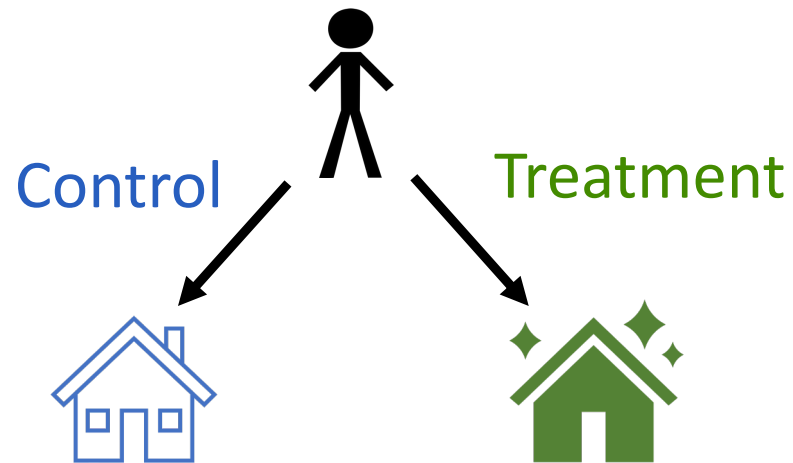# Marketplace Experimentation: Interference, Inference, and Decisions

Hannah Li (MIT)

Simons Workshop on Quantifying Uncertainty

Joint work with Ramesh Johari and Gabriel Weintraub (Stanford)
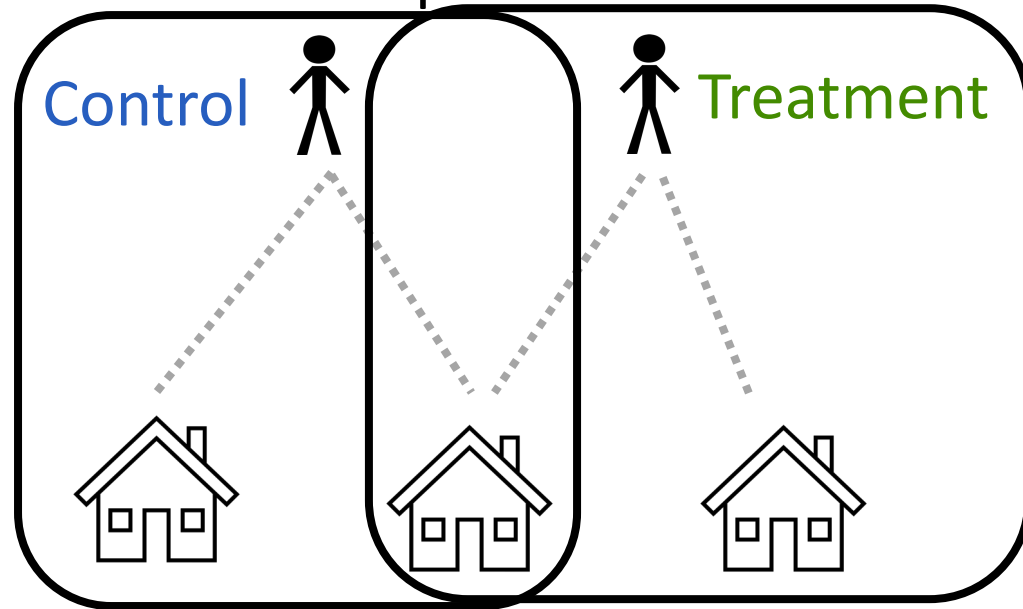
# Decision-making in online marketplaces



Control          Treatment

"If we show higher quality photos, do the number of bookings increase?"

- Experimentation ("A/B tests")
- Goal: estimate *Global Treatment Effect*

    *GTE* = Bookings in global treatment
    – bookings in global control

- Give intervention to some (treatment) and not others (control)
- Large platforms run > 10,000 per year

But estimates of GTE in marketplaces often **biased** due to interference!

# *Competition ⟹ Interference ⟹ Bias*

Customer-side experiment



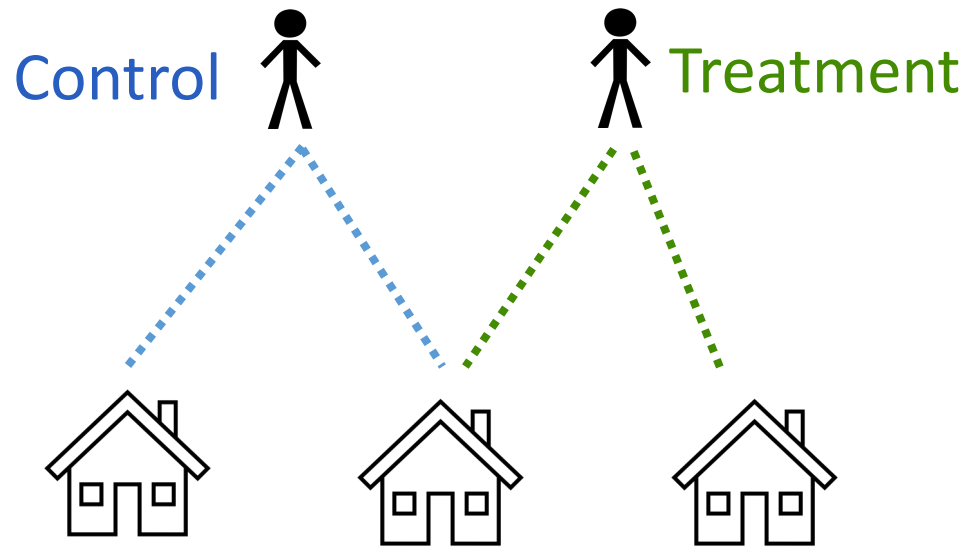Global Treatment Effect (GTE)  =  Global Treatment  −  Global Control

# Competition ⟹ Interference ⟹ Bias

Customer-side experiment



Control    Treatment
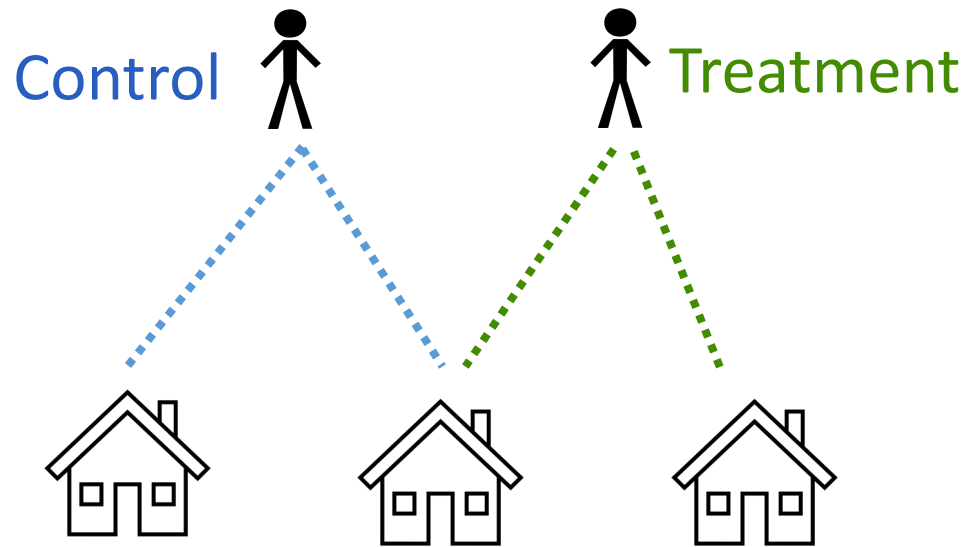
- Suppose feature makes treatment customer more likely to book than control
- Treatment customer books listing
- Reduces supply for control customer
- This instance: overestimate GTE

Global Treatment Effect (GTE) = Global Treatment − Global Control

# *Competition ⟹ Interference ⟹ Bias*

Customer-side experiment



Control    Treatment

- Suppose feature makes treatment customer more likely to book than control
- Treatment customer books listing
- Reduces supply for control customer
- This instance: overestimate GTE

More generally:

- Change a customer's booking prob. ⇒ change supply for other customers
- Change a listing's display ⇒ make other listing relatively more/less attractive

# Prior work: Interference and $\widehat{GTE}$-bias

- $\widehat{GTE}$-bias is **30% – 230%** size of GTE. [Blake and Coey '14, Fradkin '19, Holtz et al. '20, Liu et al. '21]

- Methods to reduce $\widehat{GTE}$-bias: Cluster randomization, switchback testing, and TSR. [Holtz '18, Candogan et al. '21, Sneider et al. '19, Glynn et al '20, Bojinov et al. '21, Wager and Xu '19, Ha-Thuc et al. '20, Novak et al. '20, Han et al. '21, Liu et al. '21, Bajari et al. '21, Li et al. '21, Johari et al. '22, Bright et al. '22]

- Size of bias depends on supply and demand imbalance. [Li et al. '21] [Johari et al. '22]

**This talk:** How do biases affect resulting **decisions**?
**Takeaway:** Interference creates multiple biases, fixing one bias alone can actually worsen decisions.
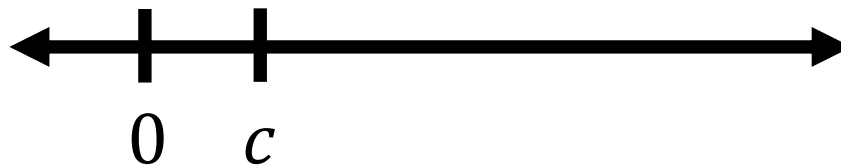
# Decision-making pipeline
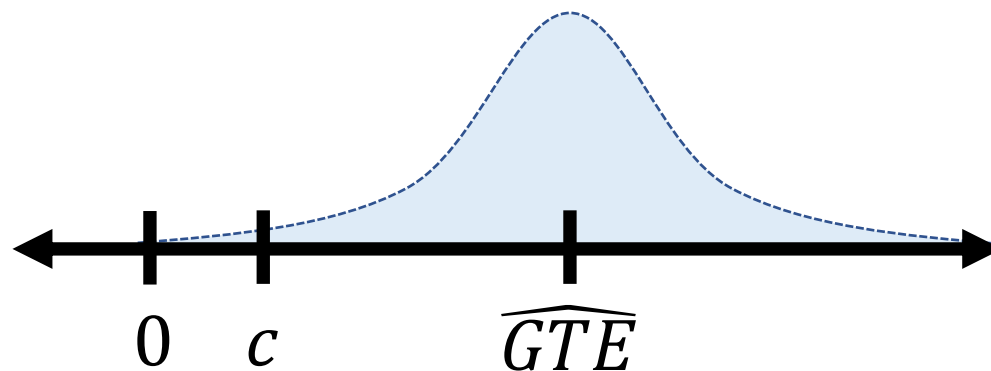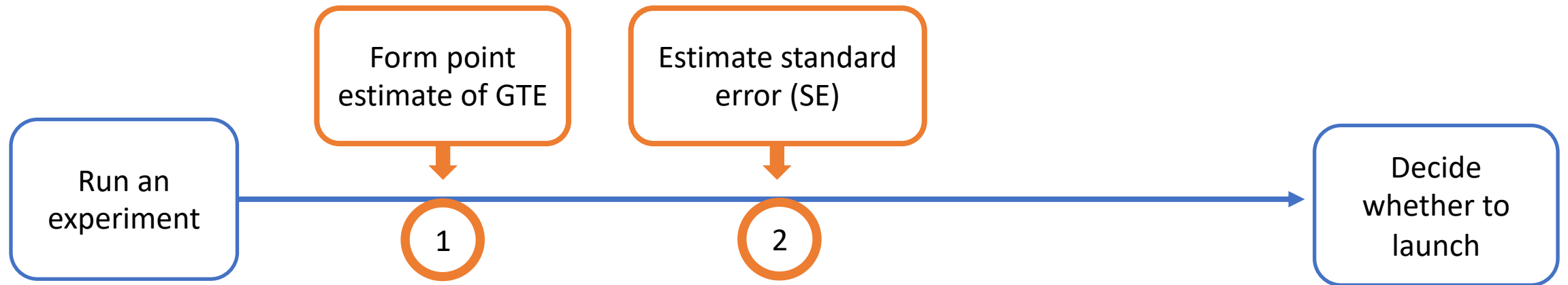
Run an experiment → Decide whether to launch

# Decision-making pipeline

Given significance level $\alpha$ and launch threshold $c$:

# Decision-making pipeline

Given significance level $\alpha$ and launch threshold $c$:

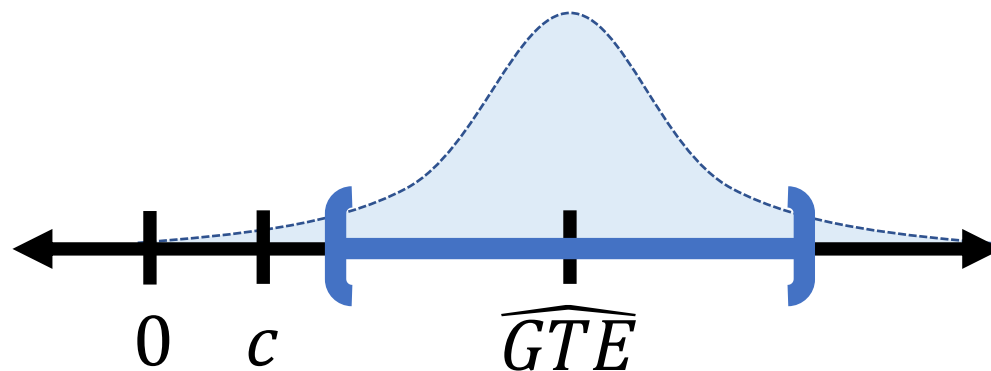# Decision-making pipeline

Given significance level $\alpha$ and launch threshold $c$:

# Decision-making pipeline

Given significance level $\alpha$ and launch threshold $c$:

# Decision-making pipeline

Given significance level $\alpha$ and launch threshold $c$:

# Decision-making pipeline

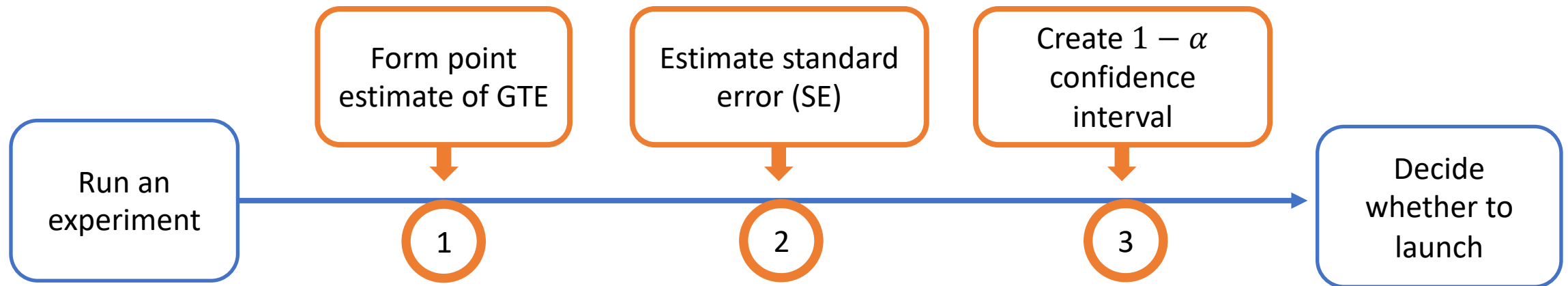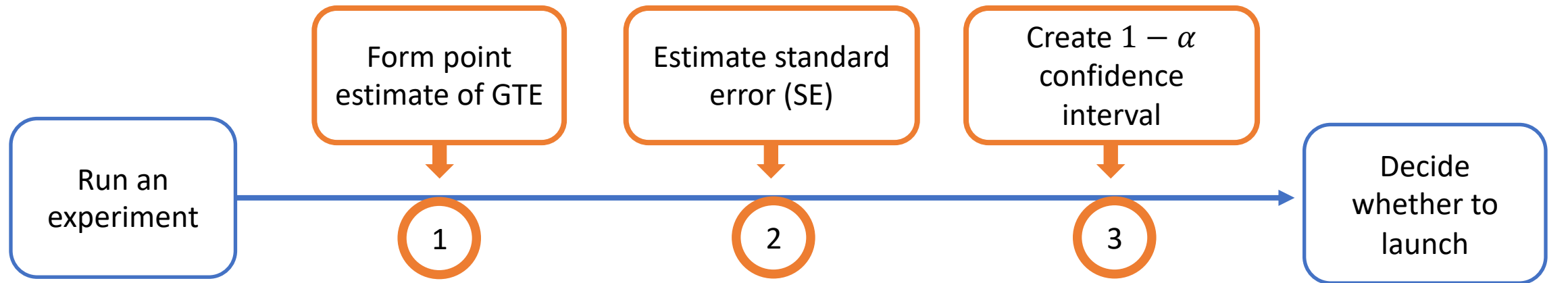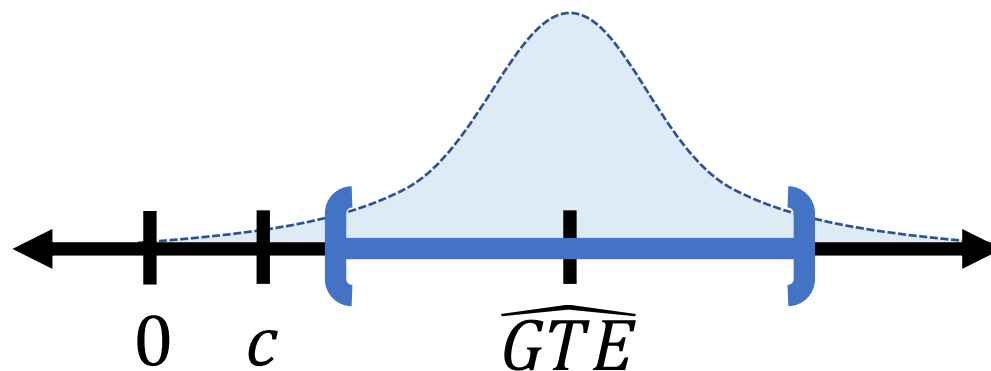Given significance level $\alpha$ and launch threshold $c$:
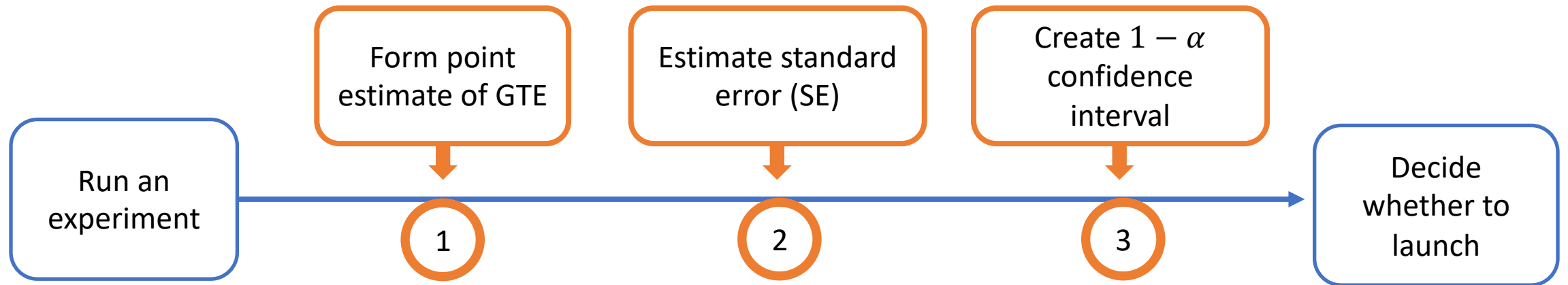


Interference can create multiple errors.

Decision heuristic:
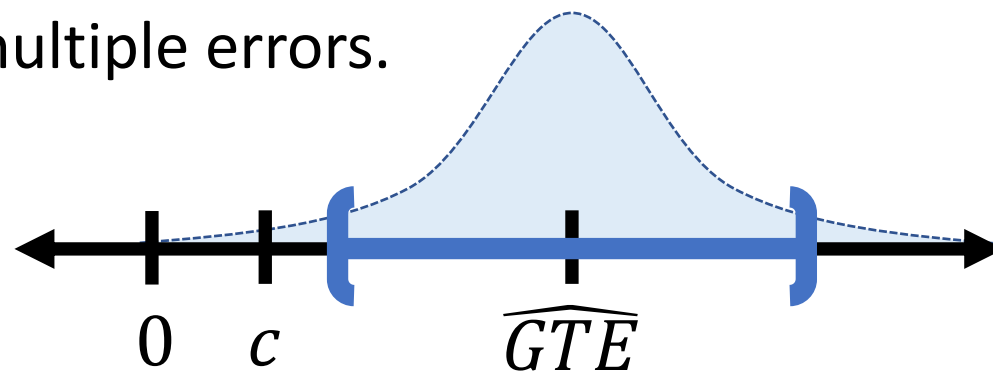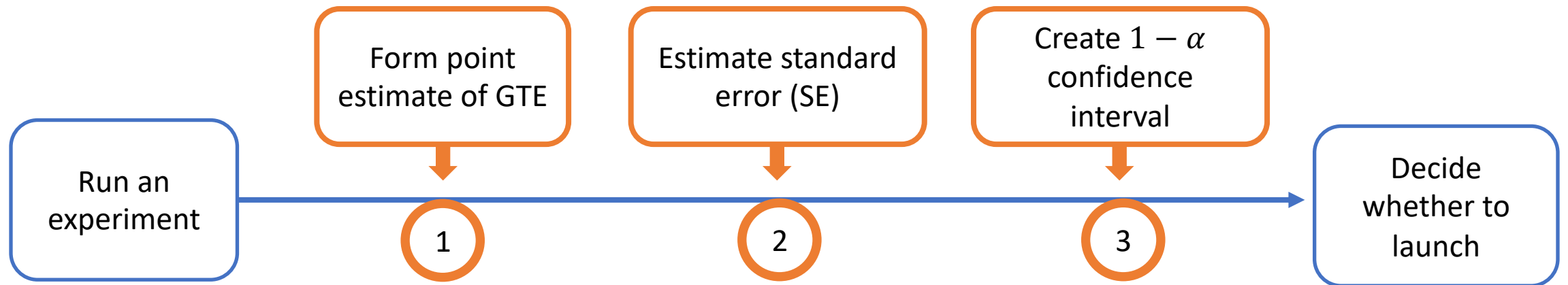Launch if lower bound of confidence interval $> c$

# Decision-making pipeline

Given significance level $\alpha$ and launch threshold $c$:



Interference can create multiple errors.

Prior work: Focuses on **(1)**

This work: Studies **(2)**, **(3)**, and impact on decisions

Decision heuristic:
Launch if lower bound of
confidence interval $> c$

# This Work

**Use a dynamic market model to study:**
1. What biases arise in $\widehat{SE}$ estimates?
2. When/how do biases in $\widehat{GTE}$ and $\widehat{SE}$ ests. affect decision-making?
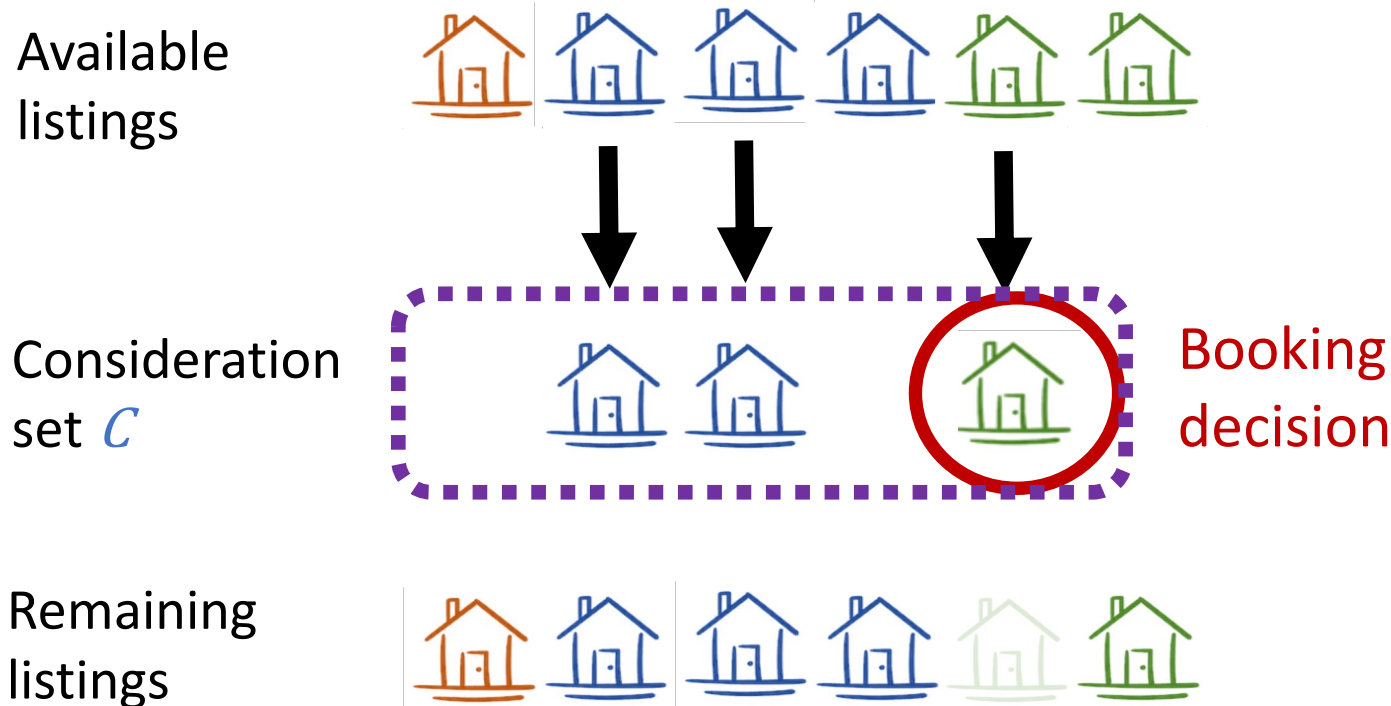
**Takeaways:**
- In a large class of interventions ("positive interventions"), $\widehat{GTE}$ and $\widehat{SE}$ - bias lead platform to launch too often.
- Two types of biases interact; fixing only one can lead to worse decisions.
- Provide a method to reduce $\widehat{SE}$-bias *and* improve decisions

# CTMC model of two-sided markets

Customers have type $\gamma \in \Gamma$. Type $\gamma$ customers arrive at rate $\Lambda_\gamma$ (Poisson).

Available listings



Consideration set $C$

Booking decision

Remaining listings

Booked listing of type $\theta$ becomes unavailable for an exponential time with parameter $\tau(\theta)$.
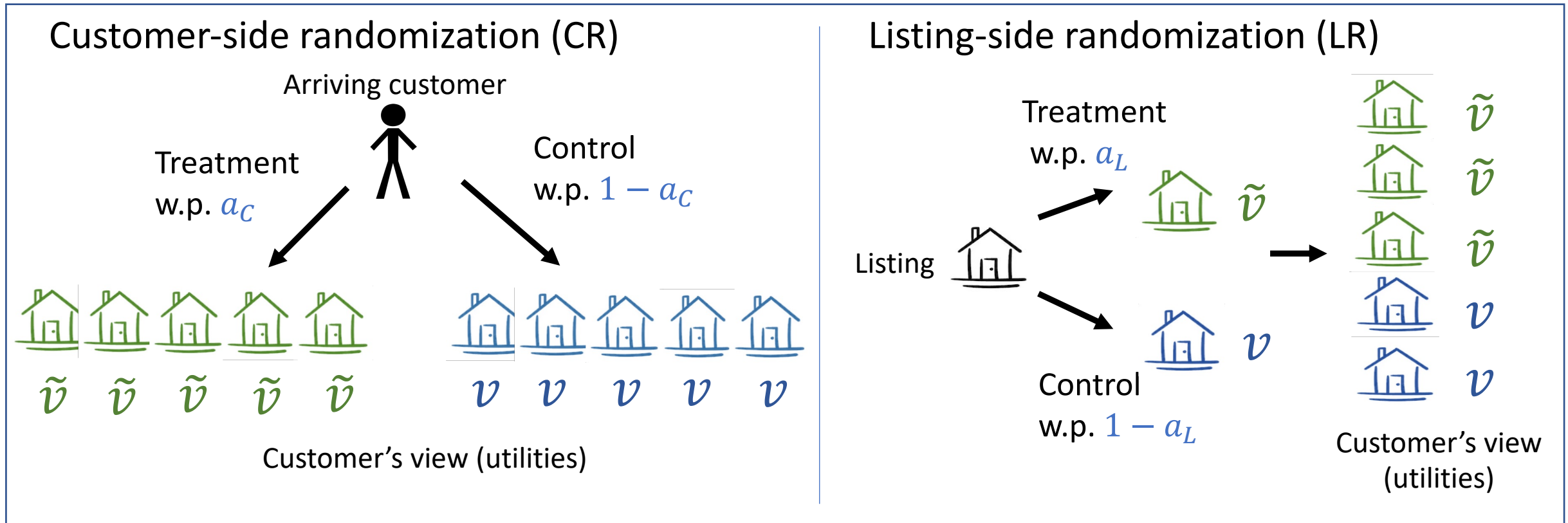
**1. Consideration set.** Includes each listing $l$ in consideration set w.p. $\alpha_\gamma(\theta_l)$ (independent across listings).

**2. Choice.** Chooses from consideration set according to multinomial logit model.

$$P_\gamma(\text{choose } l) = \frac{v_\gamma(\theta_l)}{E_\gamma + \sum_{l' \in C} v_\gamma(\theta_{l'})}$$

# Running an experiment

We focus on two common types of marketplace experiments.



## Customer-side randomization (CR)

Arriving customer

Treatment w.p. $a_C$

Control w.p. $1 - a_C$

$\tilde{v}$ $\tilde{v}$ $\tilde{v}$ $\tilde{v}$ $\tilde{v}$    $v$ $v$ $v$ $v$ $v$

Customer's view (utilities)

## Listing-side randomization (LR)

Treatment w.p. $a_L$

Listing

Control w.p. $1 - a_L$

$\tilde{v}$ $v$

Customer's view (utilities)

$$\widehat{GTE}^{CR} = \frac{\#\ Treatment\ Bookings}{a_C\ T} - \frac{\#\ Control\ Bookings}{(1 - a_C)T}$$

# Quantities of interest

Study Markov chain behavior in counterfactual worlds

$\qquad$ (global treatment, global control, experiment)

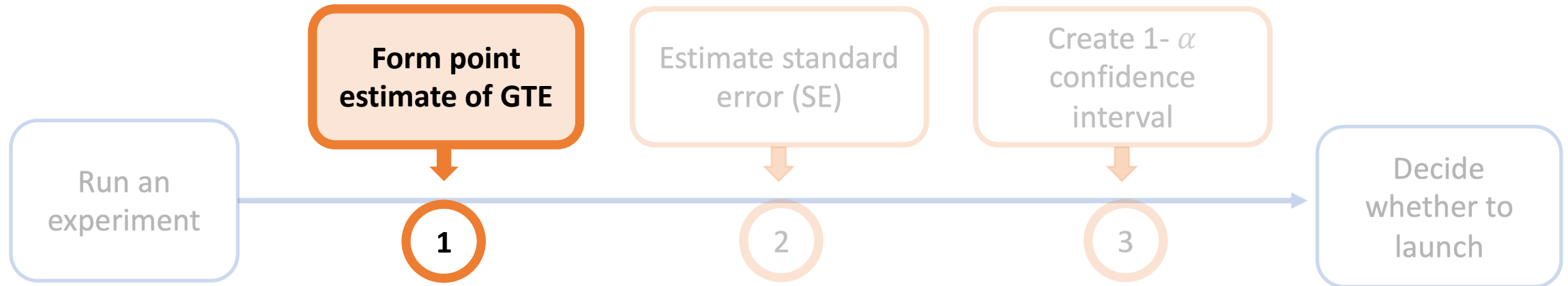- Estimand: Global Treatment Effect $GTE$ – evaluated in steady state

$$GTE = \underbrace{Q^{GT}}_{\substack{\text{Global Treatment} \\ \text{rate of booking}}} - \underbrace{Q^{GC}}_{\substack{\text{Global Control} \\ \text{rate of booking}}}$$
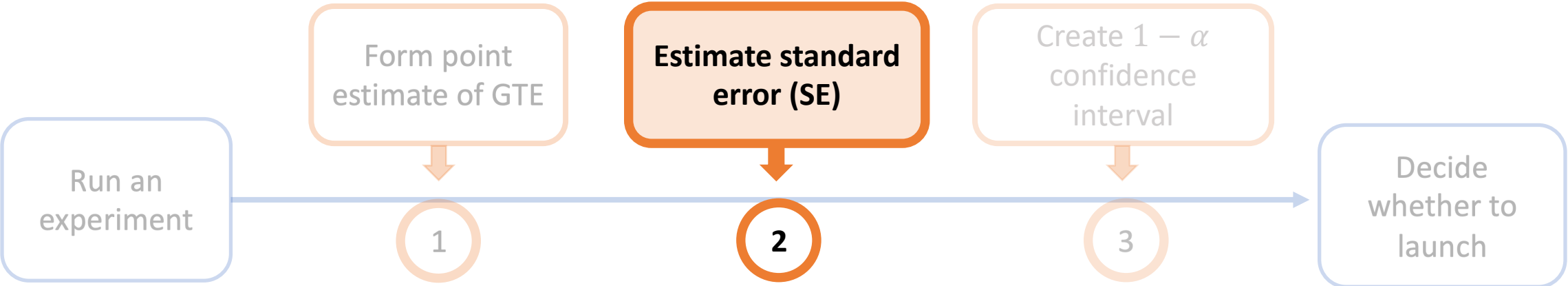
- Estimator (calculated from experiment booking rates):

$$\widehat{GTE} = \frac{\#\,Treatment\,Bookings}{a_C\,T} - \frac{\#\,Control\,Bookings}{(1-a_C)T}$$

- Standard Error $\qquad SE = \left(Var\left(\widehat{GTE}\right)\right)^{1/2}$

This talk: Focus on a class of interventions that increases utilities, denoted "positive" interventions



**Theorem (informal)** [JLLW '21]:
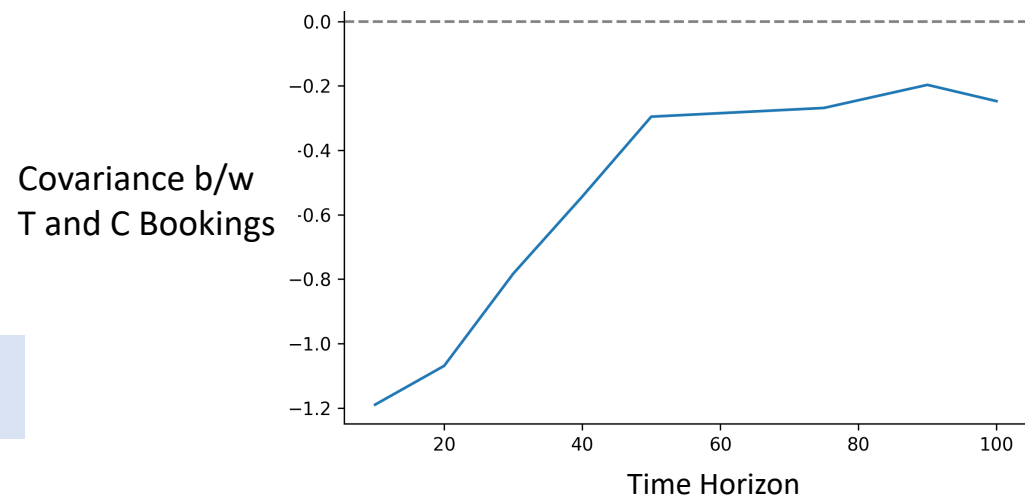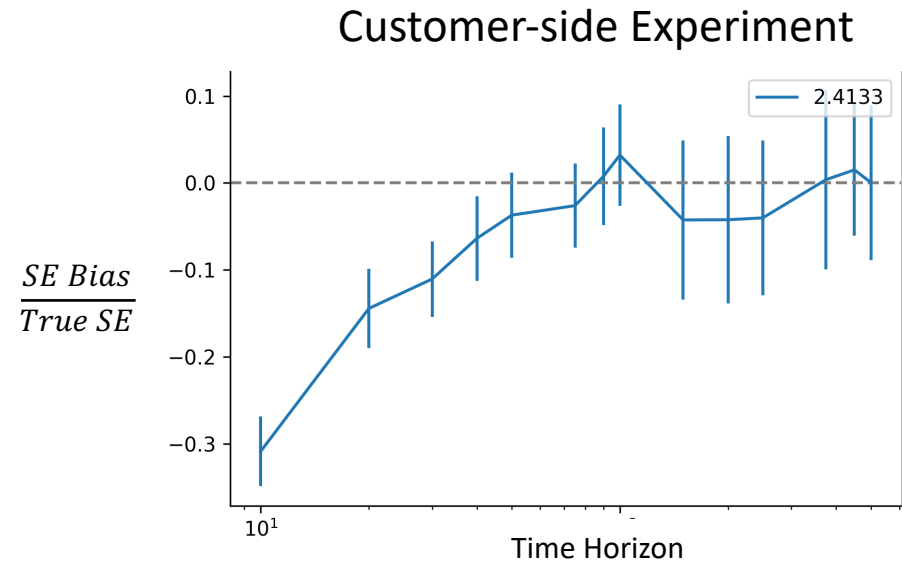For a positive intervention, $CR$ and $LR$ *overestimate* the magnitude of the $GTE$.

# (2) Inference and $SE$ estimation

- Estimate $SE = \left(Var\left(\widehat{GTE}\right)\right)^{1/2}$

- "Naive" $\widehat{SE}$ estimate: **Assume individuals are independent**

- Leads to biased estimates of $SE$

$$Var(T - C)$$
$$= Var(T) + Var(C) - Cov(T, C)$$

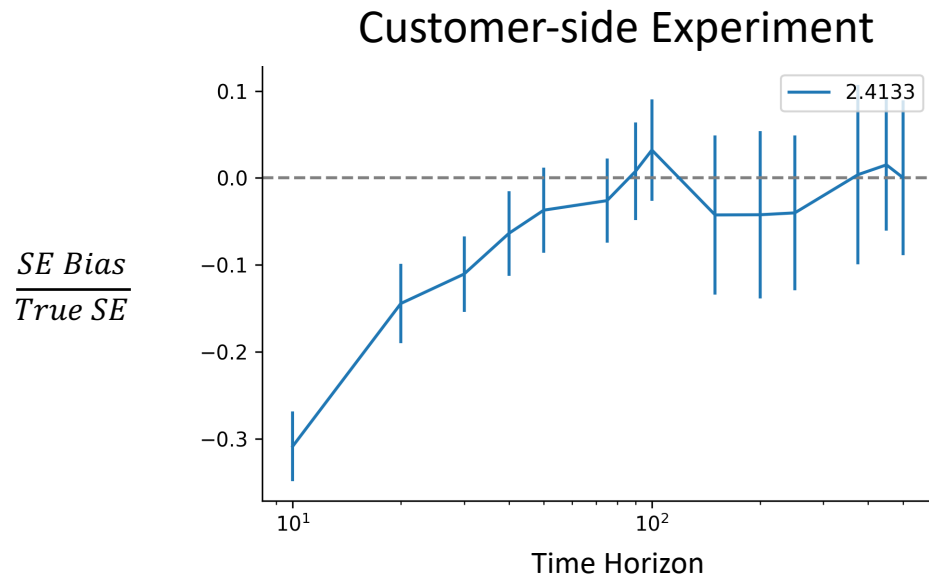- Ignores correlation between individual outcomes

**Competition $\Longrightarrow$ Interference $\Longrightarrow$ Bias**



Customer-side Experiment

$\frac{SE\ Bias}{True\ SE}$

Time Horizon



Covariance b/w T and C Bookings

Time Horizon

# Reducing $\widehat{SE}$ bias
# Method 1: Longer experiments



Customer-side Experiment

$\frac{SE\ Bias}{True\ SE}$

Time Horizon

**Theorem (informal).**

For a customer-side experiment, the bias of the "naive" $\widehat{SE}$ estimate approaches 0 as $T \to \infty$.

*Proof idea.*

System is a regenerative process.

# Reducing $\widehat{SE}$ bias
# Method 2: Block bootstrap

- Standard bootstrap: resample individuals

- Block bootstrap [Hardle et al. '03]
  - Resample "blocks" from observed time series, create "pseudo-time series"

Observed time series of bookings

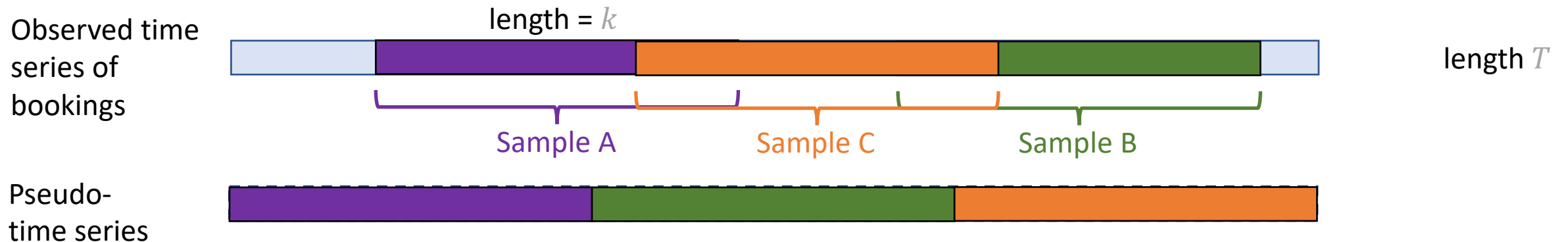length $T$

Pseudo-time series

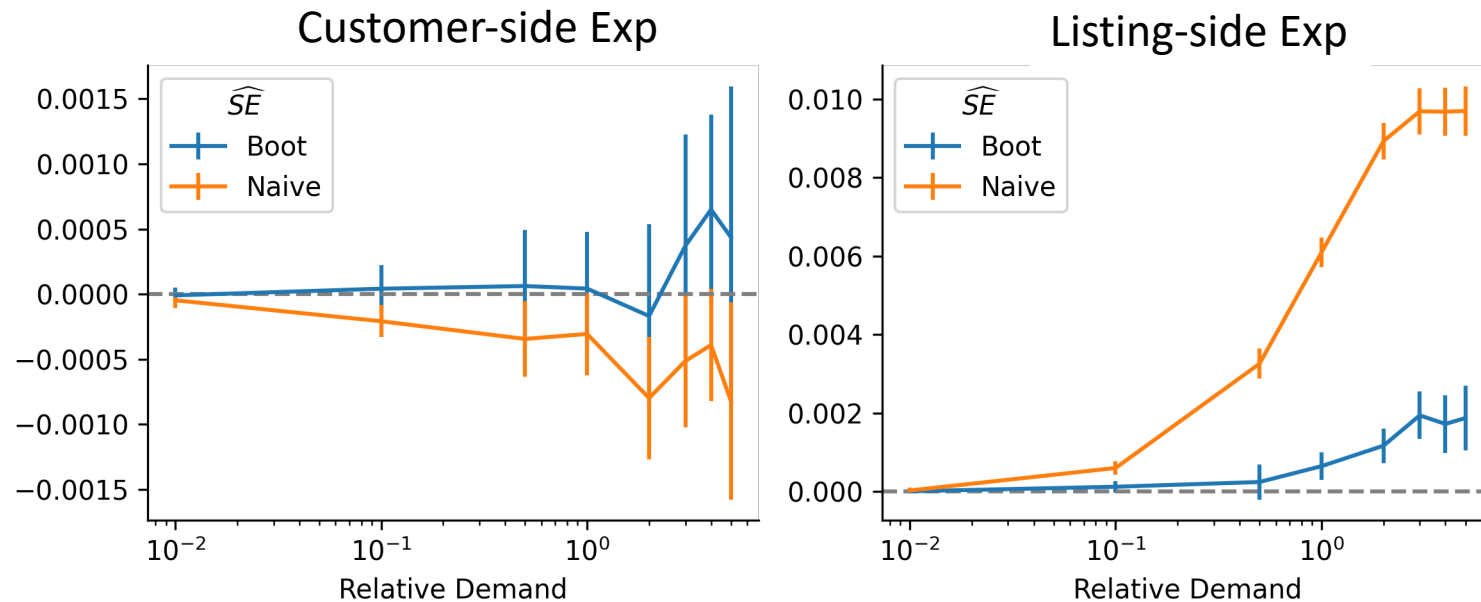# Reducing $\widehat{SE}$ bias
# Method 2: Block bootstrap

- Standard bootstrap: resample individuals

- Block bootstrap [Hardle et al. '03]

  - Resample "blocks" from observed time series, create "pseudo-time series"



- From each bootstrap run $b$ (pseudo-time series): calculate $\widehat{GTE_b}$

- Repeat $B$ times, calculate std. dev. across $\widehat{GTE_b}$ estimates $\rightarrow \widehat{SE}^{boot}$
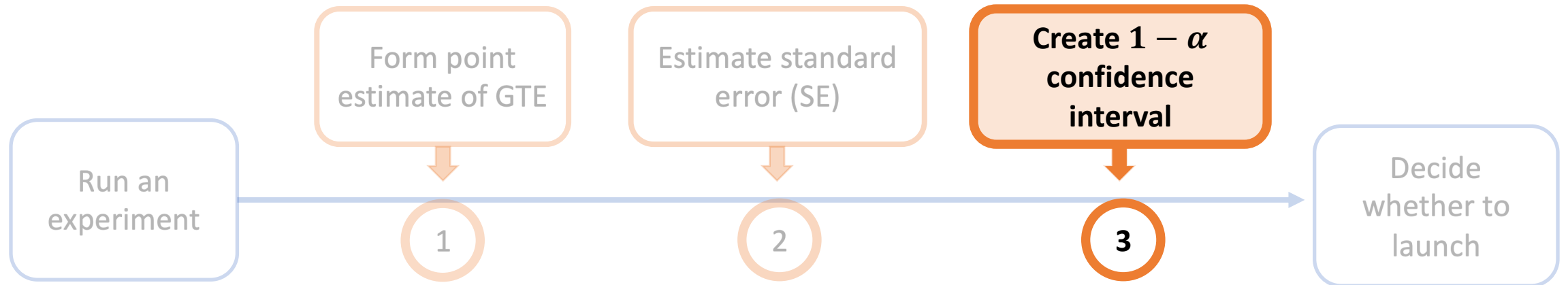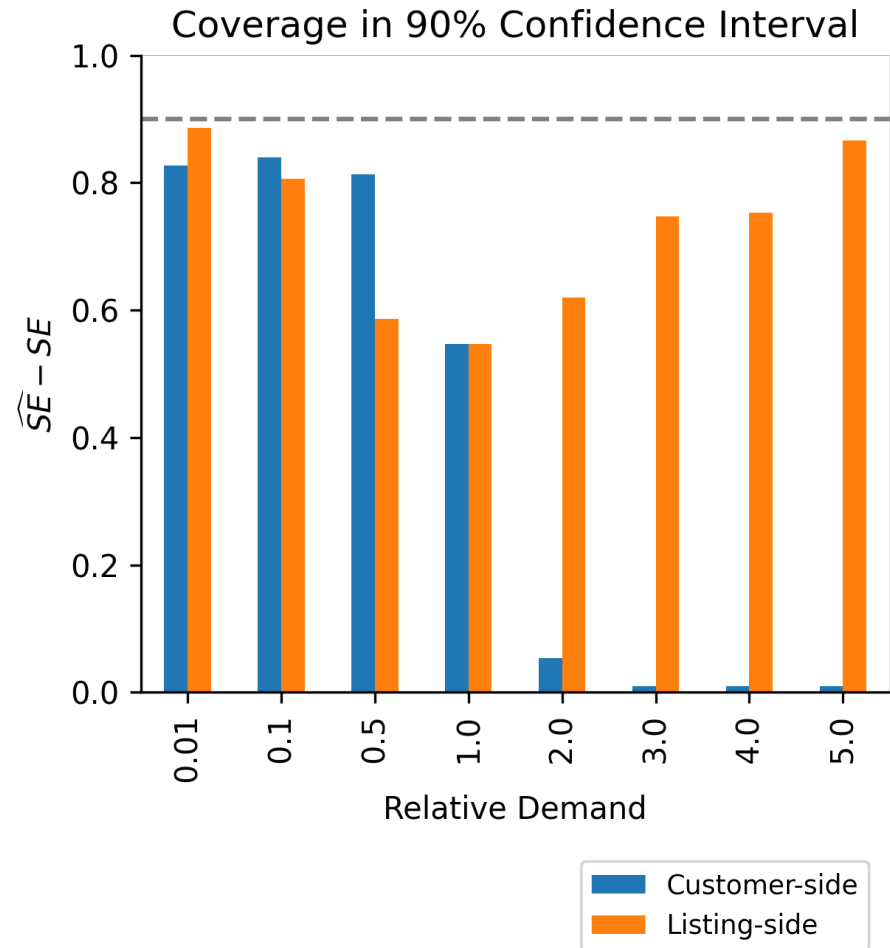
# Reducing SE bias
# Method 2: Block bootstrap



**Takeaway:** Bootstrapping can mitigate biases.

**Caveat:** Need to tune block length.

# Decision-making pipeline
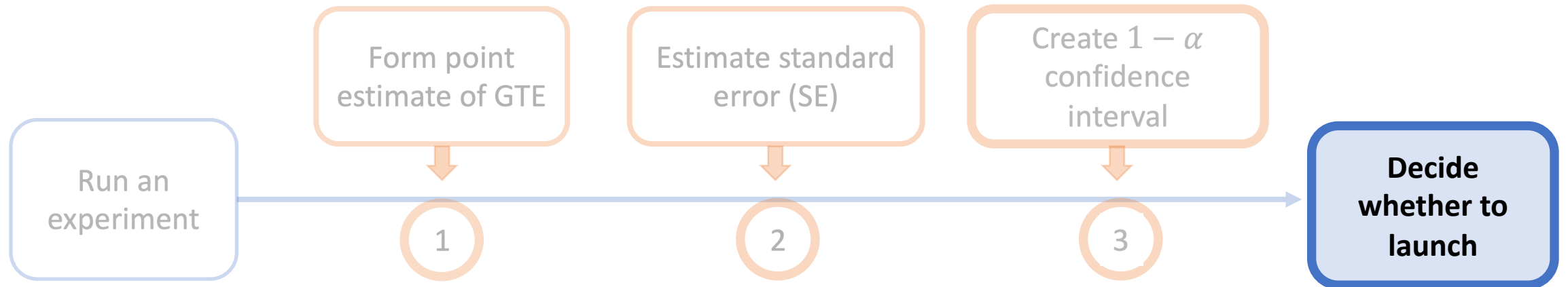
# (3) Coverage of confidence intervals



Coverage in 90% Confidence Interval

- $\widehat{GTE}$-bias shifts confidence intervals
- $\widehat{SE}$-bias changes width of intervals
- Interactions between $\widehat{GTE}$-bias and $\widehat{SE}$-bias determine coverage
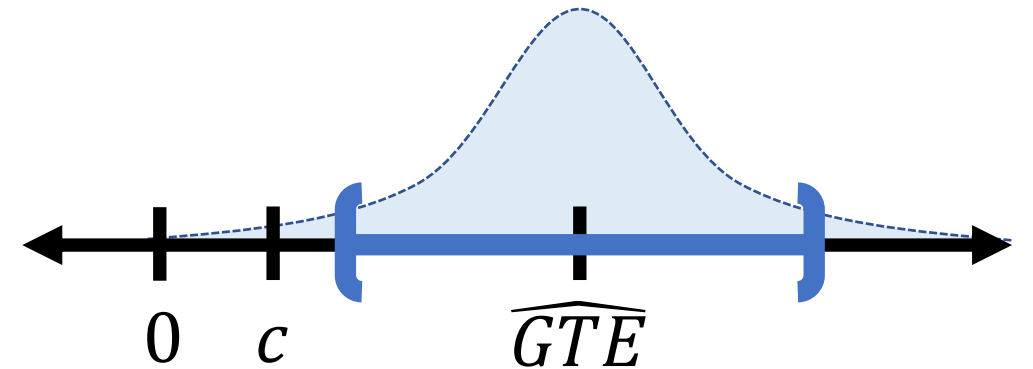
We characterize asymptotic coverage of conf. ints. as a function of $\widehat{GTE}$-bias, $\widehat{SE}$-bias, and $SE$

# Decision-making pipeline

# Implications for decision-making

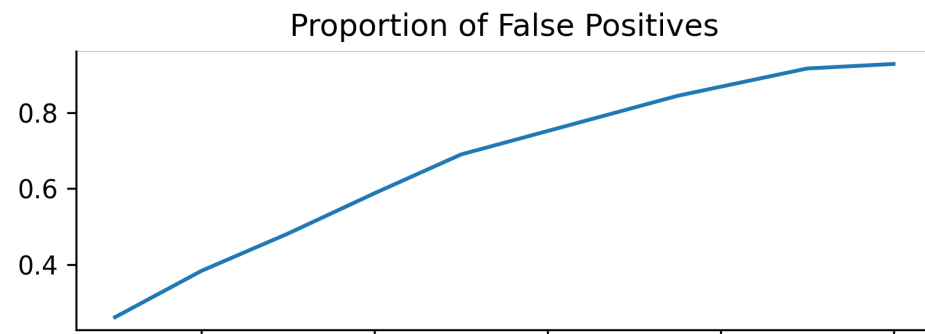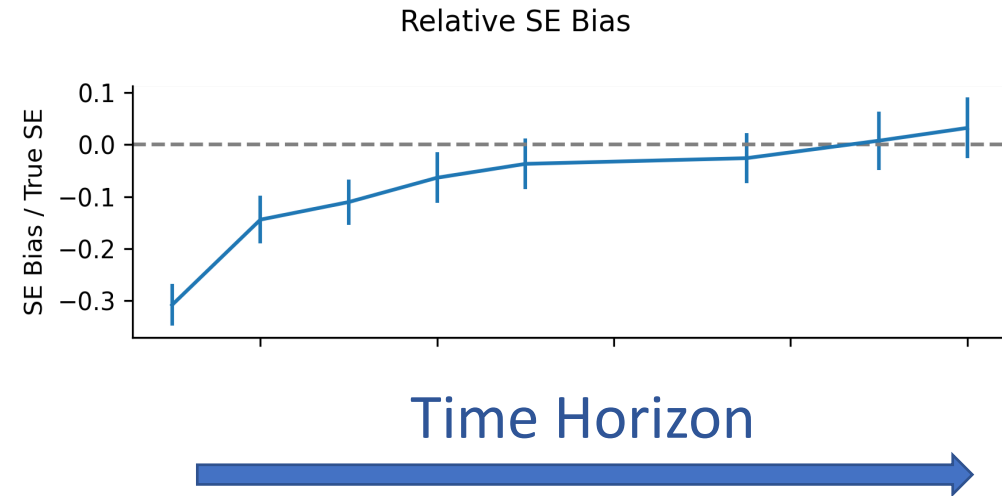| Goal: | Launch if $GTE > c$. |
|---|---|
| Decision Heuristic: | Launch if lower bound of conf int > $c$. |
| Evaluating Decision: | Decision is correct if we launch only when $GTE > c$. |



In positive interventions, we see:

1. Overestimation of $GTE$ in CR and LR experiments
2. Underestimation of $SE$ in CR experiments

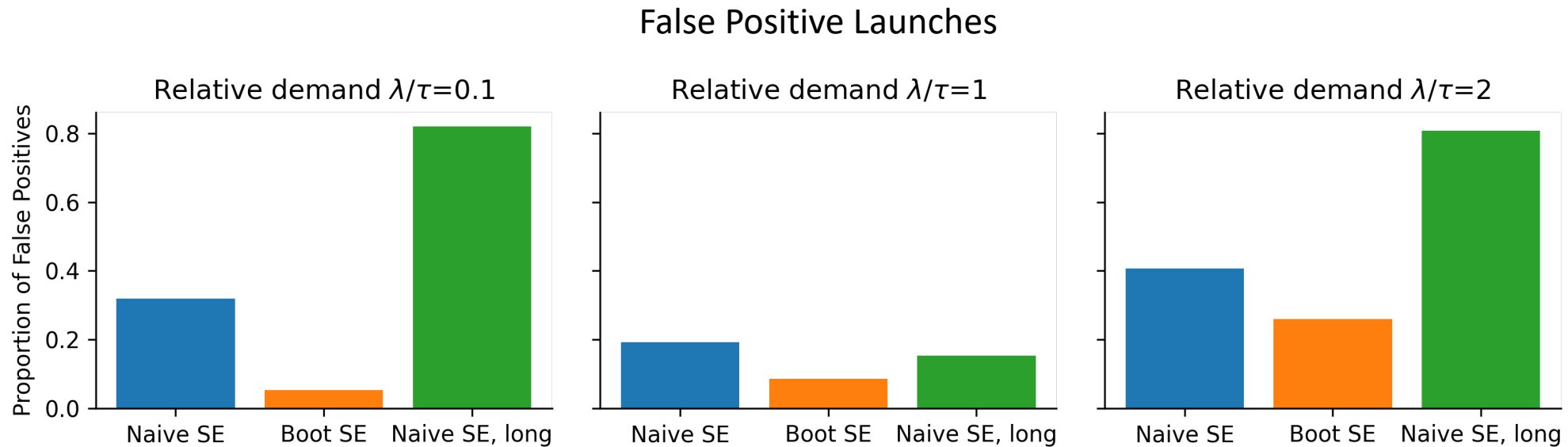Combination leads to more **false positives** (launch feature when $GTE < c$).

# Fixing bias ≠ improving decisions

- Scenario: $GTE < c$

- Any decision to launch is a "false positive"

- Implement **Method 1** for SE-bias reduction: Run longer experiment

- As time horizon increases:
  - Actual $SE$ of CR estimator decreases together with $\widehat{SE}$ bias
  - More confident about our biased $\widehat{GTE}$

# Alternative: Reduce $\widehat{SE}$ bias with bootstrap

- Scenario: $GTE < c$

- (With appropriate block length) bootstrap method reduces SE bias and reduces false positive launches



False Positive Launches

# Takeaways

- $\widehat{GTE}$ and $\widehat{SE}$-biases interact and cause incorrect decisions
- Propose two methods to reduce $\widehat{SE}$-bias
    1. Increasing time horizon – can worsen decisions
    2. Block bootstrapping – can improve decisions

Open questions
- Combining $\widehat{SE}$-bias reduction with $\widehat{GTE}$-bias reduction
- Increased attention on decisions made from experiments
- Marketplace interactions complicate many statistical methods. How do complications interact with the ways platforms utilize experiments?
    - e.g., simultaneous experiments, ramp-up experiments