

# Human-AI Collaborative Uncertainty Quantification

Sima Noorani  
University of Pennsylvania

November 17th, 2025

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



**AI has become powerful**

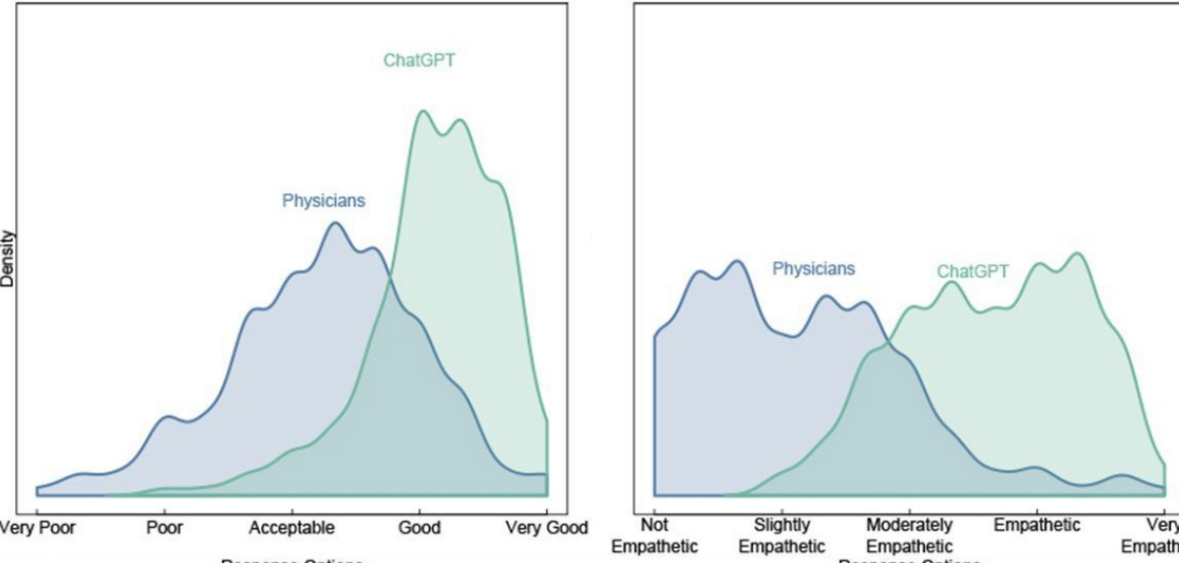
# AI has become powerful

UC San Diego

Menu UC SAN DIEGO TODAY

## 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 figure consists of two density plots. The left plot shows the distribution of response options, with the x-axis ranging from 'Very Poor' to 'Very Good' and the y-axis representing 'Density'. The 'Physicians' distribution (blue) is centered around 'Acceptable', while the 'ChatGPT' distribution (green) is centered around 'Good'. The right plot shows the distribution of empathy levels, with the x-axis ranging from 'Not Empathetic' to 'Very Empathetic' and the y-axis representing 'Density'. The 'Physicians' distribution (blue) is centered around 'Moderately Empathetic', while the 'ChatGPT' distribution (green) is centered around 'Empathetic'.

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

Empathy Level	Physicians Density	ChatGPT Density
Not Empathetic	High	Low
Slightly Empathetic	High	Low
Moderately Empathetic	High	Low
Empathetic	Low	High
Very Empathetic	Very Low	High

# AI has become powerful

Menu UC San Diego TODAY

## 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

Density

Response Options

Very Poor Poor Acceptable Good Very Good

Not Empathetic Slightly Empathetic Moderately Empathetic Empathetic Very Empathetic

EXPLORE TIME SUBSCRIBE

JUL 2, 2025 10:11 AM ET

## Microsoft's AI Is Better Than Doctors at Diagnosing Disease

HEALTH AI

Kilito Chan—Getty Images

by [Alice Park](#)  
SENIOR CORRESPONDENT

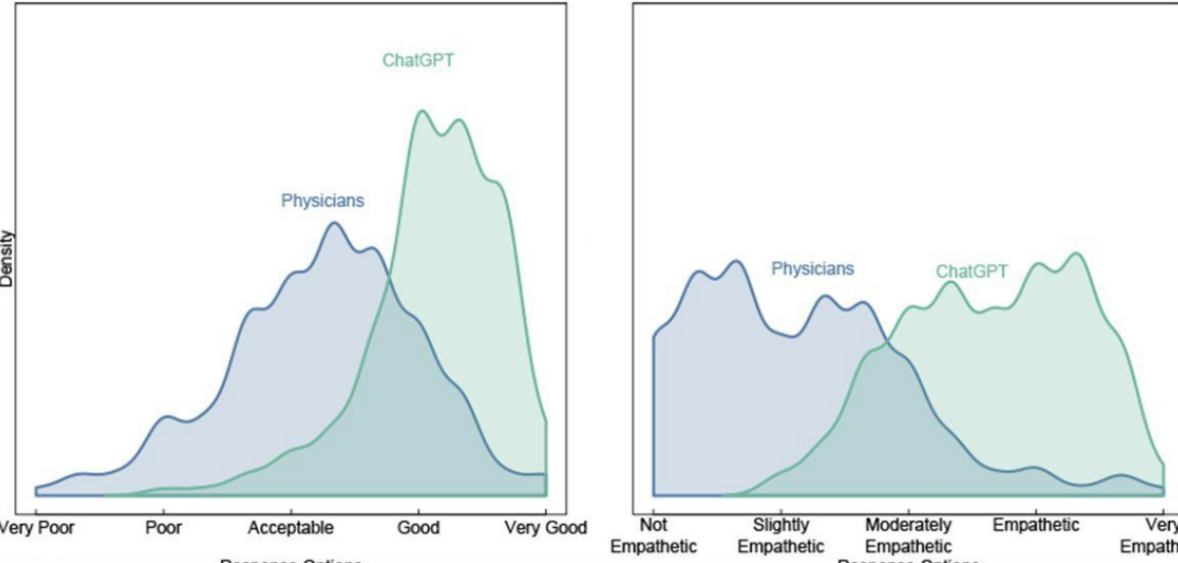
**M**edicine may be a combination of art and science, but Microsoft just showed that much of both can be learned

# AI has become powerful

UC San Diego TODAY

## 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



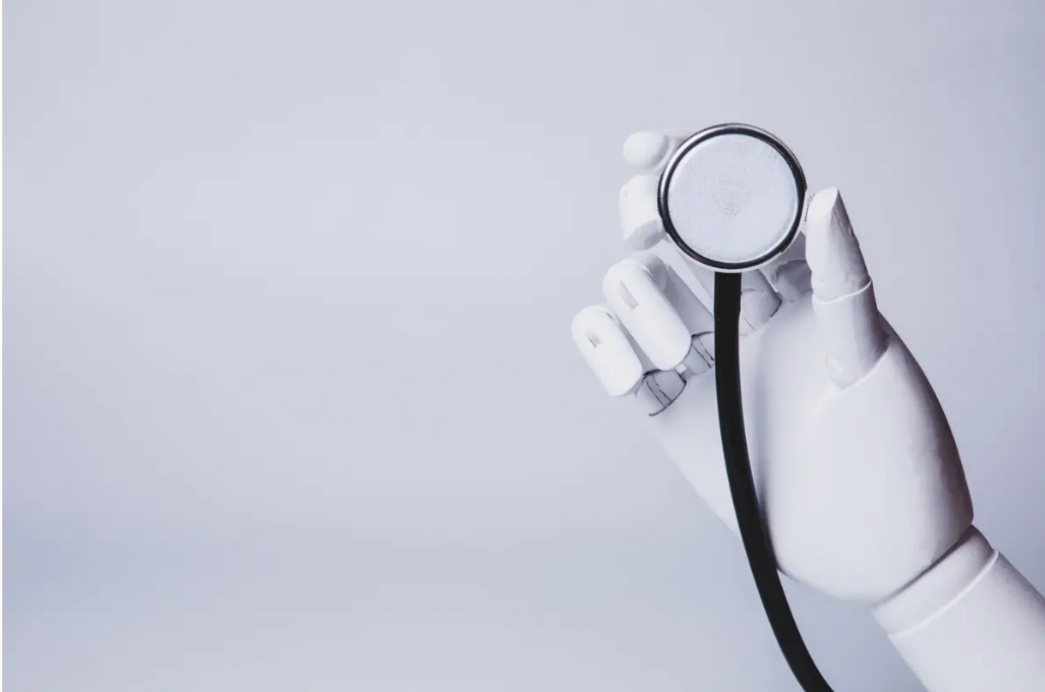
The figure consists of two density plots. The left plot shows the distribution of response quality options, with ChatGPT (green) peaking at 'Good' and Physicians (blue) peaking at 'Acceptable'. The right plot shows the distribution of empathy levels, with ChatGPT (green) peaking at 'Moderately Empathetic' and Physicians (blue) peaking at 'Slightly Empathetic'.

EXPLORE TIME SUBSCRIBE

JUL 2, 2025 10:11 AM ET

## Microsoft's AI Is Better Than Doctors at Diagnosing Disease

HEALTH AI



Kilito 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

Stanford University

StanfordReport

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

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

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 has become powerful

UC San Diego TODAY

## 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

Density

Very Poor Poor Acceptable Good Very Good

Response Options

Physicians ChatGPT

Not Empathetic Slightly Empathetic Moderately Empathetic Empathetic Very Empathetic

Response Options

Physicians ChatGPT

EXPLORE TIME SUBSCRIBE

JUL 2, 2025 10:11 AM ET

## Microsoft's AI Is Better Than Doctors at Diagnosing Disease

HEALTH AI

Kilito Chan—Getty Images

by [Alice Park](#)  
SENIOR CORRESPONDENT

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

Stanford University

StanfordReport

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

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

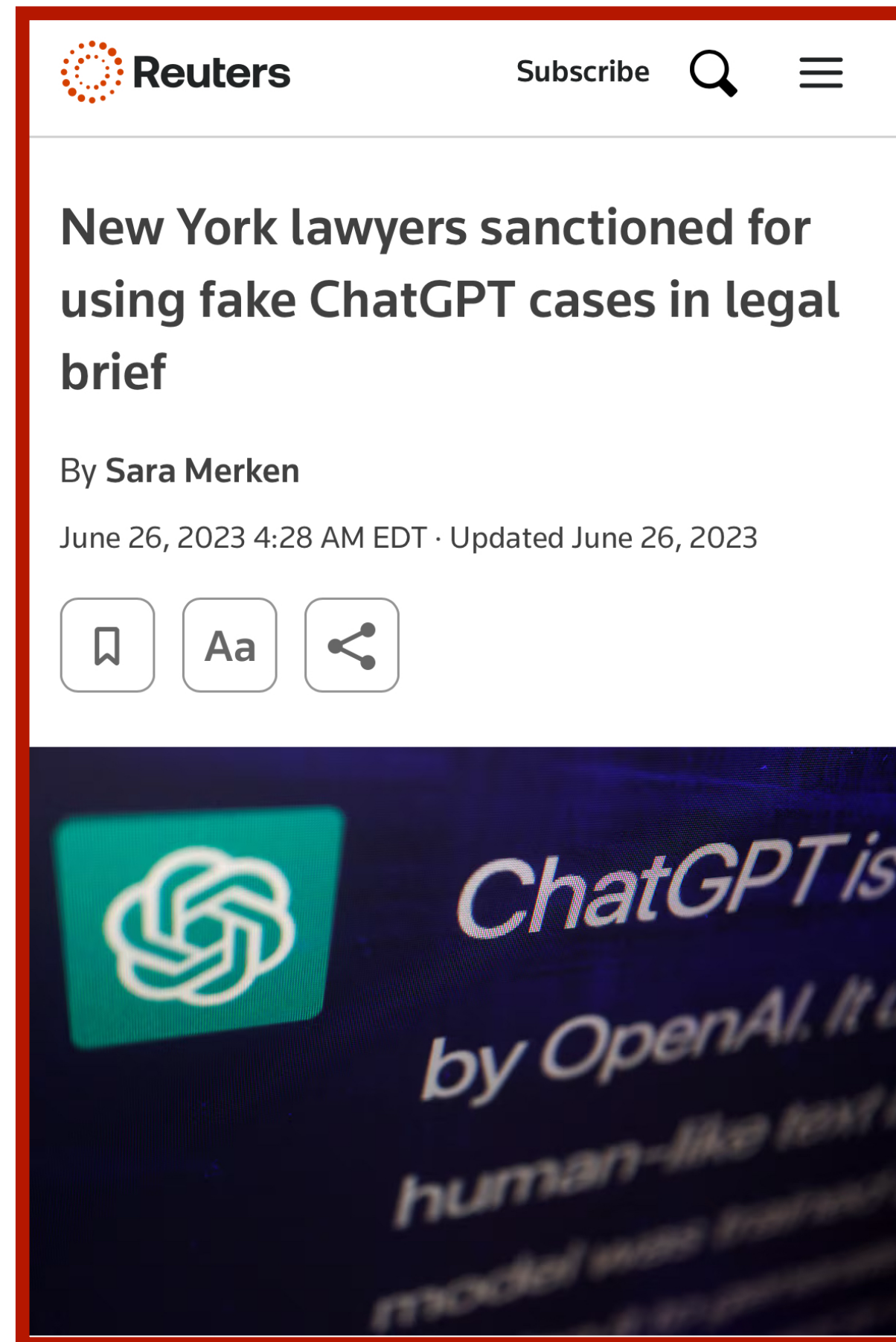
In hopes of creating better access to medical care, Stanford researchers have trained an algorithm to diagnose skin cancer.

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

**Question:** Why not replace human decision makers with AI already?

**Question:** Why not replace human decision makers with AI already?

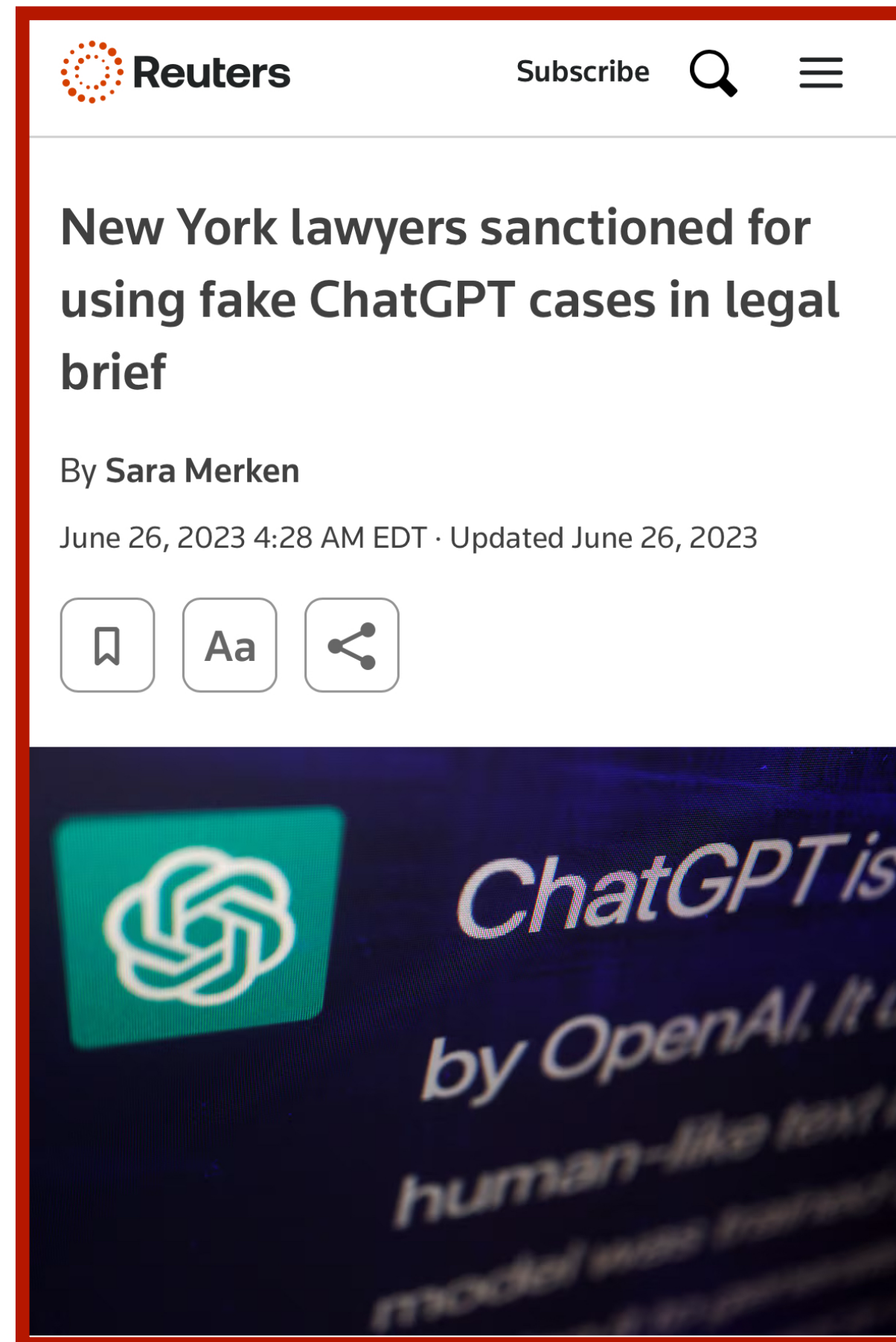
# Question: Why not replace human decision makers with AI already?



The image is a screenshot of a news article from Reuters. At the top left is the Reuters logo, followed by a 'Subscribe' button, a search icon, and a menu icon. The main headline reads 'New York lawyers sanctioned for using fake ChatGPT cases in legal brief'. Below the headline, it says 'By Sara Merken' and 'June 26, 2023 4:28 AM EDT · Updated June 26, 2023'. There are three icons below the text: a bookmark icon, an 'Aa' font size icon, and a share icon. The bottom portion of the screenshot shows a blurred image of a blue banner with the OpenAI logo and the text 'ChatGPT is by OpenAI. It is human-like text model...'.

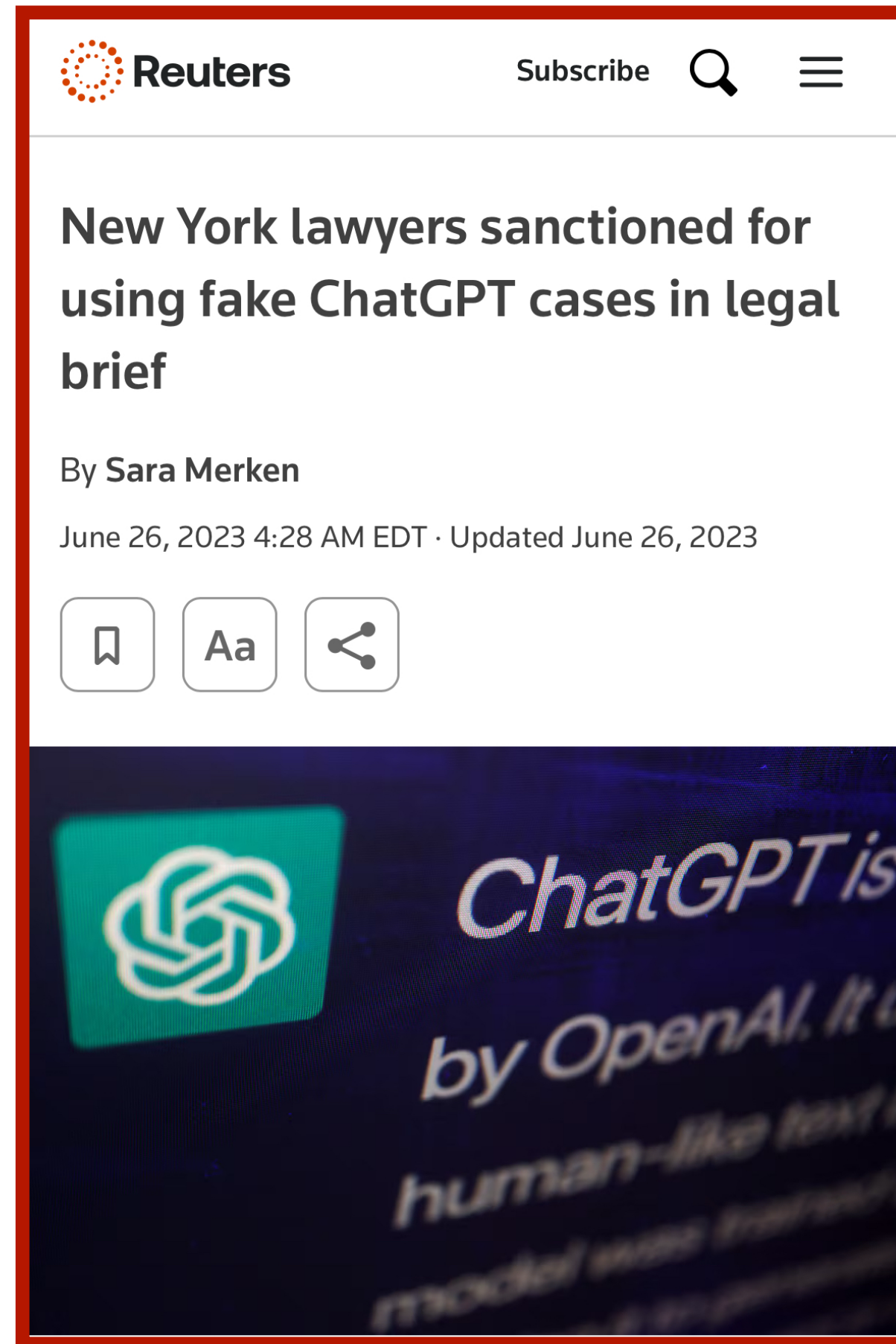
# Question: Why not replace human decision makers with AI already?



## Hallucinations



# Question: Why not replace human decision makers with AI already?




## Hallucinations

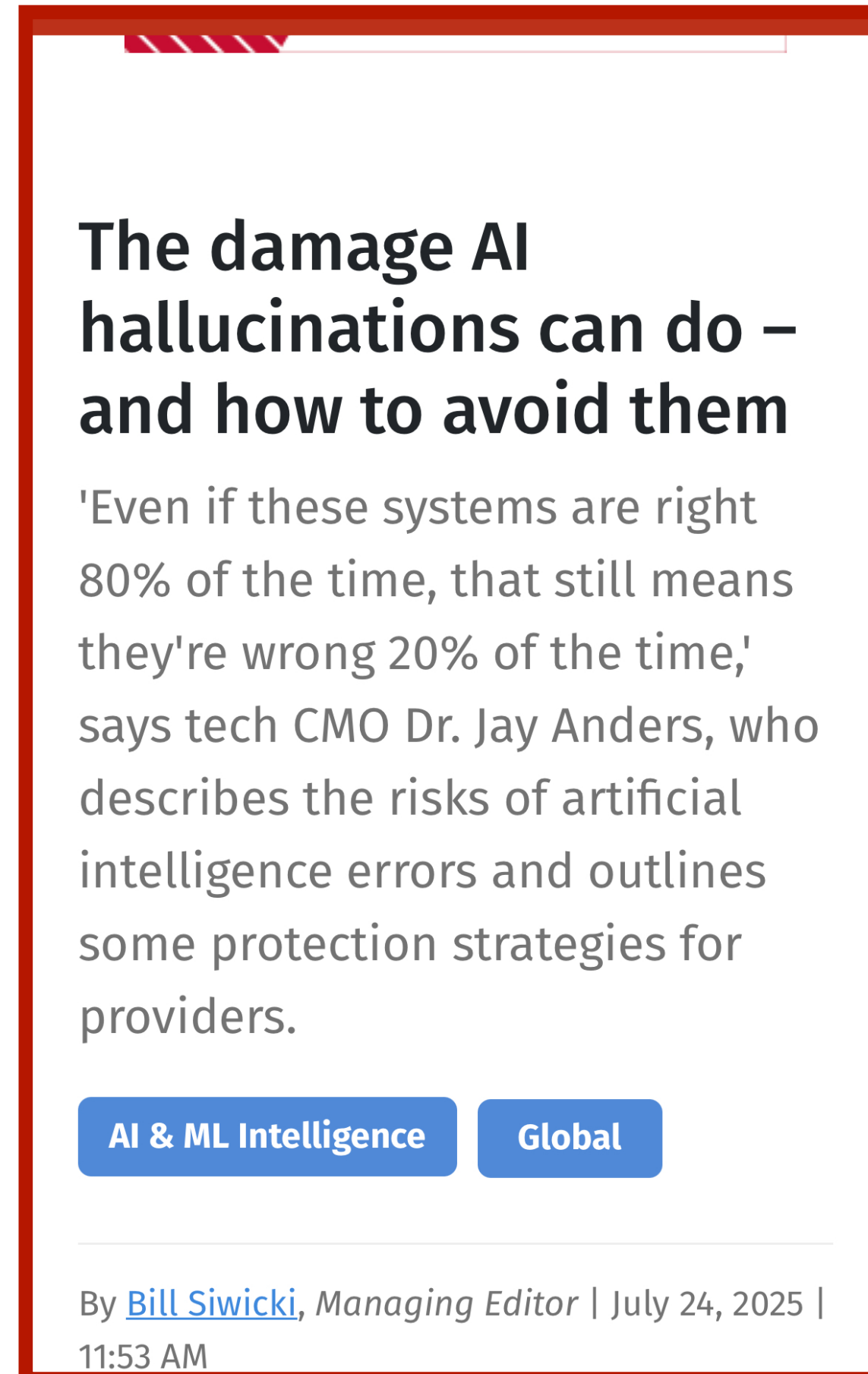
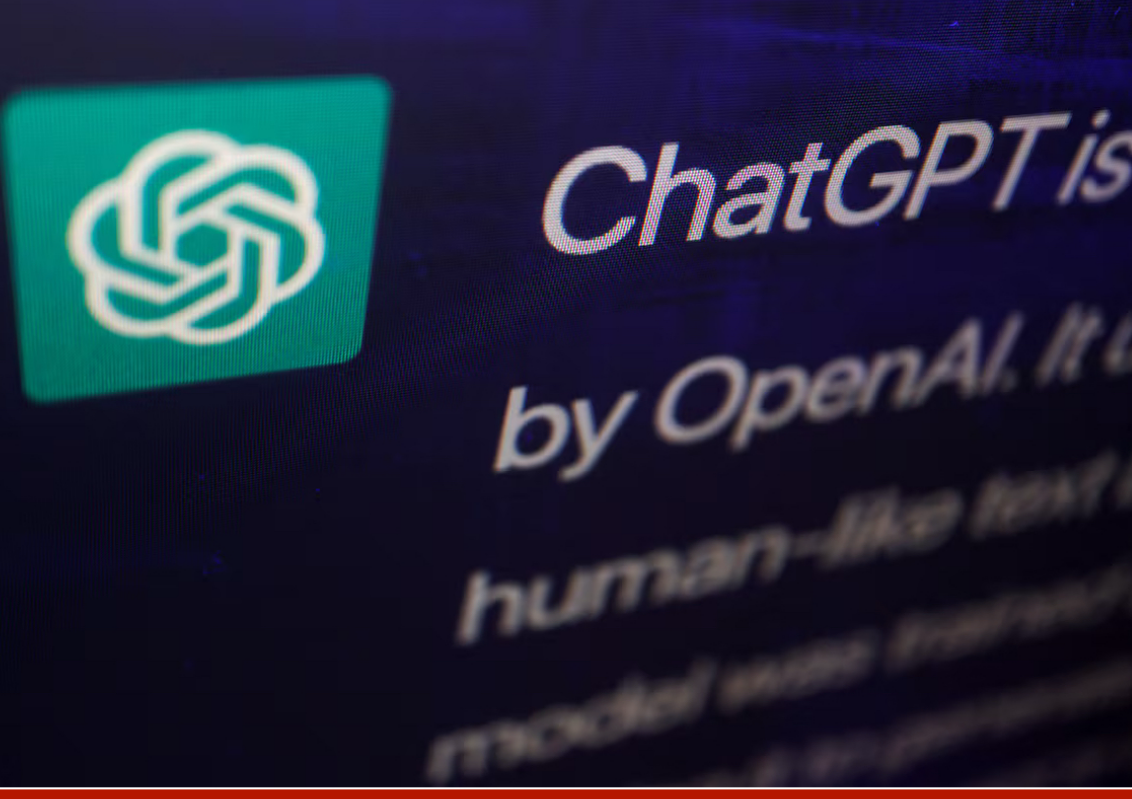


Reuters Subscribe  

### 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



### The damage AI hallucinations can do – and how to avoid them

'Even if these systems are right 80% of the time, that still means 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.

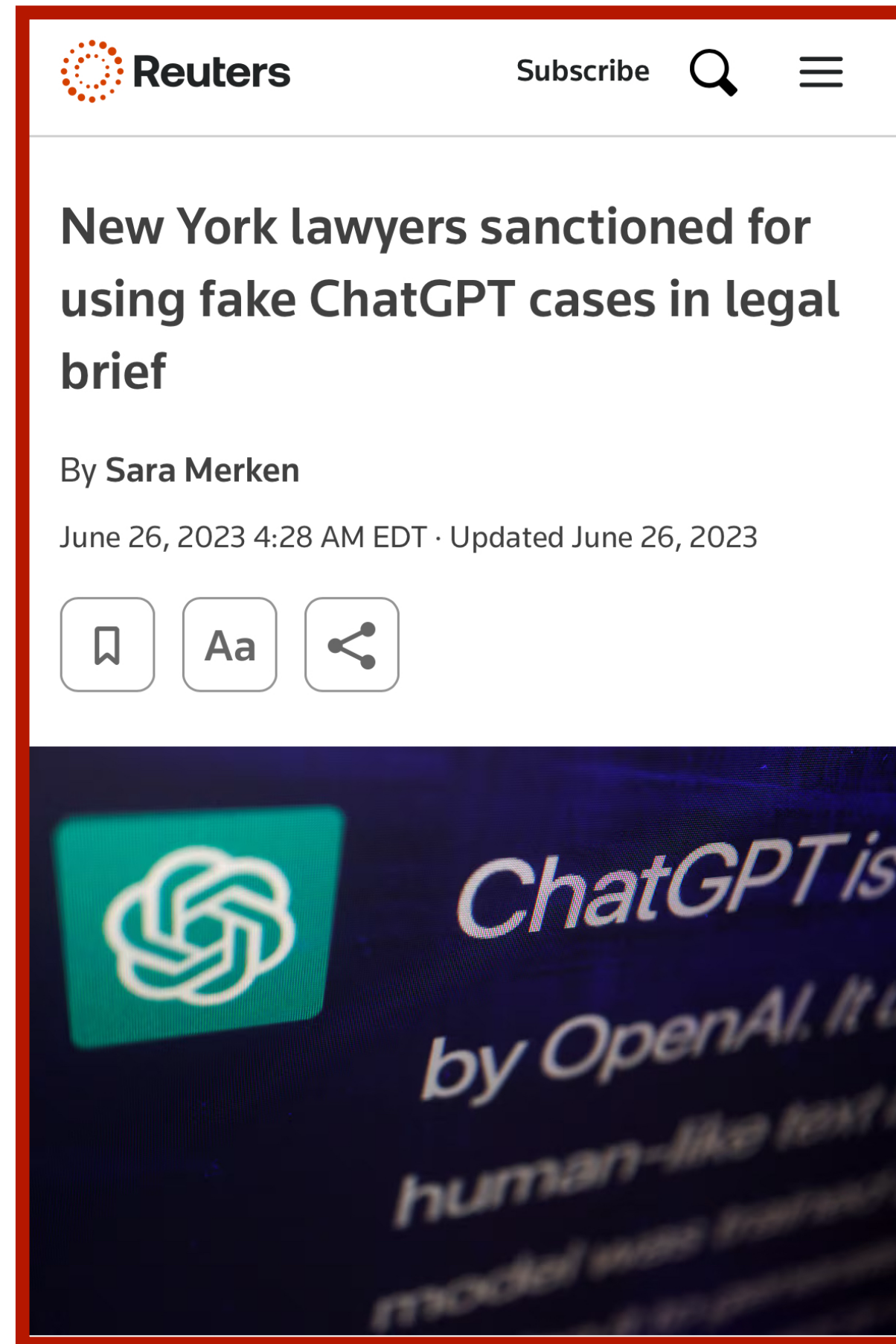
[AI & ML Intelligence](#) [Global](#)



---

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

# Question: Why not replace human decision makers with AI already?




## Hallucinations

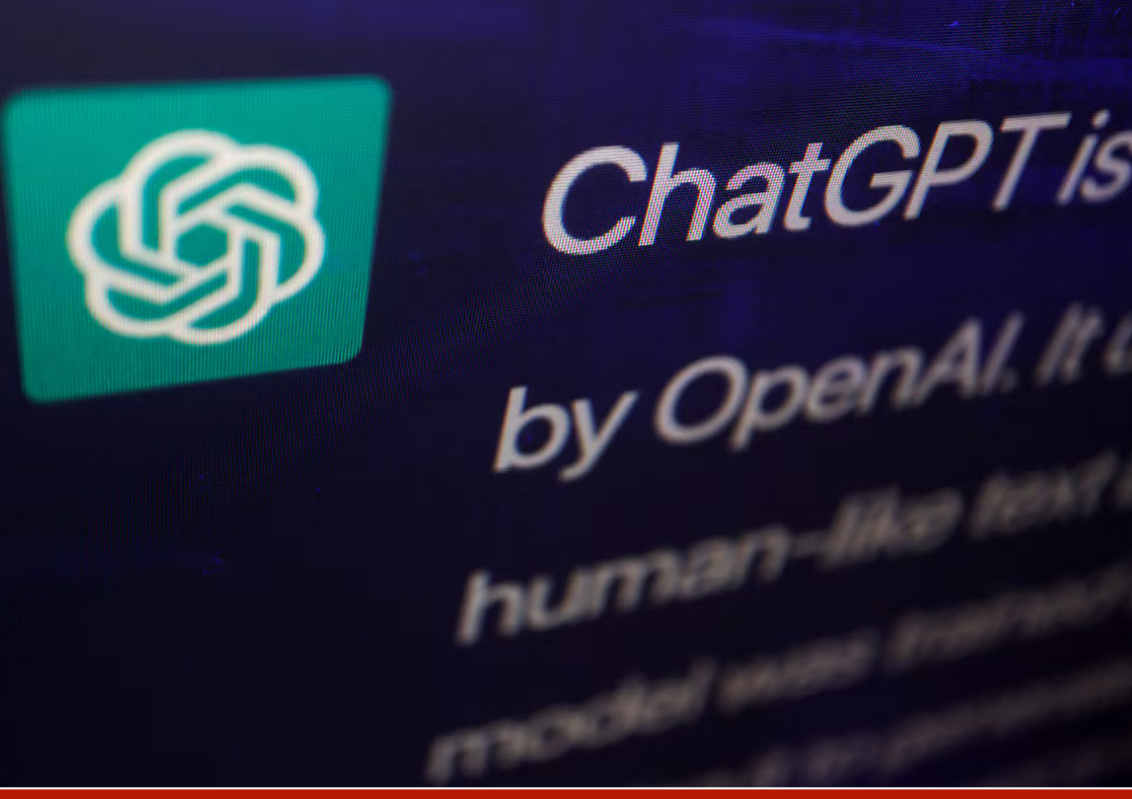


Reuters Subscribe  

### 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



### The damage AI

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.

describes the risks of artificial intelligence errors and outlines some protection strategies for providers.

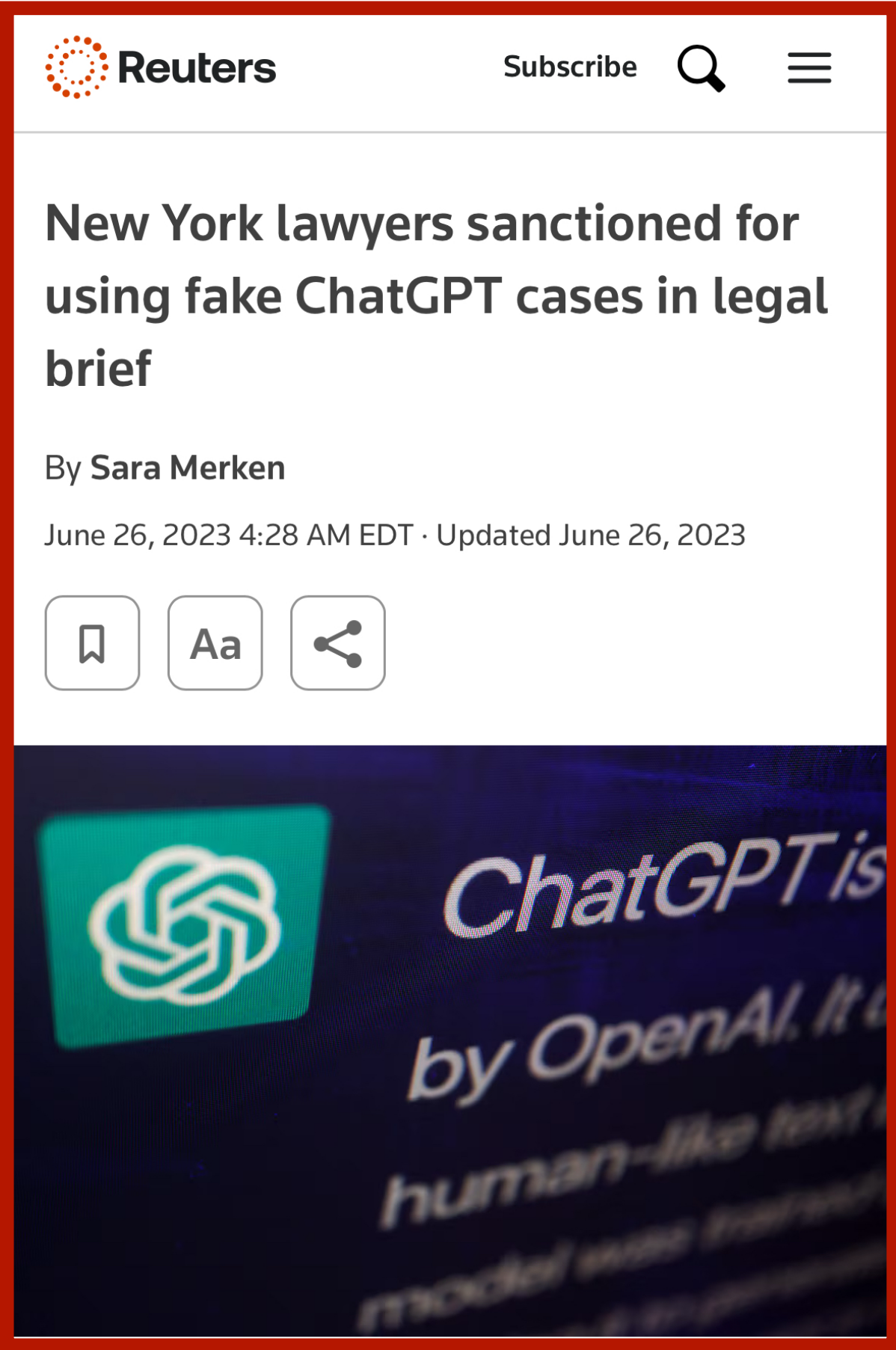
AI & ML Intelligence

Global

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

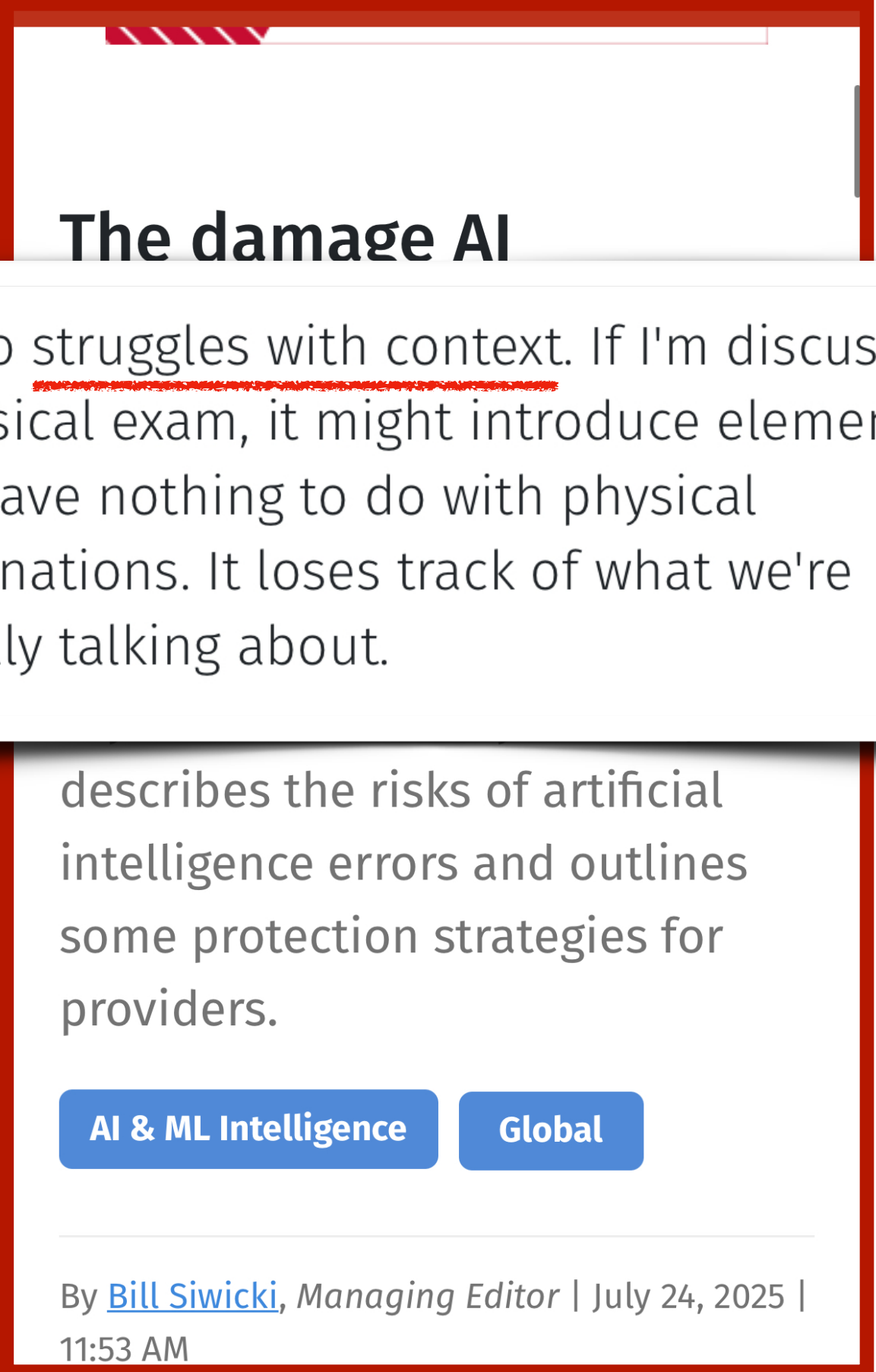
# Question: Why not replace human decision makers with AI already?

## Hallucinations



The screenshot shows a Reuters article. At the top left is the Reuters logo and a 'Subscribe' button. The article title is 'New York lawyers sanctioned for using fake ChatGPT cases in legal brief'. The author is Sara Merken, and the date is June 26, 2023. Below the title are icons for bookmarking, font size adjustment, and sharing. The bottom part of the image shows a blurred background with the OpenAI logo and the text 'ChatGPT is by OpenAI. It's human-like text model'.

## Struggles with long-term context



The screenshot shows an article titled 'The damage AI'. A text box highlights a paragraph: '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.' Below the text box, the article continues: 'describes the risks of artificial intelligence errors and outlines some protection strategies for providers.' There are two blue buttons: 'AI & ML Intelligence' and 'Global'. At the bottom, it says 'By Bill Siwicki, Managing Editor | July 24, 2025 | 11:53 AM'.

# Question: Why not replace human decision makers with AI already?

## Hallucinations

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    🔄

ChatGPT is by OpenAI. It's human-like text model

## Struggles with long-term context

### The damage AI

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.

describes the risks of artificial intelligence errors and outlines some protection strategies for providers.

AI & ML Intelligence    Global

By [Bill Siwicky](#), 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    Ƴ    ✉

📄 Add Axios on Google

Illustration: Allie Carl/Axios

# Question: Why not replace human decision makers with AI already?

## Hallucinations

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    🔄

ChatGPT is by OpenAI. It's human-like text model

## Struggles with long-term context

### The damage AI

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.

describes the risks of artificial intelligence errors and outlines some protection strategies for providers.

AI & ML Intelligence    Global

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

## Unaware when uncertain

AXIOS Pittsburgh    🔍    👤

Aug 27, 2025 - News

### AI is overconfident even when wrong, says report

Ryan Deto

f    X    in    Ƴ    ✉

📄 Add Axios on Google

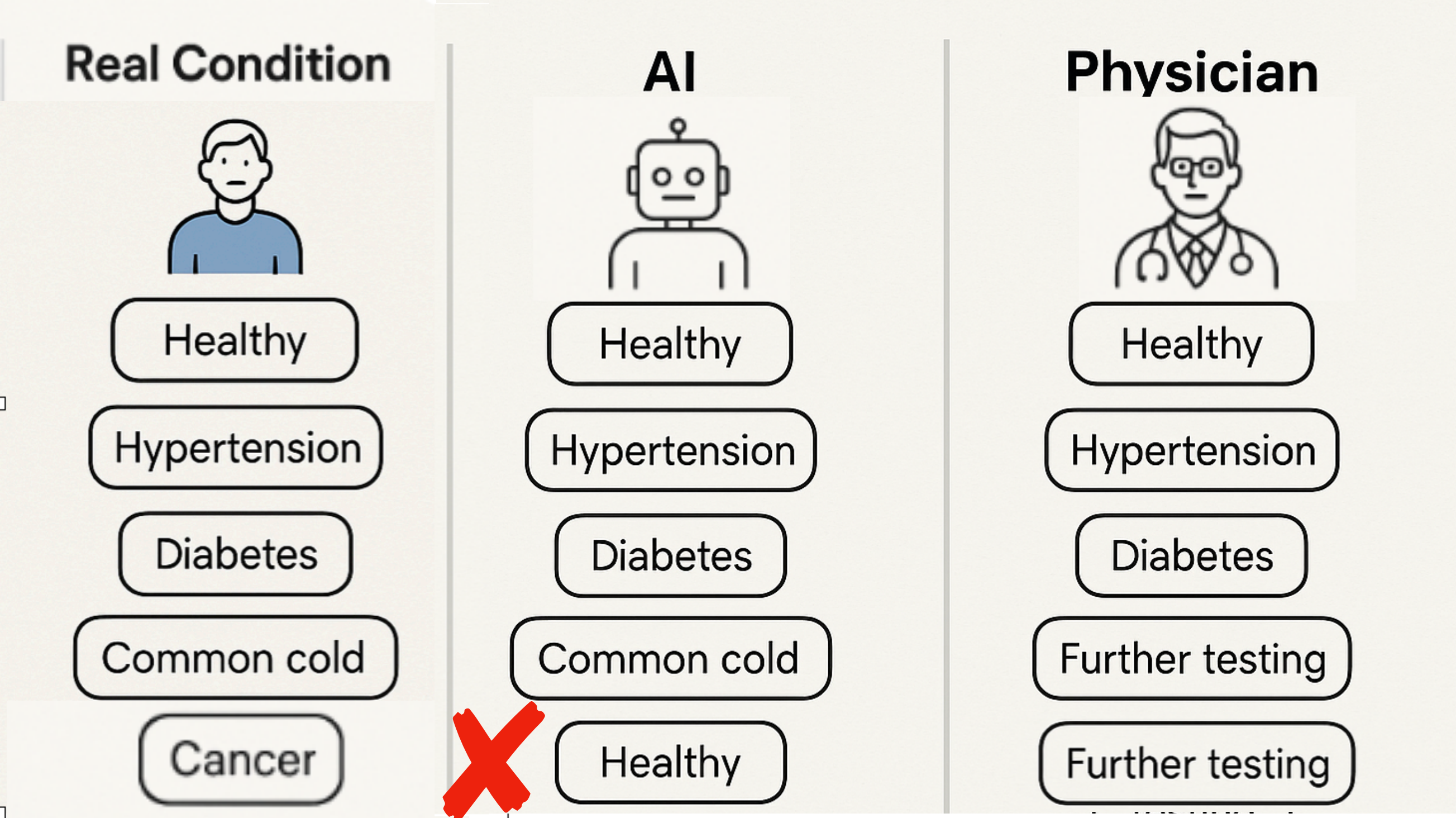
Illustration: Allie Carl/Axios

# Question: Why not replace human decision makers with AI already?

Hallucinations

Struggles with long-term context

Unaware when uncertain

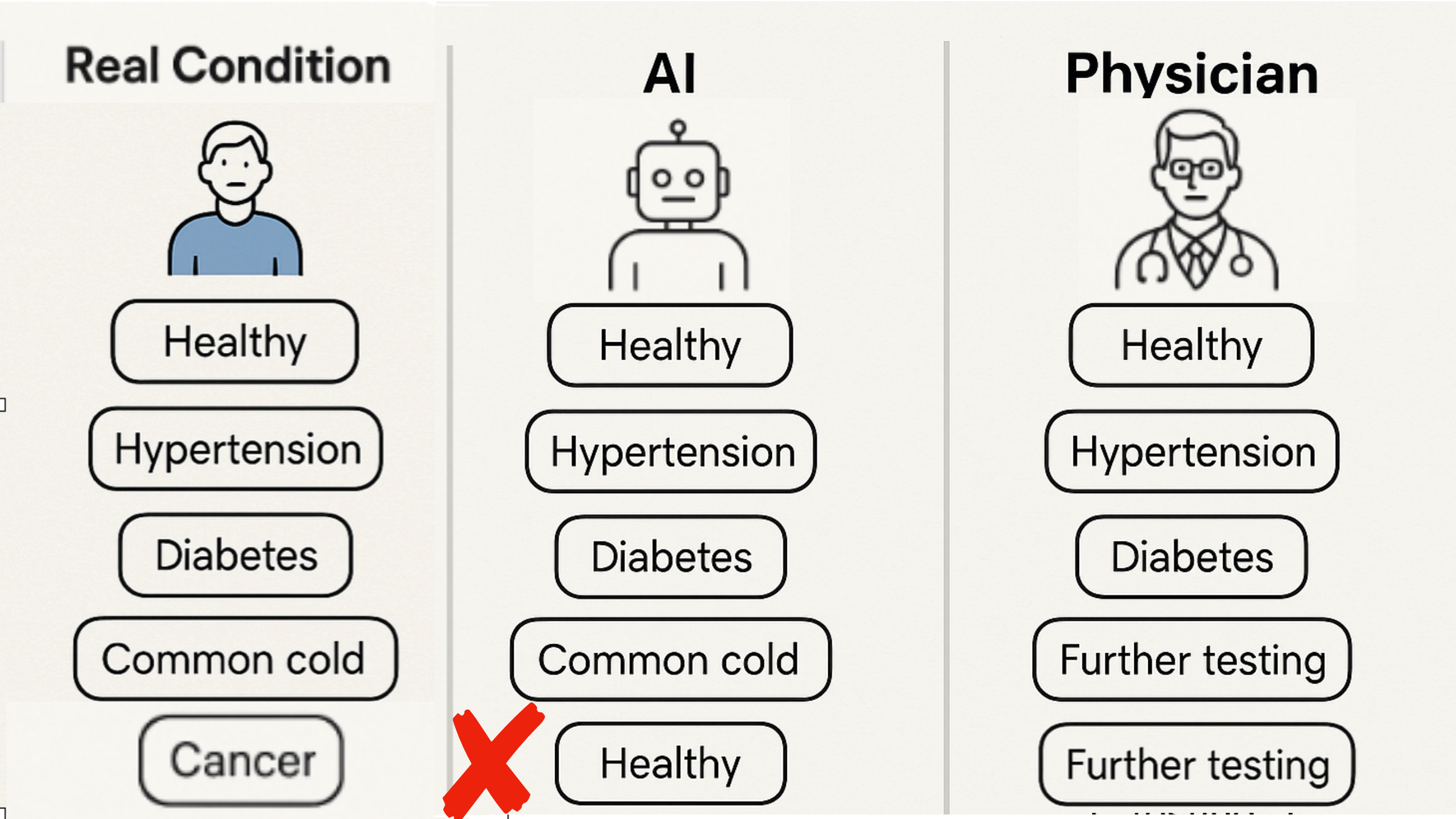


# Question: Why not replace human decision makers with AI already?

**Hallucinations**

**Struggles with long-term context**

**Unaware when uncertain**

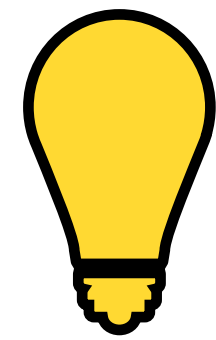


**Reliability matters!**

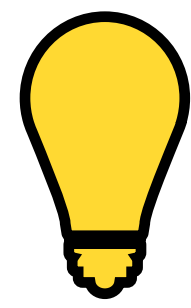
**Robust decision making requires more than predictive accuracy!**

**Question:** Why not replace human decision makers with AI already?

**Question:** Why not replace human decision makers with AI already?



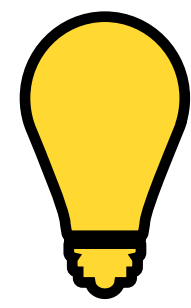
Human and AI should **coexist** in the **decision making** pipeline



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

AI

Handle large unstructured data

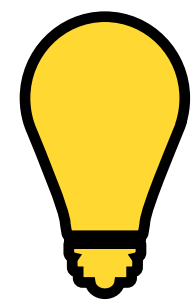


Human and AI should **coexist** in the **decision making** pipeline

AI

Handle large unstructured data

Excel at patten extraction



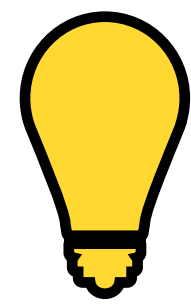
**Human and AI should coexist in the decision making pipeline**

**AI**

**Handle large unstructured data**

**Excel at patten extraction**

**Offer statistical accuracy**



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

**AI**

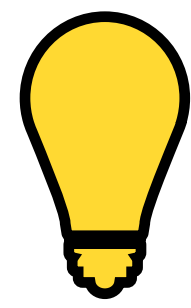
**Handle large unstructured data**

**Excel at patten extraction**

**Offer statistical accuracy**



**Domain Knowledge**



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

**AI**

**Handle large unstructured data**

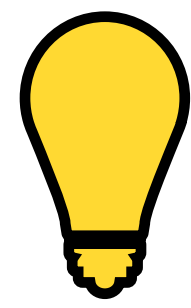
**Excel at patten extraction**

**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**

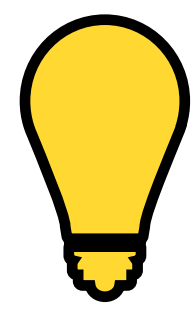
**Offer statistical accuracy**



**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**

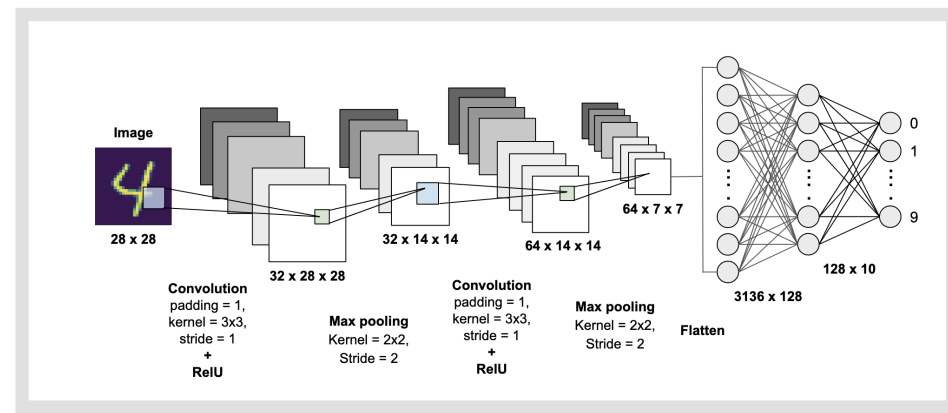
**Reason and act in  
the physical world**

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

# Decision making pipeline

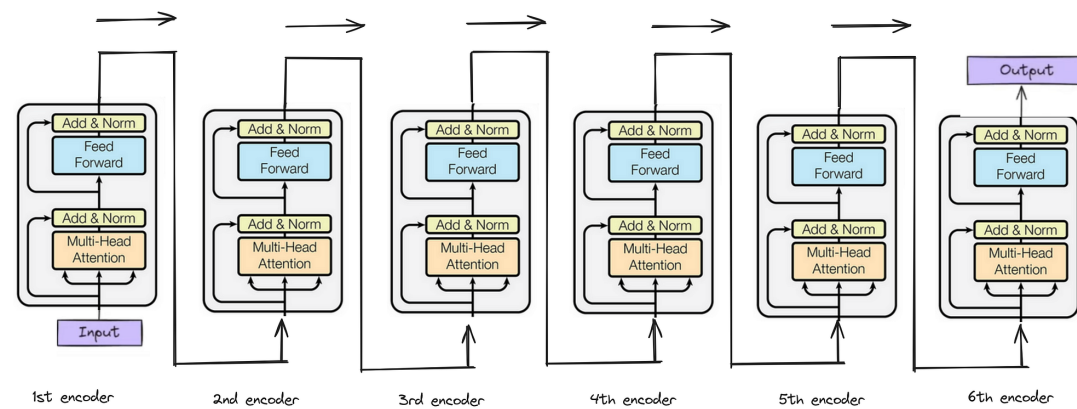


# Decision making pipeline

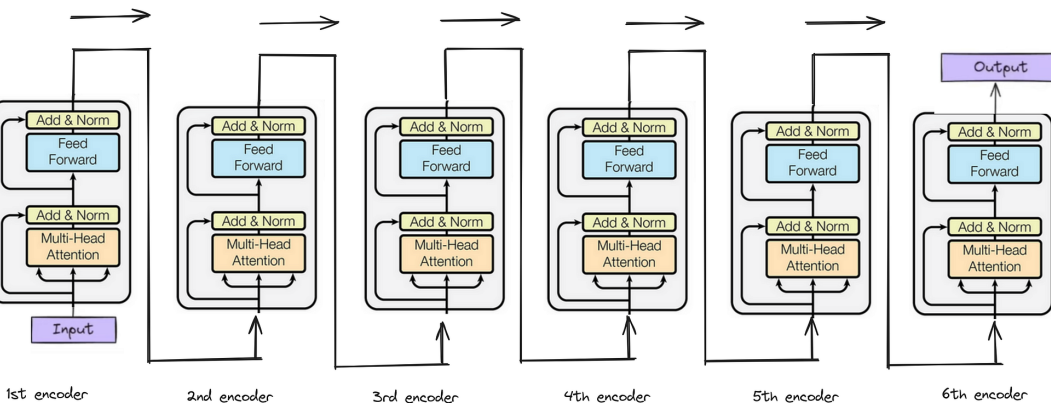
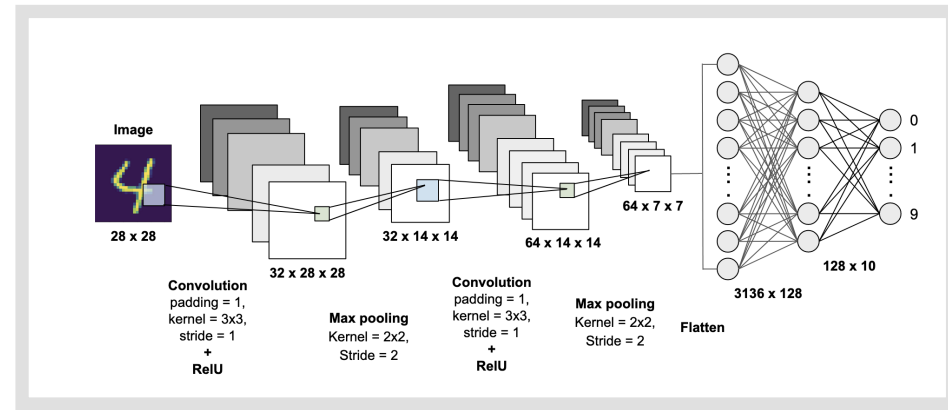


- (Heavily) Pre-trained

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



# Decision making pipeline

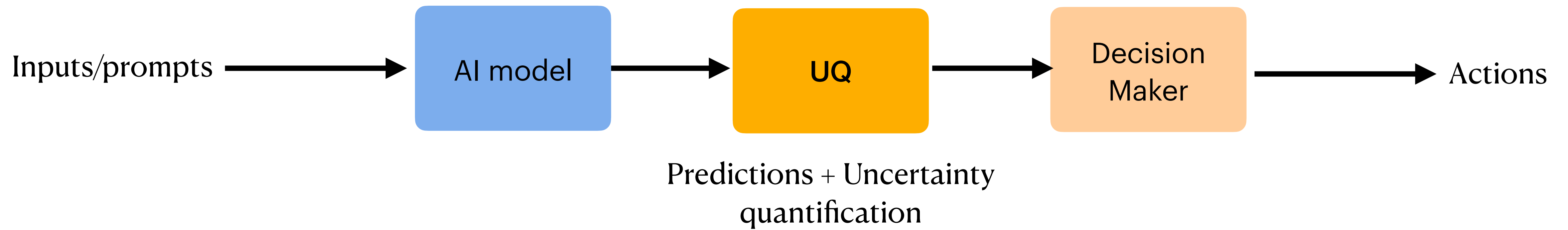


- (Heavily) Pre-trained
- Possibility of little tweaks using extra data (fine-tuning)

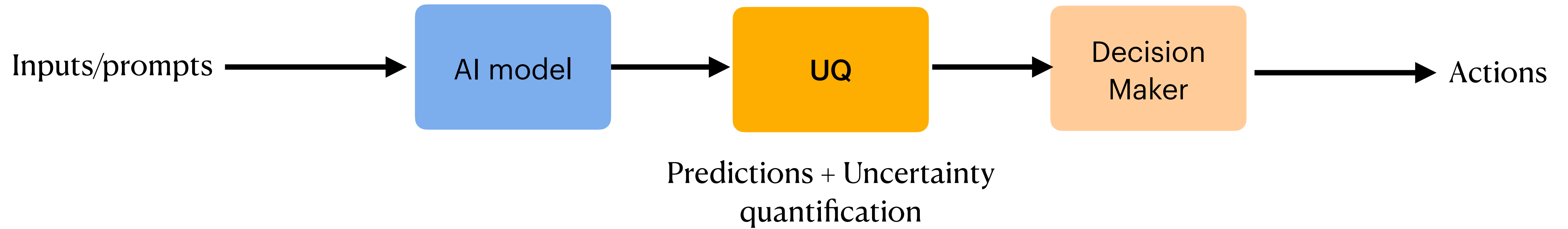
- Automated
- Humans



# Robust decision making pipeline



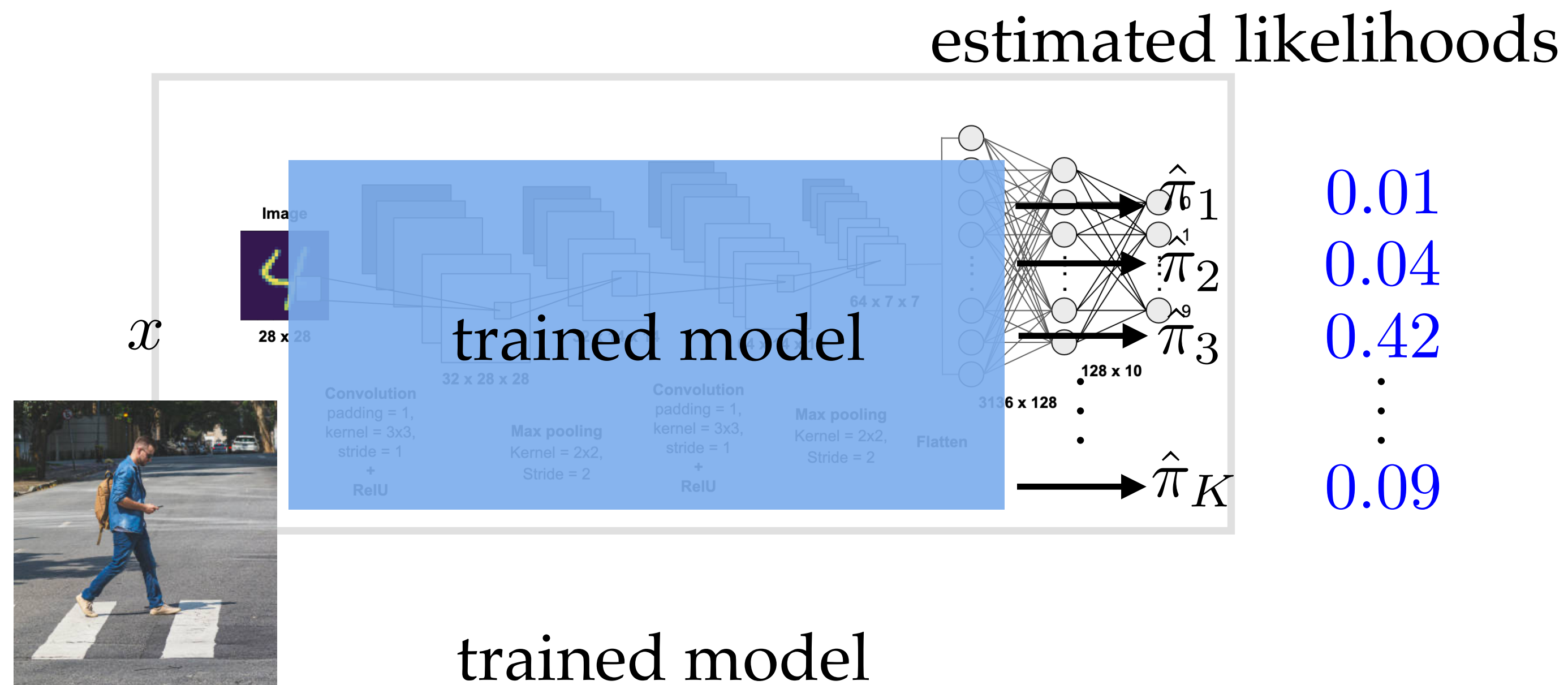
# Robust decision making pipeline



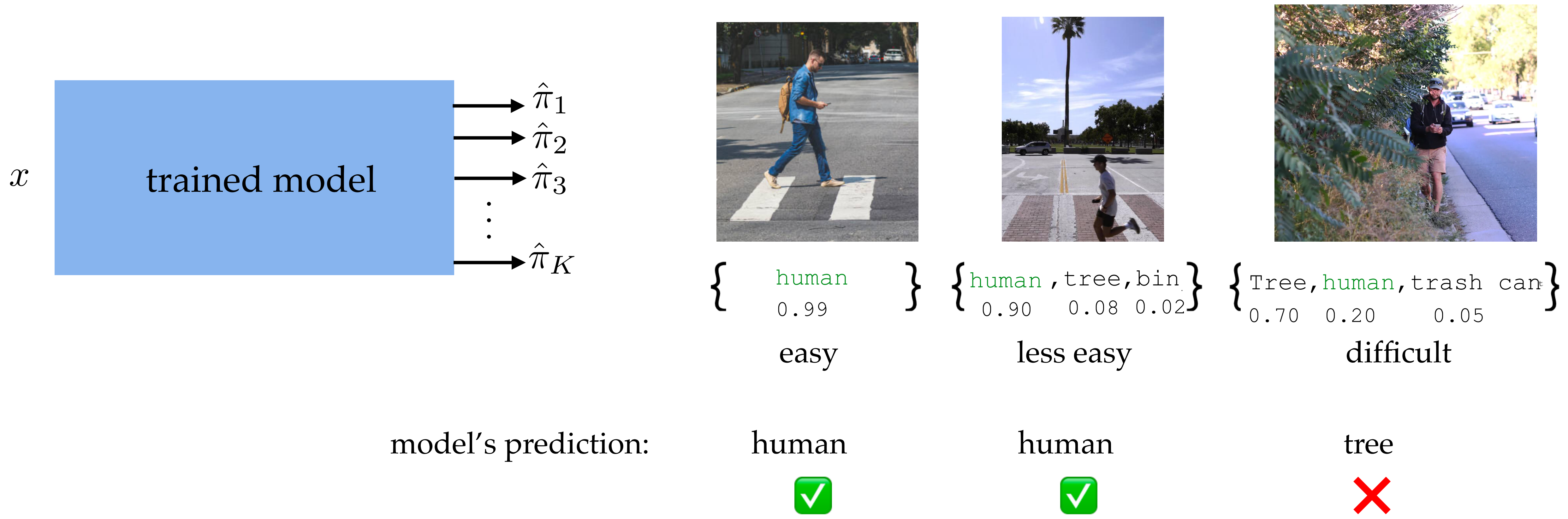
**Robust Decision making requires precise Uncertainty Quantification**

# **(Quick) Review of Conformal Prediction**

# Conformal Prediction



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

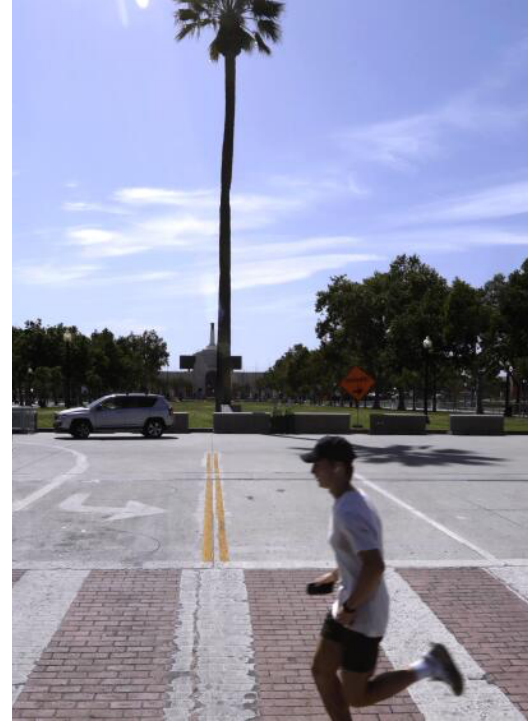


- the trained model provides estimated likelihoods (a notion of **uncertainty**)
- these likelihoods are informative but not always correct

$x$



{ 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



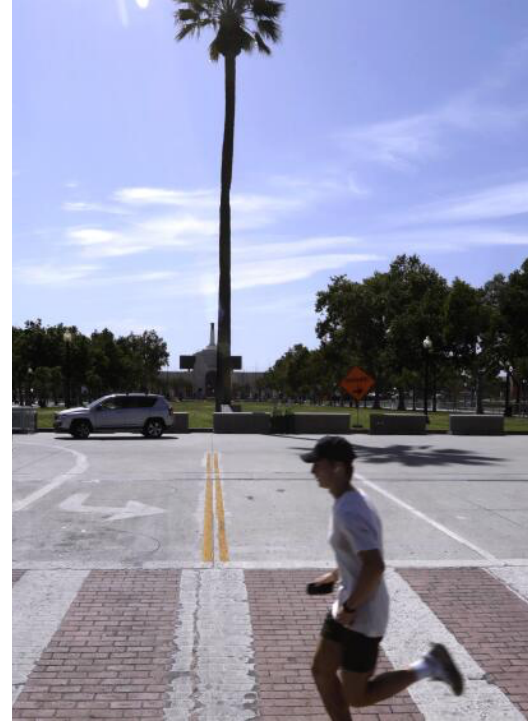
human



tree



$x$



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

$C(x)$

(model's conformal prediction)

$\{ \text{human} \}$



$\{ \text{human} \}$

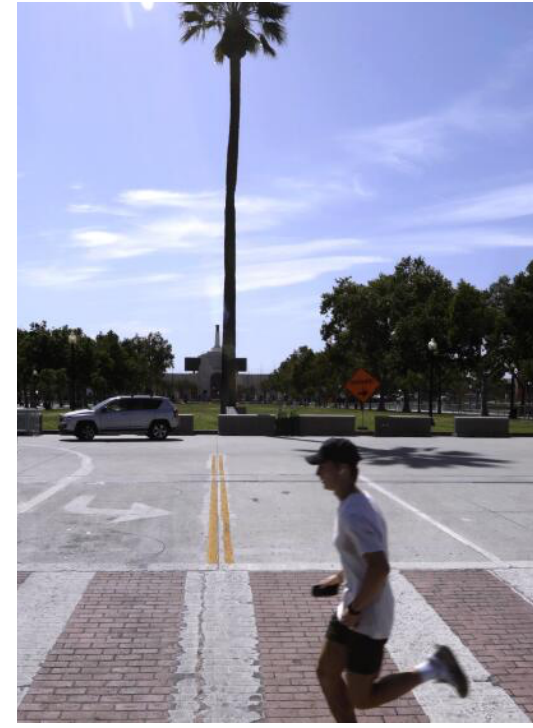


$\{ \text{tree}, \text{human} \}$



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 }

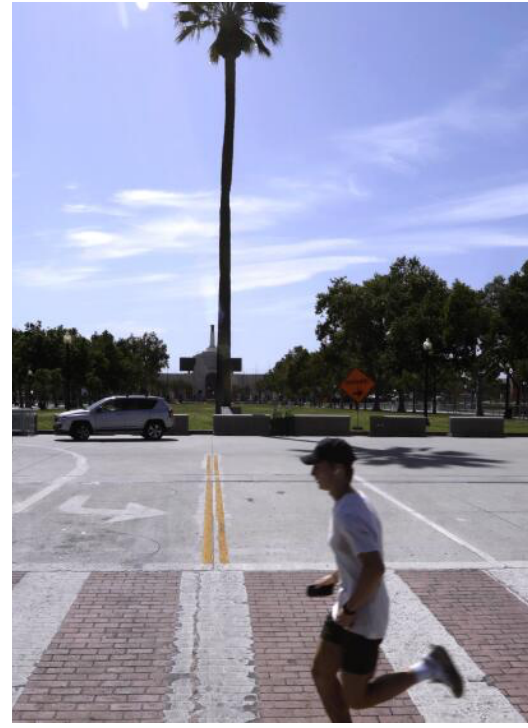


### Marginal Coverage

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

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

$x$



{ human }  
0.99

{ human, tree, bin }  
0.90 0.08 0.02

{ Tree, human, trash can }  
0.70 0.20 0.05

$C(x)$

{ human }

{ human }

{ tree, human }

(model's conformal prediction)



### Marginal Coverage

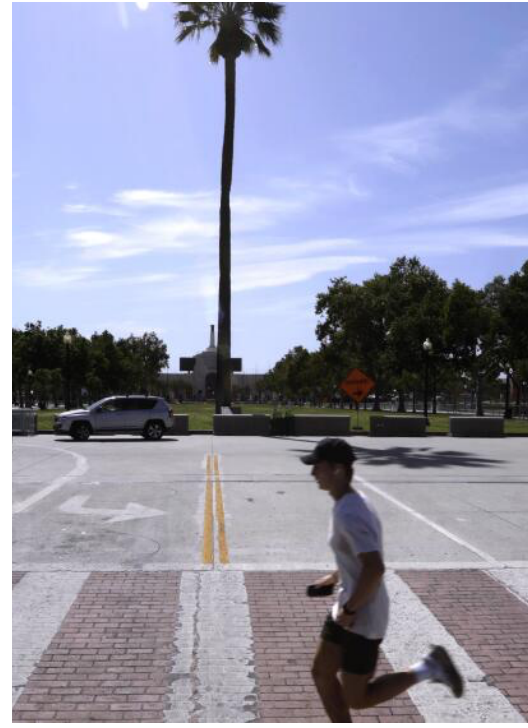
$$\Pr \{Y_{\text{test}} \in C(X_{\text{test}})\} \geq 1 - \alpha$$

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

the threshold



$x$



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

$C(x)$

(model's conformal prediction)

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



$|C(x)|$

size

1

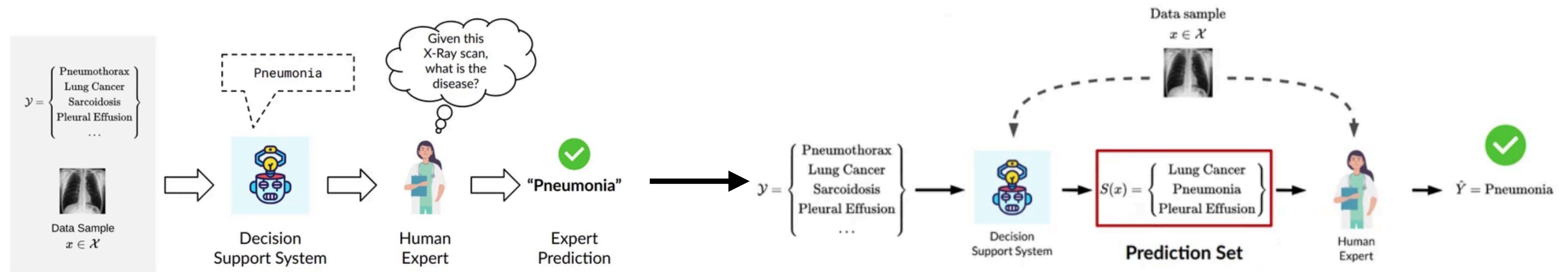
1

2

difficulty of the input  $\uparrow$  models' uncertainty about the label  $\uparrow$  size of the prediction set  $\uparrow$

**Prediction sets are useful for decision makers**

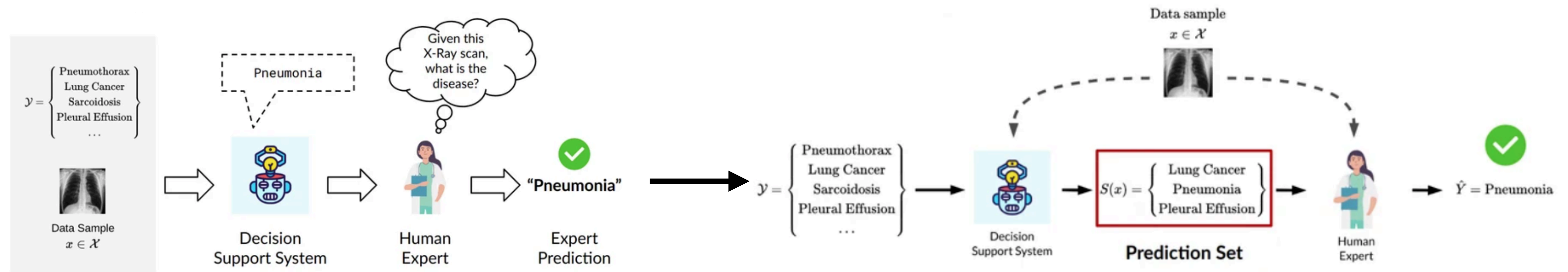
# Prediction sets are useful for decision makers



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.

# Prediction sets are useful for decision makers



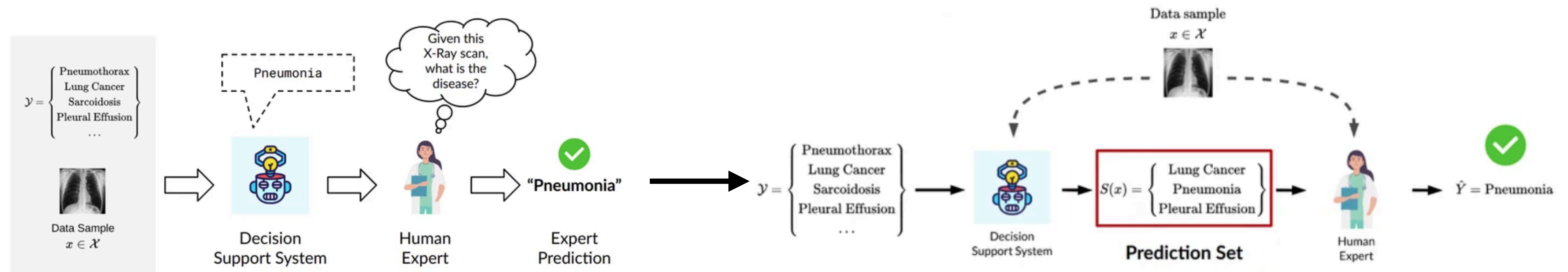
Straitouri et al., Improving Expert Predictions with Conformal Prediction, ICML, 2023.

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

“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.

“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**

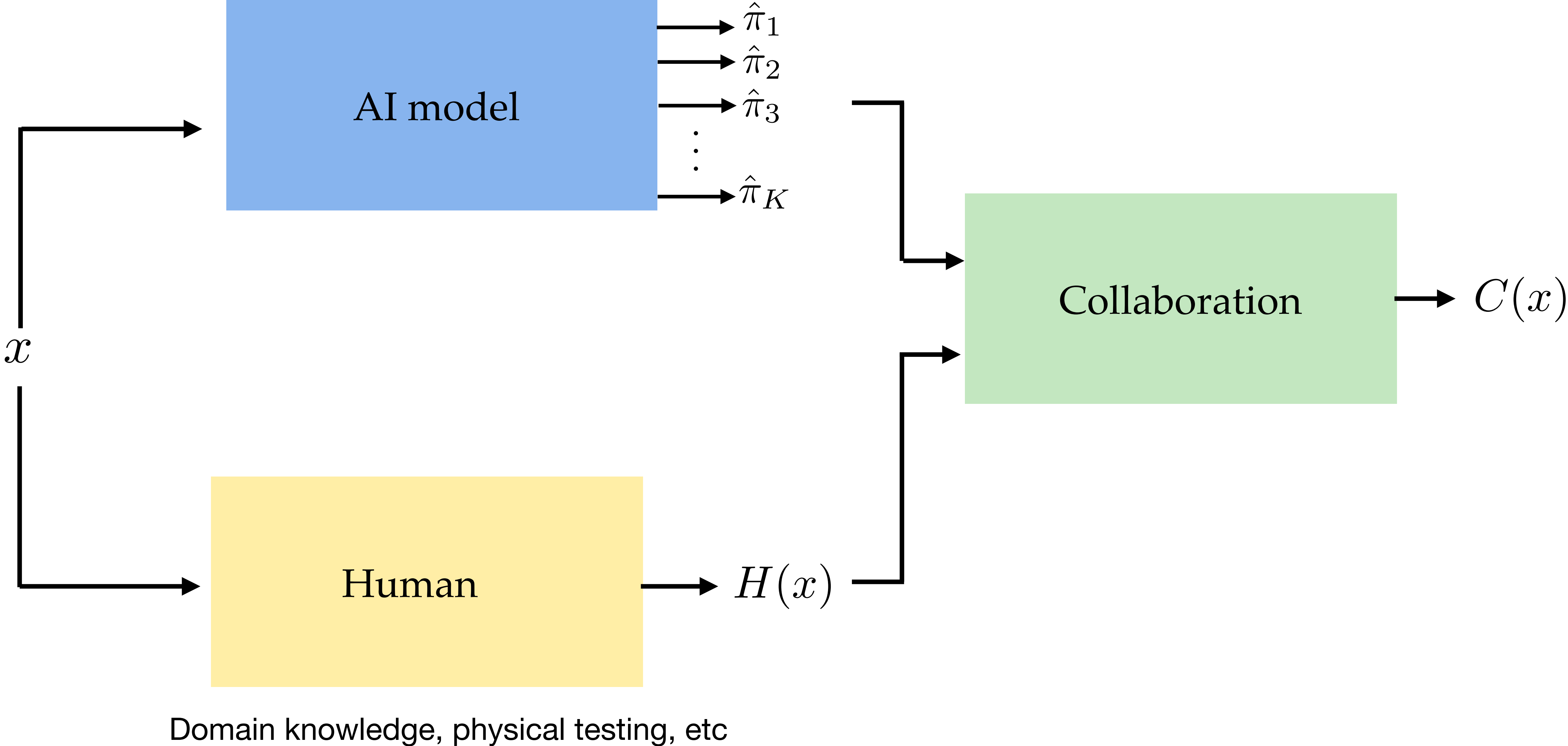
UQ



**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?**

# Collaborative Uncertainty Quantification pipeline

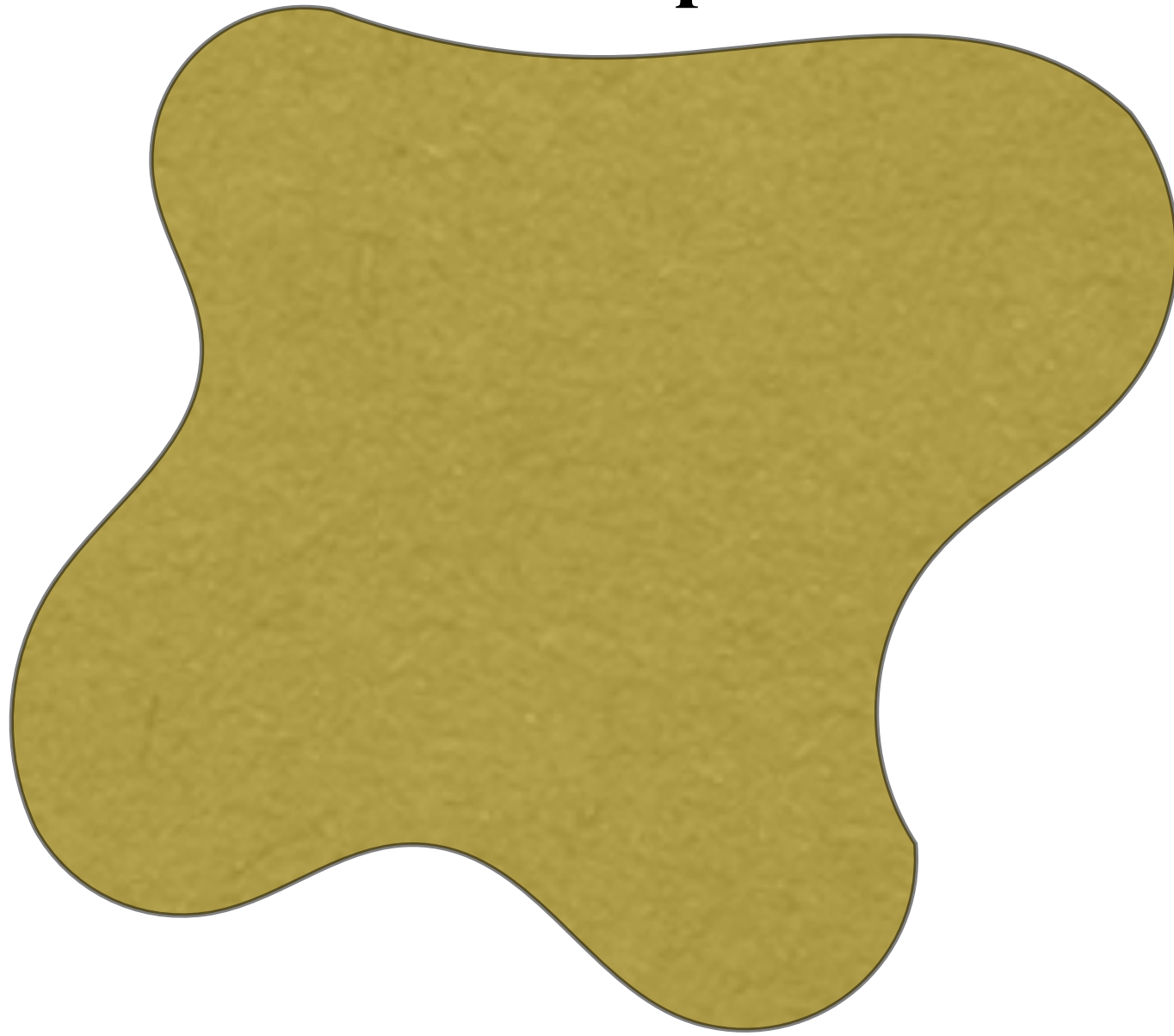


**Question:** what constitutes a good collaboration?

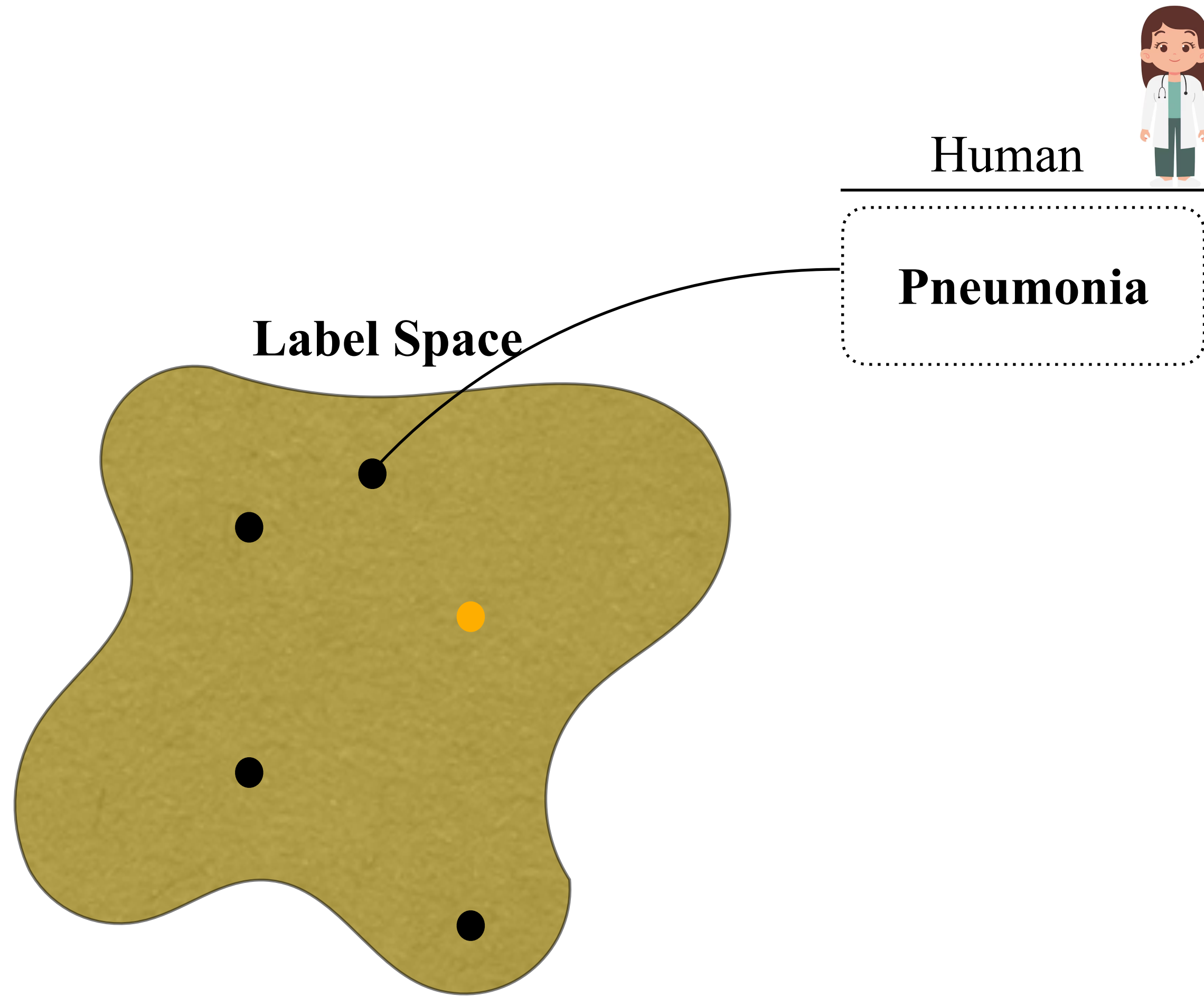
**Question:** what constitutes a good collaboration?

**Question:** what constitutes a good collaboration?

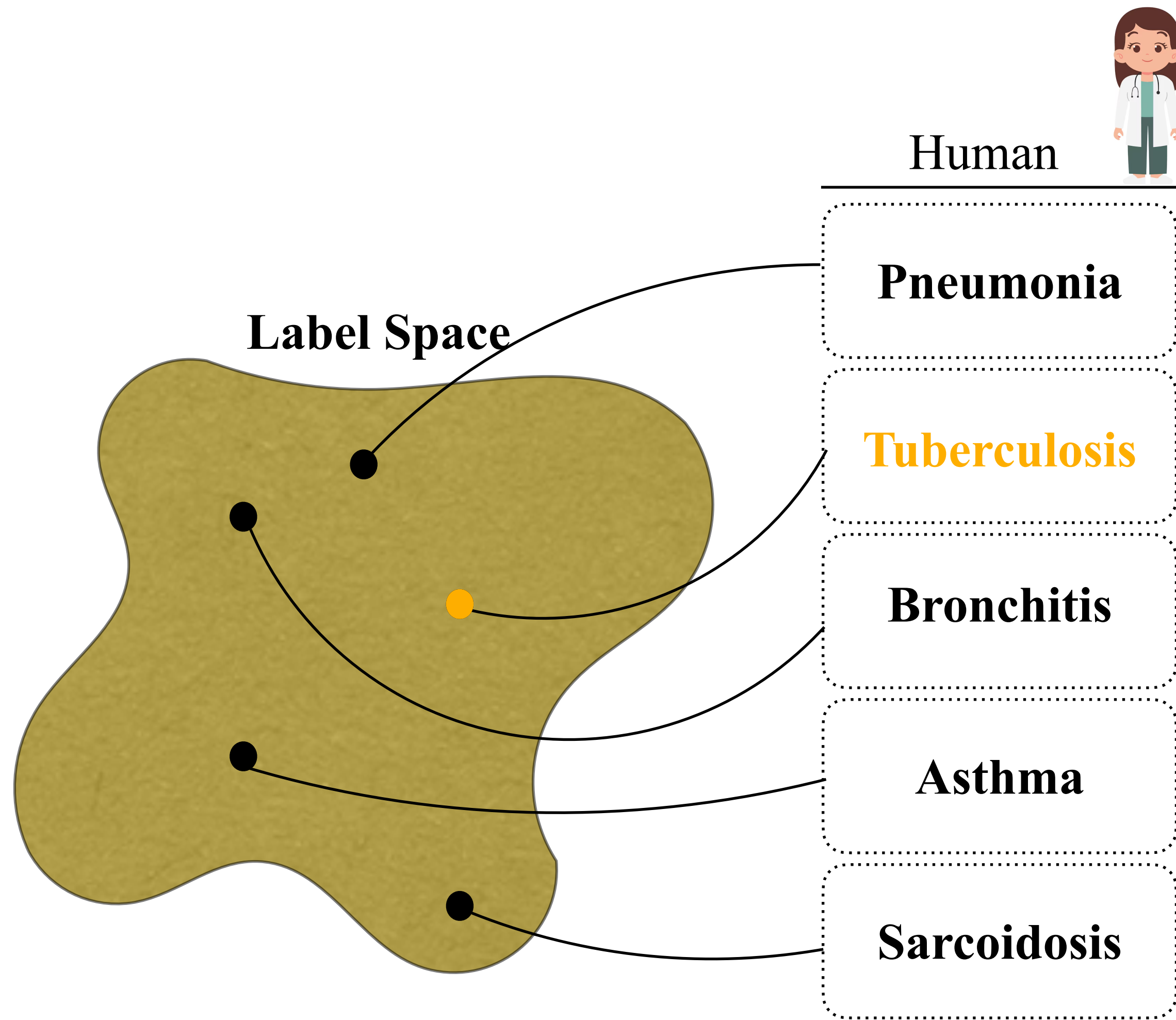
**Label Space**



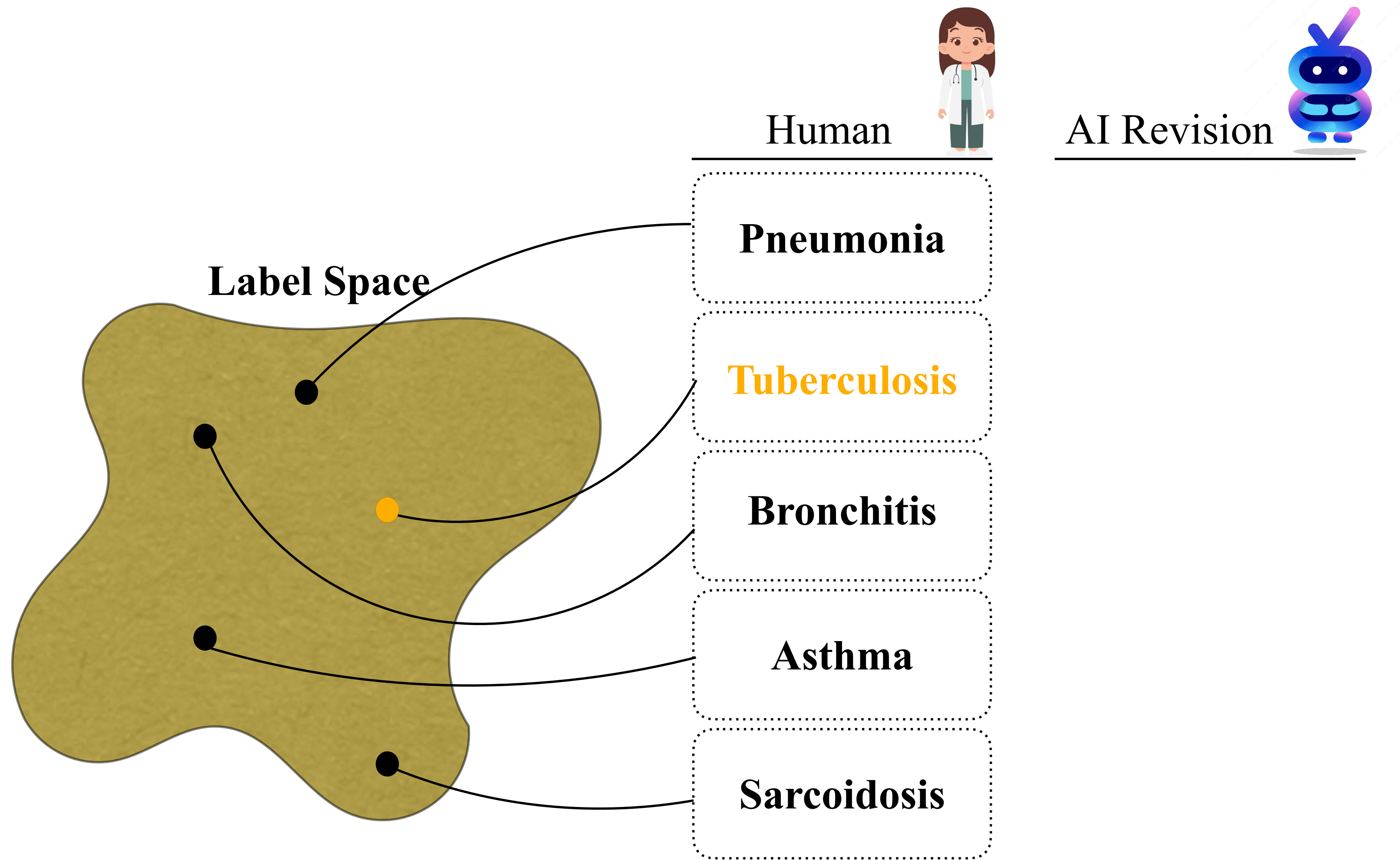
# Question: what constitutes a good collaboration?



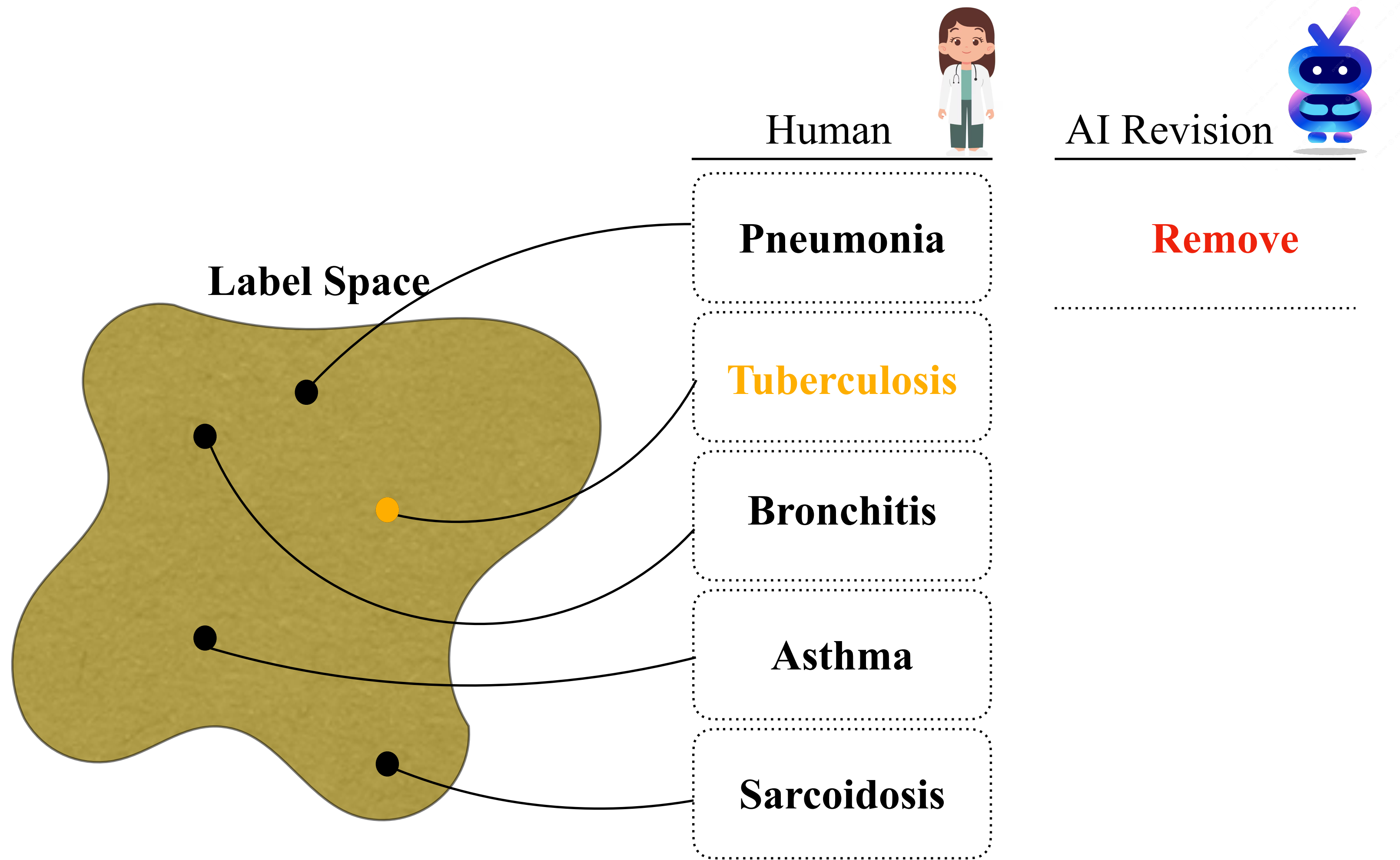
# Question: what constitutes a good collaboration?



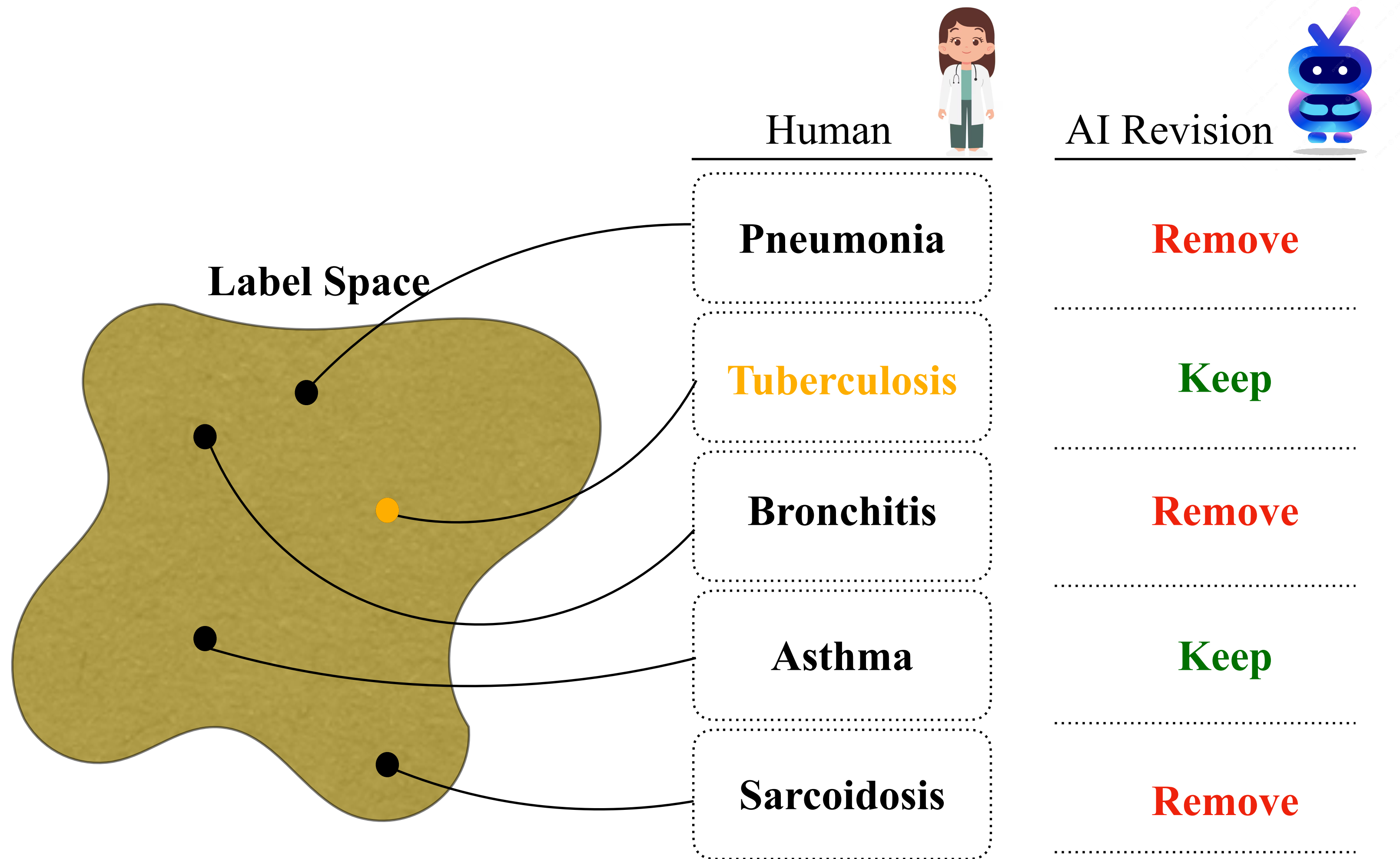
# Question: what constitutes a good collaboration?



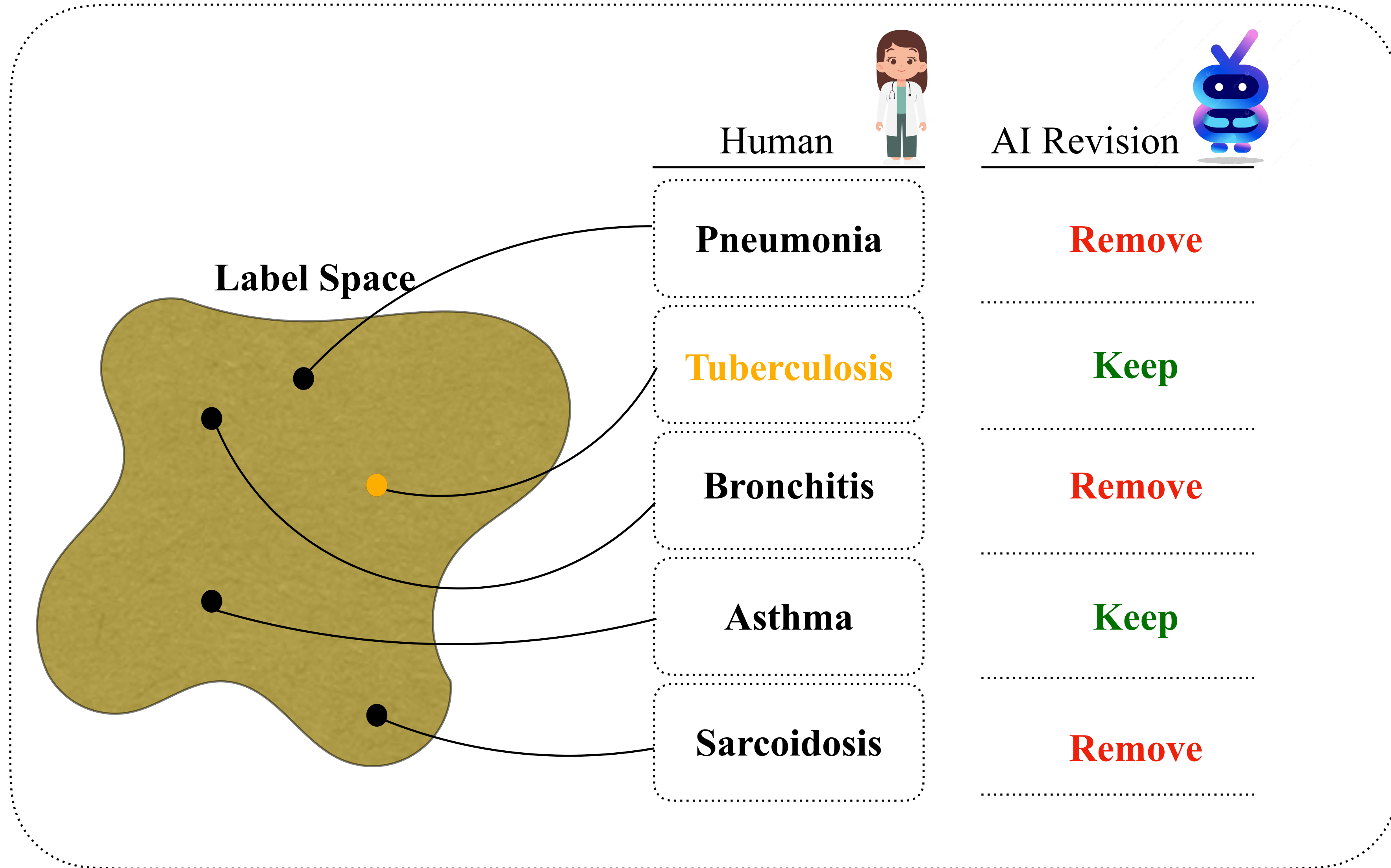
# Question: what constitutes a good collaboration?



# Question: what constitutes a good collaboration?

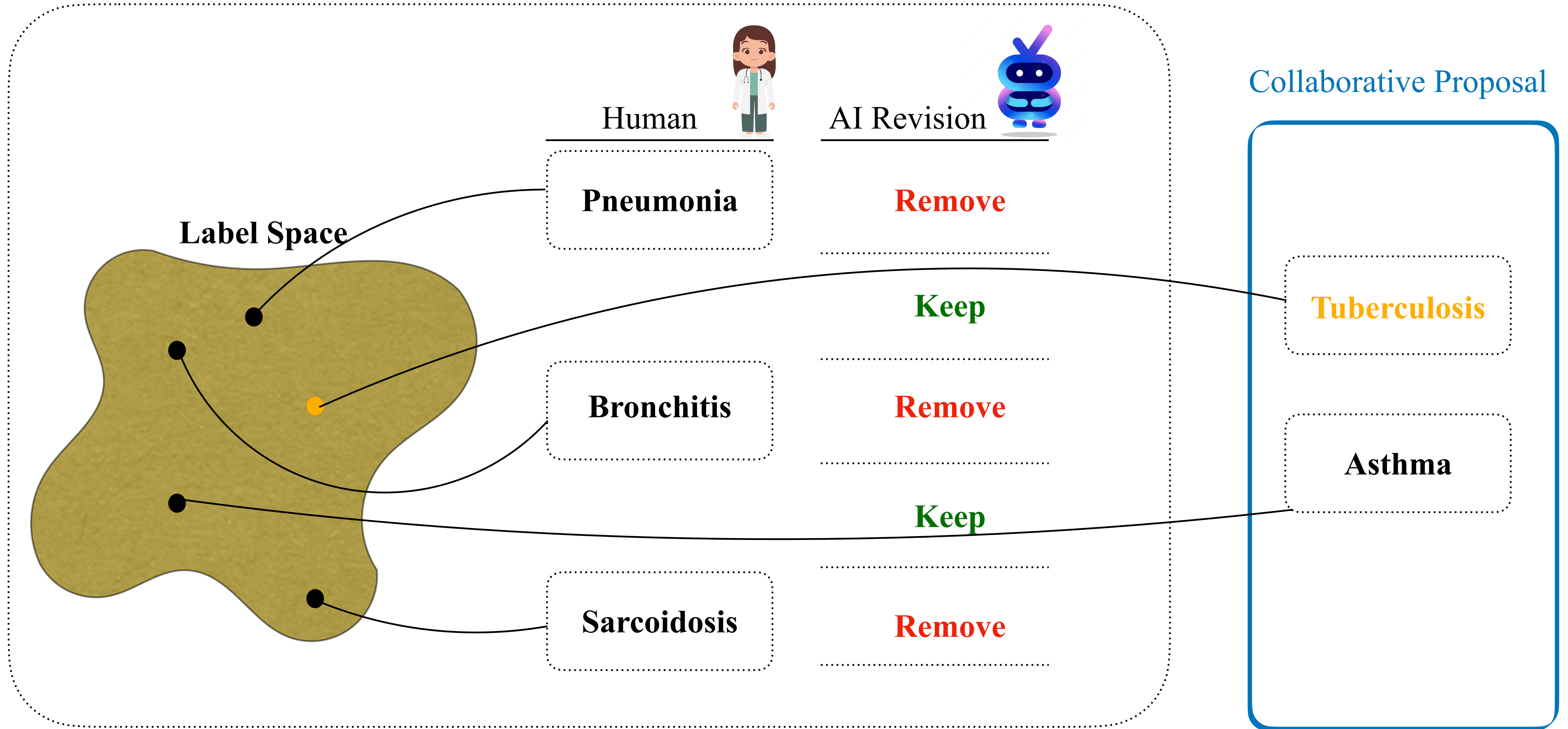


# Question: what constitutes a good collaboration?

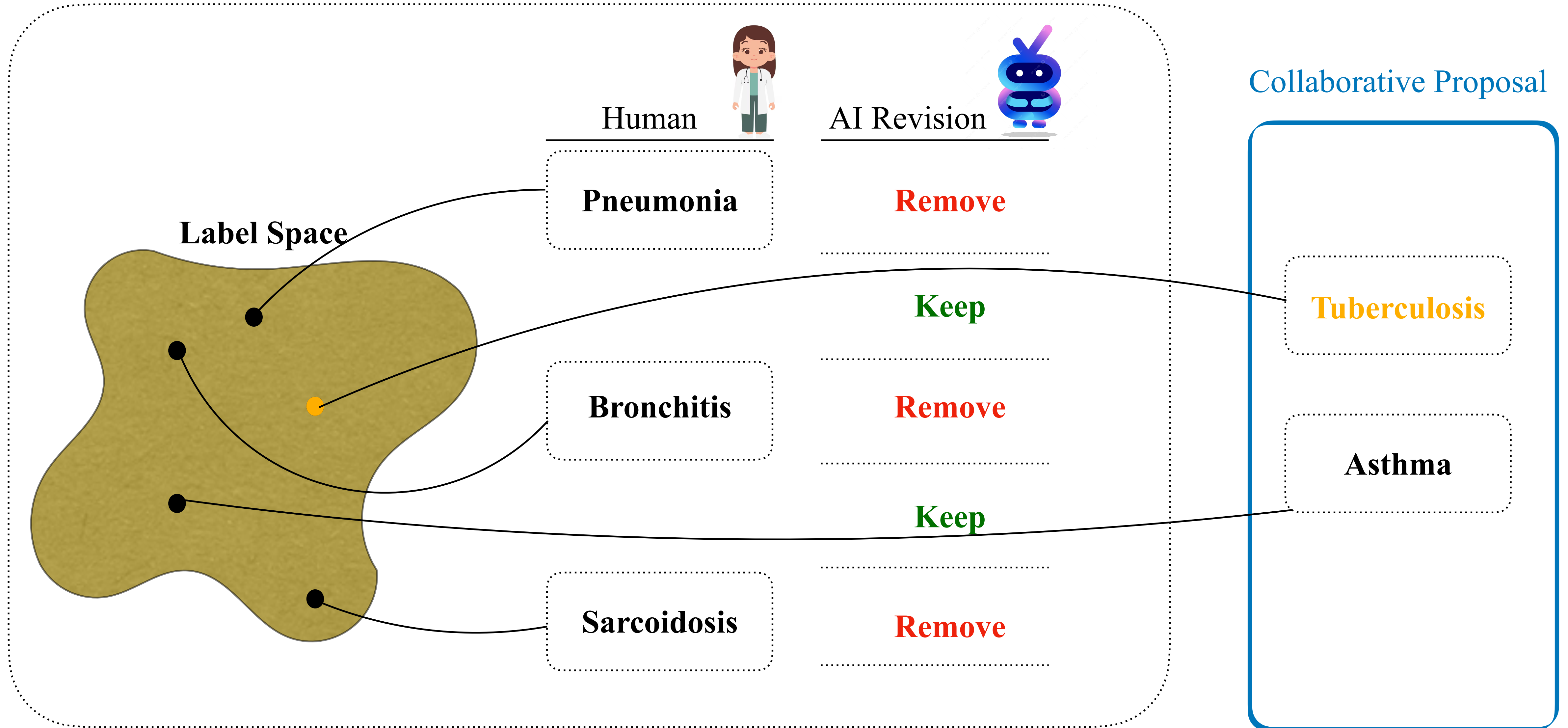


## Collaborative Proposal

# Question: what constitutes a good collaboration?



# Question: what constitutes a good collaboration?



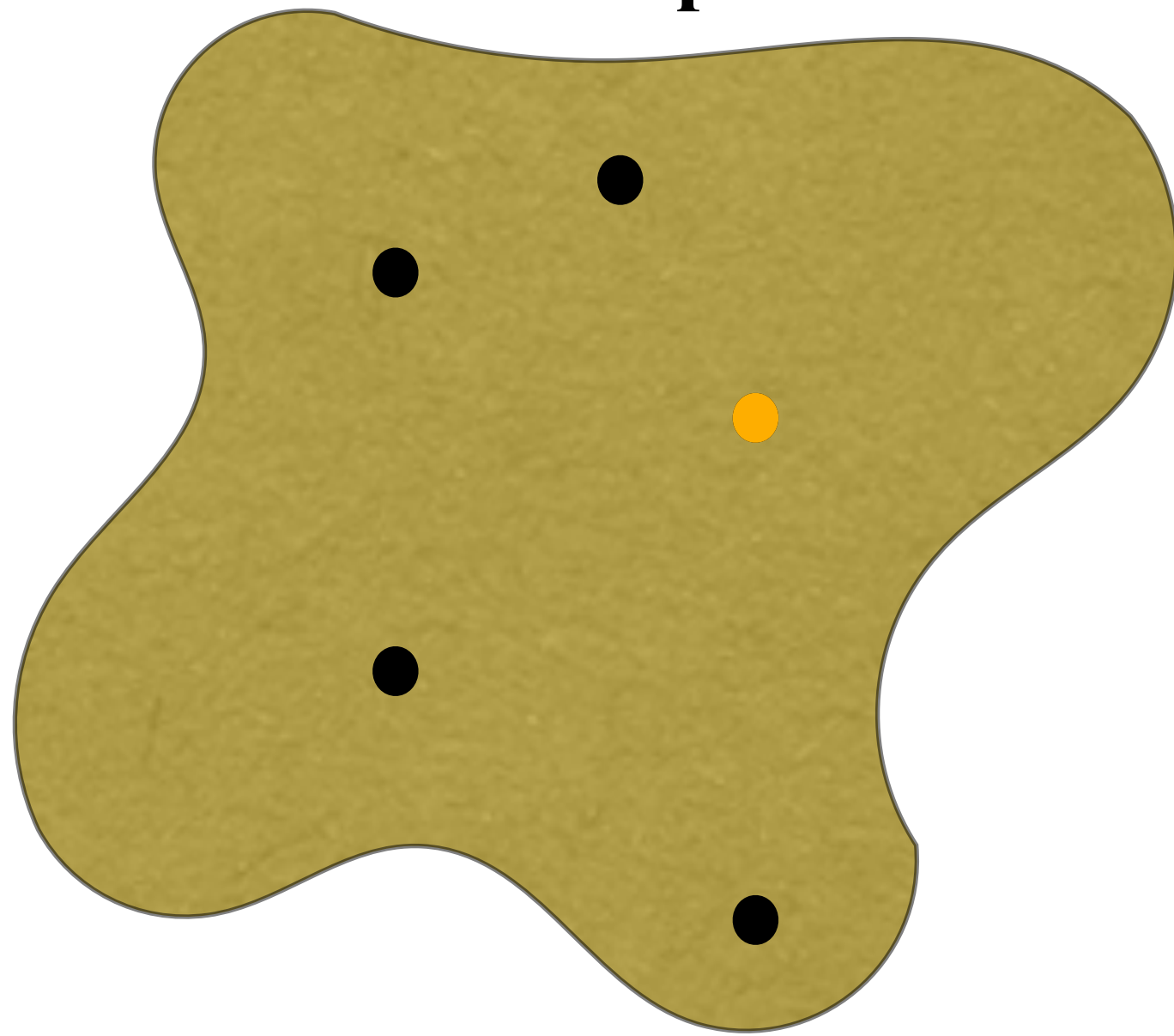
**First: dont cause harm!**

# Question: what constitutes a good collaboration?

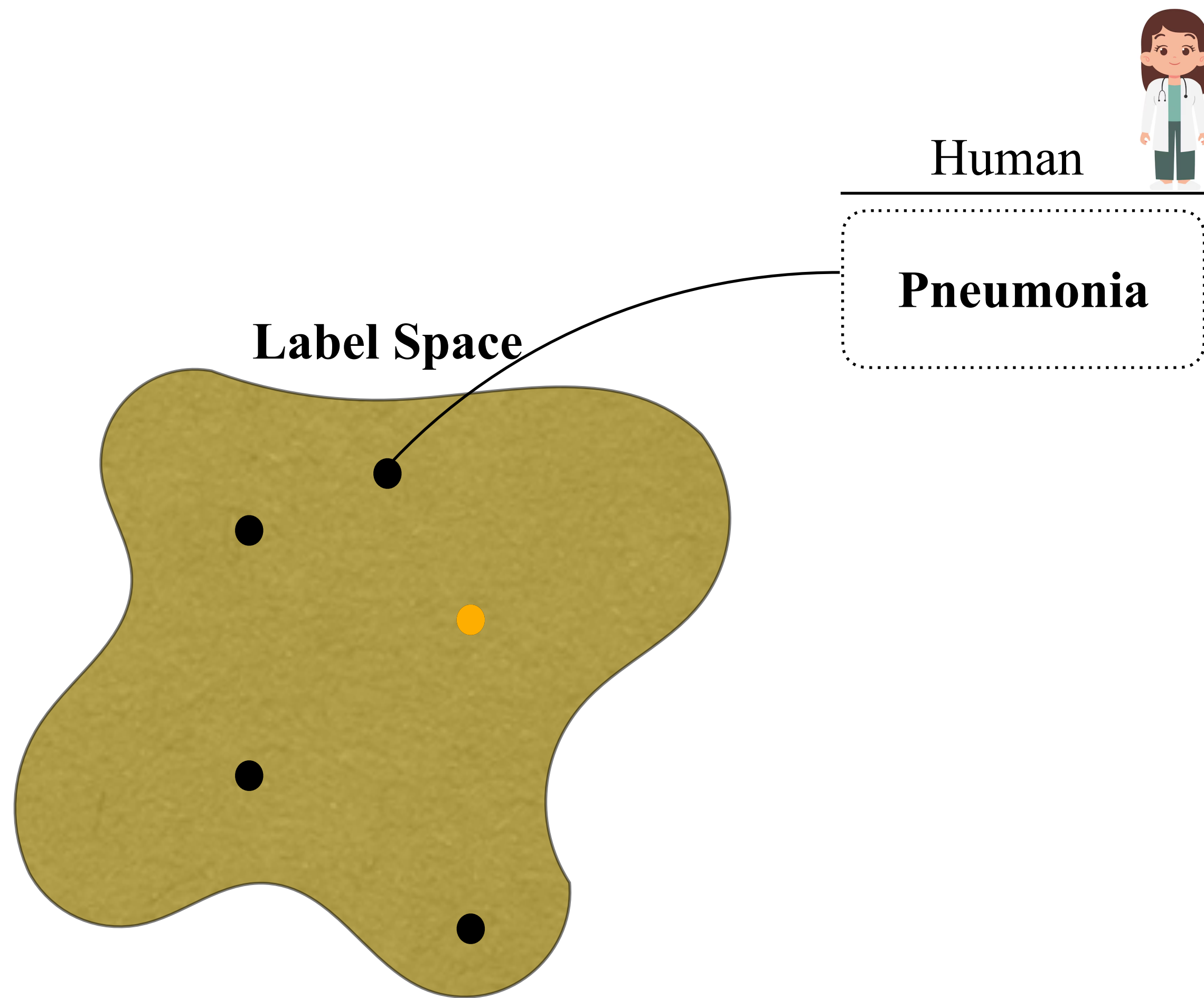
Human



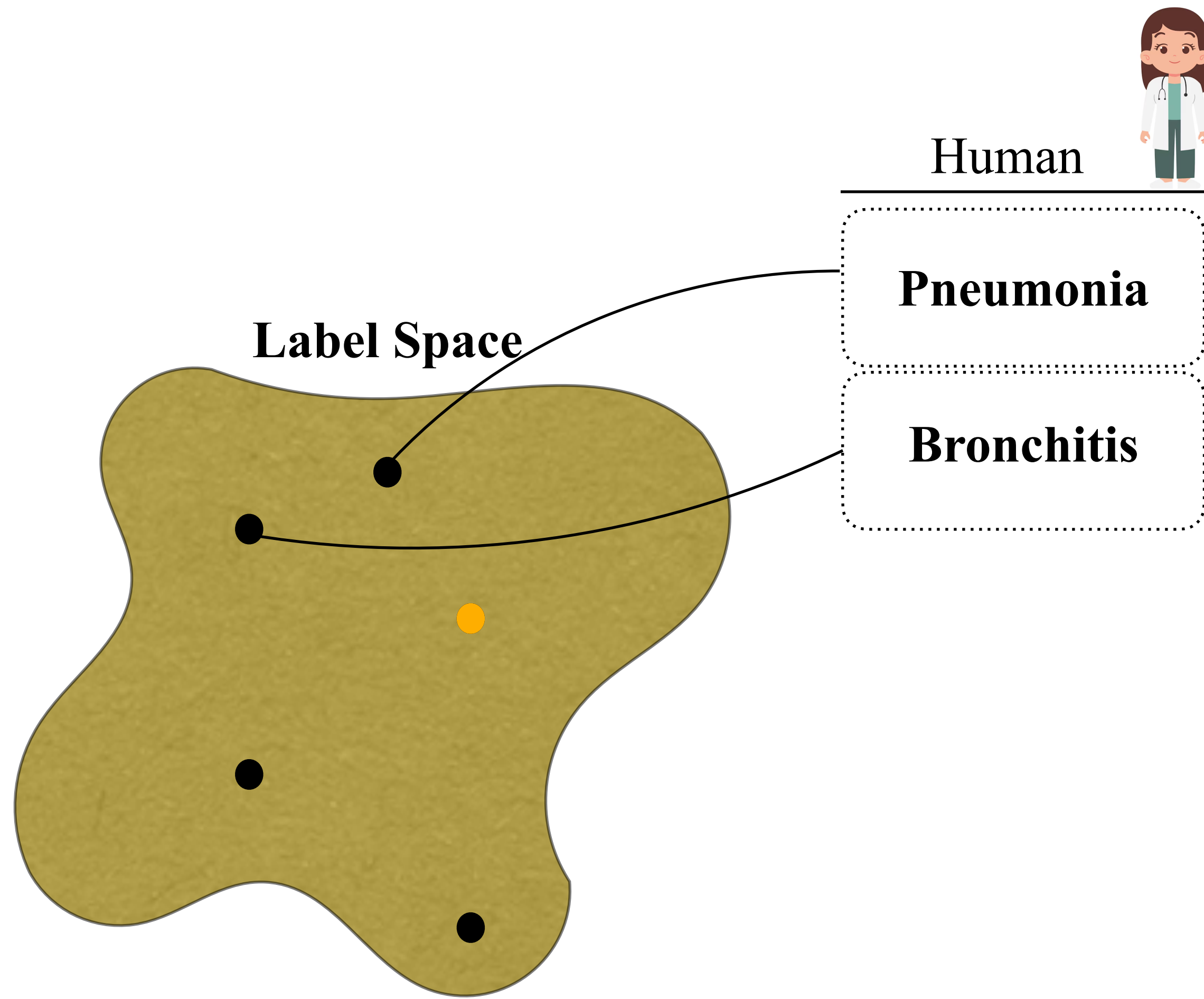
Label Space



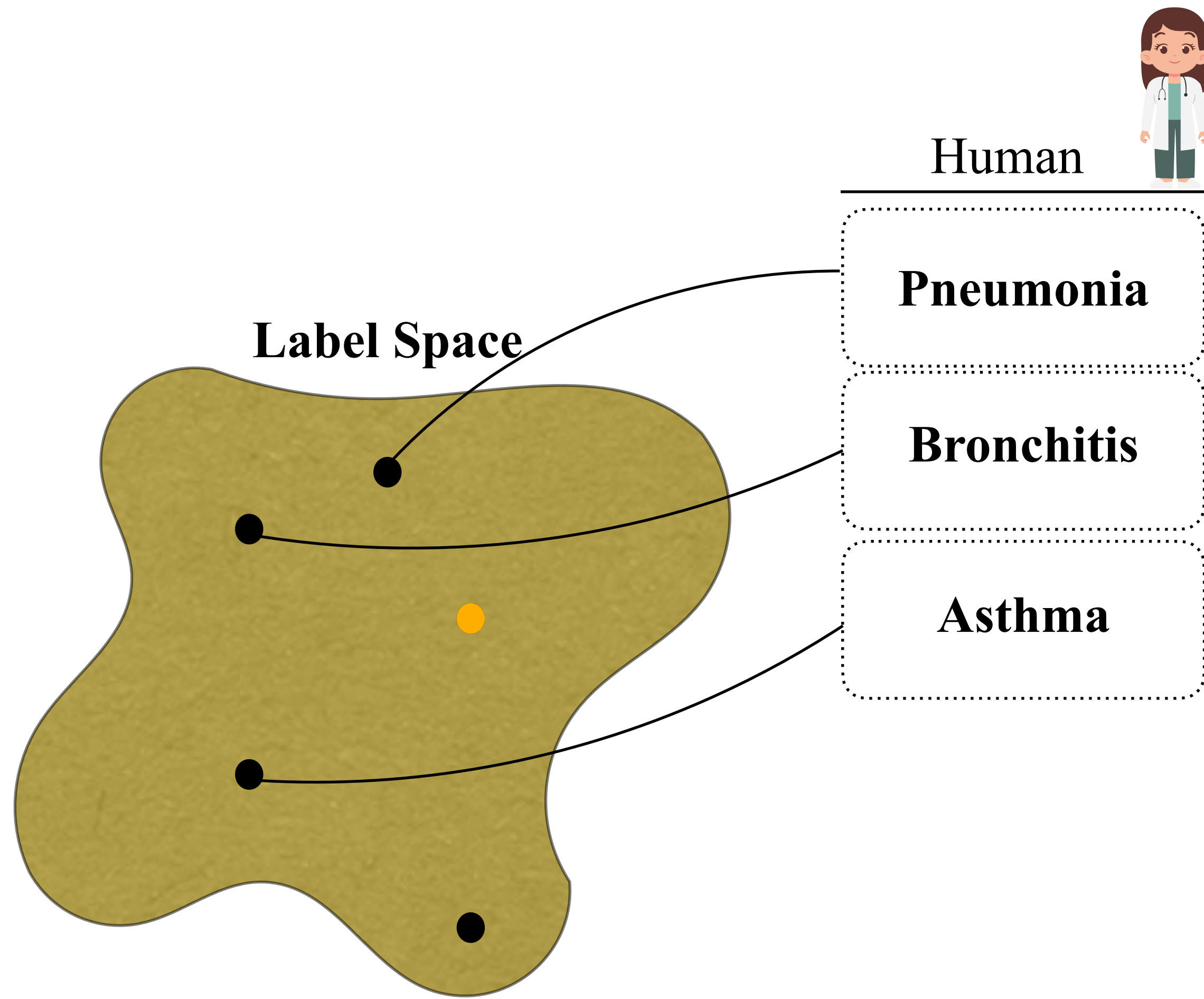
# Question: what constitutes a good collaboration?



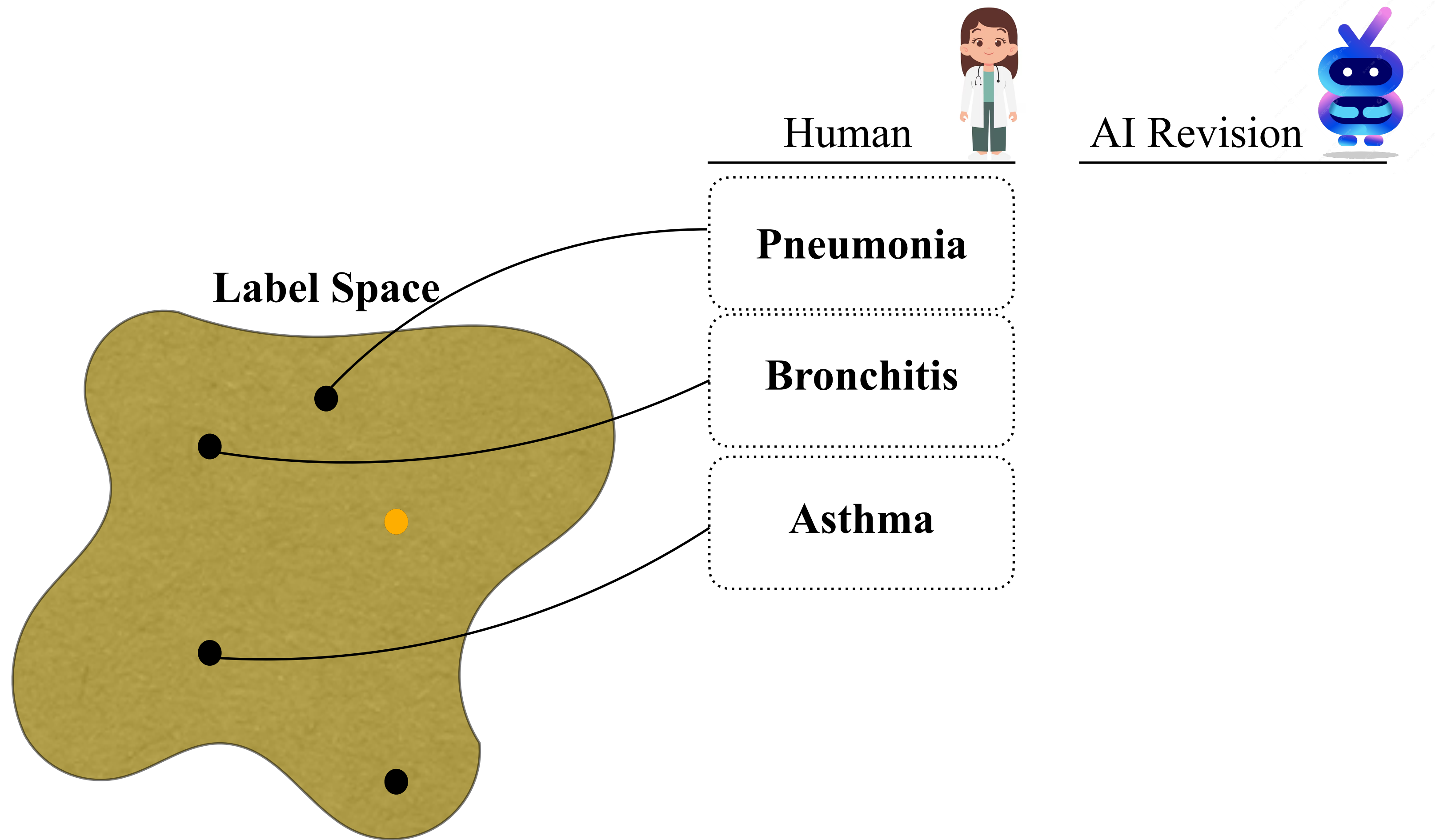
# Question: what constitutes a good collaboration?



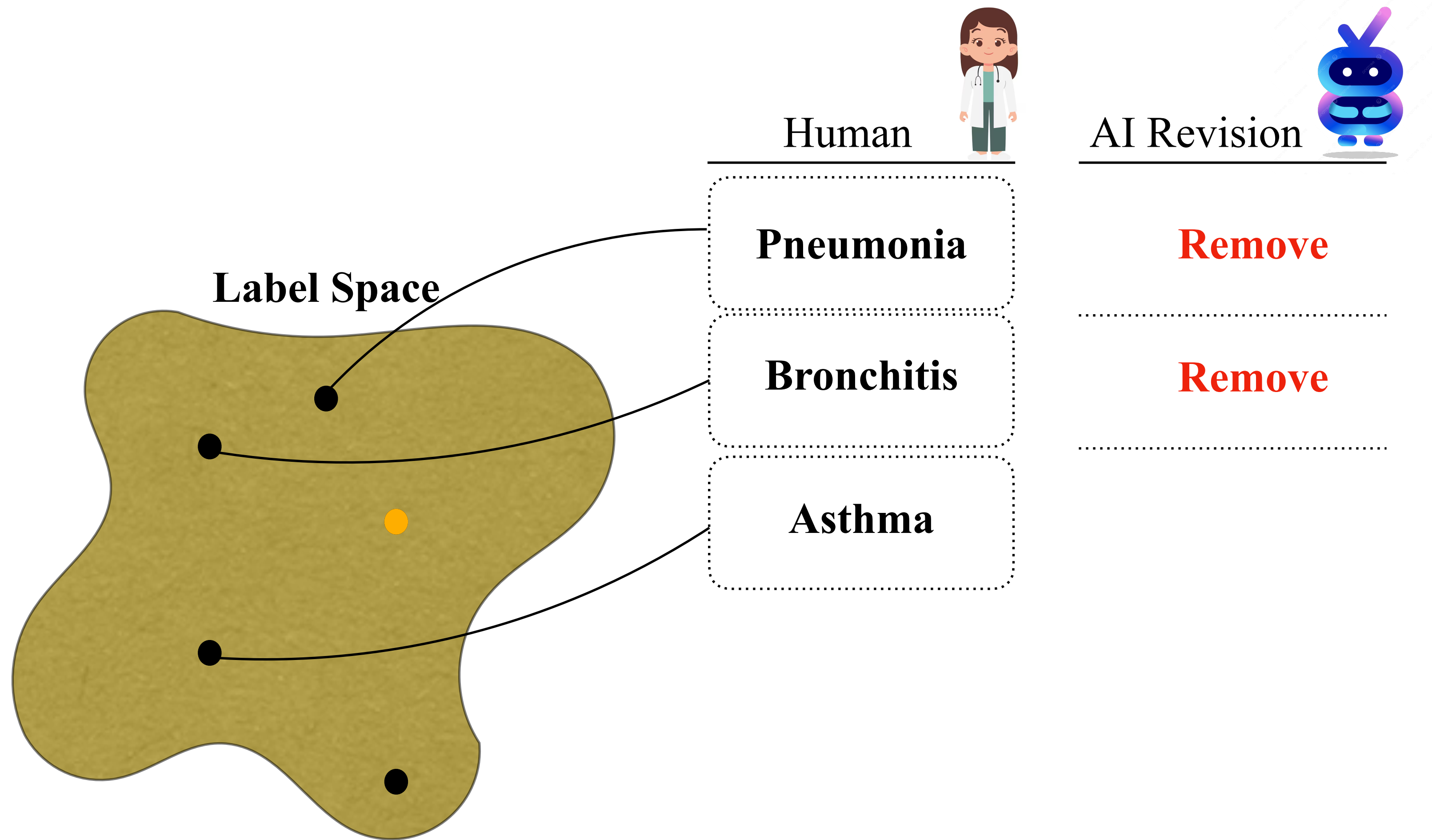
# Question: what constitutes a good collaboration?



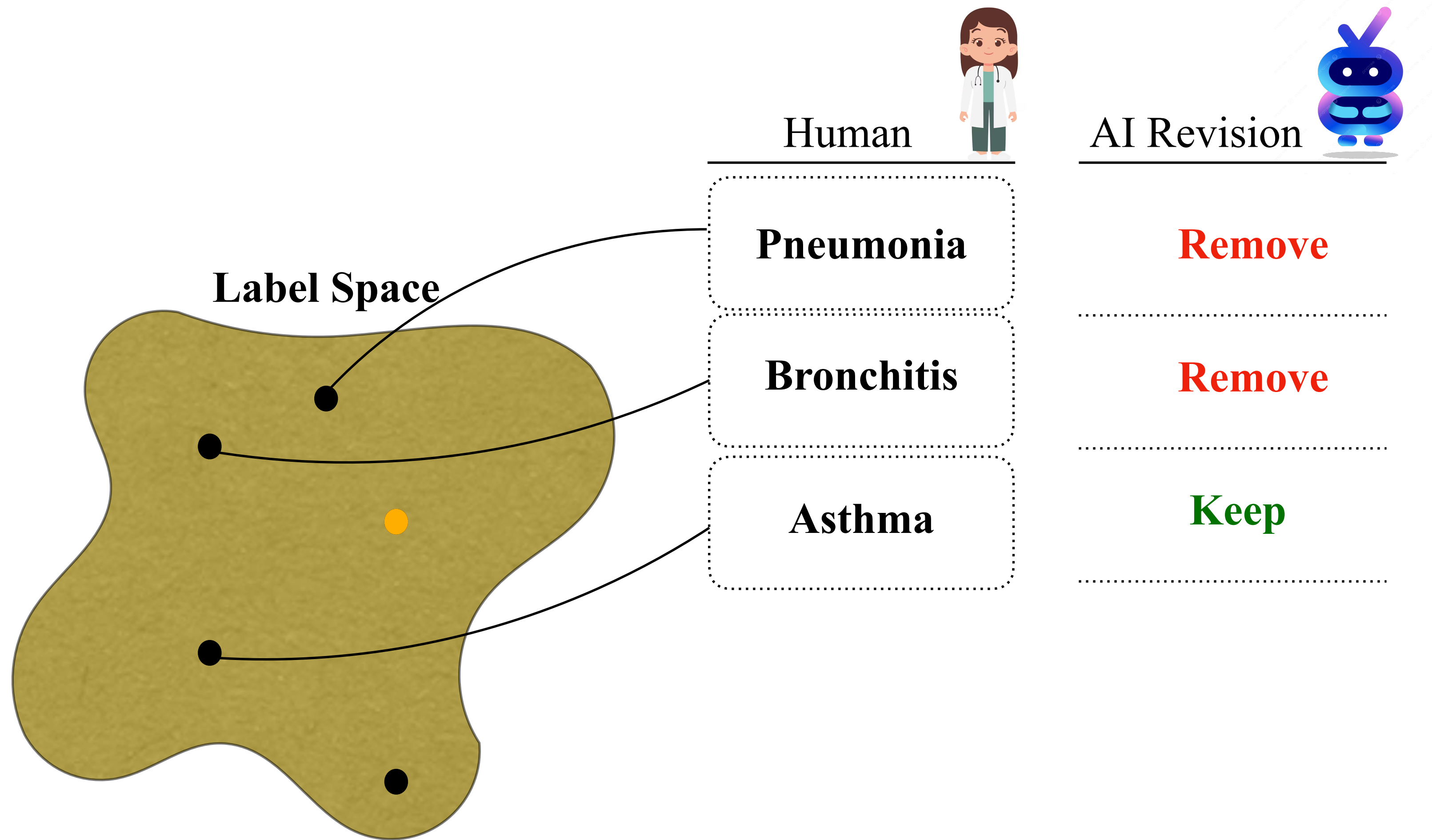
# Question: what constitutes a good collaboration?



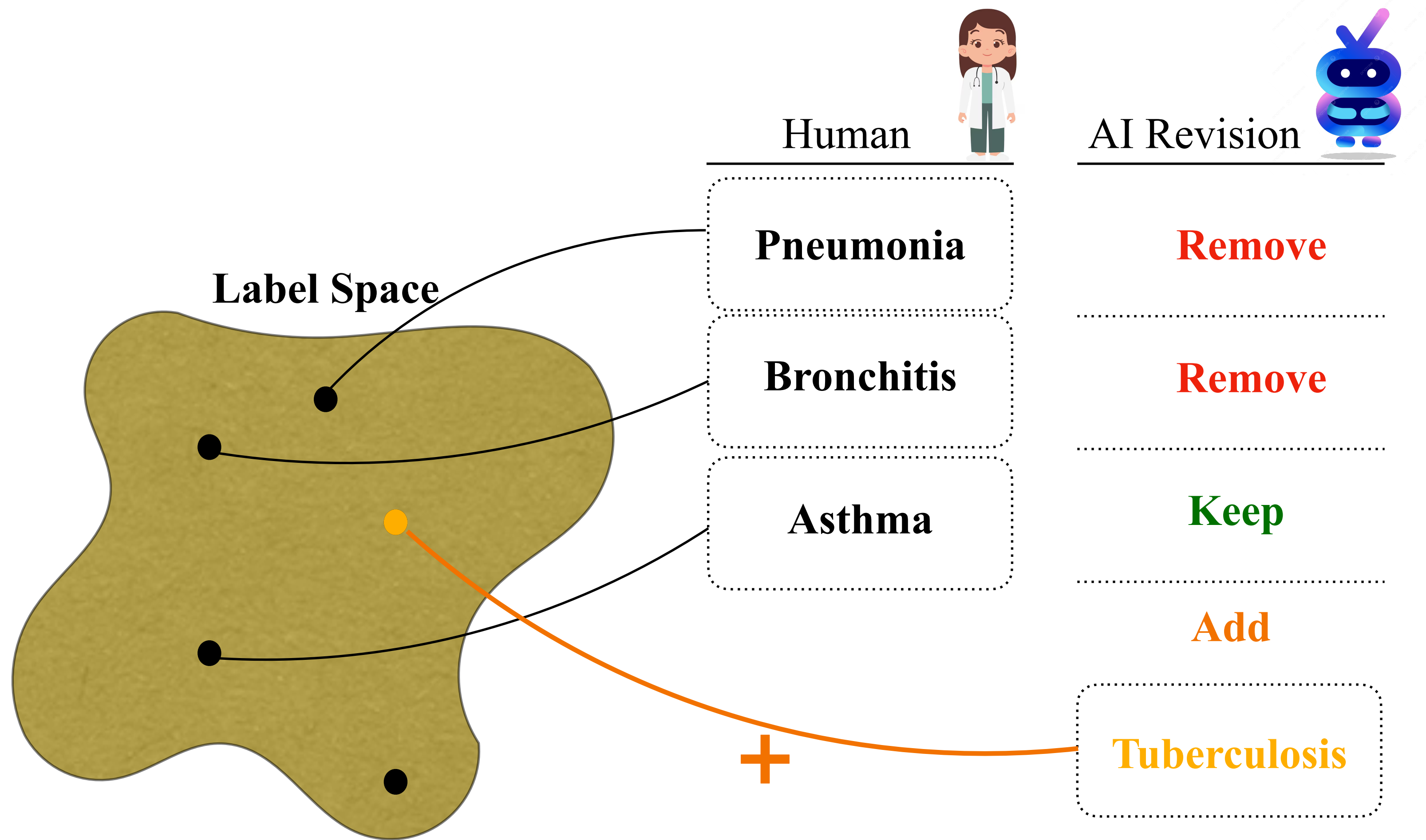
# Question: what constitutes a good collaboration?



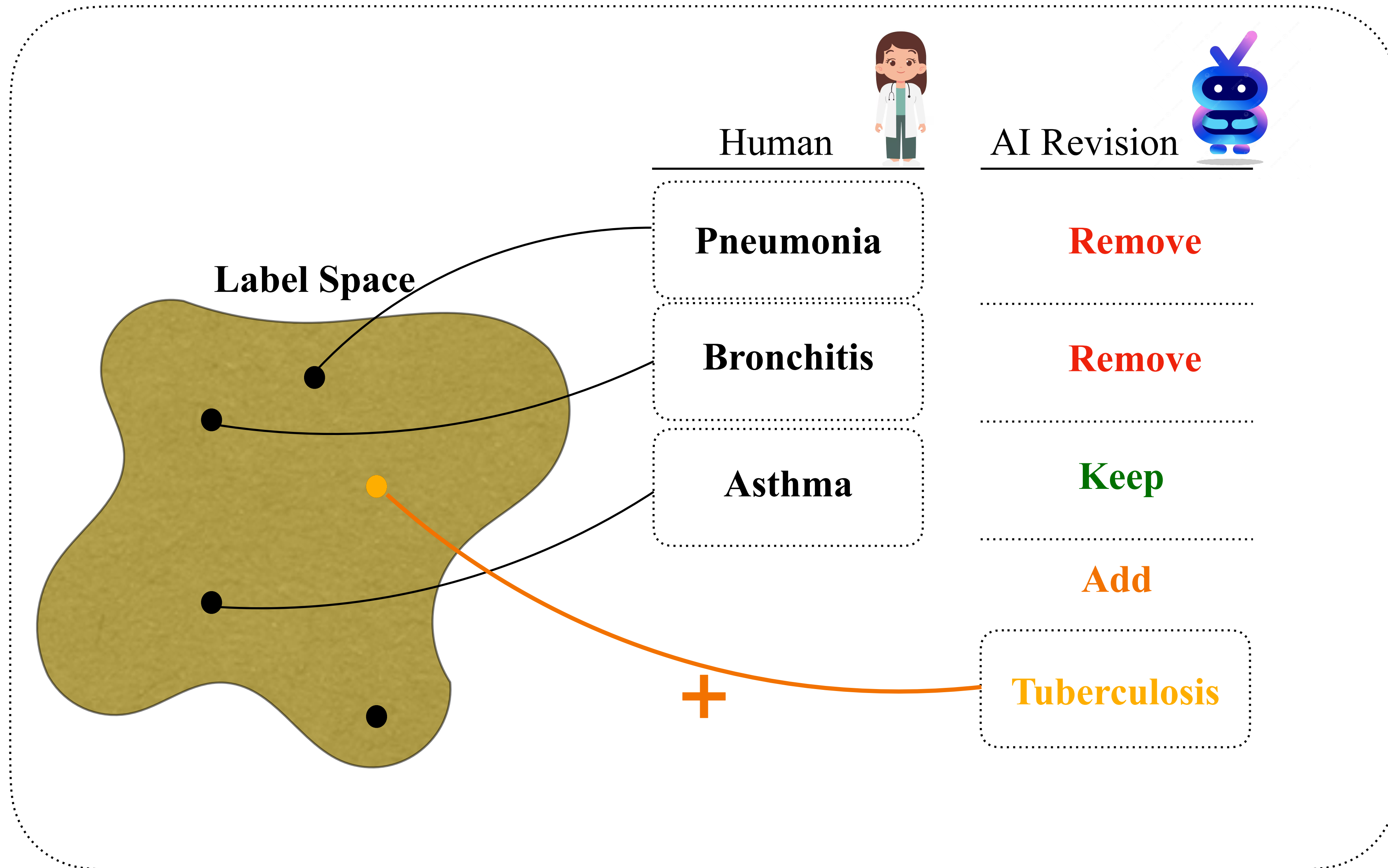
# Question: what constitutes a good collaboration?



# Question: what constitutes a good collaboration?

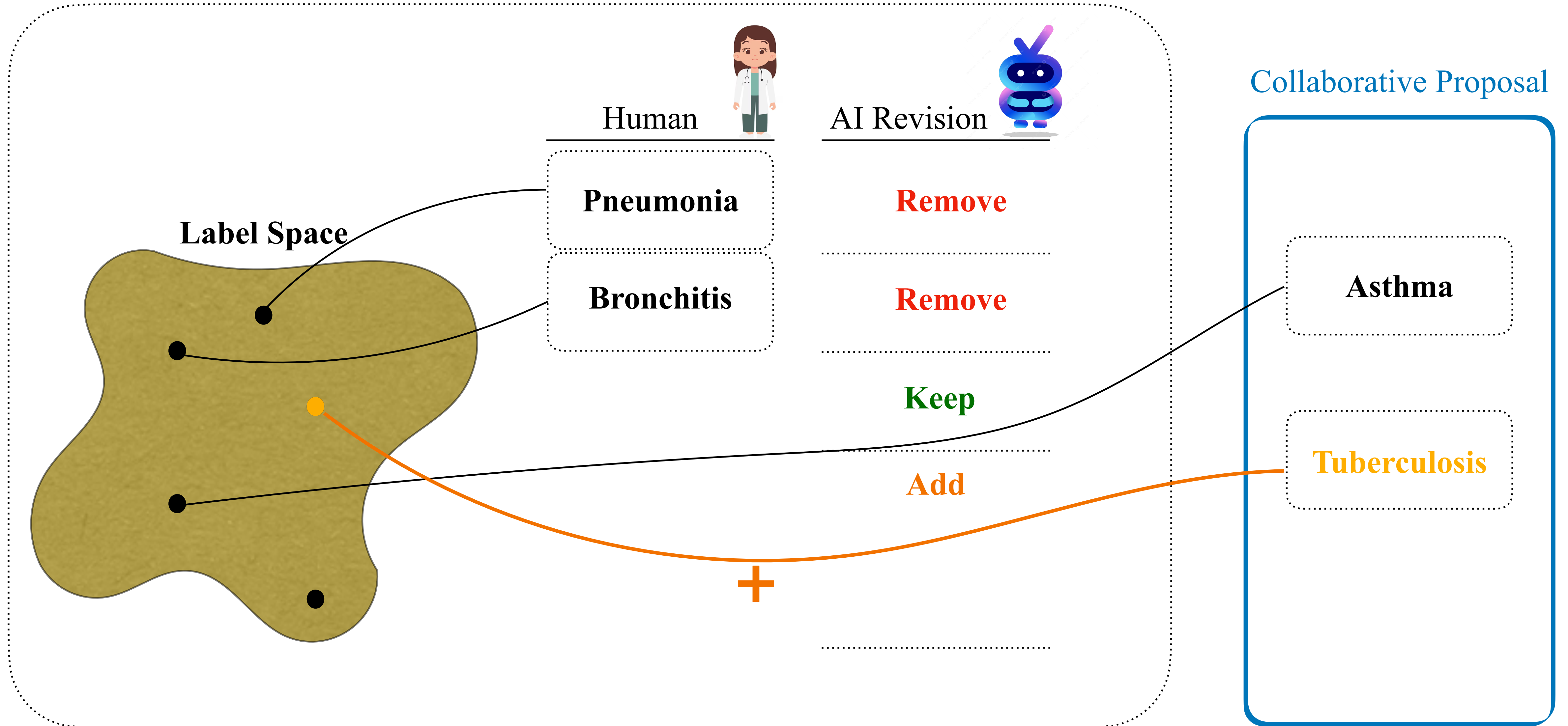


# Question: what constitutes a good collaboration?

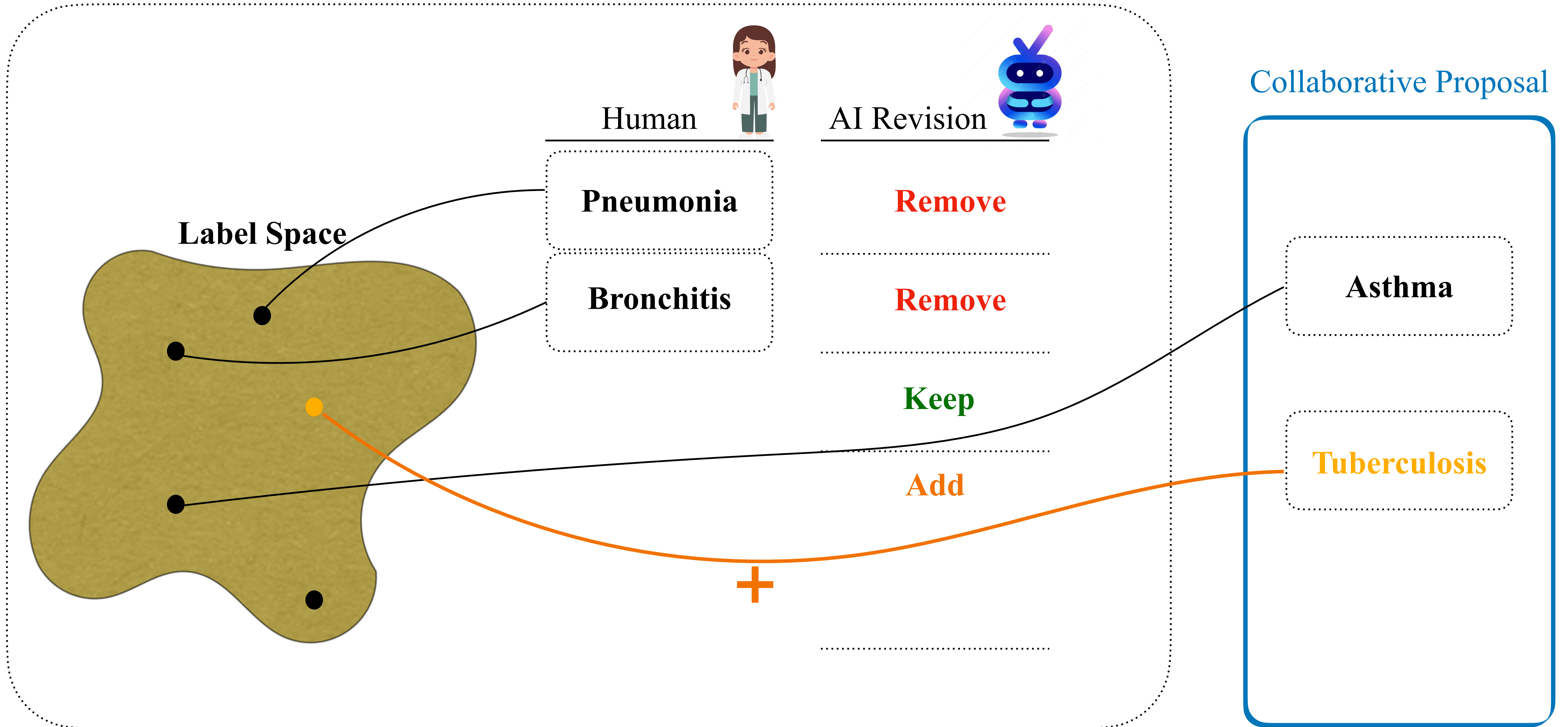


Collaborative Proposal

# Question: what constitutes a good collaboration?



# Question: what constitutes a good collaboration?



**Second: Add value!**

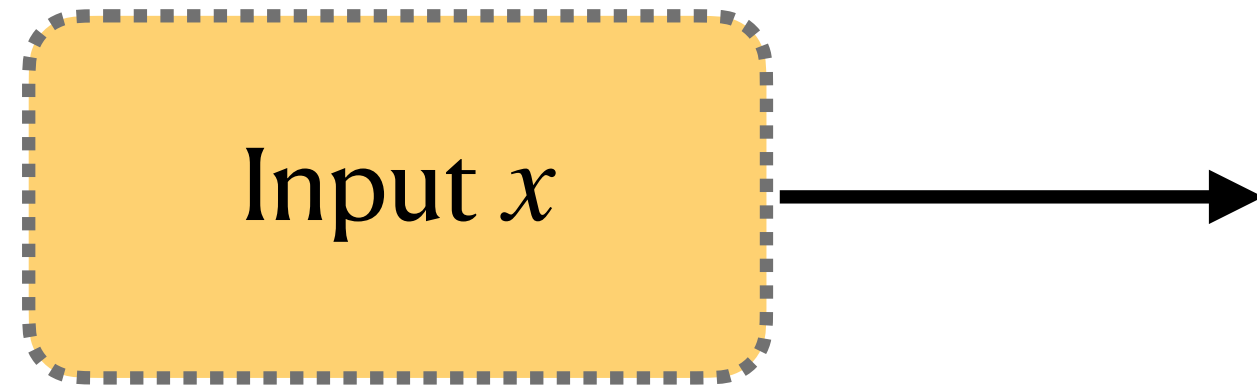
# **Two fundamentals of collaboration**

# Two fundamentals of collaboration

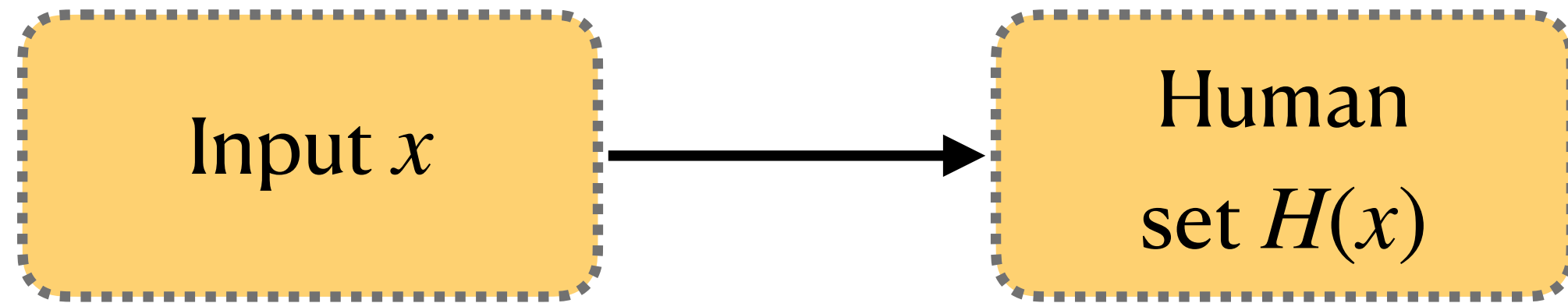


Input  $x$

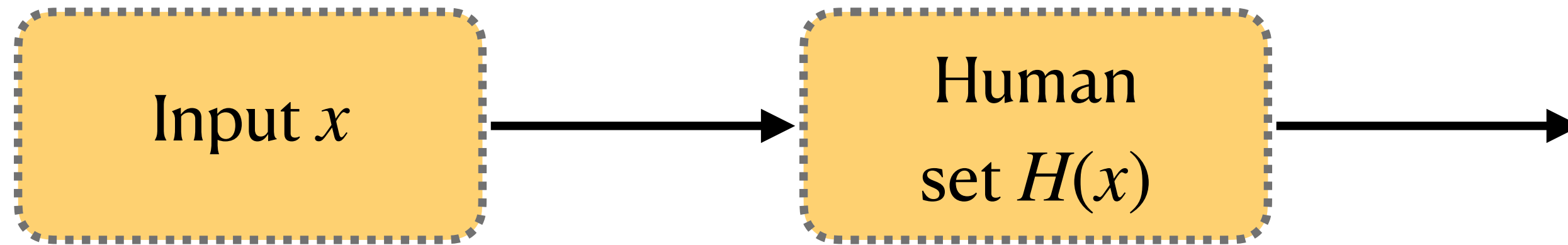
# Two fundamentals of collaboration



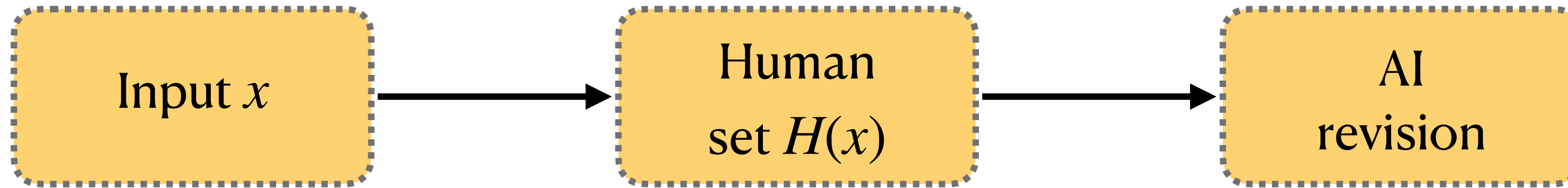
## Two fundamentals of collaboration



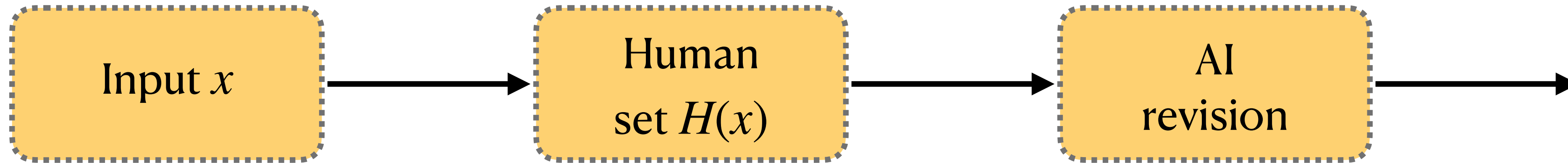
# Two fundamentals of collaboration



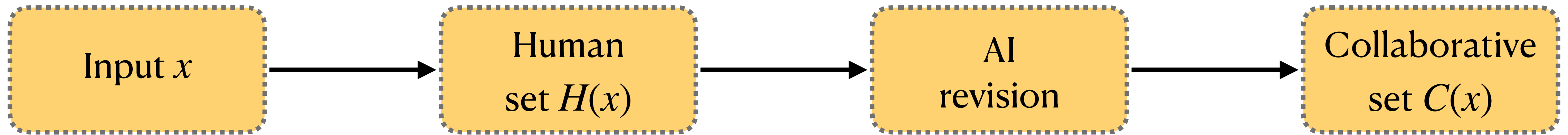
## Two fundamentals of collaboration



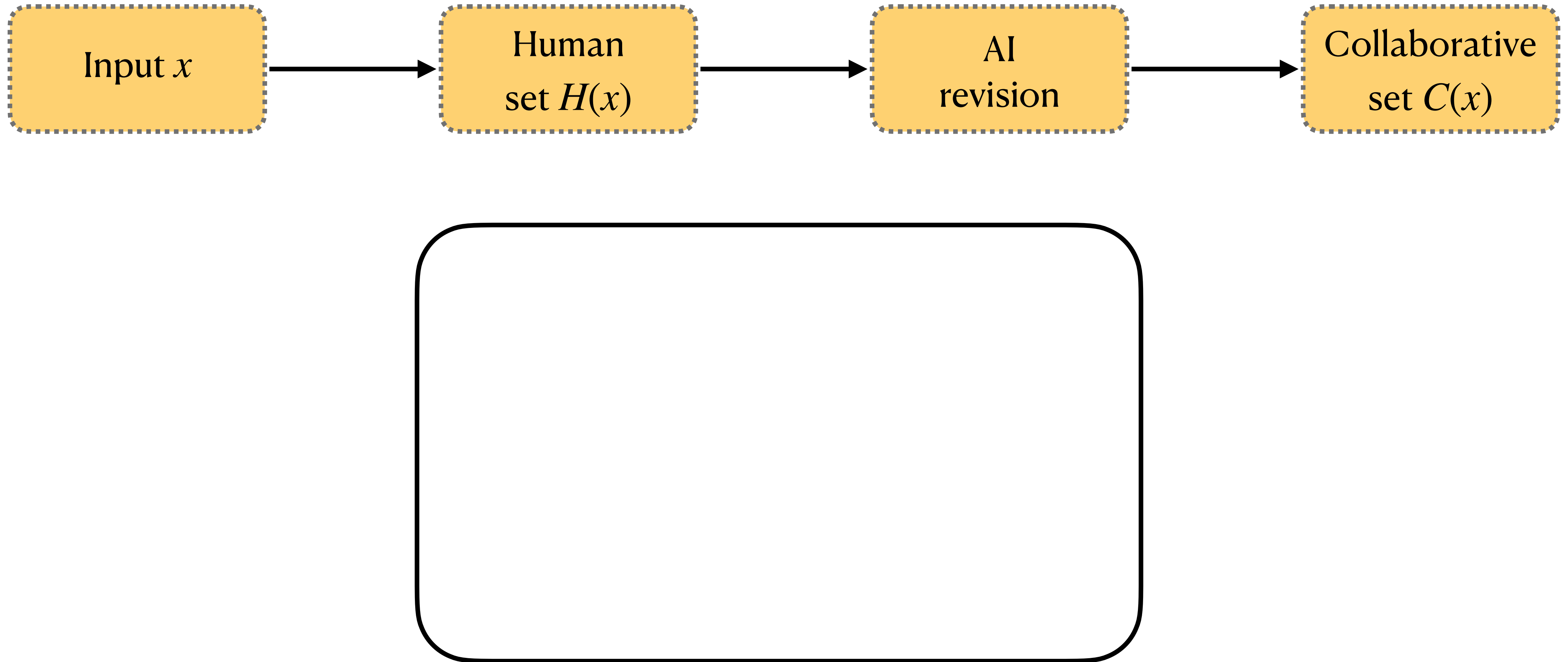
## Two fundamentals of collaboration



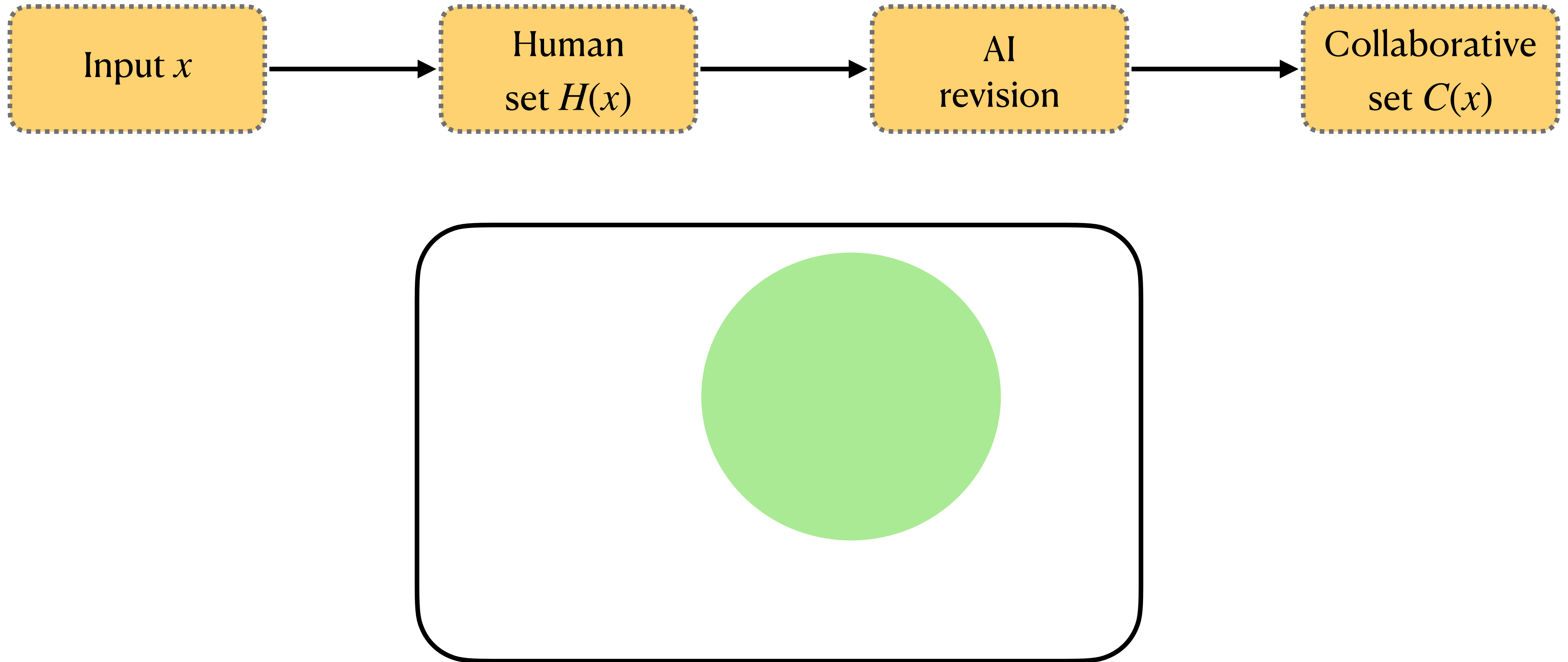
## Two fundamentals of collaboration



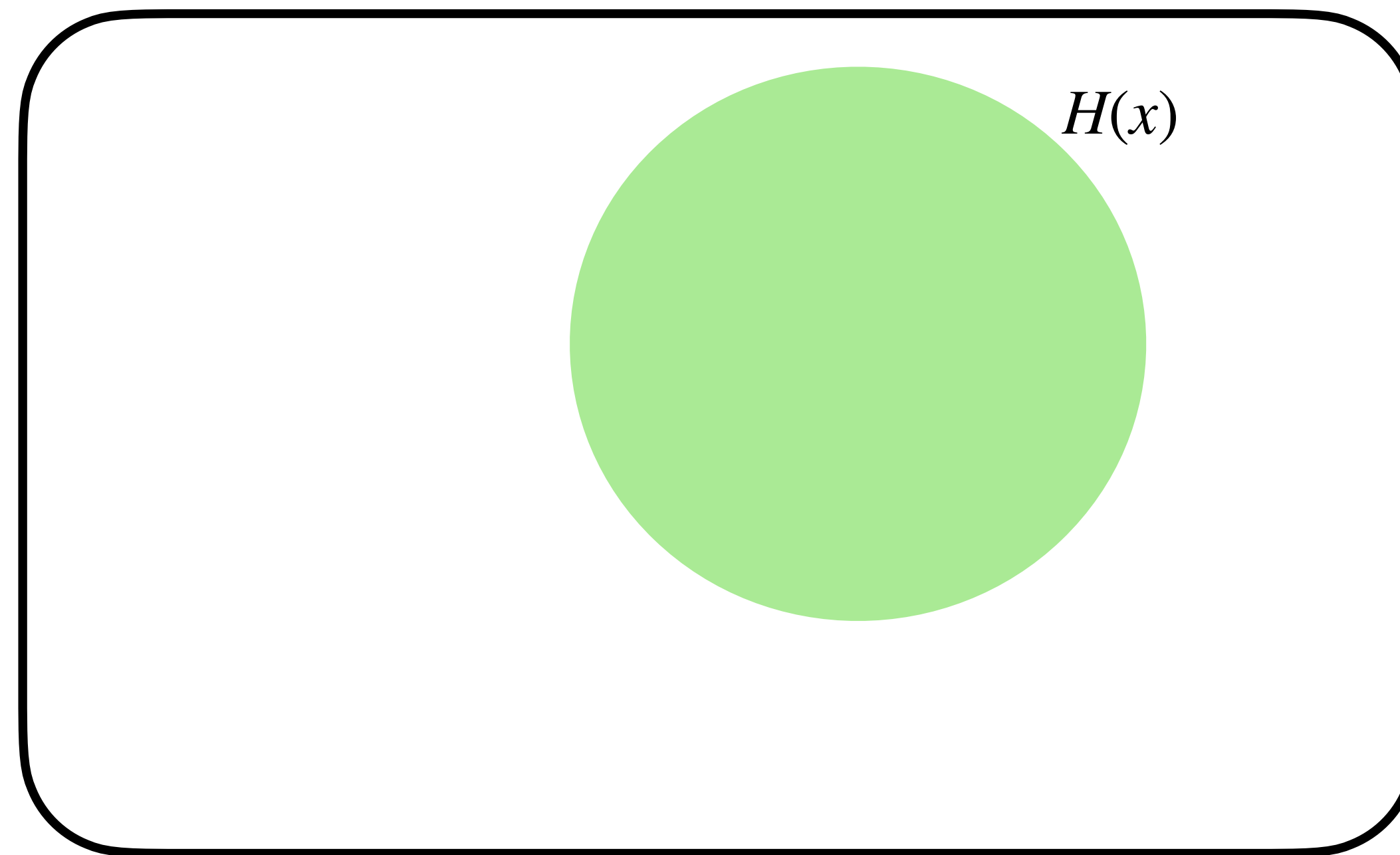
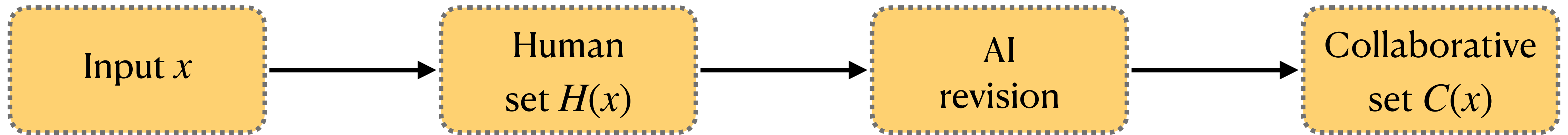
# Two fundamentals of collaboration



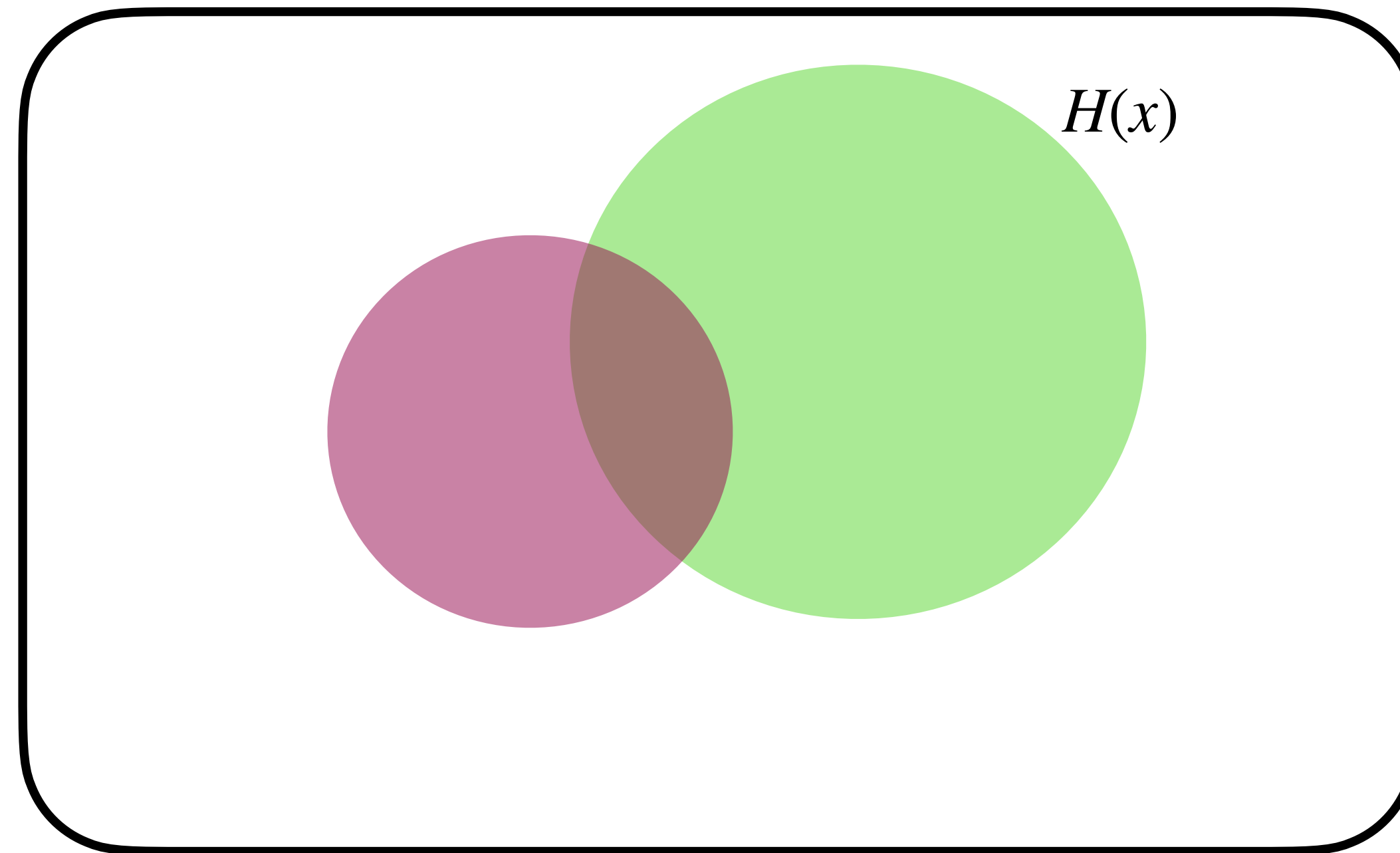
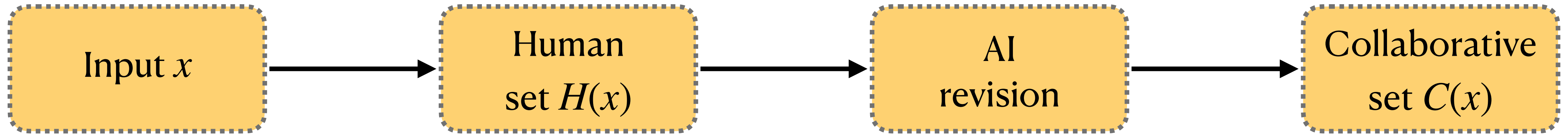
# Two fundamentals of collaboration



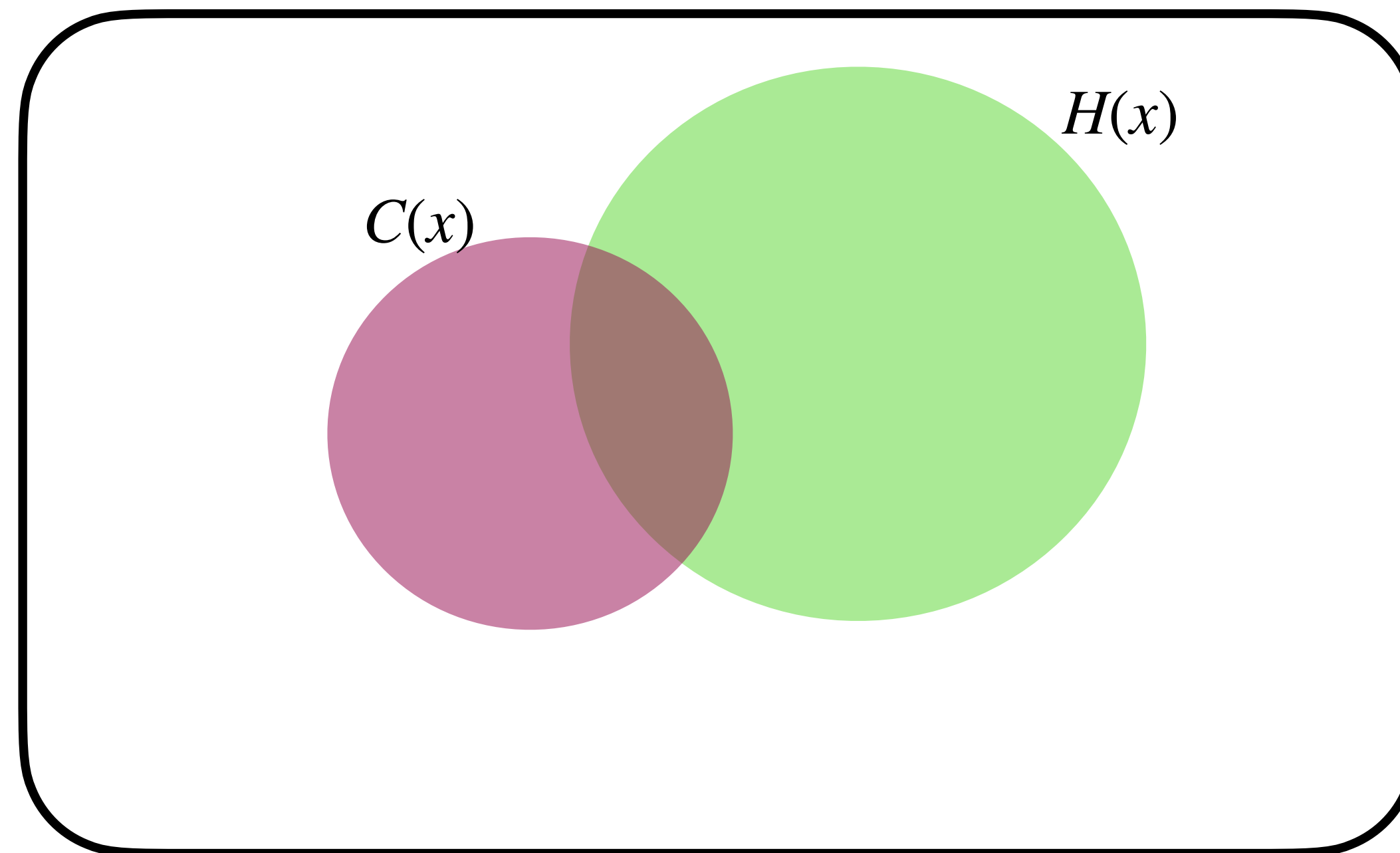
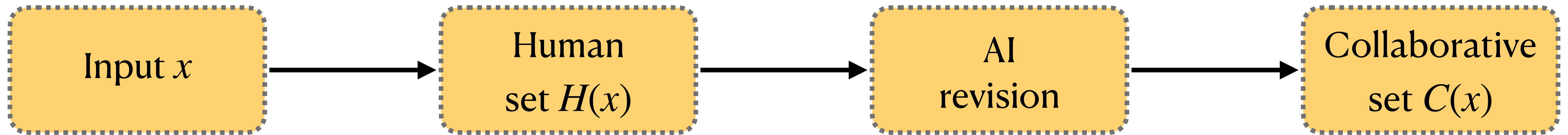
# Two fundamentals of collaboration



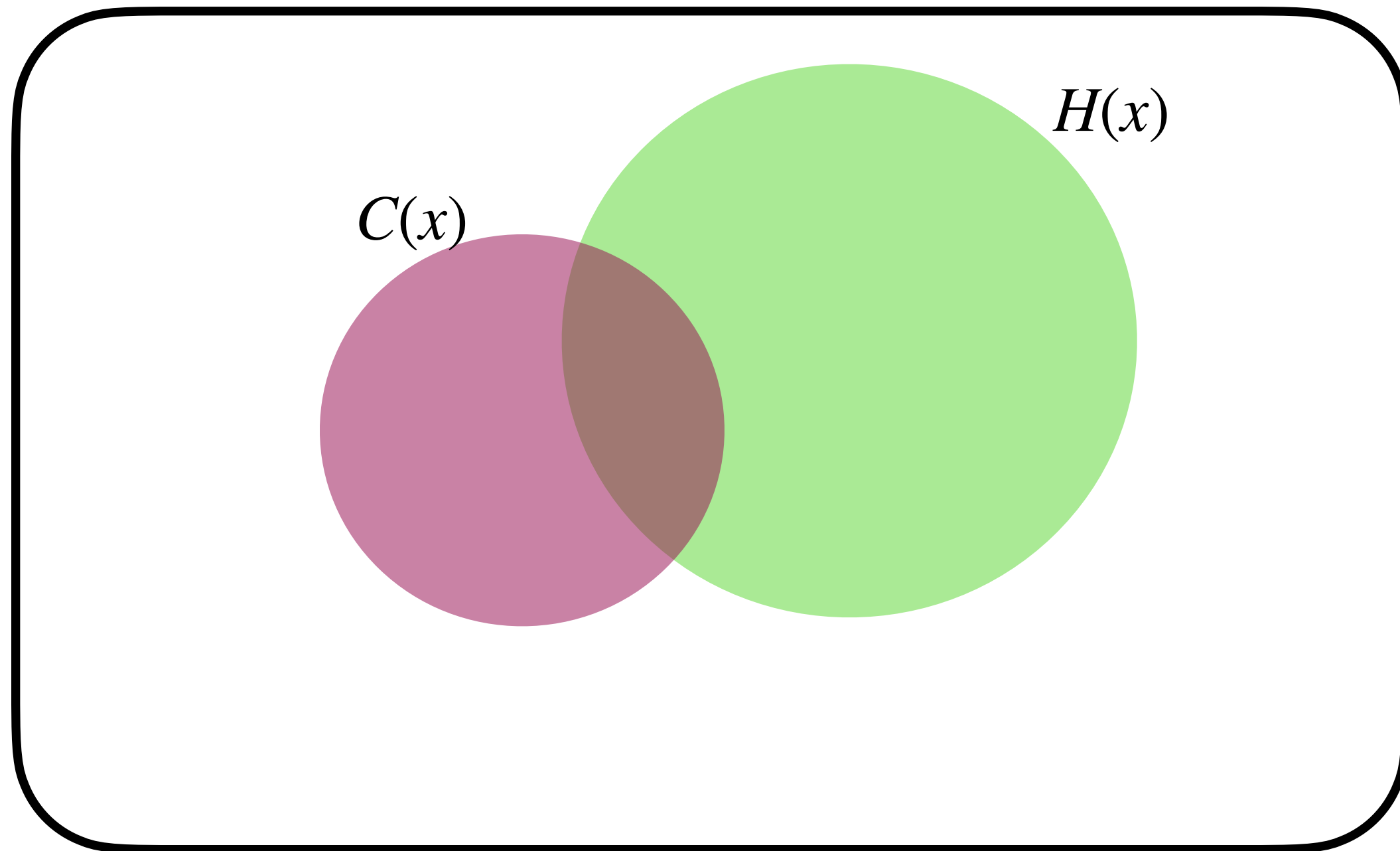
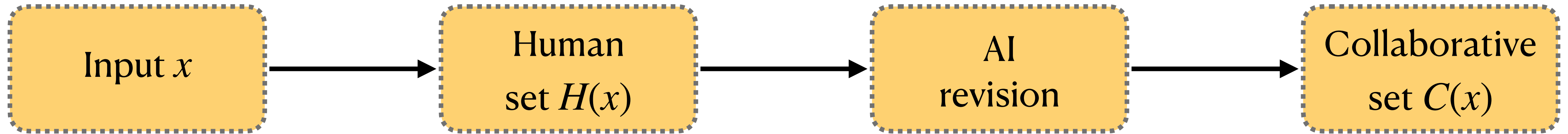
# Two fundamentals of collaboration



# Two fundamentals of collaboration

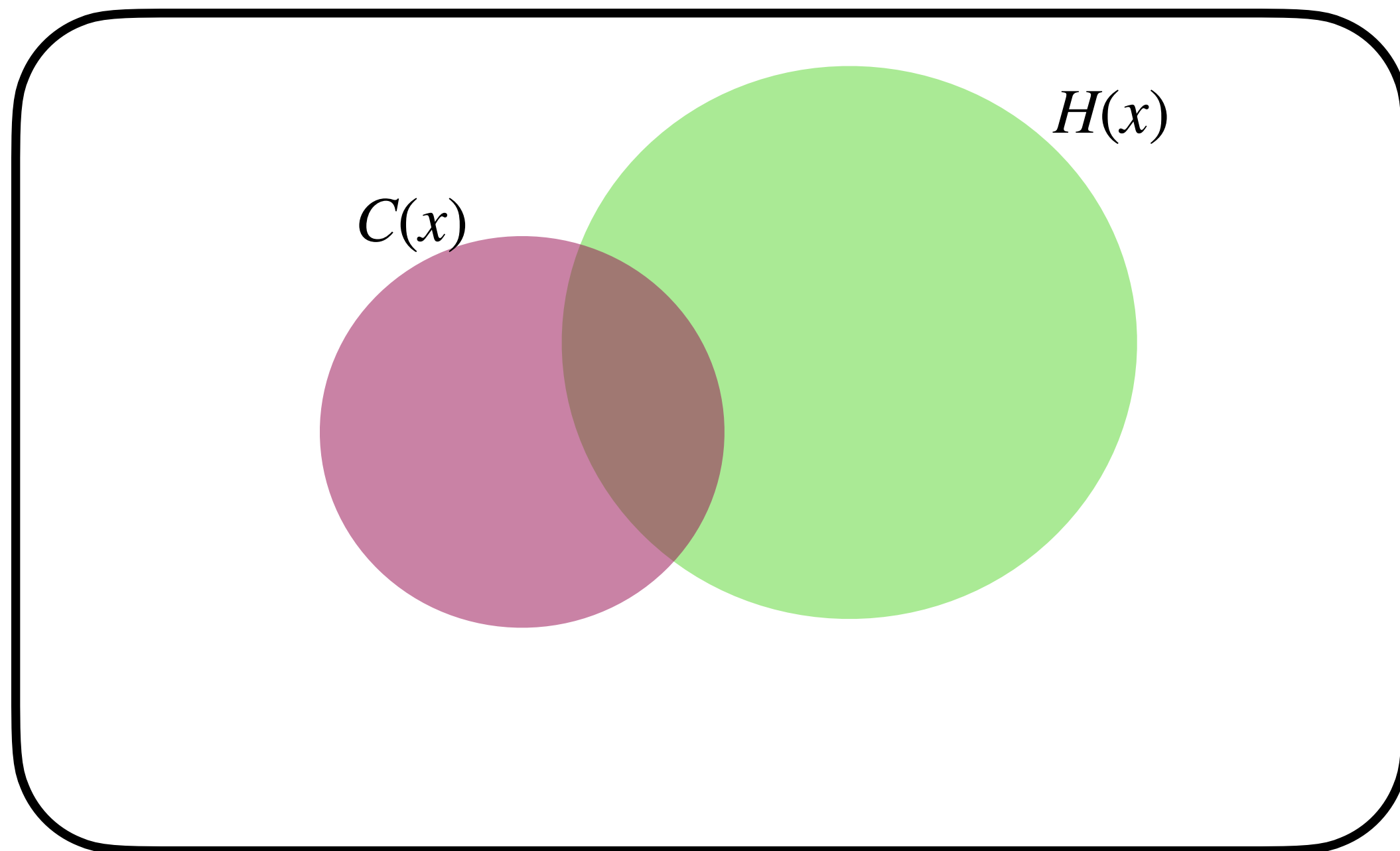
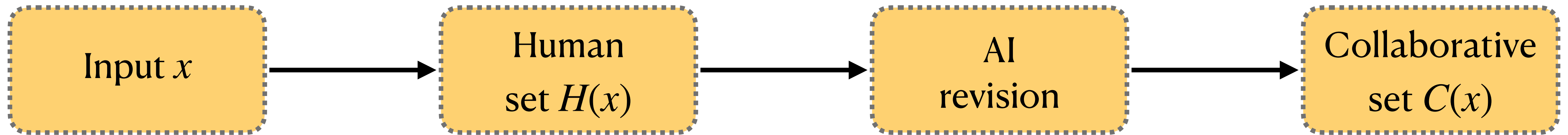


# Two fundamentals of collaboration



**Counterfactual Harm**

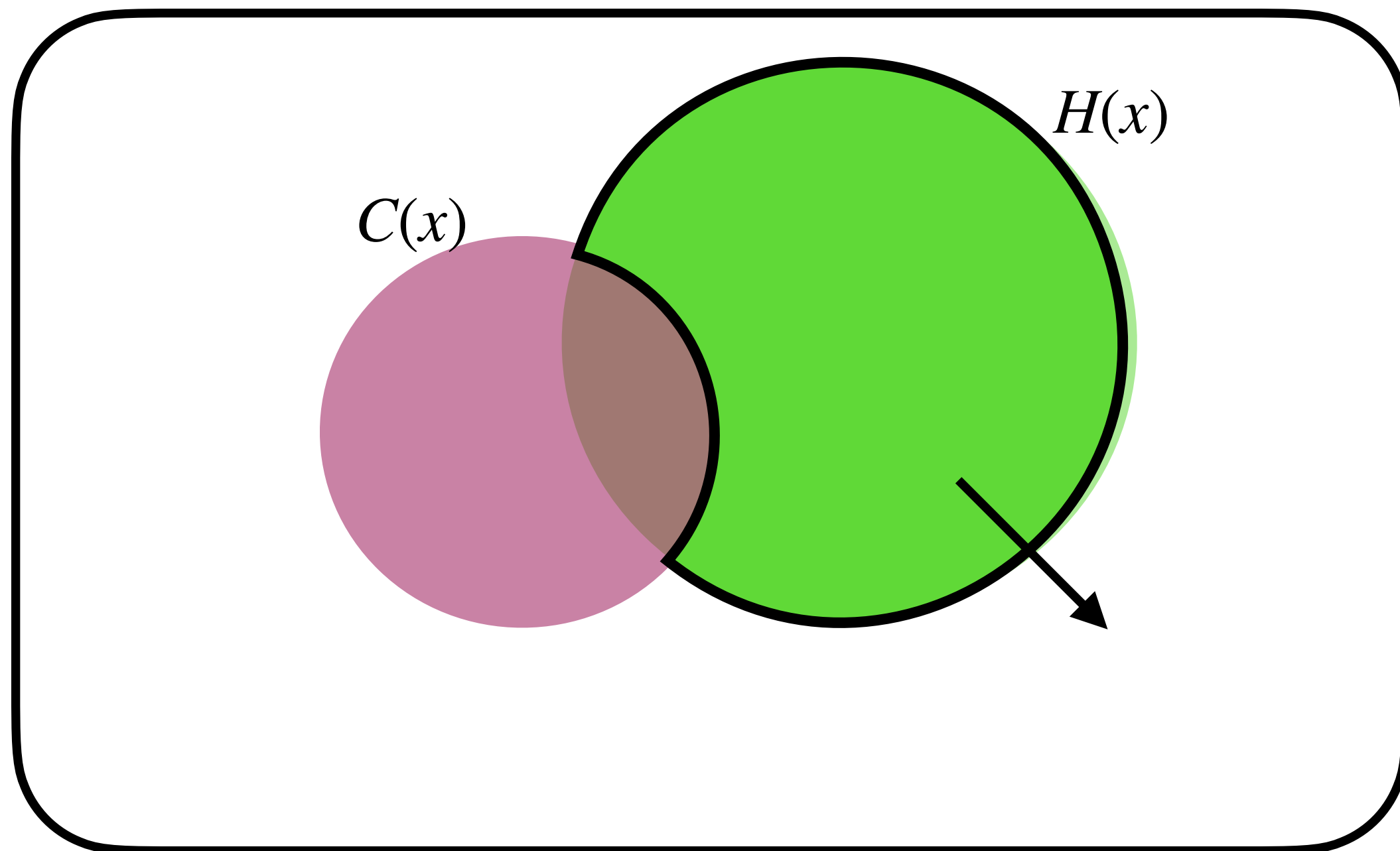
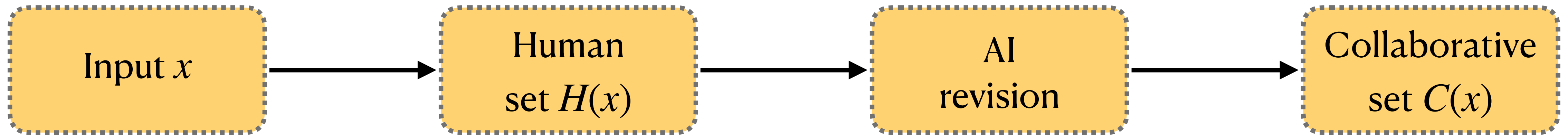
## Two fundamentals of collaboration



### Counterfactual Harm

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

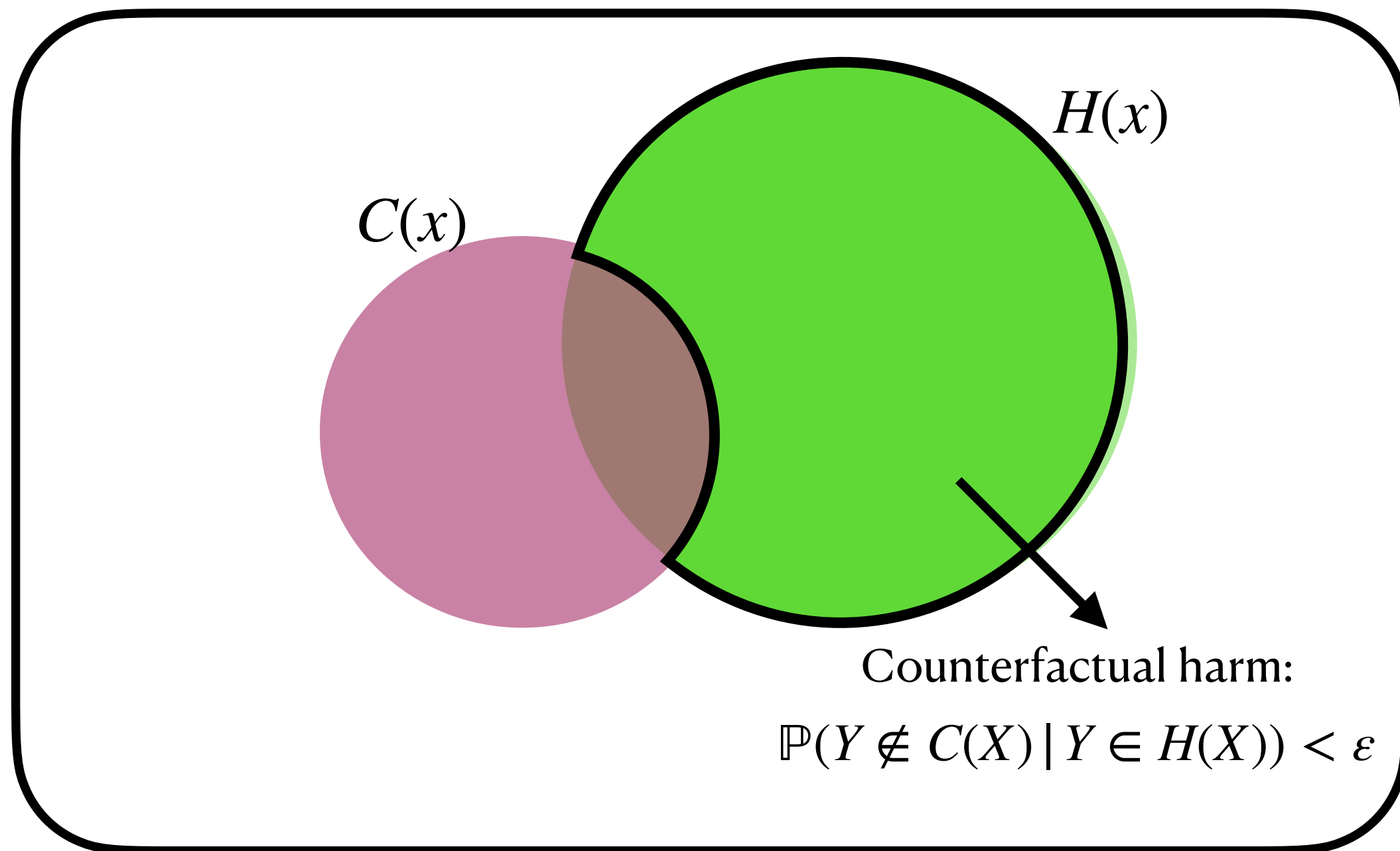
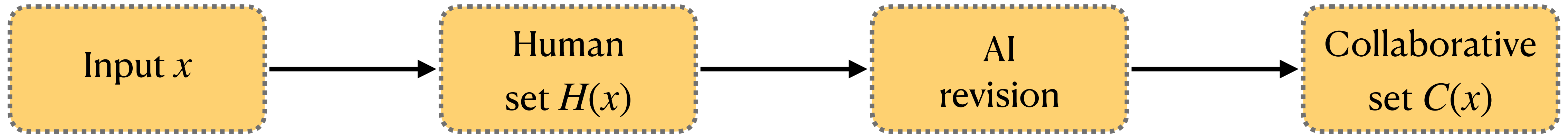
# Two fundamentals of collaboration



## Counterfactual Harm

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

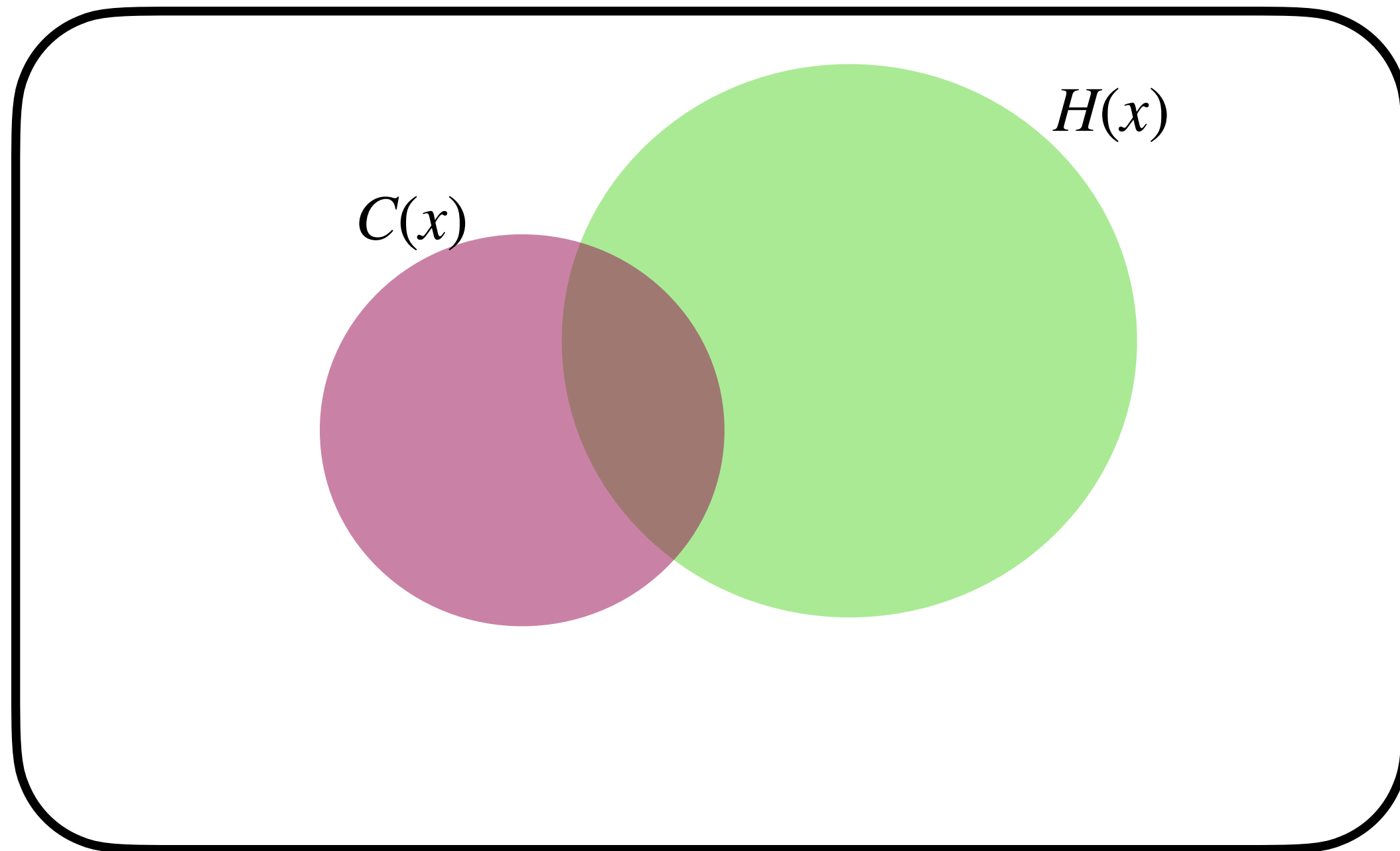
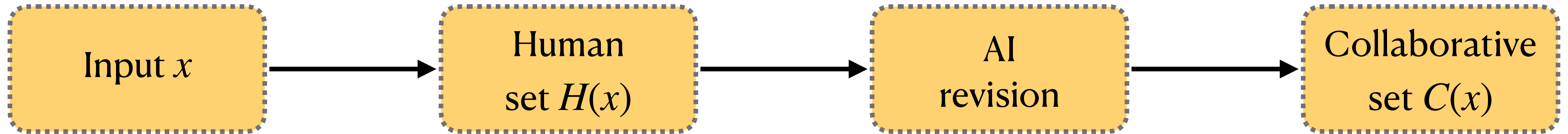
# Two fundamentals of collaboration



## Counterfactual Harm

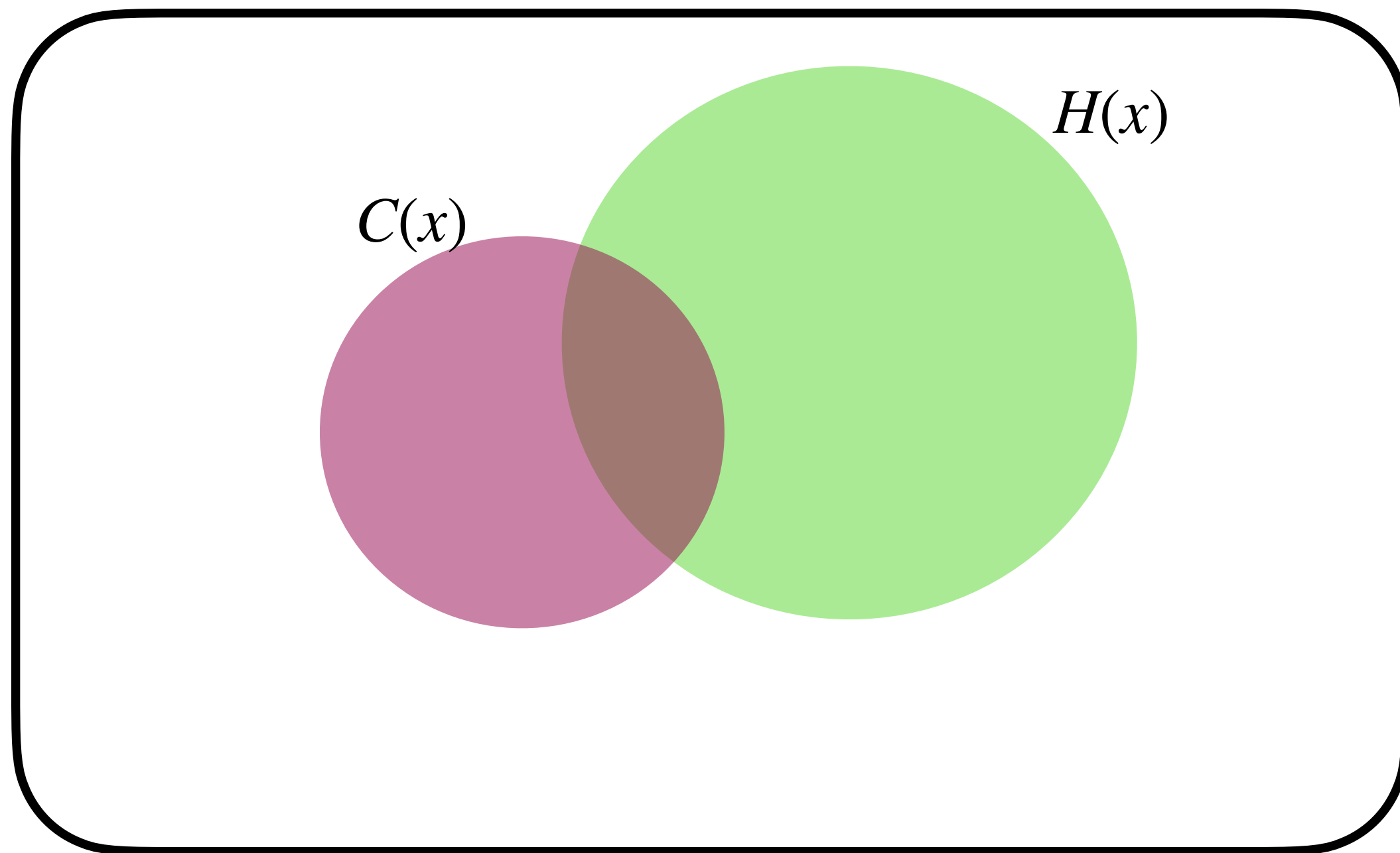
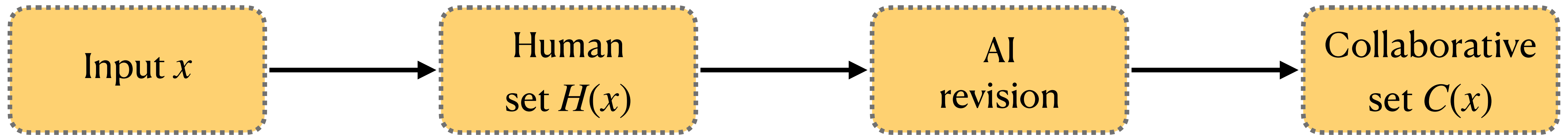
$$\mathbb{P}(Y \notin C(X) | Y \in H(X)) < \varepsilon$$

# Two fundamentals of collaboration



**Complementarity**

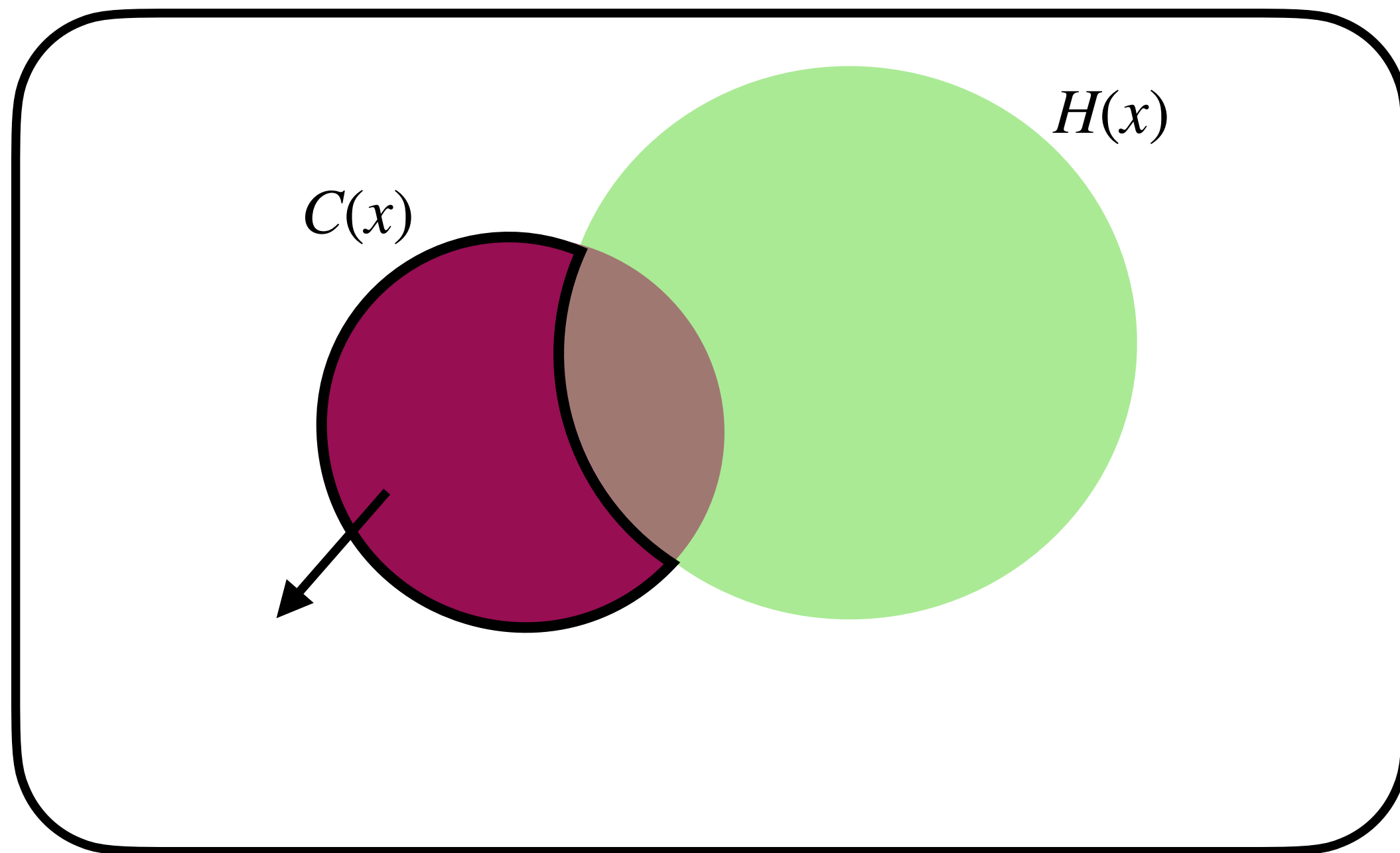
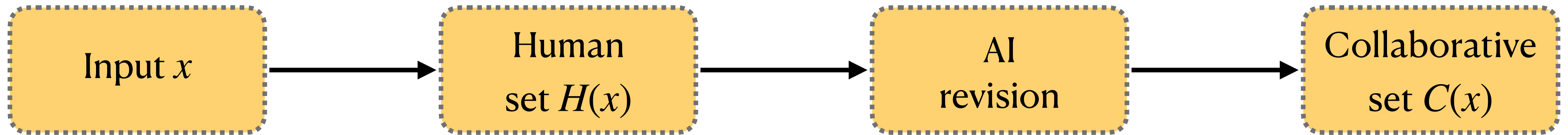
## Two fundamentals of collaboration



### Complementarity

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

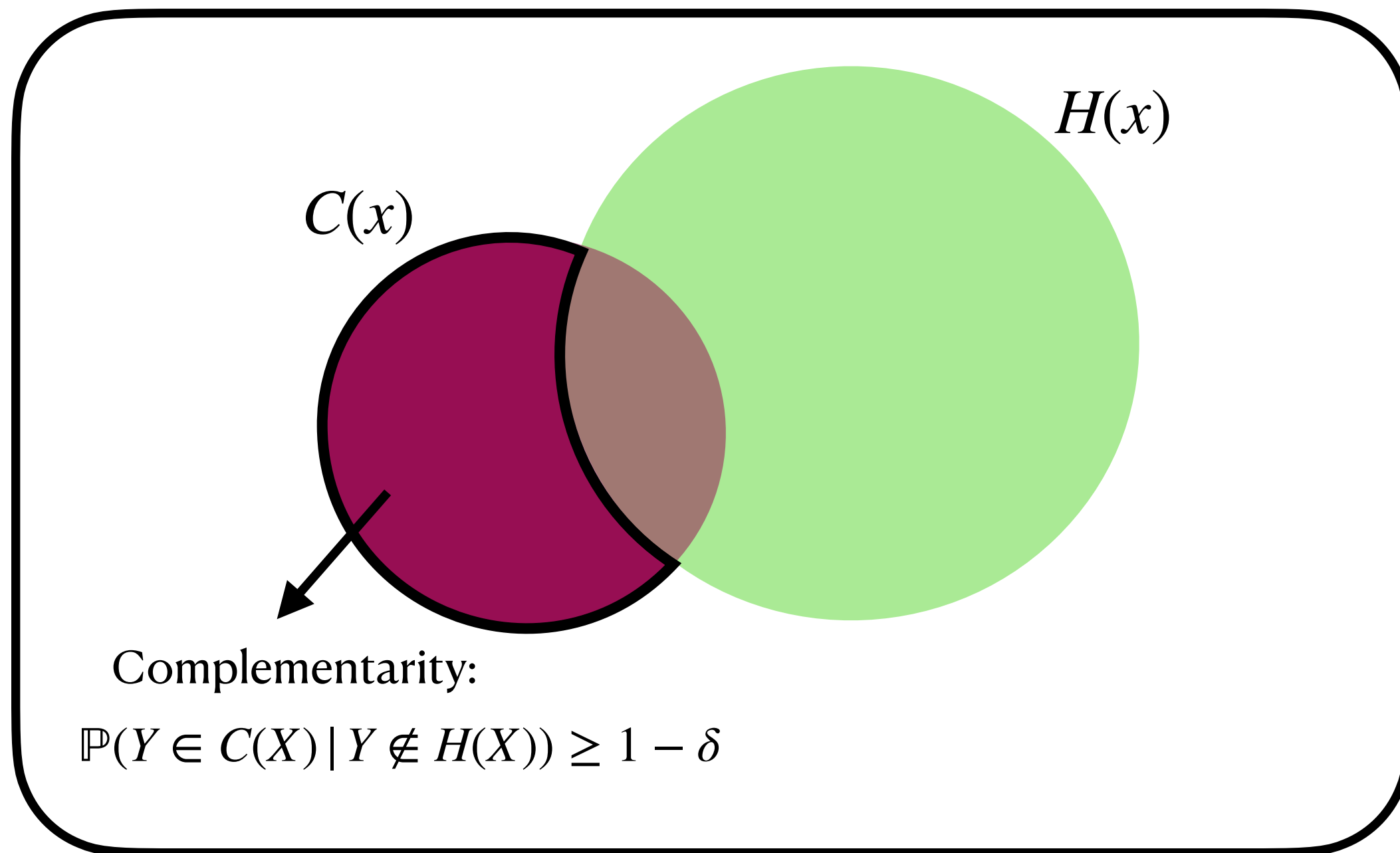
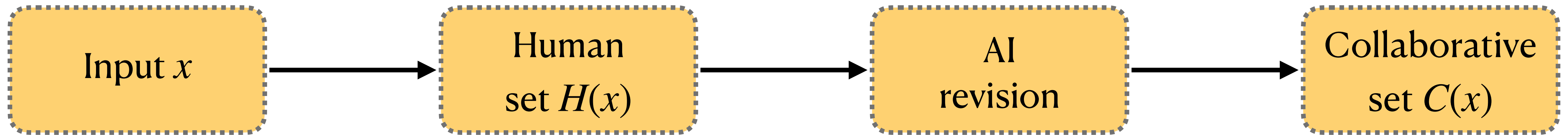
# Two fundamentals of collaboration



## Complementarity

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

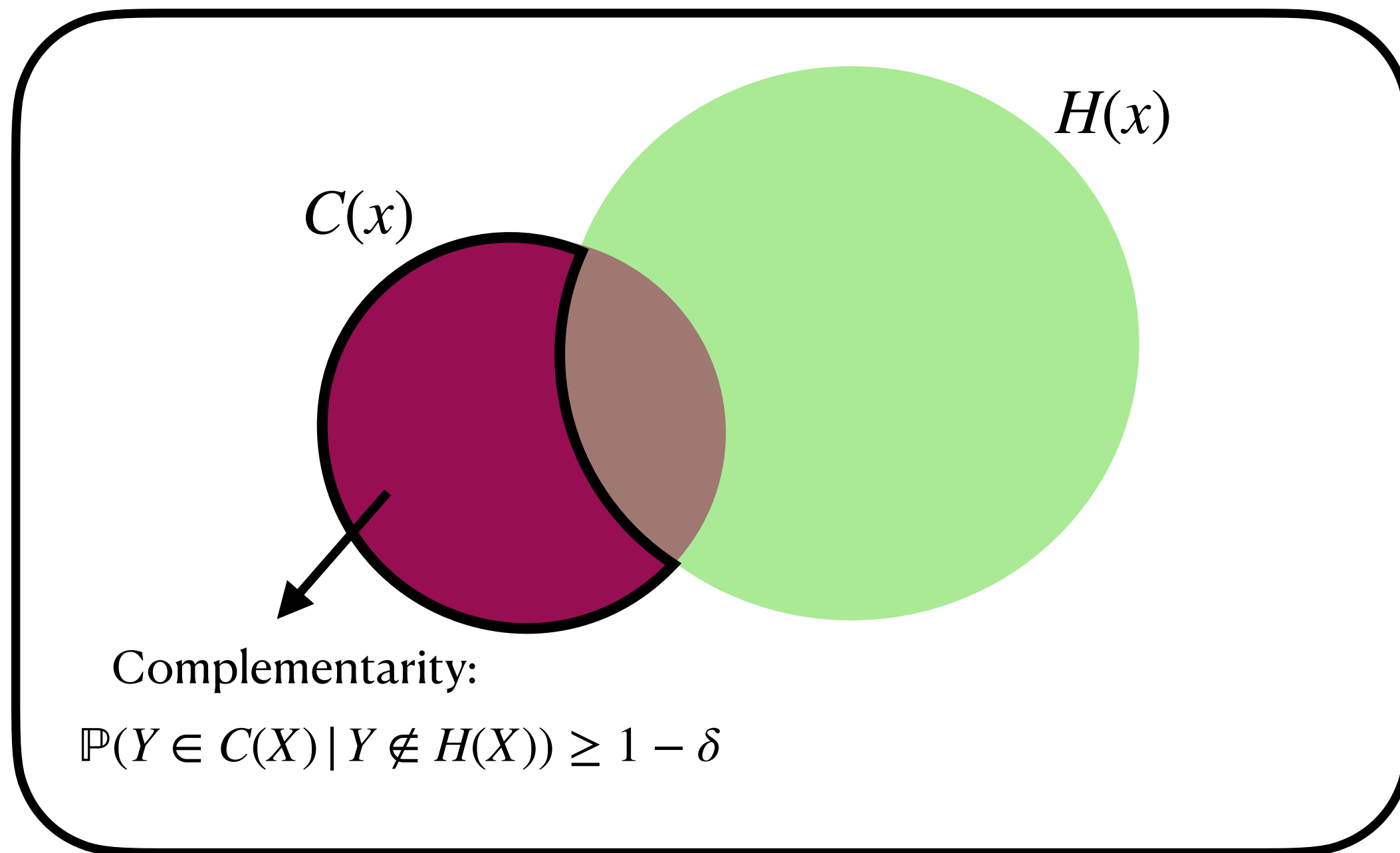
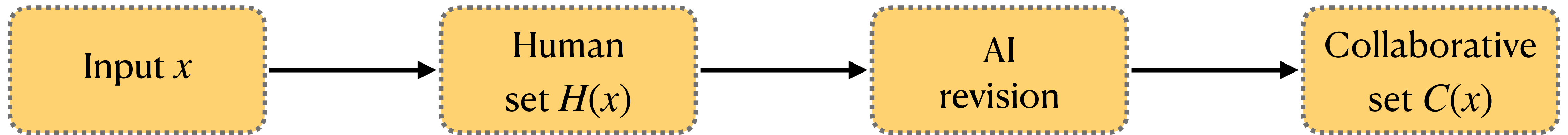
# Two fundamentals of collaboration



## Complementarity

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

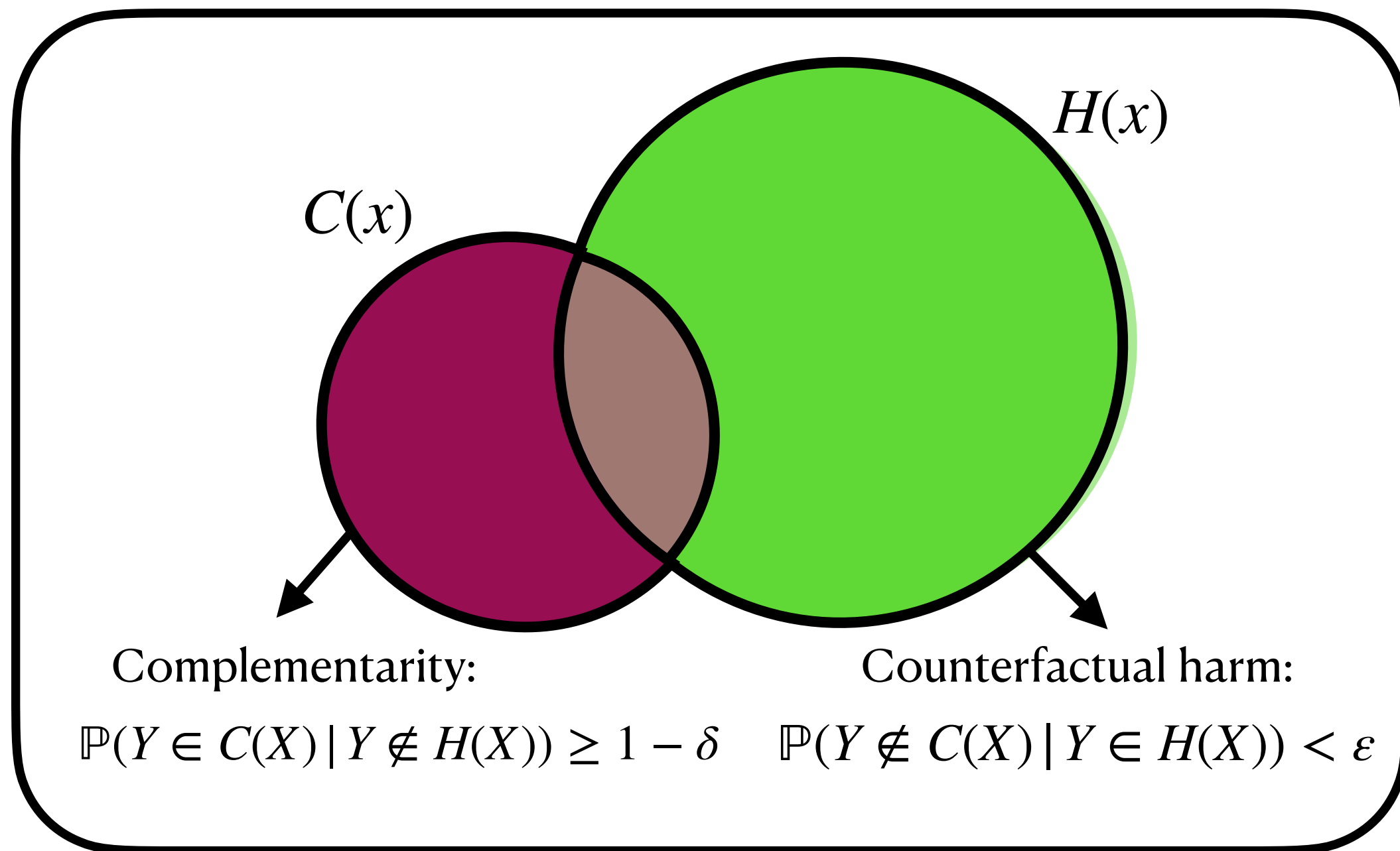
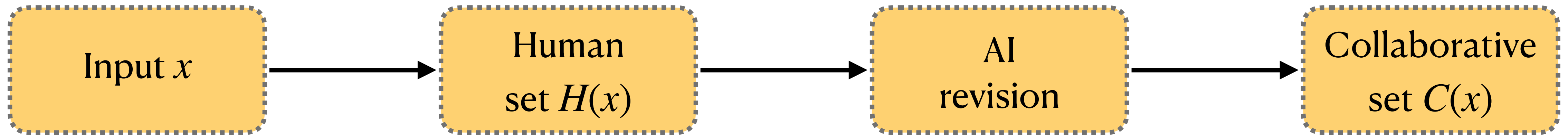
# Two fundamentals of collaboration



## 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?

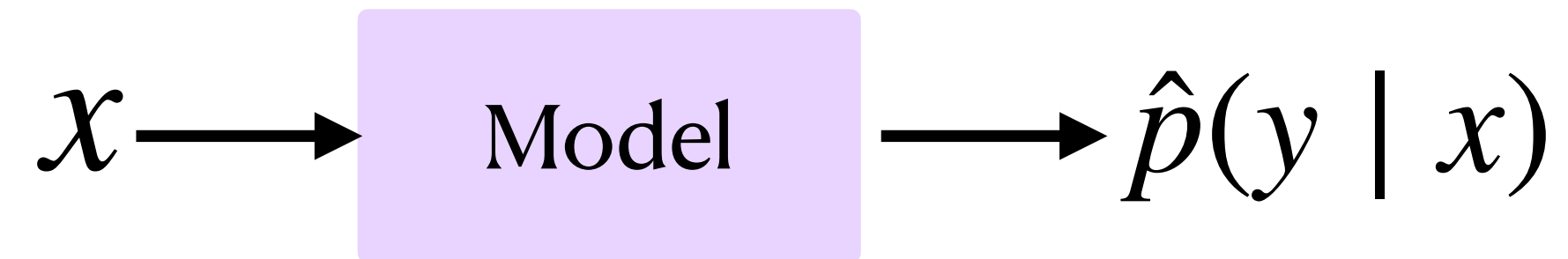
2 **Question:** How to debias the thresholds?

**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

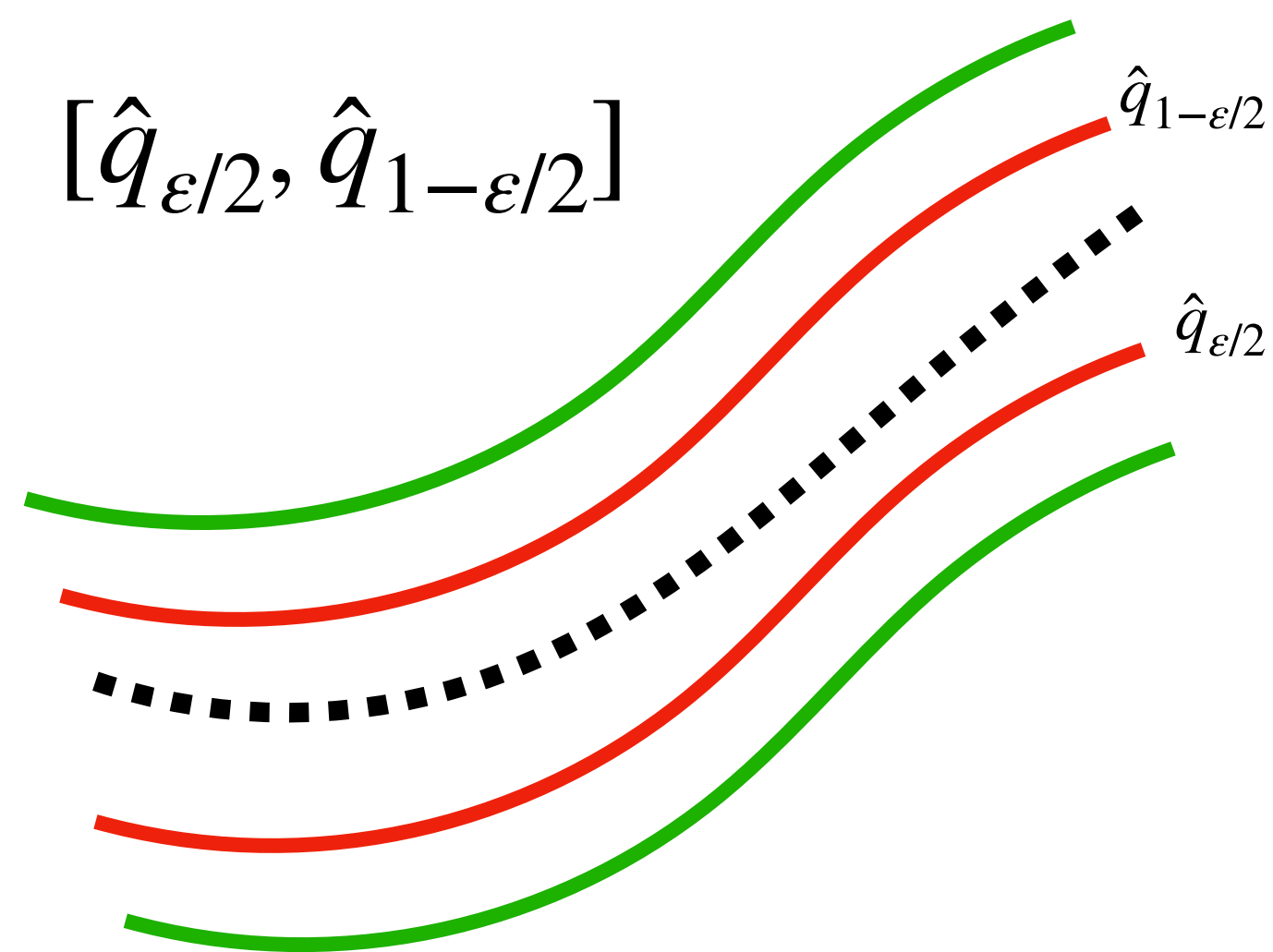


$$\hat{s}(x, y) = 1 - \hat{p}(y \mid x)$$

## 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

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

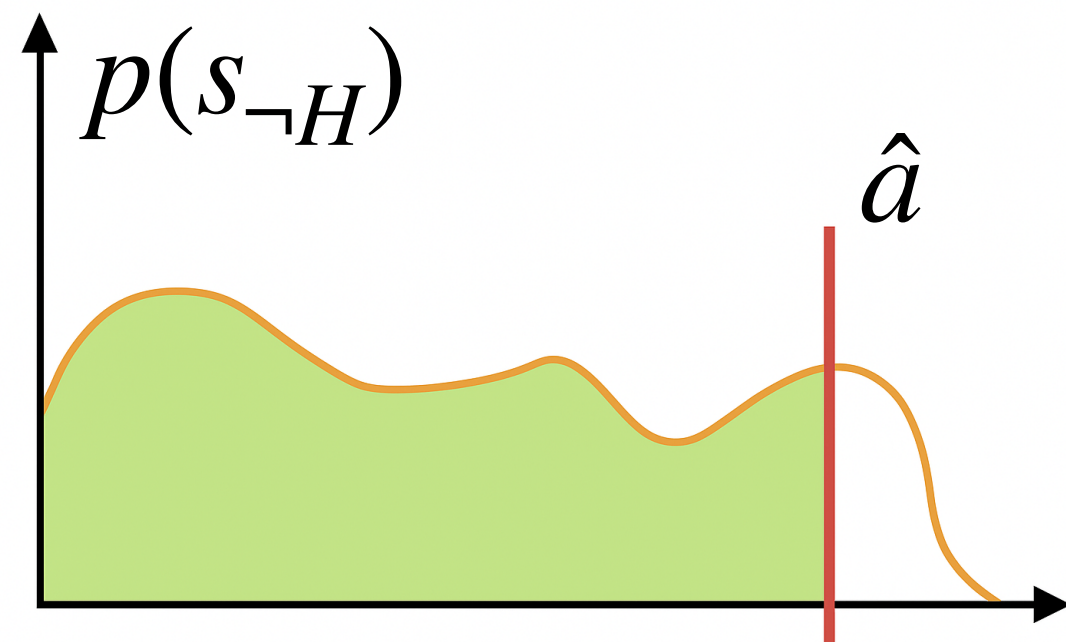
**Assume Exchangeability**    What should a and b be?

## 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?

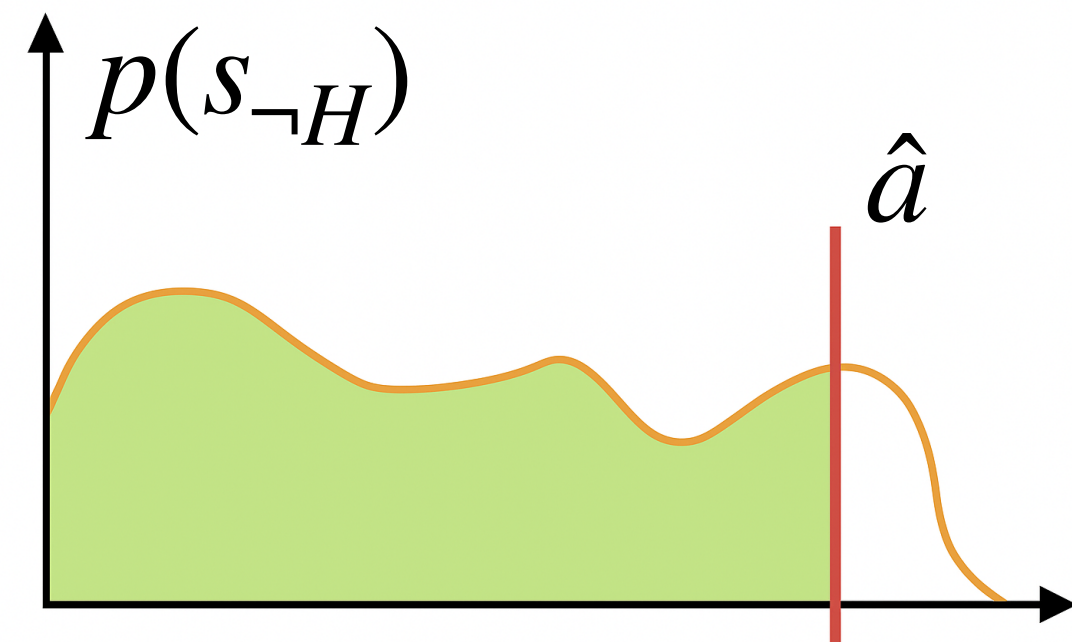


## Question: How to debias the thresholds?

**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?



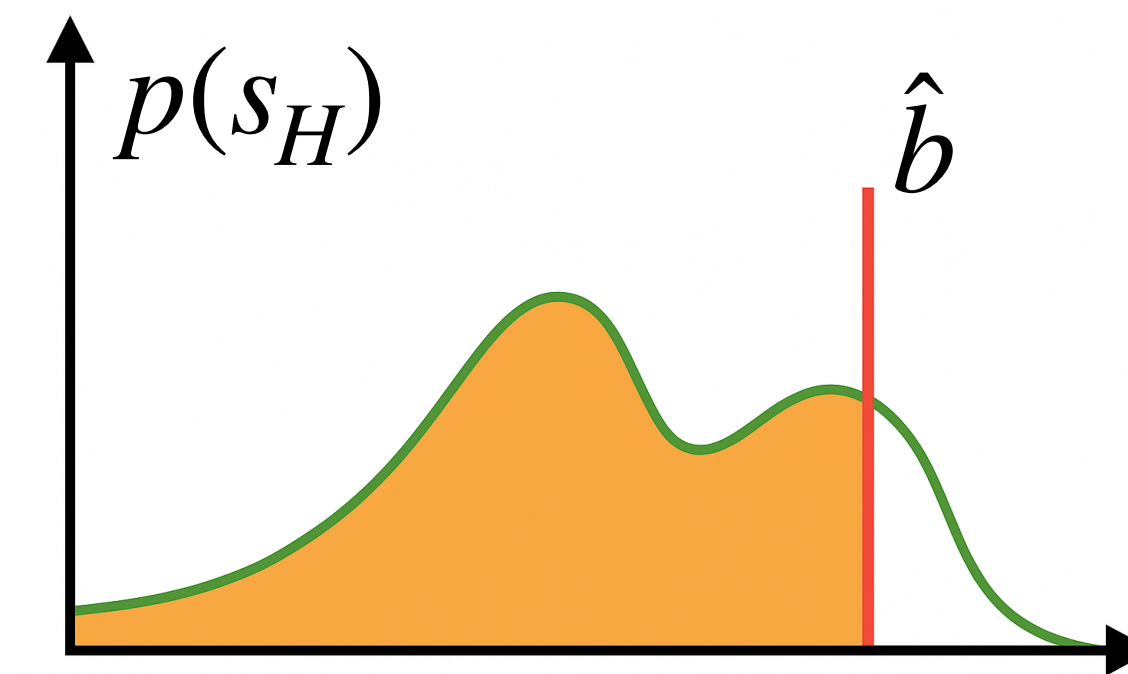
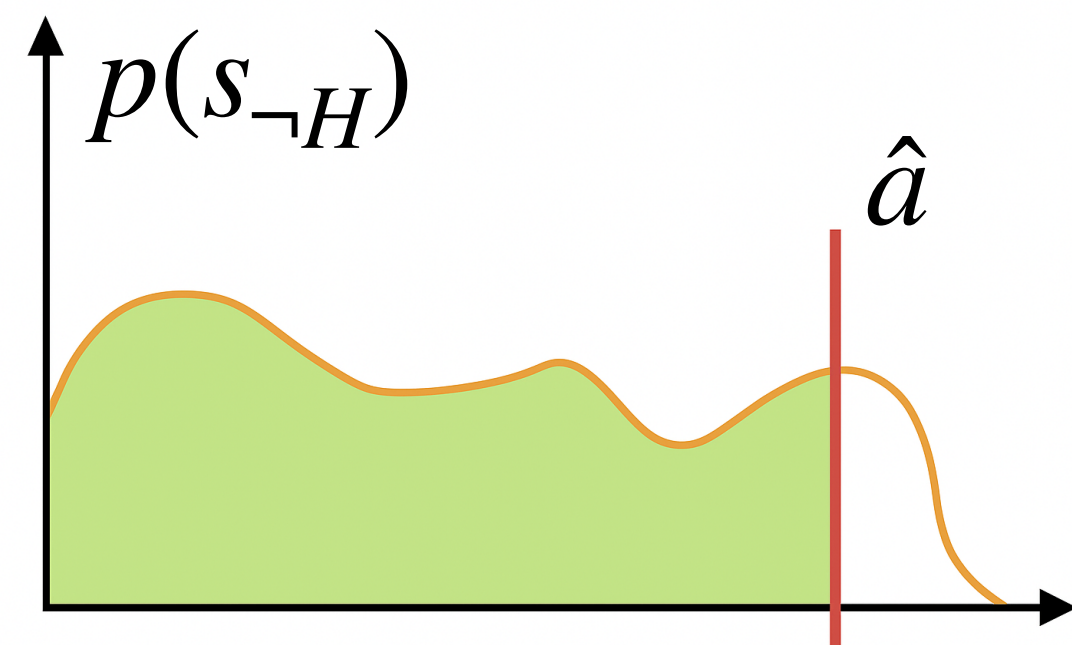
$$\hat{a} = Q_{1-\delta}(\{s_i : Y_i \notin H(X_i)\} \cup \{\infty\})$$

## 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?



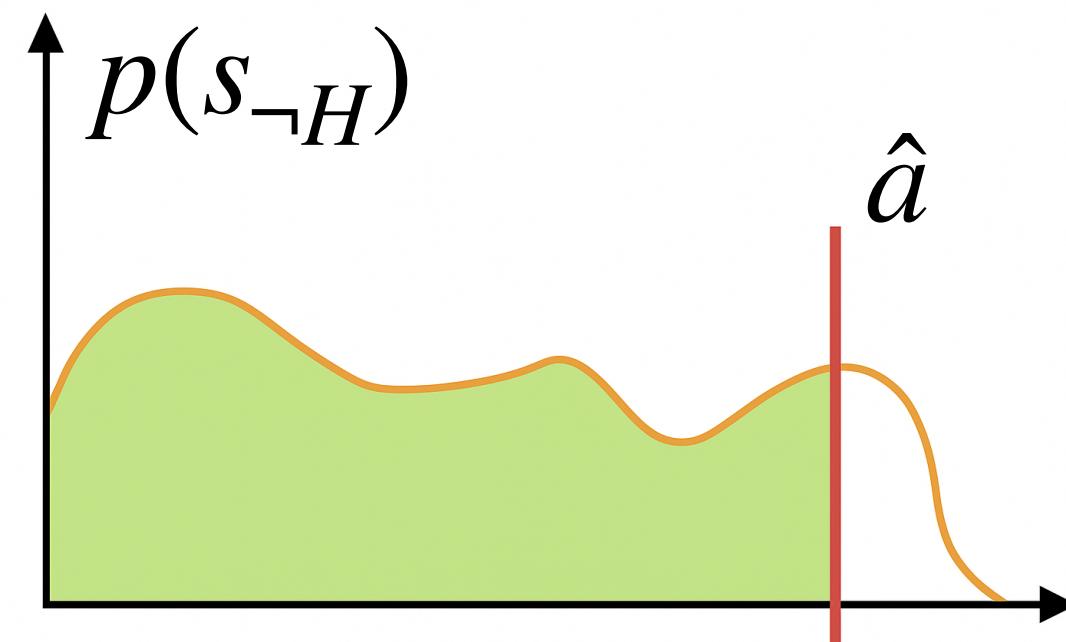
$$\hat{a} = Q_{1-\delta}(\{s_i : Y_i \notin H(X_i)\} \cup \{\infty\})$$

## Question: How to debias the thresholds?

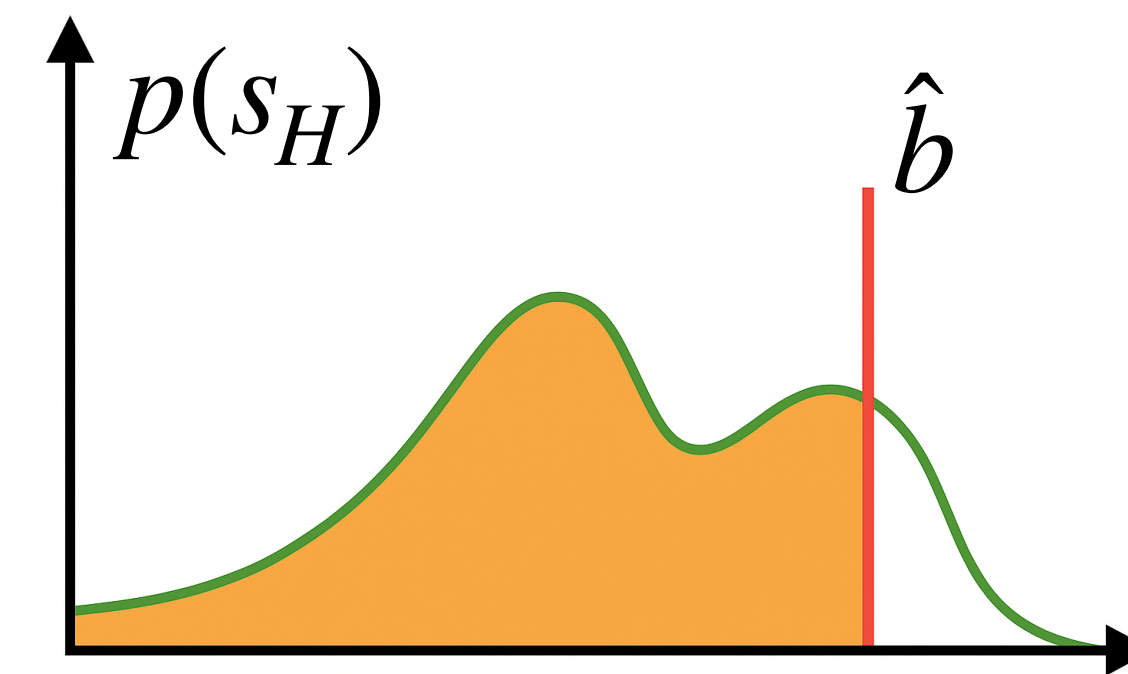
**Theorem:** The optimal solution to HACO is of the form

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

**Assume Exchangeability** What should a and b be?



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



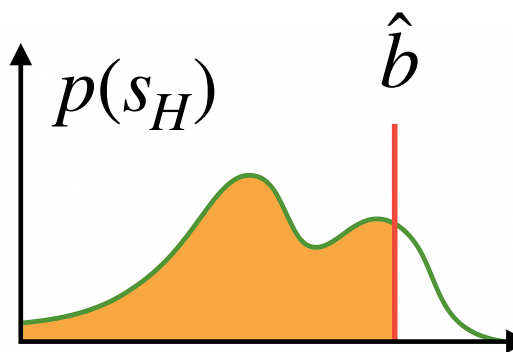
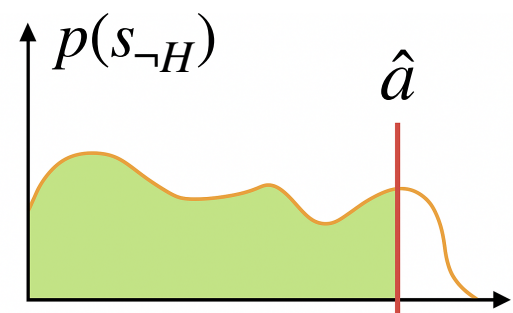
$$\hat{b} = Q_{1-\varepsilon}(\{s_i : Y_i \in H(X_i)\} \cup \{\infty\})$$

## 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**

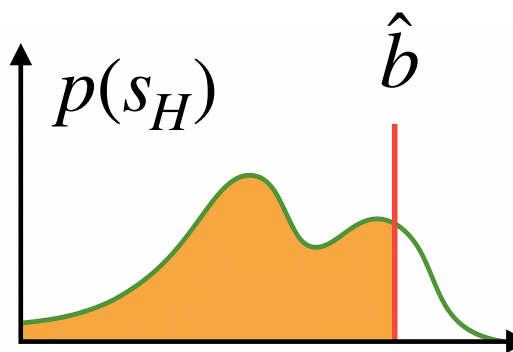
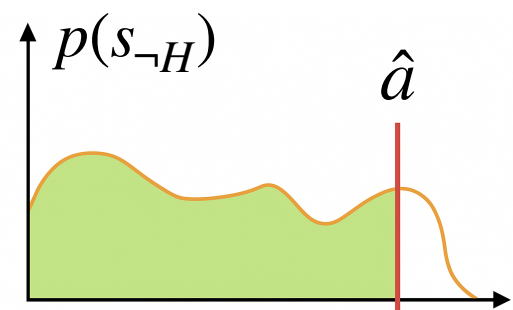


## Question: How to debias the thresholds?

**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**



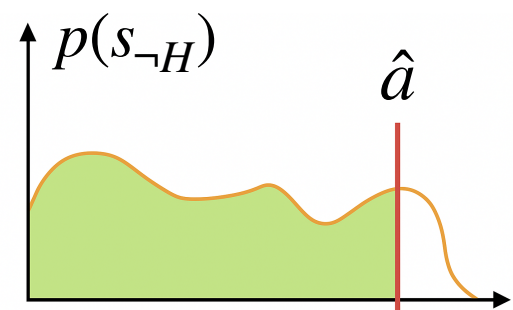
For a **new** test example  $(x_{test}, y_{test})$  :

## Question: How to debias the thresholds?

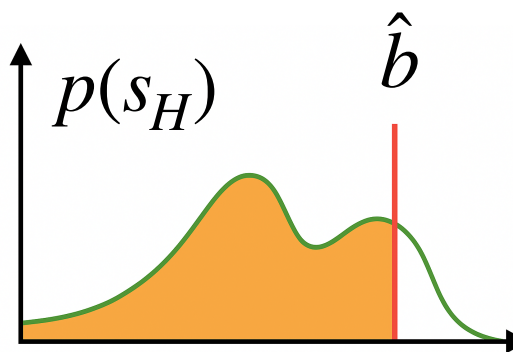
**Theorem:** The optimal solution to HACO is of the form

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

**Assume Exchangeability**



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



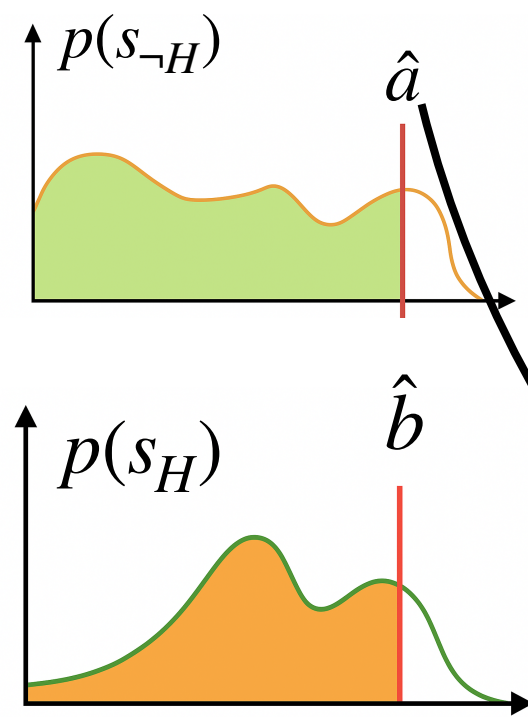
$$C(x_{test}) = \left\{ y \mid s(x_{test}, y) \leq \right\}.$$

## 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  $(x_{test}, y_{test})$ :

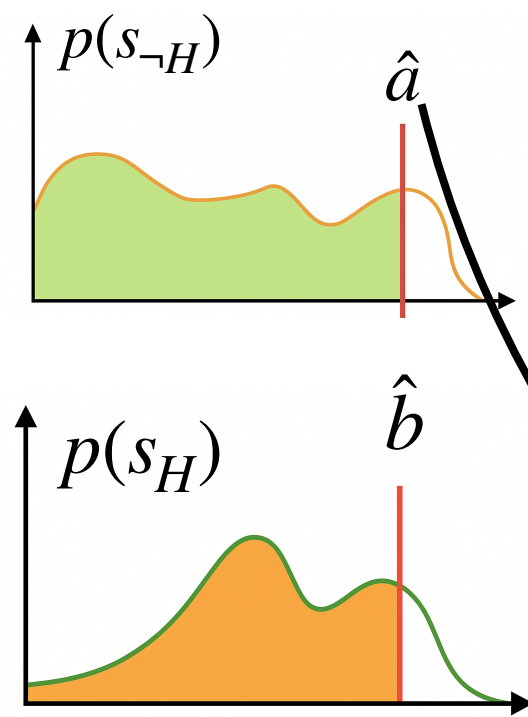
$$C(x_{test}) = \left\{ y \mid s(x_{test}, y) \leq \right\}.$$

## 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  $(x_{test}, y_{test})$  :

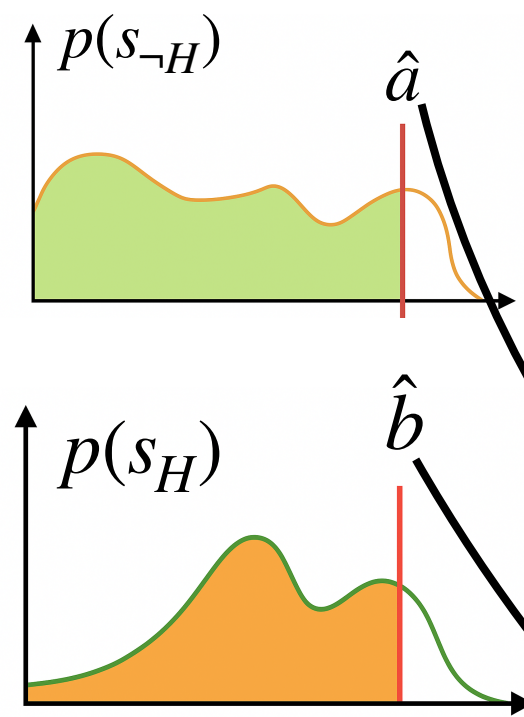
$$C(x_{test}) = \left\{ y \mid s(x_{test}, y) \leq \hat{a} \mathbf{1}_{\{y \notin H(x_{test})\}} \right\}.$$

## 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  $(x_{test}, y_{test})$  :

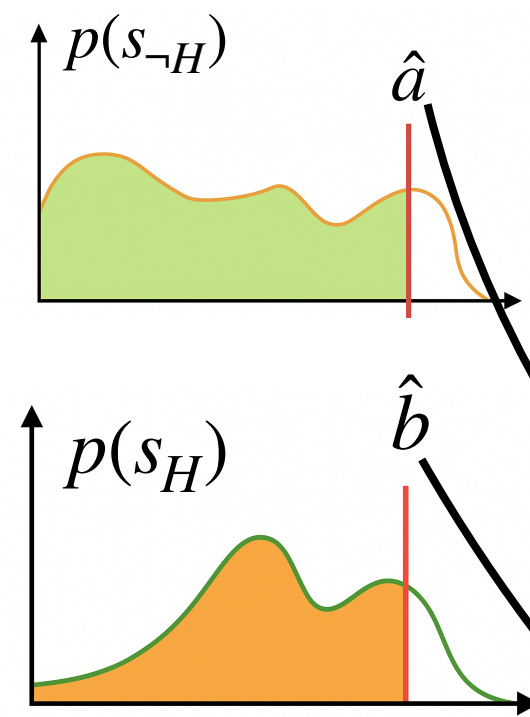
$$C(x_{test}) = \left\{ y \mid s(x_{test}, y) \leq \hat{a} \mathbf{1}_{\{y \notin H(x_{test})\}} + \hat{b} \mathbf{1}_{\{y \in H(x_{test})\}} \right\}.$$

## 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

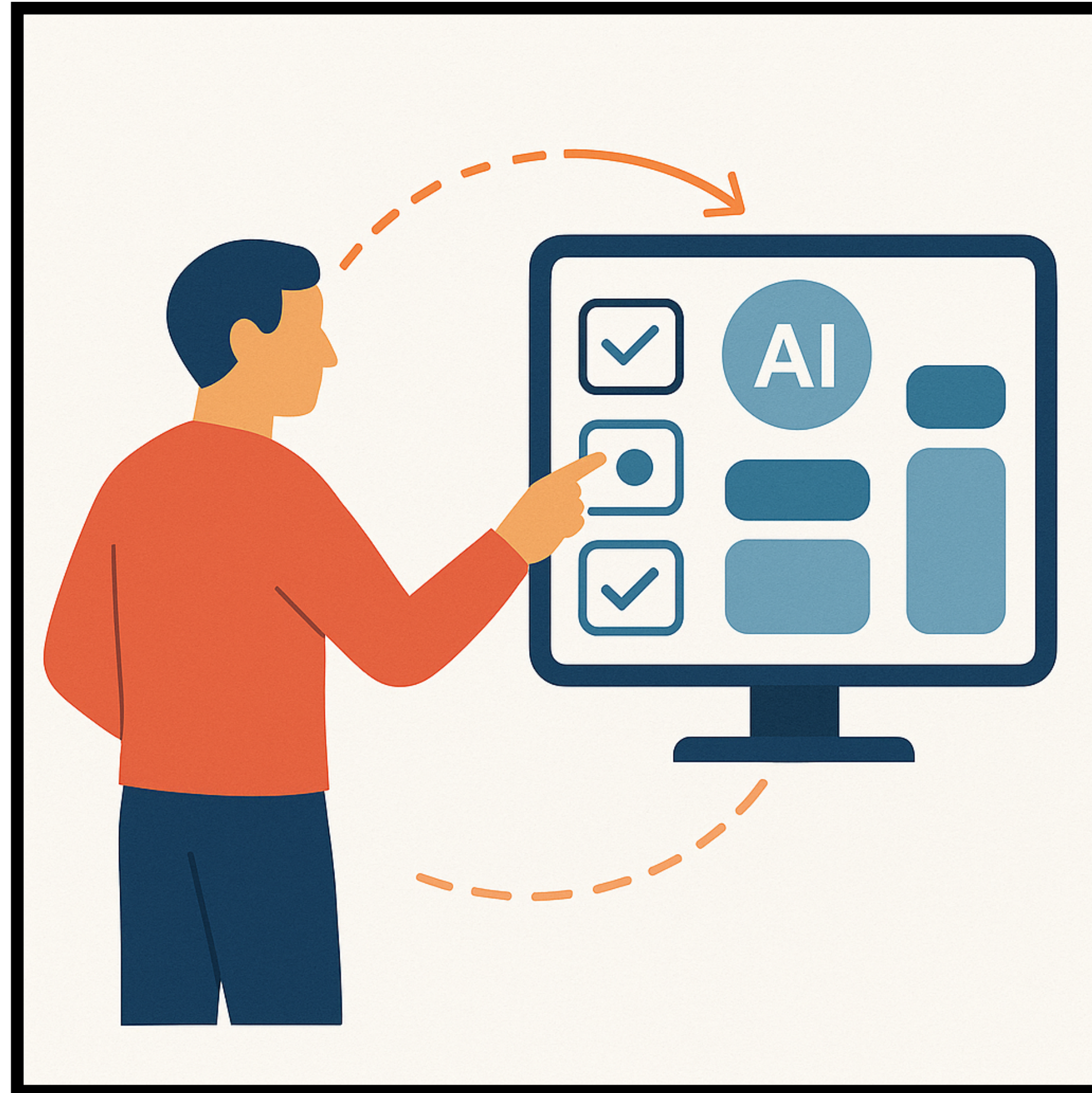
$$1 - \varepsilon \leq \Pr[Y_{\text{test}} \in C(X_{\text{test}}) \mid Y_{\text{test}} \in H(X_{\text{test}})] < 1 - \varepsilon + \frac{1}{n_1 + 1}$$

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

# ~~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

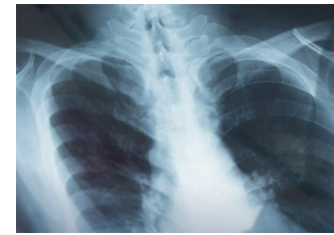
At round  $t$



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

At round  $t$

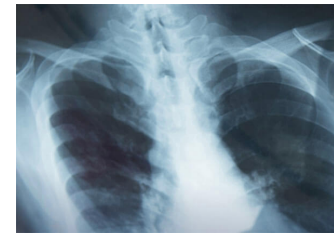
input  $x_t$



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

At round  $t$

input  $x_t$



$H(x_t)$



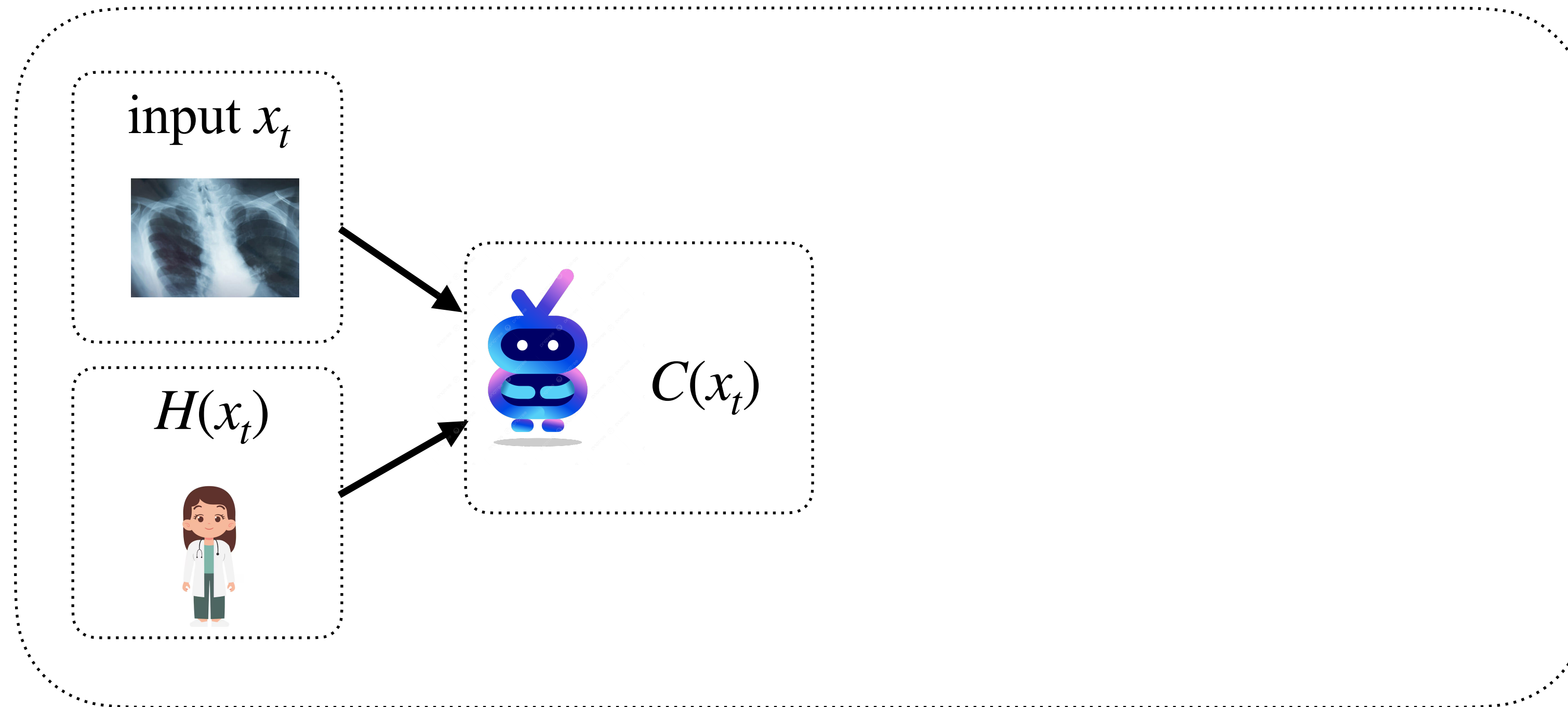
# Question: How to debias the thresholds in the **online** setting

At round  $t$



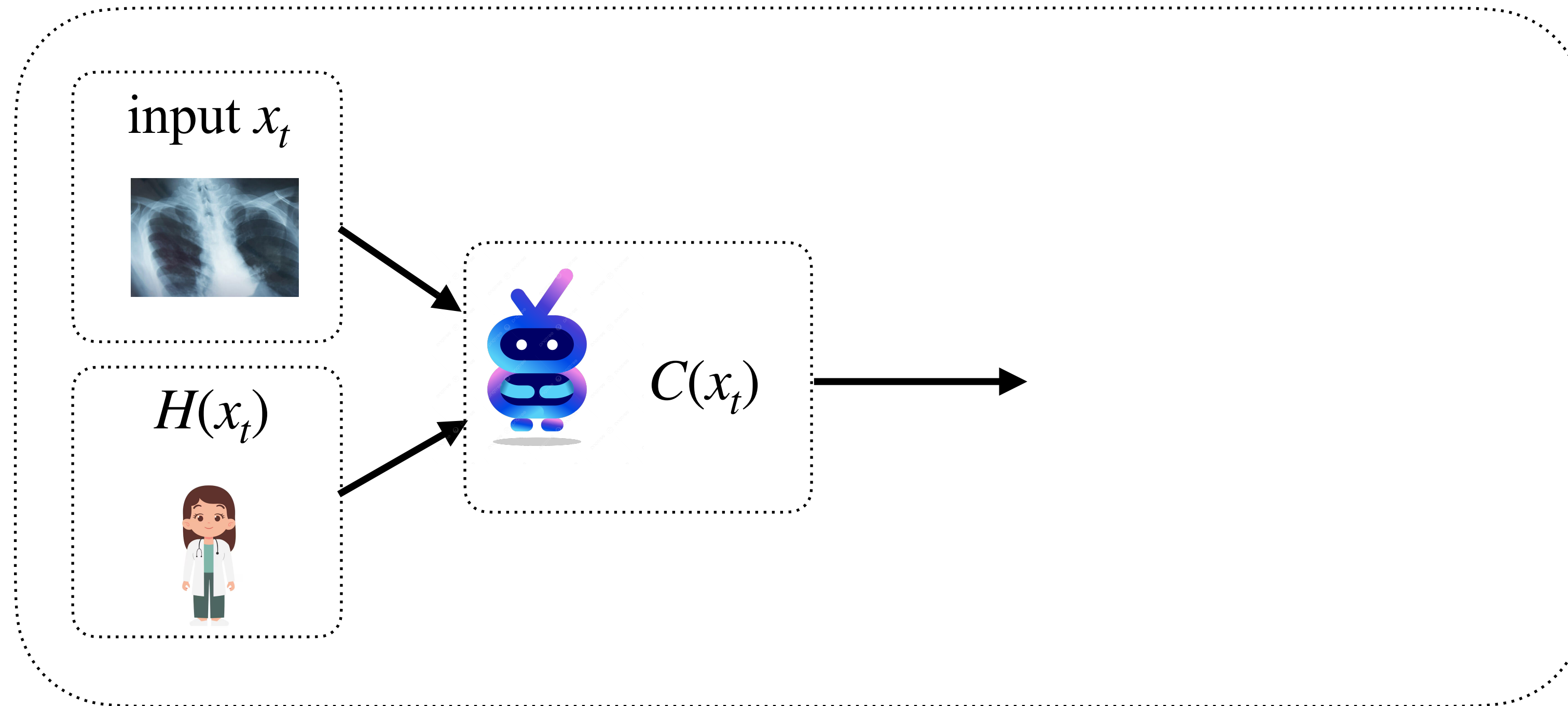
# Question: How to debias the thresholds in the **online** setting

At round  $t$



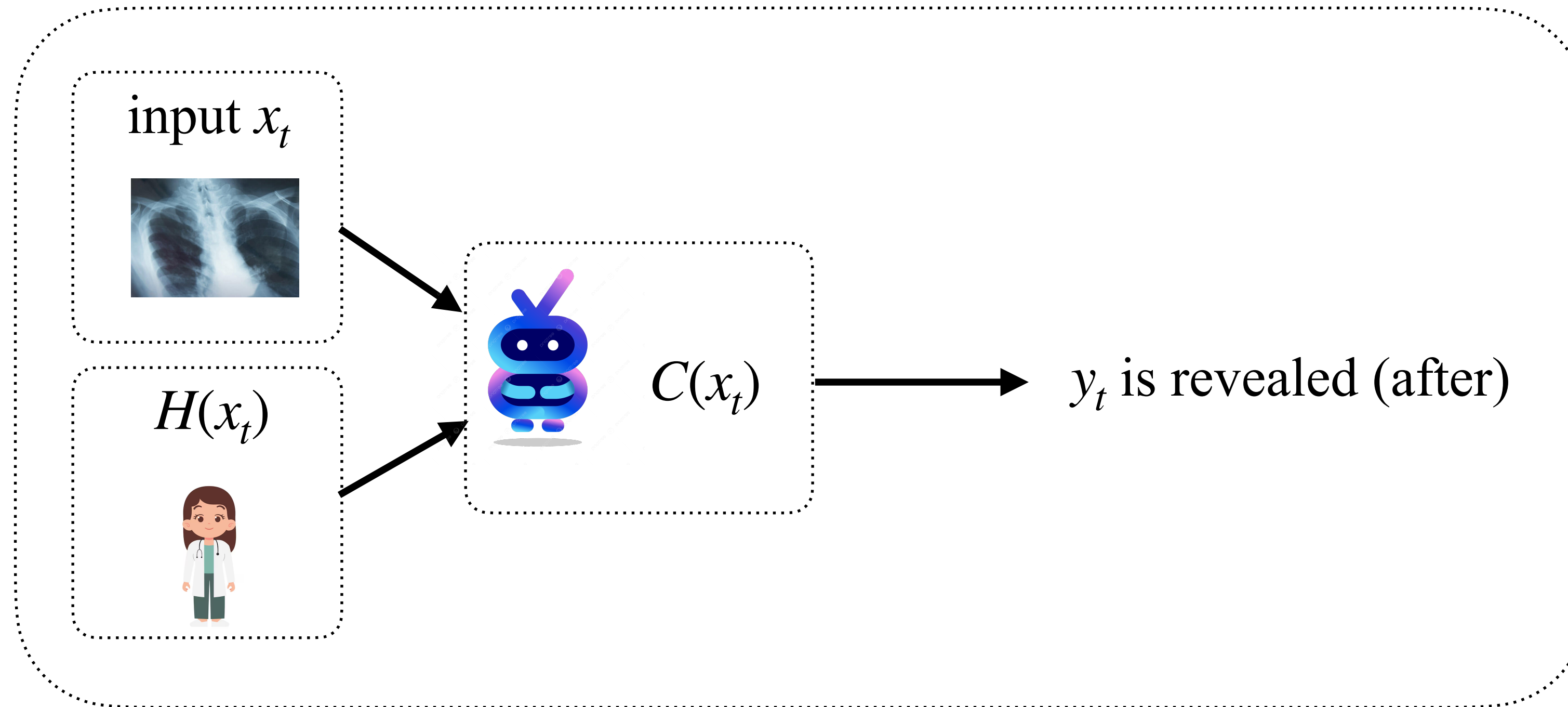
# Question: How to debias the thresholds in the **online** setting

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  $s(x, y) \in [0, 1]$ :

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

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

# Experiments

**We consider three modalities of data**

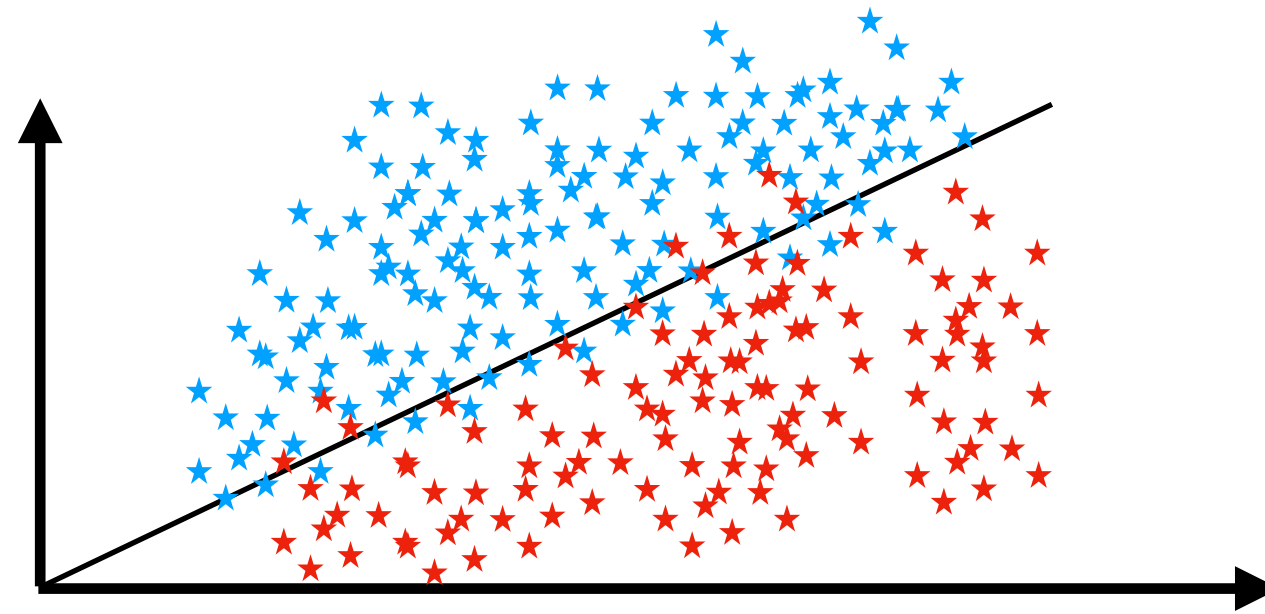
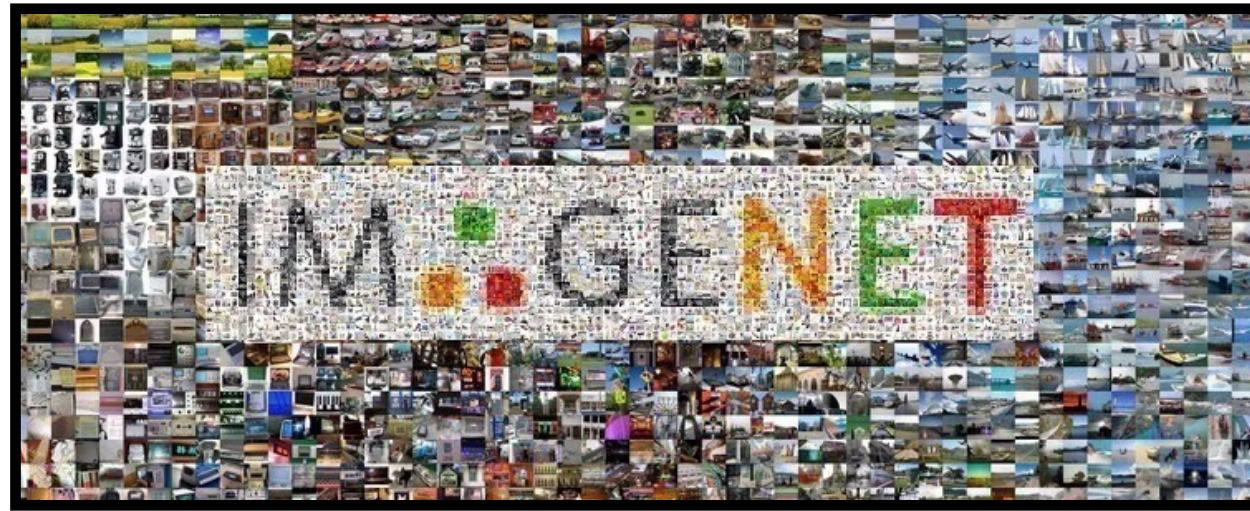
# Experiments

**We consider three modalities of data**



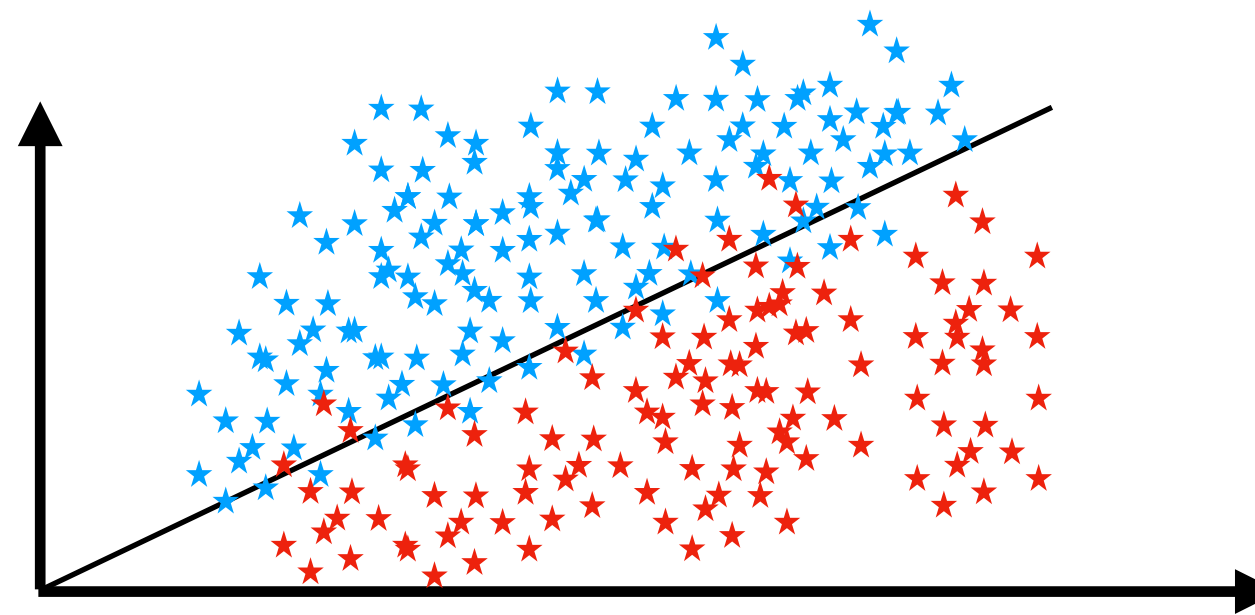
# Experiments

We consider three modalities of data



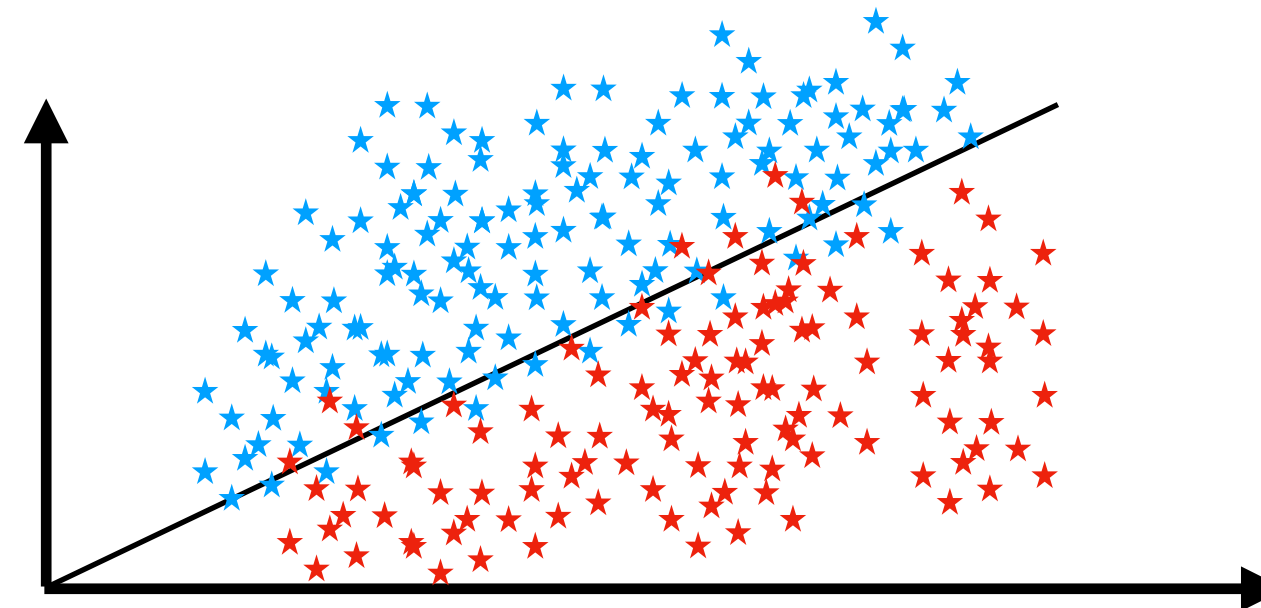
# Experiments

We consider three modalities of data



# Experiments

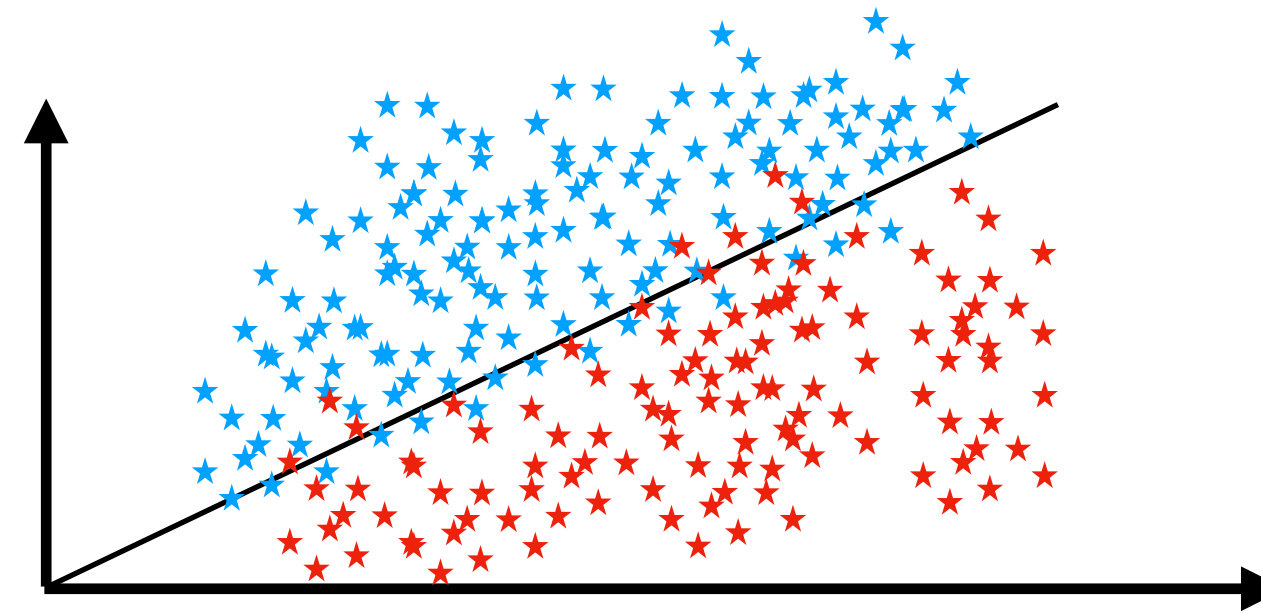
We consider three modalities of data



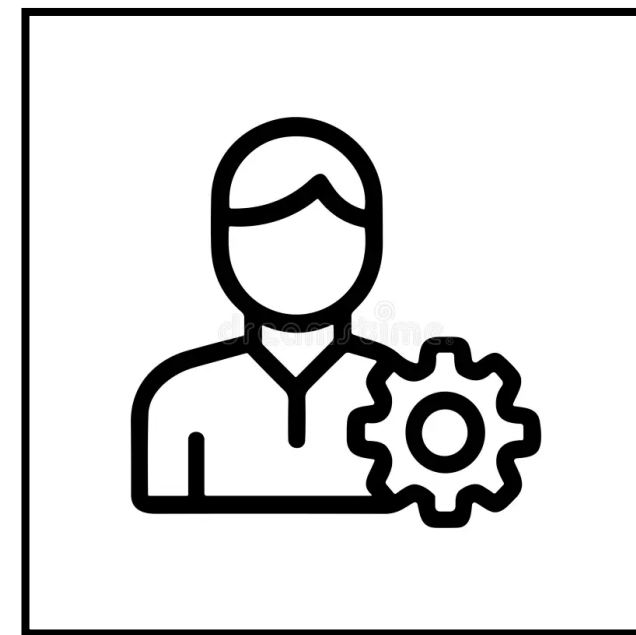
**Baselines**

# Experiments

We consider three modalities of data



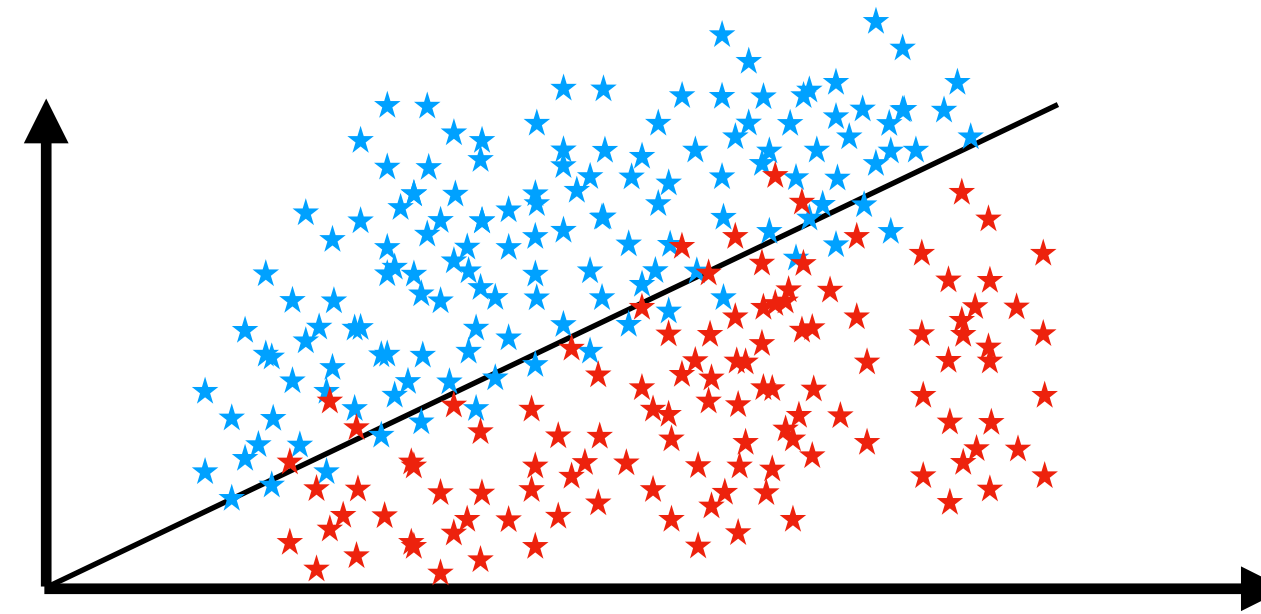
## Baselines



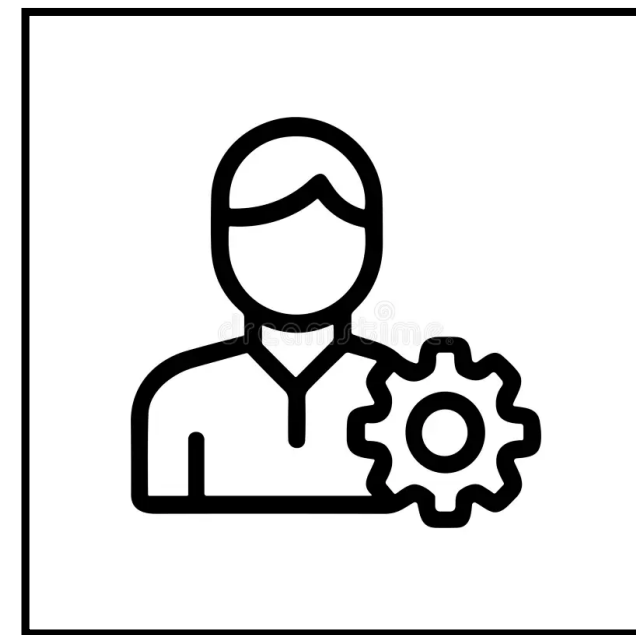
**Human Alone**

# Experiments

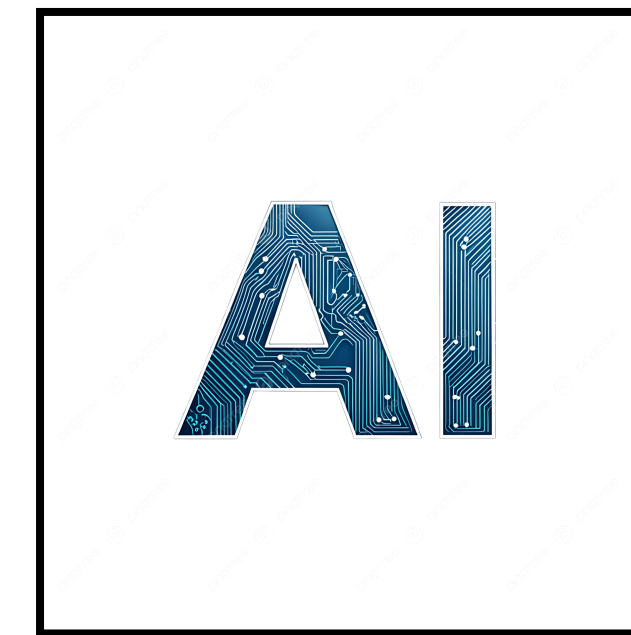
We consider three modalities of data



## Baselines

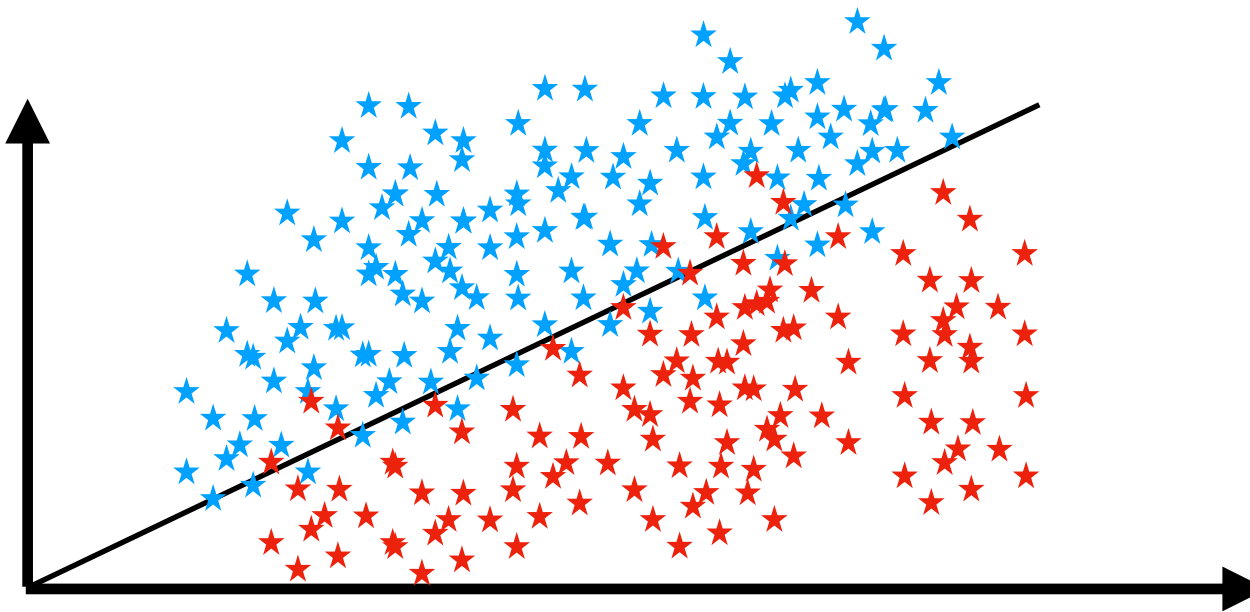


**Human Alone**



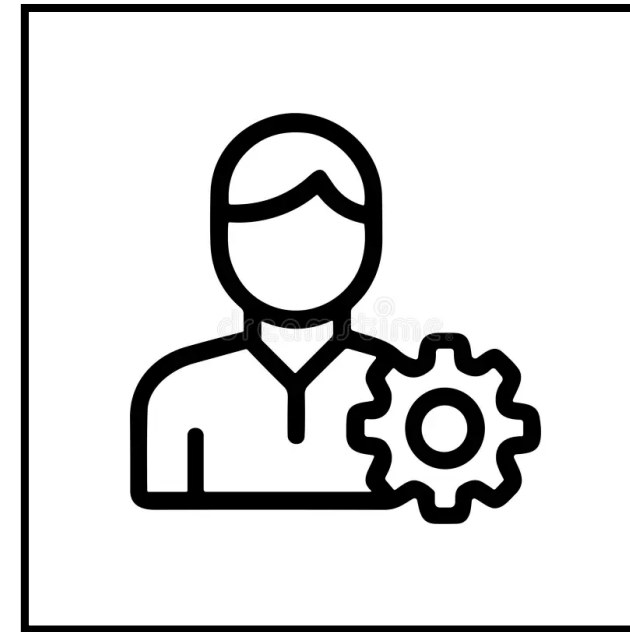
**AI Alone**

**We consider three modalities of data**

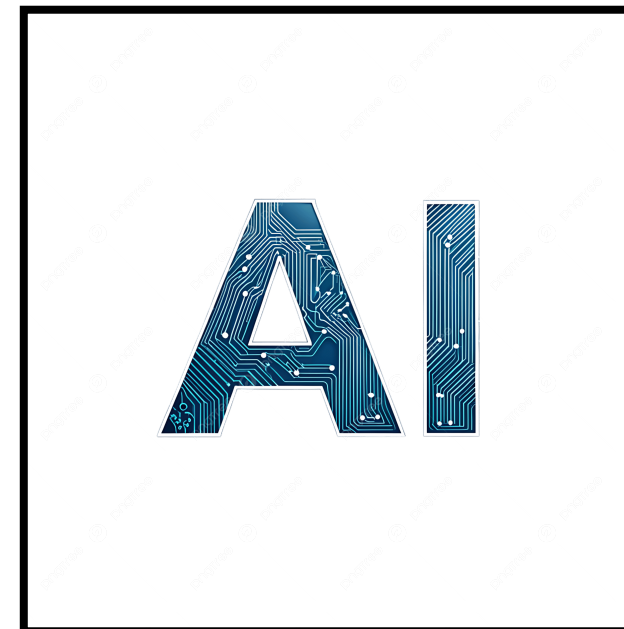


# Experiments

## Baselines



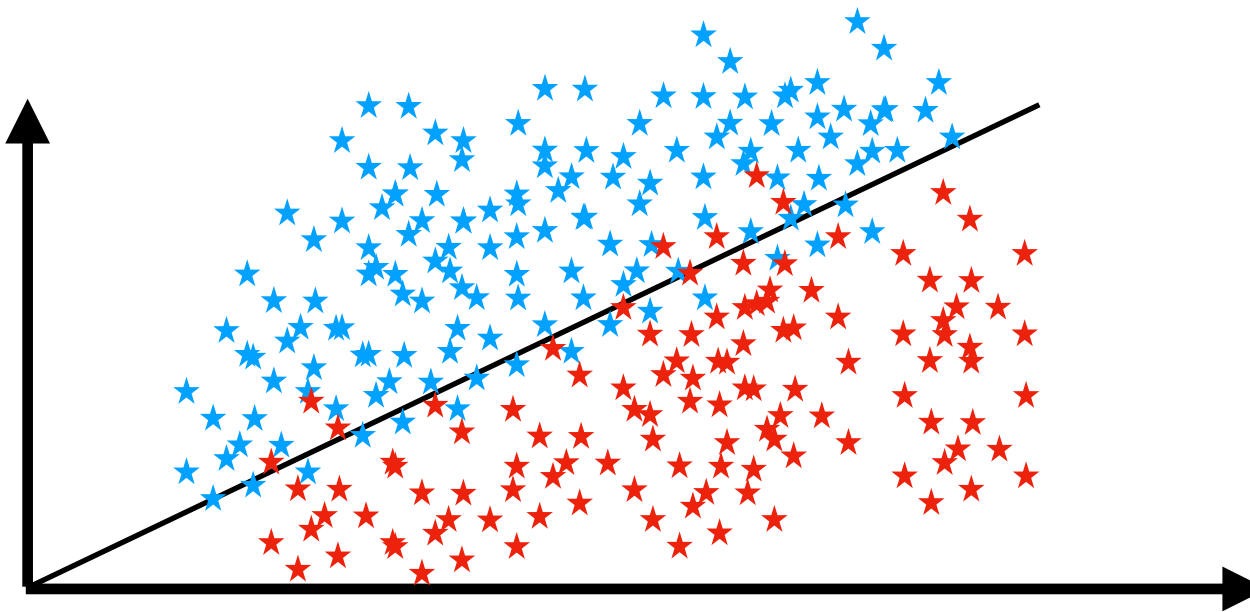
## Human Alone



## AI Alone

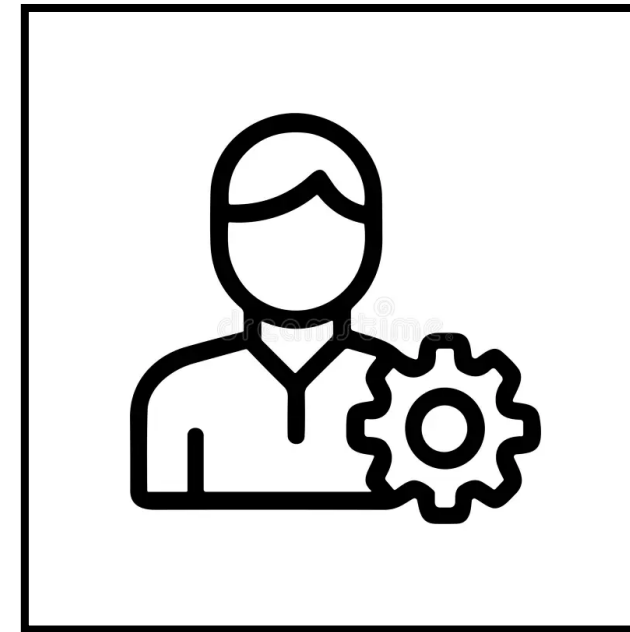
## Evaluation Metrics

**We consider three modalities of data**

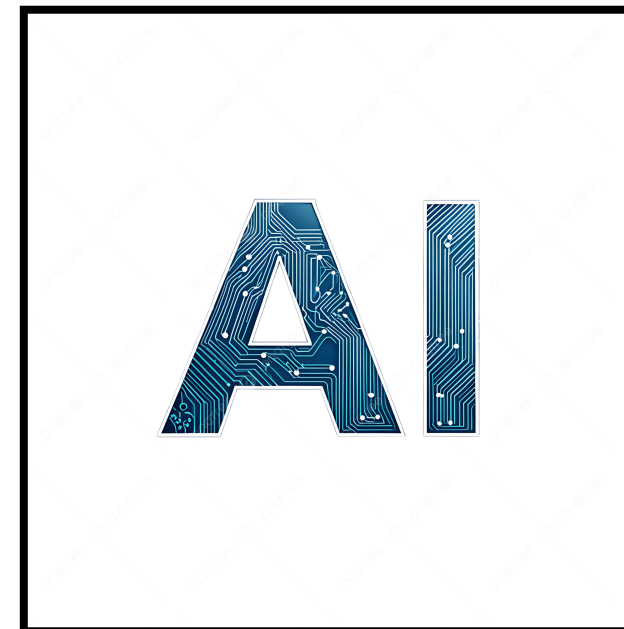


## Experiments

### Baselines

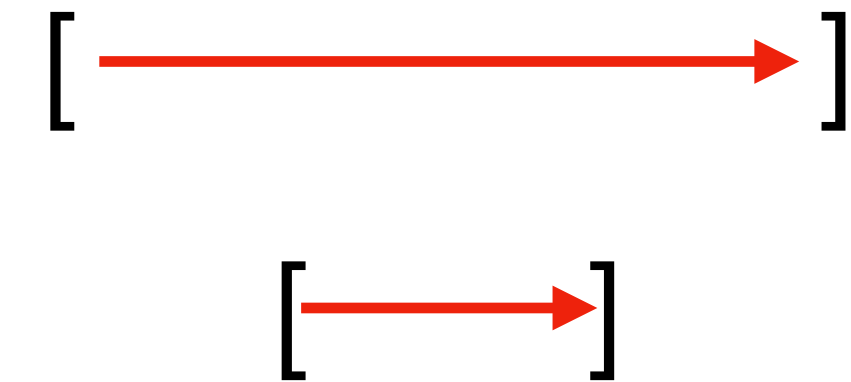


### Human Alone



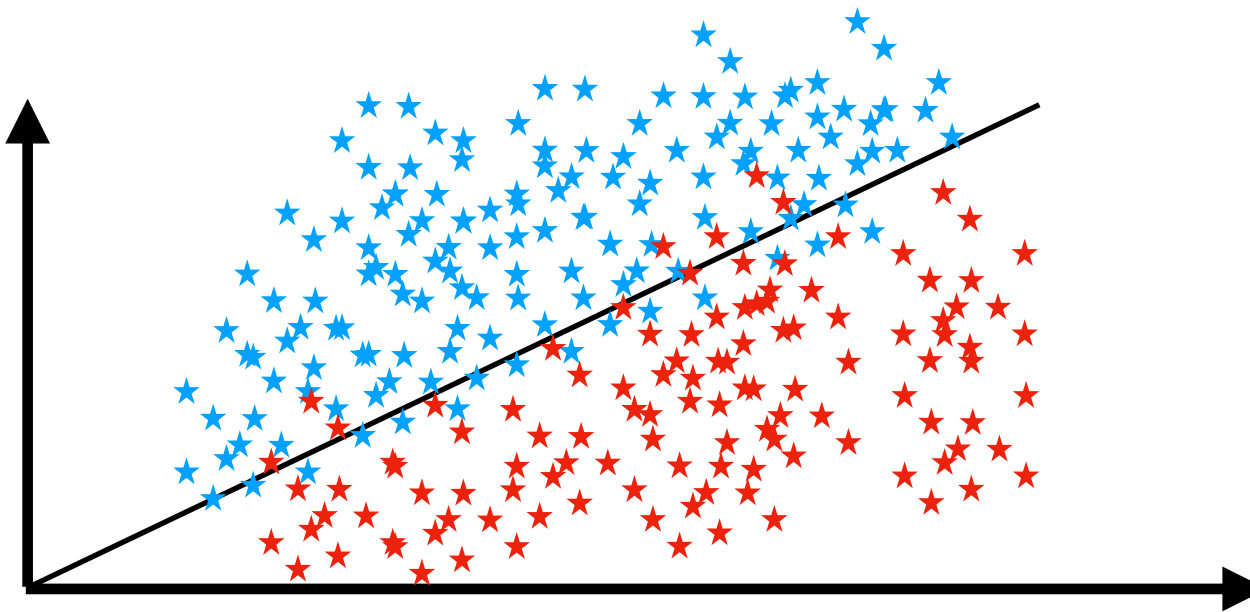
### AI Alone

### Evaluation Metrics



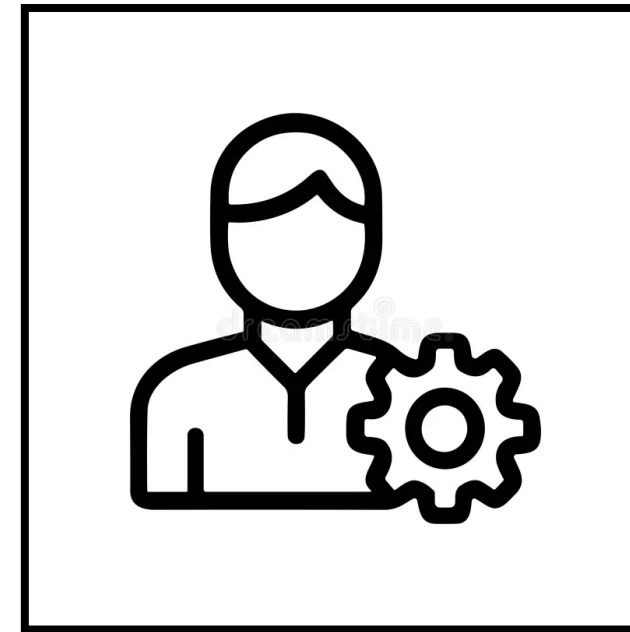
### Set Size

**We consider three modalities of data**

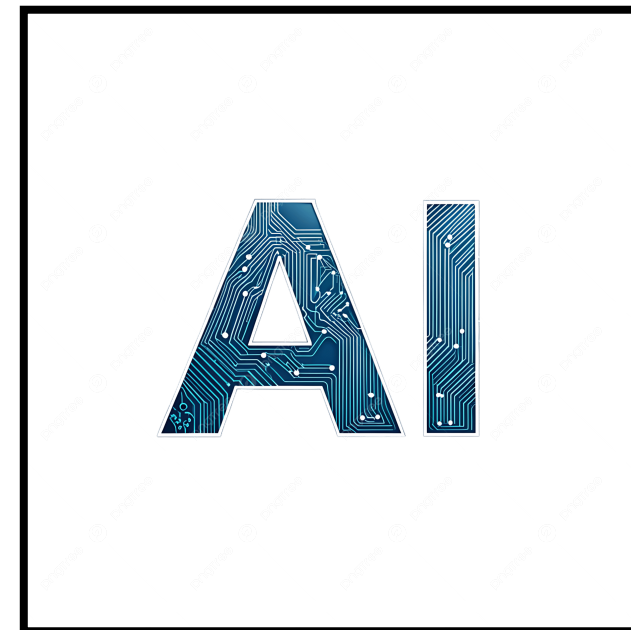


# Experiments

## Baselines

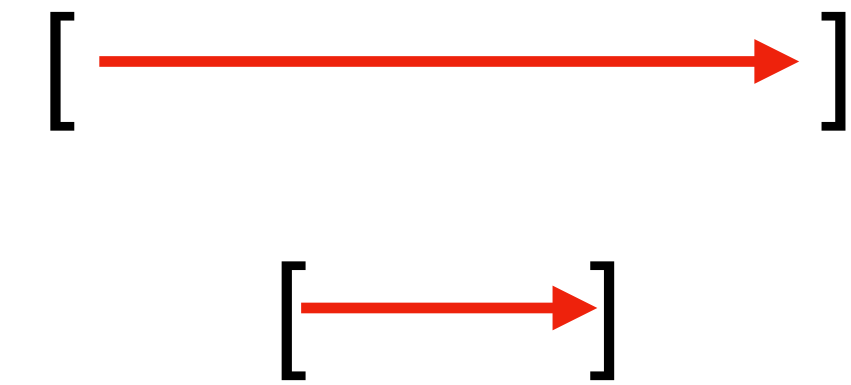


## Human Alone

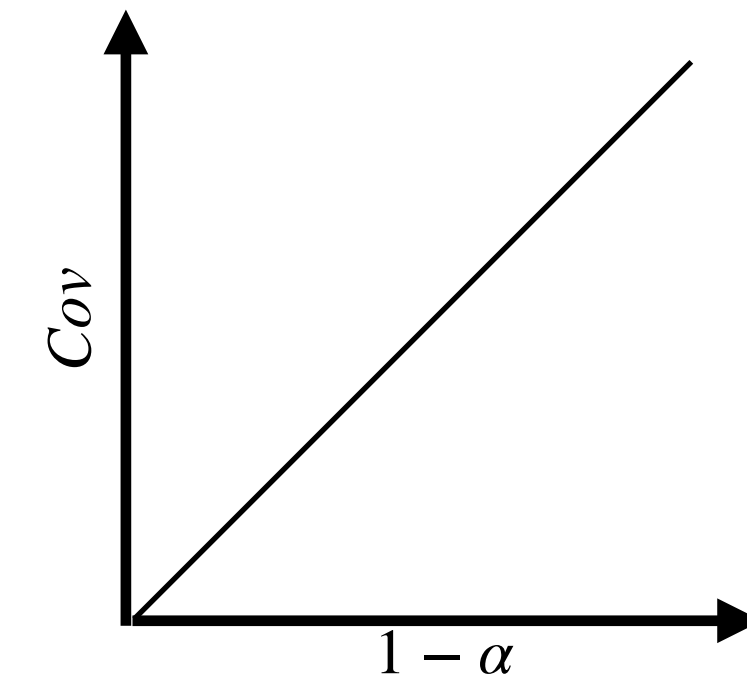


## AI Alone

## Evaluation Metrics



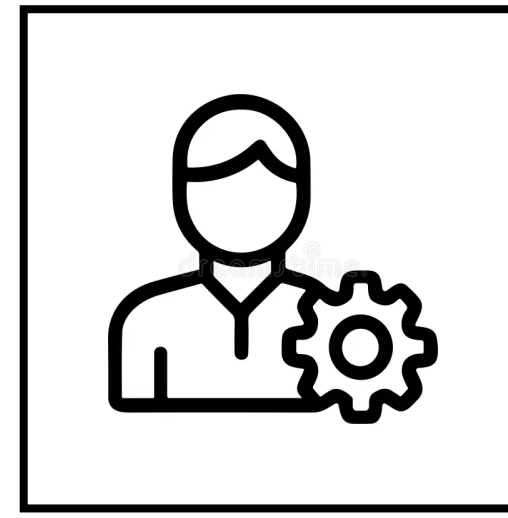
## Set Size



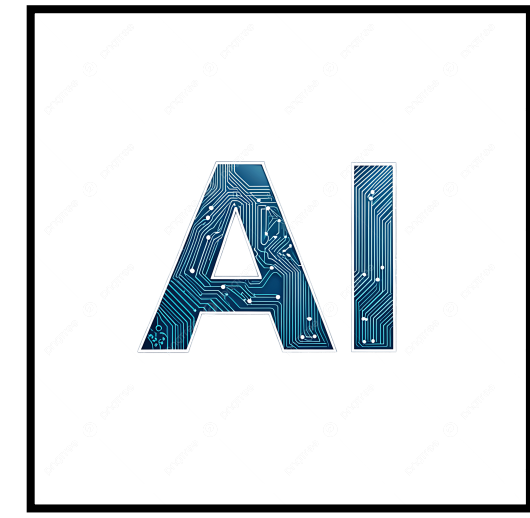
## Marginal Coverage



**Classification**



**Crowdsourcing**

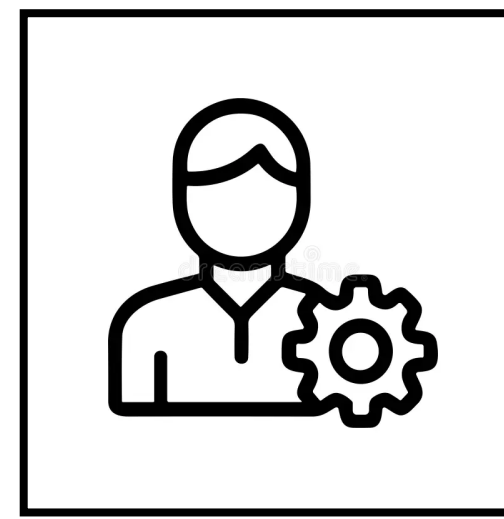


**AlexNet**

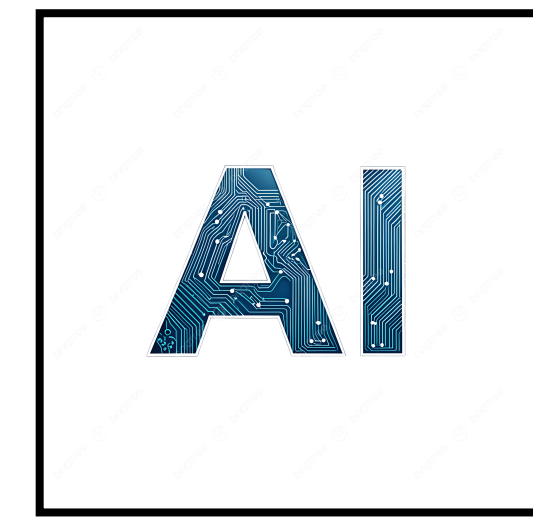
---



**Classification**



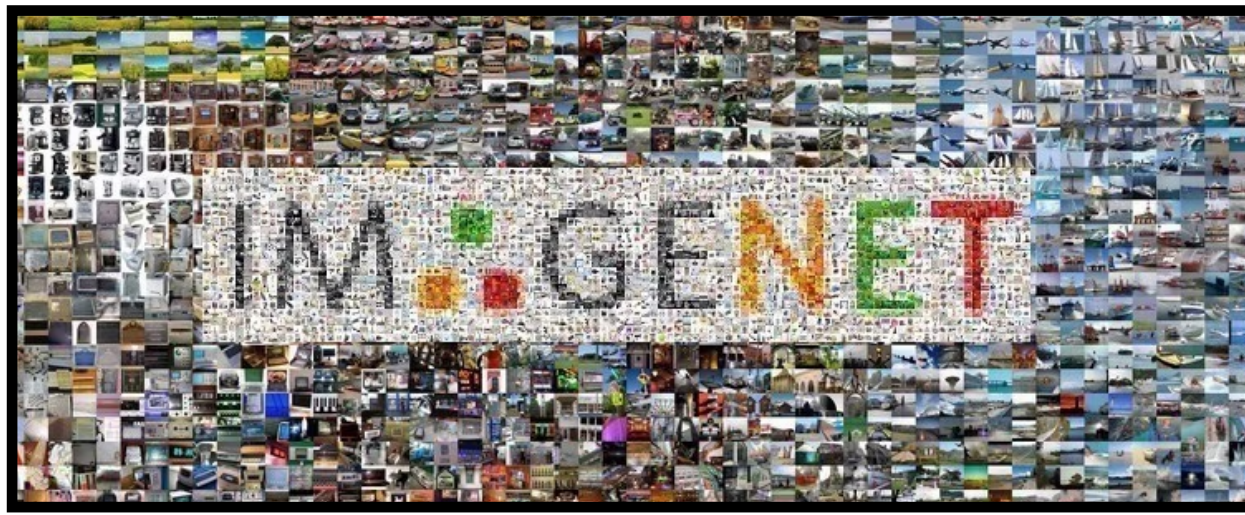
**Crowdsourcing**



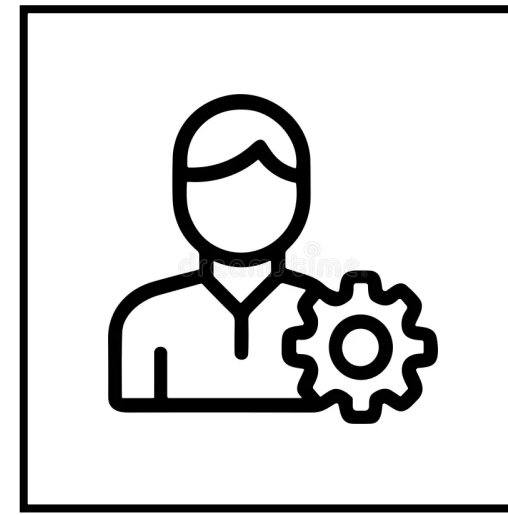
**AlexNet**

**Offline Setting**

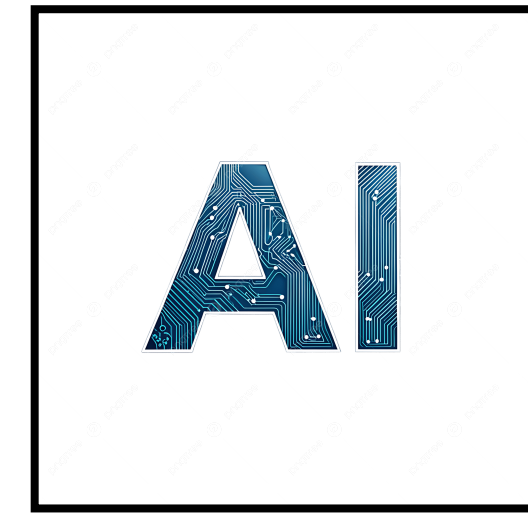
$\omega = 125$								
Strategy	Human Alone		CUP				AI Alone	
	Coverage	Size	Coverage	Size	$\epsilon$	$\delta$	Coverage	Size
Top-2	$0.8008 \pm 0.0090$	$2.00 \pm 0.00$	<b><math>0.9022 \pm 0.0083</math></b>	<b><math>1.49 \pm 0.04</math></b>	0.05	0.70	$0.9072 \pm 0.0138$	$1.65 \pm 0.07$
Top-1	$0.7245 \pm 0.0103$	$1.00 \pm 0.00$	<b><math>0.8823 \pm 0.0134</math></b>	$1.36 \pm 0.07$	0.05	0.70	$0.8828 \pm 0.0140$	$1.48 \pm 0.05$



**Classification**



**Crowdsourcing**



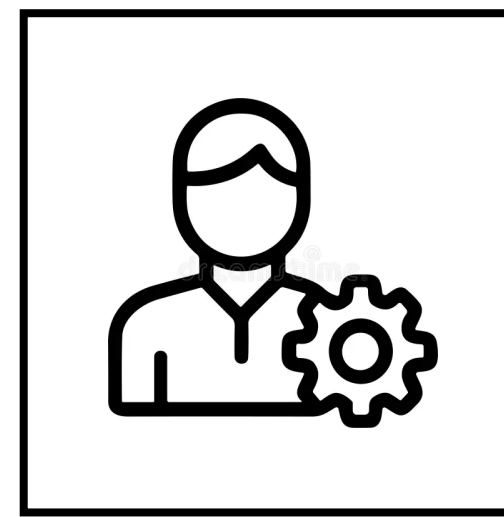
**AlexNet**

---

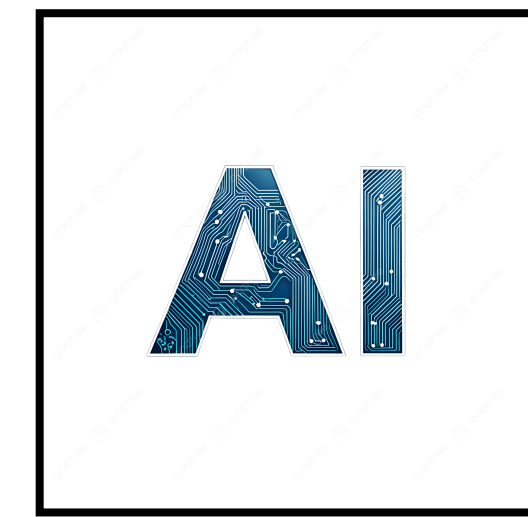
**Online Setting** - “Human Adaptation to AI”



**Classification**

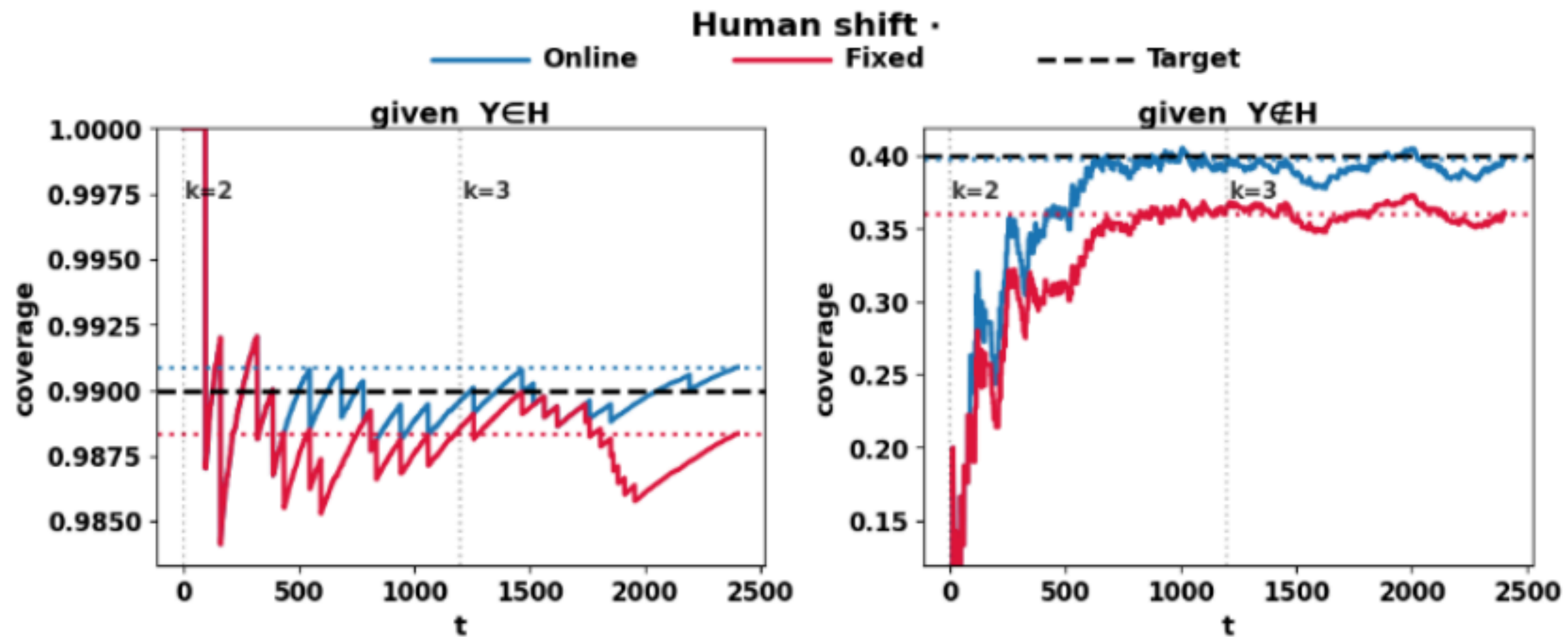


**Crowdsourcing**



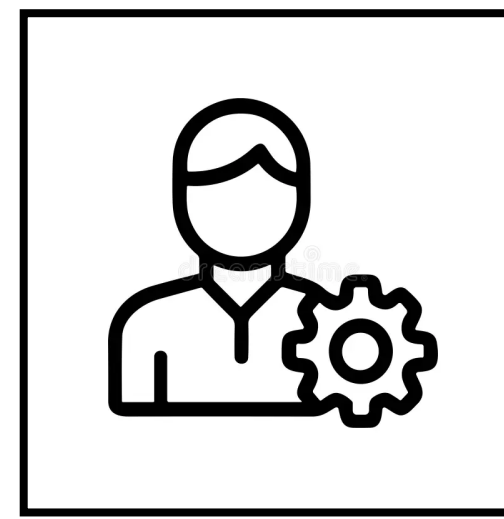
**AlexNet**

**Online Setting - “Human Adaptation to AI”**

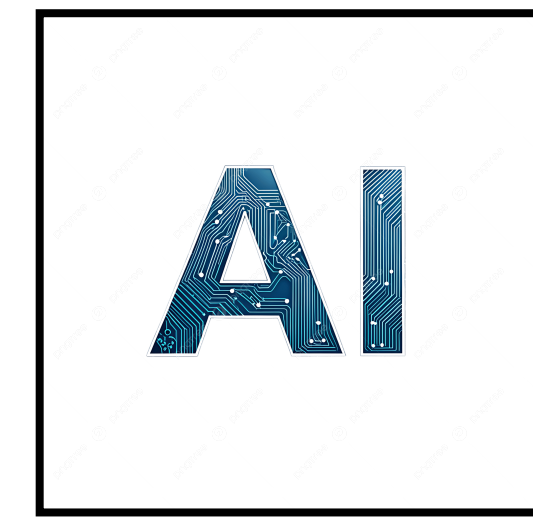




**Classification**

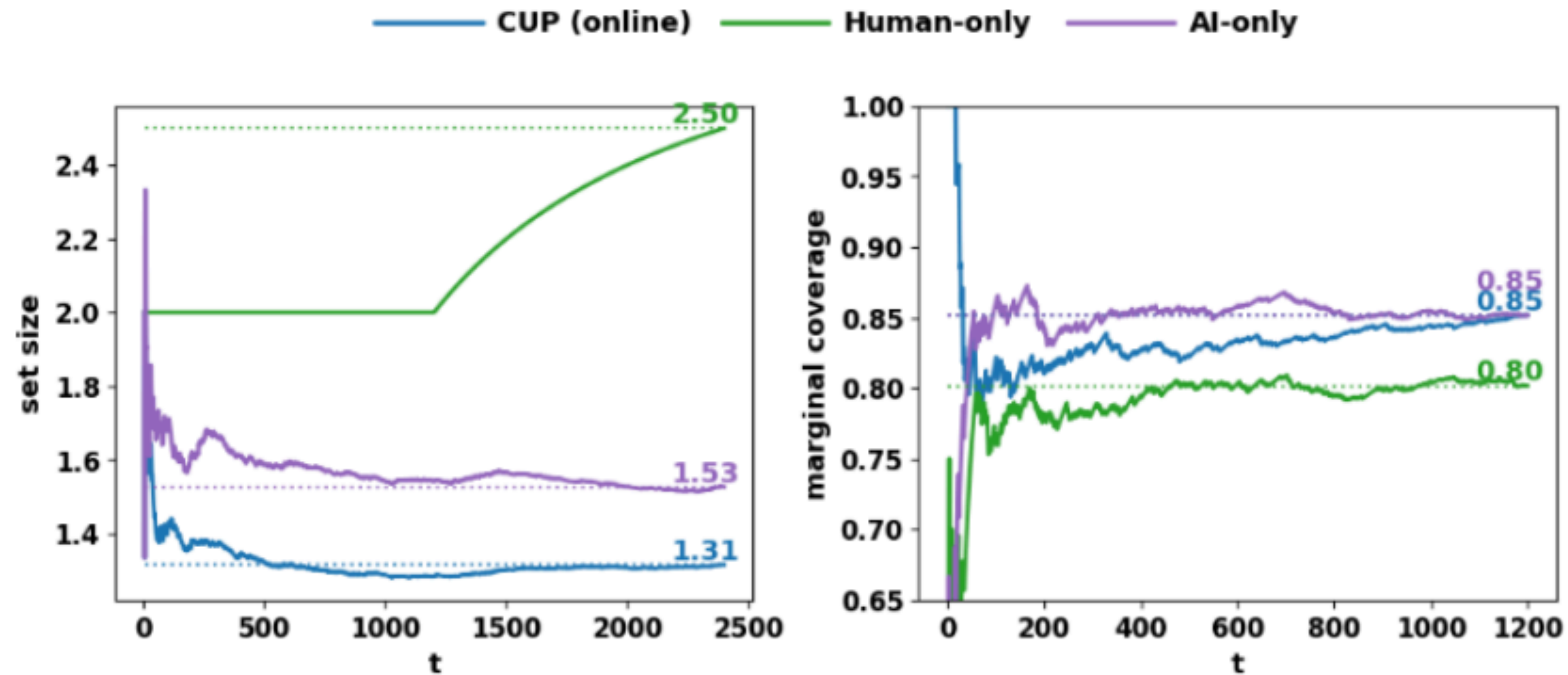


**Crowdsourcing**



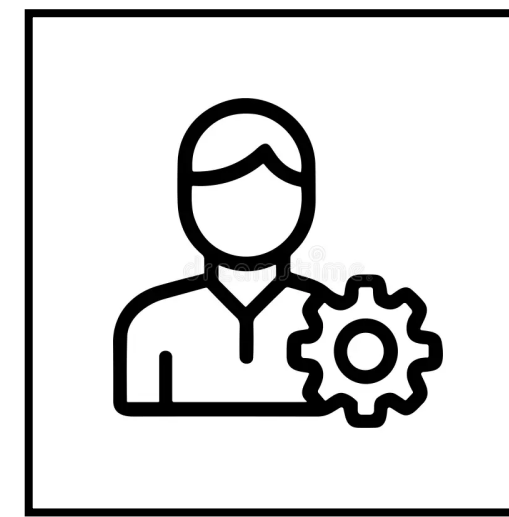
**AlexNet**

**Online Setting - “Human Adaptation to AI”**

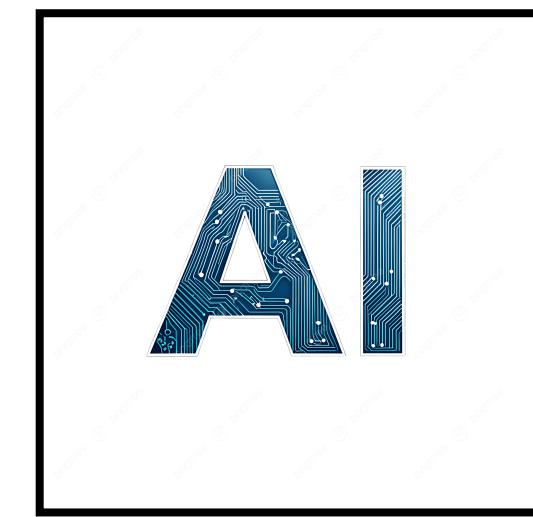




## **Text Based Medical Diagnosis**



## **Rule-based diagnostic system**

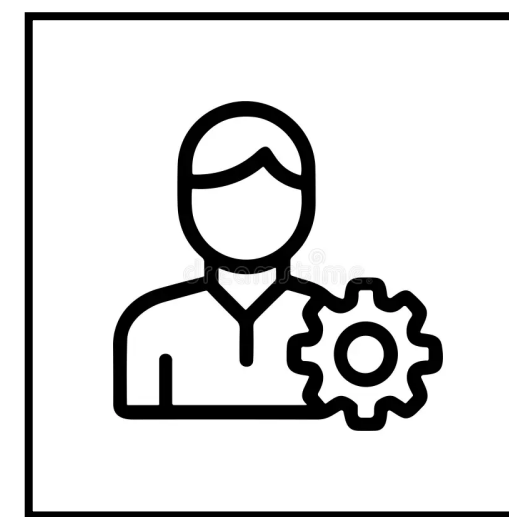


## **LLMs**

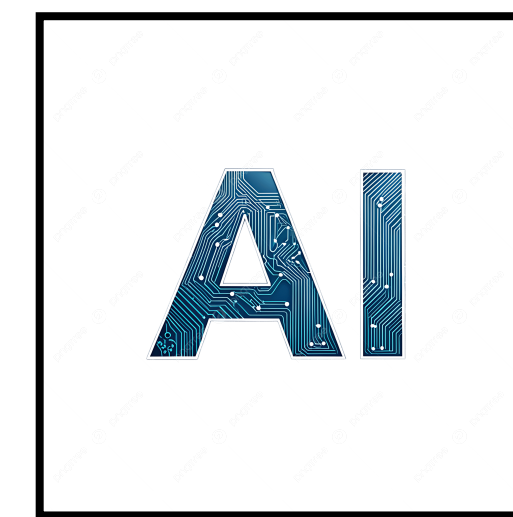
---



## **Text Based Medical Diagnosis**



## **Rule-based diagnostic system**



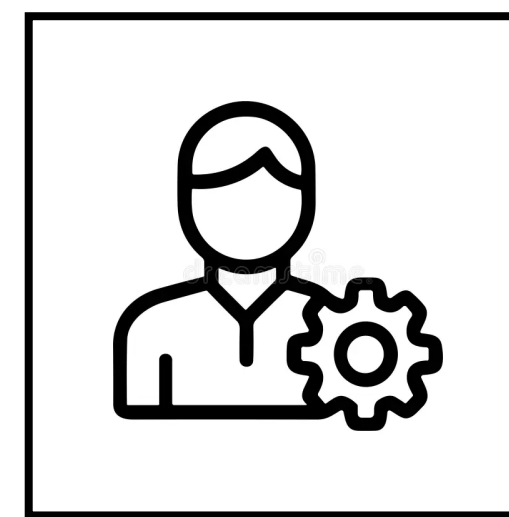
## **LLMs**

---

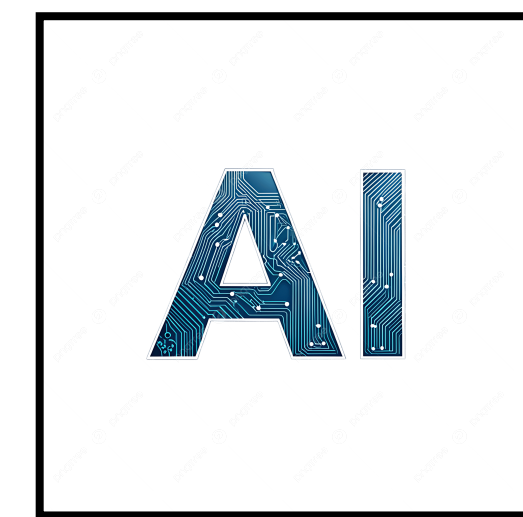
**Offline Setting**



## Text Based Medical Diagnosis



## Rule-based diagnostic system



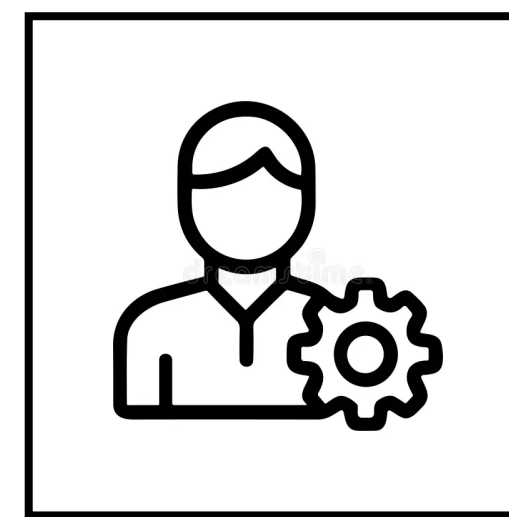
## LLMs

### Offline Setting

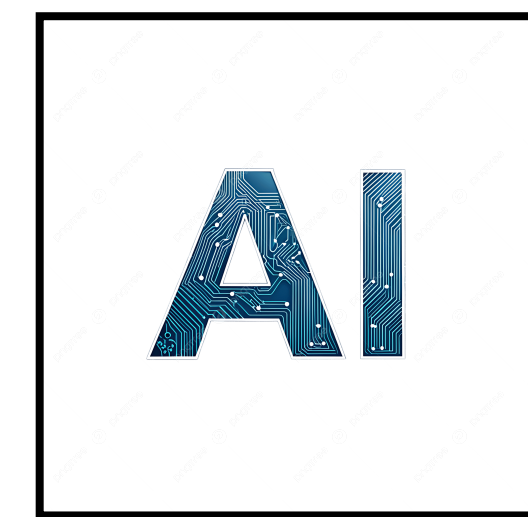
Strategy	Human	GPT-4o			GPT-5		
	C/S	CUP C/S	$(\epsilon, \delta)$	AI C/S	CUP C/S	$(\epsilon, \delta)$	AI C/S
Top-1	0.71 / 1.00	0.90 / 2.84	(0.02, 0.70)	0.88 / 4.64	0.91 / 1.59	(0.02, 0.70)	0.91 / 1.76
Top-2	0.87 / 1.95	0.93 / 3.14	(0.01, 0.45)	0.90 / 9.12	0.93 / 1.65	(0.02, 0.45)	0.93 / 1.95



## Text Based Medical Diagnosis



## Rule-based diagnostic system



## LLMs

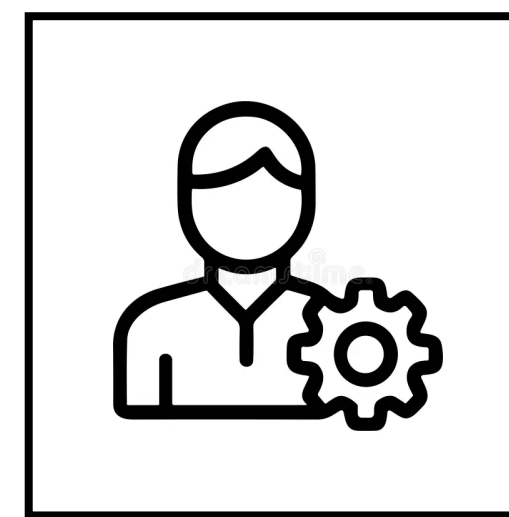
### Offline Setting

Strategy	Human	GPT-4o			GPT-5		
	C/S	CUP C/S	$(\epsilon, \delta)$	AI C/S	CUP C/S	$(\epsilon, \delta)$	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

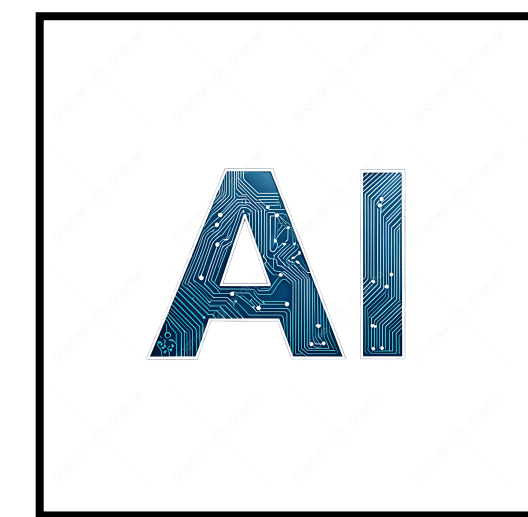
**AI's quality affects the overall collaboration quality!**



## **Text Based Medical Diagnosis**



## **Rule-based diagnostic system**

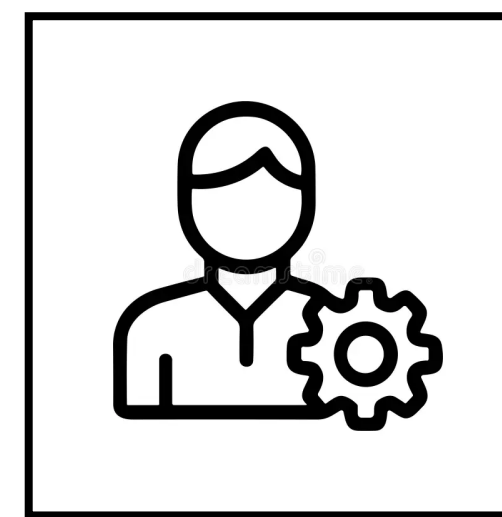


## **LLMs**

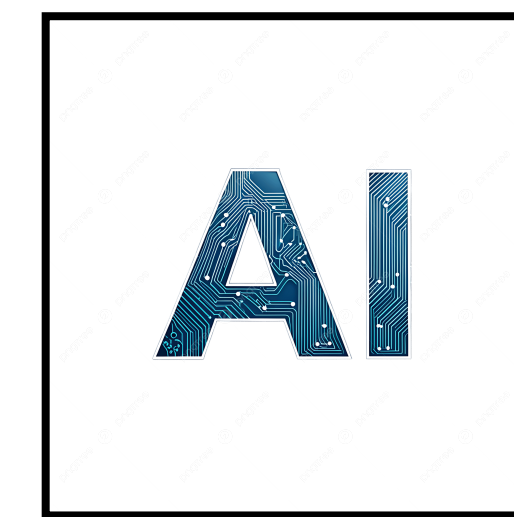
**Online Setting - "age shift"**



## Text Based Medical Diagnosis

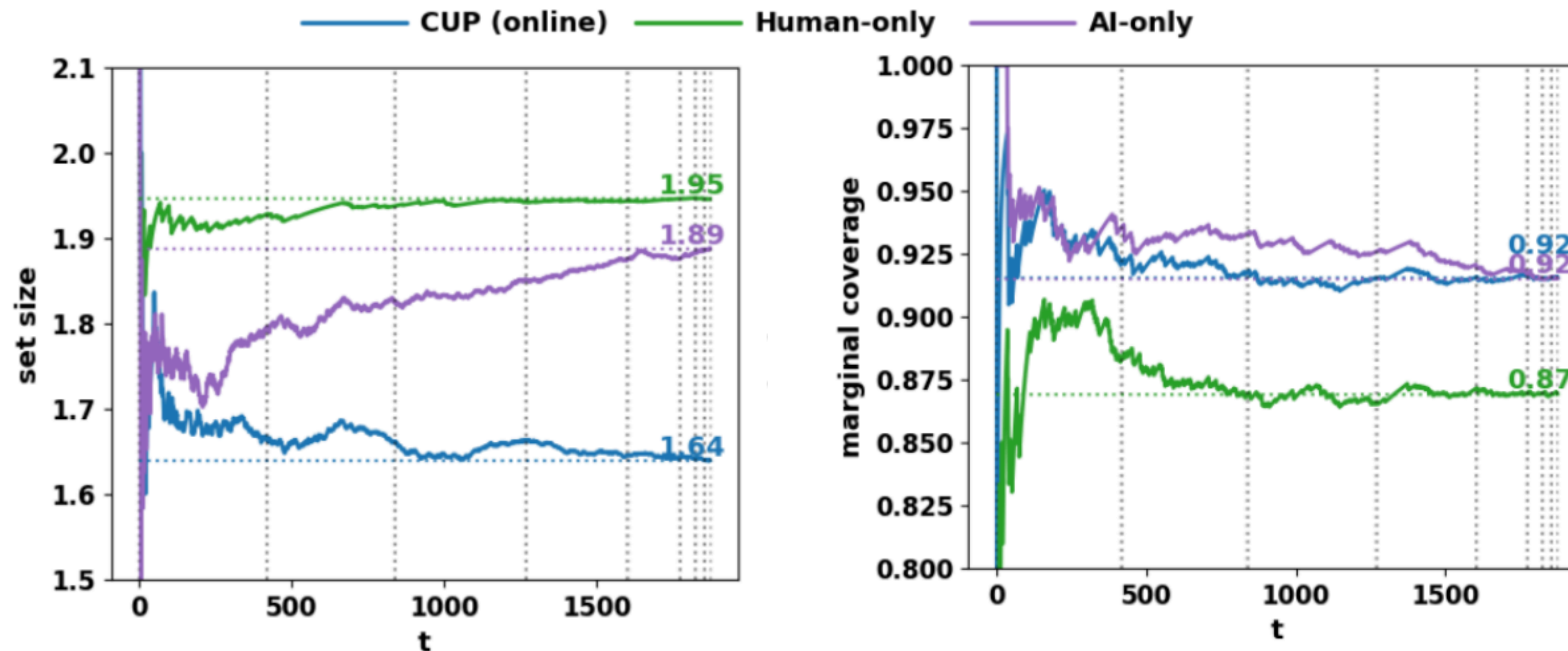


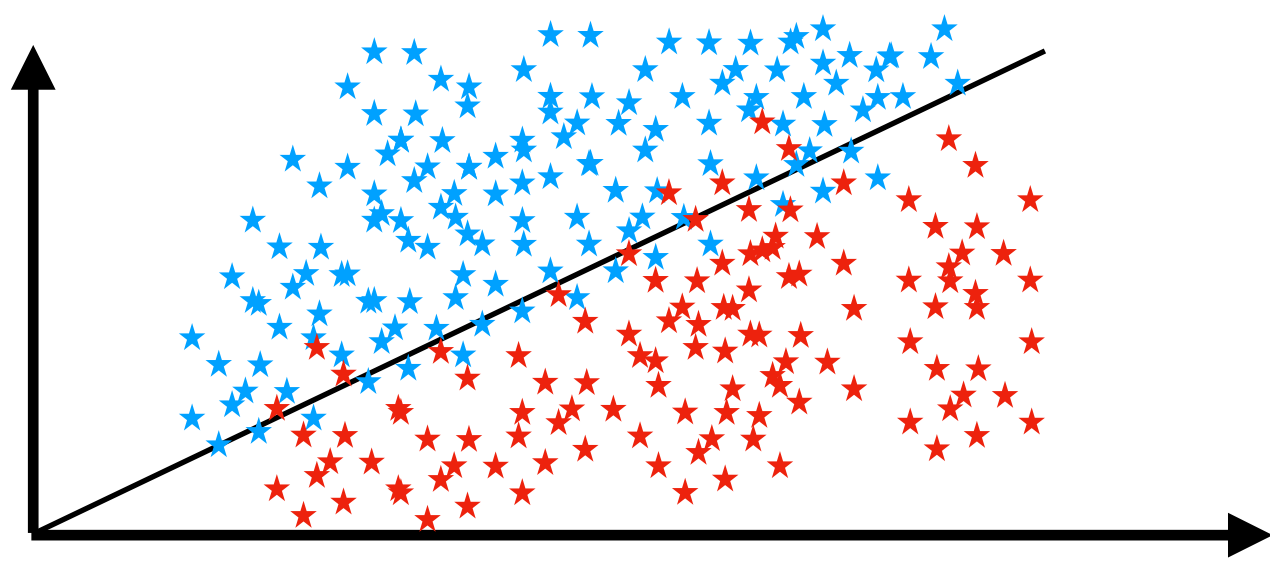
## Rule-based diagnostic system



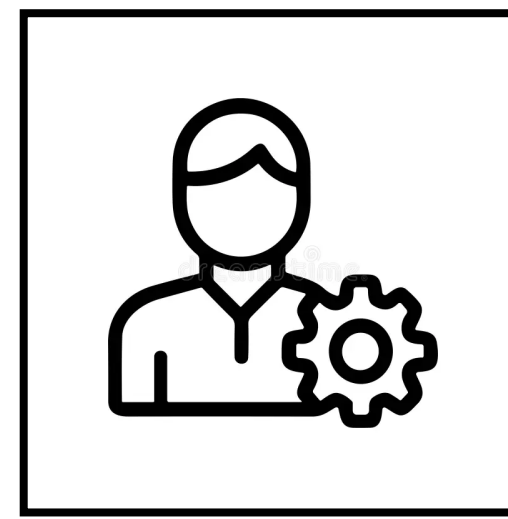
## LLMs

### Online Setting - "age shift"

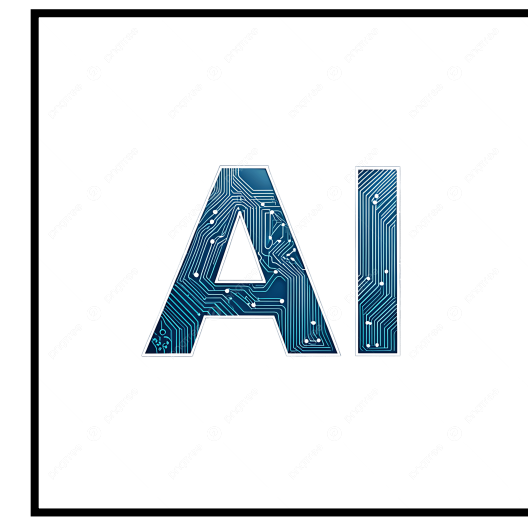




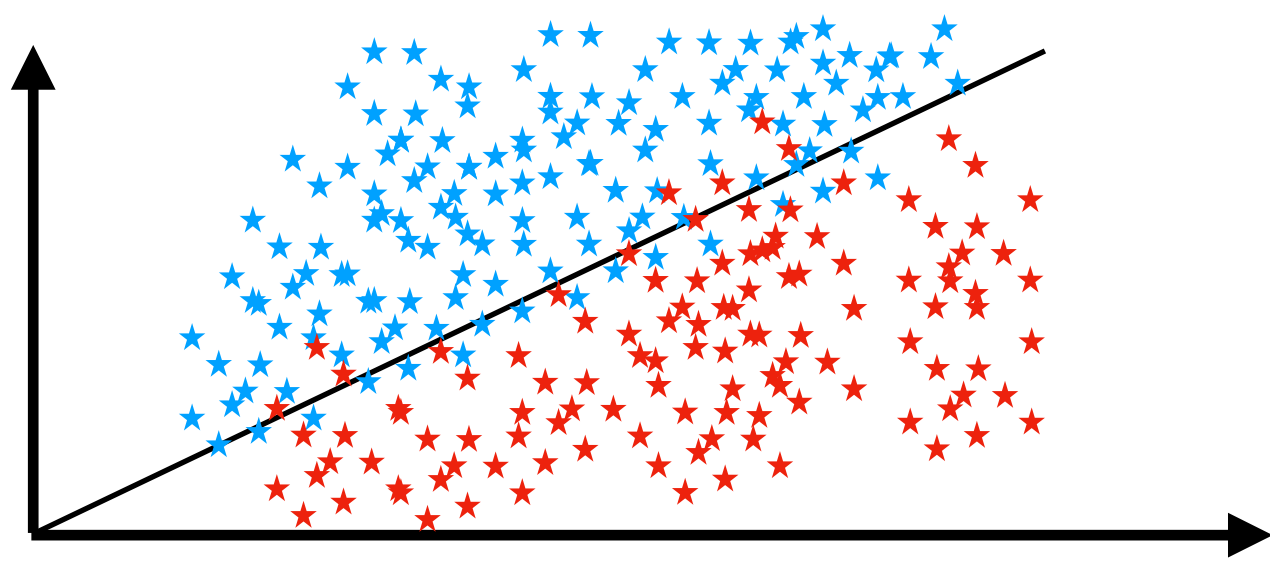
**Regression**



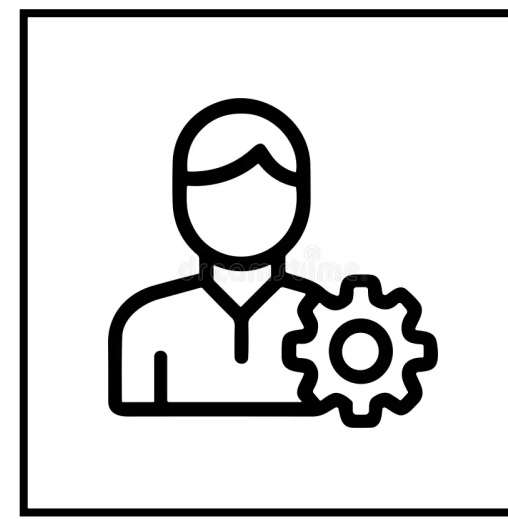
**Synthetic**



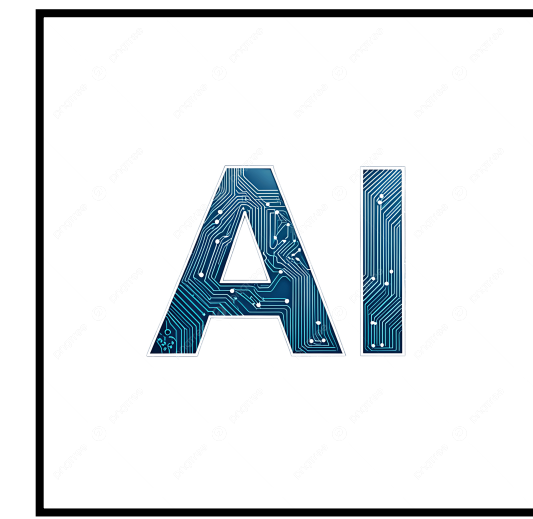
**MLP**



**Regression**



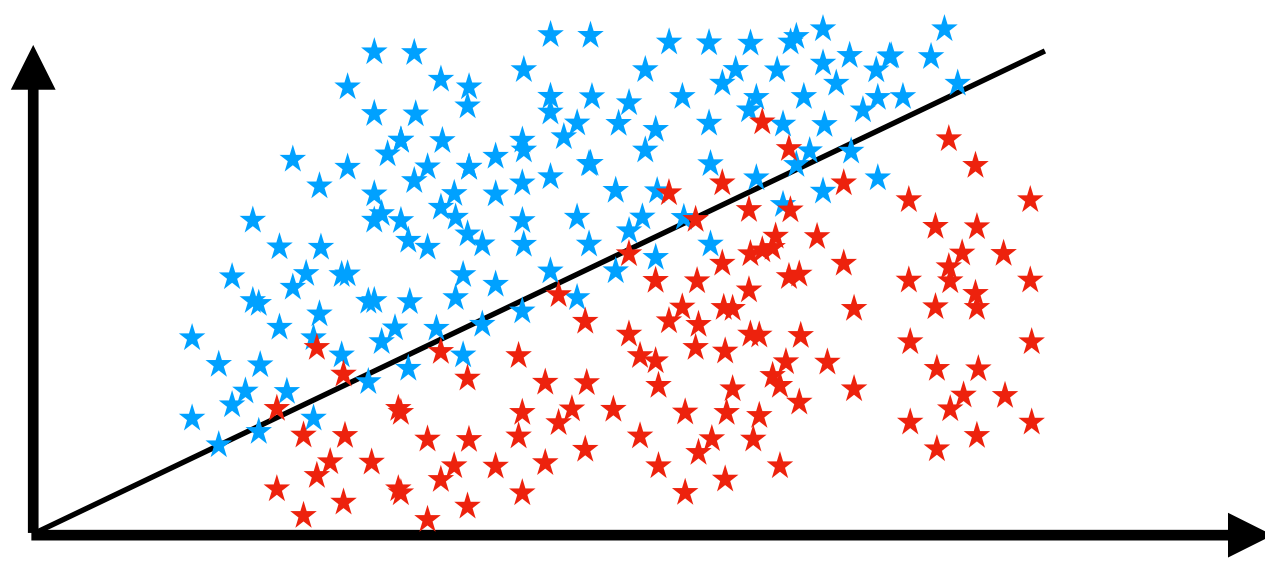
**Synthetic**



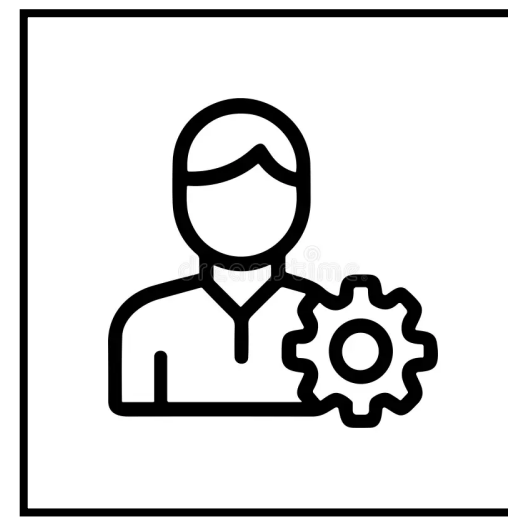
**MLP**

**Offline Setting**

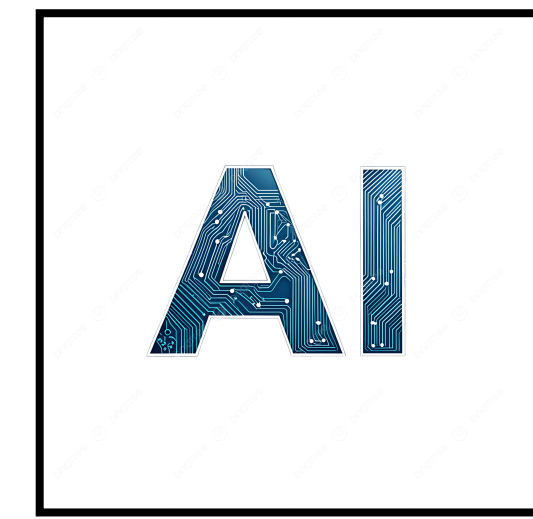
<b>Human A</b>				<b>Human B</b>			
<b>Human A C/S</b>	<b>CUP C/S</b>	$(\epsilon, 1 - \delta)$	<b>AI C/S</b>	<b>Human B C/S</b>	<b>CUP C/S</b>	$(\epsilon, 1 - \delta)$	<b>AI C/S</b>
0.760 / 0.581	0.862 / 0.380	(0.10, 0.70)	0.862 / 0.394	0.872 / 0.618	0.948 / 0.528	(0.05, 0.90)	0.948 / 0.608
0.760 / 0.581	0.825 / 0.326	(0.15, 0.70)	0.825 / 0.337	0.872 / 0.618	0.953 / 0.558	(0.05, 0.95)	0.953 / 0.588



**Regression**



**Synthetic**

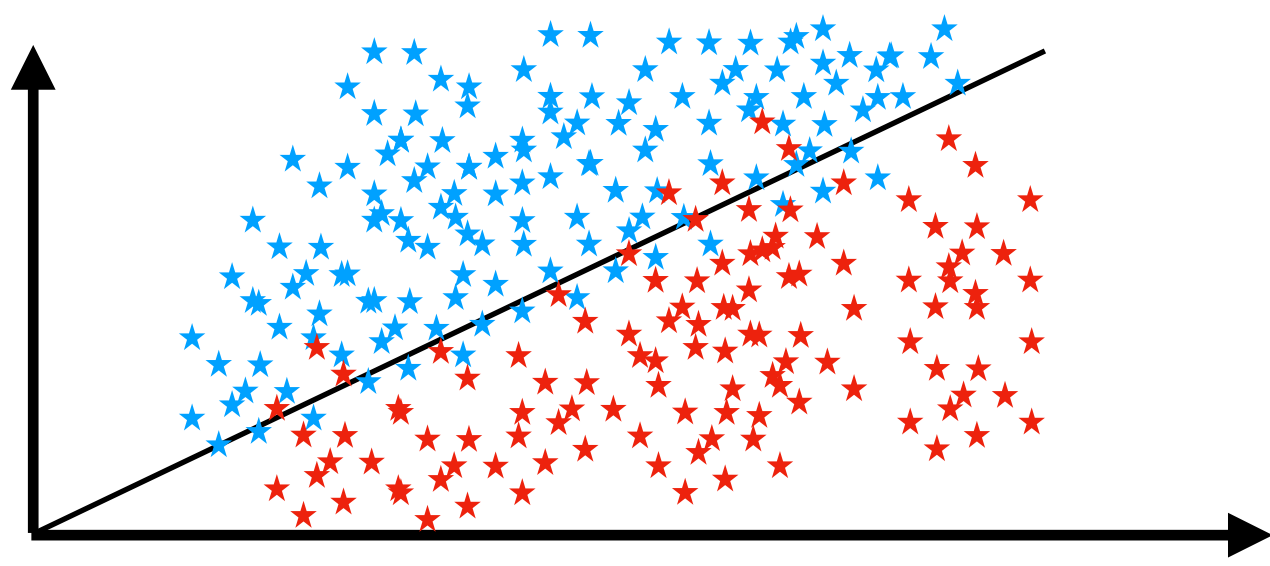


**MLP**

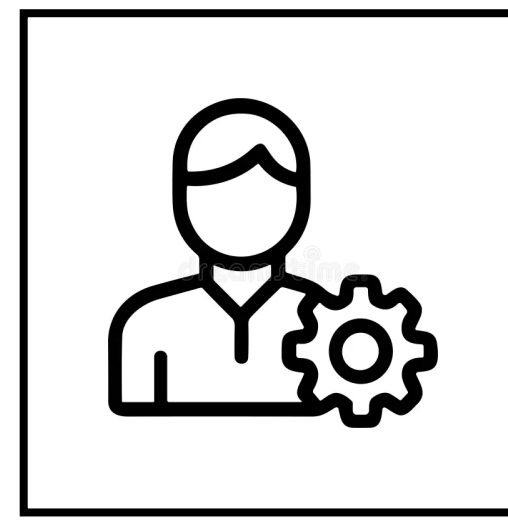
**Offline Setting**

<b>Human A</b>				<b>Human B</b>			
<b>Human A C/S</b>	<b>CUP C/S</b>	$(\epsilon, 1 - \delta)$	<b>AI C/S</b>	<b>Human B C/S</b>	<b>CUP C/S</b>	$(\epsilon, 1 - \delta)$	<b>AI C/S</b>
0.760 / 0.581	0.862 / 0.380	(0.10, 0.70)	0.862 / 0.394	0.872 / 0.618	0.948 / 0.528	(0.05, 0.90)	0.948 / 0.608
0.760 / 0.581	0.825 / 0.326	(0.15, 0.70)	0.825 / 0.337	0.872 / 0.618	0.953 / 0.558	(0.05, 0.95)	0.953 / 0.588

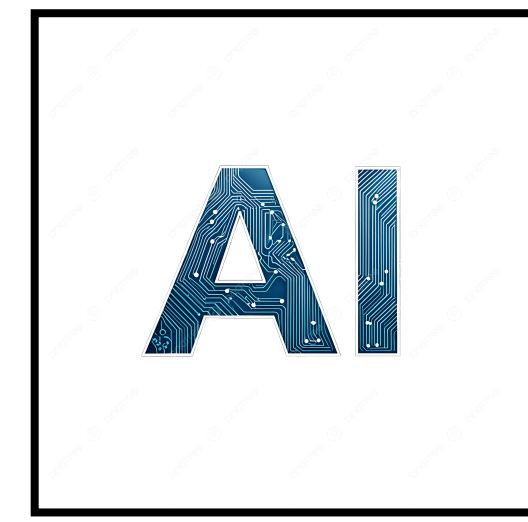
**Human's quality affects the overall collaboration quality!**



**Regression**



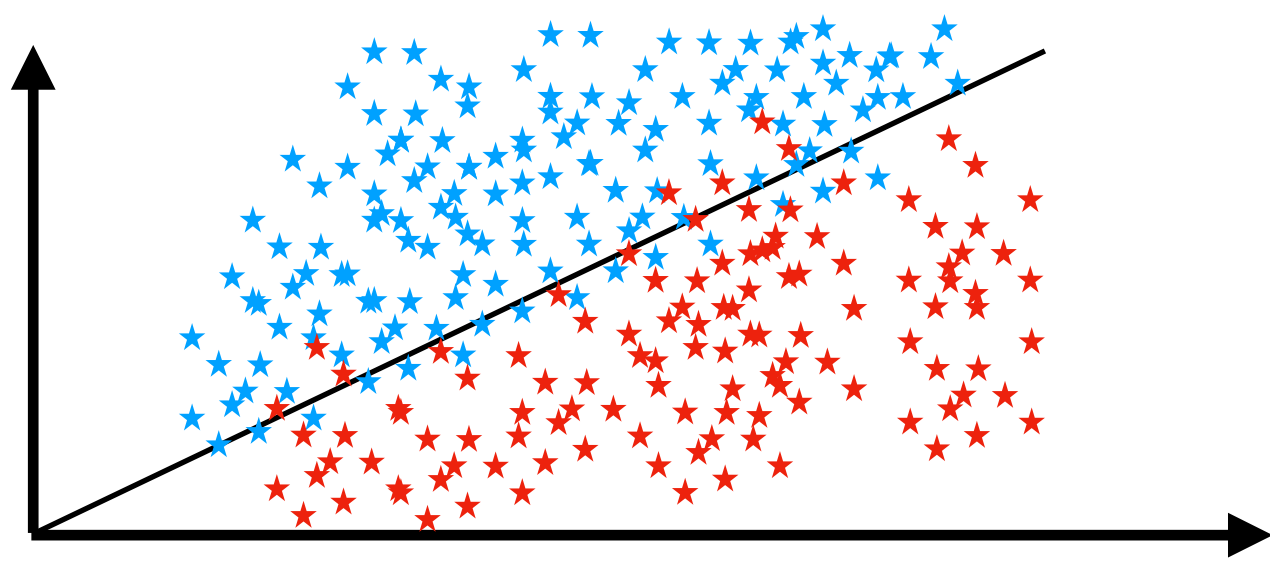
**Synthetic**



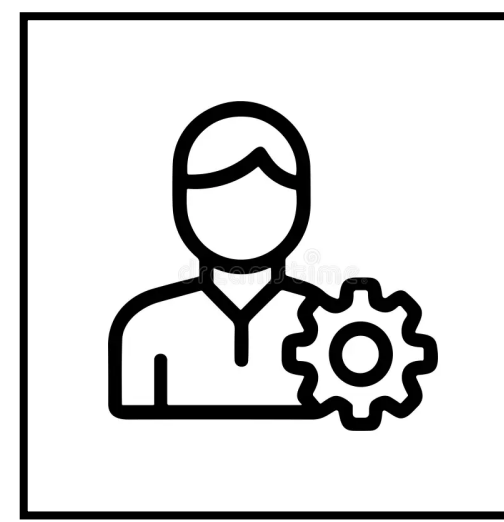
**MLP**

---

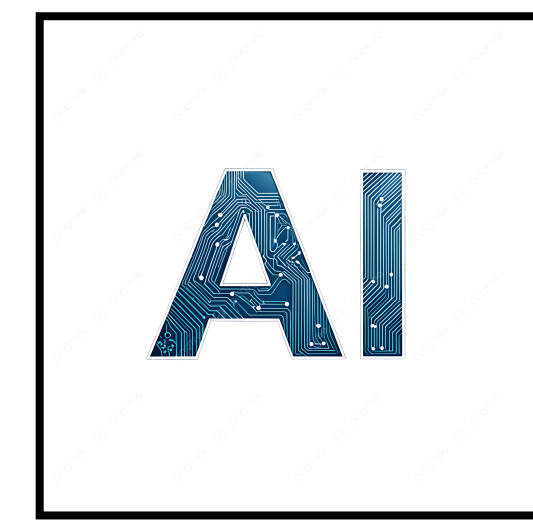
**Online Setting - “demographic shift”**



**Regression**

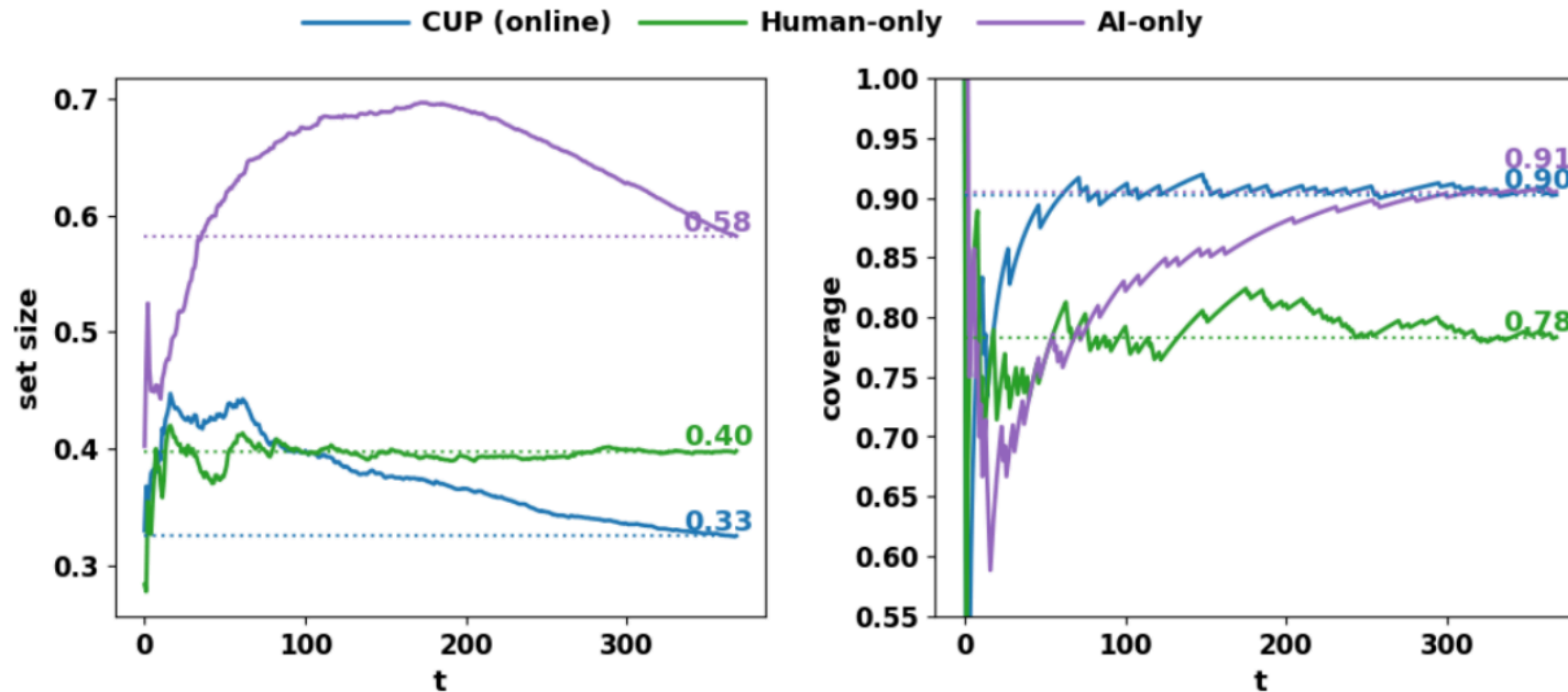


**Synthetic**



**MLP**

**Online Setting - “demographic shift”**



**Thank You!**