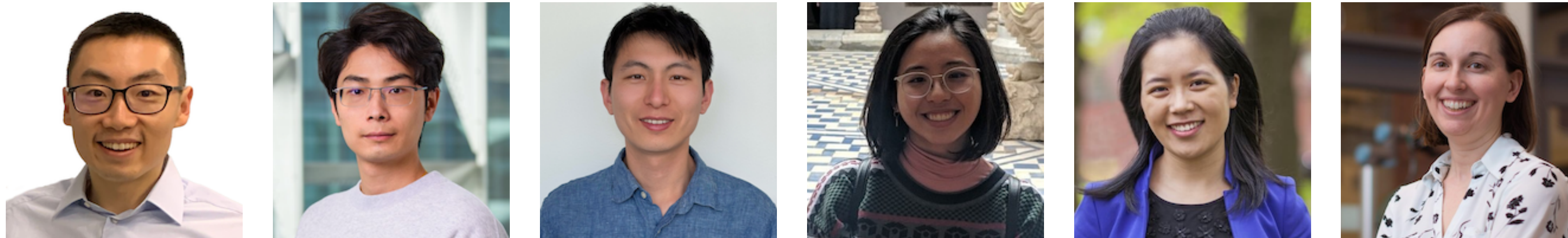


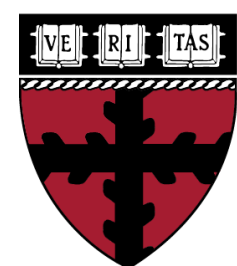
Enhancing Preference-based Linear Bandits via Human Response Time

Shen Li^{*1}, Yuyang Zhang^{*2}, Zhaolin Ren², Claire Liang¹, Na Li², Julie A. Shah¹



¹ MIT ² Harvard

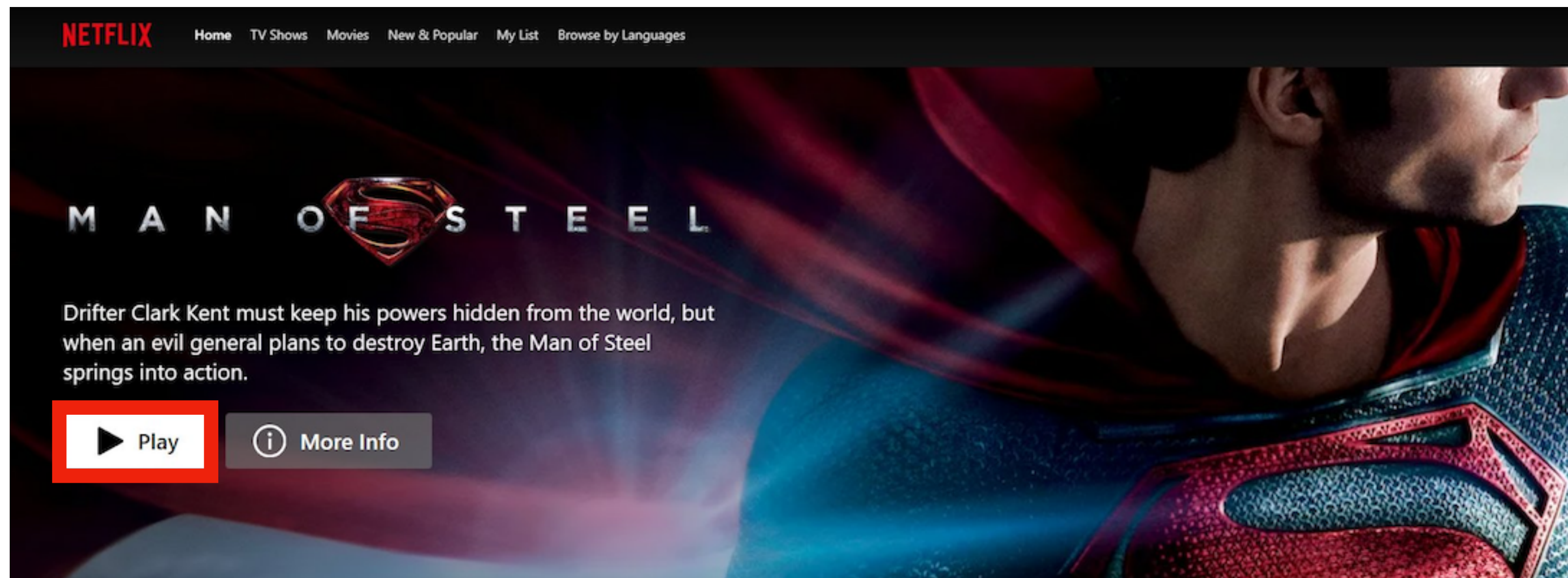
* First two authors have equal contribution



Harvard John A. Paulson
School of Engineering
and Applied Sciences



Binary Choices Are Widely Used to Align AI Systems With Human Preferences

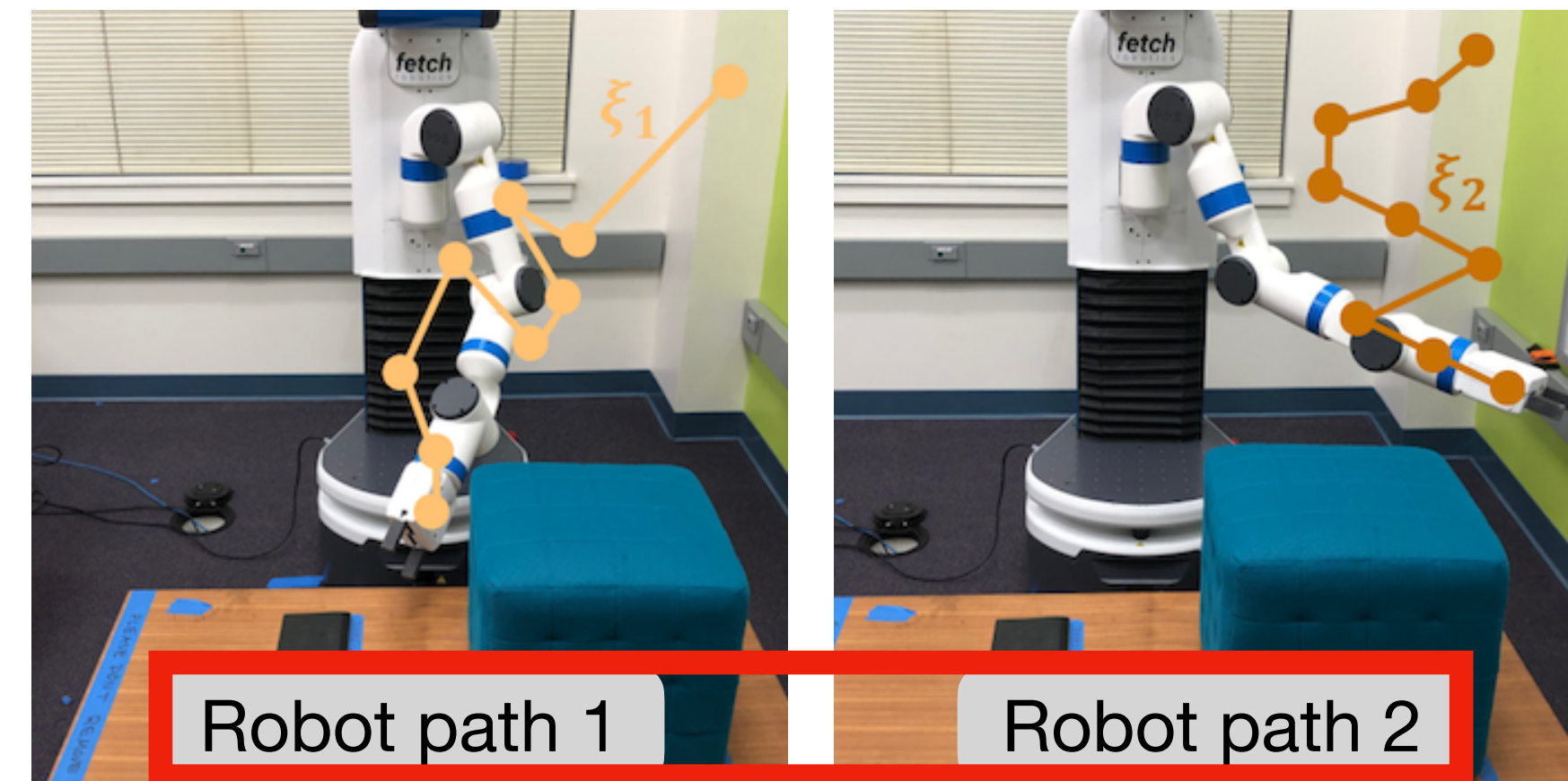


Q: why do our eyebrows and eyelashes stop growing after a certain length?

A: Hair has different cycles[1]. It grows for a certain amount of time before falling out and starting fresh. For eyelashes, the cycle is much shorter (a few weeks) whereas eyebrows grow more slowly and have a cycle closer to two months. The hair on your head has a cycle of several years, allowing it to grow much longer. Some people have longer arm or leg hair because they have a faster speed of growth or a longer cycle. Look it up, it's really interesting

[1] <https://www.philipkingsley.co.uk> > hair-science > hair-g.

Plausible answer? Yes No Not sure



(Bıyık et al., 2022)

Choices Provide Limited Information About Preference Strengths. Our Work Resolves This by Incorporating Response Times.

Which one would you like during the poster session?



Long response time



Weak preference

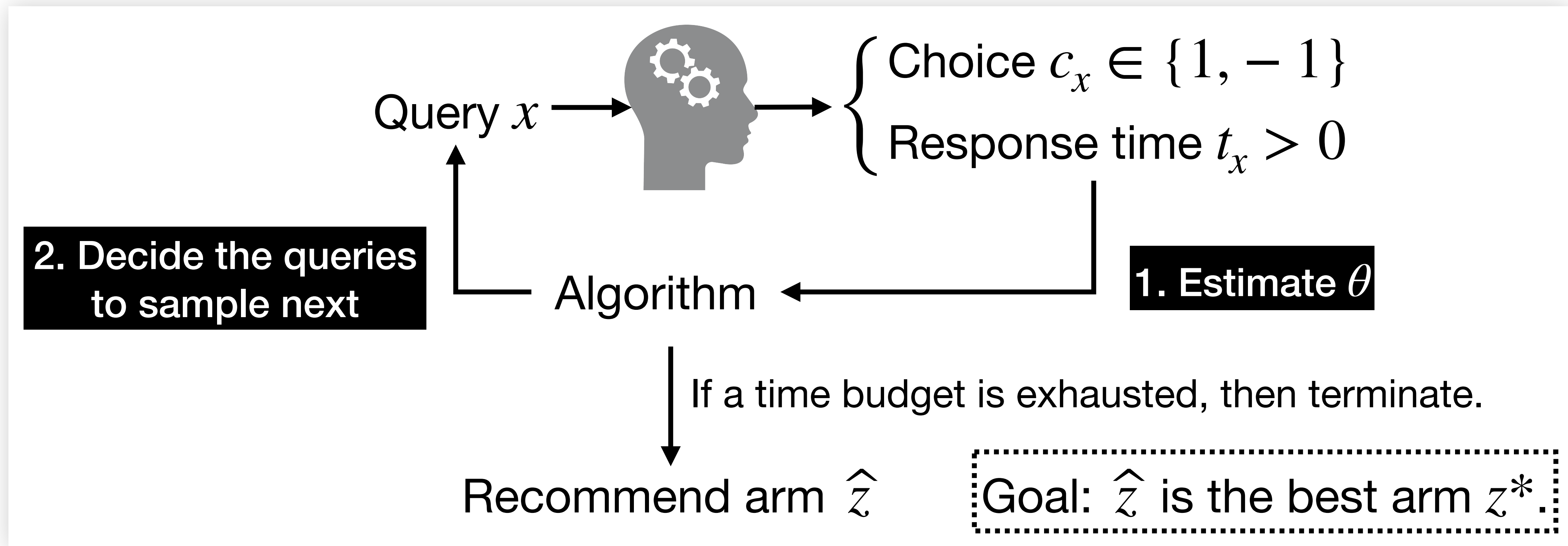
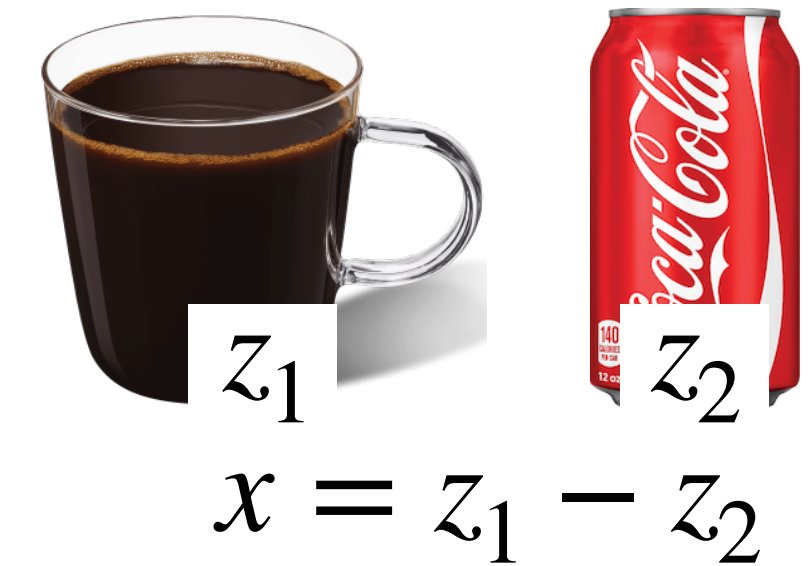
Short response time



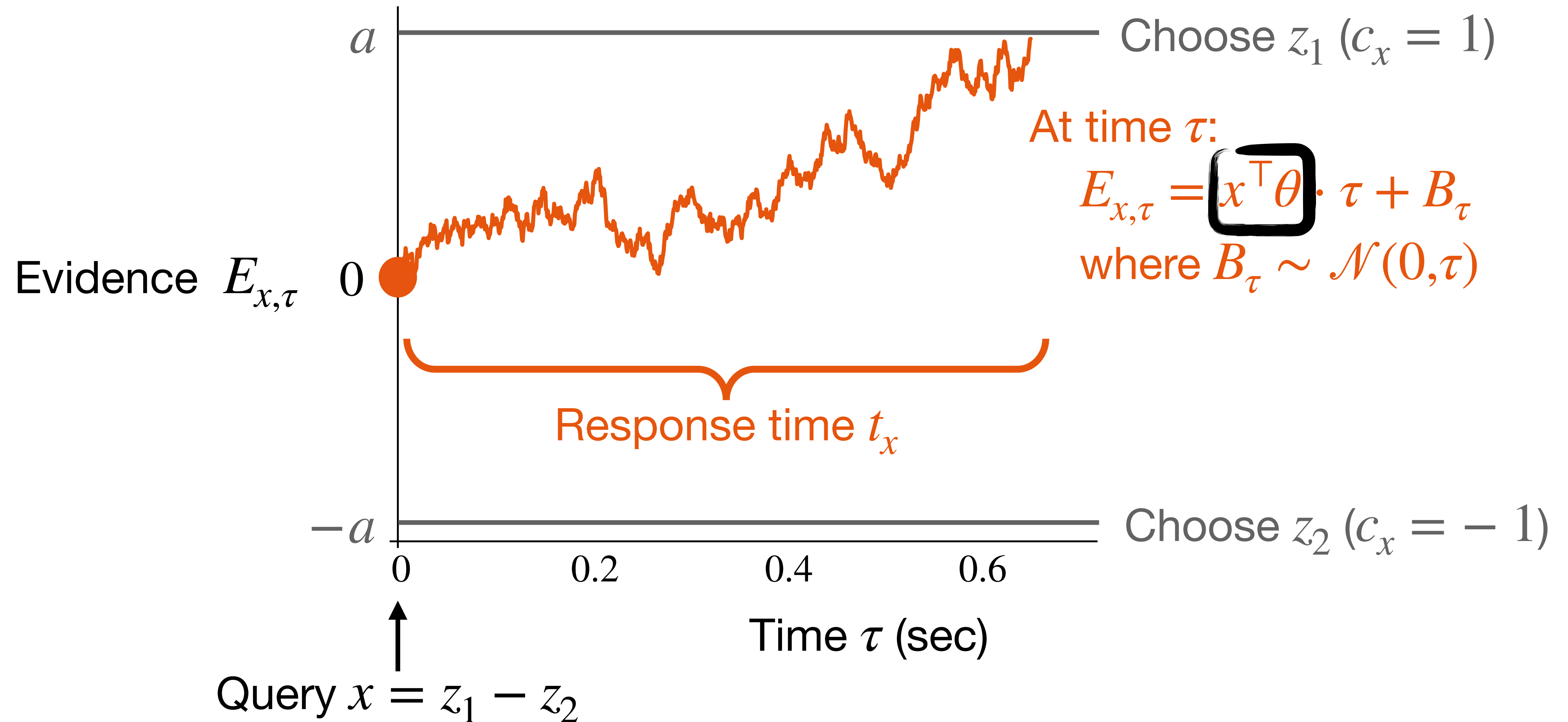
Strong preference

Problem Formulation: Linear Bandit With Binary Choice and Response Time Feedback

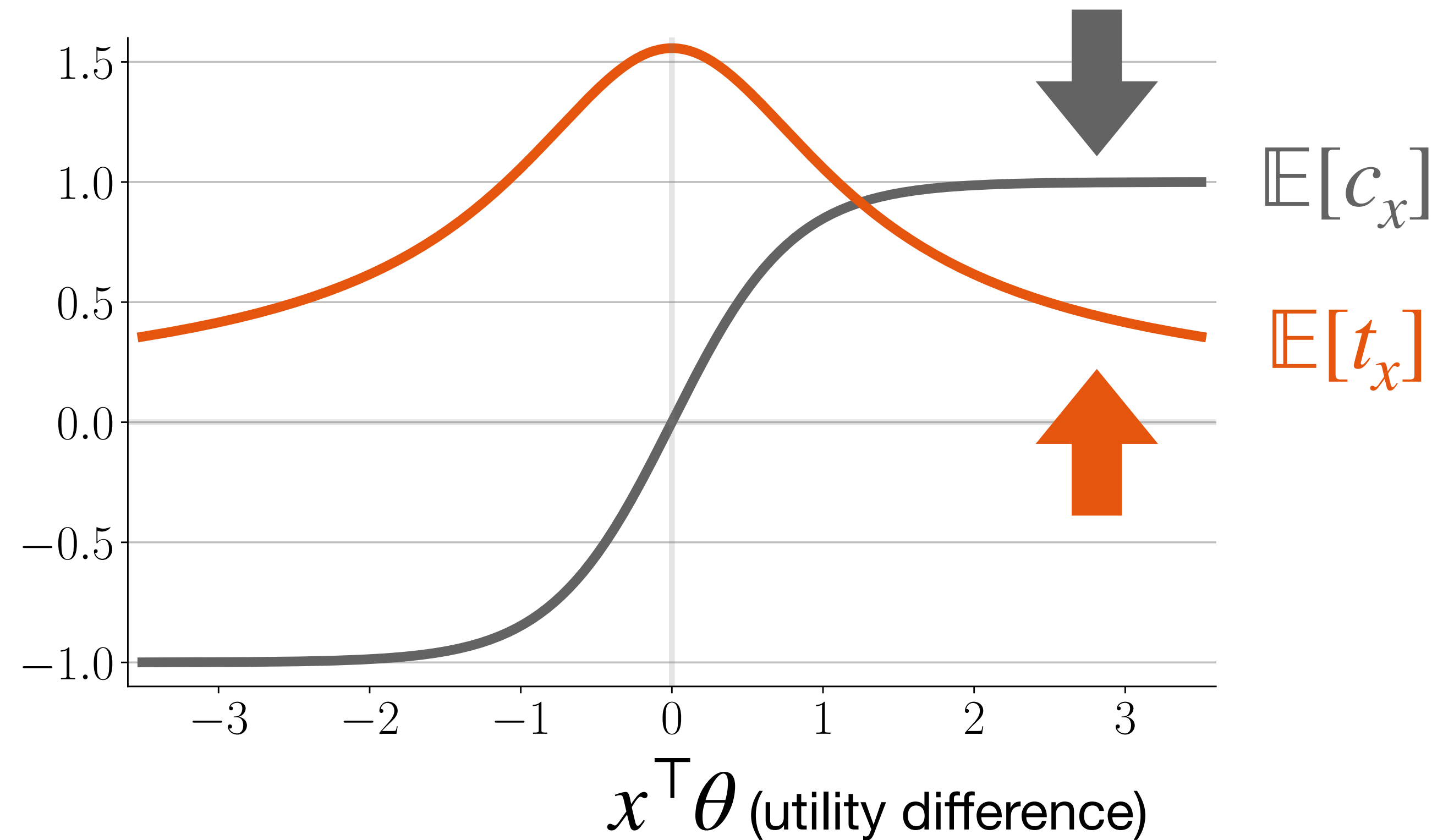
- Human preference: $\theta \in \mathbb{R}^d$
- Each arm: z with a utility $z^\top \theta$
- Each query: $x := z_1 - z_2$ with a utility difference $x^\top \theta$



EZ-Diffusion Model Links Utility Differences, Choices, and Response Times



The Magnitude of Utility Difference Is Proportional to Expected Choice and Inversely Proportional to Expected Response Time



$x^T \theta$ (utility difference)

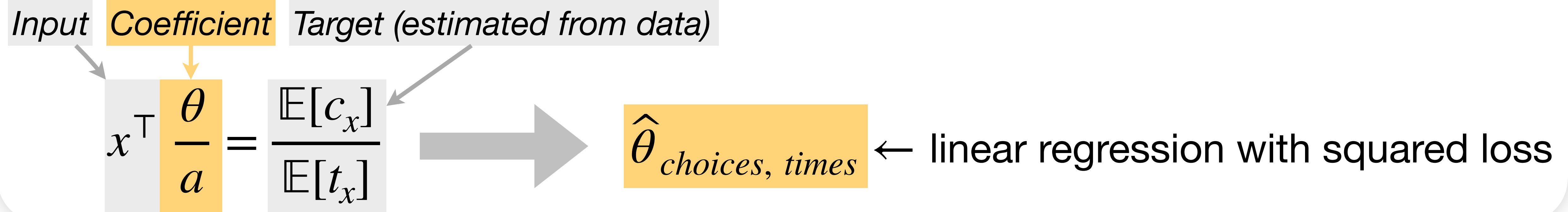
Strong pref Weak pref Strong pref



Our Novel θ Estimator Uses Choices and Response Times, While Prior Methods Only Use Choices

Given a dataset with i.i.d. choices and response times from various queries:

Our estimator uses both choices and response times:



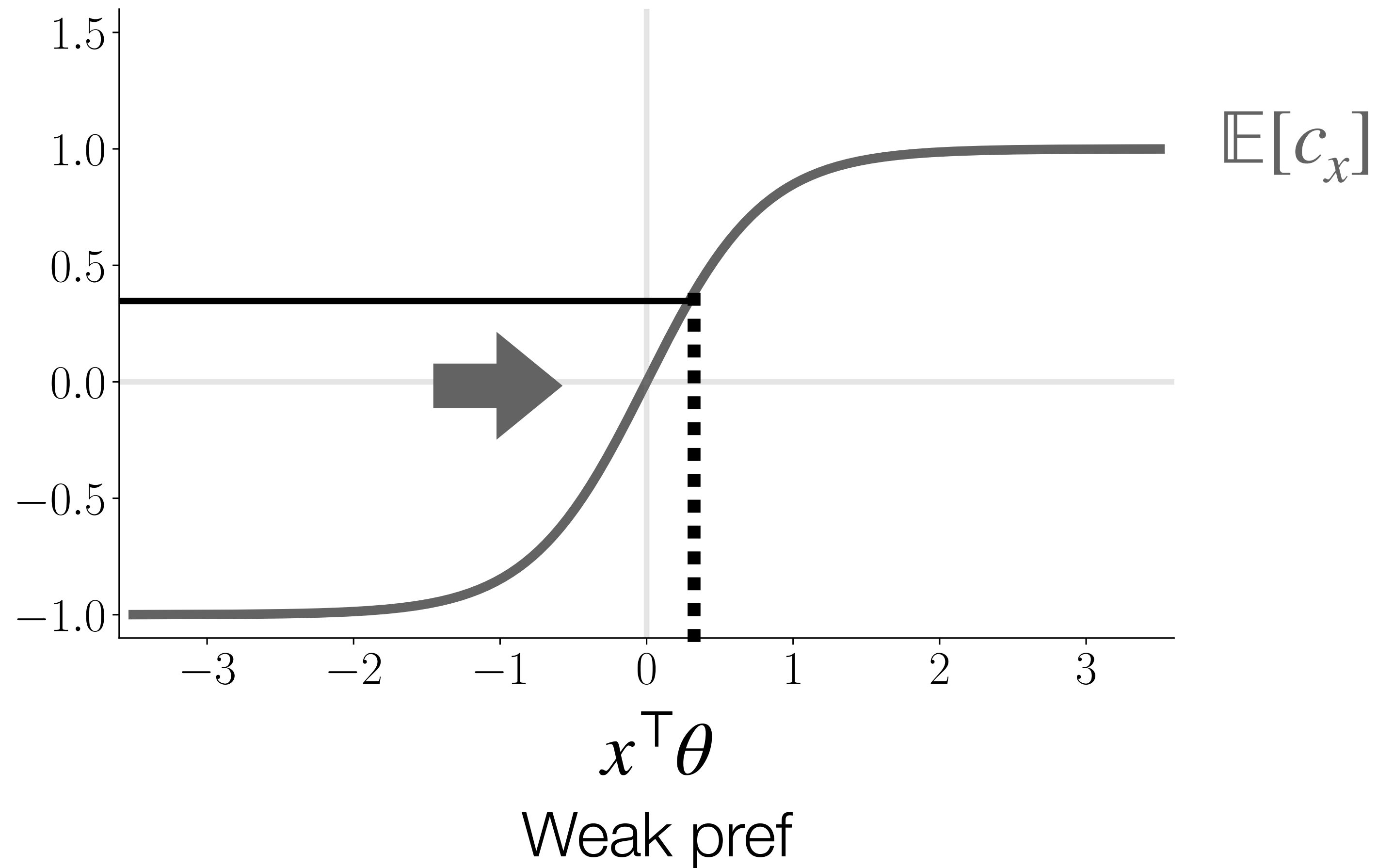
Prior methods only use choices:

$$\mathbb{P}[c_x = 1] = \frac{1}{1 + \exp(-c_x \cdot x^\top \cdot 2a\theta)}$$

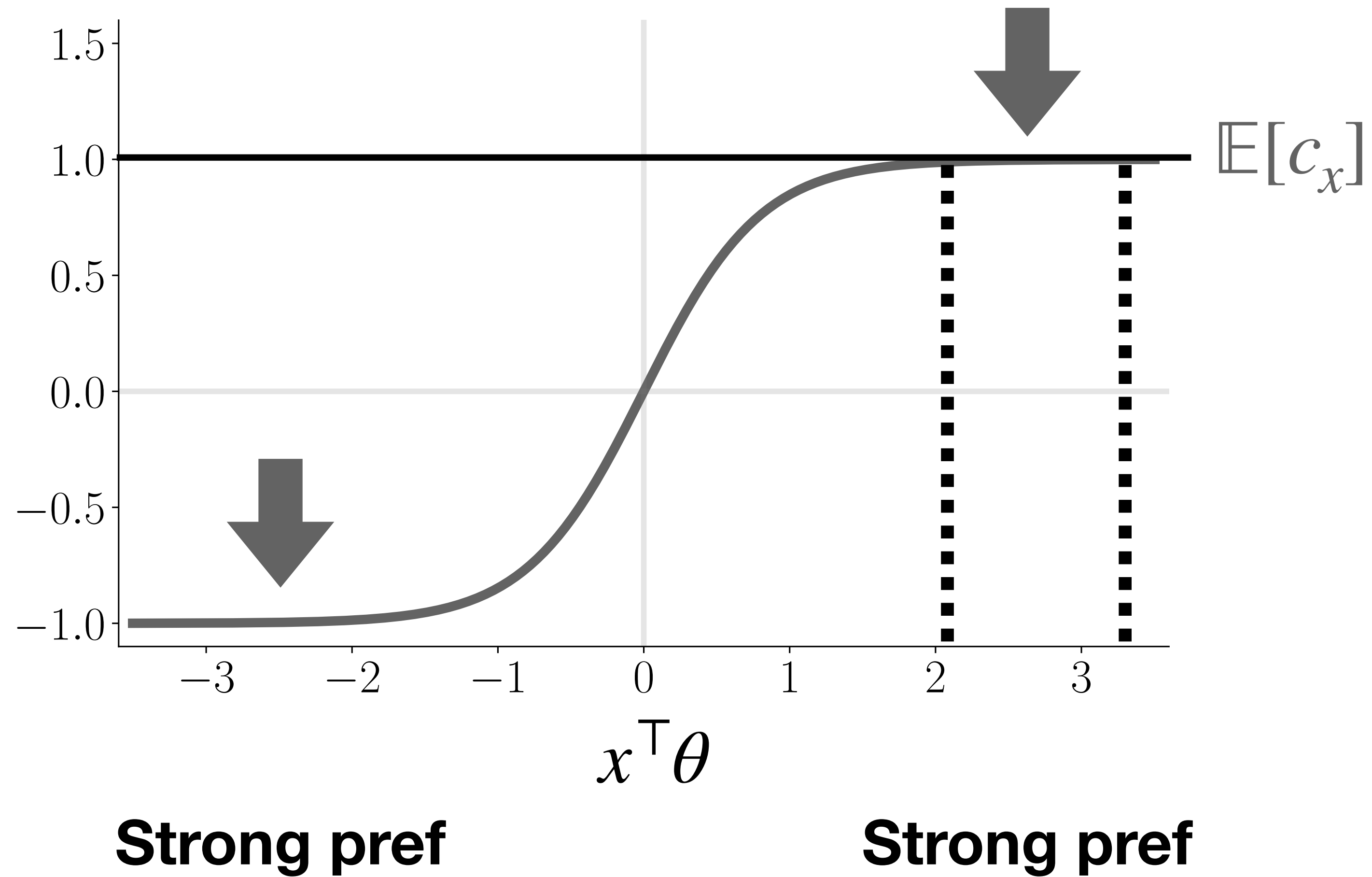
$\hat{\theta}_{choices}$ ← logistic regression

(Same as the Bradley Terry model)

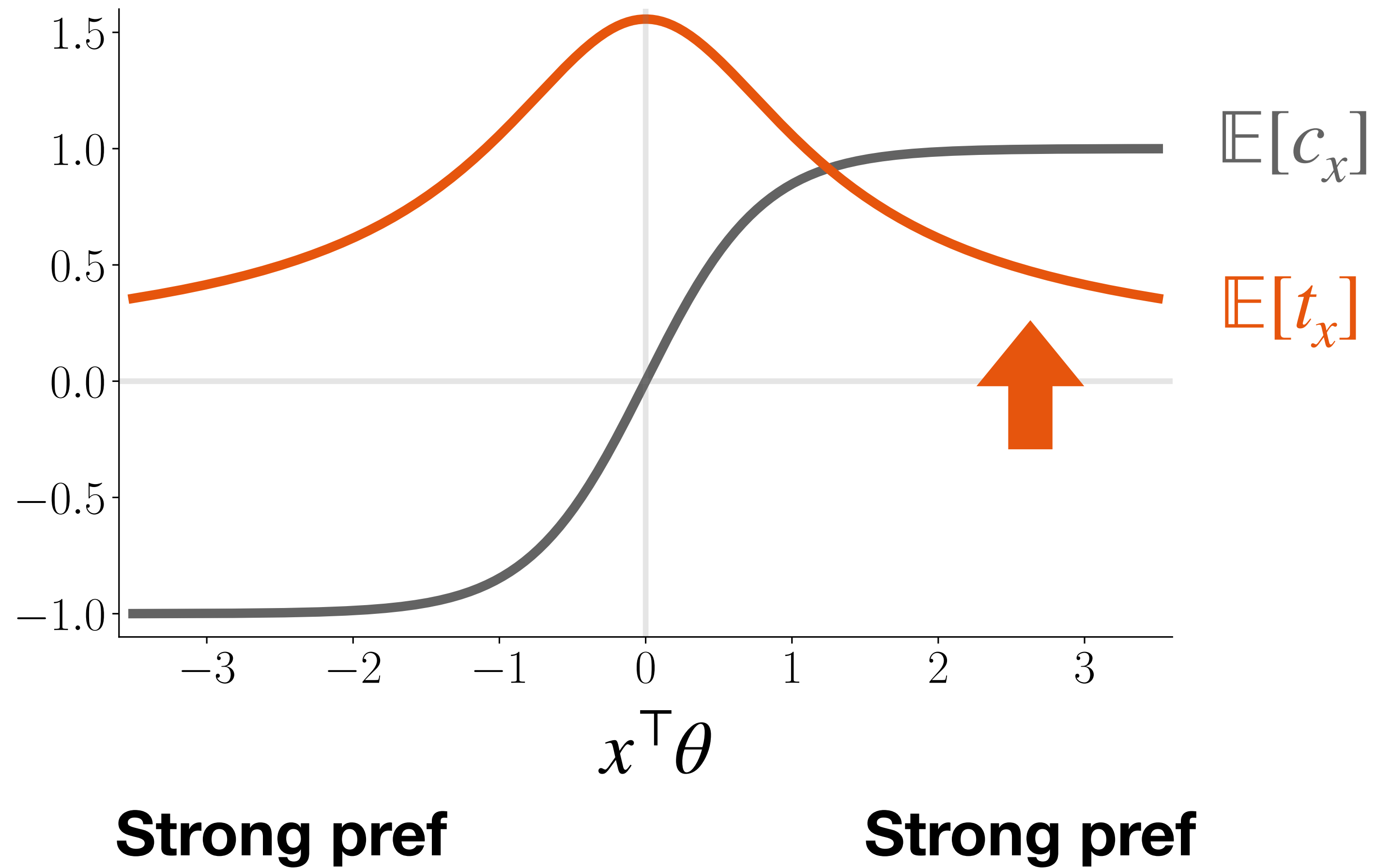
Intuitively, for Queries With *Strong* Preferences, Response Times Provide Information That Complements Choices



Intuitively, for Queries With *Strong* Preferences, Response Times Provide Information That Complements Choices



Intuitively, for Queries With *Strong* Preferences, Response Times Provide Information That Complements Choices

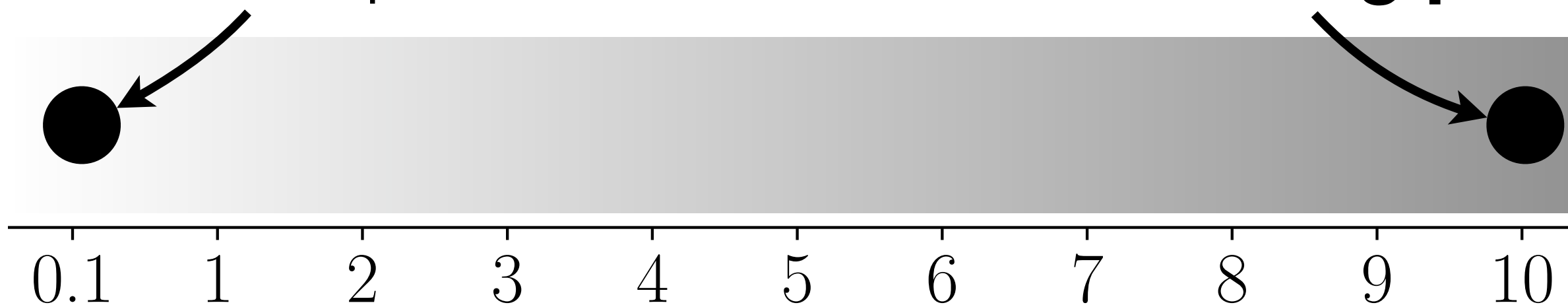


To Verify This Insight, We Use a Synthetic Problem to Compare Estimator Performance

Scaling factor c that scales each arm z to $c \cdot z$

A bandit problem with weak pref

A bandit problem with strong pref



1. Computes an experimental design:



2. Sample 50 queries and gather feedback

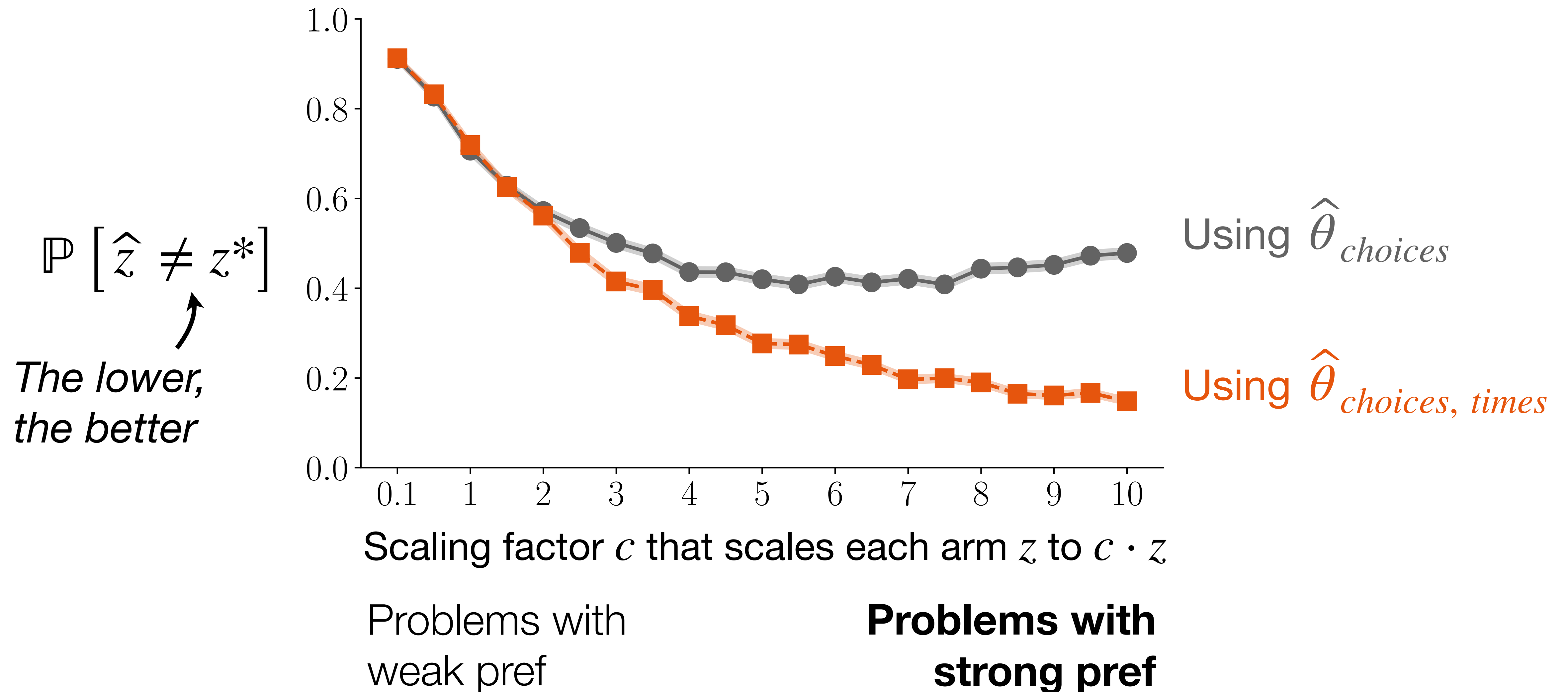


3. Estimate $\hat{\theta}_{choices, times}$ or $\hat{\theta}_{choices}$ and recommend \hat{z}



Performance measure: $\mathbb{P} [\hat{z} \neq z^*]$

Empirical Result Confirms That for Queries With *Strong* Preferences, Response Times Provide Information That Complements Choices



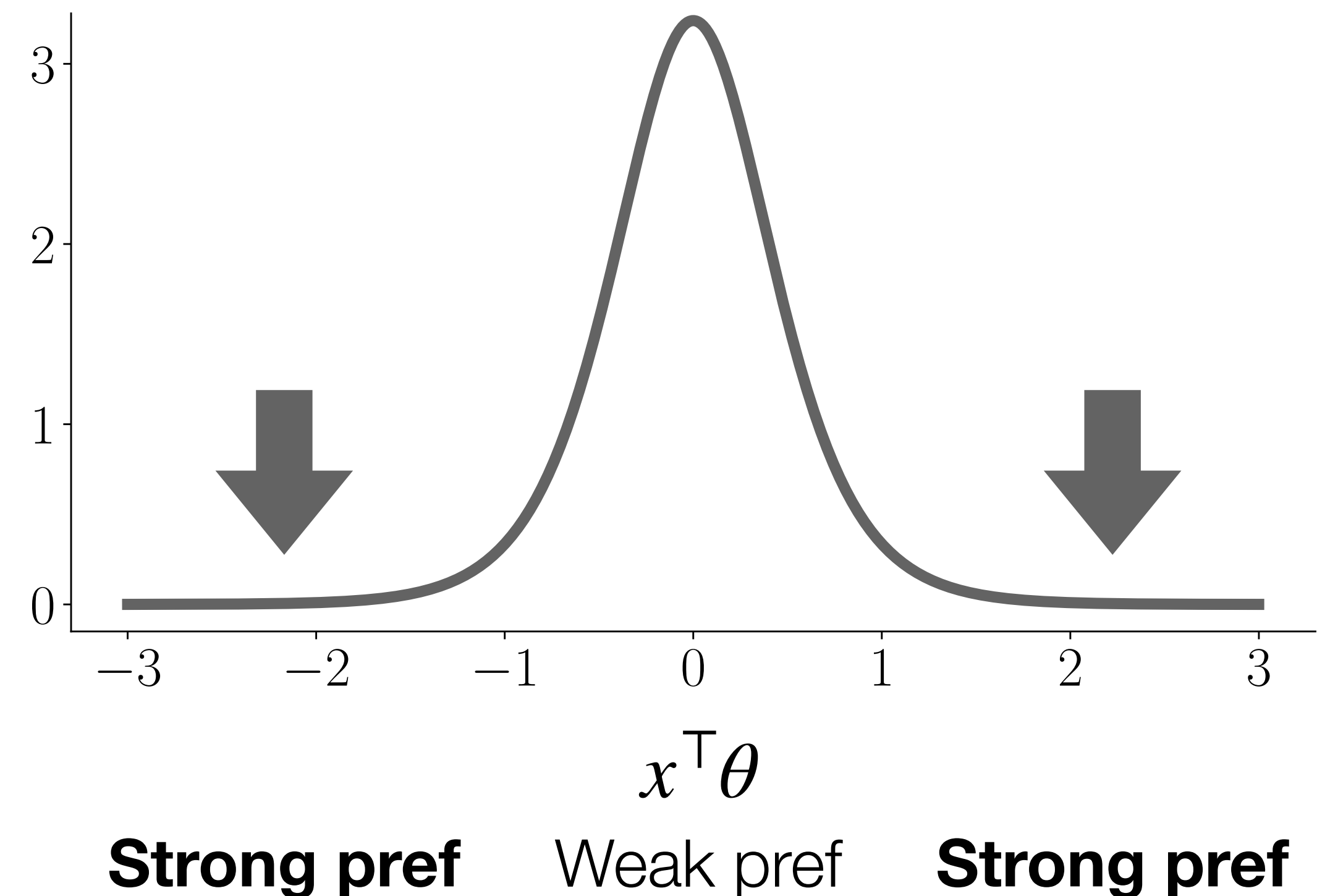
Asymptotic Variances Shows That for Queries With *Weak* Preferences, Choices Provide a Limited Amount of Information

Given a fixed dataset that contains n choices and response times for each query in \mathcal{X} , then, for each arm z , the utility estimation error satisfies $\sqrt{n} \left(z^\top \hat{\theta} - z^\top \theta \right) \xrightarrow{D} \mathcal{N} \left(0, A\text{Var}_z \right)$.

If using $\hat{\theta}_{choices}$, then

$$A\text{Var}_z = z^\top \left(\sum_{x \in \mathcal{X}} a^2 \text{Var}[c_x] x x^\top \right)^{-1} z$$

The highlighted term
(The higher, the better)



Asymptotic Variances Confirms That for Queries With *Strong* Preferences, Response Times Provide Information That Complements Choices

Given a fixed dataset that contains n choices and response times for each query in \mathcal{X} , then, for each arm z , the utility estimation error satisfies $\sqrt{n} \left(z^\top \hat{\theta} - z^\top \theta \right) \xrightarrow{D} \mathcal{N} \left(0, A\text{Var}_z \right)$.

If using $\hat{\theta}_{\text{choices}}$, then

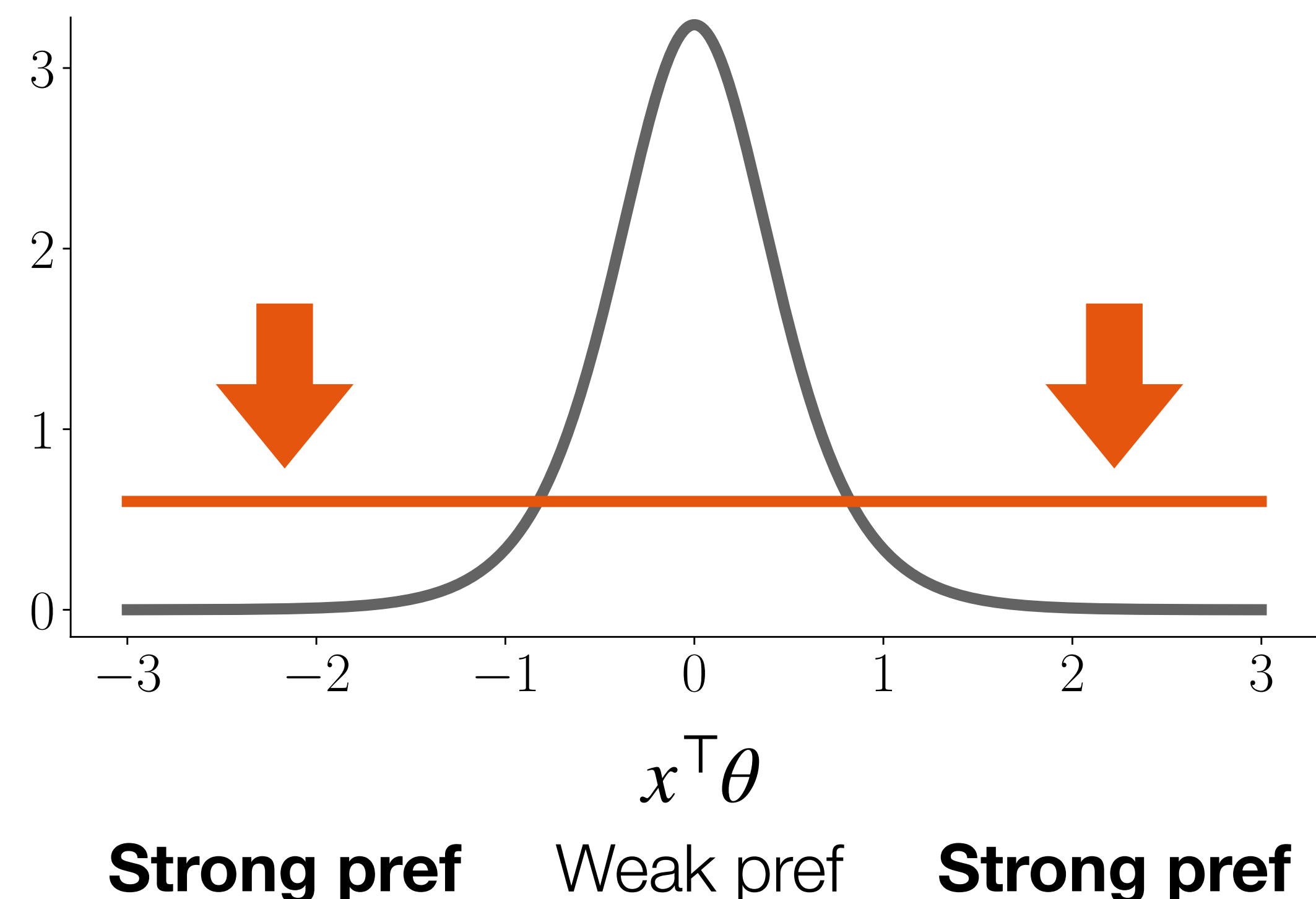
$$A\text{Var}_z = z^\top \left(\sum_{x \in \mathcal{X}} a^2 \text{Var}[c_x] x x^\top \right)^{-1} z$$

If using $\hat{\theta}_{\text{choices, times}}$, then

$$A\text{Var}_z \leq z^\top \left(\sum_{x \in \mathcal{X}} \min_{\tilde{x} \in \mathcal{X}} \mathbb{E}[t_{\tilde{x}}] x x^\top \right)^{-1} z$$

The highlighted term
(The higher, the better)

Example: assume each $x^\top \theta \in [-3, 3]$:



This Plot Does Not Provide Definitive Conclusions for Comparing Response Times and Choices for Queries With *Weak* Preferences

Given a fixed dataset that contains n choices and response times for each query in \mathcal{X} , then, for each arm z , the utility estimation error satisfies $\sqrt{n} \left(z^\top \hat{\theta} - z^\top \theta \right) \xrightarrow{D} \mathcal{N} \left(0, A\text{Var}_z \right)$.

If using $\hat{\theta}_{\text{choices}}$, then

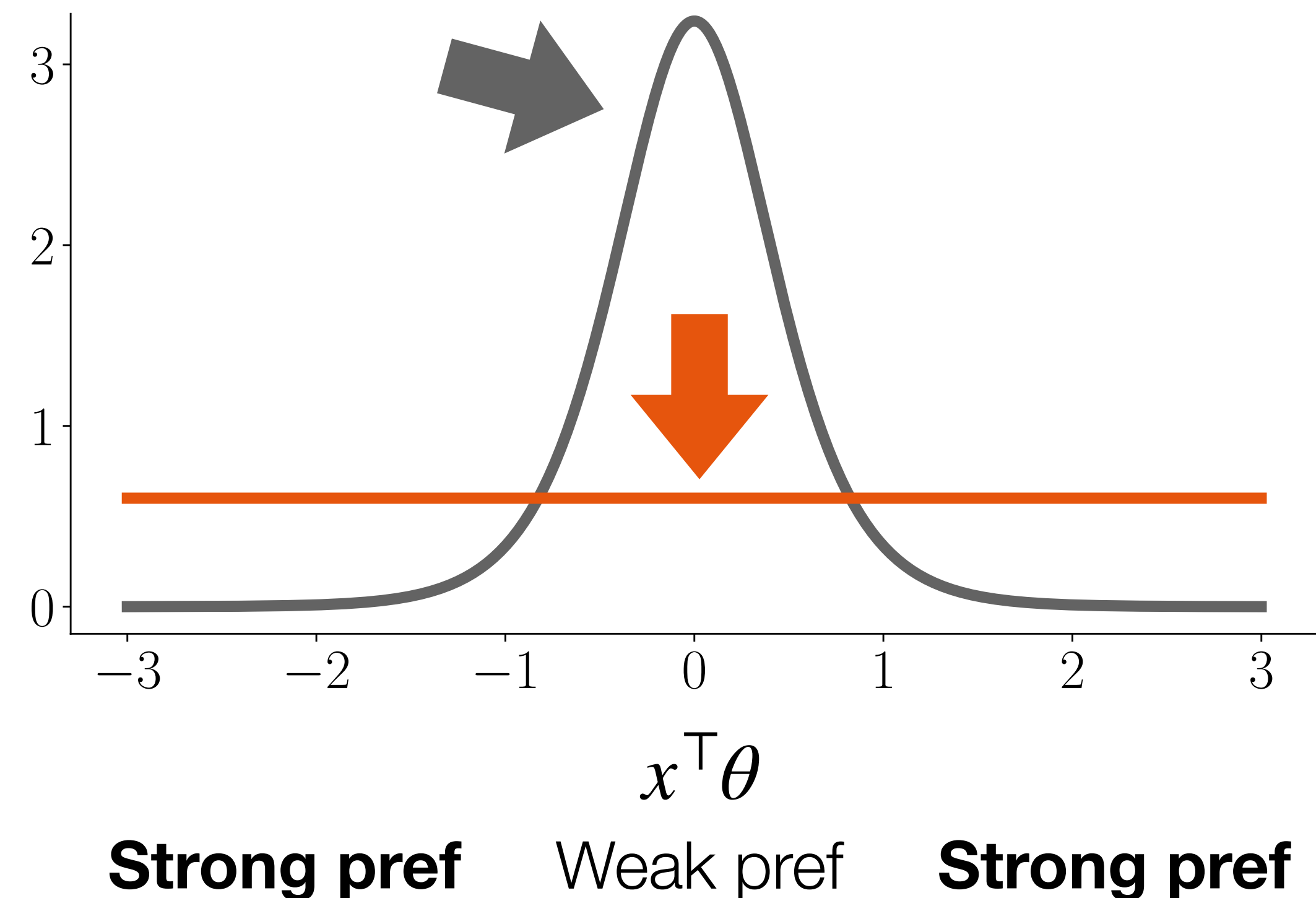
$$A\text{Var}_z = z^\top \left(\sum_{x \in \mathcal{X}} a^2 \text{Var}[c_x] x x^\top \right)^{-1} z$$

If using $\hat{\theta}_{\text{choices, times}}$, then

$$A\text{Var}_z \leq z^\top \left(\sum_{x \in \mathcal{X}} \min_{\tilde{x} \in \mathcal{X}} \mathbb{E}[t_{\tilde{x}}] x x^\top \right)^{-1} z$$

The highlighted term (The higher, the better)

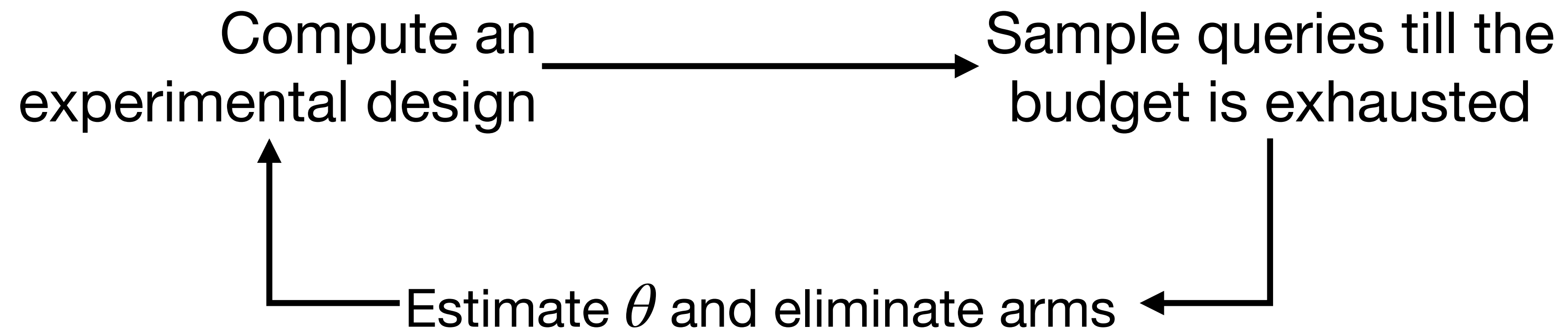
Example: assume each $x^\top \theta \in [-3, 3]$:



Integrating Both Estimators Into the Generalized Successive Elimination Algorithm to Identify the Best Arm Within a Fixed Time Budget

- Split the total budget evenly into multiple phases.

- For each phase:



- Recommended the remaining arm \hat{z} .

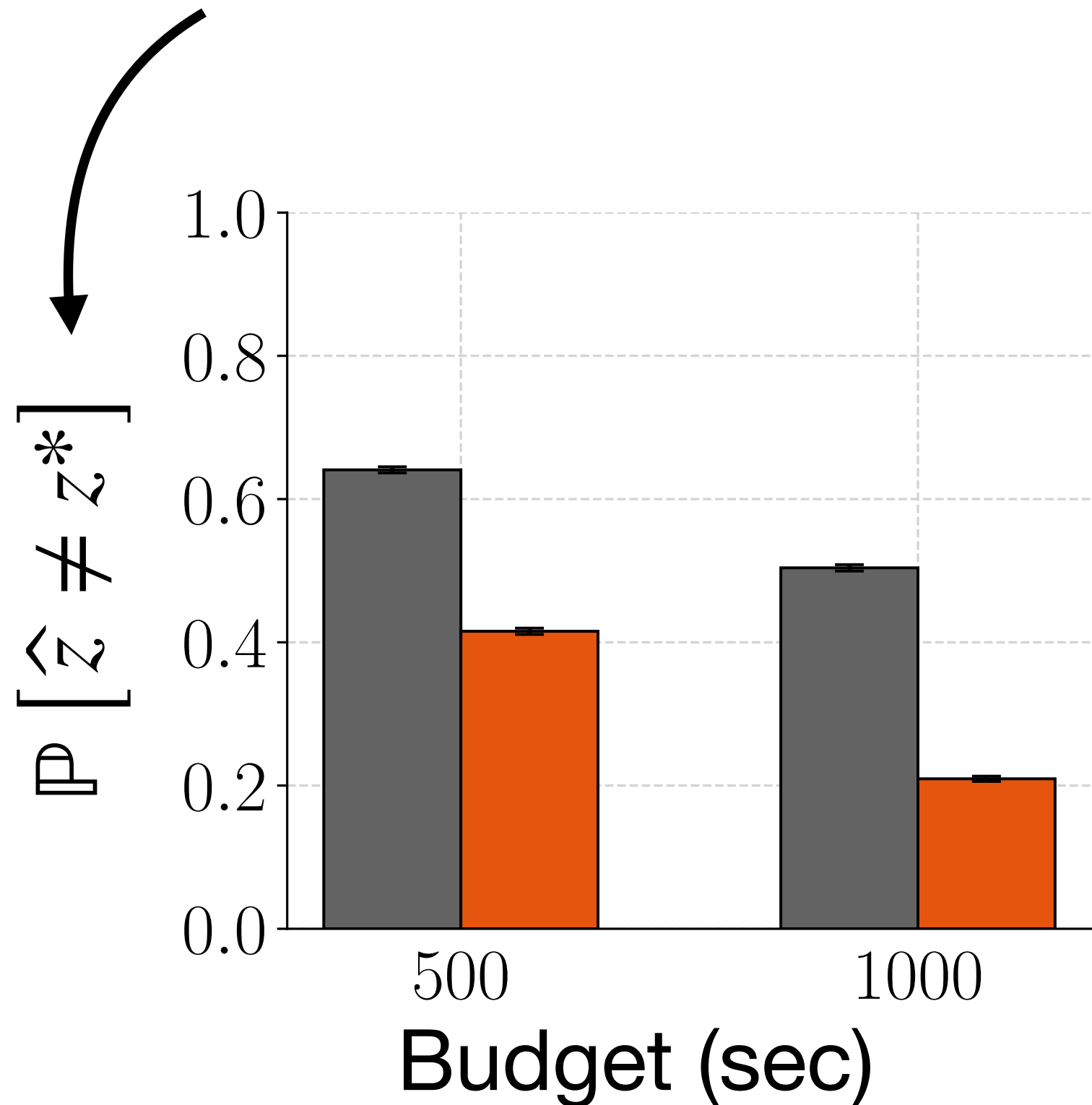
Performance measure: $\mathbb{P} \left[\hat{z} \neq z^* \right]$

Empirical Result of Bandit Learning Shows That Incorporating Response Times Reduces Learning Errors

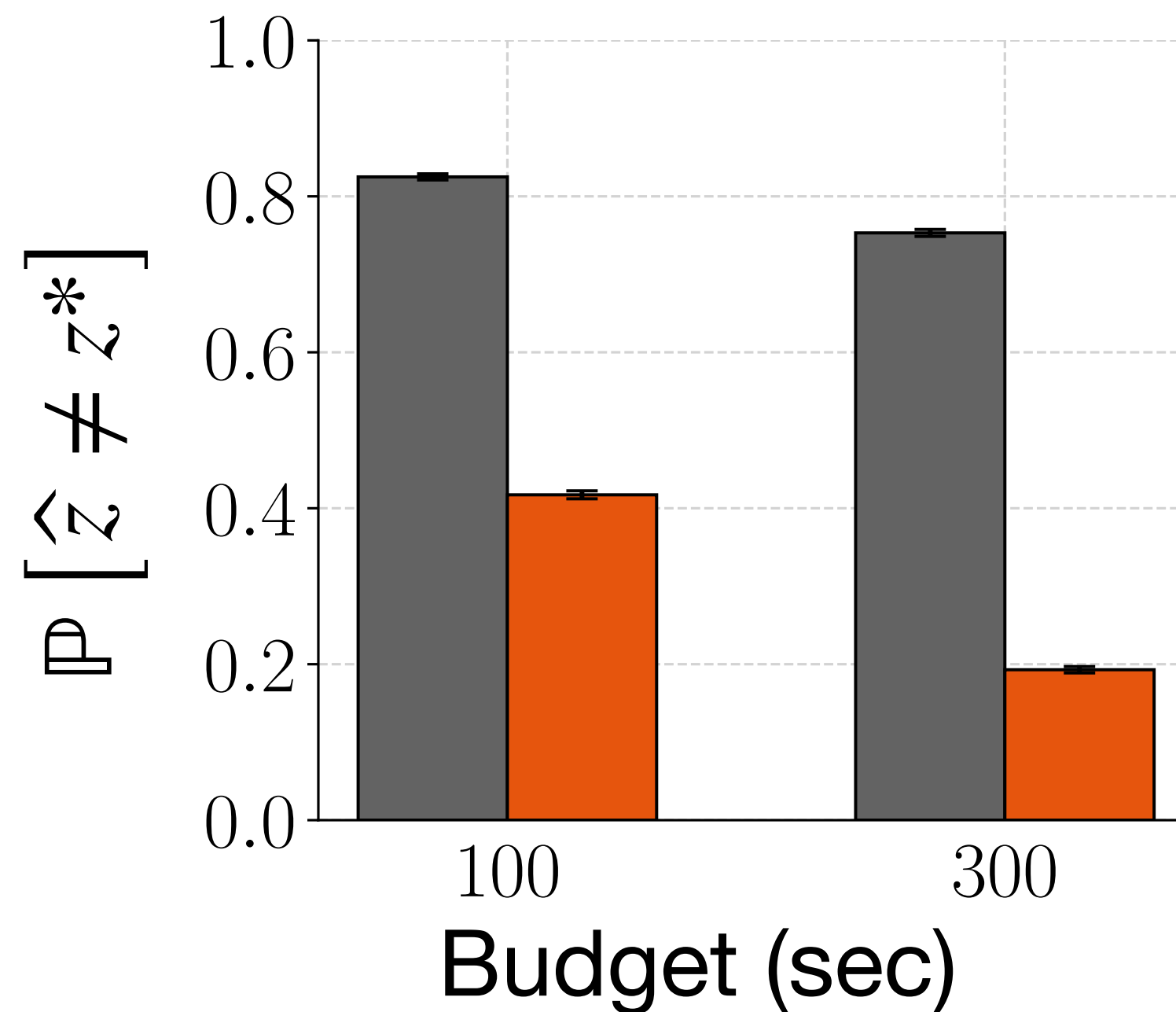
Generalized Successive Elimination with $\hat{\theta}_{choices}$

Generalized Successive Elimination with $\hat{\theta}_{choices, times}$

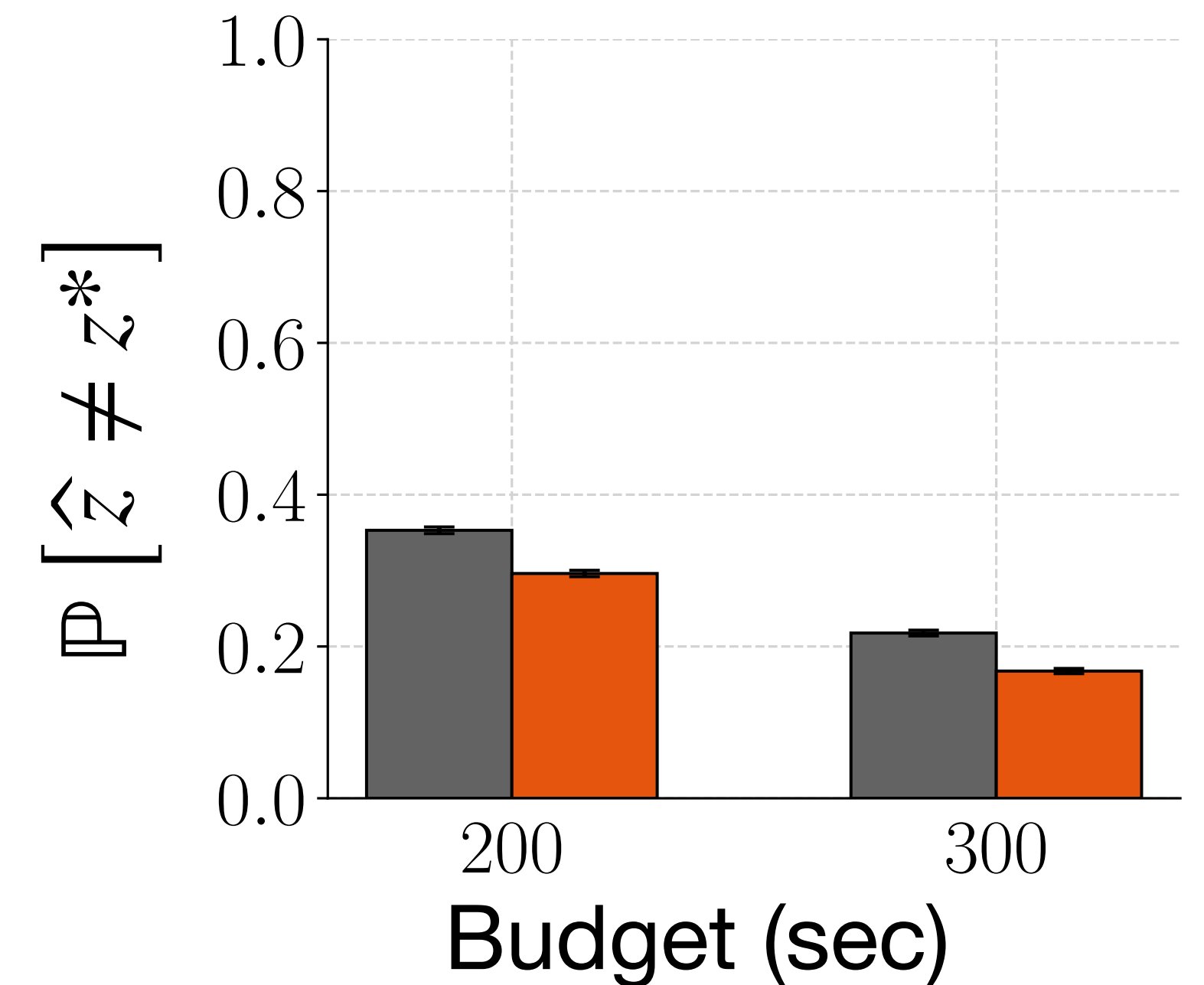
The lower, the better.



(Smith and Krajbich, 2018)



(Clithero, 2018)



(Krajbich, et al., 2010)

Key Contributions: The *First* to Use Response Times for Preference Learning

- A utility estimator using both choices and response times.
- An insight: response times from **queries with strong preferences** provide extra information that complements choices.

Full paper:



Poster: 4:30-7:30

@ East Exhibit Hall A-C #4901

(Shen is on faculty job market)

