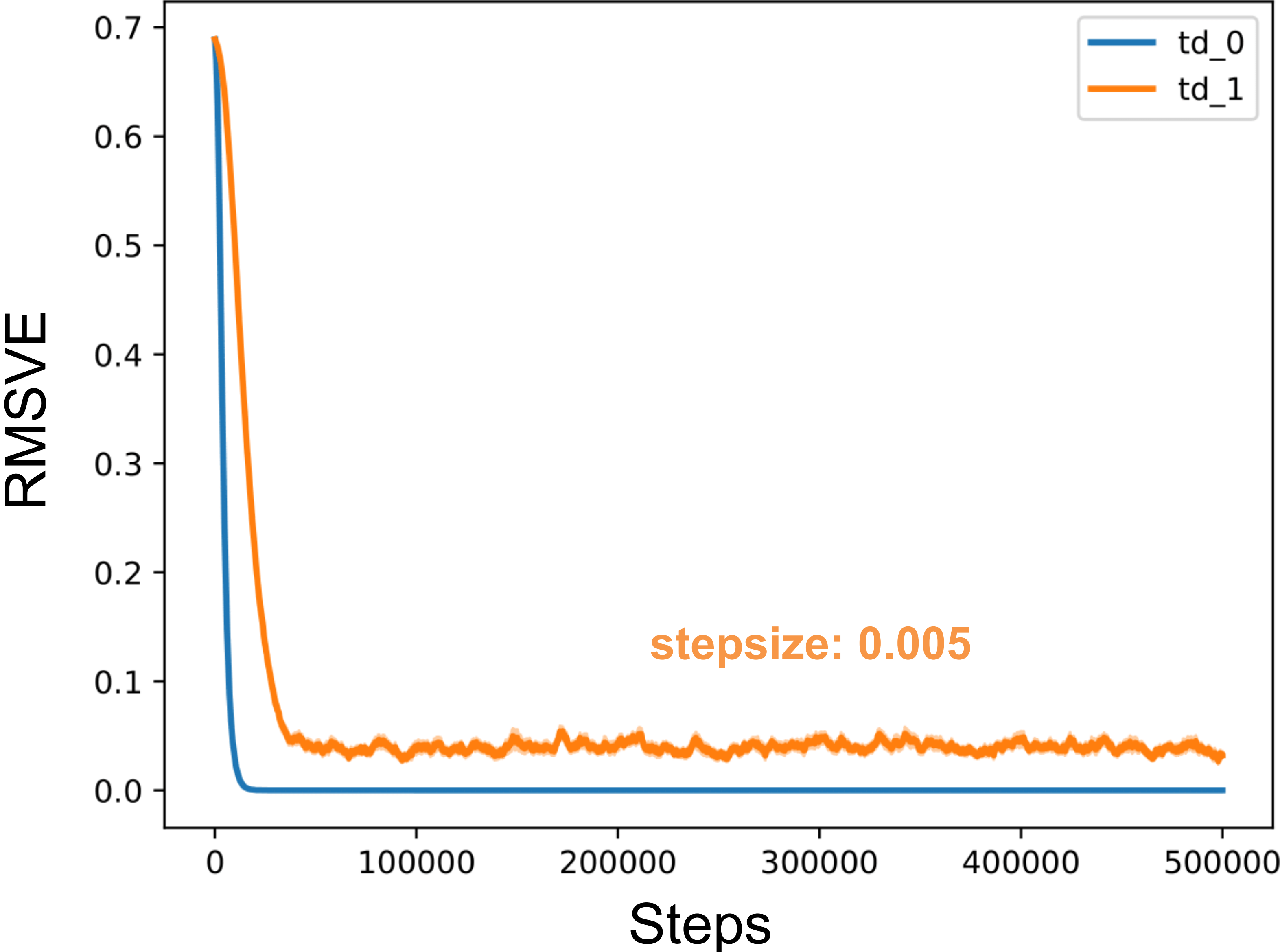




# 50k steps



# 500k steps



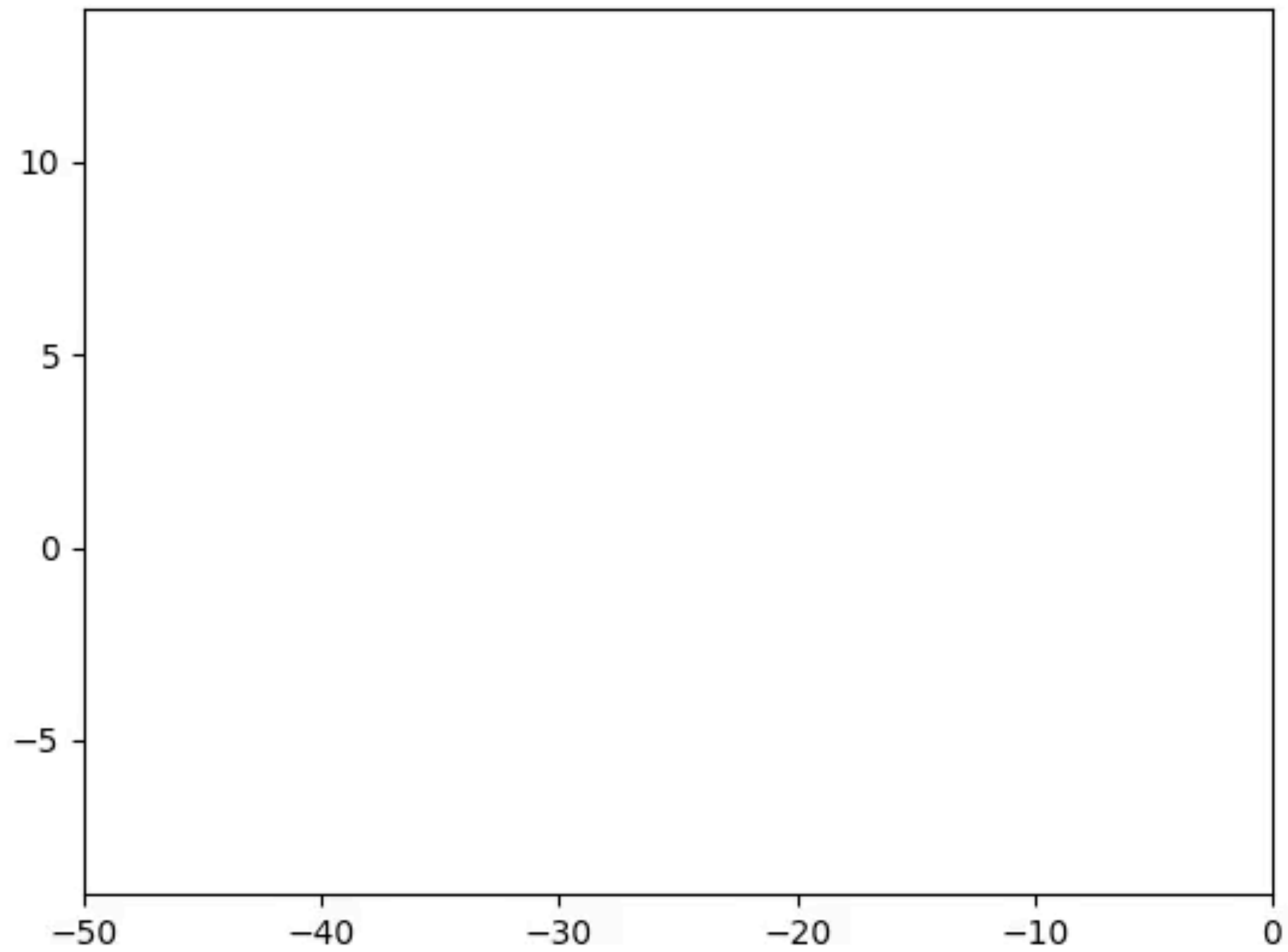
# Roadmap

- Importance Sampling
- **Off-policy TD(0) isr placement**
- IS variance
- Gradient-TD placements

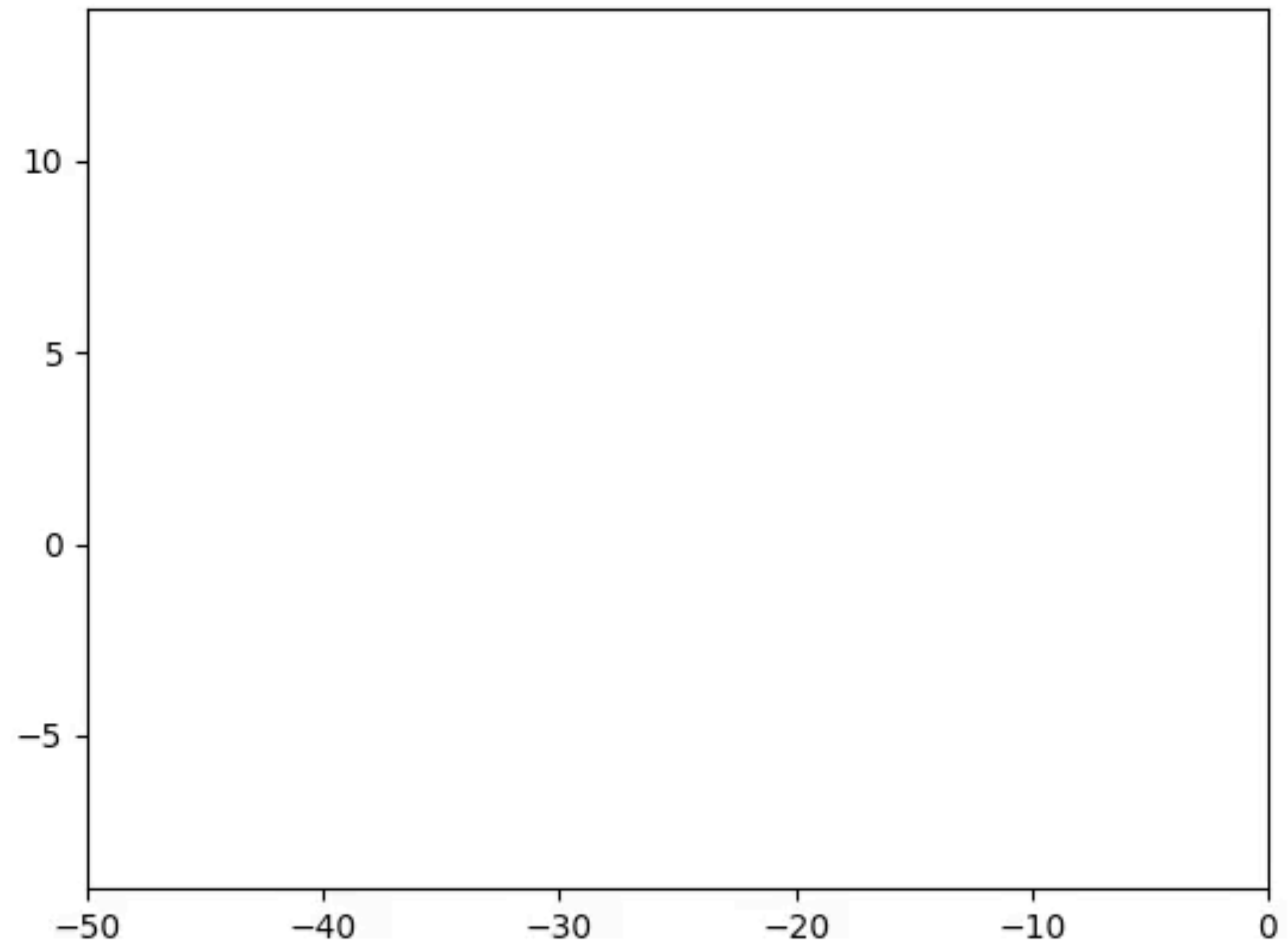
# Roadmap

- Importance Sampling
- Off-policy TD(0) isr placement
- **IS variance**
- Gradient-TD placements

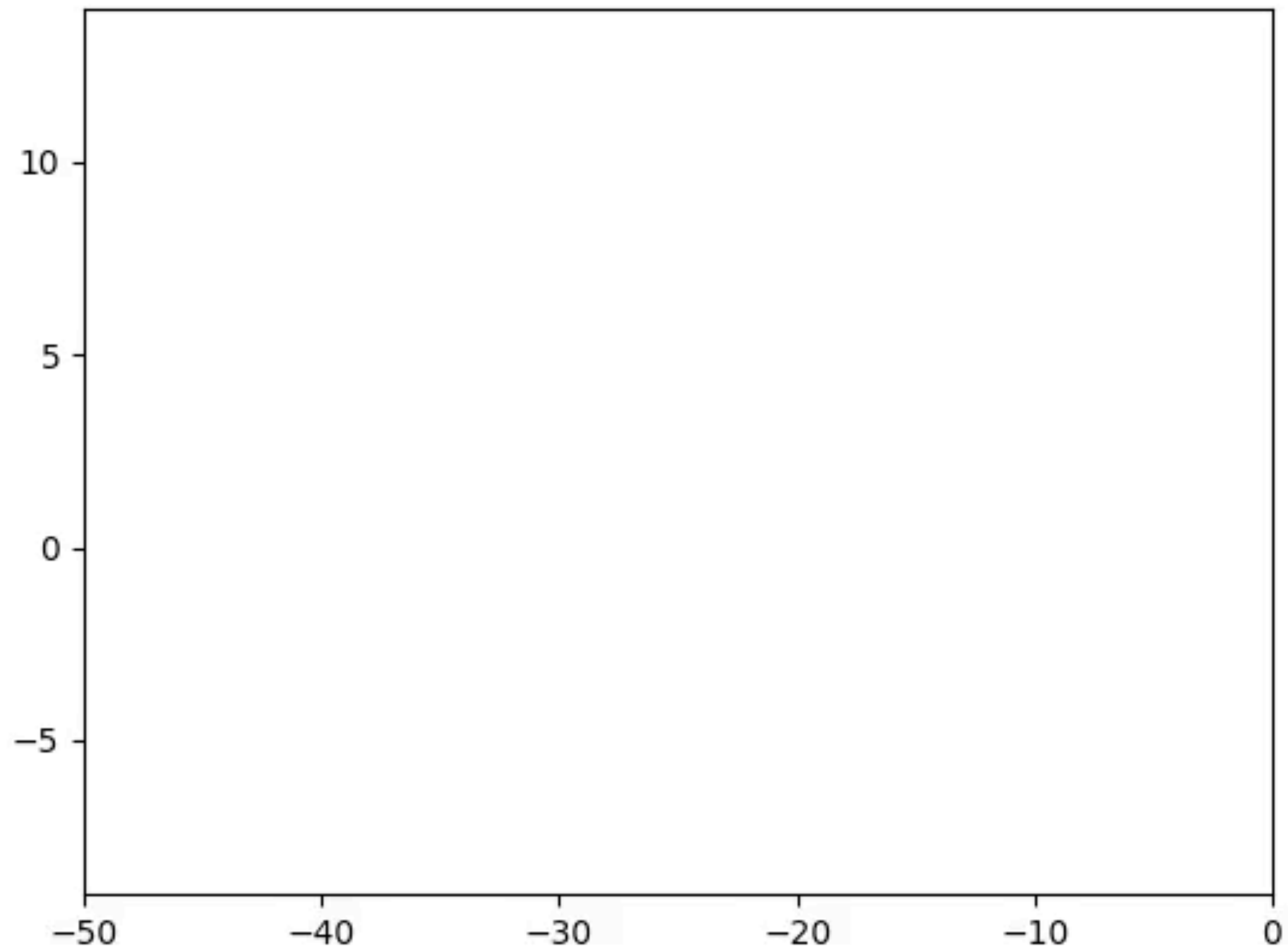
$$x \sim b$$



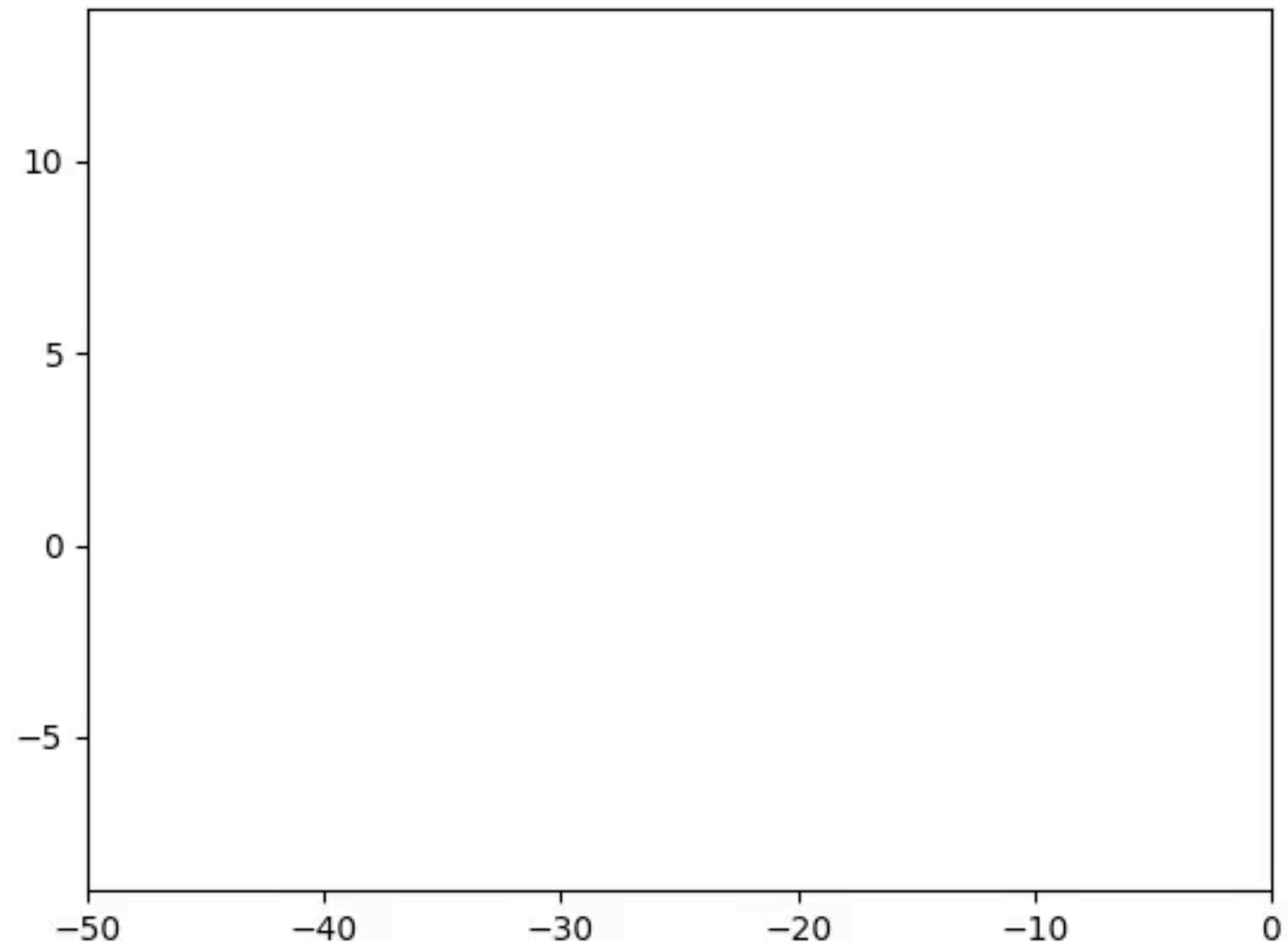
$$\rho(x)x$$



$$x \sim b$$



$$\rho(x)x$$



# Control Variates

$$X^* = X + c(Y - \mathbb{E}_b[Y])$$

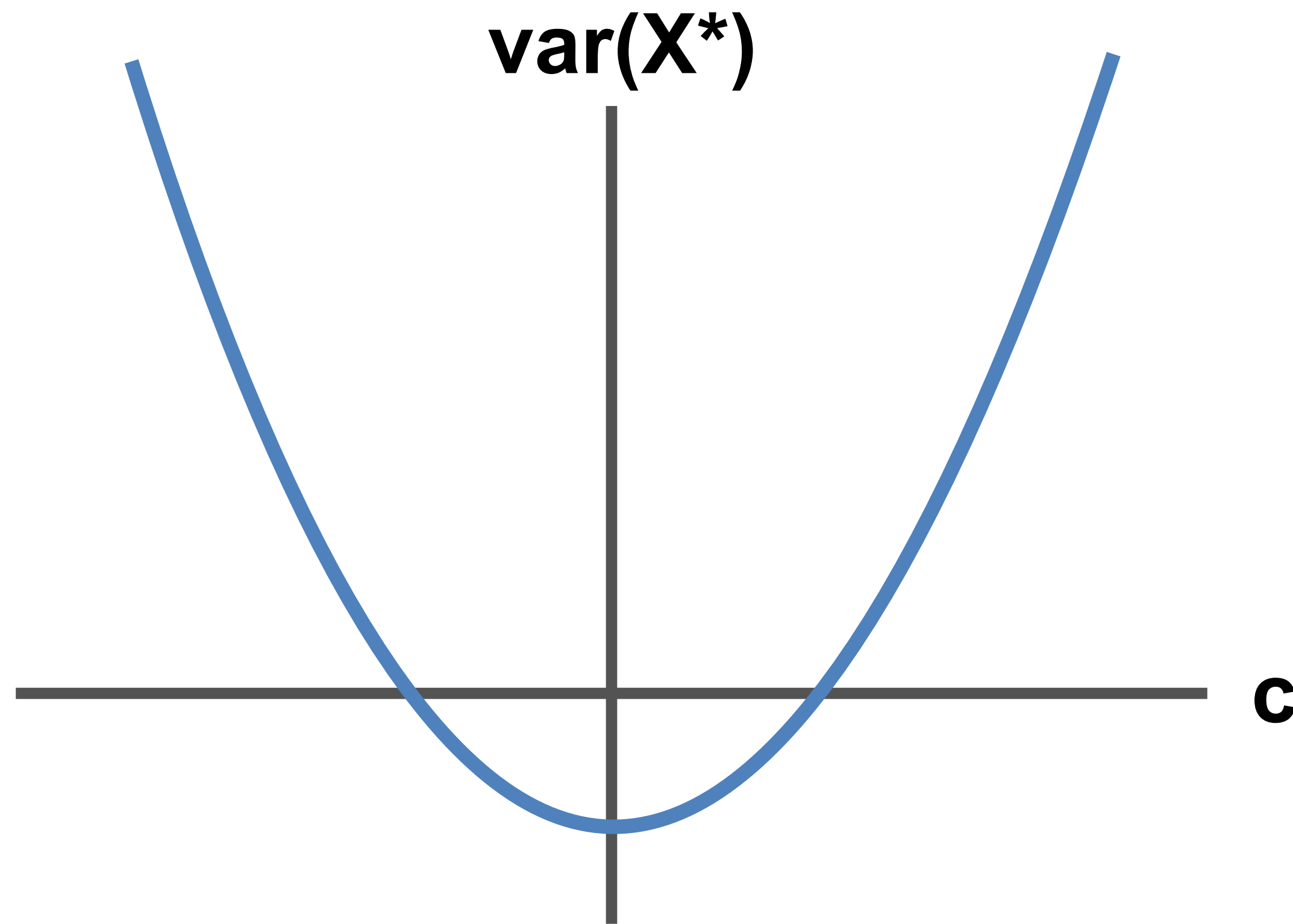
# Control Variates

$$X^* = X + c(Y - \mathbb{E}_b[Y])$$

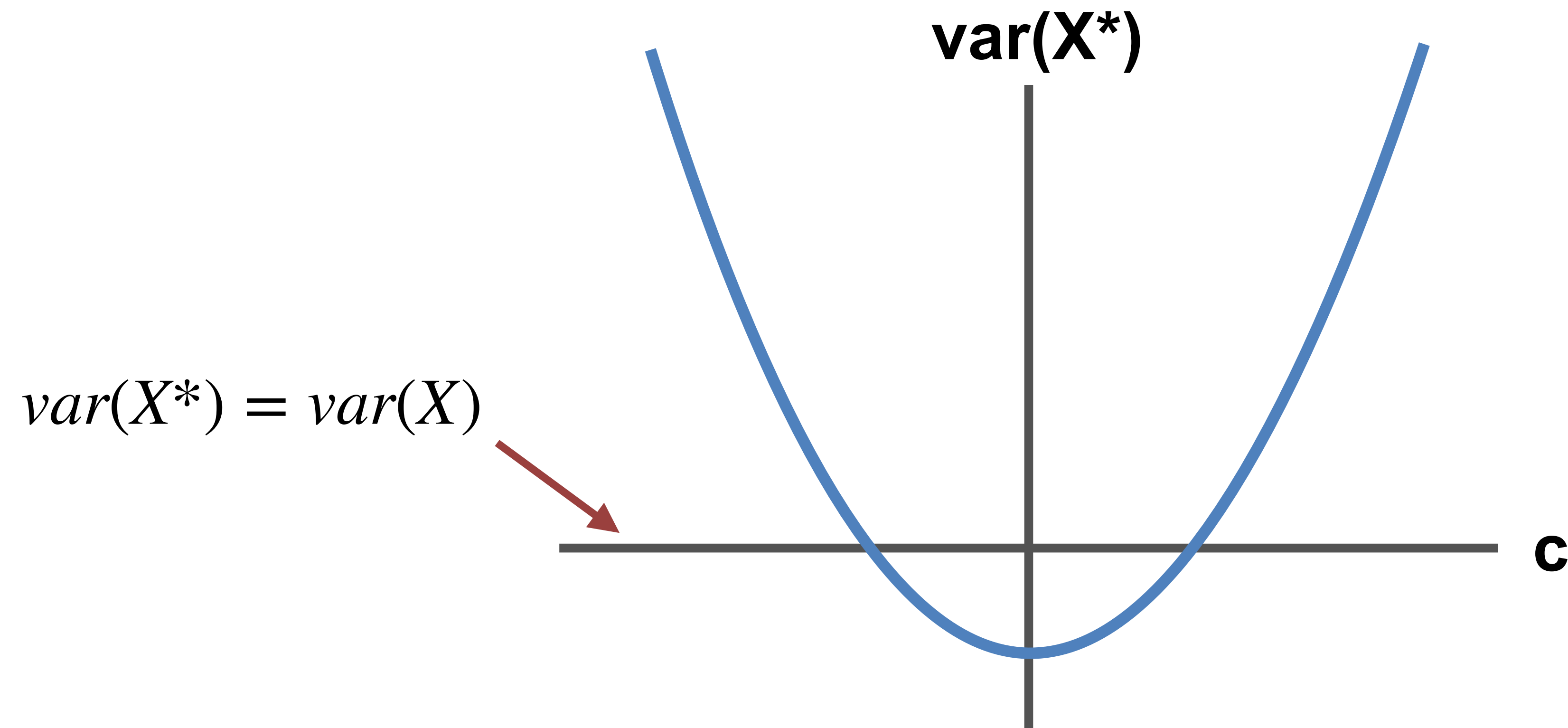
Variance  
Control



$$X^* = X + c(Y - \mathbb{E}_b[Y])$$



$$X^* = X + c(Y - \mathbb{E}_b[Y])$$



$$X^* = X + c(Y - \mathbb{E}_b[Y])$$

$$X \doteq \rho(r + \gamma v') - v$$

$$Y \doteq \rho v$$

$$c \doteq -1$$

$$X^* = X + c(Y - \mathbb{E}_b[Y])$$

$$\delta^* = \delta + (-1)(\rho v - \mathbb{E}_b[\rho v])$$

$$X \doteq \rho(r + \gamma v') - v$$

$$Y \doteq \rho v$$

$$c \doteq -1$$

$$X^* = X + c(Y - \mathbb{E}_b[Y])$$

$$\begin{aligned}\delta^* &= \delta + (-1)(\rho v - \mathbb{E}_b[\rho v]) \\ &= \delta + (-1)(\rho v - \mathbb{E}_b[\rho v])\end{aligned}$$

$$X \doteq \rho(r + \gamma v') - v$$

$$Y \doteq \rho v$$

$$c \doteq -1$$

$$X^* = X + c(Y - \mathbb{E}_b[Y])$$

$$\begin{aligned}\delta^* &= \delta + (-1)(\rho v - \mathbb{E}_b[\rho v]) \\ &= \rho(r + \gamma v') - v + (-1)(\rho v - \mathbb{E}_b[\rho v])\end{aligned}$$

$$X \doteq \rho(r + \gamma v') - v$$

$$Y \doteq \rho v$$

$$c \doteq -1$$

$$X^* = X + c(Y - \mathbb{E}_b[Y])$$

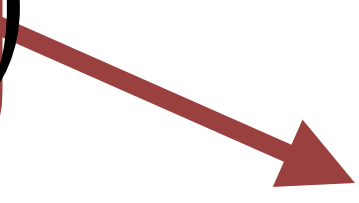
$$X \doteq \rho(r + \gamma v') - v$$

$$Y \doteq \rho v$$

$$c \doteq -1$$

$$\delta^* = \delta + (-1)(\rho v - \mathbb{E}_b[\rho v])$$

$$= \rho(r + \gamma v') - v + (-1)(\rho v - \mathbb{E}_b[\rho v])$$


$$\mathbb{E}_\pi[v] = v$$

$$X^* = X + c(Y - \mathbb{E}_b[Y])$$

$$\begin{aligned}\delta^* &= \delta + (-1)(\rho v - \mathbb{E}_b[\rho v]) \\ &= \rho(r + \gamma v') - v + (-1)(\rho v - \mathbb{E}_b[\rho v])\end{aligned}$$

$$X \doteq \rho(r + \gamma v') - v$$

$$Y \doteq \rho v$$

$$c \doteq -1$$



$$X^* = X + c(Y - \mathbb{E}_b[Y])$$

$$\begin{aligned}\delta^* &= \delta + (-1)(\rho v - \mathbb{E}_b[\rho v]) \\ &= \rho(r + \gamma v') - v + (-1)(\rho v - \mathbb{E}_b[\rho v]) \\ &= \rho(r + \gamma v') - v + v - \rho v\end{aligned}$$

$$X \doteq \rho(r + \gamma v') - v$$

$$Y \doteq \rho v$$

$$c \doteq -1$$

$$X^* = X + c(Y - \mathbb{E}_b[Y])$$

$$\begin{aligned}\delta^* &= \delta + (-1)(\rho v - \mathbb{E}_b[\rho v]) \\ &= \rho(r + \gamma v') - v + (-1)(\rho v - \mathbb{E}_b[\rho v]) \\ &= \rho(r + \gamma v') - v + v - \rho v \\ &= \rho(r + \gamma v') - \rho v\end{aligned}$$

$$X \doteq \rho(r + \gamma v') - v$$

$$Y \doteq \rho v$$

$$c \doteq -1$$

$$X^* = X + c(Y - \mathbb{E}_b[Y])$$

$$\begin{aligned}\delta^* &= \delta + (-1)(\rho v - \mathbb{E}_b[\rho v]) \\ &= \rho(r + \gamma v') - v + (-1)(\rho v - \mathbb{E}_b[\rho v]) \\ &= \rho(r + \gamma v') - v + v - \rho v \\ &= \rho(r + \gamma v') - \rho v \\ &= \rho(r + \gamma v' - v)\end{aligned}$$

$$X \doteq \rho(r + \gamma v') - v$$

$$Y \doteq \rho v$$

$$c \doteq -1$$

$$X^* = X + c(Y - \mathbb{E}_b[Y])$$

$$\begin{aligned}\delta^* &= \delta + (-1)(\rho v - \mathbb{E}_b[\rho v]) \\ &= \rho(r + \gamma v') - v + (-1)(\rho v - \mathbb{E}_b[\rho v]) \\ &= \rho(r + \gamma v') - v + v - \rho v \\ &= \rho(r + \gamma v') - \rho v \\ &= \rho(r + \gamma v' - v) \\ &= \delta^+\end{aligned}$$

$$X \doteq \rho(r + \gamma v') - v$$

$$Y \doteq \rho v$$

$$c \doteq -1$$

# Roadmap

- Importance Sampling
- Off-policy TD(0) isr placement
- **IS variance**
- Gradient-TD placements

# Roadmap

- Importance Sampling
- Off-policy TD(0) isr placement
- IS variance
- **Gradient-TD placements**

# Gradient-TD Update Equations

$$\delta = \rho_t [r_{t+1} + \gamma(w_t^\top x_{t+1}) - (w_t^\top x_t)]$$

$$z_t \leftarrow \rho_{t-1}(\gamma\lambda z_{t-1} + x_t)$$

$$h_{t+1} \leftarrow h_t + \alpha_h [\delta z_t - (h_t^\top x_t)x_t]$$

## TDC

$$w_{t+1} \leftarrow w_t + \alpha [\delta z_t - \rho_t \gamma (1 - \lambda) (h_t^\top z_t) x_{t+1}]$$

## GTD2

$$w_{t+1} \leftarrow w_t + \alpha [(h_t^\top x_t)x_t - \rho_t \gamma (1 - \lambda) (h_t^\top z_t)x_{t+1}]$$

# Gradient-TD Update Equations

$$\delta = \rho_t [r_{t+1} + \gamma(w_t^\top x_{t+1}) - (w_t^\top x_t)]$$

$$z_t \leftarrow \rho_{t-1}(\gamma\lambda z_{t-1} + x_t)$$

$$h_{t+1} \leftarrow h_t + \alpha_h [\delta z_t - (h_t^\top x_t)x_t]$$

## TDC

$$w_{t+1} \leftarrow w_t + \alpha [\delta z_t - \rho_t \gamma (1 - \lambda) (h_t^\top z_t) x_{t+1}]$$

## GTD2

$$w_{t+1} \leftarrow w_t + \alpha [(h_t^\top x_t)x_t - \rho_t \gamma (1 - \lambda) (h_t^\top z_t)x_{t+1}]$$



# Gradient-TD Update Equations

$$\delta = \rho_t [r_{t+1} + \gamma(w_t^\top x_{t+1}) - (w_t^\top x_t)]$$

$$z_t \leftarrow \rho_{t-1}(\gamma\lambda z_{t-1} + x_t)$$

$$h_{t+1} \leftarrow h_t + \alpha_h [\delta z_t - (h_t^\top x_t)x_t]$$

## TDC

$$w_{t+1} \leftarrow w_t + \alpha [\delta z_t - \rho_t \gamma (1 - \lambda)(h_t^\top z_t)x_{t+1}]$$

## GTD2

$$w_{t+1} \leftarrow w_t + \alpha [(h_t^\top x_t)x_t - \rho_t \gamma (1 - \lambda)(h_t^\top z_t)x_{t+1}]$$

# Gradient-TD Update Equations

$$\delta = \rho_t [r_{t+1} + \gamma(w_t^\top x_{t+1}) - (w_t^\top x_t)]$$

$$z_t \leftarrow \rho_{t-1}(\gamma\lambda z_{t-1} + x_t)$$

$$h_{t+1} \leftarrow h_t + \alpha_h [\delta z_t - (h_t^\top x_t)x_t]$$

## TDC

$$w_{t+1} \leftarrow w_t + \alpha [\delta z_t - \rho_t \gamma(1 - \lambda)(h_t^\top z_t)x_{t+1}]$$

## GTD2

$$w_{t+1} \leftarrow w_t + \alpha [(h_t^\top x_t)x_t - \rho_t \gamma(1 - \lambda)(h_t^\top z_t)x_{t+1}]$$

**0: correct everything**  
**1: correct as little as possible**

**a:**  $\nabla_h$   
**b:**  $\delta_h$   
**c:**  $\delta_w$

**TDC**

**GTD2**

**a:**  $\nabla_h$   
**b:**  $\delta_h$   
**c:**  $\nabla_w$

**0: correct everything**  
**1: correct as little as possible**

**a:**  $\nabla_h$   
**b:**  $\delta_h$   
**c:**  $\delta_w$

## TDC

$$h_{t+1} \leftarrow h_t + \alpha_h [\delta z_t - (h_t^\top x_t) x_t]$$


## GTD2

**a:**  $\nabla_h$   
**b:**  $\delta_h$   
**c:**  $\nabla_w$

0: correct everything  
1: correct as little as possible

**a:**  $\nabla_h$   
**b:**  $\delta_h$   
**c:**  $\delta_w$

### TDC

$$h_{t+1} \leftarrow h_t + \alpha_h [\delta z_t - \rho_t (h_t^\top x_t) x_t]$$


### GTD2

**a:**  $\nabla_h$   
**b:**  $\delta_h$   
**c:**  $\nabla_w$

**0: correct everything**  
**1: correct as little as possible**

**a:**  $\nabla_h$   
**b:**  $\delta_h$   
**c:**  $\delta_w$

## TDC

$$h_{t+1} \leftarrow h_t + \alpha_h [\delta z_t - \rho_t (h_t^\top x_t) x_t]$$
$$\delta = \rho_t [r_{t+1} + \gamma (w_t^\top x_{t+1}) - (w_t^\top x_t)]$$

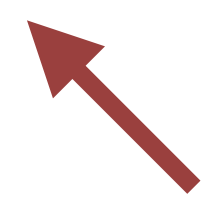
## GTD2

**a:**  $\nabla_h$   
**b:**  $\delta_h$   
**c:**  $\nabla_w$

0: correct everything  
1: correct as little as possible

- a:**  $\nabla_h$
- b:**  $\delta_h$
- c:**  $\delta_w$

### TDC

$$h_{t+1} \leftarrow h_t + \alpha_h [\delta z_t - \rho_t (h_t^\top x_t) x_t]$$
$$\delta = \rho_t [r_{t+1} + \gamma (w_t^\top x_{t+1})] - (w_t^\top x_t)$$


### GTD2

- a:**  $\nabla_h$
- b:**  $\delta_h$
- c:**  $\nabla_w$

**0: correct everything**  
**1: correct as little as possible**

**a:**  $\nabla_h$   
**b:**  $\delta_h$   
**c:**  $\delta_w$

## TDC

$$h_{t+1} \leftarrow h_t + \alpha_h [\delta z_t - \rho_t (h_t^\top x_t) x_t]$$
$$\delta = \rho_t [r_{t+1} + \gamma (w_t^\top x_{t+1})] - (w_t^\top x_t)$$
$$\delta = \rho_t [r_{t+1} + \gamma (w_t^\top x_{t+1}) - (w_t^\top x_t)]$$

## GTD2


**a:**  $\nabla_h$   
**b:**  $\delta_h$   
**c:**  $\nabla_w$



**0: correct everything**  
**1: correct as little as possible**

**a:**  $\nabla_h$   
**b:**  $\delta_h$   
**c:**  $\delta_w$

## TDC

$$h_{t+1} \leftarrow h_t + \alpha_h [\delta z_t - \rho_t (h_t^\top x_t) x_t]$$
$$\delta = \rho_t [r_{t+1} + \gamma (w_t^\top x_{t+1})] - (w_t^\top x_t)$$
$$\delta = \rho_t [r_{t+1} + \gamma (w_t^\top x_{t+1})] - (w_t^\top x_t)$$


## GTD2

**a:**  $\nabla_h$   
**b:**  $\delta_h$   
**c:**  $\nabla_w$

0: correct everything

1: correct as little as possible

**a:**  $\nabla_h$

**b:**  $\delta_h$

**c:**  $\delta_w$

## TDC

$$h_{t+1} \leftarrow h_t + \alpha_h [\delta z_t - \rho_t (h_t^\top x_t) x_t]$$

$$\delta = \rho_t [r_{t+1} + \gamma (w_t^\top x_{t+1})] - (w_t^\top x_t)$$

$$\delta = \rho_t [r_{t+1} + \gamma (w_t^\top x_{t+1})] - (w_t^\top x_t)$$

## GTD2

$$h_{t+1} \leftarrow h_t + \alpha_h [\delta z_t - (h_t^\top x_t) x_t]$$

**a:**  $\nabla_h$

**b:**  $\delta_h$

**c:**  $\nabla_w$

**0: correct everything**  
**1: correct as little as possible**


**a:**  $\nabla_h$   
**b:**  $\delta_h$   
**c:**  $\delta_w$

## TDC

$$h_{t+1} \leftarrow h_t + \alpha_h [\delta z_t - \rho_t (h_t^\top x_t) x_t]$$
$$\delta = \rho_t [r_{t+1} + \gamma (w_t^\top x_{t+1})] - (w_t^\top x_t)$$
$$\delta = \rho_t [r_{t+1} + \gamma (w_t^\top x_{t+1})] - (w_t^\top x_t)$$

## GTD2

**a:**  $\nabla_h$   
**b:**  $\delta_h$   
**c:**  $\nabla_w$

$$h_{t+1} \leftarrow h_t + \alpha_h [\delta z_t - \rho_t (h_t^\top x_t) x_t]$$


0: correct everything

1: correct as little as possible

**a:**  $\nabla_h$

**b:**  $\delta_h$

**c:**  $\delta_w$

## TDC

$$h_{t+1} \leftarrow h_t + \alpha_h [\delta z_t - \rho_t (h_t^\top x_t) x_t]$$

$$\delta = \rho_t [r_{t+1} + \gamma (w_t^\top x_{t+1})] - (w_t^\top x_t)$$

$$\delta = \rho_t [r_{t+1} + \gamma (w_t^\top x_{t+1})] - (w_t^\top x_t)$$

## GTD2

**a:**  $\nabla_h$

**b:**  $\delta_h$

**c:**  $\nabla_w$

$$h_{t+1} \leftarrow h_t + \alpha_h [\delta z_t - \rho_t (h_t^\top x_t) x_t]$$

$$\delta = \rho_t [r_{t+1} + \gamma (w_t^\top x_{t+1}) - (w_t^\top x_t)]$$

**0: correct everything**  
**1: correct as little as possible**


**a:**  $\nabla_h$   
**b:**  $\delta_h$   
**c:**  $\delta_w$

## TDC

$$h_{t+1} \leftarrow h_t + \alpha_h [\delta z_t - \rho_t (h_t^\top x_t) x_t]$$
$$\delta = \rho_t [r_{t+1} + \gamma (w_t^\top x_{t+1})] - (w_t^\top x_t)$$
$$\delta = \rho_t [r_{t+1} + \gamma (w_t^\top x_{t+1})] - (w_t^\top x_t)$$

## GTD2

**a:**  $\nabla_h$   
**b:**  $\delta_h$   
**c:**  $\nabla_w$

$$h_{t+1} \leftarrow h_t + \alpha_h [\delta z_t - \rho_t (h_t^\top x_t) x_t]$$
$$\delta = \rho_t [r_{t+1} + \gamma (w_t^\top x_{t+1})] - (w_t^\top x_t)$$


0: correct everything

1: correct as little as possible

**a:**  $\nabla_h$

**b:**  $\delta_h$

**c:**  $\delta_w$

## TDC

$$h_{t+1} \leftarrow h_t + \alpha_h [\delta z_t - \rho_t (h_t^\top x_t) x_t]$$

$$\delta = \rho_t [r_{t+1} + \gamma (w_t^\top x_{t+1})] - (w_t^\top x_t)$$

$$\delta = \rho_t [r_{t+1} + \gamma (w_t^\top x_{t+1})] - (w_t^\top x_t)$$

## GTD2

**a:**  $\nabla_h$

**b:**  $\delta_h$

**c:**  $\nabla_w$

$$h_{t+1} \leftarrow h_t + \alpha_h [\delta z_t - \rho_t (h_t^\top x_t) x_t]$$

$$\delta = \rho_t [r_{t+1} + \gamma (w_t^\top x_{t+1})] - (w_t^\top x_t)$$

$$w_{t+1} \leftarrow w_t + \alpha [ (h_t^\top x_t) x_t - \rho_t \gamma (1 - \lambda) (h_t^\top z_t) x_{t+1} ]$$

**0: correct everything**  
**1: correct as little as possible**


**a:**  $\nabla_h$   
**b:**  $\delta_h$   
**c:**  $\delta_w$

## TDC

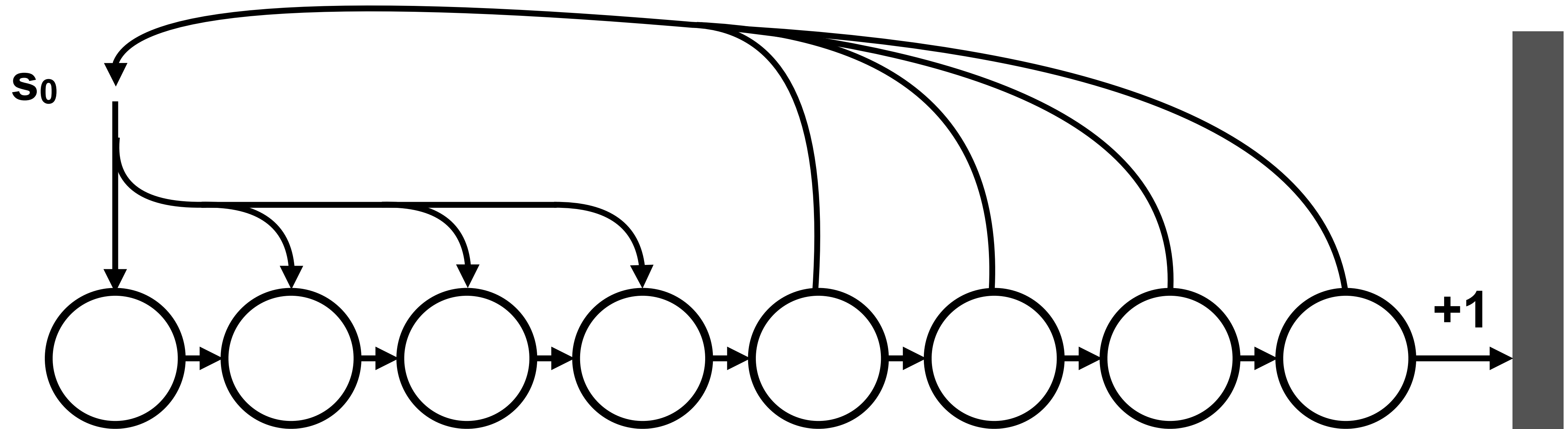
$$h_{t+1} \leftarrow h_t + \alpha_h [\delta z_t - \rho_t (h_t^\top x_t) x_t]$$
$$\delta = \rho_t [r_{t+1} + \gamma (w_t^\top x_{t+1})] - (w_t^\top x_t)$$
$$\delta = \rho_t [r_{t+1} + \gamma (w_t^\top x_{t+1})] - (w_t^\top x_t)$$

## GTD2

**a:**  $\nabla_h$   
**b:**  $\delta_h$   
**c:**  $\nabla_w$

$$h_{t+1} \leftarrow h_t + \alpha_h [\delta z_t - \rho_t (h_t^\top x_t) x_t]$$
$$\delta = \rho_t [r_{t+1} + \gamma (w_t^\top x_{t+1})] - (w_t^\top x_t)$$
$$w_{t+1} \leftarrow w_t + \alpha [\rho_t (h_t^\top x_t) x_t - \rho_t \gamma (1 - \lambda) (h_t^\top z_t) x_{t+1}]$$


# Collision



**Target**

**Behavior**

**Right: 100%**

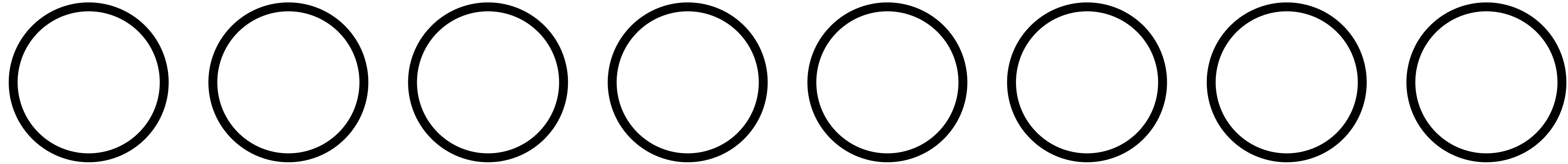
**Right: 50%**

**Retreat: 0%**

**Retreat: 50%**



# Tabular Collision

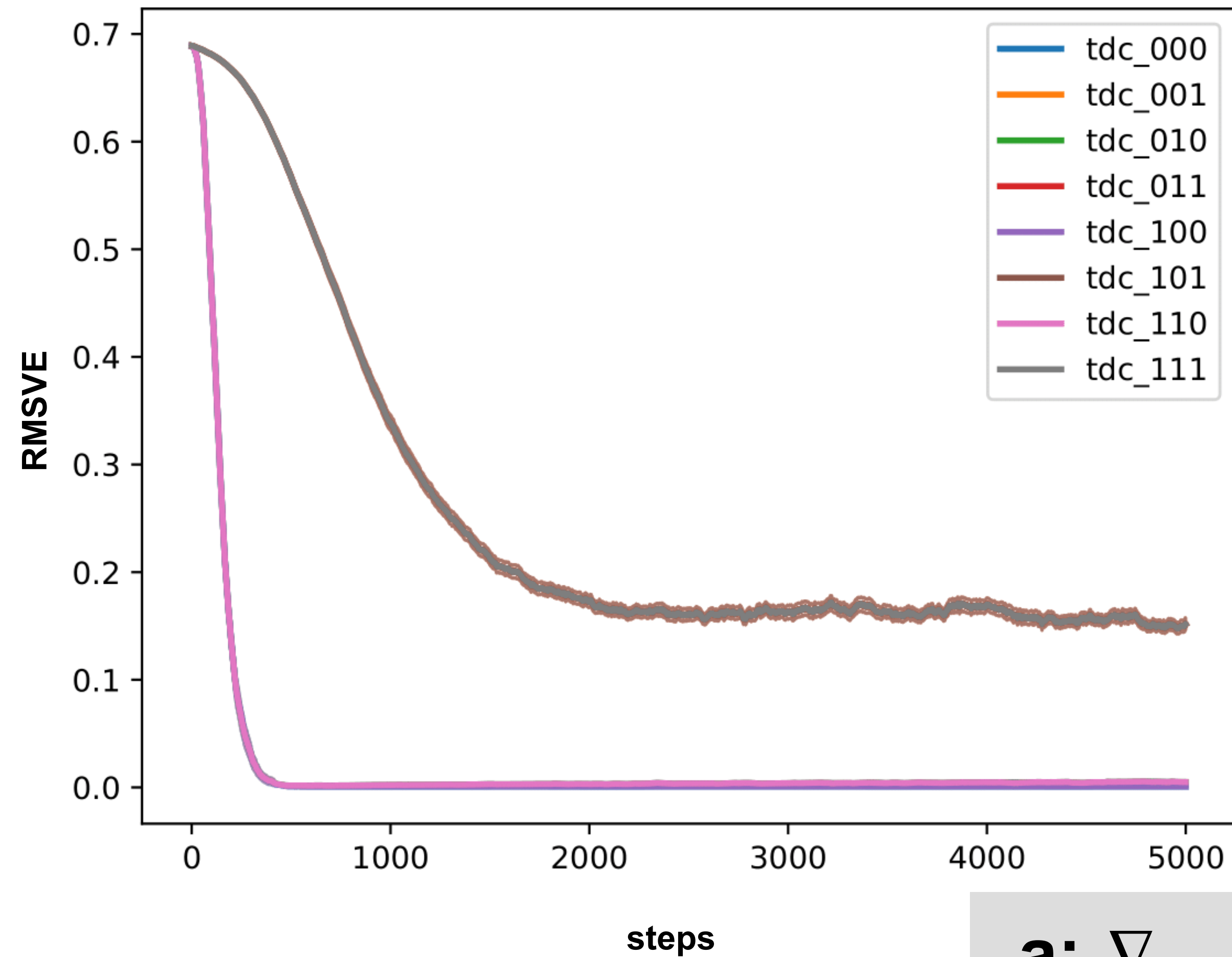


1	0	0	0	0	0	0	0
0	1	0	0	0	0	0	0
0	0	1	0	0	0	0	0
0	0	0	1	0	0	0	0
0	0	0	0	1	0	0	0
0	0	0	0	0	1	0	0
0	0	0	0	0	0	1	0
0	0	0	0	0	0	0	1

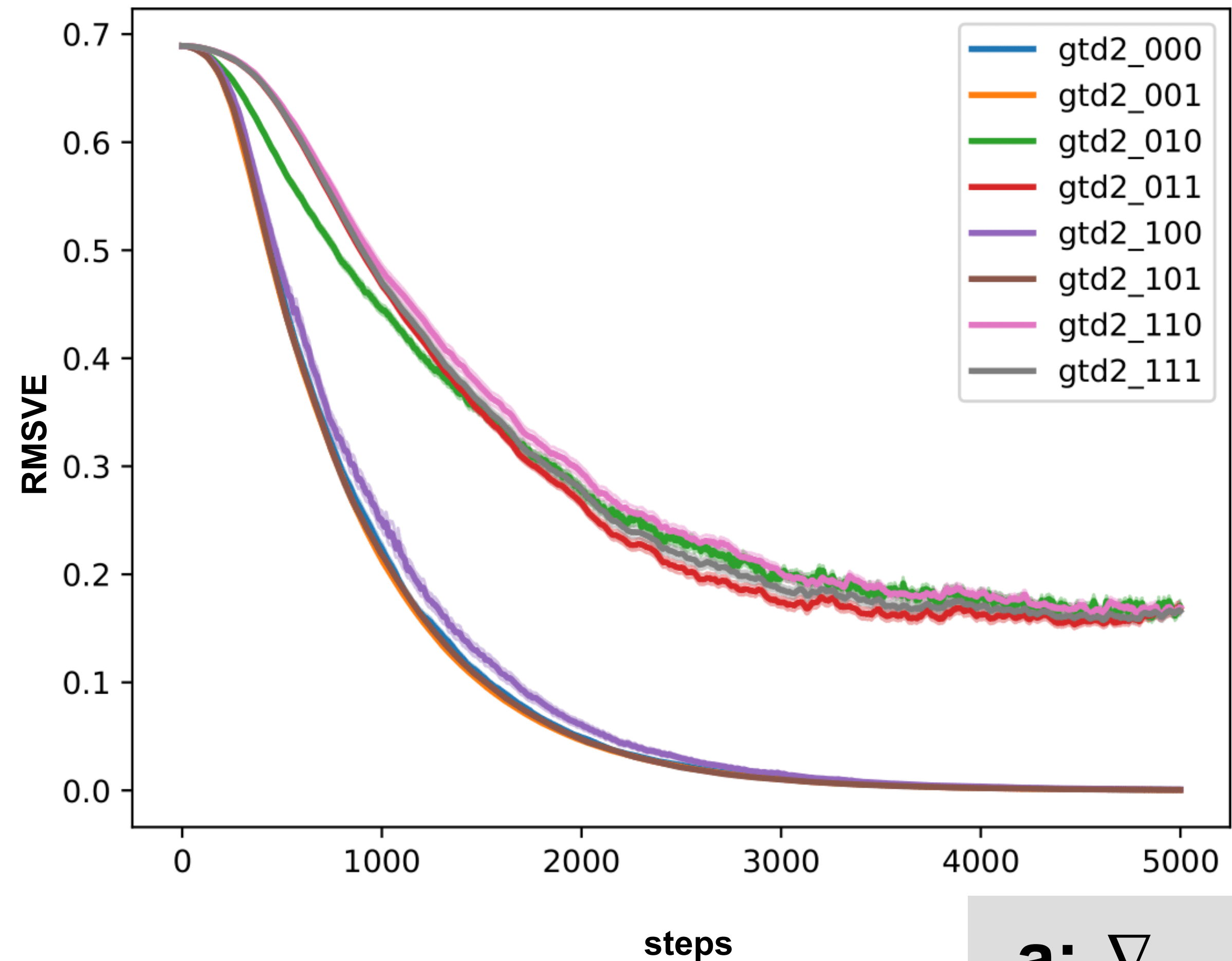
# TDC

0: correct everything  
1: correct as little as possible

# GTD2

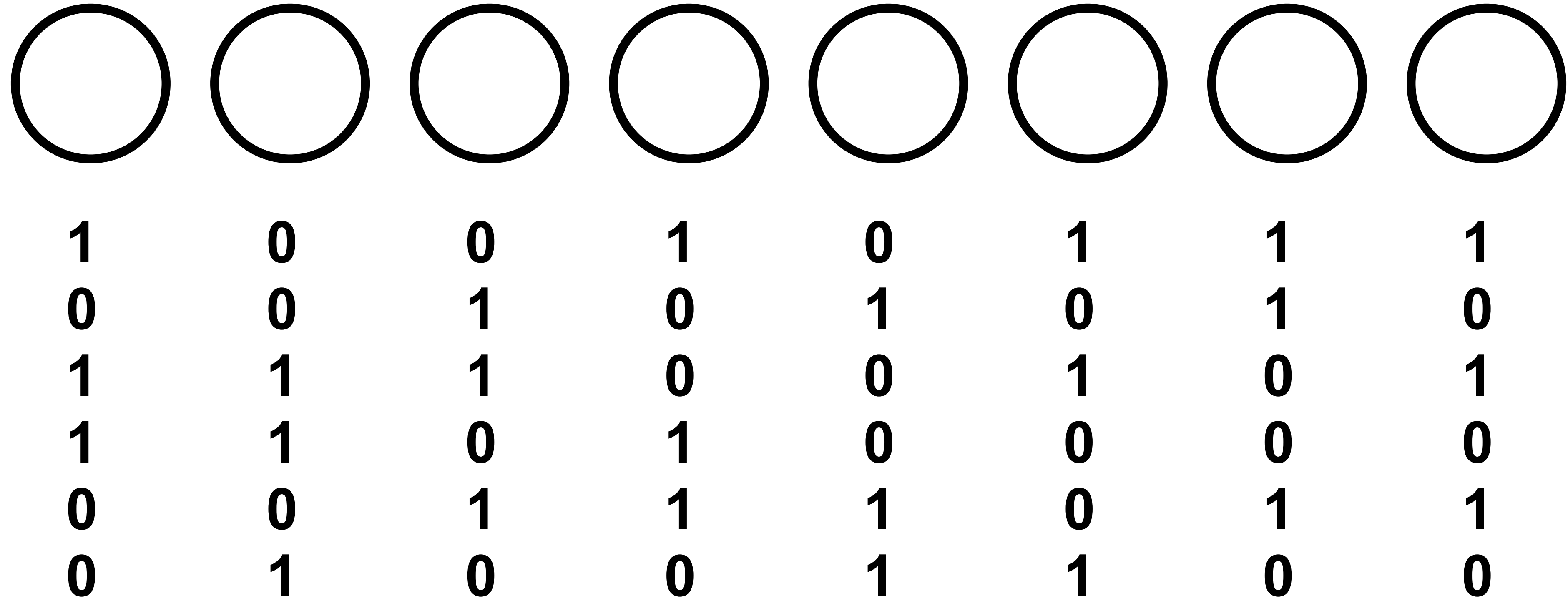


**a:**  $\nabla_h$   
**b:**  $\delta_h$   
**c:**  $\delta_w$



**a:**  $\nabla_h$   
**b:**  $\delta_h$   
**c:**  $\nabla_w$

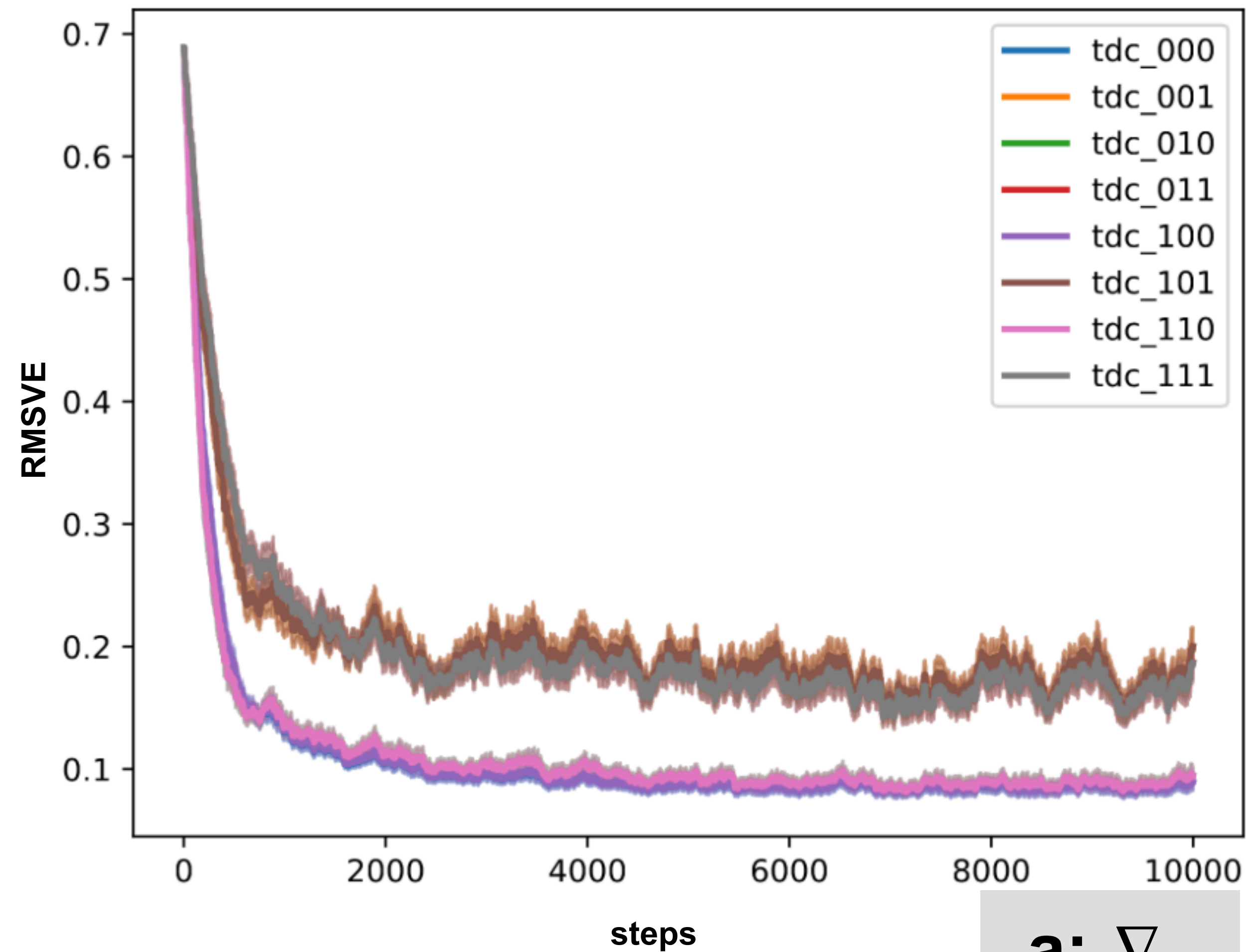
# Binary Encoder Collision



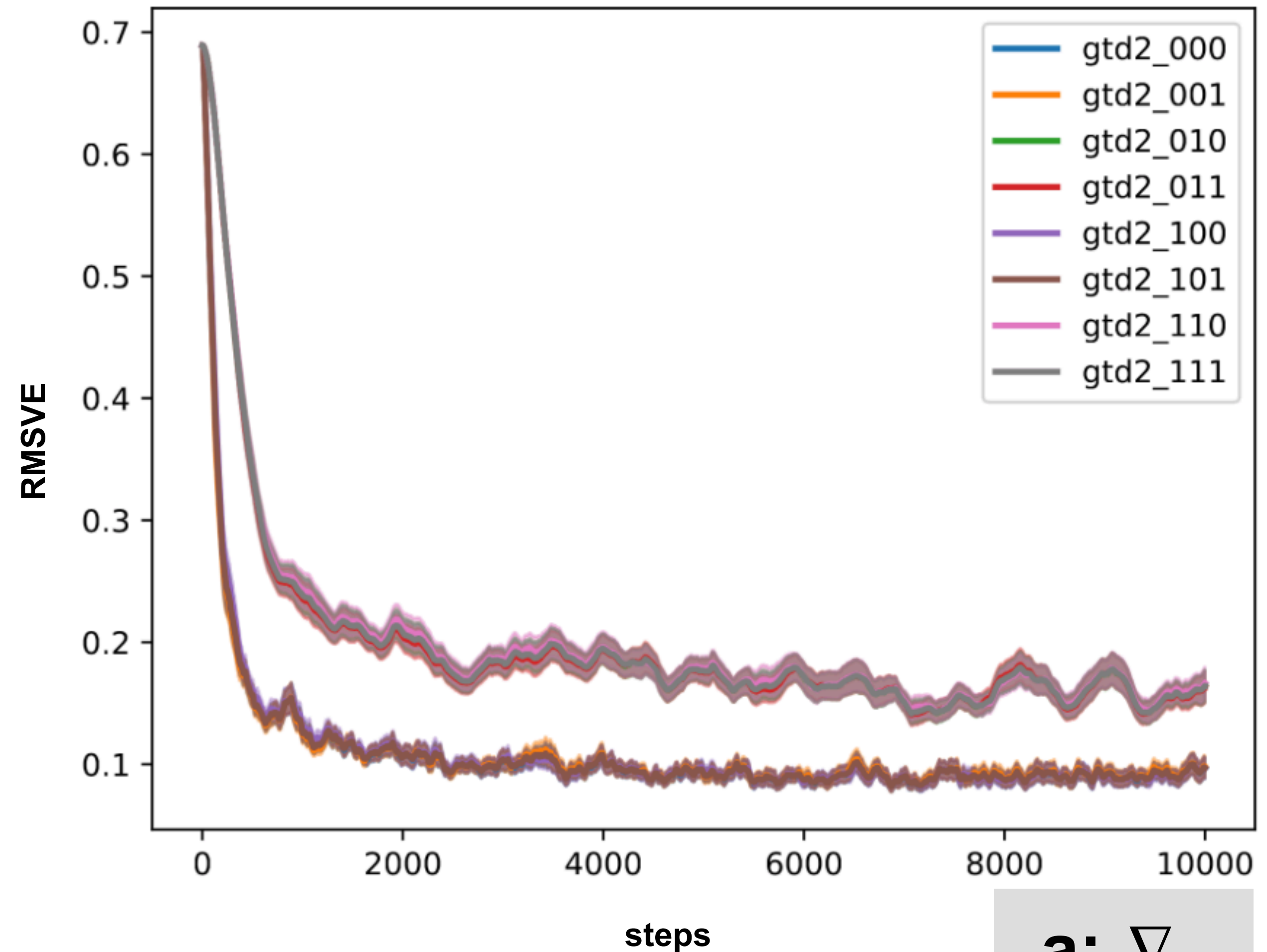
# TDC

0: correct everything  
1: correct as little as possible

# GTD2



**a:**  $\nabla_h$   
**b:**  $\delta_h$   
**c:**  $\delta_w$

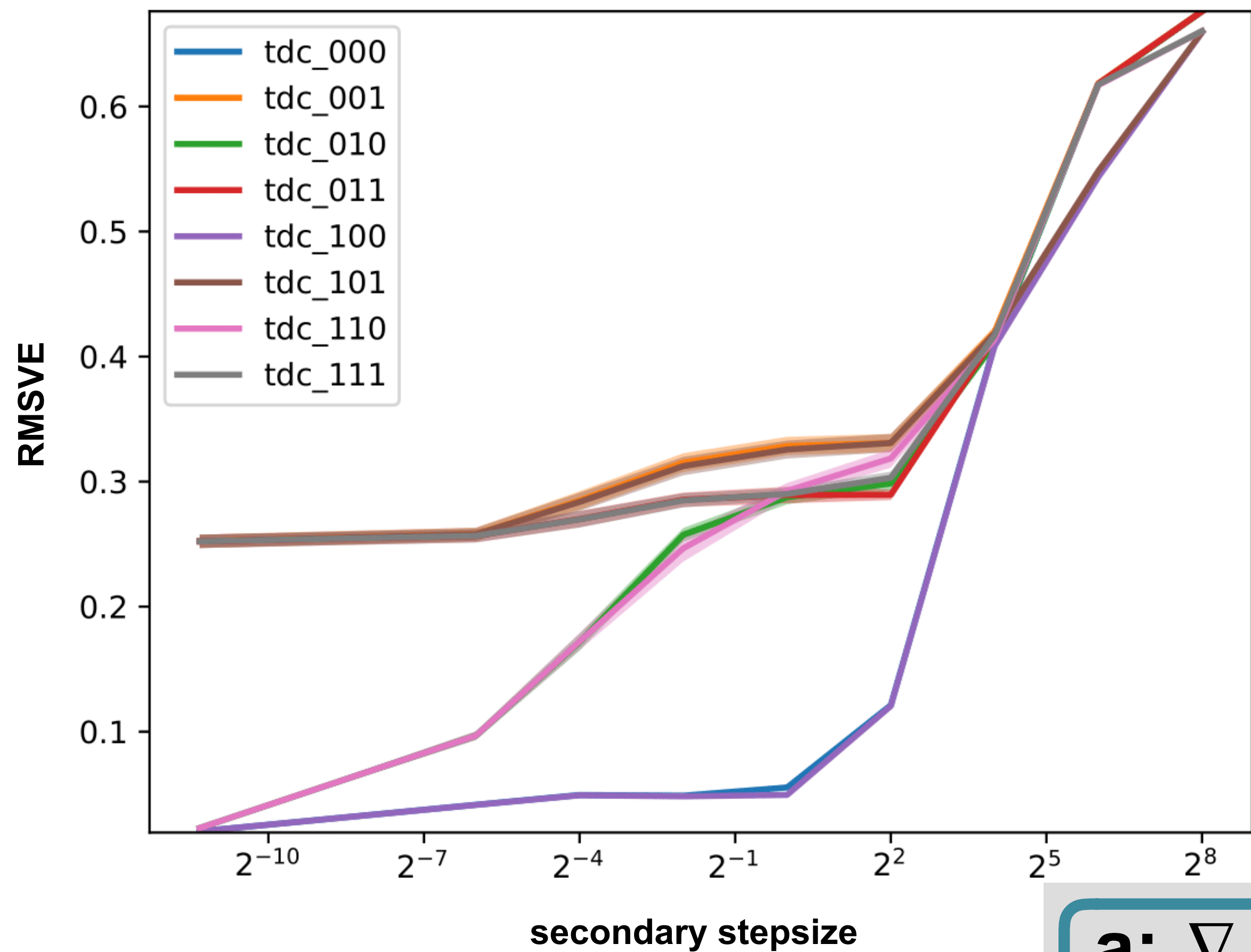


**a:**  $\nabla_h$   
**b:**  $\delta_h$   
**c:**  $\nabla_w$

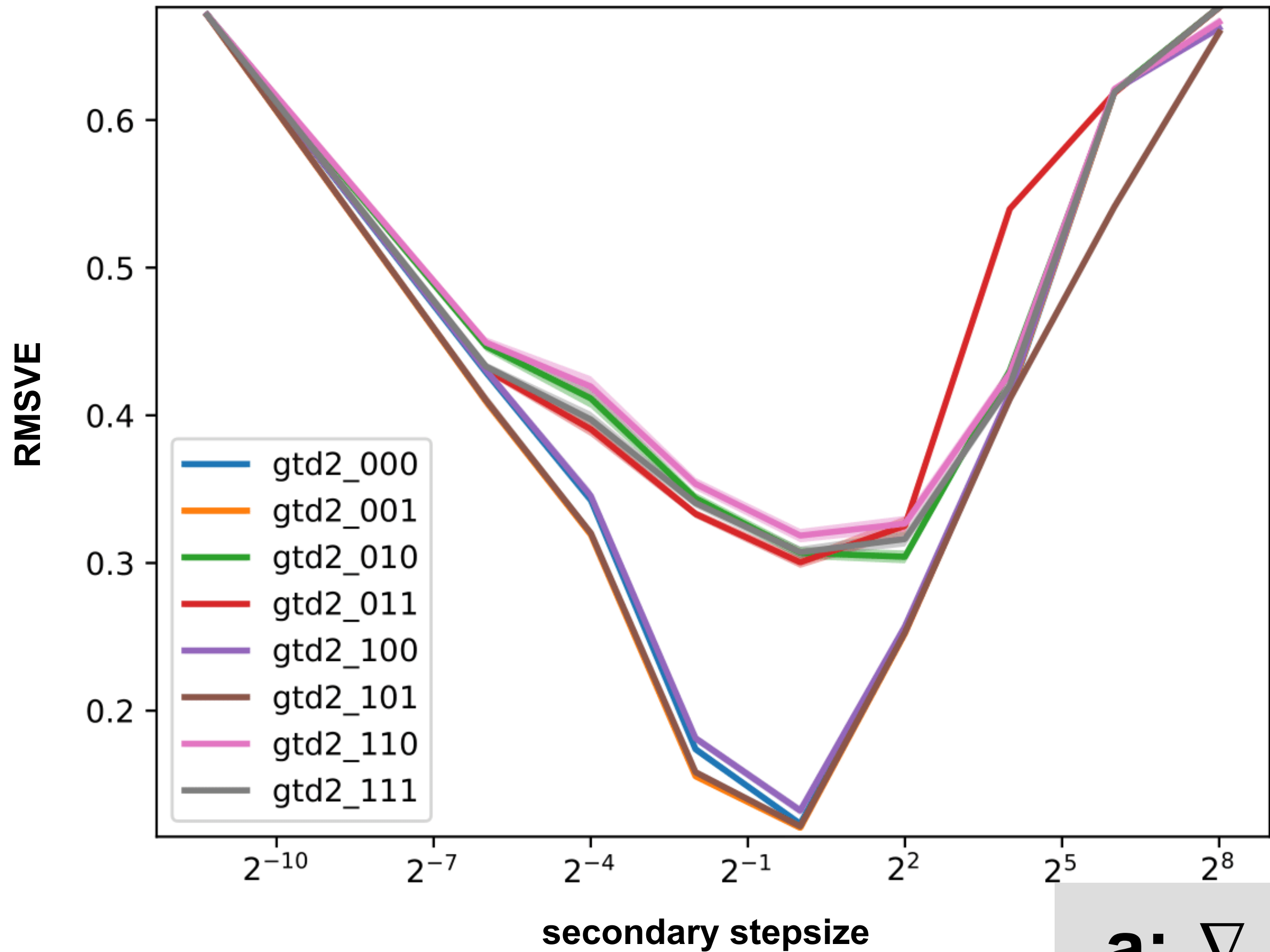
# TDC

0: correct everything  
1: correct as little as possible

# GTD2

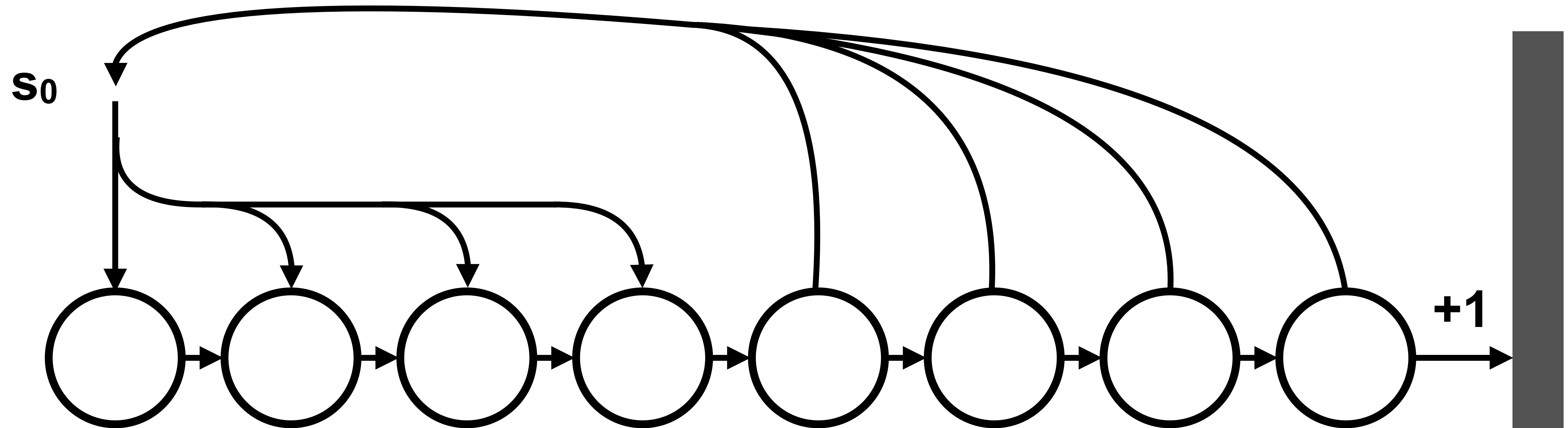


**a:**  $\nabla_h$   
**b:**  $\delta_h$   
**c:**  $\delta_w$



**a:**  $\nabla_h$   
**b:**  $\delta_h$   
**c:**  $\nabla_w$

# Collision



**Target**

**Behavior**

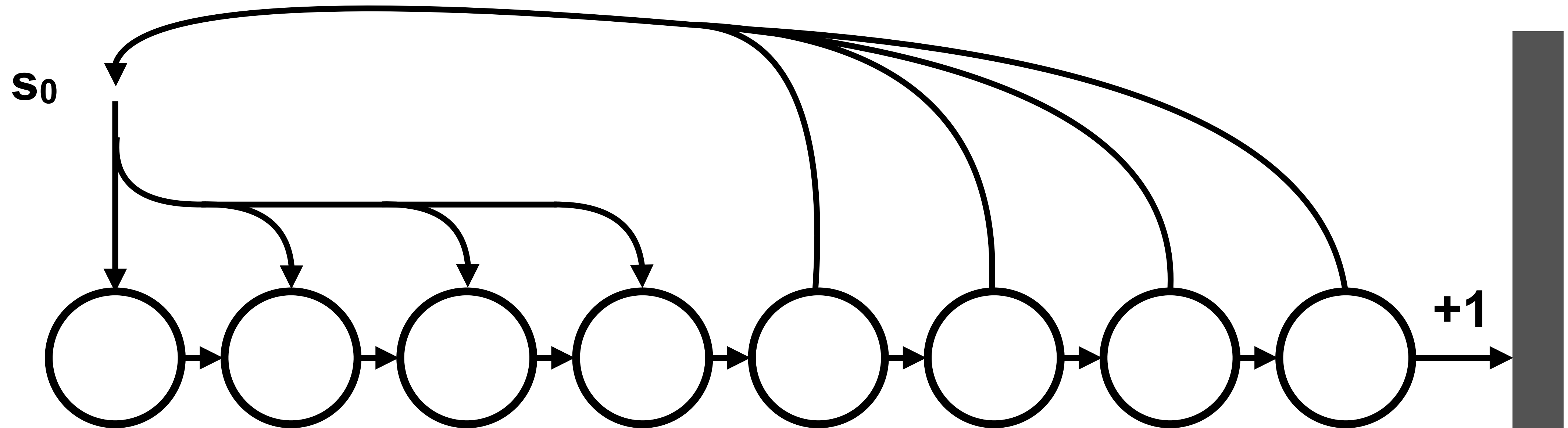
**Right: 100%**

**Right: 50%**

**Retreat: 0%**

**Retreat: 50%**

# Collision



**Target**

**Behavior**

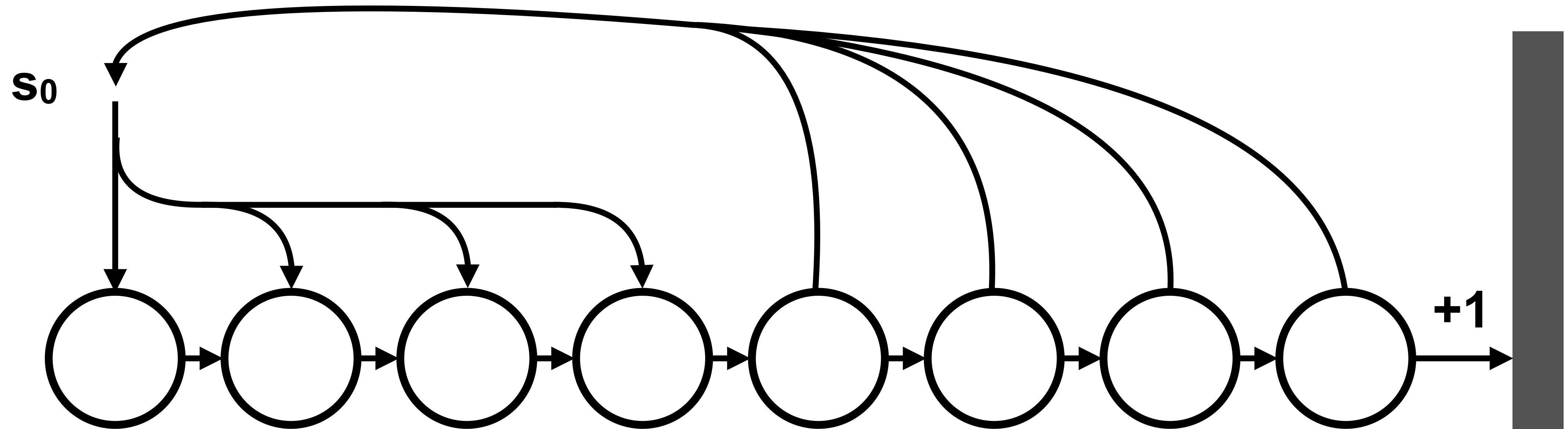
**Right: 100%**

**Right: 50%**

**Retreat: 0%**

**Retreat: 50%**

# Collision



**Target**

**Behavior**

**Right: 100%**

**Right: 25%**

**Retreat: 0%**

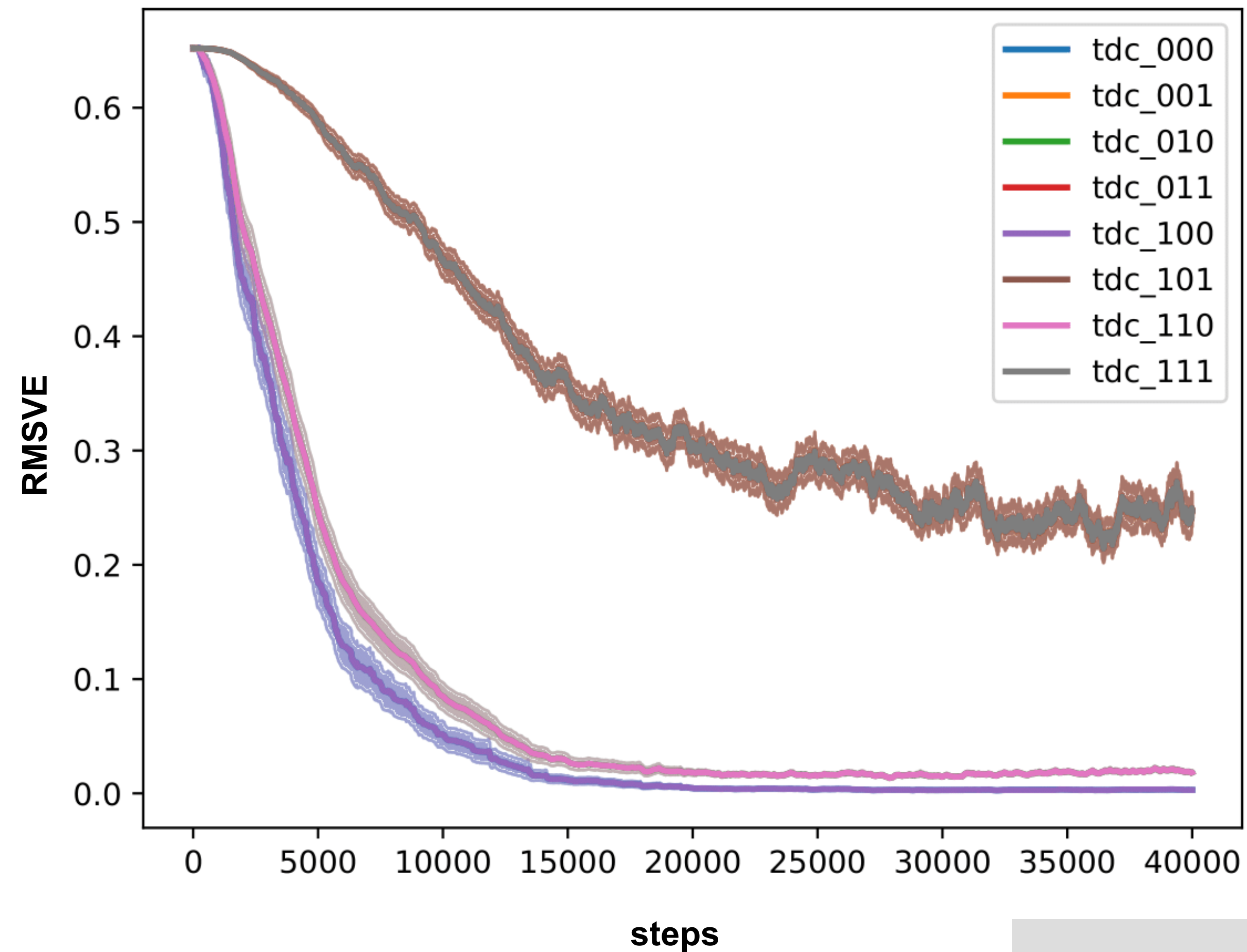
**Retreat: 75%**



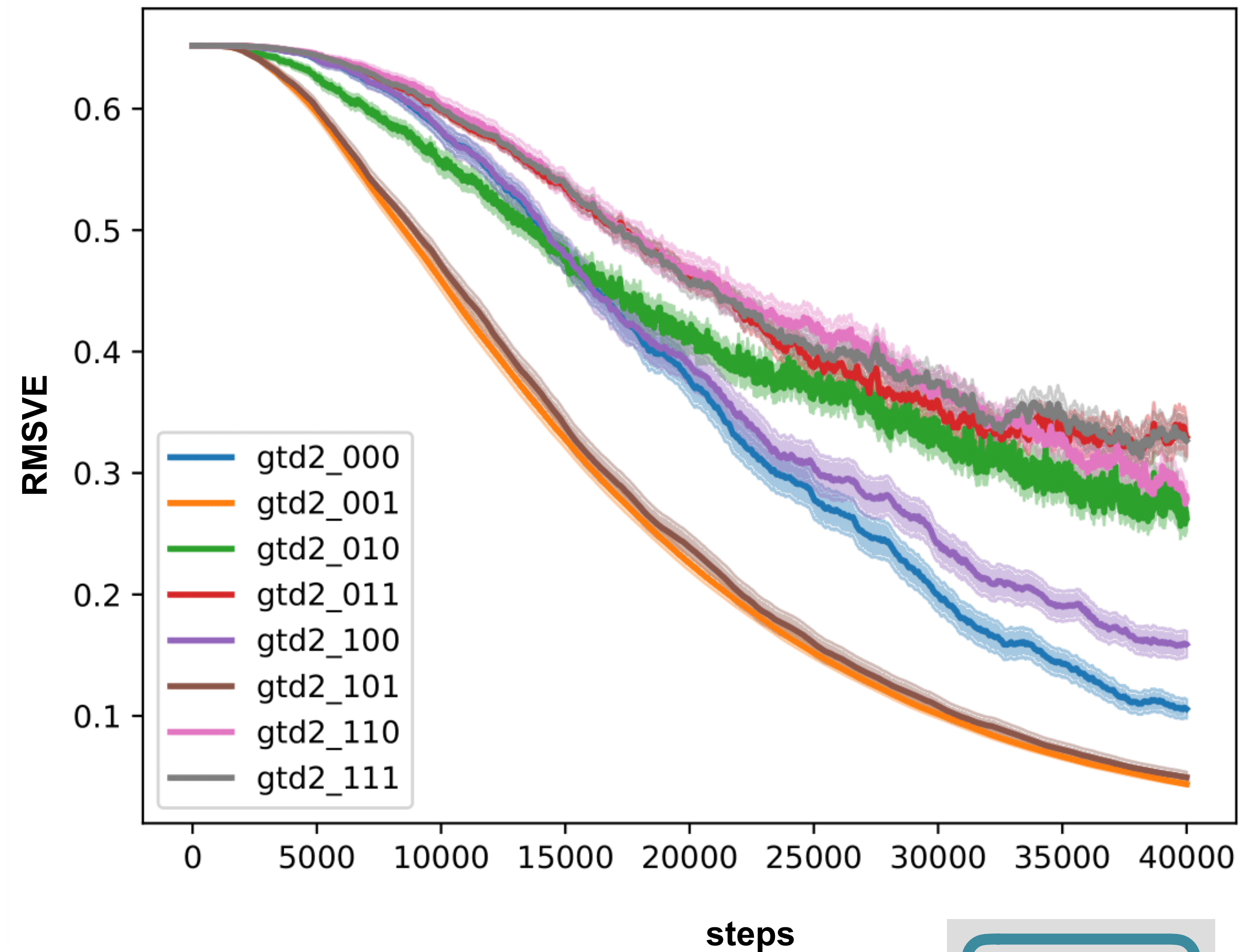
# TDC

0: correct everything  
1: correct as little as possible

# GTD2



**a:**  $\nabla_h$   
**b:**  $\delta_h$   
**c:**  $\delta_w$

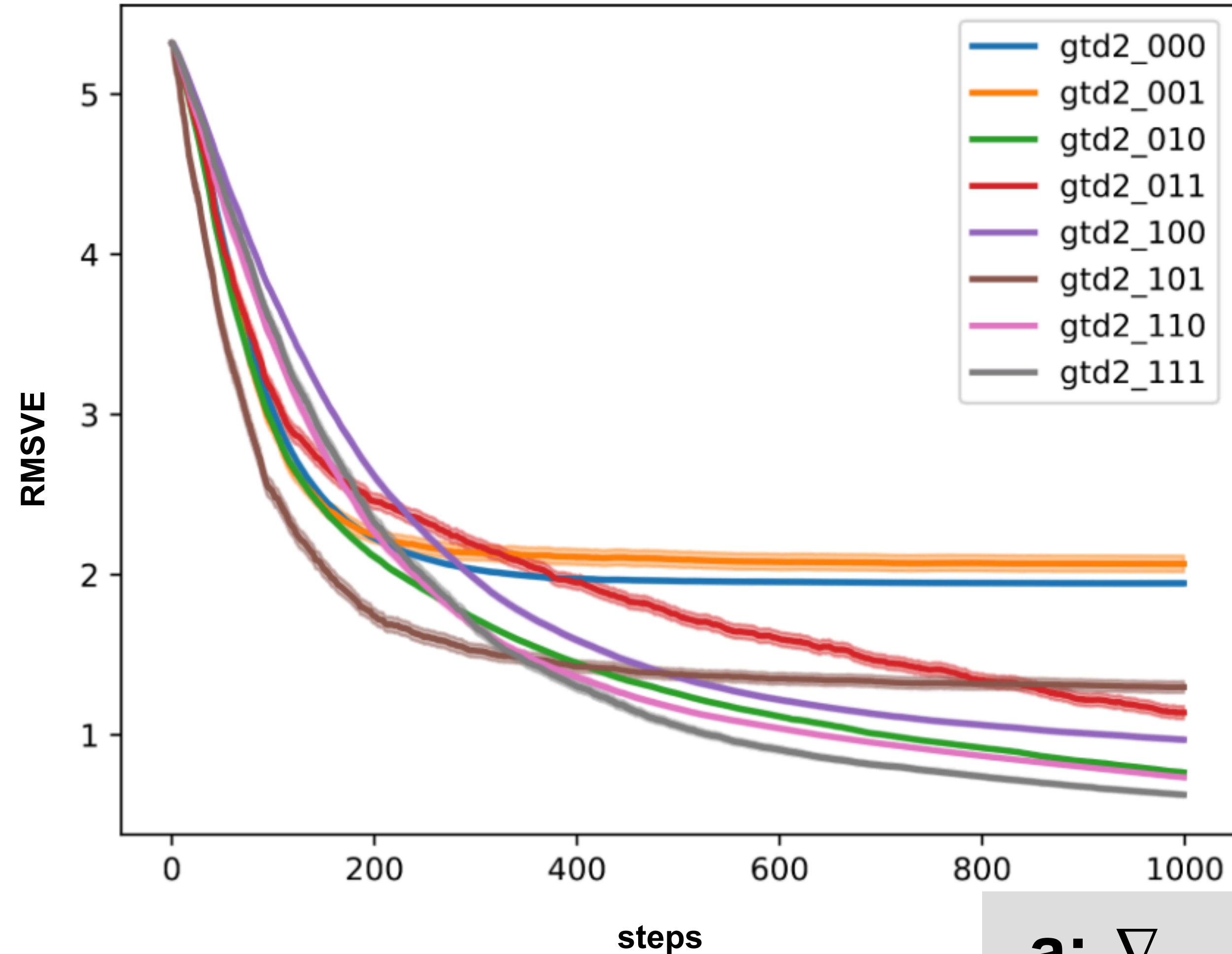
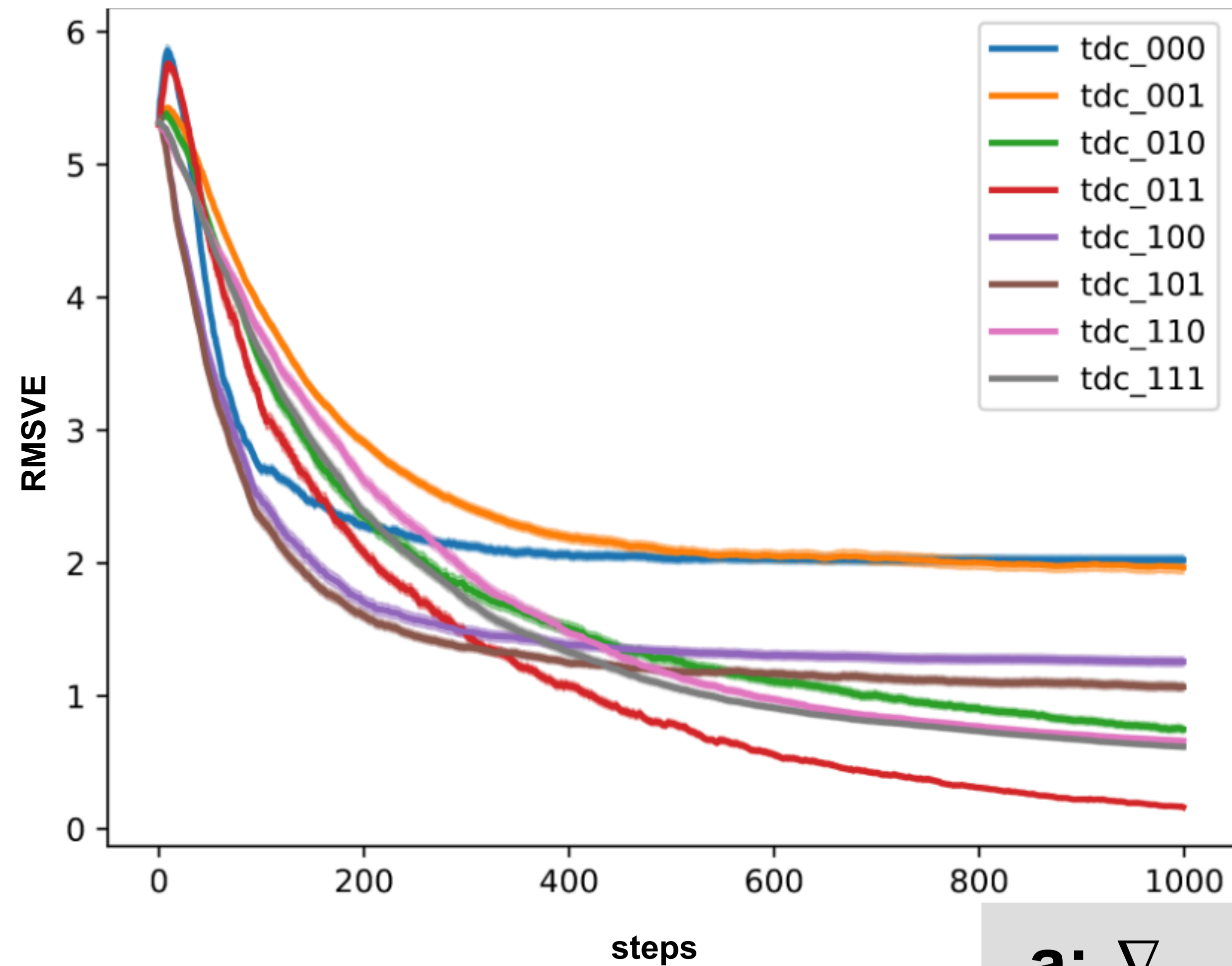


**a:**  $\nabla_h$   
**b:**  $\delta_h$   
**c:**  $\nabla_w$

# TDC

0: correct everything  
1: correct as little as possible

# GTD2



## Baird's Counterexample

Baird, L. C. (1995).

**a:**  $\nabla_h$

**b:**  $\delta_h$

**c:**  $\delta_w$

**a:**  $\nabla_h$

**b:**  $\delta_h$

**c:**  $\nabla_w$

**Thanks for your time**