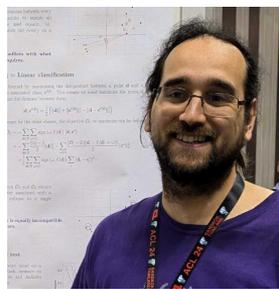


# Data complexity is not Data uncertainty



Aman Sinha<sup>1</sup>



Timothee Mickus<sup>2</sup>



Raul Vazquez<sup>2</sup>

<sup>1</sup>Université de Lorraine, Nancy, France

<sup>2</sup>University of Helsinki, Finland

# Bonjour! Hello! I'm Aman Sinha

- ▶ I'm a final year PhD Student at Université de Lorraine, Nancy and ICANS Strasbourg

# Bonjour! Hello! I'm Aman Sinha

- ▶ I'm a final year PhD Student at Université de Lorraine, Nancy and ICANS Strasbourg
- ▶ My thesis supervisors are:



**Marianne Clausel (IECL)**



**Mathieu Constant (ATILF)**



**Xavier Coubez (ICANS)**

# Bonjour! Hello! I'm Aman Sinha

- ▶ I'm a final year PhD Student at Université de Lorraine, Nancy and ICANS Strasbourg
- ▶ My thesis supervisors are:



Marianne Clausel (IECL)



Mathieu Constant (ATILF)



Xavier Coubez (ICANS)

- ▶ My thesis topic broadly related to **Medical language understanding**.

# Bonjour! Hello! I'm Aman Sinha

- ▶ I'm a final year PhD Student at Université de Lorraine, Nancy and ICANS Strasbourg
- ▶ My thesis supervisors are:



Marianne Clausel (IECL)



Mathieu Constant (ATILF)



Xavier Coubez (ICANS)

- ▶ My thesis topic broadly related to **Medical language understanding**.
  - ▶ medical data sources, multi-modality, and structures.

# Bonjour! Hello! I'm Aman Sinha

- ▶ I'm a final year PhD Student at Université de Lorraine, Nancy and ICANS Strasbourg
- ▶ My thesis supervisors are:



Marianne Clausel (IECL)



Mathieu Constant (ATILF)



Xavier Coubez (ICANS)

- ▶ My thesis topic broadly related to **Medical language understanding**.
  - ▶ medical data sources, multi-modality, and structures.
  - ▶ interpretability of medical models

# Bonjour! Hello! I'm Aman Sinha

- ▶ I'm a final year PhD Student at Université de Lorraine, Nancy and ICANS Strasbourg
- ▶ My thesis supervisors are:



Marianne Clausel (IECL)



Mathieu Constant (ATILF)



Xavier Coubez (ICANS)

- ▶ My thesis topic broadly related to **Medical language understanding**.
  - ▶ medical data sources, multi-modality, and structures.
  - ▶ interpretability of medical models
  - ▶ uncertainty & reliability of medical models.

# Bonjour! Hello! I'm Aman Sinha

- ▶ I'm a final year PhD Student at Université de Lorraine, Nancy and ICANS Strasbourg
- ▶ My thesis supervisors are:



Marianne Clausel (IECL)



Mathieu Constant (ATILF)



Xavier Coubez (ICANS)

- ▶ My thesis topic broadly related to **Medical language understanding**.
  - ▶ medical data sources, multi-modality, and structures.
  - ▶ interpretability of medical models
  - ▶ uncertainty & reliability of medical models.

# Bonjour! Hello! I'm Aman Sinha

- ▶ I'm a final year PhD Student at Université de Lorraine, Nancy and ICANS Strasbourg
- ▶ My thesis supervisors are: Marianne Clausel (IECL), Mathieu Constant (ATILF), and Xavier Coubez (ICANS)
- ▶ My thesis topic broadly related to **Medical language understanding**.

Today, I will be talking about the interaction of data uncertainty and data complexity.

# Bonjour! Hello! I'm Aman Sinha

- ▶ I'm a final year PhD Student at Université de Lorraine, Nancy and ICANS Strasbourg
- ▶ My thesis supervisors are: Marianne Clausel (IECL), Mathieu Constant (ATILF), and Xavier Coubez (ICANS)
- ▶ My thesis topic broadly related to **Medical language understanding**.

Today, I will be talking about the interaction of data uncertainty and data complexity.

## Your Model is Overconfident, and Other Lies We Tell Ourselves

**Timothee Mickus**        **Aman Sinha**         **Raúl Vázquez**   
 University of Helsinki       Université de Lorraine       ICANS Strasbourg  
firstname.lastname@<sup>1</sup>helsinki.fi,<sup>2</sup>univ-lorraine.fr

**Context**

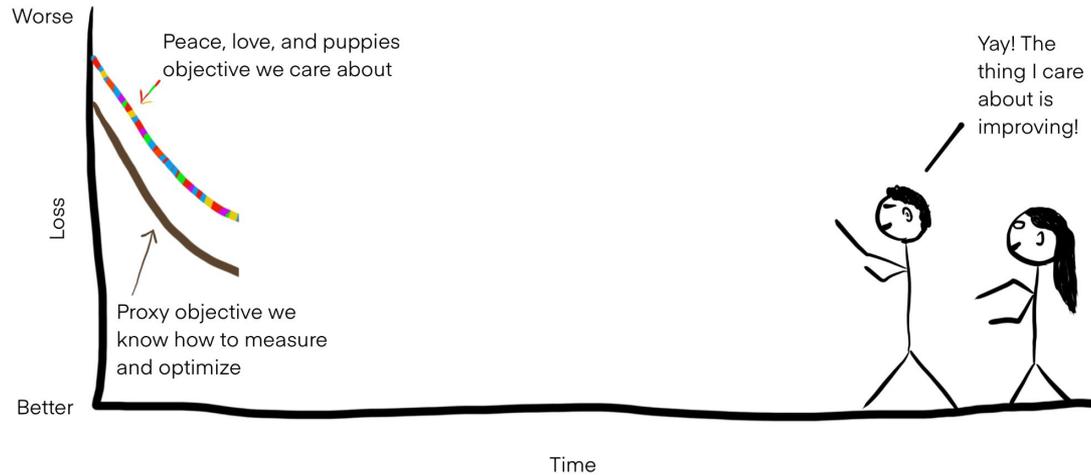
**Motivation**

---

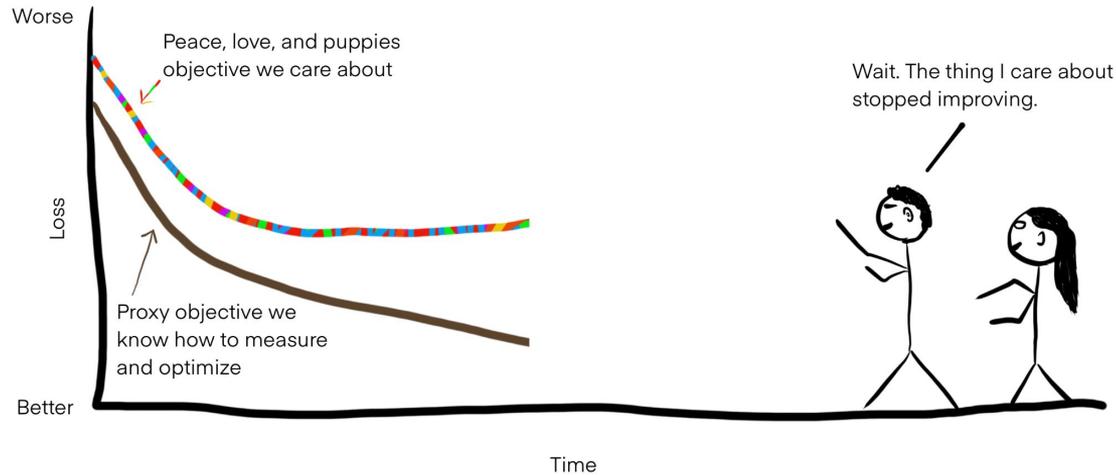




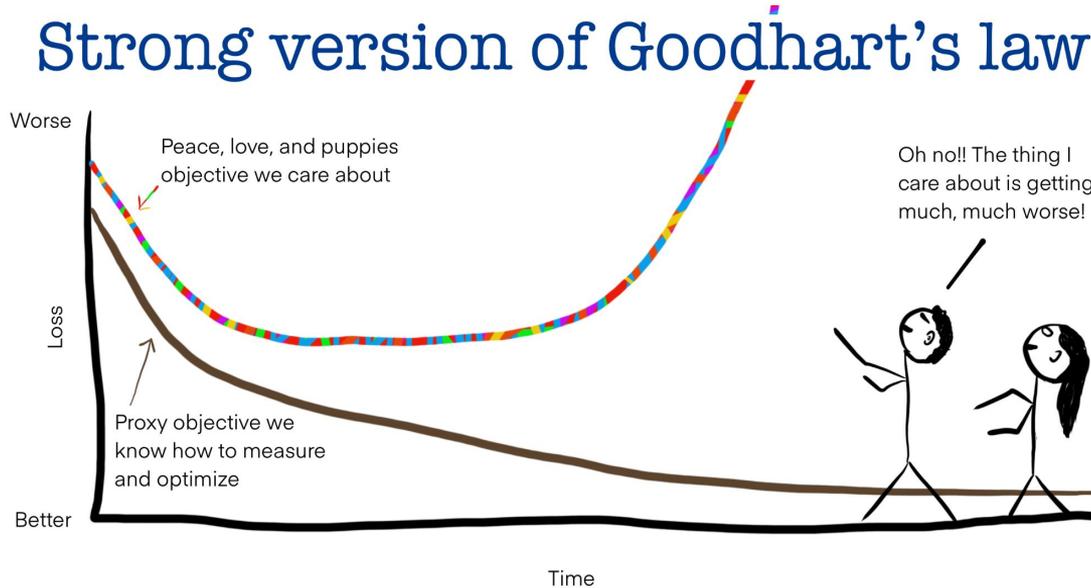
## Well-aligned phase



## Overfitting / Goodhart's law



## Strong version of Goodhart's law



*“When a measure becomes a target, if it is effectively optimized, then the thing it is designed to measure will grow worse.”*

- ▶ We expect any model to be good in terms of *performance*, *interpretation*, and *calibration*, as a wholesome behavior.

# Motivation - I

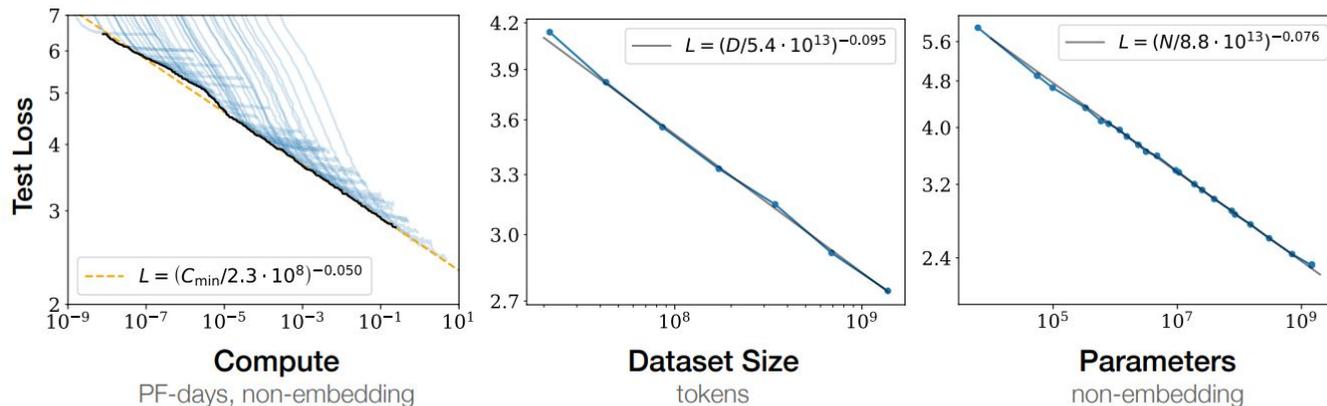
- ▶ We expect any model to be good in terms of *performance*, *interpretation*, and *calibration*, as a wholesome behavior.



From Greek mythology, Cerberus, often referred to as the hound of Hades

# Motivation - I

- ▶ We expect any model to be good in terms of *performance*, *interpretation*, and *calibration*, as a wholesome behavior.



**Figure 1** Language modeling performance improves smoothly as we increase the model size, dataset size, and amount of compute<sup>2</sup> used for training. For optimal performance all three factors must be scaled up in tandem. Empirical performance has a power-law relationship with each individual factor when not bottlenecked by the other two.

# Motivation - I

- ▶ We expect any model to be good in terms of *performance*, *interpretation*, and *calibration*, as a wholesome behavior.

Decision Tree  
(if/else)

LLMs



```
def recommend_activity_marseille(interest, preference):  
    if interest == "outdoors":  
        if preference == "land":  
            return "Calanques Hiking"  
        if preference == "sea":  
            return "Boat Tour"  
  
    if interest == "culture":  
        if preference == "modern":  
            return "MuCEM Museum"  
        if preference == "historic":  
            return "Old Port Walk"
```



# Motivation - I

- ▶ We expect any model to be good in terms of *performance*, *interpretation*, and *calibration*, as a wholesome behavior.

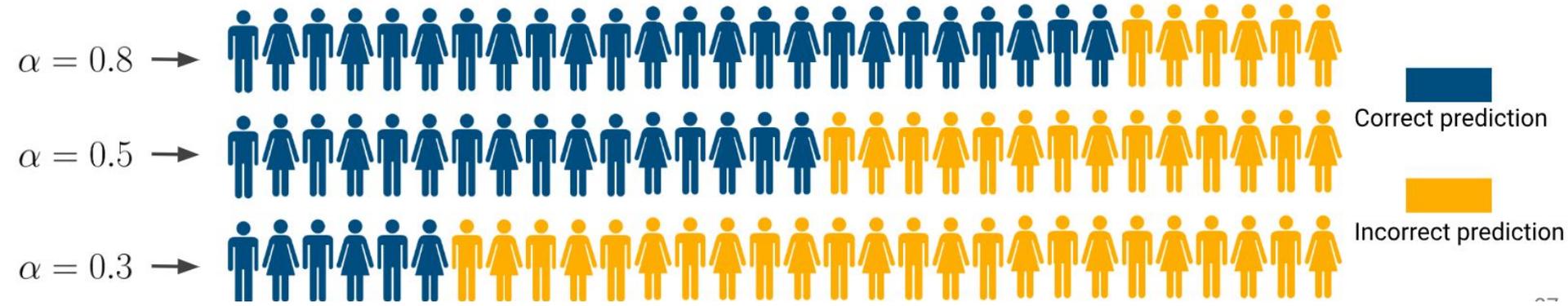
**What does being calibrated imply?**

# 💡 What does being *calibrated* imply?

► If a model is calibrated (aka **uncertainty-aware** or **reliable**),

$$P(\text{model is correct} \mid \text{confidence is } \alpha) = \alpha$$

It means,  $\alpha$ -fraction of the predicted classes with  $\alpha$  **confidence** should be correct.



# Motivation - I

- We expect any model to be good in terms of *performance*, *interpretation*, and *calibration*, as a wholesome behavior.

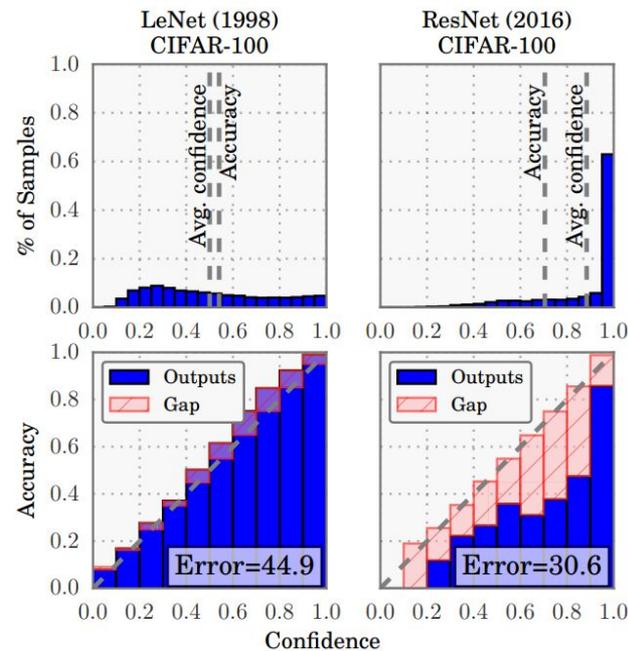


Figure 1. Confidence histograms (top) and reliability diagrams (bottom) for a 5-layer LeNet (left) and a 110-layer ResNet (right) on CIFAR-100. Refer to the text below for detailed illustration.

Modern neural networks exhibit a strange phenomenon: **probabilistic error and miscalibration worsen** even as classification error is reduced [Guo et al 2017].

# Motivation - I

- ▶ We expect any model to be good in terms of *performance*, *interpretation*, and *calibration*, as a wholesome behavior.



आईसीडी १० E११ का मुख्य कारण क्या है ?  
(Quelle est la cause principale du code CIM-10 E11 ?)

आईसीडी १० ई११ का मुख्य कारण पर्यावरणीय क्षरण है।  
(La cause principale du code CIM-10 E11 est la dégradation de l'environnement.)

Diabetes Type 2  
Classification and external resources

Universal blue circle symbol for diabetes. <sup>[1]</sup>

ICD- 10	E11 . <a href="#">↗</a>
ICD- 9	250.00 <a href="#">↗</a> , 250.02 <a href="#">↗</a>
OMIM	125853 <a href="#">↗</a>
Disease-DB	3661 <a href="#">↗</a>
Medline Plus	000313 <a href="#">↗</a>
eMedicine	article/117853 <a href="#">↗</a>
M.E.S.H.	D003924 <a href="#">↗</a>

- ▶ What makes an example unclear?

- ▶ What makes an example unclear?

Uncertainty  
quantification   $P(y^* | x, D) = \int P(y^* | x, \theta) P(\theta | D) d\theta$

- ▶ What makes an example unclear?

Uncertainty quantification  
$$P(y^* | x, D) = \int \underbrace{P(y^* | x, \Theta)}_{\text{data}} \underbrace{P(\Theta | D)}_{\text{model}} d\Theta$$

# Motivation - II

- ▶ What makes an example unclear?

Uncertainty quantification  $\rightarrow$   $P(y^* | x, D) = \int \underbrace{P(y^* | x, \Theta)}_{\text{data}} \underbrace{P(\Theta | D)}_{\text{model}} d\Theta$   $\rightarrow$  **intractable**

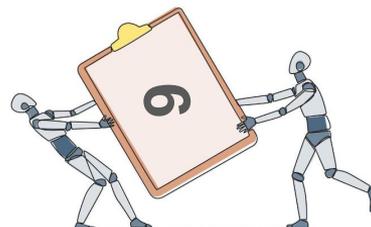
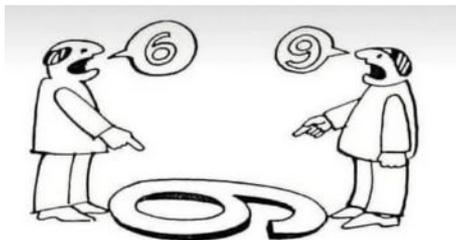
# Motivation - II

- ▶ What makes an example unclear?

Uncertainty quantification  $\rightarrow$   $P(y^* | x, D) = \int \underbrace{P(y^* | x, \Theta)}_{\text{data}} \underbrace{P(\Theta | D)}_{\text{model}} d\Theta$   $\rightarrow$  **intractable**

$\downarrow$

“inherent ambiguity or noise of the language”



- ▶ When is an example unclear?

► When is an example unclear?

to amongst others attention slips. Crucially, HLV assumes humans usually provide their best judgements, and variation emerges due to, e.g., ambiguity of the instance, uncertainty of the annotator, genuine disagreement, or simply the fact that multiple options are correct. Aggregation obfuscates this real-world complexity.

Plank et al., 2022

## ► When is an example unclear?

to amongst others attention slips. Crucially, HLV assumes humans usually provide their best judgements, and variation emerges due to, e.g., ambiguity of the instance, uncertainty of the annotator, genuine disagreement, or simply the fact that multiple options are correct. Aggregation obfuscates this real-world complexity.

Plank et al., 2022

- **Aleatoric Uncertainty** It is also known as data uncertainty, which refers to the uncertainty inherent in data due to its randomness or noise. This type of uncertainty is irreducible, meaning it cannot be eliminated through model improvements or tuning. It can arise from a variety of sources, such as noisy observations, overlapping classes, ground truth errors, inherent randomness, or other factors that are not entirely predictable.

Hu et al., 2023

## ► When is an example unclear?

to amongst others attention slips. Crucially, HLV assumes humans usually provide their best judgments, and variation emerges due to, e.g., ambiguity of the instance, uncertainty of the annotator, genuine disagreement, or simply the fact that multiple options are correct. Aggregation obfuscates this real-world complexity.

Plank et al., 2022

- **Aleatoric Uncertainty** It is also known as data uncertainty, which refers to the uncertainty inherent in data due to its randomness or noise. This type of uncertainty is irreducible, meaning it cannot be eliminated through model improvements or tuning. It can arise from a variety of sources, such as noisy observations, overlapping classes, ground truth errors, inherent randomness, or other factors that are not entirely predictable.

Hu et al., 2023

**Easy examples:** (Low  $PD_{Val.}$ , Low  $PD_{Train.}$ ). Such examples are often visually typical members of their class and the predicted label nearly always matches the ground truth.

**Looks like a different class:** (Low  $PD_{Val.}$ , High  $PD_{Train.}$ ). In the validation set, there is a clear (and nearly always incorrect) classification for such an input, but it is difficult to connect such inputs to other examples of their ground truth class during training. Mislabeled examples are of this kind, as are visually confusing images which at first appear to show something else.

**Ambiguous unless the label is given:** (High  $PD_{Val.}$ , Low  $PD_{Train.}$ ). These examples are difficult to connect to their predicted class in the validation split but easy to connect to their ground truth class during training. These points may, for example, visually resemble both their own class and another class. They are likely to be misclassified.

**Ambiguous:** (High  $PD_{Val.}$ , High  $PD_{Train.}$ ). These examples may be corrupted or show an example of a rare sub-class. Predictions for these inputs can depend strongly on the random seed used for training and initialization.

Baldock et al., 2021

## **Complexity, uncertainty, human variation**

---

# Complexity, uncertainty, human variation

- ▶ Same root causes

► Same root causes

Tab. 1 shows examples from *WinoGrande* belonging to the different regions defined above. *easy-to-learn* examples are straightforward for the model, as well as for humans. In contrast, most *hard-to-learn* and some *ambiguous* examples could be challenging for humans (see green highlights in Tab. 1), which might explain why the model shows lower **confidence** on them. These cate-

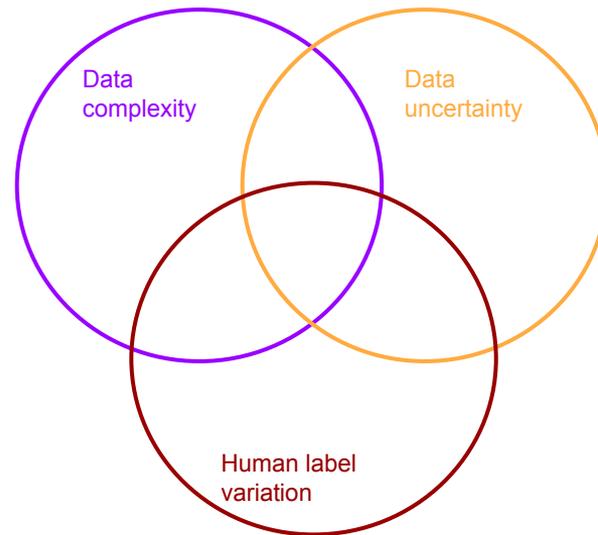
Swayamdipta et al., 2020

# Complexity, uncertainty, human variation

## ► Same root causes

Tab. 1 shows examples from *WinoGrande* belonging to the different regions defined above. *easy-to-learn* examples are straightforward for the model, as well as for humans. In contrast, most *hard-to-learn* and some *ambiguous* examples could be challenging for humans (see green highlights in Tab. 1), which might explain why the model shows lower **confidence** on them. These cate-

Swayamdipta et al., 2020

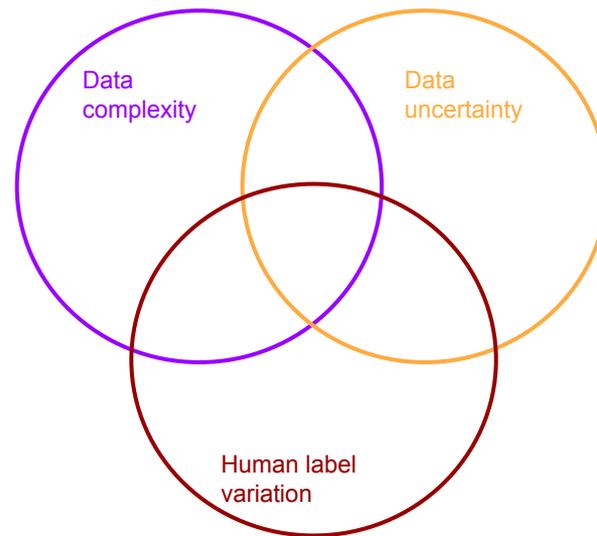


# Complexity, uncertainty, human variation

## ► Same root causes

Tab. 1 shows examples from *WinoGrande* belonging to the different regions defined above. *easy-to-learn* examples are straightforward for the model, as well as for humans. In contrast, most *hard-to-learn* and some *ambiguous* examples could be challenging for humans (see green highlights in Tab. 1), which might explain why the model shows lower **confidence** on them. These cate-

Swayamdipta et al., 2020



**Are these truly similar concepts?**

# Experimental Protocol

---

Context	Hypothesis	Old Labels majority and individual labels	New Labels	Source	Type
With the sun rising, a person is gliding with a huge parachute attached to them.	The person is falling to safety with the parachute	<b>Entailment</b> E E N N	<b>Entailment</b> E <sup>(50)</sup> N <sup>(50)</sup>	SNLI	Low agreements
A woman in a tan top and jeans is sitting on a bench wearing headphones.	A woman is listening to music.	<b>Entailment</b> E E N N E	<b>Neutral</b> N <sup>(93)</sup> E <sup>(7)</sup>	SNLI	Majority changed
A group of guys went out for a drink after work, and sitting at the bar was a real a 6 foot blonde with a fabulous face and figure to match.	The men didn't appreciate the figure of the blonde woman sitting at the bar.	<b>Contradiction</b> C N N C C	<b>Contradiction</b> C <sup>(56)</sup> N <sup>(44)</sup>	MNLI	Low agreements
In the other sight he saw Adrin's hands cocking back a pair of dragon-hammered pistols.	He had spotted Adrin preparing to fire his pistols.	<b>Neutral</b> N E N N E	<b>Entailment</b> E <sup>(94)</sup> N <sup>(5)</sup> C <sup>(1)</sup>	MNLI	Majority changed

Table 3: Examples from ChaosNLI-S and ChaosNLI-M development set. ‘Old Labels’ is the 5 label annotations from original dataset. ‘New Labels’ refers to the newly collected 100 label annotations. Superscript indicates the frequency of the label.

Context	Hypothesis	Old Labels majority and individual labels	New Labels	Source	Type
With the sun rising, a person is gliding with a huge parachute attached to them.	The person is falling to safety with the parachute	<b>Entailment</b> E E E N N	<b>Entailment</b> E <sup>(50)</sup> N <sup>(50)</sup>	SNLI	Low agreements
A woman in a tan top and jeans is sitting on a bench wearing headphones.	A woman is listening to music.	<b>Entailment</b> E E N N E	<b>Neutral</b> N <sup>(93)</sup> E <sup>(7)</sup>	SNLI	Majority changed
A group of guys went out for a drink after work, and sitting at the bar was a real a 6 foot blonde with a fabulous face and figure to match.	The men didn't appreciate the figure of the blonde woman sitting at the bar.	<b>Contradiction</b> C N N C C	<b>Contradiction</b> C <sup>(56)</sup> N <sup>(44)</sup>	MNLI	Low agreements
In the other sight he saw Adrin's hands cocking back a pair of dragon-hammered pistols.	He had spotted Adrin preparing to fire his pistols.	<b>Neutral</b> N E N N E	<b>Entailment</b> E <sup>(94)</sup> N <sup>(5)</sup> C <sup>(1)</sup>	MNLI	Majority changed

Table 3: Examples from ChaosNLI-S and ChaosNLI-M development set. ‘Old Labels’ is the 5 label annotations from original dataset. ‘New Labels’ refers to the newly collected 100 label annotations. Superscript indicates the frequency of the label.

Re-annotation of SNLI, MNLI,  $\alpha$ NLI

Context	Hypothesis	Old Labels majority and individual labels	New Labels	Source	Type
With the sun rising, a person is gliding with a huge parachute attached to them.	The person is falling to safety with the parachute	<b>Entailment</b> E E E N N	<b>Entailment</b> E <sup>(50)</sup> N <sup>(50)</sup>	SNLI	Low agreements
A woman in a tan top and jeans is sitting on a bench wearing headphones.	A woman is listening to music.	<b>Entailment</b> E E N N E	<b>Neutral</b> N <sup>(93)</sup> E <sup>(7)</sup>	SNLI	Majority changed
A group of guys went out for a drink after work, and sitting at the bar was a real a 6 foot blonde with a fabulous face and figure to match.	The men didn't appreciate the figure of the blonde woman sitting at the bar.	<b>Contradiction</b> C N N C C	<b>Contradiction</b> C <sup>(56)</sup> N <sup>(44)</sup>	MNLI	Low agreements
In the other sight he saw Adrin's hands cocking back a pair of dragon-hammered pistols.	He had spotted Adrin preparing to fire his pistols.	<b>Neutral</b> N E N N E	<b>Entailment</b> E <sup>(94)</sup> N <sup>(5)</sup> C <sup>(1)</sup>	MNLI	Majority changed

Table 3: Examples from ChaosNLI-S and ChaosNLI-M development set. ‘Old Labels’ is the 5 label annotations from original dataset. ‘New Labels’ refers to the newly collected 100 label annotations. Superscript indicates the frequency of the label.

Re-annotation of SNLI, MNLI,  $\alpha$ NLI

- 100 annotators per datapoint

Context	Hypothesis	Old Labels majority and individual labels	New Labels	Source	Type
With the sun rising, a person is gliding with a huge parachute attached to them.	The person is falling to safety with the parachute	<b>Entailment</b> E E N N	<b>Entailment</b> E <sup>(50)</sup> N <sup>(50)</sup>	SNLI	Low agreements
A woman in a tan top and jeans is sitting on a bench wearing headphones.	A woman is listening to music.	<b>Entailment</b> E E N N E	<b>Neutral</b> N <sup>(93)</sup> E <sup>(7)</sup>	SNLI	Majority changed
A group of guys went out for a drink after work, and sitting at the bar was a real a 6 foot blonde with a fabulous face and figure to match.	The men didn't appreciate the figure of the blonde woman sitting at the bar.	<b>Contradiction</b> C N N C C	<b>Contradiction</b> C <sup>(56)</sup> N <sup>(44)</sup>	MNLI	Low agreements
In the other sight he saw Adrin's hands cocking back a pair of dragon-hammered pistols.	He had spotted Adrin preparing to fire his pistols.	<b>Neutral</b> N E N N E	<b>Entailment</b> E <sup>(94)</sup> N <sup>(5)</sup> C <sup>(1)</sup>	MNLI	Majority changed

Table 3: Examples from ChaosNLI-S and ChaosNLI-M development set. ‘Old Labels’ is the 5 label annotations from original dataset. ‘New Labels’ refers to the newly collected 100 label annotations. Superscript indicates the frequency of the label.

## Re-annotation of SNLI, MNLI, $\alpha$ NLI

- ▶ 100 annotators per datapoint
- ▶ Provides distribution over labels

# Dataset : ChaosNLI

Context	Hypothesis	Old Labels majority and individual labels	New Labels	Source	Type
With the sun rising, a person is gliding with a huge parachute attached to them.	The person is falling to safety with the parachute	<b>Entailment</b> E E N N	<b>Entailment</b> E <sup>(50)</sup> N <sup>(50)</sup>	SNLI	Low agreements
A woman in a tan top and jeans is sitting on a bench wearing headphones.	A woman is listening to music.	<b>Entailment</b> E E N N E	<b>Neutral</b> N <sup>(93)</sup> E <sup>(7)</sup>	SNLI	Majority changed
A group of guys went out for a drink after work, and sitting at the bar was a real a 6 foot blonde with a fabulous face and figure to match.	The men didn't appreciate the figure of the blonde woman sitting at the bar.	<b>Contradiction</b> C N N C C	<b>Contradiction</b> C <sup>(56)</sup> N <sup>(44)</sup>	MNLI	Low agreements
In the other sight he saw Adrin's hands cocking back a pair of dragon-hammered pistols.	He had spotted Adrin preparing to fire his pistols.	<b>Neutral</b> N E N N E	<b>Entailment</b> E <sup>(94)</sup> N <sup>(5)</sup> C <sup>(1)</sup>	MNLI	Majority changed

Table 3: Examples from ChaosNLI-S and ChaosNLI-M development set. ‘Old Labels’ is the 5 label annotations from original dataset. ‘New Labels’ refers to the newly collected 100 label annotations. Superscript indicates the frequency of the label.

Re-annotation of SNLI, MNLI,  $\alpha$ NLI

- ▶ 100 annotators per datapoint
- ▶ Provides distribution over labels

We're considering SNLI

- ▶ Each time we consider a **pool** models of trained on NLI

- ▶ Each time we consider a **pool** models of trained on NLI

## <1B models

	$N_{\text{param}}$	$\theta_i$		$N_{\text{param}}$	$\theta_i$
BU_L-2_H-128_A-2	4385920	$m_1$	BU_L-2_H-512_A-8	22458880	$m_{13}$
BU_L-4_H-128_A-2	4782464	$m_2$	BU_L-4_H-512_A-8	28763648	$m_{14}$
BU_L-6_H-128_A-2	5179008	$m_3$	BU_L-6_H-512_A-8	35068416	$m_{15}$
BU_L-8_H-128_A-2	5575552	$m_4$	BU_L-2_H-768_A-12	38603520	$m_{16}$
BU_L-10_H-128_A-2	5972096	$m_5$	BU_L-8_H-512_A-8	41373184	$m_{17}$
BU_L-12_H-128_A-2	6368640	$m_6$	BU_L-10_H-512_A-8	47677952	$m_{18}$
BU_L-2_H-256_A-4	9591040	$m_7$	BU_L-4_H-768_A-12	52779264	$m_{19}$
BU_L-4_H-256_A-4	11170560	$m_8$	BU_L-12_H-512_A-8	53982720	$m_{20}$
BU_L-6_H-256_A-4	12750080	$m_9$	BU_L-6_H-768_A-12	66955008	$m_{21}$
BU_L-8_H-256_A-4	14329600	$m_{10}$	BU_L-8_H-768_A-12	81130752	$m_{22}$
BU_L-10_H-256_A-4	15909120	$m_{11}$	BU_L-10_H-768_A-12	95306496	$m_{23}$
BU_L-12_H-256_A-4	17488640	$m_{12}$	BU_L-12_H-768_A-12	109482240	$m_{24}$

- Each time we consider a **pool** models of trained on NLI

<1B models

	$N_{\text{param}}$	$\theta_i$		$N_{\text{param}}$	$\theta_i$
BU_L-2_H-128_A-2	4385920	$m_1$	BU_L-2_H-512_A-8	22458880	$m_{13}$
BU_L-4_H-128_A-2	4782464	$m_2$	BU_L-4_H-512_A-8	28763648	$m_{14}$
BU_L-6_H-128_A-2	5179008	$m_3$	BU_L-6_H-512_A-8	35068416	$m_{15}$
BU_L-8_H-128_A-2	5575552	$m_4$	BU_L-2_H-768_A-12	38603520	$m_{16}$
BU_L-10_H-128_A-2	5972096	$m_5$	BU_L-8_H-512_A-8	41373184	$m_{17}$
BU_L-12_H-128_A-2	6368640	$m_6$	BU_L-10_H-512_A-8	47677952	$m_{18}$
BU_L-2_H-256_A-4	9591040	$m_7$	BU_L-4_H-768_A-12	52779264	$m_{19}$
BU_L-4_H-256_A-4	11170560	$m_8$	BU_L-12_H-512_A-8	53982720	$m_{20}$
BU_L-6_H-256_A-4	12750080	$m_9$	BU_L-6_H-768_A-12	66955008	$m_{21}$
BU_L-8_H-256_A-4	14329600	$m_{10}$	BU_L-8_H-768_A-12	81130752	$m_{22}$
BU_L-10_H-256_A-4	15909120	$m_{11}$	BU_L-10_H-768_A-12	95306496	$m_{23}$
BU_L-12_H-256_A-4	17488640	$m_{12}$	BU_L-12_H-768_A-12	109482240	$m_{24}$

1B models



- Each time we consider a **pool** models of trained on NLI

<1B models

	$N_{\text{param}}$	$\theta_i$		$N_{\text{param}}$	$\theta_i$
BU_L-2_H-128_A-2	4385920	$m_1$	BU_L-2_H-512_A-8	22458880	$m_{13}$
BU_L-4_H-128_A-2	4782464	$m_2$	BU_L-4_H-512_A-8	28763648	$m_{14}$
BU_L-6_H-128_A-2	5179008	$m_3$	BU_L-6_H-512_A-8	35068416	$m_{15}$
BU_L-8_H-128_A-2	5575552	$m_4$	BU_L-2_H-768_A-12	38603520	$m_{16}$
BU_L-10_H-128_A-2	5972096	$m_5$	BU_L-8_H-512_A-8	41373184	$m_{17}$
BU_L-12_H-128_A-2	6368640	$m_6$	BU_L-10_H-512_A-8	47677952	$m_{18}$
BU_L-2_H-256_A-4	9591040	$m_7$	BU_L-4_H-768_A-12	52779264	$m_{19}$
BU_L-4_H-256_A-4	11170560	$m_8$	BU_L-12_H-512_A-8	53982720	$m_{20}$
BU_L-6_H-256_A-4	12750080	$m_9$	BU_L-6_H-768_A-12	66955008	$m_{21}$
BU_L-8_H-256_A-4	14329600	$m_{10}$	BU_L-8_H-768_A-12	81130752	$m_{22}$
BU_L-10_H-256_A-4	15909120	$m_{11}$	BU_L-10_H-768_A-12	95306496	$m_{23}$
BU_L-12_H-256_A-4	17488640	$m_{12}$	BU_L-12_H-768_A-12	109482240	$m_{24}$

1B models



Guarantees that we have pool of heterogeneous models: we can average out model quirks

- Each time we consider a **pool** models of trained on NLI

<1B models

	$N_{\text{param}}$	$\theta_i$		$N_{\text{param}}$	$\theta_i$
BU_L-2_H-128_A-2	4385920	$m_1$	BU_L-2_H-512_A-8	22458880	$m_{13}$
BU_L-4_H-128_A-2	4782464	$m_2$	BU_L-4_H-512_A-8	28763648	$m_{14}$
BU_L-6_H-128_A-2	5179008	$m_3$	BU_L-6_H-512_A-8	35068416	$m_{15}$
BU_L-8_H-128_A-2	5575552	$m_4$	BU_L-2_H-768_A-12	38603520	$m_{16}$
BU_L-10_H-128_A-2	5972096	$m_5$	BU_L-8_H-512_A-8	41373184	$m_{17}$
BU_L-12_H-128_A-2	6368640	$m_6$	BU_L-10_H-512_A-8	47677952	$m_{18}$
BU_L-2_H-256_A-4	9591040	$m_7$	BU_L-4_H-768_A-12	52779264	$m_{19}$
BU_L-4_H-256_A-4	11170560	$m_8$	BU_L-12_H-512_A-8	53982720	$m_{20}$
BU_L-6_H-256_A-4	12750080	$m_9$	BU_L-6_H-768_A-12	66955008	$m_{21}$
BU_L-8_H-256_A-4	14329600	$m_{10}$	BU_L-8_H-768_A-12	81130752	$m_{22}$
BU_L-10_H-256_A-4	15909120	$m_{11}$	BU_L-10_H-768_A-12	95306496	$m_{23}$
BU_L-12_H-256_A-4	17488640	$m_{12}$	BU_L-12_H-768_A-12	109482240	$m_{24}$

1B models



Guarantees that we have pool of heterogeneous models: we can average out model quirks

(let's keep it simple, and **we will talk about 1B models**)

- ▶ Contrasting values of metrics for human label variation / data difficulty / data uncertainty.

Human label variation	Data complexity	Data uncertainty
HUMAN-BASED	MODEL-BASED	
	<i>With reference</i>	<i>Without reference</i>
Entropy [Nie et al., 2020]	Early-exit [Baldock et al., 2021]	CP set size [Vovk et al., 2005]
Dissensus	Confidence [Swayamdipta et al., 2020]	Model entropy
	Early acquisition	Dissensus across model decisions
	Avg. acc. across models	Entropy across model decisions
	Avg. acc. across training	

- ▶ Contrasting values of metrics for human label variation / data difficulty / data uncertainty.

Human label variation	Data complexity	Data uncertainty
HUMAN-BASED	MODEL-BASED	
	<i>With reference</i>	<i>Without reference</i>
<b>Entropy</b> [Nie et al., 2020]	<b>Early-exit</b> [Baldock et al., 2021]	<b>CP set size</b> [Vovk et al., 2005]
Dissensus	Confidence [Swayamdipta et al., 2020]	Model entropy
	Early acquisition	Dissensus across model decisions
	Avg. acc. across models	Entropy across model decisions
	Avg. acc. across training	

Diversity in the label distribution, better accounting for both dominant and minority labels.

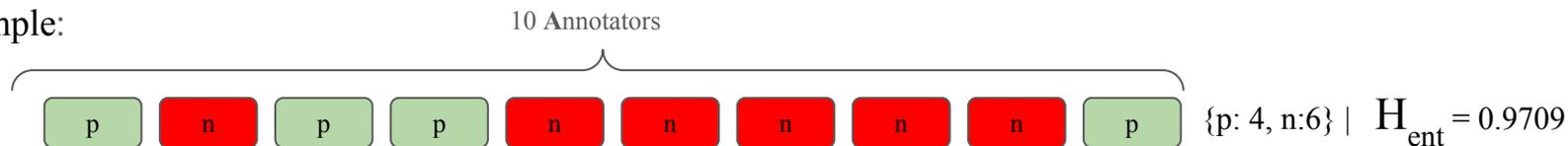
$$H_{\text{ent}} = - \sum_{y_i \in Y} \Pr_H(y_i | \mathbf{x}) \log \Pr_H(y_i | \mathbf{x})$$

**NB:** Human-based

Diversity in the label distribution, better accounting for both dominant and minority labels.

$$H_{\text{ent}} = - \sum_{y_i \in Y} \Pr_H(y_i | x) \log \Pr_H(y_i | x)$$

example:



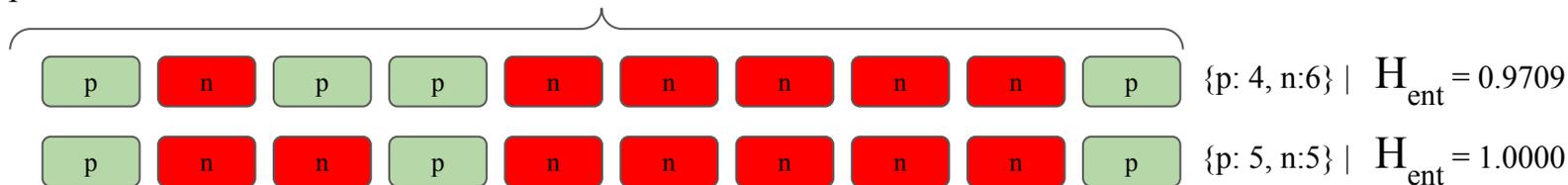
**NB:** Human-based

Diversity in the label distribution, better accounting for both dominant and minority labels.

$$H_{\text{ent}} = - \sum_{y_i \in Y} \Pr_H(y_i | x) \log \Pr_H(y_i | x)$$

example:

10 Annotators



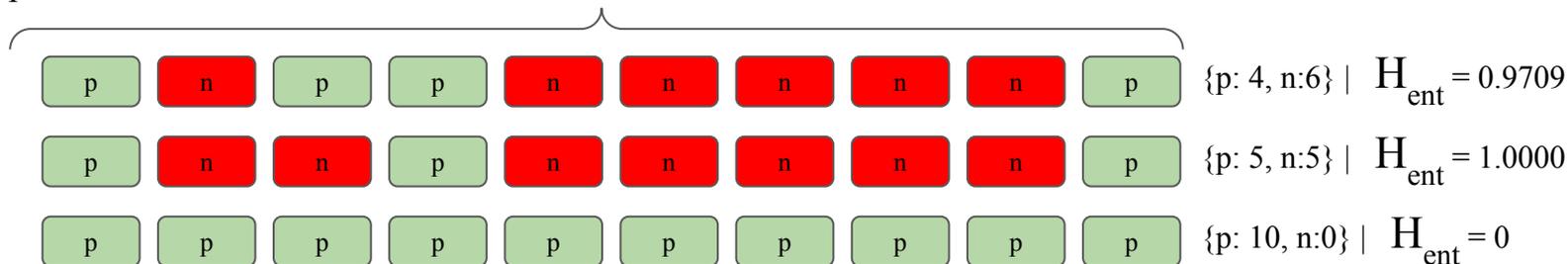
**NB:** Human-based

Diversity in the label distribution, better accounting for both dominant and minority labels.

$$H_{\text{ent}} = - \sum_{y_i \in Y} \Pr_H(y_i | x) \log \Pr_H(y_i | x)$$

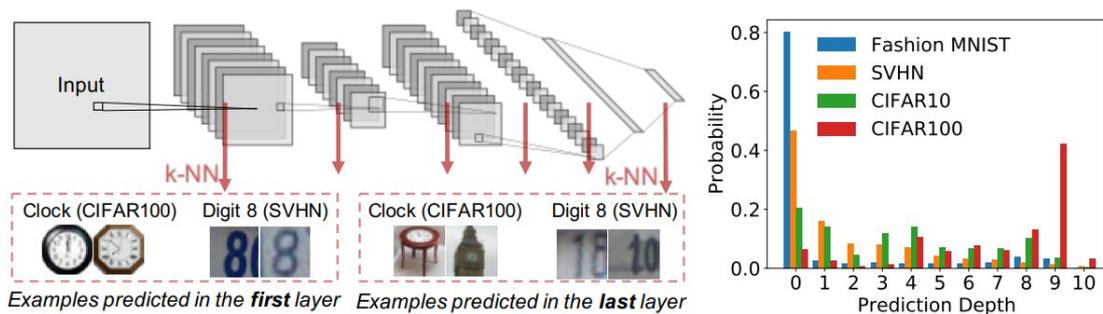
example:

10 Annotators



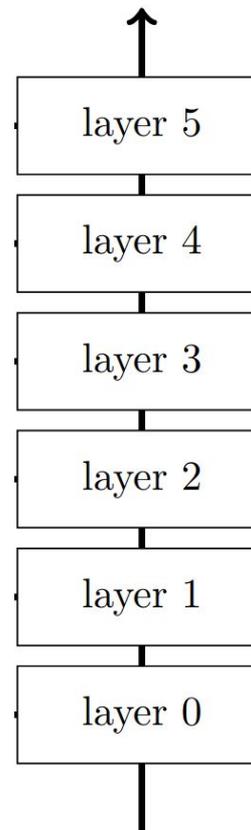
**NB:** Human-based

# Early-exit [Baldock et al., 2021]



*“Deep models use fewer layers to (effectively) determine the prediction for easy examples and more layers for hard examples”*

Given a deep learning model with layers of the same shape

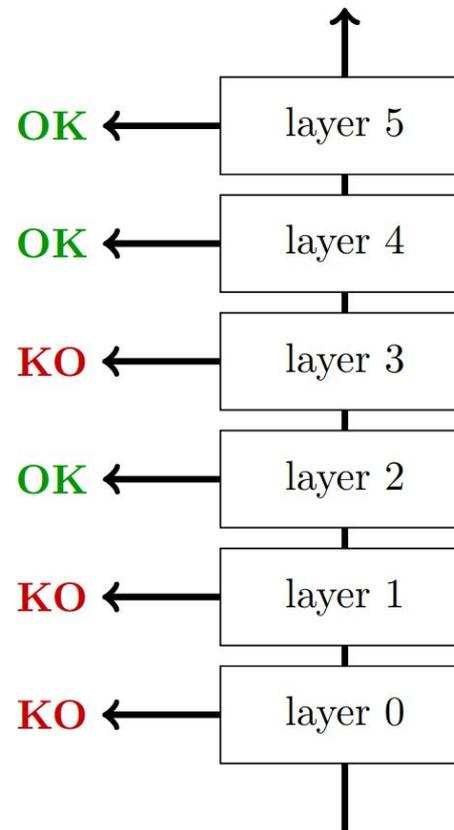


**NB:** Model-based, uses a reference

Given a deep learning model with layers of the same shape

- ▶ get the predictions at every layer

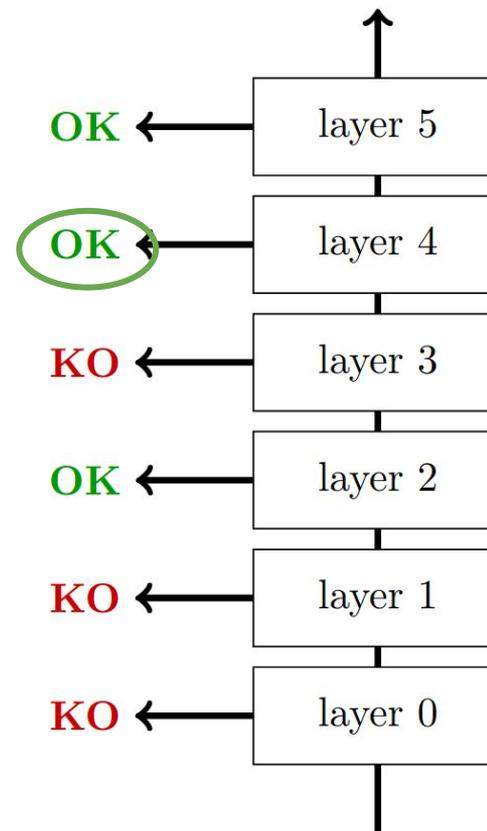
**NB:** Model-based, uses a reference



Given a deep learning model with layers of the same shape

- ▶ get the predictions at every layer
- ▶ select the layer where you start getting consistently correct predictions

**NB:** Model-based, uses a reference

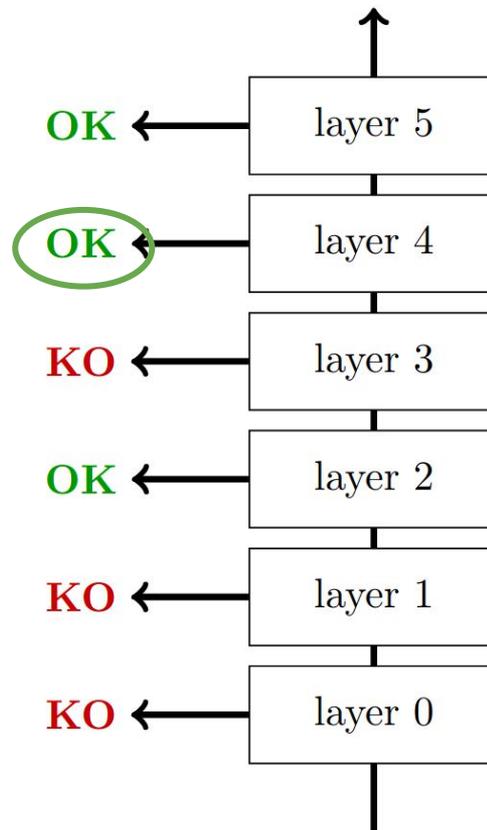


Given a deep learning model with layers of the same shape

- ▶ get the predictions at every layer
- ▶ select the layer where you start getting consistently correct predictions

Examples for which this layer is lower are *easier*

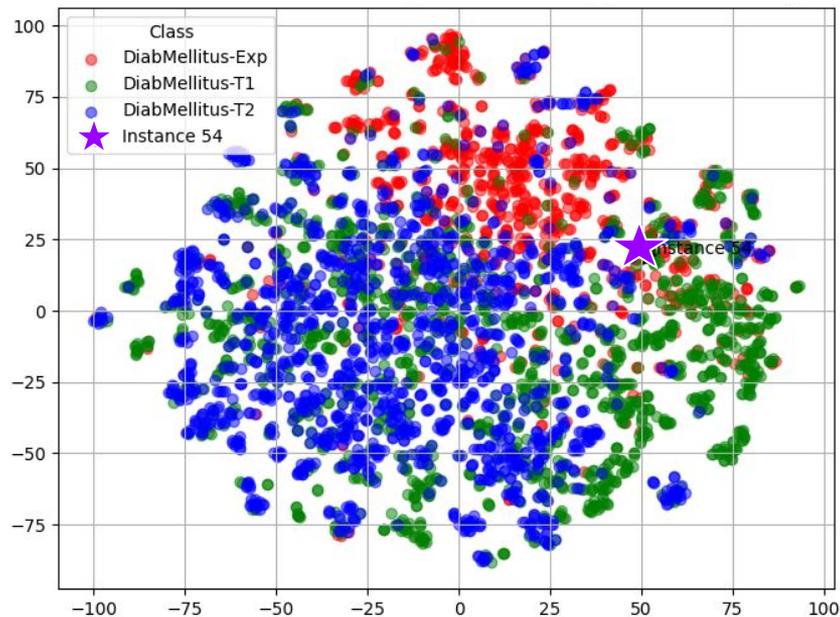
**NB:** Model-based, uses a reference



Conformal prediction considers the set of labels one must select to statistically guarantee that they contain the right answer (with a given risk)

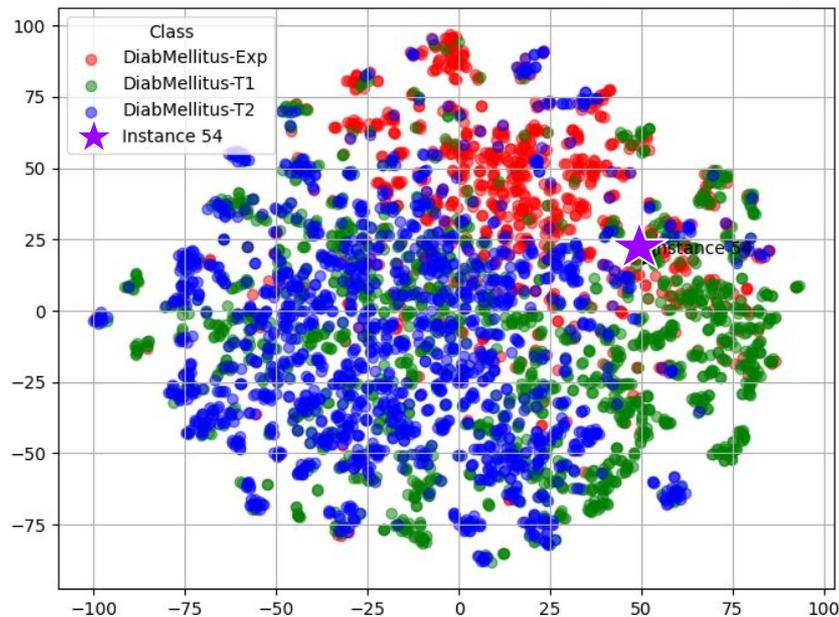
Conformal prediction considers the set of labels one must select to statistically guarantee that they contain the right answer (with a given risk)

example:



Conformal prediction considers the set of labels one must select to statistically guarantee that they contain the right answer (with a given risk)

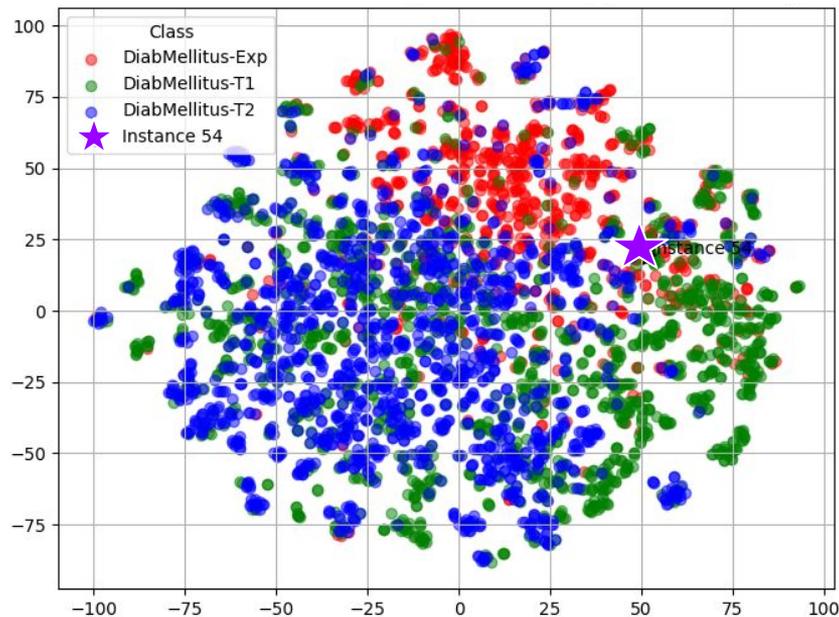
example:



Is that a Diab-Exp? Is that a Diab-T1?  
Is my classifier confident enough to  
decide on a single answer?  
How ok am I with my classifier being  
wrong?

Conformal prediction considers the set of labels one must select to statistically guarantee that they contain the right answer (with a given risk)

example:

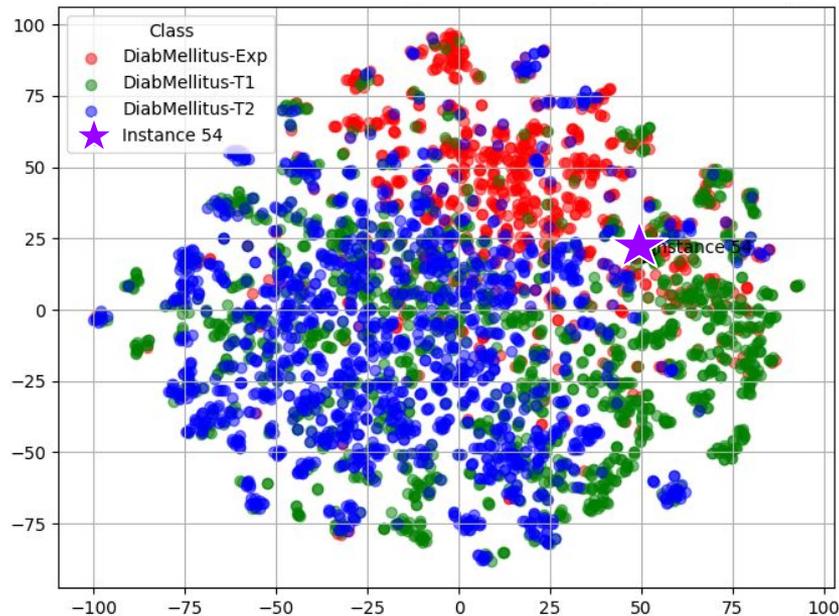


In practice:

$$\hat{P}(Y|\star) \rightarrow \begin{matrix} \text{Diab}_{\text{Exp}}, & \text{Diab}_{\text{T1}}, & \text{Diab}_{\text{T2}} \\ 0.45 & 0.35 & 0.20 \end{matrix}$$

Conformal prediction considers the set of labels one must select to statistically guarantee that they contain the right answer (with a given risk)

example:



In practice:

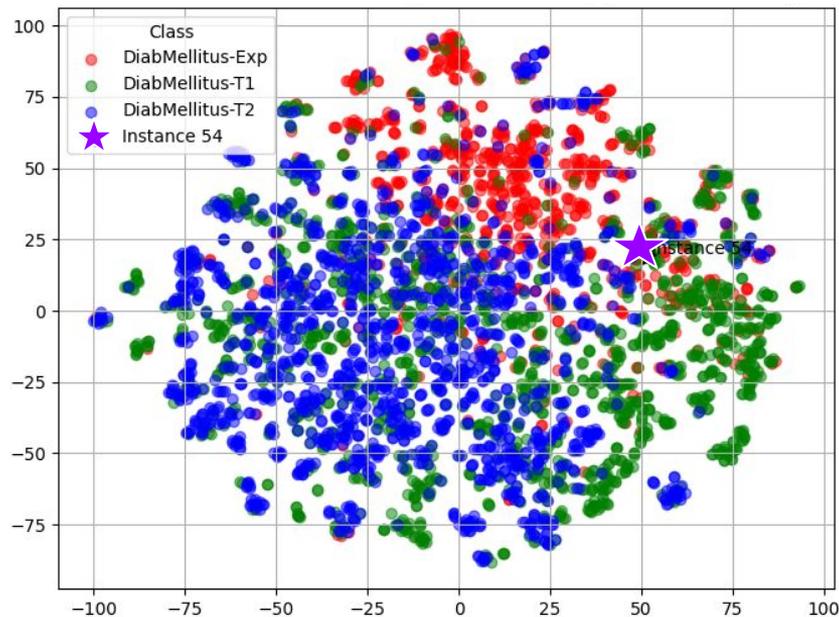
$$\hat{P}(Y|\star) \rightarrow \begin{matrix} \text{Diab}_{\text{Exp}}, & \text{Diab}_{\text{T1}}, & \text{Diab}_{\text{T2}} \\ 0.45 & 0.35 & 0.20 \end{matrix}$$

↓

$$U(\{\text{Diab}_{\text{Exp}}\}, \hat{P}, u_{f1}) = 0.5$$

Conformal prediction considers the set of labels one must select to statistically guarantee that they contain the right answer (with a given risk)

example:



In practice:

$$\hat{P}(Y|\star) \rightarrow \begin{matrix} \text{Diab}_{\text{Exp}}, & \text{Diab}_{\text{T1}}, & \text{Diab}_{\text{T2}} \\ 0.45 & 0.35 & 0.20 \end{matrix}$$

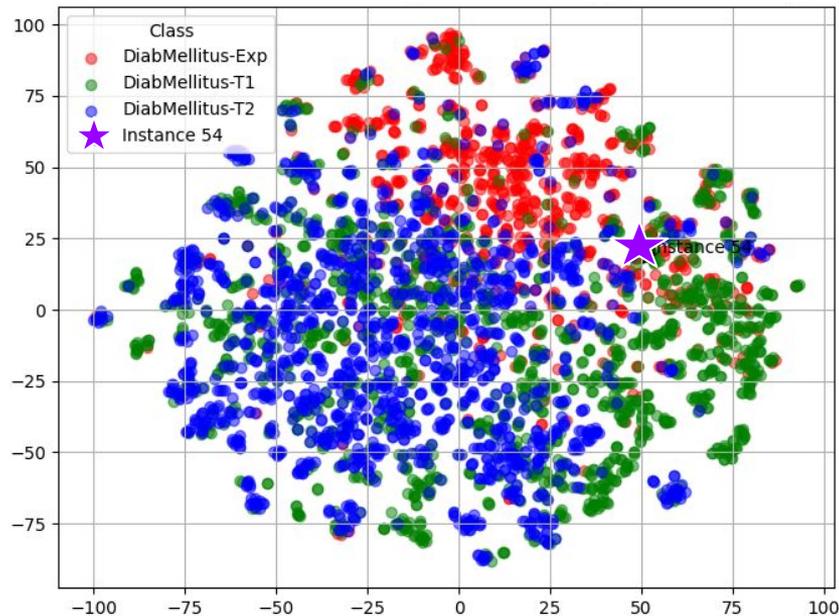


$$U(\{\text{Diab}_{\text{Exp}}\}, \hat{P}, u_{f1}) = 0.5$$

$$U(\{\text{Diab}_{\text{Exp}}, \text{Diab}_{\text{T1}}\}, \hat{P}, u_{f1}) = 0.6$$

Conformal prediction considers the set of labels one must select to statistically guarantee that they contain the right answer (with a given risk)

example:



In practice:

$$\hat{P}(Y|\star) \rightarrow \begin{matrix} \text{Diab}_{\text{Exp}}, & \text{Diab}_{\text{T1}}, & \text{Diab}_{\text{T2}} \\ 0.45 & 0.35 & 0.20 \end{matrix}$$



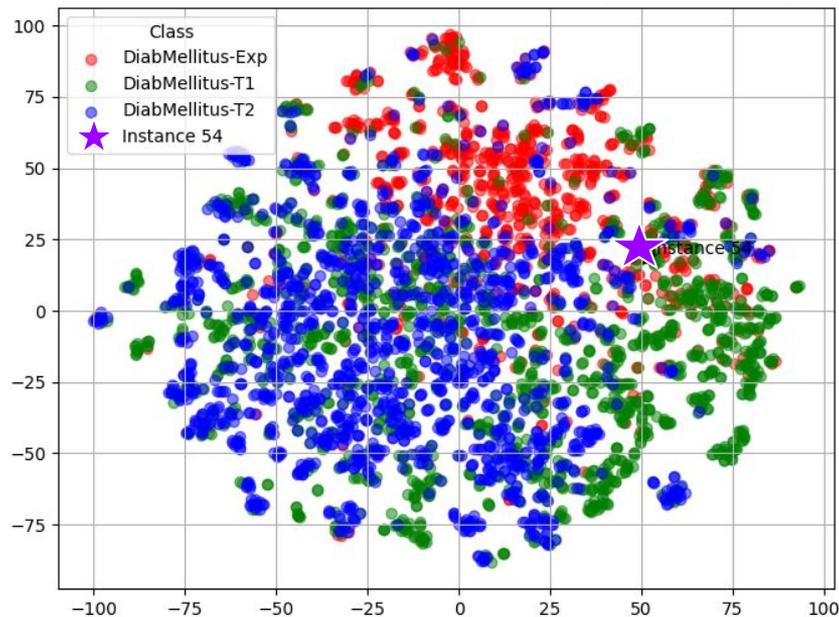
$$U(\{\text{Diab}_{\text{Exp}}\}, \hat{P}, u_{f1}) = 0.5$$

$$U(\{\text{Diab}_{\text{Exp}}, \text{Diab}_{\text{T1}}\}, \hat{P}, u_{f1}) = 0.6$$

$$U(\{\text{Diab}_{\text{Exp}}, \text{Diab}_{\text{T1}}, \text{Diab}_{\text{T2}}\}, \hat{P}, u_{f1}) = 0.465$$

Conformal prediction considers the set of labels one must select to statistically guarantee that they contain the right answer (with a given risk)

example:



In practice:

$$\hat{P}(Y|\star) \rightarrow \begin{matrix} \text{Diab}_{\text{Exp}}, & \text{Diab}_{\text{T1}}, & \text{Diab}_{\text{T2}} \\ 0.45 & 0.35 & 0.20 \end{matrix}$$



$$U(\{\text{Diab}_{\text{Exp}}\}, \hat{P}, u_{f1}) = 0.5$$

$$U(\{\text{Diab}_{\text{Exp}}, \text{Diab}_{\text{T1}}\}, \hat{P}, u_{f1}) = 0.6$$

$$U(\{\text{Diab}_{\text{Exp}}, \text{Diab}_{\text{T1}}, \text{Diab}_{\text{T2}}\}, \hat{P}, u_{f1}) = 0.465$$

In our case:

1. the number of labels we need to select is an indicator of uncertainty

**NB:** Model-based, no reference needed

## Results

---

# Human-based vs. model-based indicators

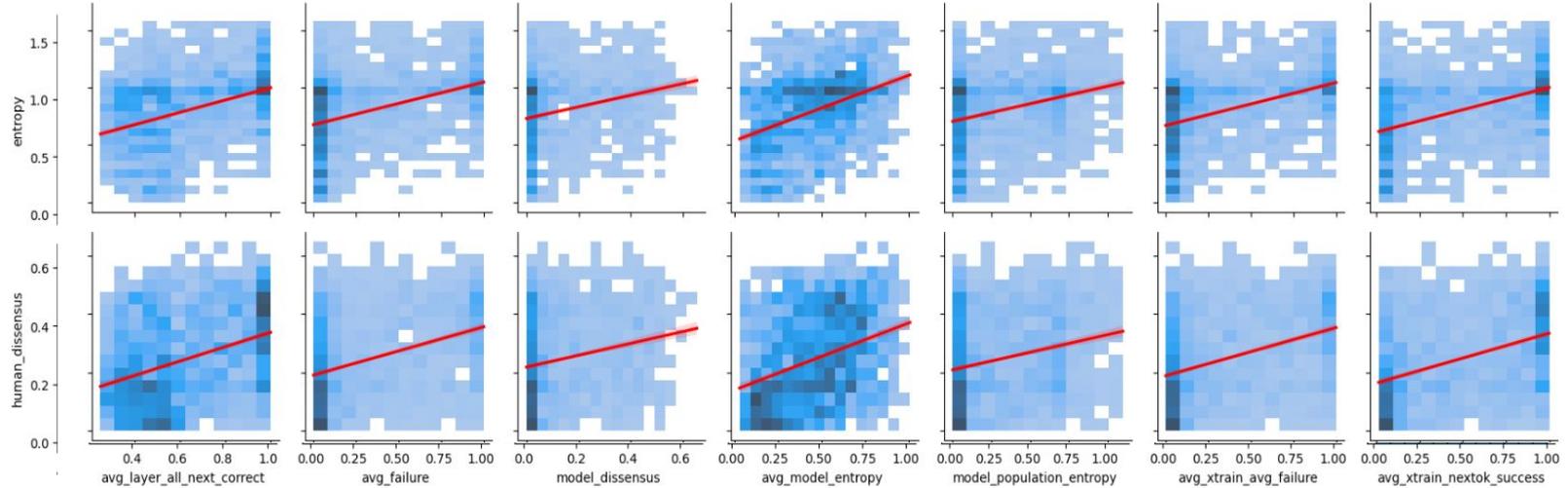
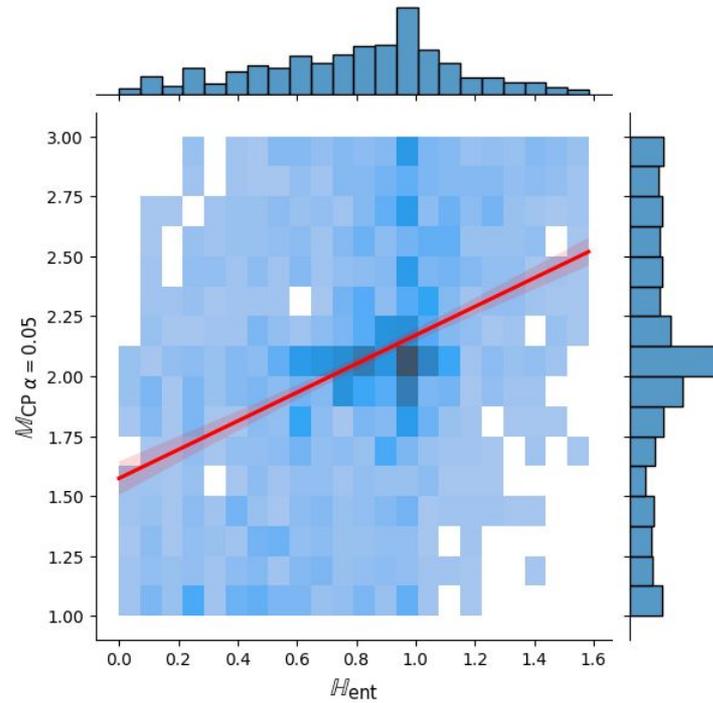
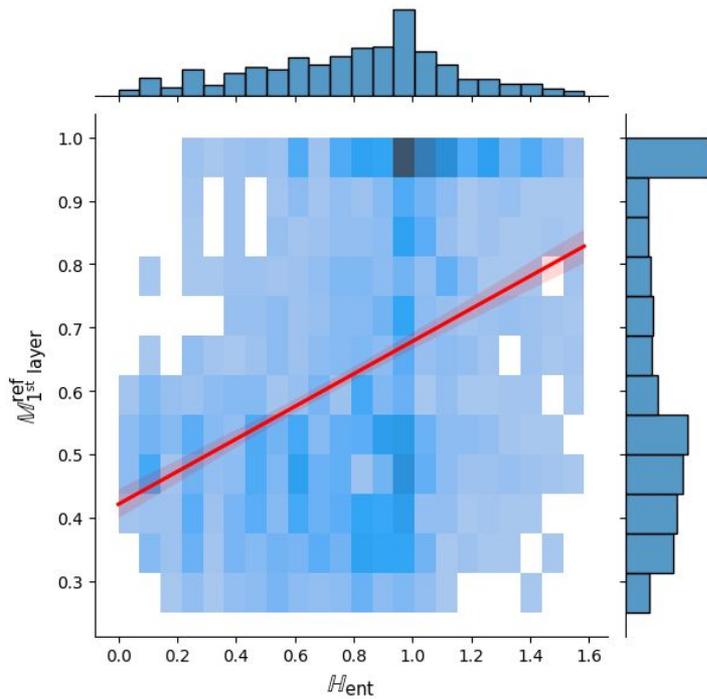
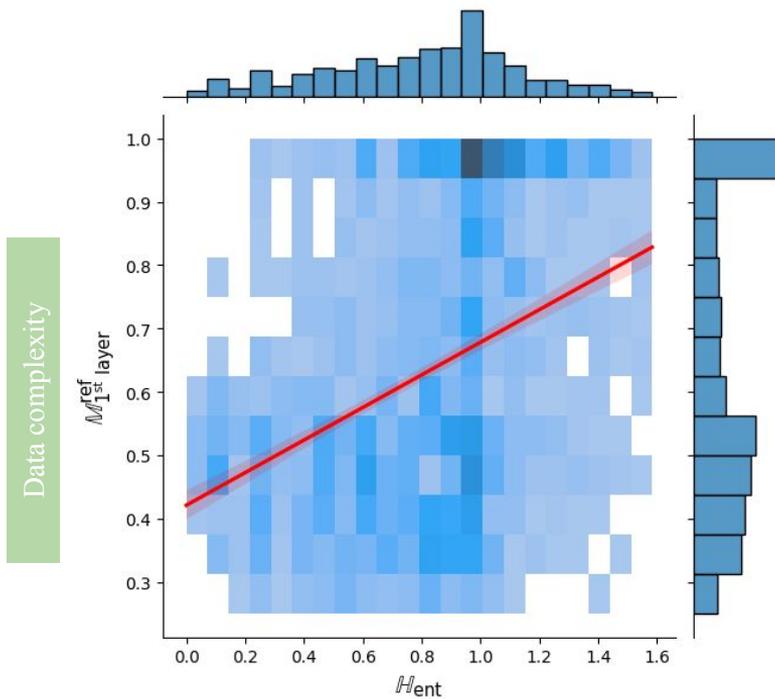


Figure : Interaction between human-based and model-based indicators

# Human-based vs. model-based indicators

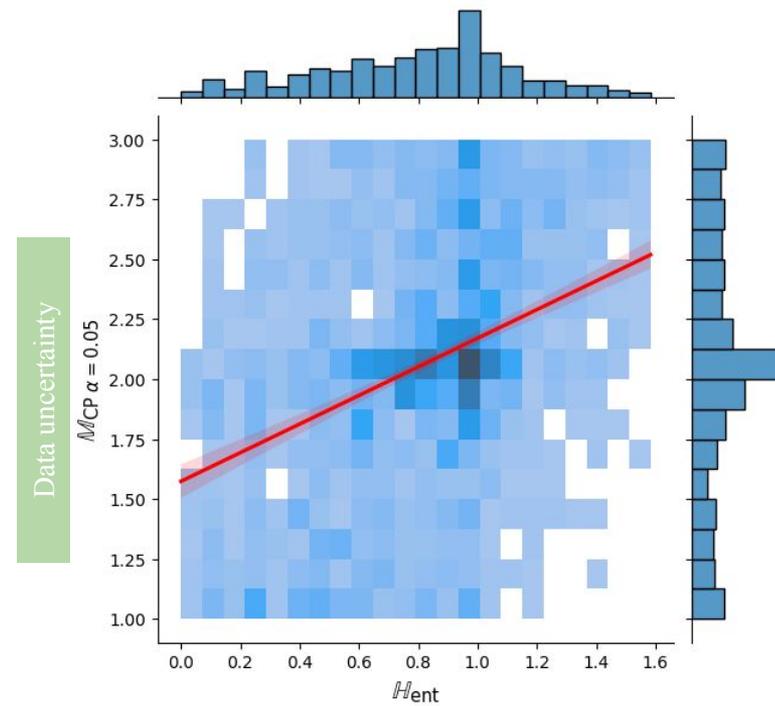


# Human-based vs. model-based indicators



Human label variation

Residual (R2 Variance) = 0.1313



Human label variation

Residual (R2 Variance) = 0.0766

# there is actually some order to this chaos

		Human label variation	
		$\mathbb{H}_{\text{ent}}$	$\mathbb{H}_{\text{dis}}$
Data uncertainty	$M_{\text{dis}}$	0.1947	0.1772
	$M_{\text{ent}}$	0.2183	0.1970
	$M_{\text{avg ent}}$	0.2811	0.2398
	$M_{\text{CP } \alpha=0.05}$	0.2767	0.2315
	$M_{\text{CP } \alpha=0.1}$	0.2819	0.2393
	$M_{\text{CP } \alpha=0.2}$	0.2482	0.2157
Data complexity	$M_{\text{fail}}^{\text{ref}}$	0.3497	0.3330
	$M_{\text{1st layer}}^{\text{ref}}$	0.3624	0.3387
	$M_{\text{1st ckpt}}^{\text{ref}}$	0.3682	0.3443
	$M_{\text{avg ckpt}}^{\text{ref}}$	0.3477	0.3274
	$M_{\text{avg ckpt } p}^{+\text{ref}}$	0.3670	0.3428

# there is actually some order to this chaos

- In practice, we observe (**low**) correlations throughout

		Human label variation	
		$\mathbb{H}_{\text{ent}}$	$\mathbb{H}_{\text{dis}}$
Data uncertainty	$M_{\text{dis}}$	0.1947	0.1772
	$M_{\text{ent}}$	0.2183	0.1970
	$M_{\text{avg ent}}$	0.2811	0.2398
	$M_{\text{CP } \alpha=0.05}$	0.2767	0.2315
	$M_{\text{CP } \alpha=0.1}$	0.2819	0.2393
	$M_{\text{CP } \alpha=0.2}$	0.2482	0.2157
Data complexity	$M_{\text{fail}}^{\text{ref}}$	0.3497	0.3330
	$M_{\text{1st layer}}^{\text{ref}}$	0.3624	0.3387
	$M_{\text{1st ckpt}}^{\text{ref}}$	0.3682	0.3443
	$M_{\text{avg ckpt}}^{\text{ref}}$	0.3477	0.3274
	$M_{\text{avg ckpt } p}^{+\text{ref}}$	0.3670	0.3428

# there is actually some order to this chaos

- ▶ In practice, we observe (**low**) correlations throughout
- ▶ whether an indicator **uses a reference** drives the correlation up

		Human label variation	
		$\mathbb{H}_{\text{ent}}$	$\mathbb{H}_{\text{dis}}$
Data uncertainty	$M_{\text{dis}}$	0.1947	0.1772
	$M_{\text{ent}}$	0.2183	0.1970
	$M_{\text{avg ent}}$	0.2811	0.2398
	$M_{\text{CP } \alpha=0.05}$	0.2767	0.2315
	$M_{\text{CP } \alpha=0.1}$	0.2819	0.2393
	$M_{\text{CP } \alpha=0.2}$	0.2482	0.2157
Data complexity	$M_{\text{fail}}^{\text{ref}}$	0.3497	0.3330
	$M_{\text{1st layer}}^{\text{ref}}$	0.3624	0.3387
	$M_{\text{1st ckpt}}^{\text{ref}}$	0.3682	0.3443
	$M_{\text{avg ckpt}}^{\text{ref}}$	0.3477	0.3274
	$M_{\text{avg ckpt } p}^{+\text{ref}}$	0.3670	0.3428

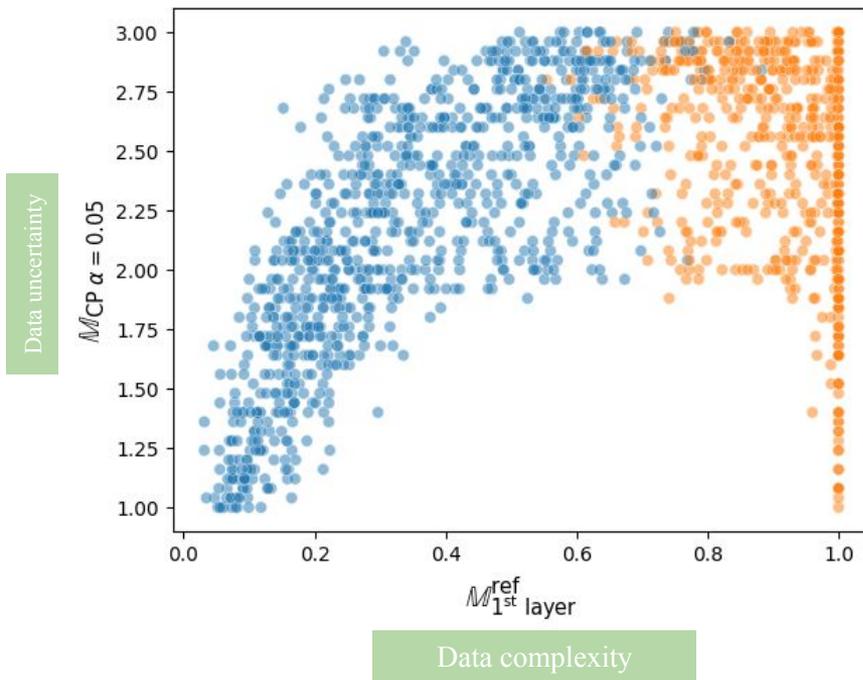
# there is actually some order to this chaos

- ▶ In practice, we observe (low) correlations throughout
- ▶ whether an indicator **uses a reference** drives the correlation up
- ▶ correlations are much higher when comparing two metrics within a group of indicators

		Human label variation	
		$\mathbb{H}_{ent}$	$\mathbb{H}_{dis}$
Data uncertainty	$M_{dis}$	0.1947	0.1772
	$M_{ent}$	0.2183	0.1970
	$M_{avg\ ent}$	0.2811	0.2398
	$M_{CP\ \alpha=0.05}$	0.2767	0.2315
	$M_{CP\ \alpha=0.1}$	0.2819	0.2393
	$M_{CP\ \alpha=0.2}$	0.2482	0.2157
Data complexity	$M_{fail}^{ref}$	0.3497	0.3330
	$M_{1^{st}\ layer}^{ref}$	0.3624	0.3387
	$M_{1^{st}\ ckpt}^{ref}$	0.3682	0.3443
	$M_{avg\ ckpt}^{ref}$	0.3477	0.3274
	$M_{avg\ ckpt\ p}^{+ref}$	0.3670	0.3428

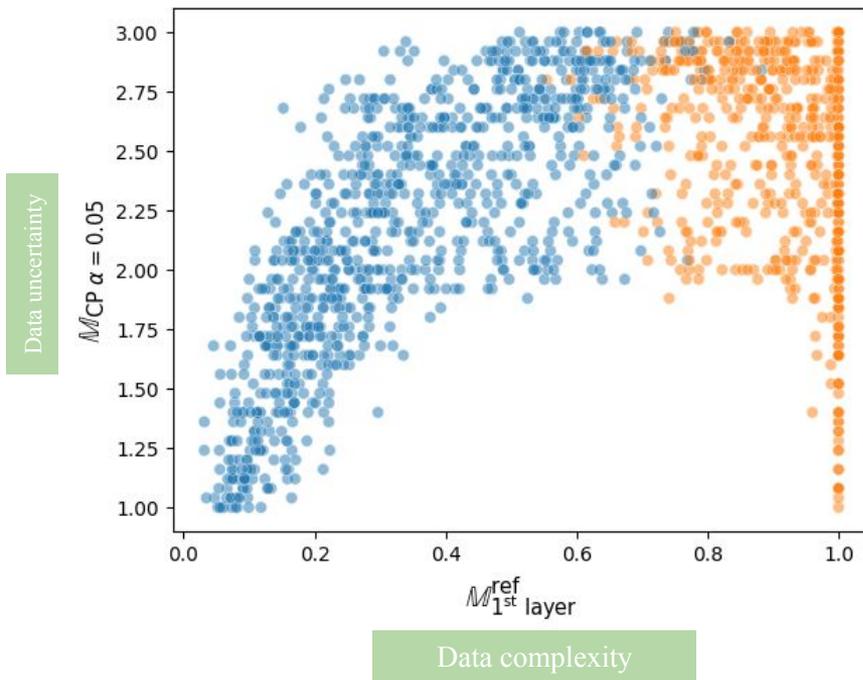
# Reference-free vs. reference-dependent indicators

► U-shaped curve



# Reference-free vs. reference-dependent indicators

- ▶ U-shaped curve
- ▶ **orange**: models tend to fail; **blue**: they tend to succeed



# More generally

	$M_{\text{fail}}^{\text{ref}}$	$M_{1^{\text{st}} \text{ layer}}^{\text{ref}}$	$M_{1^{\text{st}} \text{ ckpt}}^{\text{ref}}$	$M_{\text{avg ckpt}}^{\text{ref}}$	$M_{\text{avg ckpt } p}^{\text{ref}}$
$M_{\text{dis}}$	-0.7761	-0.7593	-0.7737	-0.7136	-0.6838
$M_{\text{ent}}$	-0.7131	-0.7075	-0.7174	-0.6539	-0.6140
$M_{\text{avg ent}}$	-0.5615	-0.5264	-0.5283	-0.5303	-0.5111
$M_{\text{CP } \alpha=0.05}$	-0.3670	-0.3633	-0.3515	-0.3389	-0.3037
$M_{\text{CP } \alpha=0.1}$	-0.4761	-0.4565	-0.4575	-0.4453	-0.4182
$M_{\text{CP } \alpha=0.2}$	-0.6427	-0.5836	-0.5967	-0.6156	-0.6116

when models tend to fail

	$M_{\text{fail}}^{\text{ref}}$	$M_{1^{\text{st}} \text{ layer}}^{\text{ref}}$	$M_{1^{\text{st}} \text{ ckpt}}^{\text{ref}}$	$M_{\text{avg ckpt}}^{\text{ref}}$	$M_{\text{avg ckpt } p}^{\text{ref}}$
$M_{\text{dis}}$	0.9536	0.8928	0.9159	0.9003	0.8934
$M_{\text{ent}}$	0.9464	0.8891	0.9107	0.8980	0.8962
$M_{\text{avg ent}}$	0.8803	0.8955	0.9116	0.8694	0.9315
$M_{\text{CP } \alpha=0.05}$	0.7748	0.7939	0.7971	0.7759	0.8546
$M_{\text{CP } \alpha=0.1}$	0.8546	0.8601	0.8816	0.8491	0.9103
$M_{\text{CP } \alpha=0.2}$	0.8996	0.8982	0.9320	0.8799	0.9166

when models don't

# More generally

	$M_{\text{fail}}^{\text{ref}}$	$M_{1^{\text{st}} \text{ layer}}^{\text{ref}}$	$M_{1^{\text{st}} \text{ ckpt}}^{\text{ref}}$	$M_{\text{avg ckpt}}^{\text{ref}}$	$M_{\text{avg ckpt } p}^{\text{ref}}$
$M_{\text{dis}}$	-0.7761	-0.7593	-0.7737	-0.7136	-0.6838
$M_{\text{ent}}$	-0.7131	-0.7075	-0.7174	-0.6539	-0.6140
$M_{\text{avg ent}}$	-0.5615	-0.5264	-0.5283	-0.5303	-0.5111
$M_{\text{CP } \alpha=0.05}$	-0.3670	-0.3633	-0.3515	-0.3389	-0.3037
$M_{\text{CP } \alpha=0.1}$	-0.4761	-0.4565	-0.4575	-0.4453	-0.4182
$M_{\text{CP } \alpha=0.2}$	-0.6427	-0.5836	-0.5967	-0.6156	-0.6116

when models tend to fail

	$M_{\text{fail}}^{\text{ref}}$	$M_{1^{\text{st}} \text{ layer}}^{\text{ref}}$	$M_{1^{\text{st}} \text{ ckpt}}^{\text{ref}}$	$M_{\text{avg ckpt}}^{\text{ref}}$	$M_{\text{avg ckpt } p}^{\text{ref}}$
$M_{\text{dis}}$	0.9536	0.8928	0.9159	0.9003	0.8934
$M_{\text{ent}}$	0.9464	0.8891	0.9107	0.8980	0.8962
$M_{\text{avg ent}}$	0.8803	0.8955	0.9116	0.8694	0.9315
$M_{\text{CP } \alpha=0.05}$	0.7748	0.7939	0.7971	0.7759	0.8546
$M_{\text{CP } \alpha=0.1}$	0.8546	0.8601	0.8816	0.8491	0.9103
$M_{\text{CP } \alpha=0.2}$	0.8996	0.8982	0.9320	0.8799	0.9166

when models don't

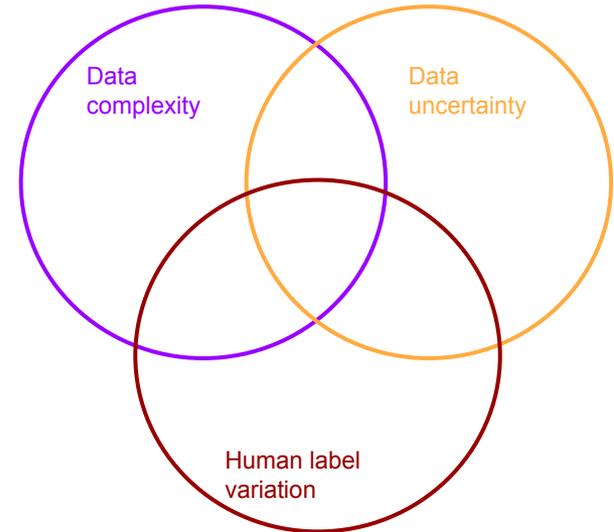
Reference-free model-based indicators systematically conflate model successes and model failures

## Takeaways

---

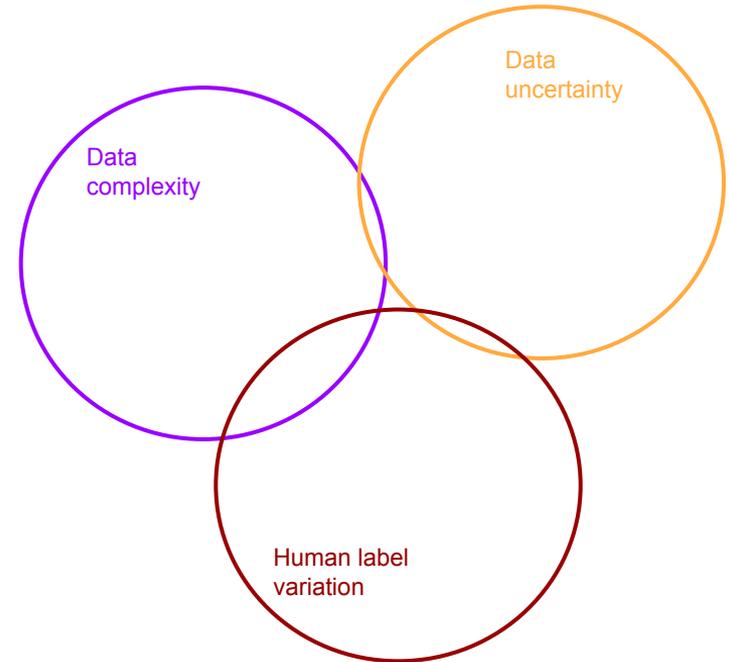
# So what?

- ▶ Model-based indicators **align poorly** with human-based indicators!



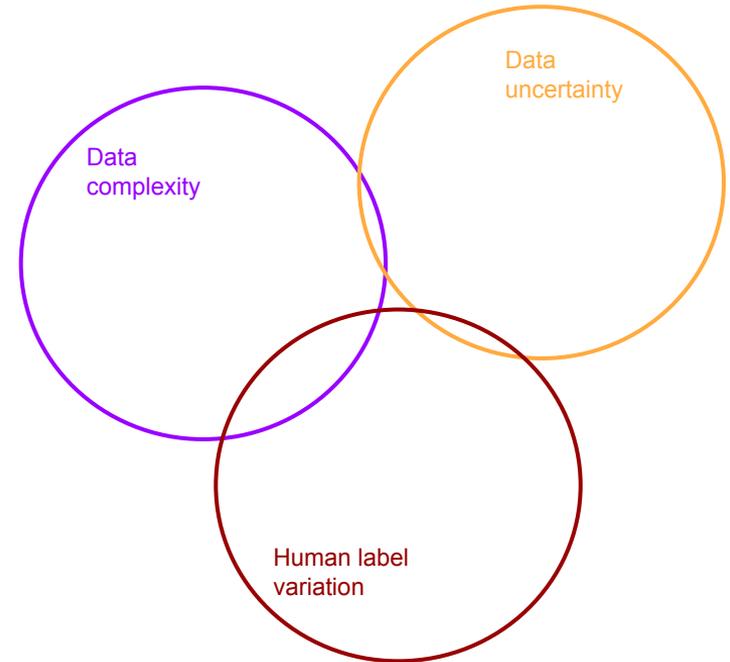
# So what?

- ▶ Model-based indicators **align poorly** with human-based indicators!
- ▶ If *what was easy/difficult for humans was also easy/difficult for models*, why does **data complexity** more aligned with **human disagreement** than **data uncertainty**.



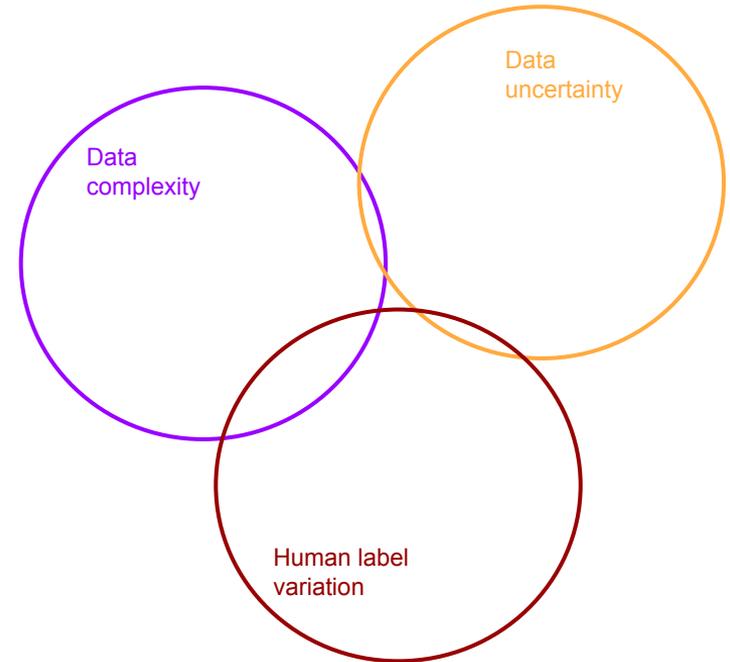
# So what?

- ▶ Model-based indicators **align poorly** with human-based indicators!
- ▶ If *what was easy/difficult for humans was also easy/difficult for models*, why does **data complexity** more aligned with **human disagreement** than **data uncertainty**.
- ▶ We need to **rethink practices such as active learning** (where difficulty of samples is taken in account during training)



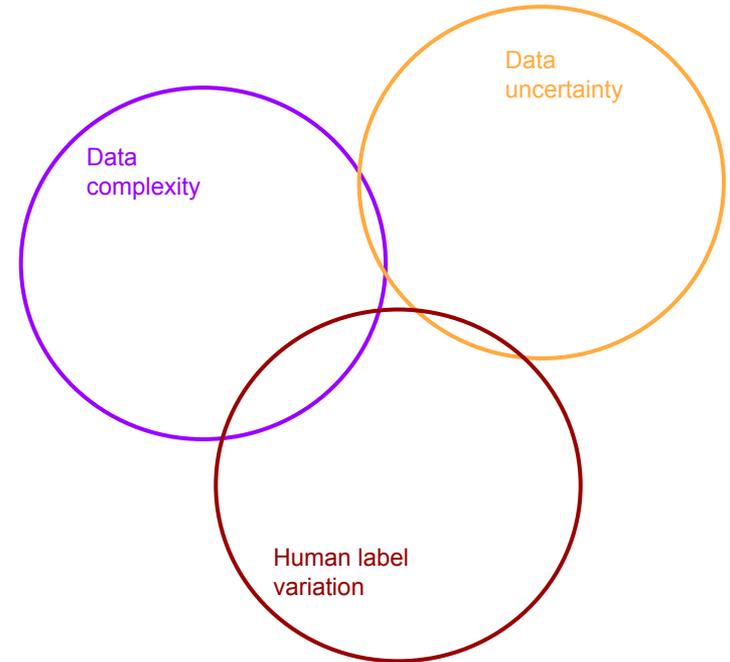
# So what?

- ▶ Model-based indicators **align poorly** with human-based indicators!
- ▶ If *what was easy/difficult for humans was also easy/difficult for models*, why does **data complexity** more aligned with **human disagreement** than **data uncertainty**.
- ▶ We need to **rethink practices such as active learning** (where difficulty of samples is taken in account during training)
- ▶ not all framework are created equal (out of the box **CP has flaws**)



# So what?

- ▶ Model-based indicators **align poorly** with human-based indicators!
- ▶ If *what was easy/difficult for humans was also easy/difficult for models*, why does **data complexity** more aligned with **human disagreement** than **data uncertainty**.
- ▶ We need to **rethink practices such as active learning** (where difficulty of samples is taken in account during training)
- ▶ not all framework are created equal (out of the box **CP has flaws**)
- ▶ Linguistic ambiguity remains distinct for model-based indicators



# Appendix - I

- **Low residual** : Non-linear relationship between human and model based indicators.

	$\mathbb{H}_{\text{dis}}$	$\mathbb{H}_{\text{ent}}$
$M_{\text{dis}}$	0.0314	0.0379
$M_{\text{avg ent}}$	0.0575	0.0790
$M_{\text{ent}}$	0.0388	0.0477
$M_{\text{CP } \alpha=0.05}$	0.0536	0.0766
$M_{\text{CP } \alpha=0.1}$	0.0573	0.0795
$M_{\text{CP } \alpha=0.2}$	0.0465	0.0616
$M_{\text{fail}}^{\text{ref}}$	0.1109	0.1223
$M_{\text{1st layer}}^{\text{ref}}$	0.1147	0.1313
$M_{\text{1st ckpt}}^{\text{ref}}$	0.1186	0.1356
$M_{\text{avg ckpt}}^{\text{ref}}$	0.1072	0.1209
$M_{\text{avg ckpt } p}^{+\text{ref}}$	0.1175	0.1347

Table : Proportion of explained variance (R2) of linear regressions predicting a model-based indicator from a human-based indicator.

	$M_{\text{fail}}^{\text{ref}}$	$M_{\text{1st layer}}^{\text{ref}}$	$M_{\text{1st ckpt}}^{\text{ref}}$	$M_{\text{avg ckpt}}^{\text{ref}}$	$M_{\text{avg ckpt } p}^{\text{ref}}$
$M_{\text{dis}}$	0.5154	0.4966	0.5029	0.5035	0.5048
$M_{\text{ent}}$	0.5168	0.4984	0.5047	0.5073	0.5127
$M_{\text{avg ent}}$	0.5292	0.5419	0.5468	0.5292	0.5560
$M_{\text{CP } \alpha=0.05}$	0.4860	0.4958	0.4972	0.4916	0.5286
$M_{\text{CP } \alpha=0.1}$	0.5216	0.5290	0.5353	0.5241	0.5527
$M_{\text{CP } \alpha=0.2}$	0.5232	0.5338	0.5437	0.5188	0.5338

Table : Spearman correlation between reference dependent and reference-free indicators

**Thank you for your attention!**

**Any questions?**

Interested in my work ?



Want to read our full paper ?



*I am looking for postdoc opportunities in NLP, healthcare, LLM Evaluations*