

Google DeepMind

Evaluation and calibration of AI models with *uncertain* ground truth

David Stutz

Jul 11 2023

Outline

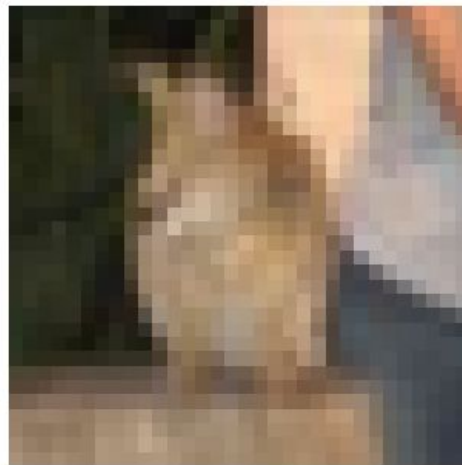
- ❑ Motivation through annotator disagreement
- ❑ Statistical model
- ❑ Measuring uncertainty
- ❑ Evaluating AI models
- ❑ Case study in dermatology:
 - ❑ Results
 - ❑ Bonus: calibration
- ❑ Conclusion and outlook

Paper: arxiv.org/abs/2307.02191

Evaluating AI systems under uncertain ground truth: a case study in dermatology

David Stutz^{*,@,1}, Ali Taylan Cemgil^{*,@,1}, Abhijit Guha Roy^{*,@,2}, Tatiana Matejovicova^{*,1}, Melih Barsbey^{*,1,3}, Patricia Strachan², Mike Schaeckermann², Jan Freyberg², Rajeev Rikhye², Beverly Freeman², Javier Perez Matos², Umesh Telang², Dale R. Webster², Yuan Liu², Greg S. Corrado², Yossi Matias², Pushmeet Kohli¹, Yun Liu^{2,+}, Arnaud Doucet^{1,+} and Alan Karthikesalingam^{2,+}

¹Google DeepMind, ²Google, ³Bogazici University, ^{*}Equal first authors, [@]Corresponding authors, ⁺Equal last authors



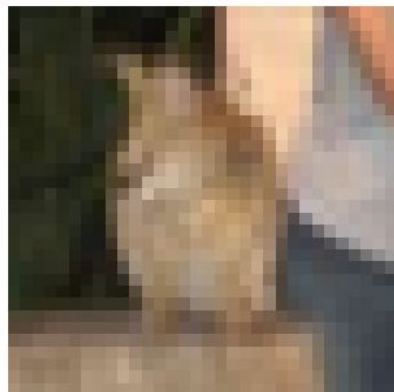
“Bird”, “cat”, or “frog”?



“Hemangioma” or “Melanoma”?
Benign or cancer?

Standard evaluation of supervised models

Observation



Correct/good prediction?



AI prediction

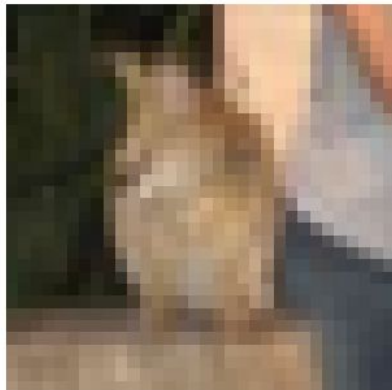
“bird”

Standard evaluation of supervised models

Unknown
true label

?

Observation



Correct/good prediction?

AI prediction

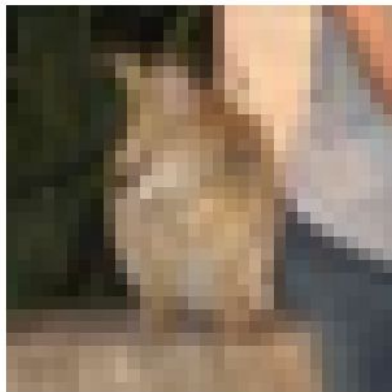
“bird”

Standard evaluation of supervised models

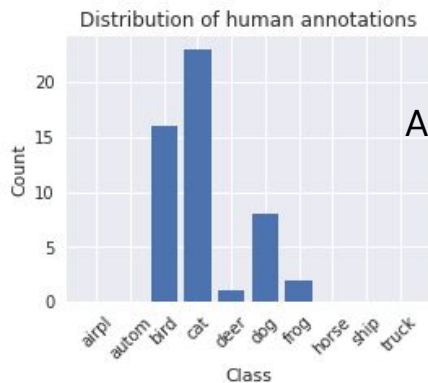
Unknown
true label

?

Observation



Annotations



← ? →

AI prediction

“bird”

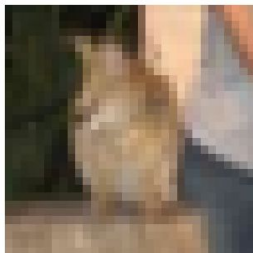
Annotators disagree!

Standard evaluation of supervised models

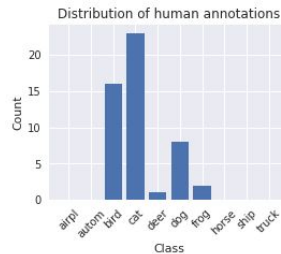
Unknown
true label

?

Observation



Annotations



Majority vote

"cat"

AI prediction

"bird"

Incorrect!



Standard evaluation of supervised models

Unknown
true label

?

Observation



Annotations

- b¹**: {*Pyogenic granuloma* (Low)} {*Hemangioma* (Med)}
{*Melanoma* (High)}
- b²** {*Angiokeratoma of skin* (Low)} {*Atypical Nevus* (Med)}
- b³**: {*Hemangioma* (Med)} {*Melanocytic Nevus* (Low),
Melanoma (High), *O/E - ecchymoses present* (Low)}
- b⁴**: {*Hemangioma* (Med), *Melanoma* (High), *Skin Tag* (Low)}
- b⁵**: {*Melanoma* (High)}
- b⁶**: {*Hemangioma* (Med)} {*Melanoma* (High)} {*Melanocytic Nevus* (Low)}

Conditions, Low/Med/High risk conditions

← ? →

AI prediction

“Hemangioma”

Majority vote non-trivial

Standard evaluation of supervised models

Unknown
true label

?

Observation



Annotations

b¹: {*Pyogenic granuloma* (Low)} {*Hemangioma* (Med)}
{*Melanoma* (High)}
b² {*Angiokeratoma of skin* (Low)} {*Atypical Nevus* (Med)}
b³: {*Hemangioma* (Med)} {*Melanocytic Nevus* (Low),
Melanoma (High), *O/E - ecchymoses present* (Low)}
b⁴: {*Hemangioma* (Med), *Melanoma* (High), *Skin Tag* (Low)}
b⁵: {*Melanoma* (High)}
b⁶: {*Hemangioma* (Med)} {*Melanoma* (High)} {*Melanocytic
Nevus* (Low)}

AI prediction *set*

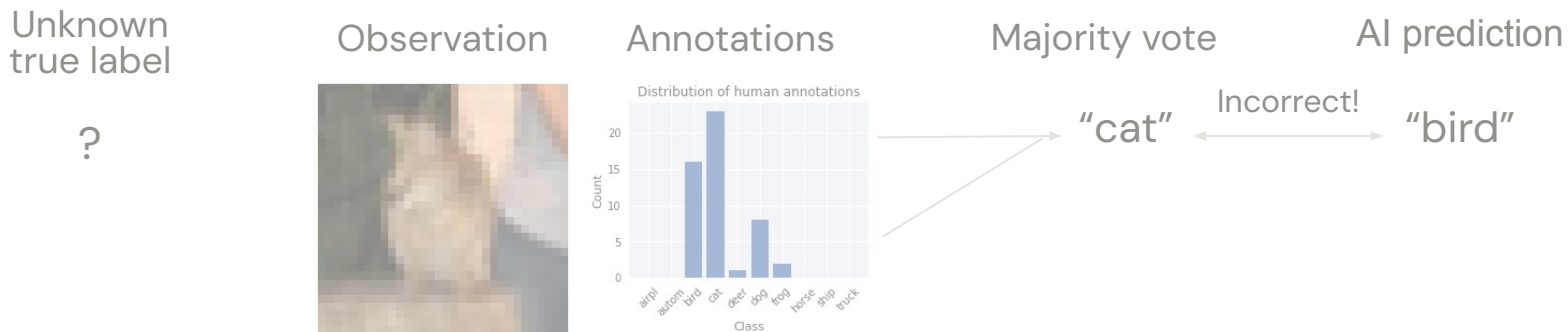
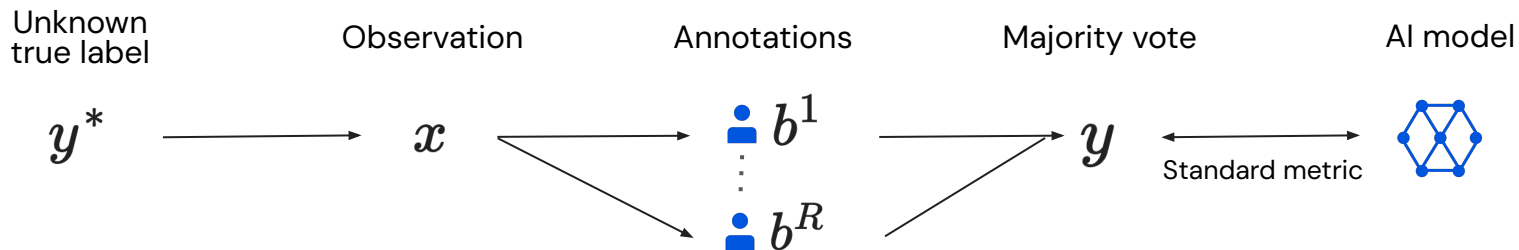
?

“**Hemangioma**”

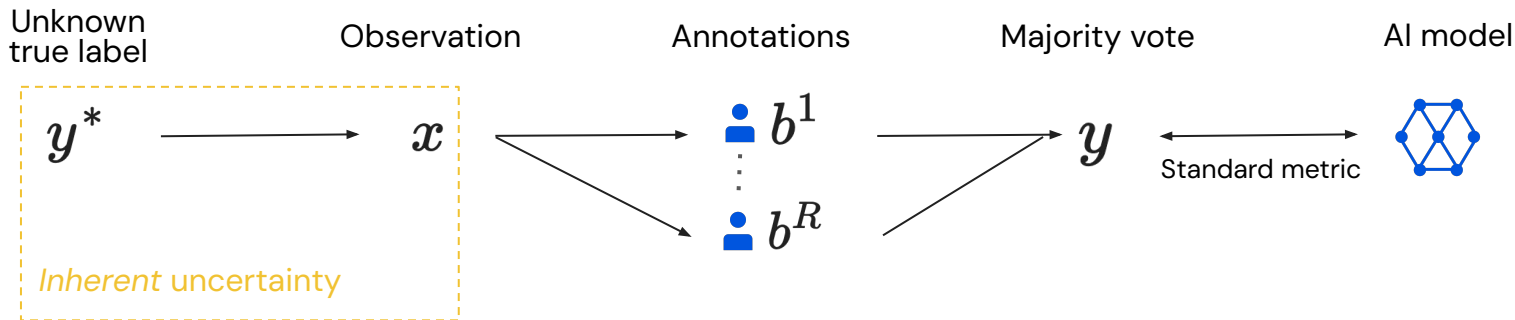
“**Atypical
Nevus**”

“**Melanocytic
Nevus**”

Standard evaluation of supervised models



Inherent uncertainty

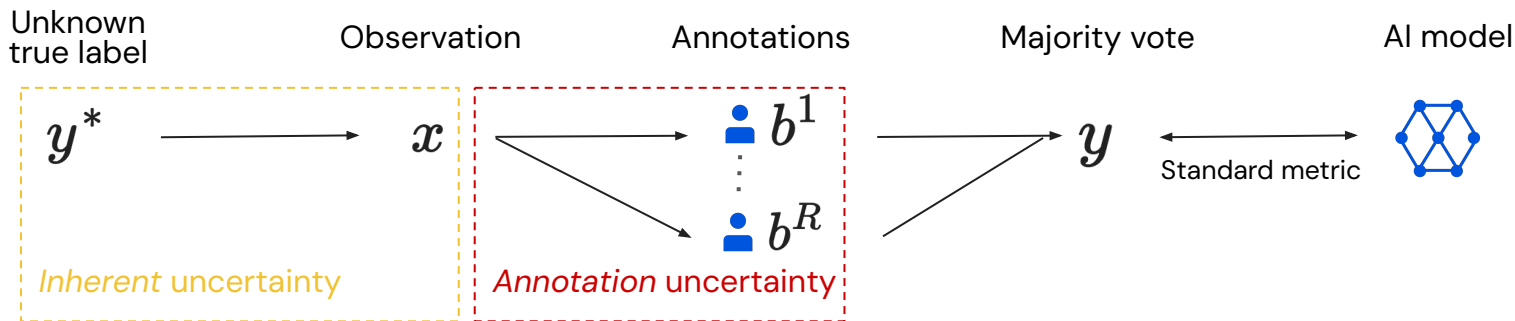


Inherent uncertainty = limited observational information:
(typically called data uncertainty)

- Low-resolution images in image recognition (e.g., CIFAR10)
- Single 2D view in 3D reconstruction
- Missing meta information or no option to question the patient in health
- ...

TL;DR: $p(y^*|x)$ is not one-hot and has high entropy!

Annotation uncertainty

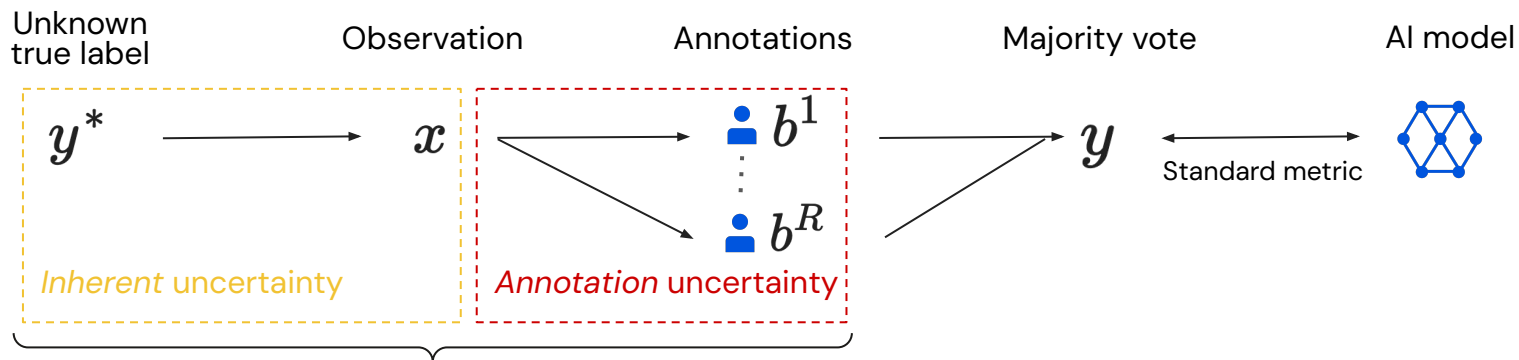


Annotation uncertainty = uncertainty induced through human annotators:

- Subjective tasks
- Inexperience of annotators
- Insufficient training of annotators
- Inappropriate annotation tool
- Different biases or background from annotators

TL;DR: annotation is difficult.

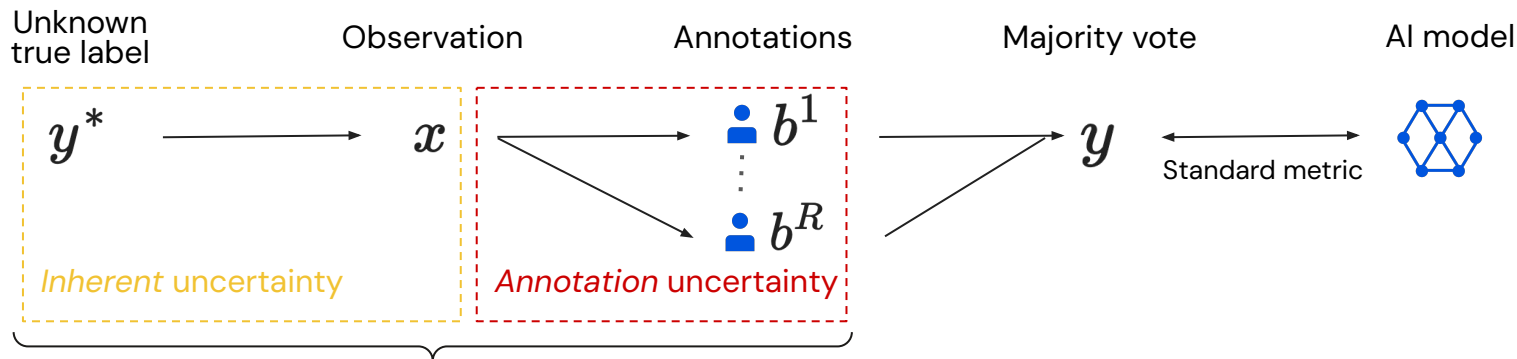
Ground truth uncertainty



Ground truth uncertainty = inherent + annotation uncertainty

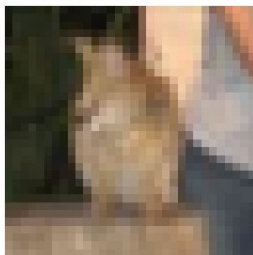
- We observe both through annotation **disagreement**

Ground truth uncertainty

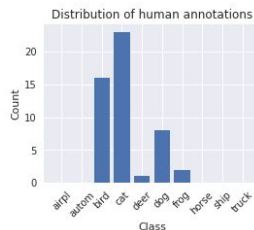


Ground truth uncertainty = inherent + annotation uncertainty

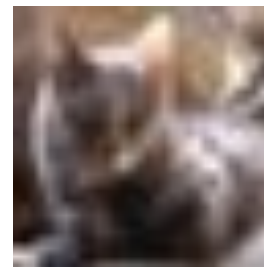
- We observe both through annotation disagreement
- Usually we cannot disentangle between inherent and annotation uncertainty



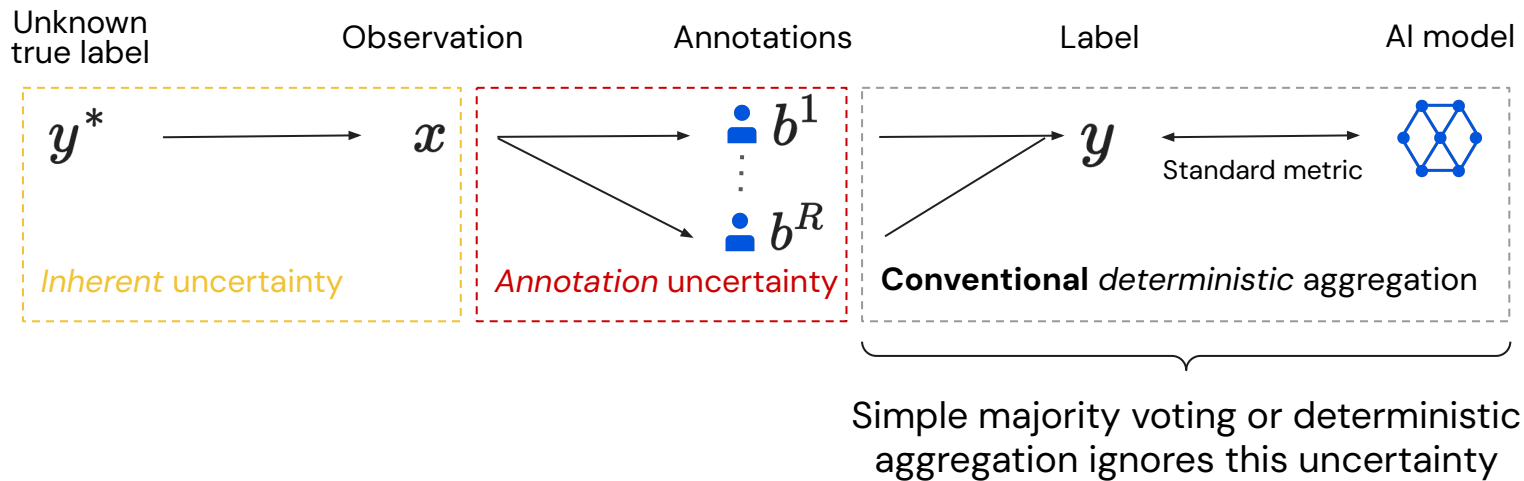
Inherent
uncertainty



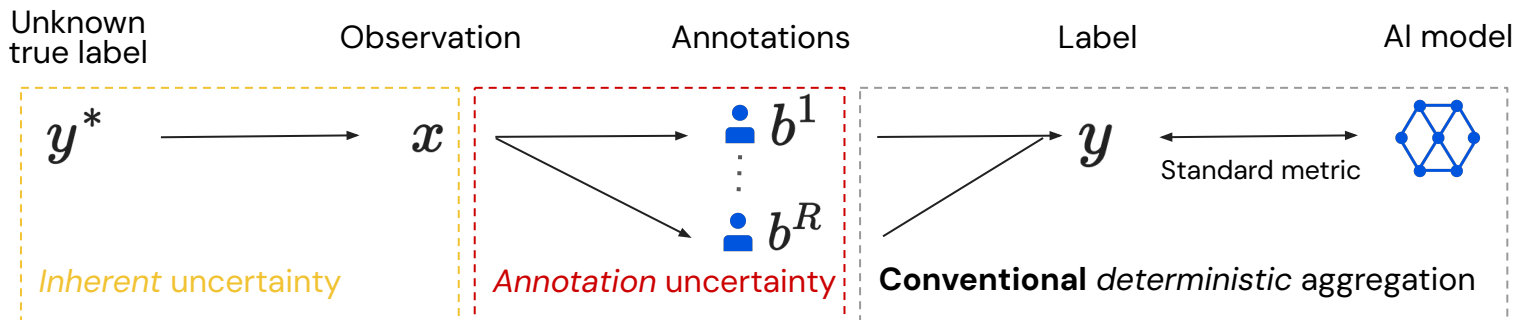
Annotation
uncertainty



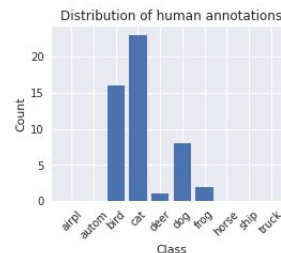
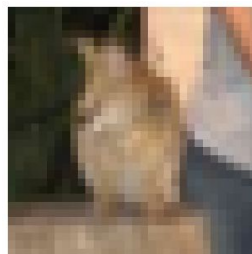
Deterministic aggregation ignores uncertainty



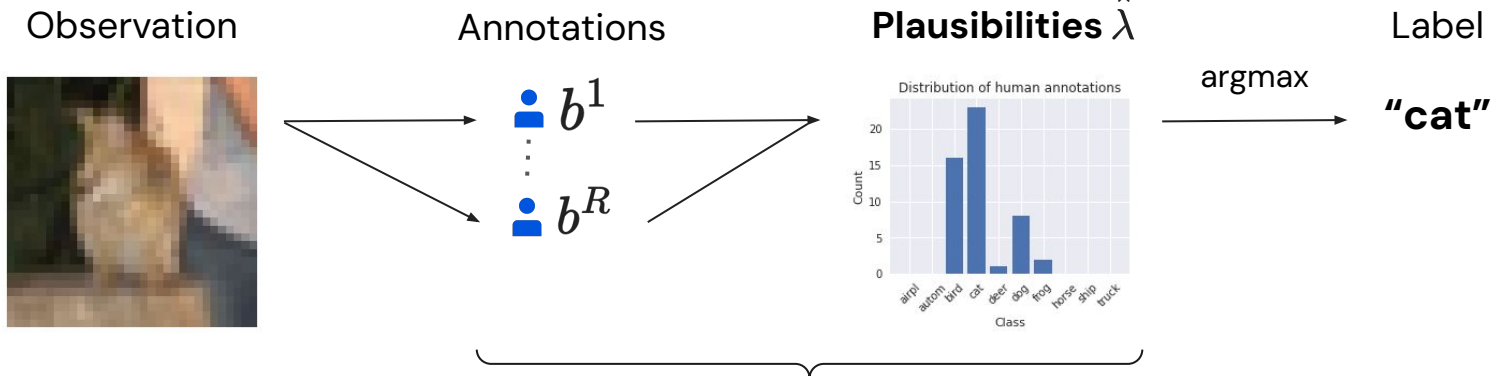
Deterministic aggregation ignores uncertainty



- Deterministic aggregation:
- Might evaluate against the wrong labels
 - Ignores large parts of the annotators
 - Does not quantify uncertainty on top of metrics



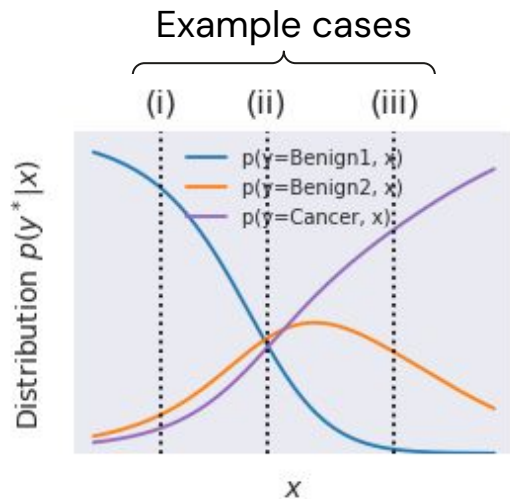
Introducing *plausibilities*



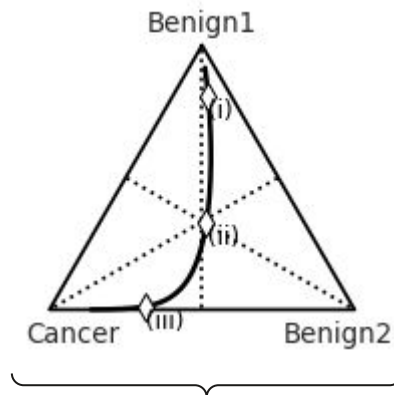
Deterministic aggregation approximates posterior $p(y^*|x)$ using a point estimate $\hat{\lambda}$

- “Plausibilities” = how *plausible* is a label given the annotations
- In this talk: categorical distributions over classes

Plausibilities on one-dimensional toy example

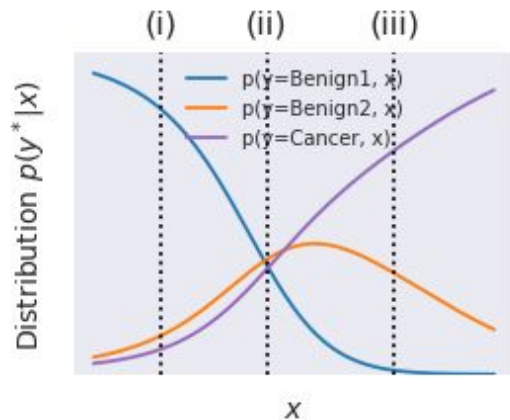


True posterior on
3-simplex

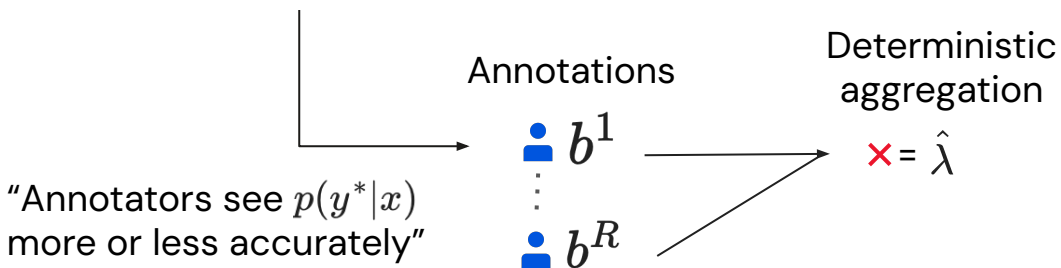
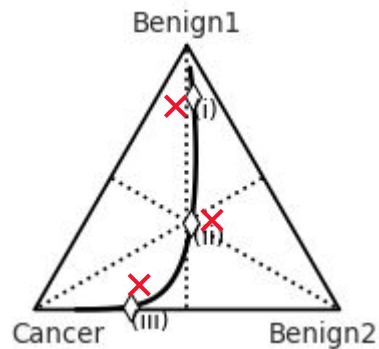


Inherent uncertainty =
location on simplex

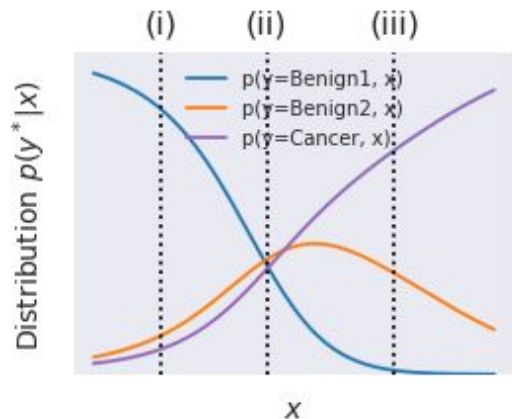
Point estimates from deterministic aggregation



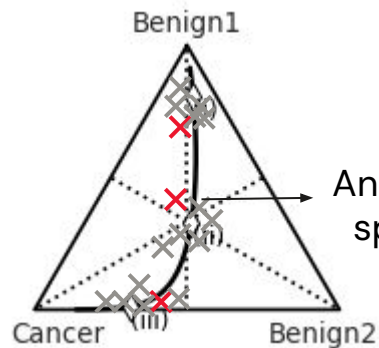
True posterior
 $\diamond = \lambda^* = p(y^*|x)$
on 3-simplex



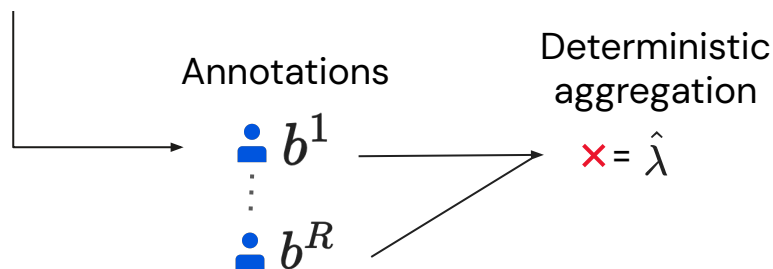
Variation in plausibilities through re-annotating



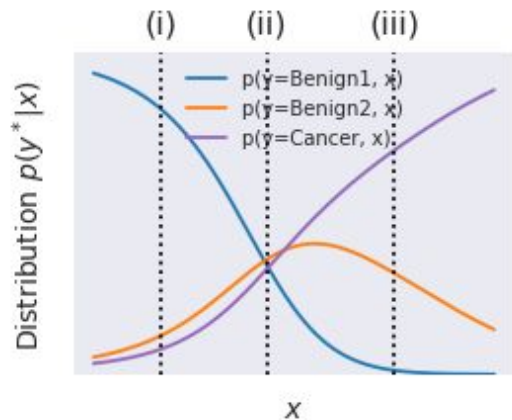
True posterior
 $\diamond = \lambda^* = p(y^*|x)$
on 3-simplex



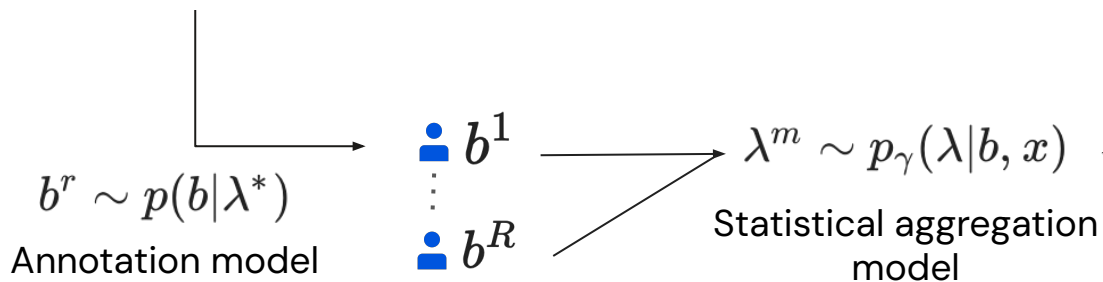
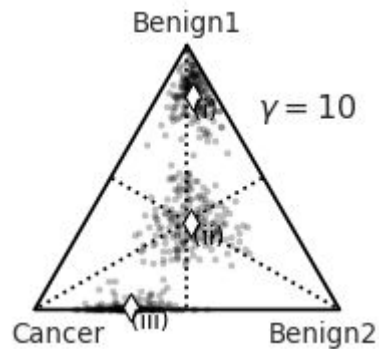
Different annotators yield
different plausibilities $\hat{\lambda}$



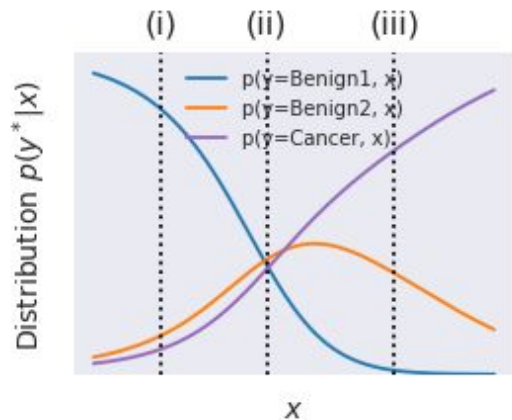
Statistical aggregation



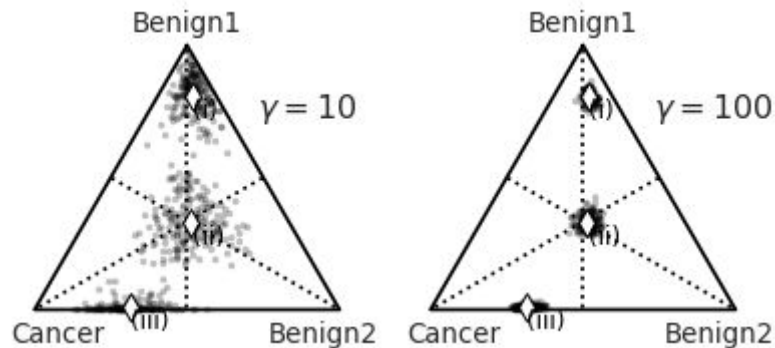
True posterior
 $\diamond = \lambda^* = p(y^* | x)$
on 3-simplex



Annotator *reliability* in statistical aggregation



True posterior
 $\diamond = \lambda^* = p(y^*|x)$
 on 3-simplex



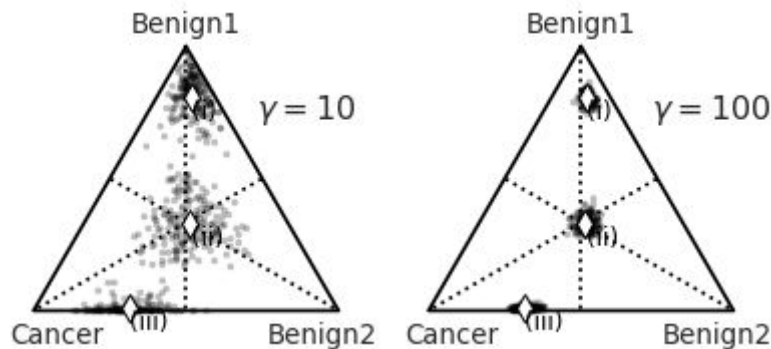
$b^r \sim p(b|\lambda^*)$

- b^1
- \vdots
- b^R

$\lambda^m \sim p_\gamma(\lambda|b, x)$

Reliability $\gamma =$ lower or higher prior trust in annotators

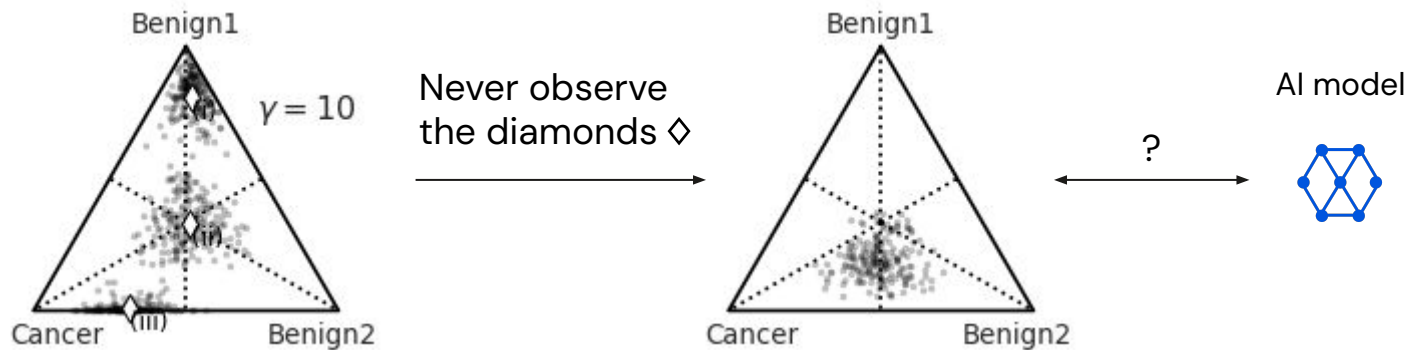
Conclusion: plausibilities on toy example



Ground truth uncertainty on the simplex:

- Location of plausibilities on simplex = inherent uncertainty
- Spread of plausibilities = annotation uncertainty

Conclusion: plausibilities on toy example



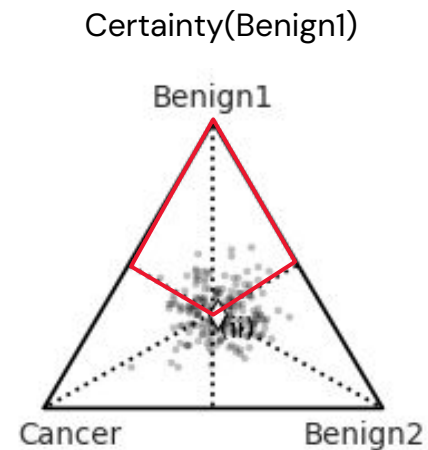
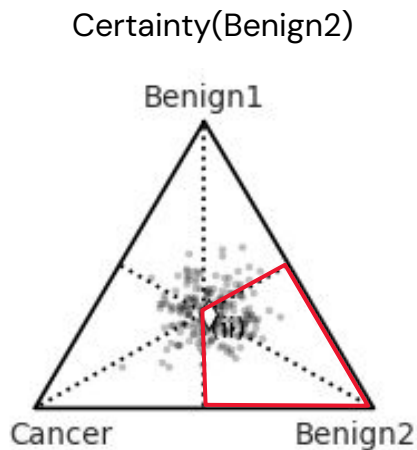
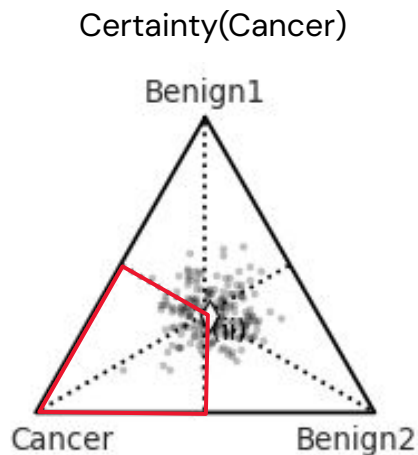
In practice, we only observe annotation disagreement:

- Cannot disentangle annotation and inherent uncertainty
 - Often impossible to deterministically derive a good approximation $\hat{\lambda}$
- How to *measure* this uncertainty and to *evaluate* AI models?

Measuring annotation uncertainty

- How *certain* is it that y is the top-1 label?

$$\text{Certainty}(y; b, x) = \mathbb{E}_{p(\lambda|b,x)} \left[\delta[y = \arg \max_j \lambda_j] \right]$$

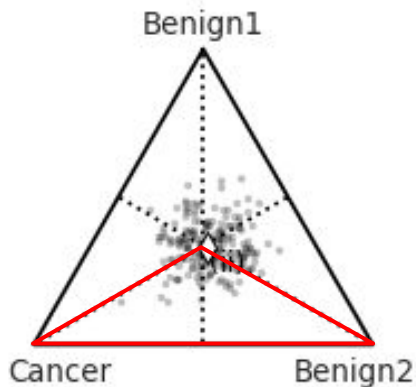


Measuring annotation uncertainty

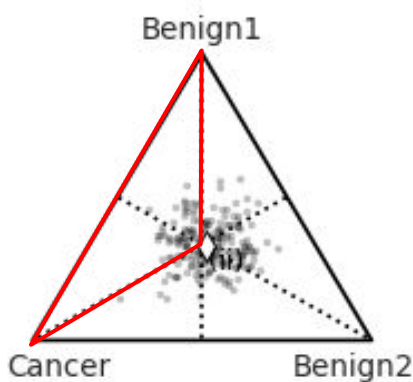
- Can also quantify certainty of label sets Y :

$$\text{Certainty}(Y; b, x) = \mathbb{E}_{p(\lambda|b,x)} [\delta[Y = \text{top}_k(\lambda)]]$$

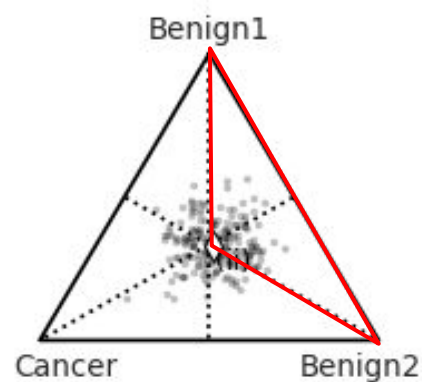
Certainty({Cancer, Benign2})



Certainty({Cancer, Benign2})



Certainty({Benign1, Benign2})



Measuring annotation uncertainty

- How *certain* is it that y is the top-1 label?

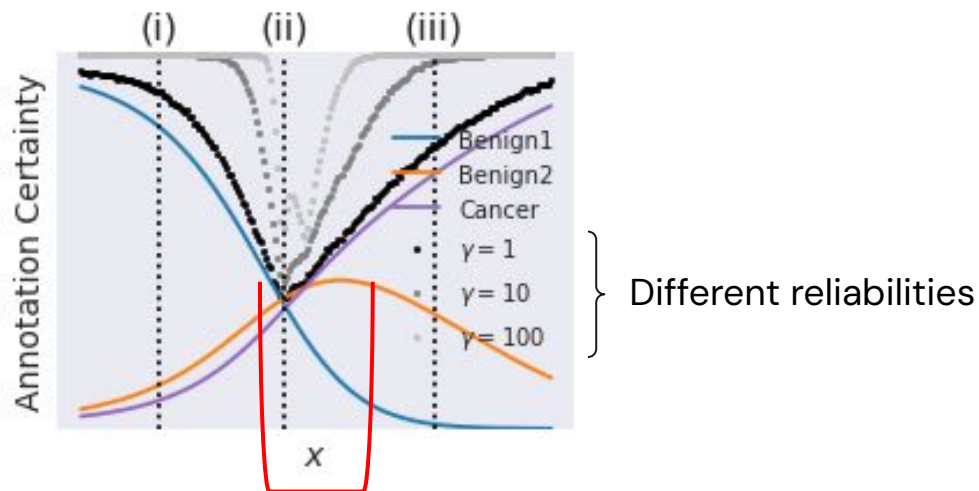
$$\text{Certainty}(y; b, x) = \mathbb{E}_{p(\lambda|b,x)} \left[\delta[y = \arg \max_j \lambda_j] \right]$$

- What is the highest certainty across labels?

$$\text{AnnotationCertainty}(b, x) = \max_y \text{Certainty}(y; b, x)$$

Measuring annotation uncertainty

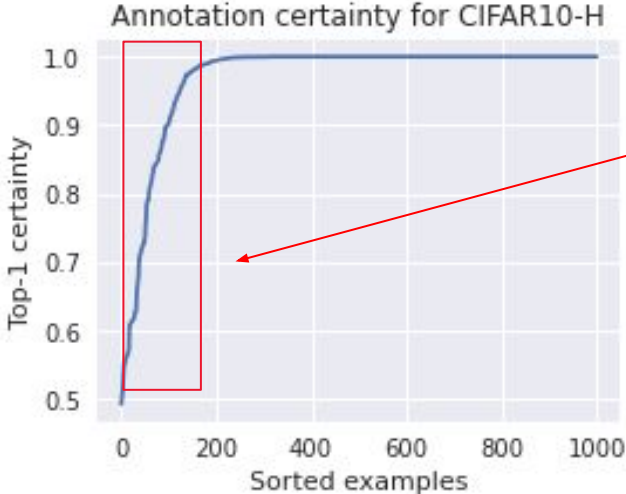
- Annotation certainty on toy example for different reliabilities γ :



Top-1 label unclear uncertain irrespective
of how much we trust our annotators

Measuring annotation uncertainty

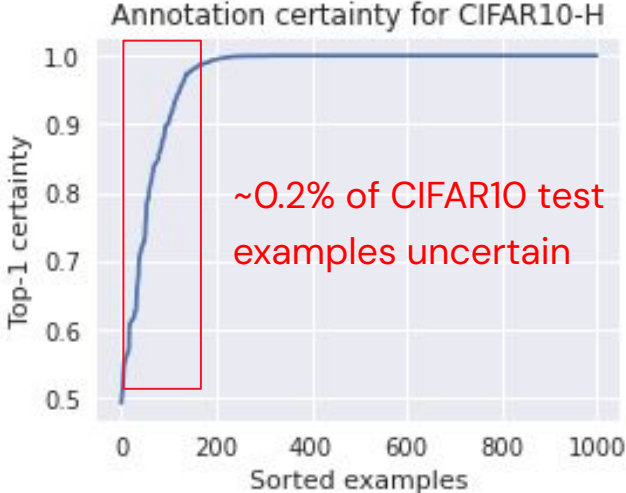
- Annotation certainty on CIFAR10 using annotations from CIFAR10-H:



- 178 examples with annotation certainty < 0.99
- This is ~0.2% of all CIFAR10 test examples

Measuring annotation uncertainty

- Annotation certainty on CIFAR10 using annotations from CIFAR10-H:



Papers with Code leaderboard:

μ 2Net	99.49	2022
ViT-L/16	99.42	2020
CaiT-M-36 U 224	99.4	2021
CvT-W24	99.39	2021
BiT-L	99.37	2019
ViT-B	99.3	2022

Improvements within 0.2%

Uncertainty-adjusted (top-k) accuracy

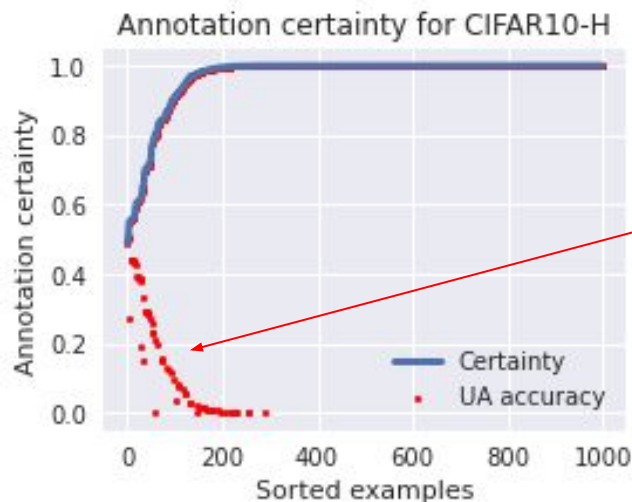
Given a model that yields a top-k prediction set $C_{\text{top-}k}(x)$:

$$\text{UA-Accuracy}_{\text{top-}k} = \mathbb{E}_{p(x)} \mathbb{E}_{p(\lambda|b,x)} \left[\delta[\arg \max_j \lambda_j \in C_{\text{top-}k}(x)] \right]$$

Uncertainty-adjusted (top-k) accuracy

Given a model that yields a top-k prediction set $C_{\text{top-}k}(x)$:

$$\text{UA-Accuracy}_{\text{top-}k} = \mathbb{E}_{p(x)} \mathbb{E}_{p(\lambda|b,x)} \left[\delta[\arg \max_j \lambda_j \in C_{\text{top-}k}(x)] \right]$$



- $C_{\text{top-}k}(x)$ = original CIFAR10 labels ($k = 1$)
- Annotations taken from CIFAR10-H
- Even CIFAR10 labels perform poorly on uncertain examples!

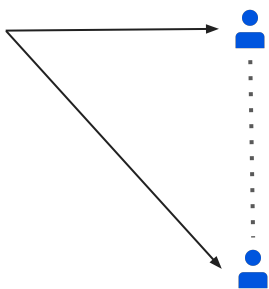
Case study: dermatology

Task: predict dermatological conditions from (multiple, consumer-grade) images.

Observation



Annotations



- b¹**: {*Pyogenic granuloma* (Low)} {*Hemangioma* (Med)} {*Melanoma* (High)}
- b²**: {*Angiokeratoma of skin* (Low)} {*Atypical Nevus* (Med)}
- b³**: {*Hemangioma* (Med)} {*Melanocytic Nevus* (Low), *Melanoma* (High),
O/E - ecchymoses present (Low)}
- b⁴**: {*Hemangioma* (Med), *Melanoma* (High), *Skin Tag* (Low)}
- b⁵**: {*Melanoma* (High)}
- b⁶**: {*Hemangioma* (Med)} {*Melanoma* (High)} {*Melanocytic Nevus* (Low)}

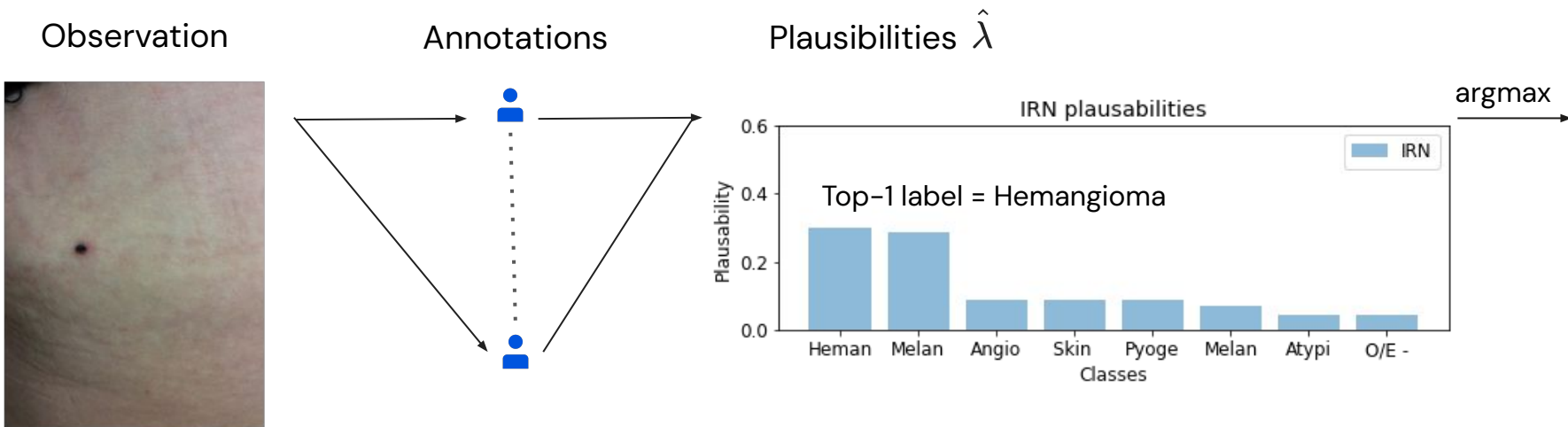
Conditions, Low/Med/High risk conditions

Partial rankings to model differential diagnoses

Case study: deterministic aggregation using IRN

Task: predict dermatological conditions from images.

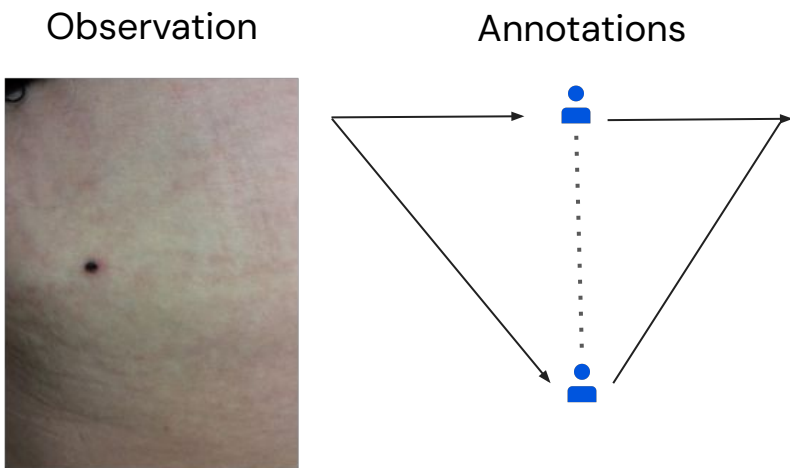
- Inverse rank normalization (IRN) to aggregate annotators' differential diagnoses.



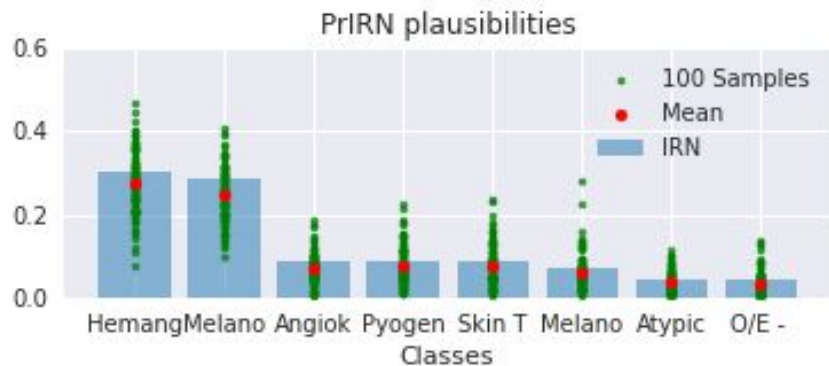
Case study: statistical aggregation using PrIRN

Task: predict dermatological conditions from images.

- Plackett-Luce or probabilistic IRN (PrIRN) to model $p(\lambda|b)$



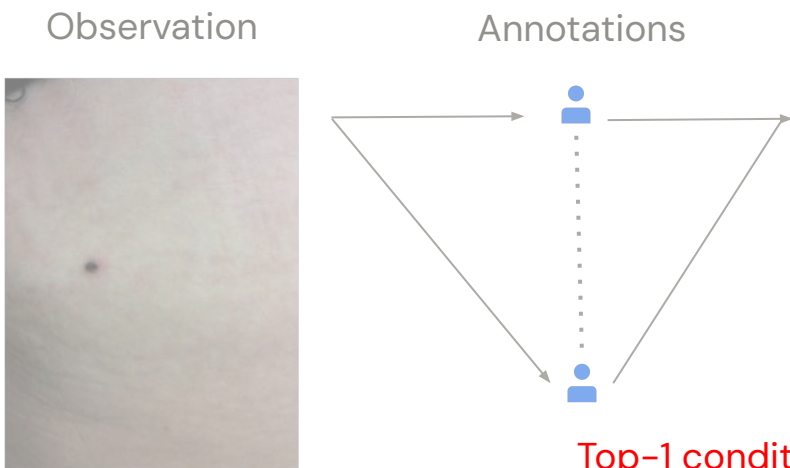
Plausibilities $\lambda^m \sim p_\gamma(\lambda|b)$



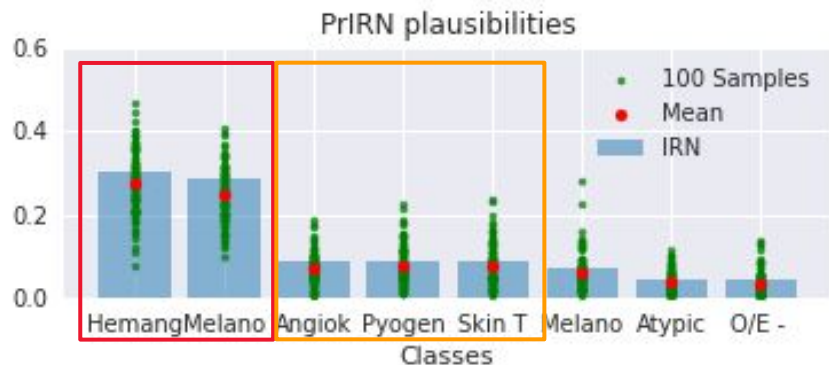
Case study: statistical aggregation using PrIRN

Task: predict dermatological conditions from images.

- Plackett-Luce or probabilistic IRN to model $p(\lambda|b)$



Plausibilities $\lambda^m \sim p_\gamma(\lambda|b)$

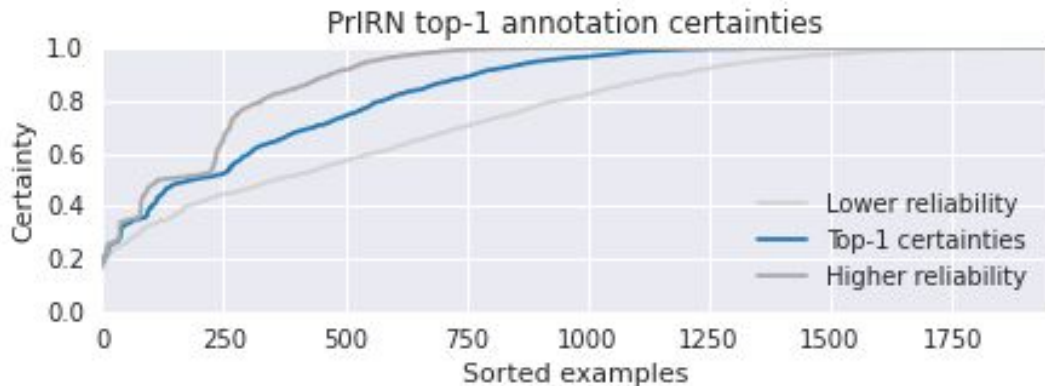


Top-1 condition changes easily
= low annotation certainty

3rd, 4th, 5h conditions also
change easily

High annotation uncertainty

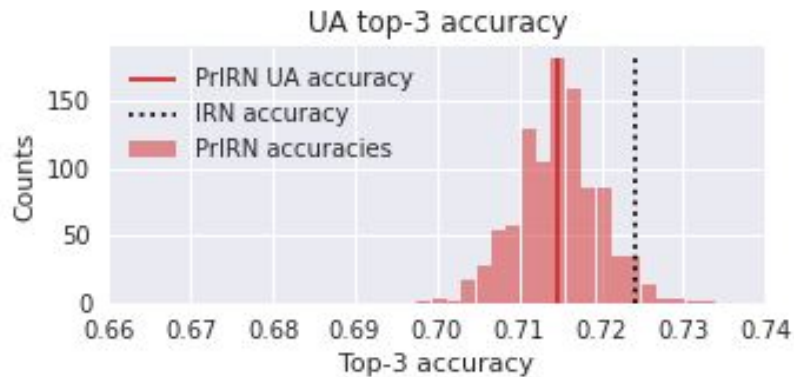
- Significant portions of cases with high annotation uncertainty:



- In discussions with dermatologists often attributed to inherent uncertainty

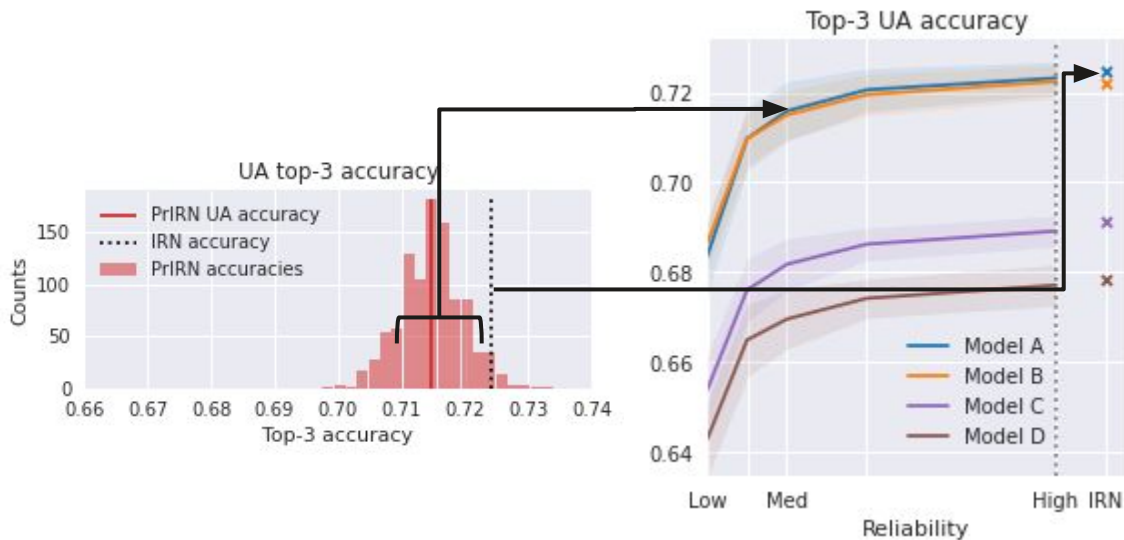
Uncertainty-adjusted top-3 accuracy

- Across cases / per plausibility:



→ Significant variation in top-3 accuracy

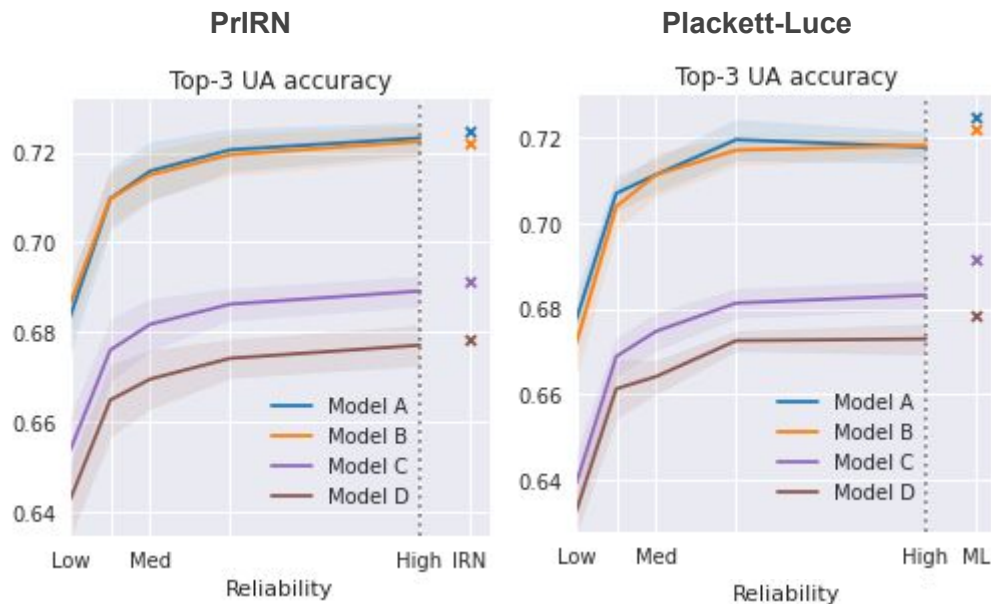
Evaluation across annotator reliabilities



- UA accuracy varies significantly by reliability
- IRN implicitly evaluates infinite annotator reliability
- Large spread/uncertainty in accuracies (shaded)

Alternative statistical aggregation methods

- Alternative statistical aggregation models exhibit different results:



→ Aggregation is a mode choice usually not made explicit!

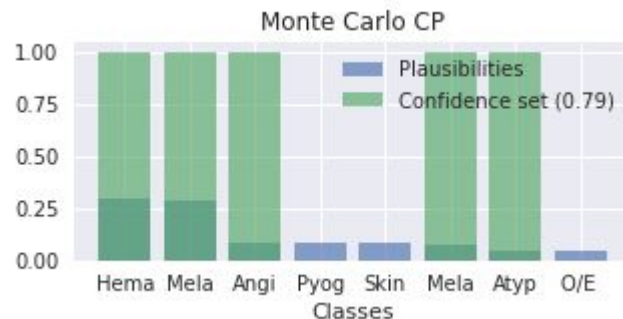
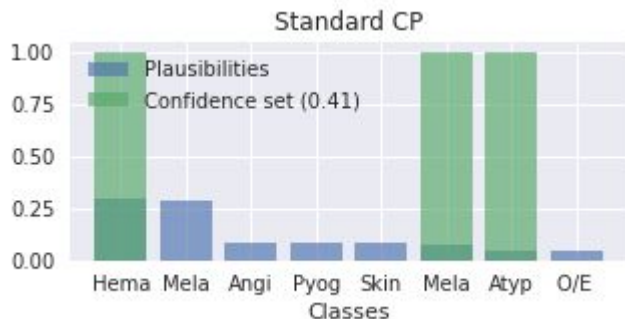
Bonus: calibration with uncertain ground truth

Calibration usually based on ground truth labels on a calibration/validation set:

- Conformal prediction uses ground truth labels to calibrate a softmax threshold τ
- Threshold used to predict confidence sets of classes at test time instead of the top-k:

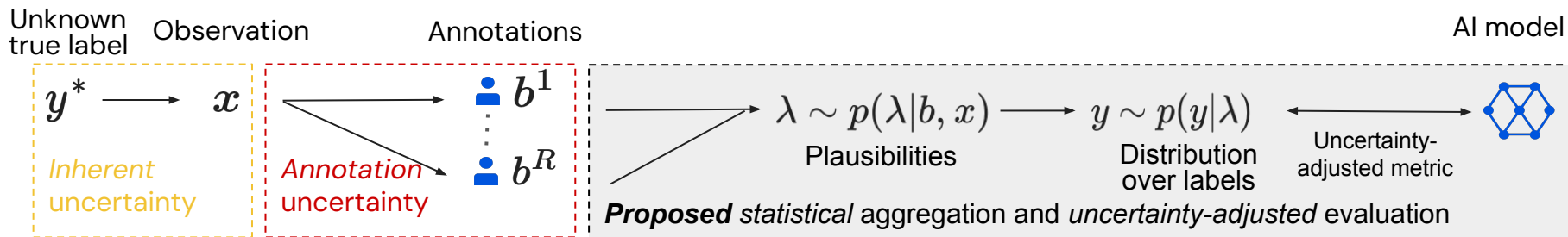
$$C_{\text{top-}k}(x) \longrightarrow C_{\text{CP}}(x) := \{k \in [K] : k\text{-th softmax} \geq \tau\}$$

- We propose *Monte Carlo* conformal prediction to address this issue, improving uncertainty-adjusted performance

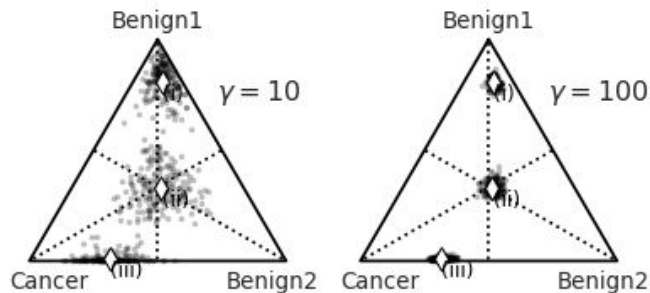


Conclusion

Proposed a statistical framework for dealing with ground truth uncertainty:



- Ground truth uncertainty = inherent + annotation uncertainty (location + spread of plausibilities)
- *Annotation certainty* explicitly measures annotation uncertainty
- Uncertainty-adjusted metrics to evaluate and evaluate models



More details: arxiv.org/abs/2307.02191 | dstutz@google.com