

Dynamic Remote Patient Monitoring under Alert Fatigue

Jimmy Qin, UT Dallas

Presented at Advances in Quantitative Medical Care, IMSI
Joint work with Wanting Yang and Jingwei Zhang, CUHK-SZ

Remote Patient Monitoring (RPM)

- Remote patient monitoring (RPM) uses digital devices to collect patient health data from a patient's home and transmit it to healthcare providers

Remote Patient Monitoring (RPM)

- Remote patient monitoring (RPM) uses digital devices to collect patient health data from a patient's home and transmit it to healthcare providers
- Some examples of devices
 - Blood pressure monitor
 - Glucose meter



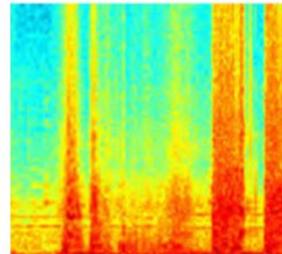
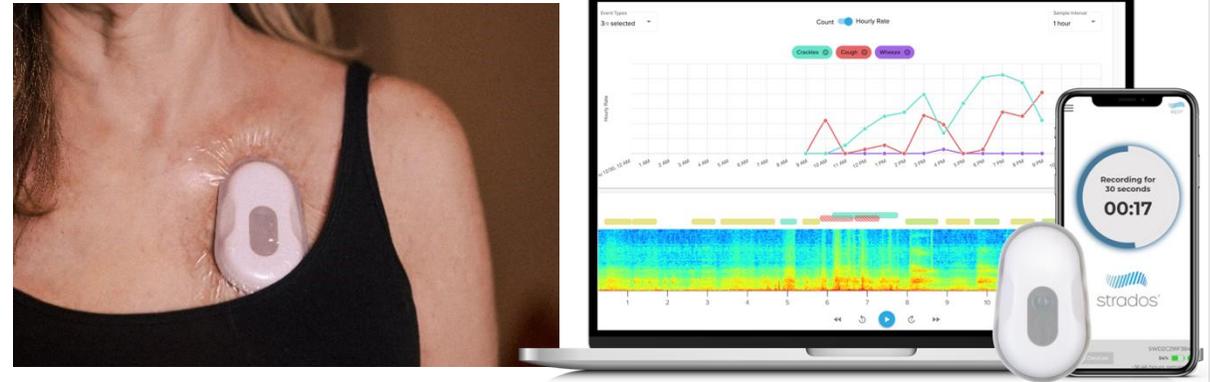
Remote Patient Monitoring (RPM)

- Remote patient monitoring (RPM) uses digital devices to collect patient health data from a patient's home and transmit it to healthcare providers
- Some examples of devices
 - Blood pressure monitor
 - Glucose meter
 - Internal cardiac devices

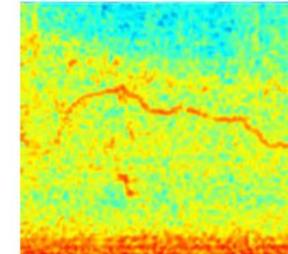


Remote Patient Monitoring (RPM)

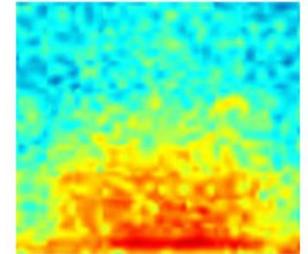
- Remote patient monitoring (RPM) uses digital devices to collect patient health data from a patient's home and transmit it to healthcare providers
- Some examples of devices
 - Blood pressure monitor
 - Glucose meter
 - Internal cardiac devices
 - Remote respiratory monitoring
 - ...



Cough



Wheeze



Crackles

Remote Patient Monitoring: benefits

- Early detection of issues and real-time monitoring
- Enhanced patient engagement
- Increased convenience and access
- Reduced costs
- Reduced clinical burden

Remote glucose monitor



Lessons we learnt from current practice

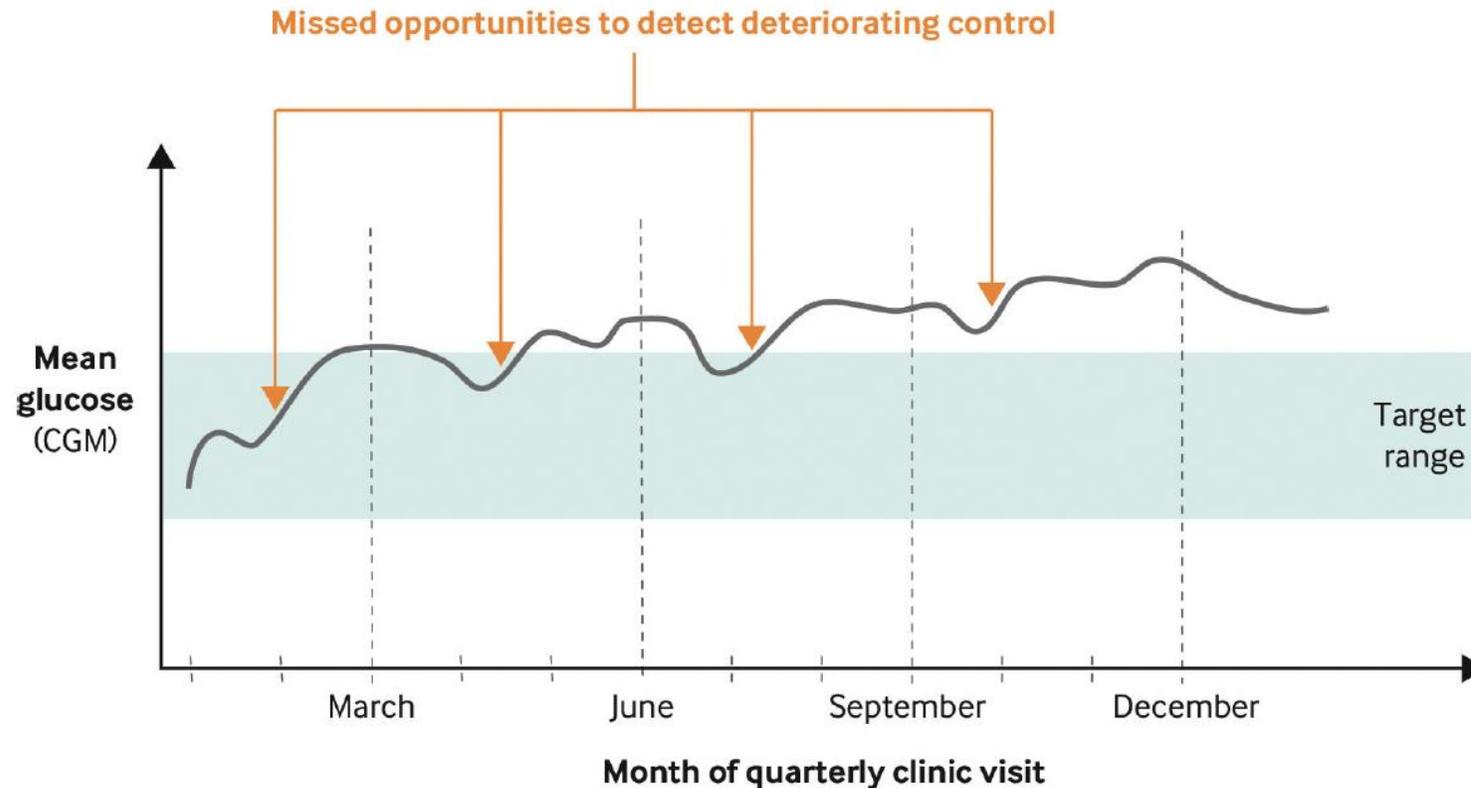


FIGURE 1. Glucose Management with Fixed-Cadence Visits (Hypothetical)

Glucose management improves after fixed-cadence visits, but not enough to make up for episodes of deterioration occurring between visits. The trend illustrated is consistent with an observed decline in postdiagnosis glucose management.

CGM = continuous glucose monitor.

Source: Lucile Packard Children's Hospital Stanford

NEJM Catalyst (catalyst.nejm.org) © Massachusetts Medical Society

Lessons we learnt from current practice

Why the Most “Accurate” Glucose Monitors Are Failing Some Users > For diabetics who rely on them, lived experience trumps metrics

BY GWENDOLYN RAK | 09 DEC 2025 | 5 MIN READ | 

Faulty glucose sensors used by diabetes patients linked to 7 deaths, hundreds of medical issues

The issue involves approximately 3 million glucose sensors in the U.S.

Continuous Glucose Monitors Can Overestimate Blood Sugar Levels, Study Finds

NEWS ARTICLE

3/4/2025

TUESDAY, March 4, 2025 (HealthDay News) – Continuous blood glucose monitors have been promoted as potentially life-changing for people with [diabetes](#) – allowing real-time updates on blood sugar levels without the need for repeated finger pricks.

But a new small-scale study suggests these devices might not be as accurate as many believe, and could lead some to mismanage their diets.

Continuous glucose monitors (CGMs) appear to overestimate blood sugar levels in healthy adults, according to findings published recently in the [American Journal of Clinical Nutrition](#).

The monitors consistently reported elevated blood sugar levels two to four times more often than finger prick tests, which are the gold standard for blood glucose testing, researchers report.



Related Topics

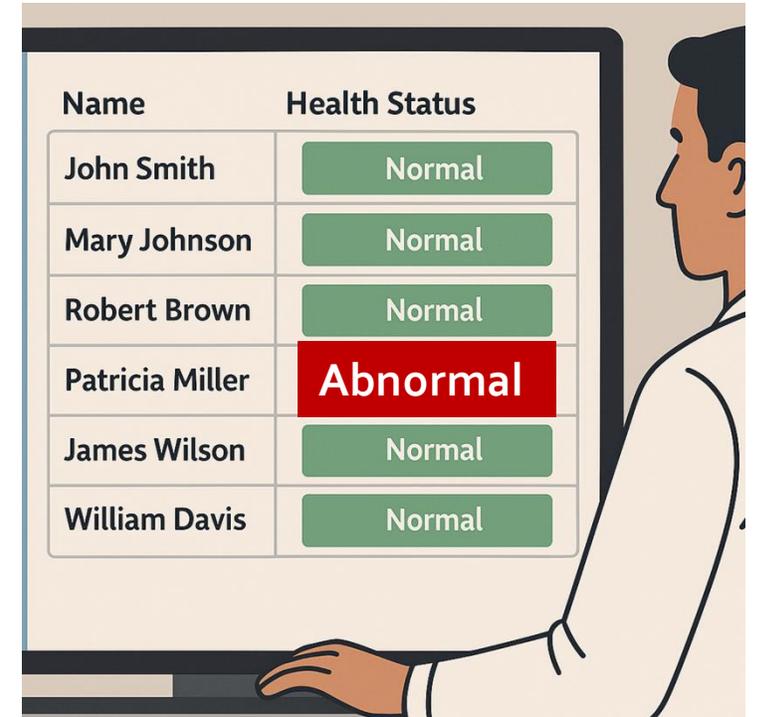
[Diabetes: Misc.](#)

FDA Warns Public About Defective Glucose Monitors

Certain sensors have been linked to seven deaths and more than 700 serious injuries

Research question: optimal strategy to alert patients

- If a healthcare provider can observe the remote monitoring signals of all patients under their care, then what is the optimal strategy to alert the patients and bring them to the hospital?



An illustration of a healthcare provider in a white coat looking at a computer monitor. The monitor displays a table with patient names and their health status. The table has two columns: 'Name' and 'Health Status'. The rows are: John Smith (Normal), Mary Johnson (Normal), Robert Brown (Normal), Patricia Miller (Abnormal), James Wilson (Normal), and William Davis (Normal). The 'Abnormal' status for Patricia Miller is highlighted in a red box.

Name	Health Status
John Smith	Normal
Mary Johnson	Normal
Robert Brown	Normal
Patricia Miller	Abnormal
James Wilson	Normal
William Davis	Normal

Model setup: finite-horizon discounted POMDP model

■ State:

1. True health states are binary (*healthy* or *sick*) and unobservable
2. Belief of being healthy $b \in [0,1]$
 - Remote monitoring can have false positive error α and false negative error β (known)
3. Disease naturally evolves according to the transition matrix P (known)

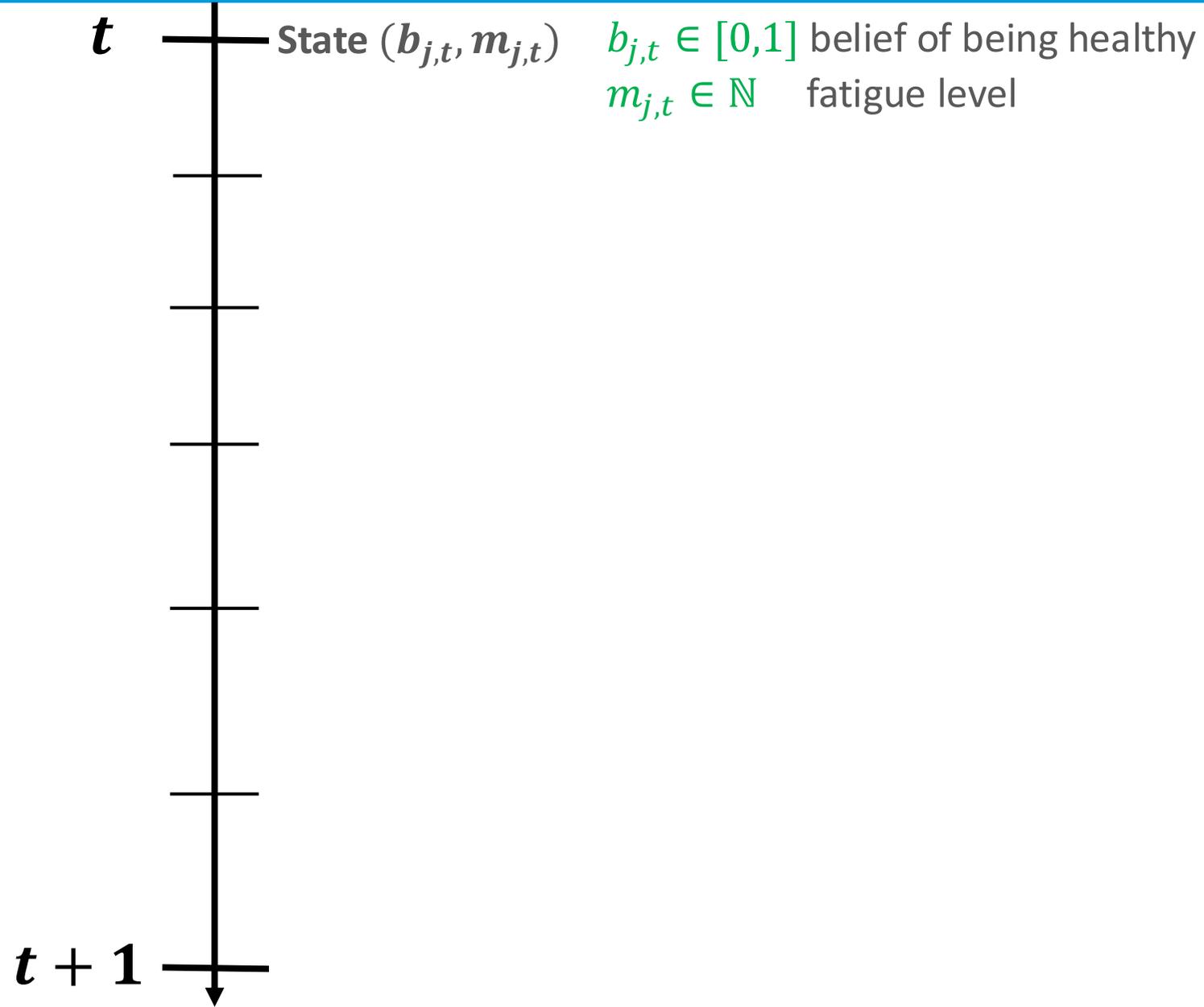
■ Action: Alert/Schedule at most one patient among n patients

■ After being alerted/scheduled, the patient may

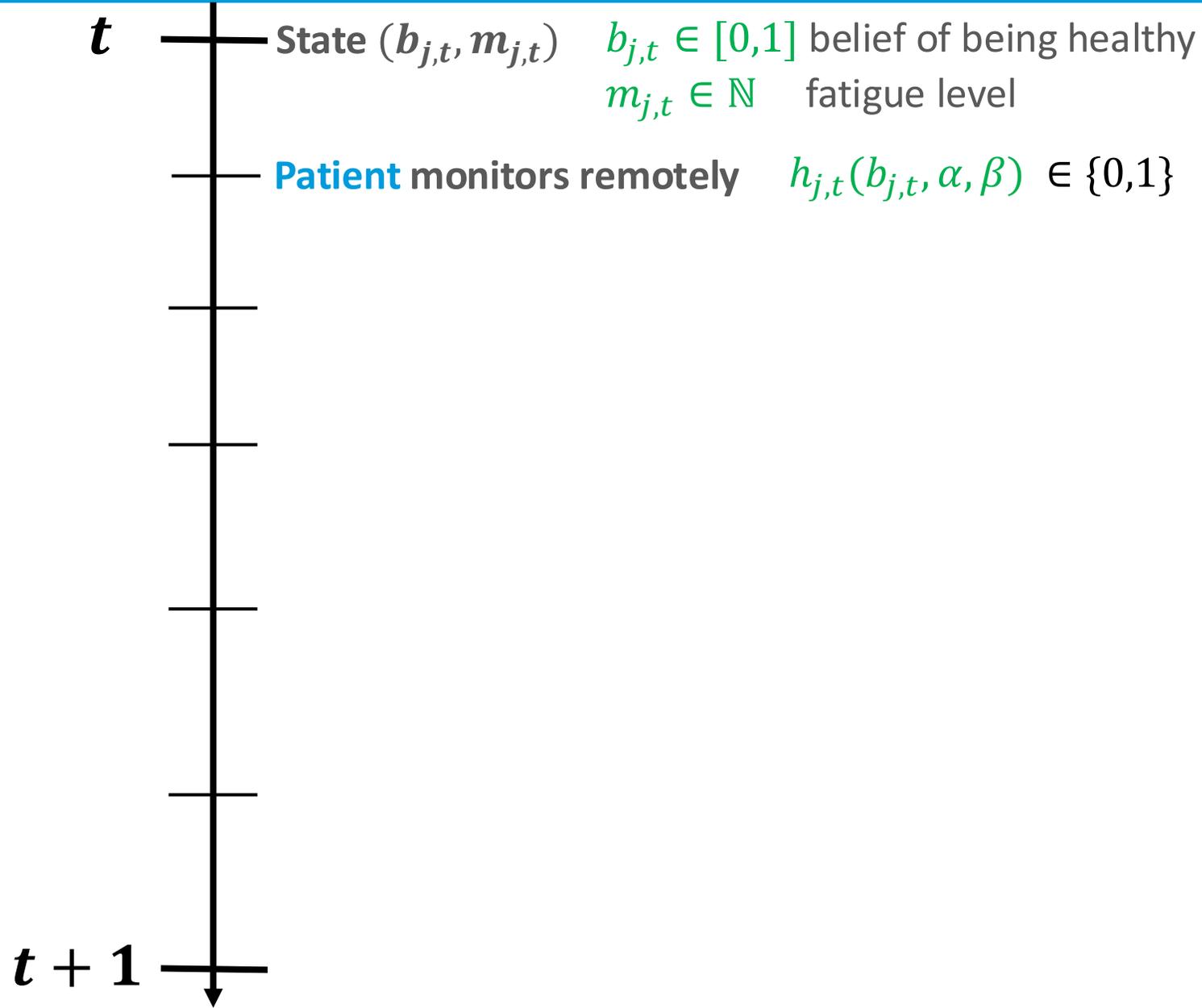
1. No-show with probability $\phi(m)$ where m is the fatigue level
2. Show-up with probability $1 - \phi(m)$
 - Diagnosis: reveal true health state
 - Treatment

■ Objective: maximizing the aggregate health benefits measured in quality adjusted life years (QALYs) for the patient cohort

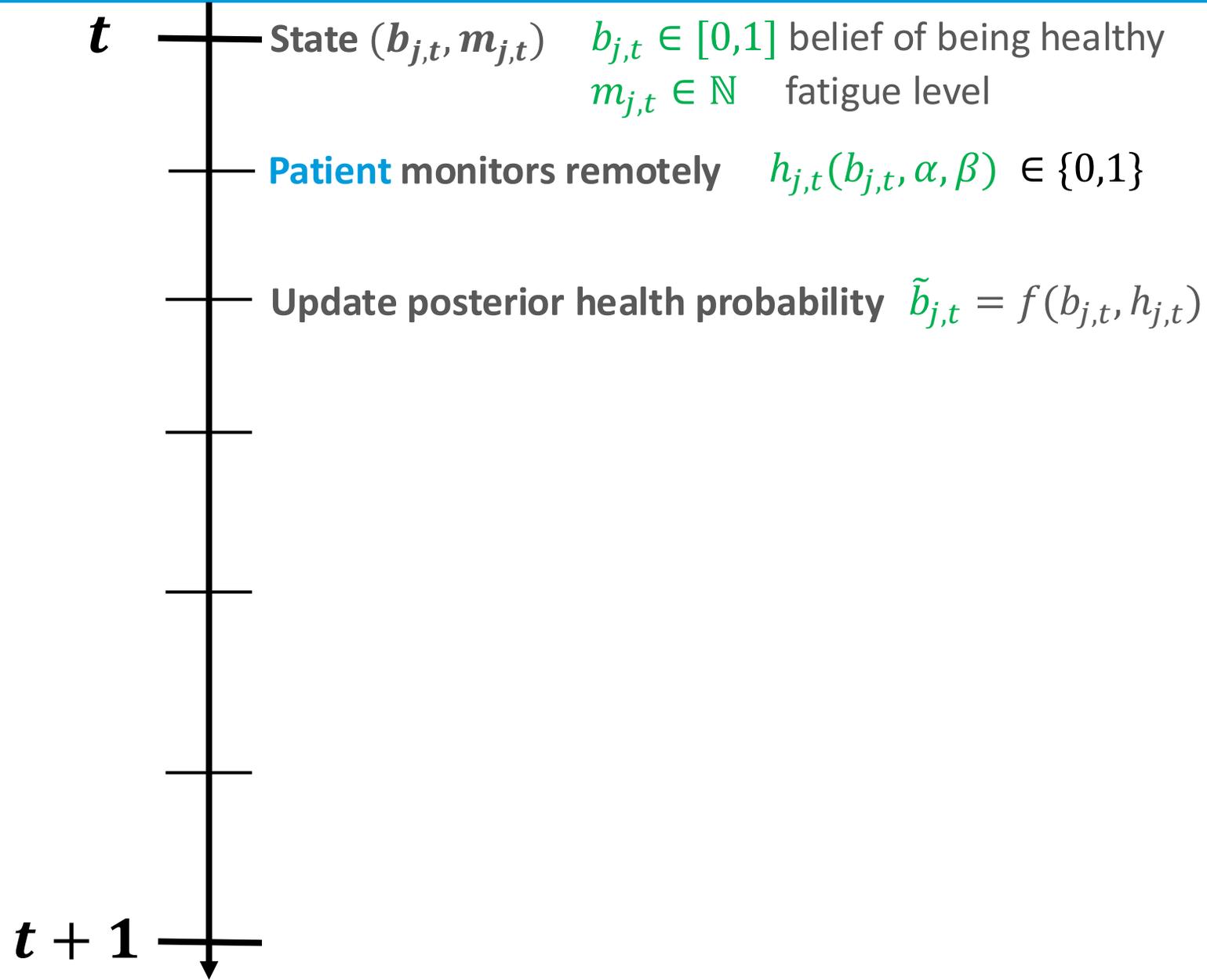
Model: sequence of events in period t



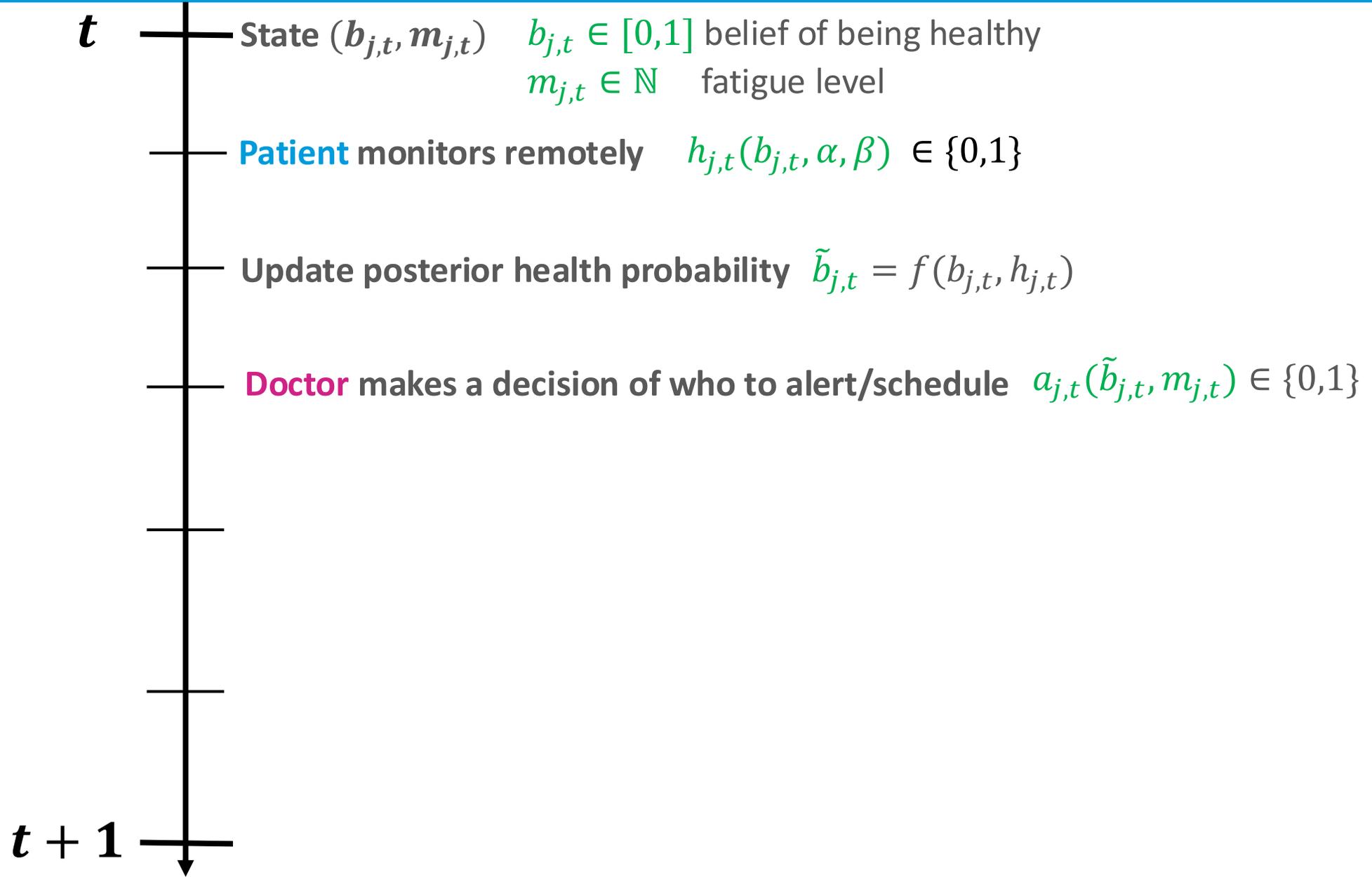
Model: sequence of events in period t



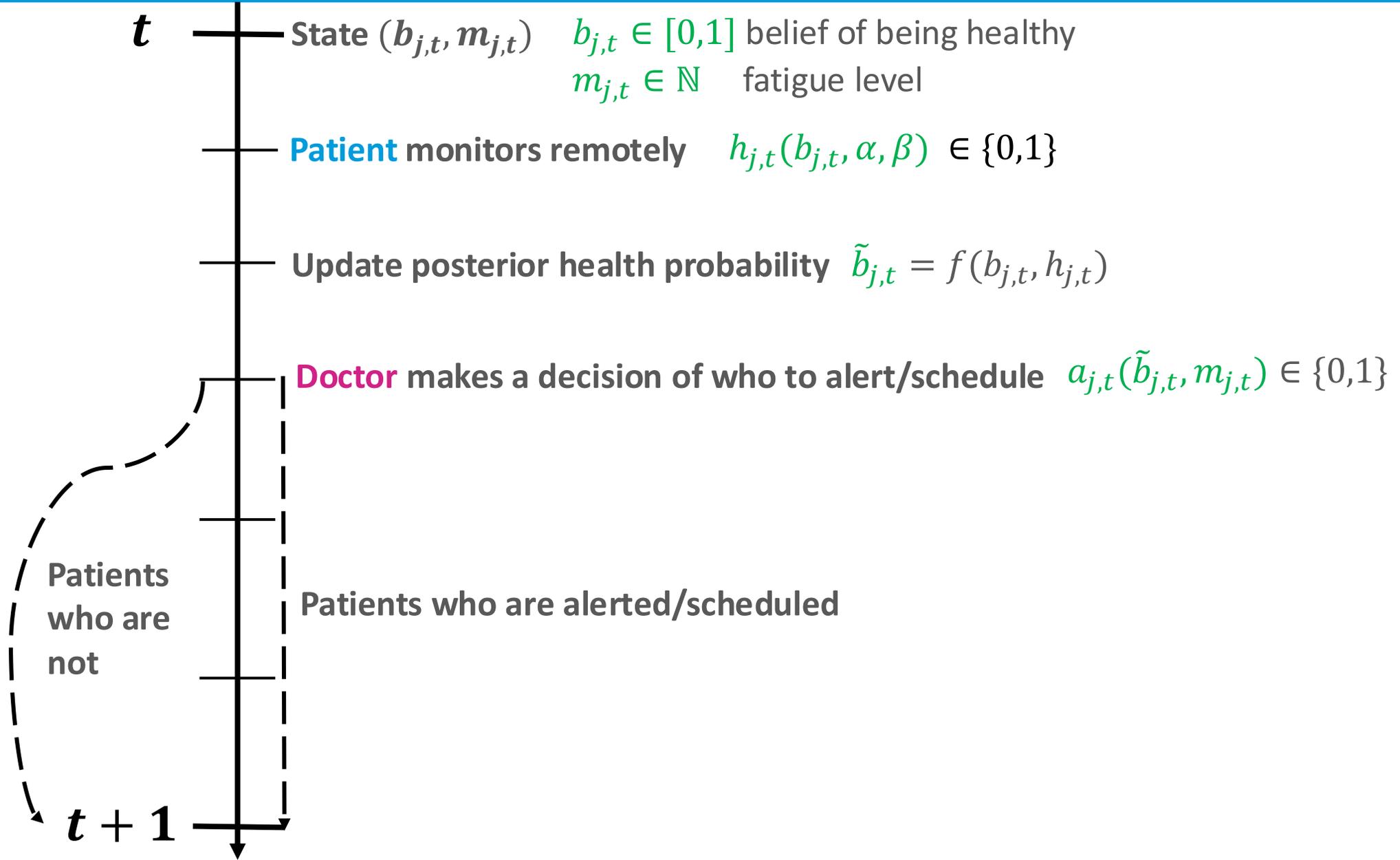
Model: sequence of events in period t



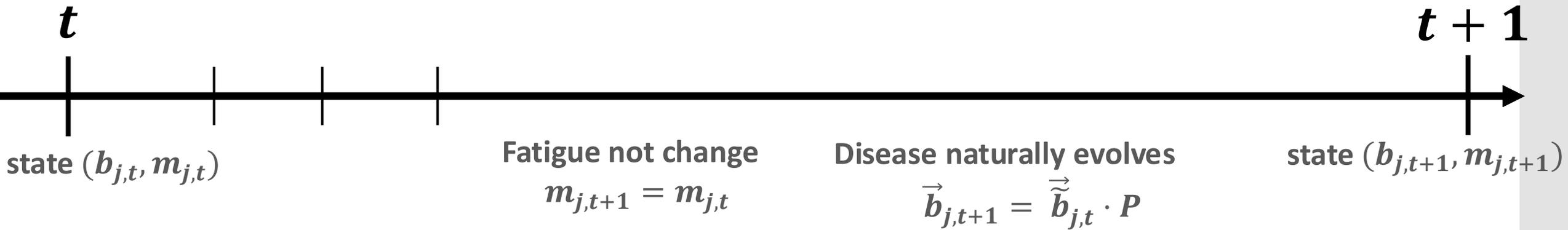
Model: sequence of events in period t



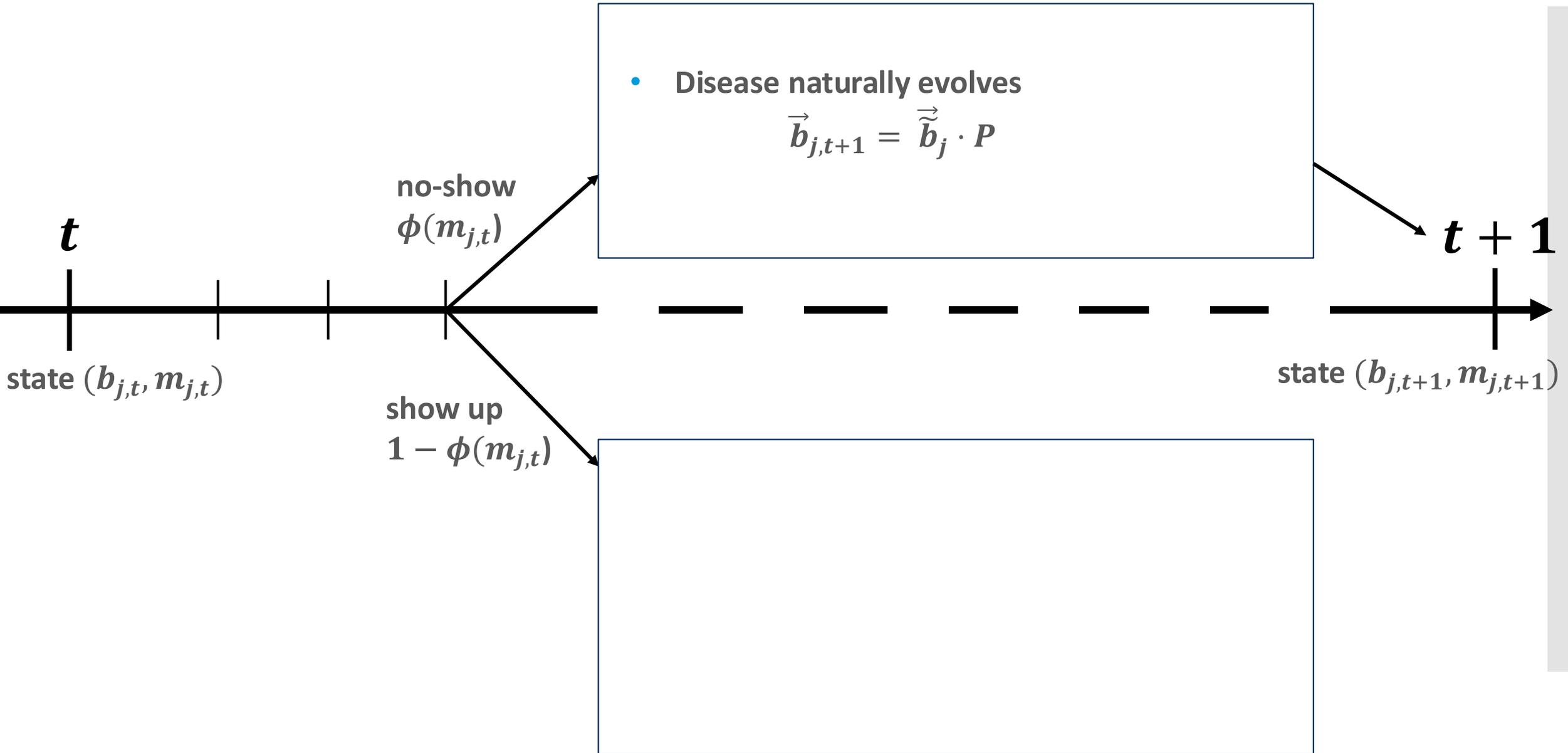
Model: sequence of events in period t



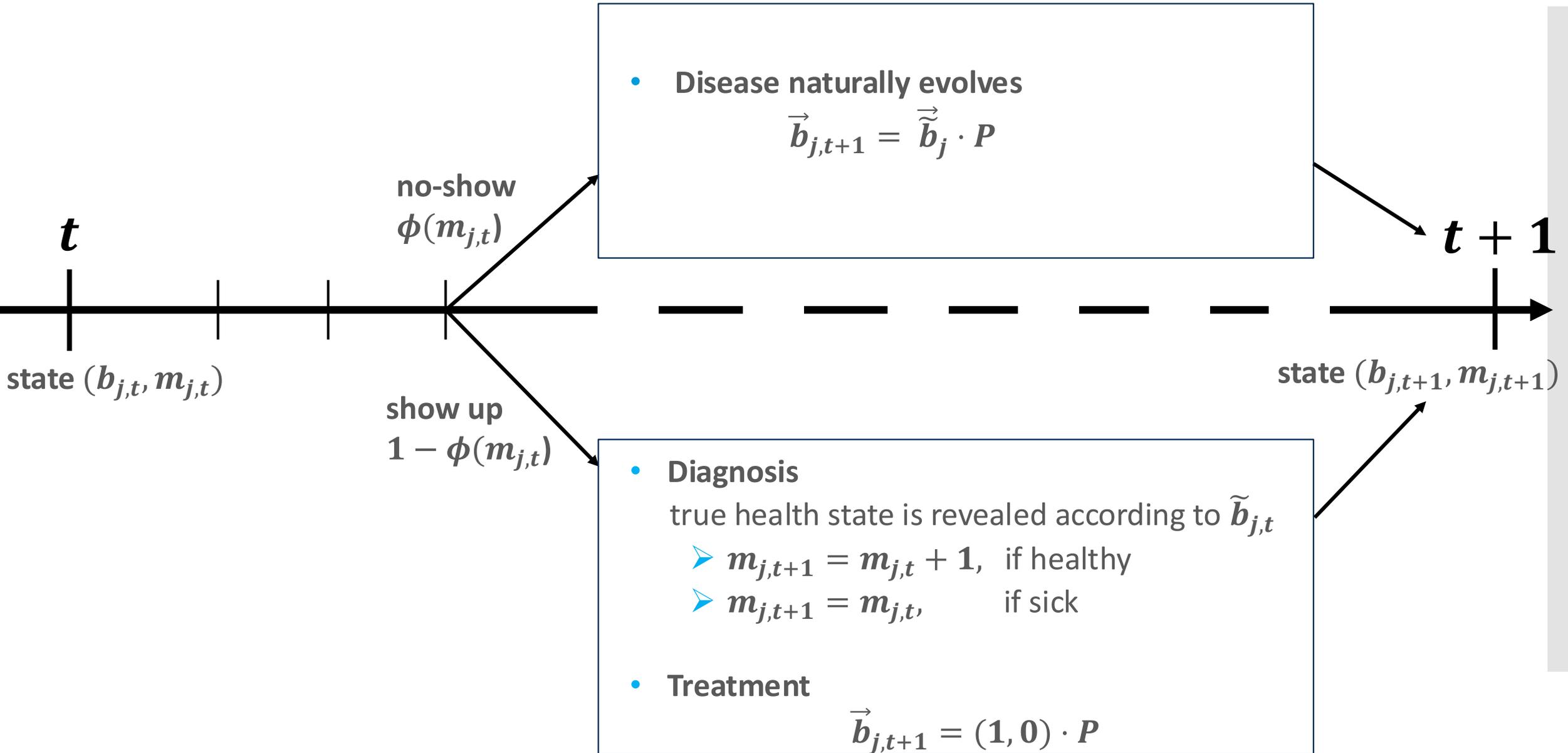
Model: patients who are not alerted/scheduled



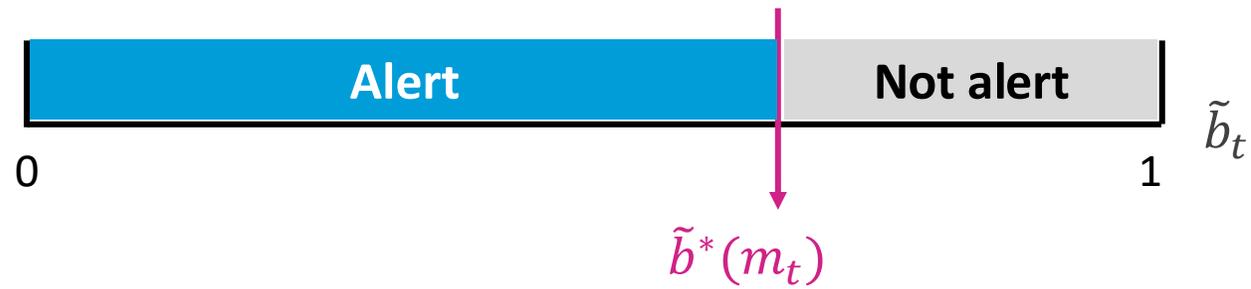
Model: patients who are alerted/scheduled



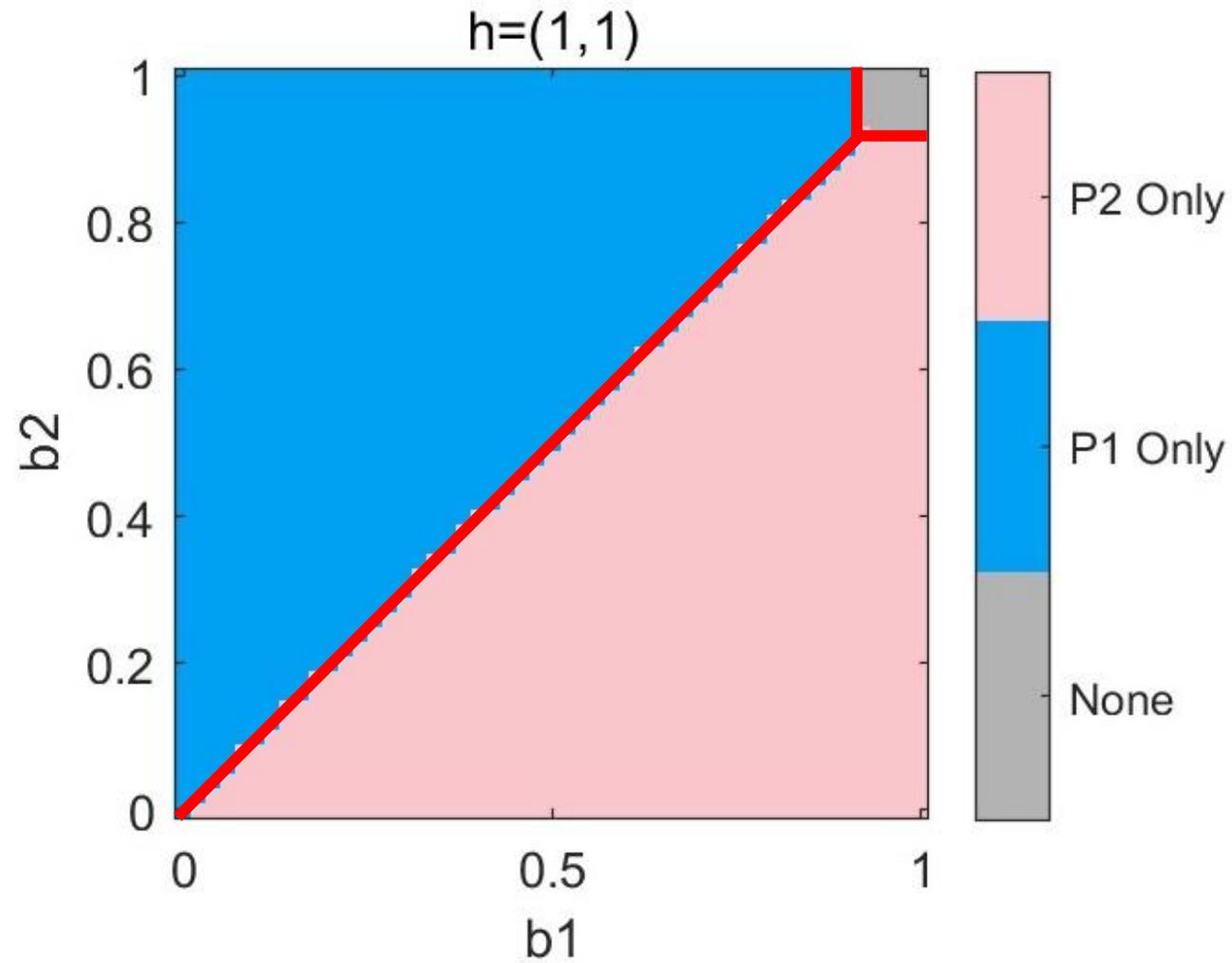
Model: patients who are alerted/scheduled



Optimal policy: single-patient case



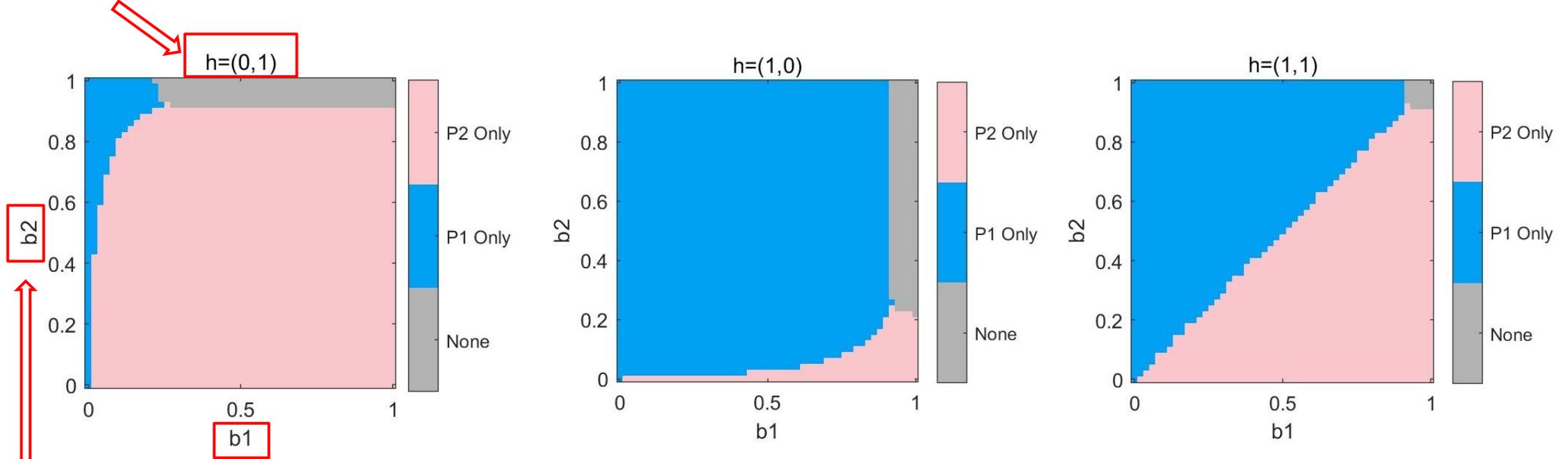
Optimal policy: two-patient case



Optimal policy: two-patient case

Decision Regions at $m_1 = 2, m_2 = 2$

Monitoring Signals



The belief of patient 1 being healthy

The belief of patient 2 being healthy

A quick look at a few more results

- The threshold policy can be generalized to N -patient case
 - Given system states (\mathbf{m}, \mathbf{h}) , the N dimensional belief space (b_1, \dots, b_n) can be partitioned into $N + 1$ unique optimal decision regions: not alert anyone, or alert patient i for $i \in \{1, \dots, N\}$

A quick look at a few more results

- The threshold policy can be generalized to N -patient case
 - Given system states (\mathbf{m}, \mathbf{h}) , the N dimensional belief space (b_1, \dots, b_n) can be partitioned into $N + 1$ unique optimal decision regions: not alert anyone, or alert patient i for $i \in \{1, \dots, N\}$
- Lagrangian truncation policy for larger capacity
 - Relax the constraint $\sum_i a_{i,t} \leq C$ (where $C \geq 1$) by introducing a nonnegative Lagrangian multiplier λ_t and decouple into the N subproblems
 - Solve each subproblem with the optimal Lagrangian multiplier λ^* and truncate if exceeding the capacity C
 - Asymptotically optimal and perform well in small systems

A quick look at a few more results

- The threshold policy can be generalized to N -patient case
 - Given system states (\mathbf{m}, \mathbf{h}) , the N dimensional belief space (b_1, \dots, b_n) can be partitioned into $N + 1$ unique optimal decision regions: not alert anyone, or alert patient i for $i \in \{1, \dots, N\}$
- Lagrangian truncation policy for larger capacity
 - Relax the constraint $\sum_i a_{i,t} \leq C$ (where $C \geq 1$) by introducing a nonnegative Lagrangian multiplier λ_t and decouple into the N subproblems
 - Solve each subproblem with the optimal Lagrangian multiplier λ^* and truncate if exceeding the capacity C
 - Asymptotically optimal and perform well in small systems
- Some numerical findings
 - RPM brings about 5% increase on the QALY than without RPM
 - Ignoring alert fatigue (i.e., endogenous no-show behavior) is equivalent to increasing the false positive rate of the RPM device by 15%

Thank you!

A Bayesian Decision-Theoretic Framework for Adaptive Learning in Phase I Cancer Trials

Shouhao Zhou
Pennsylvania State University

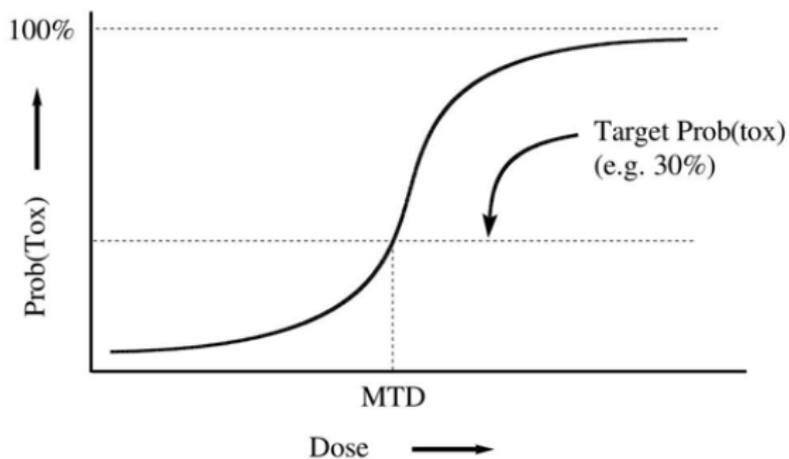
A joint work with Chenqi Fu (Penn State), Xinying Fang (Penn State),
David Madigan (Northeastern), and J Jack Lee (MD Anderson)

IMSI: Advances in Quantitative Medical Care

February 4, 2026

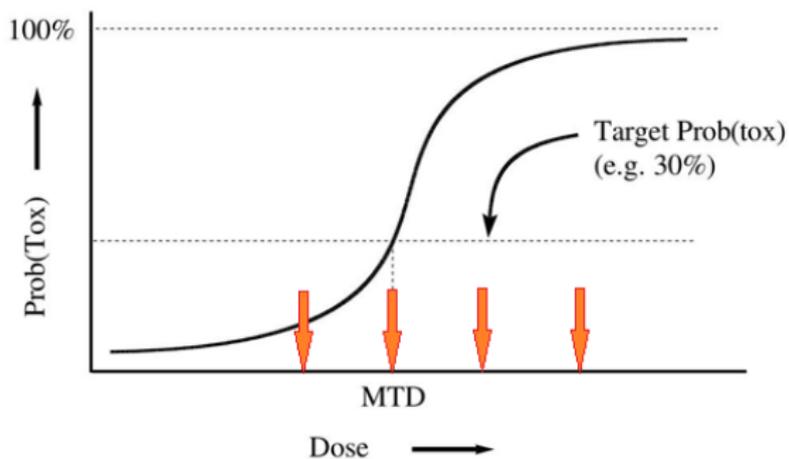
PHASE I CLINICAL TRIALS

Figure: Dose-Toxicity Curve



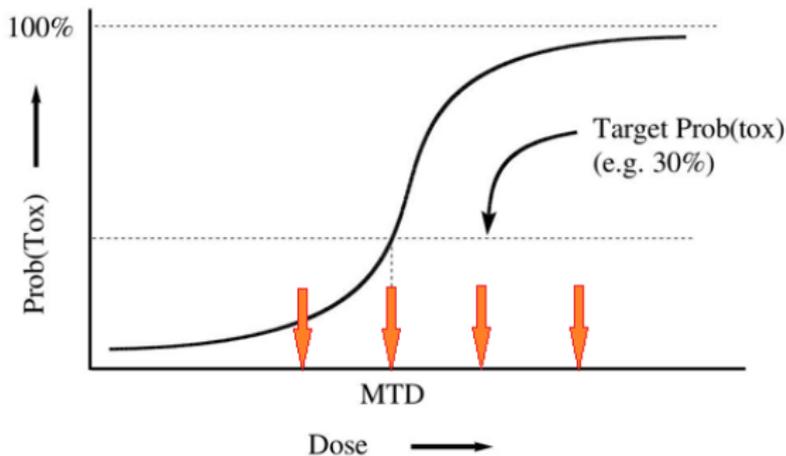
PHASE I CLINICAL TRIALS

Figure: Dose-finding in phase I trials



PHASE I CLINICAL TRIALS

Figure: Dose-finding in phase I trials

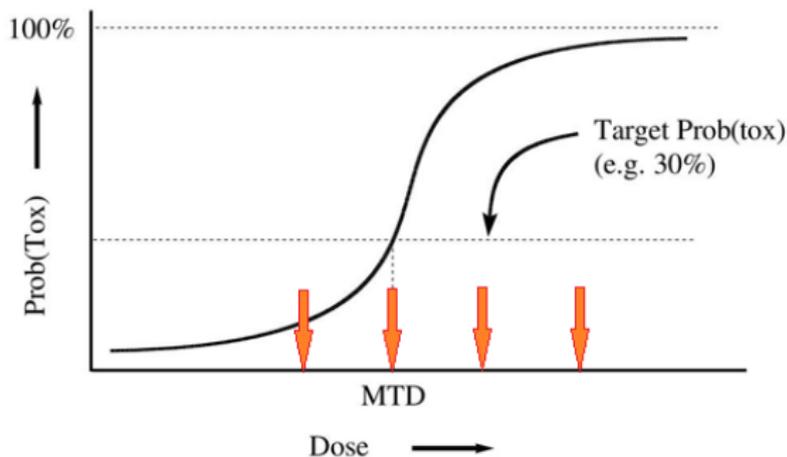


Challenges in phase I trial design:

- Sequential data
- Binary outcomes
- Small sample size

PHASE I CLINICAL TRIALS

Figure: Dose-finding in phase I trials



Phase I cancer trials:

- 1,500+ new trials initiated annually in U.S.
- Average \$8.1M cost, 32m duration (Sertkaya et al, 2024)

TABLE OF CONTENTS

- A Bayesian Decision Theoretic Adaptive Learning Framework
 - Motivation
 - Predictive Bayes Factors
 - Posterior Predictive (PoP) Design
 - Theoretical Property
 - Empirical Performance + Implementation Tools

PHASE I CLINICAL TRIALS

Primary goal: Given several doses, to assess dose limiting toxicities (**DLT**) and estimate the maximum tolerated dose (**MTD**).

- Algorithmic designs
 - Dose transition based on a set of prespecified rules/algorithm
 - Examples: 3+3 design; interval 3+3 design (Liu et al., 2020)
- Model-based designs
 - Dose transition based on an updated dose-toxicity model using accrued data
 - Examples: Continual Reassessment Method (CRM) design (O'Quigley et al., 1990), BMA-CRM design (Yin and Yuan, 2009)
- Model-assisted / Interval-based designs
 - Dose transition based on statistical decision model, and dose transition rules pretabulated similar to the algorithmic designs
 - Examples: Bayesian Optimal INterval (BOIN) design (Liu and Yuan, 2015), Keyboard design (Yan et al., 2017), **Posterior Predictive (PoP) design** (Fu et al., 2025)

PRINCIPAL OF INTERVAL-BASED TRANSITION

Given a target DLT rate ϕ , if true toxicity probability p_j at current dose level j is known, there are 3 possible **Oracle** decisions for next patient cohort assignment.

1. Retain, if $p_j = \phi$; $H_{0j} : p_j = \phi$
2. Escalate, if $p_j < \phi$; $H_{1j} : p_j < \phi$
3. Deescalate, if $p_j > \phi$. $H_{2j} : p_j > \phi$

PRINCIPAL OF INTERVAL-BASED TRANSITION

Given a target DLT rate ϕ , if true toxicity probability p_j at current dose level j is known, there are 3 possible **Oracle** decisions for next patient cohort assignment.

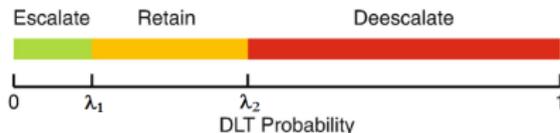
1. Retain, if $p_j = \phi$; $H_{0j} : p_j = \phi$
2. Escalate, if $p_j < \phi$; $H_{1j} : p_j < \phi$
3. Deescalate, if $p_j > \phi$. $H_{2j} : p_j > \phi$

PRINCIPAL OF INTERVAL-BASED TRANSITION

Given a target DLT rate ϕ , if true toxicity probability p_j at current dose level j is known, there are 3 possible **Oracle** decisions for next patient cohort assignment.

1. Retain, if $p_j = \phi$; $H_{0j} : p_j = \phi$
2. Escalate, if $p_j < \phi$; $H_{1j} : p_j < \phi$
3. Deescalate, if $p_j > \phi$. $H_{2j} : p_j > \phi$

Without knowing p_j , empirically **escalate** if $\hat{p}_j \leq \lambda_1$ or **deescalate** if $\hat{p}_j \geq \lambda_2$



Revise the decision rule given $\{n_j, \hat{p}_j\}$ for next patient cohort assignment,

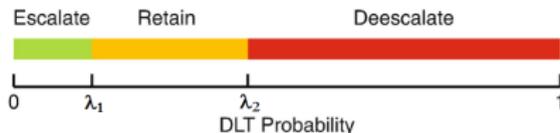
1. Retain, if $\hat{p}_j = \phi$;
 2. Escalate, if $\hat{p}_j < \lambda_1(\phi, n_j)$;
 3. Deescalate, if $\hat{p}_j > \lambda_2(\phi, n_j)$.
- $\hat{p}_j = y_j/n_j$: the observed toxicity rate at dose j
 - y_j : the number of patients experienced DLT at dose j
 - n_j : the number of patients treated at dose j

BAYESIAN OPTIMAL INTERVAL (BOIN) DESIGN

Given a target DLT rate ϕ , if true toxicity probability p_j at current dose level j is known, there are 3 possible **Oracle** decisions for next patient cohort assignment.

1. Retain, if $p_j = \phi$; $H_{0j} : p_j = \phi$
2. Escalate, if $p_j < \phi$; $H_{1j} : p_j < \phi$
3. Deescalate, if $p_j > \phi$. $H_{2j} : p_j > \phi$

Without knowing p_j , empirically **escalate** if $\hat{p}_j \leq \lambda_1$ or **deescalate** if $\hat{p}_j \geq \lambda_2$

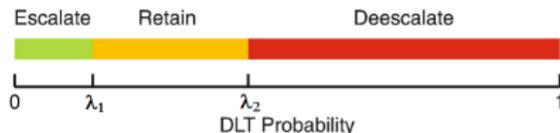


Revise the decision rule given \hat{p}_i for next patient cohort assignment,

1. Retain, if $\hat{p}_j = \phi$; $H_{0j} : p_j = \phi$
 2. Escalate, if $\hat{p}_j < \lambda_1$; $H_{1j} : p_j = \phi_1$
 3. Deescalate, if $\hat{p}_j > \lambda_2$. $H_{2j} : p_j = \phi_2$
- $\hat{p}_j = y_j/n_j$: the observed toxicity rate at dose j
 - y_j : the number of patients experienced DLT at dose j
 - n_j : the number of patients treated at dose j

BAYESIAN OPTIMAL INTERVAL (BOIN) DESIGN

Decision rule: Escalate if $\hat{p}_j \leq \lambda_1$ or Deescalate if $\hat{p}_j \geq \lambda_2$



- The decision for dose transition could be **wrong** if $\hat{p}_j \neq p_j$
 - Escalate/retain when the current dose is above the MTD
 - Deescalate/retain when the current dose is below the MTD
 - Escalate/deescalate when the current dose is the MTD
- Recommended by [FDA *Fit-for-Purpose Initiative*](#) (December, 2021)
- Optimal in a sense to minimize the risk of incorrect decision

STATISTICAL HYPOTHESIS TESTING

The rationale of hypothesis testing in BOIN:

$$H_{0j} : p_j = \phi \quad (= 25\%)$$

$$H_{1j} : p_j = \phi_1 \quad (= 15\%)$$

$$H_{2j} : p_j = \phi_2 \quad (= 35\%)$$

1. Retain, if $\hat{p}_j = \phi$;

2. Escalate, if $\hat{p}_j < \lambda_1$;

3. Deescalate, if $\hat{p}_j > \lambda_2$.

- ϕ_1 is the highest toxicity rate that is subtherapeutic
- ϕ_2 is the lowest toxicity rate that is overly toxic
- **Arbitrarily**, set $\phi_1 = 0.6 \times \phi$ and $\phi_2 = 1.4 \times \phi$

For example, if $\phi = 25\%$,

$$\phi_1 = 0.6 \times \phi = 0.6 \times 25\% = 15\%$$

$$\phi_2 = 1.4 \times \phi = 1.4 \times 25\% = 35\%$$

- **Lack** of theoretical justification on the *local* choice of ϕ_1 and ϕ_2 .
- Decision boundaries $\lambda_1 = \lambda_1(\phi_1)$ and $\lambda_2 = \lambda_2(\phi_2)$ determined by Bayesian decision theory, **dependent** on hypothesis boundaries ϕ_1, ϕ_2 .

STATISTICAL HYPOTHESIS TESTING

The rationale of hypothesis testing in BOIN:

$$H_{0j} : p_j = \phi (= 25\%)$$

$$H_{1j} : p_j = \phi_1 (= 15\%)$$

$$H_{2j} : p_j = \phi_2 (= 35\%)$$

? \mapsto

$$H_{0j} : p_j = \phi (= 25\%)$$

$$H_{1j} : p_j < \phi (< 25\%)$$

$$H_{2j} : p_j > \phi (> 25\%)$$

- ϕ_1 is the highest toxicity rate that is subtherapeutic
- ϕ_2 is the lowest toxicity rate that is overly toxic
- **Arbitrarily**, set $\phi_1 = 0.6 \times \phi$ and $\phi_2 = 1.4 \times \phi$

For example, if $\phi = 25\%$,

$$\phi_1 = 0.6 \times \phi = 0.6 \times 25\% = 15\%$$

$$\phi_2 = 1.4 \times \phi = 1.4 \times 25\% = 35\%$$

- Lack of **theoretical justification** on the *local* choice of ϕ_1 and ϕ_2 .
- Decision boundaries $\lambda_1 = \lambda_1(\phi_1)$ and $\lambda_2 = \lambda_2(\phi_2)$ determined by Bayesian decision theory, **dependent** on hypothesis boundaries ϕ_1, ϕ_2 .
- May converge to **suboptimal** dose even with infinite sample size.

BAYES FACTORS AND LINDLEY'S PARADOX

Consider two competing Bayesian models $M_k(\mathbf{y})$ and $M_{k'}(\mathbf{y})$ with parameter priors $\pi_k(\theta_k)$ and $\pi_{k'}(\theta_{k'})$ for hypotheses H_k vs $H_{k'}$.

Bayes factor comparing M_k to $M_{k'}$ is defined as

$$BF_{k,k'} = \frac{p(\mathbf{y} | M_k)}{p(\mathbf{y} | M_{k'})} = \frac{\int \prod p(y_i | \theta_k, M_k) \pi_k(\theta_k) d\theta_k}{\int \prod p(y_i | \theta_{k'}, M_{k'}) \pi_{k'}(\theta_{k'}) d\theta_{k'}},$$

which measures the relative evidence provided by the observed data $\mathbf{y} = \{y_1, \dots, y_n\}$ in favor of M_k versus $M_{k'}$.

- Natural interpretability of model (hypothesis) preference

Lindley's Paradox

As the sample size n increases:

- Classical hypothesis tests may strongly reject a null model.
- The Bayes factor may simultaneously favor the null model,

particularly when diffuse or weakly informative priors are used.

BAYES FACTORS AND LINDLEY'S PARADOX

Consider two competing Bayesian models $M_k(\mathbf{y})$ and $M_{k'}(\mathbf{y})$ with parameter priors $\pi_k(\theta_k)$ and $\pi_{k'}(\theta_{k'})$ for hypotheses H_k vs $H_{k'}$.

Bayes factor comparing M_k to $M_{k'}$ is defined as

$$BF_{k,k'} = \frac{p(\mathbf{y} | M_k)}{p(\mathbf{y} | M_{k'})} = \frac{\int \prod p(y_i | \theta_k, M_k) \pi_k(\theta_k) d\theta_k}{\int \prod p(y_i | \theta_{k'}, M_{k'}) \pi_{k'}(\theta_{k'}) d\theta_{k'}},$$

which measures the relative evidence provided by the observed data $\mathbf{y} = \{y_1, \dots, y_n\}$ in favor of M_k versus $M_{k'}$.

- Natural interpretability of model (hypothesis) preference

Lindley's Paradox

As the sample size n increases:

- Classical hypothesis tests may strongly reject a null model.
- The Bayes factor may simultaneously favor the null model, particularly when diffuse or weakly informative priors are used.

PREDICTIVE BAYES FACTORS

For Bayesian models $M_k(\mathbf{y})$ with posterior predictive density $p_k(y|\mathbf{y})$:

Theorem 1

Under mild regularity conditions,

$$\sum_i \log p_k(y_i|\mathbf{y}) + \hat{b}_k$$

is an asymptotic unbiased estimator of $n \cdot E_{\tilde{y}} \log p_k(\tilde{y}|\mathbf{y})$, where the term $\hat{b}_k = -tr\{J_{n,k}^{-1}(\hat{\theta}^k)I_{n,k}(\hat{\theta}^k)\}$ corrects the over-estimation bias for the ‘double use’ of data \mathbf{y} .

Therefore, we present **predictive Bayes factor (PrBF)** (Zhou & Madigan, 2026+):

$$PrBF_{k,k'} := \frac{\prod_i p_k(y_i|\mathbf{y})}{\prod_i p_{k'}(y_i|\mathbf{y})} \cdot \frac{\exp(\hat{b}_k)}{\exp(\hat{b}_{k'})}$$

as a measure of the weight of predictive sample evidence in favor of $M_k(\mathbf{y})$ compared with $M_{k'}(\mathbf{y})$.

PROPERTIES OF PREDICTIVE BAYES FACTORS

- Compare *posterior predictive* Bayesian models rather than *prior predictive* models.
- Rectify the over-estimation error with **asymptotic unbiased** estimator.
- Outstanding **finite sample** performance.
- The error correction in $PrBF$ holds for **mis-specified models**.
- Inherit the property of **coherence** (Hu and Johnson, 2009).

$$PrBF_{k,k''} = PrBF_{k,k'} \times PrBF_{k',k''}$$

- **Avoid the degeneration** of the integrated likelihood.
- **Elucidate the Lindley's paradox**.
- **Reduce the sensitivity** to variations in the prior distribution.
- **Simplify** the computation and stabilize the estimation.

CALIBRATION OF PREDICTIVE BAYES FACTORS

Table: The interpretation of the posterior probability in favor of *predictive* model $\mathcal{M}_1(\mathbf{y})$ using the predictive Bayes factor (*PrBF*).

$PrBF_{12}$	$\mathcal{M}_1(\mathbf{y})$ WEIGHT	EVIDENCE
1 to 3	50% to 75%	NOT WORTH MORE THAN A BARE MENTION
3 to 19	75% to 95%	POSITIVE
19 to 99	95% to 99%	STRONG
> 99	> 99%	VERY STRONG

POSTERIOR PREDICTIVE DESIGN

The statistical hypothesis for dose transition:

BOIN

$$H_{0j} : p_j = \phi \quad (= 25\%)$$

$$H_{1j} : p_j = \phi_1 \quad (= 15\%)$$

$$H_{2j} : p_j = \phi_2 \quad (= 35\%)$$

versus

PoP

$$H_{0j}^* : p_j = \phi \quad (= 25\%)$$

$$H_{1j}^* : p_j \neq \phi \quad (\neq 25\%)$$

For simplified Oracle hypotheses H_{0j}^* and H_{1j}^* ,¹

$$PrBF_{0,1} = e (n_j + 2)^{n_j} \left(\frac{\phi}{y_j + 1} \right)^{y_j} \left(\frac{1 - \phi}{n_j - y_j + 1} \right)^{n_j - y_j} .$$

Transition rule for PoP design:

- Retain, *iff* $PrBF_{0,1} \geq C$.
- Exclude, *iff* $PrBF_{0,1} < E$.

¹Under H_{0j}^* , $PrBF_{0,1} \rightarrow e$; Under H_{1j}^* , $PrBF_{0,1} \rightarrow 0$.

OPTIMAL BAYESIAN DECISIONS

Three potential actions:

- retain the dose: \mathcal{R}
- transit but not eliminate the dose: \mathcal{T}
- transit and eliminate the dose: \mathcal{E}

Loss function $L(a(y_j), H)$:

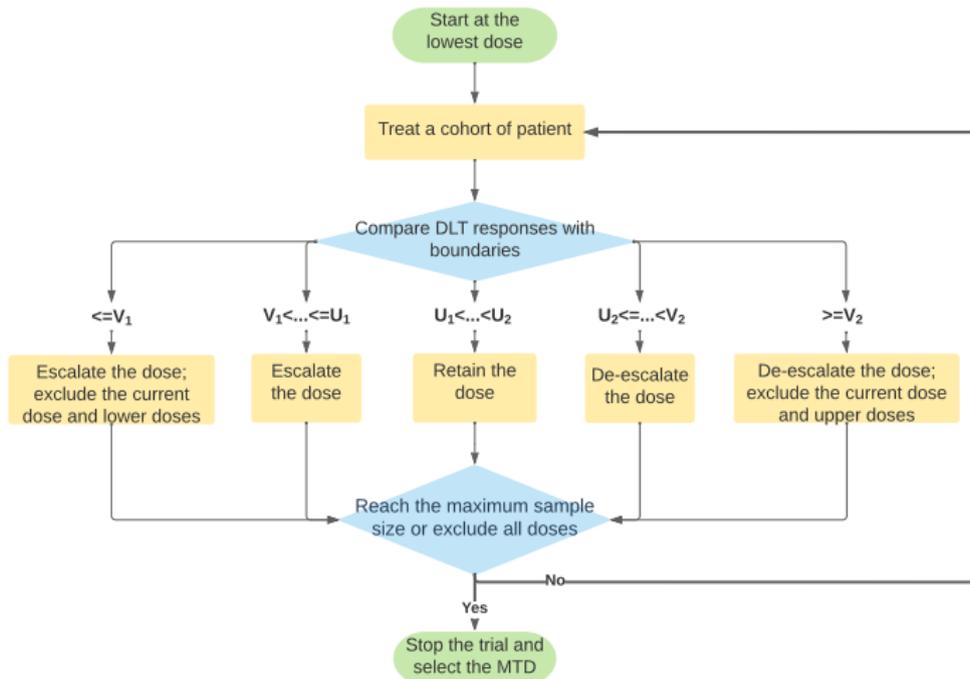
	\mathcal{R}	\mathcal{T}	\mathcal{E}
$H_0 : \pi_j = \phi$	0	b_1	1
$H_1 : \pi_j \neq \phi$	b_2	b_3	0

Theorem 2: Global Optimality

PoP design minimizes the risk of incorrect decision of dose assignment

$$R(a) = P(H_0) \int L(a(y), H_0) dy + P(H_1) \int L(a(y), H_1) dy$$

when $C = \frac{b_2 - b_3}{b_1}$ and $E = \frac{b_3}{1 - b_1}$.



Simplicity / Transparency

Under the PoP design, dose allocation can be pretabulated similar to algorithmic designs for easy implementation.

COMPARISON OF DECISION BOUNDARIES

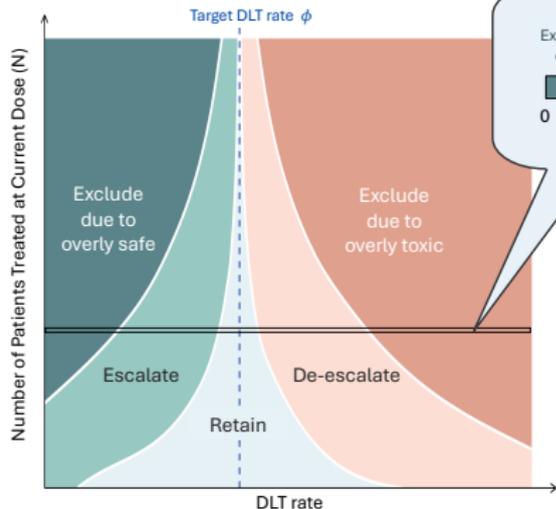
$\phi_0 = 0.25$	Number of patients treated at the current dose									
	3	6	9	12	15	18	21	24	27	30
PoP Design										
Escalation if no. of DLT \leq	0	0	1	2	3	3	4	5	6	6
De-escalation if no. of DLT \geq	1	2	3	4	5	5	6	7	8	9
BOIN Designs										
Escalation if no. of DLT \leq	0	1	1	2	2	3	4	4	5	5
De-escalation if no. of DLT \geq	1	2	3	4	5	6	7	8	9	9

Safety / Efficiency

Under the PoP design, dose transition is conservative when information is limited, but efficient to exploit data for dose allocation.

Theorem 3: Consistency

Under the PoP design, dose allocation and selection almost surely converge to the target MTD.



Example of Decision Boundaries ($\phi=0.25$)

	Number of patients treated at current dose (N)										
	N	3	6	9	12	15	18	21	24	27	30
Overly safe exclusion if no. of DLT $\leq V_1$	V_1	-	-	-	-	0	0	1	1	2	2
Escalation if no. of DLT $\leq U_1$	U_1	0	0	1	2	2	3	4	4	5	6
De-escalation if no. of DLT $\geq U_2$	U_2	2	3	3	4	5	6	7	7	8	9
Overly toxic exclusion if no. of DLT $\geq V_2$	V_2	3	5	6	7	8	9	11	12	13	14

RATE OF CONVERGENCE

Denote L_n and U_n the lower and upper PoP transition boundaries.

Theorem 4: Convergence rate

For $C < e$, we have the boundaries of PoP design satisfy $|L_n - \phi| \leq kn^{-1/2}$ and $|U_n - \phi| \leq kn^{-1/2}$ as $n \rightarrow \infty$, where $k = \sqrt{2\phi(1-\phi)(1-\log C)}$.

Both L_n and U_n converge to ϕ at the **root- n rate** $O(n^{-1/2})$, which matches the **optimal** rate attainable for learning a probability parameter in binomial models under regularity conditions, implying that the PoP design is highly **efficient** in identifying MTD.

Coherence

The PoP design is (long-memory) coherent in the sense that the design will never escalate the dose when the observed toxicity rate \hat{p}_j at the current dose is higher than the target toxicity rate ϕ ; and will never deescalate the dose when the observed toxicity rate \hat{p}_j at the current dose is lower than the target toxicity rate ϕ .

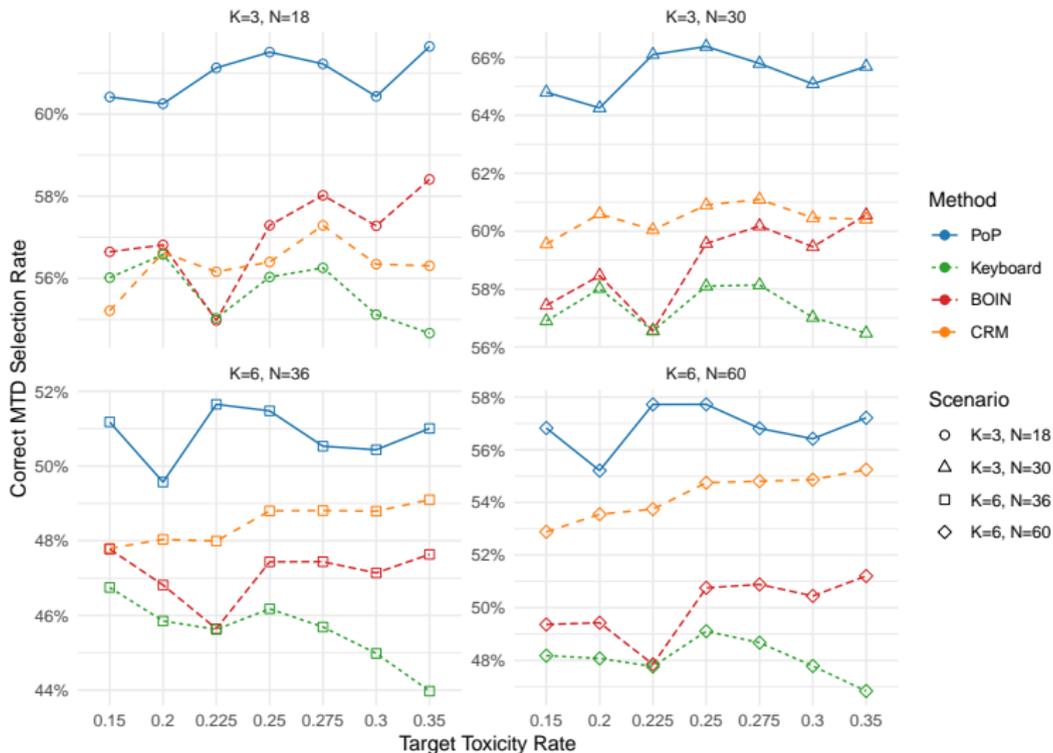
Example: for a target toxicity rate $\phi = 30\%$,

- if $y_j/n_j = 1/3 > \phi$, PoP design will *never* escalate dose;
- if $y_j/n_j = 0/3 < \phi$, PoP design will *never* deescalate dose.

SIMULATION STUDY

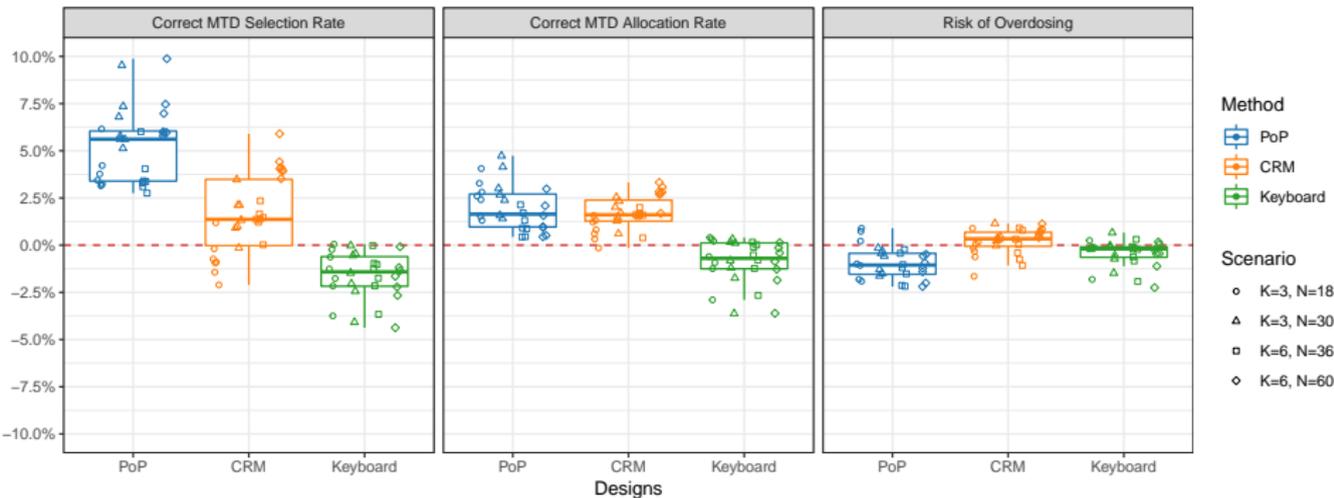
- Total sample size $n = 18$ ($K=3$), 30 ($K=3$), 36 ($K=6$) or 60 ($K=6$), target DLT rate $\phi = \underline{0.15}$ to 0.35 .
- 10,000 dose-toxicity scenarios randomly generated using the pseudo-uniform algorithm (Clertant and O'Quigley, 2017)
- 20,000 trials under each scenario
 1. percentage of correct selection (PCS) of true MTD
 2. percentage of correct patients allocated (PCA) to MTD
 3. risk of overdosing
 4. average sample size.
- CRM: skeleton based on indifference-interval (Lee and Cheung, 2009).
- BOIN: set boundaries as $0.6 \times \phi$ and $1.4 \times \phi$ (recommended) and assigned equal prior probability to the hypotheses
- Keyboard: use the default values for key width of 0.1 (recommended)
- PoP: $C = 2.5$, $E = 5/24$ (based on $b_1 = 0.2$, $b_2 = 2/3$, $b_3 = 1/6$).
- All of the methods incorporate accelerated titration steps.

PERCENTAGE OF CORRECT SELECTION (PCS)



PoP, KEYBOARD, AND CRM VS BOIN

OC Comparison of Design Methods



Take-home message:

- Predictive Bayes factor is posterior predictive-based, powerful for general-purpose Bayesian hypothesis testing, model selection, or ensemble inference.
- Bayesian decision-theoretic framework is efficient in optimizing adaptive learning in Phase I cancer trials, projecting more than \$100M cost savings with the minimum 2.5% MTD selection improvement.

IMPLEMENTATION AND SOFTWARE

- PoP has been implemented in an ongoing phase I/II study (NCT06541262) with adaptive dose assignment applied for ten enrolled pediatric oncology patients.
- R package *PoPdesign* on CRAN
- R Shiny https://xinying-fang.shinyapps.io/PoPdesign_Shiny/

PoPdesign: Posterior Predictive (PoP) Design for Phase I Clinical Trials

Cheng Fu¹, Xinying Fang¹, J. Jack Lee², Shouhao Zhou¹

¹Division of Biostatistics and Bioinformatics, Department of Public Health Sciences, Pennsylvania State University

²Department of Biostatistics, The University of Texas MD Anderson Cancer Center

PoPdesign R package

Navigation: Trial Setting | Simulation | Trial Protocol | Select MTD | Reference

Design Flow Chart

Decision Table

Download flowchart

Number of doses: 5 Starting dose level: 1

Cohort size: 3 Number of cohort: 10

Target Toxicity Probability α : 0.25

Risk Cutoff: 0.5

Get Decision Table

KEY REFERENCES

- Fu, C., Zhou, S., & Lee, J. J. (2025). Posterior predictive design for Phase I clinical trials. *Journal of the American Statistical Association*. 1-12. doi: 10.1080/01621459.2025.2484044.
- Zhou, S., & Madigan, D. (2025+). Predictive Bayes factors. Under review.
- Zhou, S., Fu, C., Fang, X., & Lee, J. J. (2025+). PoP: Next-Generation Model-Assisted Phase I Design with Global Optimality. Under review.

Thank you!

contact: Shouhao Zhou ✉ szhou1@pennstatehealth.psu.edu



An Integrated Approach for Claim Outcome Prediction and Denial Appeal Generation in Health Insurance

Daiwen Zhang (dwzhang@umich.edu)

IMSI Advances in Quantitative Medical Care Workshop
February 4, 2026



General Hospital Billing Process

Treatment Ordered

Treatment Provided



Prior Authorization Request
[upon payor's approval]



Payment Request
[upon payor's approval]



Burden and Costs

- HHS OIG reports in 2018-2022 on Medicare Advantage (> 30M enrollments):
 - 5% of prior authorization requests and 9.5% of payment requests are denied (2018. OEI-09-19-00350)
 - 1% of denials were appealed, 75% of appealed denials were overturned (2014-2016. OEI-09-16-00410)
 - 13% of prior authorization denials and 18% of payment denials met Medicare Coverage/Billing Rules (2022. OEI-09-18-00260. <https://oig.hhs.gov/report/all/>)
- Physicians and their staff spend 13 hours each week on Prior Authorizations
 - 2024 AMA Prior Authorization Physician Survey, <https://www.ama-assn.org/system/files/prior-authorization-survey.pdf>
 - 89 FR 8758 Final Rule <https://www.federalregister.gov/d/2024-00895/p-2124>



Scope

- Project at a hospital in US; summary on the technical aspects..
- We focus mainly on three aspects:
 1. Identify key factors that could affect the outcome of claims through data analytics
 2. ML models for predicting the outcome of claims;
 3. LLMs for generating communication letters (appeal letters).
- The goals are to provide insights to claim outcomes, while reducing the operational costs in denial management by automated processing, and hopefully improve the overall approval rate through customized letter generation.



Source of Data

1. Claims Administrative/Financial Records (payor information, charge details, etc.)
2. EHR FHIR Database
 - Patient demographics information & observation data
 - Clinical documents including Dx report and references
3. Communication history with the payor (letters and other messages/notices/packages)
4. Public Datasets
 - Location-based Health Risk Indices



Task 1: Exploratory Data Analytics

- Sankey Diagram: Payor -> Specialty -> Outcome
- Pearson's chi-square test on various factors
- Distinctions between different specialties

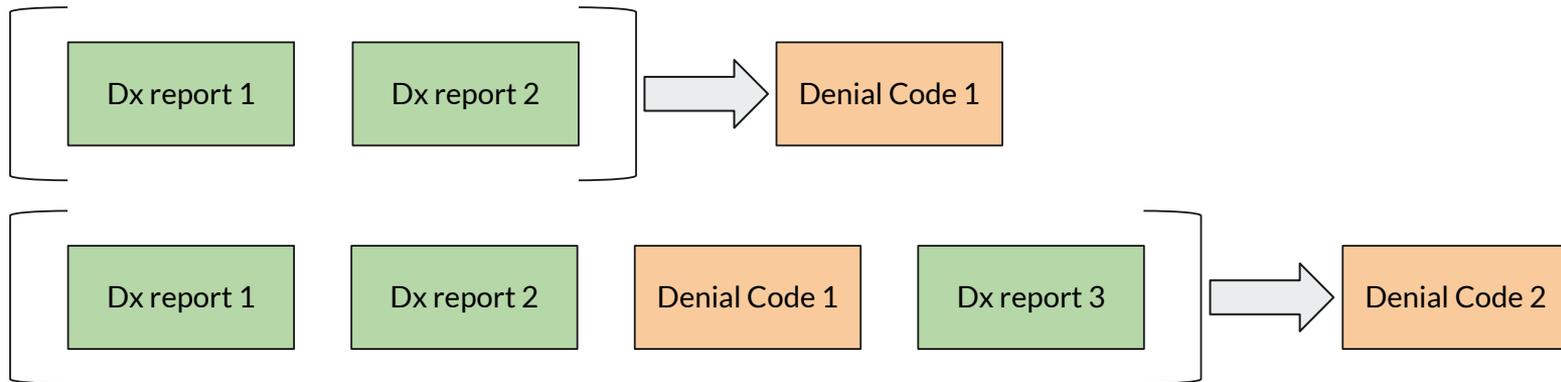


Task 2: Predictive NLP Models

- Predicting the outcome of claims with both structured tabular data and unstructured plain-text clinical documents.
- Clinical texts (training set) are used to fine-tune a BERT-like encoder model with classification head layer (BioClinical-ModernBERT-base).
- The fine-tuned model then encodes the text inputs into vector embedding.
- The tabular data together with the text embedding and/or the output logits are then fed into a Gradient Boosting Tree model (LightGBM).

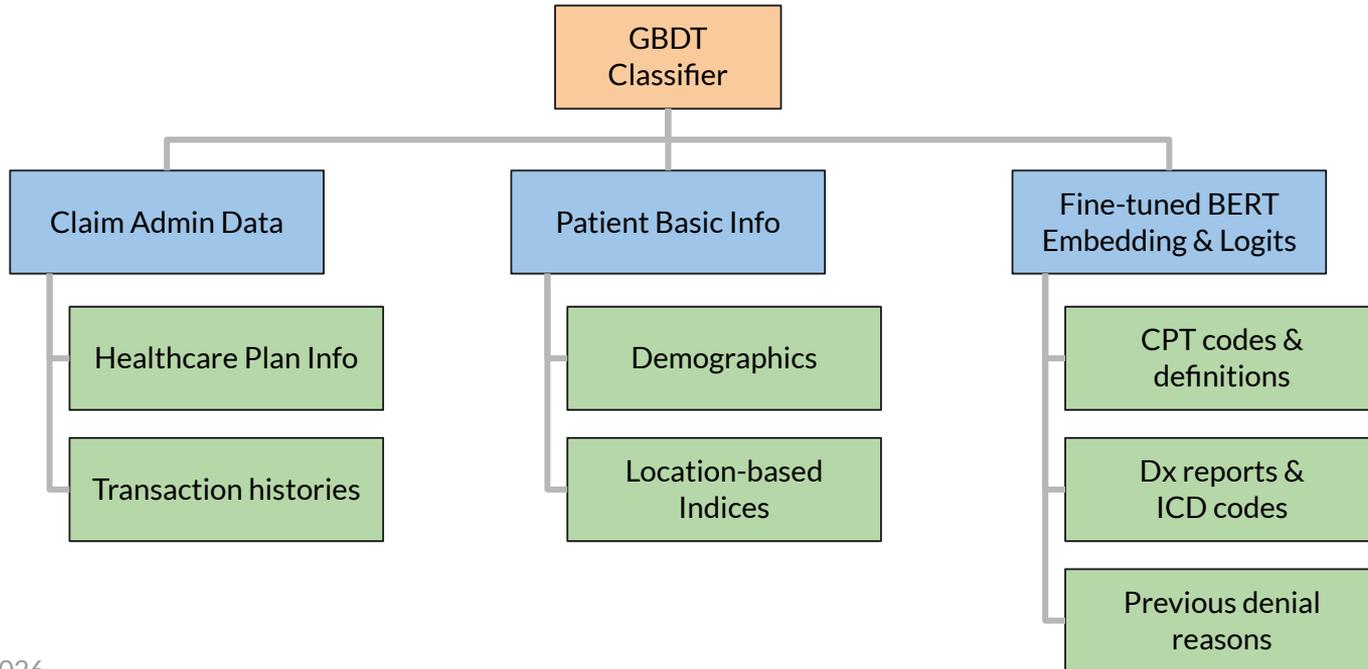


Predictive Model Architecture





Predictive Model Architecture





Preliminary Results Summary

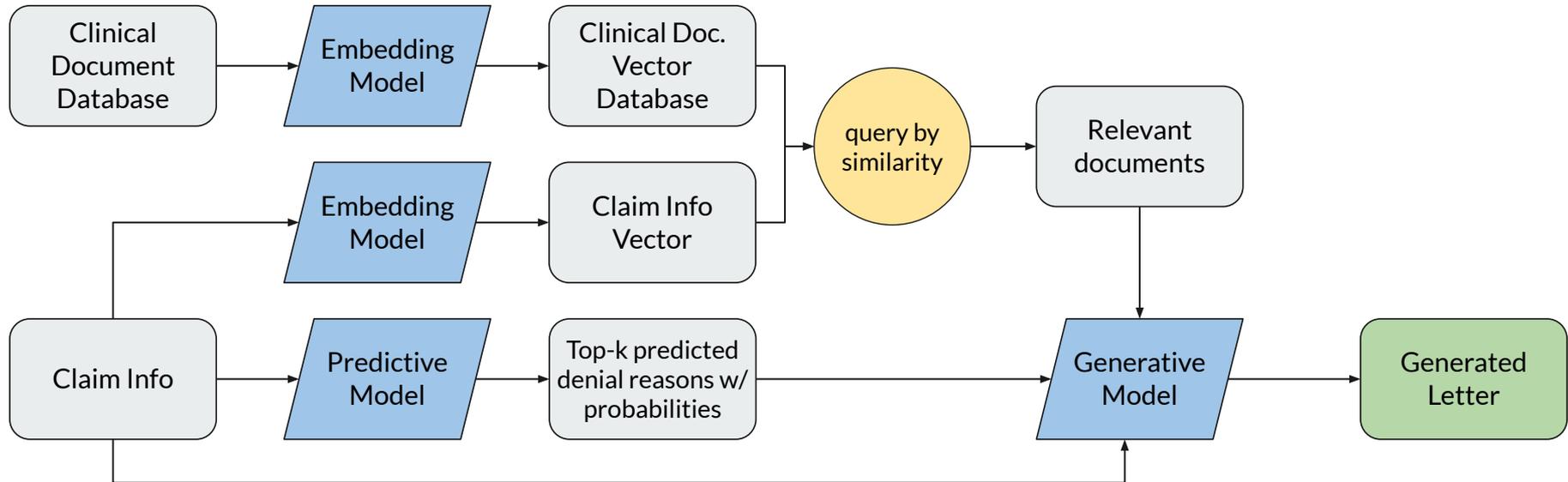
1. **Binary Classification: whether the next outcome is [denial] or [approval]**
 - Highly Imbalanced data
 - **Accuracy**: ~ 0.9; **AU-ROC**: 0.8 - 0.85
 - **AU-PRC** for [approval]: > 0.9
 - **AU-PRC** for [denial]: 0.4 - 0.8, performance varies by specialty
 - **Tabular + Text** provides 5% improvement over **Tabular-only**
2. **Binary Classification: whether the insurer paid a certain proportion of the expected amount**
 - 75% quartile: **All four metrics** > 0.9
 - 50% quartile: **All four metrics** > 0.9
 - 25% quartile: **Accuracy**: ~0.8; **AU-ROC**: ~0.85; **AU-PRC** for [below]: ~0.6



Preliminary Results Summary

1. Multi-class classification: predicting the next denial code
 - 90-300 possible denial codes
 - Output: Softmax
 - Performance vary by denial code, better on more common codes
 - **Weight Average F1**: >0.9; **Macro Average F1**: ~ 0.1
 - **Top-2 predictions** improved the recall without losing too much precision

Task 3: RAG System for Letter Generation





Task 3: RAG System for Letter Generation

- Inputs:
 1. Basic claims and patient information
 2. Top-2 denial reasons given of the predictive model
 3. Vector database of clinical documents
- The LLM is prompted to summarize the the selected documents and develop strategies to address the given potential denial reasons when generating the output letter.
- The relevance and accuracy are improved with more explicit instructions



Potential Future Directions

- BERT embedding vs. GPT embedding
- Grouping the denial codes to reduce the number of classes
- More interpretable models
- Systematic Evaluation on the Relevance and Accuracy of the RAG Document Retrieval
- Systematic Evaluation on the RAG generation with human feedback



Thank You!

Optimizing Systems of Risk-Appropriate Maternal Care

Abel Sapirstien¹ Meghan Meredith² Lauren N. Steimle¹

¹Industrial and Systems Engineering
Georgia Institute of Technology

²Center for Surgical & Transplant Applied Research
NYU Langone Health

Advances in Quantitative Medical Care, February 4th, 2026

Severe maternal morbidity is a persistent public health challenge.

"unexpected outcomes from labor & delivery that can result in significant short- or long-term health consequences" [3]

Severe maternal morbidity is a persistent public health challenge.

"unexpected outcomes from labor & delivery that can result in significant short- or long-term health consequences" [3]

SMM in the US:

- ~60,000 non-fatal events annually [3]
- Highest among high-income countries [5]
- 80% of SMM events are *preventable* [12]
- Modifiable risk factors: cesarean delivery, gestational diabetes [7]

Maternity ward are closing, increasing travel times to delivery.

Where do closures occur?

Across urban and rural hospitals

Rural hospitals in rural states have seen the most closures [9]

Most rural hospitals (52.4%) no longer have maternity wards [9]

Maternity ward are closing, increasing travel times to delivery.

Where do closures occur?

Across urban and rural hospitals

Rural hospitals in rural states have seen the most closures [9]

Most rural hospitals (52.4%) no longer have maternity wards [9]

What drives closures?

Maternity ward are closing, increasing travel times to delivery.

Where do closures occur?

Across urban and rural hospitals

Rural hospitals in rural states have seen the most closures [9]

Most rural hospitals (52.4%) no longer have maternity wards [9]

What drives closures?

High fixed costs associated with maternity wards.

Aging population \implies fewer births

Maternity ward are closing, increasing travel times to delivery.

Where do closures occur?

Across urban and rural hospitals

Rural hospitals in rural states have seen the most closures [9]

Most rural hospitals (52.4%) no longer have maternity wards [9]

What drives closures?

High fixed costs associated with maternity wards.

Aging population \implies fewer births

More rural births are paid by Medicaid

Maternity ward are closing, increasing travel times to delivery.

Where do closures occur?

Across urban and rural hospitals

Rural hospitals in rural states have seen the most closures [9]

Most rural hospitals (52.4%) no longer have maternity wards [9]

What drives closures?

High fixed costs associated with maternity wards.

Aging population \implies fewer births

More rural births are paid by Medicaid

Of low-volume rural hospitals:

41.7% with $<$ births than *financial viability* [8]

29.9% with $<$ births than *clinical safety* [8]

Maternity ward are closing, increasing travel times to delivery.

Where do closures occur?

Across urban and rural hospitals

Rural hospitals in rural states have seen the most closures [9]

Most rural hospitals (52.4%) no longer have maternity wards [9]

What drives closures?

High fixed costs associated with maternity wards.

Aging population \implies fewer births

More rural births are paid by Medicaid

Of low-volume rural hospitals:

41.7% with $<$ births than *financial viability* [8]

29.9% with $<$ births than *clinical safety* [8]

As a recent Health Affairs article concluded:

"Absent substantive changes, obstetric unit closures will keep happening, families will suffer, and no one should be surprised." [10]

Regionalization strategically assigns a level of care to each hospital.

Different levels of care (\mathcal{L}) available at each $h \in \mathcal{H}$

Level 1



uncomplicated pregnancies

Level 2



some specialty care

Level 3+



specialty, subspecialty, NICU

Regionalization strategically assigns a level of care to each hospital.

Different levels of care (\mathcal{L}) available at each $h \in \mathcal{H}$

Level 1



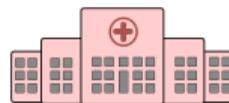
uncomplicated pregnancies

Level 2



some specialty care

Level 3+



specialty, subspecialty, NICU

\mathcal{H} distributed across region, \mathcal{C} seeking care distributed across same region

Regionalization strategically assigns a level of care to each hospital.

Different levels of care (\mathcal{L}) available at each $h \in \mathcal{H}$

Level 1



uncomplicated pregnancies

Level 2



some specialty care

Level 3+



specialty, subspecialty, NICU

\mathcal{H} distributed across region, \mathcal{C} seeking care distributed across same region

Two Stage Modeling Framework:

1st stage: Policy maker assigns levels:

$$x_{hl} \in \{0, 1\} \quad h \in \mathcal{H}, l \in \mathcal{L}$$

2nd stage: Each community $c \in \mathcal{C}$ seeks care based on behavior and assigned levels.

Regionalization strategically assigns a level of care to each hospital.

Different levels of care (\mathcal{L}) available at each $h \in \mathcal{H}$

Level 1



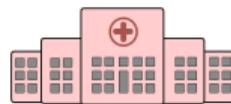
uncomplicated pregnancies

Level 2



some specialty care

Level 3+



specialty, subspecialty, NICU

\mathcal{H} distributed across region, \mathcal{C} seeking care distributed across same region

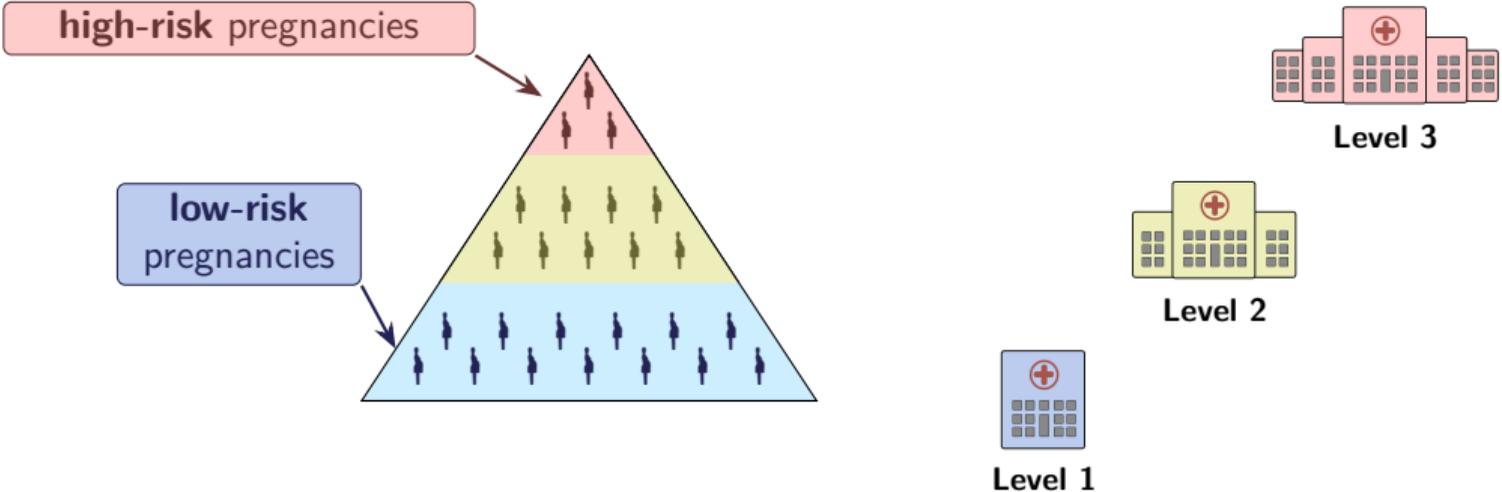
Two Stage Modeling Framework:

1st stage: Policy maker assigns levels:

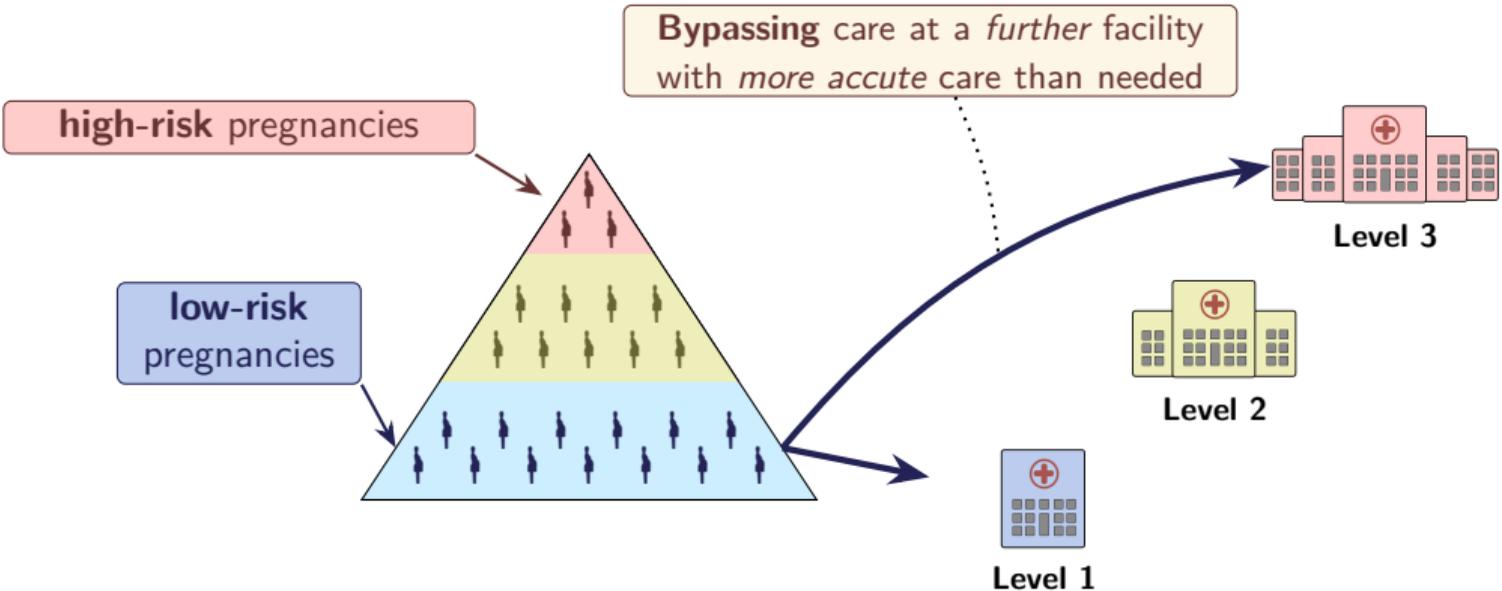
$$x_{hl} \in \{0, 1\} \quad h \in \mathcal{H}, l \in \mathcal{L}$$

2nd stage: Each community $c \in \mathcal{C}$ seeks care based on behavior and assigned levels.

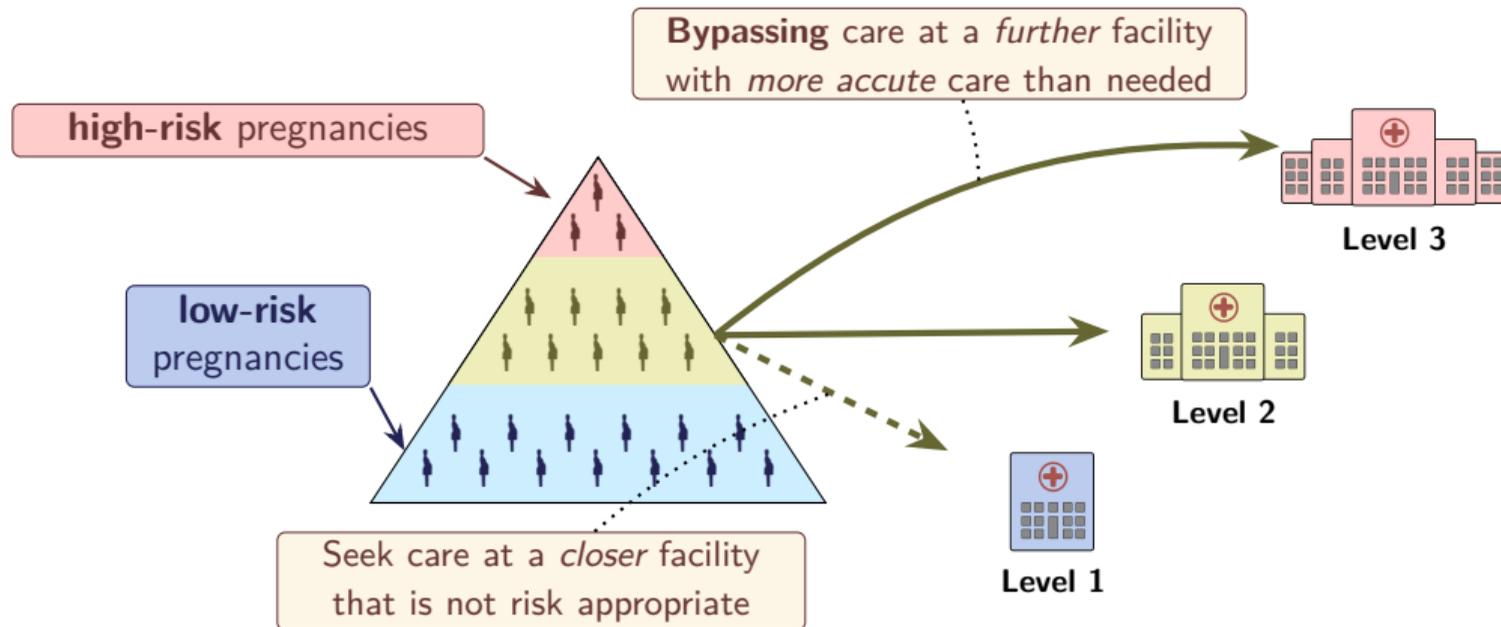
Patient Behavior: Bypassing and Appropriate Referrals



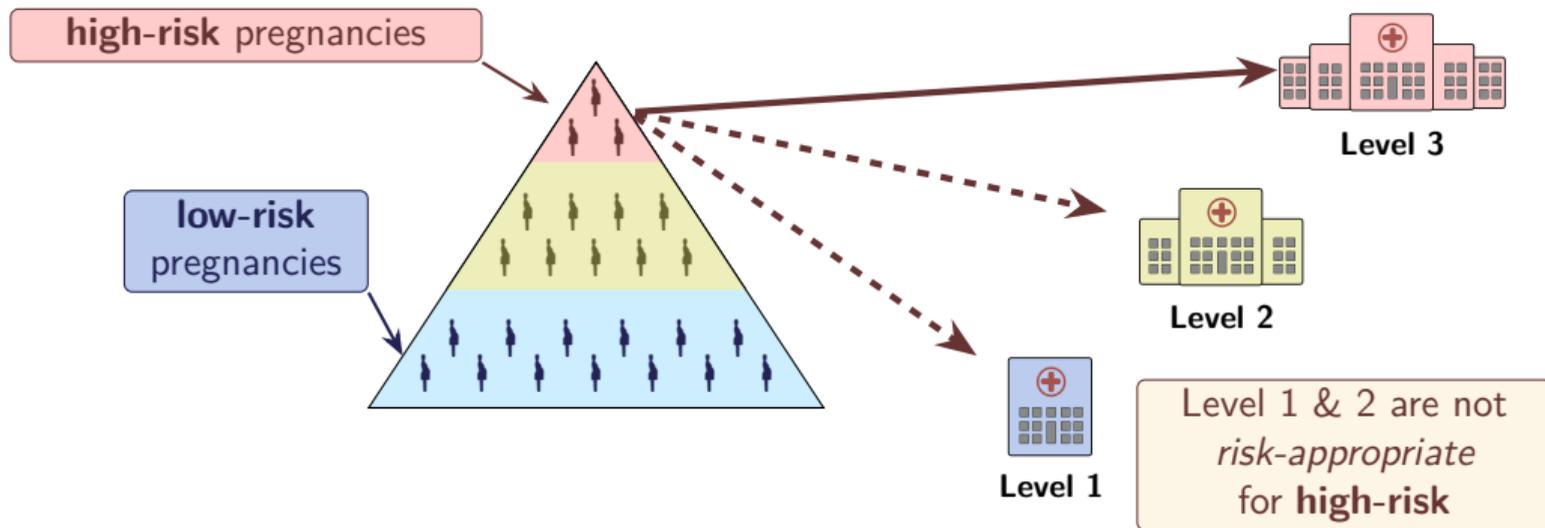
Patient Behavior: Bypassing and Appropriate Referrals



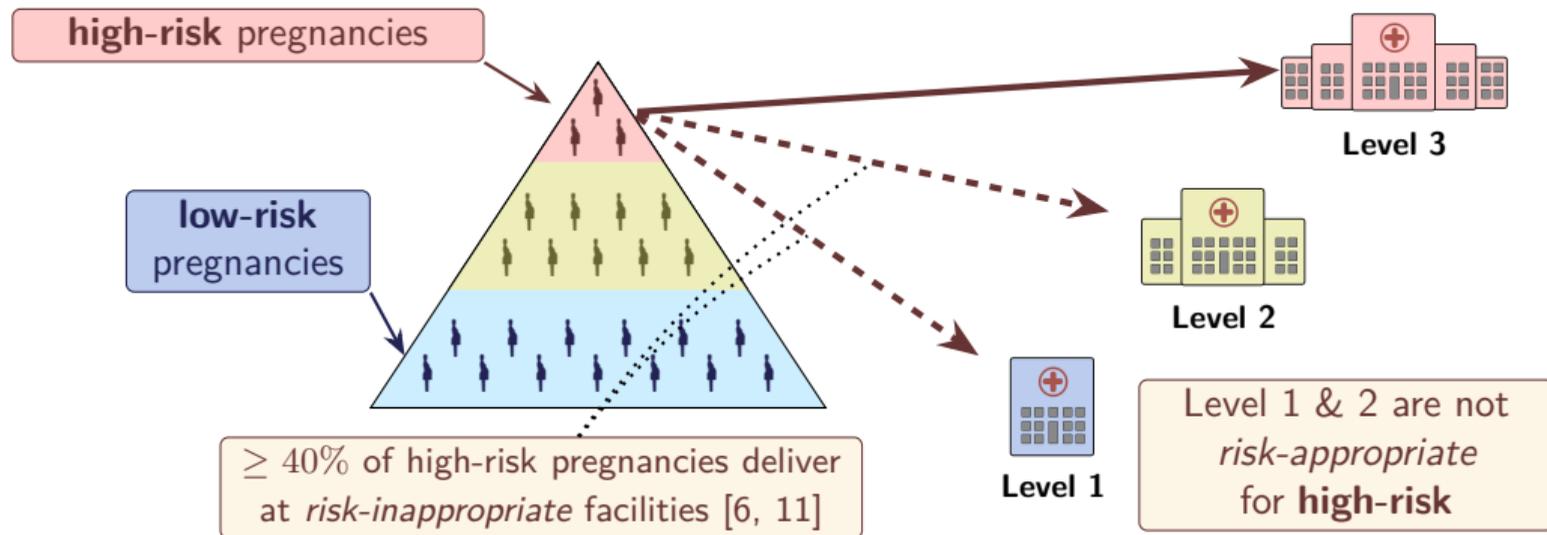
Patient Behavior: Bypassing and Appropriate Referrals



Patient Behavior: Bypassing and Appropriate Referrals



Patient Behavior: Bypassing and Appropriate Referrals



Multi-Objective Model for Systems of Regionalized Maternity Care (MOMS)

Goals: Minimize travel-time and minimize high-risk pregnancies delivered at low-acuity hospitals.

Multi-Objective Model for Systems of Regionalized Maternity Care (MOMS)

Goals: Minimize travel-time and minimize high-risk pregnancies delivered at low-acuity hospitals.

Our Model Minimize expected risk, while being ϵ -close to the travel-time optimal policy

Let \mathcal{P} denote a set of *regionalization* policies, which assign a *level of care* to each hospital

Multi-Objective Model for Systems of Regionalized Maternity Care (MOMS)

Goals: Minimize travel-time and minimize high-risk pregnancies delivered at low-acuity hospitals.

Our Model Minimize expected risk, while being ϵ -close to the travel-time optimal policy

Let \mathcal{P} denote a set of *regionalization* policies, which assign a *level of care* to each hospital

minimize $\mathbb{E}[\text{mismatch between patient needs and received care}]$
 $p_1 \in \mathcal{P}$

subject to

travel-time $_{p_1} \leq (1 + \epsilon) \min_{p_2 \in \mathcal{P}} \text{travel-time}_{p_2}$ for each $c \in \mathcal{C}$,

volume at open hospital \geq sufficient volume,

patient bypassing behavior

Multi-Objective Model for Systems of Regionalized Maternity Care (MOMS)

Goals: Minimize travel-time and minimize high-risk pregnancies delivered at low-acuity hospitals.

Our Model Minimize expected risk, while being ϵ -close to the travel-time optimal policy

Let \mathcal{P} denote a set of *regionalization* policies, which assign a *level of care* to each hospital

minimize $\mathbb{E}[\text{mismatch between patient needs and received care}]$
 $p_1 \in \mathcal{P}$

subject to

travel-time $_{p_1} \leq (1 + \epsilon) \min_{p_2 \in \mathcal{P}} \text{travel-time}_{p_2}$ for each $c \in \mathcal{C}$,

volume at open hospital \geq sufficient volume,

patient bypassing behavior

Rationale: Easily explainable, interpretable for non specialists, quantifiable tradeoff.

Multi-Objective Model for Systems of Regionalized Maternity Care (MOMS)

Goals: Minimize travel-time and minimize high-risk pregnancies delivered at low-acuity hospitals.

Our Model Minimize expected risk, while being as close to the travel-time optimal policy

Let \mathcal{P} denote

Technical Challenges:

- Risk function depends on matching between needs and available (non-convex)
- Not amenable to traditional decomposition, due to incomplete recourse

Rationale: Easily explainable, interpretable for non specialists, quantifiable tradeoff.

Multi-Objective Model for Systems of Regionalized Maternity Care (MOMS)

Goals: Minimize travel-time and minimize high-risk pregnancies delivered at low-acuity hospitals.

Our Model: Minimize expected risk, while being as close to the travel-time optimal policy

Let \mathcal{P} denote

Technical Challenges:

- Risk function depends on matching between needs and available (non-convex)
- Not amenable to traditional decomposition, due to incomplete recourse

Our Solution Approach:

- new *portfolio*-based linearization, with exponentially many variables
 - each portfolio captures the set of hospitals at which a community may seek care
- ⇒ hierarchical branch-price-and-cut to solve state-size instances

Rationale: Easily explainable, interpretable for non specialists, quantifiable tradeoff.

Case Study: Regionalization in Iowa

Why Iowa?

Many low-volume rural hospitals

23 % rural births vs 8% nationwide [4]

State actively undergoing regionalization

Data Sources:

Travel times from OSMnx [1]

Census block groups and data [13]

Need distribution across levels [6]

Estimates of sufficient volume [6]

Bypassing rates estimated as 25%
(10-55%) [2]

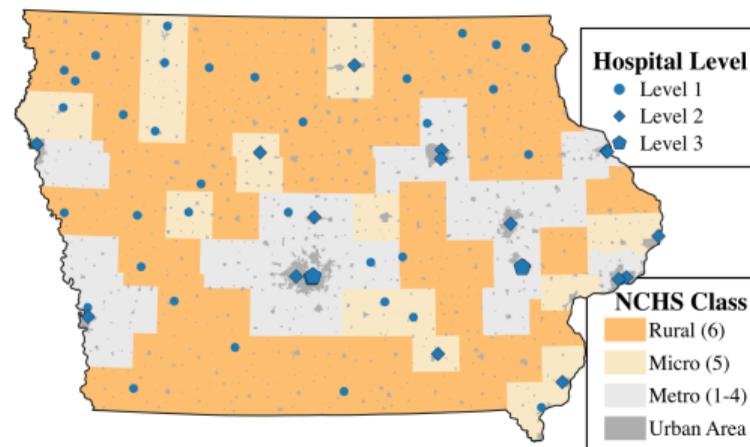


Figure: Iowa has 36 level 1 facilities, 18 level 2 facilities, and 3 level 3 facilities. Data from Iowa Health and Human Services.

Generated policies may improve access while reducing travel times.

Performance Metrics	Status Quo	Generated Policies	
		Time then Risk	Risk then Time
<i>Hospitals by level in Solutions</i>			
Low Volume	15	0	0
Level 1	36	25	0
Level 2	18	11	0
Level 3	3	11	16
<i>Risk-inappropriate care per 1000 births</i>			
Level 2	65.2	59.1	0.0
Level 3	168.4	107.6	0.0
Overall	12.3	10.0	0.0
<i>Distance to care median (IQR) (mins)</i>			
Overall	22.4(11.8-38.2)	14.2(6.0-28.2)	13.1(5.8-38.5)
Mean	25.8	18.3	25.0

Generated policies may improve access while reducing travel times.

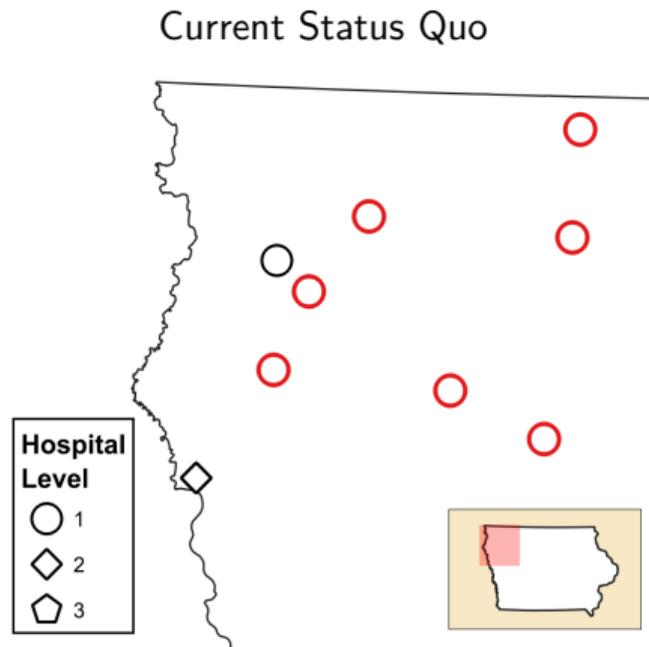
Performance Metrics	Status Quo	Generated Policies	
		Time then Risk	Risk then Time
<i>Hospitals by level in Solutions</i>			
Low Volume	15	0	0
Level 1	36	25	0
Level 2	18	11	0
Level 3	3	11	16
<i>Risk-inappropriate care per 1000 births</i>			
Level 2	65.2	59.1	0.0
Level 3	168.4	107.6	0.0
Overall	12.3	10.0	0.0
<i>Distance to care median (IQR) (mins)</i>			
Overall	22.4(11.8-38.2)	14.2(6.0-28.2)	13.1(5.8-38.5)
Mean	25.8	18.3	25.0

Conclusion: status quo is not on the optimality frontier between travel-time and risk appropriateness.

Observed System Design: Consolidation and Promotion

At time-minimizing policies we observe:

1. Close a subset of **low-volume** hospitals

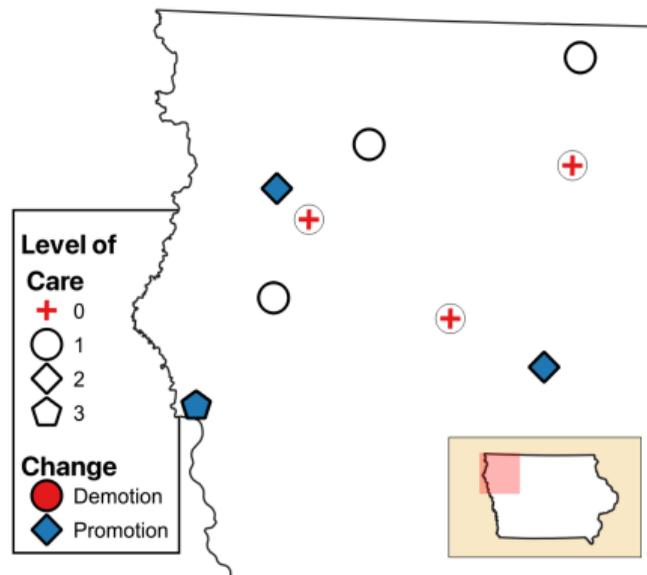


Observed System Design: Consolidation and Promotion

At time-minimizing policies we observe:

1. Close a subset of **low-volume** hospitals
2. Patient volume is **consolidated**
3. **Promote** hospital to higher-acuity

Relative to Status Quo

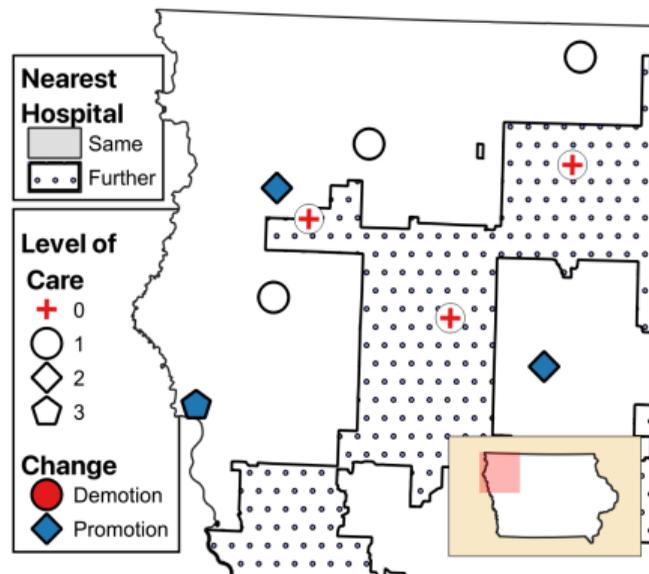


Observed System Design: Consolidation and Promotion

At time-minimizing policies we observe:

1. Close a subset of **low-volume** hospitals
2. Patient volume is **consolidated**
3. **Promote** hospital to higher-acuity

Relative to Status Quo

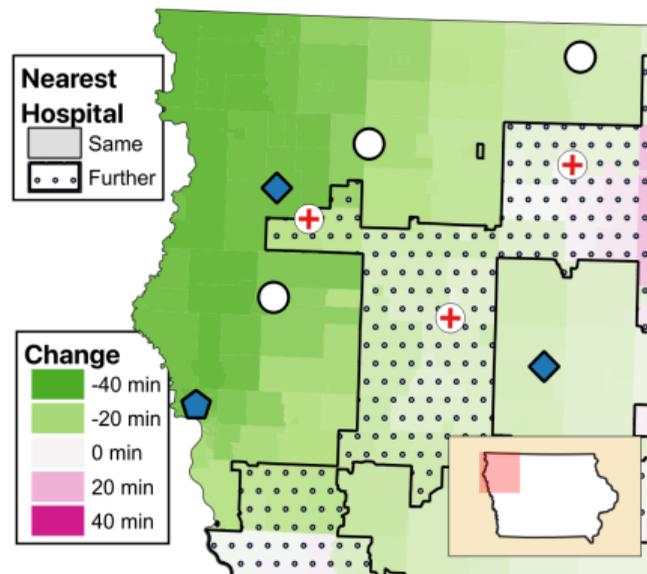


Observed System Design: Consolidation and Promotion

At time-minimizing policies we observe:

1. Close a subset of **low-volume** hospitals
2. Patient volume is **consolidated**
3. **Promote** hospital to higher-acuity

Relative to Status Quo



Among patient population needing level 1 care

Observed System Design: Consolidation and Promotion

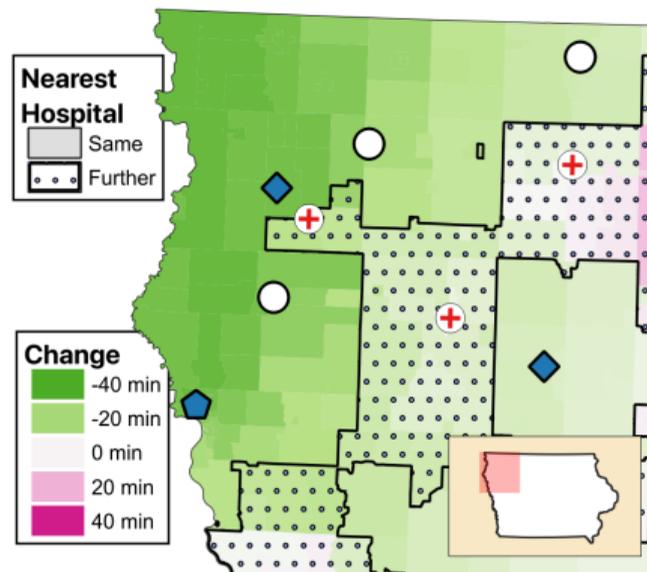
At time-minimizing policies we observe:

1. Close a subset of **low-volume** hospitals
2. Patient volume is **consolidated**
3. **Promote** hospital to higher-acuity

Average travel time decreases

- even among patients who need level 1 care
- driven by closer care for bypassing patients

Relative to Status Quo



Among patient population needing level 1 care

Observed System Design: Consolidation and Promotion

At time-minimizing policies we observe:

1. Close a subset of **low-volume** hospitals
2. Patient volume is **consolidated**
3. **Promote** hospital to higher-acuity

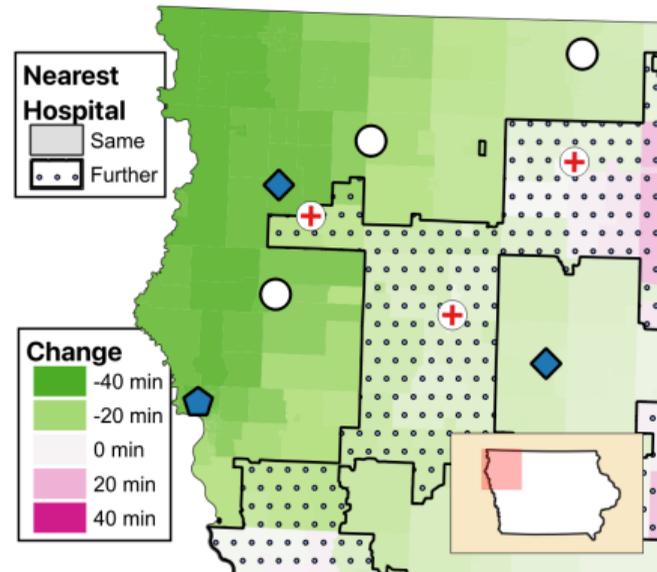
Average travel time decreases

- even among patients who need level 1 care
- driven by closer care for bypassing patients

Risk-appropriateness increases:

- closest hospital may provide higher-acuity care
- benefit for patients not appropriate referred

Relative to Status Quo



Among patient population needing level 1 care

In summary

- Context driven model of hierarchical regionalization of maternity care wards
- *New linearization with portfolios of care-seeking behavior, branch-price-and-cut*
- Identify that current *status-quo* in Iowa is not on the travel-time - risk-appropriateness frontier
- Observed consolidate-and-promote in distance optimal policies for Iowa

In summary

- Context driven model of hierarchical regionalization of maternity care wards
- *New linearization with portfolios of care-seeking behavior, branch-price-and-cut*
- Identify that current *status-quo* in Iowa is not on the travel-time - risk-appropriateness frontier
- Observed consolidate-and-promote in distance optimal policies for Iowa

Current / Future Work

- Identify parameter regions under-which consolidation-and-promotion occurs
- Integrate the impact of available antenatal care into referral rates
- Consider spatial-heterogeneity in risk and payment models

In summary

- Context driven model of hierarchical regionalization of maternity care wards
- *New linearization with portfolios of care-seeking behavior, branch-price-and-cut*
- Identify that current *status-quo* in Iowa is not on the travel-time - risk-appropriateness frontier
- Observed consolidate-and-promote in distance optimal policies for Iowa

Current / Future Work

- Identify parameter regions under-which consolidation-and-promotion occurs
- Integrate the impact of available antenatal care into referral rates
- Consider spatial-heterogeneity in risk and payment models

Thank you!

Questions: asapirstein3@gatech.edu

References I

- [1] Geoff Boeing. Modeling and analyzing urban networks and amenities with OSMnx. *Geographical Analysis*, 57(4):567–577, 2025.
- [2] Margaret Carrel, Barbara C. Keino, Nicole L. Novak, Kelli K. Ryckman, and Stephanie Radke. Bypassing of nearest labor & delivery unit is contingent on rurality, wealth, and race. *Birth*, 50(1):5–10, 2023.
- [3] CDC. Severe maternal morbidity, May 2024.
- [4] Centers for Disease Control and Prevention, National Center for Health Statistics. National Vital Statistics System, Natality on CDC WONDER Online Database. Data are from the Natality Records 2016-2022, as compiled from data provided by the 57 vital statistics jurisdictions through the Vital Statistics Cooperative Program, 2022. Accessed on Feb 13, 2024 at 10:56:11 AM.
- [5] Yingxi Chen, Meredith S. Shiels, Tarsicio Uribe-Leitz, Rose L. Molina, Wayne R. Lawrence, Neal D. Freedman, and Christian C. Abnet. Pregnancy-related deaths in the us, 2018-2022. *JAMA Network Open*, 8(4):e254325, April 2025.

References II

- [6] Sarah Rae Easter, Julian N. Robinson, M. Kathryn Menard, Andreea A. Creanga, Xinling Xu, Sarah E. Little, and Brian T. Bateman. Potential effects of regionalized maternity care on u.s. hospitals. *Obstetrics & Gynecology*, 134(3):545–552, sept 2019.
- [7] Dorothy A Fink, Deborah Kilday, Zhun Cao, Kelly Larson, Adrienne Smith, Craig Lipkin, Raymond Perigard, Richelle Marshall, Taryn Deirmenjian, Ashley Finke, et al. Trends in maternal mortality and severe maternal morbidity during delivery-related hospitalizations in the united states, 2008 to 2021. *JAMA Network Open*, 6(6):e2317641–e2317641, 2023.
- [8] Katy B Kozhimannil, Julia D Interrante, Lindsay K Admon, and Bridget L Basile Ibrahim. Rural hospital administrators' beliefs about safety, financial viability, and community need for offering obstetric care. In *JAMA Health Forum*, volume 3, pages e220204–e220204. American Medical Association, 2022.
- [9] Katy B. Kozhimannil, Julia D. Interrante, Caitlin Carroll, Emily C. Sheffield, Alyssa H. Fritz, Alecia J. McGregor, and Sara C. Handley. Obstetric care access at rural and urban hospitals in the united states. *JAMA*, December 2024.
- [10] Katy Backes Kozhimannil, Emily C. Sheffield, and Julia D. Interrante. Millions of women don't have access to maternity care - and the number is growing. *Health Affairs Forefront*, 2026.

References III

- [11] Godwin K Osei-Poku, Julia C Prentice, Sarah Rae Easter, and Hafsatou Diop. Delivery at an inadequate level of maternal care is associated with severe maternal morbidity. *American Journal of Obstetrics and Gynecology*, 231(5):546–e1, 2024.
- [12] Susanna Trost, Jennifer Beauregard, Gyan Chandra, Fanny Njie, Jasmine Berry, Alyssa Harvey, and David A Goodman. Pregnancy-related deaths: data from maternal mortality review committees in 36 US states, 2017–2019. *Education*, 45(10):1–0, 2022.
- [13] U.S. Census Bureau. Centers of Population. Retrieved from <https://www.census.gov/geographies/reference-files/time-series/geo/centers-population.html>, 2021.

Efficient Inference Using Large Language Models with Limited Human Data: Fine-Tuning then Rectification

Lei Wang (UW), Zikun Ye (UW), **Jinglong Zhao (BU)**

IMSI

Feb 4, 2026

Assessing Quality of Care in Telehealth

- ▶ Nurse auditors ensuring the care meets quality standards

Assessing Quality of Care in Telehealth

- ▶ Nurse auditors ensuring the care meets quality standards
- ▶ Abundant digital footprints **X**

Patient Intake Log

Timestamp: 08:45 AM

Patient: AAA

I've had a throbbing headache behind my left eye for 3 days. It feels worse when I look at my phone. I took some Tylenol, but it didn't really help. I'm also feeling a bit nauseous and the light in my kitchen is making it worse.

Doctor's Telehealth Notes

Timestamp: 03:53 PM

Provider: Dr. BBB

Duration: 12 minutes

Patient presents with unilateral, throbbing cephalalgia with associated photophobia and nausea. No history of aura reported. Onset was gradual. Neurological screening (self-reported via camera) shows no focal deficits. Diagnosis: Acute migraine without aura. Advised darkened room and hydration.

Follow-up Instructions

Timestamp: 04:12 PM

Patient instructed to rest in a quiet, dark room and continue monitoring. Advised to seek emergency care for visual changes described as a "curtain dropping," sudden weakness, fever, or onset of the "worst headache of life." Prescription for Sumatriptan transmitted to pharmacy.

Assessing Quality of Care in Telehealth

- ▶ Nurse auditors ensuring the care meets quality standards
- ▶ Abundant digital footprints **X**
- ▶ Costly nurse auditor evaluations **Y**

Patient Intake Log

Timestamp: 08:45 AM

Patient: AAA

I've had a throbbing headache behind my left eye for 3 days. It feels worse when I look at my phone. I took some Tylenol, but it didn't really help. I'm also feeling a bit nauseous and the light in my kitchen is making it worse.

Doctor's Telehealth Notes

Timestamp: 03:53 PM

Provider: Dr. BBB

Duration: 12 minutes

Patient presents with unilateral, throbbing cephalalgia with associated photophobia and nausea. No history of aura reported. Onset was gradual. Neurological screening (self-reported via camera) shows no focal deficits. Diagnosis: Acute migraine without aura. Advised darkened room and hydration.

Follow-up Instructions

Timestamp: 04:12 PM

Patient instructed to rest in a quiet, dark room and continue monitoring. Advised to seek emergency care for visual changes described as a "curtain dropping," sudden weakness, fever, or onset of the "worst headache of life." Prescription for Sumatriptan transmitted to pharmacy.

Assessing Quality of Care in Telehealth

- ▶ Nurse auditors ensuring the care meets quality standards
- ▶ Abundant digital footprints **X**
- ▶ Costly nurse auditor evaluations **Y**

Patient Intake Log

Timestamp: 08:45 AM

Patient: AAA

I've had a throbbing headache behind my left eye for 3 days. It feels worse when I look at my phone. I took some Tylenol, but it didn't really help. I'm also feeling a bit nauseous and the light in my kitchen is making it worse.

Doctor's Telehealth Notes

Timestamp: 03:53 PM

Provider: Dr. BBB

Duration: 12 minutes

Patient presents with unilateral, throbbing cephalalgia with associated photophobia and nausea. No history of aura reported. Onset was gradual. Neurological screening (self-reported via camera) shows no focal deficits. Diagnosis: Acute migraine without aura. Advised darkened room and hydration.

Follow-up Instructions

Timestamp: 04:12 PM

Patient instructed to rest in a quiet, dark room and continue monitoring. Advised to seek emergency care for visual changes described as a "curtain dropping," sudden weakness, fever, or onset of the "worst headache of life." Prescription for Sumatriptan transmitted to pharmacy.

Nurse Auditor Evaluation: 85/100

Assessing Quality of Care in Telehealth

- ▶ Abundant digital footprints \mathbf{X}
- ▶ Nurse auditor evaluations Y
- ▶ LLM evaluations $f(\mathbf{X})$

Patient Intake Log

Timestamp: 08:45 AM

Patient: AAA

I've had a throbbing headache behind my left eye for 3 days. It feels worse when I look at my phone. I took some Tylenol, but it didn't really help. I'm also feeling a bit nauseous and the light in my kitchen is making it worse.

Doctor's Telehealth Notes

Timestamp: 03:53 PM

Provider: Dr. BBB

Duration: 12 minutes

Patient presents with unilateral, throbbing cephalalgia with associated photophobia and nausea. No history of aura reported. Onset was gradual. Neurological screening (self-reported via camera) shows no focal deficits. Diagnosis: Acute migraine without aura. Advised darkened room and hydration.

Follow-up Instructions

Timestamp: 04:12 PM

Patient instructed to rest in a quiet, dark room and continue monitoring. Advised to seek emergency care for visual changes described as a "curtain dropping," sudden weakness, fever, or onset of the "worst headache of life." Prescription for Sumatriptan transmitted to pharmacy.

Nurse Auditor Evaluation: 85/100

LLM Evaluation: 90/100

Assessing Quality of Care in Telehealth

- ▶ Abundant digital footprints \mathbf{X}
- ▶ Nurse auditor evaluations Y : high-fidelity but costly
- ▶ LLM evaluations $f(\mathbf{X})$: low-fidelity but near-zero cost

Patient Intake Log

Timestamp: 08:45 AM

Patient: AAA

I've had a throbbing headache behind my left eye for 3 days. It feels worse when I look at my phone. I took some Tylenol, but it didn't really help. I'm also feeling a bit nauseous and the light in my kitchen is making it worse.

Doctor's Telehealth Notes

Timestamp: 03:53 PM

Provider: Dr. BBB

Duration: 12 minutes

Patient presents with unilateral, throbbing cephalalgia with associated photophobia and nausea. No history of aura reported. Onset was gradual. Neurological screening (self-reported via camera) shows no focal deficits. Diagnosis: Acute migraine without aura. Advised darkened room and hydration.

Follow-up Instructions

Timestamp: 04:12 PM

Patient instructed to rest in a quiet, dark room and continue monitoring. Advised to seek emergency care for visual changes described as a "curtain dropping," sudden weakness, fever, or onset of the "worst headache of life." Prescription for Sumatriptan transmitted to pharmacy.

Nurse Auditor Evaluation: 85/100

LLM Evaluation: 90/100

Using LLMs as Human Surrogates

Using LLMs as Human Surrogates

- ▶ What do we have:
 - ▶ (1) small labeled data (2) large unlabeled data (3) LLM predictor

Using LLMs as Human Surrogates

- ▶ What do we have:
 - ▶ (1) small labeled data (2) large unlabeled data (3) LLM predictor
- ▶ Two existing approaches: fine-tuning and rectification

Using LLMs as Human Surrogates

- ▶ What do we have:
 - ▶ (1) small labeled data (2) large unlabeled data (3) LLM predictor
- ▶ Two existing approaches: fine-tuning and rectification
- ▶ Fine-tuning: use small labeled data to improve LLM

Using LLMs as Human Surrogates

- ▶ What do we have:
 - ▶ (1) small labeled data (2) large unlabeled data (3) LLM predictor
- ▶ Two existing approaches: fine-tuning and rectification
- ▶ Fine-tuning: use small labeled data to improve LLM
 - ▶ Marginal value of one additional labeled data point decreases

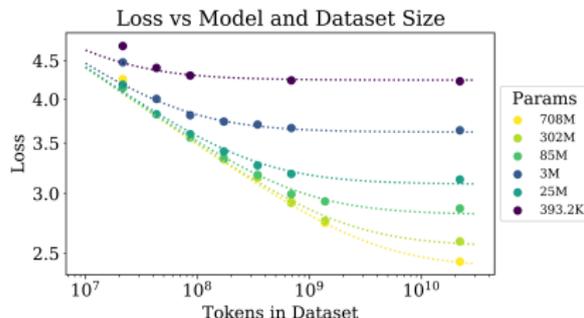


Figure 4 in [Kaplan et al. '20]

Using LLMs as Human Surrogates

- ▶ What do we have:
 - ▶ (1) small labeled data (2) large unlabeled data (3) LLM predictor
- ▶ Two existing approaches: fine-tuning and rectification
- ▶ Fine-tuning: use small labeled data to improve LLM
 - ▶ Marginal value of one additional labeled data point decreases
- ▶ Rectification: use small labeled data to quantify errors of LLM

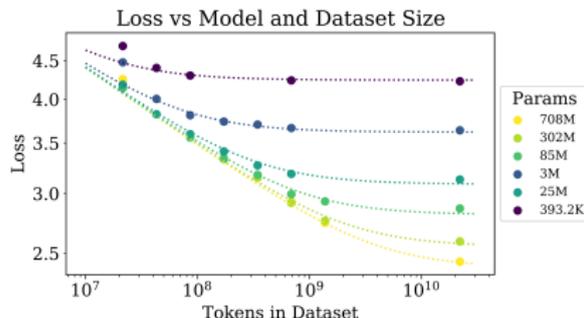


Figure 4 in [Kaplan et al. '20]

Using LLMs as Human Surrogates

- ▶ What do we have:
 - ▶ (1) small labeled data (2) large unlabeled data (3) LLM predictor
- ▶ Two existing approaches: fine-tuning and rectification
- ▶ Fine-tuning: use small labeled data to improve LLM
 - ▶ Marginal value of one additional labeled data point decreases
- ▶ Rectification: use small labeled data to quantify errors of LLM
 - ▶ Marginal value of one additional labeled data point decreases

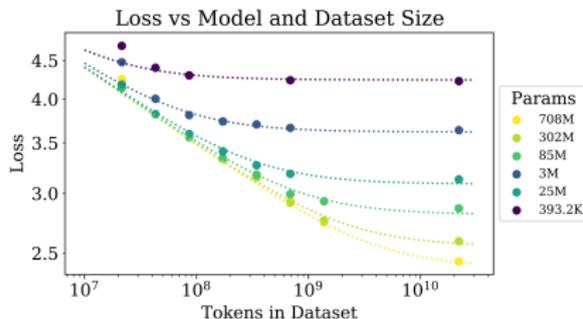


Figure 4 in [Kaplan et al. '20]

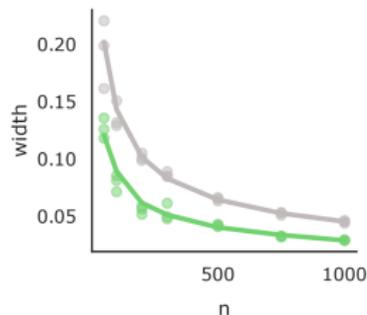


Figure 1 in [Angelopoulos et al. '23]

Using LLMs as Human Surrogates

- ▶ What do we have:
 - ▶ (1) small labeled data (2) large unlabeled data (3) LLM predictor
- ▶ Two existing approaches: fine-tuning and rectification
- ▶ Fine-tuning: use small labeled data to improve LLM
 - ▶ Marginal value of one additional labeled data point decreases
- ▶ Rectification: use small labeled data to quantify errors of LLM
 - ▶ Marginal value of one additional labeled data point decreases
- ▶ This paper: “fine-tuning then rectification”
 - ▶ Taking advantage of the rapid initial gains of both methods

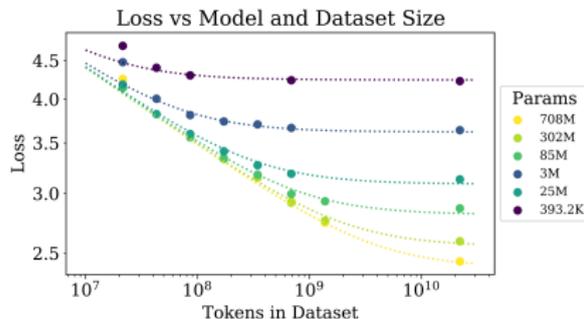


Figure 4 in [Kaplan et al. '20]

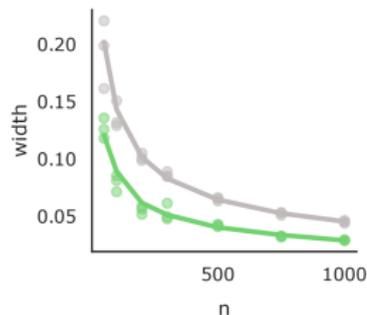


Figure 1 in [Angelopoulos et al. '23]

Fine-Tuning then Rectification

- ▶ Partition the labeled samples into two subsets:
the first used to fine-tune LLM, the second used for rectification

Fine-Tuning then Rectification

- ▶ Partition the labeled samples into two subsets:
the first used to fine-tune LLM, the second used for rectification
- ▶ One rectification method: prediction-powered inference (PPI)

Fine-Tuning then Rectification

- ▶ Partition the labeled samples into two subsets:
the first used to fine-tune LLM, the second used for rectification
- ▶ One rectification method: prediction-powered inference (PPI)

Example (Rectification)

Sample	Nurse Auditor Evaluation	LLM Evaluation
1	85	90
2	80	85
3	83	88
4	89	94
5	?	99

Fine-Tuning then Rectification

- ▶ Partition the labeled samples into two subsets:
the first used to fine-tune LLM, the second used for rectification
- ▶ One rectification method: prediction-powered inference (PPI)

Example (Rectification)

Sample	Nurse Auditor Evaluation	LLM Evaluation
1	85	90
2	80	85
3	83	88
4	89	94
5	94	99

Fine-Tuning then Rectification

- ▶ Partition the labeled samples into two subsets:
the first used to fine-tune LLM, the second used for rectification
- ▶ One rectification method: prediction-powered inference (PPI)

Example (Rectification)

Sample	Nurse Auditor Evaluation	LLM Evaluation
1	85	90
2	80	85
3	83	88
4	89	94
5	94	99

Example (Different LLMs)

Fine-Tuning then Rectification

- ▶ Partition the labeled samples into two subsets:
the first used to fine-tune LLM, the second used for rectification
- ▶ One rectification method: prediction-powered inference (PPI)

Example (Rectification)

Sample	Nurse Auditor Evaluation	LLM Evaluation
1	85	90
2	80	85
3	83	88
4	89	94
5	94	99

Example (Different LLMs)

- ▶ LLM A: always over-predict by 5

Fine-Tuning then Rectification

- ▶ Partition the labeled samples into two subsets:
the first used to fine-tune LLM, the second used for rectification
- ▶ One rectification method: prediction-powered inference (PPI)

Example (Rectification)

Sample	Nurse Auditor Evaluation	LLM Evaluation
1	85	90
2	80	85
3	83	88
4	89	94
5	94	99

Example (Different LLMs)

- ▶ LLM A: always over-predict by 5
- ▶ LLM B: sometimes over-predict by 1, sometimes under-predict by 1

Fine-Tuning then Rectification

- ▶ Partition the labeled samples into two subsets: the first used to fine-tune LLM, the second used for rectification
- ▶ One rectification method: prediction-powered inference (PPI)
- ▶ Change fine-tuning objective: minimize variance instead of MSE

Example (Rectification)

Sample	Nurse Auditor Evaluation	LLM Evaluation
1	85	90
2	80	85
3	83	88
4	89	94
5	94	99

Example (Fine-Tuning a Biased LLM)

- ▶ LLM A: always over-predict by 5
- ▶ LLM B: sometimes over-predict by 1, sometimes under-predict by 1

This Talk...

- ▶ A fine-tuning then rectification framework
- ▶ Optimal sample allocation between fine-tuning and rectification
- ▶ Data-driven procedure to guide sample allocation
- ▶ Empirical demonstration using Wine Reviews dataset

The Fine-Tuning then Rectification Framework

- ▶ $\{(\mathbf{X}_i, Y_i)\}_{i=1}^n$: small labeled data sampled from \mathcal{F}

The Fine-Tuning then Rectification Framework

- ▶ $\{(\mathbf{X}_i, Y_i)\}_{i=1}^n$: small labeled data sampled from \mathcal{F}
 - ▶ Digital footprints with nurse auditor evaluations

The Fine-Tuning then Rectification Framework

- ▶ $\{(\mathbf{X}_i, Y_i)\}_{i=1}^n$: small labeled data sampled from \mathcal{F}
 - ▶ Digital footprints with nurse auditor evaluations
- ▶ $\{\tilde{\mathbf{X}}_j\}_{j=1}^m$: large unlabeled data sampled from \mathcal{F}_X

The Fine-Tuning then Rectification Framework

- ▶ $\{(\mathbf{X}_i, Y_i)\}_{i=1}^n$: small labeled data sampled from \mathcal{F}
 - ▶ Digital footprints with nurse auditor evaluations
- ▶ $\{\tilde{\mathbf{X}}_j\}_{j=1}^m$: large unlabeled data sampled from \mathcal{F}_X
 - ▶ Digital footprints only

The Fine-Tuning then Rectification Framework

- ▶ $\{(\mathbf{X}_i, Y_i)\}_{i=1}^n$: small labeled data sampled from \mathcal{F}
 - ▶ Digital footprints with nurse auditor evaluations
- ▶ $\{\tilde{\mathbf{X}}_j\}_{j=1}^m$: large unlabeled data sampled from \mathcal{F}_X
 - ▶ Digital footprints only
- ▶ $f(\cdot)$: LLM predictor, independent of both data

The Fine-Tuning then Rectification Framework

- ▶ $\{(\mathbf{X}_i, Y_i)\}_{i=1}^n$: small labeled data sampled from \mathcal{F}
 - ▶ Digital footprints with nurse auditor evaluations
- ▶ $\{\tilde{\mathbf{X}}_j\}_{j=1}^m$: large unlabeled data sampled from \mathcal{F}_X
 - ▶ Digital footprints only
- ▶ $f(\cdot)$: LLM predictor, independent of both data
 - ▶ Replace the classification head by a regression head

The Fine-Tuning then Rectification Framework

- ▶ $\{(\mathbf{X}_i, Y_i)\}_{i=1}^n$: small labeled data sampled from \mathcal{F}
 - ▶ Digital footprints with nurse auditor evaluations
- ▶ $\{\tilde{\mathbf{X}}_j\}_{j=1}^m$: large unlabeled data sampled from \mathcal{F}_X
 - ▶ Digital footprints only
- ▶ $f(\cdot)$: LLM predictor, independent of both data
 - ▶ Replace the classification head by a regression head
- ▶ $f^{(s)}(\cdot)$: LLM predictor fine-tuned on s samples of the small labeled data

Fine-Tuning then Rectification: Mean Estimation

- ▶ $\{(\mathbf{X}_i, Y_i)\}_{i=1}^n$: small labeled data sampled from \mathcal{F}
 - ▶ Digital footprints with nurse auditor evaluations
- ▶ $\{\tilde{\mathbf{X}}_j\}_{j=1}^m$: large unlabeled data sampled from \mathcal{F}_X
 - ▶ Digital footprints only
- ▶ $f(\cdot)$: LLM predictor, independent of both data
 - ▶ Replace the classification head by a regression head
- ▶ $f^{(s)}(\cdot)$: LLM predictor fine-tuned on s samples of the small labeled data
- ▶ Estimator based on fine-tuned predictor $f^{(s)}(\cdot)$

$$\hat{\mu} = \frac{1}{n-s} \sum_{i=1}^{n-s} (Y_i - f^{(s)}(\mathbf{X}_i)) + \frac{1}{m} \sum_{j=1}^m f^{(s)}(\tilde{\mathbf{X}}_j)$$

Variance of Estimator

- ▶ Estimator based on fine-tuned predictor $f^{(s)}(\cdot)$

$$\hat{\mu} = \frac{1}{n-s} \sum_{i=1}^{n-s} (Y_i - f^{(s)}(\mathbf{x}_i)) + \frac{1}{m} \sum_{j=1}^m f^{(s)}(\tilde{\mathbf{x}}_j)$$

Variance of Estimator

- ▶ Estimator based on fine-tuned predictor $f^{(s)}(\cdot)$

$$\hat{\mu} = \frac{1}{n-s} \sum_{i=1}^{n-s} (Y_i - f^{(s)}(\mathbf{x}_i)) + \frac{1}{m} \sum_{j=1}^m f^{(s)}(\tilde{\mathbf{x}}_j)$$

- ▶ Variance of $\hat{\mu}$

$$\text{Var}(\hat{\mu}) = \frac{1}{n-s} \text{Var}\left(Y - f^{(s)}(\mathbf{X})\right) + \frac{1}{m} \text{Var}\left(f^{(s)}(\mathbf{X})\right)$$

Variance of Estimator

- ▶ Estimator based on fine-tuned predictor $f^{(s)}(\cdot)$

$$\hat{\mu} = \frac{1}{n-s} \sum_{i=1}^{n-s} (Y_i - f^{(s)}(\mathbf{X}_i)) + \frac{1}{m} \sum_{j=1}^m f^{(s)}(\tilde{\mathbf{X}}_j)$$

- ▶ When m , sample size of the large unlabeled data, is large

$$\begin{aligned} \text{Var}(\hat{\mu}) &= \frac{1}{n-s} \text{Var}\left(Y - f^{(s)}(\mathbf{X})\right) + \frac{1}{m} \text{Var}\left(f^{(s)}(\mathbf{X})\right) \\ &\approx \frac{1}{n-s} \text{Var}\left(Y - f^{(s)}(\mathbf{X})\right) \end{aligned}$$

Variance of Estimator

- ▶ Estimator based on fine-tuned predictor $f^{(s)}(\cdot)$

$$\hat{\mu} = \frac{1}{n-s} \sum_{i=1}^{n-s} (Y_i - f^{(s)}(\mathbf{X}_i)) + \frac{1}{m} \sum_{j=1}^m f^{(s)}(\tilde{\mathbf{X}}_j)$$

- ▶ When m , sample size of the large unlabeled data, is large

$$\begin{aligned} \text{Var}(\hat{\mu}) &= \frac{1}{n-s} \text{Var}\left(Y - f^{(s)}(\mathbf{X})\right) + \frac{1}{m} \text{Var}\left(f^{(s)}(\mathbf{X})\right) \\ &\approx \frac{1}{n-s} \text{Var}\left(Y - f^{(s)}(\mathbf{X})\right) \end{aligned}$$

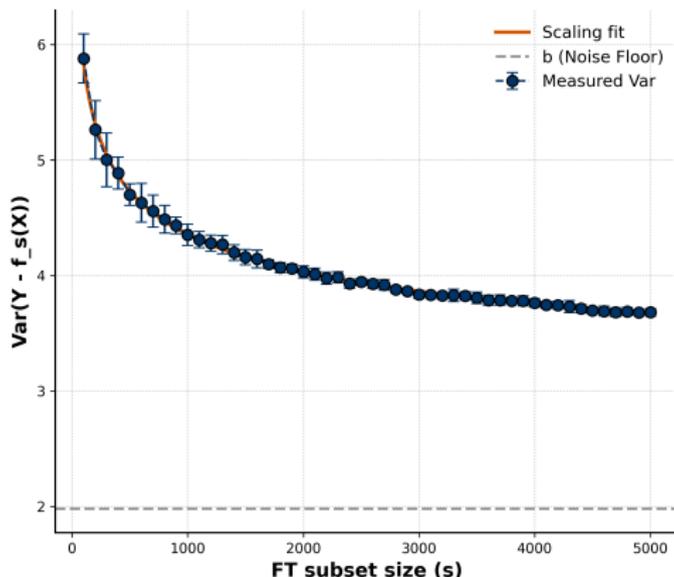
- ▶ Use $\text{Var}(Y - f^{(s)}(\mathbf{X}))$ as fine-tuning objective

Optimal Sample Allocation

► Scaling law

$$\text{Var}\left(Y - f^{(s)}(\mathbf{X})\right) = as^{-\alpha} + b,$$

where $\alpha, a > 0$, $b \geq 0$ are task-specific constants



Optimal Sample Allocation

- ▶ Scaling law

$$\text{Var}\left(Y - f^{(s)}(\mathbf{X})\right) = as^{-\alpha} + b,$$

where $\alpha, a > 0$, $b \geq 0$ are task-specific constants

- ▶ Optimal sample allocation

$$\min_{s \in (0, n)} \frac{\text{Var}\left(Y - f^{(s)}(\mathbf{X})\right)}{n - s} = \min_{s \in (0, n)} \frac{as^{-\alpha} + b}{n - s}$$

Optimal Sample Allocation

- ▶ Scaling law

$$\text{Var}\left(Y - f^{(s)}(\mathbf{X})\right) = as^{-\alpha} + b,$$

where $\alpha, a > 0$, $b \geq 0$ are task-specific constants

- ▶ Optimal sample allocation

$$\min_{s \in (0, n)} \frac{\text{Var}(Y - f^{(s)}(\mathbf{X}))}{n - s} = \min_{s \in (0, n)} \frac{as^{-\alpha} + b}{n - s}$$

Theorem

The optimal sample allocation s^ is given as the unique solution to*

$$\alpha ans^{-\alpha-1} - (\alpha + 1)as^{-\alpha} - b = 0$$

Optimal Sample Allocation

- ▶ Scaling law

$$\text{Var}\left(Y - f^{(s)}(\mathbf{X})\right) = as^{-\alpha} + b,$$

where $\alpha, a > 0$, $b \geq 0$ are task-specific constants

- ▶ Optimal sample allocation

$$\min_{s \in (0, n)} \frac{\text{Var}(Y - f^{(s)}(\mathbf{X}))}{n - s} = \min_{s \in (0, n)} \frac{as^{-\alpha} + b}{n - s}$$

- ▶ Estimate $\hat{\alpha}, \hat{a}, \hat{b}$ during fine-tuning

Theorem

The optimal sample allocation s^ is given as the unique solution to*

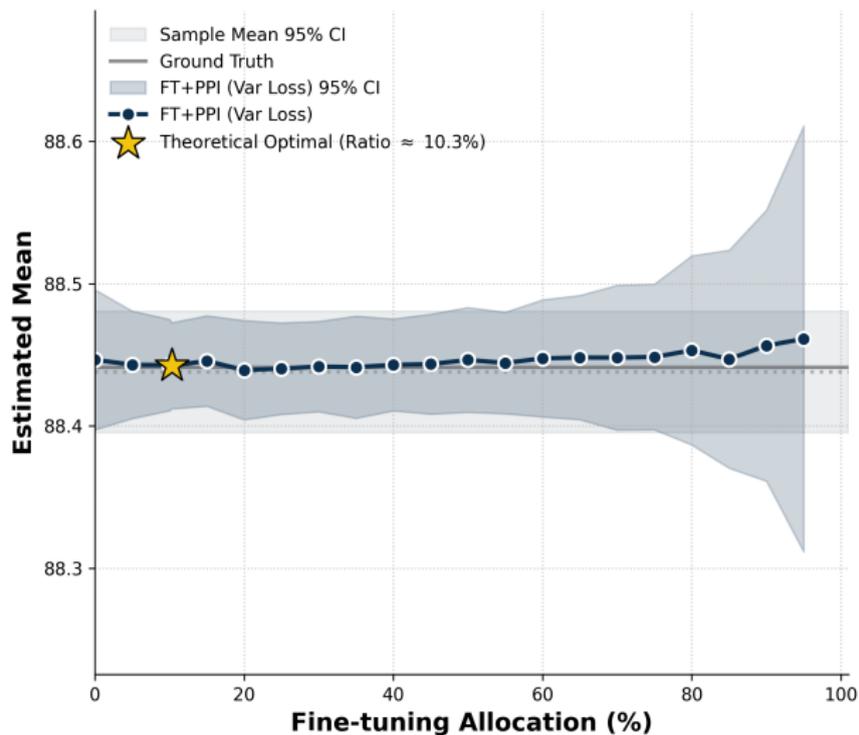
$$\alpha ans^{-\alpha-1} - (\alpha + 1)as^{-\alpha} - b = 0$$

Empirical Demonstration Using Wine Reviews Data

Table: Examples from the Wine Reviews dataset

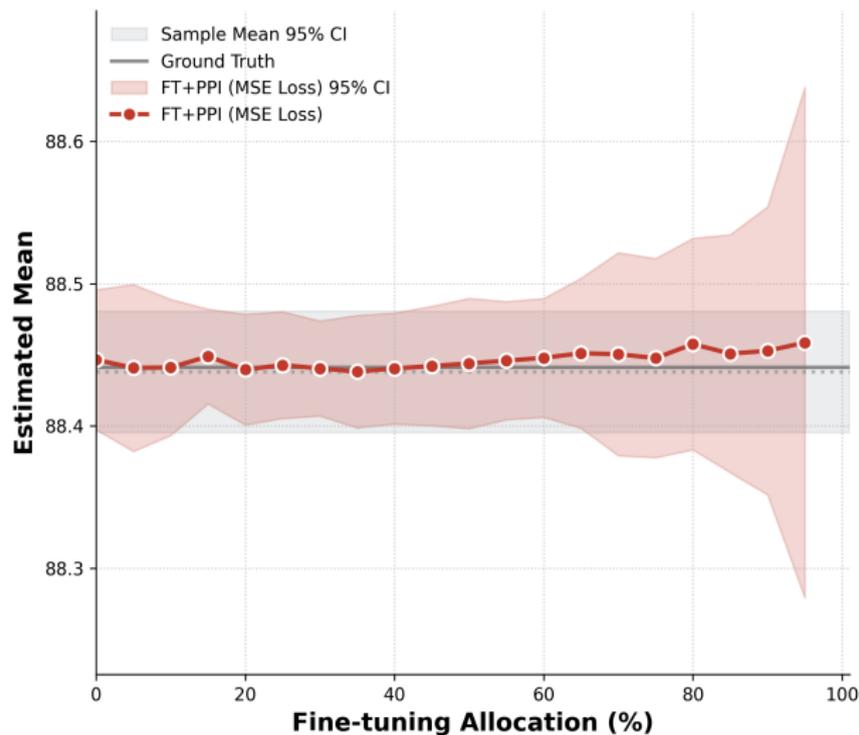
Review (X)	This is a walk backward after the impressive 2012. Almost impenetrably black, the flavors converge around espresso and bitter chocolate, yet the tannins have a green edge. The wine simply feels flat in the mouth with no life to it.
Rating (Y)	85
Review (X)	This is one of the great Rieslings from the Wachau, a wonderful panoply of ripe, tropical fruit, pierced with flint, spice and minerality. It is rich and opulent, while never losing sight of the core tautness of a fine Riesling.
Rating (Y)	95

Variance-Based vs. MSE-Based Fine-Tuning



Variance-based fine-tuning

Variance-Based vs. MSE-Based Fine-Tuning



MSE-based fine-tuning

Thank You!

- ▶ A fine-tuning then rectification framework
- ▶ Optimal sample allocation between fine-tuning and rectification
- ▶ Data-driven procedure to guide sample allocation
- ▶ Empirical demonstration using Wine Reviews dataset



reBandit: Random Effects based Online RL algorithm for Reducing Cannabis Use

Susobhan Ghosh
Harvard University

**Joint work with Yongyi Guo, Pei-Yao Hung, Lara Coughlin, Erin E. Bonar, Inbal Nahum-Shani,
Maureen Walton, Susan A. Murphy**



Harvard John A. Paulson
School of Engineering
and Applied Sciences



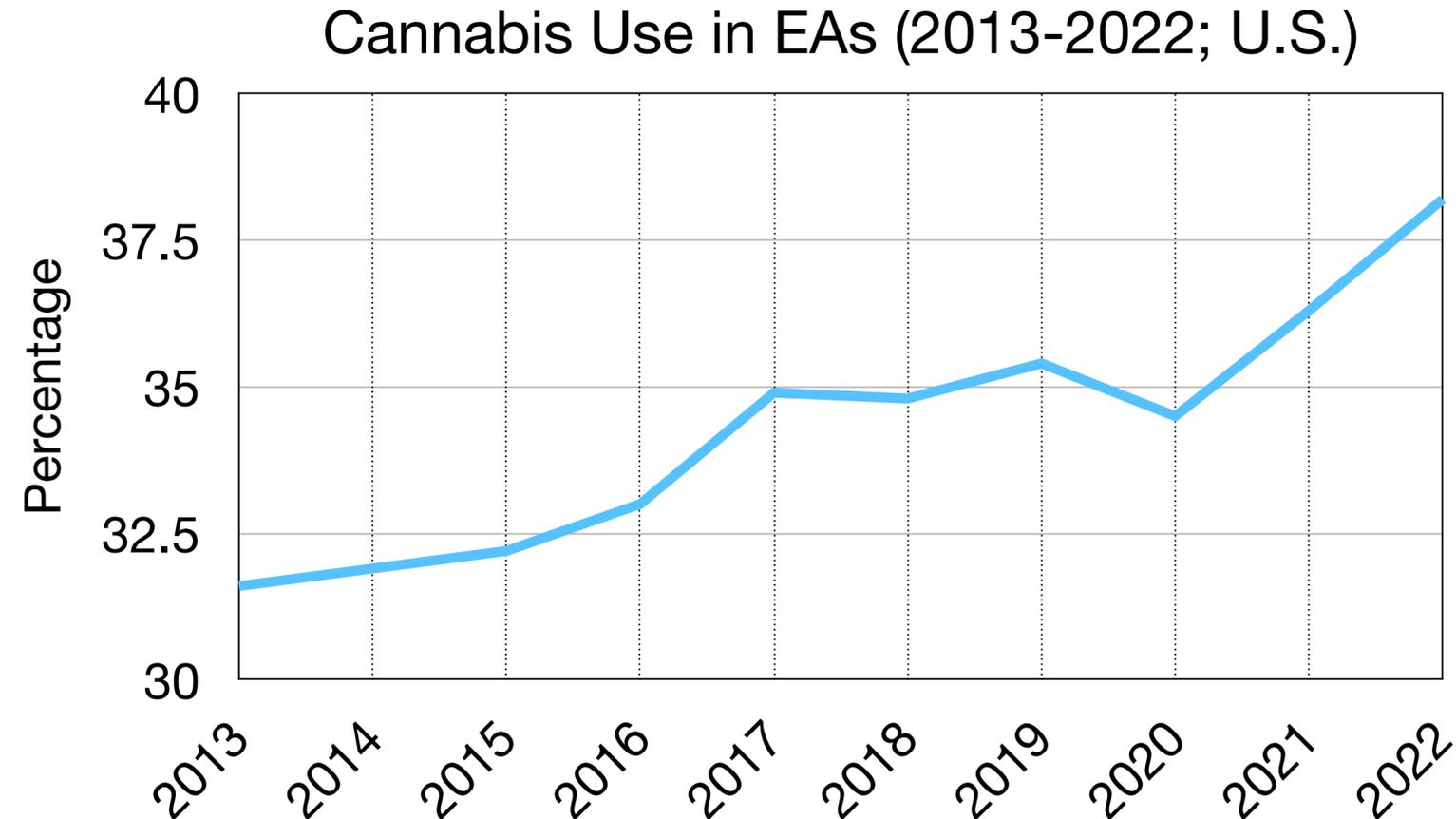
NIH/NIDA P50DA054039
NIH/NIBIB and OD P41EB028242
NIH/NIDCR UH3DE028723

Introduction: Cannabis Use

- **Cannabis legalization** in the US: Public perception of cannabis being **less risky** than in prior decades

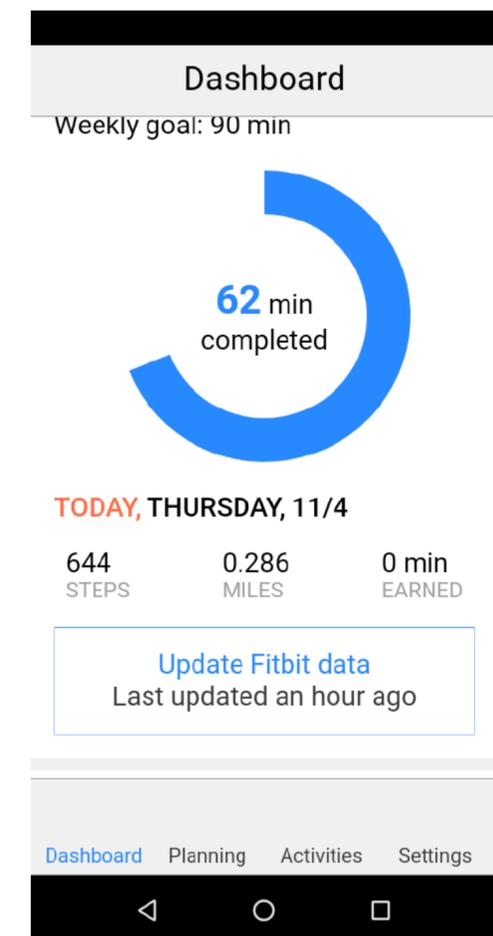
Introduction: Cannabis Use

- **Cannabis legalization** in the US: Public perception of cannabis being **less risky** than in prior decades
- Cannabis use - more prevalent among **emerging adults** (EAs) (age 18-25) than any other age group



Digital Interventions and Decision Making Algorithms

- Optimize intervention delivery to improve some short term outcome, which is on the pathway to longer term behavioral outcome



HeartSteps (*Liao et al, 2019*)

MiWaves Pilot Study and reBandit

March to May 2024

PURPOSE

Just-In-Time Adaptive Intervention (JITAI) to help reduce cannabis use amongst emerging adults (EAs) (ages 18-25)



Photo by مهدي کردی on Unsplash

INTERVENTION

Participants are delivered digital interventions through smartphone prompts at most twice a day.

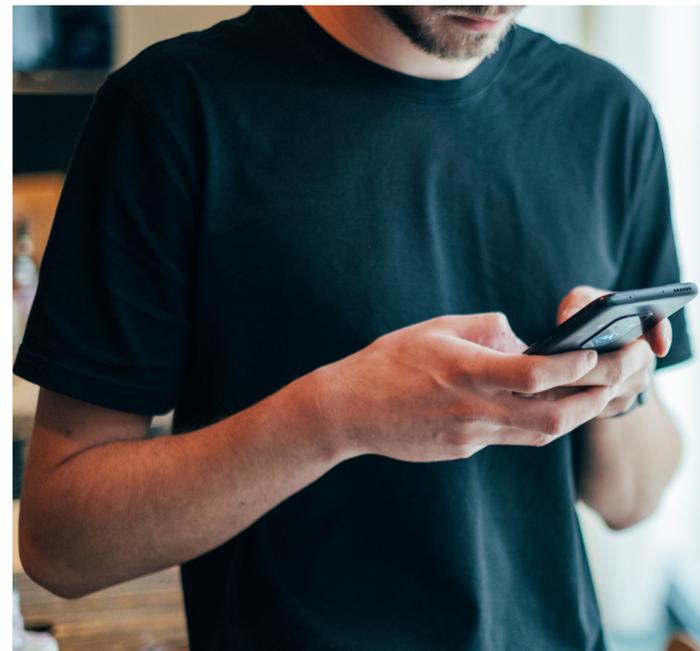
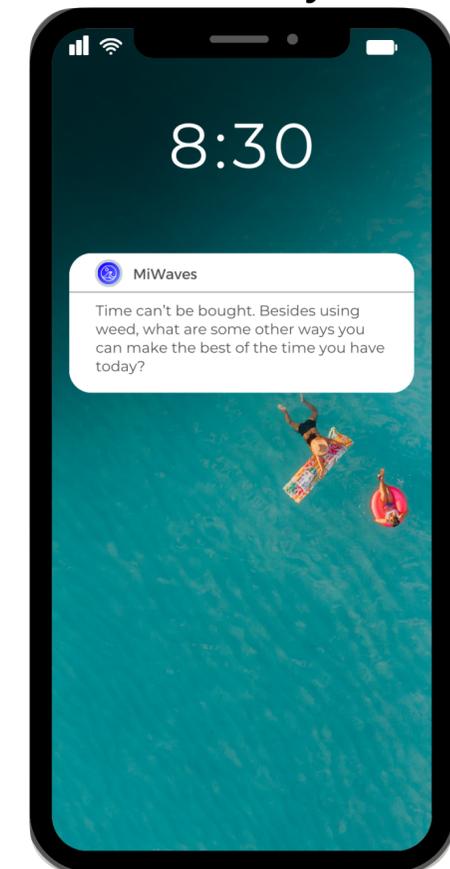


Photo by Jonas Leupe on Unsplash

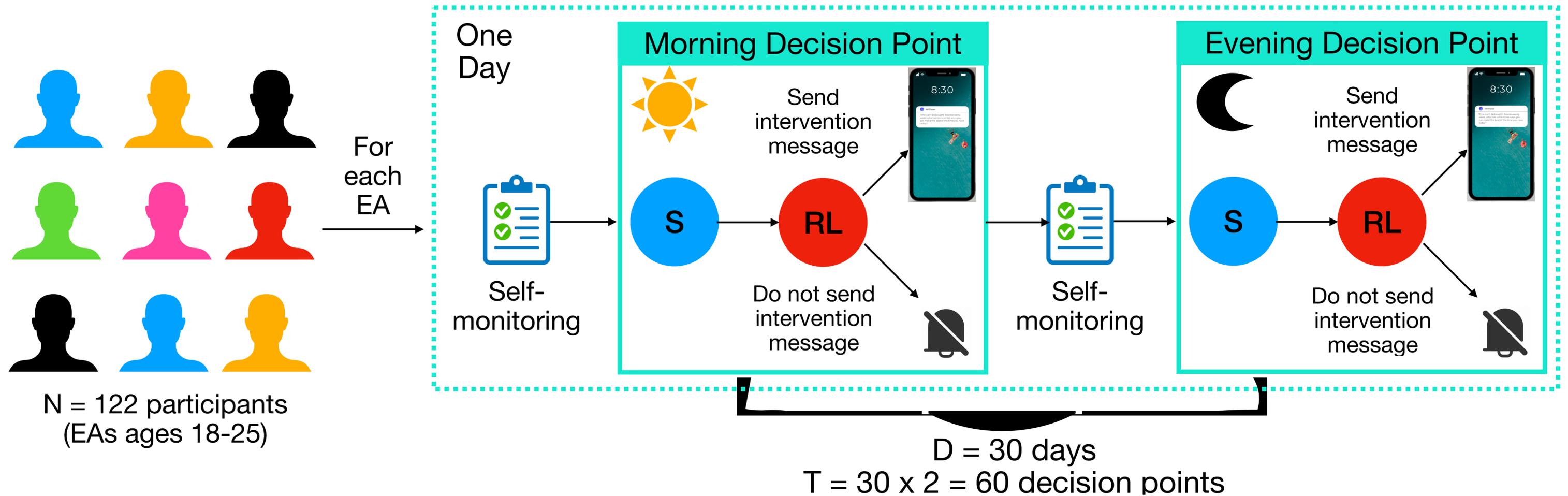
OUR CONTRIBUTION

Autonomous online RL algorithm called **reBandit** that personalizes likelihood of intervention prompt delivery



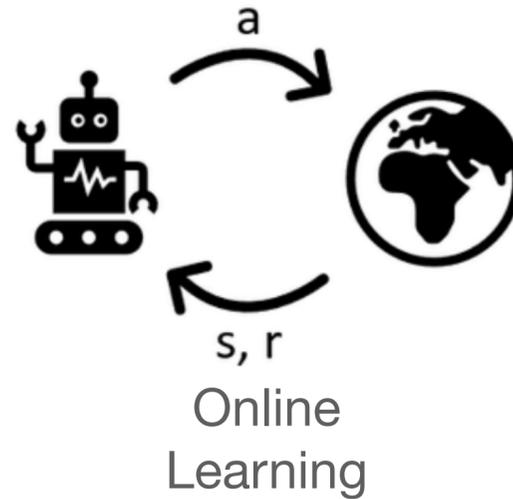


MiWaves Pilot Study: Overview

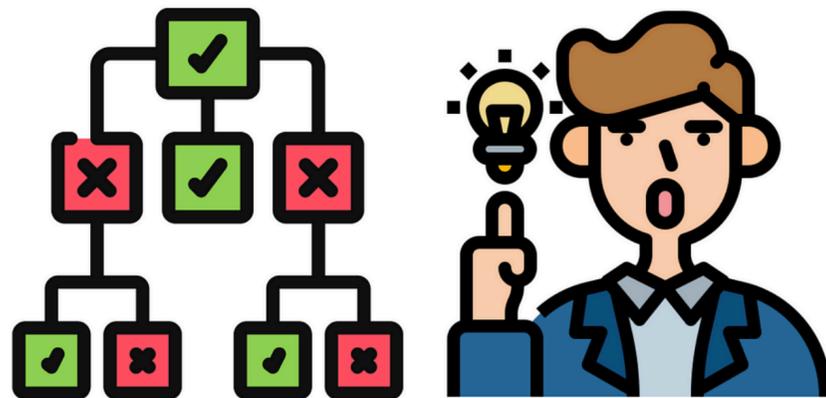


Goal: Reduce cannabis-use by improving participant's app engagement, behavior change by making participants self-monitor regularly

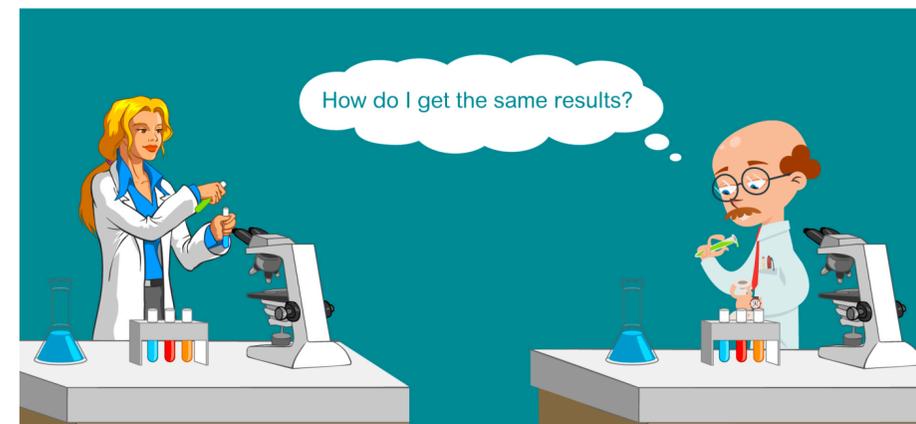
Decision-Making in Digital Interventions: Challenges



Autonomy and Stability



Interpretability



Reproducibility

Our contributions

Our contributions

- Developed reBandit, an online Bayesian RL algorithm

Our contributions

- Developed reBandit, an online Bayesian RL algorithm
- Autonomous deployment of reBandit - released as a software package/API
 - Publicly available and deployed live on-field
 - Autonomous updates
 - Real-time monitoring

Our contributions

- Developed reBandit, an online Bayesian RL algorithm
- Autonomous deployment of reBandit - released as a software package/API
 - Publicly available and deployed live on-field
 - Autonomous updates
 - Real-time monitoring
- Simulation environment to benchmark variants of RL algorithm

Our contributions

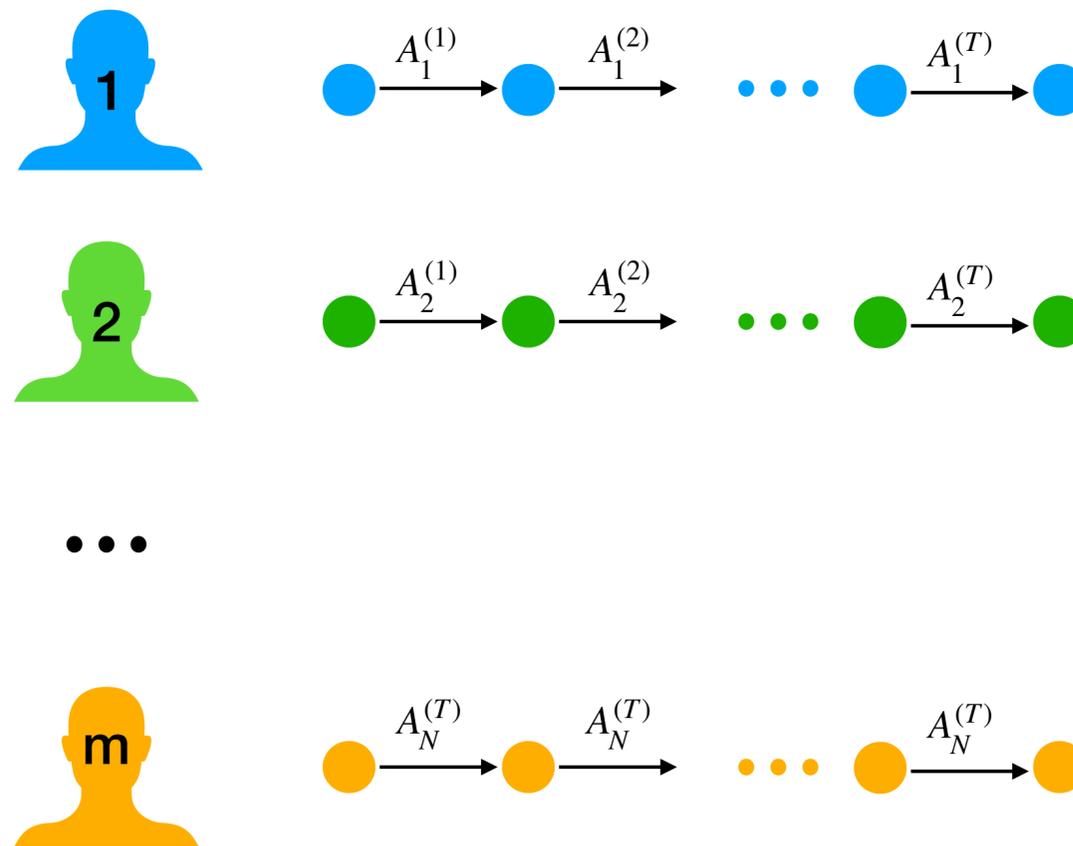
- Developed reBandit, an online Bayesian RL algorithm
- Autonomous deployment of reBandit - released as a software package/API
 - Publicly available and deployed live on-field
 - Autonomous updates
 - Real-time monitoring
- Simulation environment to benchmark variants of RL algorithm
- Bayesian priors to warm start reBandit

Our contributions

- Developed reBandit, an online Bayesian RL algorithm
- Autonomous deployment of reBandit - released as a software package/API
 - Publicly available and deployed live on-field
 - Autonomous updates
 - Real-time monitoring
- Simulation environment to benchmark variants of RL algorithm
- Bayesian priors to warm start reBandit

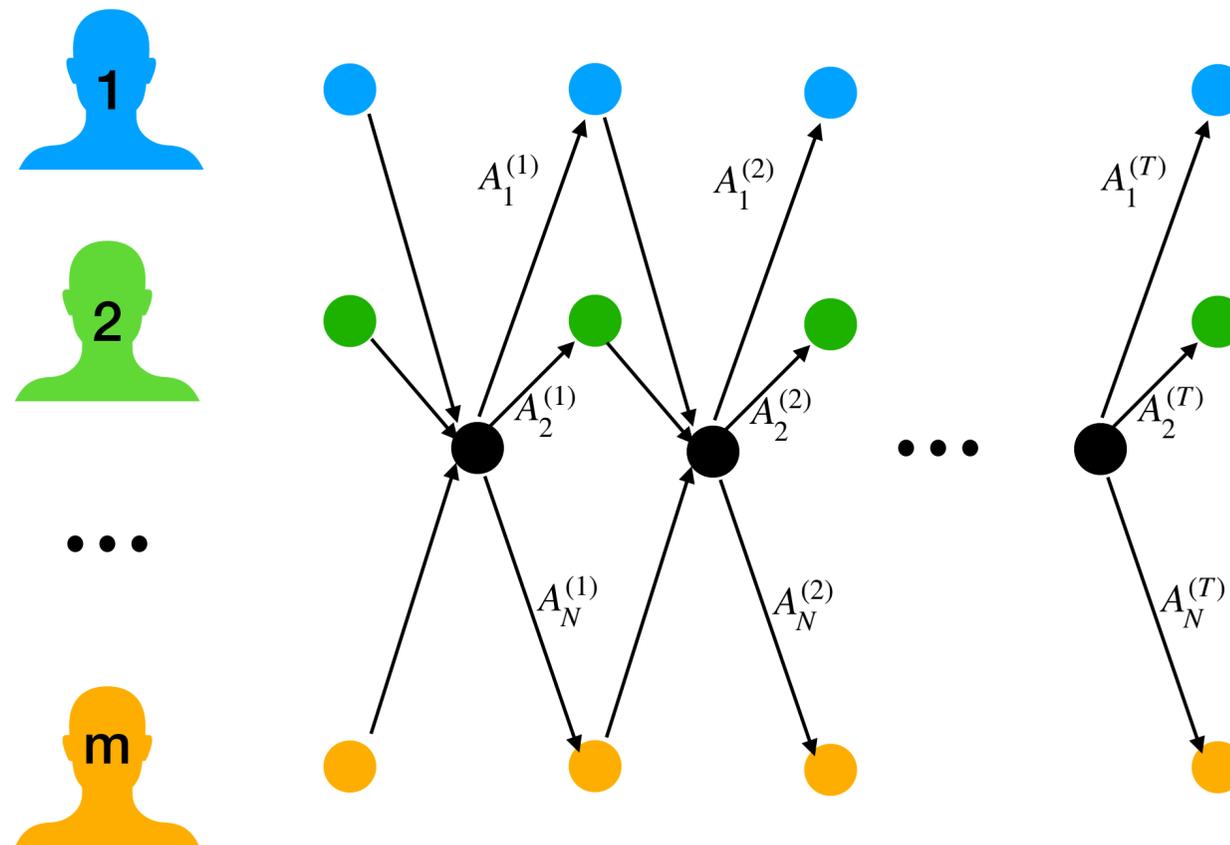
Fully Personalized Approach

- Separate model for each participant
- Requires **lot of data** to personalize
- **Only 60 data points per participant in MiWaves - not feasible!**



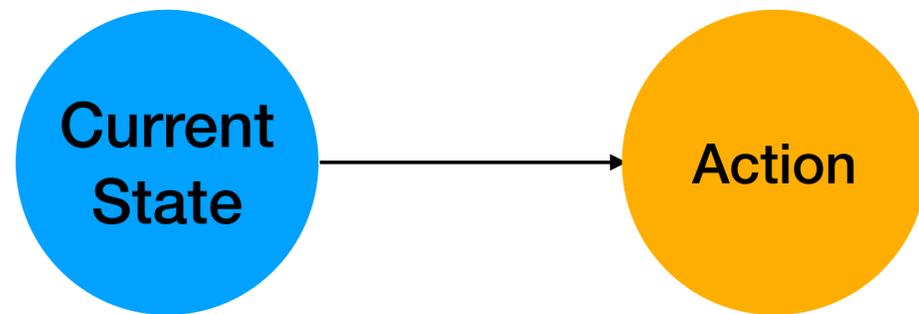
Fully Pooled Approach

- One model for the entire population of participants
- **Population heterogeneity** - if substantial, may personalize **poorly**
- **Can we adaptively pool and personalize?**



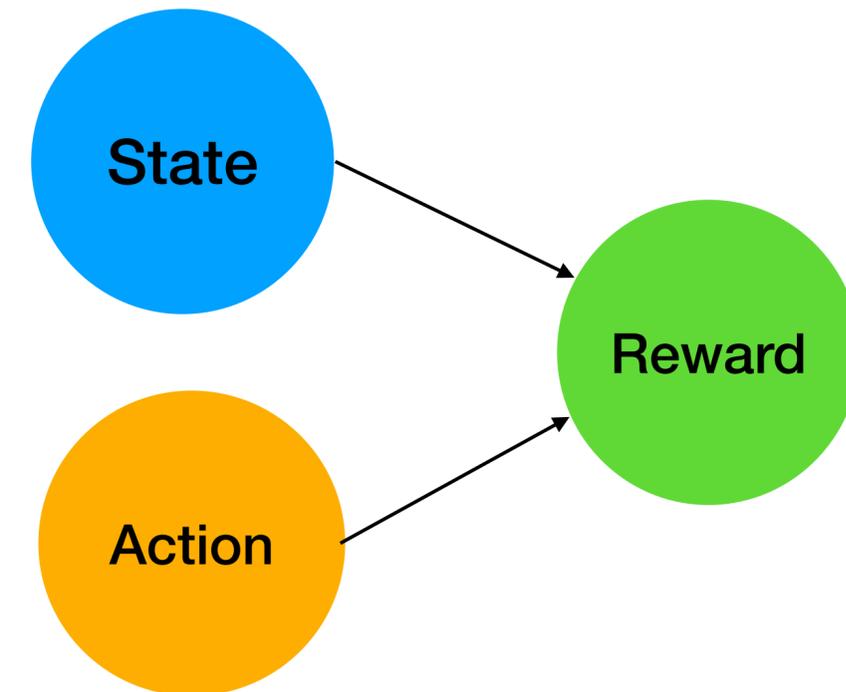
Online RL: Two Components

- Optimization algorithm or action-selection procedure



Function that outputs action to take, given some current state

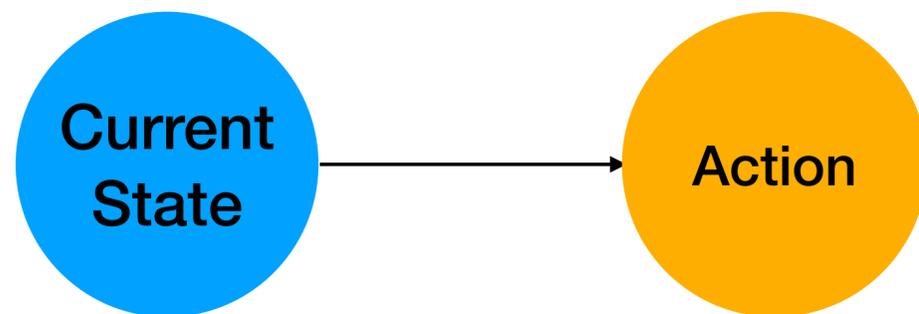
- Learning algorithm



Function that models the reward given some state and action

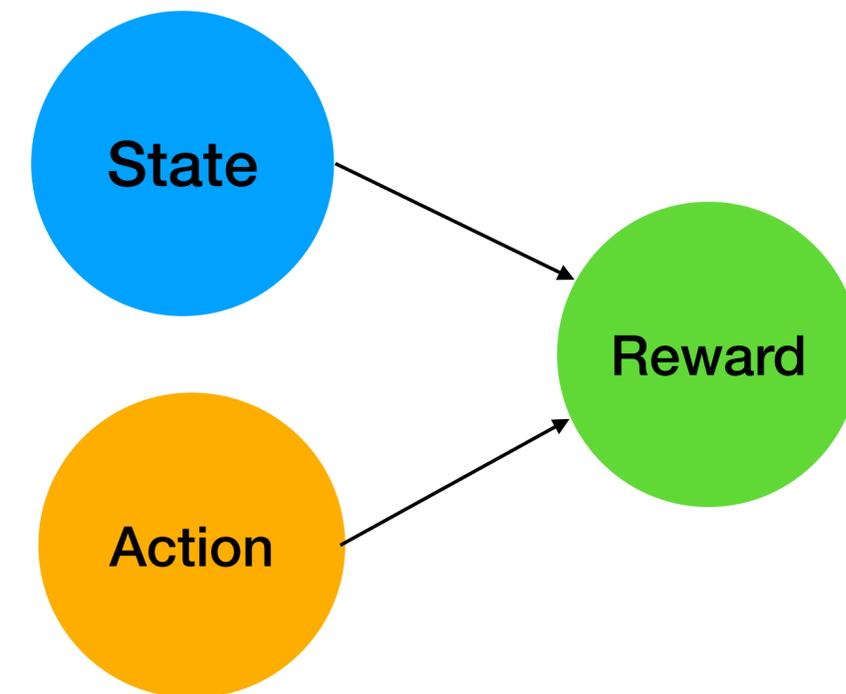
Online RL: Two Components

- Optimization algorithm or action-selection procedure



Function that outputs action to take, given some current state

- Learning algorithm



Function that models the reward given some state and action

MiWaves RL Framework

- **Time Horizon:** 30 days, 2 data points per day **60**
- **States:** $S = \{S_1, S_2, S_3\}$ - Binary features **8**
- **Actions:** Binary, i.e. $A = \{0, 1\}$ - show or do not show intervention. **2**
- **Reward:** $R \in \{0, 1, 2, 3\}$ increases linearly with engagement

MiWaves RL Framework

- **Time Horizon:** 30 days, 2 data points per day **60**
- **States:** $S = \{S_1, S_2, S_3\}$ - Binary features **8**
- **Actions:** Binary, i.e. $A = \{0, 1\}$ - show or do not show intervention. **2**
- **Reward:** $R \in \{0, 1, 2, 3\}$ increases linearly with engagement

8 x 2 = 16 weights

Noisy Environment

Reward model

Mixed Effects Model: m users

- Using a mixed effects model, we approximate the reward as:

$$R_i^{(t)} = \Phi_{it}^T \theta_i + \epsilon_i^{(t)}$$

$$\theta_i = \theta_{pop} + u_i$$

$$\Phi_{it}^T = \Phi_{it}^T(s, a)$$

By definition, $u_i \sim \mathcal{N}(\mathbf{0}, \Sigma_u)$ $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma_\epsilon^2 \mathbf{I}_{mt})$

- Assuming gaussian prior : $\theta_{pop} \sim \mathcal{N}(\mu_{prior}, \Sigma_{prior})$

Reward model

Mixed Effects Model: m users

- Using a mixed effects model, we approximate the reward as:

$$R_i^{(t)} = \Phi_{it}^T \theta_i + \epsilon_i^{(t)}$$

$$\theta_i = \theta_{pop} + u_i$$

$$\Phi_{it}^T = \Phi_{it}^T(s, a)$$

Hyper-parameters

By definition, $u_i \sim \mathcal{N}(\mathbf{0}, \Sigma_u)$ $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma_\epsilon^2 \mathbf{I}_{mt})$

- Assuming gaussian prior : $\theta_{pop} \sim \mathcal{N}(\mu_{prior}, \Sigma_{prior})$

Posterior Update

Posteriors for parameters θ

$$\mu_{post} = \left(\tilde{\Sigma}_{\theta}^{-1} + \frac{1}{\sigma_{\epsilon}^2} \mathbf{A} \right)^{-1} \left(\tilde{\Sigma}_{\theta}^{-1} \mu_{\theta} + \frac{1}{\sigma_{\epsilon}^2} \mathbf{B} \right)$$

$$\Sigma_{post} = \left(\tilde{\Sigma}_{\theta}^{-1} + \frac{1}{\sigma_{\epsilon}^2} \mathbf{A} \right)^{-1}$$

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_m \end{bmatrix}$$

where $\mu_{\theta} = \mathbb{J}_{m,1} \otimes \mu_{prior}$ $\tilde{\Sigma}_{\theta} = \mathbb{J}_{m,m} \otimes \Sigma_{prior} + \mathbb{I}_m \otimes \Sigma_u$

$$\mathbf{A} = \text{BlockDiag}(\mathbf{A}_1, \dots, \mathbf{A}_m) \quad \mathbf{A}_i = \sum_{\tau=1}^t \Phi_{i\tau} \Phi_{i\tau}^T$$

$$\mathbf{B}^T = [\mathbf{B}_1^T \ \dots \ \mathbf{B}_m^T] \quad \mathbf{B}_i = \sum_{\tau=1}^t \Phi_{i\tau} R_i^{(\tau)}$$

Posterior Update

$$\mu_{post} = \left(\tilde{\Sigma}_{\theta}^{-1} + \frac{1}{\sigma_{\epsilon}^2} \mathbf{A} \right)^{-1} \left(\tilde{\Sigma}_{\theta}^{-1} \mu_{\theta} + \frac{1}{\sigma_{\epsilon}^2} \mathbf{B} \right)$$

$$\Sigma_{post} = \left(\tilde{\Sigma}_{\theta}^{-1} + \frac{1}{\sigma_{\epsilon}^2} \mathbf{A} \right)^{-1}$$

where $\mu_{\theta} = \mathbb{J}_{m,1} \otimes \mu_{prior}$ $\tilde{\Sigma}_{\theta} = \mathbb{J}_{m,m} \otimes \Sigma_{prior} + \mathbb{I}_m \otimes \Sigma_u$

$$\mathbf{A} = \text{BlockDiag}(\mathbf{A}_1, \dots, \mathbf{A}_m) \quad \mathbf{A}_i = \sum_{\tau=1}^t \Phi_{i\tau} \Phi_{i\tau}^T$$

$$\mathbf{B}^T = [\mathbf{B}_1^T \ \dots \ \mathbf{B}_m^T] \quad \mathbf{B}_i = \sum_{\tau=1}^t \Phi_{i\tau} R_i^{(\tau)}$$

Posteriors for parameters θ

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_m \end{bmatrix}$$

Estimates

Hyper-parameter update

- Updating the variance terms - by maximizing the (log) marginal likelihood of observed rewards, marginalized over the parameters θ

$$\Sigma_u, \sigma_\epsilon^2 = \operatorname{argmax} l(\Sigma_u, \sigma_\epsilon^2; \mathcal{H})$$

History

subject to the constraints: $\Sigma_u \succ 0$ and $\sigma_\epsilon^2 > 0$

- Solve online using gradient descent, use optimization tricks like using Cholesky decomposition and positive eigenvalue checks

Experimental Results

500 simulated trials

Using SARA based simulation testbed

- reBandit vs fully pooled vs random:

reBandit consistently performs better than fully pooled and random action selection across all simulation environments, across all metrics.

Learning Algorithm	Minimal TE		Low TE		High TE	
	Mean TRPU	SE	Mean TRPU	SE	Mean TRPU	SE
reBandit	128.535	0.176	129.441	0.170	132.245	0.164
fully pooled	127.773	0.162	129.097	0.166	132.209	0.164
random	127.830	0.162	128.970	0.170	131.050	0.168

Environments with no habituation effect

TE: Treatment Effect; TRPU: Total Reward Per User; SE: Standard Error

MiWaves Pilot Study Results

- 77% engagement rate
- After-study analysis paper - recently submitted
- Resampling analysis for RL algorithm personalization - currently ongoing
- HCI-centric analysis of the intervention - recently published

Thank You



Personal Website



Link to reBandit paper

INTERACTION LEVEL	SHORT LENGTH	LONG LENGTH
A (acknowledge the message)	You are the artist and the future is your canvas. When you think about your life in 6 months from now, what do you hope for?	Do you ever find yourself burying or suppressing your emotions? Talking your feelings out with someone you trust or writing a private journal entry can help you free those emotions.
B (participant requested to visit external resource)	What's your favorite song? Learn more about how music is beneficial to your mental health: https://www.youtube.com/watch?v=zJ2YGLuzGfo	Trying new things doesn't have to be expensive. When you're tight on cash, check out this list to see if any of these cheap and easy hobbies seem interesting: https://www.buzzfeed.com/tomvellner/cheap-easy-hobbies
C (requires input from participant)	Fun fact: even just five mins of physical activity can be beneficial. How do you get your body moving? <hr/>	You have the power to achieve anything you put your mind to. From this list, what are some ways you are interested in building a plan to create the future that you hope for yourself? (A) Writing goals (B) Attending therapy (C) Budget money (D) Establish healthy routines



ROAD 3

Mitigating Observed and Unobserved Confounding in Observational Data

Vasiliki Stoumpou, Georgios Antonios Margonis,
Dimitris Bertsimas

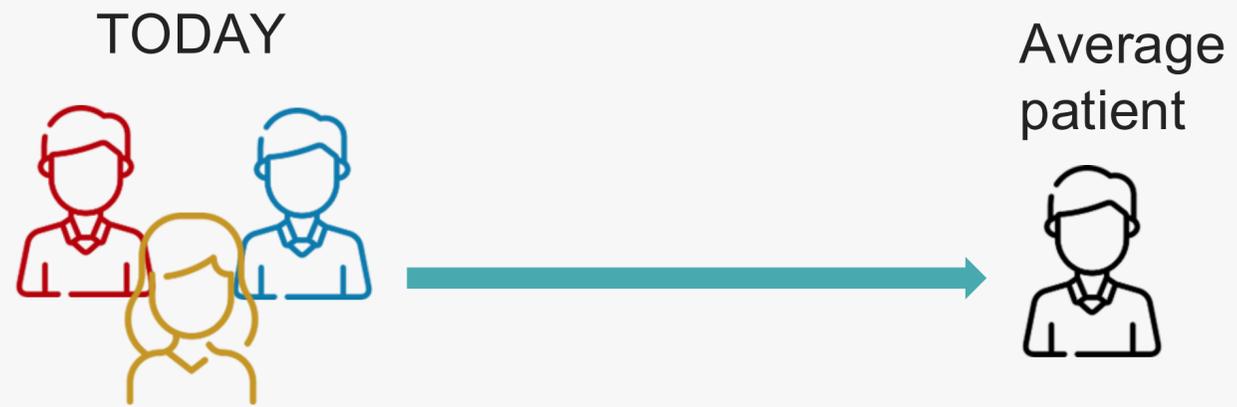


Precision Medicine: The Future of Healthcare

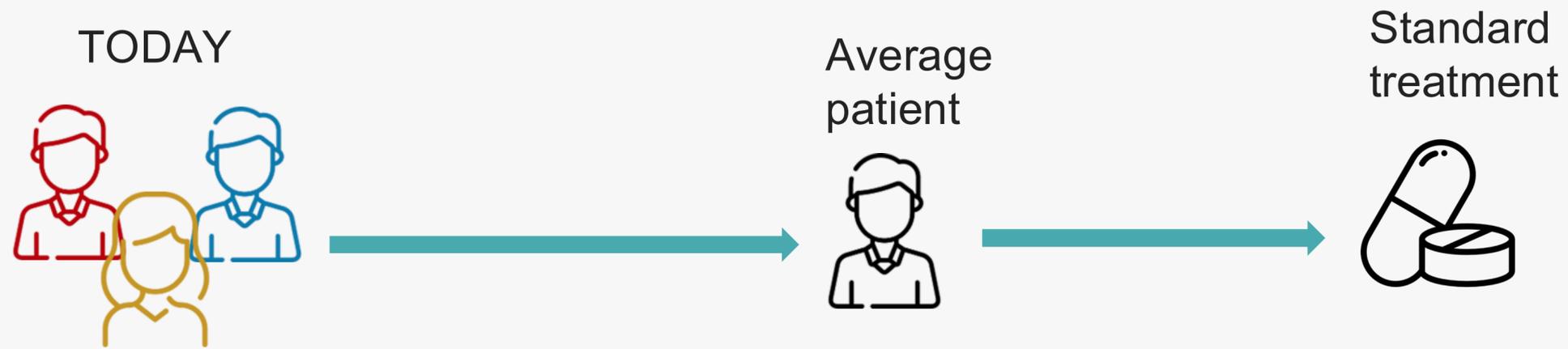
TODAY



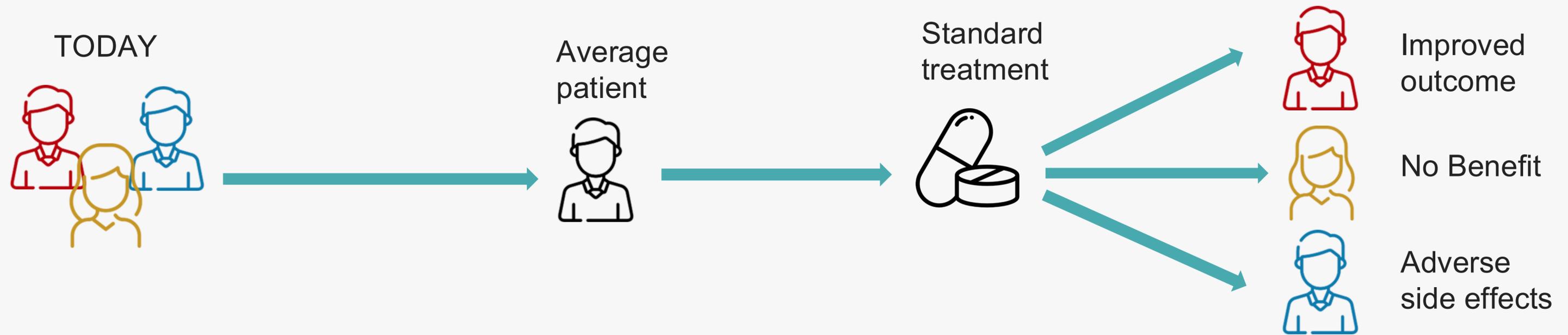
Precision Medicine: The Future of Healthcare



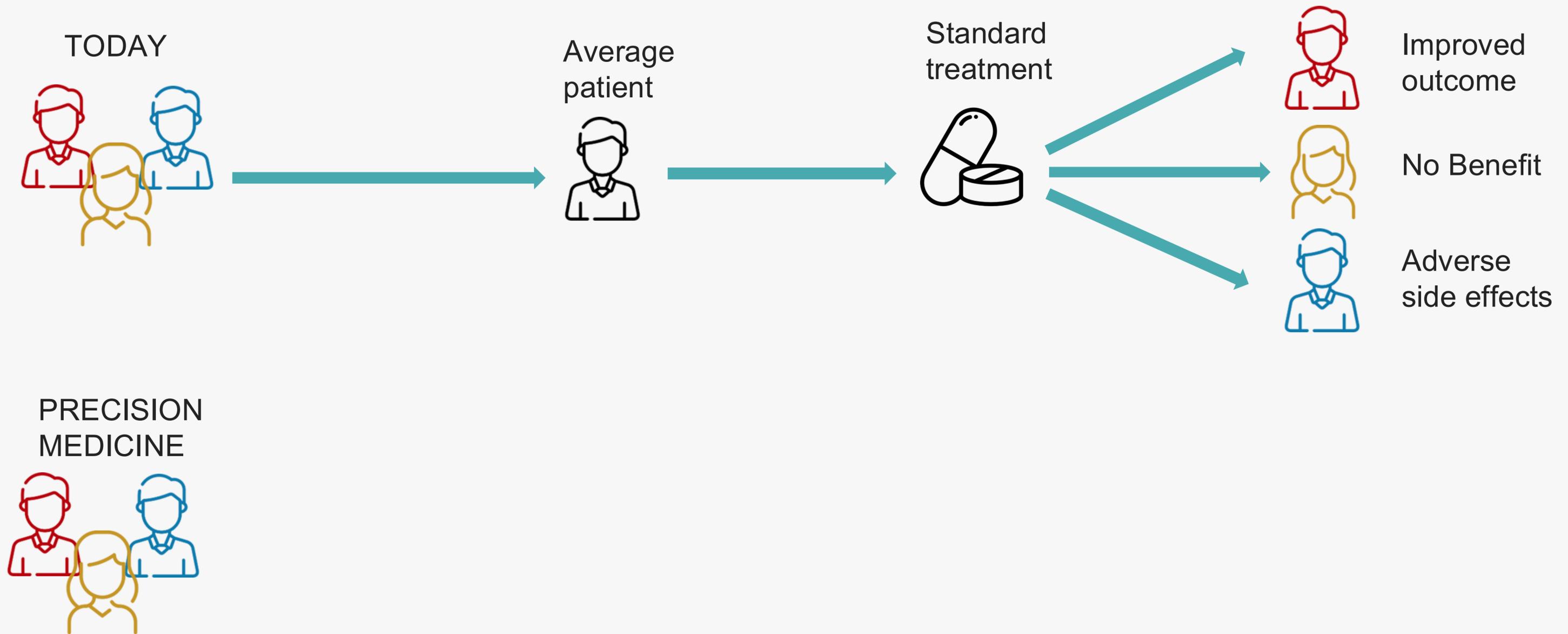
Precision Medicine: The Future of Healthcare



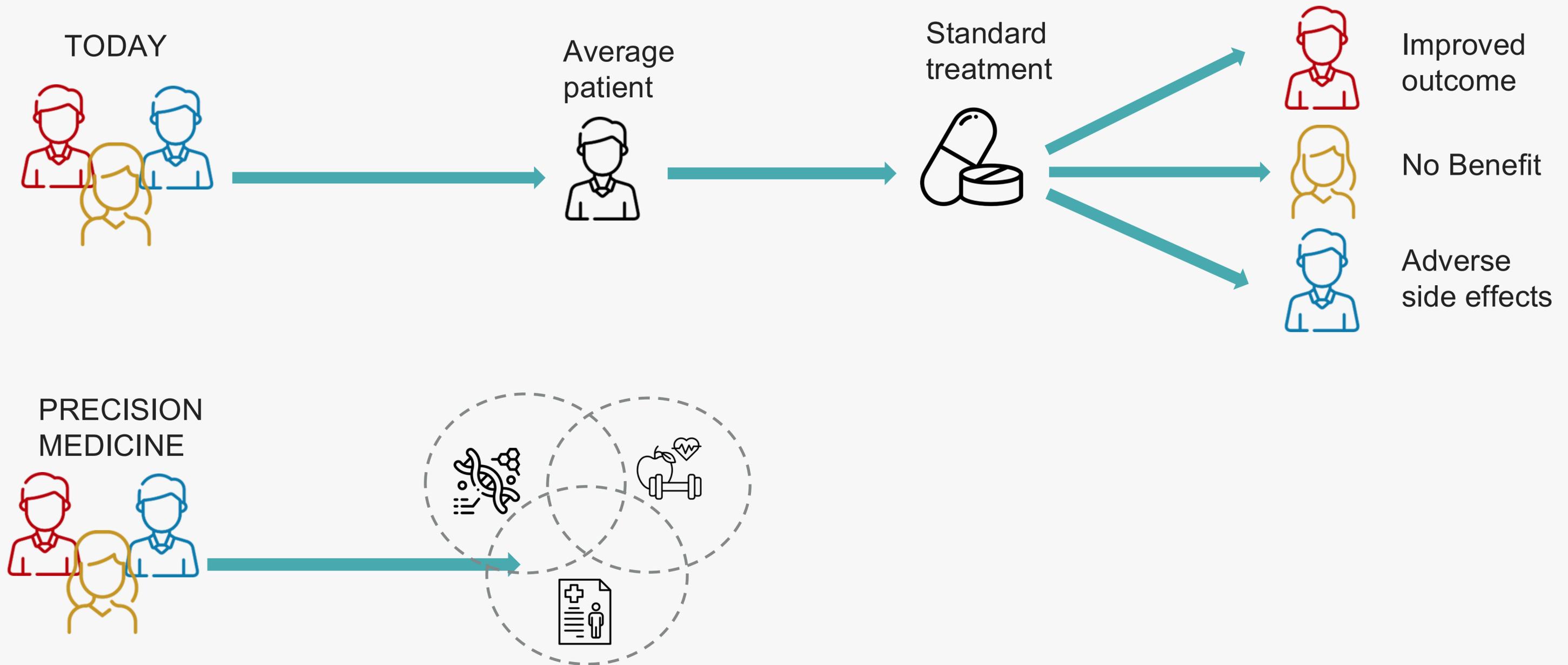
Precision Medicine: The Future of Healthcare



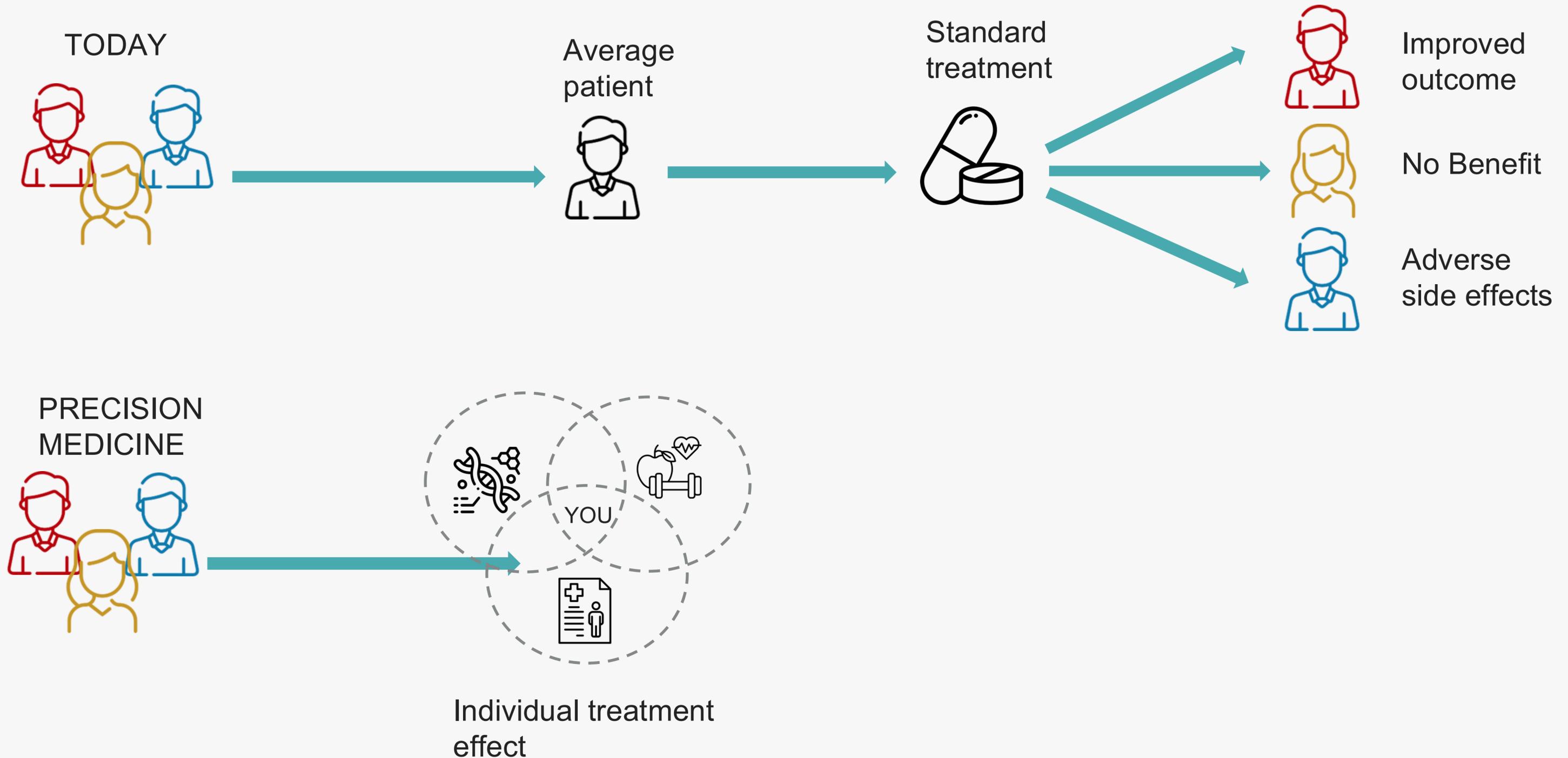
Precision Medicine: The Future of Healthcare



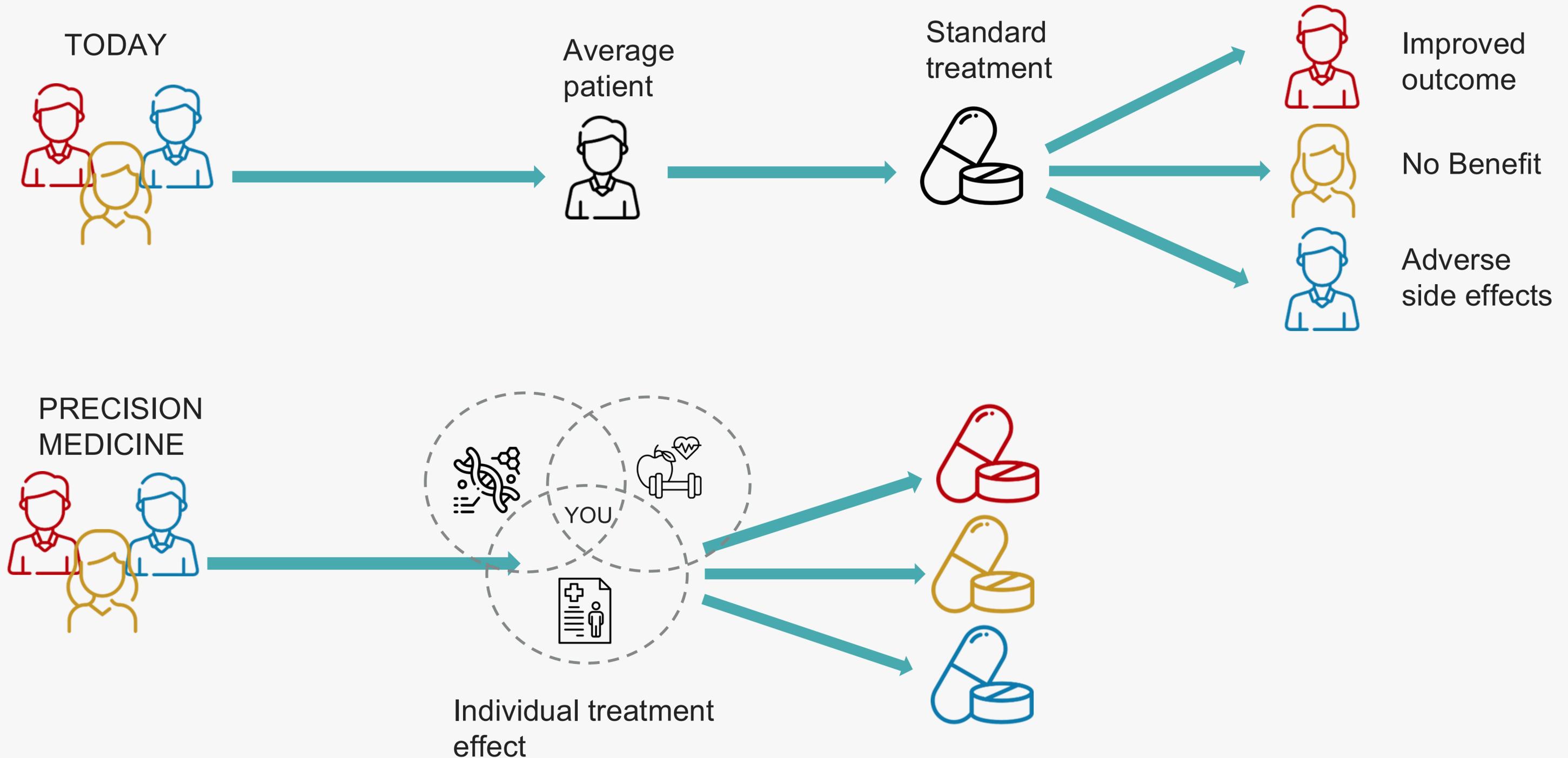
Precision Medicine: The Future of Healthcare



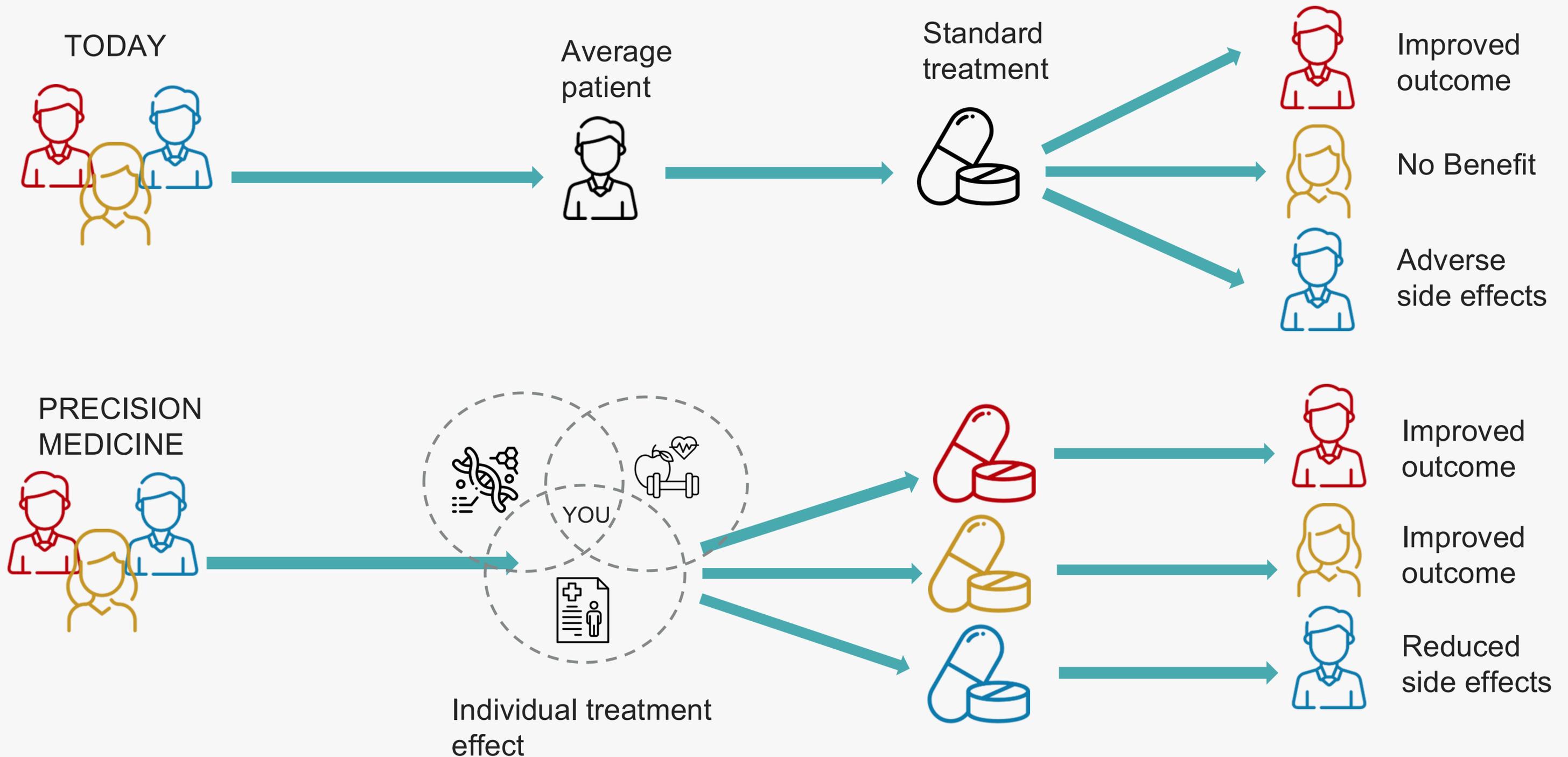
Precision Medicine: The Future of Healthcare



Precision Medicine: The Future of Healthcare

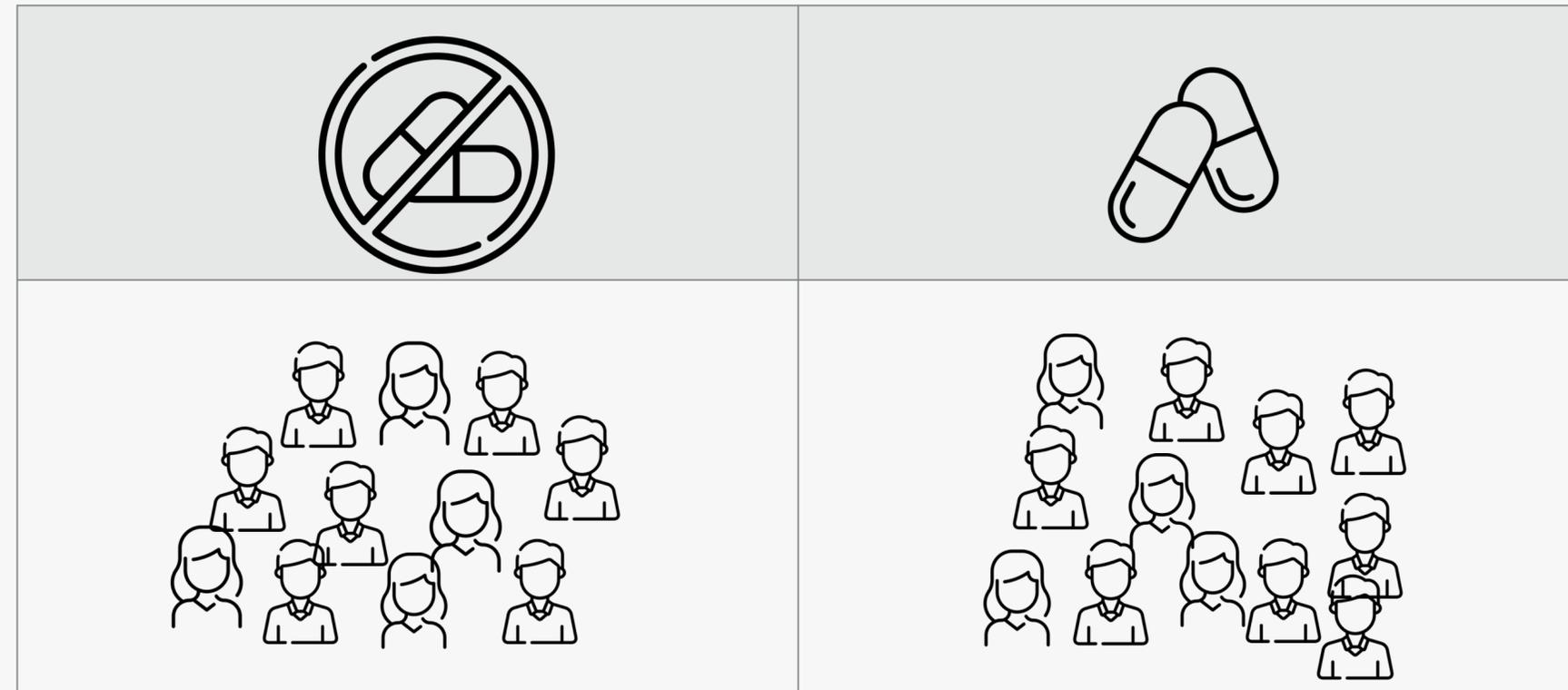


Precision Medicine: The Future of Healthcare



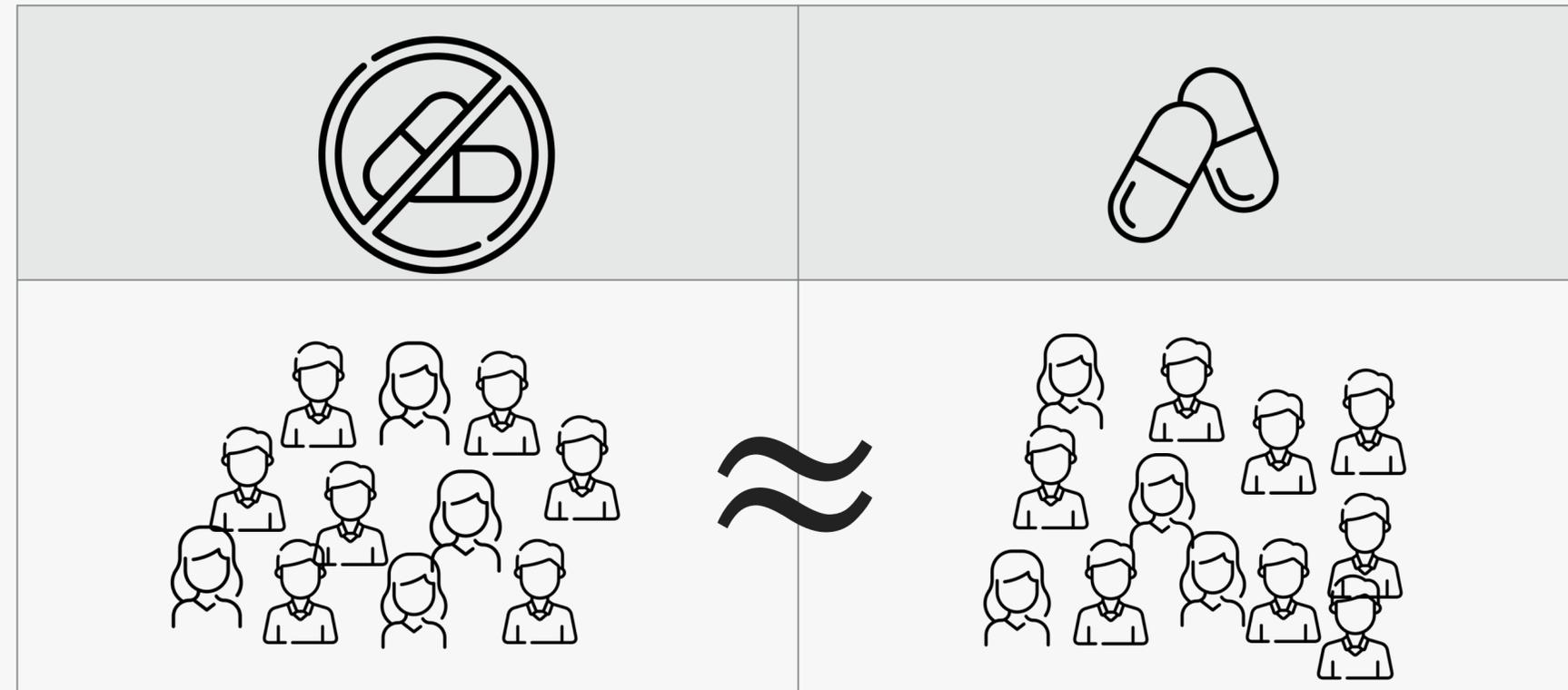
How are treatment effects estimated?

Randomized
Clinical Trial
(RCT)



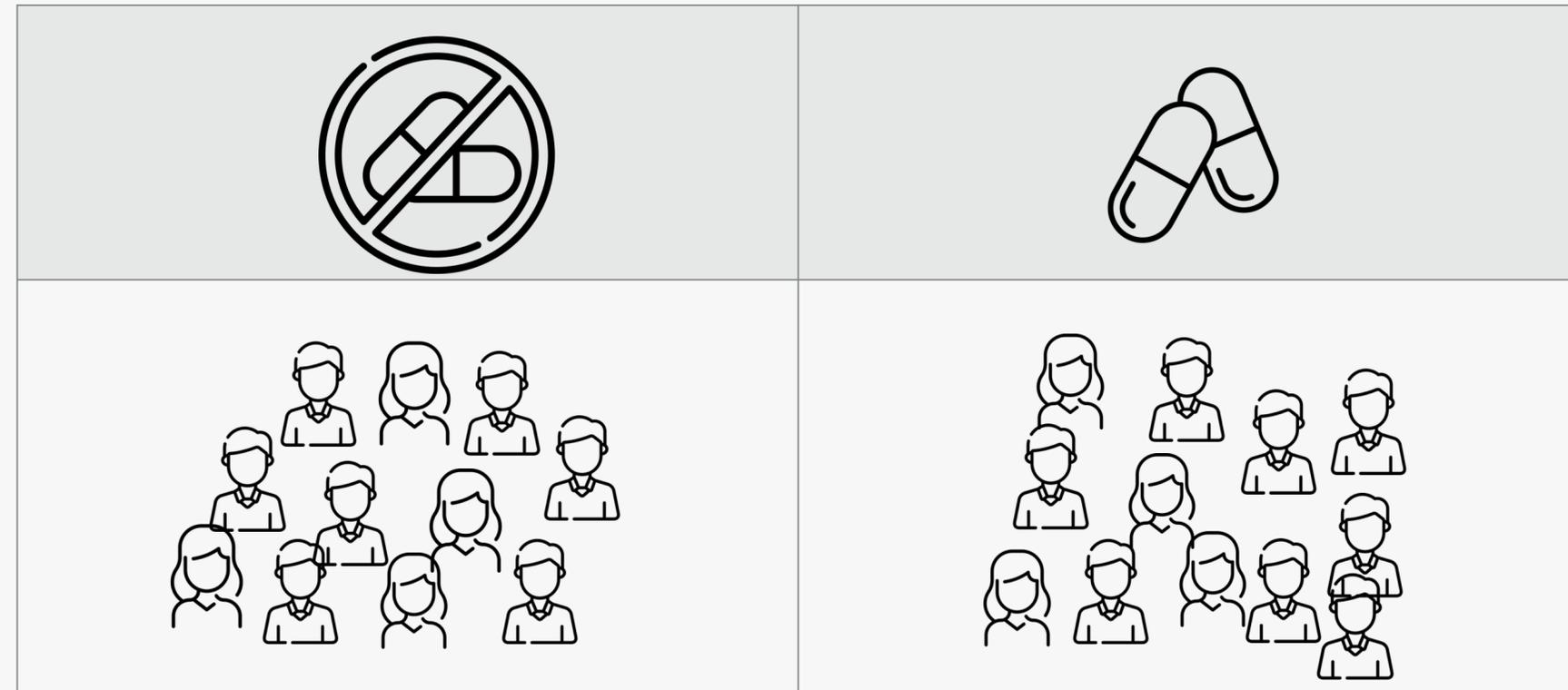
How are treatment effects estimated?

Randomized
Clinical Trial
(RCT)



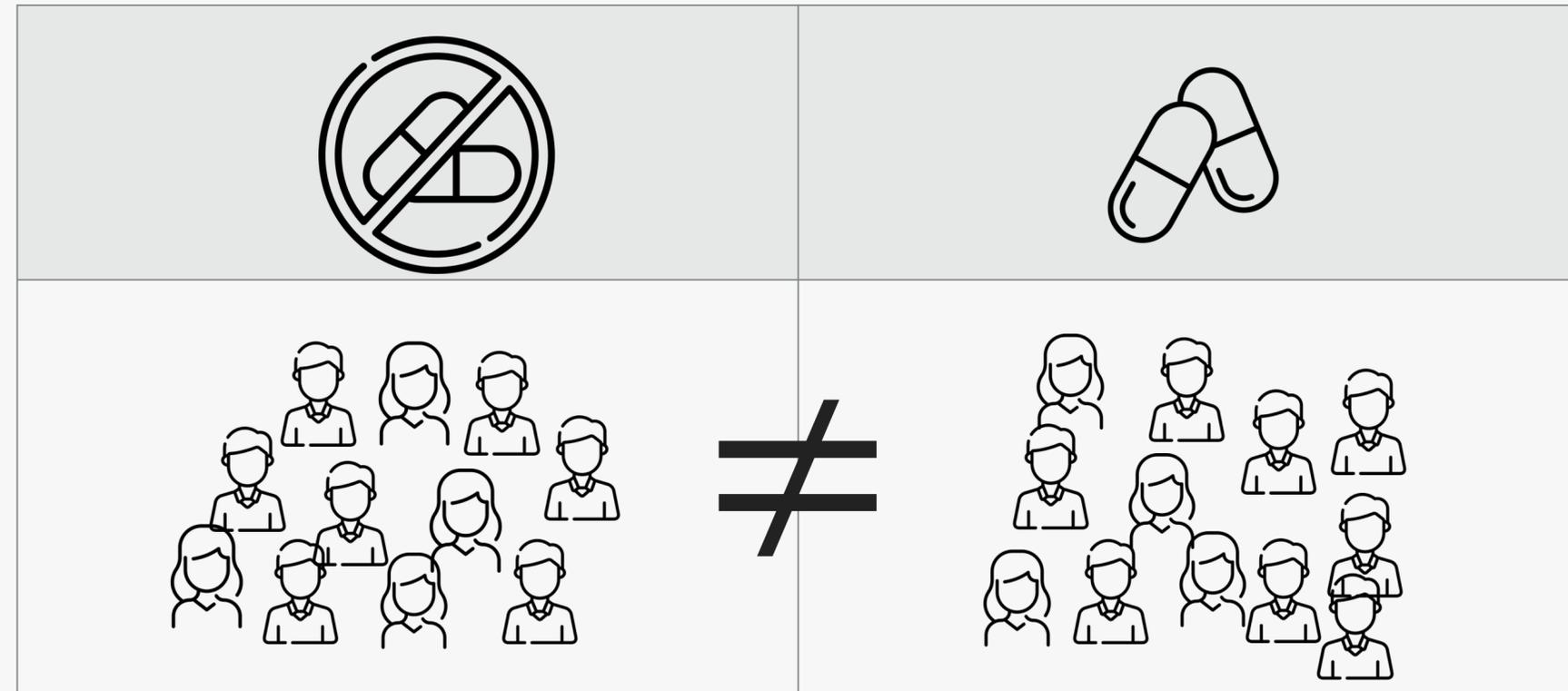
Are RCTs the golden standard?

Observational data



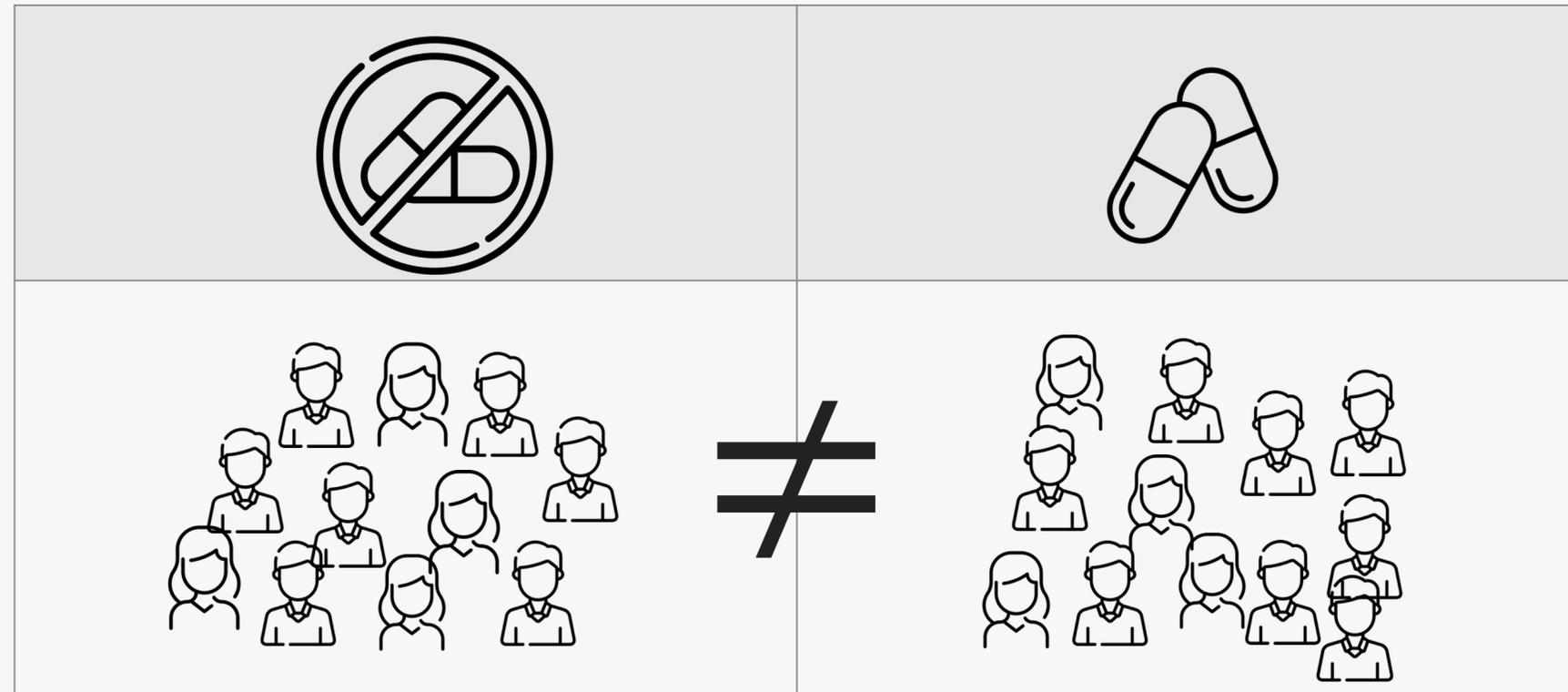
Are RCTs the golden standard?

Observational data



Are RCTs the golden standard?

Observational data



Type of confounding	What is it?	Examples	Solutions
Observed	Can be explained by differences in observed features	Age, Size of the tumor, Sex etc.	Covariate adjustment Propensity matching Propensity weighting...
Unobserved	Due to factors that we cannot explain	Mutations	Rely on ignorability assumption

Are Observational data that problematic?

Observational data: Patients with GIST
Gastrointestinal Stromal Tumor

Type of cancer in the digestive tract

Are Observational data that problematic?

Observational data: Patients with GIST
Gastrointestinal Stromal Tumor

Type of cancer in the digestive tract

Treatment of interest: Imatinib
(blocks enzymes that promote cancer growth)

Are Observational data that problematic?

Observational data: Patients with GIST
Gastrointestinal Stromal Tumor

Type of cancer in the digestive tract

Treatment of interest: Imatinib
(blocks enzymes that promote cancer growth)

We fit a Cox proportional hazards model
where the only feature is the treatment (Rx,
imatinib)

Are Observational data that problematic?

Observational data: Patients with GIST
Gastrointestinal Stromal Tumor

Type of cancer in the digestive tract

Treatment of interest: Imatinib
(blocks enzymes that promote cancer growth)

We fit a Cox proportional hazards model
where the only feature is the treatment (Rx,
imatinib)

Data type	Feature	Coefficient (b)	Hazard Ratio (HR=exp(b))

Are Observational data that problematic?

Observational data: Patients with GIST
Gastrointestinal Stromal Tumor

Type of cancer in the digestive tract

Treatment of interest: Imatinib
(blocks enzymes that promote cancer growth)

We fit a Cox proportional hazards model
where the only feature is the treatment (Rx,
imatinib)

Data type	Feature	Coefficient (b)	Hazard Ratio (HR=exp(b))

HR < 1: 

HR > 1: 

HR=0.46: the treatment has 46% the hazard
of control

Are Observational data that problematic?

Observational data: Patients with GIST
Gastrointestinal Stromal Tumor

Type of cancer in the digestive tract

Treatment of interest: Imatinib
(blocks enzymes that promote cancer growth)

We fit a Cox proportional hazards model
where the only feature is the treatment (Rx,
imatinib)

Data type	Feature	Coefficient (b)	Hazard Ratio (HR=exp(b))
RCT	Rx	-0.77	0.46

HR < 1:

HR > 1:

HR=0.46: the treatment has 46% the hazard
of control

Are Observational data that problematic?

Observational data: Patients with GIST
Gastrointestinal Stromal Tumor

Type of cancer in the digestive tract

Treatment of interest: Imatinib
(blocks enzymes that promote cancer growth)

We fit a Cox proportional hazards model
where the only feature is the treatment (Rx,
imatinib)

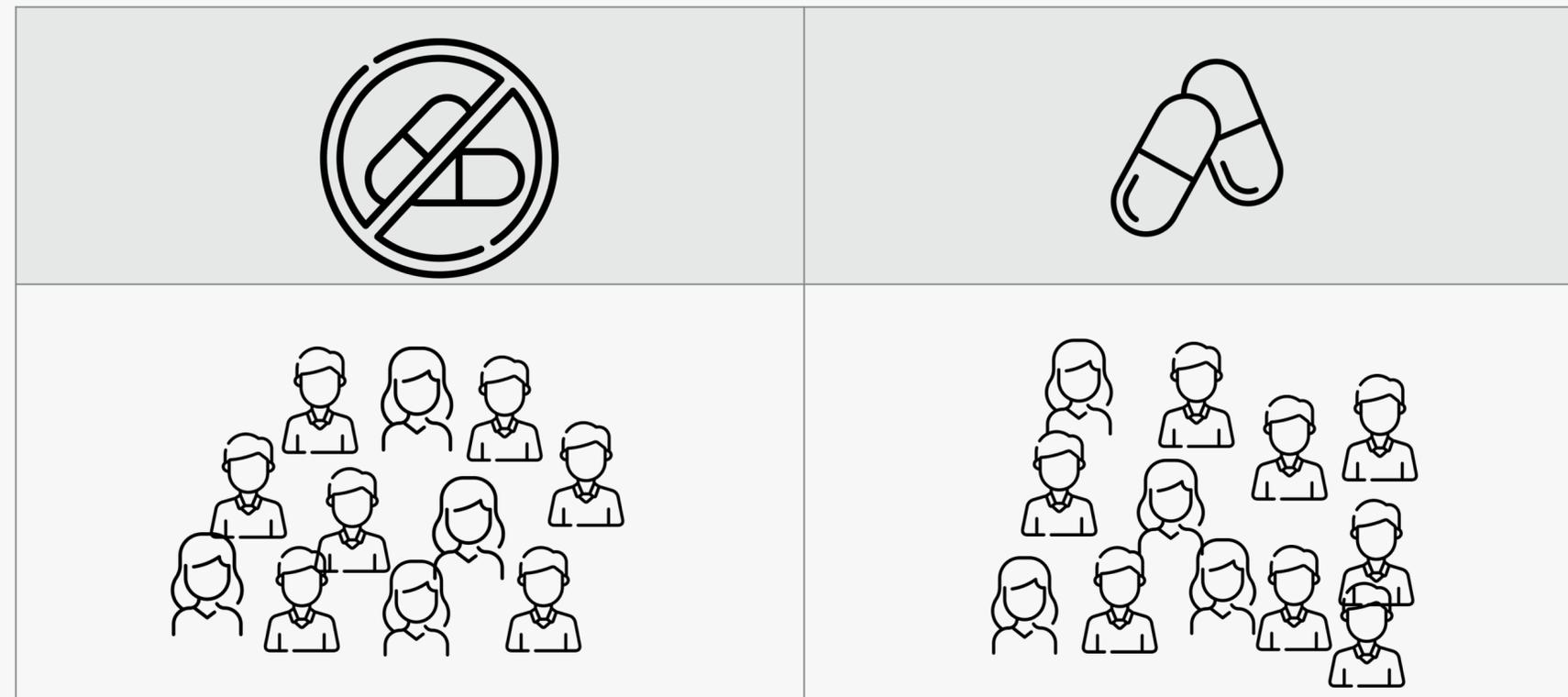
Data type	Feature	Coefficient (b)	Hazard Ratio (HR=exp(b))
RCT	Rx	-0.77	0.46
Observational	Rx	0.90	2.47

HR < 1:

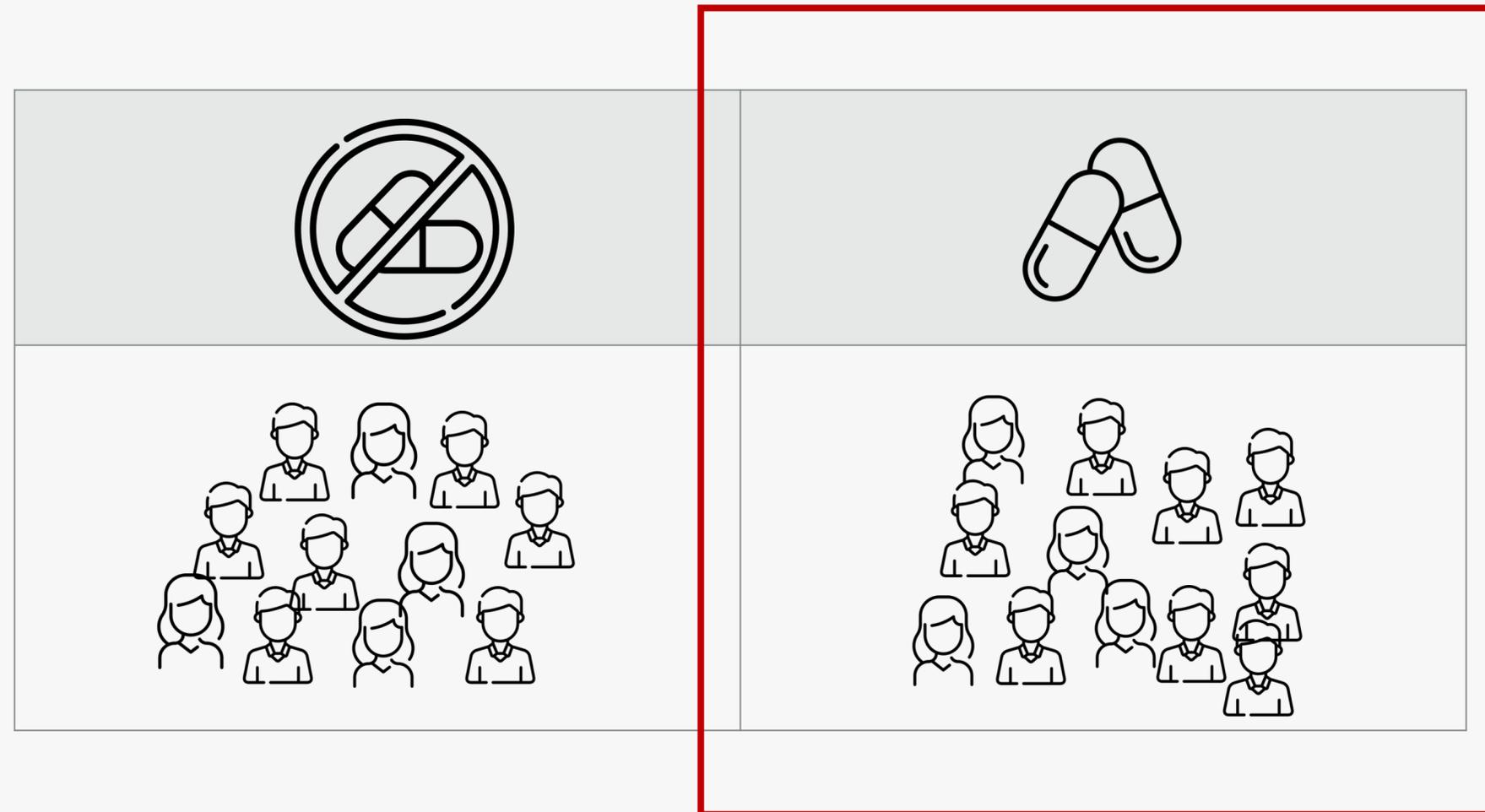
HR > 1:

HR=0.46: the treatment has 46% the hazard
of control

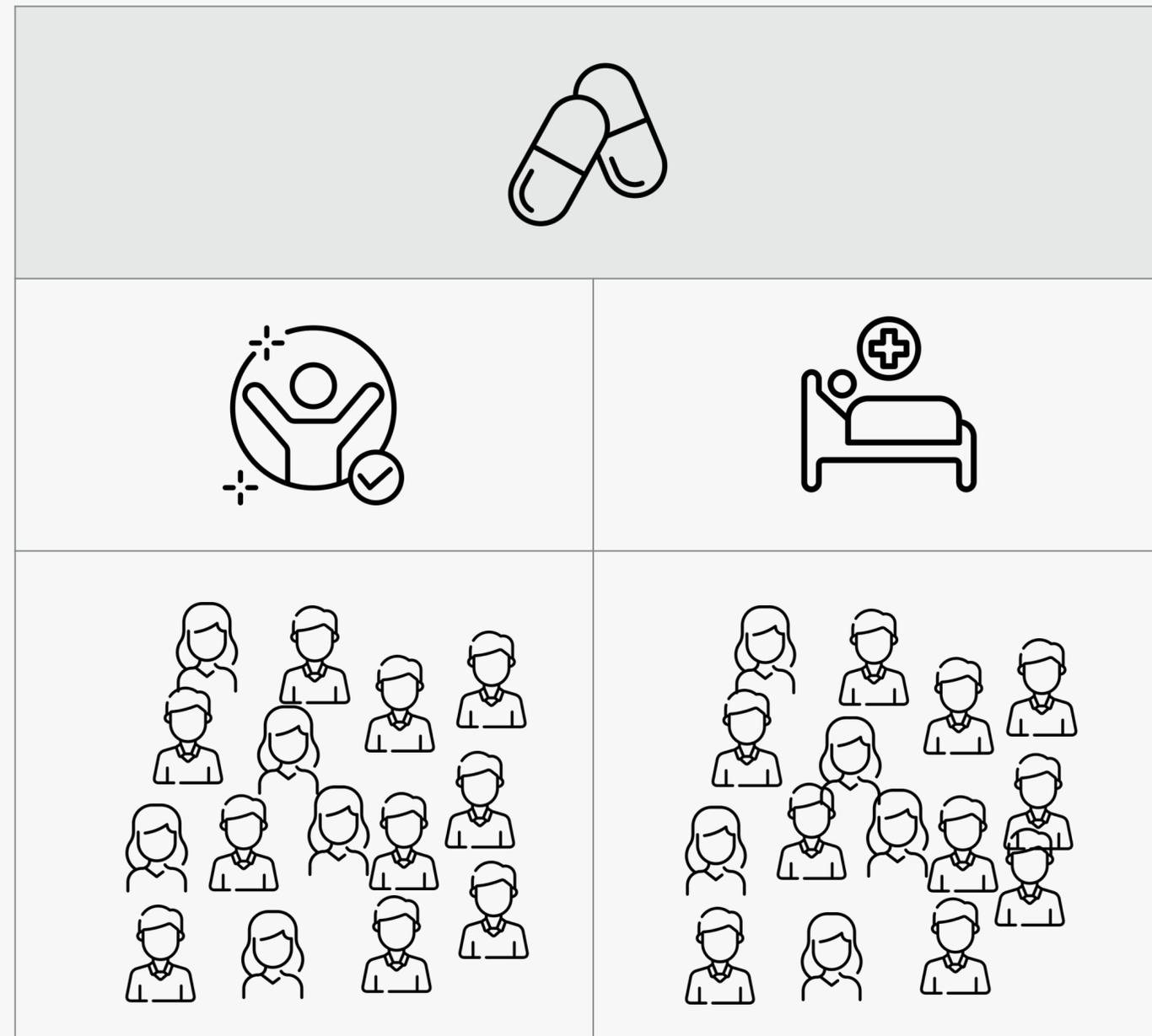
Taking a step back: What is one way to think about unobserved factors?



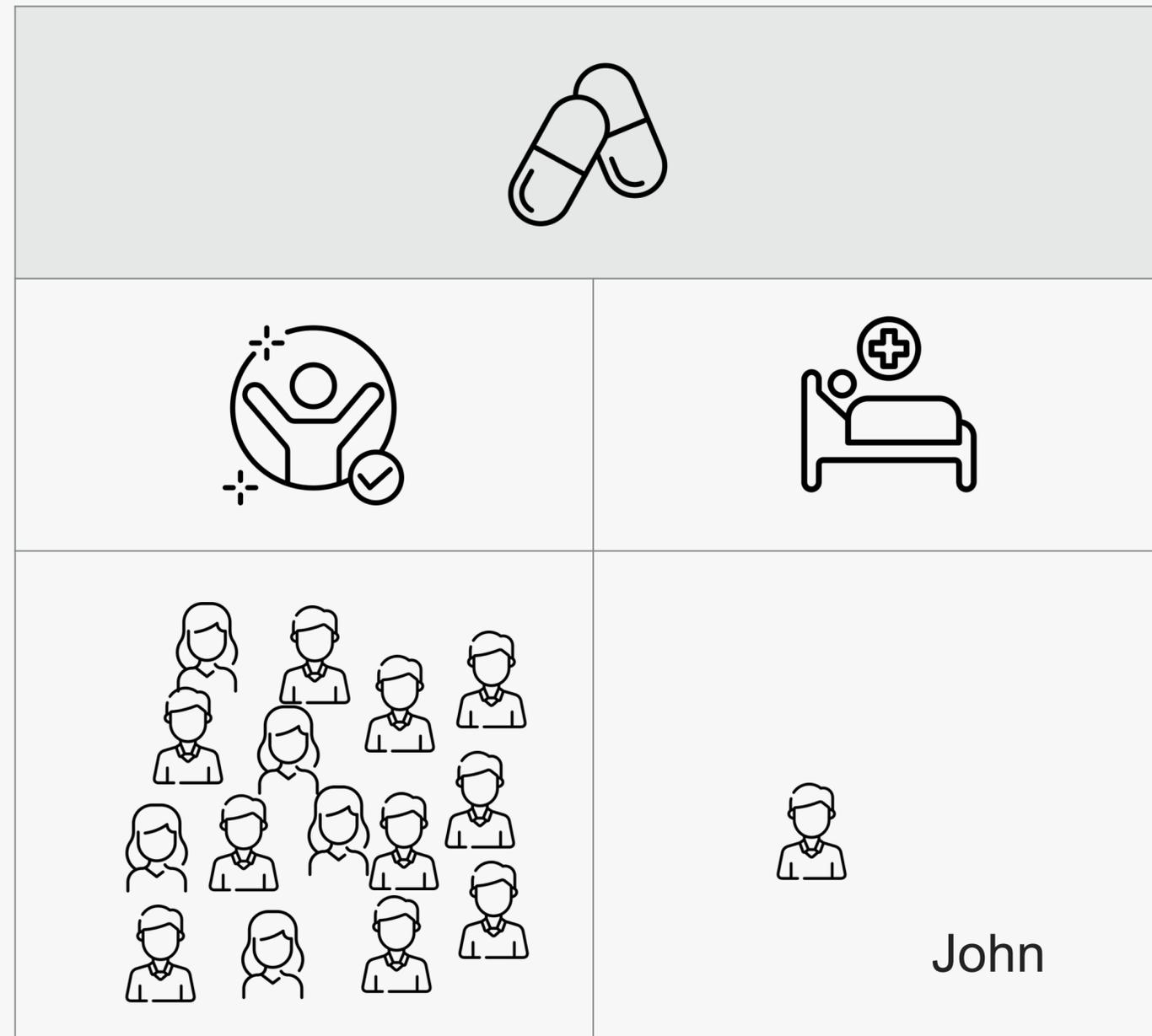
Taking a step back: What is one way to think about unobserved factors?



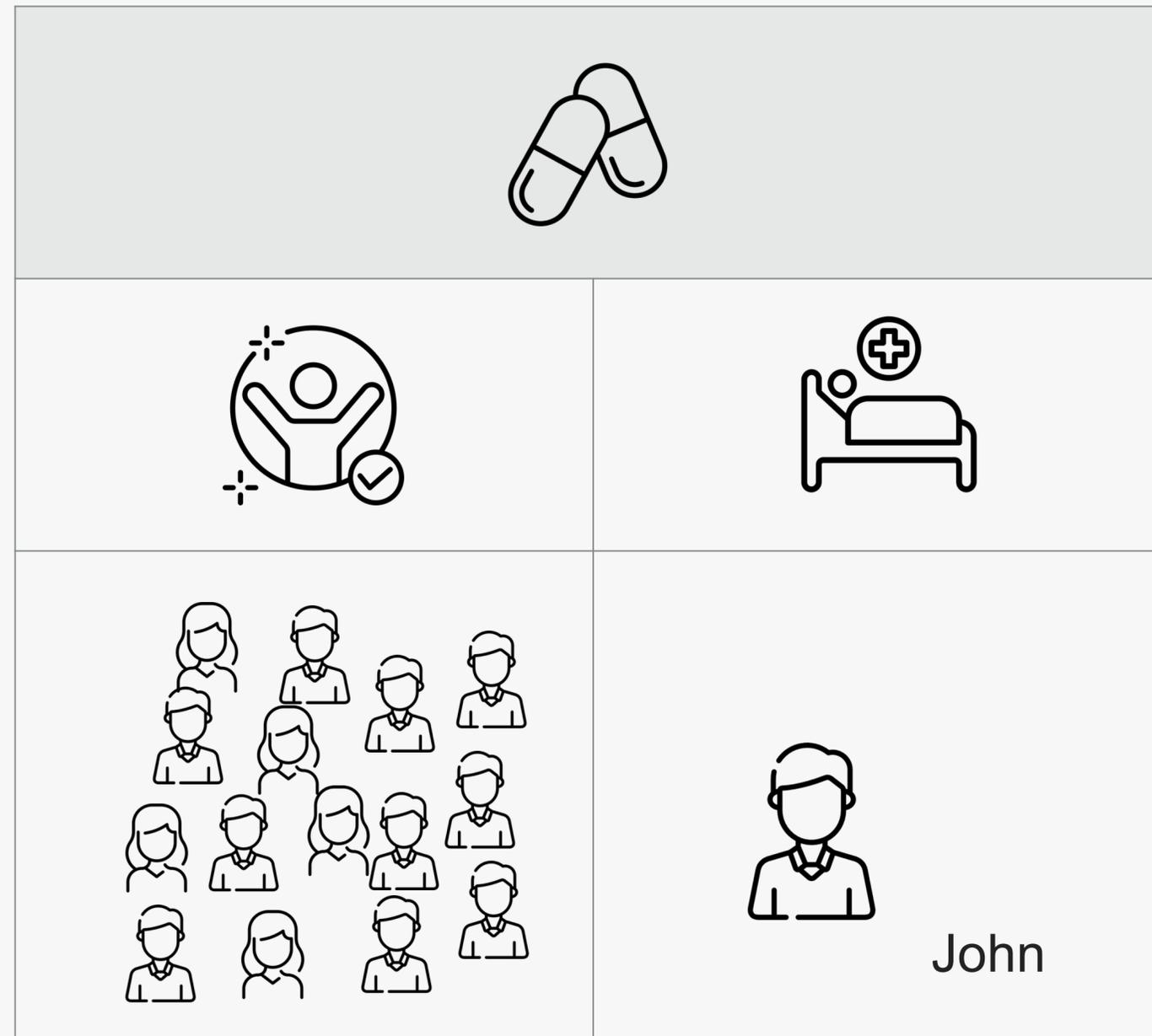
Taking a step back: What is one way to think about unobserved factors?



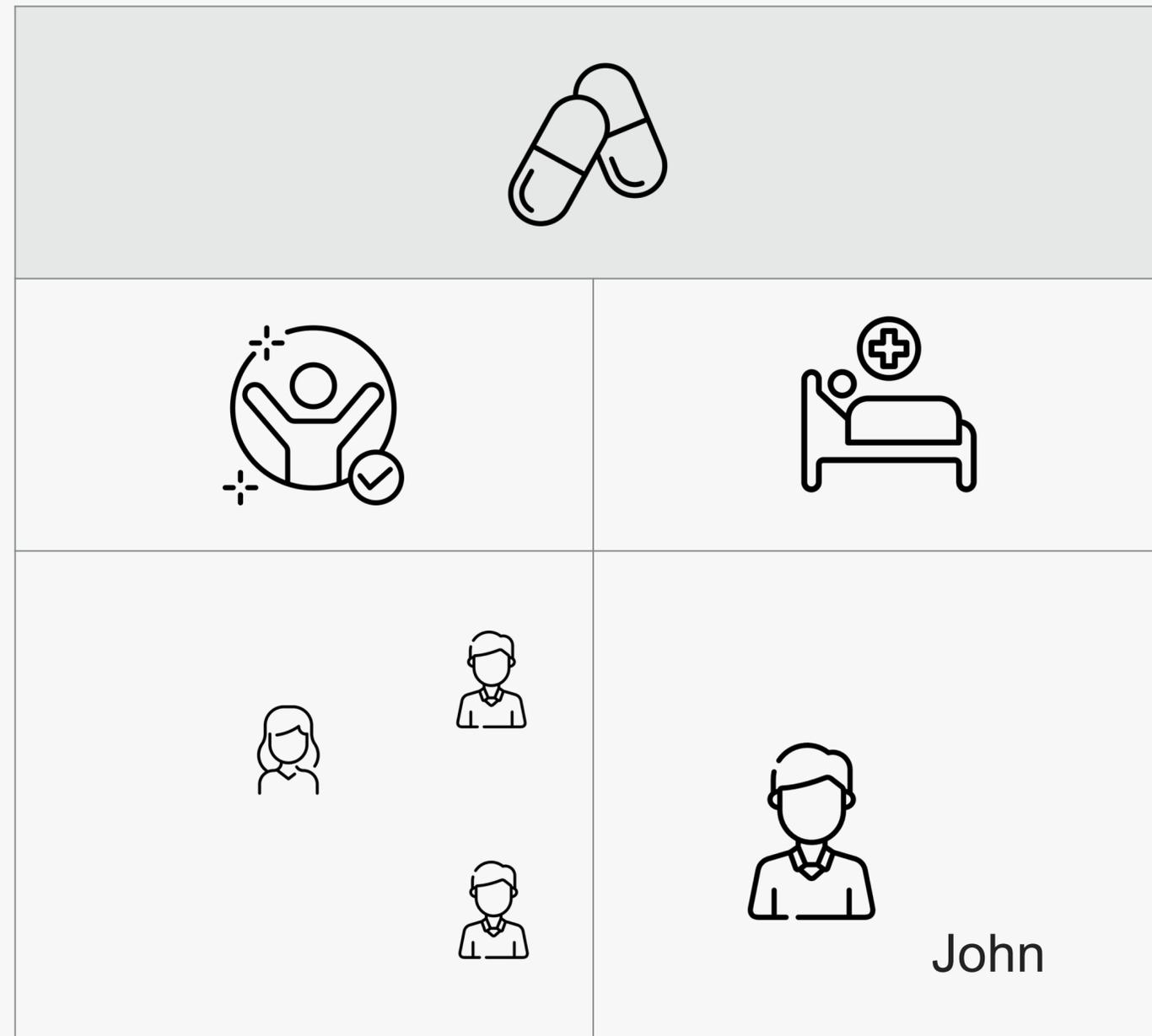
Taking a step back: What is one way to think about unobserved factors?



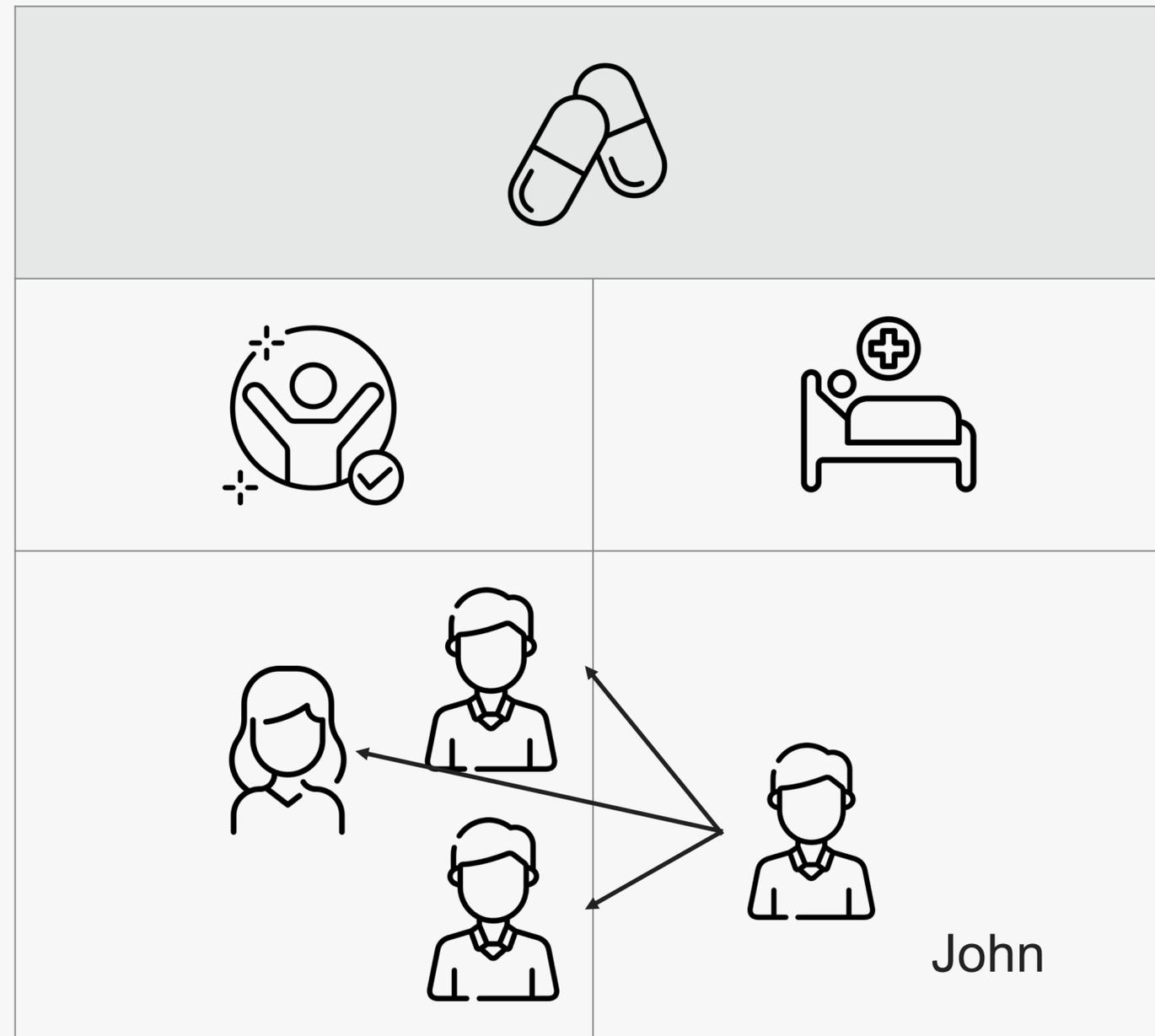
Taking a step back: What is one way to think about unobserved factors?



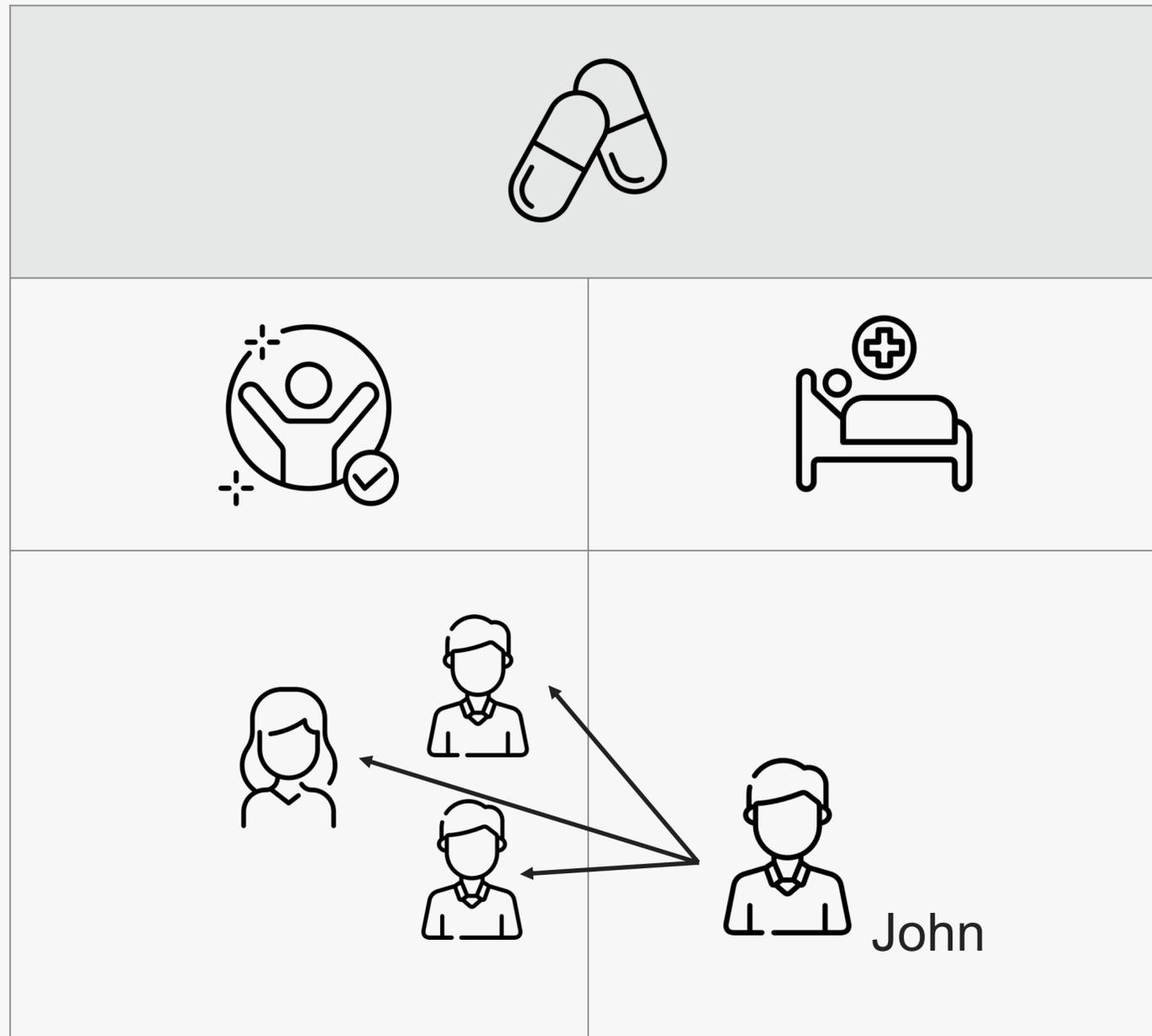
Taking a step back: What is one way to think about unobserved factors?



Taking a step back: What is one way to think about unobserved factors?

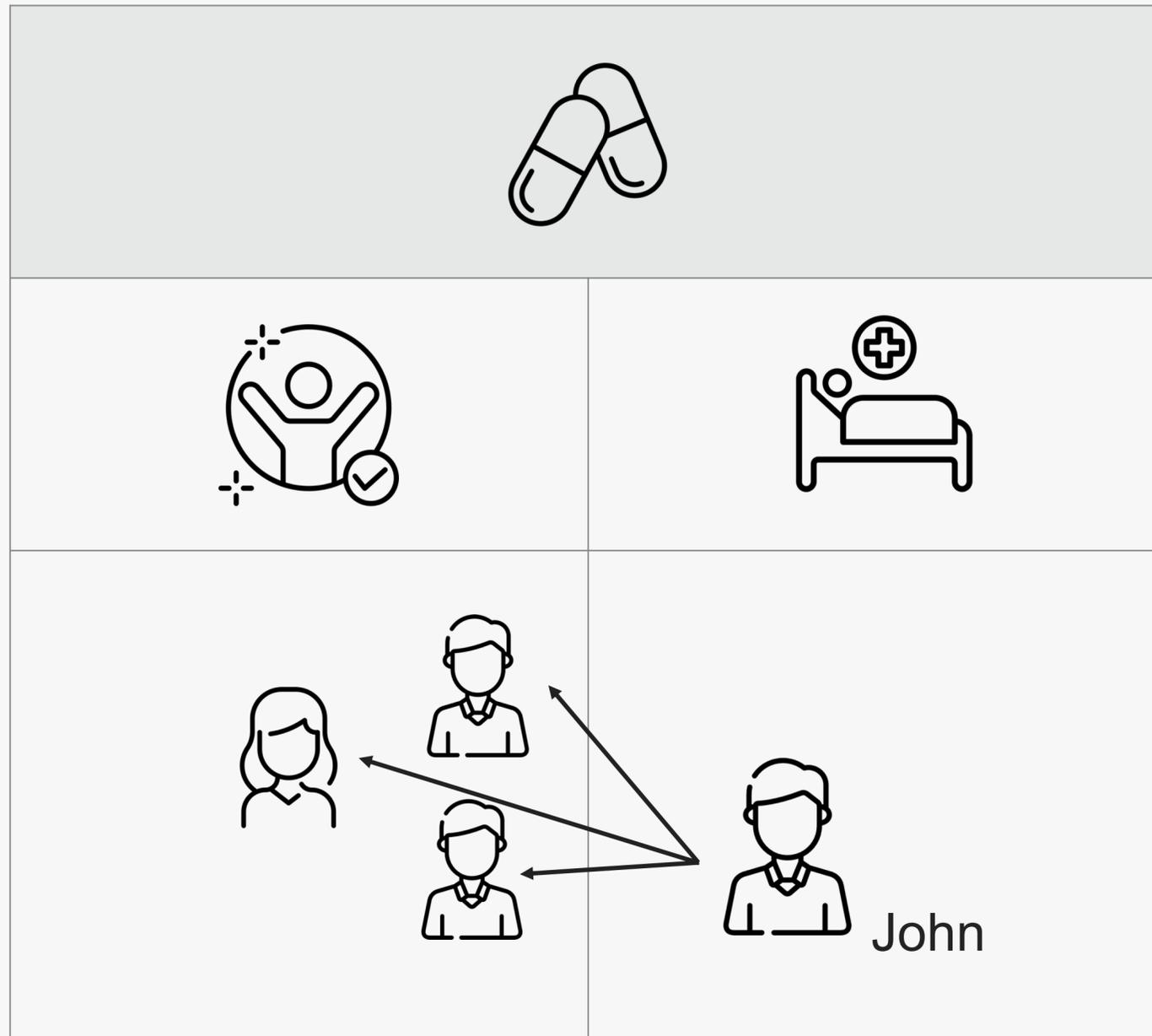


Taking a step back: What is one way to think about unobserved factors?



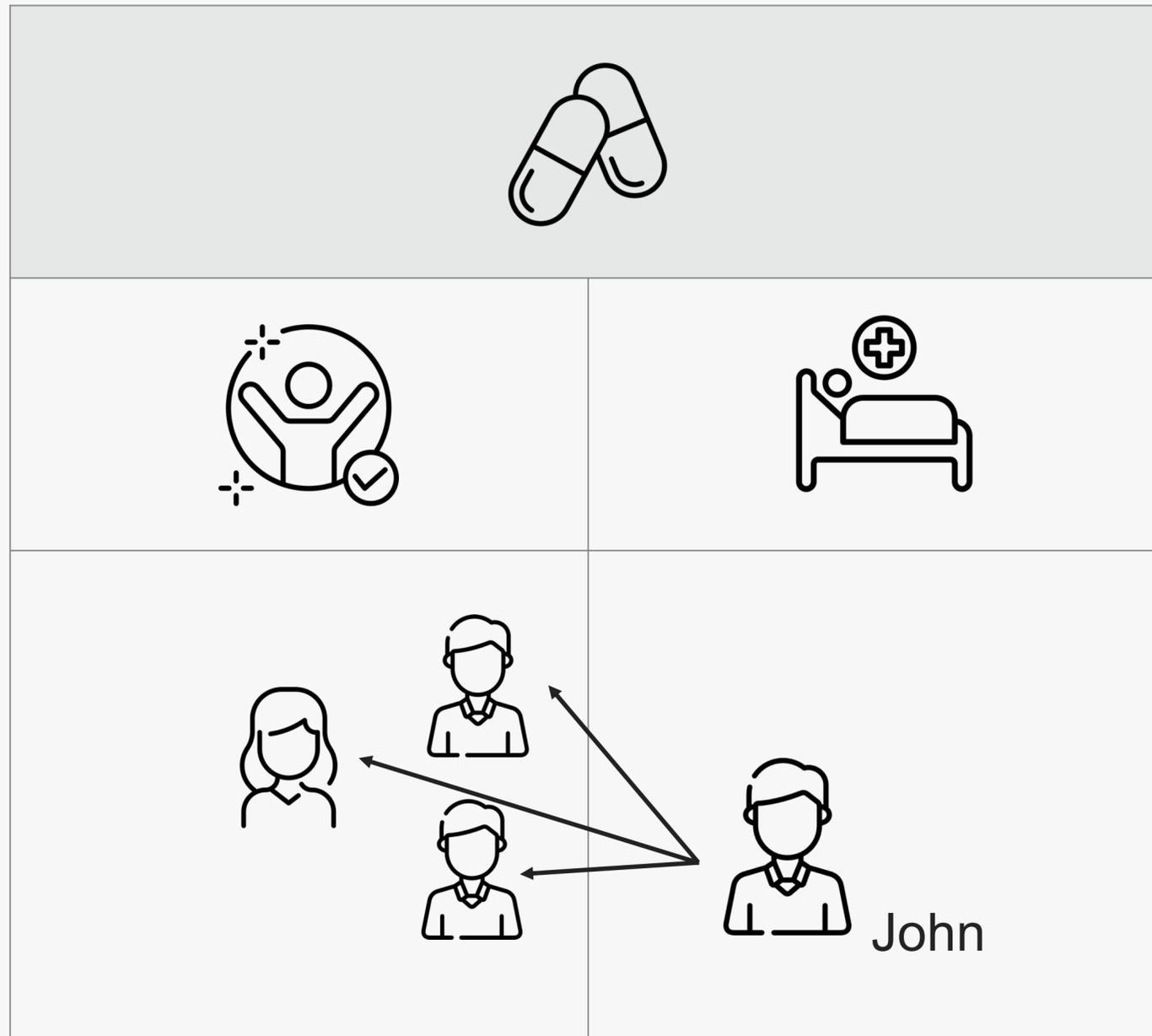
- John's nearest neighbors with a good outcome in the treated cohort are Mary, Peter and Tom
- Similar age
- Got the treatment for approximately the same time-period
- Their tumor is located in the same area

Taking a step back: What is one way to think about unobserved confounders?



- But John had a recurrence, and Mary, Peter and Tom didn't
- There must be a hidden factor that differentiated John from the rest
- We are calculating this unobserved feature for John: \tilde{U}_{John}
- $\tilde{U}_{\text{John}} < 0$ unusual frailty
- $\tilde{U}_{\text{John}} > 0$ unusual robustness
- All widely used balancing approaches can be extended

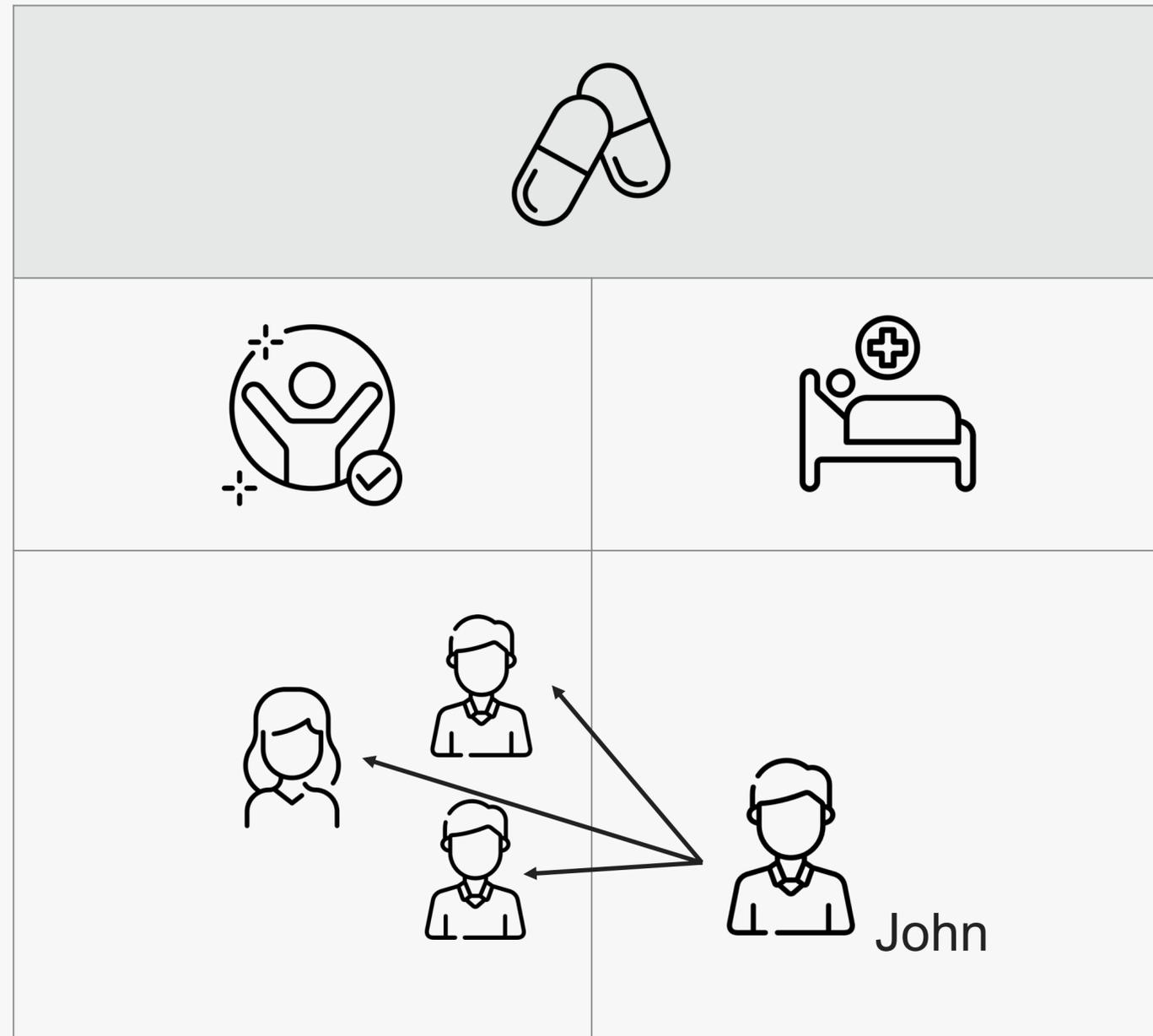
Calculating U



$$Y_{John} = f_t(X_{John}) + g_t(V_{John}) + \epsilon$$

- $f_t(X)$: treatment interaction with observed factors X
- $g_t(V)$: treatment interaction with unobserved factors V
- ϵ : error

Calculating U



$$Y_{John} = f_t(X_{John}) + g_t(V_{John}) + \epsilon$$

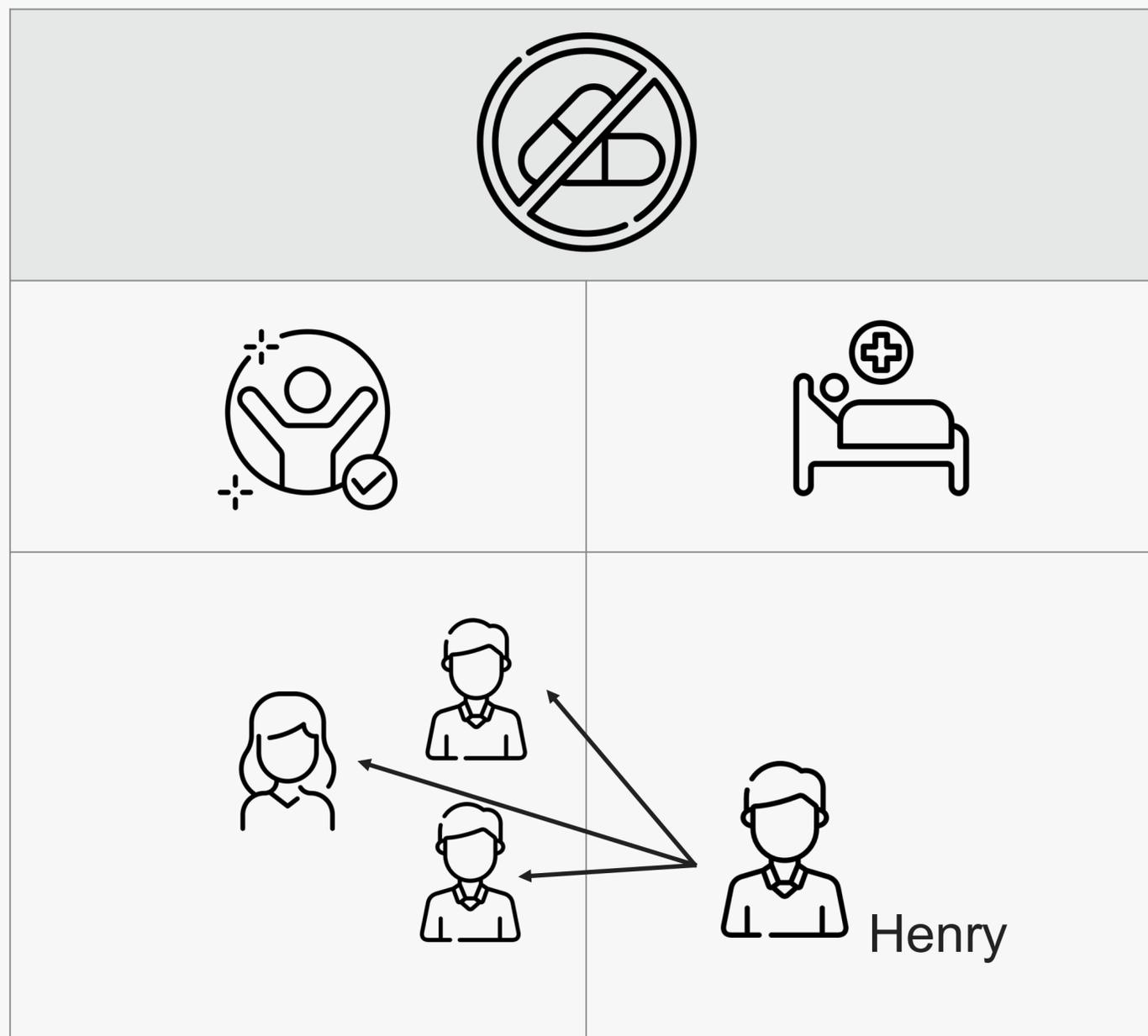
- $f_t(X)$: treatment interaction with observed factors X
- $g_t(V)$: treatment interaction with unobserved factors V
- ϵ : error

$$U_{John} = Y_{John} - \frac{1}{|\mathcal{N}(John)|} \sum_{j \in \mathcal{N}(John)} Y_j$$

$$\approx g_t(V_{John}) - E[g_t(V_{John})]$$

$\mathcal{N}(John)$: John's neighbors with worse outcomes

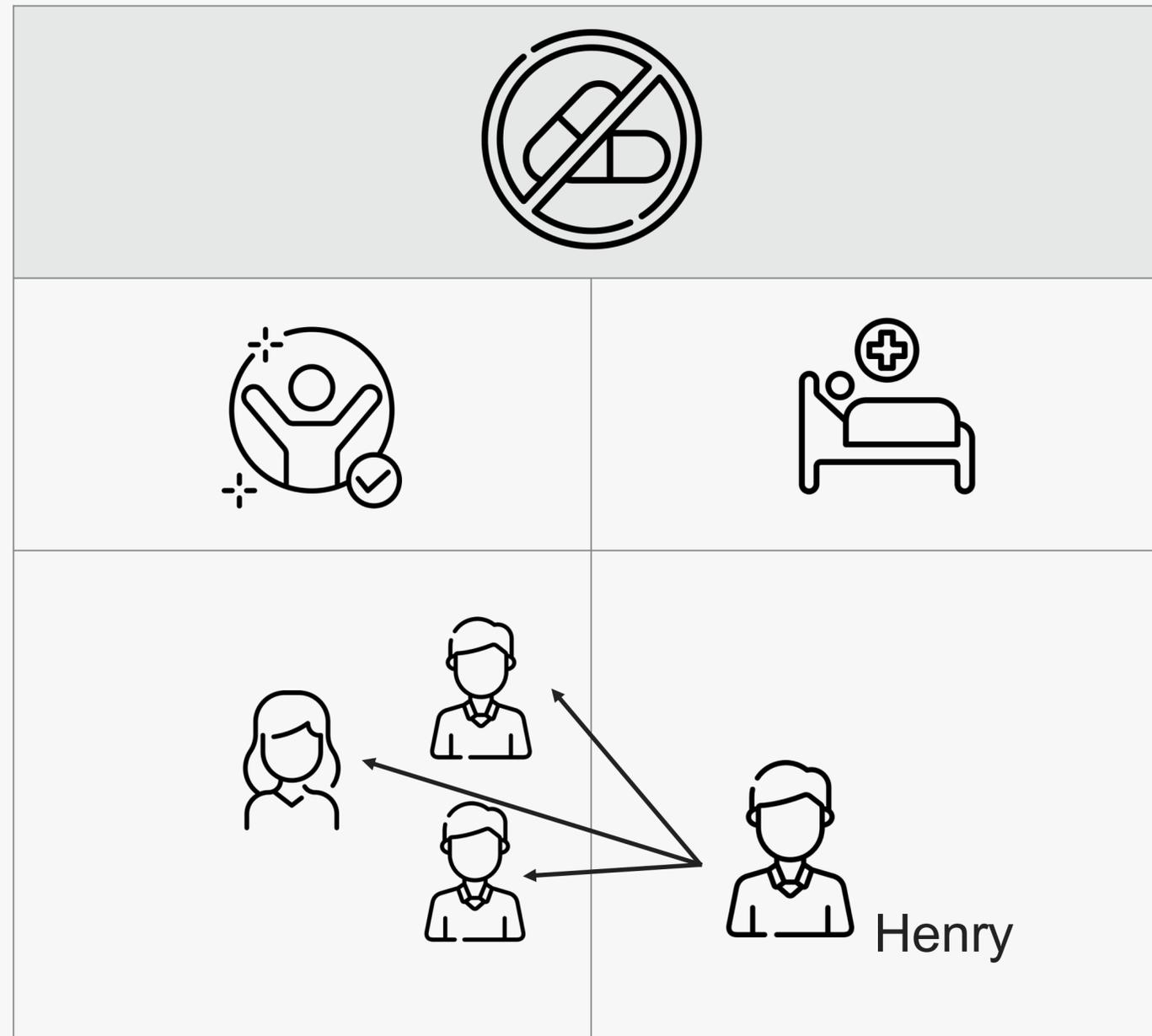
Calculating U



$$U_{Henry} = Y_{Henry} - \frac{1}{|\mathcal{N}(Henry)|} \sum_{j \in \mathcal{N}(Henry)} Y_j$$

$$\approx g_u(V_{Henry}) - E[g_u(V_{Henry})]$$

Calculating U



$$U_{Henry} = Y_{Henry} - \frac{1}{|\mathcal{N}(Henry)|} \sum_{j \in \mathcal{N}(Henry)} Y_j$$

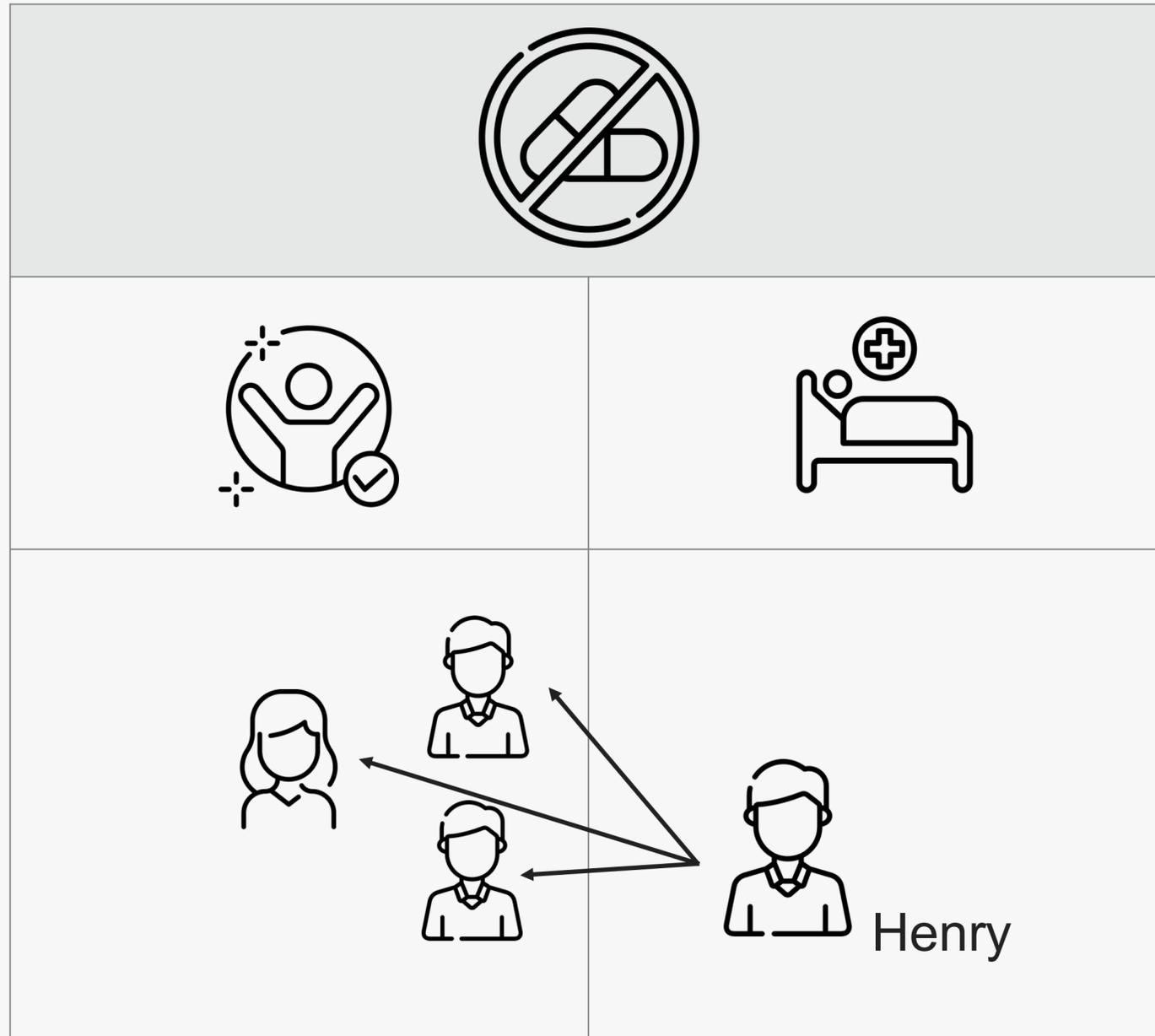
$$\approx g_u(V_{Henry}) - E[g_u(V_{Henry})]$$

- We perform the same approach for all the **untreated** patients too
- We normalize **U** so that

$$0 \leq \tilde{U}_i \leq 1 \quad \forall \text{ patient } i: U_i \geq 0$$

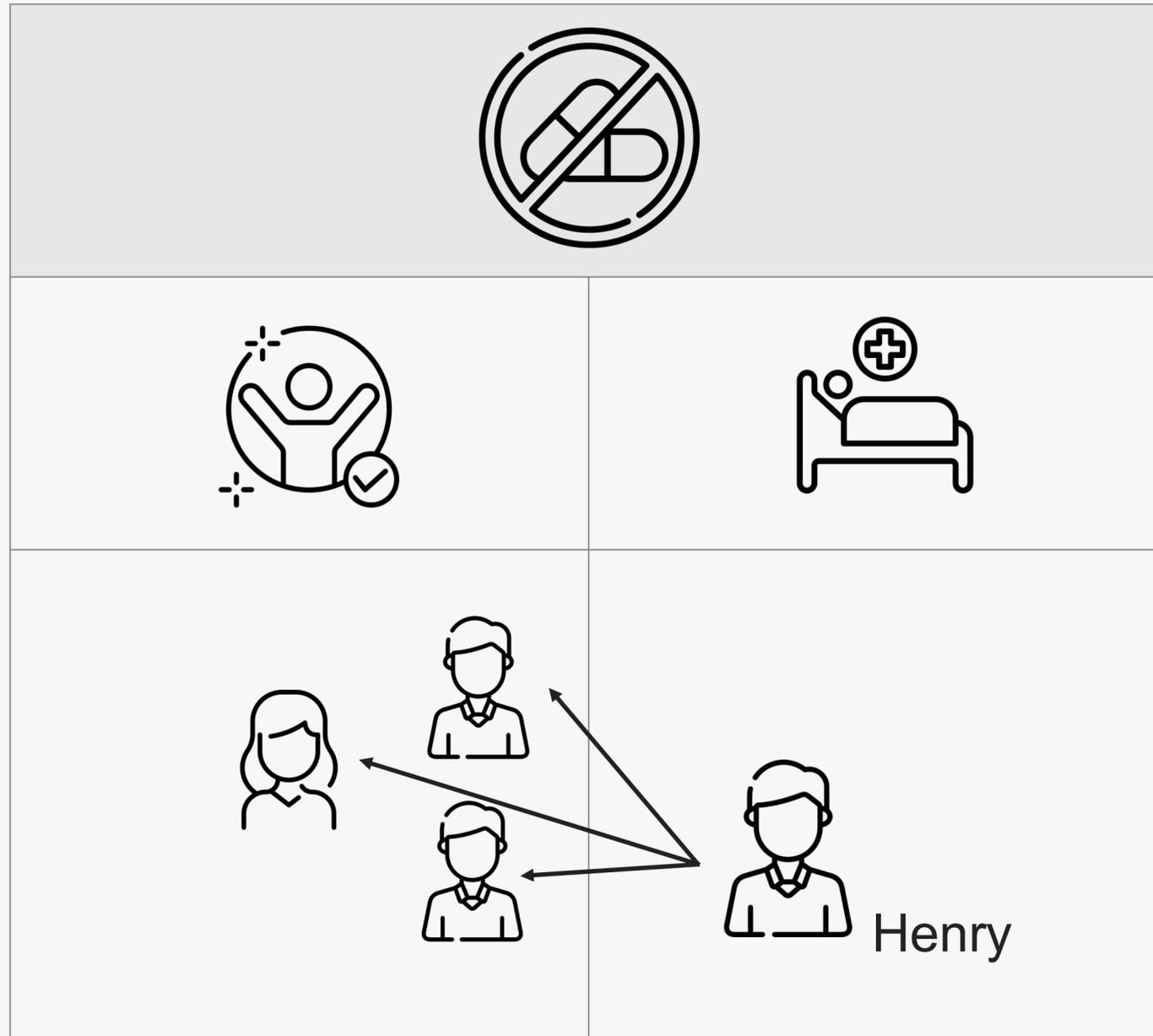
$$-1 \leq \tilde{U}_i \leq 0 \quad \forall \text{ patient } i: U_i \leq 0$$
- A "Top 10% Survivor" in the treatment arm is comparable to a "Top 10% Survivor" in the control arm

Calculating U



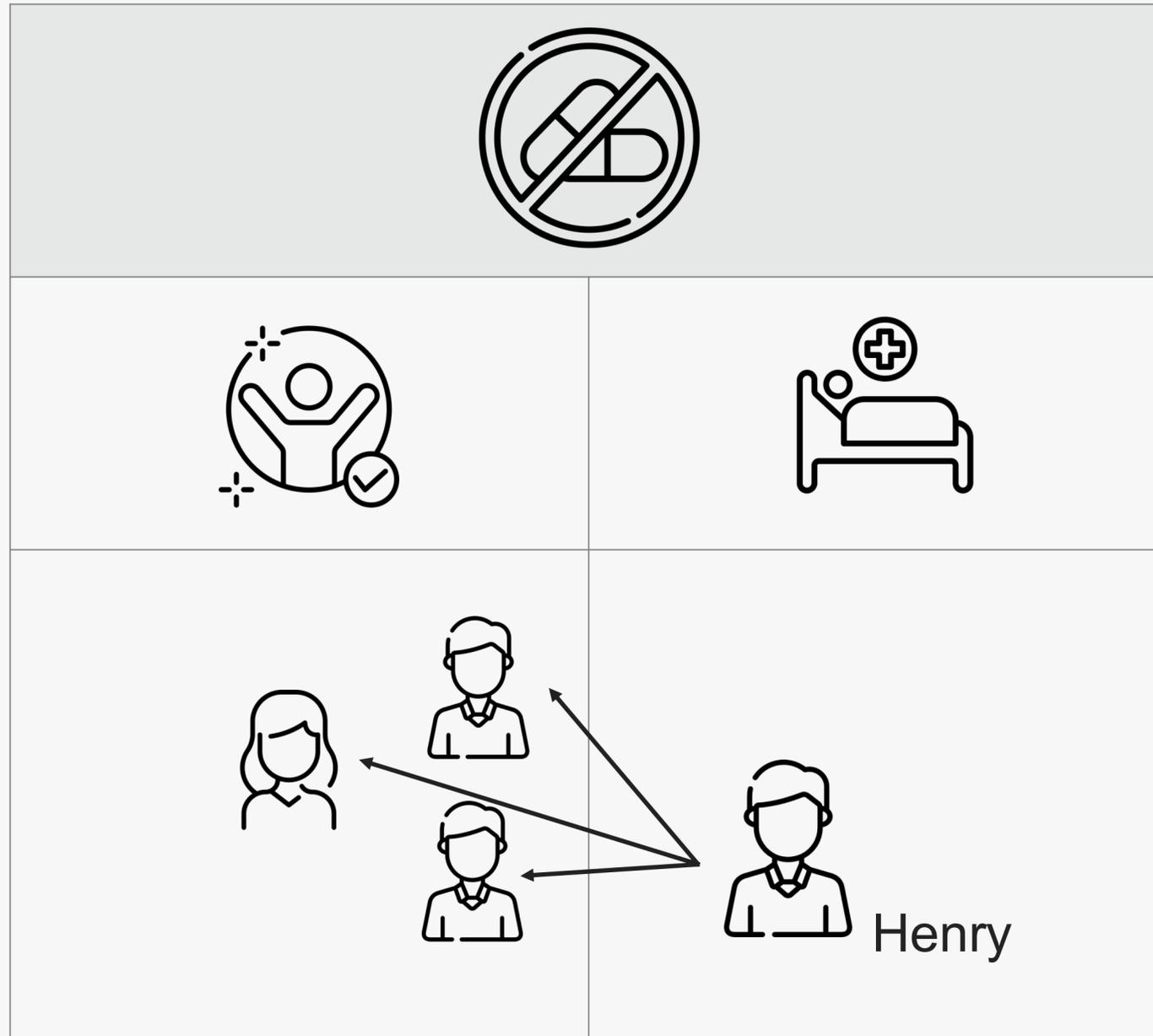
$$U_{Henry} = Y_{Henry} - \frac{1}{|\mathcal{N}(John)|} \sum_{j \in \mathcal{N}(Henry)} Y_j$$

Calculating U



$$U_{Henry} = Y_{Henry} - \frac{1}{|\mathcal{N}(John)|} \sum_{j \in \mathcal{N}(Henry)} Y_j$$

Calculating U



$$U_{Henry} = Y_{Henry} - \frac{1}{|\mathcal{N}(John)|} \sum_{j \in \mathcal{N}(Henry)} Y_j$$

- The Problem:

We can't use raw survival time because of censoring

- The Solution:

Restricted Mean Survival Time (RMST)
– Area under the survival curve up to τ years

- The Calculation:

Jackknife Pseudo-Observations

$$Y_{Henry} = n\hat{\mu} - (n-1)\hat{\mu}_{(-i)}$$

where $\hat{\mu} = \frac{1}{n} \sum_{j=1}^n Y_j$

Balancing treatment arms

Method	Standard Approach (X)	Augmented Approach ($X + U$)
1. Prognostic Matching ¹		

1. Bertsimas, Dimitris, Angelos Georgios Koulouras, and Georgios Antonios Margonis. "The ROAD to precision medicine." *npj Digital Medicine* 7.1 (2024): 307.

Balancing treatment arms

Method	Standard Approach (X)	Augmented Approach ($X + U$)
1. Prognostic Matching¹	Matches patients based on (X) within predicted risk buckets: $m(X)$	

1. Bertsimas, Dimitris, Angelos Georgios Koulouras, and Georgios Antonios Margonis. "The ROAD to precision medicine." *npj Digital Medicine* 7.1 (2024): 307.

Balancing treatment arms

Method	Standard Approach (X)	Augmented Approach ($X + U$)
1. Prognostic Matching¹	Matches patients based on (X) within predicted risk buckets: $m(X)$	Matches on risk & calibrates weights to align residual means: $E[U]_{control} \approx E[U]_{treat}$

1. Bertsimas, Dimitris, Angelos Georgios Koulouras, and Georgios Antonios Margonis. "The ROAD to precision medicine." *npj Digital Medicine* 7.1 (2024): 307.

Balancing treatment arms

Method	Standard Approach (X)	Augmented Approach ($X + U$)
1. Prognostic Matching ¹	Matches patients based on (X) within predicted risk buckets: $m(X)$	Matches on risk & calibrates weights to align residual means: $E[U]_{control} \approx E[U]_{treat}$
2. Entropy Balancing ²		

1. Bertsimas, Dimitris, Angelos Georgios Koulouras, and Georgios Antonios Margonis. "The ROAD to precision medicine." *npj Digital Medicine* 7.1 (2024): 307.

2. Hainmueller, Jens. "Entropy balancing for causal effects: A multivariate reweighting method to produce balanced samples in observational studies." *Political analysis* 20.1 (2012): 25-46.

Balancing treatment arms

Method	Standard Approach (X)	Augmented Approach ($X + U$)
1. Prognostic Matching¹	Matches patients based on (X) within predicted risk buckets: $m(X)$	Matches on risk & calibrates weights to align residual means: $E[U]_{control} \approx E[U]_{treat}$
2. Entropy Balancing²	Optimizes weights to force observed moments to match: $\sum w_i X_i = \bar{X}_{treat}$	

1. Bertsimas, Dimitris, Angelos Georgios Koulouras, and Georgios Antonios Margonis. "The ROAD to precision medicine." *npj Digital Medicine* 7.1 (2024): 307.
2. Hainmueller, Jens. "Entropy balancing for causal effects: A multivariate reweighting method to produce balanced samples in observational studies." *Political analysis* 20.1 (2012): 25-46.

Balancing treatment arms

Method	Standard Approach (X)	Augmented Approach ($X + U$)
1. Prognostic Matching¹	Matches patients based on (X) within predicted risk buckets: $m(X)$	Matches on risk & calibrates weights to align residual means: $E[U]_{control} \approx E[U]_{treat}$
2. Entropy Balancing²	Optimizes weights to force observed moments to match: $\sum w_i X_i = \bar{X}_{treat}$	Adds constraint to force residual moments to match: $\sum w_i [X_i, U_i] = [\bar{X}, \bar{U}]_{treat}$

1. Bertsimas, Dimitris, Angelos Georgios Koulouras, and Georgios Antonios Margonis. "The ROAD to precision medicine." *npj Digital Medicine* 7.1 (2024): 307.

2. Hainmueller, Jens. "Entropy balancing for causal effects: A multivariate reweighting method to produce balanced samples in observational studies." *Political analysis* 20.1 (2012): 25-46.

Balancing treatment arms

Method	Standard Approach (X)	Augmented Approach ($X + U$)
1. Prognostic Matching¹	Matches patients based on (X) within predicted risk buckets: $m(X)$	Matches on risk & calibrates weights to align residual means: $E[U]_{control} \approx E[U]_{treat}$
2. Entropy Balancing²	Optimizes weights to force observed moments to match: $\sum w_i X_i = \bar{X}_{treat}$	Adds constraint to force residual moments to match: $\sum w_i [X_i, U_i] = [\bar{X}, \bar{U}]_{treat}$
3. IPTW³		

1. Bertsimas, Dimitris, Angelos Georgios Koulouras, and Georgios Antonios Margonis. "The ROAD to precision medicine." *npj Digital Medicine* 7.1 (2024): 307.
2. Hainmueller, Jens. "Entropy balancing for causal effects: A multivariate reweighting method to produce balanced samples in observational studies." *Political analysis* 20.1 (2012): 25-46.
3. Robins, James M., Miguel Angel Hernan, and Babette Brumback. "Marginal structural models and causal inference in epidemiology." *Epidemiology* 11.5 (2000): 550-560.

Balancing treatment arms

Method	Standard Approach (X)	Augmented Approach ($X + U$)
1. Prognostic Matching¹	Matches patients based on (X) within predicted risk buckets: $m(X)$	Matches on risk & calibrates weights to align residual means: $E[U]_{control} \approx E[U]_{treat}$
2. Entropy Balancing²	Optimizes weights to force observed moments to match: $\sum w_i X_i = \bar{X}_{treat}$	Adds constraint to force residual moments to match: $\sum w_i [X_i, U_i] = [\bar{X}, \bar{U}]_{treat}$
3. IPTW³	Weights by propensity of observed history: $\pi = P(T X)$	

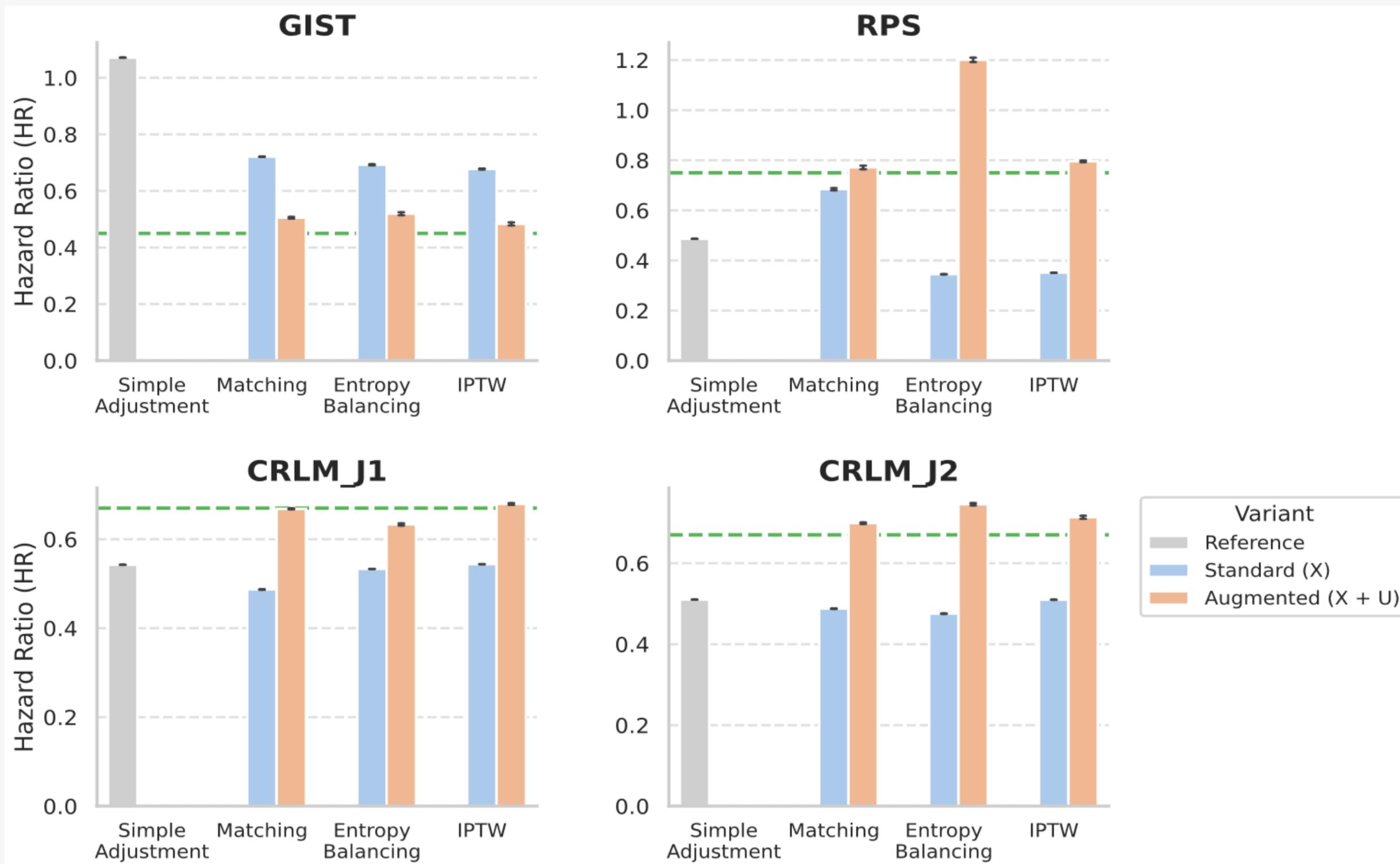
1. Bertsimas, Dimitris, Angelos Georgios Koulouras, and Georgios Antonios Margonis. "The ROAD to precision medicine." *npj Digital Medicine* 7.1 (2024): 307.
2. Hainmueller, Jens. "Entropy balancing for causal effects: A multivariate reweighting method to produce balanced samples in observational studies." *Political analysis* 20.1 (2012): 25-46.
3. Robins, James M., Miguel Angel Hernan, and Babette Brumback. "Marginal structural models and causal inference in epidemiology." *Epidemiology* 11.5 (2000): 550-560.

Balancing treatment arms

Method	Standard Approach (X)	Augmented Approach ($X + U$)
1. Prognostic Matching¹	Matches patients based on (X) within predicted risk buckets: $m(X)$	Matches on risk & calibrates weights to align residual means: $E[U]_{control} \approx E[U]_{treat}$
2. Entropy Balancing²	Optimizes weights to force observed moments to match: $\sum w_i X_i = \bar{X}_{treat}$	Adds constraint to force residual moments to match: $\sum w_i [X_i, U_i] = [\bar{X}, \bar{U}]_{treat}$
3. IPTW³	Weights by propensity of observed history: $\pi = P(T X)$	Weights by propensity of history & unobserved performance: $\pi = P(T X, U)$

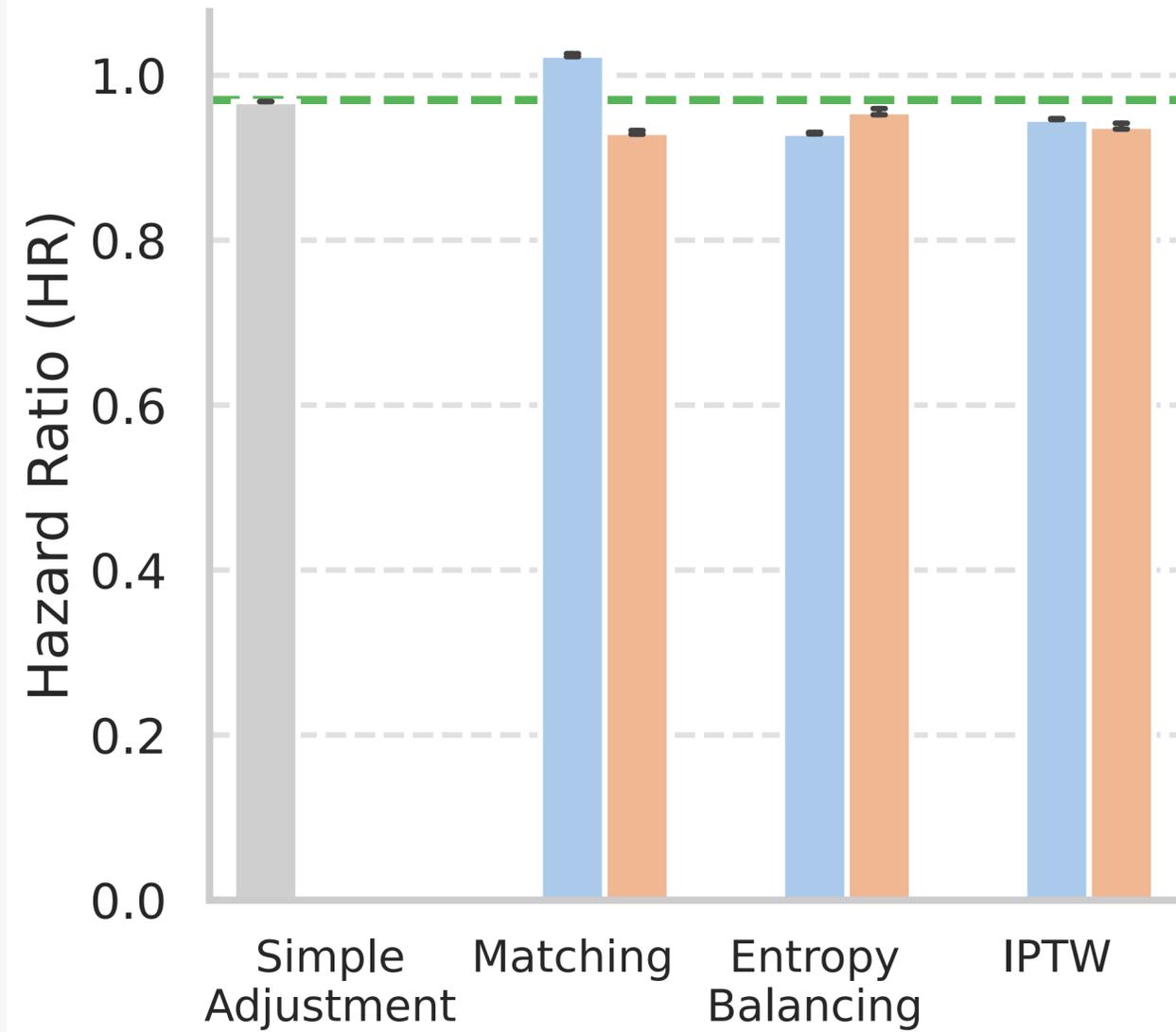
1. Bertsimas, Dimitris, Angelos Georgios Koulouras, and Georgios Antonios Margonis. "The ROAD to precision medicine." *npj Digital Medicine* 7.1 (2024): 307.
2. Hainmueller, Jens. "Entropy balancing for causal effects: A multivariate reweighting method to produce balanced samples in observational studies." *Political analysis* 20.1 (2012): 25-46.
3. Robins, James M., Miguel Angel Hernan, and Babette Brumback. "Marginal structural models and causal inference in epidemiology." *Epidemiology* 11.5 (2000): 550-560.

Observational data

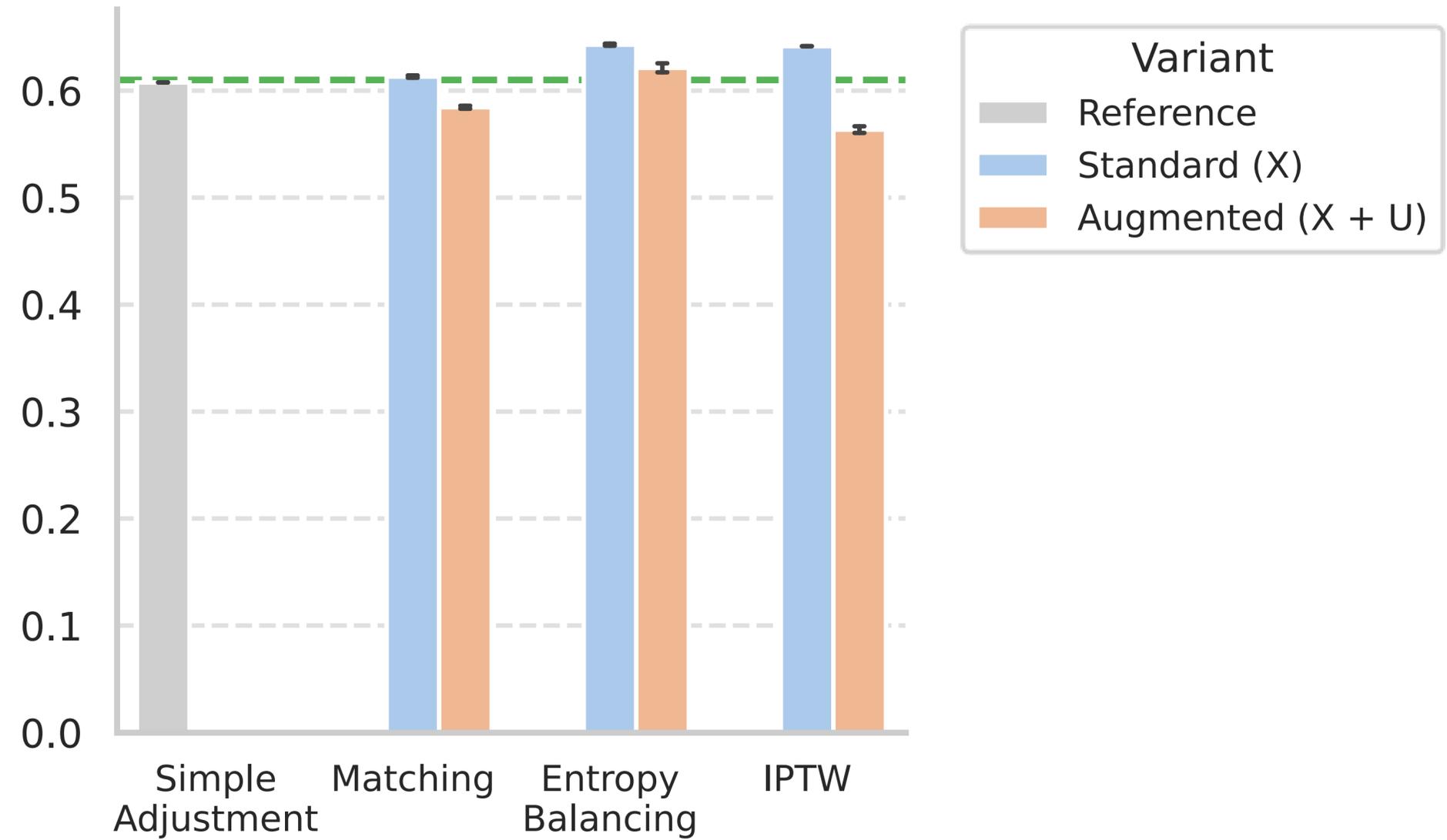


RCT data

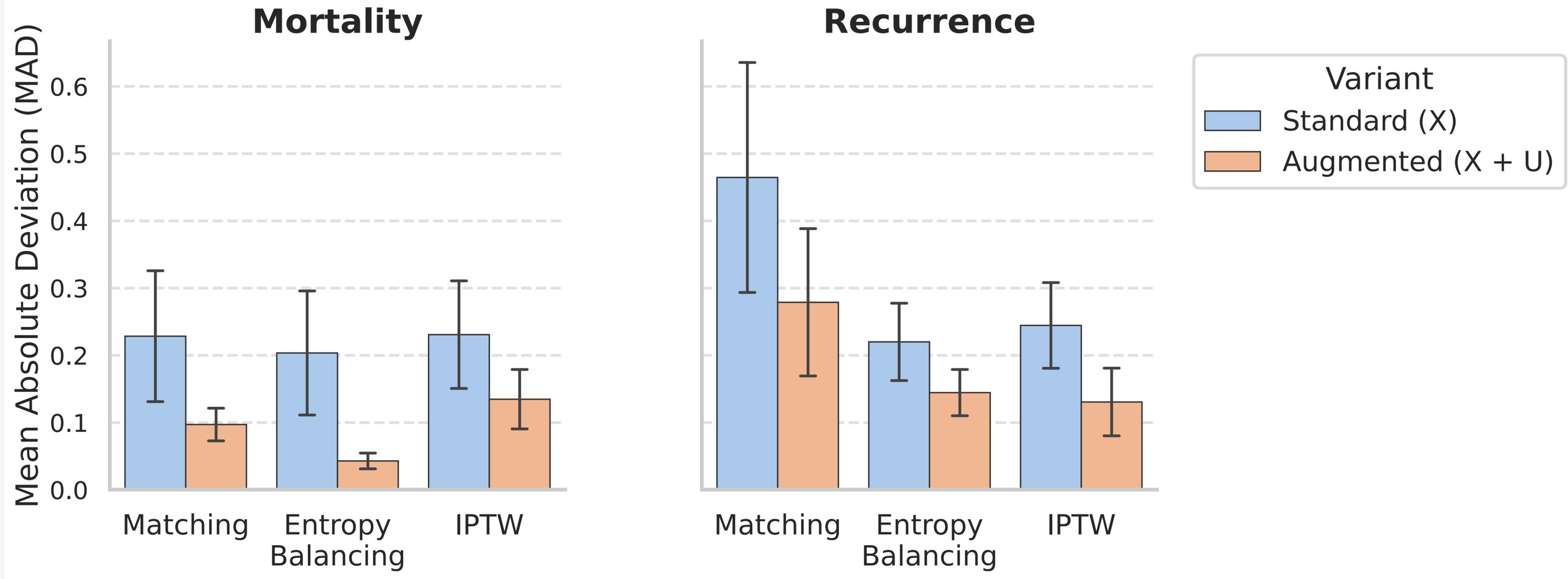
STRASS



BRT



Consistency across 6 centers



Main contributions

Introducing a way to **quantify** the **unobserved prognostic factors** (hidden information that affects the final outcomes)

Main contributions

Introducing a way to **quantify** the **unobserved prognostic factors** (hidden information that affects the final outcomes)

Easy adjustment using state of the art balancing approaches

Main contributions

Introducing a way to **quantify** the **unobserved prognostic factors** (hidden information that affects the final outcomes)

Easy adjustment using state of the art balancing approaches

Strong results compared to traditional techniques that **ignore** unobserved confounding



Appendix



Hazard Ratio

For example, imagine that X_1 is a treatment variable, with values $X_1 = 1$ for treatment and $X_1 = 0$ for control.

The hazard at time t for treatment:

$$\begin{aligned} h(t|X_1 = 1) &= h_0(t)\exp(b_1 * 1) \\ &= h_0(t)\exp(b_1) \end{aligned}$$

and for control:

$$\begin{aligned} h(t|X_1 = 0) &= h_0(t)\exp(b_1 * 0) \\ &= h_0(t)\exp(0) \\ &= h_0(t) \end{aligned}$$

We can compare the hazards for treatment and control at time t as a *hazard ratio (HR)*:

$$\begin{aligned} HR &= \frac{h(t|X_1 = 1)}{h(t|X_1 = 0)} \\ &= \frac{h_0(t)\exp(b_1)}{h_0(t)} \\ &= \exp(b_1) \end{aligned}$$

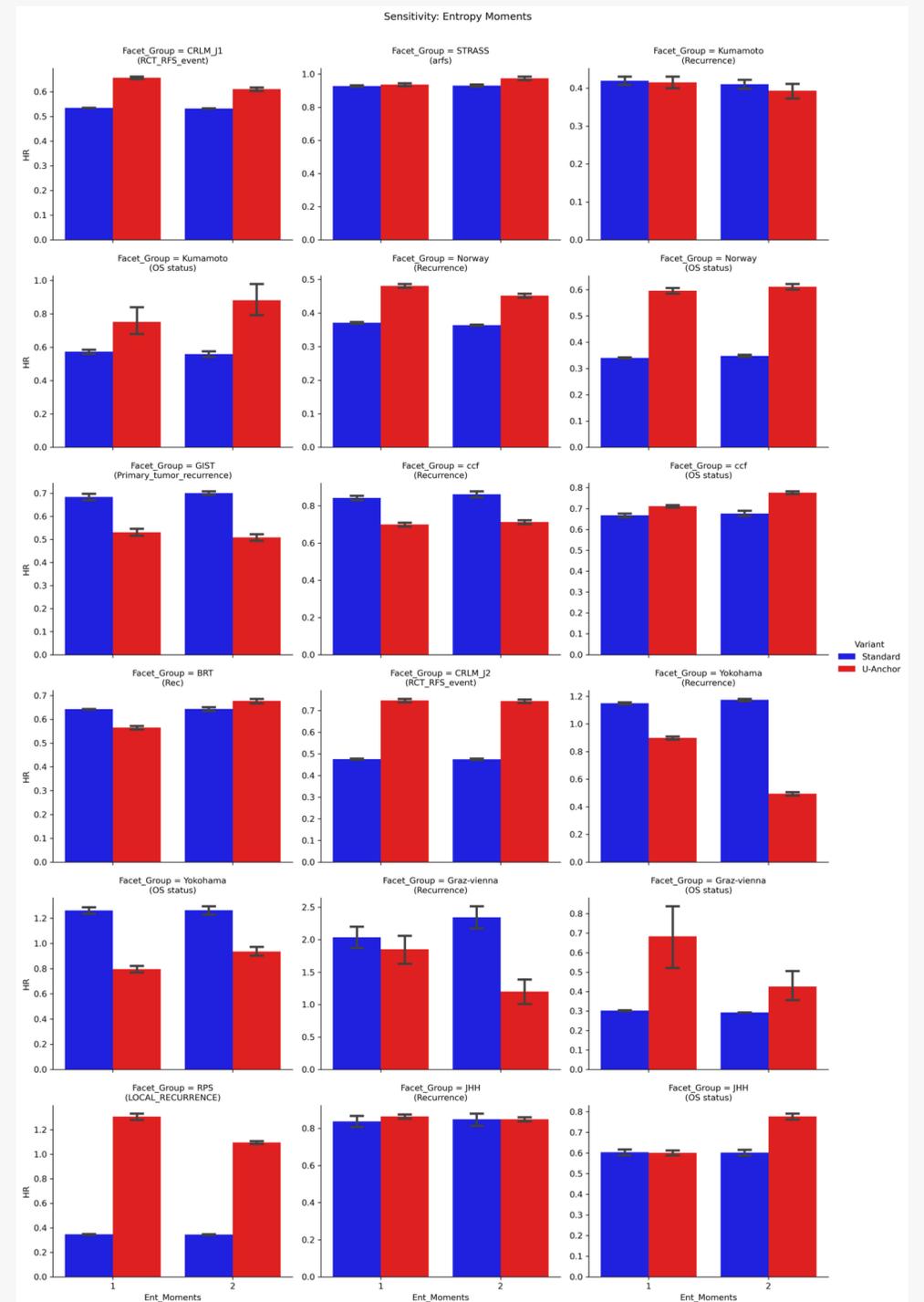
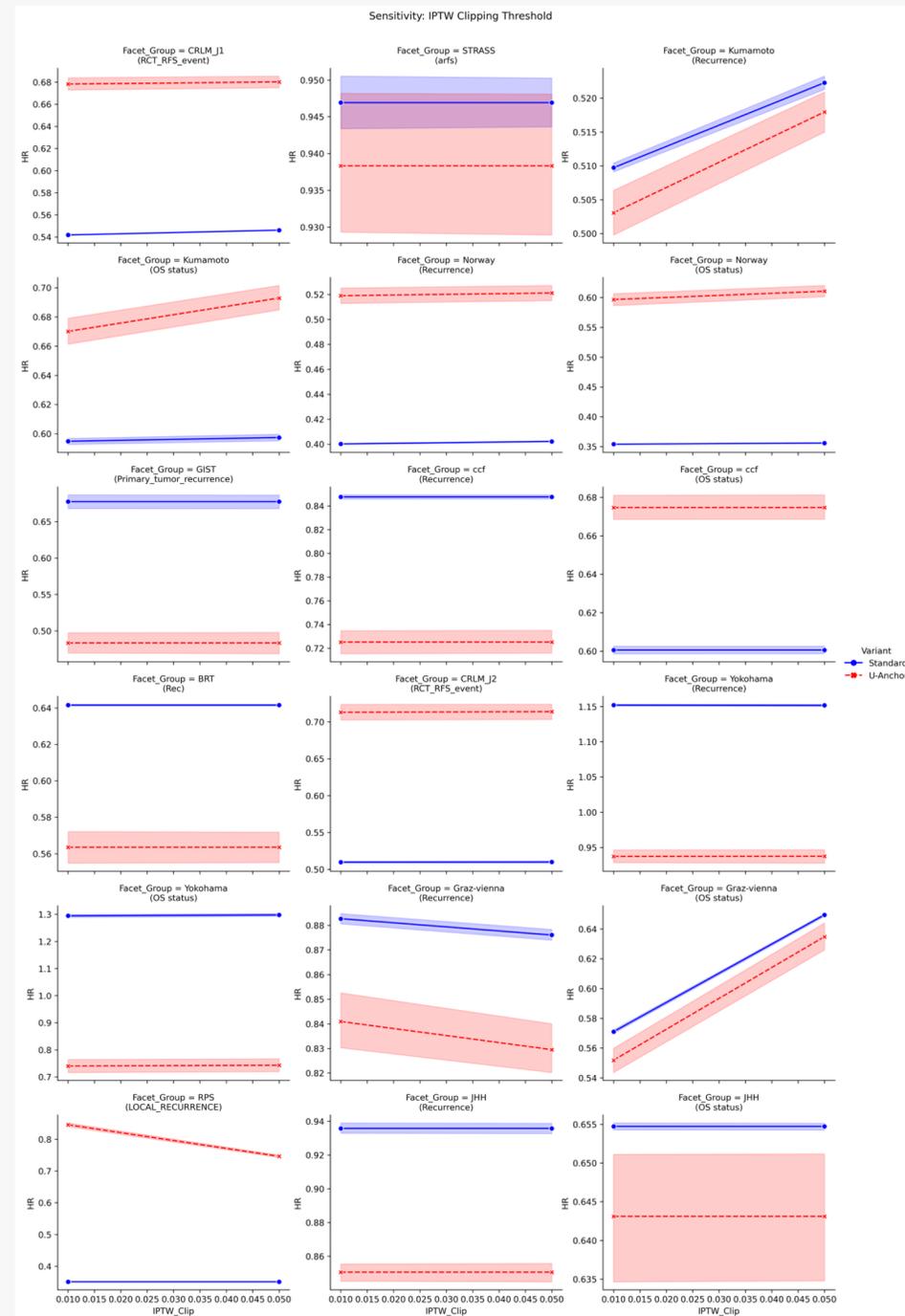
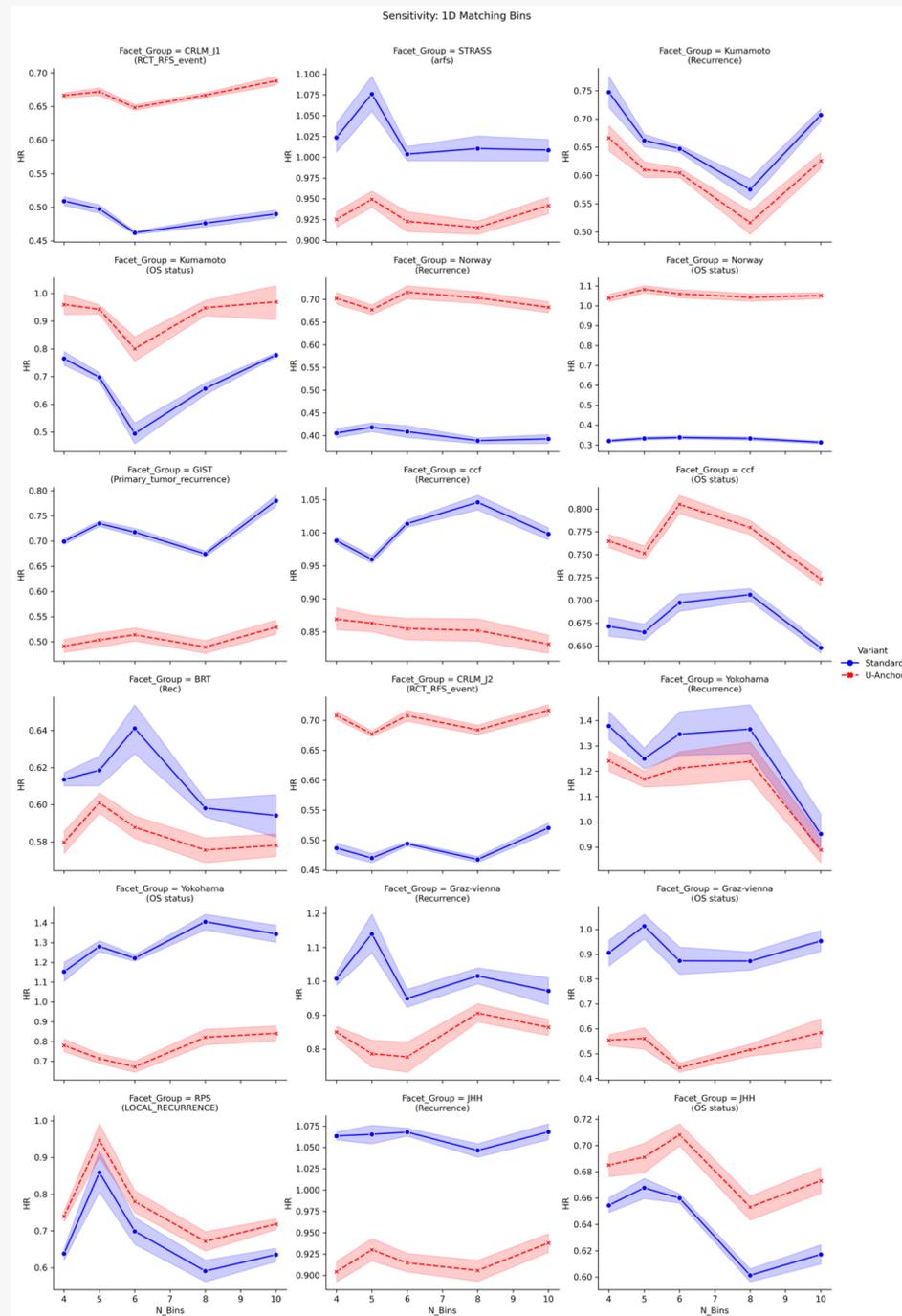
Thus, $\exp(b_1)$ is the hazard ratio comparing the hazard for treatment to controls.

- $HR = .25$ means that treatment has $\frac{1}{4}$ (25%) the hazard of control, or a 75% decrease. With a lower hazard rate, treatment will have fewer expected events and thus better survival.
- $HR = 2$ here means that treatment has twice the hazard of control, or a 100% increase, and thus worse survival.

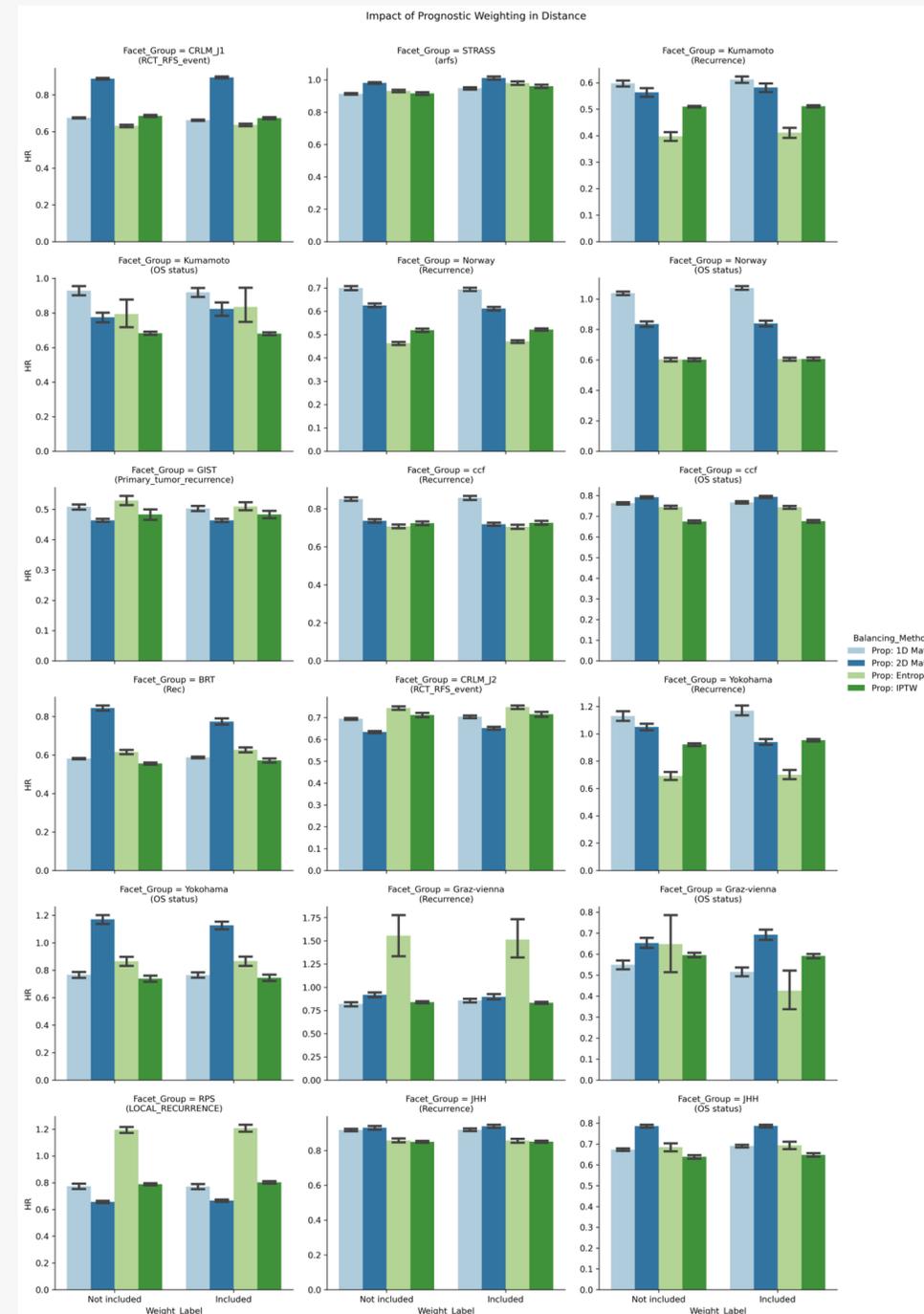
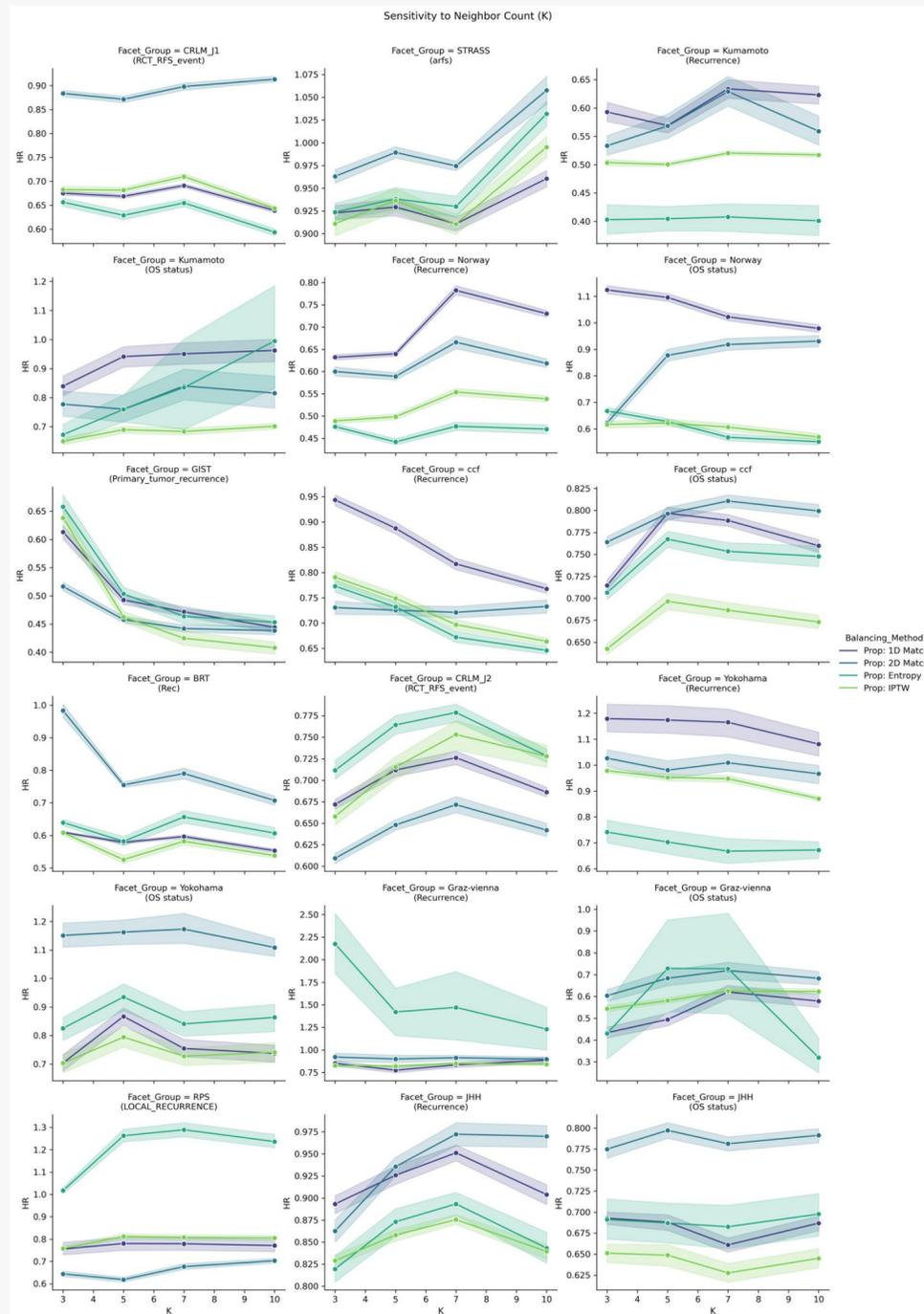
In general, $\exp(b_1)$ expresses the hazard ratio for a 1-unit increase in the associated covariate.

b_1 itself is the log-hazard ratio.

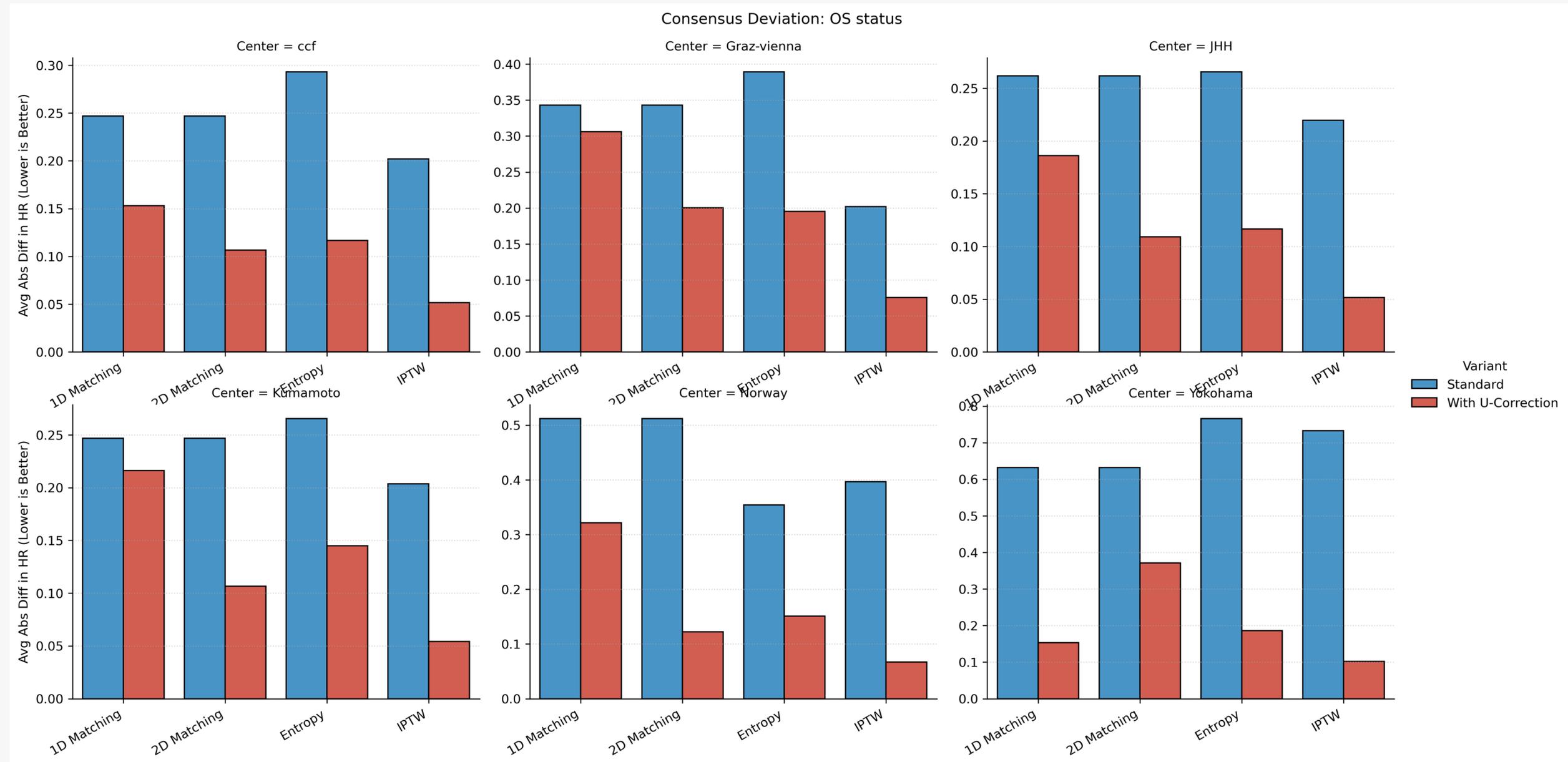
Sensitivity Analysis



Sensitivity Analysis



Pairwise Plots



Pairwise Plots

