

Uncertainty Quantification for Generative and Collaborative AI

Sima Noorani

University of Pennsylvania

Nov 20th, 2025

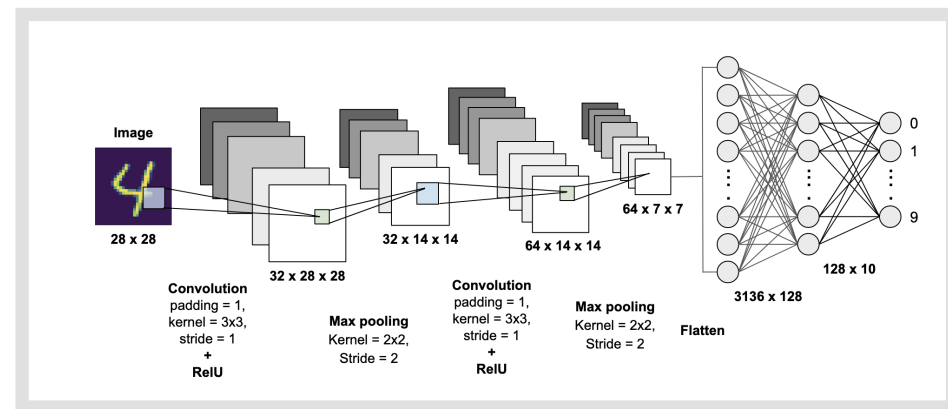
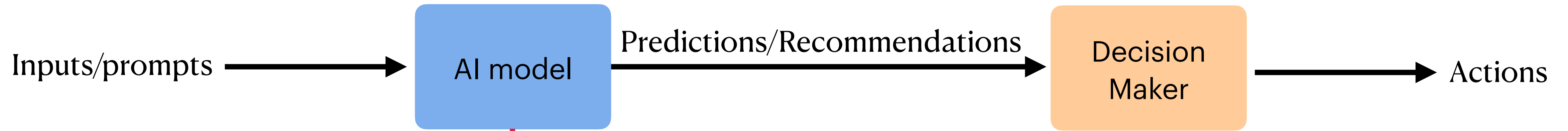
Joint work with: Shayan Kiyani, George Pappas, Hamed Hassani



AI-powered Decision making pipeline

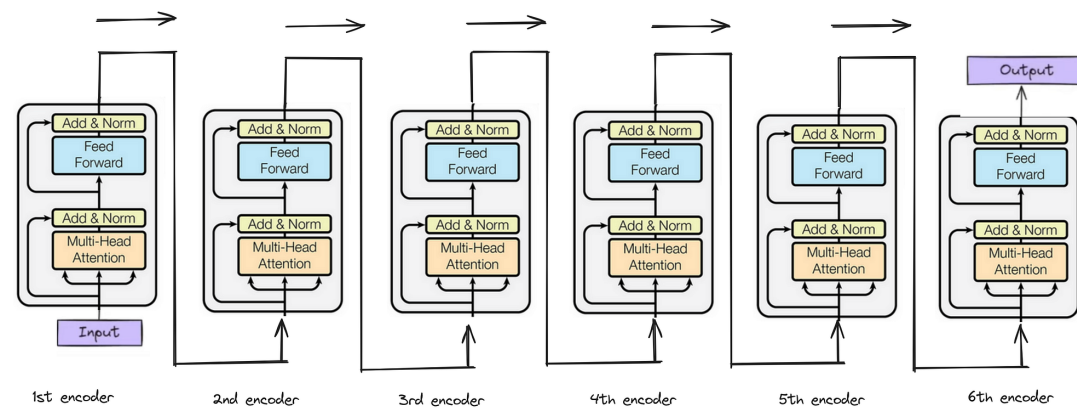


AI-powered Decision making pipeline

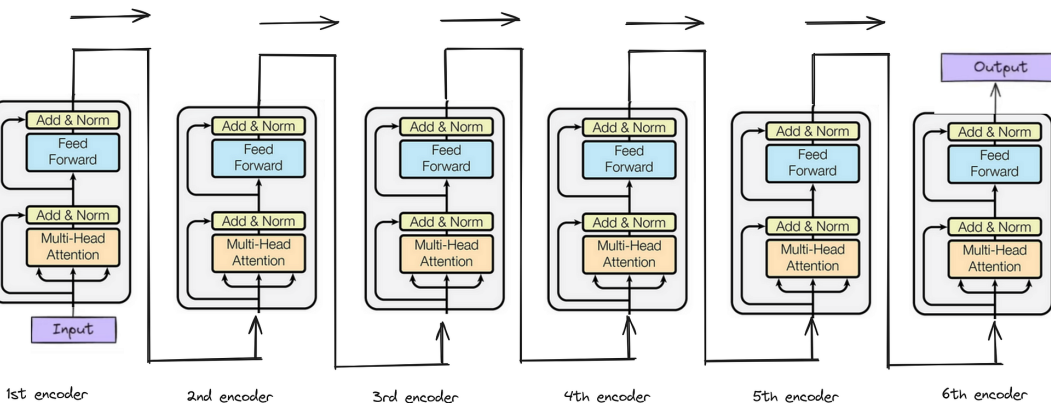
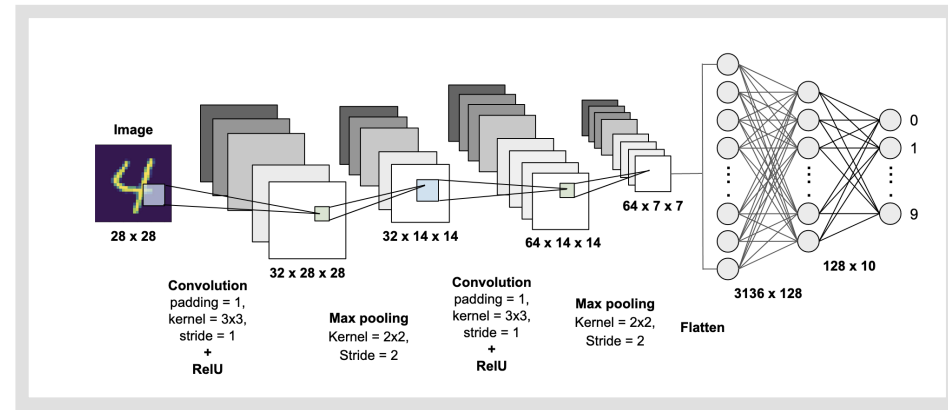


- (Heavily) Pre-trained

- Possibility of little tweaks using extra data (fine-tuning)



AI-powered Decision making pipeline

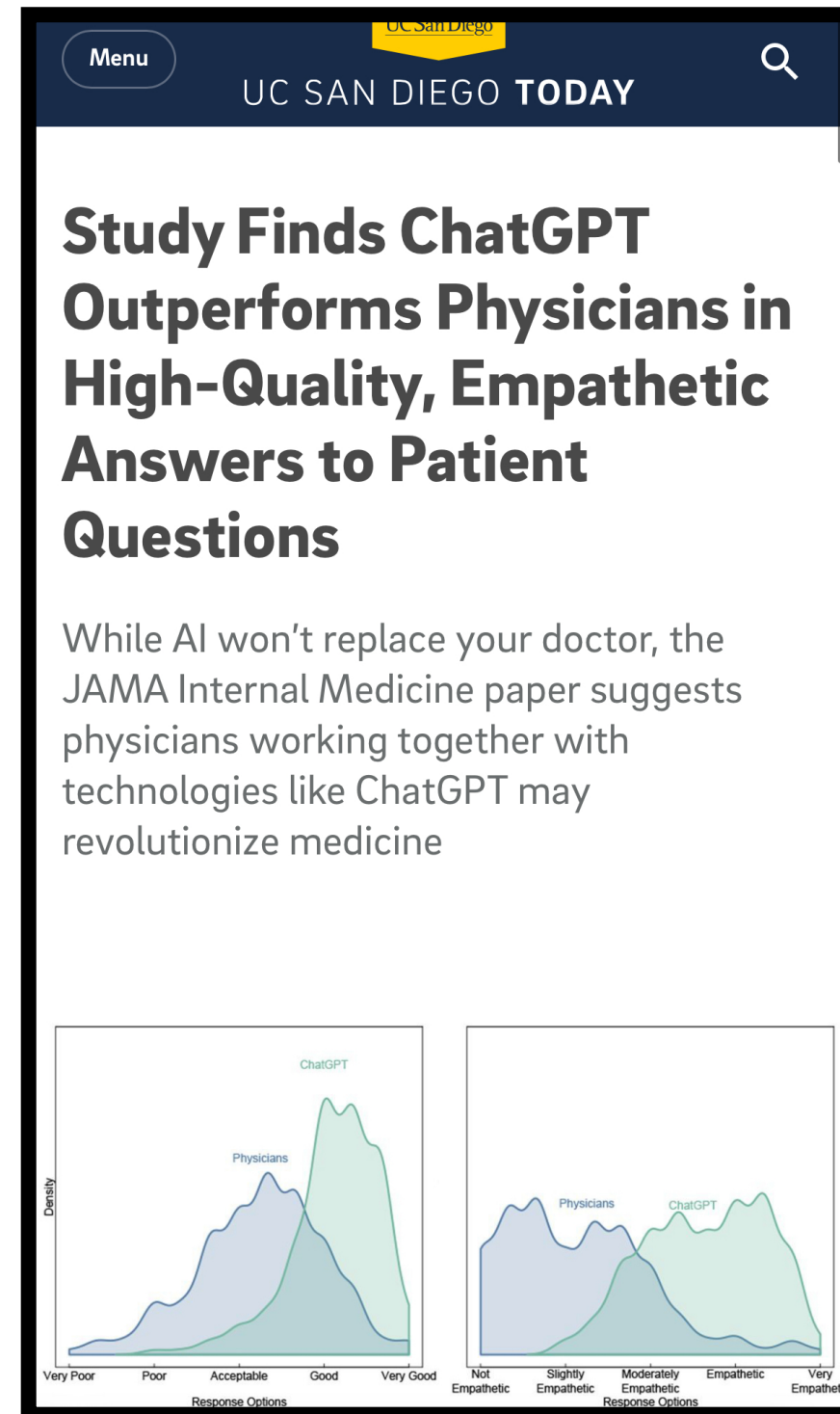
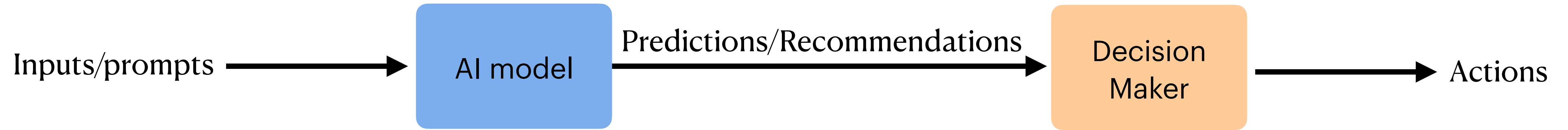


- (Heavily) Pre-trained
- Possibility of little tweaks using extra data (fine-tuning)

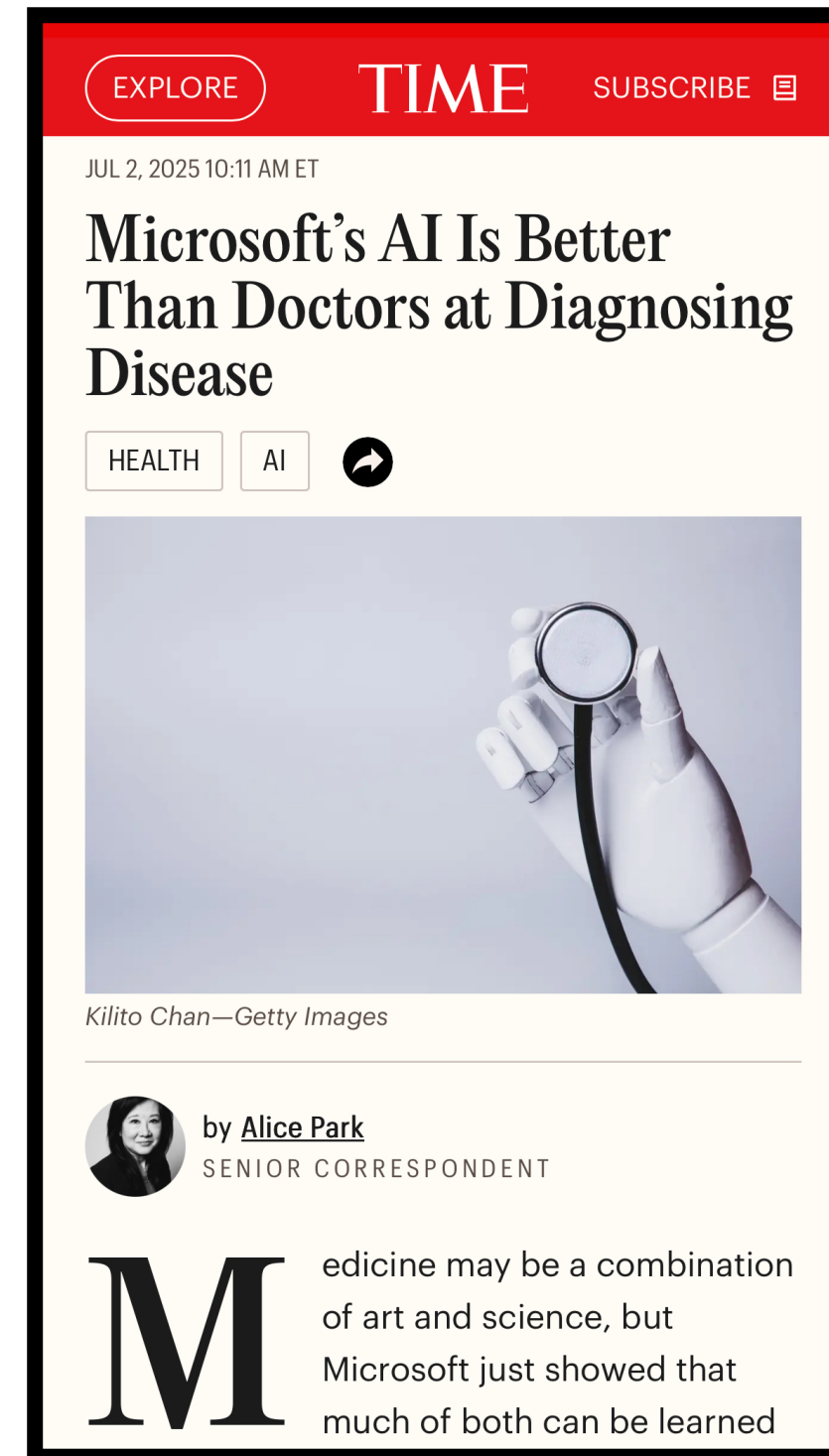
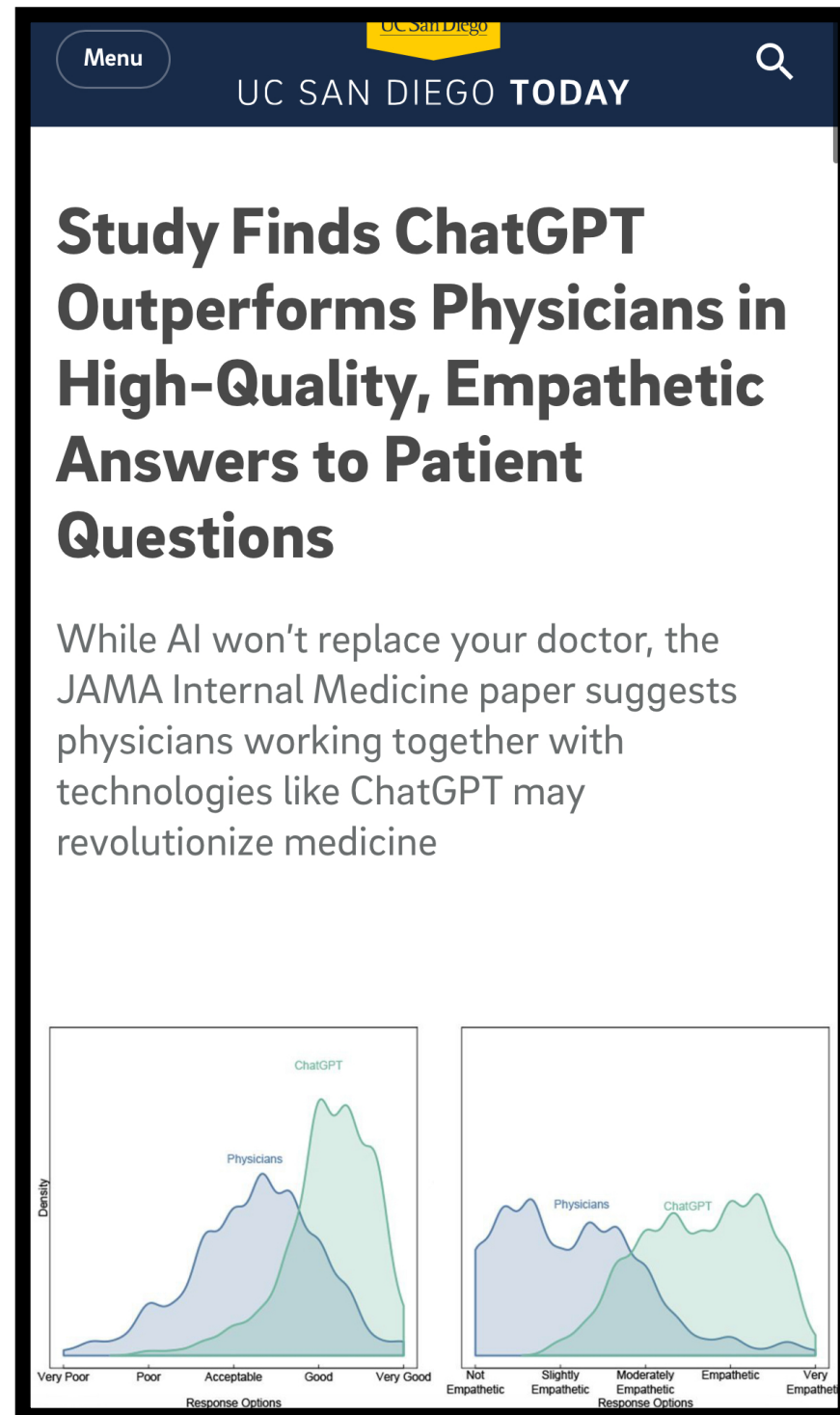
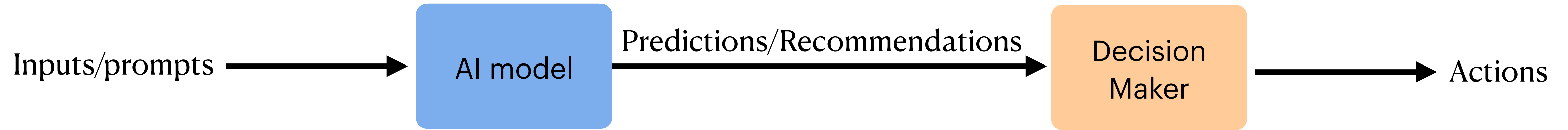
- Automated
- Humans



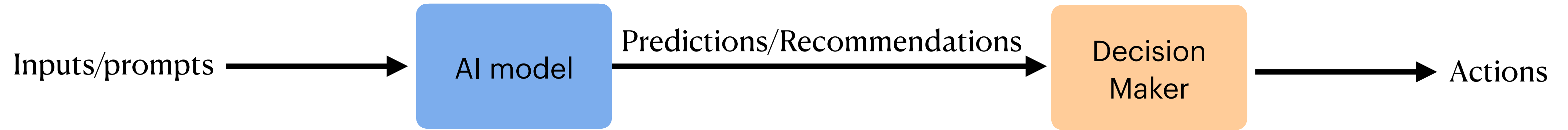
AI-powered Decision making pipeline



AI-powered Decision making pipeline



AI-powered Decision making pipeline



Study Finds ChatGPT Outperforms Physicians in High-Quality, Empathetic Answers to Patient Questions

While AI won't replace your doctor, the JAMA Internal Medicine paper suggests physicians working together with technologies like ChatGPT may revolutionize medicine

Response Option	Physicians Density	ChatGPT Density
Very Poor	Low	Very Low
Poor	Low	Very Low
Acceptable	High	Low
Good	Low	High
Very Good	Low	Low

Response Option	Physicians Density	ChatGPT Density
Not Empathetic	High	Low
Slightly Empathetic	High	Low
Moderately Empathetic	Low	High
Empathetic	Low	High
Very Empathetic	Low	Low

Microsoft's AI Is Better Than Doctors at Diagnosing Disease

by Alice Park
SENIOR CORRESPONDENT

Medicine may be a combination of art and science, but Microsoft just showed that much of both can be learned

Deep learning algorithm does as well as dermatologists in identifying skin cancer

January 25th, 2017 | 6 min read
Science & Engineering

In hopes of creating better access to medical care, Stanford researchers have trained an algorithm to diagnose skin cancer.

It's scary enough making a doctor's appointment to see if a strange mole could be

AI-powered Decision making pipeline



Study Finds ChatGPT Outperforms Physicians in High-Quality, Empathetic Answers to Patient Questions

While AI won't replace your doctor, the JAMA Internal Medicine paper suggests physicians working together with technologies like ChatGPT may revolutionize medicine

The image shows two line graphs. The left graph plots 'Density' on the y-axis against 'Response Options' on the x-axis, ranging from 'Very Poor' to 'Very Good'. It shows two curves: a blue curve for 'Physicians' and a green curve for 'ChatGPT'. The ChatGPT curve is shifted further to the right (towards 'Good' and 'Very Good') compared to the Physicians curve. The right graph plots 'Density' on the y-axis against 'Response Options' on the x-axis, ranging from 'Not Empathetic' to 'Very Empathetic'. It shows two curves: a blue curve for 'Physicians' and a green curve for 'ChatGPT'. The ChatGPT curve is shifted further to the right (towards 'Moderately Empathetic' and 'Very Empathetic') compared to the Physicians curve.

Microsoft's AI Is Better Than Doctors at Diagnosing Disease

HEALTH AI

Kililo Chan—Getty Images

by [Alice Park](#)
SENIOR CORRESPONDENT

Medicine may be a combination of art and science, but Microsoft just showed that much of both can be learned

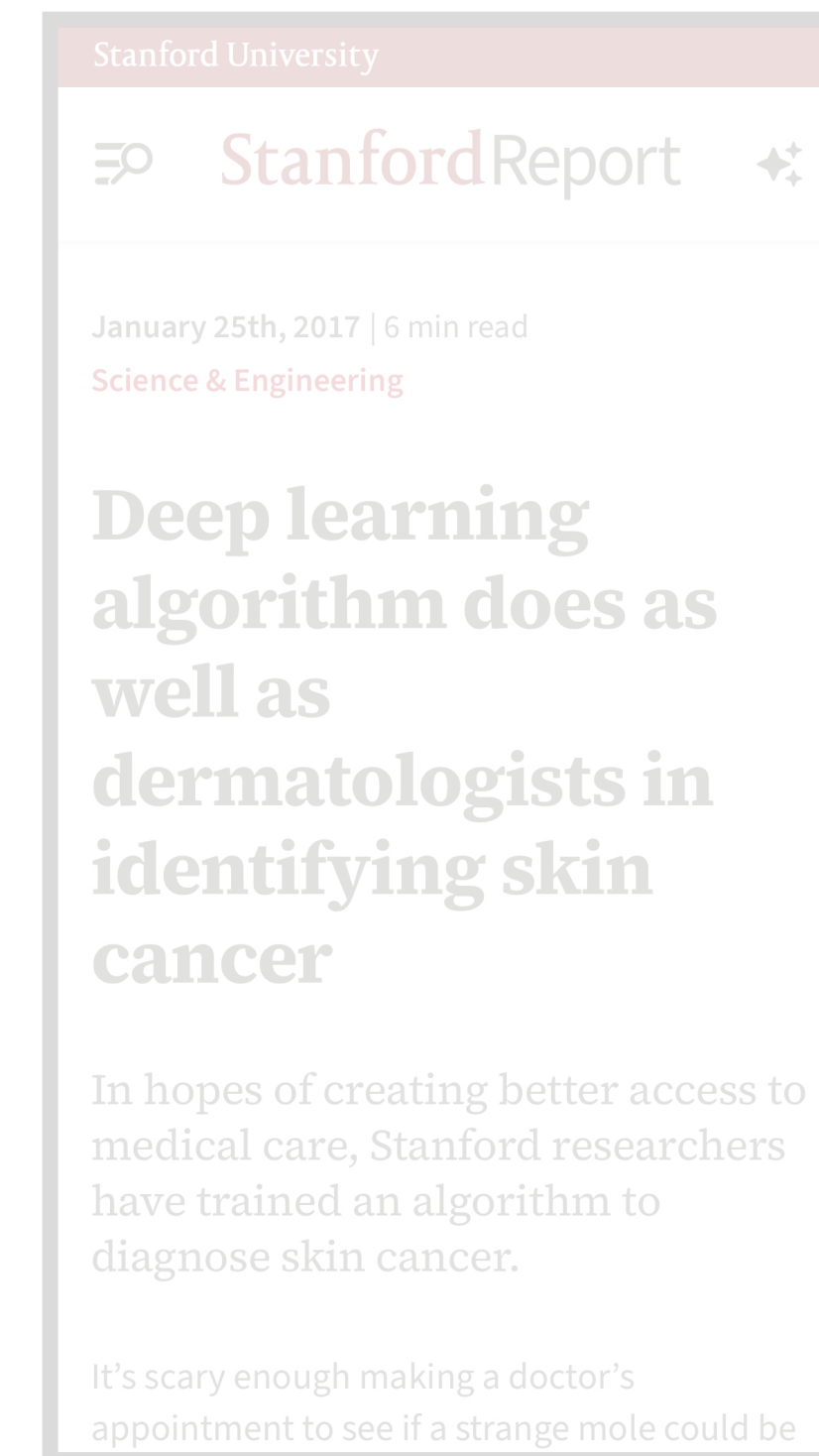
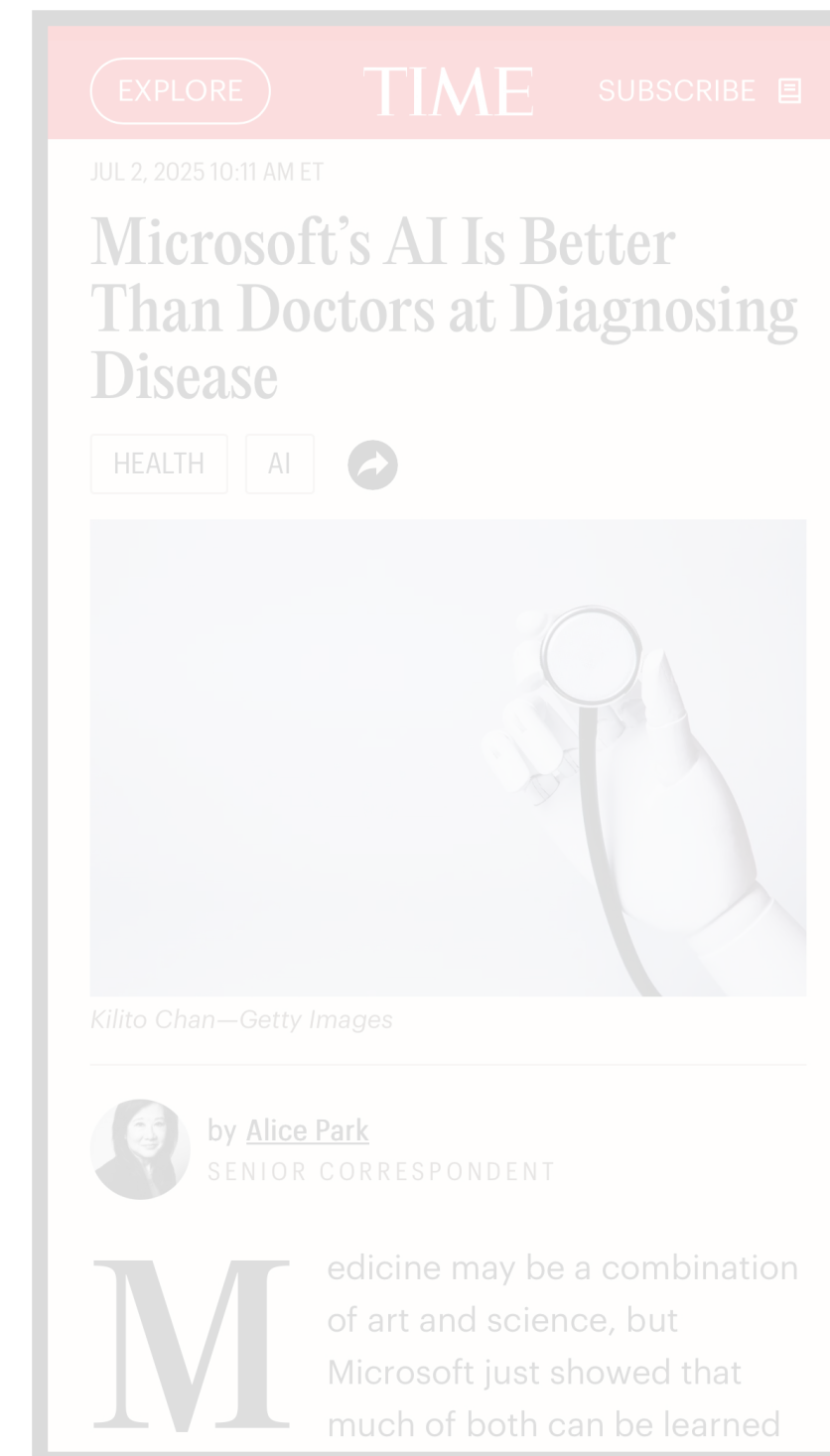
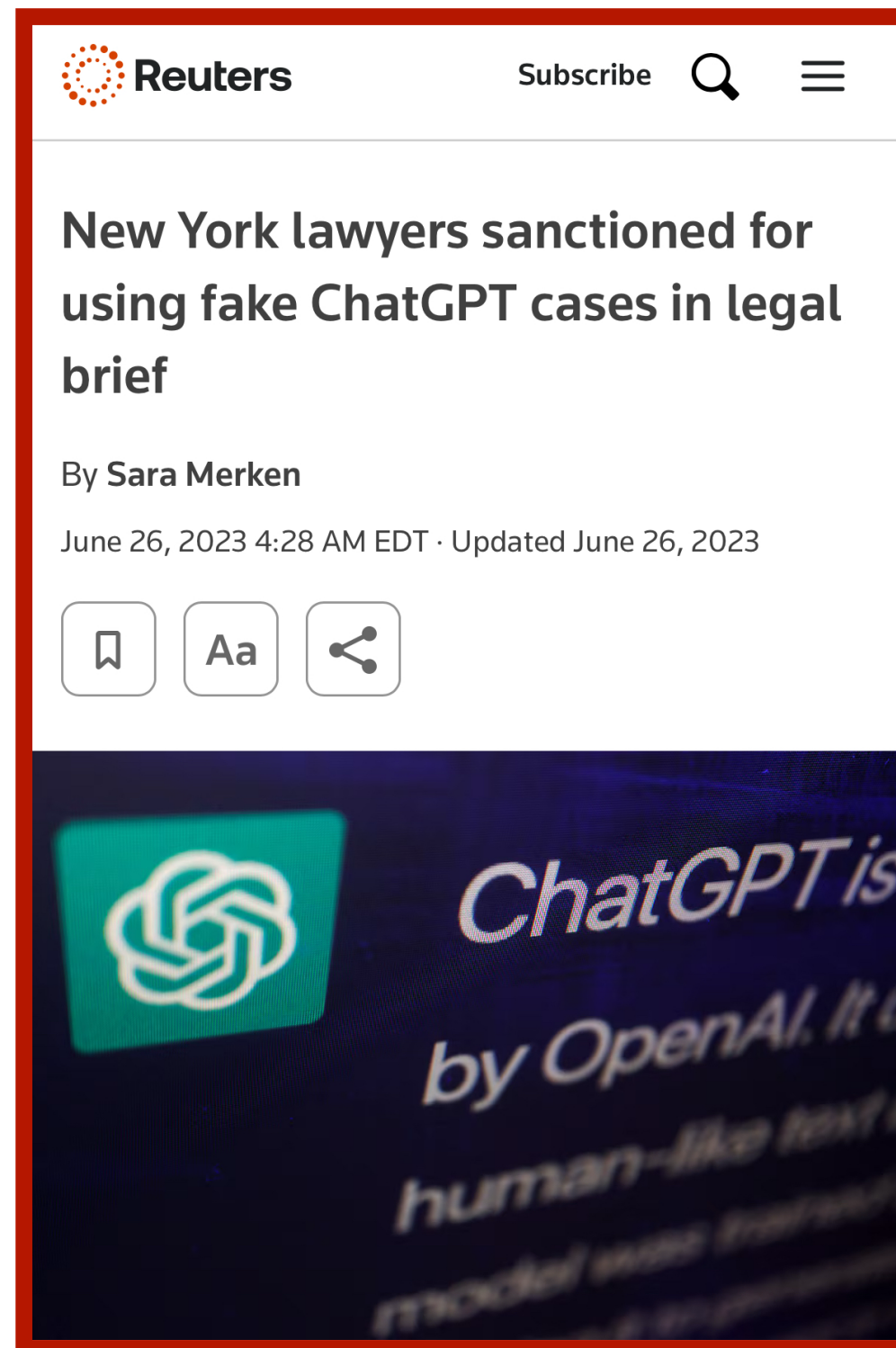
Deep learning algorithm does as well as dermatologists in identifying skin cancer

January 25th, 2017 | 6 min read
Science & Engineering

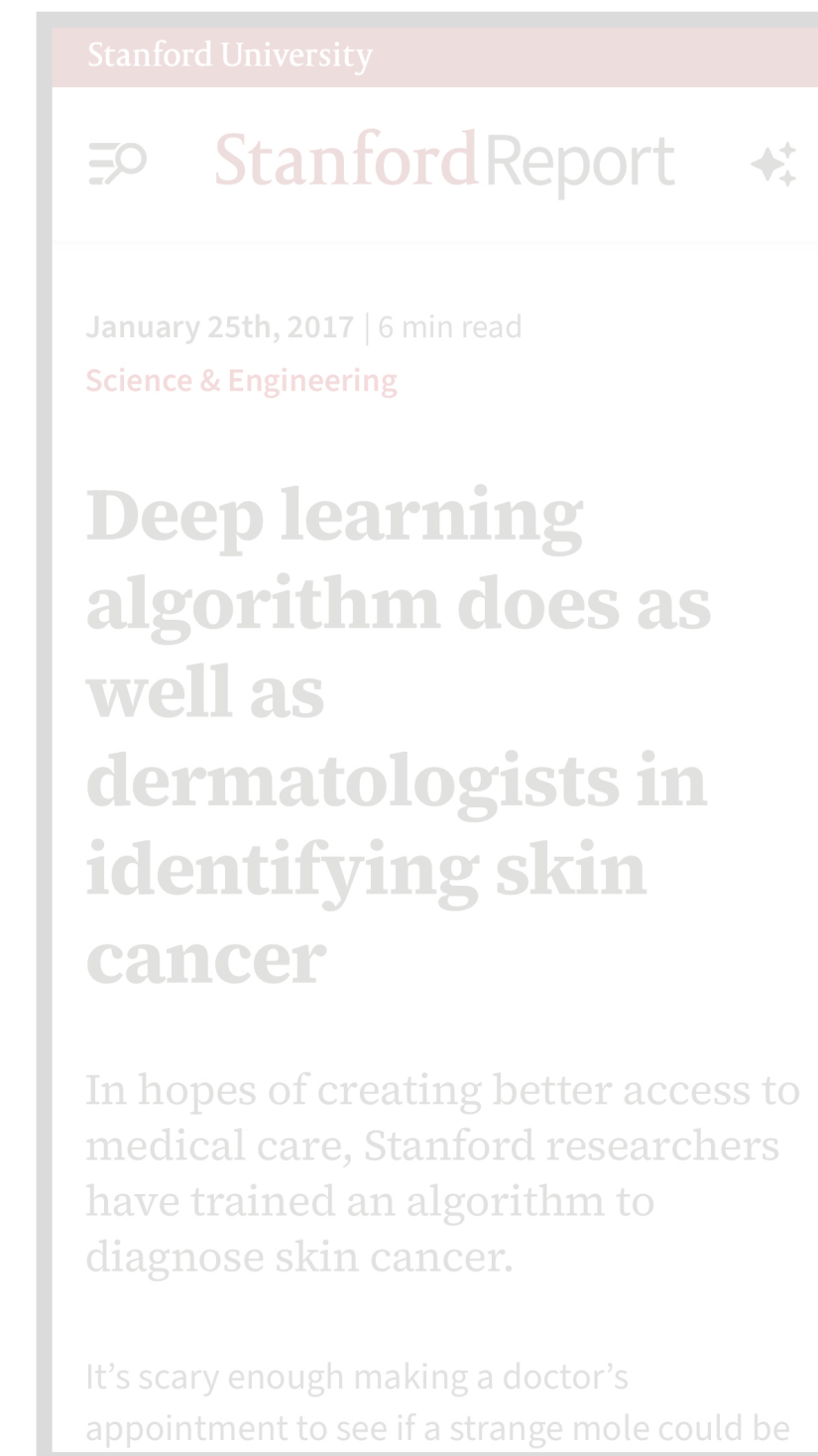
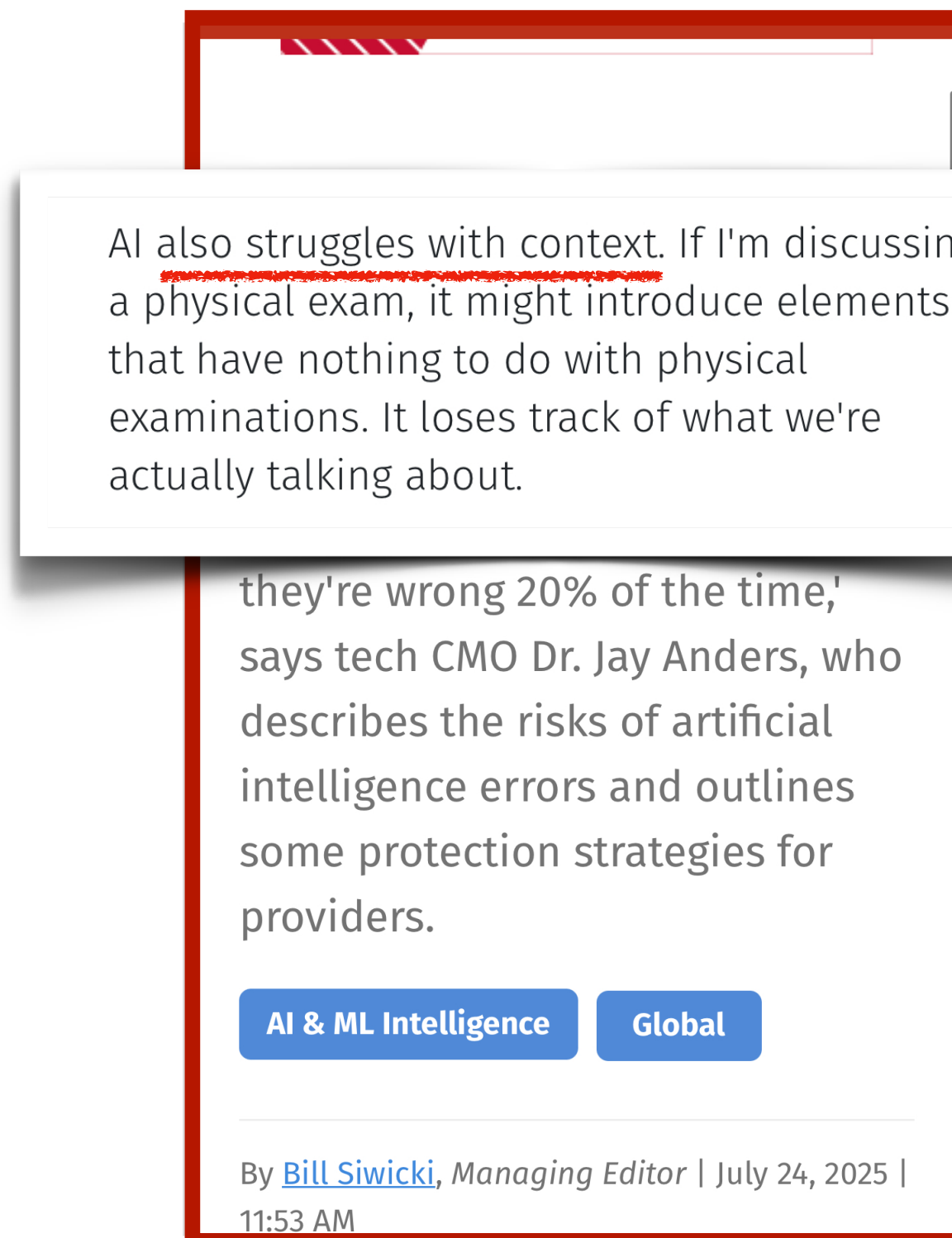
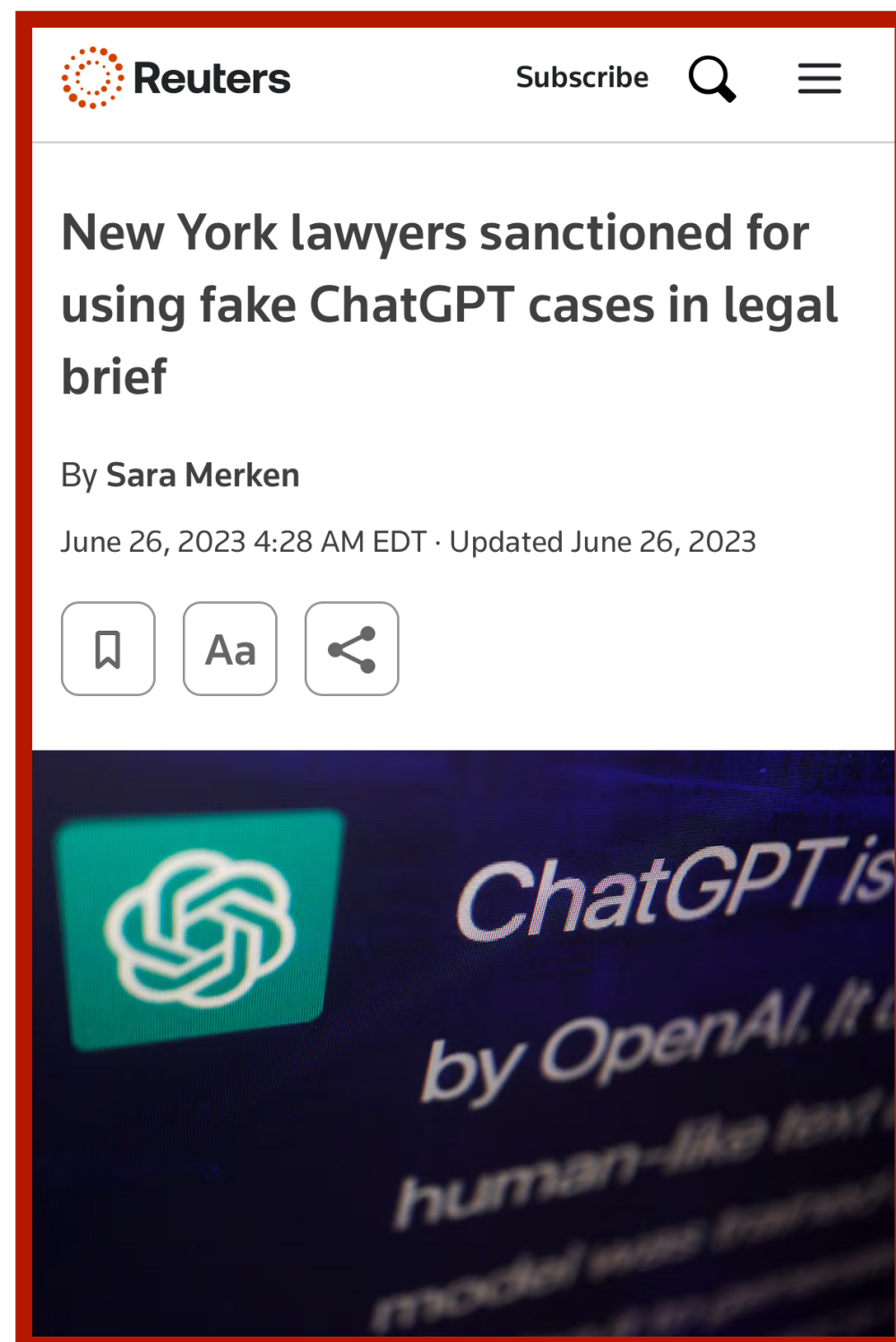
In hopes of creating better access to medical care, Stanford researchers have trained an algorithm to diagnose skin cancer.

It's scary enough making a doctor's appointment to see if a strange mole could be

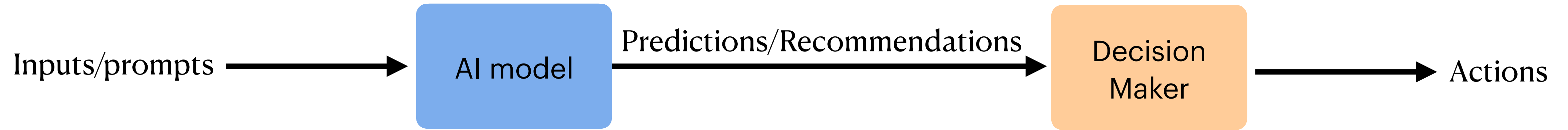
AI-powered Decision making pipeline



AI-powered Decision making pipeline



AI-powered Decision making pipeline



Reuters Subscribe 🔍 ☰

New York lawyers sanctioned for using fake ChatGPT cases in legal brief

By Sara Merken

June 26, 2023 4:28 AM EDT · Updated June 26, 2023

🔖 Aa 🔄

AI also struggles with context. If I'm discussing a physical exam, it might introduce elements that have nothing to do with physical examinations. It loses track of what we're actually talking about.

they're wrong 20% of the time,' says tech CMO Dr. Jay Anders, who describes the risks of artificial intelligence errors and outlines some protection strategies for providers.

AI & ML Intelligence Global

By [Bill Siwicki](#), Managing Editor | July 24, 2025 | 11:53 AM

AXIOS Pittsburgh 🔍 👤

Aug 27, 2025 - News

AI is overconfident even when wrong, says report

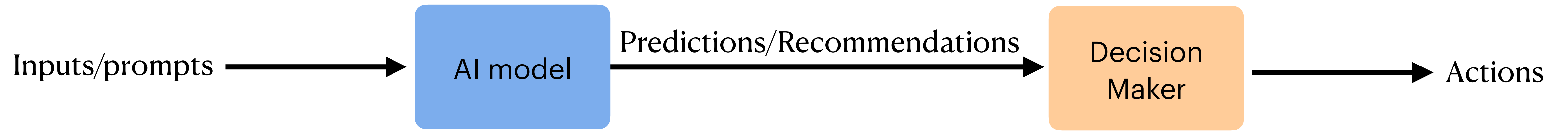
Ryan Deto

f X in T ✉

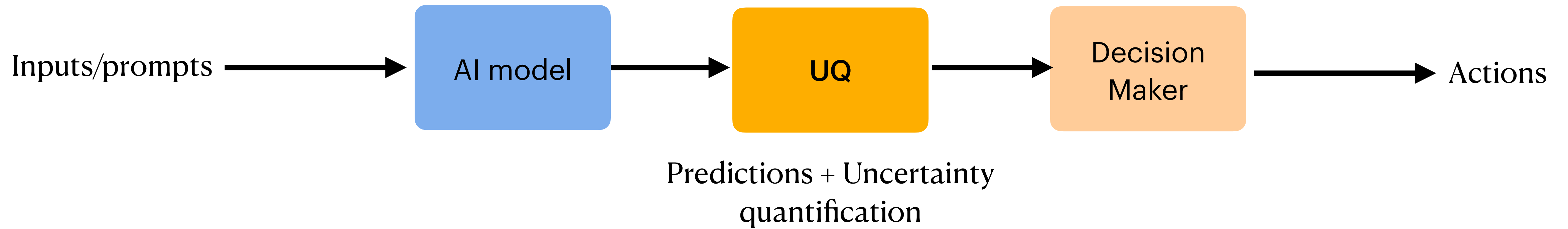
📄 Add Axios on Google

Illustration: Allie Carl/Axios

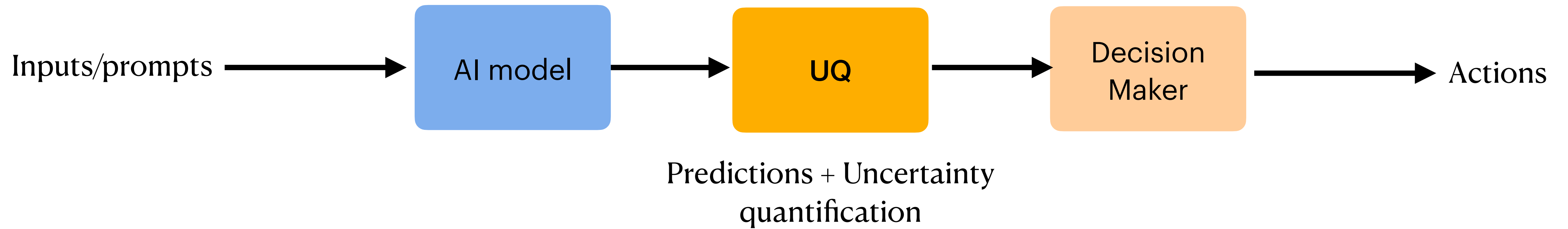
AI-powered Decision making pipeline



AI-powered Decision making pipeline

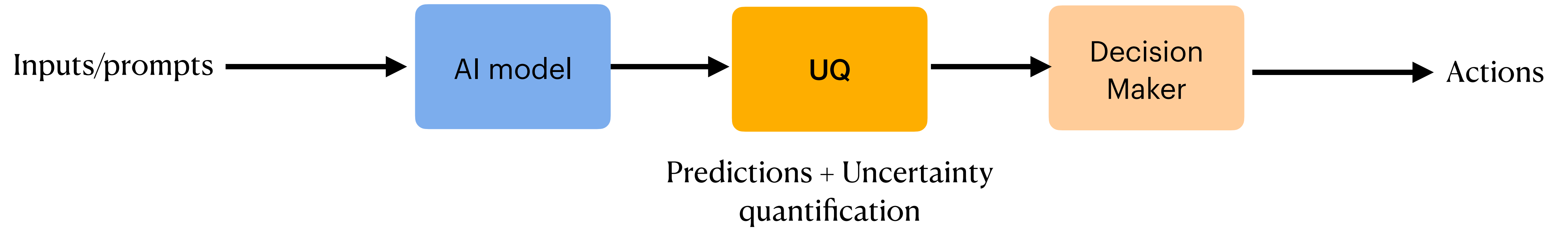


AI-powered Decision making pipeline



Robust Decision making requires precise Uncertainty Quantification

AI-powered Decision making pipeline



Part I: CP for Generative Models

Conformal Prediction Beyond the Seen: A Missing Mass Perspective for Uncertainty Quantification in Generative Models

Sima Noorani^{*1}, Shayan Kiyani^{*1}, George Pappas¹, and Hamed Hassani¹

¹University of Pennsylvania

Part II: Collaborative CP

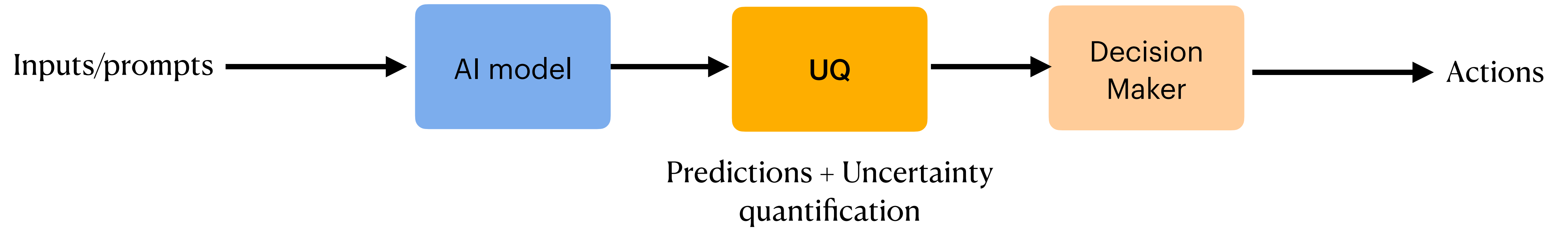
Human AI Collaborative Uncertainty Quantification

Sima Noorani^{*1}, Shayan Kiyani^{*1}, George Pappas¹, and Hamed Hassani¹

¹University of Pennsylvania

Arxiv - Oct '25

AI-powered Decision making pipeline



Part I: CP for Generative Models

Conformal Prediction Beyond the Seen: A Missing Mass Perspective for Uncertainty Quantification in Generative Models

Sima Noorani^{*1}, Shayan Kiyani^{*1}, George Pappas¹, and Hamed Hassani¹

¹University of Pennsylvania

Part II: Collaborative CP

Human AI Collaborative Uncertainty Quantification

Sima Noorani^{*1}, Shayan Kiyani^{*1}, George Pappas¹, and Hamed Hassani¹

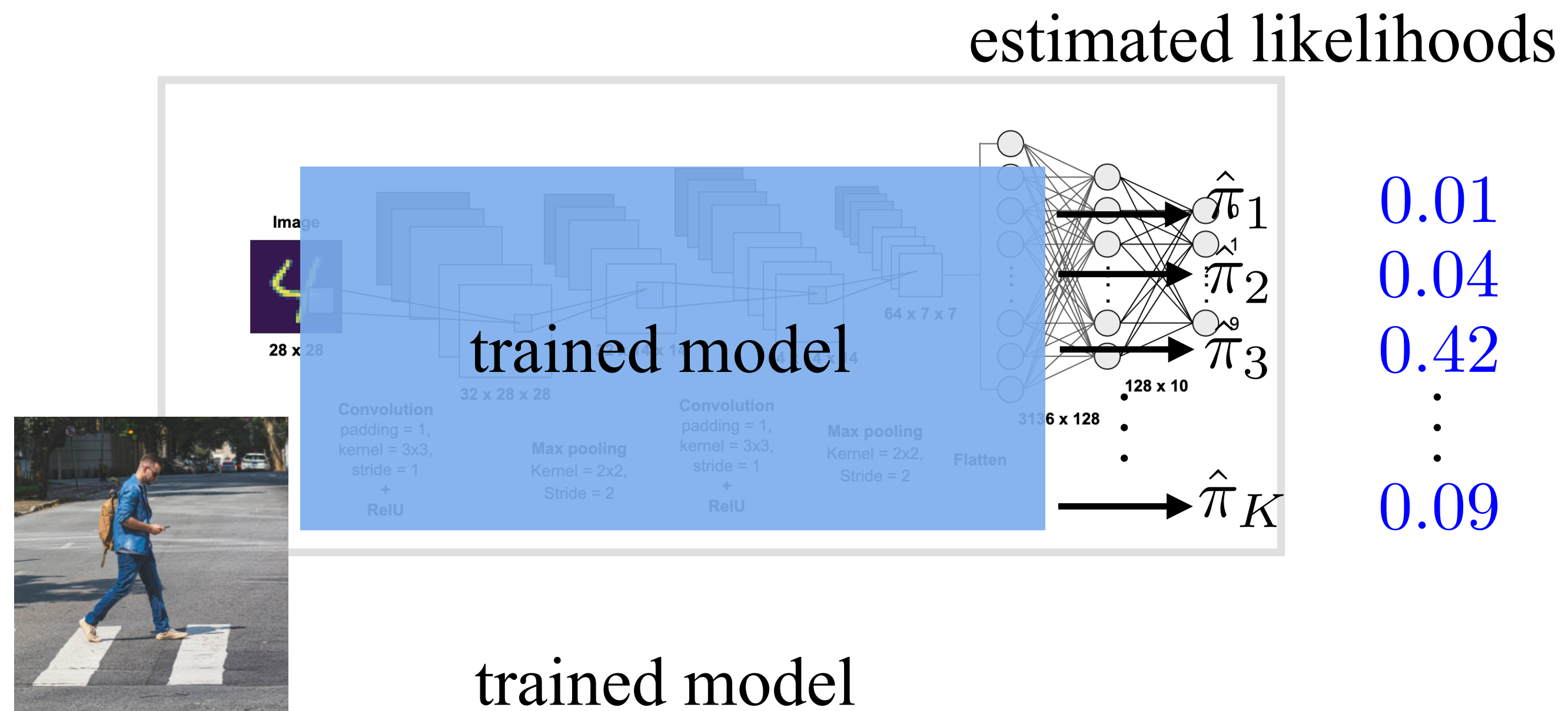
¹University of Pennsylvania

Arxiv - Oct '25

How to do uncertainty quantification in the generative setting (LLMs) ?

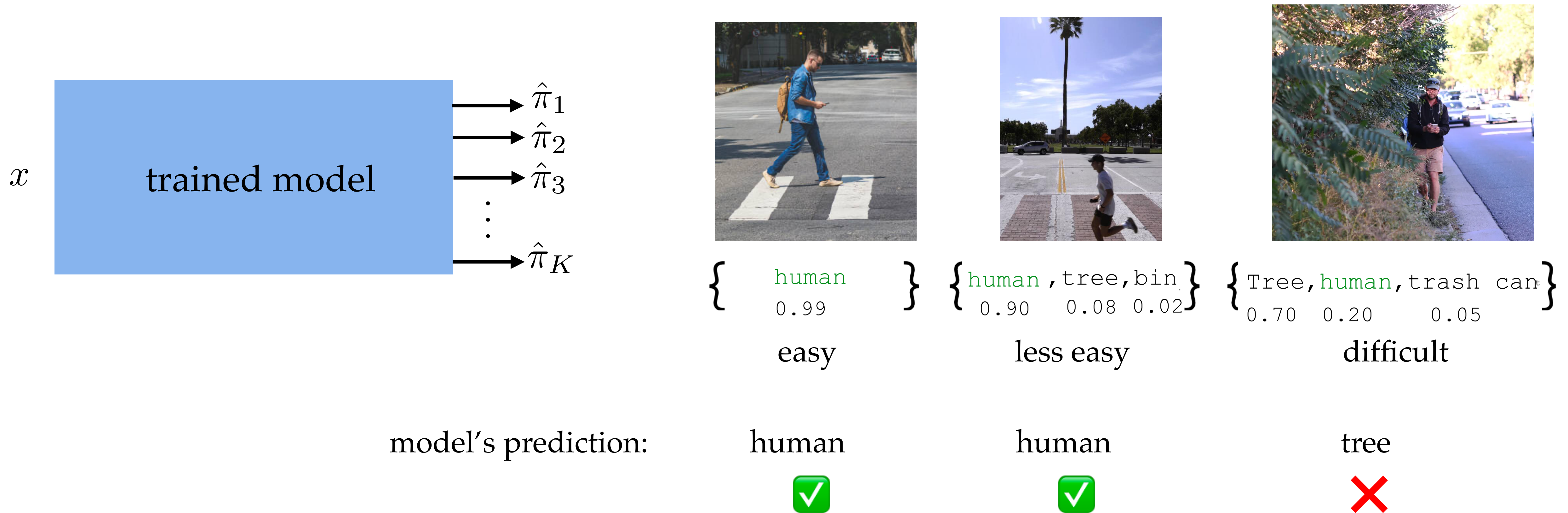
(Quick) Review of Conformal Prediction

Conformal Prediction



- prediction is then based on the class that has maximum likelihood

Conformal Prediction



➔ the trained model provides estimated likelihoods (a notion of **uncertainty**)

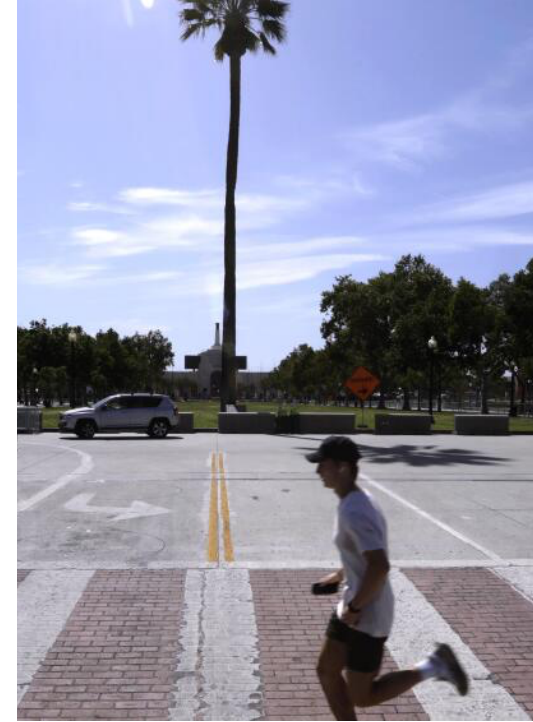
➔ these likelihoods are informative but not always correct

Conformal Prediction

x



{ human
0.99 }



{ human, tree, bin
0.90 0.08 0.02 }



{ Tree, human, trash can
0.70 0.20 0.05 }

model's prediction:

human



human

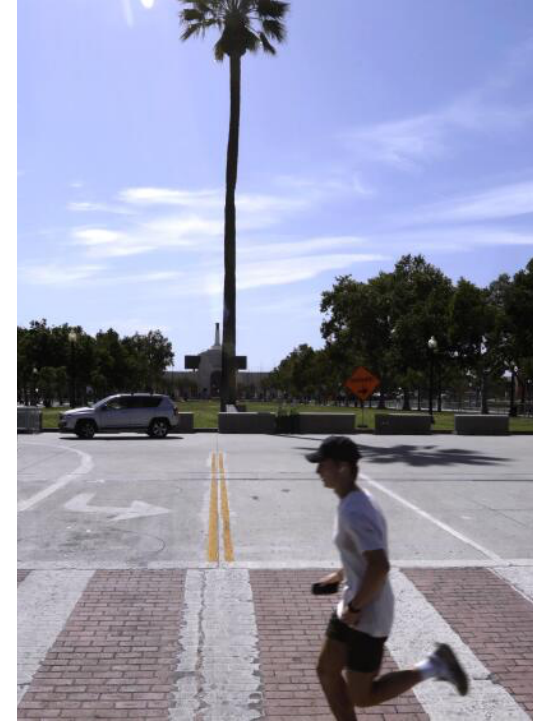


tree



Conformal Prediction

x



$\left\{ \begin{array}{l} \text{human} \\ 0.99 \end{array} \right\}$ $\left\{ \begin{array}{l} \text{human}, \text{tree}, \text{bin} \\ 0.90 \quad 0.08 \quad 0.02 \end{array} \right\}$ $\left\{ \begin{array}{l} \text{Tree}, \text{human}, \text{trash can} \\ 0.70 \quad 0.20 \quad 0.05 \end{array} \right\}$

$C(x)$

(model's conformal prediction)

$\{ \text{human} \}$



$\{ \text{human} \}$

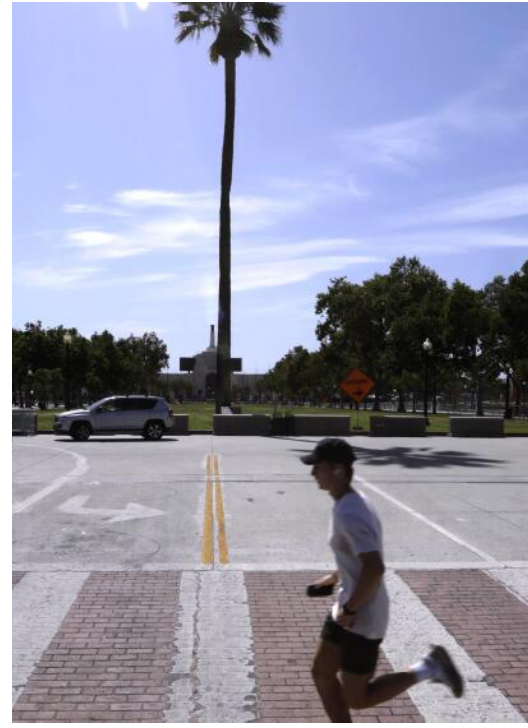


$\{ \text{tree}, \text{human} \}$



Single prediction \rightarrow set of predictions

x



{ human
0.99 }

{ human, tree, bin
0.90 0.08 0.02 }

{ Tree, human, trash can
0.70 0.20 0.05 }

$C(x)$

(model's conformal prediction)

{ human }



{ human }



{ tree, human }



Marginal Coverage

$$\Pr \{Y_{\text{test}} \in C(X_{\text{test}})\} \geq 1 - \alpha$$

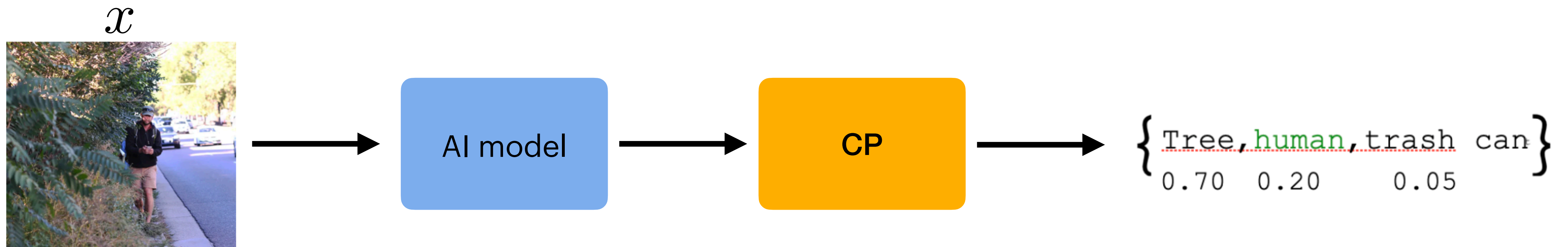
User-specified value; e.g. $1 - \alpha = 0.9$

Conformal Prediction

Marginal Coverage

$$\Pr \{Y_{\text{test}} \in C(X_{\text{test}})\} \geq 1 - \alpha$$

Classification



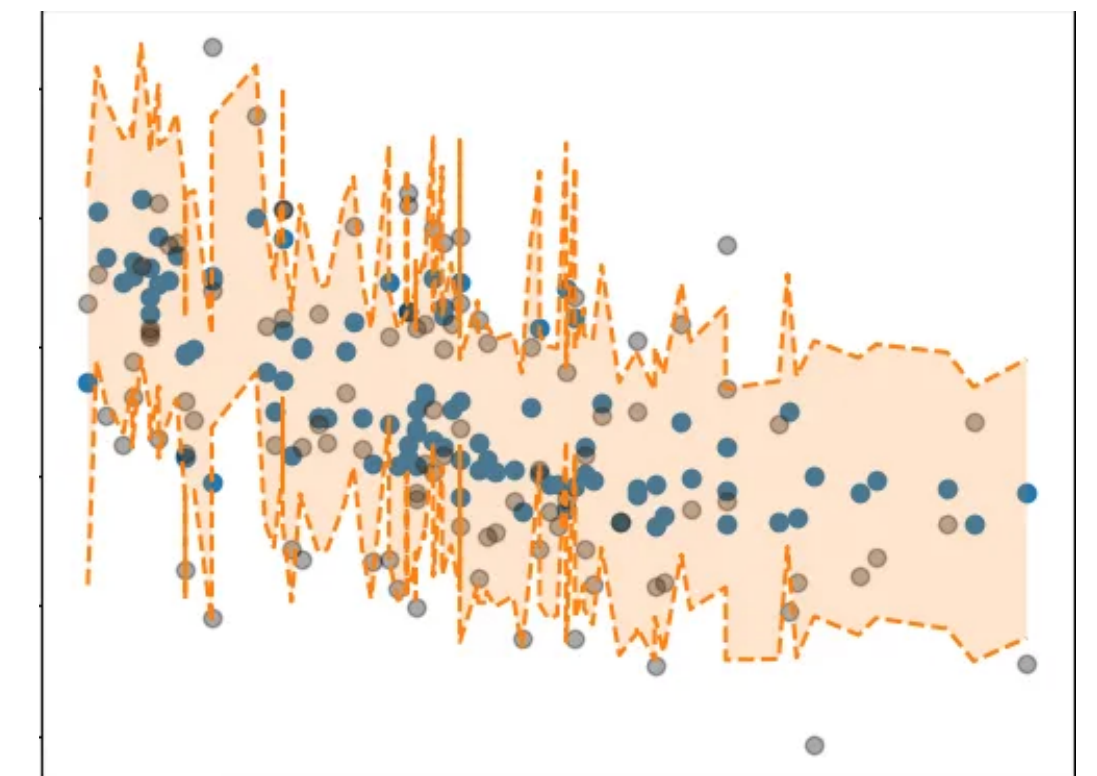
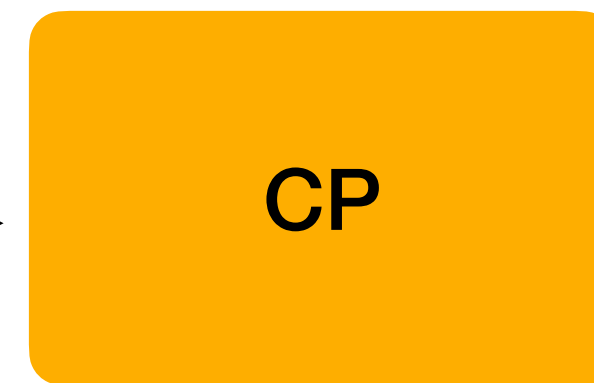
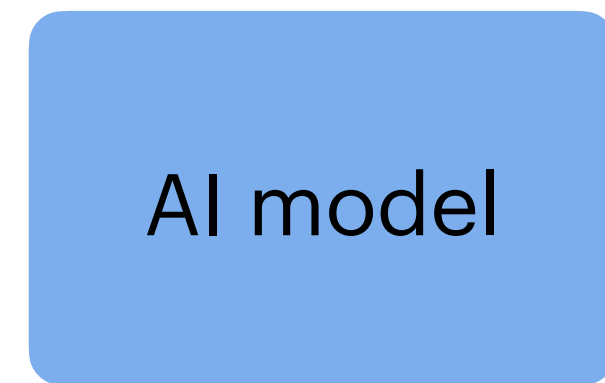
Conformal Prediction

Marginal Coverage

$$\Pr \{Y_{\text{test}} \in C(X_{\text{test}})\} \geq 1 - \alpha$$

Regression

\mathcal{X}



[10.1, 10.8]

Conformal Prediction

Marginal Coverage

$$\Pr \{Y_{\text{test}} \in C(X_{\text{test}})\} \geq 1 - \alpha$$

$$C(x) = \{y \in \mathcal{Y} : S(x, y) \leq q\}$$

the threshold



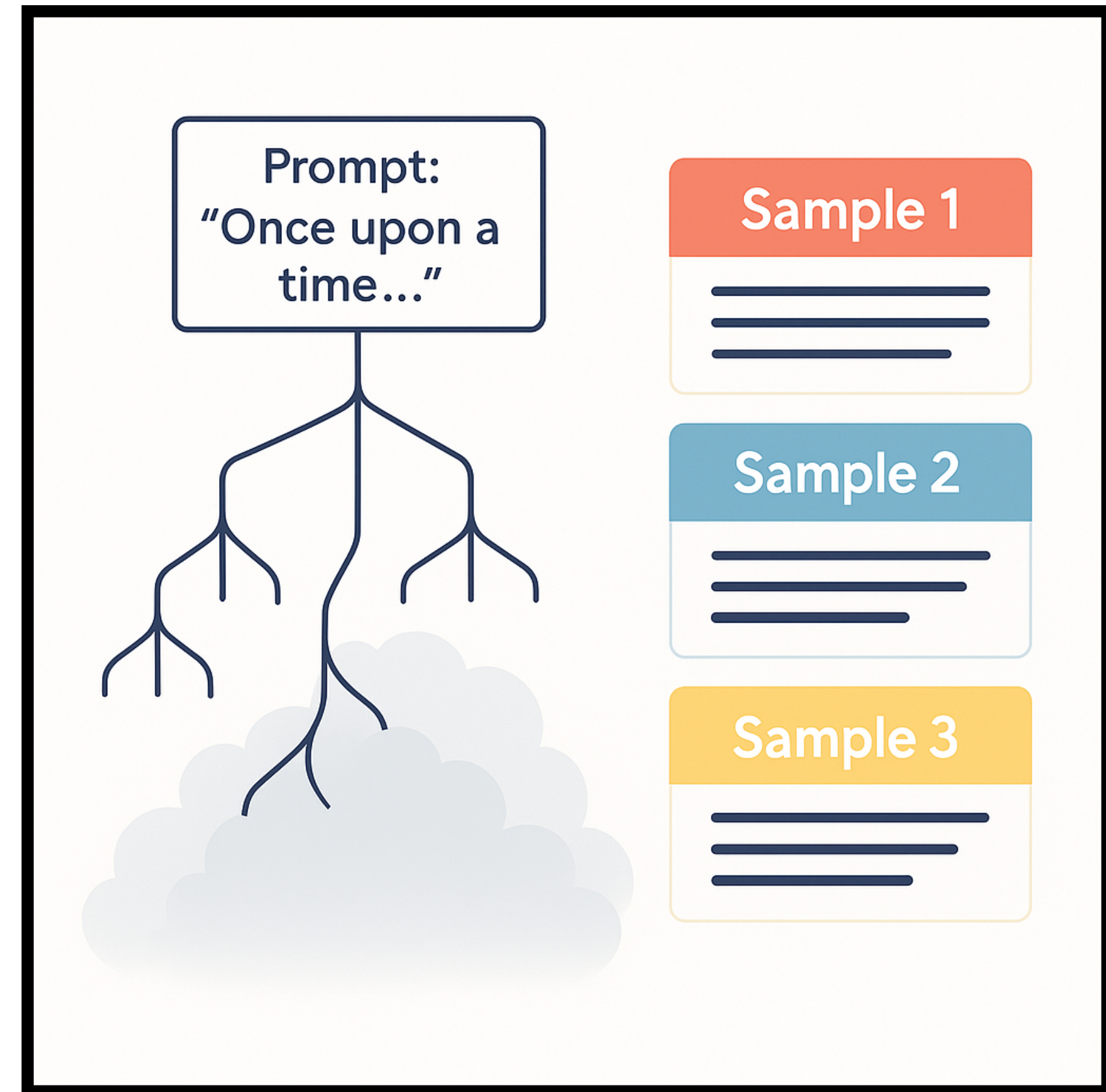
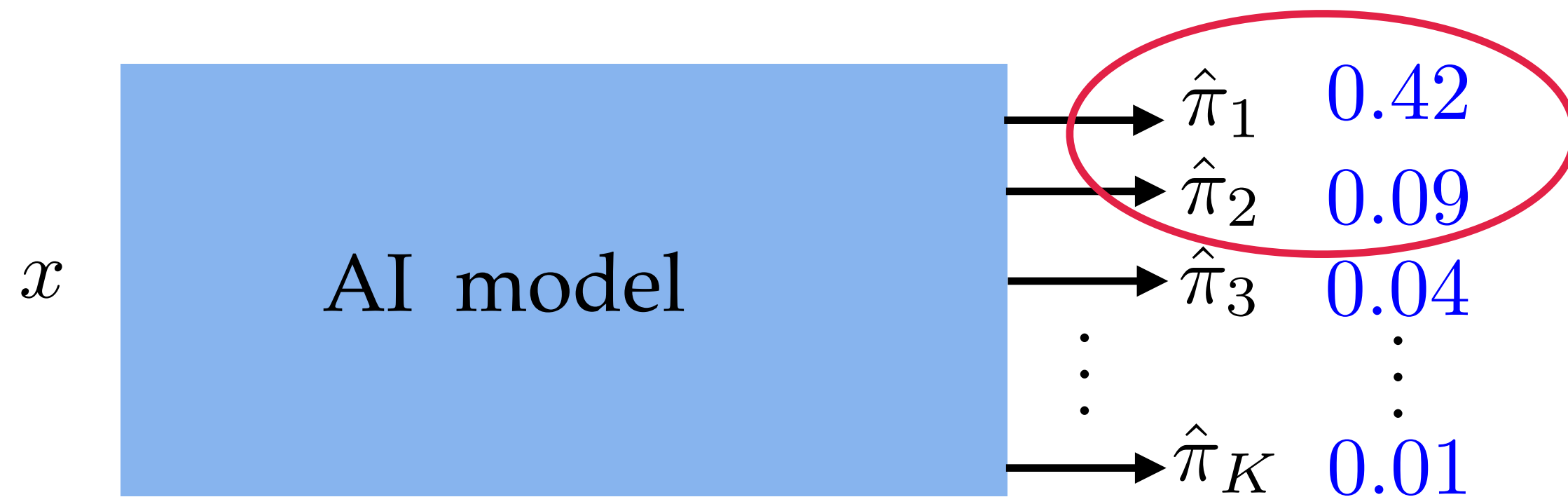
difficulty of the input ↑

models' uncertainty about the label ↑

size of the prediction set ↑

The Challenges of Extending CP to Generative Models

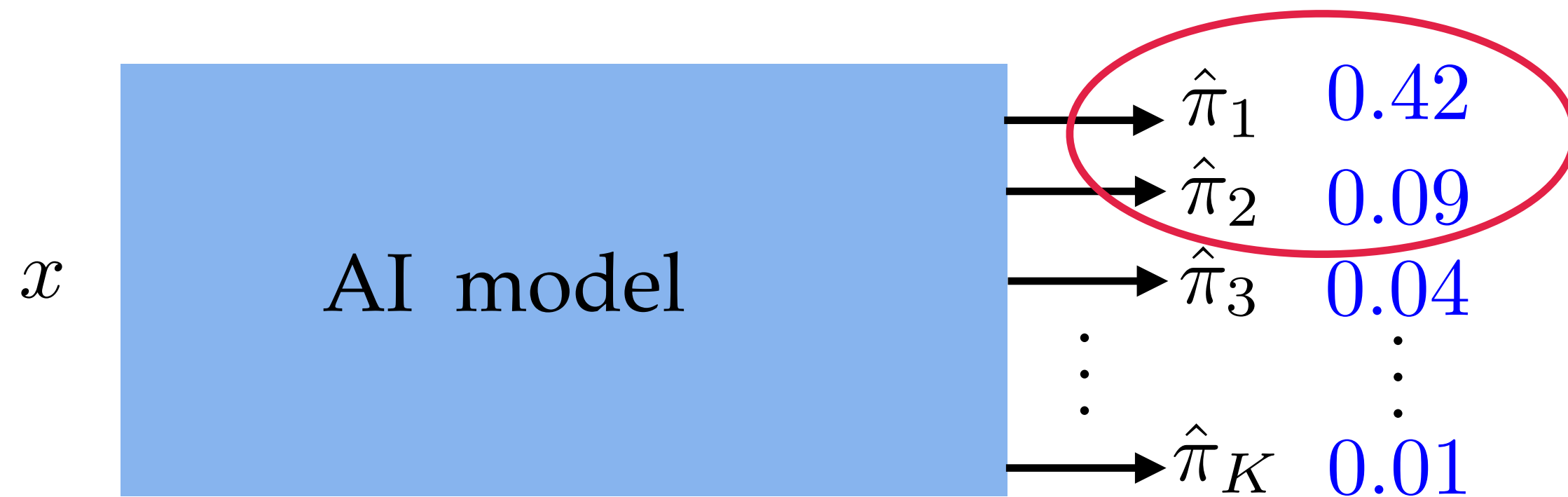
The **Challenges** of Extending CP to **Generative Models**



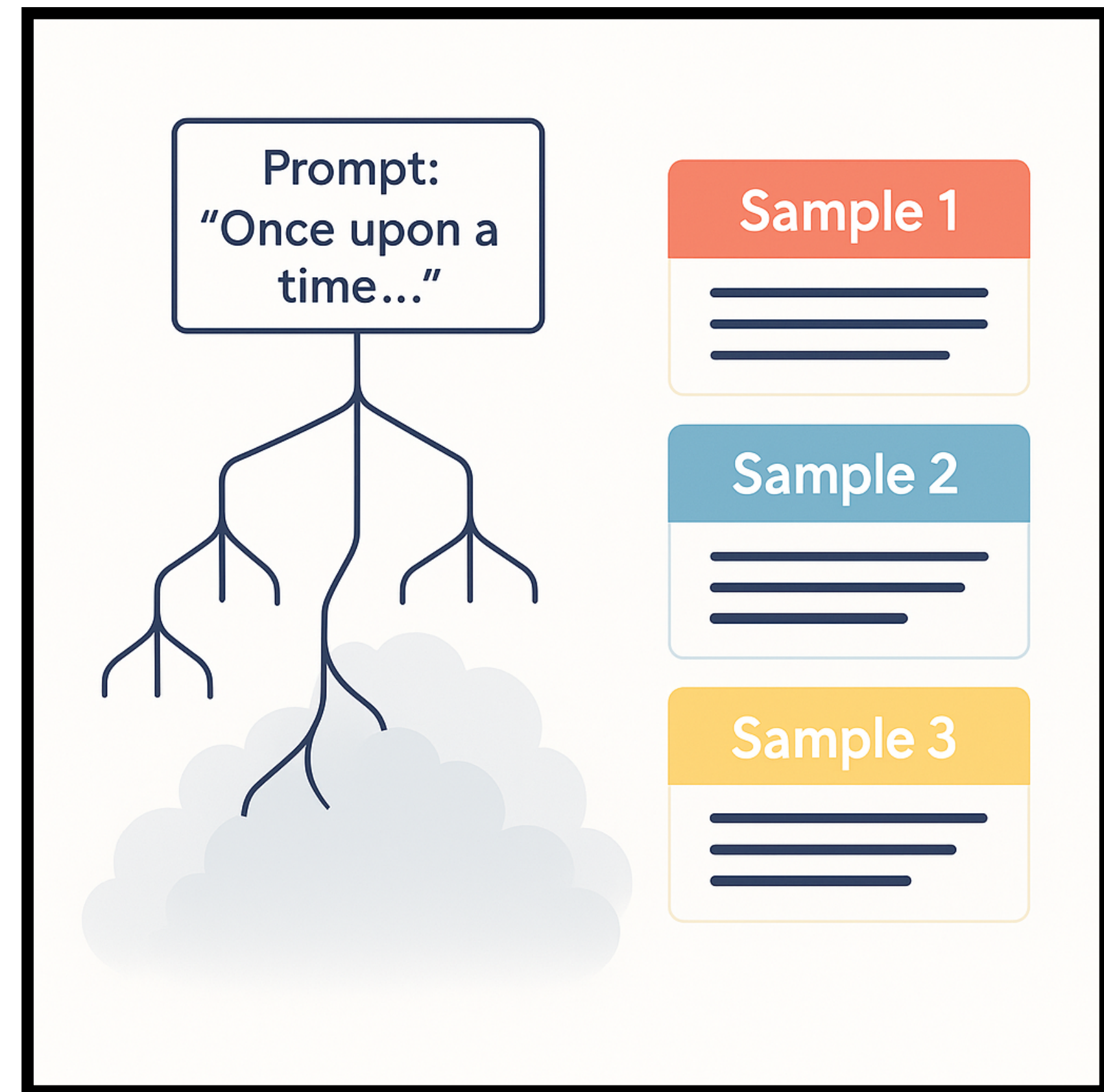
$$C(x) = \{y \in \mathcal{Y} : S(x, y) \leq q\}$$

- Enumerate? **X**
- Represent compactly? **X**

The **Challenges** of Extending CP to **Generative** Models



$$C(x) = \{y \in \mathcal{Y} : S(x, y) \leq q\}$$

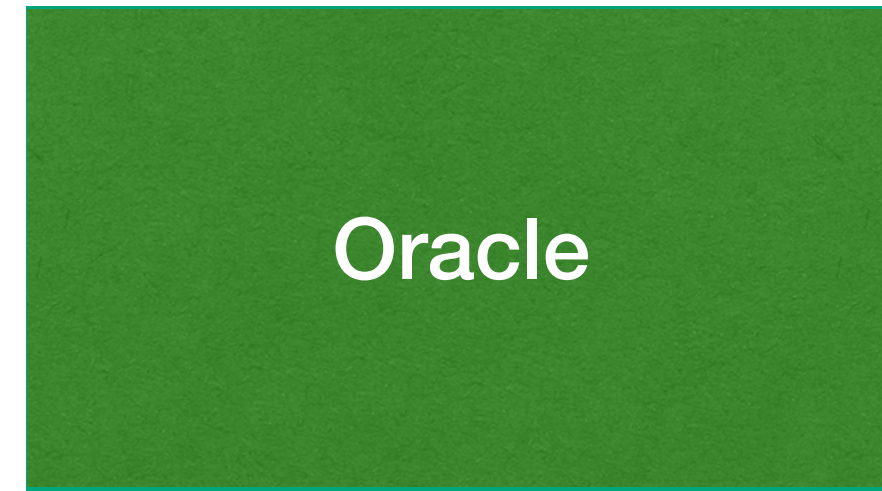


Challenge 1 : We are dealing with an unstructured + unbounded output space

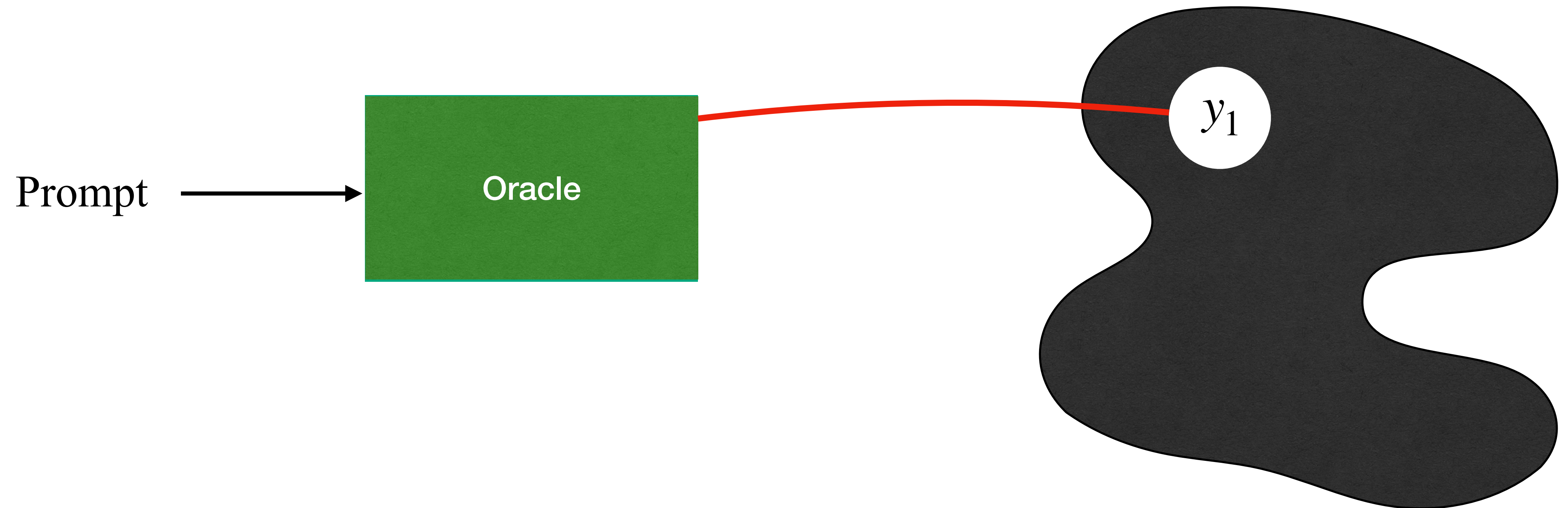
The **Challenges** of Extending CP to **Generative** Models



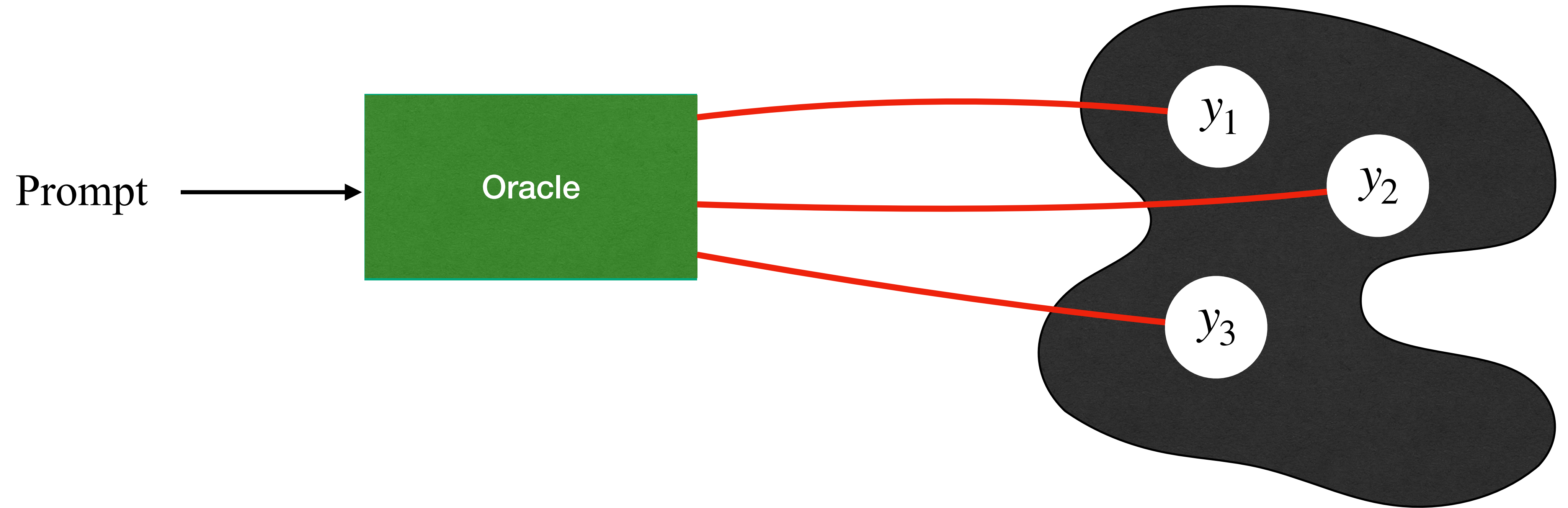
The **Challenges** of Extending CP to **Generative** Models



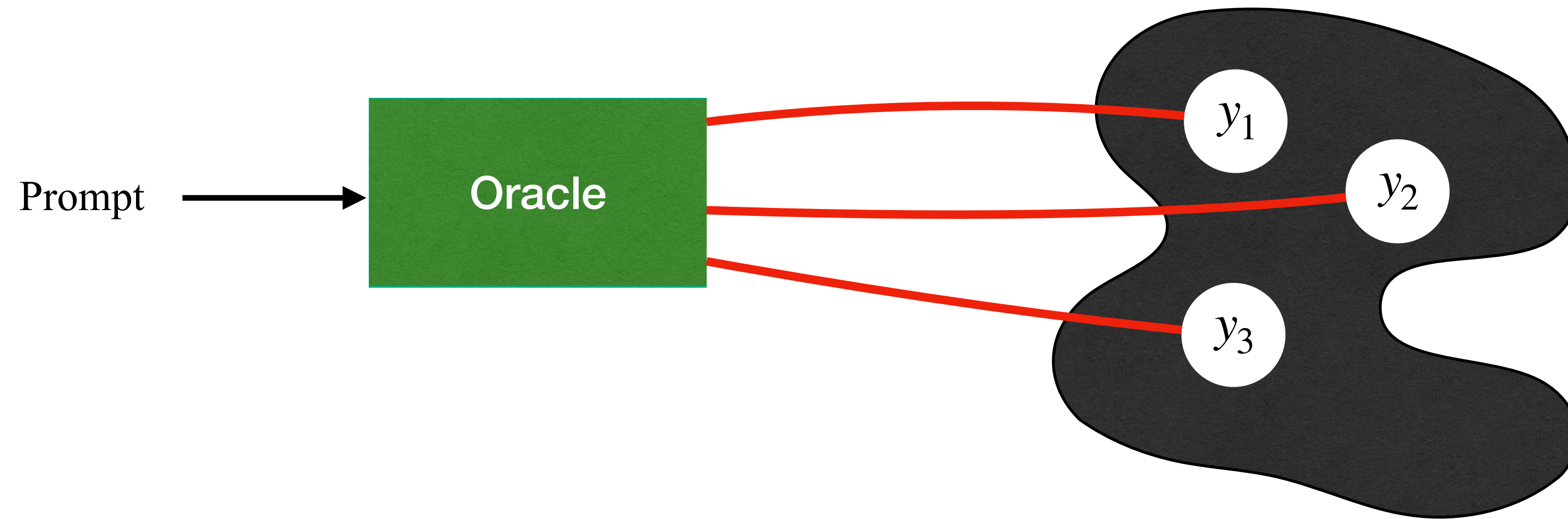
The **Challenges** of Extending CP to **Generative** Models



The **Challenges** of Extending CP to **Generative** Models

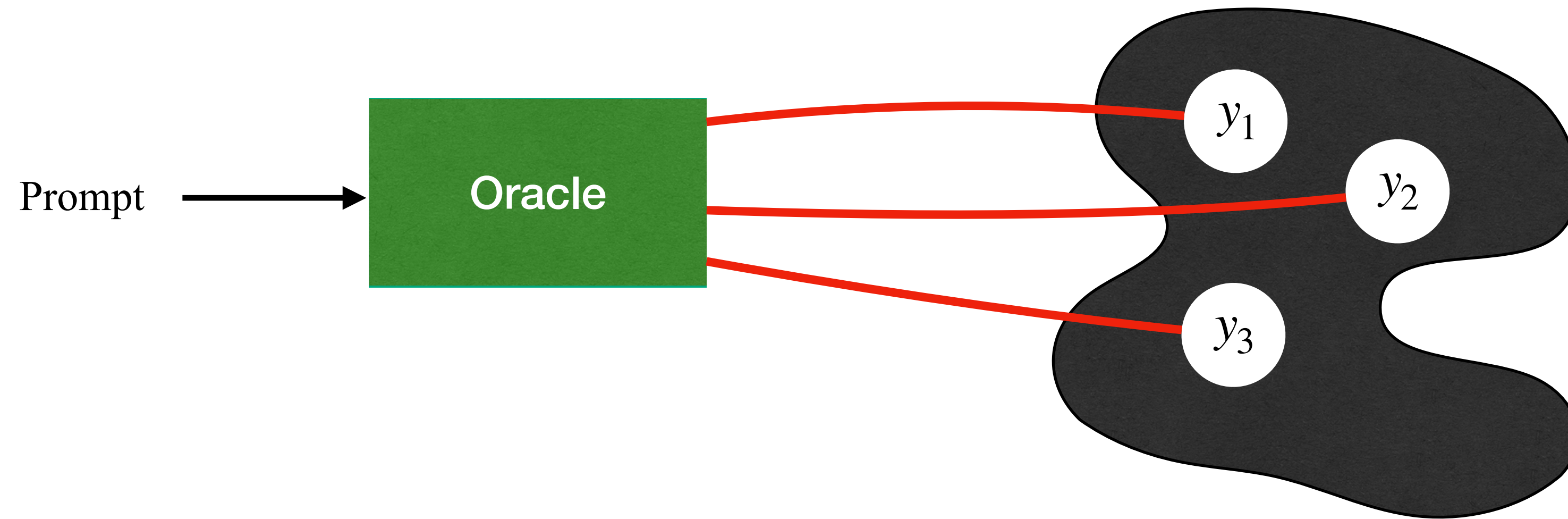


The **Challenges** of Extending CP to **Generative** Models



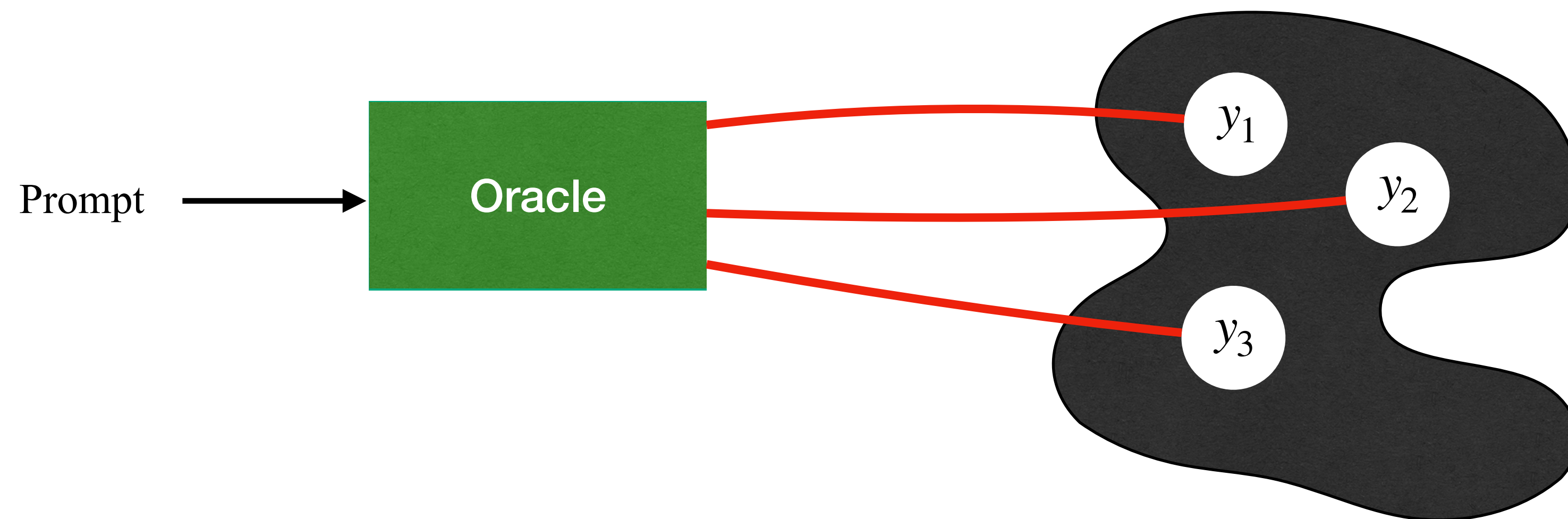
$$x \longrightarrow C(x) \subseteq \{\text{set of representative outputs}\}$$

The **Challenges** of Extending CP to **Generative** Models



$$\begin{aligned} x &\longrightarrow C(x) \subseteq \{\text{set of representative outputs}\} \\ &\subseteq \{y_1, y_2, \dots, y_t\} \end{aligned}$$

The **Challenges** of Extending CP to **Generative** Models



$$\begin{aligned} x &\longrightarrow C(x) \subseteq \{\text{set of representative outputs}\} \\ &\subseteq \{y_1, y_2, \dots, y_t\} \end{aligned}$$

Challenge 2 : query only access to the output space

Conformal Prediction for Generative Models (Setting)

$$(X, Y) \in \mathcal{X} \times \mathcal{Y}$$

$$(x, y) \sim p(x, y)$$

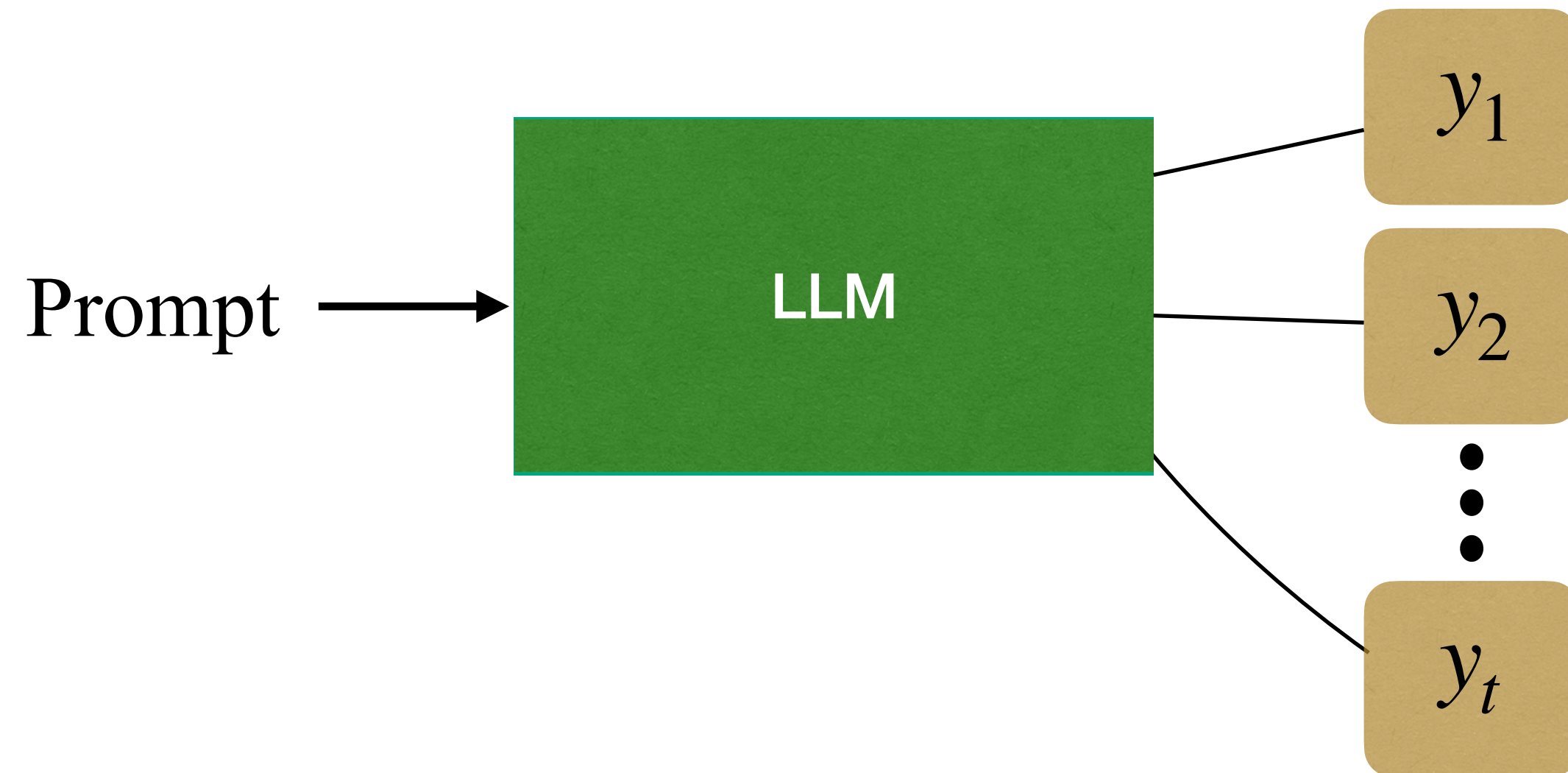
Query Oracle: $y \sim \pi(y \mid x)$

Conformal Prediction for Generative Models (Setting)

$$(X, Y) \in \mathcal{X} \times \mathcal{Y}$$

$$(x, y) \sim p(x, y)$$

$$\text{Query Oracle: } y \sim \pi(y \mid x)$$

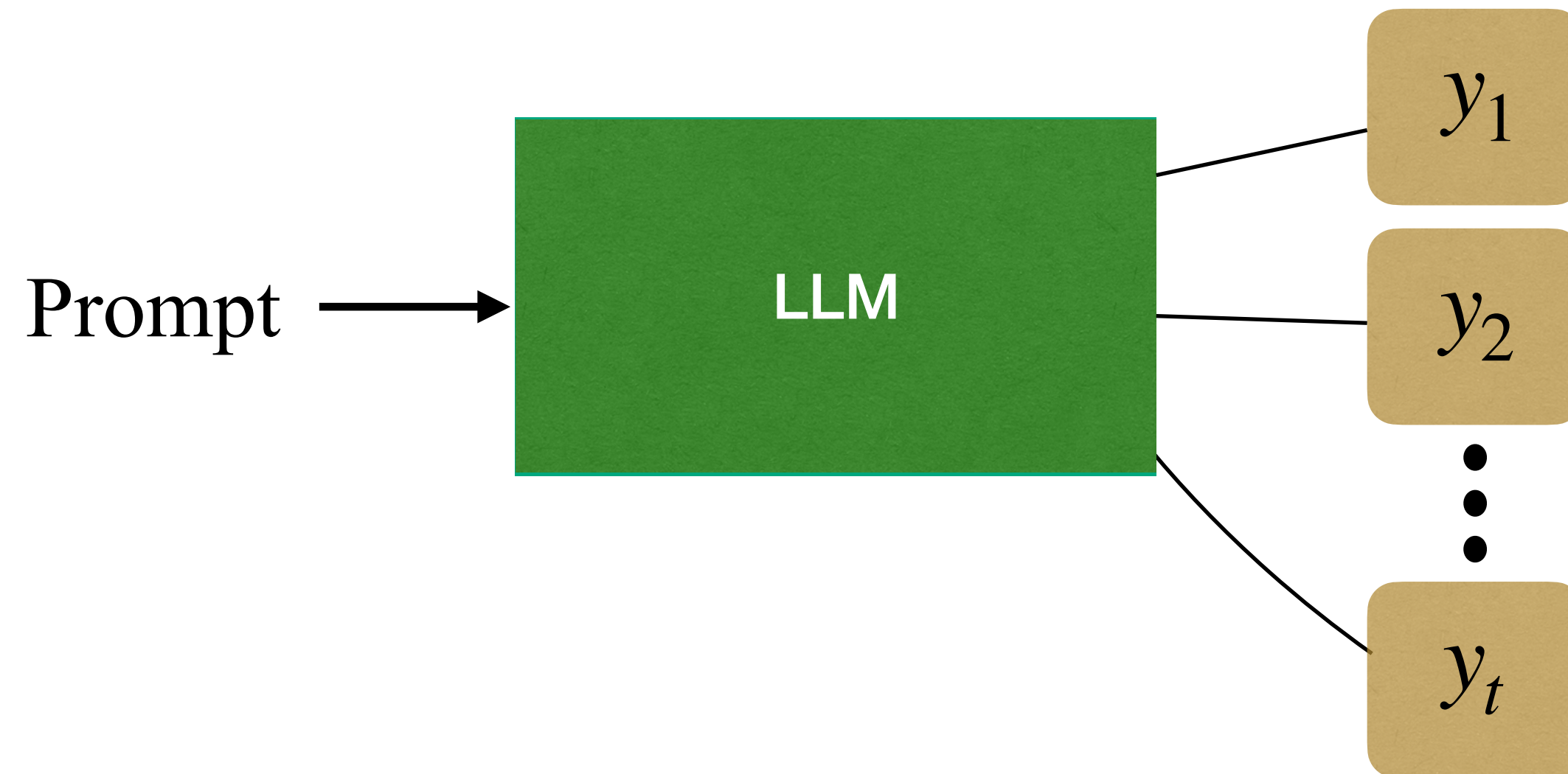


Conformal Prediction for Generative Models (Setting)

$$(X, Y) \in \mathcal{X} \times \mathcal{Y}$$

$$(x, y) \sim p(x, y)$$

$$\text{Query Oracle: } y \sim \pi(y \mid x)$$



We query the oracle for a **finite** amount of times

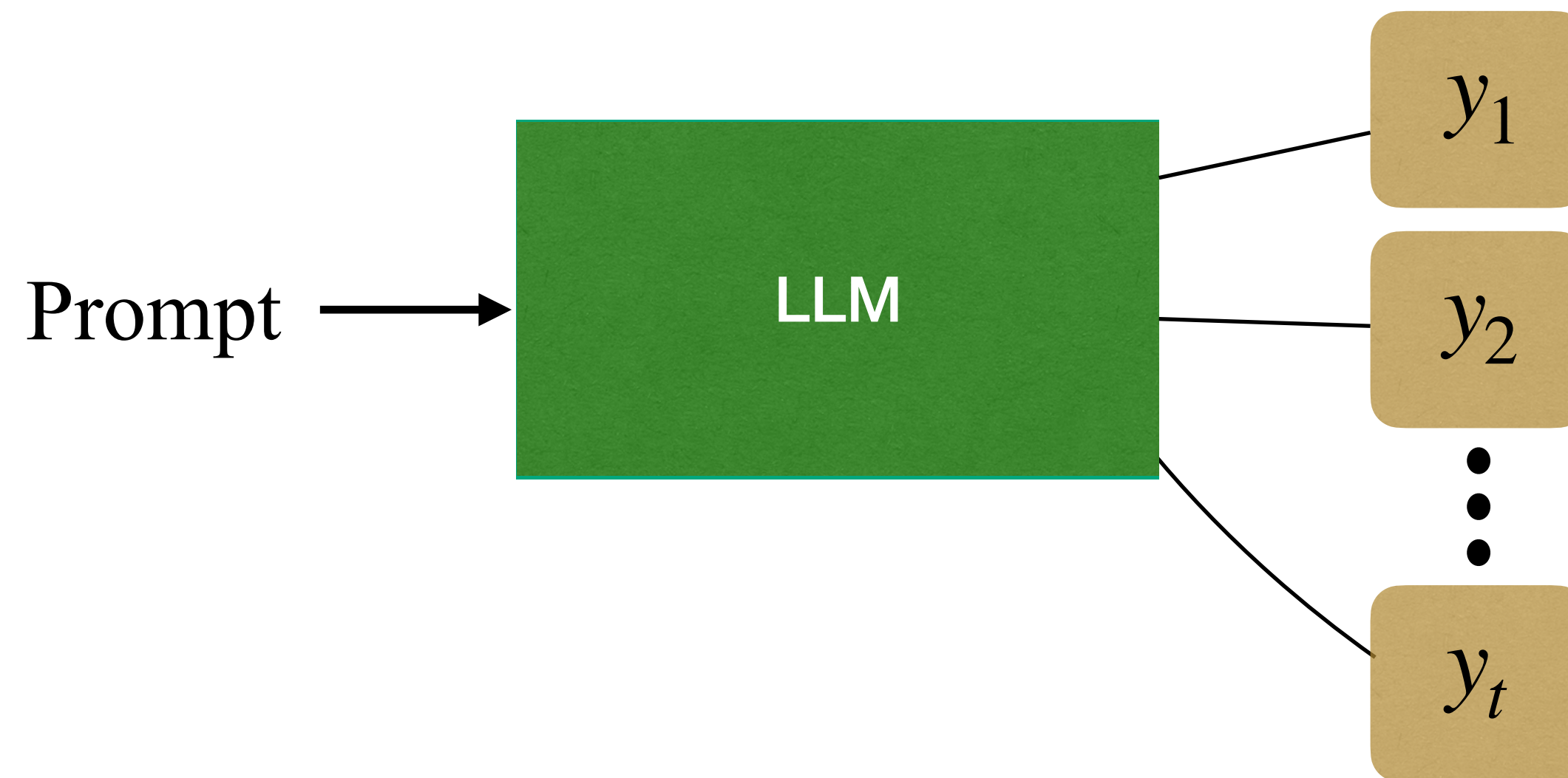
$$Z_t(x) = \{y_1^x, \dots, y_t^x\}$$

Conformal Prediction for Generative Models (Setting)

$$(X, Y) \in \mathcal{X} \times \mathcal{Y}$$

$$(x, y) \sim p(x, y)$$

$$\text{Query Oracle: } y \sim \pi(y \mid x)$$



We query the oracle for a **finite** amount of times

$$Z_t(x) = \{y_1^x, \dots, y_t^x\}$$

Caveat: What if the true label Y is not amount the queried outputs?

“Toy” Running Example

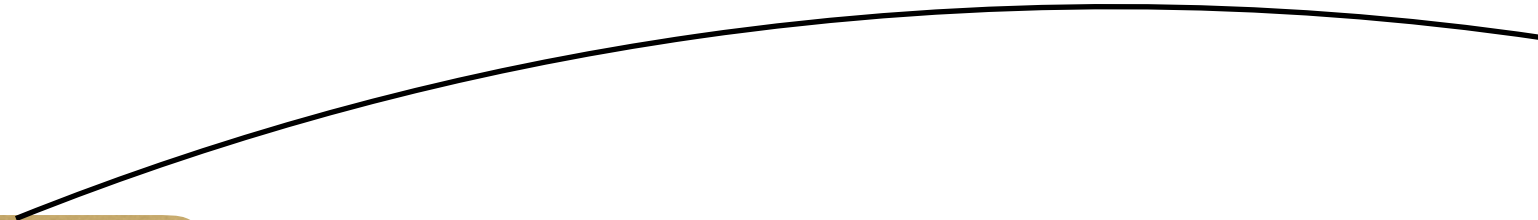
“Toy” Running Example

What is the capital of Illinois



LLM

Springfield



“Toy” Running Example

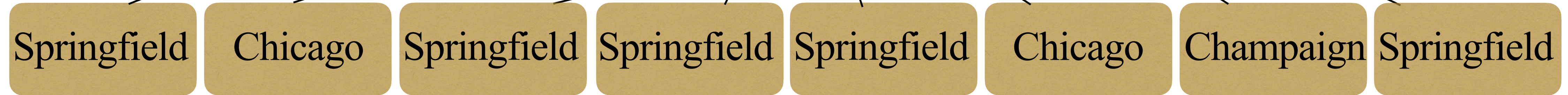
What is the capital of Illinois



Generate 8 responses, $t = 8$



Case 1

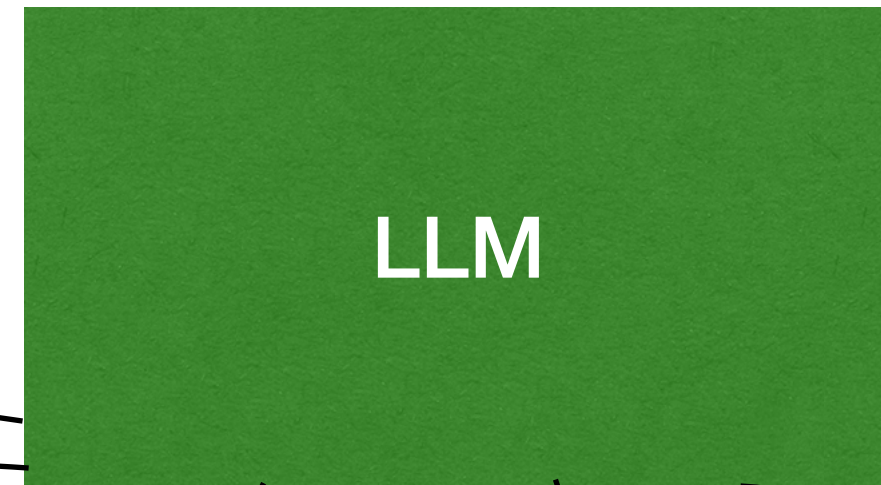


“Toy” Running Example

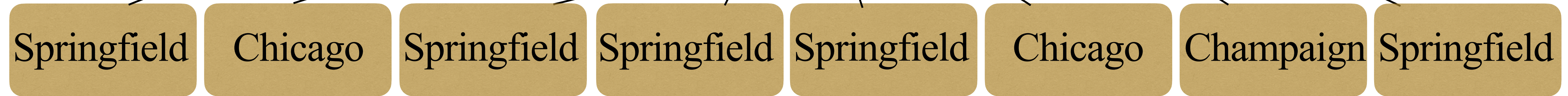
What is the capital of Illinois



Generate 8 responses, $t = 8$



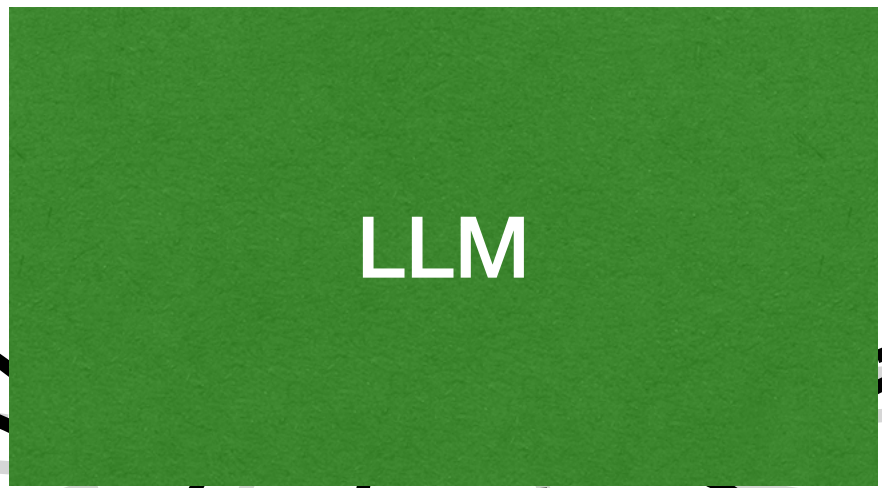
Case 1



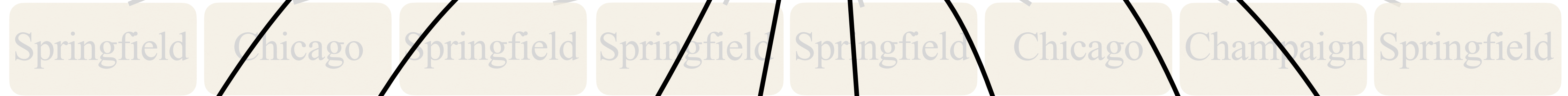
→ $C(x) = \{\text{Springfield, Chicago}\}$

“Toy” Running Example

What is the capital of Illinois



Case 1



Case 2



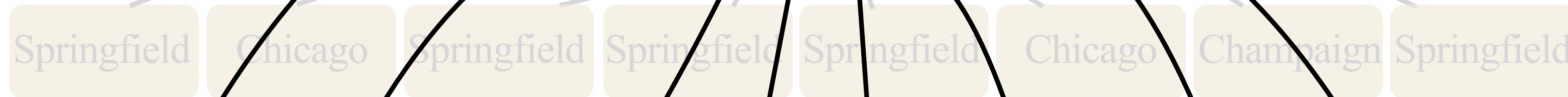
→ $C(x) = \{\text{Springfield, Chicago}\}$

“Toy” Running Example

What is the capital of Illinois



Case 1



Case 2



→ $C(x) = \{\text{Springfield, Chicago}\}$

→ $C(x) = \text{Everything?}$

Everything Else : “EE”

We want to capture the full label space in finite queries

$$Z_t(x) = \{y_1^x, \dots, y_t^x\} \cup \mathcal{Y} \setminus Z_t(x)$$

Everything Else : “EE”

We want to capture the full label space in finite queries

$$Z_t(x) = \{y_1^x, \dots, y_t^x\} \cup \mathcal{Y} \setminus Z_t(x)$$

We have not seen!

Everything Else : “EE”

We want to capture the full label space in finite queries

$$Z_t(x) = \{y_1^x, \dots, y_t^x\} \cup \text{EE}$$

“Everything Else” in the output space

Everything Else : “EE”

We want to capture the full label space in finite queries

$$Z_t(x) = \{y_1^x, \dots, y_t^x\} \cup \text{EE}$$

$$C(x) \subseteq Z_t(x) \cup \text{EE}$$

Everything Else : “EE”

We want to capture the full label space in finite queries

$$C(x) \subseteq Z_t(x) \cup EE$$

If $EE \in C(x) \rightarrow$ **Covered**

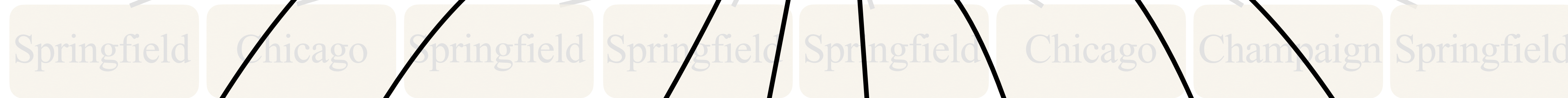
Key challenge: avoid including EE as much as possible

“Toy” Running Example

What is the capital of Illinois



Case 1



Case 2

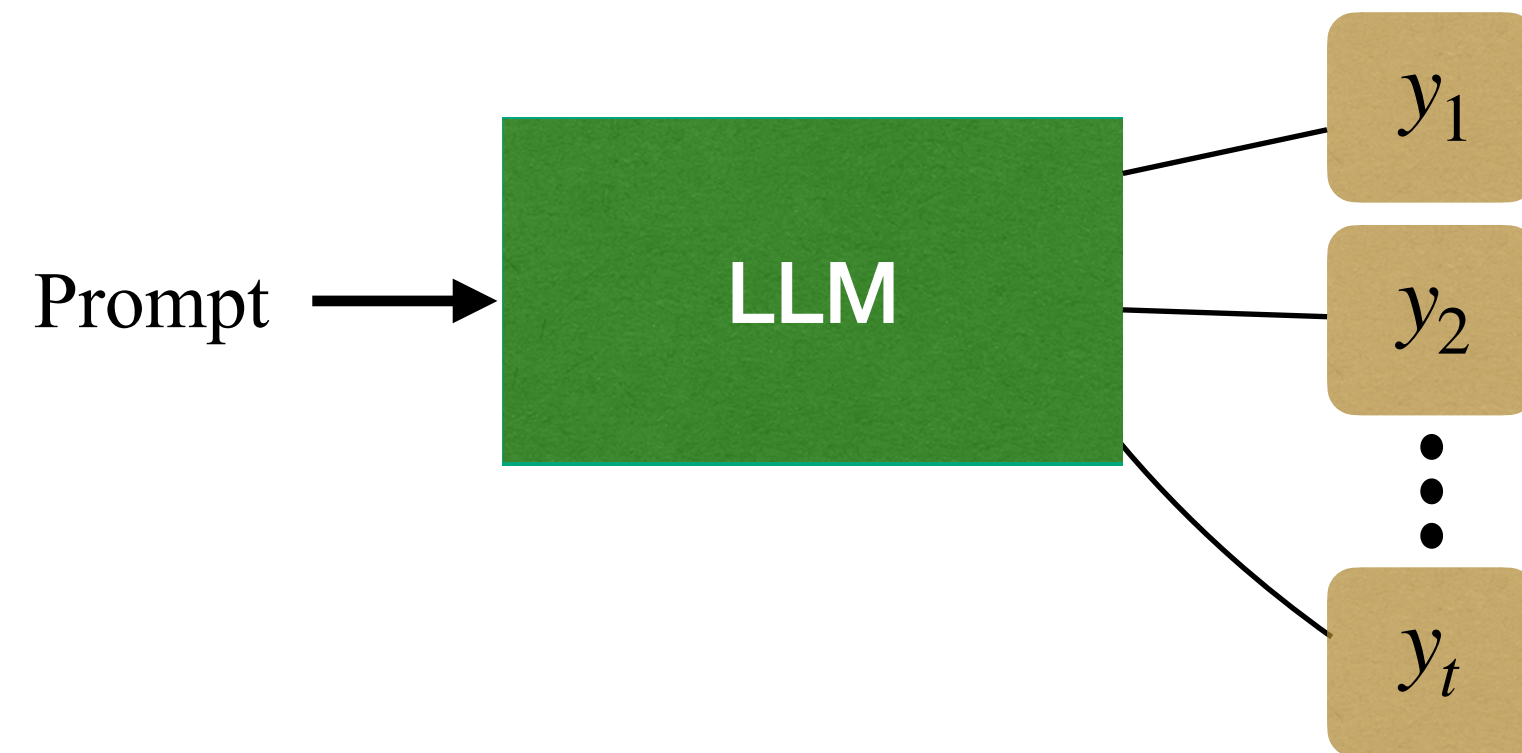


→ $C(x) = \{\text{Springfield, Chicago}\}$

→ $C(x) = \text{Everything}$
 $\{\text{generated responses}\} + \{\text{EE}\}$

So far ...

Our setting:



Trade-offs

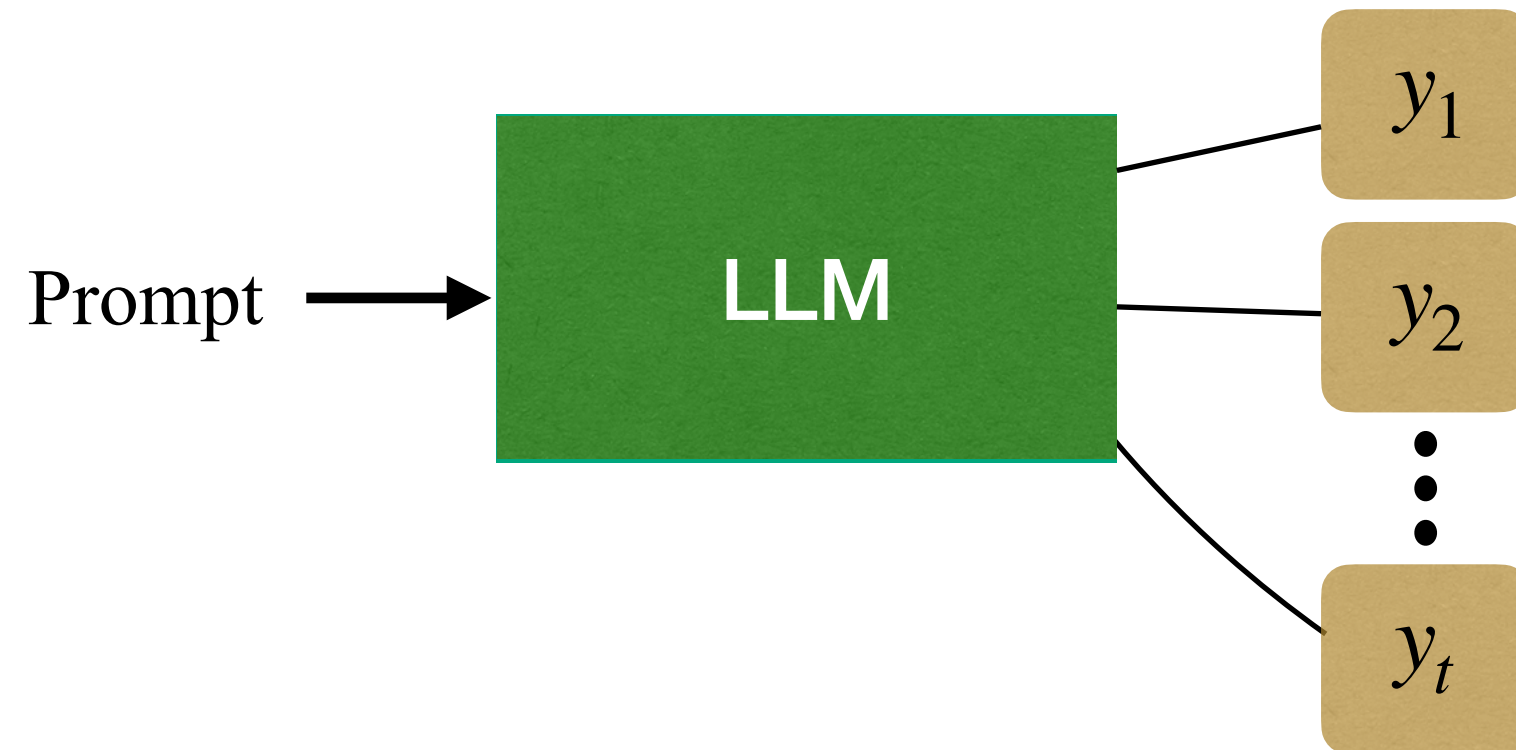
$$Z_t(x) = \{y_1^x, \dots, y_t^x\}$$

$$C(x) \subset Z_t(x) \cup \{\text{EE}\}$$

So far ...

Trade-offs

Our setting:



Finite number of queries to explore the output space (**budget**)

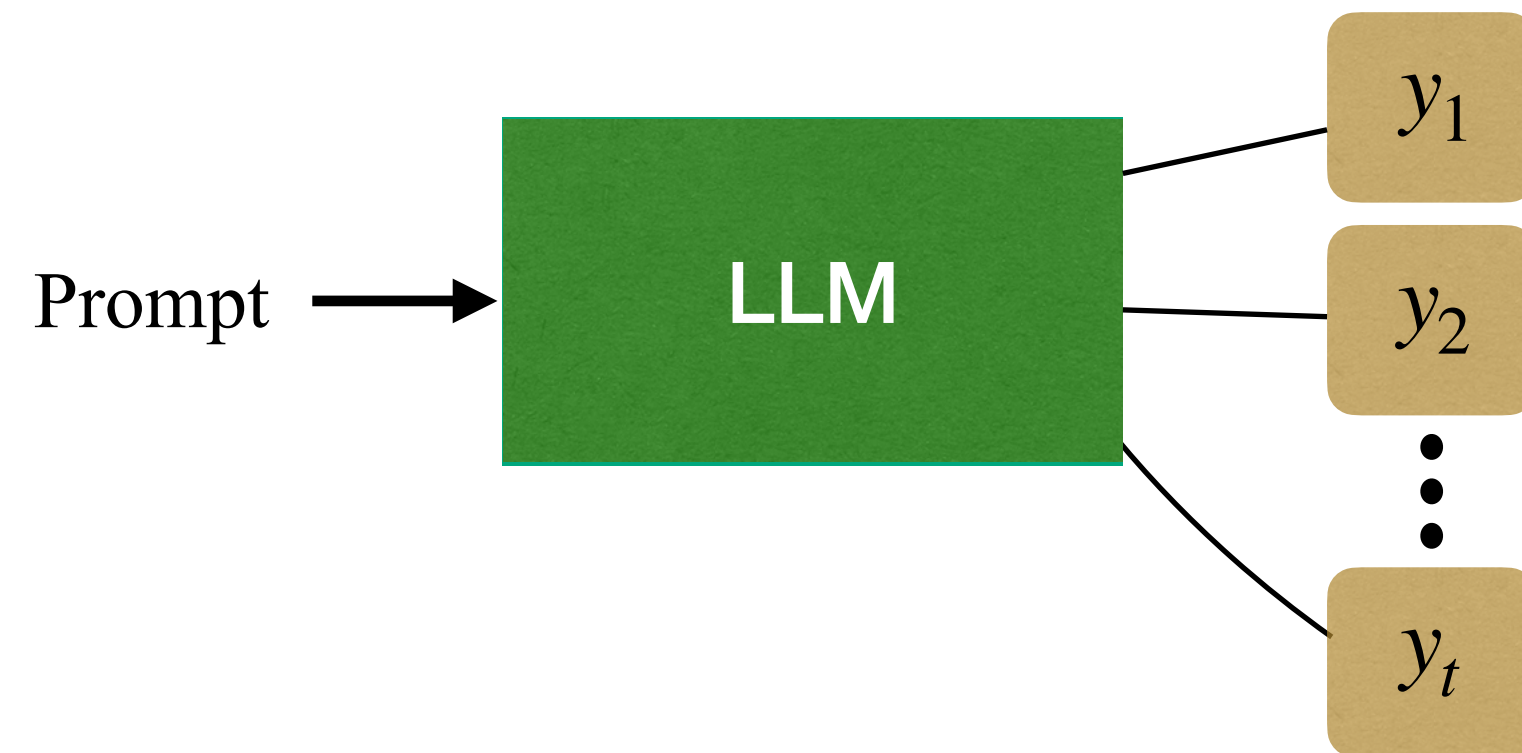
$$Z_t(x) = \{y_1^x, \dots, y_t^x\}$$

$$C(x) \subset Z_t(x) \cup \{\text{EE}\}$$

So far ...

Trade-offs

Our setting:



Finite number of queries to explore the output space (budget)

Avoid EE as much as possible (Informativeness)

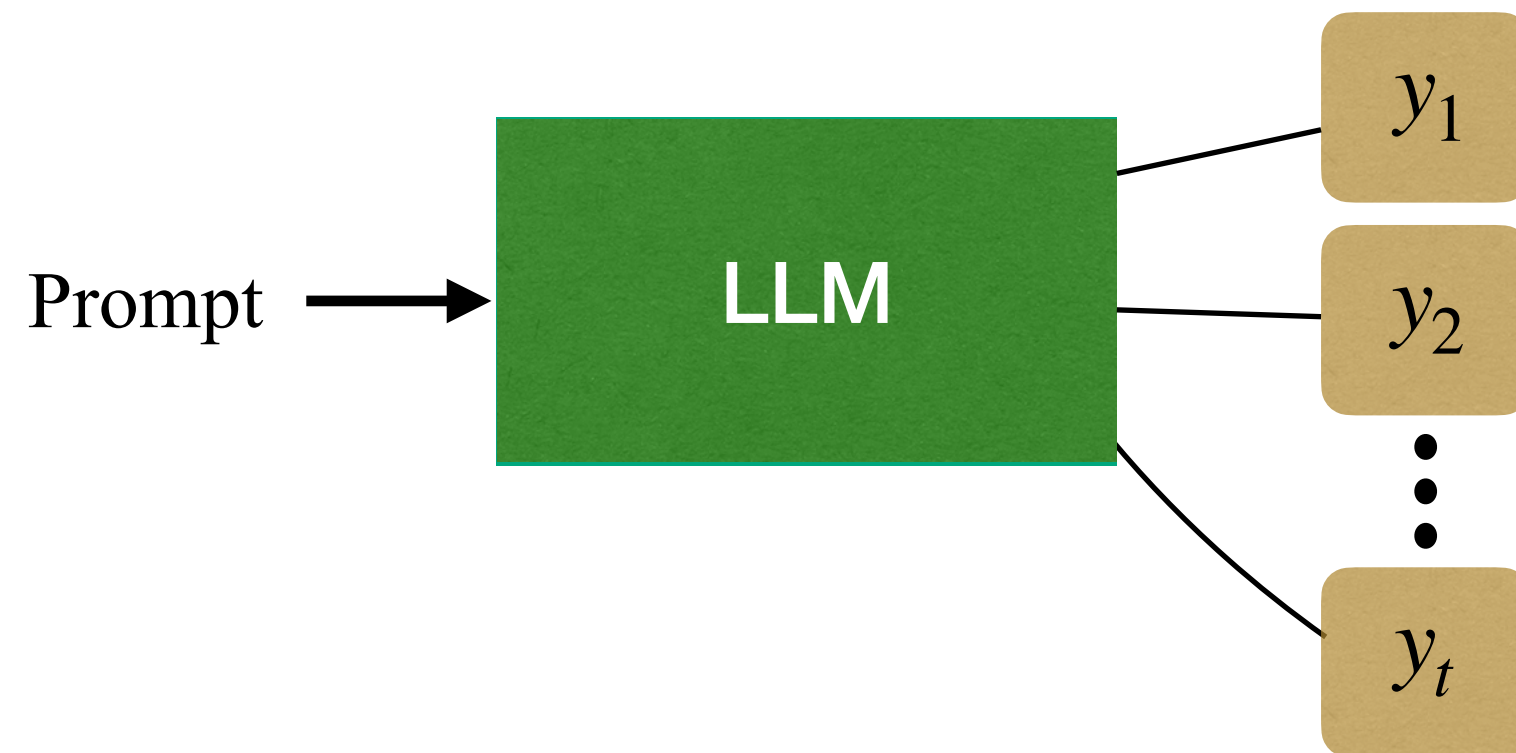
$$Z_t(x) = \{y_1^x, \dots, y_t^x\}$$

$$C(x) \subset Z_t(x) \cup \{\text{EE}\}$$

So far ...

Trade-offs

Our setting:



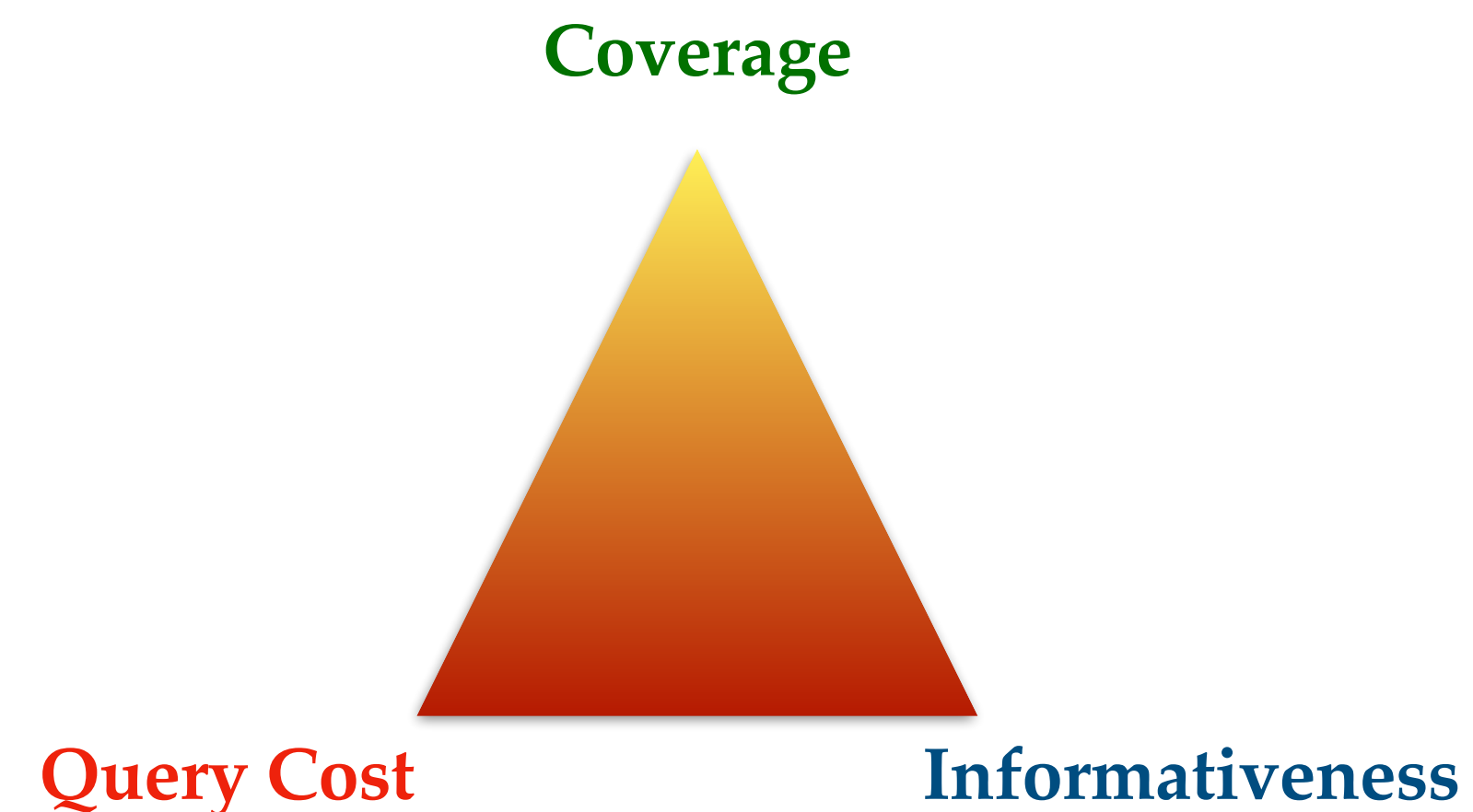
$$Z_t(x) = \{y_1^x, \dots, y_t^x\}$$

$$C(x) \subset Z_t(x) \cup \{\text{EE}\}$$

Finite number of queries to explore the output space (budget)

Avoid EE as much as possible (Informativeness)

Guarantee Coverage (coverage validity)



The Main Optimization Problem (In the Population Regime)

The Main Optimization Problem (In the Population Regime)

Consider **two** main components:

The Main Optimization Problem (In the Population Regime)

Consider **two** main components:

Query policy

$$T : \mathcal{X} \rightarrow \mathbb{N} \cup \{0\}$$

$$Z(T; x) = \{y_1^x, \dots, y_{T(x)}^x\}$$

The Main Optimization Problem (In the Population Regime)

Consider **two** main components:

Query policy

$$T : \mathcal{X} \rightarrow \mathbb{N} \cup \{0\}$$

$$Z(T; x) = \{y_1^x, \dots, y_{T(x)}^x\}$$

Set map

$$f : \mathcal{X} \times 2^{\mathcal{Y}} \rightarrow 2^{\mathcal{Y} \cup EE}$$

$$C(x) = f(x, Z(T; x))$$

The Main Optimization Problem (In the Population Regime)

$$\min_{f, T} \mathbb{E}_X \left[\lambda \mathbf{1}\{EE \in C(X)\} + \sum_{y \neq EE} \mathbf{1}\{y \in C(X)\} \right]$$

$$s.t \quad \Pr_{X, Y} [Y \in C(X)] \geq 1 - \alpha$$

$$\mathbb{E}_X [T(X)] \leq B$$

The Main Optimization Problem (In the Population Regime)

$$\begin{aligned} \min_{f, T} \quad & \mathbb{E}_X \left[\lambda \mathbf{1}\{EE \in C(X)\} + \sum_{y \neq EE} \mathbf{1}\{y \in C(X)\} \right] \\ \text{s.t.} \quad & \Pr_{X, Y} [Y \in C(X)] \geq 1 - \alpha \quad \text{Coverage} \\ & \mathbb{E}_X [T(X)] \leq B \end{aligned}$$

The Main Optimization Problem (In the Population Regime)

$$\begin{aligned} \min_{f, T} \quad & \mathbb{E}_X \left[\lambda \mathbf{1}\{EE \in C(X)\} + \sum_{y \neq EE} \mathbf{1}\{y \in C(X)\} \right] \\ \text{s.t.} \quad & \Pr_{X, Y} [Y \in C(X)] \geq 1 - \alpha \quad \text{Coverage} \\ & \mathbb{E}_X [T(X)] \leq B \quad \text{Query budget} \end{aligned}$$

The Main Optimization Problem (In the Population Regime)

$$\begin{aligned} \min_{f, T} \quad & \mathbb{E}_X \left[\begin{array}{c} \text{Informativeness} \\ \lambda \mathbf{1}\{EE \in C(X)\} + \sum_{y \neq EE} \mathbf{1}\{y \in C(X)\} \end{array} \right] \\ \text{s.t.} \quad & \Pr_{X, Y} [Y \in C(X)] \geq 1 - \alpha \quad \text{Coverage} \\ & \mathbb{E}_X [T(X)] \leq B \quad \text{Query budget} \end{aligned}$$

The Main Optimization Problem (In the Population Regime)

$$\begin{aligned} \min_{f, T} \quad & \mathbb{E}_X \left[\begin{array}{c} \text{Informativeness} \\ \lambda \mathbf{1}\{\text{EE} \in C(X)\} + \sum_{y \neq \text{EE}} \mathbf{1}\{y \in C(X)\} \end{array} \right] \\ \text{s.t.} \quad & \Pr_{X, Y} [Y \in C(X)] \geq 1 - \alpha \quad \text{Coverage} \\ & \mathbb{E}_X [T(X)] \leq B \quad \text{Query budget} \end{aligned}$$

The Main Optimization Problem (In the Population Regime)

$$\begin{aligned} \min_{f, T} \quad & \mathbb{E}_X \left[\lambda \mathbf{1}\{EE \in C(X)\} + \sum_{y \neq EE} \mathbf{1}\{y \in C(X)\} \right] \\ \text{s.t.} \quad & \Pr_{X, Y} [Y \in C(X)] \geq 1 - \alpha \\ & \mathbb{E}_X [T(X)] \leq B \end{aligned}$$

The Main Optimization Problem (In the Population Regime)

$$\begin{aligned} \min_{f, T} \quad & \mathbb{E}_X \left[\lambda \mathbf{1}\{EE \in C(X)\} + \sum_{y \neq EE} \mathbf{1}\{y \in C(X)\} \right] \\ \text{s.t.} \quad & \Pr_{X, Y} [Y \in C(X)] \geq 1 - \alpha \\ & \mathbb{E}_X [T(X)] \leq B \end{aligned}$$

1st step: Solve in the population regime (assume access to $p(y | x)$)

The Main Optimization Problem (In the Population Regime)

$$\begin{aligned} \min_{f, T} \quad & \mathbb{E}_X \left[\lambda \mathbf{1}\{EE \in C(X)\} + \sum_{y \neq EE} \mathbf{1}\{y \in C(X)\} \right] \\ \text{s.t.} \quad & \Pr_{X, Y} [Y \in C(X)] \geq 1 - \alpha \\ & \mathbb{E}_X [T(X)] \leq B \end{aligned}$$

1st step: Solve in the population regime (assume access to $p(y | x)$)

2nd step: mimic optimal solution in finite sample regime
(assume access to **oracle** $\pi(y | x)$)
(assume access to **finite calibration data**)

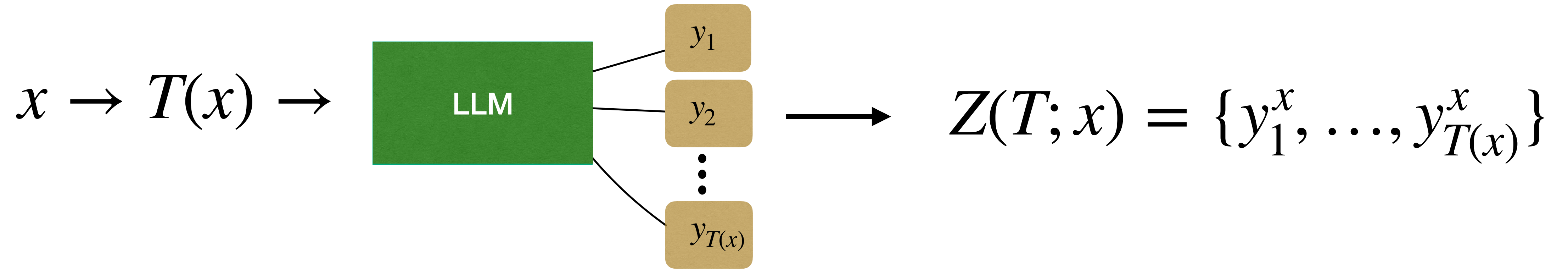
1st step: Solve in the population regime (assume access to $p(y | x)$)

$$\begin{aligned} \min_{f, T} \quad & \mathbb{E}_X \left[\lambda \mathbf{1}\{EE \in C(X)\} + \sum_{y \neq EE} \mathbf{1}\{y \in C(X)\} \right] \\ \text{s.t.} \quad & \Pr_{X, Y} [Y \in C(X)] \geq 1 - \alpha \\ & \mathbb{E}_X [T(X)] \leq B \end{aligned}$$

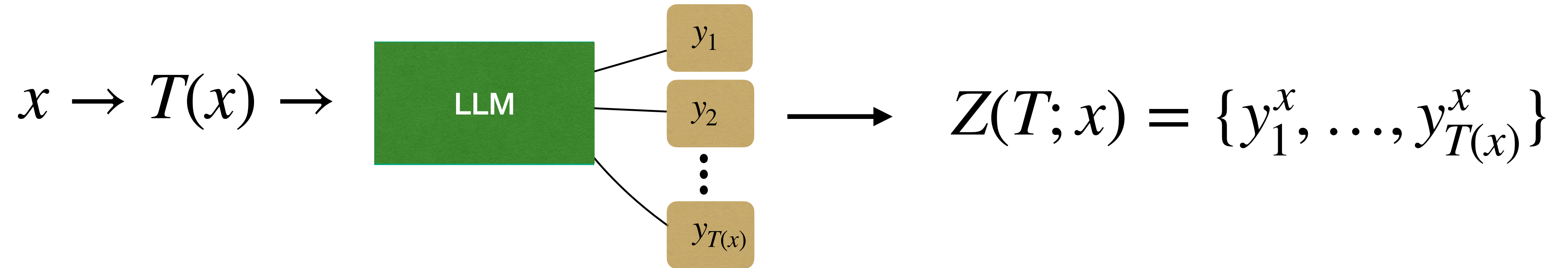
Let's ask:

What is the role of the query policy?

What is the role of the query policy?

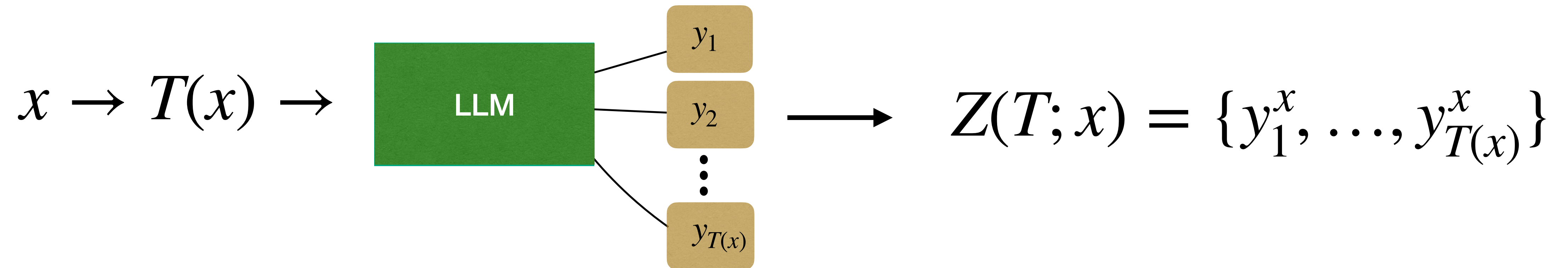


What is the role of the query policy?



Need to sample enough to not miss the correct label !

What is the role of the query policy?

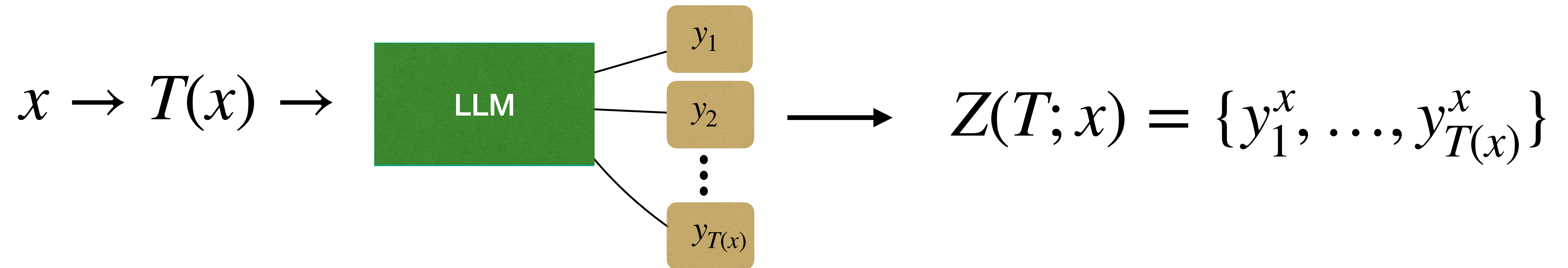


Need to sample enough to not miss the correct label !

But ...

Budget constraint: $\mathbb{E}_X[T(X)] \leq B$

What is the role of the query policy?



Need to sample enough to not miss the correct label !

$$\mathbb{E}_X[T(X)] \leq B$$

Intuition: query policy should be designed to minimize the chance of missing the correct label

A decoupled analysis - two principles

A decoupled analysis - two principles

Question 1: What is the optimal query policy to minimize the change of missing the correct label?

Query policy

$$T : \mathcal{X} \rightarrow \mathbb{N} \cup \{0\}$$
$$Z(T; x) = \{y_1^x, \dots, y_{T(x)}^x\}$$

A decoupled analysis - two principles

Question 1: What is the optimal query policy to minimize the change of missing the correct label?



Fix $T(\cdot)$

A decoupled analysis - two principles

Question 1: What is the optimal query policy to minimize the change of missing the correct label?

↓
Fix $T(\cdot)$
↓

Question 2: What is the optimal set map for constructing valid, informative sets $C(\cdot)$?

Set map

$$f: \mathcal{X} \times 2^{\mathcal{Y}} \rightarrow 2^{\mathcal{Y} \cup EE}$$

$$C(x) = f(x, Z(T; x))$$

A decoupled analysis - two principles

Question 1: What is the optimal query policy to minimize the change of missing the correct label?

Fix $T(\cdot)$

Question 2: What is the optimal set map for constructing valid, informative sets $C(\cdot)$?

Question 1: What is the optimal query policy to minimize the change of missing the correct label?

Principle 1: Optimal Query Policy

Goal: minimize the chance of missing the true label, under finite budget

$$Z_t(x) = \{y_1^x, \dots, y_t^x\}$$

Principle 1: Optimal Query Policy

Goal: minimize the chance of missing the true label, under finite budget

$$Z_t(x) = \{y_1^x, \dots, y_t^x\}$$

$$\theta(x, t) = \mathbb{P}_{Y, Z_t(x)}[Y \notin Z_t(x) \mid X = x]$$

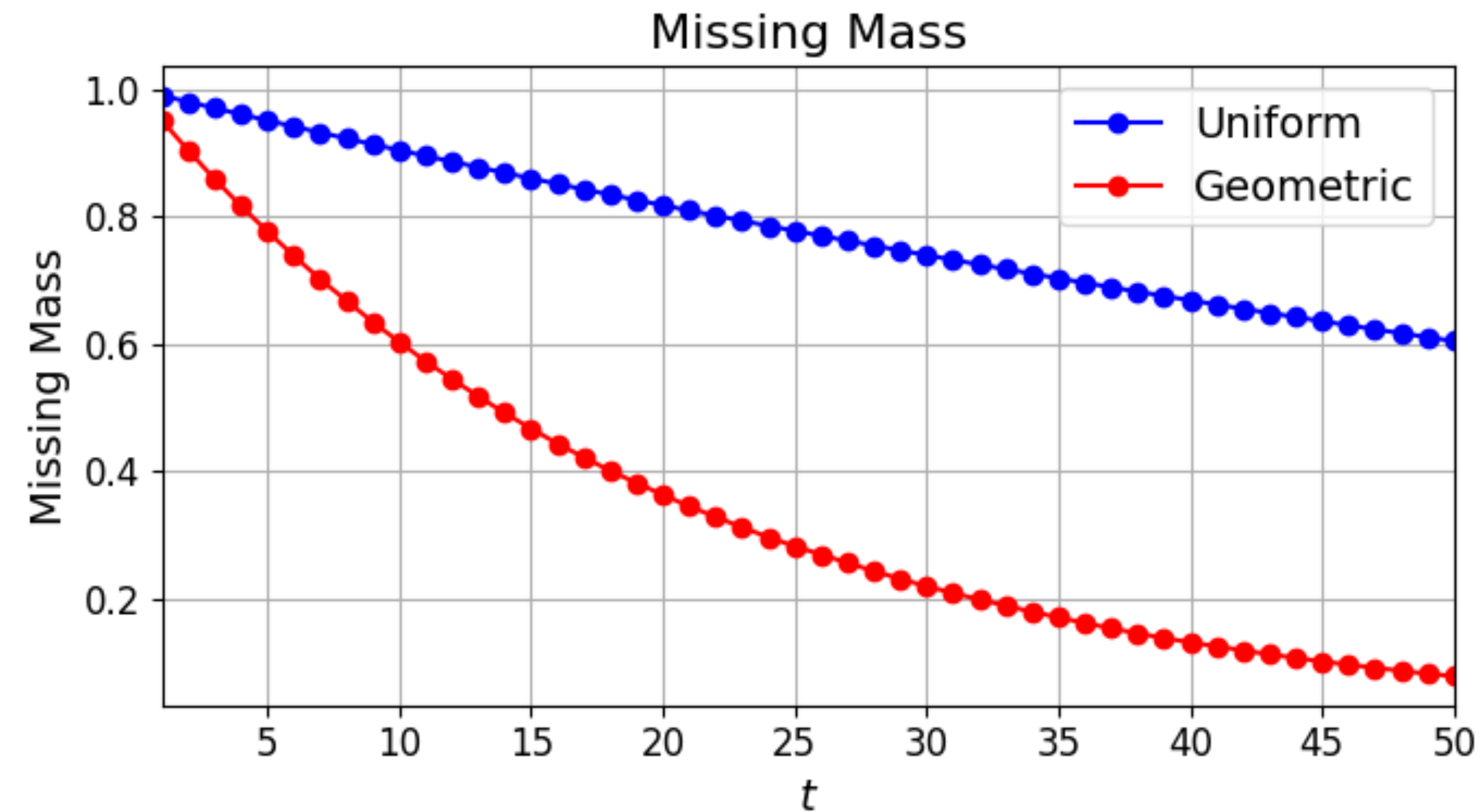
Missing Mass

Principle 1: Optimal Query Policy

Missing Mass : a natural objective

Missing Mass

$$\theta(x, t) = \theta(t) = \mathbb{P}_{Y, Z_t(x)}[Y \notin Z_t(x) \mid X = x]$$

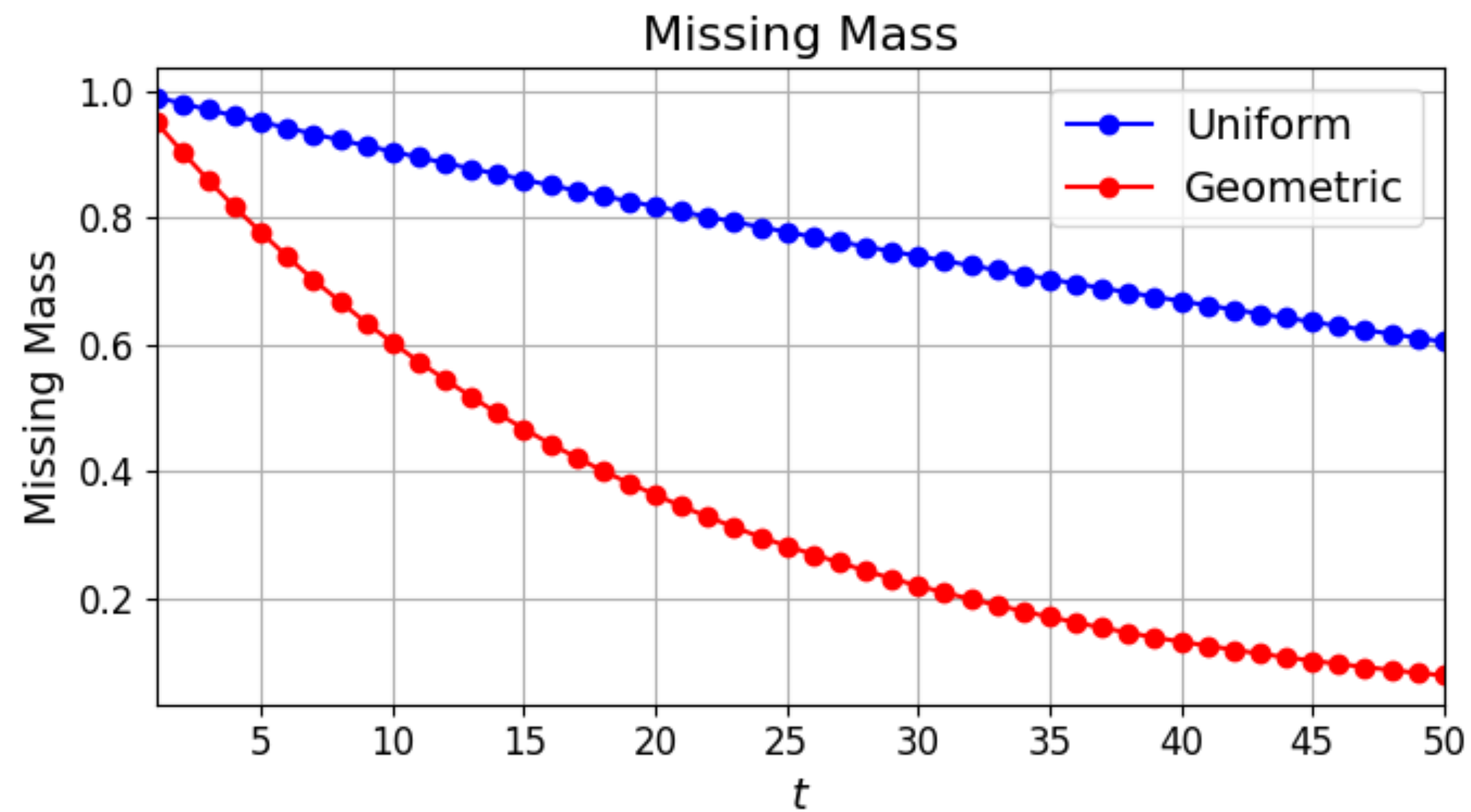


Principle 1: Optimal Query Policy

Missing Mass : a natural objective

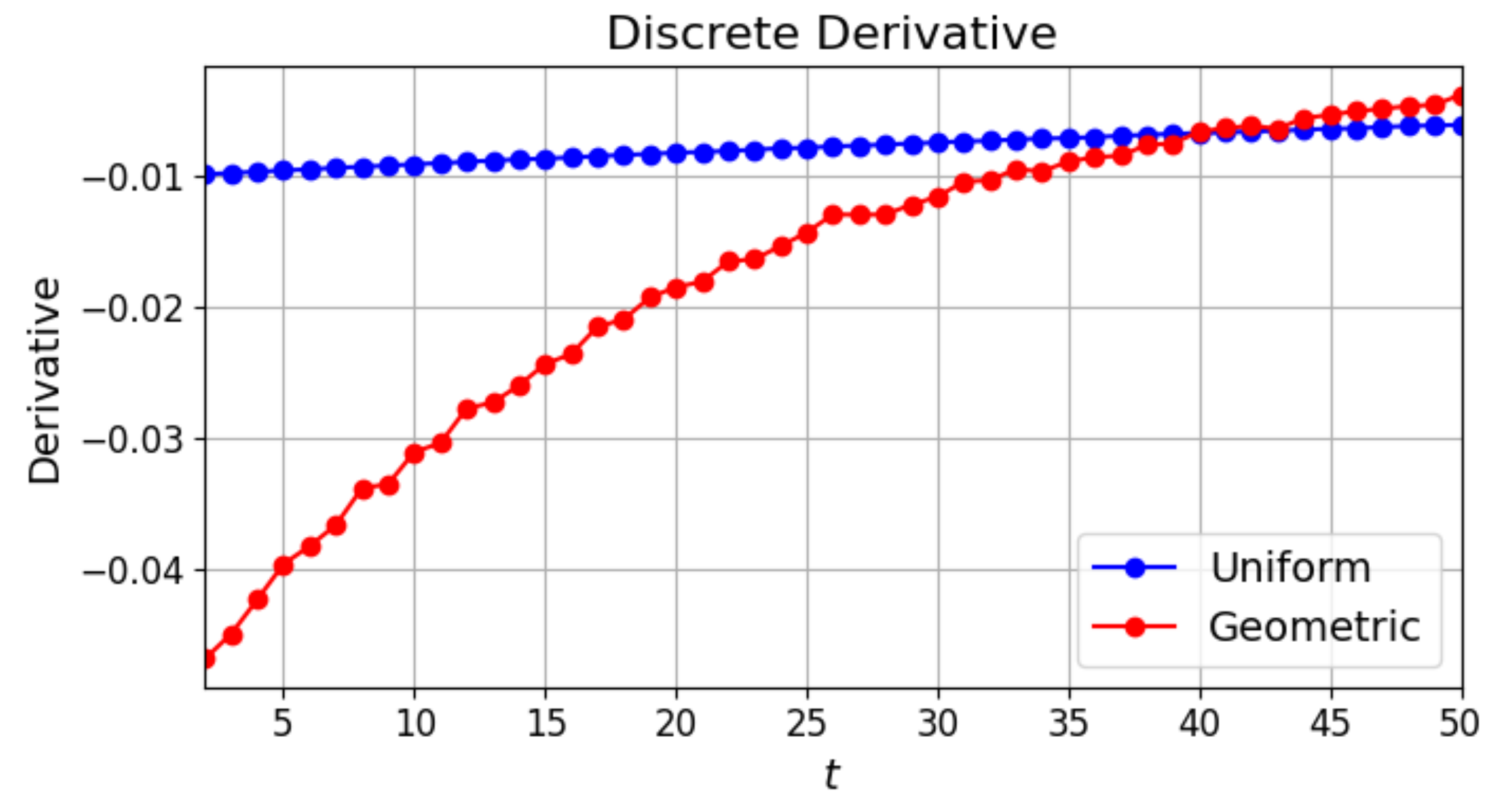
Missing Mass

$$\theta(x, t) = \theta(t) = \mathbb{P}_{Y, Z_t(x)}[Y \notin Z_t(x) \mid X = x]$$



Discrete Derivative

$$\Delta(x, t) := \theta(x, t + 1) - \theta(x, t)$$



Principle 1: Optimal Query Policy

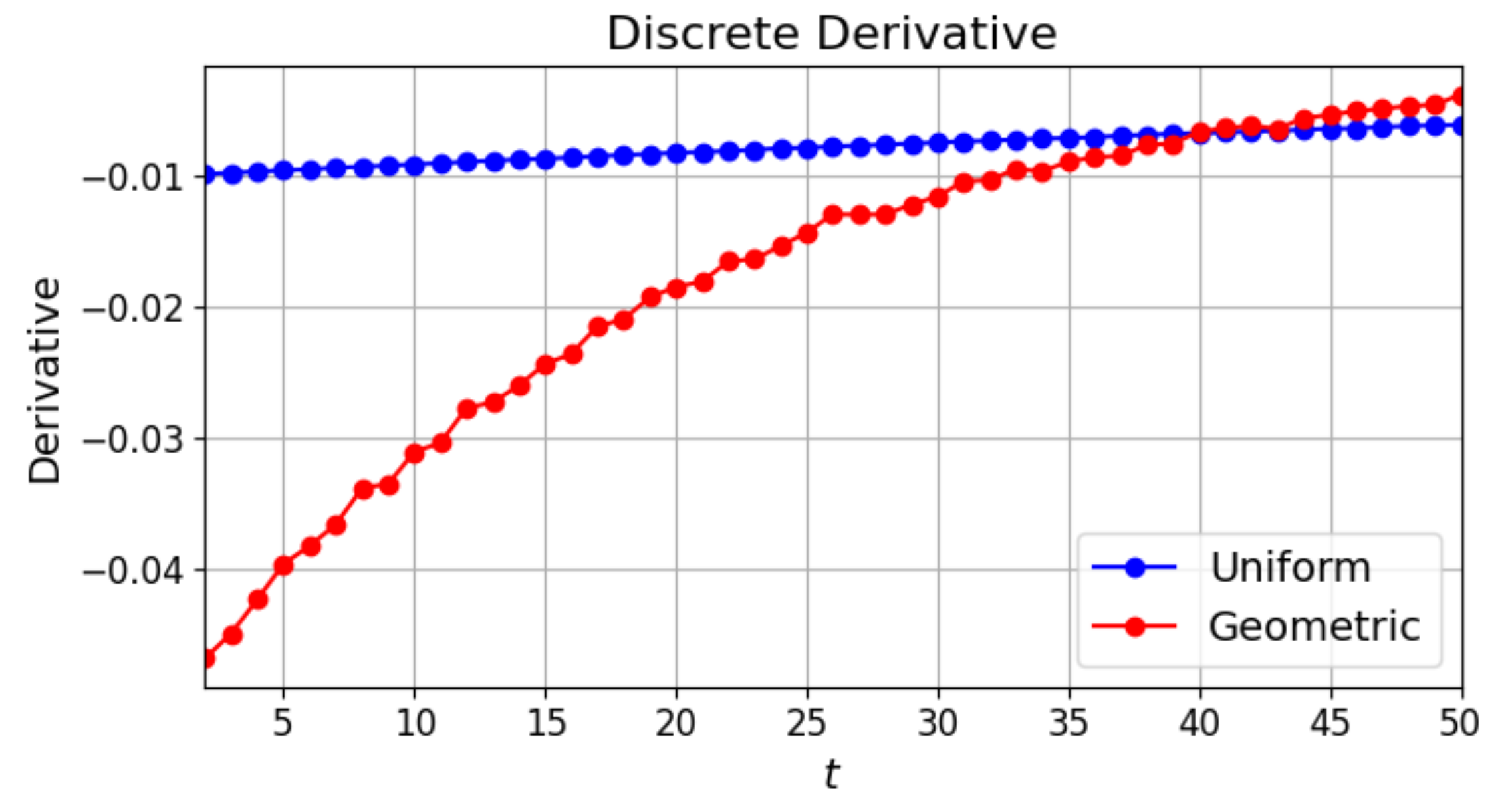
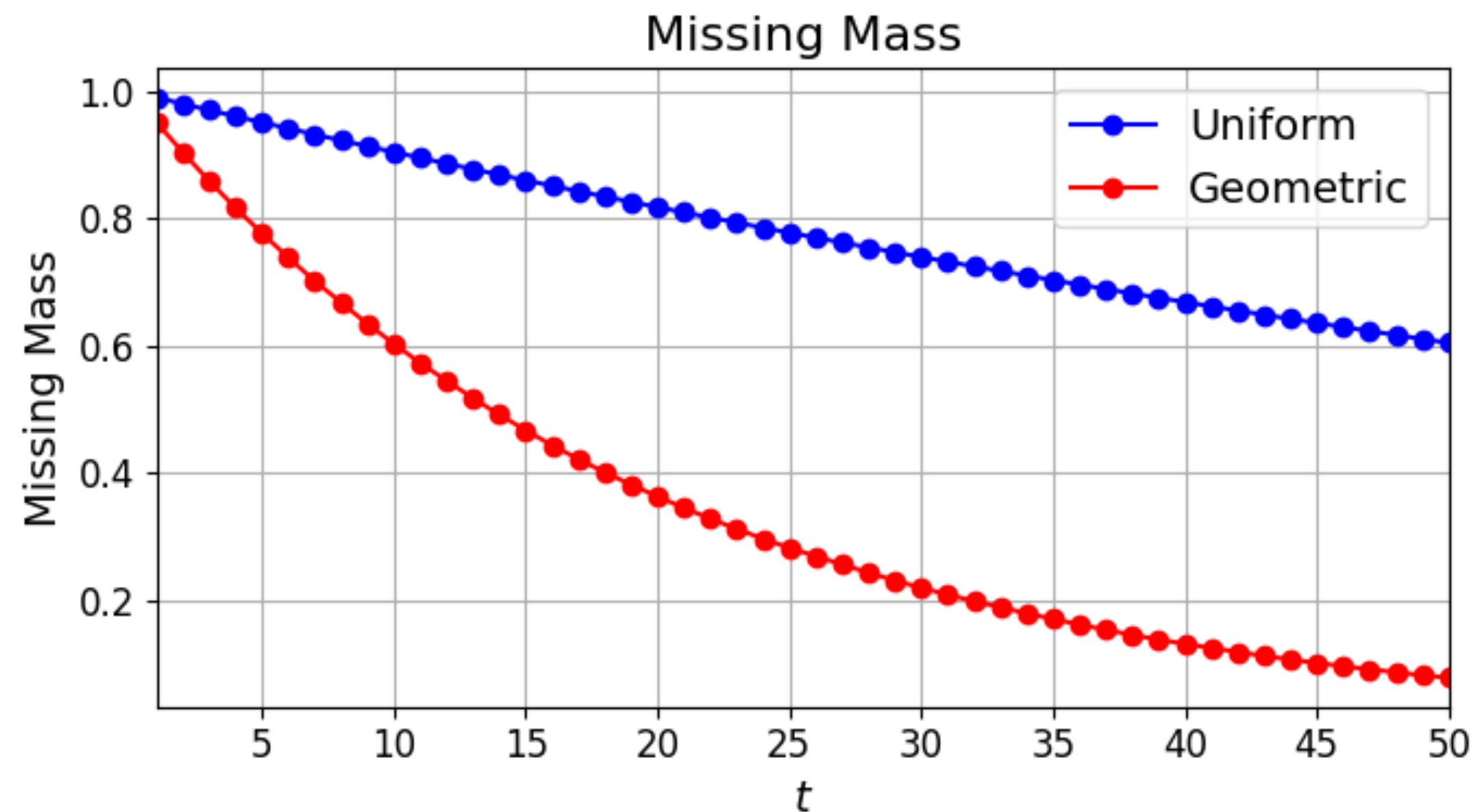
Missing Mass : a natural objective

Missing Mass

$$\theta(x, t) = \theta(t) = \mathbb{P}_{Y, Z_t(x)}[Y \notin Z_t(x) \mid X = x]$$

Discrete Derivative

$$\Delta(x, t) := \theta(x, t + 1) - \theta(x, t)$$



Missing mass decreases in t, with diminishing returns

Principle 1: Optimal Query Policy

$$\begin{aligned} \min_{T(\cdot): \mathcal{X} \rightarrow \mathbb{N} \cup \{0\}} & \quad \mathbb{E}_X[\theta(X, T(X))] \\ \text{s.t.} & \quad \mathbb{E}_X[T(X)] \leq B \end{aligned}$$

Principle 1: Optimal Query Policy

$$\begin{aligned} \min_{T(\cdot): \mathcal{X} \rightarrow \mathbb{N} \cup \{0\}} & \quad \mathbb{E}_X[\theta(X, T(X))] \\ \text{s.t.} & \quad \mathbb{E}_X[T(X)] \leq B \end{aligned}$$

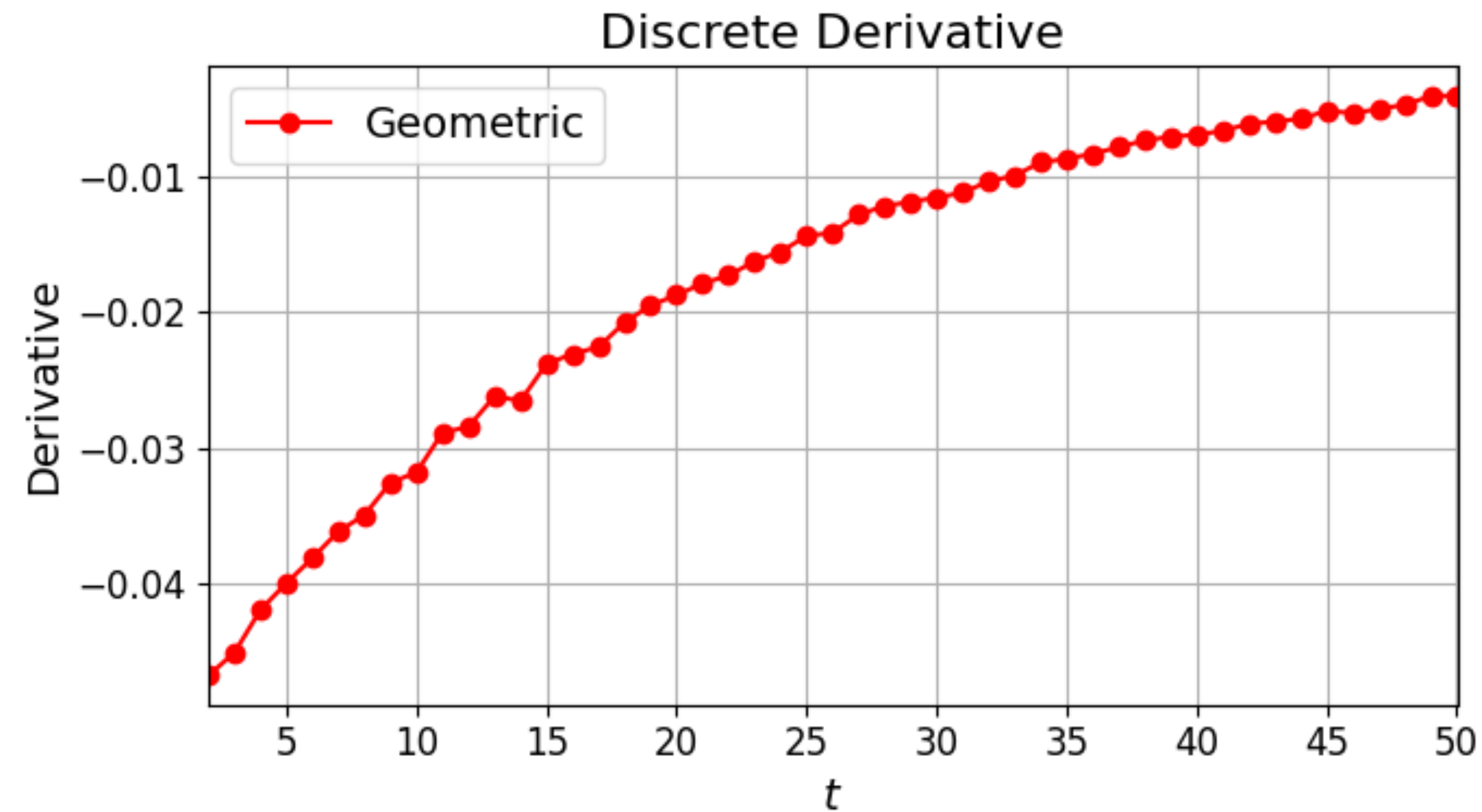
Theorem: Let $T^*(x)$ be the optimal solution. There exists $\beta^* \in \mathbb{R}$ such that for all x

$$\Delta(x, T^*(x)) = \beta^*$$

Principle 1: Optimal Query Policy

For each $x \in \mathcal{X}$, consider the **derivative of the missing mass**:

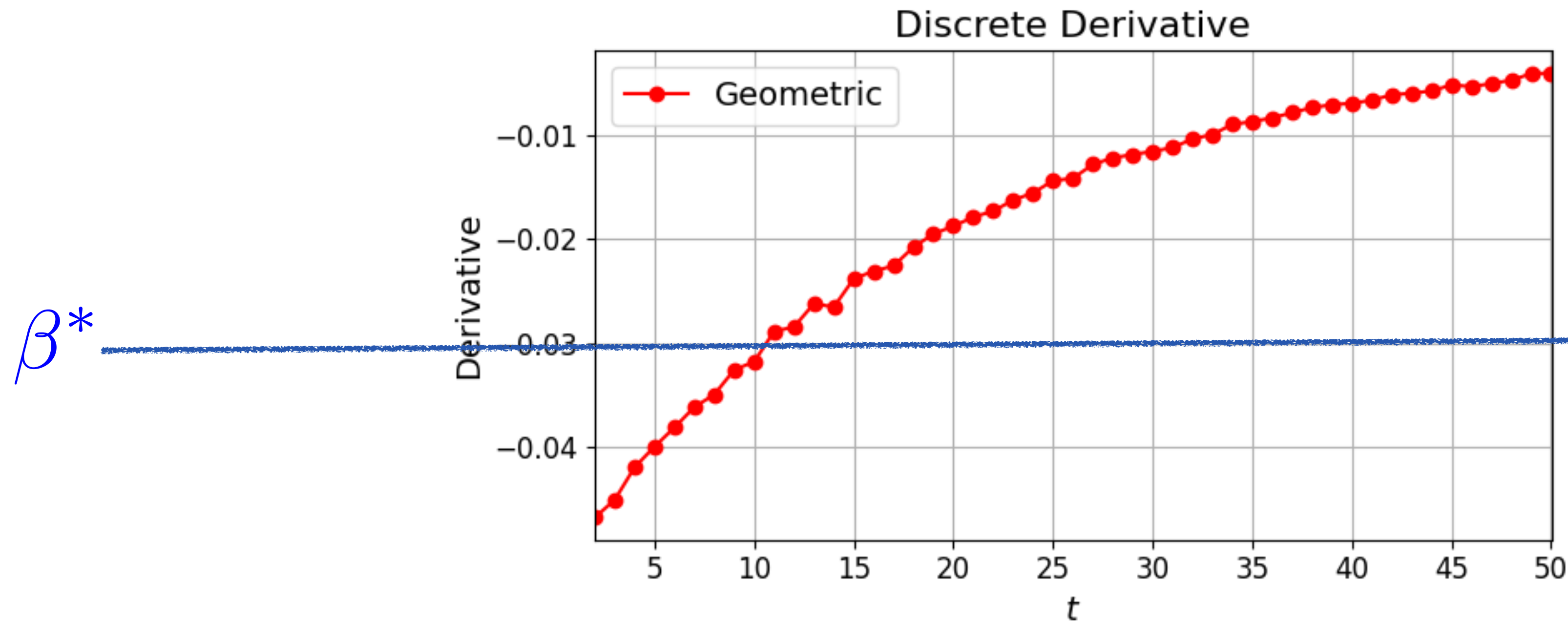
$$\Delta(x, t) := \theta(x, t + 1) - \theta(x, t)$$



Principle 1: Optimal Query Policy

For each $x \in \mathcal{X}$, consider the **derivative of the missing mass**:

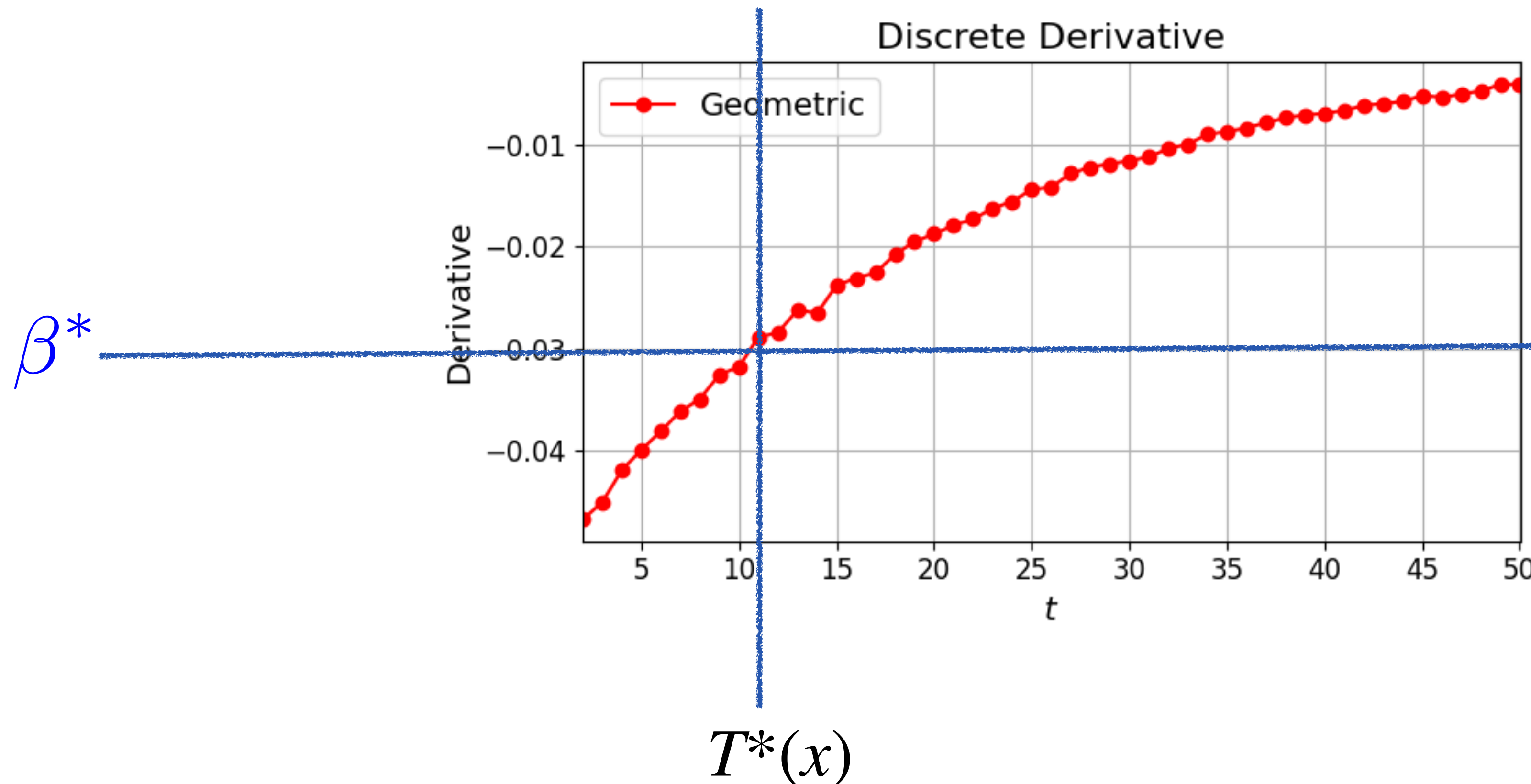
$$\Delta(x, t) := \theta(x, t + 1) - \theta(x, t)$$



Principle 1: Optimal Query Policy

For each $x \in \mathcal{X}$, consider the **derivative of the missing mass**:

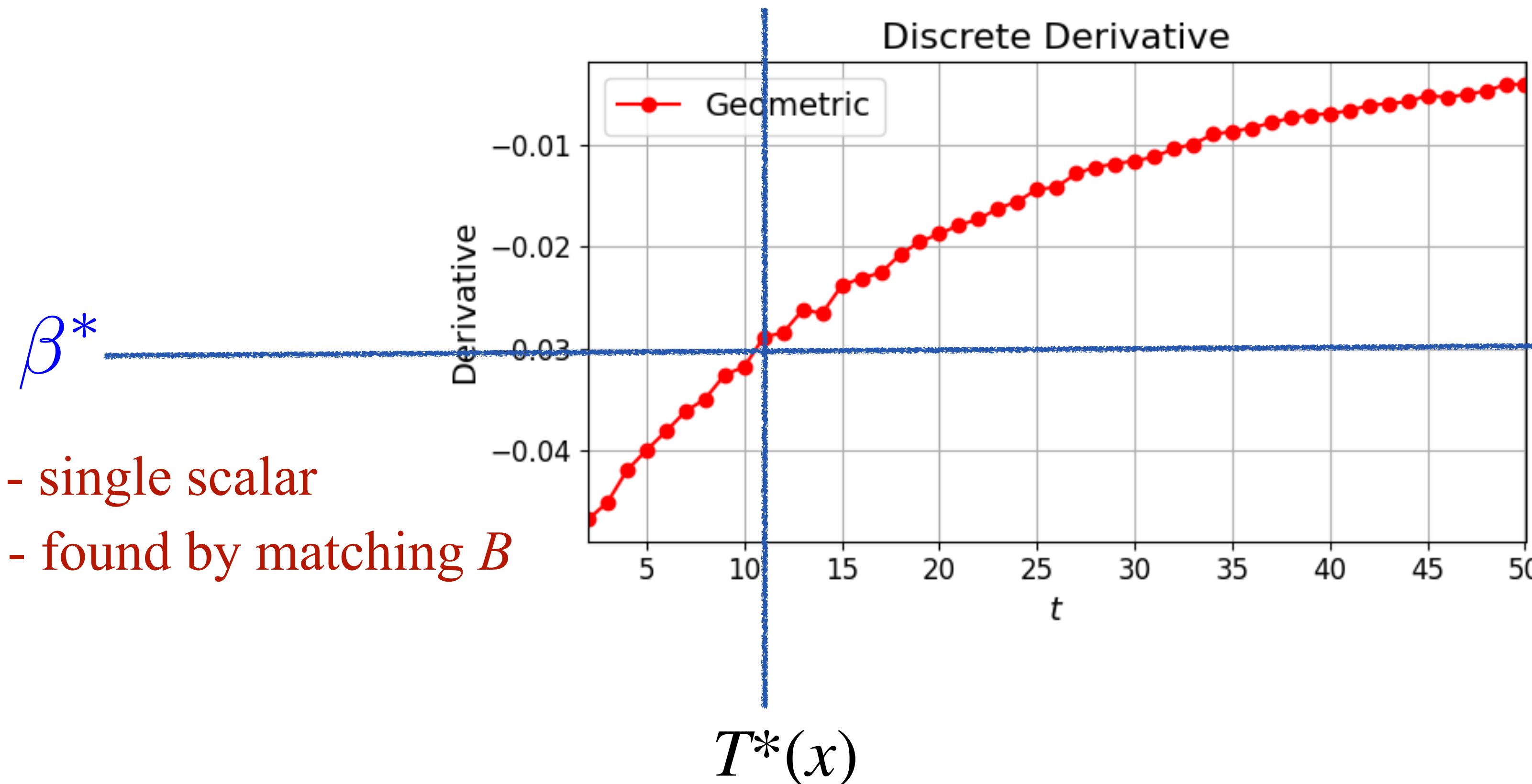
$$\Delta(x, t) := \theta(x, t + 1) - \theta(x, t)$$



Principle 1: Optimal Query Policy

For each $x \in \mathcal{X}$, consider the **derivative of the missing mass**:

$$\Delta(x, t) := \theta(x, t + 1) - \theta(x, t)$$



A decoupled analysis - two principles



Fix $T(\cdot)$

Question 2: What is the optimal set map for constructing valid, informative sets $C(\cdot)$?

Question 2: What is the optimal set map for constructing valid, informative sets $C(\cdot)$?

Principle 2: Optimal Set map

Queried samples

$$Z(x) = \{y_1^x, \dots, y_{T(x)}^x\}$$

Set map

$$f: \mathcal{X} \times 2^{\mathcal{Y}} \rightarrow 2^{\mathcal{Y} \cup EE}$$

Principle 2: Optimal Set map

Queried samples

$$Z(x) = \{y_1^x, \dots, y_{T(x)}^x\}$$

Set map

$$f: \mathcal{X} \times 2^{\mathcal{Y}} \rightarrow 2^{\mathcal{Y} \cup \mathcal{E}\mathcal{E}}$$

$$\min_{f(\cdot)} \mathbb{E}_X \left[\lambda \mathbf{1}\{\mathcal{E}\mathcal{E} \in C(X)\} + \sum_{y \neq \mathcal{E}\mathcal{E}} \mathbf{1}\{y \in C(X)\} \right]$$

$$s.t \quad \Pr_{X,Y} [Y \in C(X)] \geq 1 - \alpha$$

Principle 2: Optimal Set map

$$\begin{aligned} \min_{f(\cdot)} \quad & \mathbb{E}_X \left[\lambda \mathbf{1}\{EE \in C(X)\} + \sum_{y \neq EE} \mathbf{1}\{y \in C(X)\} \right] \\ \text{s.t} \quad & \Pr_{X,Y} [Y \in C(X)] \geq 1 - \alpha \end{aligned}$$

Theorem: Let f_λ^* be the optimal solution. There exists a scalar threshold $q^* \in \mathbb{R}^+$ s.t for every x

$$C(x) = \{y \in Z(x) \cup EE : S(x, y) \leq q^*\}$$

where with $p(EE | x) = Pr[Y \notin Z(x) | X = x]$

$$S(x, y) = \begin{cases} 1 - p(y | x), & \text{if } y \neq EE \\ 2 - p(y | x), & \text{if } y = EE \end{cases}$$

In summary ~ over population ... :

Question 1: What is the optimal query policy to minimize the change of missing the correct label?

Fix $T(\cdot)$

Question 2: What is the optimal set map for constructing valid, informative sets $C(\cdot)$?

In summary ~ over population ... :

Assume access to $p(y | x)$

Answer 1:

Find optimal set query by $\Delta(x, T^*(x)) = \beta^*$

In summary ~ over population ... :

Assume access to $p(y | x)$

Answer 1:

Find optimal set query by $\Delta(x, T^*(x)) = \beta^*$





$$Z(x) = \{y_1^x, \dots, y_{T^*(x)}^x\}$$

In summary ~ over population ... :

Assume access to $p(y | x)$

Answer 1: Find optimal set query by $\Delta(x, T^*(x)) = \beta^*$


$$Z(x) = \{y_1^x, \dots, y_{T^*(x)}^x\}$$




Answer 2: Build optimal set $C(x) = \{y \in Z(x) \cup EE : S(x, y) \leq q^*\}$

where: $S(x, y) = \begin{cases} 1 - p(y | x), & \text{if } y \neq EE \\ 2 - p(y | x), & \text{if } y = EE \end{cases}$

In summary ~~over population ...~~: In finite sample regime

Assume access to $p(y | x)$

Answer 1: Find optimal set query by $\Delta(x, T^*(x)) = \beta^*$


$$Z(x) = \{y_1^x, \dots, y_{T^*(x)}^x\}$$


Answer 2: Build optimal set $C(x) = \{y \in Z(x) \cup EE : S(x, y) \leq q^*\}$

where: $S(x, y) = \begin{cases} 1 - p(y | x), & \text{if } y \neq EE \\ 2 - p(y | x), & \text{if } y = EE \end{cases}$

In summary ~~over population ...~~:



In finite sample regime

~~Assume access to $p(y|x)$~~

Need to estimate

Answer 1:

Find optimal set query by $\Delta(x, T^*(x)) = \beta^*$


$$Z(x) = \{y_1^x, \dots, y_{T^*(x)}^x\}$$


Answer 2:

Build optimal set $C(x) = \{y \in Z(x) \cup EE : S(x, y) \leq q^*\}$

where: $S(x, y) = \begin{cases} 1 - p(y|x) & \text{if } y \neq EE \\ 2 - p(y|x) & \text{if } y = EE \end{cases}$

Our Main Algorithm (In the finite-sample regime)

Our Main Algorithm (In the finite-sample regime)

We don't have the true $p(y|x)$

Our Main Algorithm (In the finite-sample regime)

We don't have the true $p(y | x)$

But we can query $\pi(y | x)$

Our Main Algorithm (In the finite-sample regime)

We don't have the true $p(y|x)$

But we can query $\pi(y|x)$

Need to estimate $\theta(x, t)$ and $\Delta(x, t)$

Our Main Algorithm (In the finite-sample regime)

We don't have the true $p(y | x)$

But we can query $\pi(y | x)$

Need to estimate $\theta(x, t)$ and $\Delta(x, t)$

Find the right alternative for $S(x, y) = 1 - p(y | x)$

Our Main Algorithm (In the finite-sample regime)

We don't have the true $p(y | x)$

But we can query $\pi(y | x)$

Need to estimate $\theta(x, t)$ and $\Delta(x, t)$

Find the right alternative for $S(x, y) = 1 - p(y | x)$

Calibrate the scalar thresholds β^*, q^*

Our Main Algorithm (In the finite-sample regime)

We don't have the true $p(y | x)$

But we can query $\pi(y | x)$

Need to estimate $\theta(x, t)$ and $\Delta(x, t)$

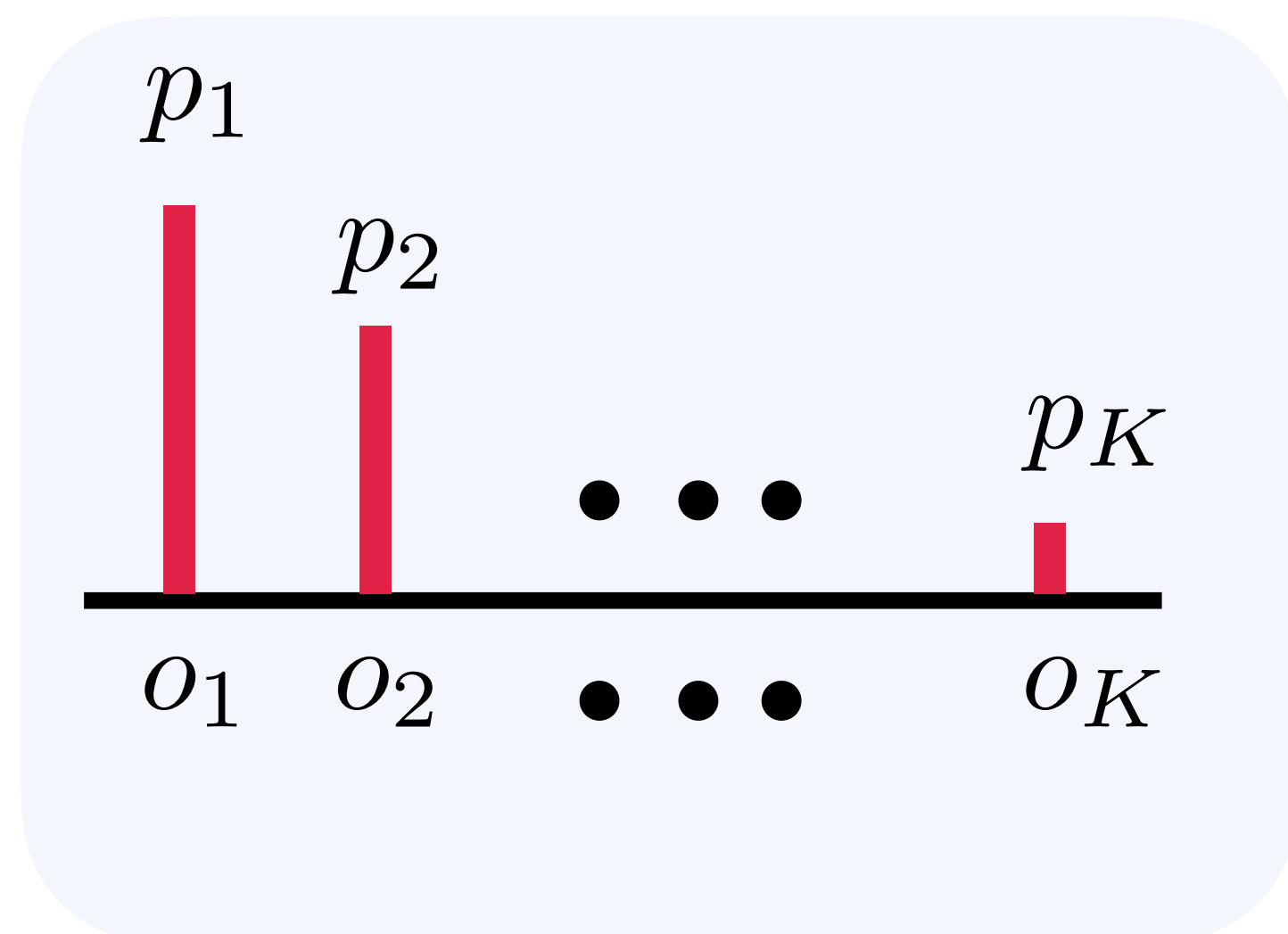
Find the right alternative for $S(x, y) = 1 - p(y | x)$

Calibrate the scalar thresholds β^*, q^*

Need to estimate $\theta(x, t)$ and $\Delta(x, t)$

Need to estimate $\theta(x, t)$ and $\Delta(x, t)$

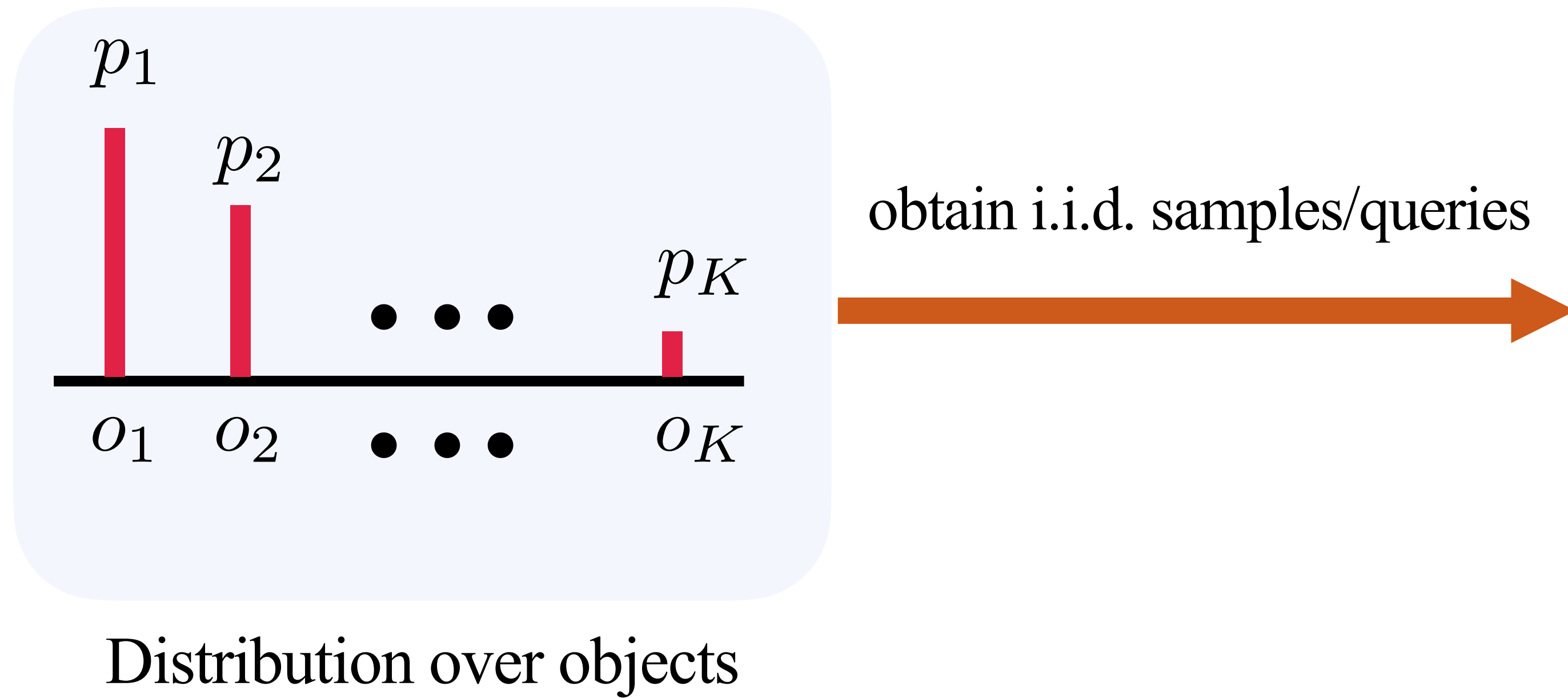
Good Turing Estimator for missing mass



Distribution over objects

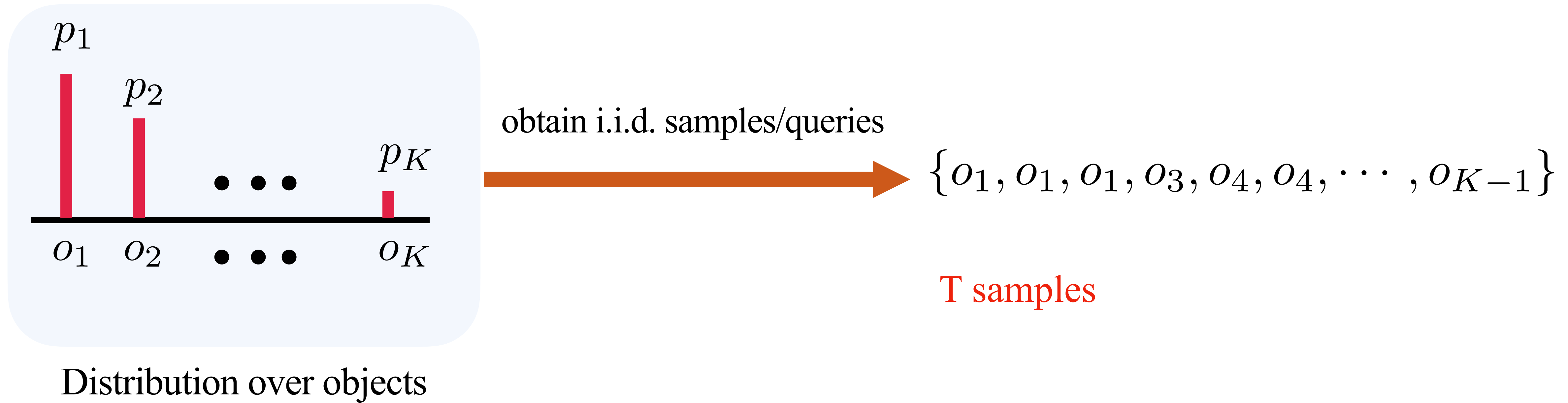
Need to estimate $\theta(x, t)$ and $\Delta(x, t)$

Good Turing Estimator for missing mass



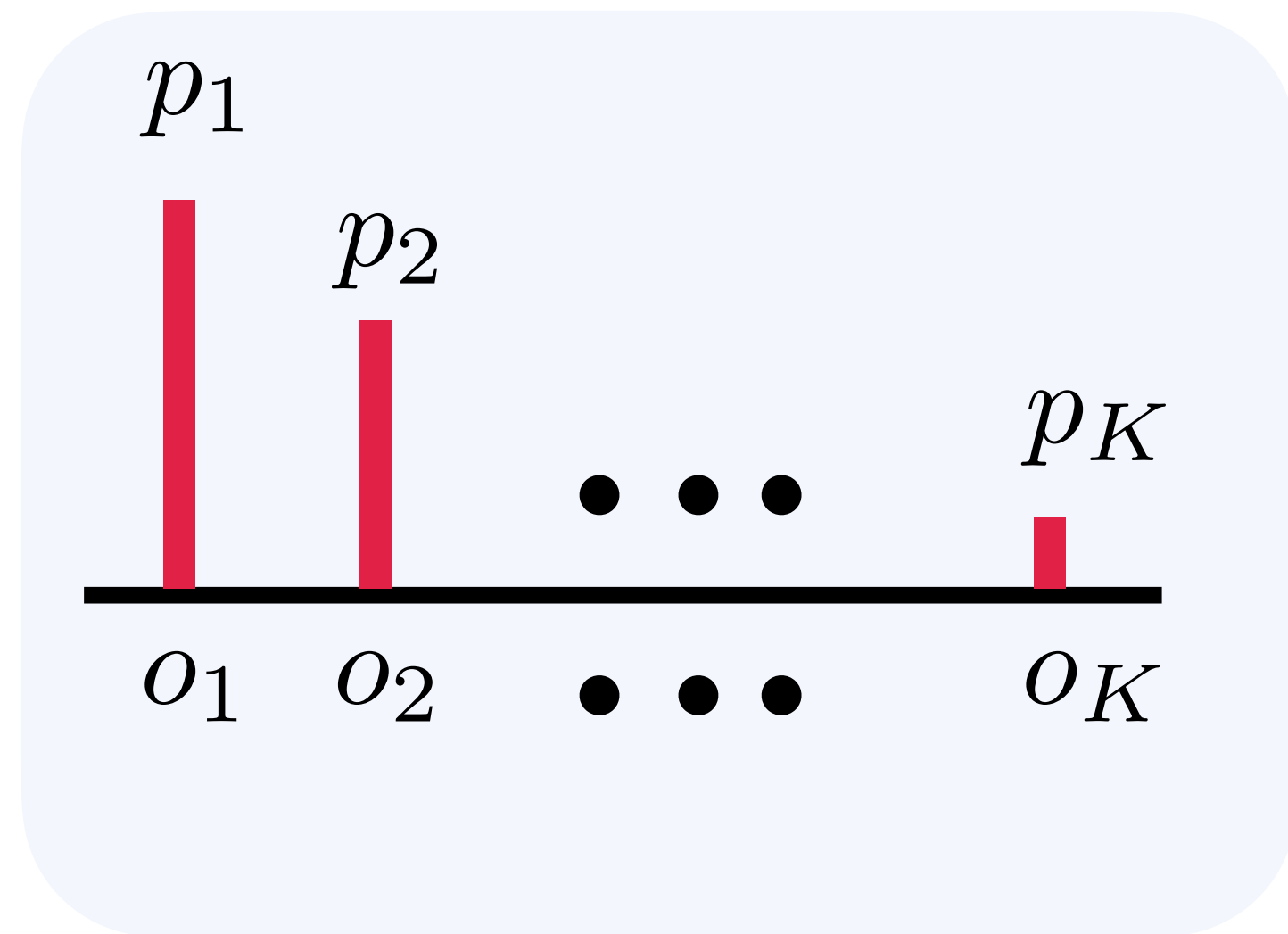
Need to estimate $\theta(x, t)$ and $\Delta(x, t)$

Good Turing Estimator for missing mass



Need to estimate $\theta(x, t)$ and $\Delta(x, t)$

Good Turing Estimator for missing mass



Distribution over objects

obtain i.i.d. samples/queries



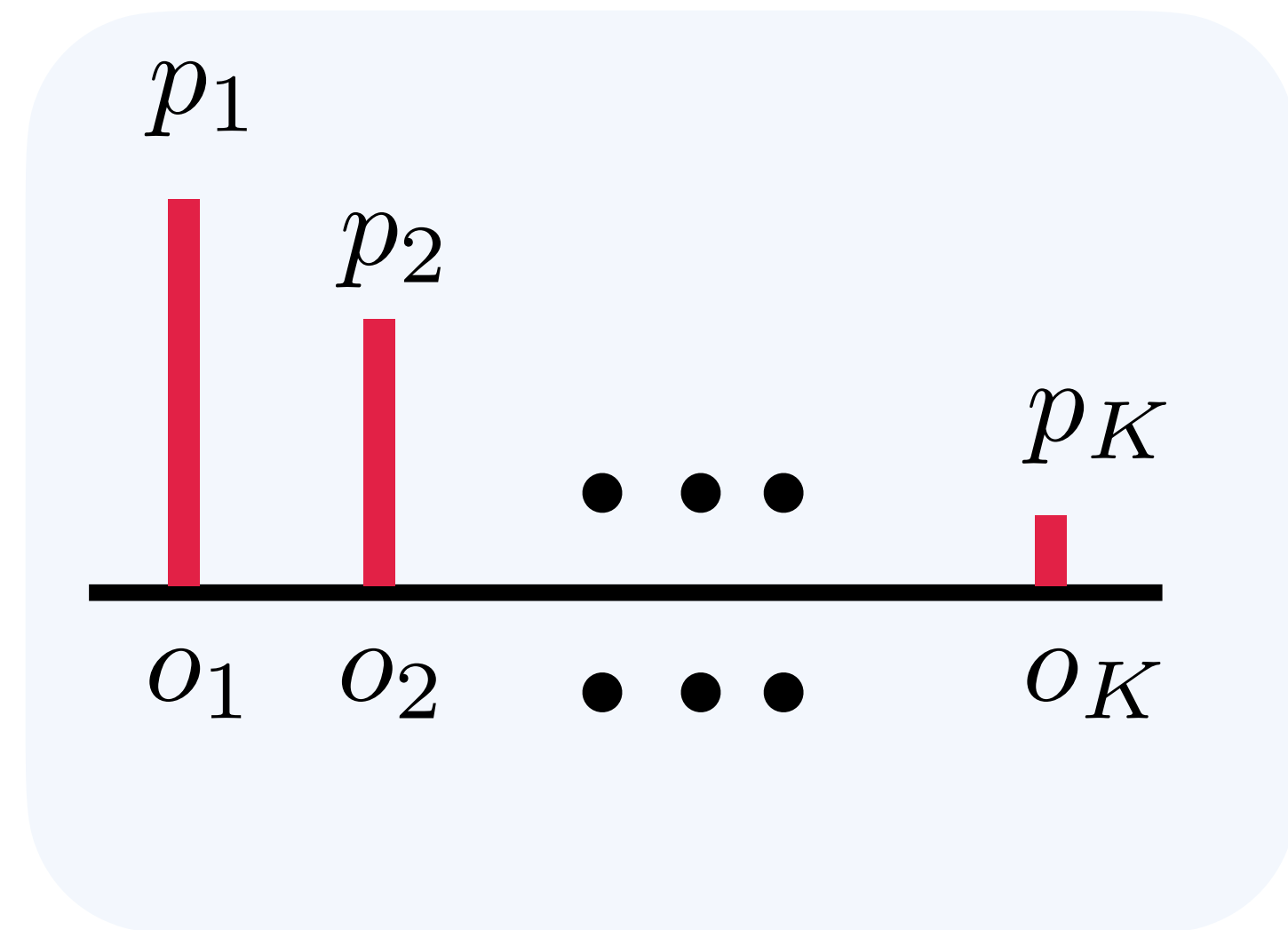
$\{o_1, o_1, o_1, o_3, o_4, o_4, \dots, o_{K-1}\}$

T samples

Some objects missing (e.g o_2)

Need to estimate $\theta(x, t)$ and $\Delta(x, t)$

Good Turing Estimator for missing mass



Distribution over objects

obtain i.i.d. samples/queries

$\{o_1, o_1, o_1, o_3, o_4, o_4, \dots, o_{K-1}\}$

T samples

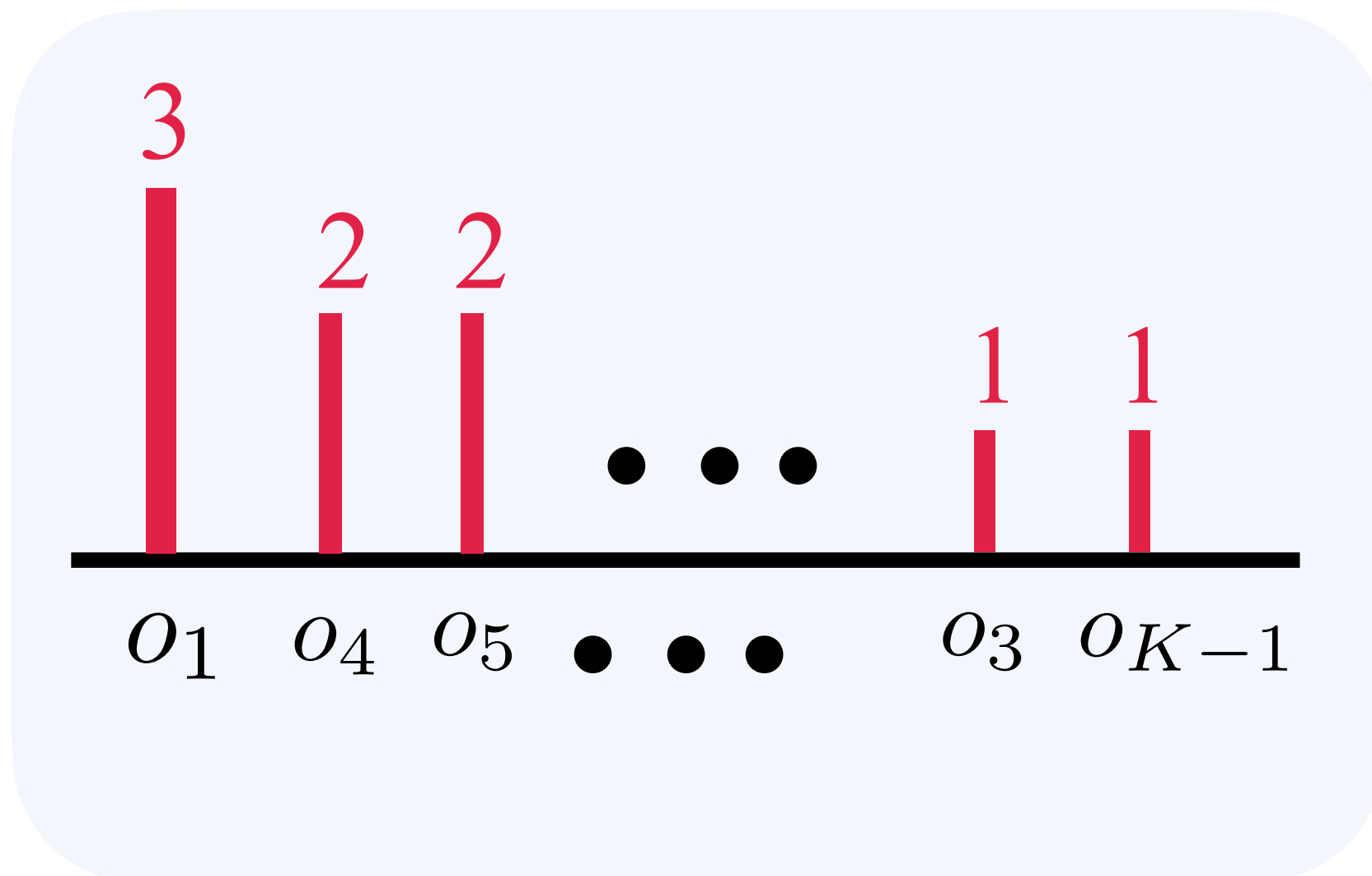
Some objects missing (e.g o_2)

Need to estimate the probability of missing objects

Need to estimate $\theta(x, t)$ and $\Delta(x, t)$

Good Turing Estimator for missing mass

$$\{o_1, o_1, o_1, o_3, o_4, o_4, \dots, o_{K-1}\}$$



Frequency of observed objects

We have $Z_t(x) = \{y_1^x, \dots, y_t^x\} \sim \pi(y | x)$

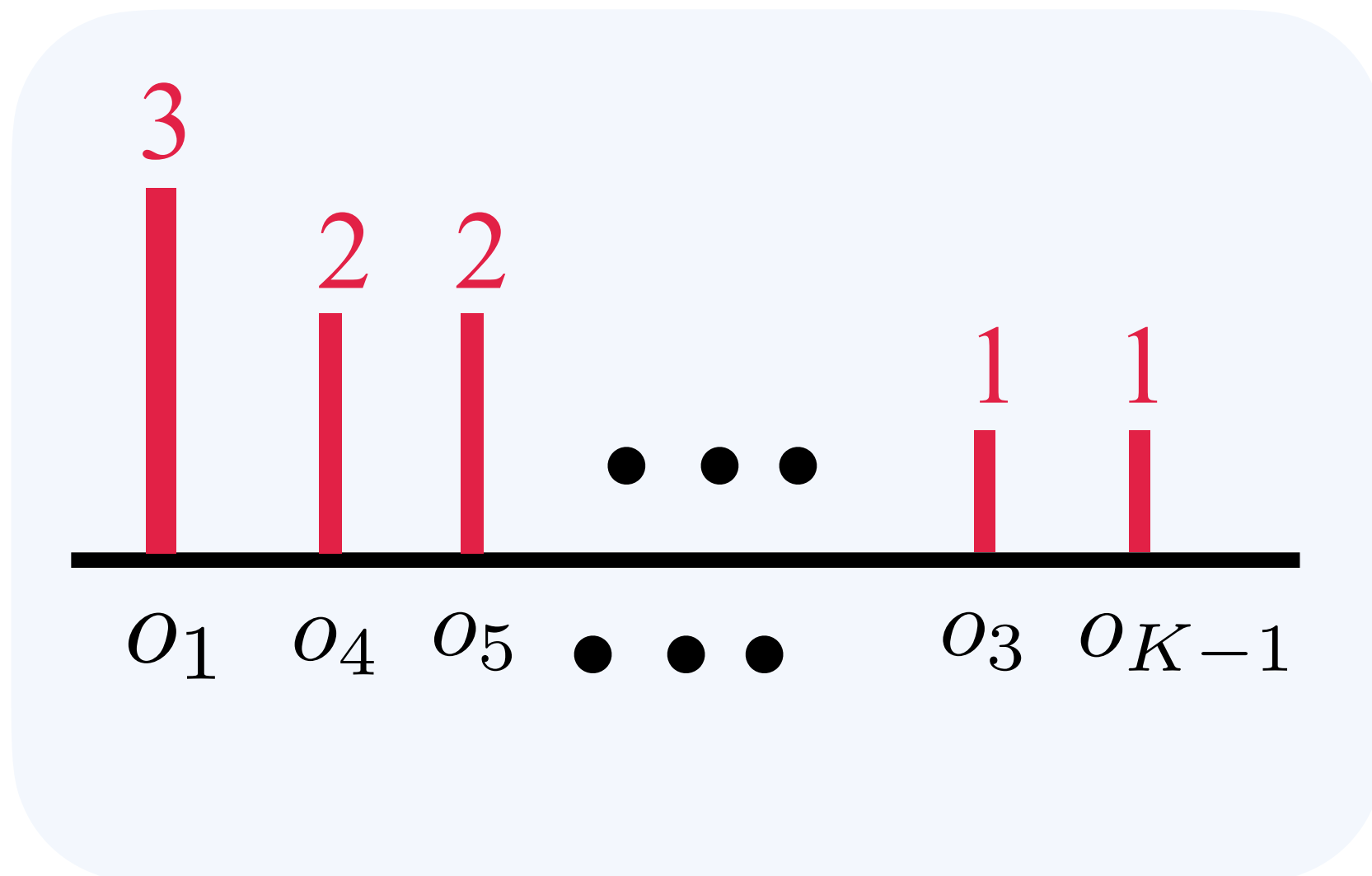
Want to estimate $\theta(x, t) = \mathbb{P}_{Y, Z_t(x)}[Y \notin Z_t(x) | X = x]$

N_r : # objects that appeared r times

Need to estimate $\theta(x, t)$ and $\Delta(x, t)$

Good Turing Estimator for missing mass

$$\{o_1, o_1, o_1, o_3, o_4, o_4, \dots, o_{K-1}\}$$



Frequency of observed objects

We have $Z_t(x) = \{y_1^x, \dots, y_t^x\} \sim \pi(y | x)$

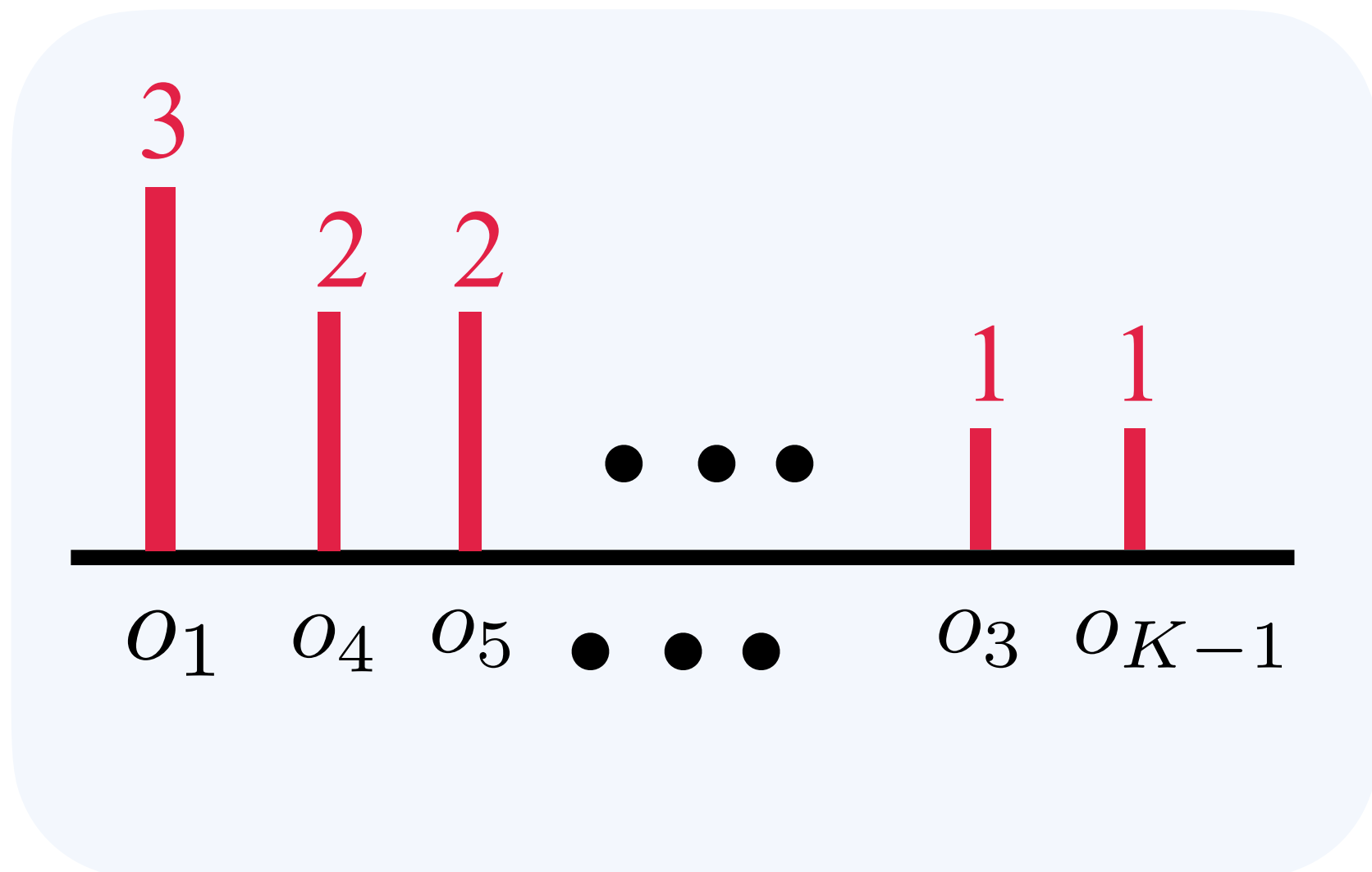
Want to estimate $\theta(x, t) = \mathbb{P}_{Y, Z_t(x)}[Y \notin Z_t(x) | X = x]$

$$N_r(x, t) = |\{y \in Z_t(x) : \#(y) = r\}|$$

Need to estimate $\theta(x, t)$ and $\Delta(x, t)$

Good Turing Estimator for missing mass

$$\{o_1, o_1, o_1, o_3, o_4, o_4, \dots, o_{K-1}\}$$



Frequency of observed objects

We have $Z_t(x) = \{y_1^x, \dots, y_t^x\} \sim \pi(y | x)$

Want to estimate $\theta(x, t) = \mathbb{P}_{Y, Z_t(x)}[Y \notin Z_t(x) | X = x]$

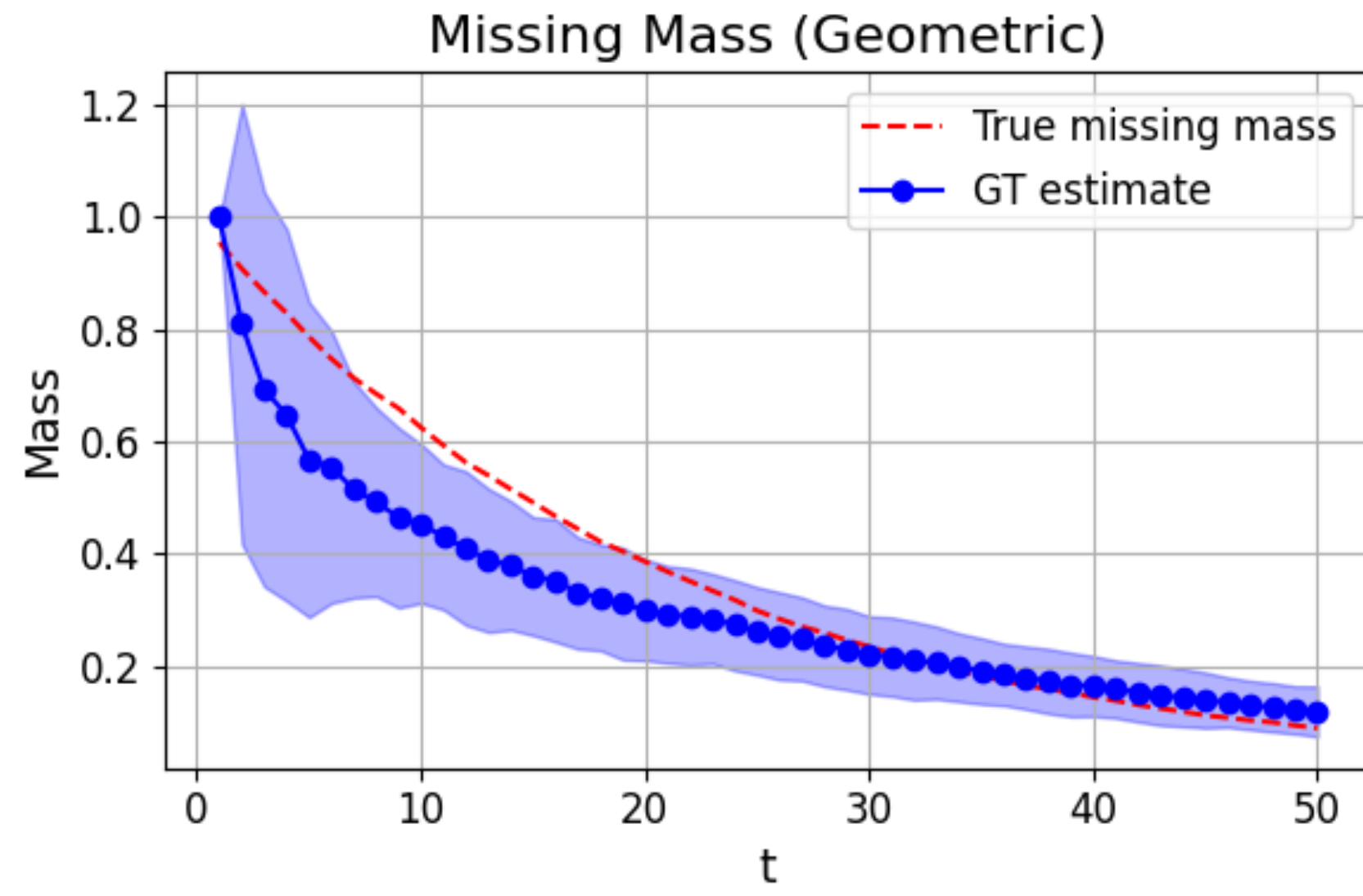
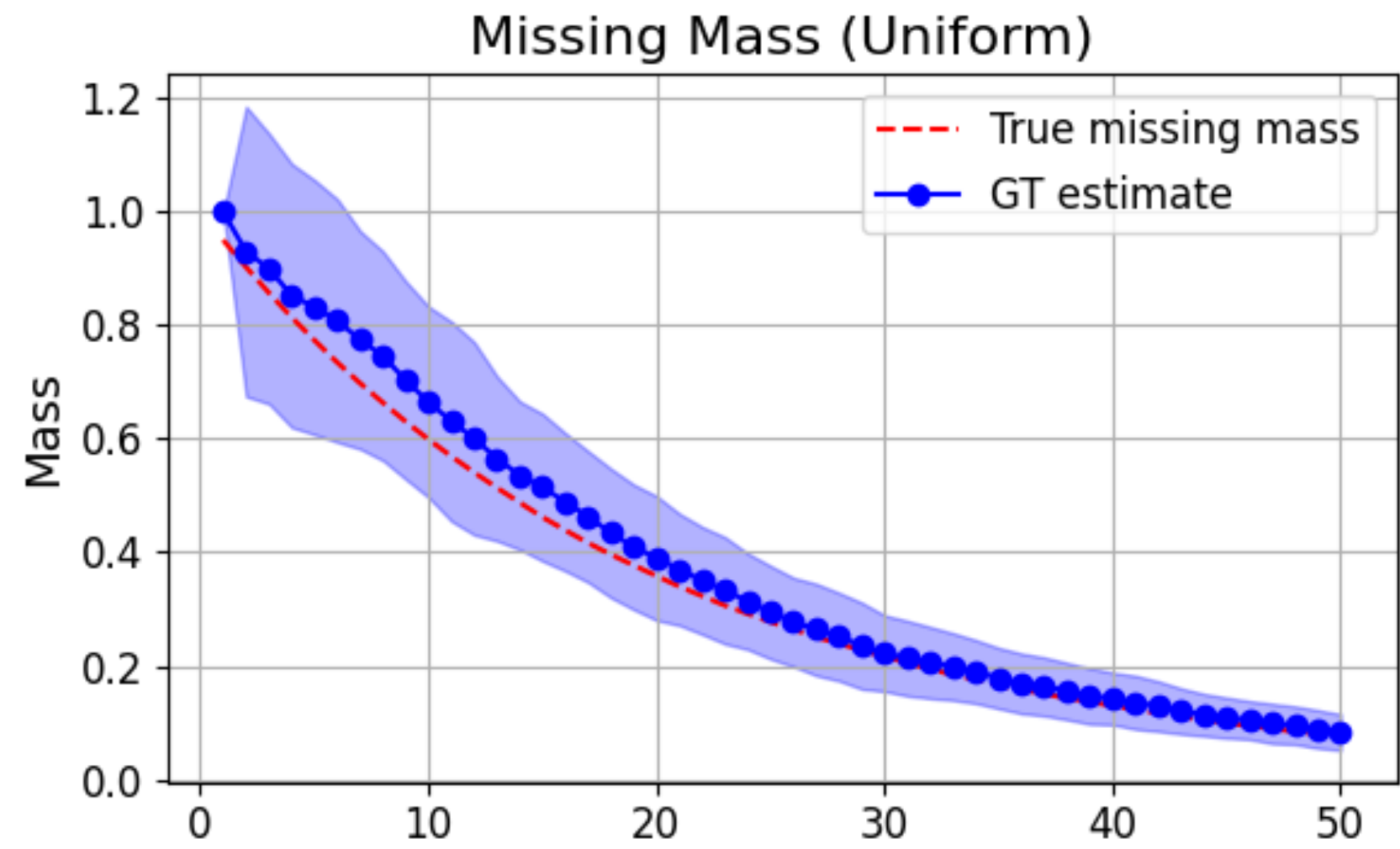
$$N_r(x, t) = |\{y \in Z_t(x) : \#(y) = r\}|$$

Good-Turing estimate for missing mass

$$\hat{\theta}(x, t) = \frac{N_1}{T}$$

Need to estimate $\theta(x, t)$ and $\Delta(x, t)$

Good Turing Estimator for missing mass



We have $Z_t(x) = \{y_1^x, \dots, y_t^x\} \sim \pi(y | x)$

Want to estimate $\theta(x, t) = \mathbb{P}_{Y, Z_t(x)}[Y \notin Z_t(x) | X = x]$

$N_r(x, t) = |\{y \in Z_t(x) : \#(y) = r\}|$

Good-Turing estimate for missing mass

$$\hat{\theta}(x, t) = \frac{N_1}{T}$$

Need to estimate $\theta(x, t)$ and $\Delta(x, t)$

Good Turing Estimator for missing mass

$$N_r(x, t) = |\{y \in Z_t(x) : \#(y) = r\}|$$

Good-Turing estimate for missing mass

$$\theta(x, t) = \mathbb{P}_{Y, Z_t(x)}[Y \notin Z_t(x) \mid X = x]$$

$$\hat{\theta}(x, t) = \frac{N_1}{T}$$

Need to estimate $\theta(x, t)$ and $\Delta(x, t)$

Good Turing Estimator for missing mass

$$N_r(x, t) = |\{y \in Z_t(x) : \#(y) = r\}|$$

Good-Turing estimate for missing mass

$$\theta(x, t) = \mathbb{P}_{Y, Z_t(x)}[Y \notin Z_t(x) \mid X = x]$$

$$\hat{\theta}(x, t) = \frac{N_1}{T}$$

Our estimate for missing mass derivative

$$\Delta(x, t) := \theta(x, t + 1) - \theta(x, t)$$

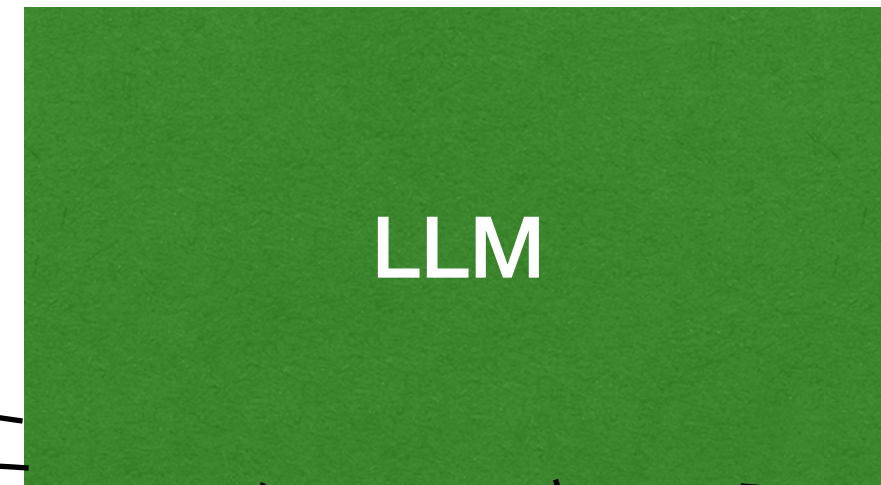
$$\hat{\Delta}(x, t) = \frac{-2N_2}{T^2}$$

“Toy” Running Example

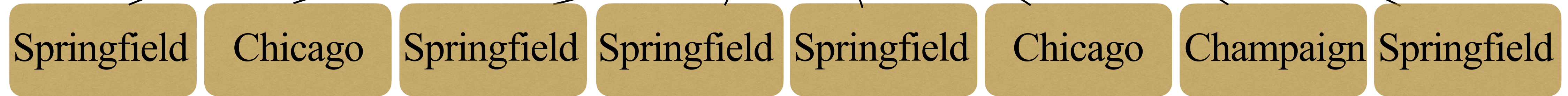
What is the capital of Illinois



Generate 8 responses, $t = 8$



Case 1



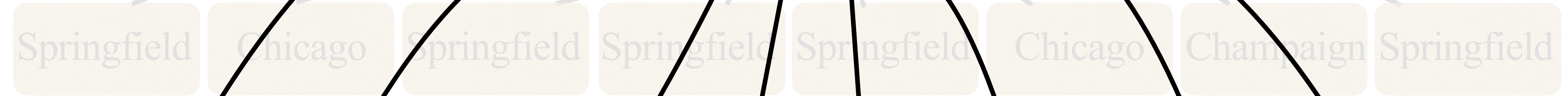
→ $\theta(x, t = 8) = \frac{1}{8}$

“Toy” Running Example

What is the capital of Illinois



Case 1



Case 2



→ $C(x) = \{\text{Springfield, Chicago}\}$

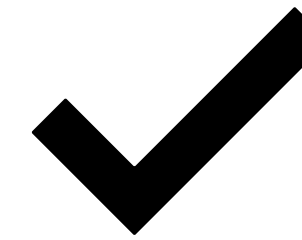
→ $\theta(x, t = 8) = 1$

Our Main Algorithm (In the finite-sample regime)

We don't have the true $p(y | x)$

But we can query $\pi(y | x)$

Need to estimate $\theta(x, t)$ and $\Delta(x, t)$



Find the right alternative for $S(x, y) = 1 - p(y | x)$

Calibrate the scalar thresholds β^*

Calibrate the scalar thresholds q^*

Calibrate the scalar thresholds β^*

Calibrate the scalar thresholds β^*

Finite sample version of principle 1

$$\Delta(x, T^*(x)) = \beta^*$$

Calibrate the scalar thresholds β^*

Finite sample version of principle 1

$$\Delta(x, T^*(x)) = \beta^*$$

For each $x \in D_{cal}$

Sample $y_{1:T(x)} \sim \pi(y | x)$ until $\hat{\Delta}(x, T(x)) \leq \beta^*$

Calibrate the scalar thresholds β^*

Finite sample version of principle 1

$$\Delta(x, T^*(x)) = \beta^*$$

For each $x \in D_{cal}$

Sample $y_{1:T(x)} \sim \pi(y | x)$ until $\hat{\Delta}(x, T(x)) \leq \beta^*$

Obtain $Z(x) = \{y_1, \dots, y_{T(x)}\}$

Calibrate the scalar thresholds β^*

Finite sample version of principle 1

$$\Delta(x, T^*(x)) = \beta^*$$

For each $x \in D_{cal}$

Sample $y_{1:T(x)} \sim \pi(y | x)$ until $\hat{\Delta}(x, T(x)) \leq \beta^*$

Obtain $Z(x) = \{y_1, \dots, y_{T(x)}\}$

Grid search

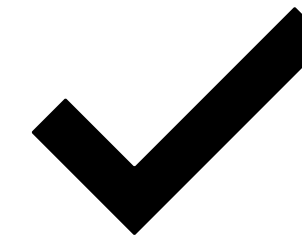
Make sure $\frac{1}{|D_{cal}|} \sum_x T(x) \leq B$

Our Main Algorithm (In the finite-sample regime)

We don't have the true $p(y | x)$

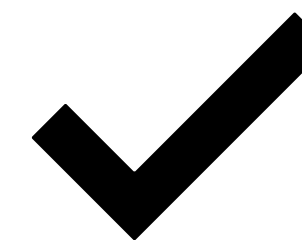
But we can query $\pi(y | x)$

Need to estimate $\theta(x, t)$ and $\Delta(x, t)$



Find the right alternative for $S(x, y) = 1 - p(y | x)$

Calibrate the scalar thresholds β^*



Calibrate the scalar thresholds q^*

Find the right alternative for $S(x, y) = 1 - p(y | x)$

Find the right alternative for $S(x, y) = 1 - p(y | x)$ and $\Delta(x, t)$

Finite sample version of **principle 2**

$$S(x, y) = \begin{cases} 1 - p(y | x), & \text{if } y \neq EE \\ 2 - p(y | x), & \text{if } y = EE \end{cases}$$

$$Z(x) = \{y_1^x, \dots, y_{T(x)}^x\}$$

Find the right alternative for $S(x, y) = 1 - p(y | x)$ and $\Delta(x, t)$

Finite sample version of **principle 2**

$$S(x, y) = \begin{cases} 1 - p(y | x), & \text{if } y \neq EE \\ 2 - p(y | x), & \text{if } y = EE \end{cases}$$

$$Z(x) = \{y_1^x, \dots, y_{T(x)}^x\}$$

Using the Good Turing estimator we can estimate:

$$\hat{p}(y | x)$$

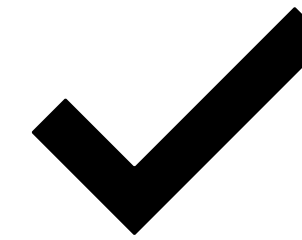
$$\hat{p}(EE | x) = \theta(x, T(x))$$

Our Main Algorithm (In the finite-sample regime)

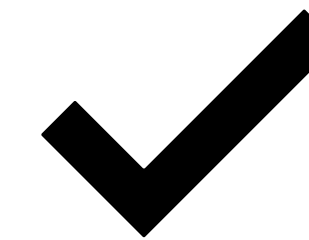
We don't have the true $p(y | x)$

But we can query $\pi(y | x)$

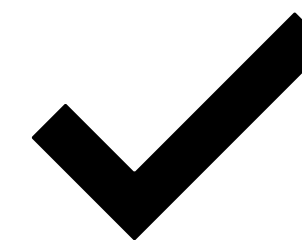
Need to estimate $\theta(x, t)$ and $\Delta(x, t)$



Find the right alternative for $S(x, y) = 1 - p(y | x)$



Calibrate the scalar thresholds β^*



Calibrate the scalar thresholds q^*

Finite Sample Algorithm

Algorithm 1 Conformal Prediction with Query Oracle (CPQ)

Input: Query oracle $\pi(y | x)$, conformity score $\hat{S}(x, y)$, calibration data $\mathcal{D}_{\text{cal}_2}$, test point x_{test} , miscoverage α , query budget B , missing-mass estimator $\hat{\Delta}(x, t)$, threshold β^*

Query Module \rightarrow Principle 1

- For each $x \in \mathcal{D}_{\text{cal}_2} \cup \{x_{\text{test}}\}$:
 - Sample $y_{1:T(x)} \sim \pi(y | x)$ until $\hat{\Delta}(x, T(x)) \leq \beta^*$. Let $Z(x) = \{y_1, \dots, y_{T(x)}\}$.

Calibration Module \rightarrow Principle 2

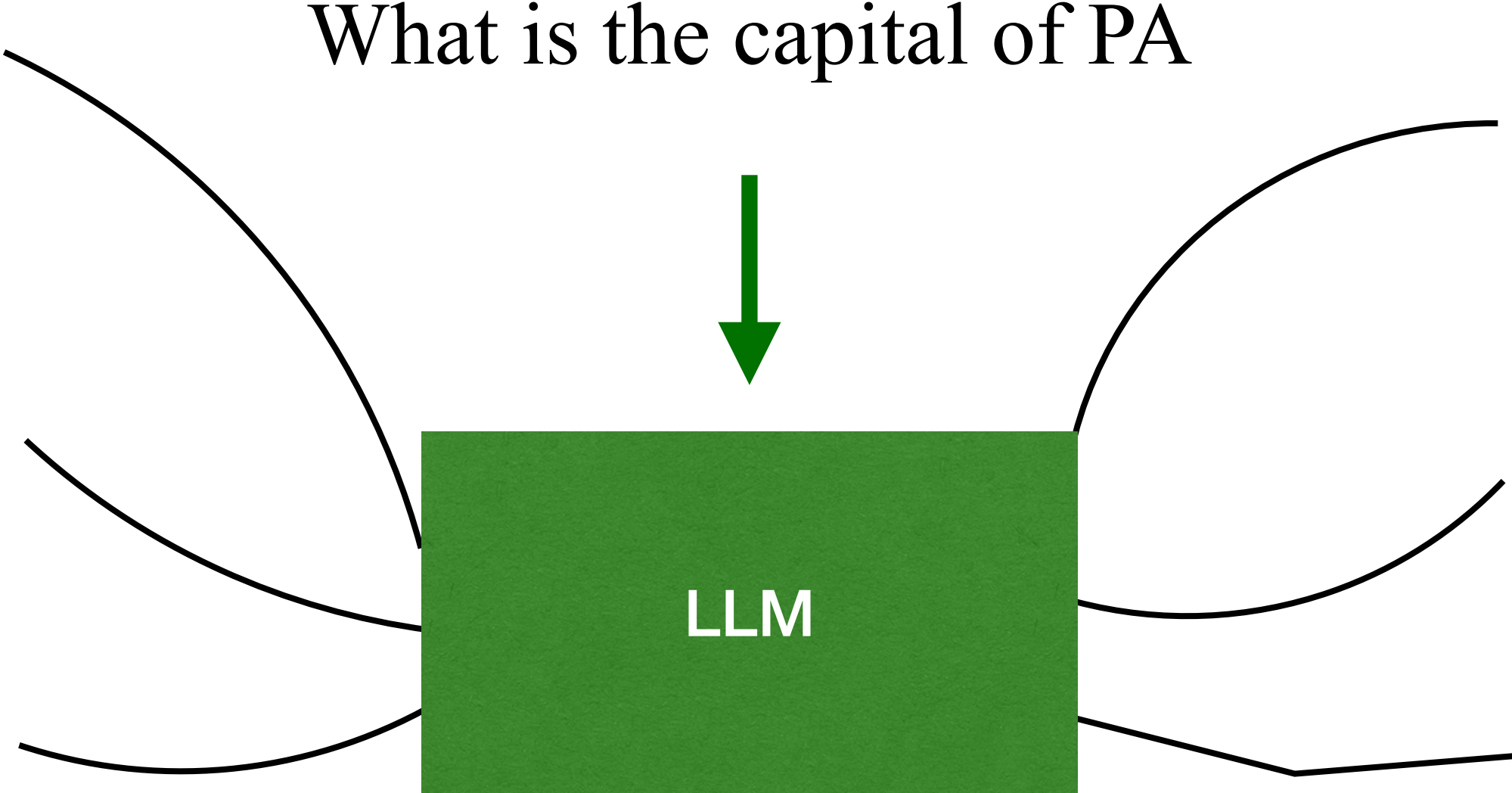
- For each $(x_i, y_i) \in \mathcal{D}_{\text{cal}_2}$ compute $s_i = \hat{S}(x_i, y_i)$.
- Set $q^* = \text{Quantile}_{1-\alpha}(s_1, \dots, s_{|\mathcal{D}_{\text{cal}_2}|}, \infty)$.

Output: $C(x_{\text{test}}) = \{y \in Z(x_{\text{test}}) \cup \{\mathbf{EE}\} : \hat{S}(x_{\text{test}}, y) \leq q^*\}$.

Let's apply our algorithm to LLMs

Let's apply our algorithm to LLMs

Clustering outputs



Clusters

Fine grained component wise comparison

We consider three methods:

Vanilla - fixed non-adaptive querying, valid but sub-optimal calibration

P1 - adds adaptive querying

P1 + P2 - optimal querying and optimal calibration -> CPQ

Fine grained component wise comparison

We consider three methods:

Vanilla - fixed non-adaptive querying, valid but sub-optimal calibration

P1 - adds adaptive querying

P1 + P2 - optimal querying and optimal calibration -> CPQ

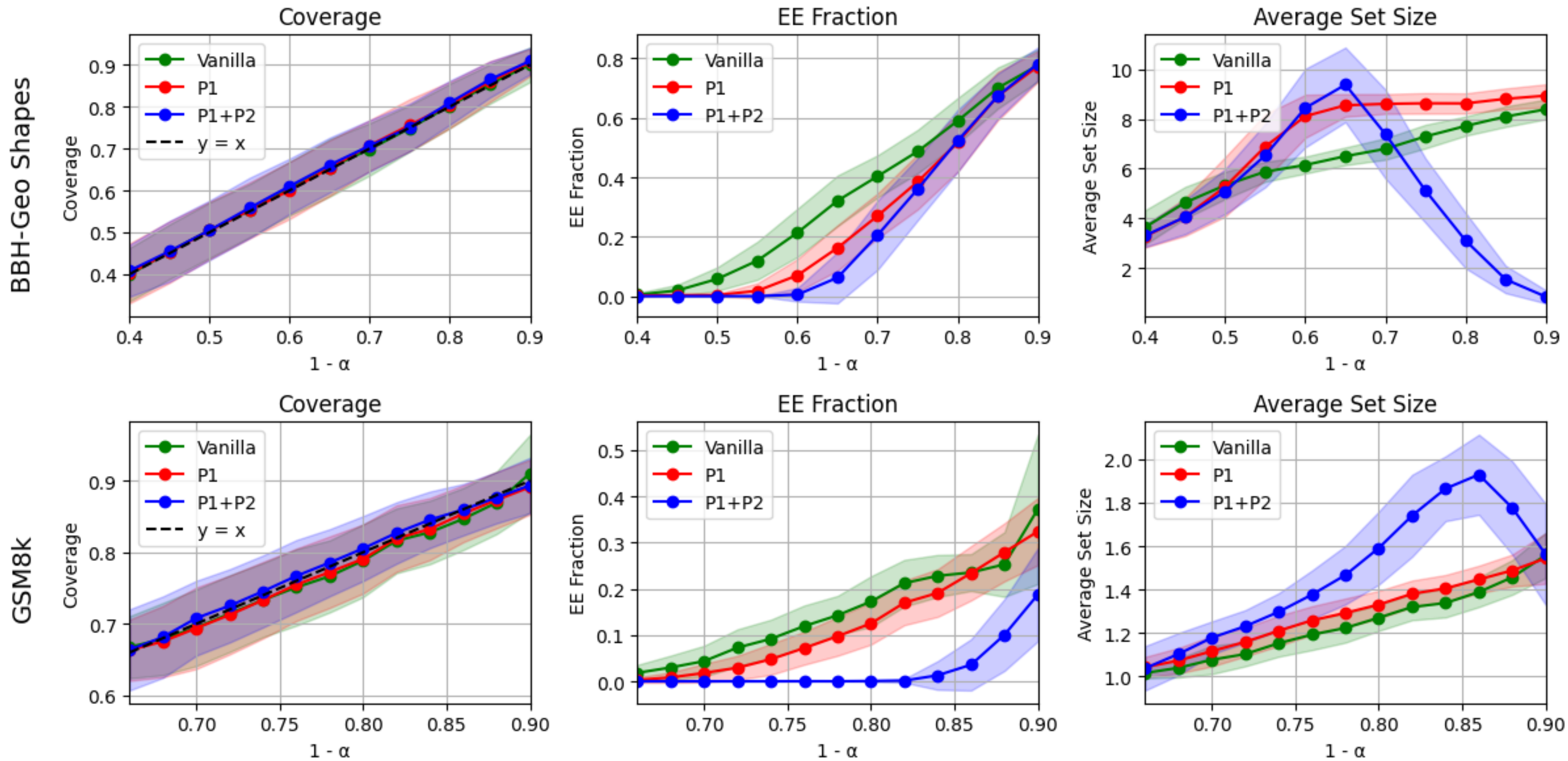
And three evaluation metrics:

Empirical Coverage

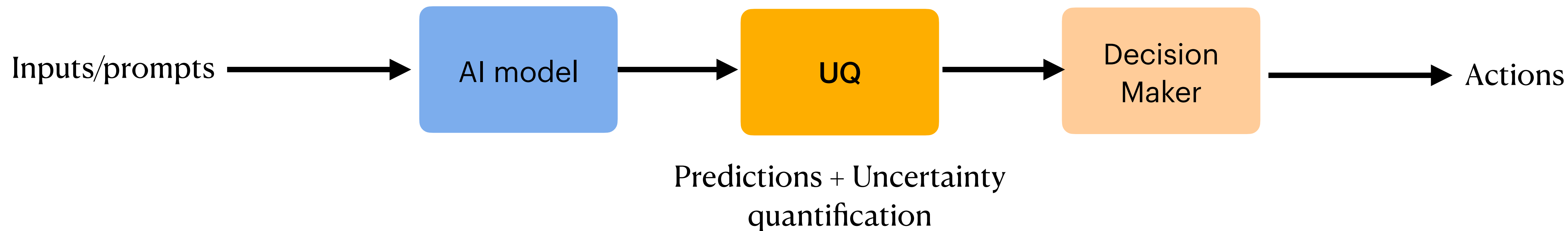
EE fraction

Prediction set size

Fine grained component wise comparison



AI-powered Decision making pipeline



Part I: CP for Generative Models

Conformal Prediction Beyond the Seen: A Missing Mass Perspective for Uncertainty Quantification in Generative Models

Sima Noorani^{*1}, Shayan Kiyani^{*1}, George Pappas¹, and Hamed Hassani¹

¹University of Pennsylvania

Part II: Collaborative CP

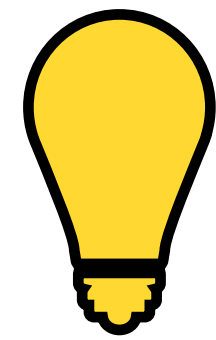
Human AI Collaborative Uncertainty Quantification

Sima Noorani^{*1}, Shayan Kiyani^{*1}, George Pappas¹, and Hamed Hassani¹

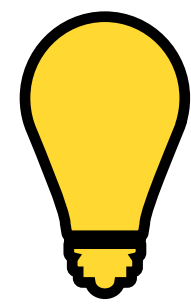
¹University of Pennsylvania

Arxiv - Oct '25

UQ when humans and AI are jointly in the loop?



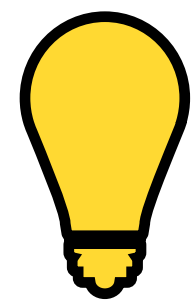
Human and AI should coexist in the decision making pipeline



Human and AI should **coexist** in the **decision making** pipeline

AI

Handle large unstructured data

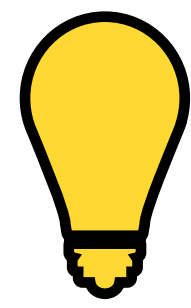


Human and AI should coexist in the decision making pipeline

AI

Handle large unstructured data

Excel at patten extraction



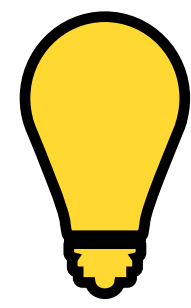
Human and AI should coexist in the decision making pipeline

AI

Handle large unstructured data

Excel at patten extraction

Offer statistical accuracy



Human and AI should **coexist** in the **decision making** pipeline

AI

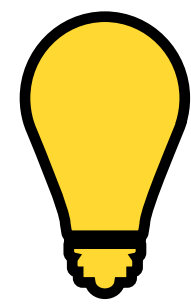
Handle large unstructured data

Excel at patten extraction

Offer statistical accuracy



Domain Knowledge



Human and AI should **coexist** in the **decision making** pipeline

AI

Handle large unstructured data

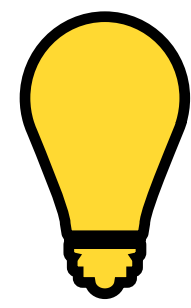
Excel at patten extraction

Offer statistical accuracy



Domain Knowledge

Persistent Memory



Human and AI should **coexist** in the **decision making** pipeline

AI

Handle large unstructured data

Excel at patten extraction

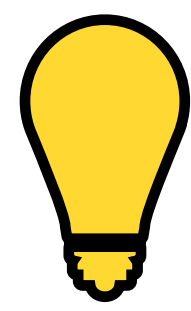
Offer statistical accuracy



Domain Knowledge

Persistent Memory

**Reason and act in
the physical world**



Human and AI should **coexist** in the **decision making** pipeline

AI

Handle large unstructured data

Excel at patten extraction

Offer statistical accuracy



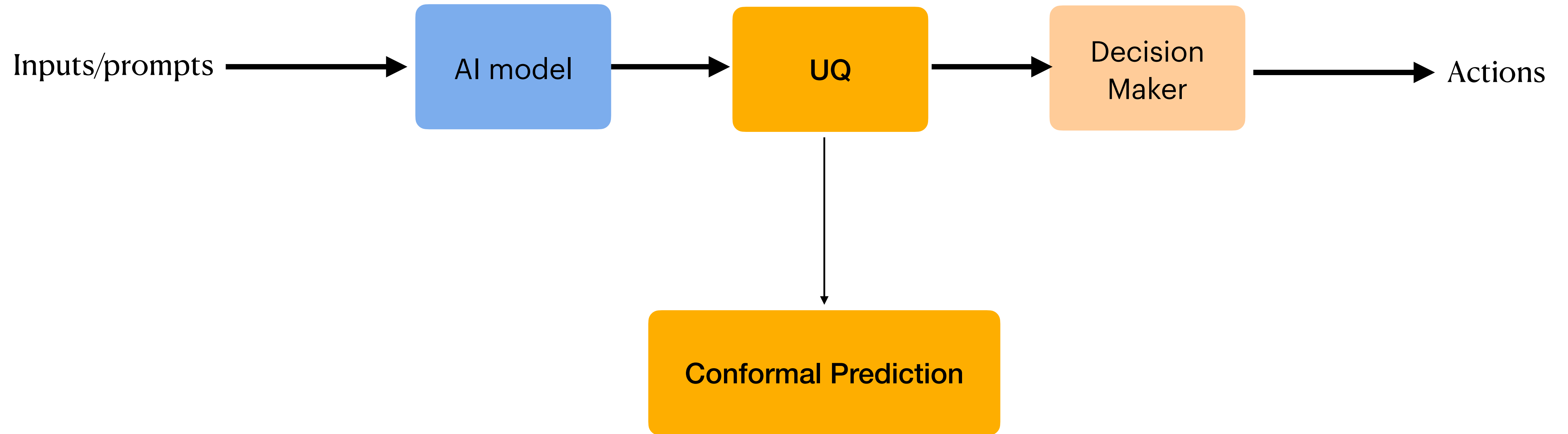
Domain Knowledge

Persistent Memory

**Reason and act in
the physical world**

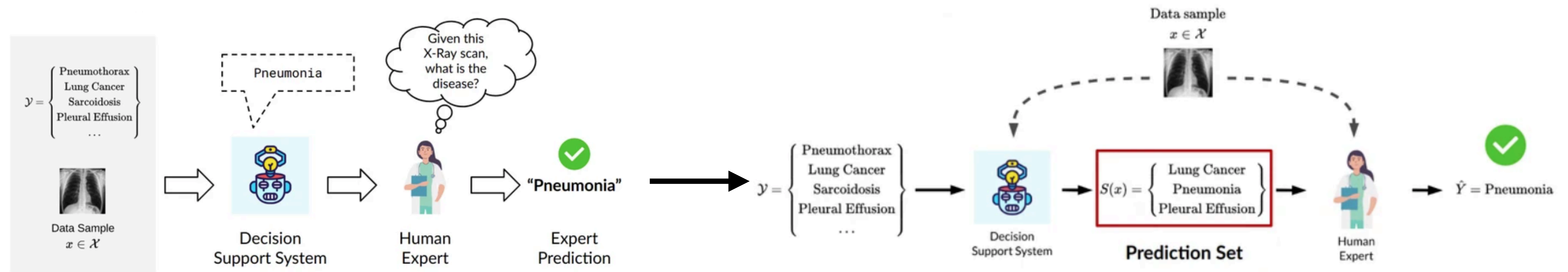
Human and AI should jointly exist in the decision making process!

AI-powered Decision making pipeline



Prediction sets are useful for decision makers

Prediction sets are useful for decision makers

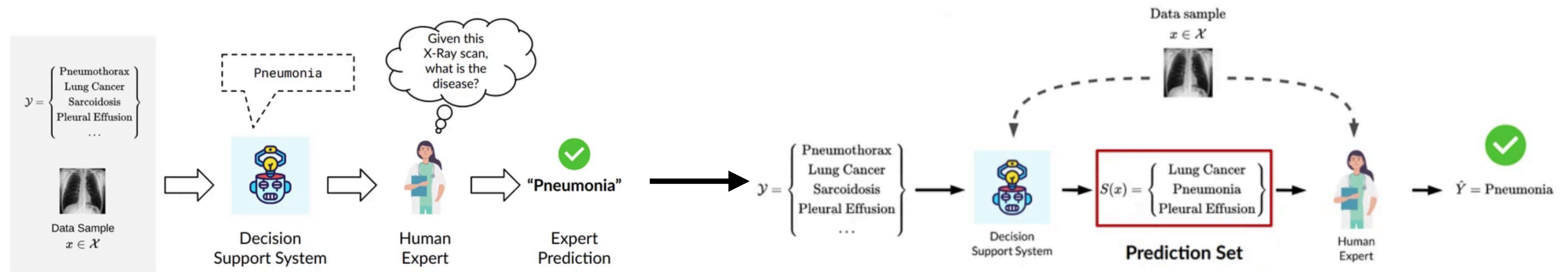


Straitouri et al., Improving Expert Predictions with Conformal Prediction, ICML, 2023.

E. Straitouri & M. Gomez-Rodriguez, *Designing Decision Support Systems using Counterfactual Prediction Sets*, ICML, 2024.

Giovanni De Toni, Nastaran Okati, Suhas Thejaswi, Eleni Straitouri, and Manuel Gomez-Rodriguez. Towards human-AI complementarity with prediction sets. In The Thirty-eighth Annual Conference on Neural Information Processing Systems, 2024.

Prediction sets are useful for decision makers



Straitouri et al., Improving Expert Predictions with Conformal Prediction, ICML, 2023.

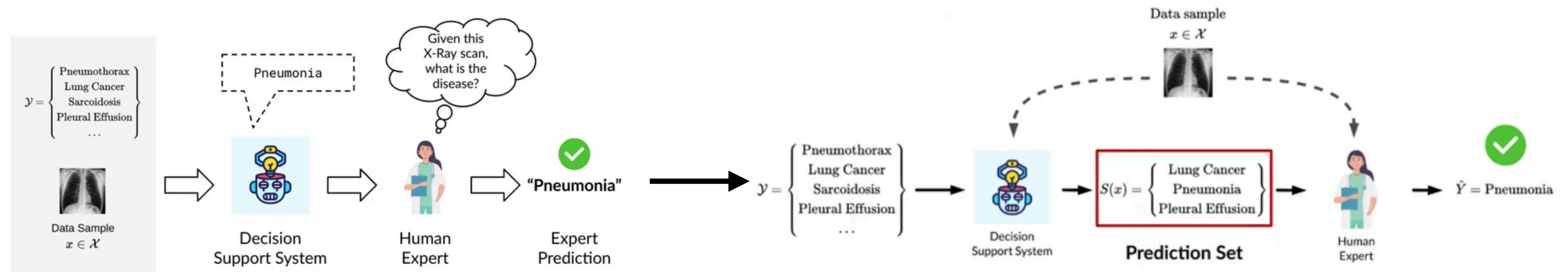
E. Straitouri & M. Gomez-Rodriguez, *Designing Decision Support Systems using Counterfactual Prediction Sets*, ICML, 2024.

Giovanni De Toni, Nastaran Okati, Suhas Thejaswi, Eleni Straitouri, and Manuel Gomez-Rodriguez. Towards human-AI complementarity with prediction sets. In The Thirty-eighth Annual Conference on Neural Information Processing Systems, 2024.

“Decision-theoretic framework for evaluating predictive uncertainty as informative signals”

Hullman et. al., *Conformal Prediction and Human Decision Making*, 2025.

Prediction sets are useful for decision makers



Straitouri et al., Improving Expert Predictions with Conformal Prediction, ICML, 2023.

E. Straitouri & M. Gomez-Rodriguez, *Designing Decision Support Systems using Counterfactual Prediction Sets*, ICML, 2024.

Giovanni De Toni, Nastaran Okati, Suhas Thejaswi, Eleni Straitouri, and Manuel Gomez-Rodriguez. Towards human-AI complementarity with prediction sets. In The Thirty-eighth Annual Conference on Neural Information Processing Systems, 2024.

“Decision-theoretic framework for evaluating predictive uncertainty as informative signals”

Hullman et. al., Conformal Prediction and Human Decision Making, 2025.

“Well designed” prediction sets are a sufficient statistic for risk averse decision making

S. Kiyani et Al., *Decision Theoretic Foundations for Conformal Prediction: Optimal Uncertainty Quantification for Risk-Averse Agents*, ICML, 2025.

Thus far ...

Thus far ...



Human and AI should coexist in the decision making pipeline

AI



Thus far ...



Human and AI should coexist in the decision making pipeline

AI



UQ is essential in the decision making pipeline

UQ



Thus far ...



Human and AI should coexist in the decision making pipeline

AI



UQ is essential in the decision making pipeline

UQ



CP is a promising tool for UQ of AI decision support systems

Thus far ...



Human and AI should coexist in the decision making pipeline

AI



UQ is essential in the decision making pipeline

UQ



CP is a promising tool for UQ of AI decision support systems

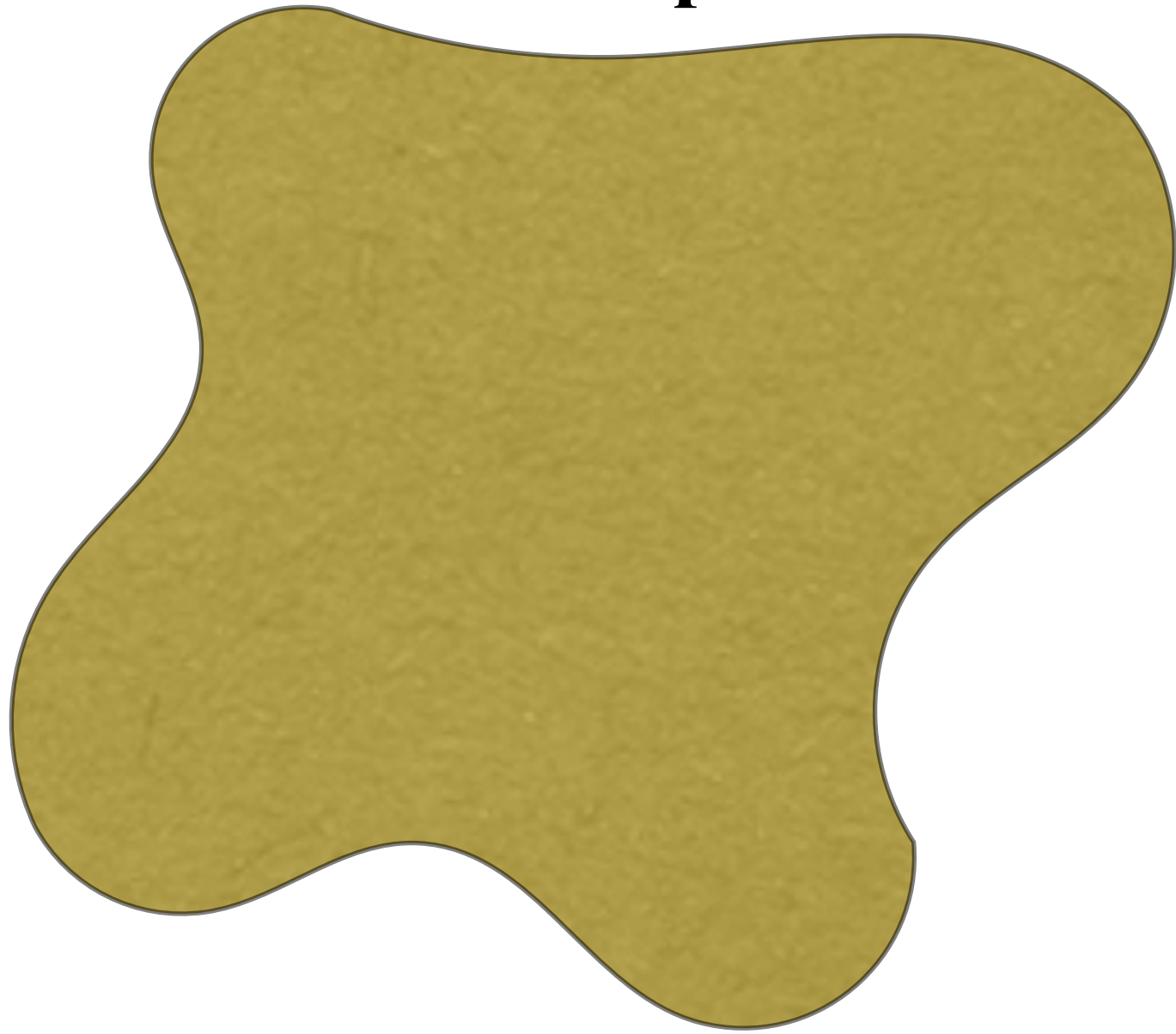
What should be the principles of UQ when Human and AI are jointly in the loop?

Question: what constitutes a good collaboration?

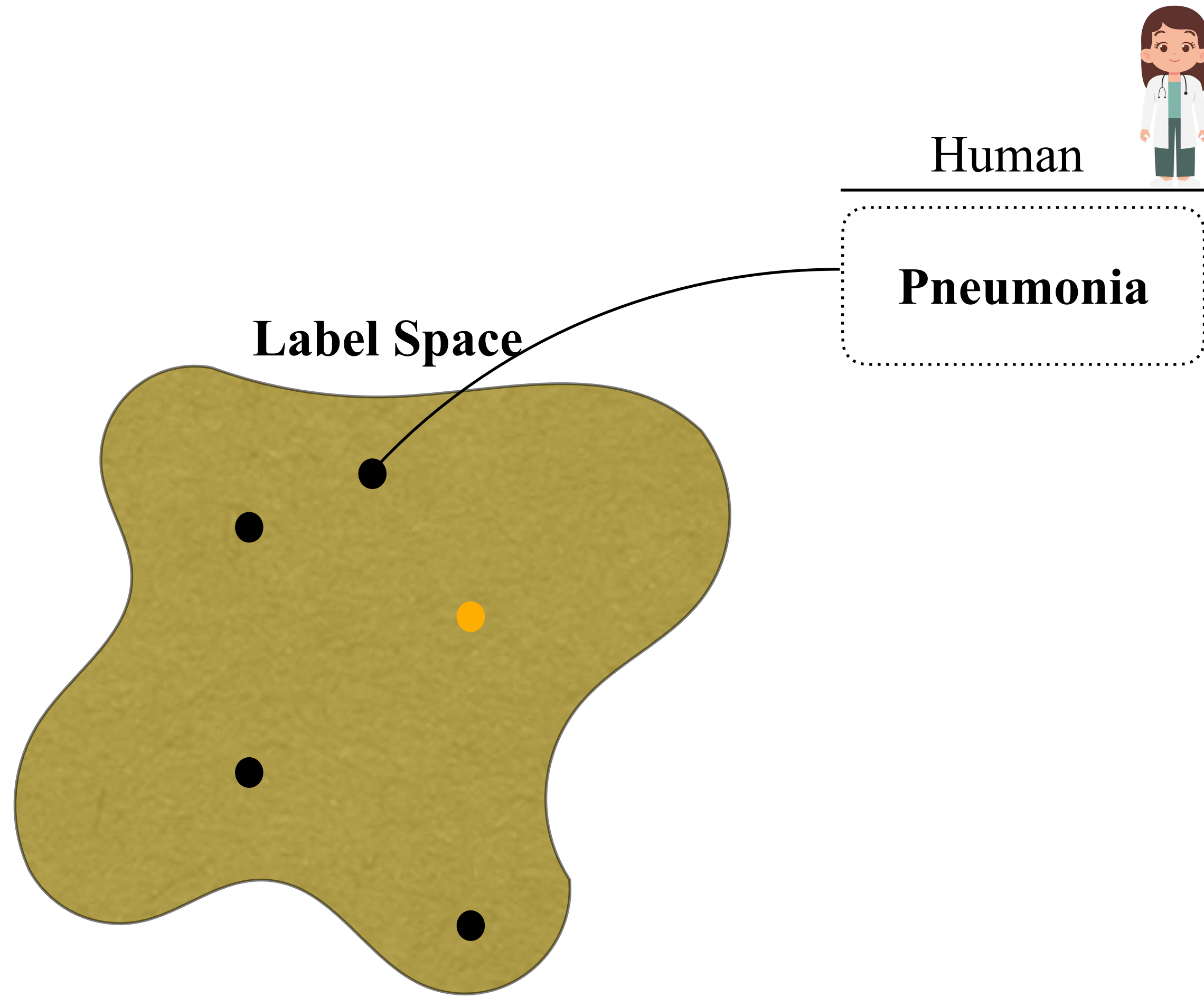
Question: what constitutes a good collaboration?

Question: what constitutes a good collaboration?

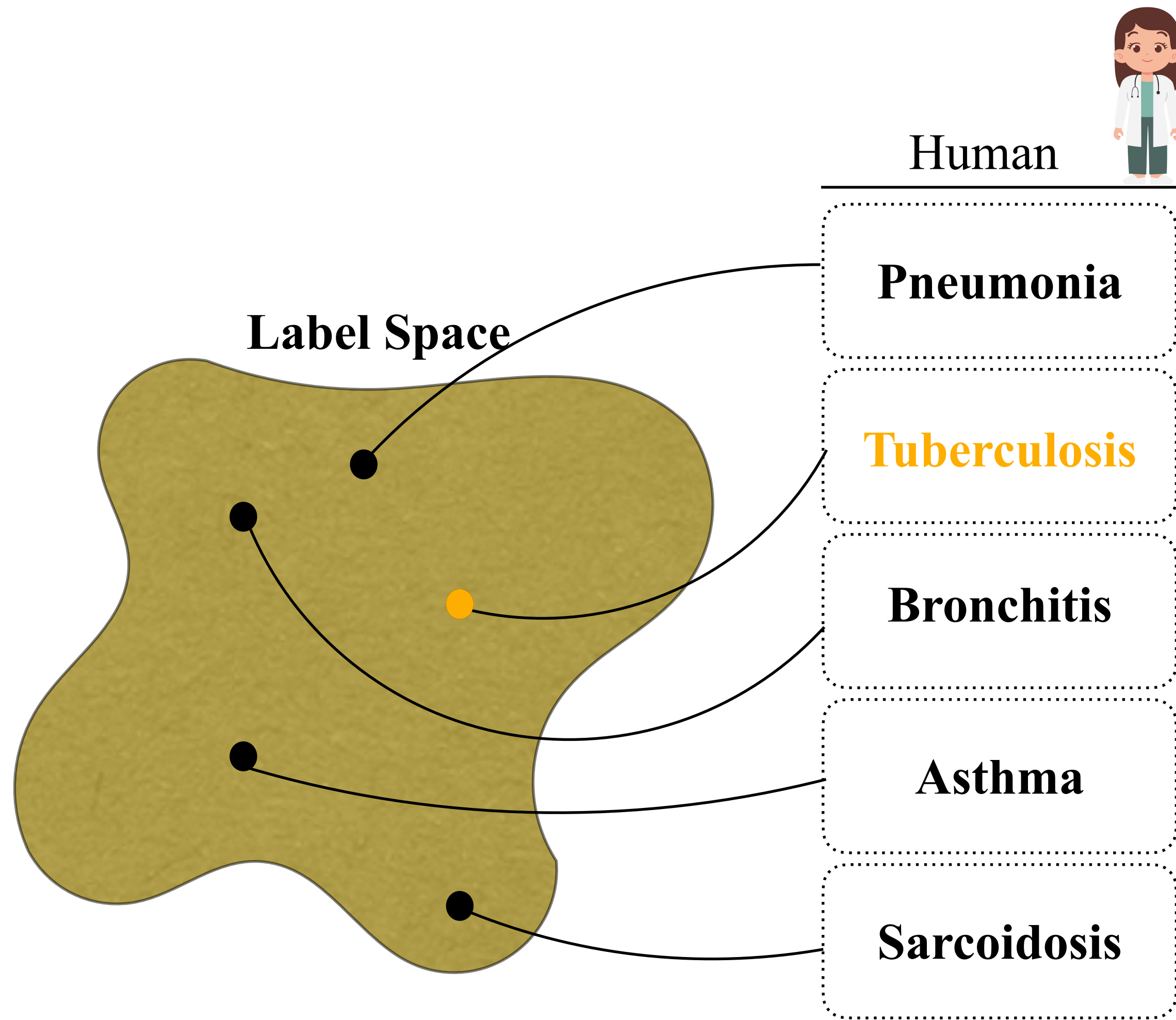
Label Space



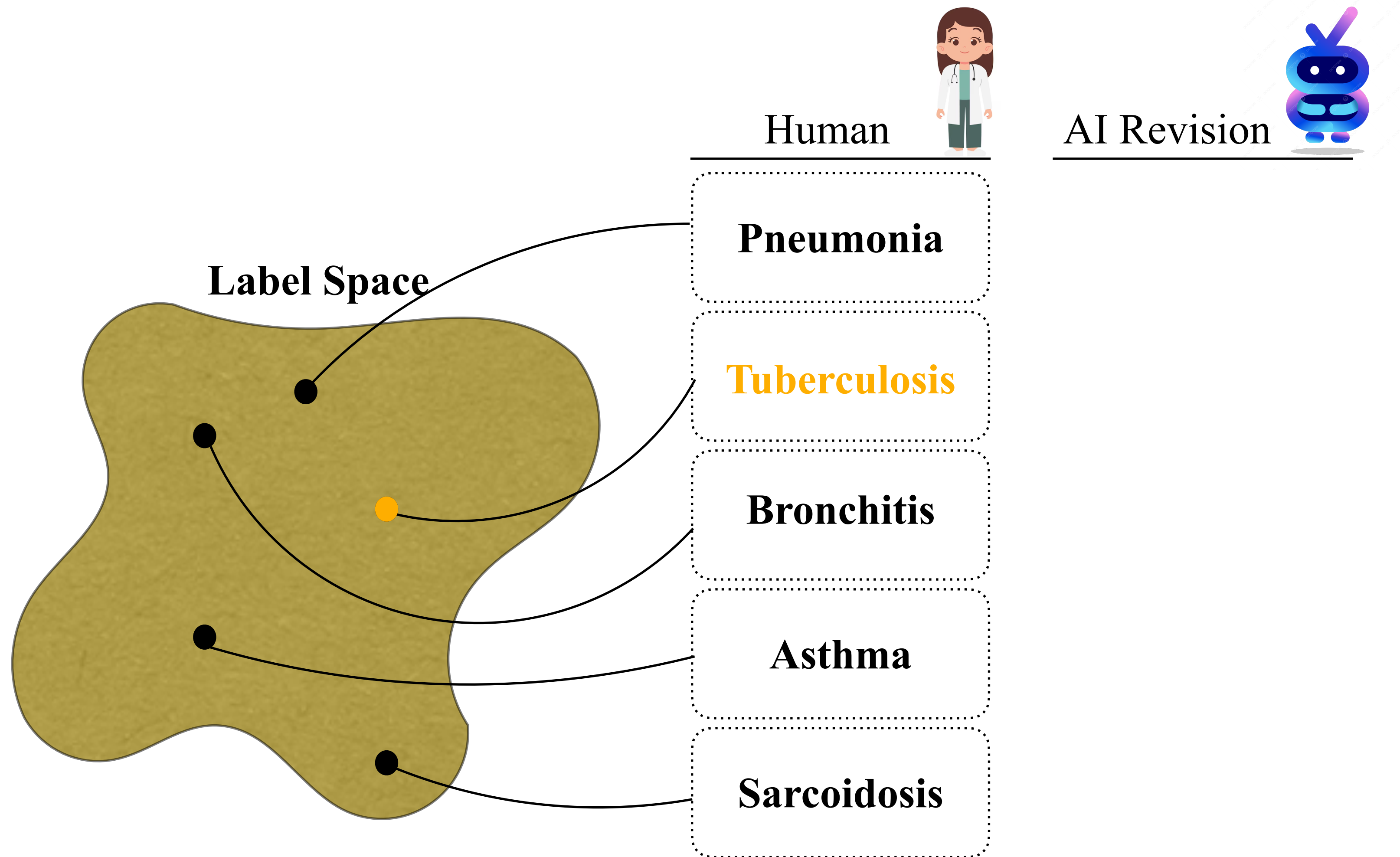
Question: what constitutes a good collaboration?



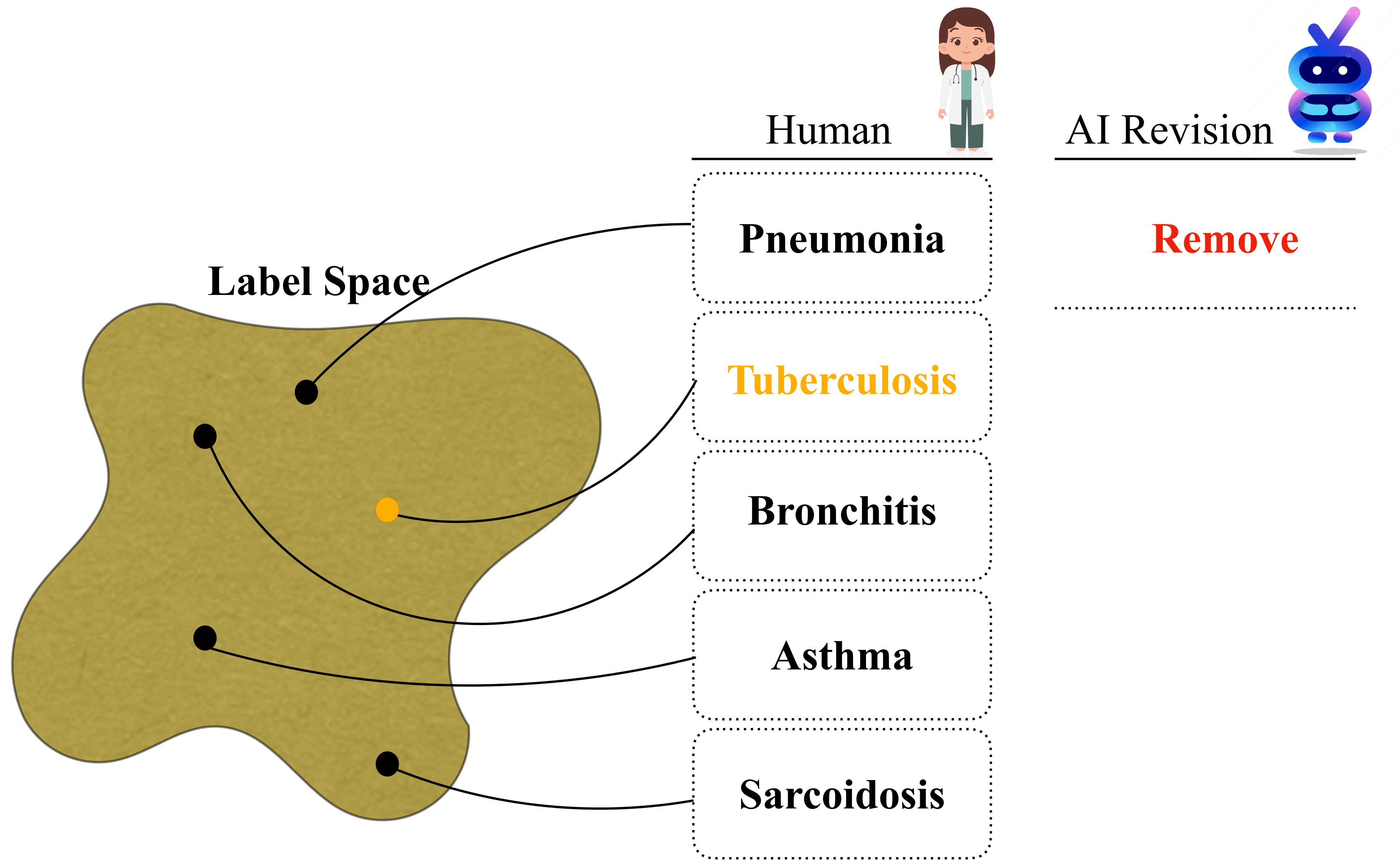
Question: what constitutes a good collaboration?



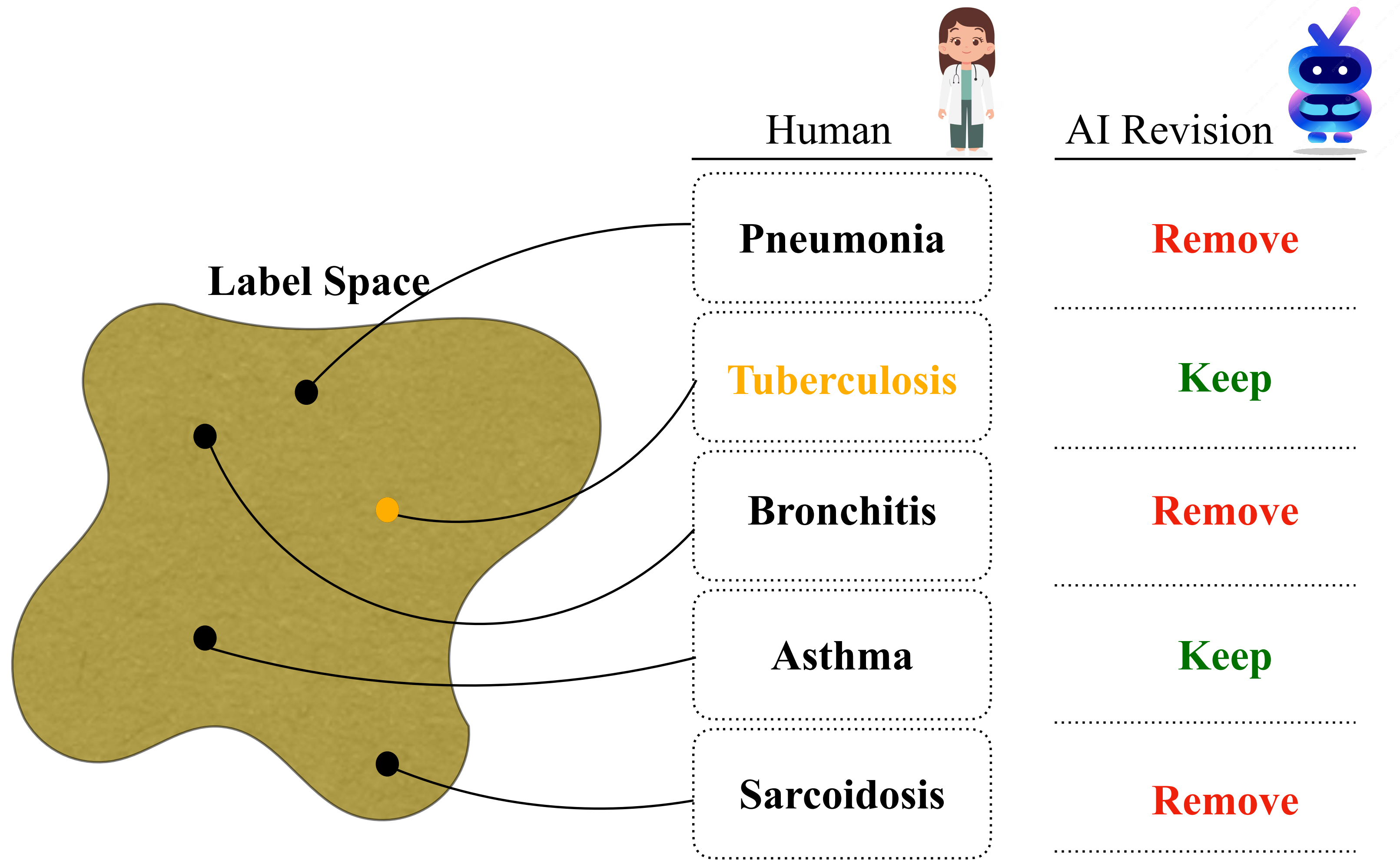
Question: what constitutes a good collaboration?



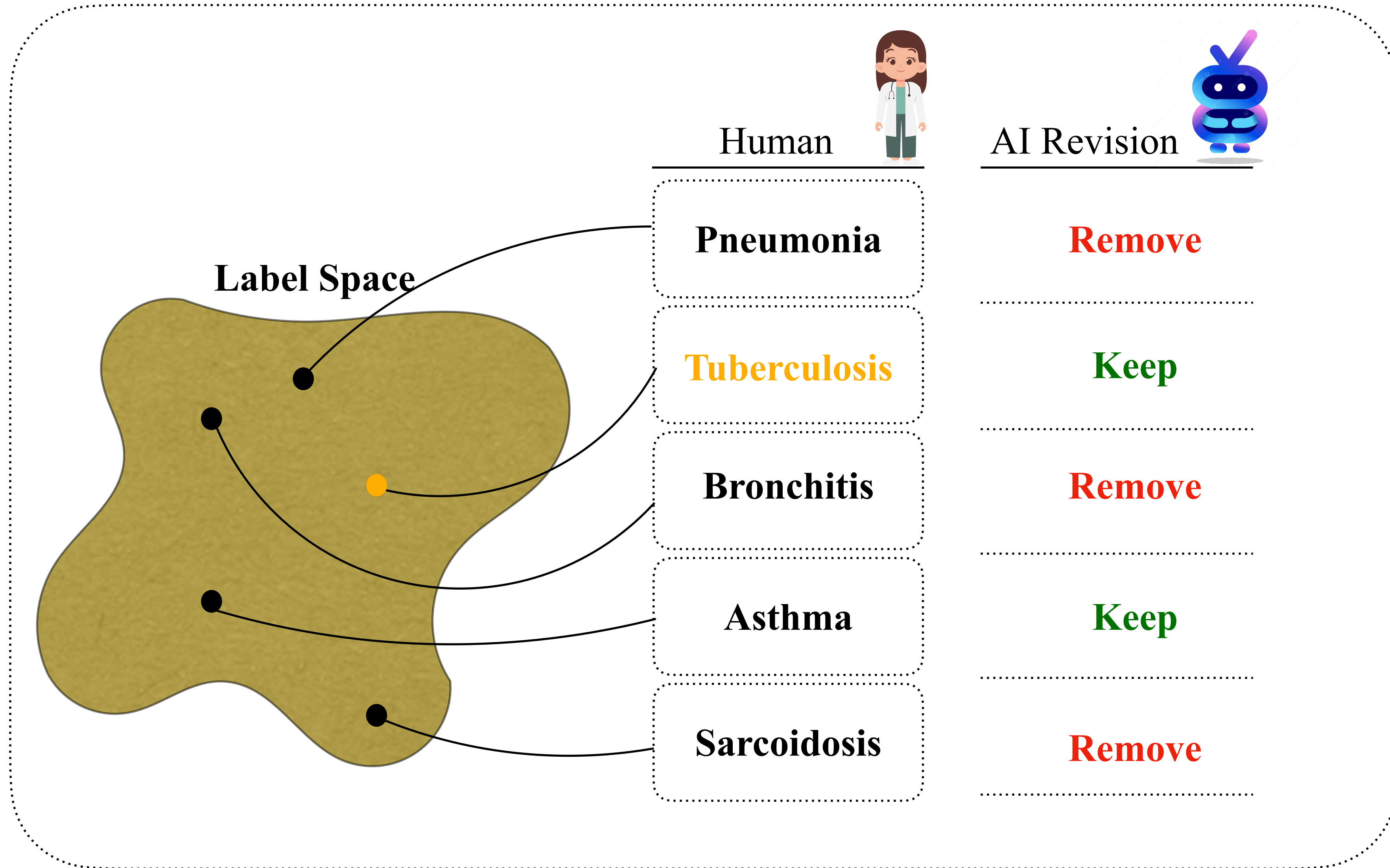
Question: what constitutes a good collaboration?



Question: what constitutes a good collaboration?

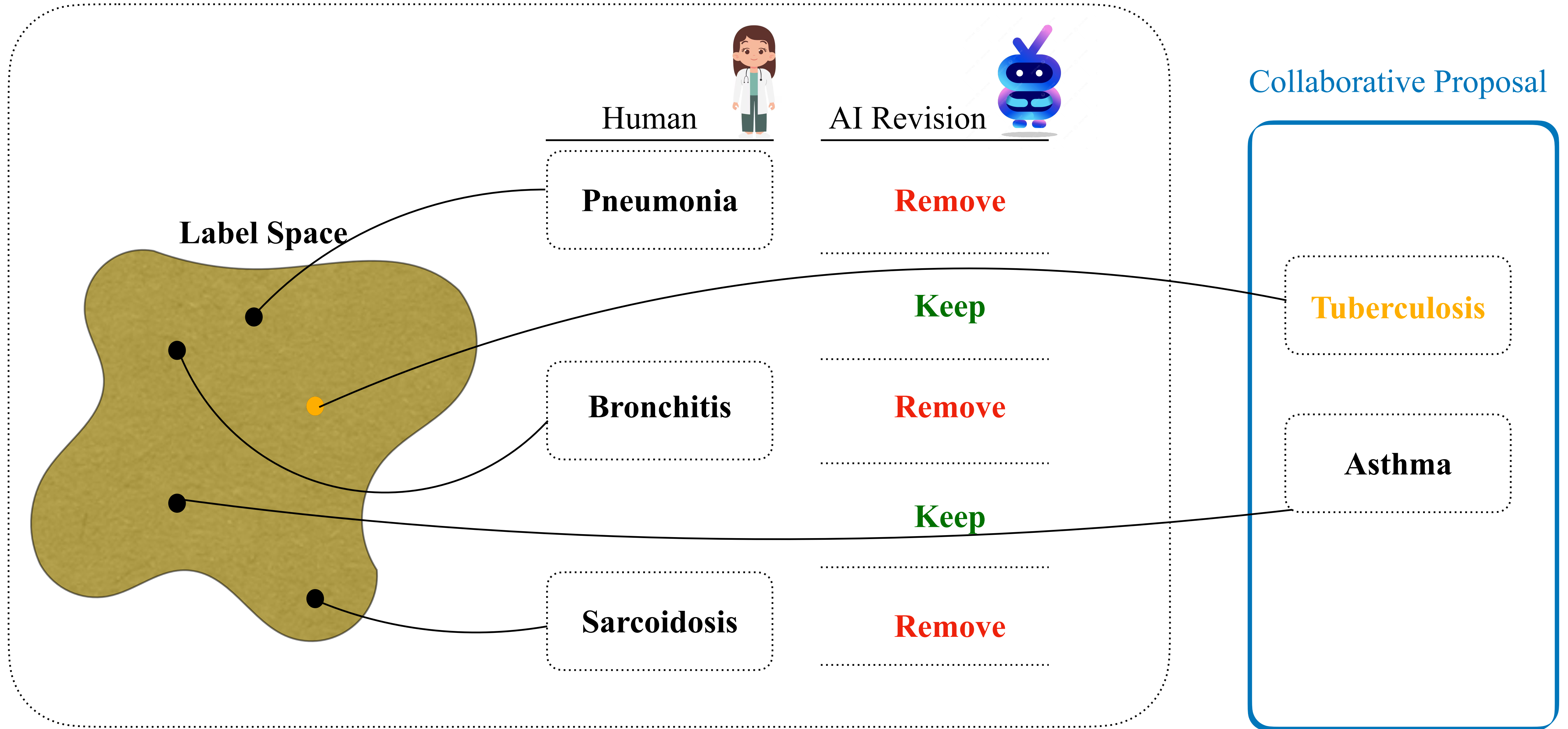


Question: what constitutes a good collaboration?

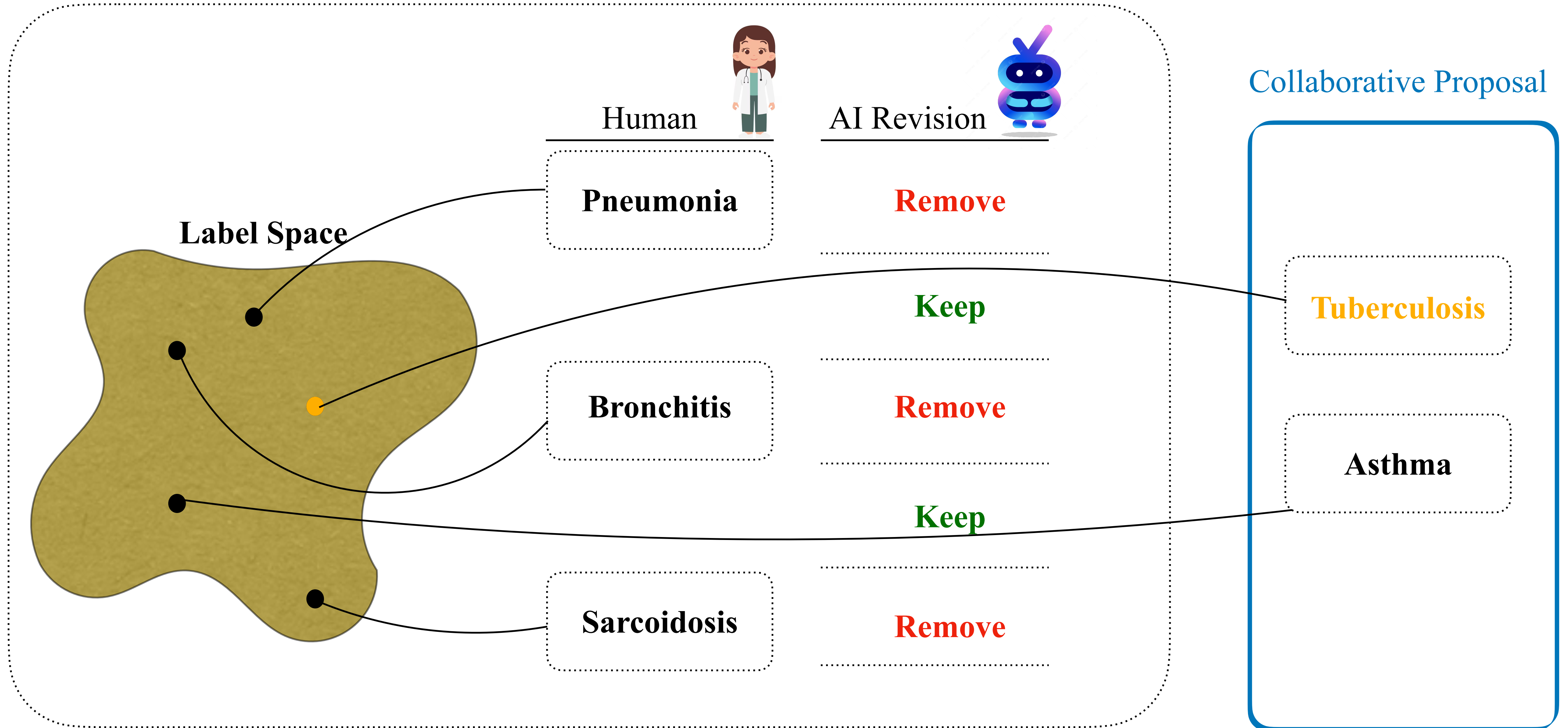


Collaborative Proposal

Question: what constitutes a good collaboration?



Question: what constitutes a good collaboration?



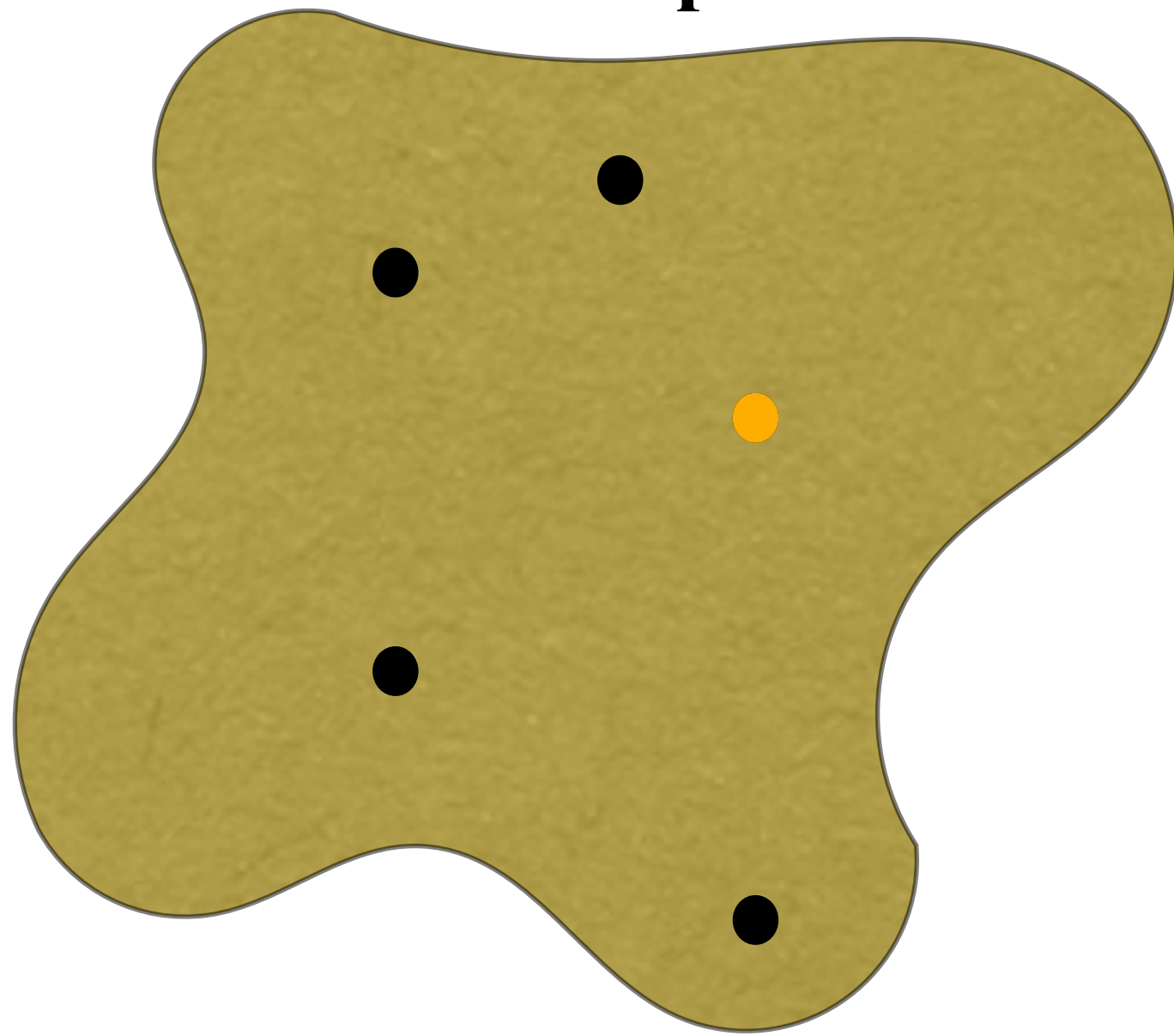
First: dont cause harm!

Question: what constitutes a good collaboration?

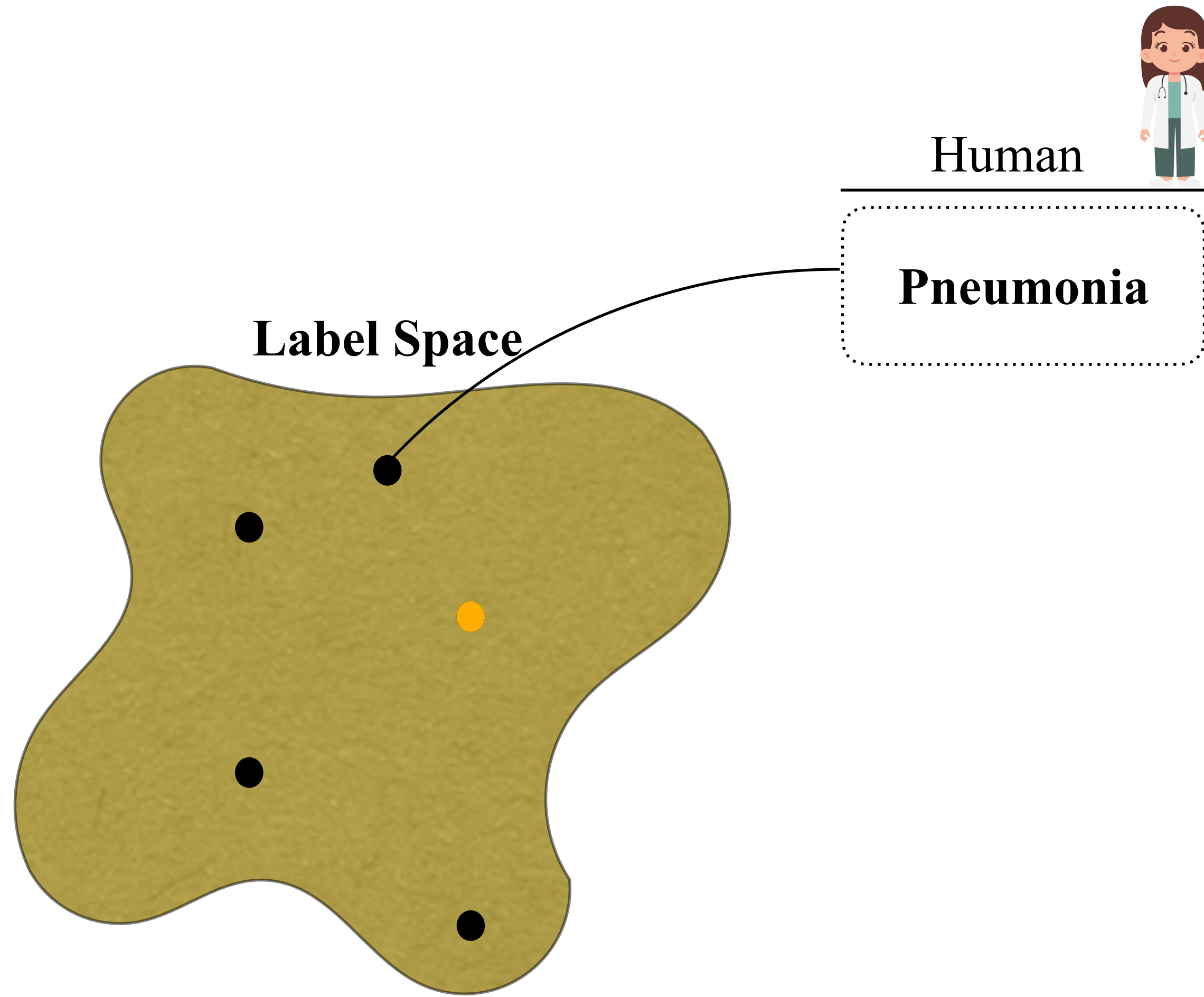
Human



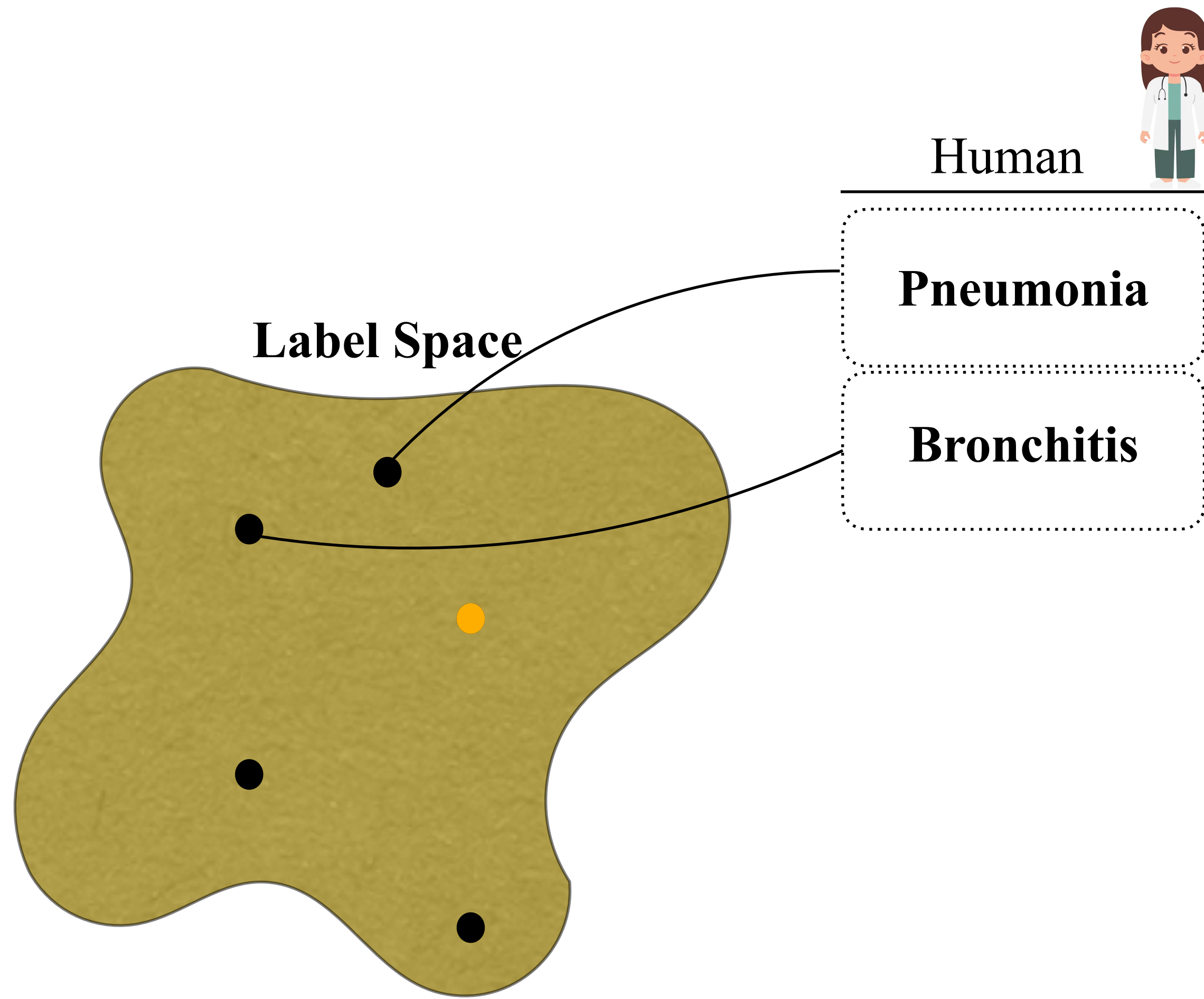
Label Space



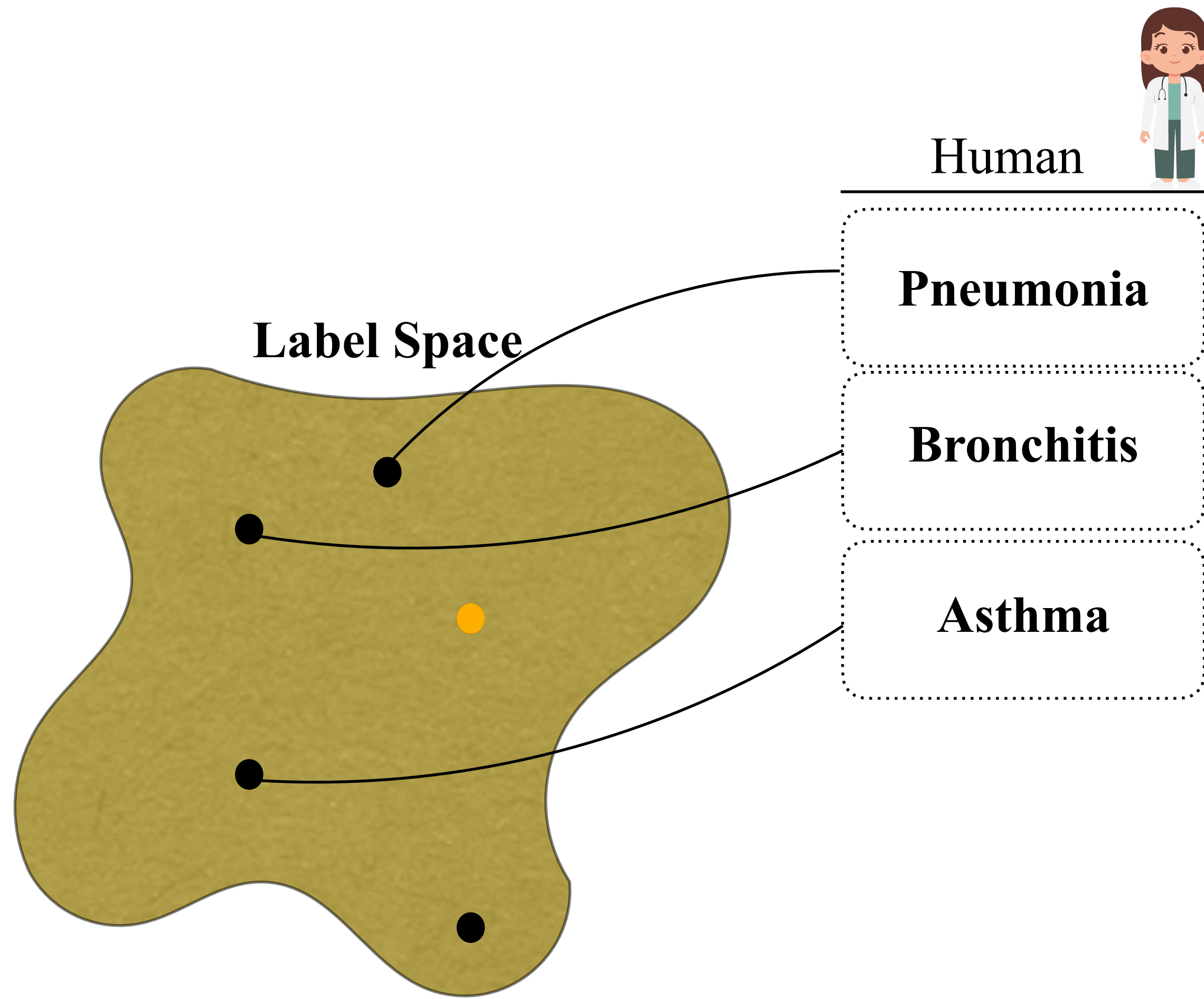
Question: what constitutes a good collaboration?



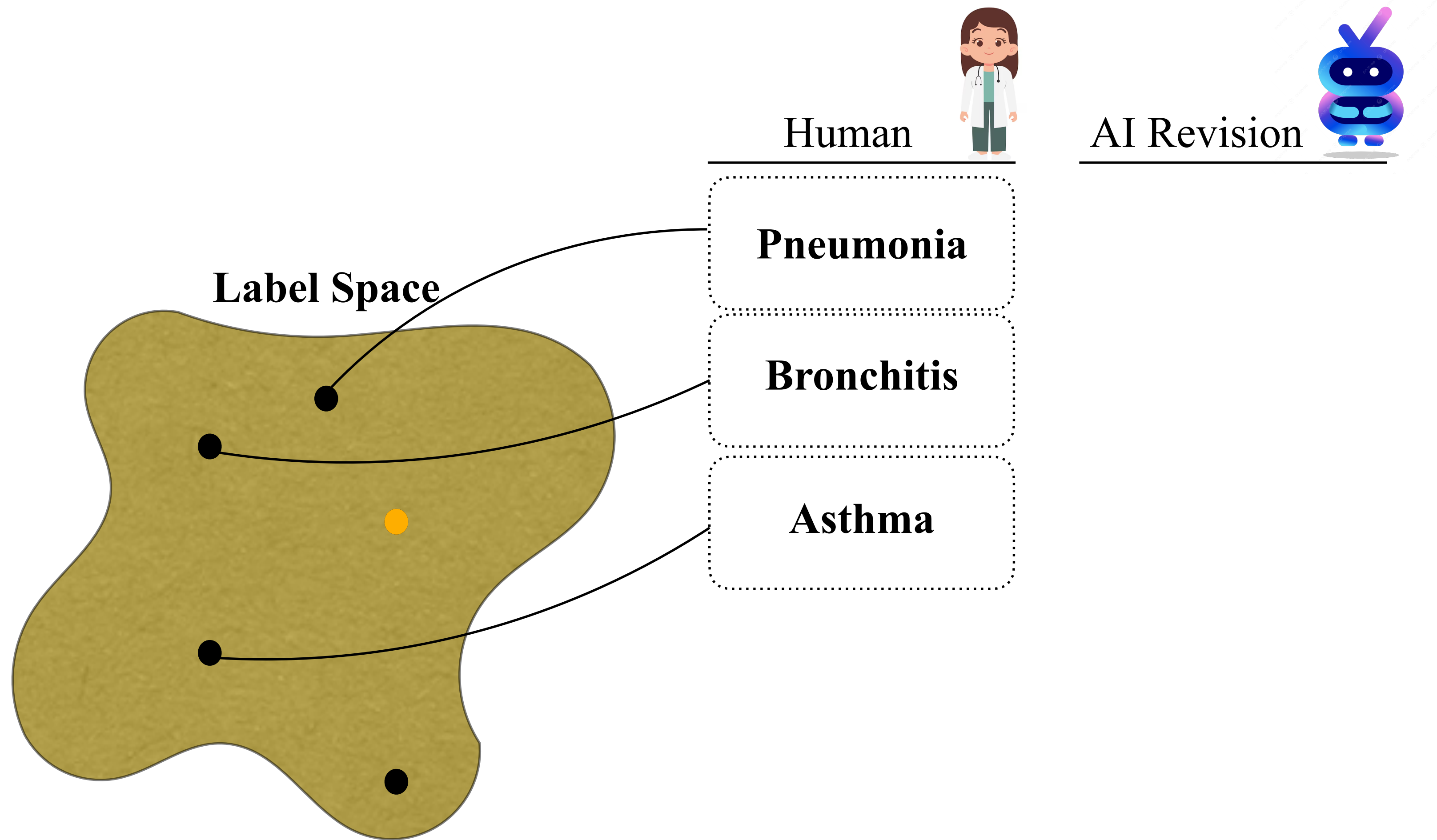
Question: what constitutes a good collaboration?



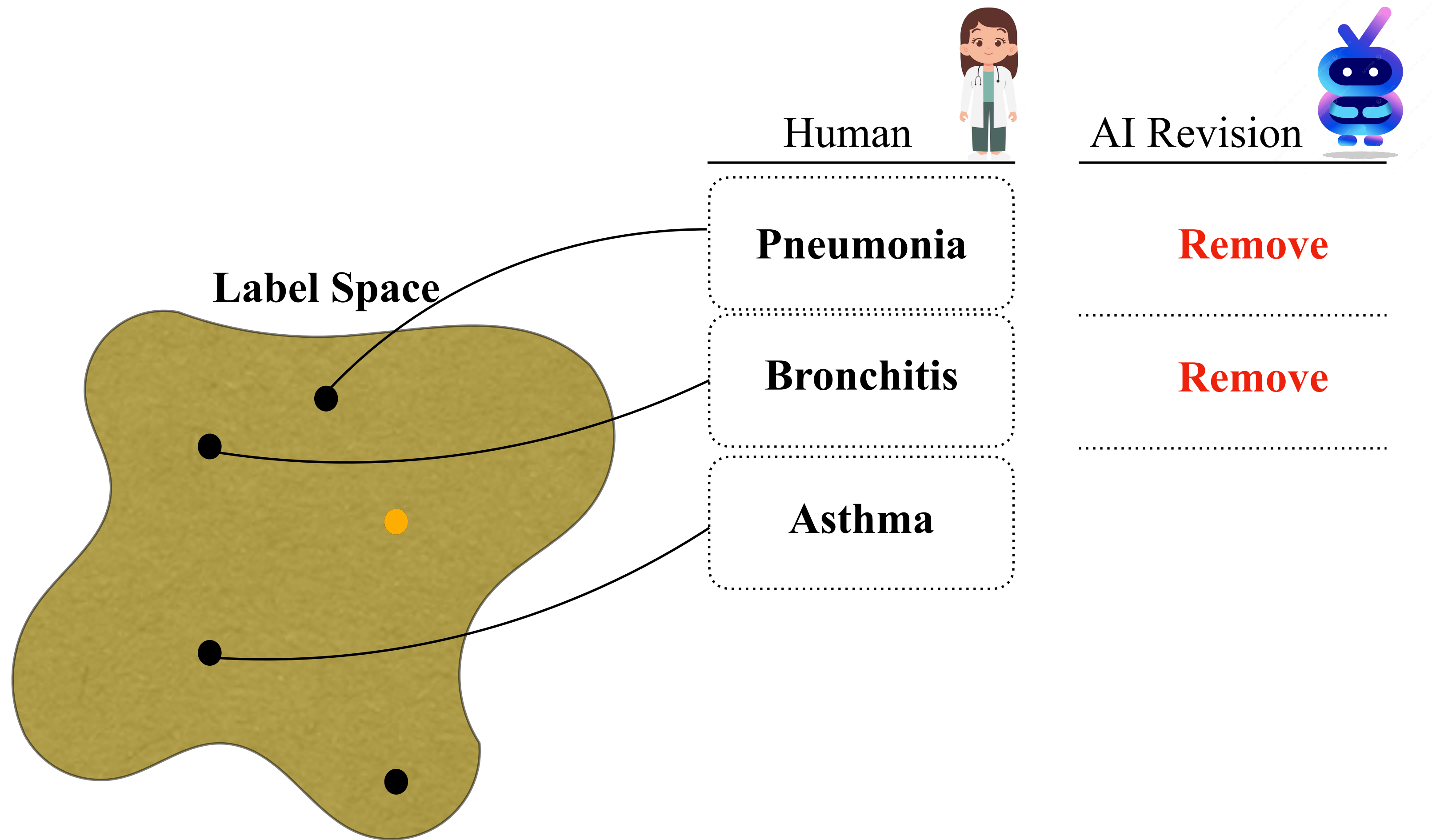
Question: what constitutes a good collaboration?



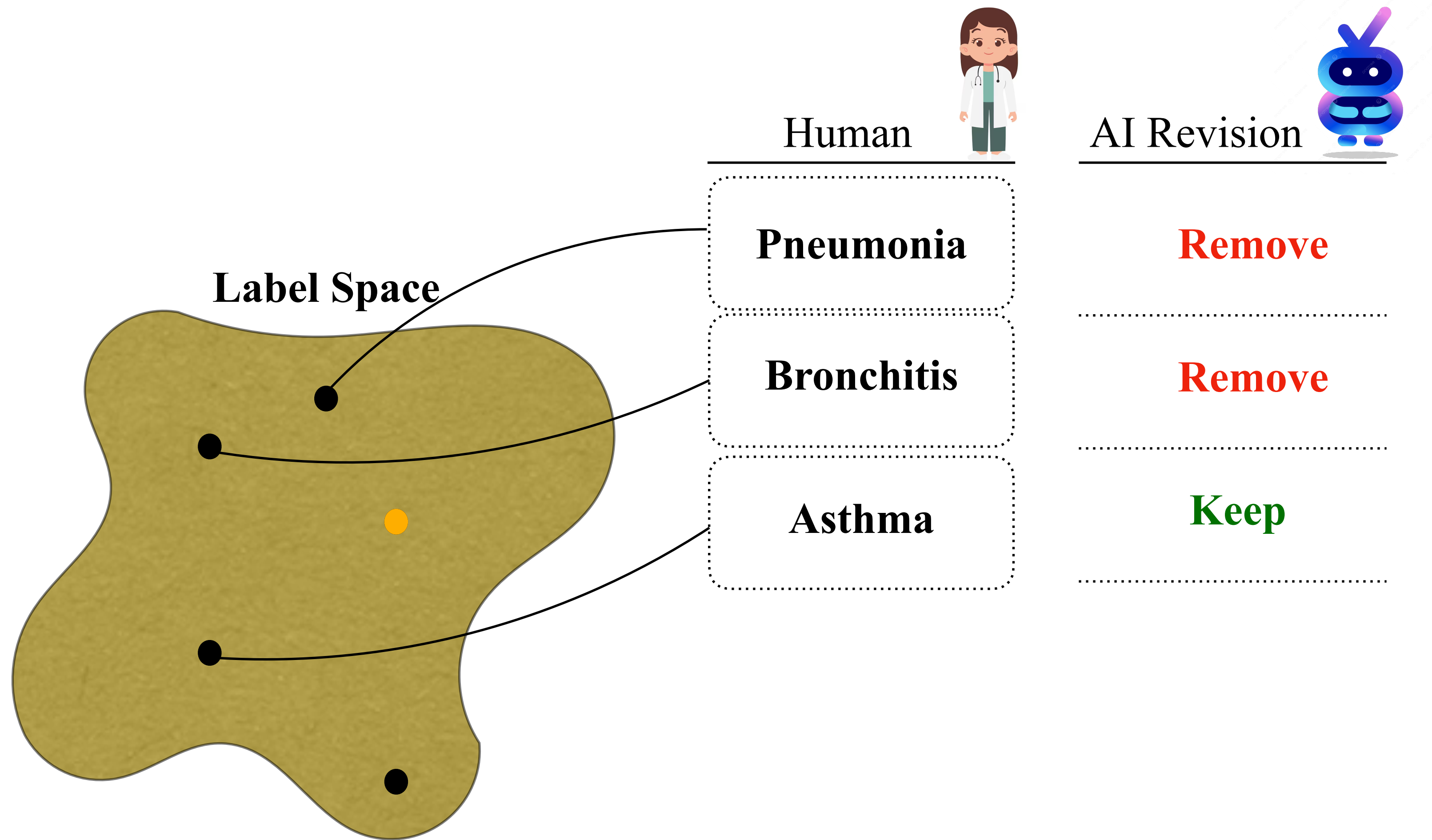
Question: what constitutes a good collaboration?



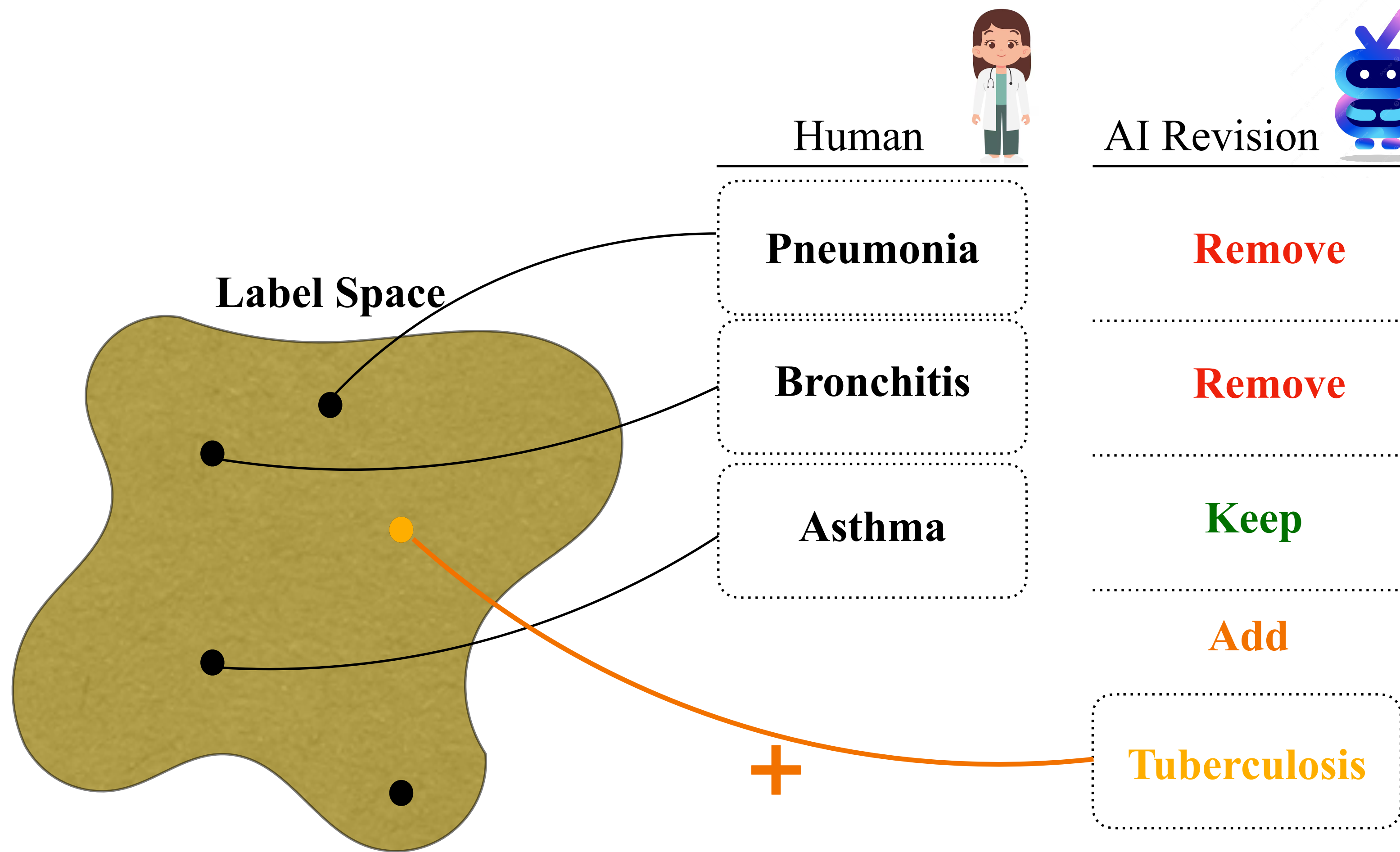
Question: what constitutes a good collaboration?



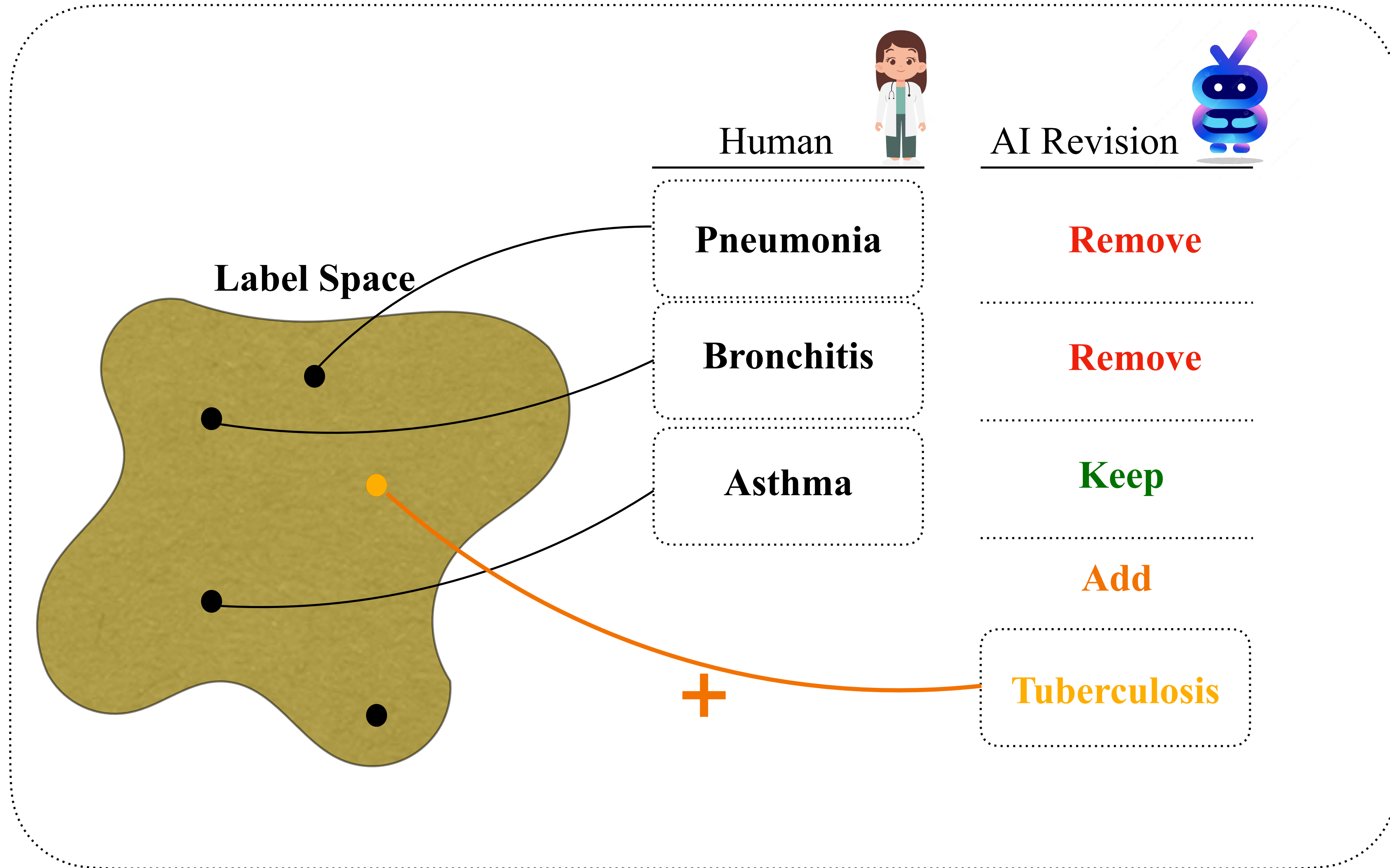
Question: what constitutes a good collaboration?



Question: what constitutes a good collaboration?

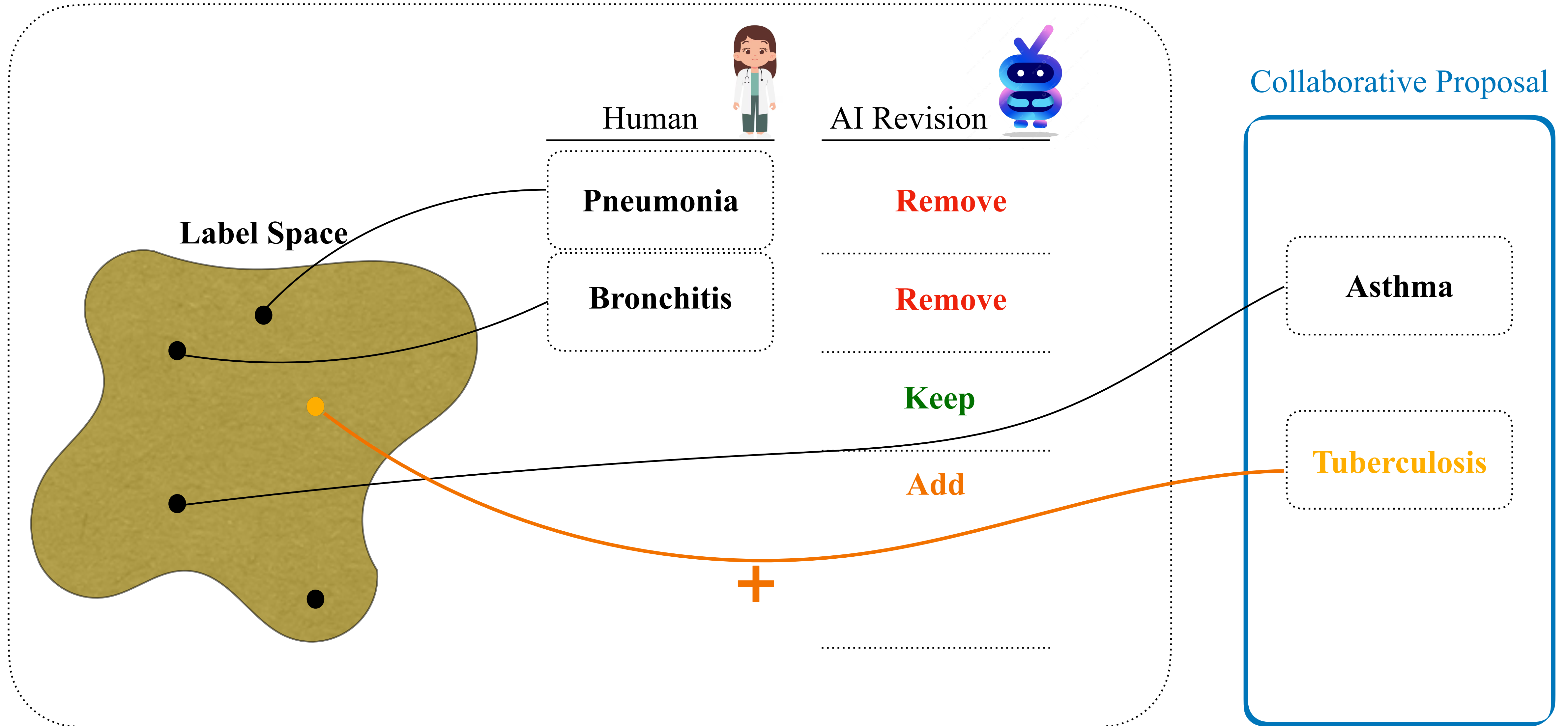


Question: what constitutes a good collaboration?

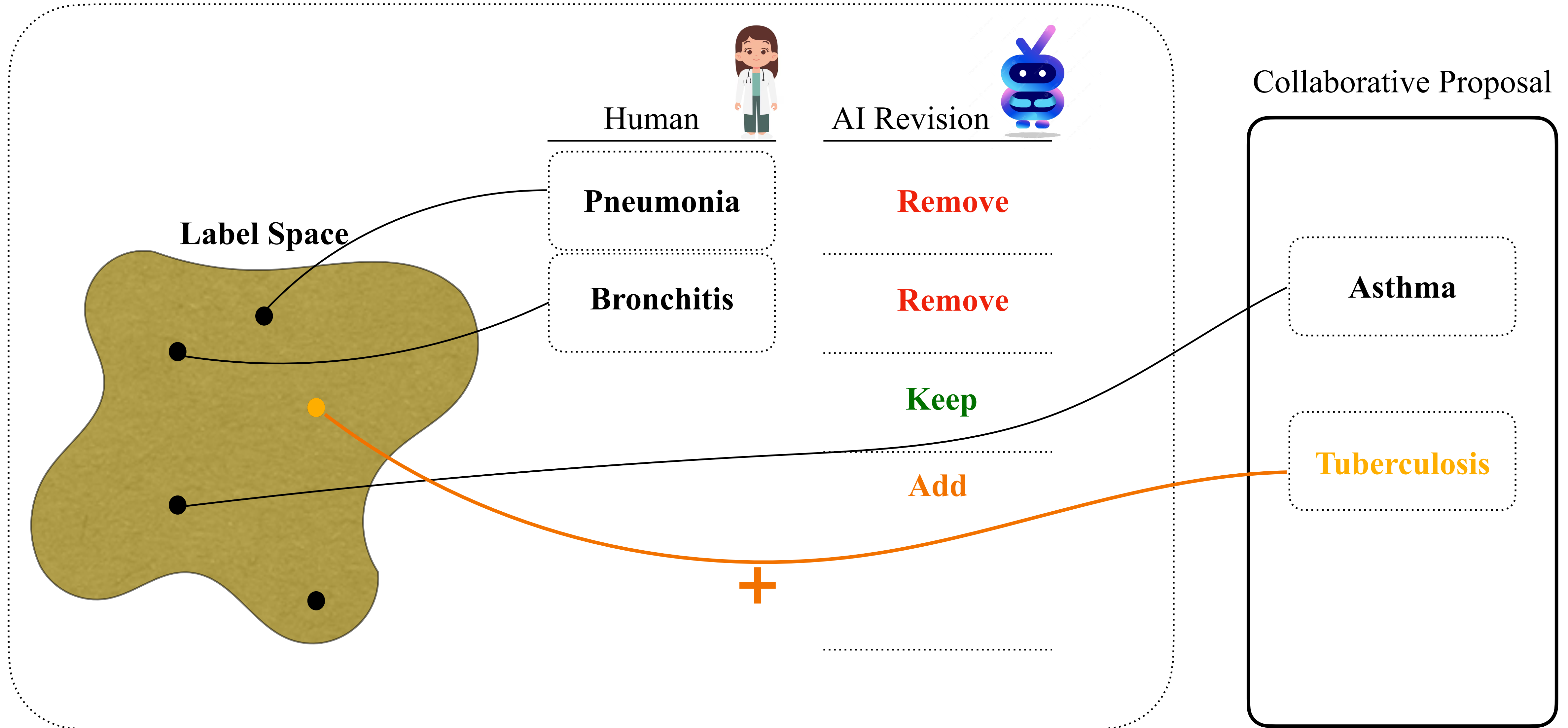


Collaborative Proposal

Question: what constitutes a good collaboration?



Question: what constitutes a good collaboration?



Second: Add value!

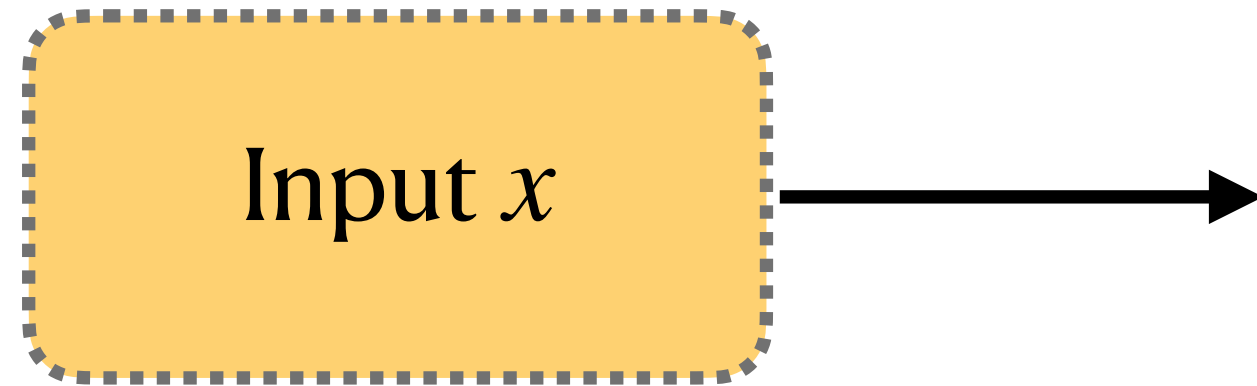
Two fundamentals of collaboration

Two fundamentals of collaboration

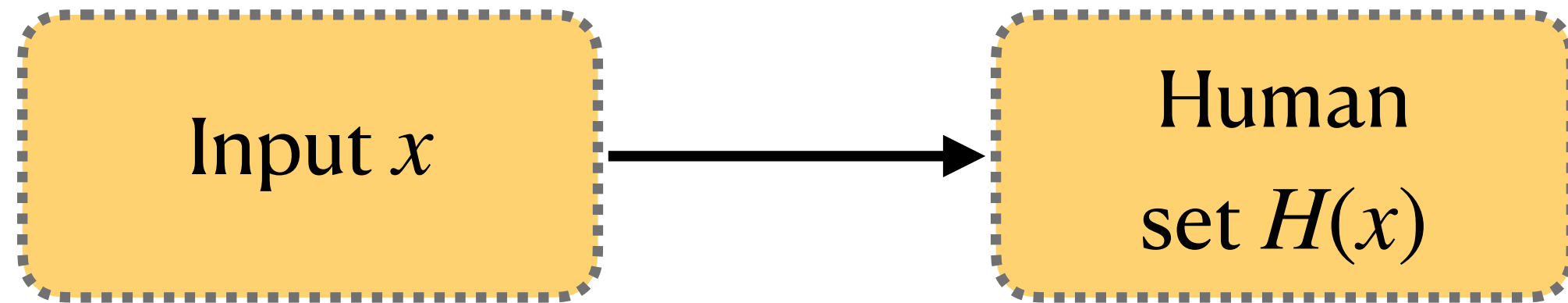


Input x

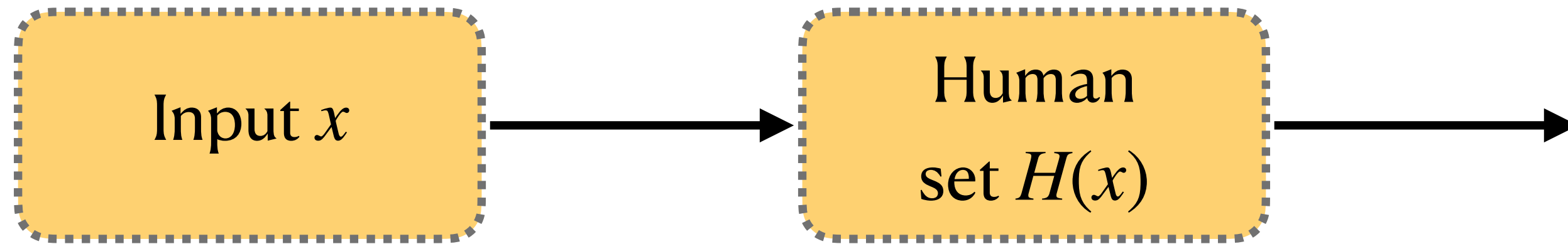
Two fundamentals of collaboration



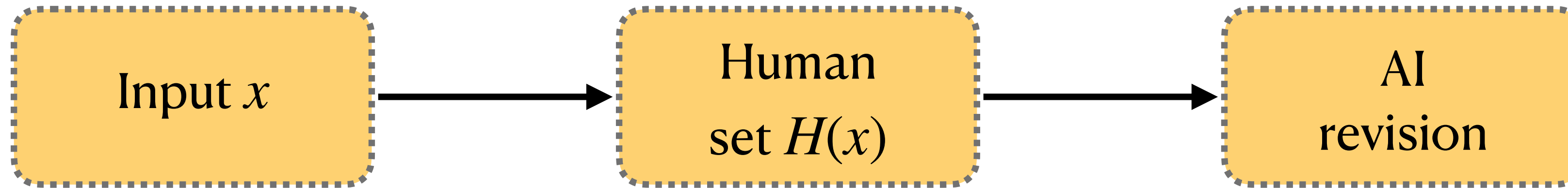
Two fundamentals of collaboration



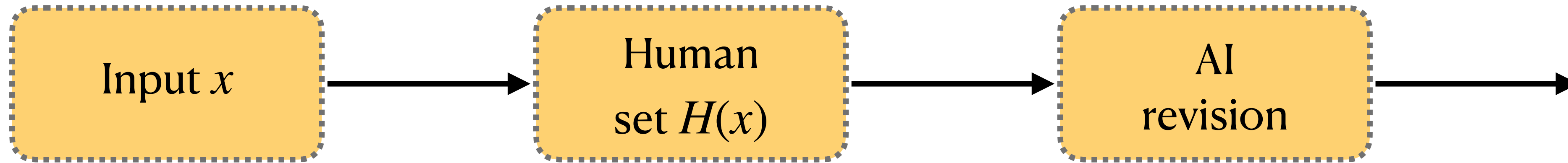
Two fundamentals of collaboration



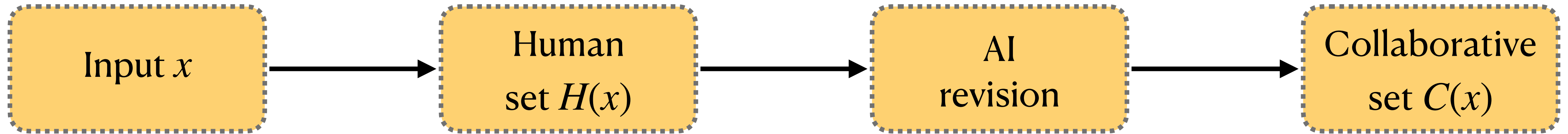
Two fundamentals of collaboration



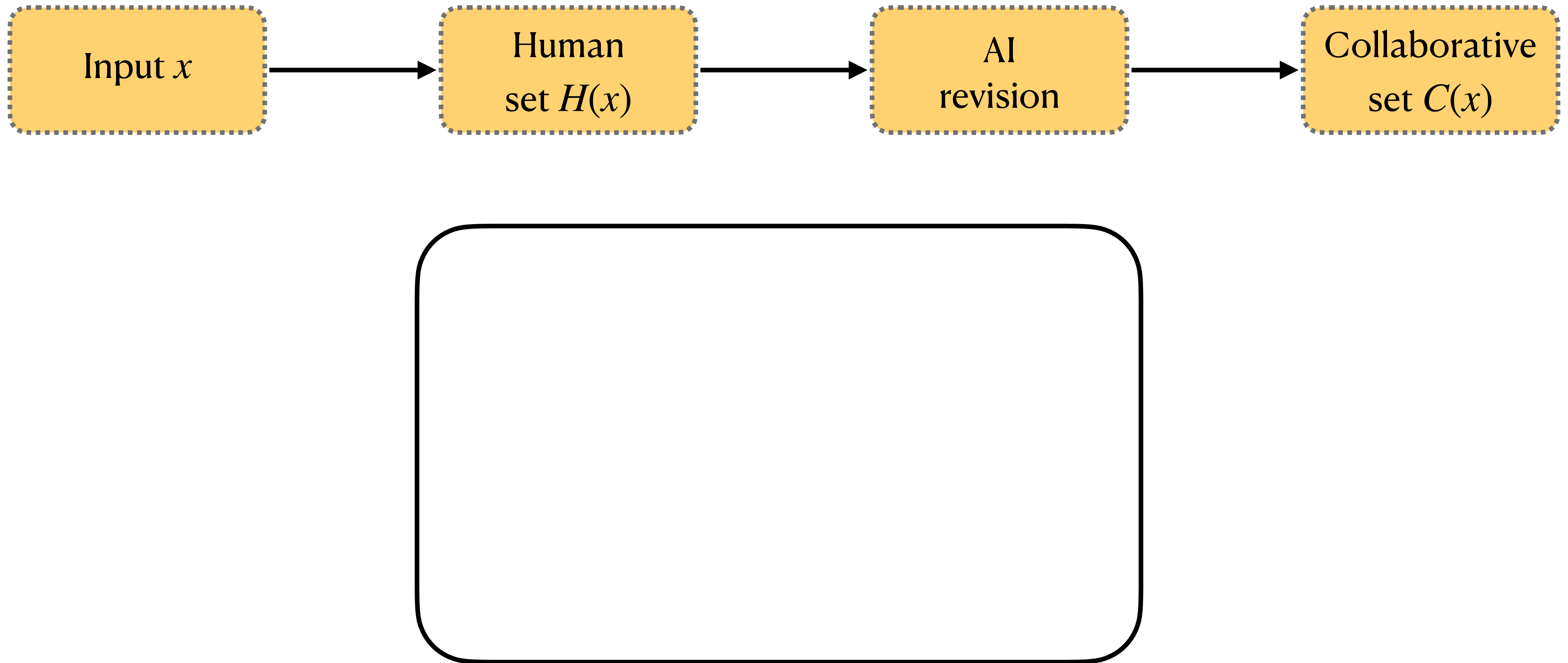
Two fundamentals of collaboration



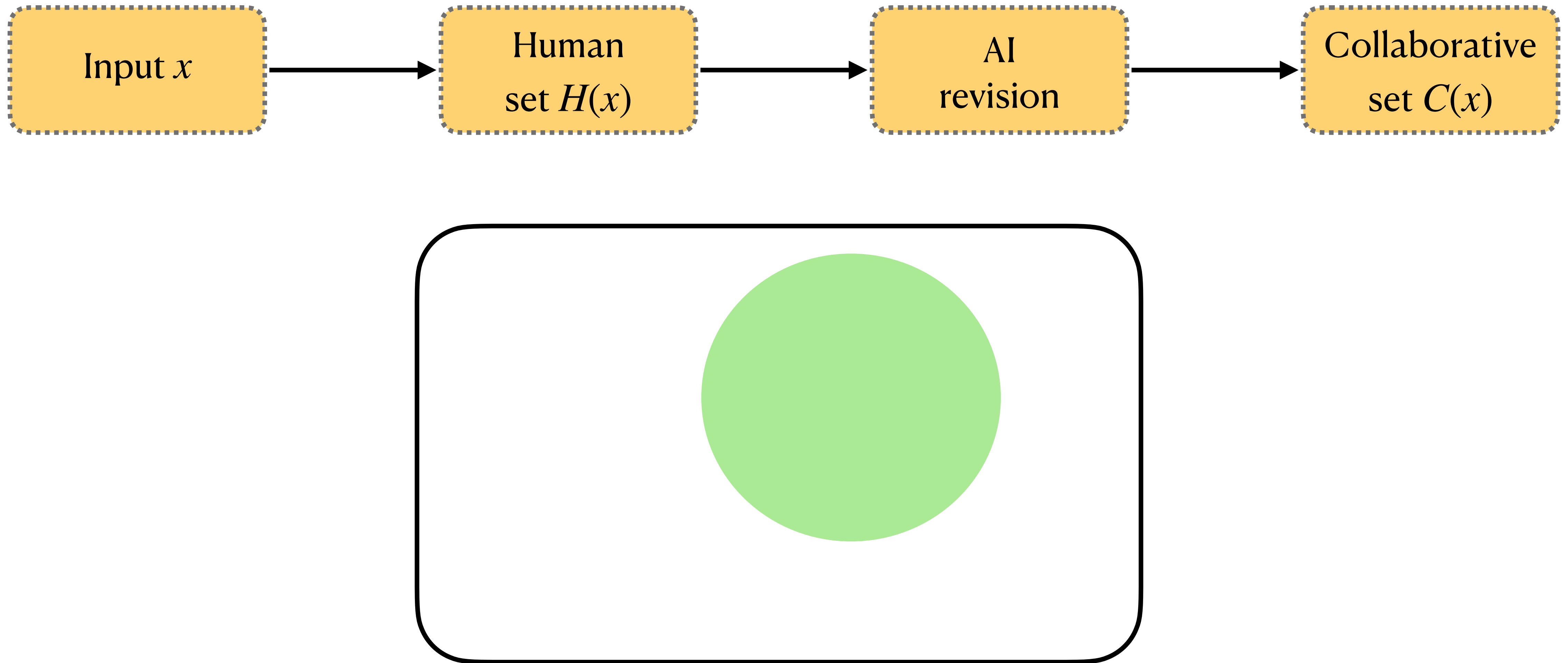
Two fundamentals of collaboration



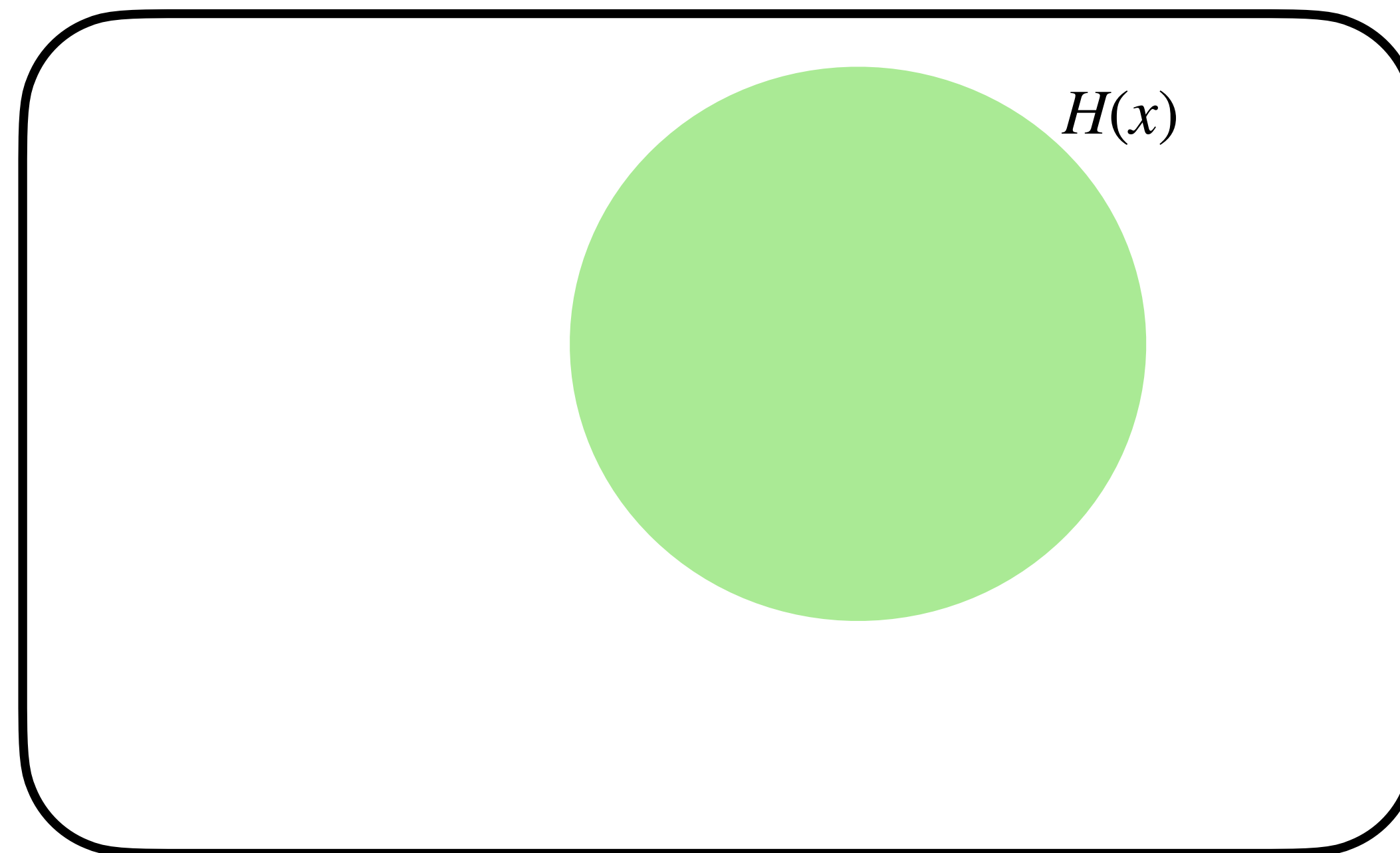
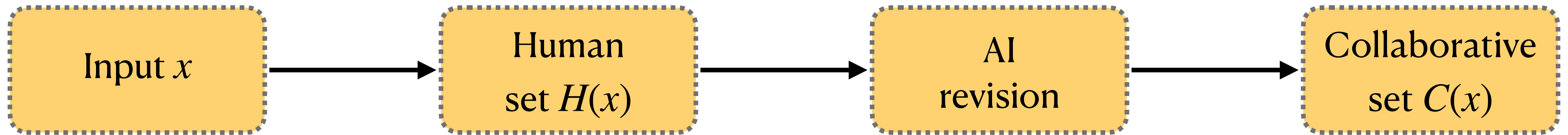
Two fundamentals of collaboration



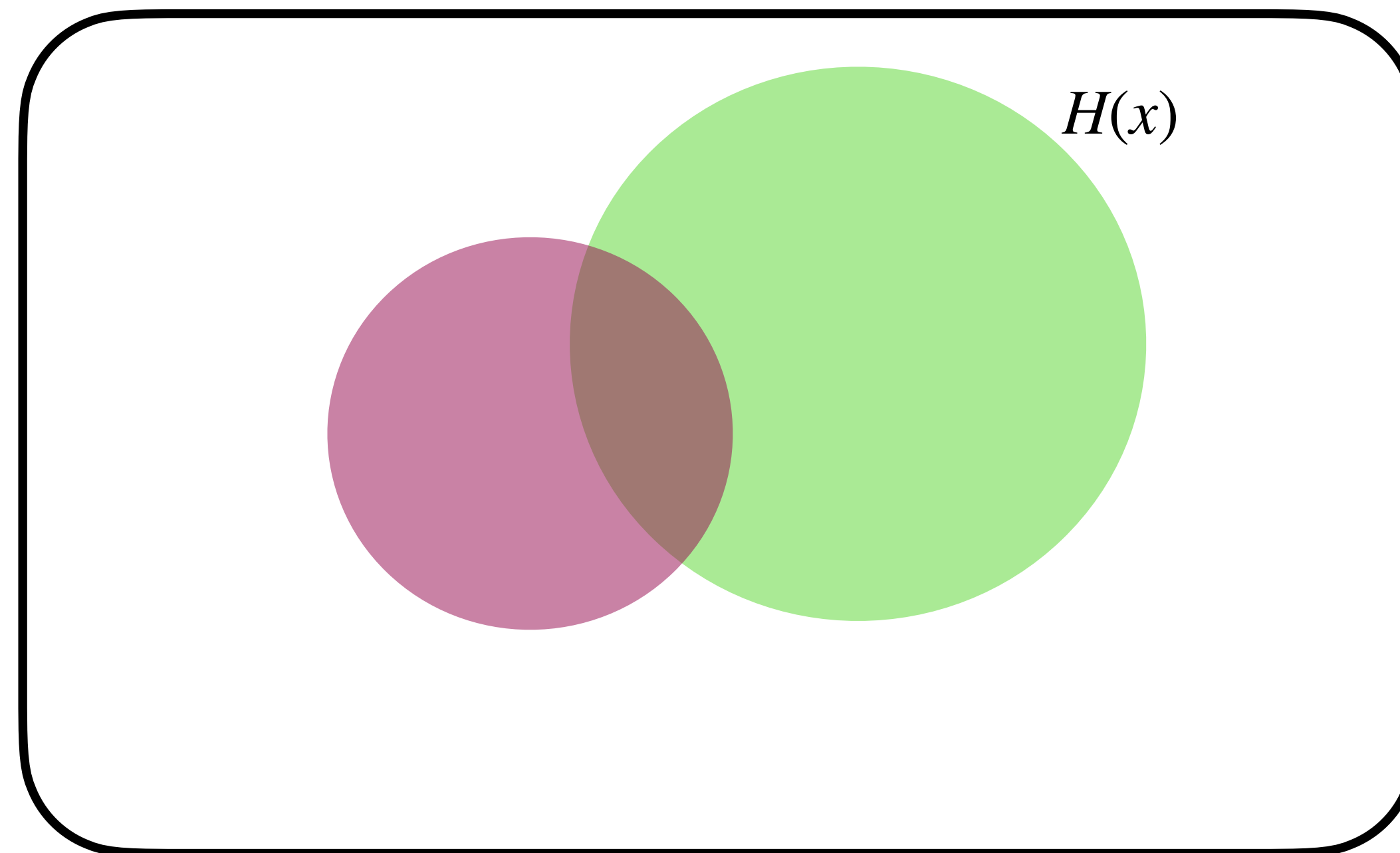
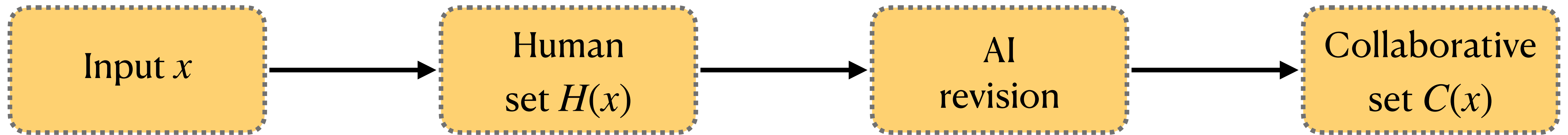
Two fundamentals of collaboration



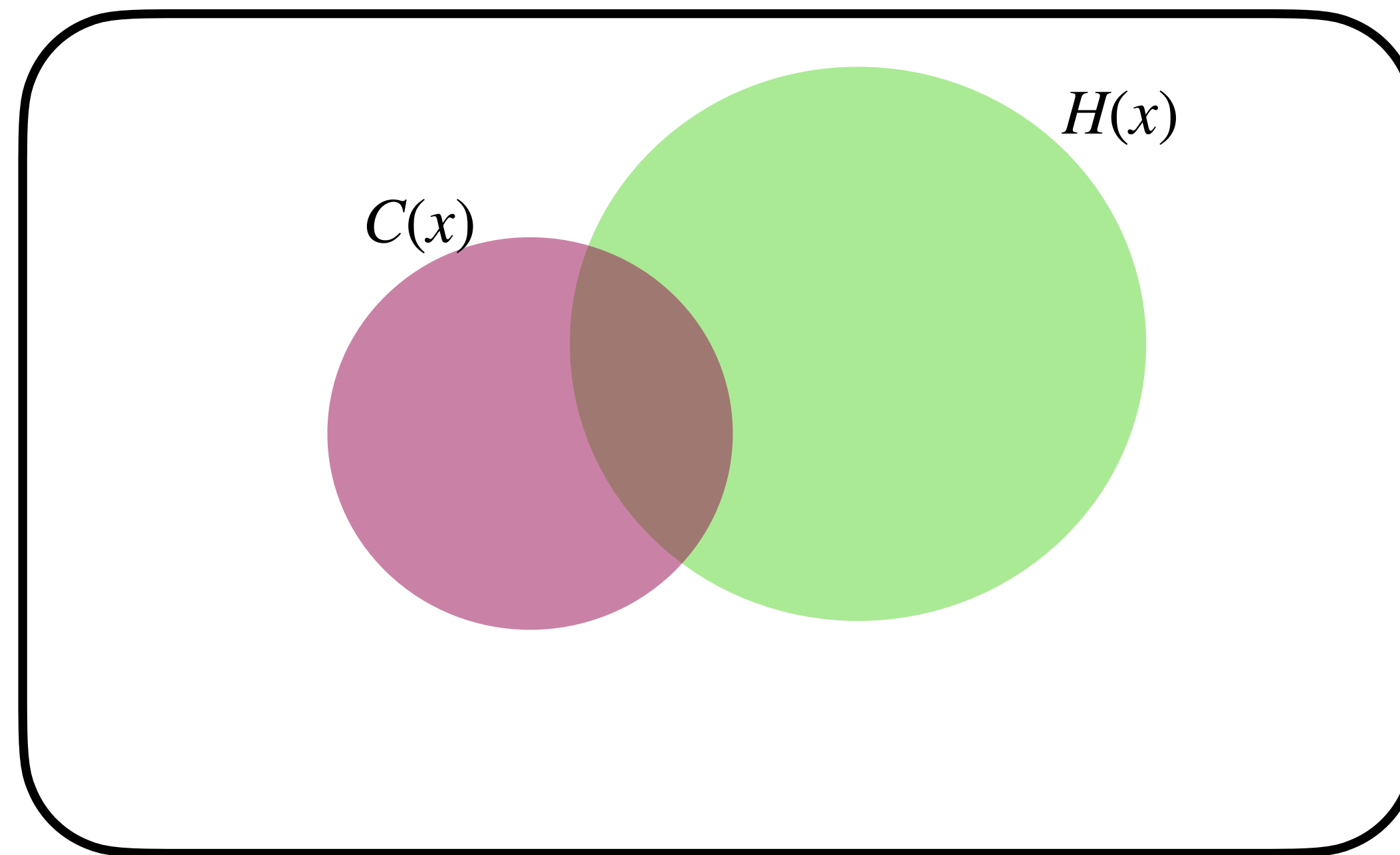
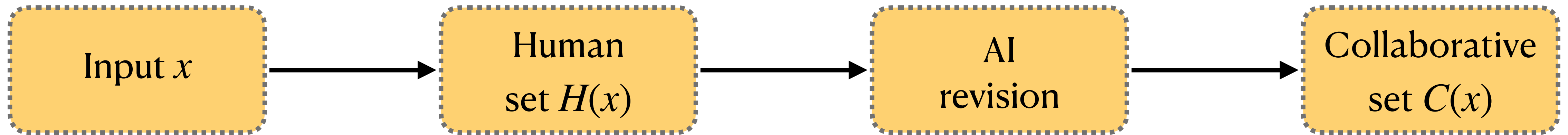
Two fundamentals of collaboration



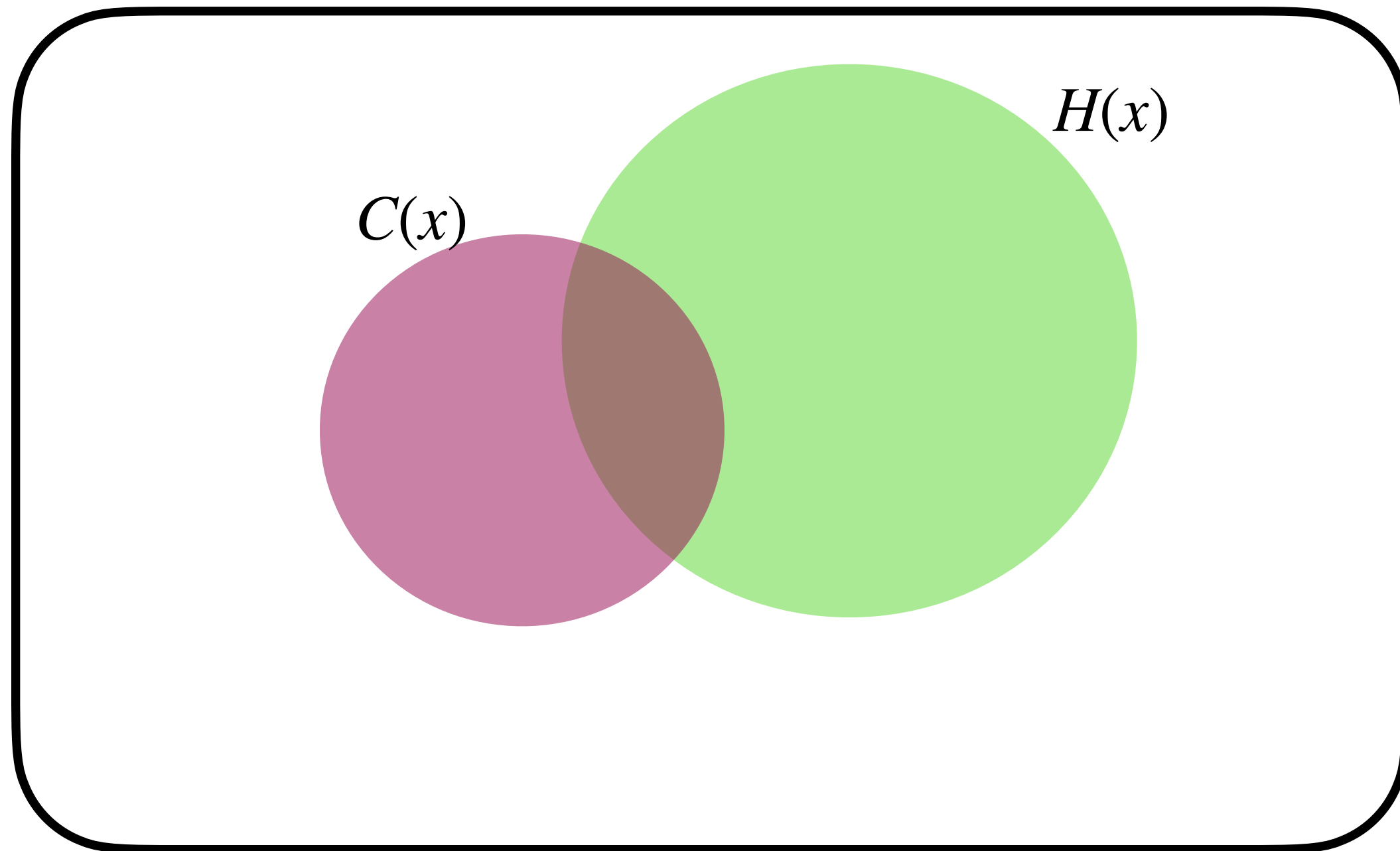
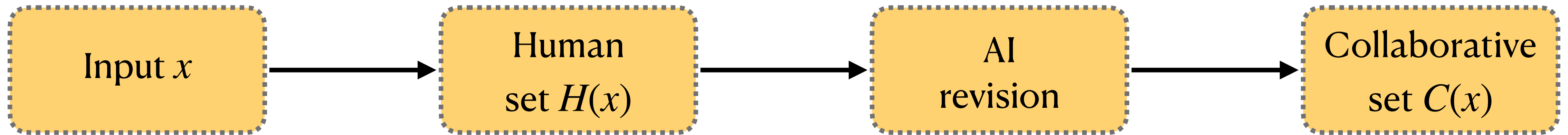
Two fundamentals of collaboration



Two fundamentals of collaboration

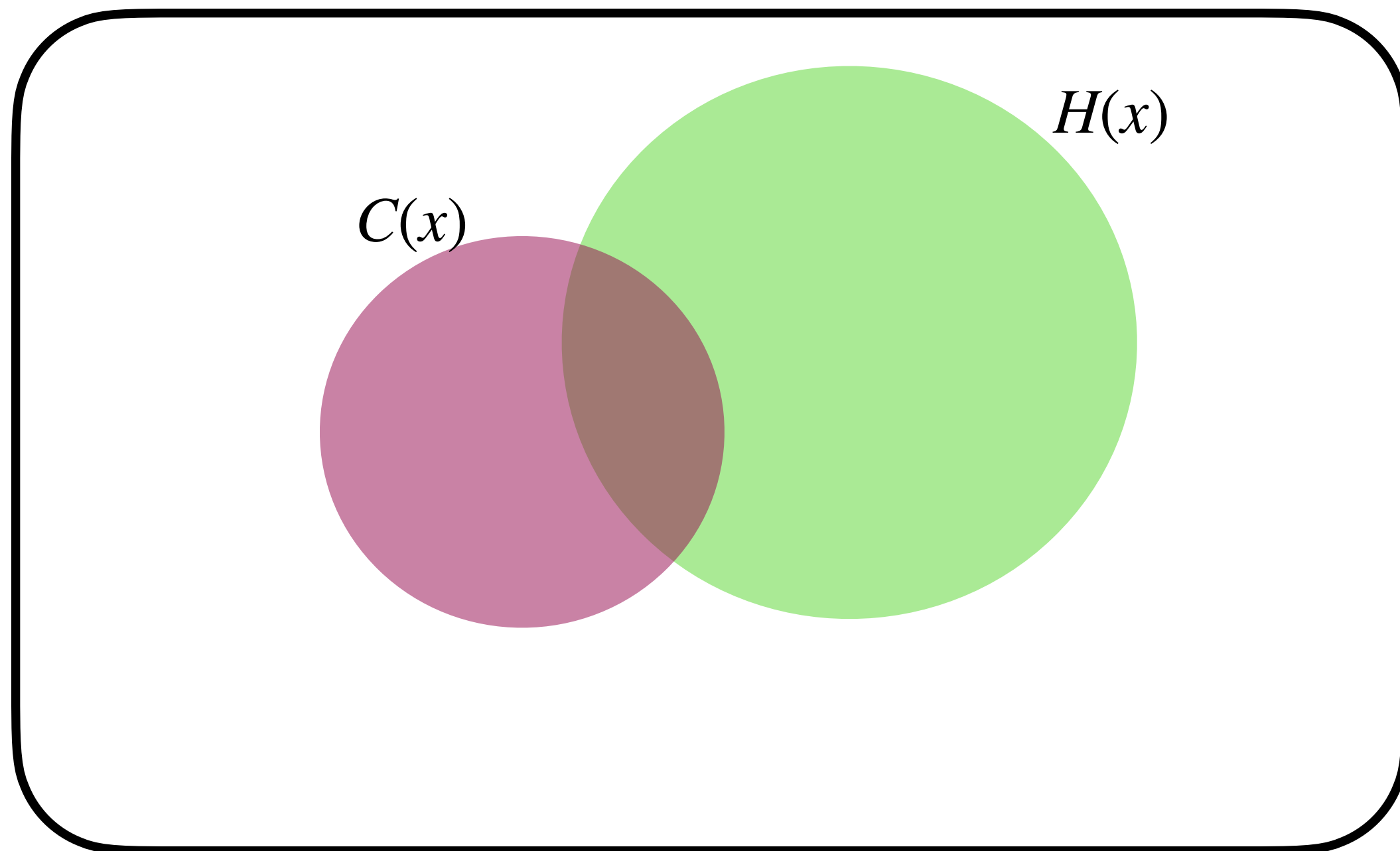
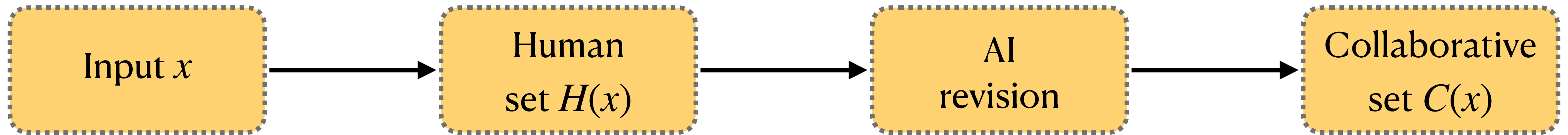


Two fundamentals of collaboration



Counterfactual Harm

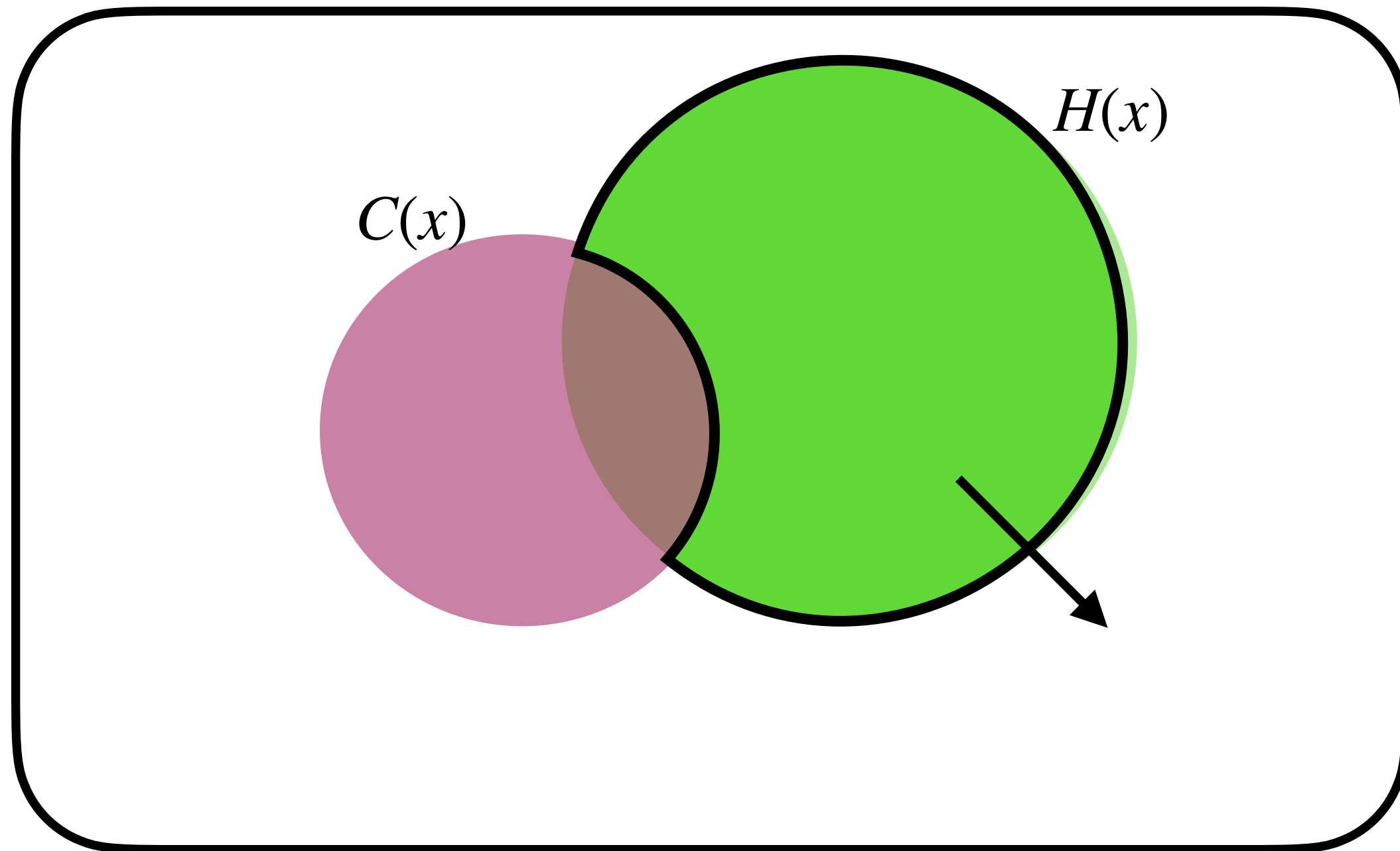
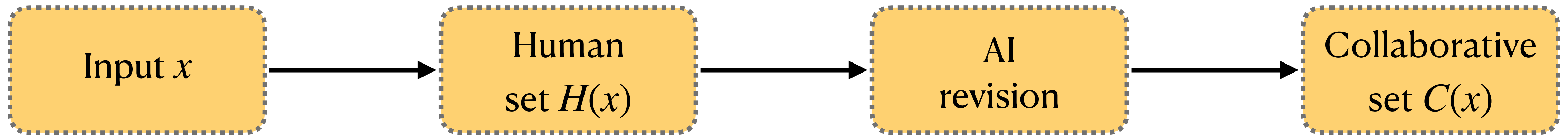
Two fundamentals of collaboration



Counterfactual Harm

$$\mathbb{P}(Y \notin C(X) \mid Y \in H(X)) < \varepsilon$$

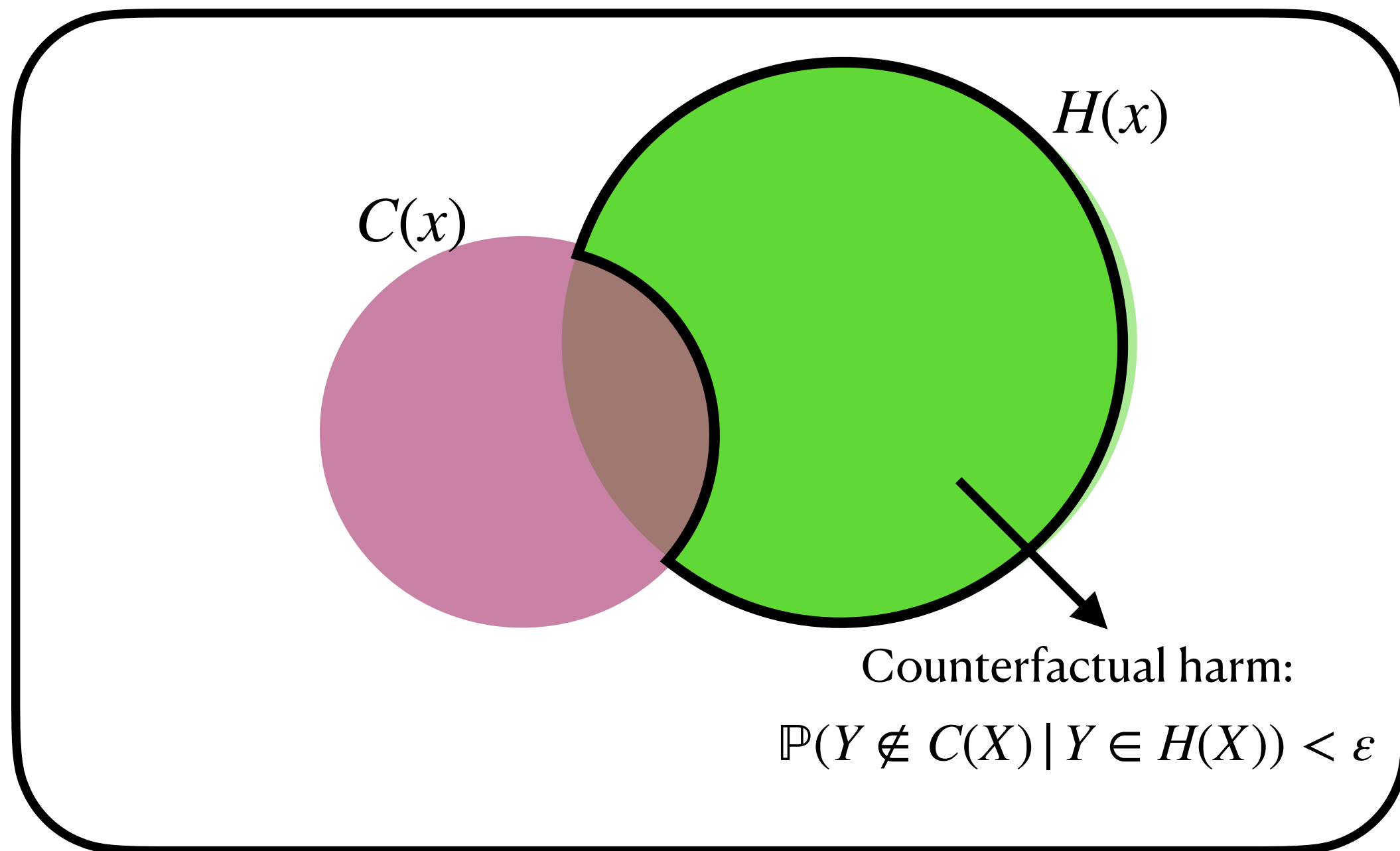
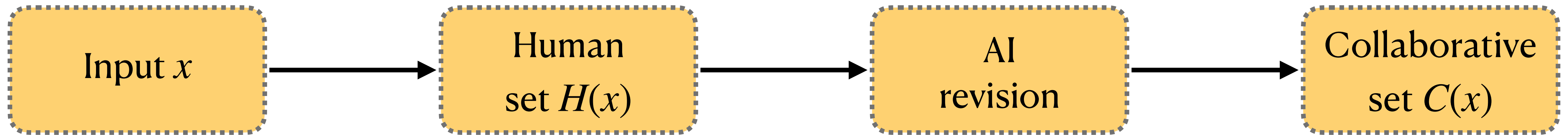
Two fundamentals of collaboration



Counterfactual Harm

$$\mathbb{P}(Y \notin C(X) \mid Y \in H(X)) < \varepsilon$$

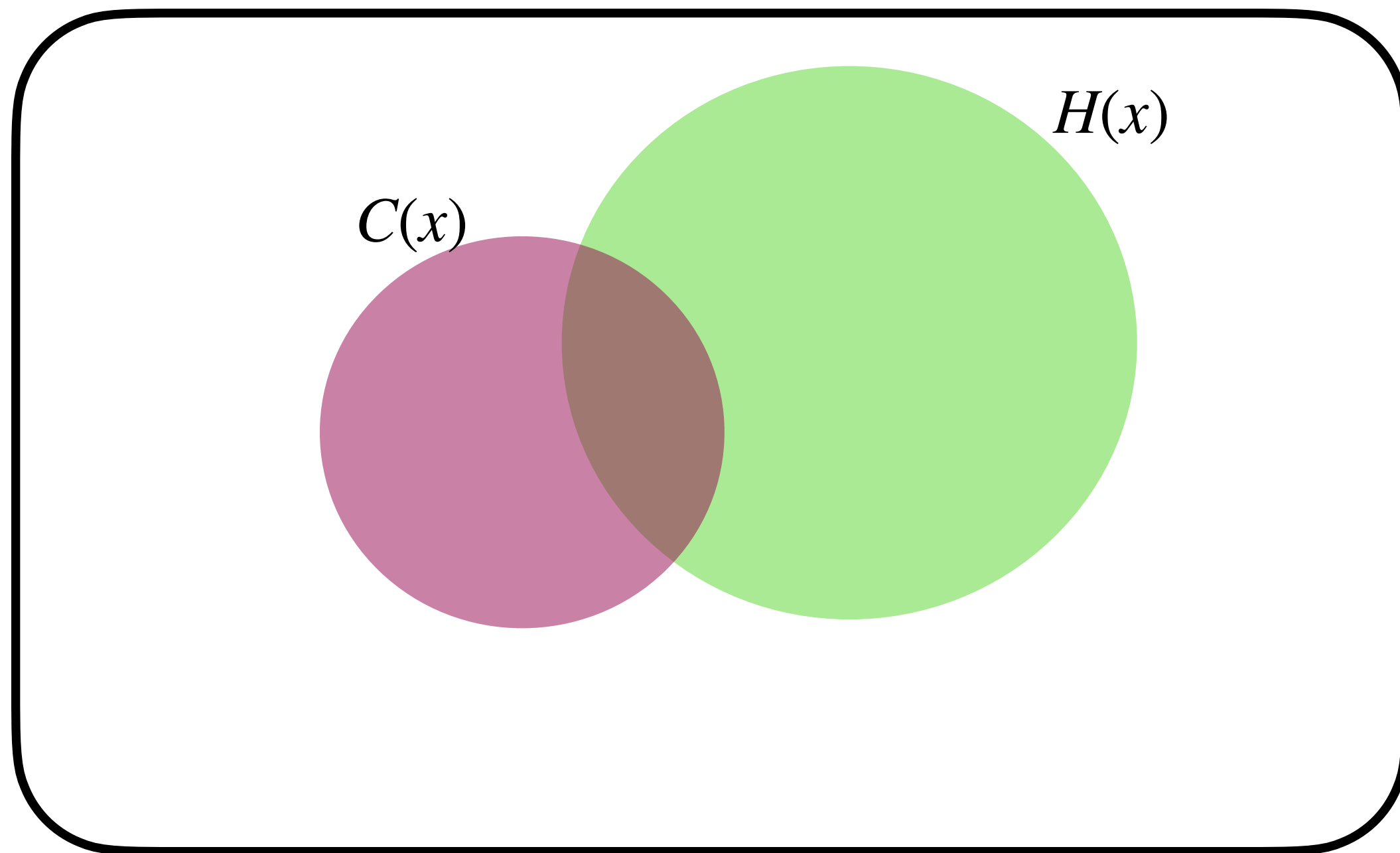
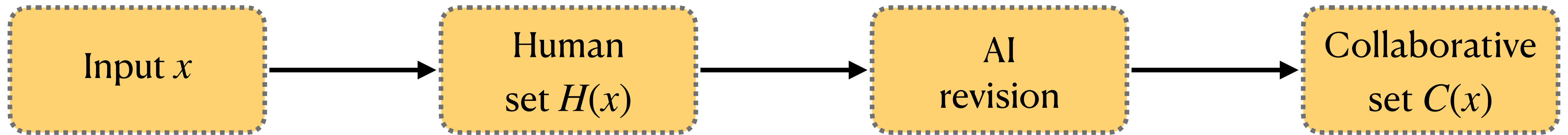
Two fundamentals of collaboration



Counterfactual Harm

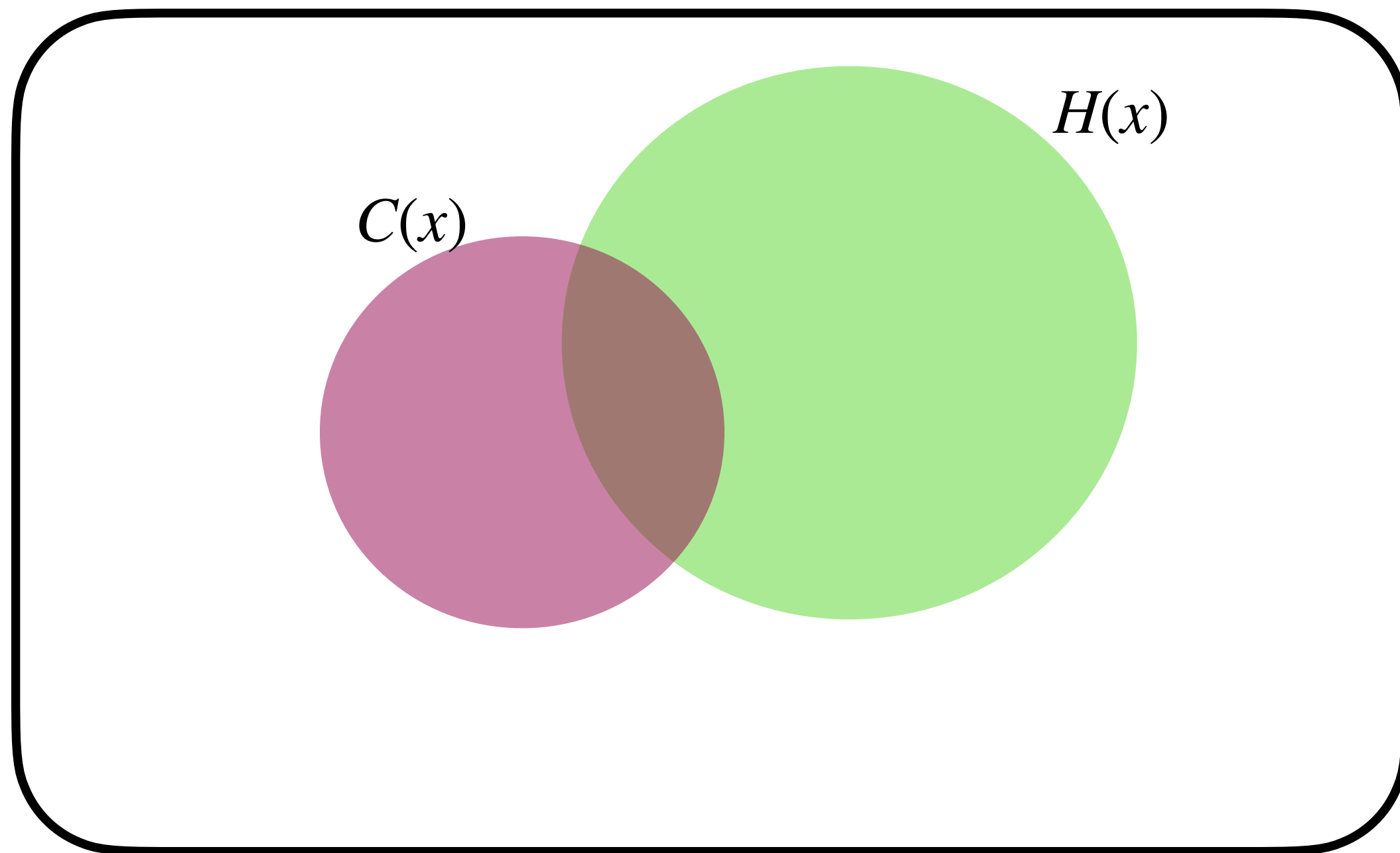
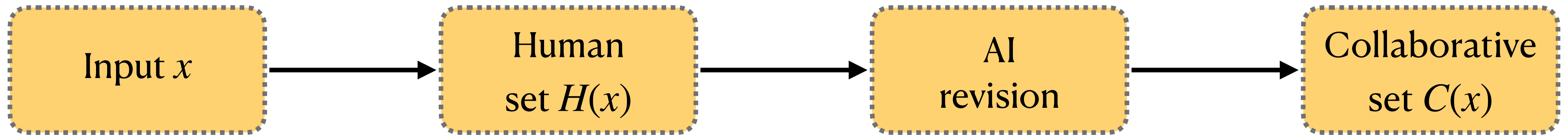
$$\mathbb{P}(Y \notin C(X) | Y \in H(X)) < \varepsilon$$

Two fundamentals of collaboration



Complementarity

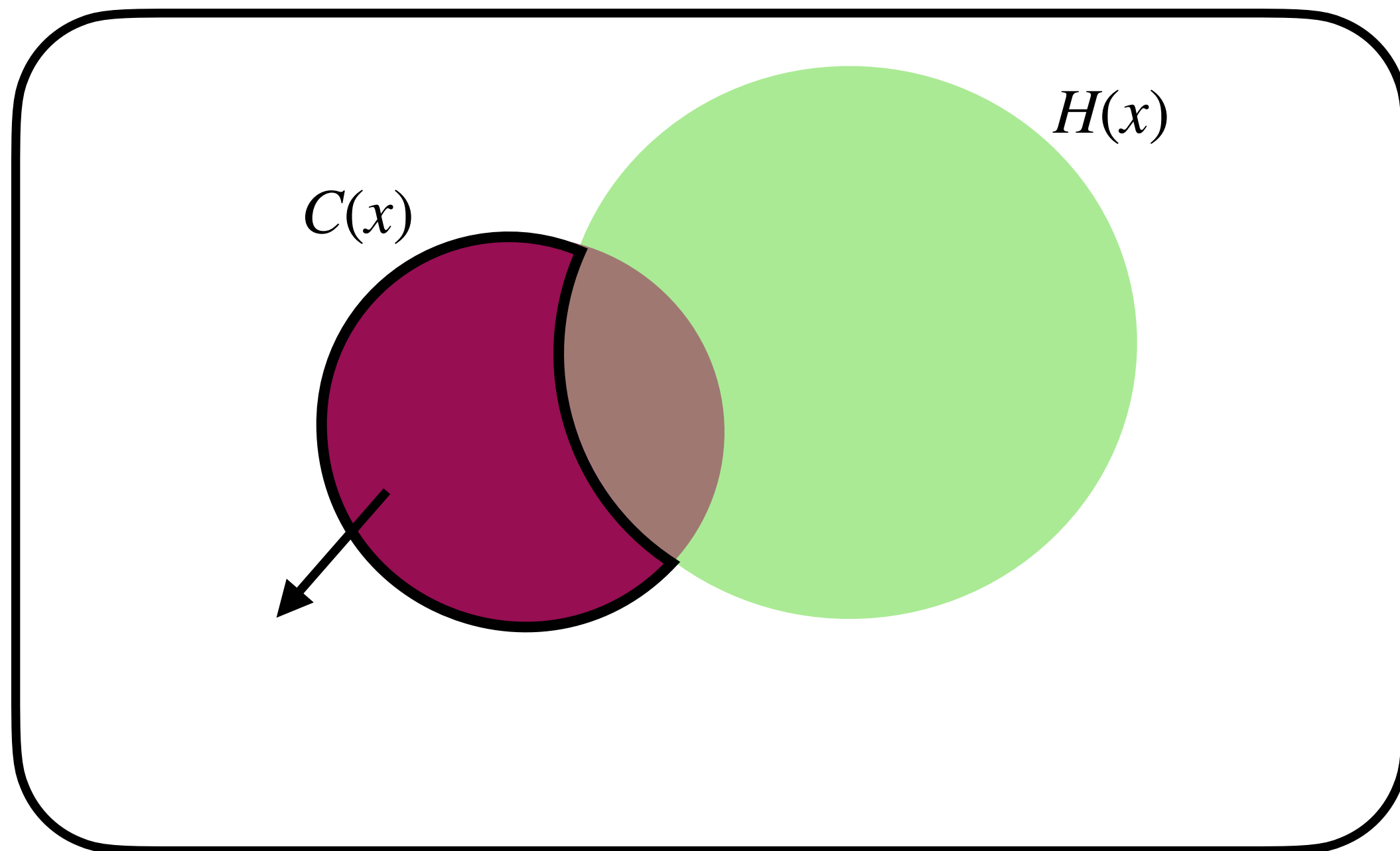
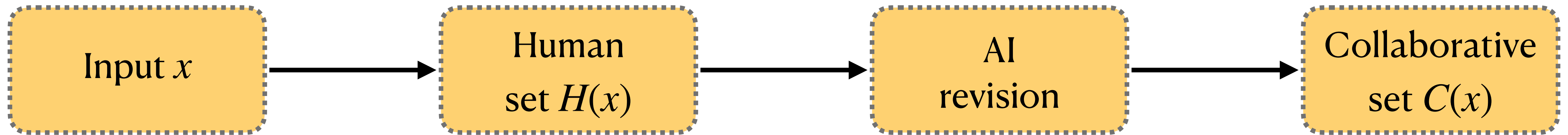
Two fundamentals of collaboration



Complementarity

$$\mathbb{P}(Y \in C(X) \mid Y \notin H(X)) \geq 1 - \delta$$

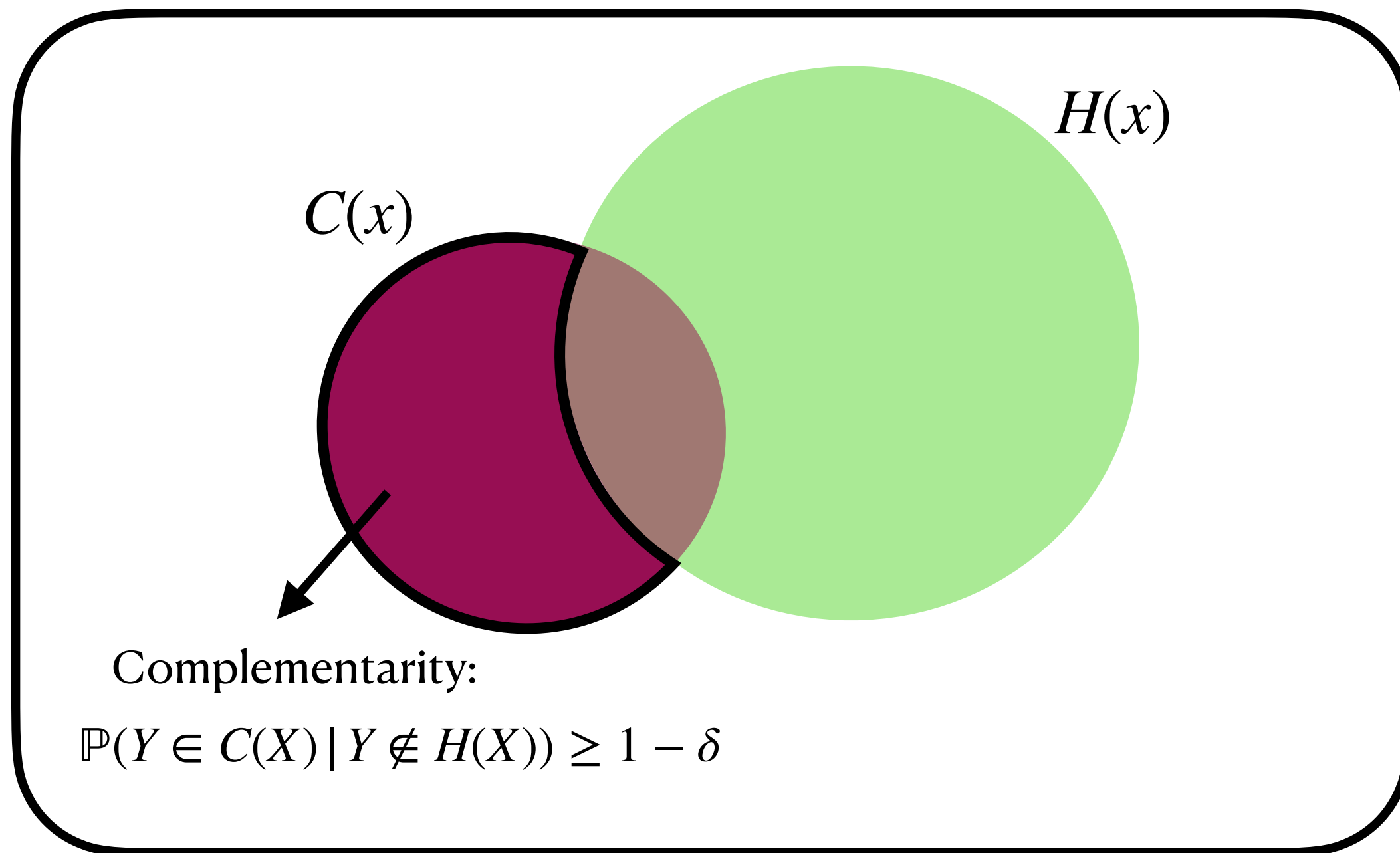
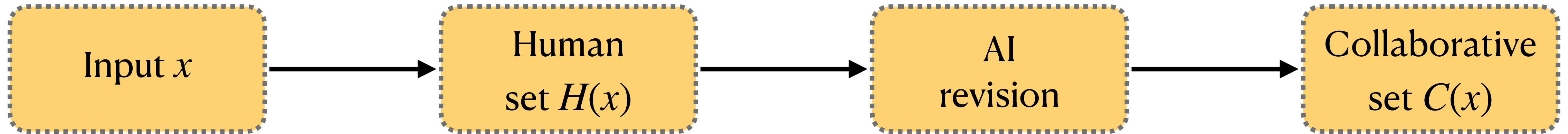
Two fundamentals of collaboration



Complementarity

$$\mathbb{P}(Y \in C(X) \mid Y \notin H(X)) \geq 1 - \delta$$

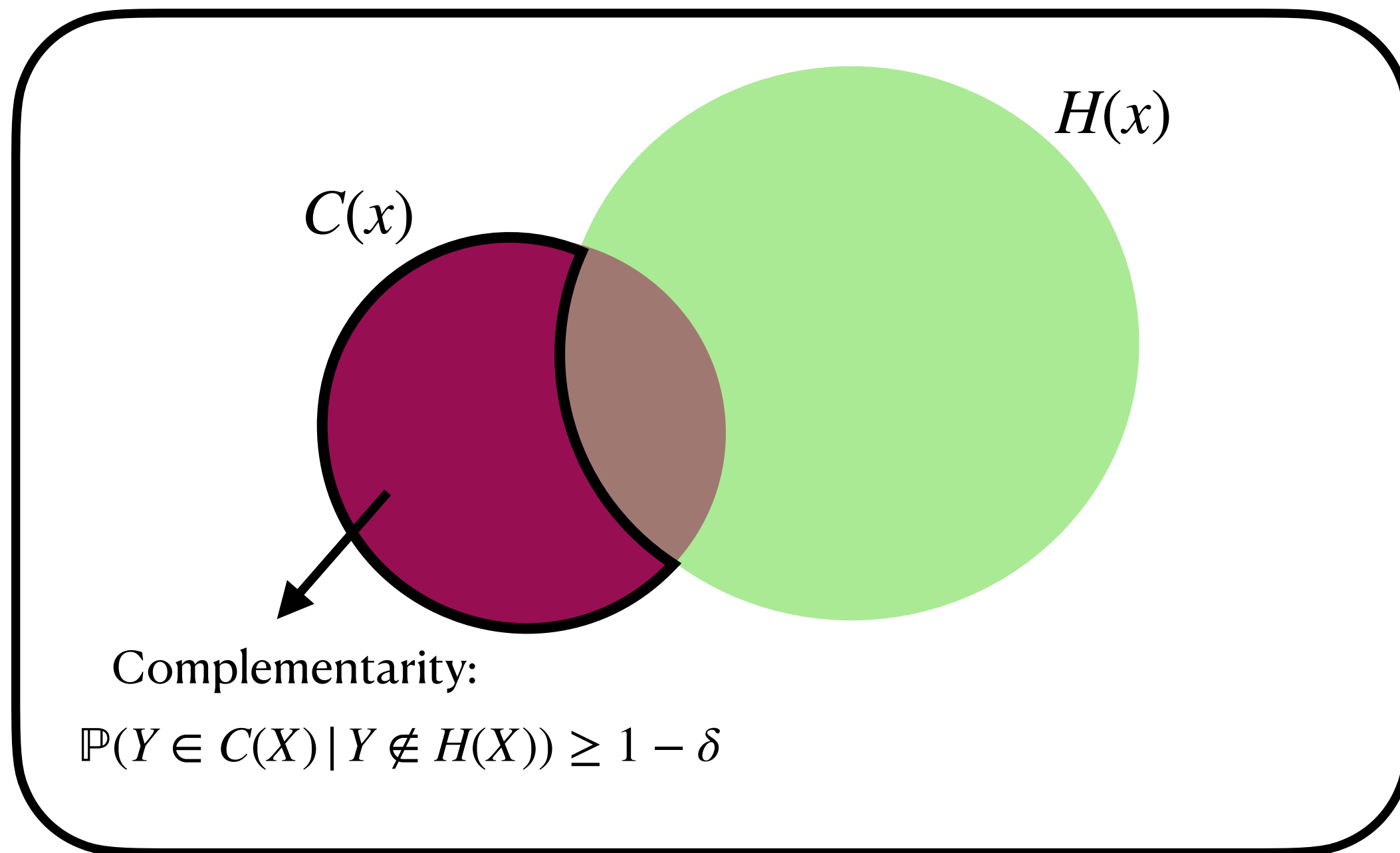
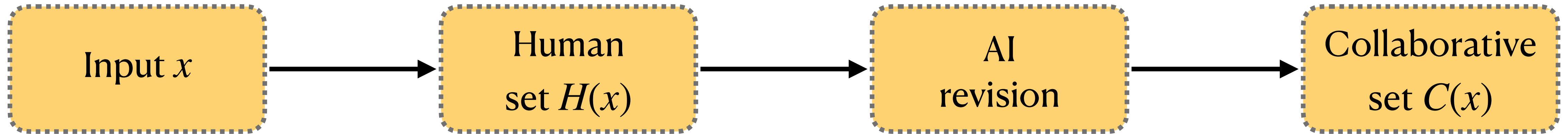
Two fundamentals of collaboration



Complementarity

$$\mathbb{P}(Y \in C(X) | Y \notin H(X)) \geq 1 - \delta$$

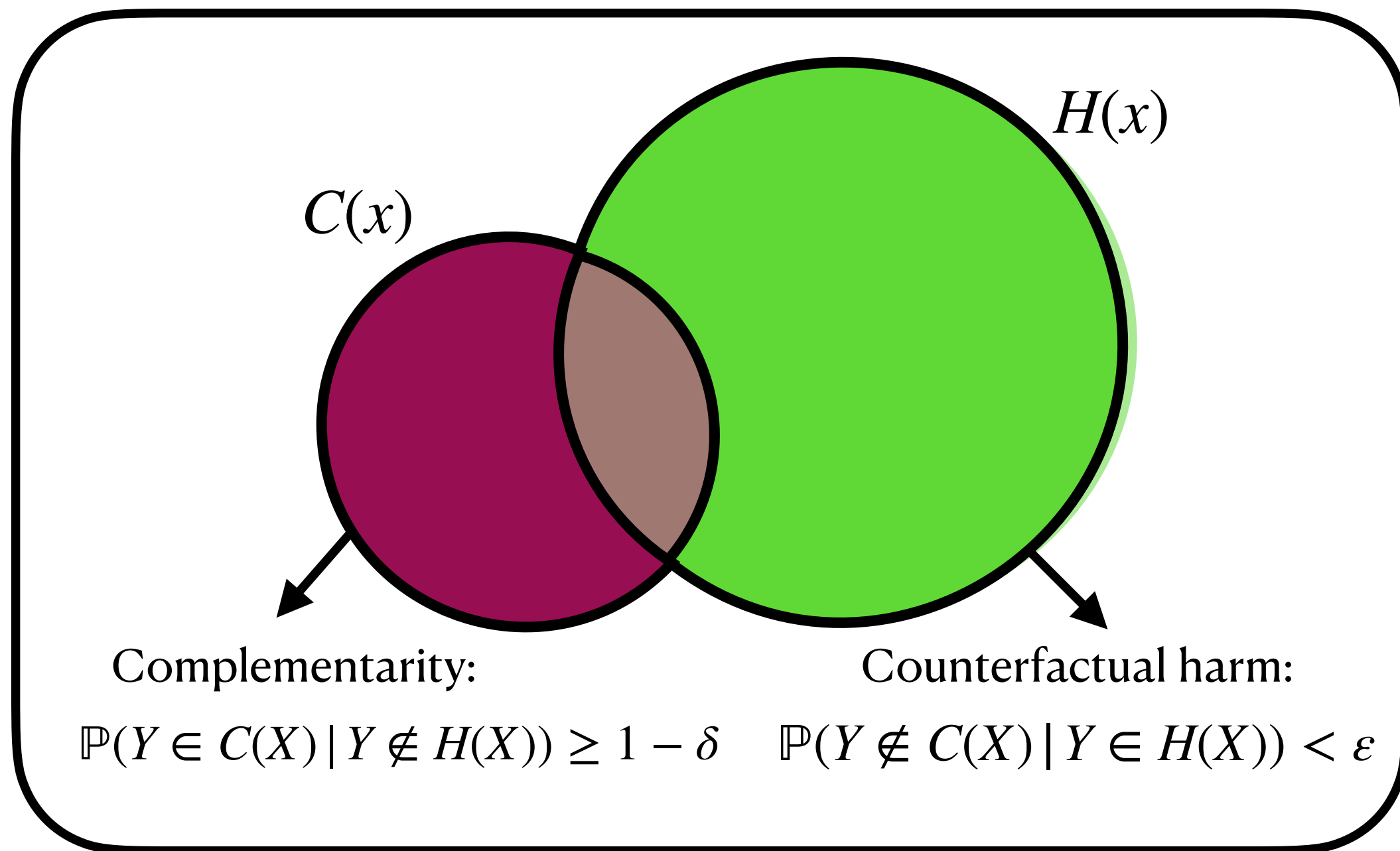
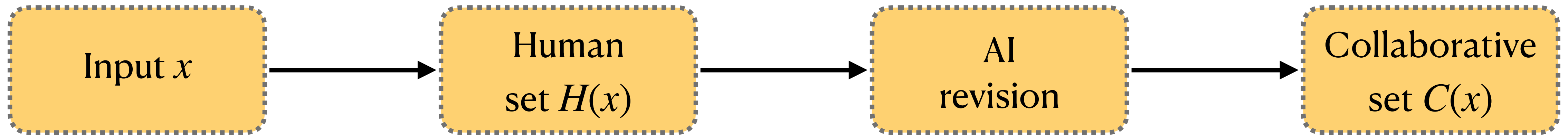
Two fundamentals of collaboration



Complementarity

$$\mathbb{P}(Y \in C(X) | Y \notin H(X)) \geq 1 - \delta$$

Two fundamentals of collaboration



Counterfactual Harm

$$\mathbb{P}(Y \notin C(X) | Y \in H(X)) < \varepsilon$$

Complementarity

$$\mathbb{P}(Y \in C(X) | Y \notin H(X)) \geq 1 - \delta$$

Question: what constitutes a good collaboration?

Question: what constitutes a good collaboration?

$$\begin{aligned} & \min_{C: \mathcal{X} \rightarrow 2^{\mathcal{Y}}} \mathbb{E} |C(X)| \\ & \text{s.t. } \mathbb{P}(Y \notin C(X) \mid Y \in H(X)) < \varepsilon, \\ & \mathbb{P}(Y \in C(X) \mid Y \notin H(X)) \geq 1 - \delta. \end{aligned}$$

Question: what constitutes a good collaboration?

$$\begin{aligned} & \min_{C: \mathcal{X} \rightarrow 2^{\mathcal{Y}}} \mathbb{E} |C(X)| \\ & \text{s.t. } \mathbb{P}(Y \notin C(X) \mid Y \in H(X)) < \varepsilon, \quad \text{Counterfactual Harm} \\ & \mathbb{P}(Y \in C(X) \mid Y \notin H(X)) \geq 1 - \delta. \end{aligned}$$

Question: what constitutes a good collaboration?

$$\begin{aligned} & \min_{C: \mathcal{X} \rightarrow 2^{\mathcal{Y}}} \mathbb{E} |C(X)| \\ & \text{s.t. } \mathbb{P}(Y \notin C(X) \mid Y \in H(X)) < \varepsilon, \quad \text{Counterfactual Harm} \\ & \mathbb{P}(Y \in C(X) \mid Y \notin H(X)) \geq 1 - \delta. \quad \text{Complementarity} \end{aligned}$$

Human-AI Collaborative Optimization

$$\begin{aligned} & \min_{C: \mathcal{X} \rightarrow 2^{\mathcal{Y}}} \mathbb{E} |C(X)| \\ & \text{s.t. } \mathbb{P}(Y \notin C(X) \mid Y \in H(X)) < \varepsilon, \\ & \mathbb{P}(Y \in C(X) \mid Y \notin H(X)) \geq 1 - \delta. \end{aligned}$$

HACO

Human-AI Collaborative Optimization

$$\begin{aligned} \min_{C: \mathcal{X} \rightarrow 2^{\mathcal{Y}}} \quad & \mathbb{E} |C(X)| \\ \text{s.t.} \quad & \mathbb{P}(Y \notin C(X) \mid Y \in H(X)) < \varepsilon, \\ & \mathbb{P}(Y \in C(X) \mid Y \notin H(X)) \geq 1 - \delta. \end{aligned}$$

HACO

Theorem: The optimal solution to HACO is of the form

Human-AI Collaborative Optimization

$$\begin{aligned} \min_{C: \mathcal{X} \rightarrow 2^{\mathcal{Y}}} \quad & \mathbb{E} |C(X)| \\ \text{s.t.} \quad & \mathbb{P}(Y \notin C(X) \mid Y \in H(X)) < \varepsilon, \\ & \mathbb{P}(Y \in C(X) \mid Y \notin H(X)) \geq 1 - \delta. \end{aligned}$$

HACO

Theorem: The optimal solution to HACO is of the form

$$C^*(x) = \left\{ y \mid 1 - p(y \mid x) \leq a^* \mathbf{1}_{\{y \notin H(x)\}} + b^* \mathbf{1}_{\{y \in H(x)\}} \right\}, \quad \forall x \in \mathcal{X}.$$

Human-AI Collaborative Optimization

Theorem: The optimal solution to HACO is of the form

$$C^*(x) = \left\{ y \mid 1 - p(y \mid x) \leq a^* \mathbf{1}_{\{y \notin H(x)\}} + b^* \mathbf{1}_{\{y \in H(x)\}} \right\}, \quad \forall x \in \mathcal{X}.$$

Human-AI Collaborative Optimization

Theorem: The optimal solution to HACO is of the form

$$C^*(x) = \left\{ y \mid 1 - p(y \mid x) \leq a^* \mathbf{1}_{\{y \notin H(x)\}} + b^* \mathbf{1}_{\{y \in H(x)\}} \right\}, \quad \forall x \in \mathcal{X}.$$

VS

Ordinary CP:

$$C(x) = \{ y \mid 1 - p(y \mid x) \leq q^* \}$$

Human-AI Collaborative Optimization

Theorem: The optimal solution to HACO is of the form

$$C^*(x) = \left\{ y \mid 1 - p(y \mid x) \leq a^* \mathbf{1}_{\{y \notin H(x)\}} + b^* \mathbf{1}_{\{y \in H(x)\}} \right\}, \quad \forall x \in \mathcal{X}.$$

Human-AI Collaborative Optimization

Theorem: The optimal solution to HACO is of the form

$$C^*(x) = \left\{ y \mid 1 - p(y \mid x) \leq a^* \mathbf{1}_{\{y \notin H(x)\}} + b^* \mathbf{1}_{\{y \in H(x)\}} \right\}, \quad \forall x \in \mathcal{X}.$$



$s(x, y)$

Human-AI Collaborative Optimization

Theorem: The optimal solution to HACO is of the form

$$C^*(x) = \left\{ y \mid s(x, y) \leq a^* \mathbf{1}_{\{y \notin H(x)\}} + b^* \mathbf{1}_{\{y \in H(x)\}} \right\}, \quad \forall x \in \mathcal{X}.$$

Human-AI Collaborative Optimization

Theorem: The optimal solution to HACO is of the form

$$C^*(x) = \left\{ y \mid s(x, y) \leq a^* \mathbf{1}_{\{y \notin H(x)\}} + b^* \mathbf{1}_{\{y \in H(x)\}} \right\}, \quad \forall x \in \mathcal{X}.$$

1

Question: How to design the score function?

Human-AI Collaborative Optimization

Theorem: The optimal solution to HACO is of the form

$$C^*(x) = \left\{ y \mid s(x, y) \leq a^* \mathbf{1}_{\{y \notin H(x)\}} + b^* \mathbf{1}_{\{y \in H(x)\}} \right\}, \quad \forall x \in \mathcal{X}.$$

1 **Question:** How to design the score function?

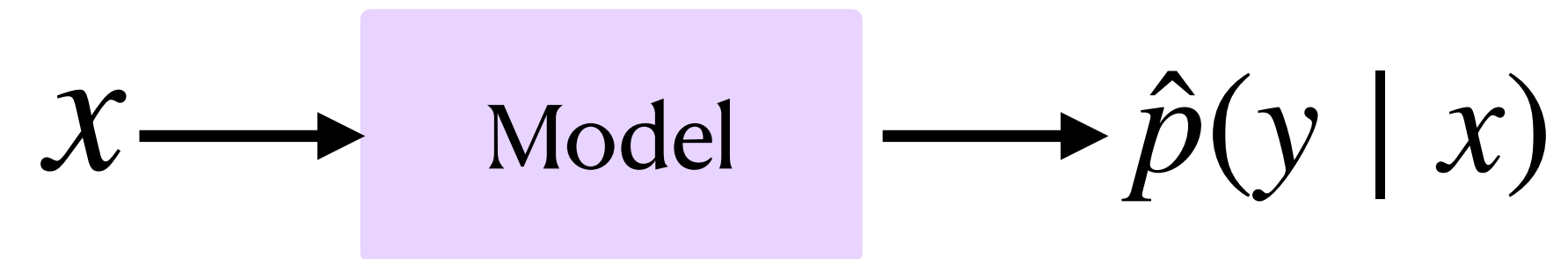
2 **Question:** How to debias the thresholds?

Question: How to design the score function?

Question: How to design the score function?

$$C^*(x) = \left\{ y \mid \underbrace{1 - p(y \mid x)}_{s(x, y)} \leq a^* \mathbf{1}_{\{y \notin H(x)\}} + b^* \mathbf{1}_{\{y \in H(x)\}} \right\}, \quad \forall x \in \mathcal{X}.$$

Classification

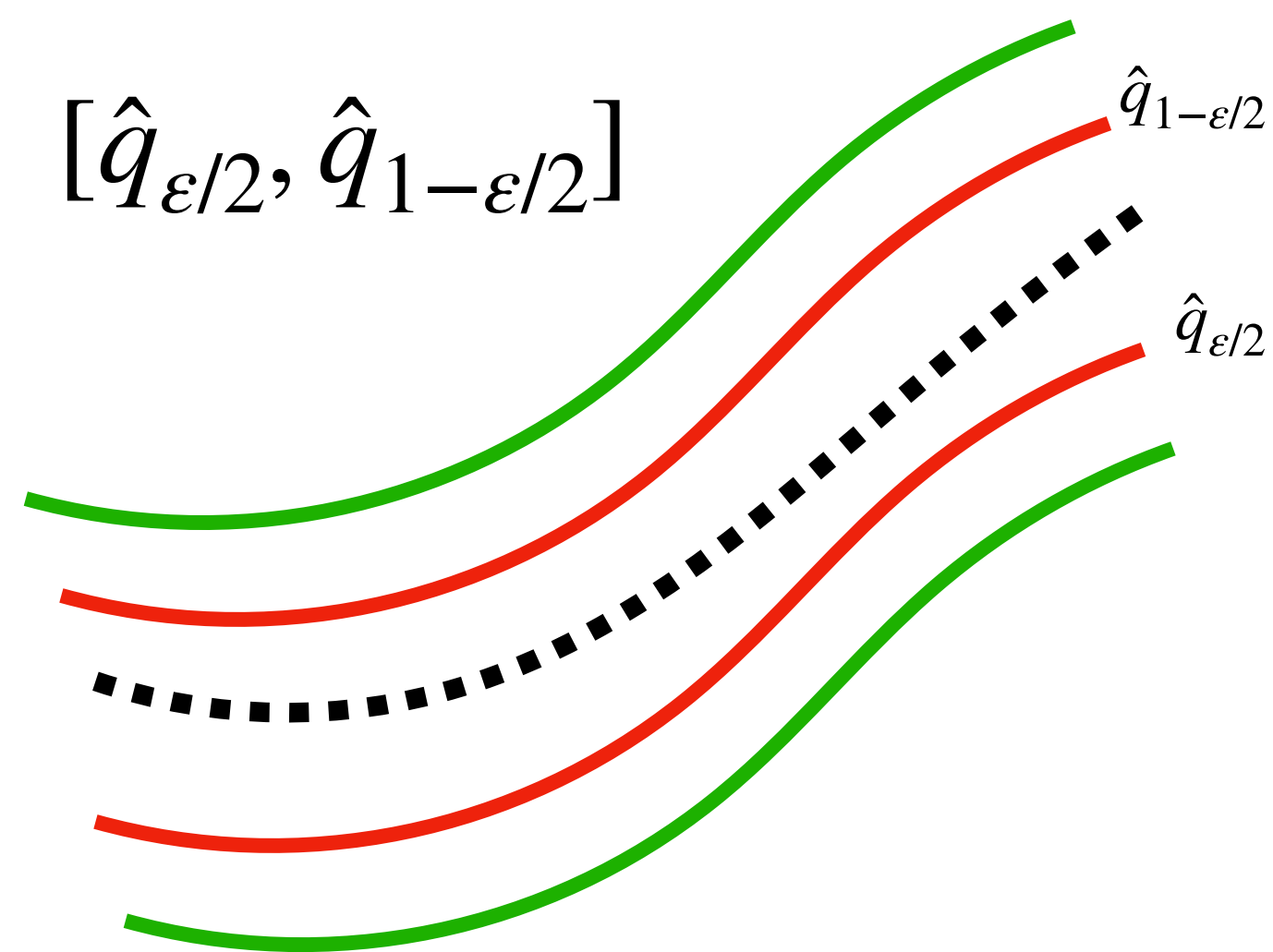


$$\hat{s}(x, y) = 1 - \hat{p}(y \mid x)$$

Question: How to design the score function?

$$C^*(x) = \left\{ y \mid \underbrace{1 - p(y \mid x)}_{s(x, y)} \leq a^* \mathbf{1}_{\{y \notin H(x)\}} + b^* \mathbf{1}_{\{y \in H(x)\}} \right\}, \quad \forall x \in \mathcal{X}.$$

Conformalized Quantile Regression



$$\hat{s}(x, y) = \max \left\{ \hat{q}_{\epsilon/2}(x) - y, y - \hat{q}_{1-\epsilon/2}(x) \right\}$$

Our two threshold structure

$$\hat{s}(x, y) = \begin{cases} \max \left\{ \hat{q}_{\epsilon/2}(x) - y, y - \hat{q}_{1-\epsilon/2}(x) \right\}, & y \in H(x), \\ \max \left\{ \hat{q}_{\delta/2}(x) - y, y - \hat{q}_{1-\delta/2}(x) \right\}, & y \notin H(x). \end{cases}$$

Human-AI Collaborative Optimization

Theorem: The optimal solution to HACO is of the form

$$C^*(x) = \left\{ y \mid \hat{s}(x, y) \leq a^* \mathbf{1}_{\{y \notin H(x)\}} + b^* \mathbf{1}_{\{y \in H(x)\}} \right\}, \quad \forall x \in \mathcal{X}.$$

1 **Question:** How to design the score function?

2 **Question:** How to debias the thresholds?

Question: How to debias the thresholds?

Question: How to debias the thresholds?

Theorem: The optimal solution to HACO is of the form

$$C^*(x) = \left\{ y \mid 1 - p(y \mid x) \leq a^* \mathbf{1}_{\{y \notin H(x)\}} + b^* \mathbf{1}_{\{y \in H(x)\}} \right\}, \quad \forall x \in \mathcal{X}.$$

Assume Exchangeability

Question: How to debias the thresholds?

Theorem: The optimal solution to HACO is of the form

$$C^*(x) = \left\{ y \mid 1 - p(y \mid x) \leq a^* \mathbf{1}_{\{y \notin H(x)\}} + b^* \mathbf{1}_{\{y \in H(x)\}} \right\}, \quad \forall x \in \mathcal{X}.$$

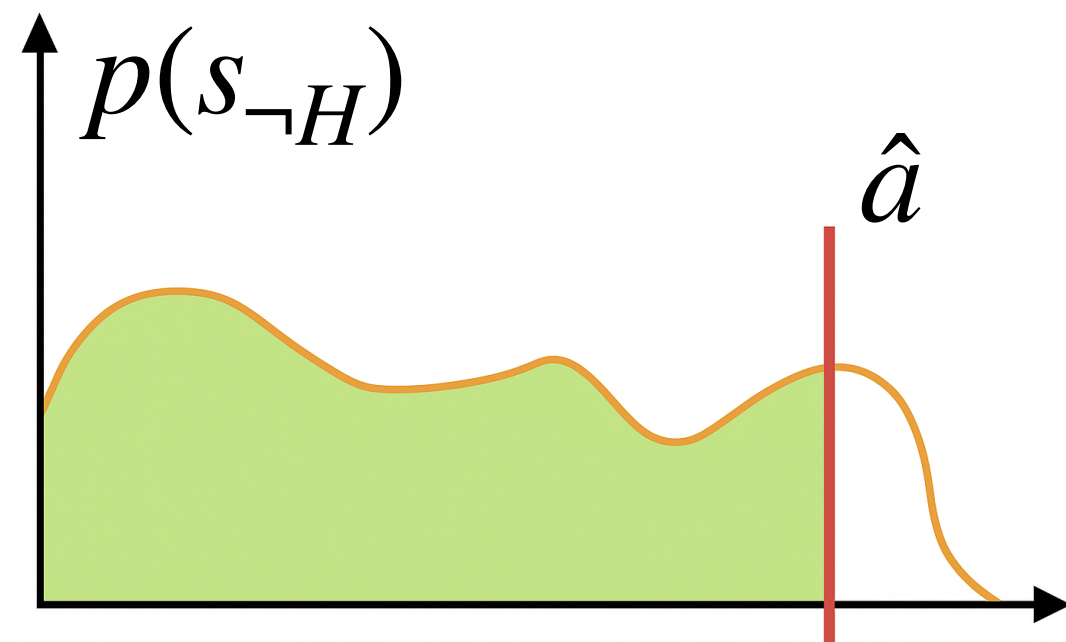
Assume Exchangeability What should a and b be?

Question: How to debias the thresholds?

Theorem: The optimal solution to HACO is of the form

$$C^*(x) = \left\{ y \mid 1 - p(y \mid x) \leq a^* \mathbf{1}_{\{y \notin H(x)\}} + b^* \mathbf{1}_{\{y \in H(x)\}} \right\}, \quad \forall x \in \mathcal{X}.$$

Assume Exchangeability What should a and b be?

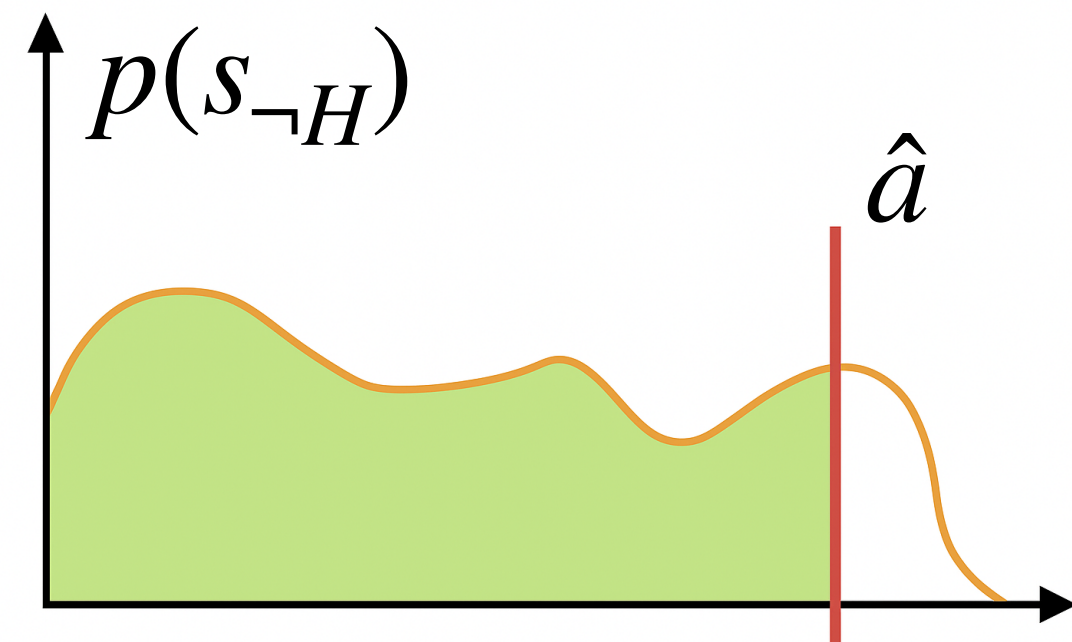


Question: How to debias the thresholds?

Theorem: The optimal solution to HACO is of the form

$$C^*(x) = \left\{ y \mid 1 - p(y \mid x) \leq a^* \mathbf{1}_{\{y \notin H(x)\}} + b^* \mathbf{1}_{\{y \in H(x)\}} \right\}, \quad \forall x \in \mathcal{X}.$$

Assume Exchangeability What should a and b be?



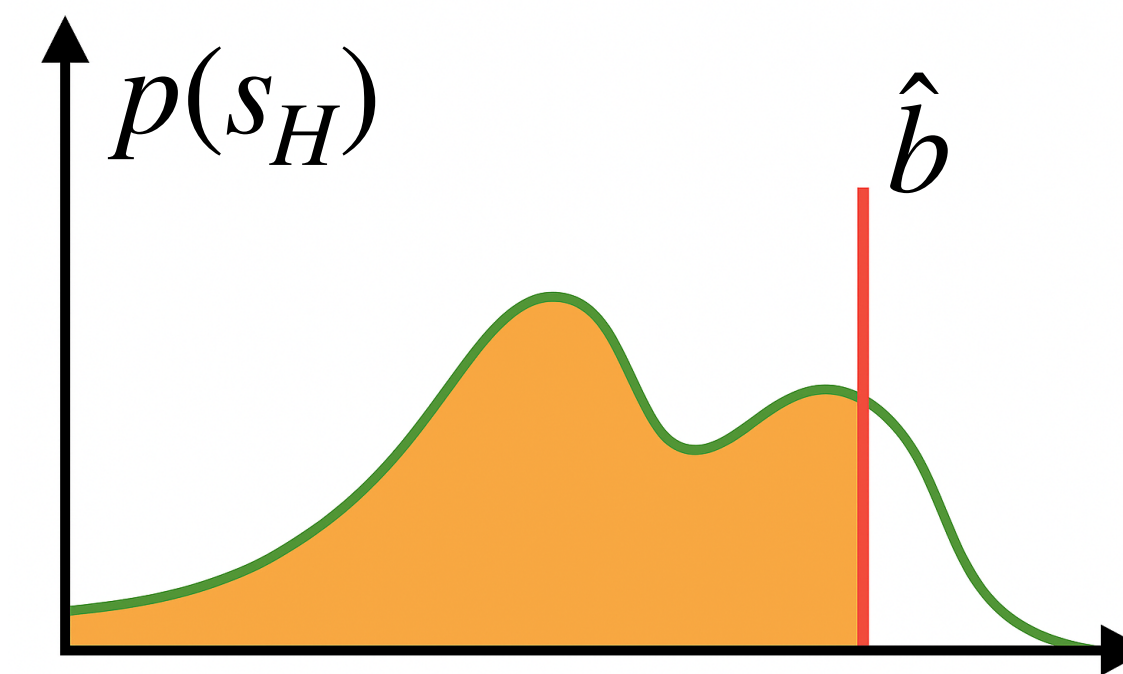
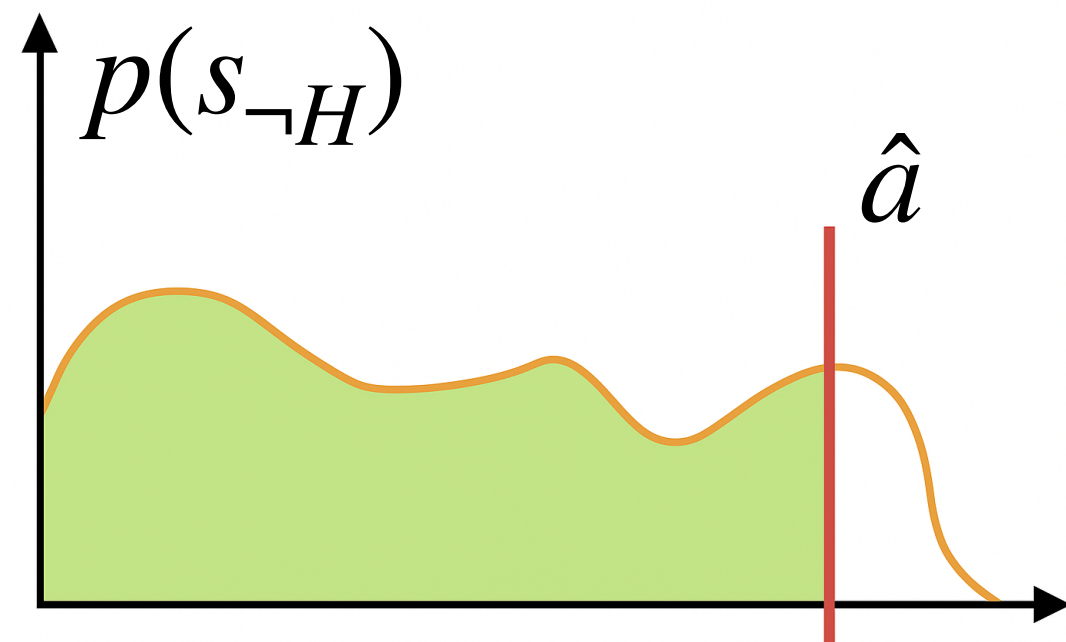
$$\hat{a} = Q_{1-\delta}(\{s_i : Y_i \notin H(X_i)\} \cup \{\infty\})$$

Question: How to debias the thresholds?

Theorem: The optimal solution to HACO is of the form

$$C^*(x) = \left\{ y \mid 1 - p(y \mid x) \leq a^* \mathbf{1}_{\{y \notin H(x)\}} + b^* \mathbf{1}_{\{y \in H(x)\}} \right\}, \quad \forall x \in \mathcal{X}.$$

Assume Exchangeability What should a and b be?



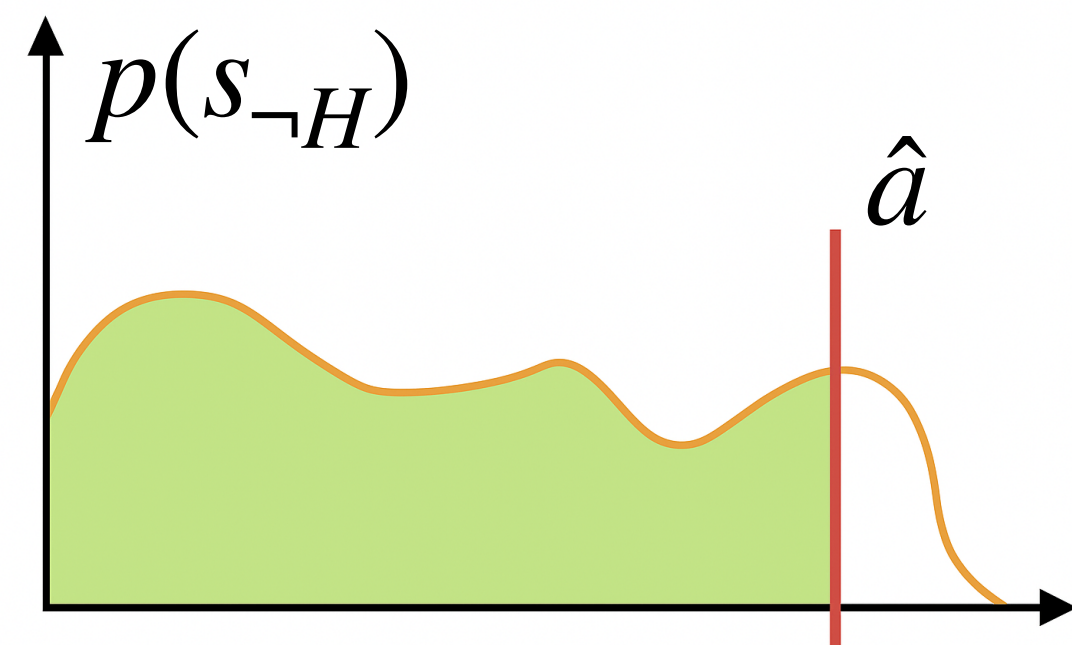
$$\hat{a} = Q_{1-\delta}(\{s_i : Y_i \notin H(X_i)\} \cup \{\infty\})$$

Question: How to debias the thresholds?

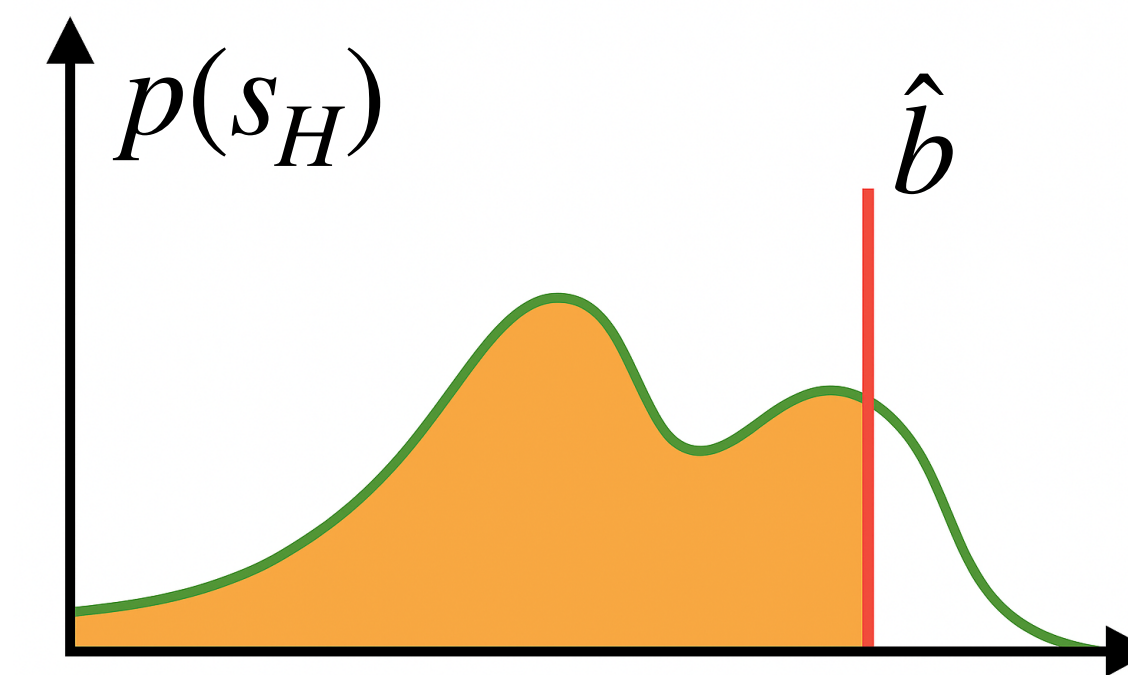
Theorem: The optimal solution to HACO is of the form

$$C^*(x) = \left\{ y \mid 1 - p(y \mid x) \leq a^* \mathbf{1}_{\{y \notin H(x)\}} + b^* \mathbf{1}_{\{y \in H(x)\}} \right\}, \quad \forall x \in \mathcal{X}.$$

Assume Exchangeability What should a and b be?



$$\hat{a} = Q_{1-\delta}(\{s_i : Y_i \notin H(X_i)\} \cup \{\infty\})$$



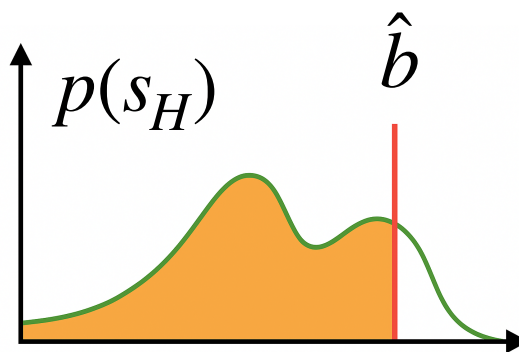
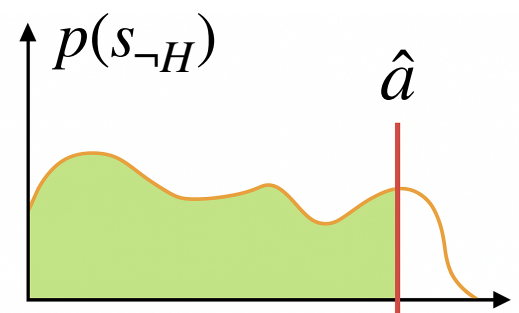
$$\hat{b} = Q_{1-\varepsilon}(\{s_i : Y_i \in H(X_i)\} \cup \{\infty\})$$

Question: How to debias the thresholds?

Theorem: The optimal solution to HACO is of the form

$$C^*(x) = \left\{ y \mid 1 - p(y \mid x) \leq a^* \mathbf{1}_{\{y \notin H(x)\}} + b^* \mathbf{1}_{\{y \in H(x)\}} \right\}, \quad \forall x \in \mathcal{X}.$$

Assume Exchangeability

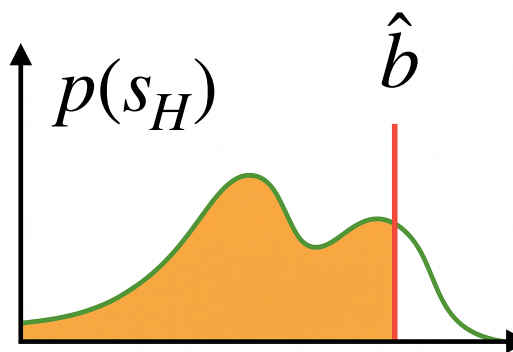
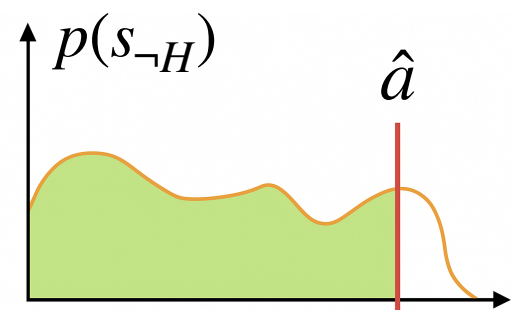


Question: How to debias the thresholds?

Theorem: The optimal solution to HACO is of the form

$$C^*(x) = \left\{ y \mid 1 - p(y \mid x) \leq a^* \mathbf{1}_{\{y \notin H(x)\}} + b^* \mathbf{1}_{\{y \in H(x)\}} \right\}, \quad \forall x \in \mathcal{X}.$$

Assume Exchangeability



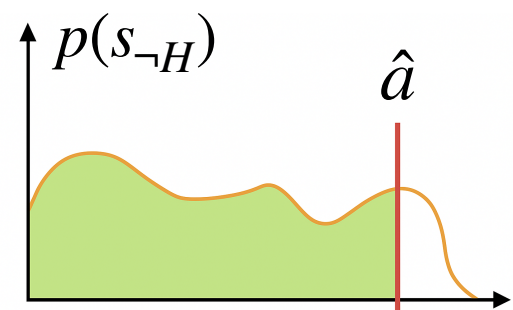
For a **new** test example (x_{test}, y_{test}) :

Question: How to debias the thresholds?

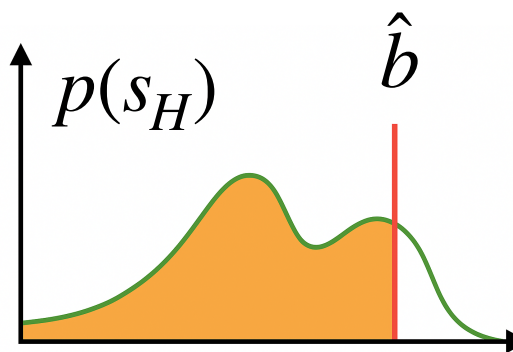
Theorem: The optimal solution to HACO is of the form

$$C^*(x) = \left\{ y \mid 1 - p(y \mid x) \leq a^* \mathbf{1}_{\{y \notin H(x)\}} + b^* \mathbf{1}_{\{y \in H(x)\}} \right\}, \quad \forall x \in \mathcal{X}.$$

Assume Exchangeability



For a **new** test example (x_{test}, y_{test}) :



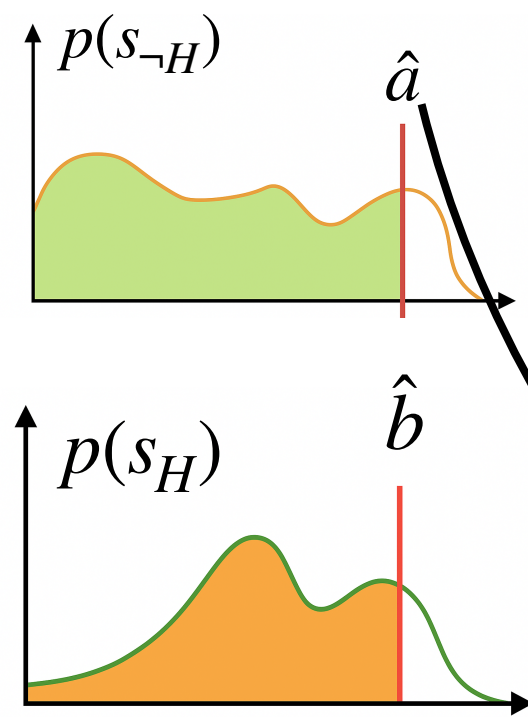
$$C(x_{test}) = \left\{ y \mid s(x_{test}, y) \leq \right\}.$$

Question: How to debias the thresholds?

Theorem: The optimal solution to HACO is of the form

$$C^*(x) = \left\{ y \mid 1 - p(y \mid x) \leq a^* \mathbf{1}_{\{y \notin H(x)\}} + b^* \mathbf{1}_{\{y \in H(x)\}} \right\}, \quad \forall x \in \mathcal{X}.$$

Assume Exchangeability



For a **new** test example (x_{test}, y_{test}) :

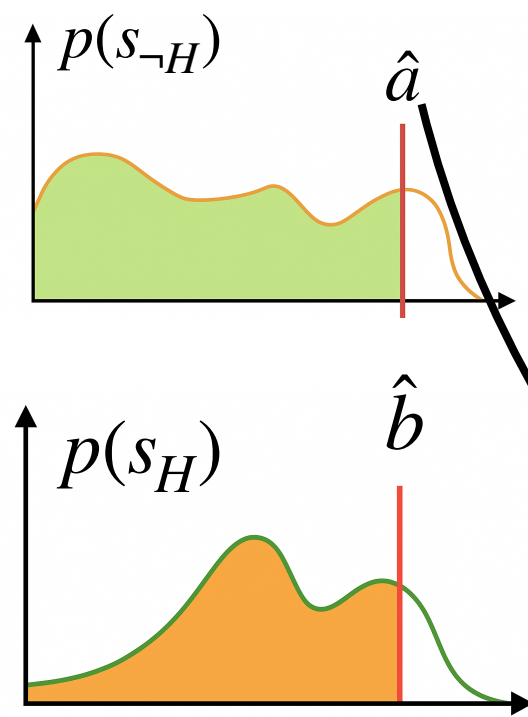
$$C(x_{test}) = \left\{ y \mid s(x_{test}, y) \leq \right\}.$$

Question: How to debias the thresholds?

Theorem: The optimal solution to HACO is of the form

$$C^*(x) = \left\{ y \mid 1 - p(y \mid x) \leq a^* \mathbf{1}_{\{y \notin H(x)\}} + b^* \mathbf{1}_{\{y \in H(x)\}} \right\}, \quad \forall x \in \mathcal{X}.$$

Assume Exchangeability



For a **new** test example (x_{test}, y_{test}) :

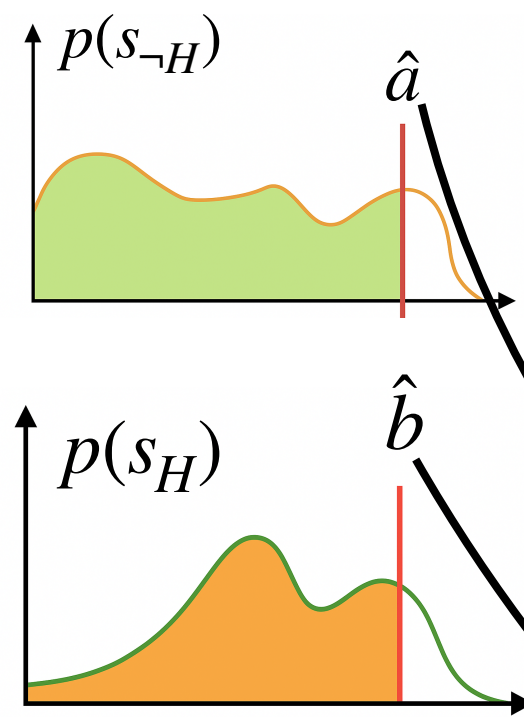
$$C(x_{test}) = \left\{ y \mid s(x_{test}, y) \leq \hat{a} \mathbf{1}_{\{y \notin H(x_{test})\}} \right\}.$$

Question: How to debias the thresholds?

Theorem: The optimal solution to HACO is of the form

$$C^*(x) = \left\{ y \mid 1 - p(y \mid x) \leq a^* \mathbf{1}_{\{y \notin H(x)\}} + b^* \mathbf{1}_{\{y \in H(x)\}} \right\}, \quad \forall x \in \mathcal{X}.$$

Assume Exchangeability



For a **new** test example (x_{test}, y_{test}) :

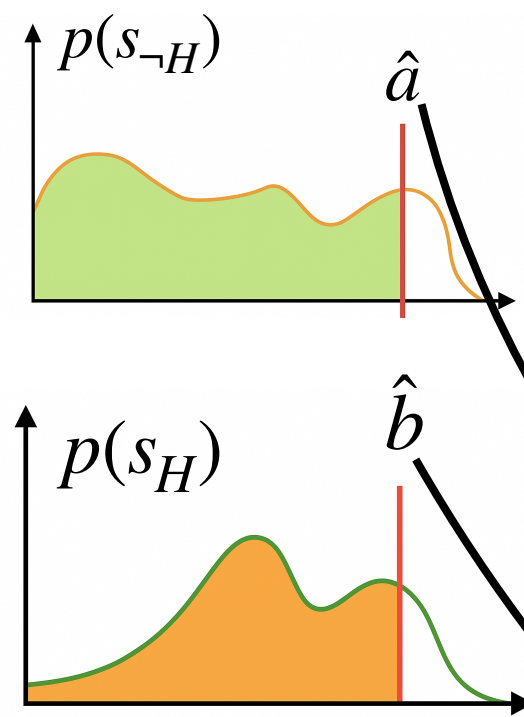
$$C(x_{test}) = \left\{ y \mid s(x_{test}, y) \leq \hat{a} \mathbf{1}_{\{y \notin H(x_{test})\}} + \hat{b} \mathbf{1}_{\{y \in H(x_{test})\}} \right\}.$$

Question: How to debias the thresholds?

Theorem: The optimal solution to HACO is of the form

$$C^*(x) = \left\{ y \mid 1 - p(y \mid x) \leq a^* \mathbf{1}_{\{y \notin H(x)\}} + b^* \mathbf{1}_{\{y \in H(x)\}} \right\}, \quad \forall x \in \mathcal{X}.$$

Assume Exchangeability



For a **new** test example (x_{test}, y_{test}) :

$$C(x_{test}) = \left\{ y \mid s(x_{test}, y) \leq \hat{a} \mathbf{1}_{\{y \notin H(x_{test})\}} + \hat{b} \mathbf{1}_{\{y \in H(x_{test})\}} \right\}.$$

Question: How to debias the thresholds?

Theorem: The optimal solution to HACO is of the form

$$C^*(x) = \left\{ y \mid 1 - p(y \mid x) \leq a^* \mathbf{1}_{\{y \notin H(x)\}} + b^* \mathbf{1}_{\{y \in H(x)\}} \right\}, \quad \forall x \in \mathcal{X}.$$

Assume Exchangeability $\hat{a} = Q_{1-\delta}(\{s_i : Y_i \notin H(X_i)\} \cup \{\infty\})$; $\hat{b} = Q_{1-\varepsilon}(\{s_i : Y_i \in H(X_i)\} \cup \{\infty\})$

$$C(x_{\text{test}}) = \left\{ y \mid s(x_{\text{test}}, y) \leq \hat{a} \mathbf{1}_{\{y \notin H(x_{\text{test}})\}} + \hat{b} \mathbf{1}_{\{y \in H(x_{\text{test}})\}} \right\}.$$

$$1 - \varepsilon \leq \Pr[Y_{\text{test}} \in C(X_{\text{test}}) \mid Y_{\text{test}} \in H(X_{\text{test}})] < 1 - \varepsilon + \frac{1}{n_1 + 1}$$

$$1 - \delta \leq \Pr[Y_{\text{test}} \in C(X_{\text{test}}) \mid Y_{\text{test}} \notin H(X_{\text{test}})] < 1 - \delta + \frac{1}{n_2 + 1}$$

Question: How to debias the thresholds?

Theorem: The optimal solution to HACO is of the form

$$C^*(x) = \left\{ y \mid 1 - p(y \mid x) \leq a^* \mathbf{1}_{\{y \notin H(x)\}} + b^* \mathbf{1}_{\{y \in H(x)\}} \right\}, \quad \forall x \in \mathcal{X}.$$

Assume Exchangeability $\hat{a} = Q_{1-\delta}(\{s_i : Y_i \notin H(X_i)\} \cup \{\infty\})$; $\hat{b} = Q_{1-\varepsilon}(\{s_i : Y_i \in H(X_i)\} \cup \{\infty\})$

$$C(x_{\text{test}}) = \left\{ y \mid s(x_{\text{test}}, y) \leq \hat{a} \mathbf{1}_{\{y \notin H(x_{\text{test}})\}} + \hat{b} \mathbf{1}_{\{y \in H(x_{\text{test}})\}} \right\}.$$

Offline Guarantees

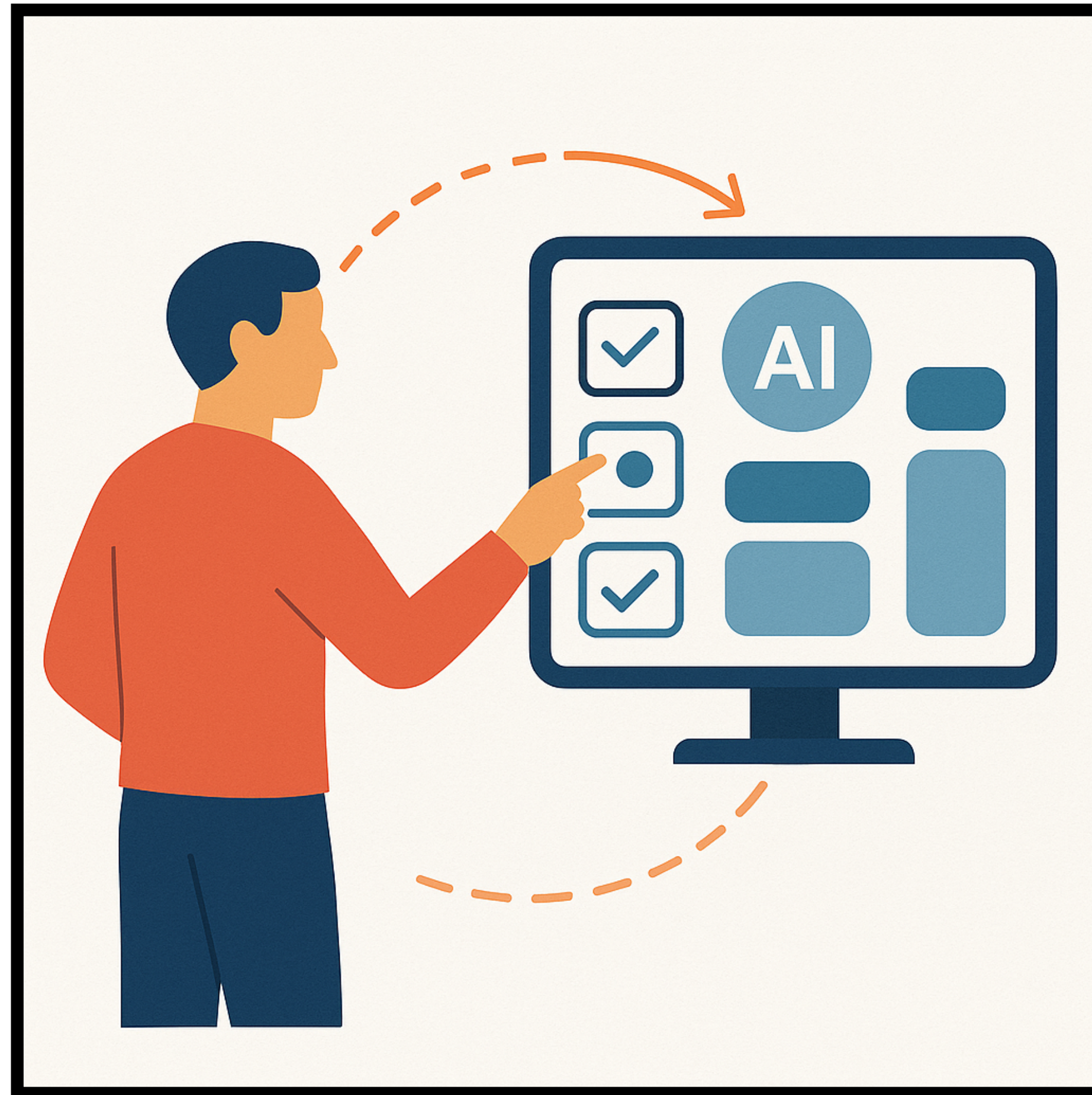
$$1 - \varepsilon \leq \Pr[Y_{\text{test}} \in C(X_{\text{test}}) \mid Y_{\text{test}} \in H(X_{\text{test}})] < 1 - \varepsilon + \frac{1}{n_1 + 1}$$

$$1 - \delta \leq \Pr[Y_{\text{test}} \in C(X_{\text{test}}) \mid Y_{\text{test}} \notin H(X_{\text{test}})] < 1 - \delta + \frac{1}{n_2 + 1}$$

~~Assume Exchangeability~~

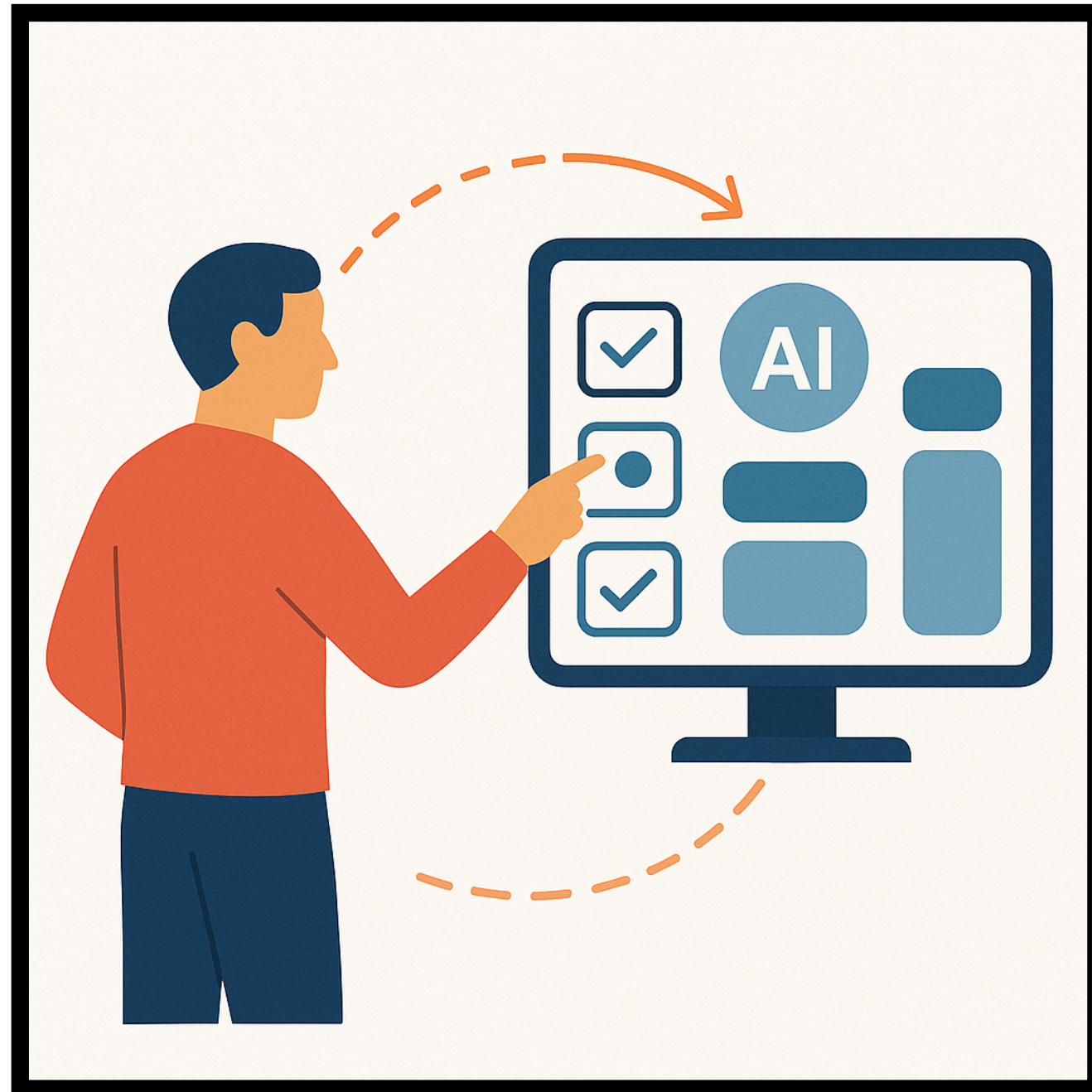
~~Assume Exchangeability~~

Exchangeability is fragile



~~Assume Exchangeability~~

Exchangeability is fragile



Question: How to debias the thresholds in the **online** setting

Question: How to debias the thresholds in the **online** setting

Question: How to debias the thresholds in the **online** setting

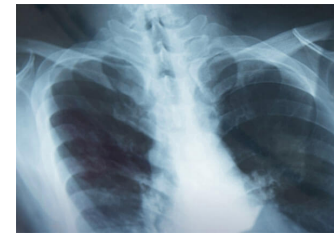
At round t



Question: How to debias the thresholds in the **online** setting

At round t

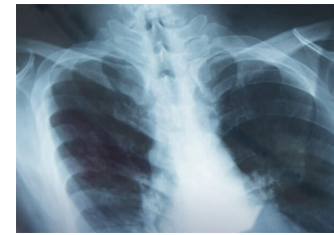
input x_t



Question: How to debias the thresholds in the **online** setting

At round t

input x_t



$H(x_t)$



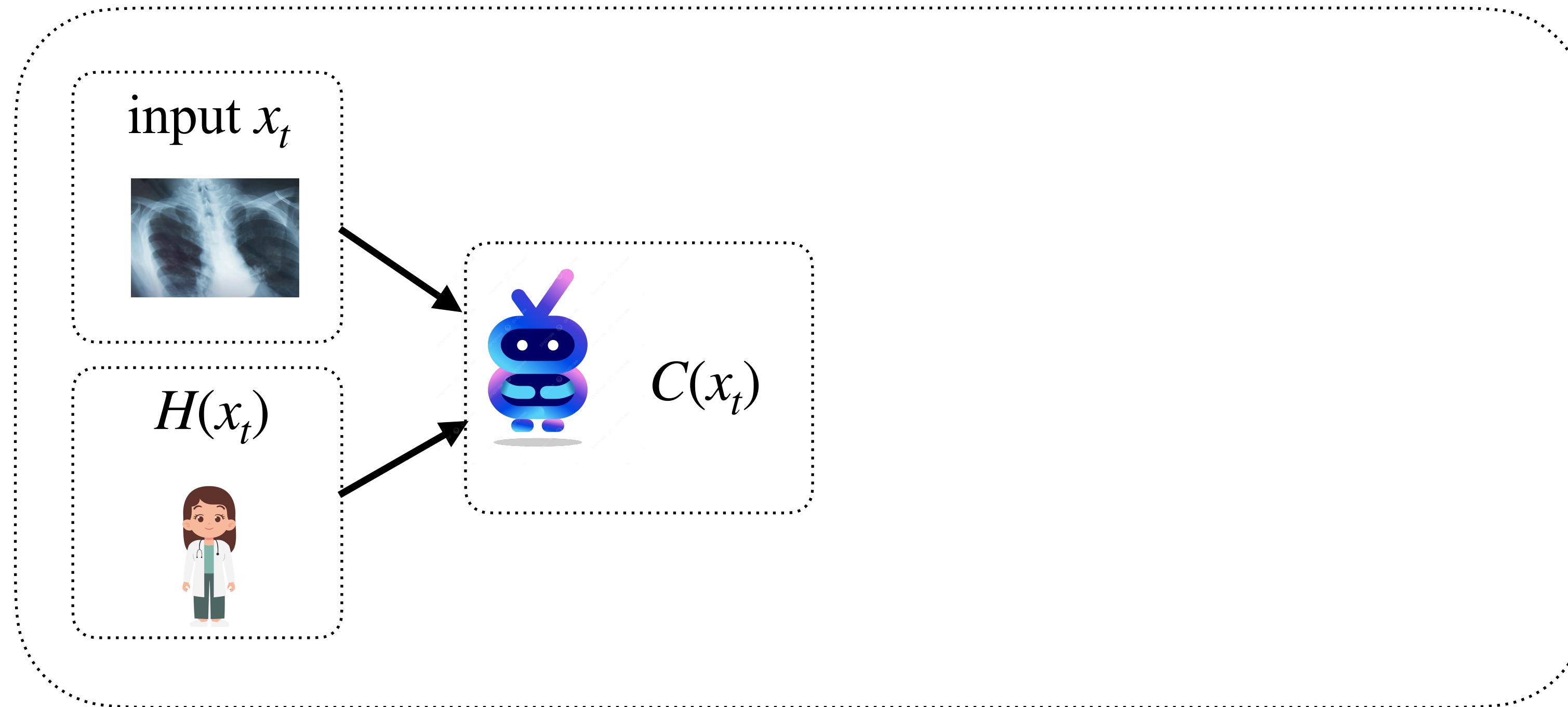
Question: How to debias the thresholds in the **online** setting

At round t



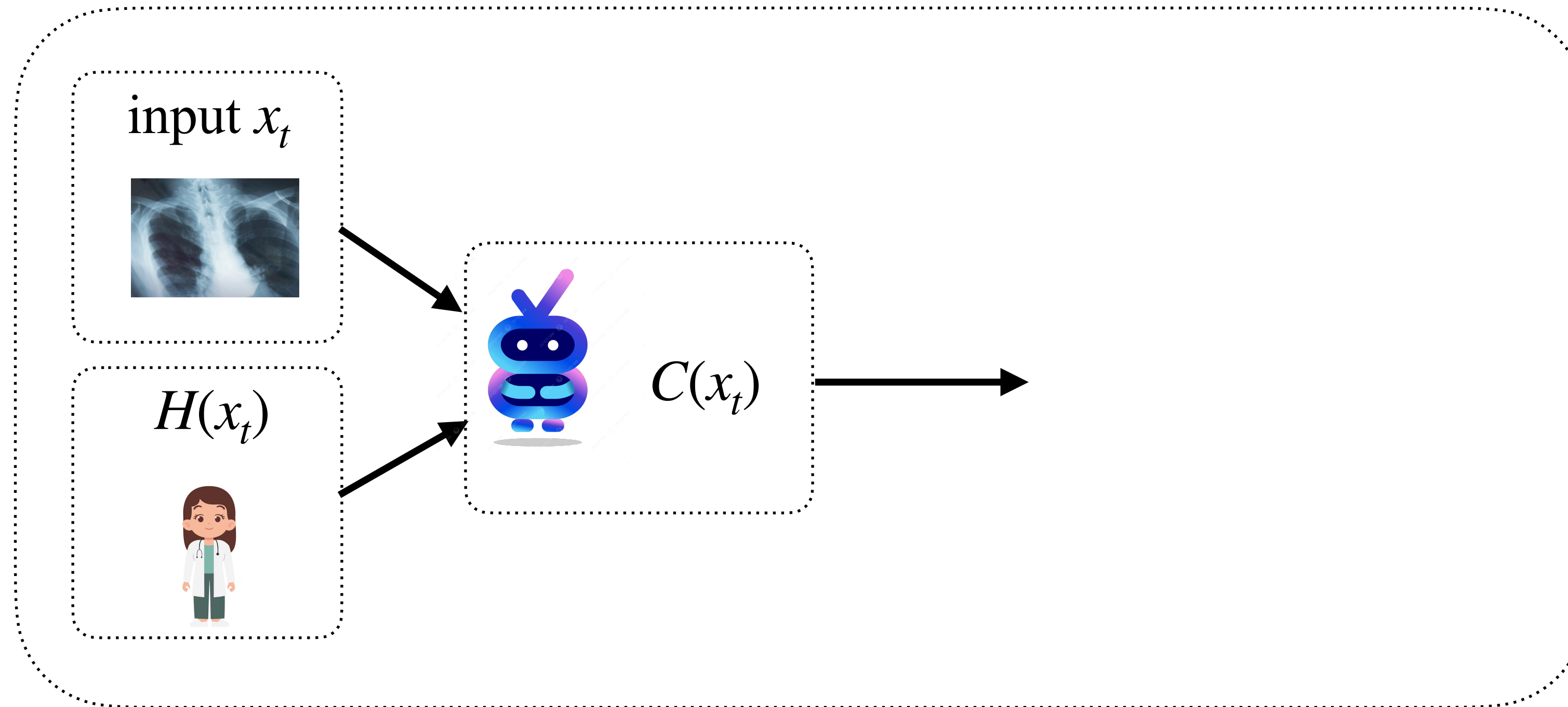
Question: How to debias the thresholds in the **online** setting

At round t



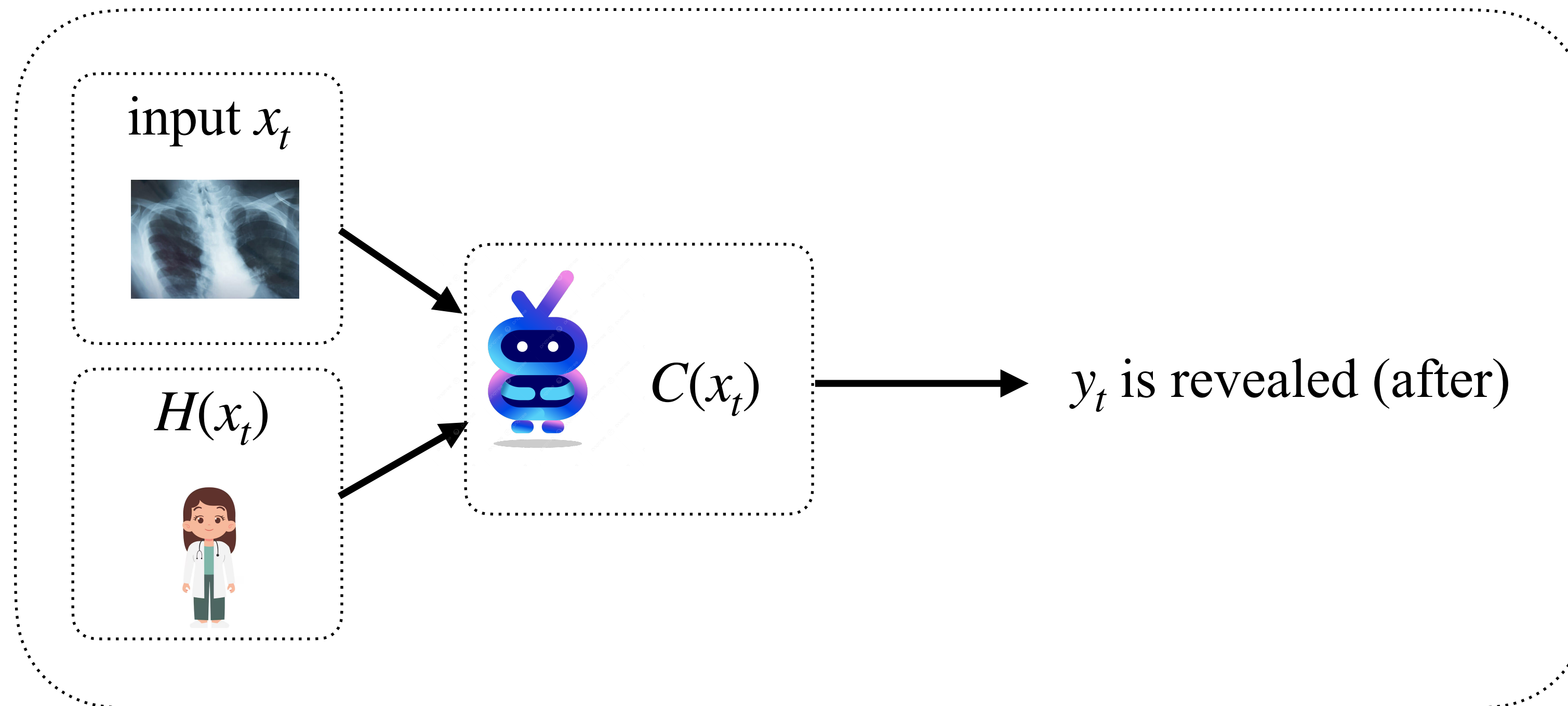
Question: How to debias the thresholds in the **online** setting

At round t



Question: How to debias the thresholds in the **online** setting

At round t



Question: How to debias the thresholds in the **online** setting

At round t input x_t ; $H(x_t)$; y_t is revealed (after)

Question: How to debias the thresholds in the **online** setting

At round t input x_t ; $H(x_t)$; y_t is revealed (after)

$$C_t(x_t) = \left\{ y \mid s(x_t, y) \leq a_t \mathbf{1}_{\{y \notin H(x_t)\}} + b_t \mathbf{1}_{\{y \in H(x_t)\}} \right\}$$

Question: How to debias the thresholds in the **online** setting

At round t input x_t ; $H(x_t)$; y_t is revealed (after)

$$C_t(x_t) = \left\{ y \mid s(x_t, y) \leq a_t \mathbf{1}_{\{y \notin H(x_t)\}} + b_t \mathbf{1}_{\{y \in H(x_t)\}} \right\}$$

We track two errors

Question: How to debias the thresholds in the **online** setting

At round t input x_t ; $H(x_t)$; y_t is revealed (after)

$$C_t(x_t) = \left\{ y \mid s(x_t, y) \leq a_t \mathbf{1}_{\{y \notin H(x_t)\}} + b_t \mathbf{1}_{\{y \in H(x_t)\}} \right\}$$

We track two errors

$$\text{err}_t^{\text{in}} := \mathbf{1}_{\{y_t \notin C_t(x_t), y_t \in H(x_t)\}}$$

Question: How to debias the thresholds in the **online** setting

At round t input x_t ; $H(x_t)$; y_t is revealed (after)

$$C_t(x_t) = \left\{ y \mid s(x_t, y) \leq a_t \mathbf{1}_{\{y \notin H(x_t)\}} + b_t \mathbf{1}_{\{y \in H(x_t)\}} \right\}$$

We track two errors

$$\text{err}_t^{\text{in}} := \mathbf{1}_{\{y_t \notin C_t(x_t), y_t \in H(x_t)\}}$$

Counterfactual Harm Error

Question: How to debias the thresholds in the **online** setting

At round t input x_t ; $H(x_t)$; y_t is revealed (after)

$$C_t(x_t) = \left\{ y \mid s(x_t, y) \leq a_t \mathbf{1}_{\{y \notin H(x_t)\}} + b_t \mathbf{1}_{\{y \in H(x_t)\}} \right\}$$

We track two errors

$$\text{err}_t^{\text{in}} := \mathbf{1}_{\{y_t \notin C_t(x_t), y_t \in H(x_t)\}}$$

Counterfactual Harm Error

$$\text{err}_t^{\text{out}} := \mathbf{1}_{\{y_t \notin C_t(x_t), y_t \notin H(x_t)\}}$$

Complementarity Error

Question: How to debias the thresholds in the **online** setting

We track two errors

$$\text{err}_t^{\text{in}} := \mathbf{1}_{\{y_t \notin C_t(x_t), y_t \in H(x_t)\}}$$

Counterfactual Harm Error

$$\text{err}_t^{\text{out}} := \mathbf{1}_{\{y_t \notin C_t(x_t), y_t \notin H(x_t)\}}$$

Complementarity Error

We update the thresholds as follows

$$\text{If } (y_t \in H(x_t)) \longrightarrow$$

Question: How to debias the thresholds in the **online** setting

We track two errors

$$\text{err}_t^{\text{in}} := \mathbf{1}_{\{y_t \notin C_t(x_t), y_t \in H(x_t)\}}$$

Counterfactual Harm Error

$$\text{err}_t^{\text{out}} := \mathbf{1}_{\{y_t \notin C_t(x_t), y_t \notin H(x_t)\}}$$

Complementarity Error

We update the thresholds as follows

$$\text{If } (y_t \in H(x_t)) \longrightarrow b_{t+1} = b_t + \eta (\mathbf{1}_{\{s(x_t, y_t) > b_t\}} - \varepsilon),$$

Question: How to debias the thresholds in the **online** setting

We track two errors

$$\text{err}_t^{\text{in}} := \mathbf{1}_{\{y_t \notin C_t(x_t), y_t \in H(x_t)\}}$$

Counterfactual Harm Error

$$\text{err}_t^{\text{out}} := \mathbf{1}_{\{y_t \notin C_t(x_t), y_t \notin H(x_t)\}}$$

Complementarity Error

We update the thresholds as follows

$$\text{If } (y_t \in H(x_t)) \longrightarrow b_{t+1} = b_t + \eta (\mathbf{1}_{\{s(x_t, y_t) > b_t\}} - \varepsilon),$$

$$a_{t+1} = a_t$$

Question: How to debias the thresholds in the **online** setting

We track two errors

$$\text{err}_t^{\text{in}} := \mathbf{1}_{\{y_t \notin C_t(x_t), y_t \in H(x_t)\}}$$

Counterfactual Harm Error

$$\text{err}_t^{\text{out}} := \mathbf{1}_{\{y_t \notin C_t(x_t), y_t \notin H(x_t)\}}$$

Complementarity Error

We update the thresholds as follows

If $(y_t \in H(x_t))$

$$b_{t+1} = b_t + \eta(\mathbf{1}_{\{s(x_t, y_t) > b_t\}} - \varepsilon),$$

$$a_{t+1} = a_t$$

If $(y_t \notin H(x_t))$

Question: How to debias the thresholds in the **online** setting

We track two errors

$$\text{err}_t^{\text{in}} := \mathbf{1}_{\{y_t \notin C_t(x_t), y_t \in H(x_t)\}}$$

Counterfactual Harm Error

$$\text{err}_t^{\text{out}} := \mathbf{1}_{\{y_t \notin C_t(x_t), y_t \notin H(x_t)\}}$$

Complementarity Error

We update the thresholds as follows

If $(y_t \in H(x_t))$

$$b_{t+1} = b_t + \eta(\mathbf{1}_{\{s(x_t, y_t) > b_t\}} - \varepsilon),$$

$$a_{t+1} = a_t$$

If $(y_t \notin H(x_t))$

$$a_{t+1} = a_t + \eta(\mathbf{1}_{\{s(x_t, y_t) > a_t\}} - \delta),$$

Question: How to debias the thresholds in the **online** setting

We track two errors

$$\text{err}_t^{\text{in}} := \mathbf{1}_{\{y_t \notin C_t(x_t), y_t \in H(x_t)\}}$$

Counterfactual Harm Error

$$\text{err}_t^{\text{out}} := \mathbf{1}_{\{y_t \notin C_t(x_t), y_t \notin H(x_t)\}}$$

Complementarity Error

We update the thresholds as follows

If $(y_t \in H(x_t))$

$$b_{t+1} = b_t + \eta(\mathbf{1}_{\{s(x_t, y_t) > b_t\}} - \varepsilon),$$

$$a_{t+1} = a_t$$

If $(y_t \notin H(x_t))$

$$a_{t+1} = a_t + \eta(\mathbf{1}_{\{s(x_t, y_t) > a_t\}} - \delta),$$

$$b_{t+1} = b_t$$

Question: How to debias the thresholds in the **online** setting

Counterfactual Harm Error

$$\text{err}_t^{\text{in}} := \mathbf{1}_{\{y_t \notin C_t(x_t), y_t \in H(x_t)\}}$$

If $(y_t \in H(x_t))$

$$b_{t+1} = b_t + \eta(\mathbf{1}_{\{s(x_t, y_t) > b_t\}} - \varepsilon),$$

$$a_{t+1} = a_t$$

Complementarity Error

$$\text{err}_t^{\text{out}} := \mathbf{1}_{\{y_t \notin C_t(x_t), y_t \notin H(x_t)\}}$$

If $(y_t \notin H(x_t))$

$$a_{t+1} = a_t + \eta(\mathbf{1}_{\{s(x_t, y_t) > a_t\}} - \delta),$$

$$b_{t+1} = b_t$$

Theorem: Assume the conformity score is bounded, i.e $s(x, y) \in [0, 1]$:

Question: How to debias the thresholds in the **online** setting

Counterfactual Harm Error

$$\text{err}_t^{\text{in}} := \mathbf{1}_{\{y_t \notin C_t(x_t), y_t \in H(x_t)\}}$$

If $(y_t \in H(x_t))$

$$b_{t+1} = b_t + \eta(\mathbf{1}_{\{s(x_t, y_t) > b_t\}} - \varepsilon),$$

$$a_{t+1} = a_t$$

Complementarity Error

$$\text{err}_t^{\text{out}} := \mathbf{1}_{\{y_t \notin C_t(x_t), y_t \notin H(x_t)\}}$$

If $(y_t \notin H(x_t))$

$$a_{t+1} = a_t + \eta(\mathbf{1}_{\{s(x_t, y_t) > a_t\}} - \delta),$$

$$b_{t+1} = b_t$$

Theorem: Assume the conformity score is bounded, i.e $s(x, y) \in [0, 1]$:

$$\left| \frac{1}{N_1(T)} \sum_{t=1}^T \text{err}_t^{\text{in}} - \varepsilon \right| \leq \frac{1 + \eta \max(\varepsilon, 1 - \varepsilon)}{\eta N_1(T)}$$

Question: How to debias the thresholds in the **online** setting

Counterfactual Harm Error

$$\text{err}_t^{\text{in}} := \mathbf{1}_{\{y_t \notin C_t(x_t), y_t \in H(x_t)\}}$$

If $(y_t \in H(x_t))$

$$b_{t+1} = b_t + \eta(\mathbf{1}_{\{s(x_t, y_t) > b_t\}} - \varepsilon),$$

$$a_{t+1} = a_t$$

Complementarity Error

$$\text{err}_t^{\text{out}} := \mathbf{1}_{\{y_t \notin C_t(x_t), y_t \notin H(x_t)\}}$$

If $(y_t \notin H(x_t))$

$$a_{t+1} = a_t + \eta(\mathbf{1}_{\{s(x_t, y_t) > a_t\}} - \delta),$$

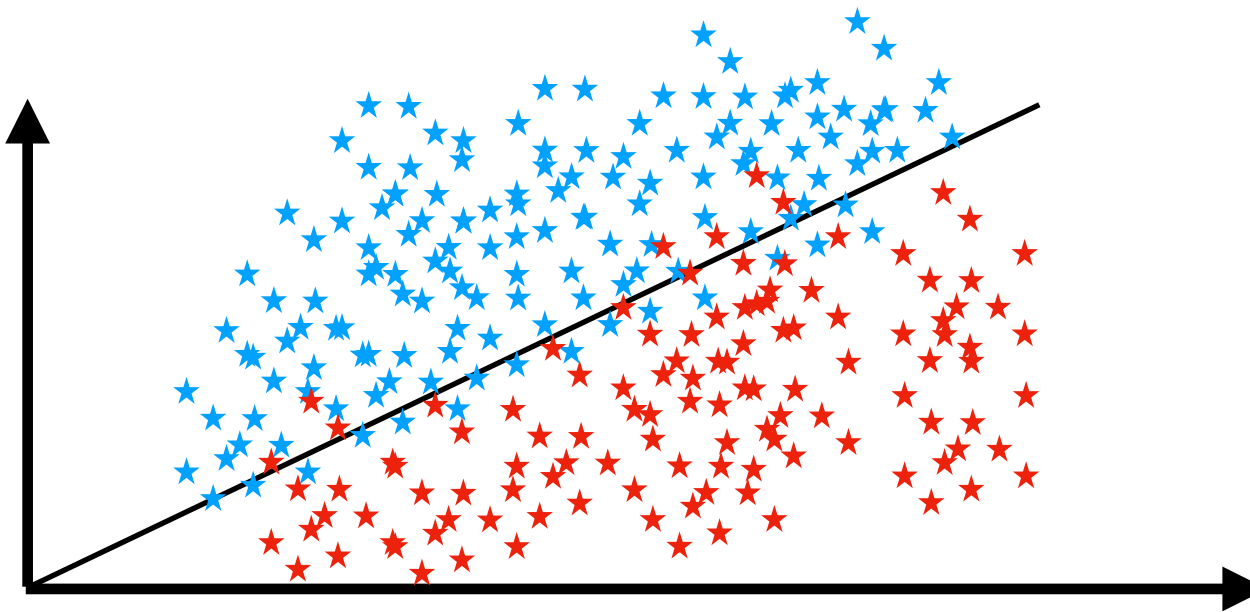
$$b_{t+1} = b_t$$

Theorem: Assume the conformity score is bounded, i.e $s(x, y) \in [0, 1]$:

$$\left| \frac{1}{N_1(T)} \sum_{t=1}^T \text{err}_t^{\text{in}} - \varepsilon \right| \leq \frac{1 + \eta \max(\varepsilon, 1 - \varepsilon)}{\eta N_1(T)}$$

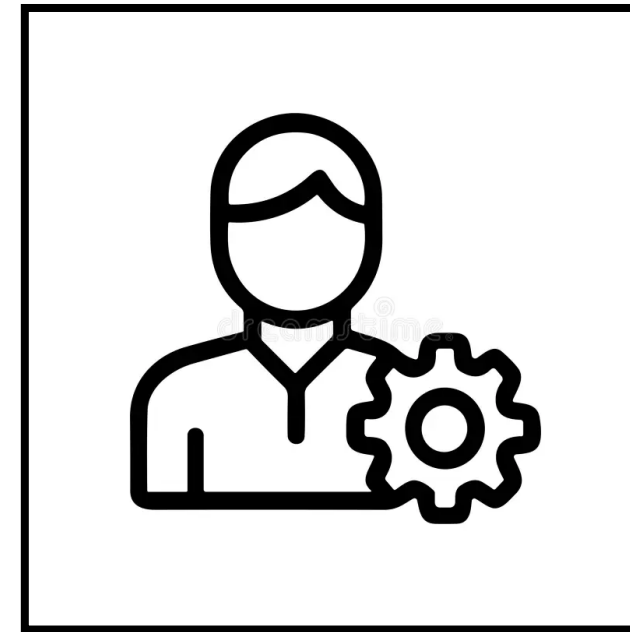
$$\left| \frac{1}{N_2(T)} \sum_{t=1}^T \text{err}_t^{\text{out}} - \delta \right| \leq \frac{1 + \eta \max(\delta, 1 - \delta)}{\eta N_2(T)}$$

We consider three modalities of data

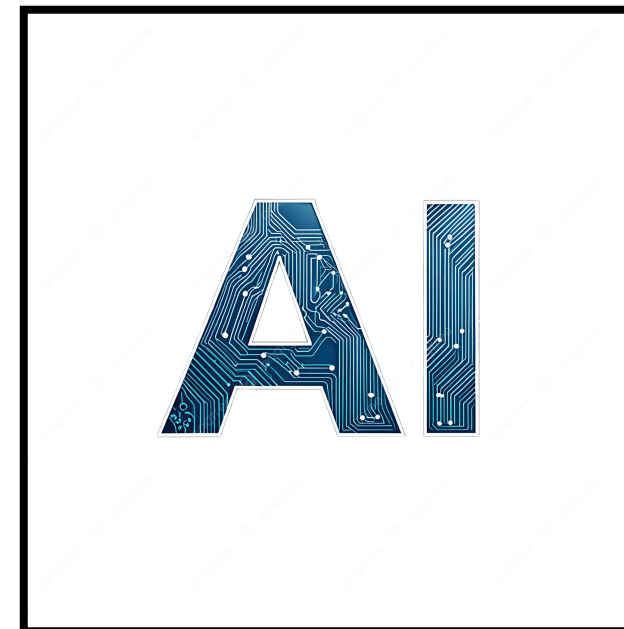


Experiments

Baselines

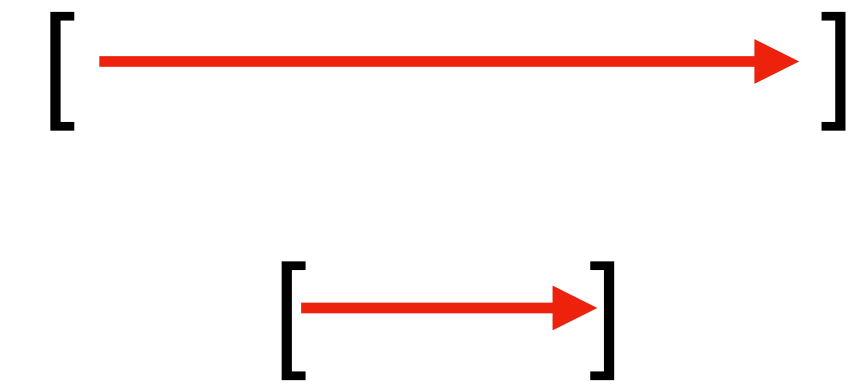


Human Alone

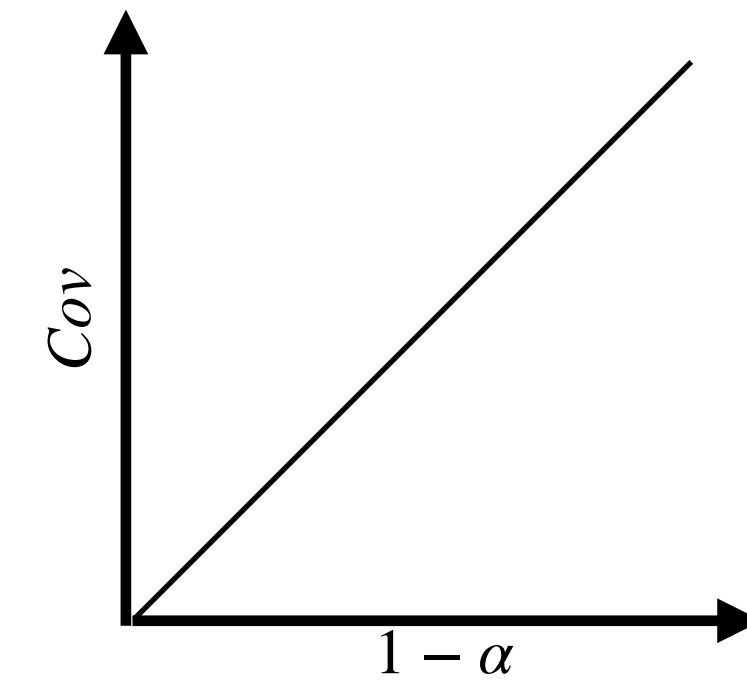


AI Alone

Evaluation Metrics



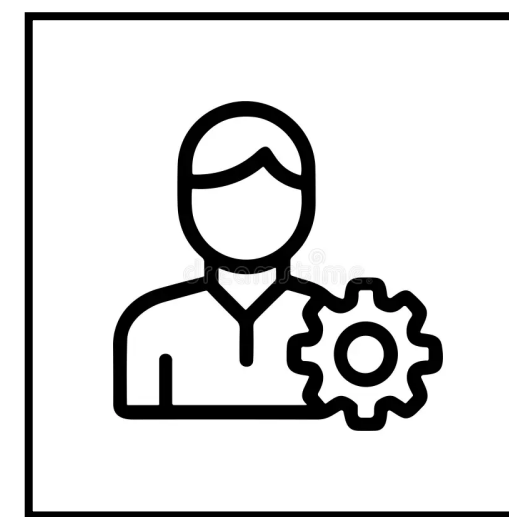
Set Size



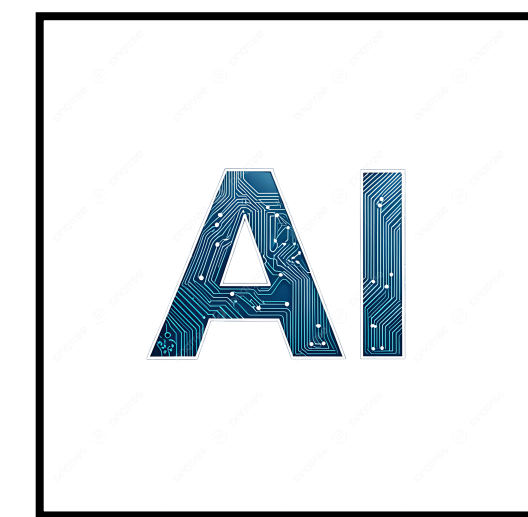
Marginal Coverage



Text Based Medical Diagnosis



Rule-based diagnostic system

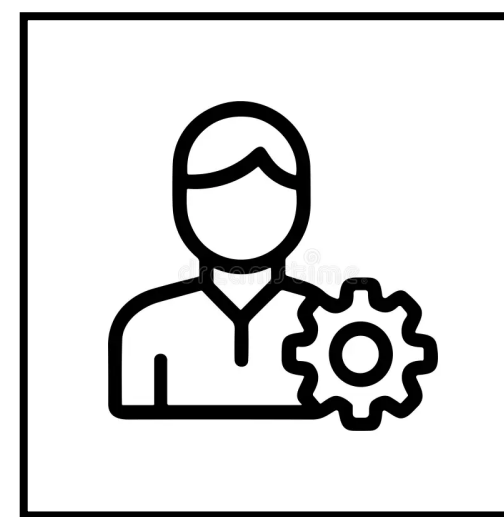


LLMs

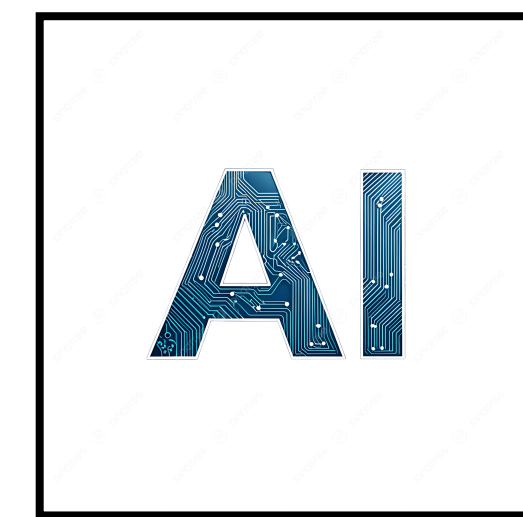
Offline Setting



Text Based Medical Diagnosis



Rule-based diagnostic system



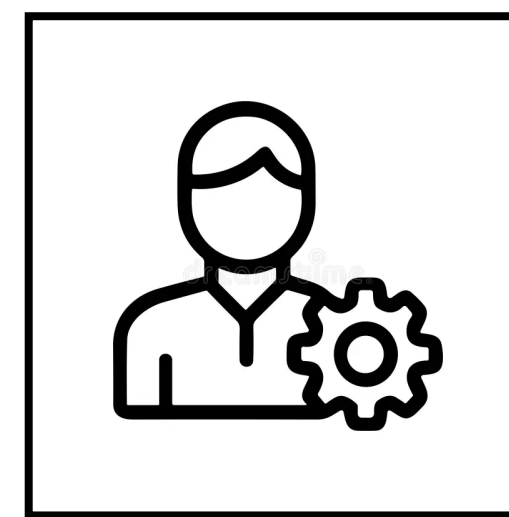
LLMs

Offline Setting

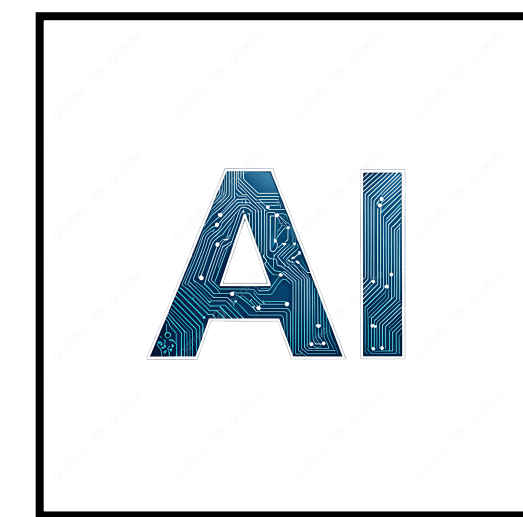
Strategy	Human	GPT-4o			GPT-5		
	C/S	CUP C/S	(ϵ, δ)	AI C/S	CUP C/S	(ϵ, δ)	AI C/S
Top-1	0.71 / 1.00	0.90 / 2.84	(0.02, 0.70)	0.88 / 4.64	0.91 / 1.59	(0.02, 0.70)	0.91 / 1.76
Top-2	0.87 / 1.95	0.93 / 3.14	(0.01, 0.45)	0.90 / 9.12	0.93 / 1.65	(0.02, 0.45)	0.93 / 1.95



Text Based Medical Diagnosis



Rule-based diagnostic system



LLMs

Offline Setting

Strategy	Human	GPT-4o			GPT-5		
	C/S	CUP C/S	(ϵ, δ)	AI C/S	CUP C/S	(ϵ, δ)	AI C/S
Top-1	0.71 / 1.00	0.90 / 2.84	(0.02, 0.70)	0.88 / 4.64	0.91 / 1.59	(0.02, 0.70)	0.91 / 1.76
Top-2	0.87 / 1.95	0.93 / 3.14	(0.01, 0.45)	0.90 / 9.12	0.93 / 1.65	(0.02, 0.45)	0.93 / 1.95

AI's quality affects the overall collaboration quality!

Thank You!