

# From Human Label Variation and Model Uncertainty to Error Detection (and Back)?

Prof. Dr. Barbara Plank  
MaiNLP lab, CIS, LMU Munich  
& IT University of Copenhagen

NLPerspectives workshop LREC-COLING 2024  
May 21, 2024

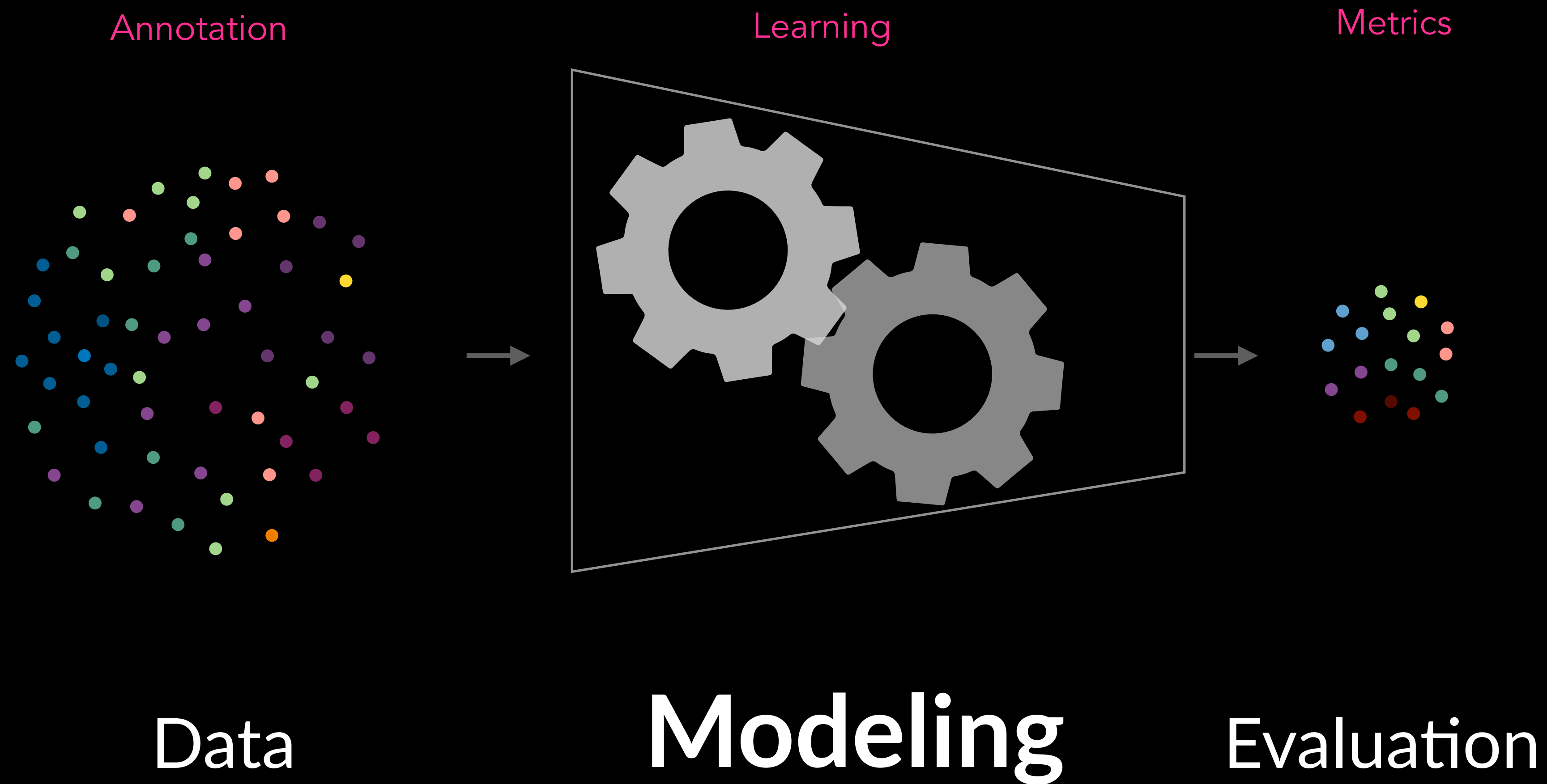


IT UNIVERSITY OF COPENHAGEN



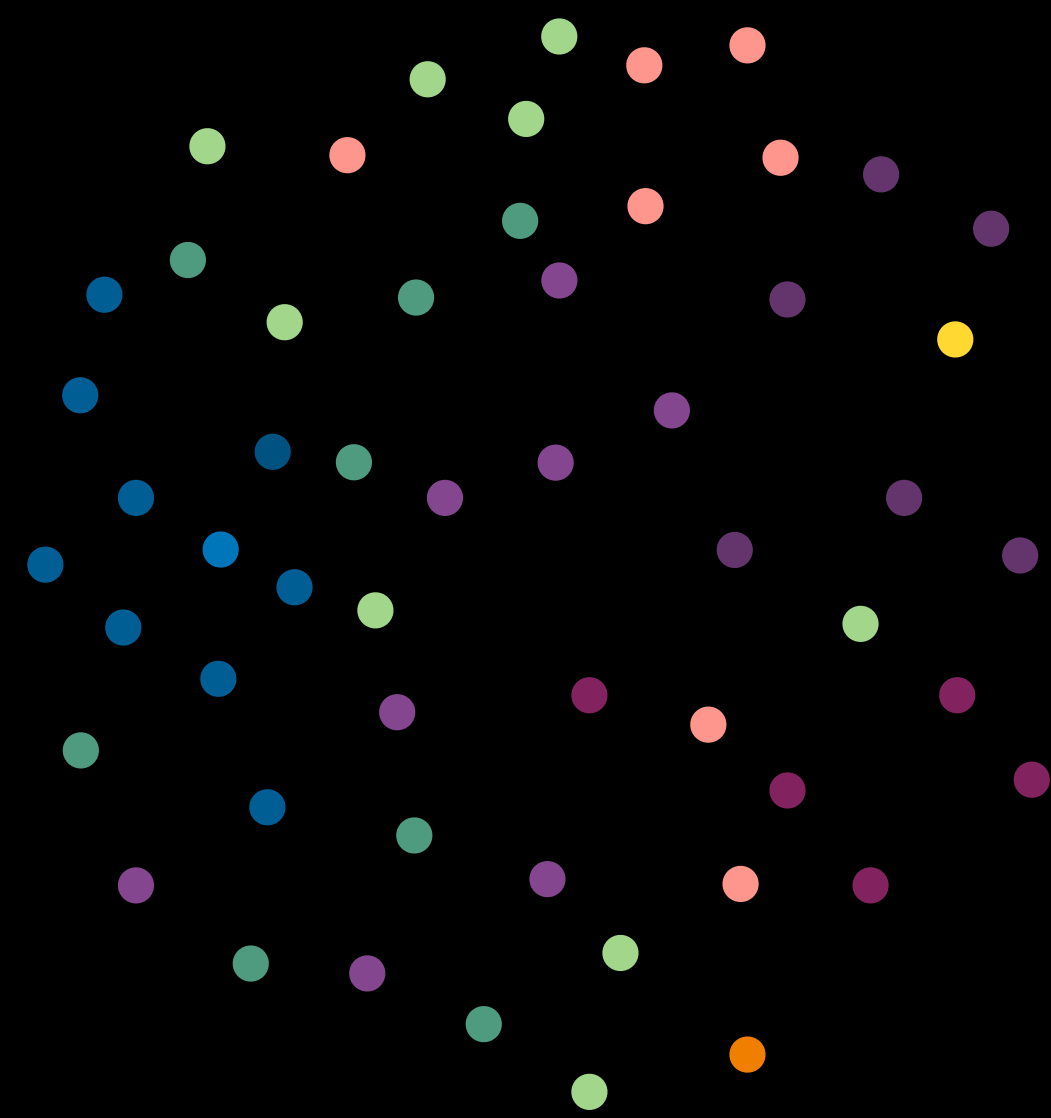
# A typical/traditional NLP/AI pipeline

---

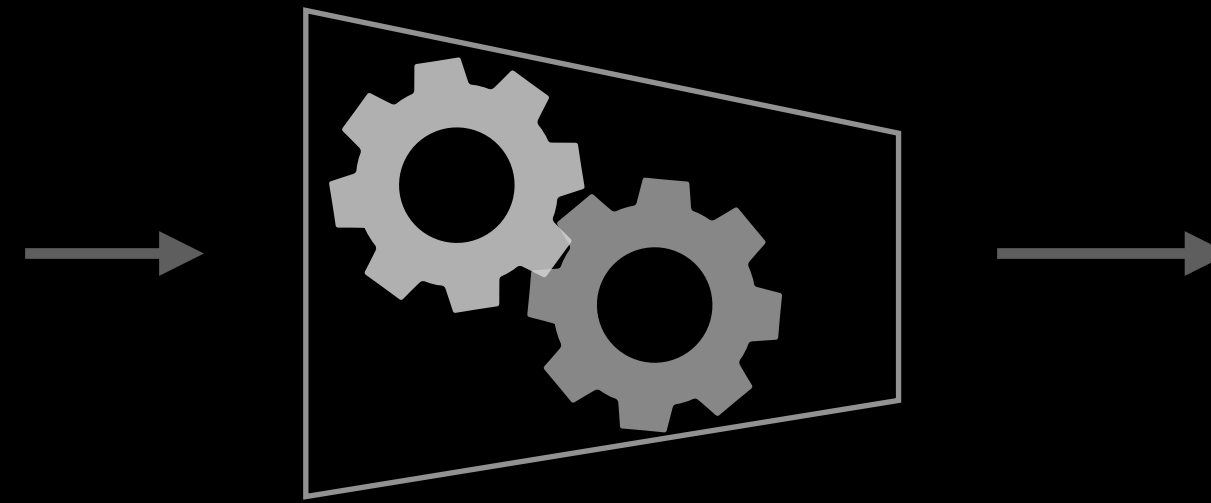


# Growing Importance of High-Quality Data and Evaluation

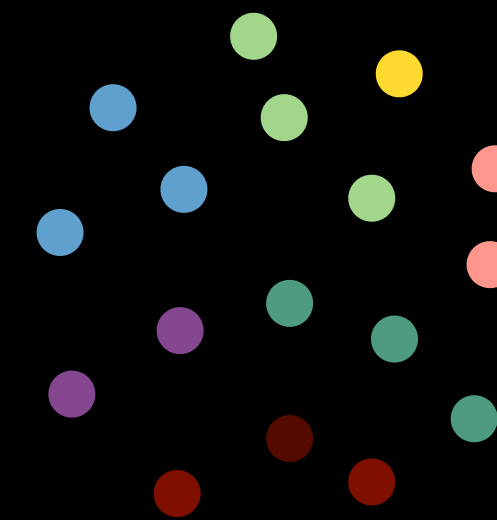
---



**Data**



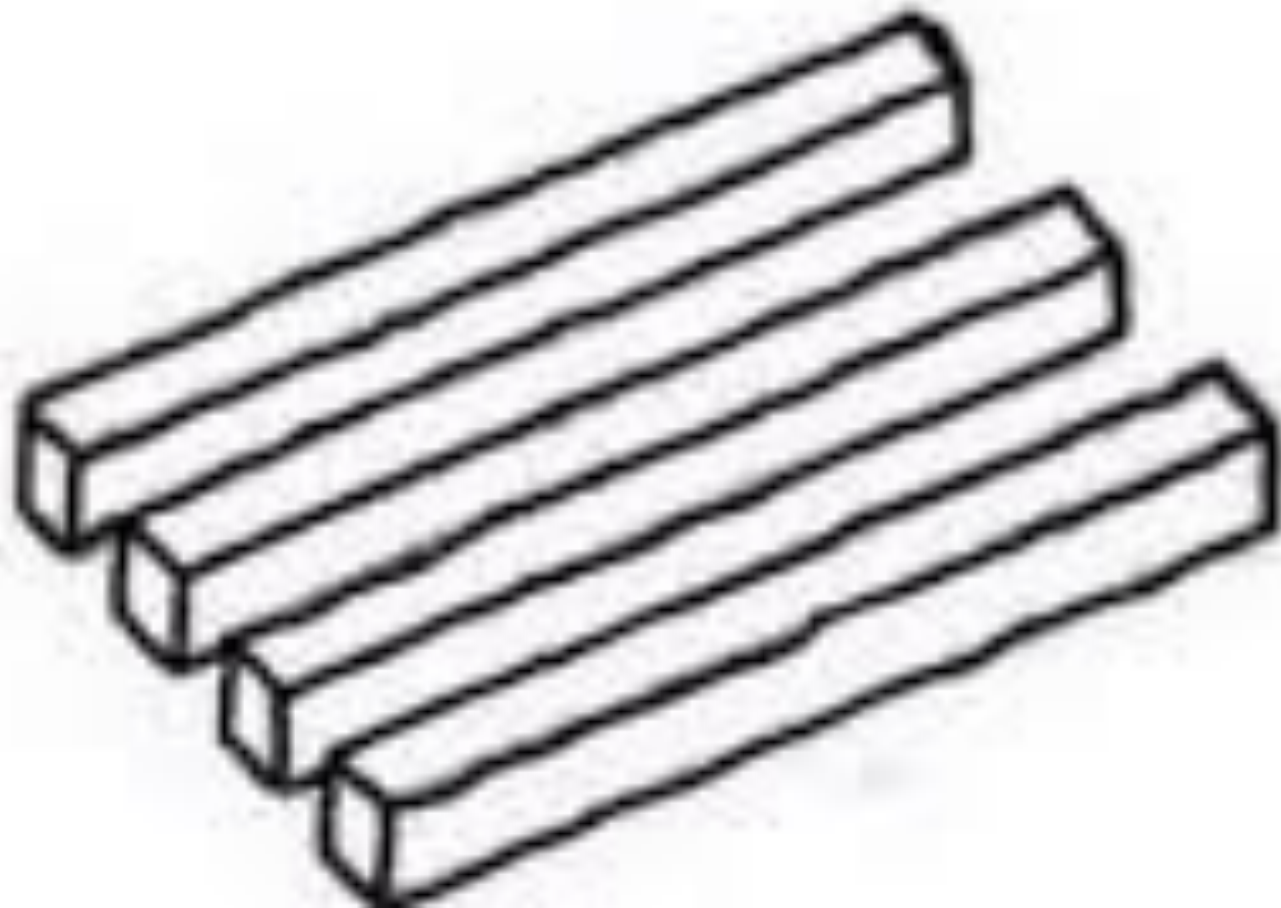
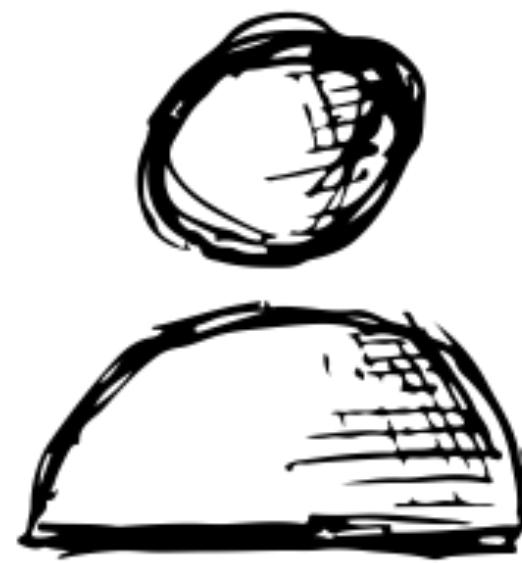
Modeling



**Evaluation**

**Four**

**No.  
Three**







**The world is beyond digital - if we  
zoom out it is rich and diverse.**

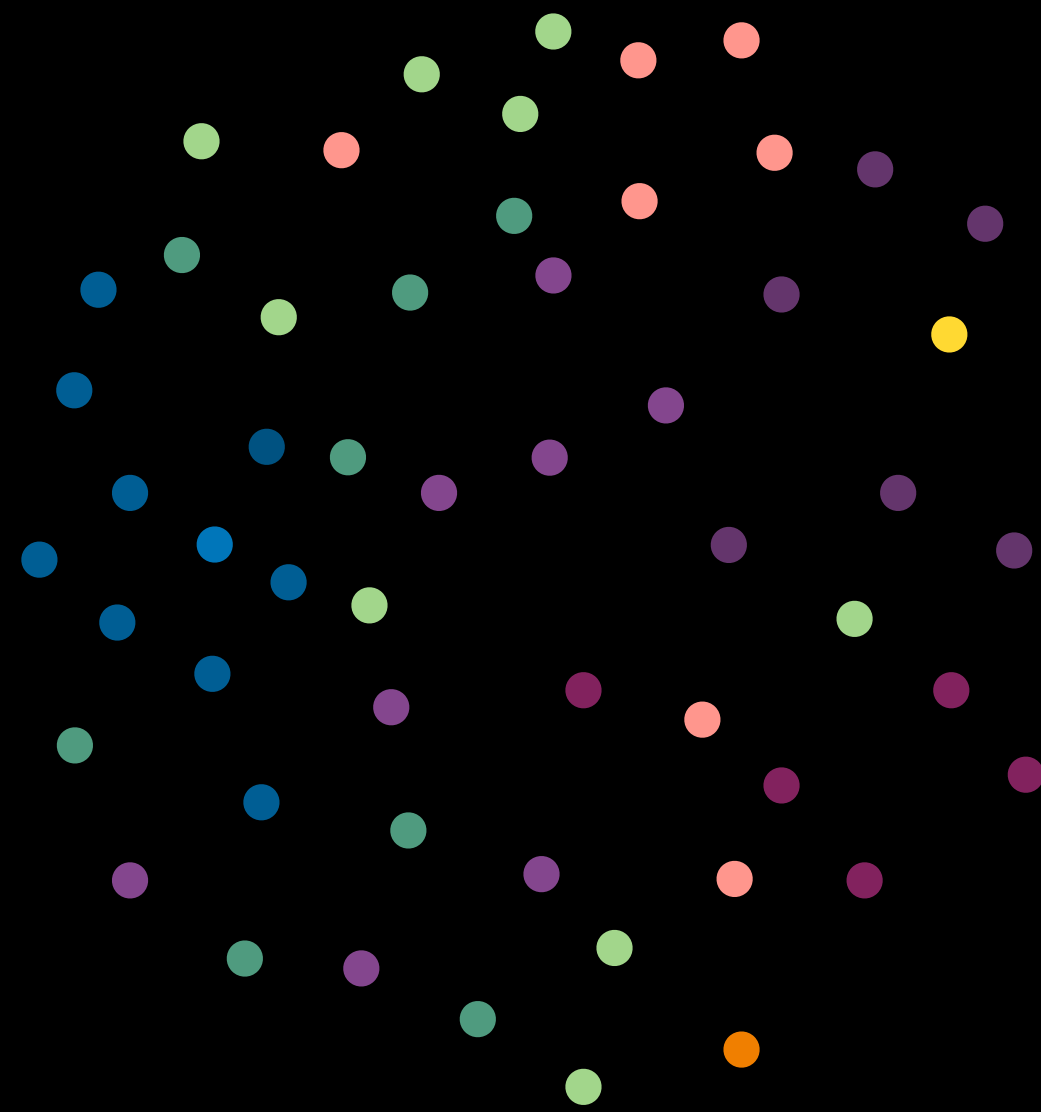


# Disagreement in human labeling is ubiquitous

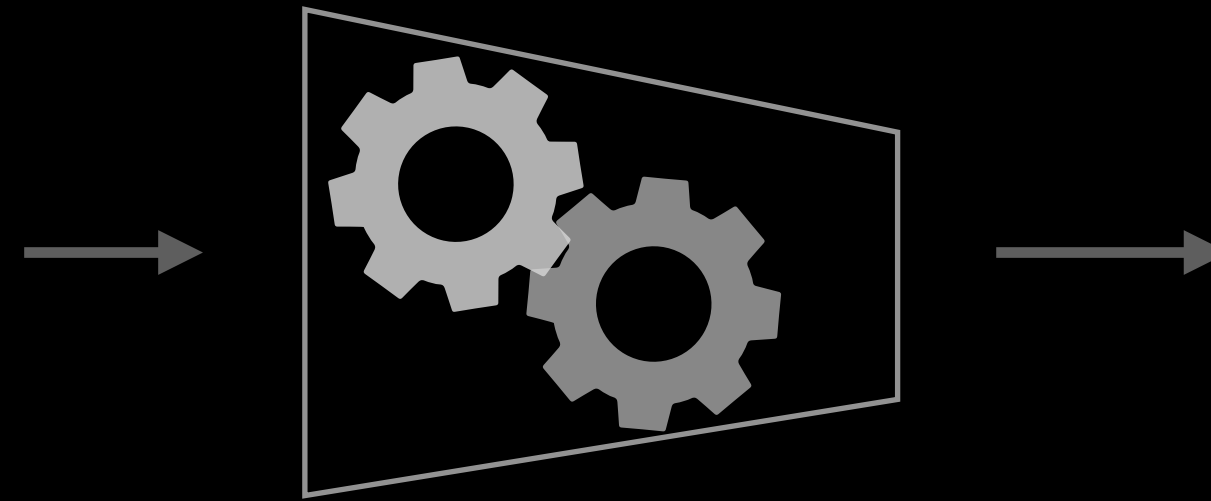
---

- It impacts **all 3 stages** of the NLP pipeline
- It is one important form of **uncertainty**

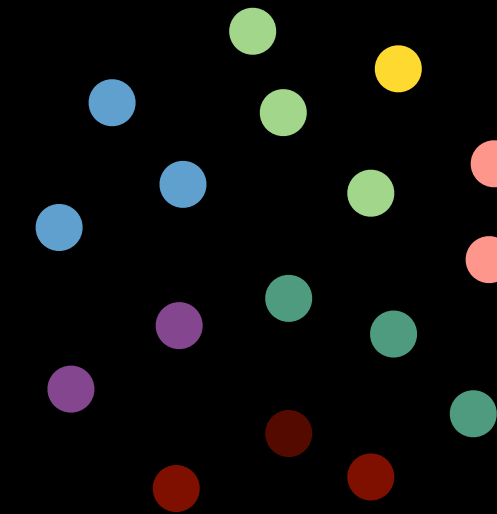
*Can we turn  
disagreement into  
advantage?*



Data



Modeling



Evaluation

**Growing Importance of  
Data Quality > Data Quantity**

# The “it” in AI models is the dataset - talk by Thom Wolf 🤗

---

## The “it” in AI models is the dataset.

Posted on June 10, 2023 by jbetker

I’ve been at OpenAI for almost a year now. In that time, I’ve trained a lot of generative models. More than anyone really has any right to train. As I’ve spent these hours observing the effects of tweaking various model configurations and hyperparameters, one thing that has struck me is the similarities in between all the training runs.

It’s becoming awfully clear to me that these models are truly approximating their datasets to an incredible degree. What that means is not only that they learn what it means to be a dog or a cat, but the interstitial frequencies between distributions that don’t matter, like what photos humans are likely to take or words humans commonly write down.

What this manifests as is – trained on the same dataset for long enough, pretty much every model with enough weights and training time converges to the same point. Sufficiently large diffusion conv-unets produce the same images as ViT generators. AR sampling produces the same images as diffusion.

This is a surprising observation! It implies that model behavior is not determined by architecture, hyperparameters, or optimizer choices. It’s determined by your dataset, nothing else. Everything else is a means to an end in efficiently delivery compute to approximating that dataset.

Then, when you refer to “Lambda”, “ChatGPT”, “Bard”, or “Claude” then, it’s not the model weights that you are referring to. It’s the dataset.



# Evidence from a talk by Sara Hooker 🦋

Model	Size (# Parameters)	Training Tokens
LaMDA ( <a href="#">Thoppilan et al., 2022</a> )	137 Billion	168 Billion
GPT-3 ( <a href="#">Brown et al., 2020</a> )	175 Billion	300 Billion
Jurassic ( <a href="#">Lieber et al., 2021</a> )	178 Billion	300 Billion
Gopher ( <a href="#">Rae et al., 2021</a> )	280 Billion	300 Billion
MT-NLG 530B ( <a href="#">Smith et al., 2022</a> )	530 Billion	270 Billion
<i>Chinchilla</i>	70 Billion	1.4 Trillion

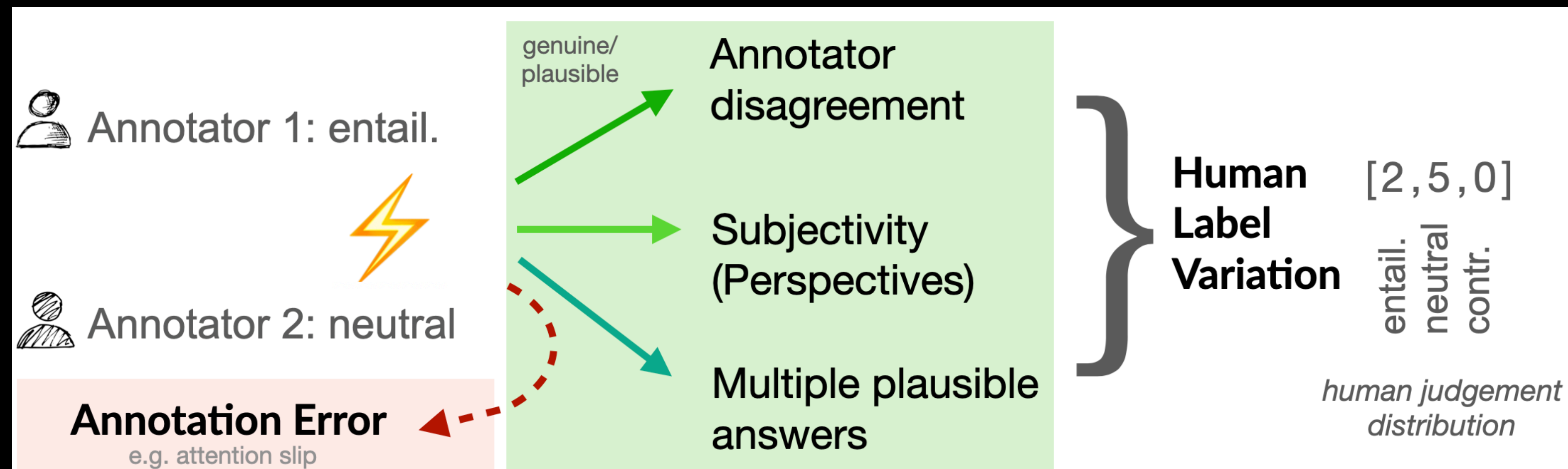
- Recent work suggests smaller amounts of higher quality data remove the need for a larger model.
- This suggest larger models may just be compensating for problems in the data pipeline.

# Roadmap: Selected Case Studies

- 1 Humans and Uncertainty: What is *Human Label Variation*?
- 2 Models and Uncertainty: *Stop Measuring Calibration When Humans Disagree*
- 3 How to detect errors? *ActiveAED*
- 4 Plausible variation or error? *VariERR*

# Disagreement or Variation?

- ▶ **Human Label Variation** (Plank, 2022 EMNLP)
  - ▶ plausible variation
  - ▶ to **reconcile different notions** in the literature (disagreement, perspectives, human uncertainty, hard cases)
  - ▶ preferred over **disagreement** as that implies two views cannot hold at the same time

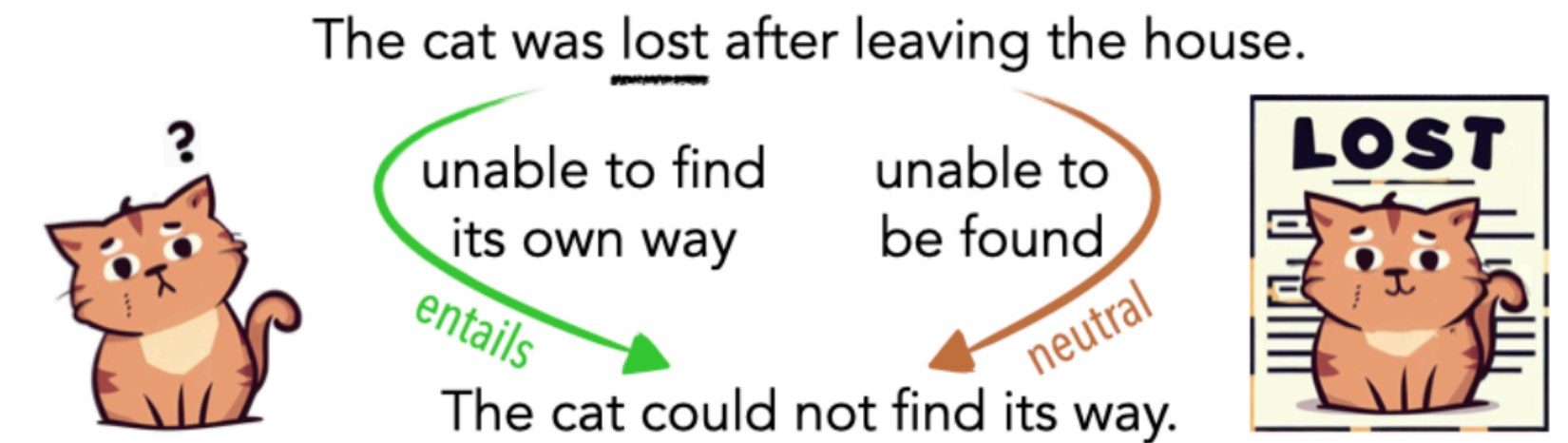


- ▶ In contrast to **errors**

# Sources of human label variation

(Basile et al., 2021)

- ▶ **Stimulus characteristics** (ambiguity, task setup and difficulty)
- ▶ **Individual differences** (incl. cultural and socio-demographics): for example in hate speech or sentiment
- ▶ **Context and attention** (Intra-coder disagreement; attention slips play a non-negligible role as well; Beigman Klebanov et al., 2008)



*Ambiguity (Example from Liu et al., 2023)*



**Examples**

# Lora Aroyo's NeurIPS 2023 keynote:



Is there a **SMILE** in this image?

YES but ...

Canada		
YES	NO	DNK
40%	40%	20%

India		
YES	NO	DNK
70%	30%	0

USA		
YES	NO	DNK
50%	0	50%





# Name the object



# Name the object



cake (53), food (19), bread (8), burger (6),  
dessert (6), snacks (3), muffin (3), pastry (3)



# Natural Language Inference: Entailment? Neutral? Contradiction?

*Context/Premise:*

*Statement/Hypothesis:*

*[E, N, C]*

A boy in an orange shirt sells fruit from a street cart.

A boy is a street vendor.

[90, 10, 0]

A women wearing a red hat and black coat.

The women is asleep.

[0, 87, 13]

People walk amonst a traffic jam in a crowded city.

The cars are zooming past the people.

[3, 15, 82]

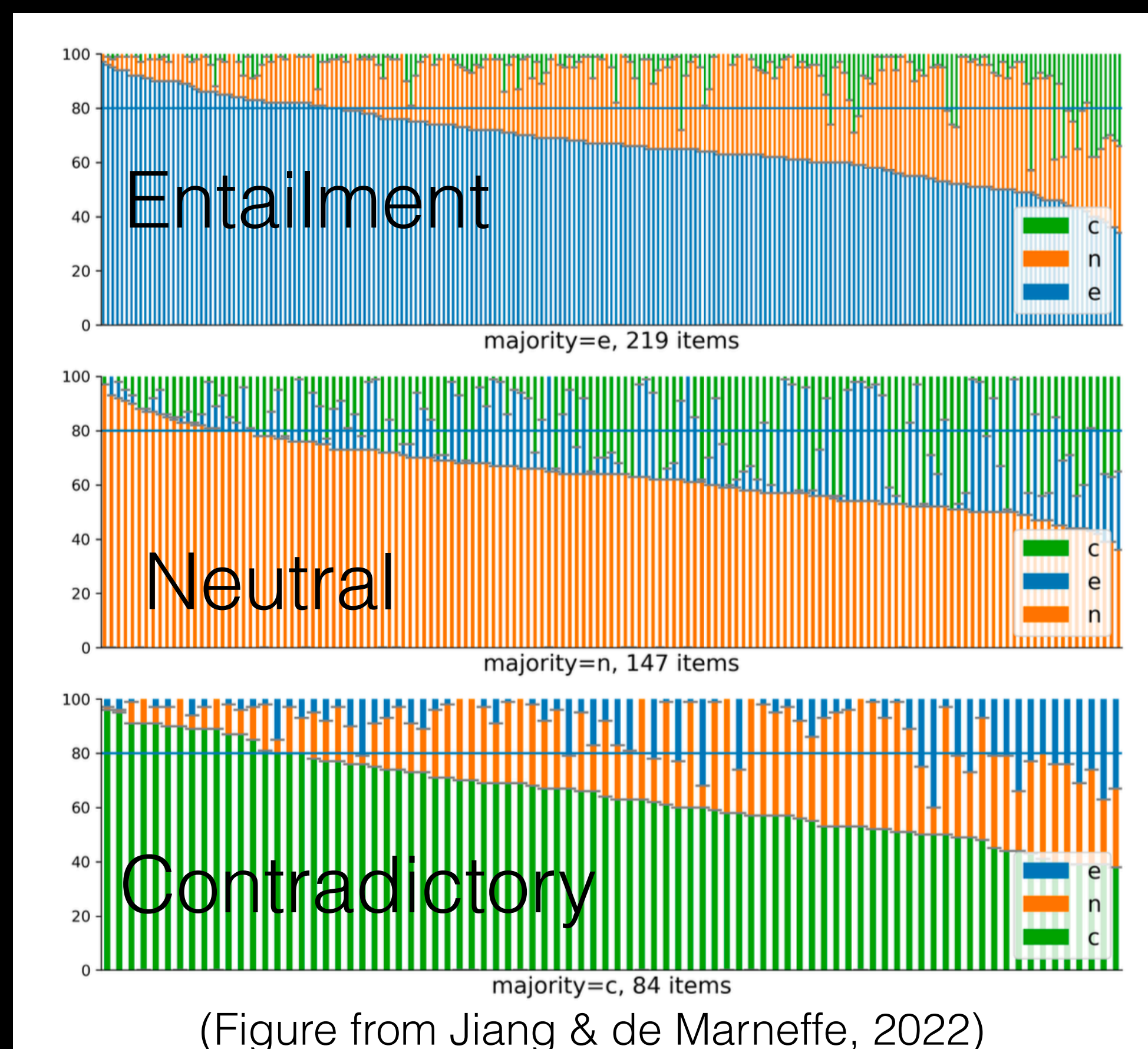
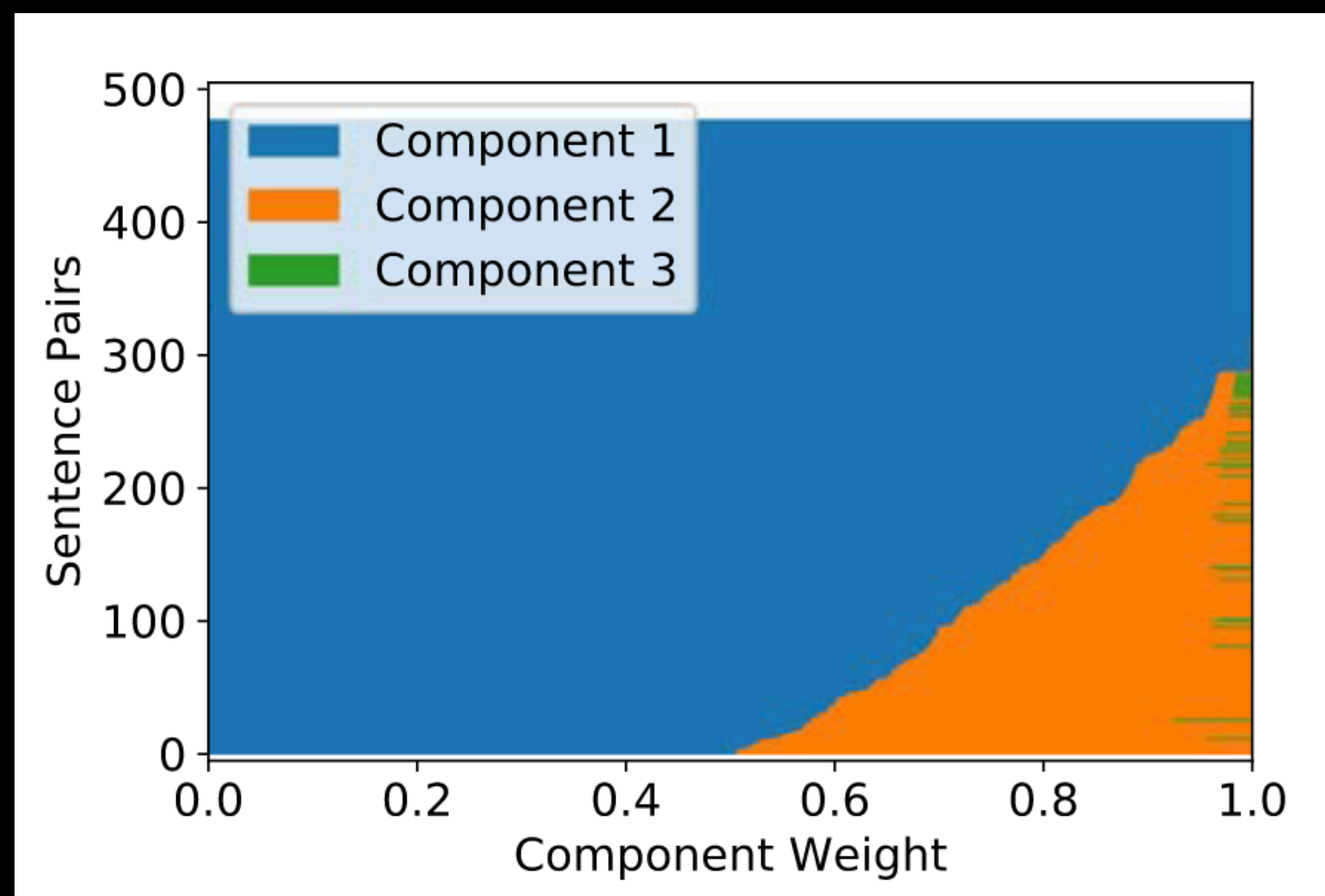
A women holding a child in a purple shirt.

The women is asleep at home.

[1, 53, 46]

# Natural Language Inference: How Frequent?

- "For 20% of the sentence pairs, there is a non-trivial second component (GMM; Pavlick & Kwiatkowski, 2019)



# More NLP task examples (to name a few):

- ▶ **Toxic language detection:** Not all text is *equally toxic* for everyone (Sap et al., 2019).  
*Subjective language tasks* (Akhtar et al, 2021; Leonardelli et al., 2021; Ceras Curry et al., 2021)

- ▶ **Understanding indirect answers** to polar questions (e.g. Damgaard et al., 2021)

Q: Hey. Everything ok?  
A: I'm just mad at my agent

? Yes

? No

? Yes, subject to some condition

? Neither Yes nor no

- ▶ **Visual Question Answering** (Jolly et al., 2021)

Q: Where is this?  
GT: road (4), outside (2), outdoors (1),  
sidewalk (1), ...





# HLV not just labels: Natural Language Generation

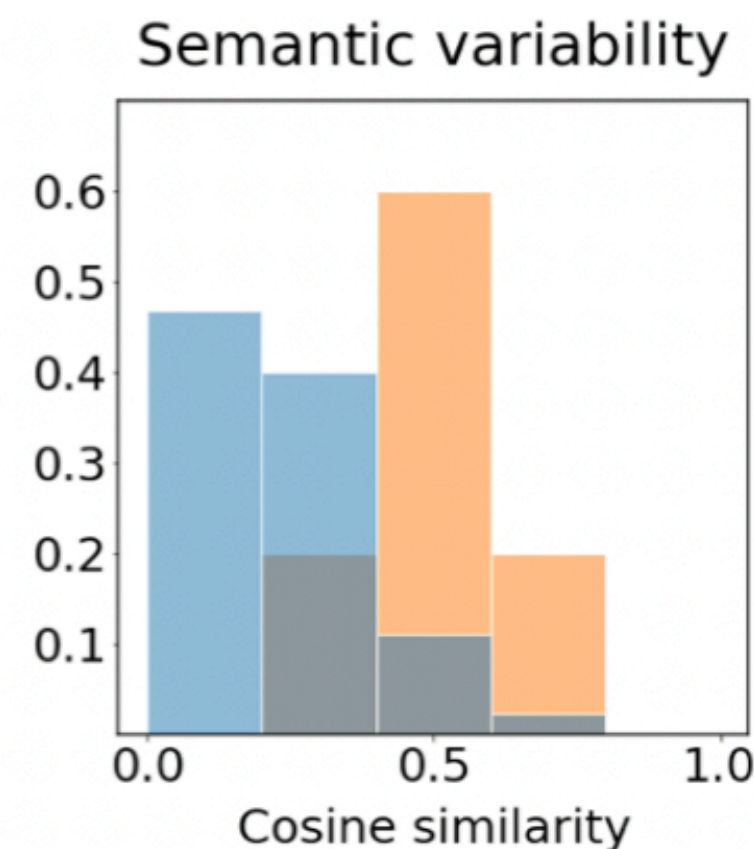
## Dialogue context

It's very dark in here. Will you turn on the light?

Okay. But our baby has fallen asleep.

Then, turn on the lamp, please.

But where's the switch?



## Humans



- Don't you know where the switch is?
- Switch is on the left side of the lamp.
- Just press the second switch on the board.
- Lamp is upon the study table and now you know where the switch is.
- I will light up the torch, so you can find the switch and turn on the lamp.

## DialoGPT-medium, nucleus $p = 0.9$



- You don't have one.
- Where's the button?
- It's on the top.
- Well, you'll want to turn it on.
- Turn it on.

(Guilianelli, Baan et al., 2023)



Is Human Label Variation So Bad? **No.**

It provides opportunities for more trustworthy, human-facing AI.



# Many open questions

---

- ▶ How does Human Label Variation interact with Socially Aware NLP?
- ▶ How much Variation is there Within-Human Label Variation?
- ▶ How does Human Label Variation relate to Model Uncertainty?
- ▶ Is All Variation Plausible?
- ▶ How to tease apart Annotation Error from Plausible Human Variation?
- ▶ Do we Need More Labels or More Cases?

# Human Label Variation: Growing interest 🎉

## Toward a Perspectivist Turn in Ground Truthing for Predictive Computing

Valerio Basile<sup>1</sup>, Federico Cabitza<sup>2</sup>, Andrea Campagner<sup>2</sup>, and Michael Fell<sup>1</sup>

## Why Don't You Do It Right? Analysing Annotators' Disagreement in Subjective Tasks

Marta Sandri  
Dept. of Humanities  
University of Pavia, Italy  
sandri.marta97@gmail.com

Elisa Leonardelli  
Fondazione Bruno Kessler  
Trento, Italy  
eleonardelli@fbk.eu

Sara Tonelli

Elisabetta Jezek

## Understanding and Predicting Human Label Variation in Natural Language Inference through Explanations

Nan-Jiang Jiang<sup>1</sup>   Chenhao Tan<sup>2</sup>   Marie-Catherine de Marneffe<sup>1,3</sup>

## DisaggregHate It Corpus: A Disaggregated Italian Dataset of Hate Speech

Marco Madeddu<sup>1</sup>, Simona Frenda<sup>1,2</sup>, Mirko Lai<sup>1,2</sup>, Viviana Patti<sup>1</sup> and Valerio Basile<sup>1</sup>

## More Labels or Cases? Assessing Label Variation in Natural Language Inference

Cornelia Gruber<sup>\*1♠</sup>   Katharina Hechinger<sup>\*1♠</sup>   Matthias Aßenmacher<sup>1,2♠</sup>  
Göran Kauermann<sup>1♠</sup>   Barbara Plank<sup>2,3♠</sup>

## Interpreting Predictive Probabilities: Model Confidence or Human Label Variation?

Joris Baan<sup>🔗</sup>, Raquel Fernández<sup>🔗</sup>, Barbara Plank<sup>🔗🔗</sup>, Wilker Aziz<sup>🔗</sup>

## EPIC: Multi-Perspective Annotation of a Corpus of Irony

Simona Frenda<sup>\*🔗</sup>, Alessandro Pedrani<sup>🔗</sup>, Valerio Basile<sup>\*</sup>, Soda Maren Lo<sup>\*</sup>,  
Alessandra Teresa Cignarella<sup>\*🔗</sup>, Raffaella Panizzon<sup>🔗</sup>, Cristina Marco<sup>🔗</sup>,  
Bianca Scarlini<sup>🔗</sup>, Viviana Patti<sup>\*</sup>, Cristina Bosco<sup>\*</sup>, Davide Bernardi<sup>🔗</sup>

## ACTOR: Active Learning with Annotator-specific Classification Heads to Embrace Human Label Variation

Xinpeng Wang and Barbara Plank

## Through the Lens of Split Vote: Exploring Disagreement, Difficulty and Calibration in Legal Case Outcome Classification

Shanshan Xu<sup>1</sup>, Santosh T.Y.S.S<sup>1</sup>, Oana Ichim<sup>2</sup>  
Barbara Plank<sup>3,4</sup>, Matthias Grabmair<sup>1</sup>

## Annotator-Centric Active Learning for Subjective NLP Tasks

Michiel van der Meer  
LIACS  
Leiden University

Neele Falk  
Institute for Natural Language Processing  
University of Stuttgart

## Consistency is Key: Disentangling Label Variation in Natural Language Processing with Intra-Annotator Agreement

Gavin Abercrombie<sup>1</sup> and Verena Rieser<sup>1,2</sup> and Dirk Hovy<sup>3</sup>

## When the Majority is Wrong: Modeling Annotator Disagreement in Subjective Tasks

Eve Fleisig<sup>†</sup>

Rediet Abebe

Dan Klein

## Which Examples Should be Multiply Annotated? Active Learning When Annotators May Disagree

Connor Baumler<sup>\*</sup>

Anna Sotnikova<sup>\*</sup>

Hal Daumé III

## The Ecological Fallacy in Annotation: Modelling Human Label Variation goes beyond Sociodemographics

Matthias Orlikowski<sup>1</sup>, Paul Röttger<sup>2</sup>, Philipp Cimiano<sup>1</sup>, and Dirk Hovy<sup>3</sup>

## Can Large Language Models Capture Dissenting Human Voices?

Noah Lee<sup>\*</sup>

Na Min An<sup>\*</sup>

James Thorne

## Wisdom of Instruction-Tuned Language Model Crowds. Exploring Model Label Variation

Flor Miriam Plaza-del-Arco, Debora Nozza, Dirk Hovy



# Roadmap: Selected Case Studies

- 1 Humans and Uncertainty: *The “Problem” of Human Label Variation*
- 2 Models and Uncertainty: *Stop Measuring Calibration When Humans Disagree*
- 3 How to detect errors? *ActiveAED*
- 4 Plausible variation or error? *VariERR*

# Stop Measuring Calibration When Humans Disagree

**Joris Baan<sup>1</sup>, Wilker Aziz<sup>1</sup>, Barbara Plank<sup>2,3,4</sup>, Raquel Fernández<sup>1</sup>**

<sup>1</sup>University of Amsterdam, <sup>2</sup>IT University of Copenhagen, <sup>3</sup>MCML Munich, <sup>4</sup>LMU Munich  
{j.s.baan,w.aziz,raquel.fernandez}@uva.nl, b.plank@lmu.de

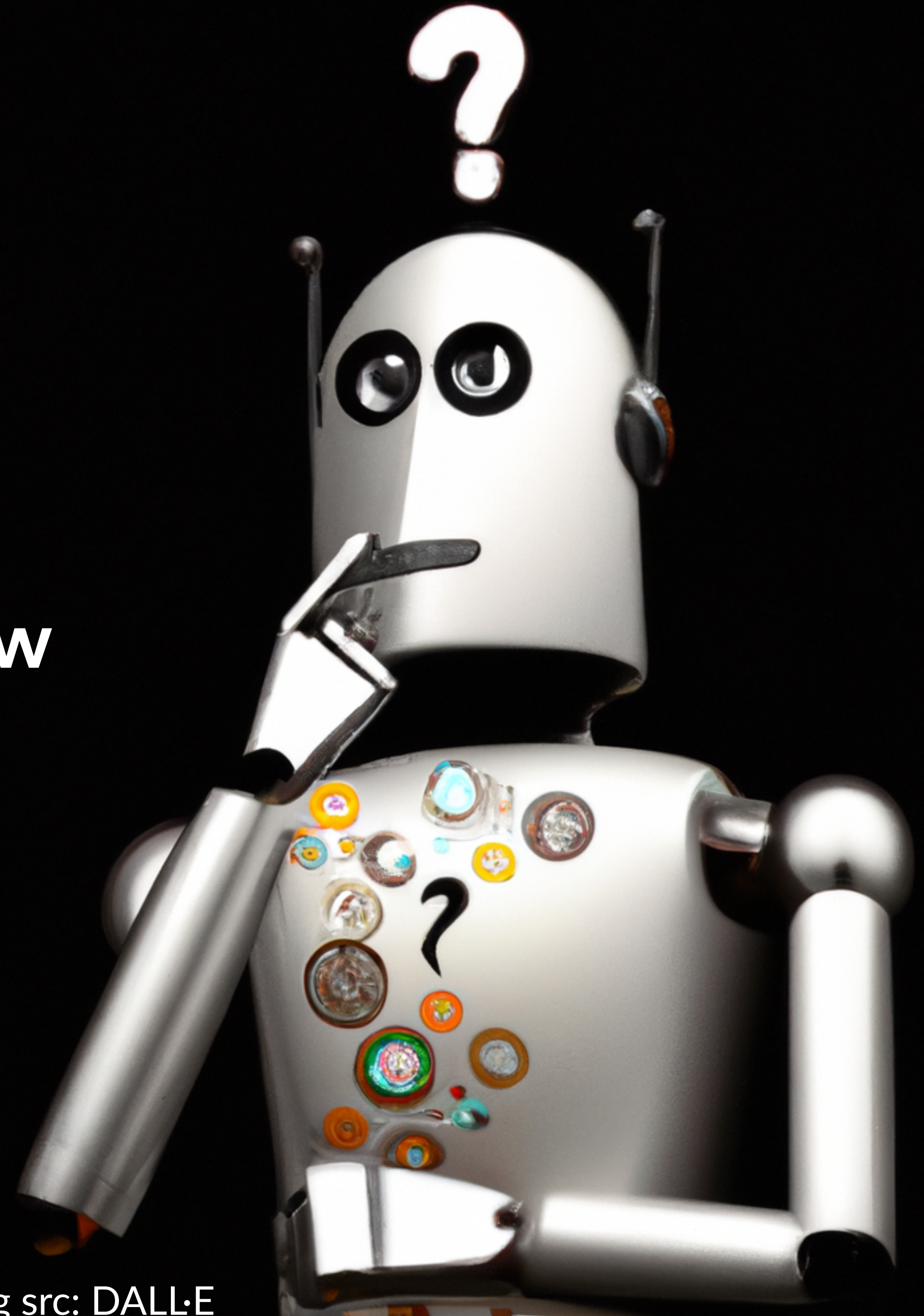


(Baan, Aziz, Plank, Fernandez, 2022 EMNLP)

**Uncertainty**



**Model Uncertainty:**  
Models don't always know  
when they don't know



Img src: DALL·E

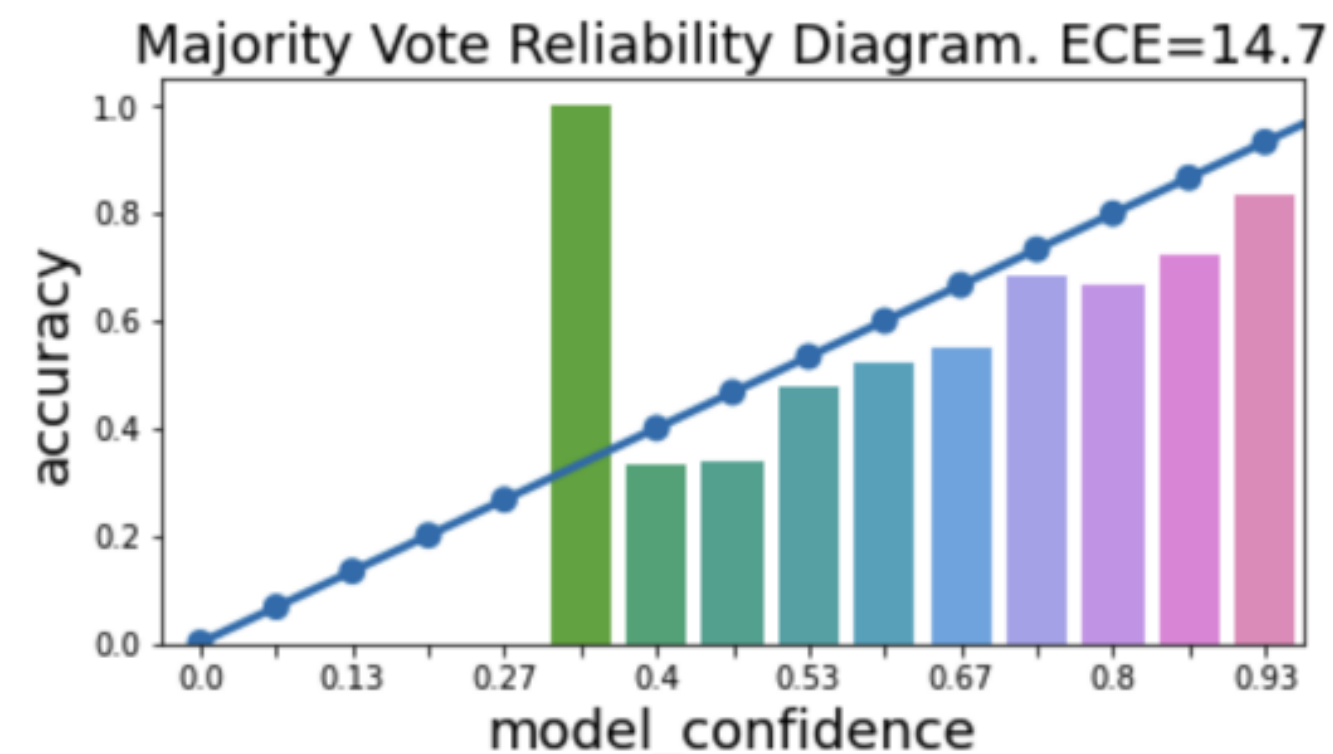




# More trustworthy models:

## Calibration & Model Uncertainty

- Calibration is a popular framework to evaluate whether a classifier knows when it does not know
- Reliability diagram to indicate how well calibrated a model is
  - ECE (expected calibration error)



- What does calibration mean when there is no ground truth?
  - We examine calibration under the lens of HLV

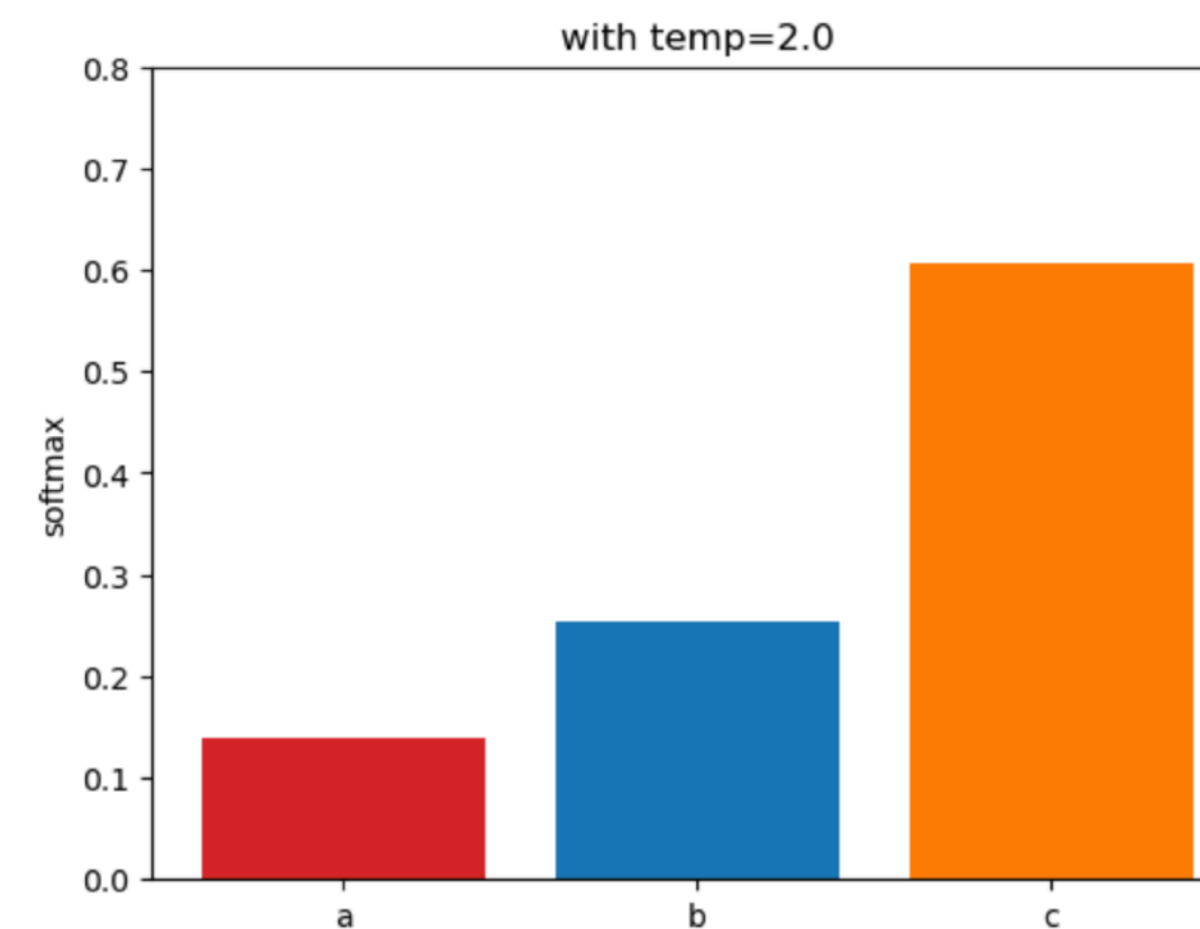
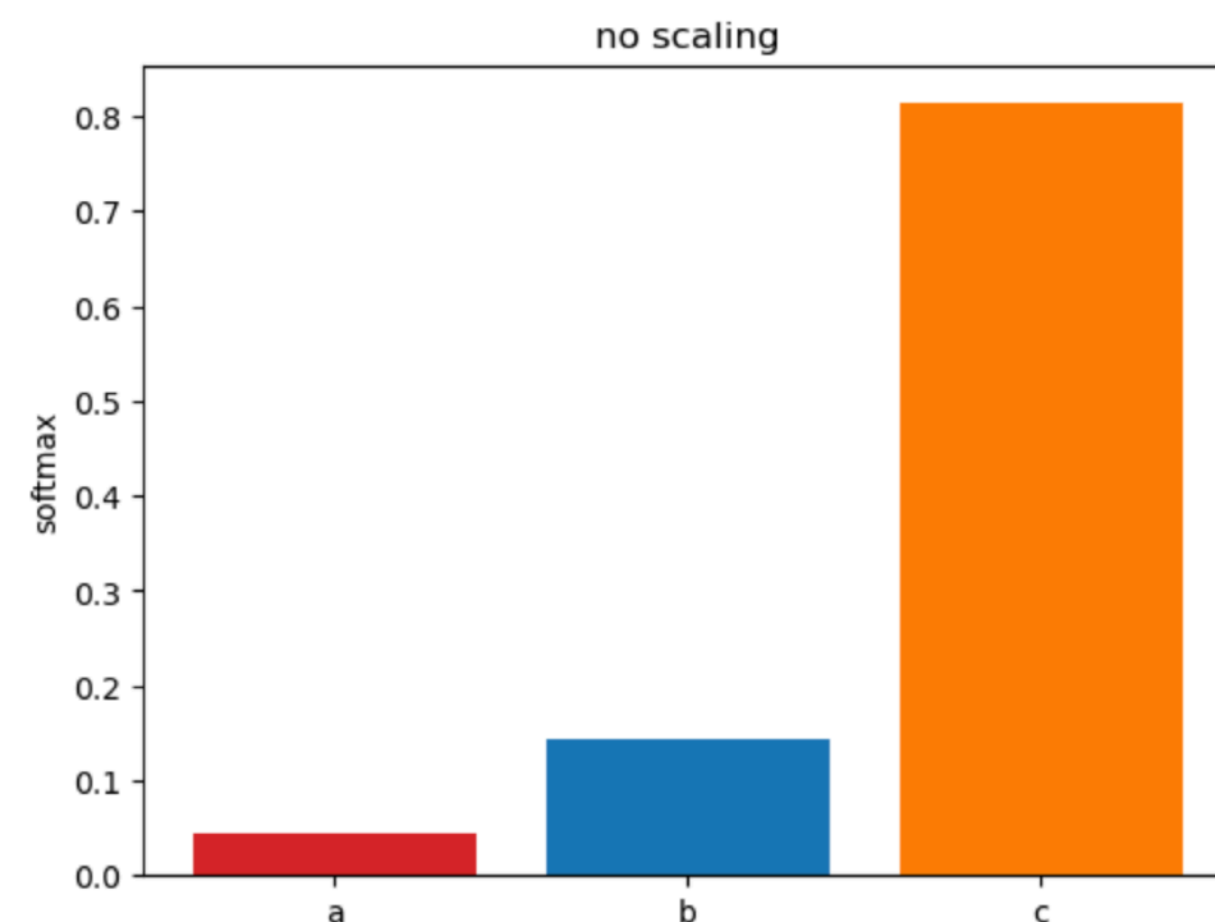


# Calibration: Temperature Scaling

- Temperature Scaling is one way to do calibration. It is a post-processing technique to improve the calibration error. It works by dividing the logits by a scalar T:

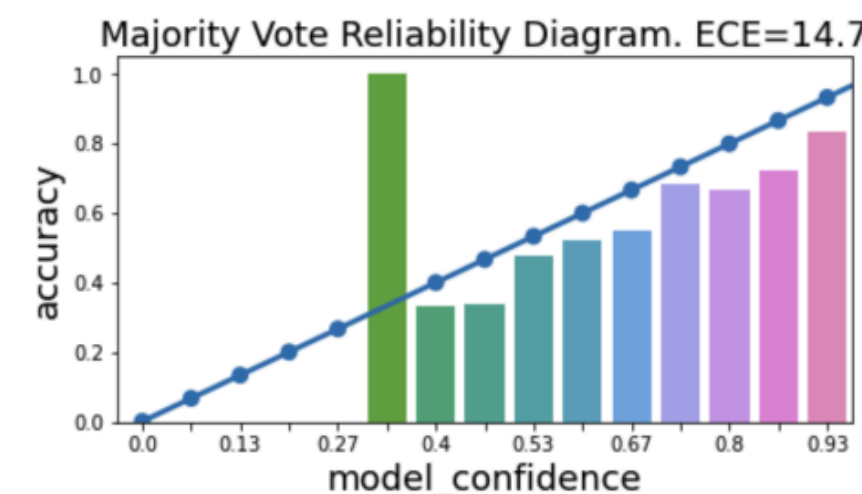
$$\text{softmax}_T = \frac{e^{z/T}}{\sum_i e^{z_i/T}}$$

- If  $T > 1.0$ , it makes the model less confident about its predictions:



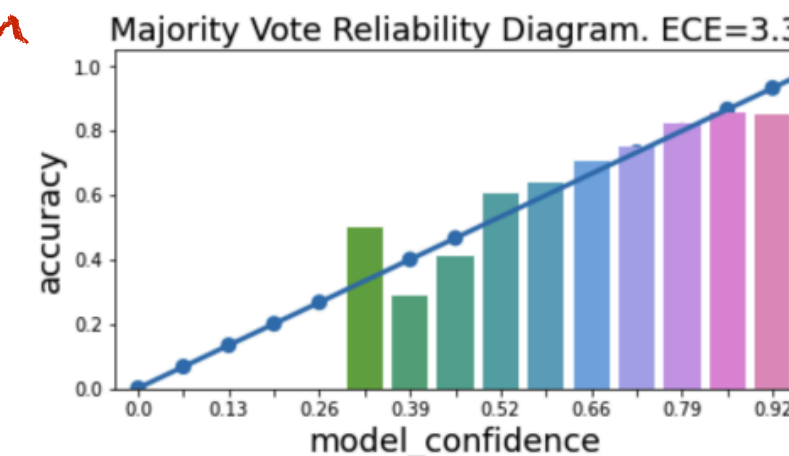
# Calibration to majority? BAD.

- ▶ Temperature Scaling can help improve ECE:



(c) ECE: Vanilla

ECE reduction  
→



(d) ECE: Temp Scaling

- ▶ However, we observe that **despite low ECE**, an oracle is still miscalibrated:

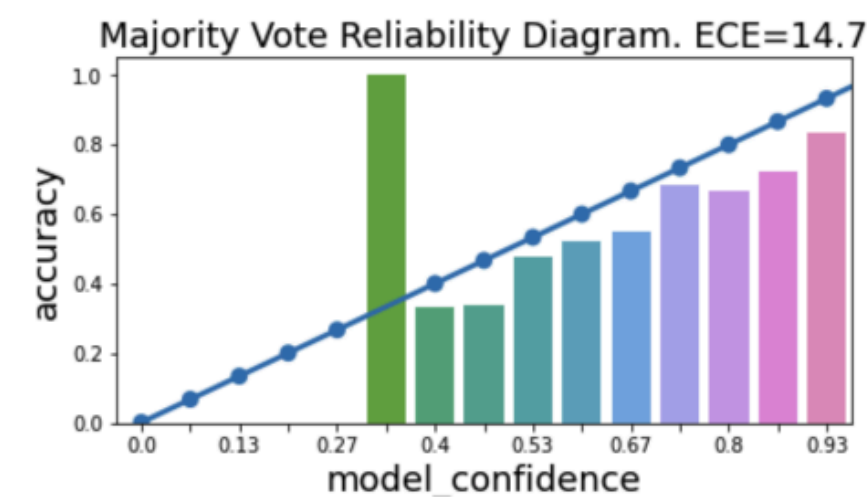
	RoBERTa	RoBERTa-TS	Oracle
Acc ↑	0.74±0.01	0.74±0.01	1.00
ECE ↓	0.14±0.01	0.03±0.01	0.25

- ▶ What can we do? Measure **Human Calibration Error (DistCE)**:
  - ▶ Total variation distance between predictive distribution and human judgement distribution (range: 0...1)

$$\text{DistCE}(x) = \text{TVD}(\mathbf{f}(x), \bar{\boldsymbol{\pi}}(x))$$

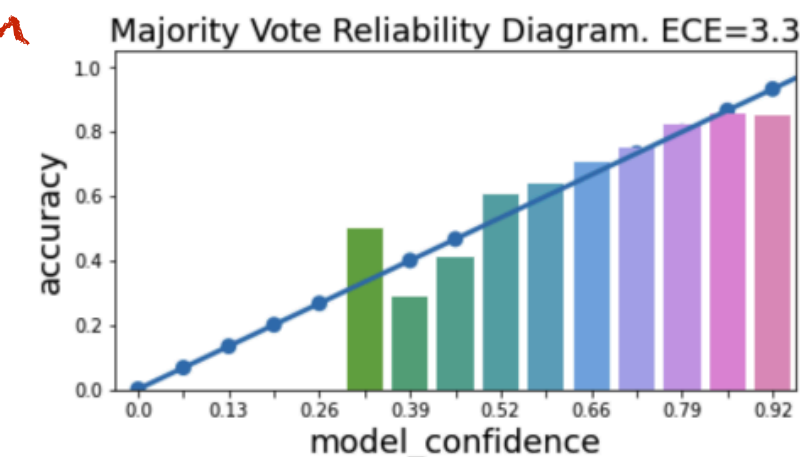
# Calibration in Light of HLV

- DistCE = instance-level analysis, enables a more fine-grained view on model calibration (Baan et al., 2022). Recall:



(c) ECE: Vanilla

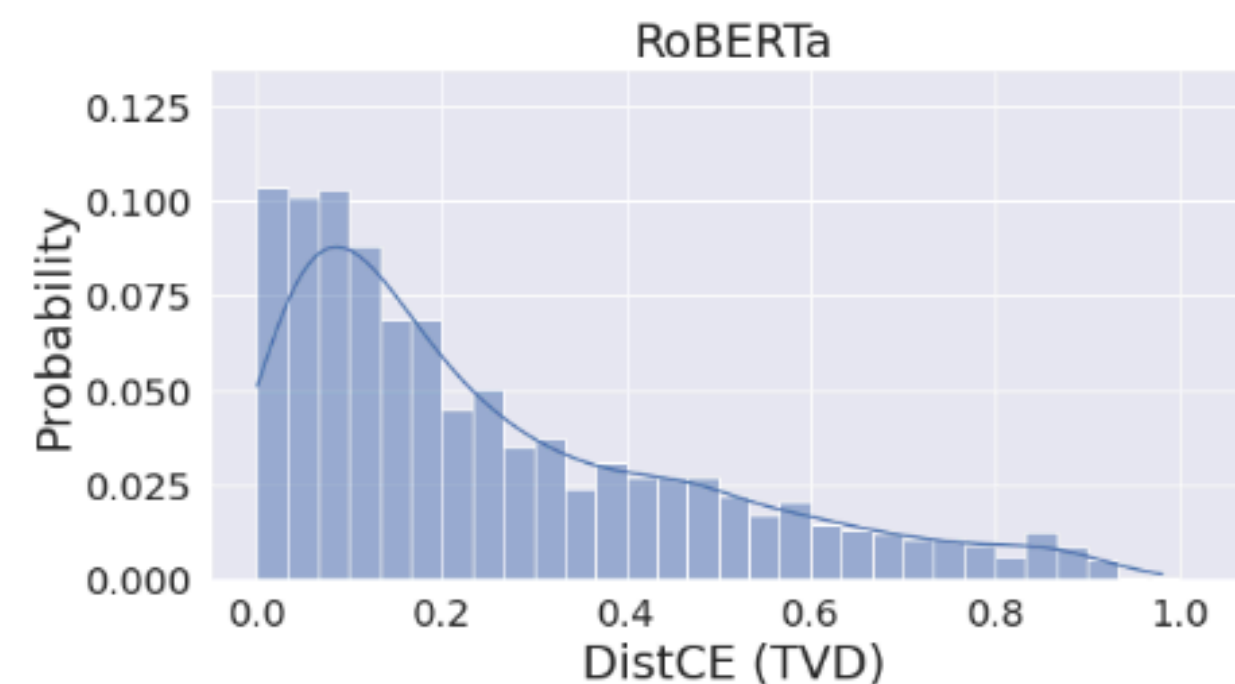
ECE reduction



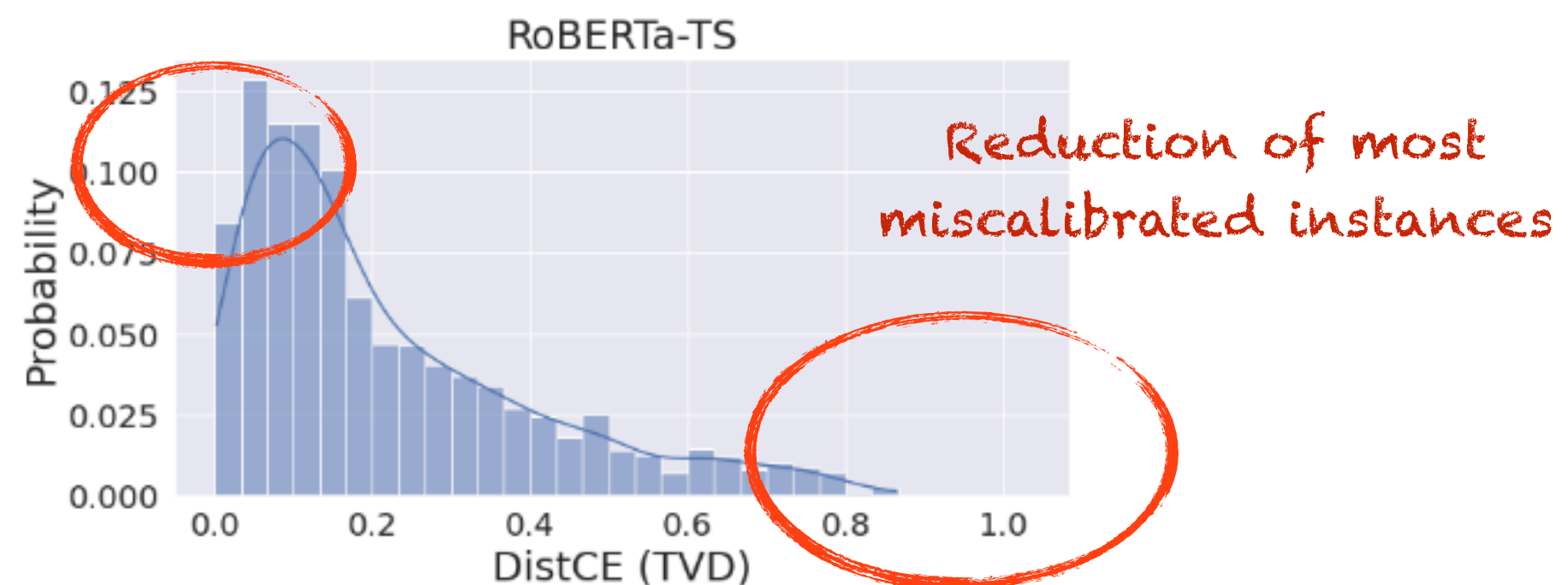
(d) ECE: Temp Scaling

- DistCE ↓ (0: perfectly calibrated to human judgement)

BUT also fewer perfectly calibrated instances!



(a) DistCE: Vanilla



(b) DistCE: Temp Scaling

# Take-home message

(Baan et al., 2022 EMNLP)

- We showed that calibration to human majority is flawed
- We suggested to look at calibration in light of HLV
  - Proposed several measures (more in the paper), incl. **Human calibration error** (DistCE), that provide us **instance-level** insights
  - More nuanced insights into model **uncertainty**
- **Limitation:** requires data with human label variation

# Roadmap: Selected Case Studies

- 1 Humans and Uncertainty: *The “Problem” of Human Label Variation*
- 2 Models and Uncertainty: *Stop Measuring Calibration When Humans Disagree*
- 3 How to detect errors? *ActiveAED*
- 4 Plausible variation or error? *VariERR*

**Existing datasets contain annotation errors**



# Data Quality

---



**Djamé..** @zehavoc · 20h

...

just found out this wonderful quote in an old paper where we described our efforts to parse the British National Corpus (100M words, back then it was huge, clusters and all) work by @Wjrigo @jenfoster, Josef van Genaboth and I  
[web.stanford.edu/group/cslipubl...](http://web.stanford.edu/group/cslipubl...)

Still applies today imho

*“Cleaning is a low-level, unglamorous task, yet crucial: The better it is done, the better the outcomes. All further layers of linguistic processing depend on the cleanliness of the data.”*

(Kilgarriff, 2007, p.149)

# Example: Sentiment (Imdb)

## Review

**\*\*SPOILERS AHEAD\*\***<br /><br />It is really unfortunate that a movie so well produced turns out to be<br /><br />such a disappointment. [...]

Lois Weber's film "Hypocrites" was and still kind of is a very bold and daring film. I enjoyed it and was very impressed by the filming and story of it. [...]

## Original Label

Positive

Negative



# Example: NER (CoNLL 2003)

## Original Annotation

1	<div>Person</div> Regula Susana Siegfried , 50 , and <div>Misc</div> Nicola <div>Person</div> Fleuchaus , 25 , were released after 71 days after a \$ 200,000 ransom was paid.
2	<div>Person</div> Laurence Courtois ( <div>Location</div> Belgium ) beat <div>Location</div> Flora <div>Person</div> Perfetti ( <div>Location</div> Italy ) 6-4 3-6 6-2
3	<div>Organization</div> Hapoel Haifa 3 <div>Org</div> Maccabi Tel Aviv 1
4	<div>Org</div> Sporting Gijon 15 4 4 7 15 22 16
5	<div>Org</div> St. Gallen 9 4 4 1 6 5 16

What to do about it?

# Annotation Error Detection (AED)

- A long-standing task (e.g. Dickinson & Meurers, 2003); recently surveyed comprehensively by Klie, Webber, Gurevych (2022)
- Typical AED methods are post-hoc processing
- We propose to combine AED with human in the loop: **Active AED**

## ActiveAED: A Human in the Loop Improves Annotation Error Detection

Leon Weber<sup>▲</sup> and Barbara Plank<sup>▲◇</sup>

<sup>▲</sup>Center for Information and Language Processing (CIS), LMU Munich, Germany

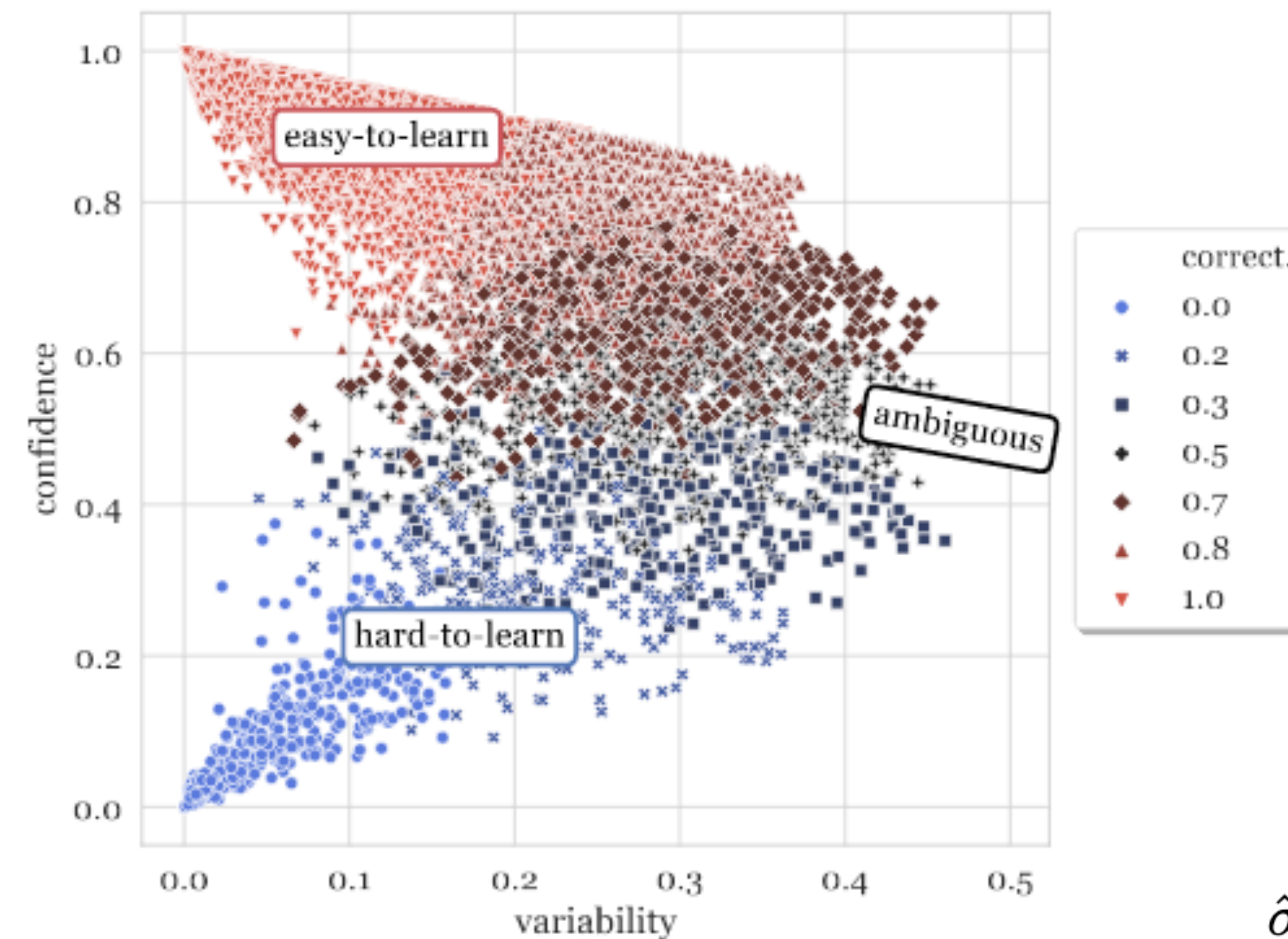
<sup>◇</sup>Munich Center for Machine Learning (MCML), Munich, Germany  
{leonweber, bplank}@cis.lmu.de



(Weber & Plank, 2023 ACL Findings)

# Dataset cartography: Training dynamics

$$\hat{\mu}_i = \frac{1}{E} \sum_{e=1}^E p_{\theta^{(e)}}(y_i^* | \mathbf{x}_i)$$



Data map for SNLI train set, based on a ROBERTA-large classifier. The x-axis shows **variability** and y-axis, the **confidence**; the colors/shapes indicate **correctness**.

$$\hat{\sigma}_i = \sqrt{\frac{\sum_{e=1}^E (p_{\theta^{(e)}}(y_i^* | \mathbf{x}_i) - \hat{\mu}_i)^2}{E}}$$

(Swayamdipta et al, 2020)



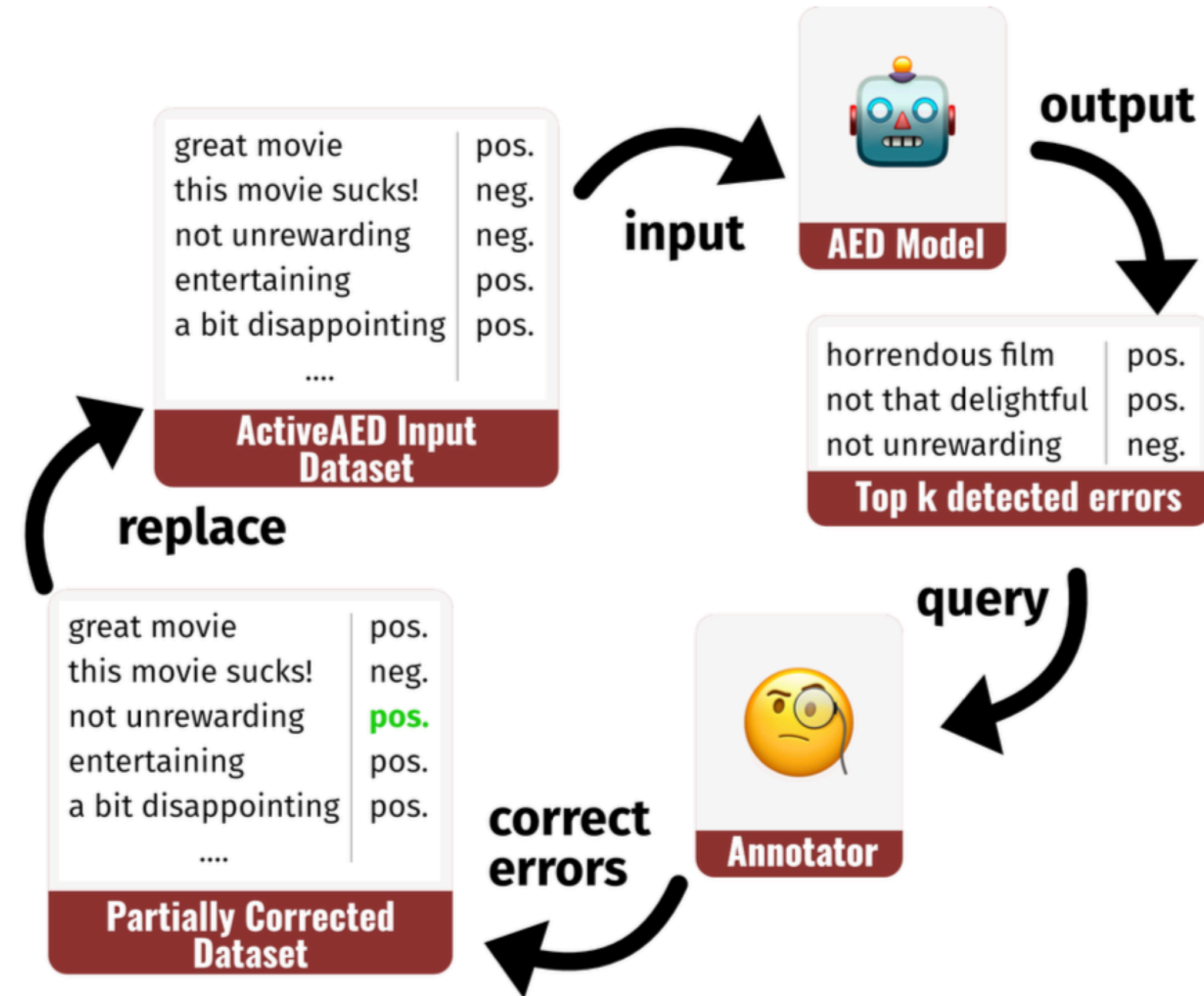
# Our solution: ActiveAED

- **ActiveAED**: Involve human annotator in pipeline, by **repeatedly querying for error corrections**
- Can be used with any **scoring-based method**. We use **Area-Under-the-Margin** (Pleiss et al. 2020)

$$s_i = \frac{1}{E} \sum_{e=1}^E \max_{y' \neq y_i} p_{\theta_e}(y'|x_i) - p_{\theta_e}(y_i|x_i)$$

- Our novel ensembling scheme merges **training-dynamics-based** and **cross-validation-based** AED for improved results

$$s_i^{train} = \frac{1}{E-1} \sum_{c \in train_i} s_{c,i}$$
$$s_i = \frac{1}{2} (s_i^{train} + s_i^{test})$$



# Main results

	ATIS	SI-Flights	IMDb	SST	GUM	CONLL-2003	SI-Companies	SI-Forex
CU	91.7±1.4	80.9±0.5	31.6±1.3	42.7±1.0	98.8±0.1	25.2±0.6	96.1±0.2	84.2±2.0
DM	97.2±0.2	79.2±2.4	30.1±3.0	47.1±1.0	99.3±0.1	30.2±0.7	97.5±0.2	80.6±0.9
AUM (p)	98.0±0.1	78.9±2.3	30.1±3.0	47.1±1.0	99.0±0.1	30.2±0.7	97.3±0.3	81.1±0.9
AUM (l)	97.3±0.4	72.6±0.3	27.5±2.5	39.6±1.3	<b>99.5±0.1</b>	29.3±0.2	97.2±0.2	66.6±1.5
ActiveAED	<b>98.6±0.1</b>	<b>86.6±0.5</b>	<b>36.6±0.1</b>	<b>53.0±0.2</b>	98.5±0.0	<b>33.3±0.2</b>	<b>99.3±0.0</b>	<b>89.7±0.6</b>
w/o active	98.7±0.1	80.3±0.6	36.0±0.4	52.9±0.4	98.4±0.0	31.7±0.4	97.9±0.1	85.5±0.6

## Original Annotation

1	<div>Person</div> Regula Susana Siegfried, 50, and <div>Misc</div> Nicola <div>Person</div> Fleuchaus, 25, were released after 71 days after a \$ 200,000 ransom was paid.
2	<div>Person</div> Laurence Courtois ( <div>Location</div> Belgium ) beat <div>Location</div> Flora <div>Person</div> Perfetti ( <div>Location</div> Italy ) 6-4 3-6 6-2
3	<div>Organization</div> Hapoel Haifa 3 <div>Org</div> Maccabi Tel Aviv 1
4	<div>Org</div> Sporting Gijon 15 4 4 7 15 22 16
5	<div>Org</div> St. Gallen 9 4 4 1 6 5 16

## Corrected Annotation

1	<div>Person</div> Regula Susana Siegfried, 50, and <div>Person</div> Nicola Fleuchaus, 25, were released after 71 days after a \$ 200,000 ransom was paid.
2	<div>Person</div> Laurence Courtois ( <div>Location</div> Belgium ) beat <div>Person</div> Flora Perfetti ( <div>Location</div> Italy ) 6-4 3-6 6-2
3	<div>Organization</div> Hapoel Haifa 3 <div>Organization</div> Maccabi Tel Aviv 1
4	<div>Organization</div> Sporting Gijon 15 4 4 7 15 22 16
5	<div>Organization</div> St. Gallen 9 4 4 1 6 5 16



# Conclusion

(Weber & Plank, 2023)

ActiveAED: 🧐 A human in the loop improves  
Annotation Error Detection.

So far studied on AED were limited to  
(discriminative) classification tasks

# DONKII: Characterizing and Detecting Errors in Instruction-Tuning Datasets

Leon Weber-Genzel<sup>▲</sup><sup>🏢</sup> and Robert Litschko<sup>▲</sup><sup>🏢</sup> and Ekaterina Artemova<sup>▲</sup><sup>\*</sup>  
and Barbara Plank<sup>▲</sup><sup>🏢</sup><sup>✉</sup>

▲ MaiNLP, Center for Information and Language Processing, LMU Munich, Germany

🏢 Munich Center for Machine Learning (MCML), Munich, Germany

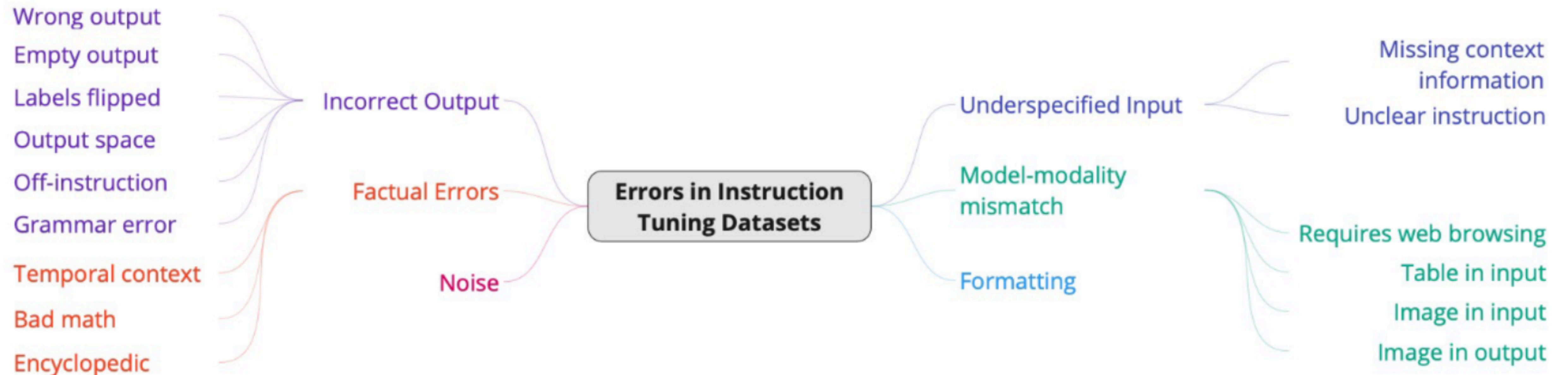
✉ Department of Computer Science, IT University of Copenhagen, Denmark

{leonweber, robert.litschko, b.plank}@lmu.de

(Weber et al., 2024 LAW@EACL 2024)



# Taxonomy of Error Types



# Examples

Error	Source	Instruction	Input	Output (shortened)
Requires web browsing	Alpaca	Summarize the New York Times article about the rise of AI.	<a href="https://www.nytimes.com/2020/06/16/opinion/ai-growth.html">www.nytimes.com/2020/06/16/opinion/ai-growth.html</a>	The New York Times article titled “The Rise of Artificial Intelligence” looks at the growing impact of AI on the world today. <...>
Bad math	Alpaca	Calculate the mean of the numbers 3, 4, 6, 10, 11		The mean of the numbers 3, 4, 6, 10, 11 is 6.4. This is calculated by adding the numbers together and then dividing by the number of values (5). $(3 + 4 + 6 + 10 + 11 = 34; 34 / 5 = 6.4)$ .
Unclear instruction	Alpaca	Sort the following 10 fruits into three categories.	Bananas, Apples, Pineapples, Plums, Cherries, Peaches, Grapes, Watermelons, Mangoes, Oranges	Category 1: Bananas, Apples, Pineapples, Plums, Cherries;

🤔 What's an error vs plausible variation?



# Roadmap: Selected Case Studies

- 1 Humans and Uncertainty: *The “Problem” of Human Label Variation*
- 2 Models and Uncertainty: *Stop Measuring Calibration When Humans Disagree*
- 3 How to detect errors? *ActiveAED*
- 4 Plausible variation or error? *VariERR*

# Motivation

- While **Human Label Variation** exists, so do **errors**.
- Annotators are inevitably prone to make errors.
- We lack both a theory and operationalizable procedures to answer the RQ:
  - Can we tease apart error from plausible human label variation?



**Error** vs. plausible **Human Label Variation**

## VARIERR NLI: Separating Annotation Error from Human Label Variation

Leon Weber-Genzel<sup>▲</sup><sup>🏢</sup>\* Siyao Peng<sup>▲</sup><sup>🏢</sup>\* Marie-Catherine de Marneffe<sup>✍</sup> Barbara Plank<sup>▲</sup><sup>🏢</sup>

<sup>▲</sup> MaiNLP, Center for Information and Language Processing, LMU Munich, Germany

<sup>🏢</sup> Munich Center for Machine Learning (MCML), Munich, Germany

<sup>✍</sup> FNRS, UCLouvain, Belgium

{leonweber, siyaopeng, bplank}@cis.lmu.de marie-catherine.demarneffe@uclouvain.be

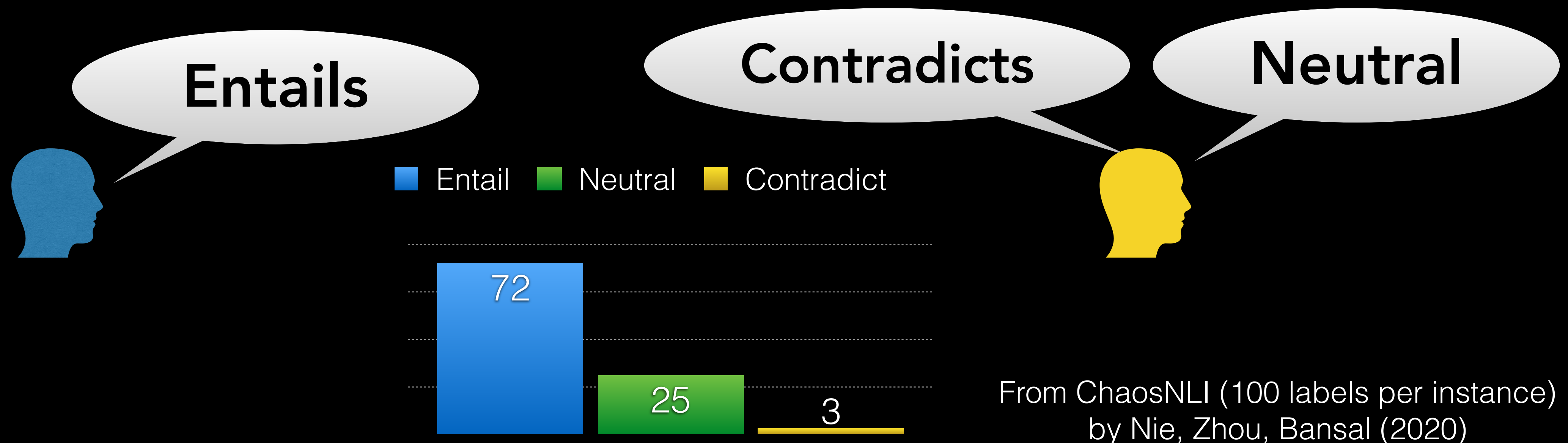
(Weber-Genzel, Peng et al., 2024 To Appear at ACL)



# Natural Language Inference

**Premise:** As he stepped across the threshold, Tommy brought the picture down with terrific force on his head.

**Hypothesis:** Tommy hurt his head bringing the picture down.



# We propose a two step-procedure: 1) Explanations

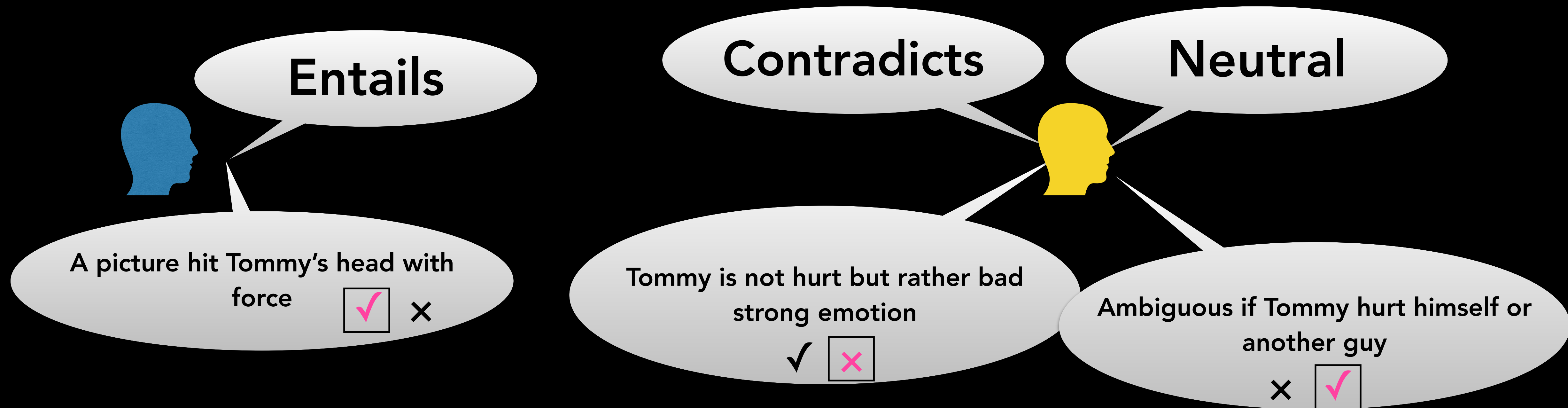
*Premise:* As he stepped across the threshold, Tommy brought the picture down with terrific force on his head.

*Hypothesis:* Tommy hurt his head bringing the picture down.



► **Ecologically valid explanations** inspired by (Jiang et al., 2023)

# We propose a two step-procedure: 2) Validations



- ▶ Another kind of validation: see your own and peer's label-explanation pairs



# ValiErr: Defining Errors



Entails



A picture hit Tommy's head with force



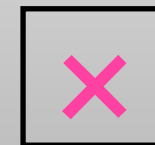
×



Contradicts



Tommy is not hurt but rather bad strong emotion



Neutral

Ambiguous if Tommy hurt himself or another guy

×



- **Self-validated:** any self-validated label-explanation pair is plausible, otherwise it is an **error**
- **Peer-validated:** A label-explanation pair is peer-validated if  $\geq 2$  annotators approved it

# Example from VariErr NLI:

*Premise:* Because marginal costs are very low, a newspaper price for preprints might be as low as 5 or 6 cents per piece.

*Hypothesis:* Newspaper preprints can cost as much as \$5.

*Label-explanation pairs:* Before:{E:1,N:2,C:1} Self-validated:{N:2} Peer-validated:{N:2,C:1}

*Label:* [N] *Errors:* [E, C]

Round 1: NLI Label & Explanation			Round 2: Validity			
L	A	Explanation	1	2	3	4
E	4	5 dollars for a piece of newspaper.	×	×	×	×
N	1	The context only mentions how low the price may be, not how high it may be.	✓	✓	✓	✓
	3	The maximum cost of newspaper preprints is not given in the context.	✓	✓	✓	✓
C	2	The context says 5 or 6 cents, not \$5.	×	×	✓	✓

(a) *id:* 72870c

**Table 1:** Sample annotations from VARIERRNLI corpus. L: Label, A: Annotator; E: Entailment, N: Neutral, C: Contradiction; ✓: ‘yes’; ×: ‘no’; ?: ‘idk’; magenta: self-judgments, black: peer-judgments, Err: label error.

# VariErr dataset

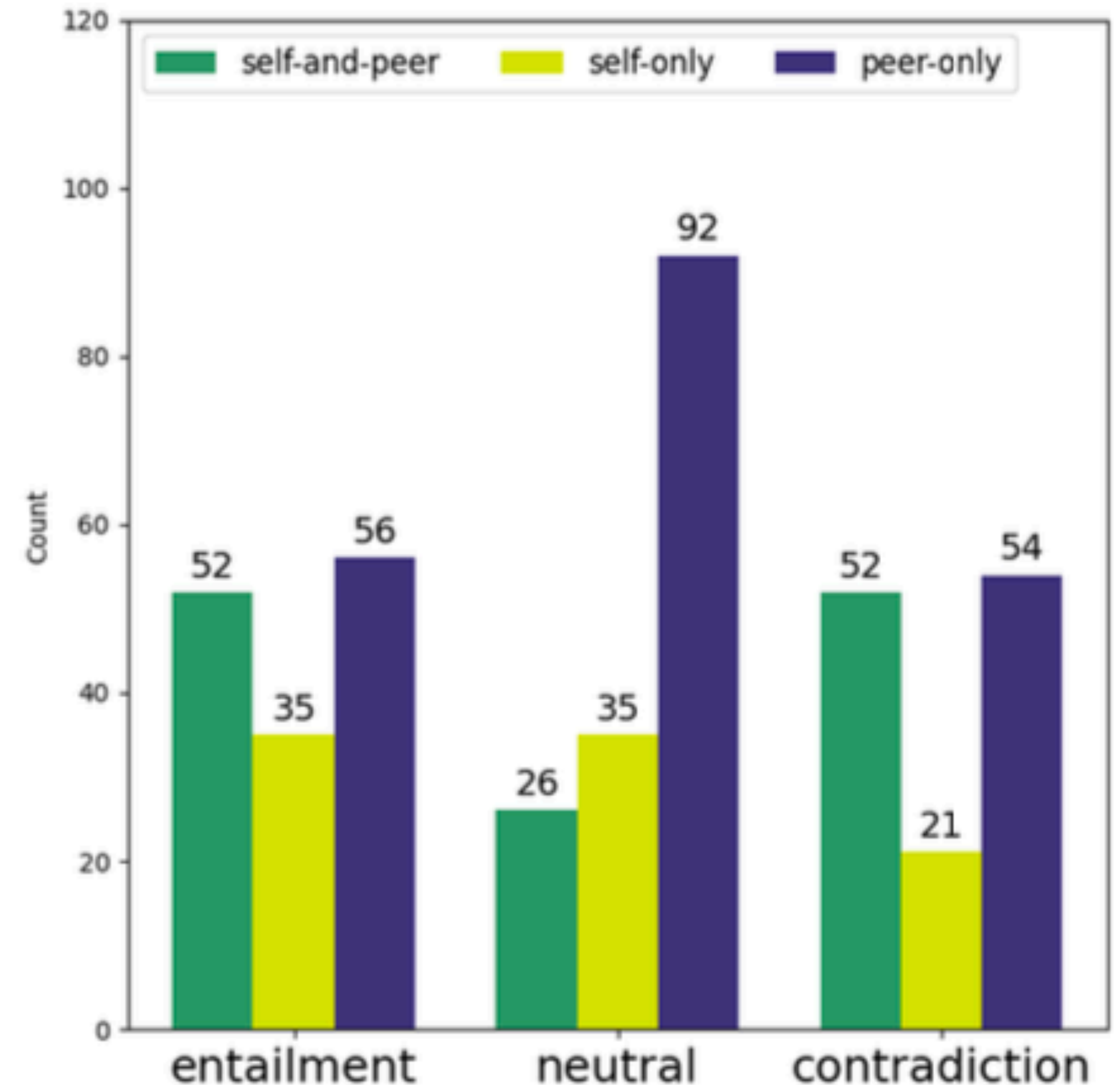
- VariErr NLI: re-annotated 500 NLI items from scratch, 1,933 label-explanation pairs
- 88.57% (1,712/1,933) are self-validated, 82.82% are peer-validated (1,601/1,933)
- Overall, 37% of items had self-identified errors (188/500)

Validation	FreqType	E	N	C	$\Sigma$	IAA
before validation	<i>repeated</i>	554	977	402	1,933	0.35
	<i>aggregated</i>	263	403	212	878	
self-validated	<i>repeated</i>	467	916	329	1,712	0.50
	<i>aggregated</i>	210	380	159	749	
peer-validated	<i>repeated</i>	446	859	296	1,601	0.69
	<i>aggregated</i>	177	335	130	642	



# Statistics on VariErr

- Number of label-explanation pairs that were rejected in phase 2
- Most Entailment and Contradiction annotations are rejected by both self- and peer-validations



# How good is Annotation Error Detection on VariErr?

- We model AED as a ranking task
- scorer to rank the list of labels with errors high
- from 500 items, give list of 878 item-label pairs to scorer
- compare ranked lists to self-flagged errors
- Metrics: Average Precision (AP), P/R/F1 of top 100 ranked labels P@100, R@100
- Five AED models: two variants of datamaps, metadata archaeology, two GPTs\* (GPT-3.5 and GPT-4)

System:

You are an expert linguistic annotator.

User:

We have collected annotations for an NLI instance together with reasons for the labels. Your task is to judge whether the reasons make sense for the label. Provide the probability (0.0 - 1.0) that the reason makes sense for the label. Give ONLY the reason and the probability, no other words or explanation. For example:

Reason: <The verbatim copy of the reason>  
Probability: <the probability between 0.0 and 1.0 that the reason makes sense for the label, without any extra commentary whatsoever; just the probability!>.

Context: {CONTEXT}

Statement: {STATEMENT}

Reason for label {LABEL}: {REASON\_1}

Reason for label {LABEL}: {REASON\_2}

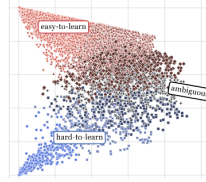
[...]

Reason for label {LABEL}: {REASON\_n}

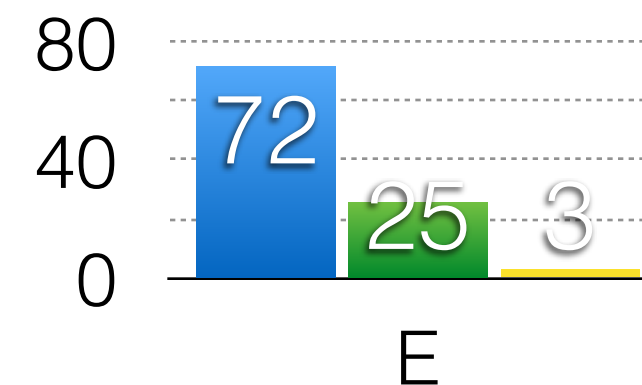
Reason {REASON\_1}

Probability:

# How good is Annotation Error Detection on VariErr?

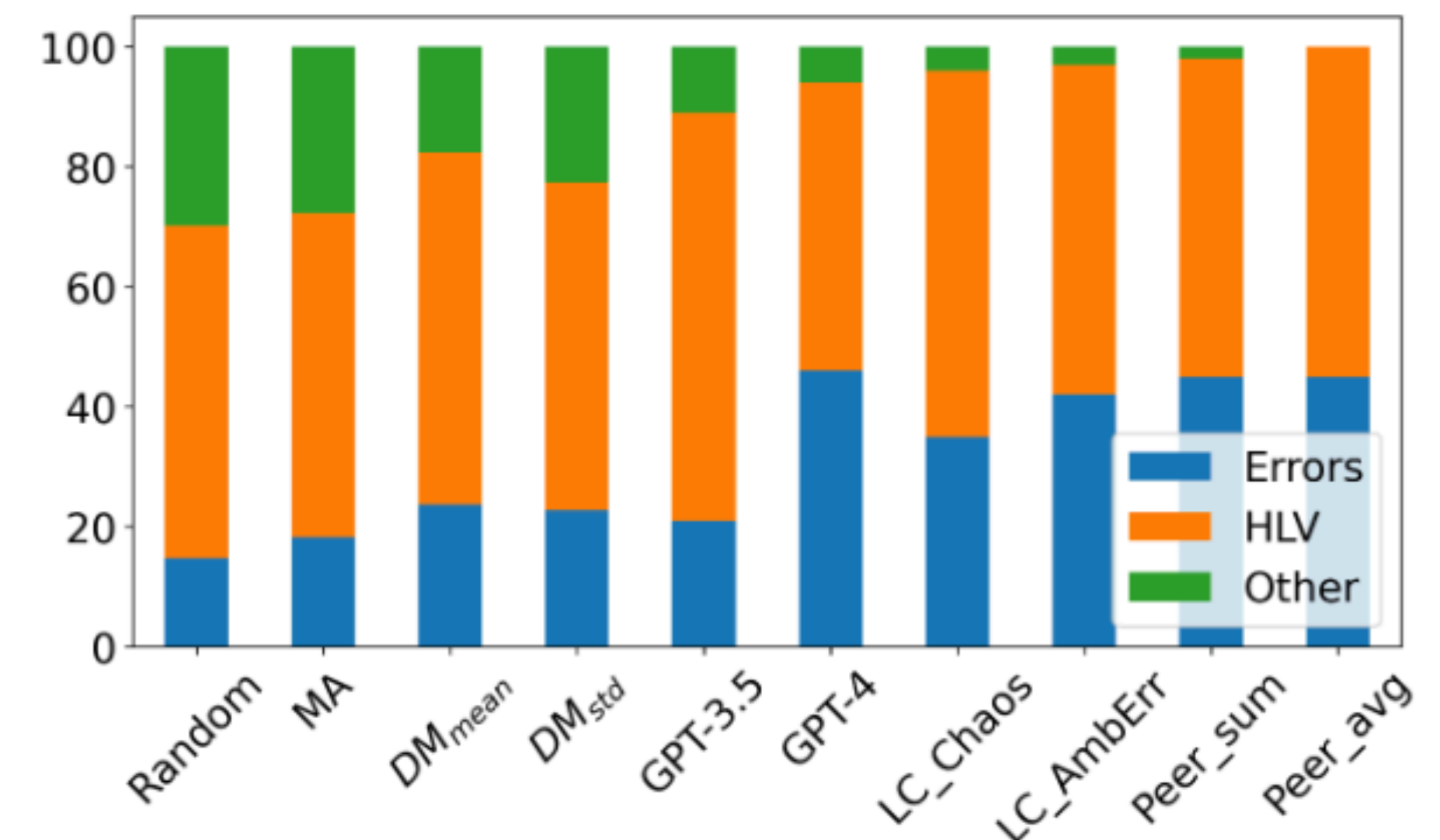
- Random baseline AP **14.7**
- Data Maps: 22 AP 
- GPTs: GPT-3.5 **17.6**, GPT-4 **31.3**
- Human label count heuristics:

- 32.5 (ChaosNLI 100 voters)
- 40.8** (4 voters)



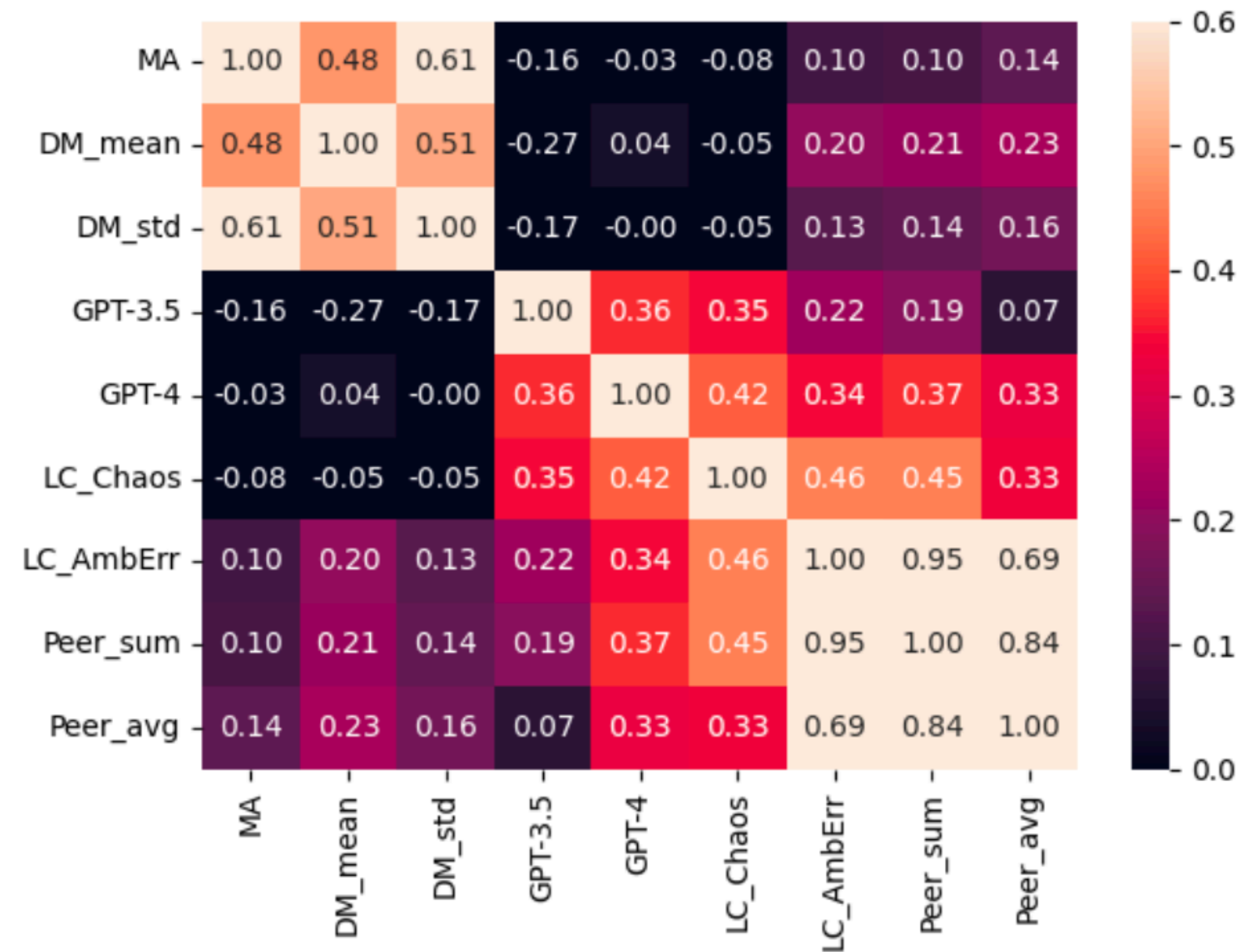
- Human heuristics outperforms GPTs, best with explanations:
- Peer heuristics from VariErr:
  - 46.5** (sum peer-validations)

- Human validation is a strong means to detect errors in data with high HLV
- Heuristics from VariErr performs (4 > 100)
- Analysis: What instances were selected?





# Complementarity



- GPT-4 correlation to LC\_chaos is 0.42

**Figure 3:** Correlations among scorer predictions.

high-stake human  
decision support  
(e.g. law)

learning from less but  
higher quality data

active learning

model uncertainty

# ***Human Label Variation***

**- many exciting connections -**

human values and LLM  
alignment

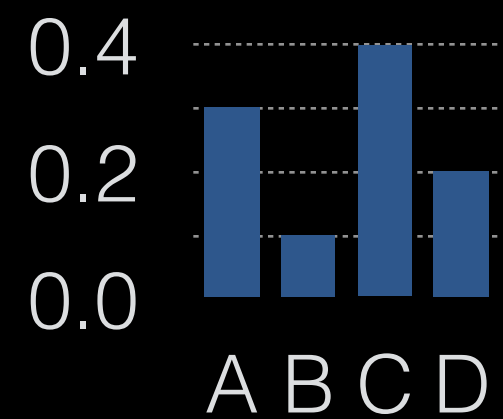
LLMs that react as  
humans do

statistics and data-  
generation process

# Take-home message



- ✓ **Human label variation** is signal (annotation errors though do exist)



- ✓ Let's **embrace** it in all stages of the AI pipelines - to not continue to model only the "mode"



- ✓ HLV will help us develop **trustworthy human-facing AI**



# From Human Label Variation and Model Uncertainty to Error Detection (and Back)?



IT UNIVERSITY OF COPENHAGEN

Thanks to my research team, collaborators and funders:



<http://mainlp.github.io>