Google DeepMind

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

David Stutz

a collaboration between Google DeepMind and Health
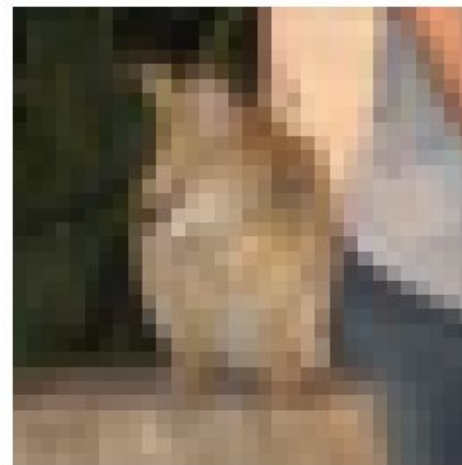
# Outline

Contents:
- ❏ Uncertainty from annotator disagreement
- ❏ Statistical framework
- ❏ Measuring uncertainty
- ❏ Evaluating AI models
- ❏ Case study in dermatology:
  - ❏ Results
  - ❏ Bonus: calibration
- ❏ Conclusion and outlook

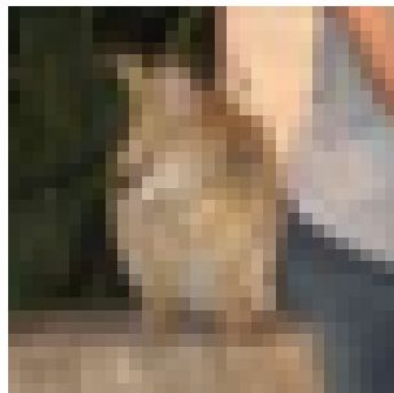Promise: you will start to question any "ground truth" labels you come across!



``Bird", "cat", or "frog"?



"Hemangioma" or "Melanoma"?
Benign or cancer?

# Standard evaluation of supervised models

Observation

AI prediction

Correct/good prediction?

**"bird"**

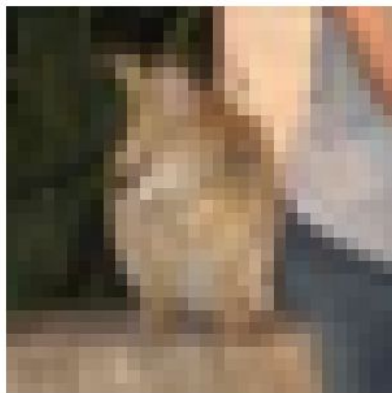# Standard evaluation of supervised models

*Unknown* true label

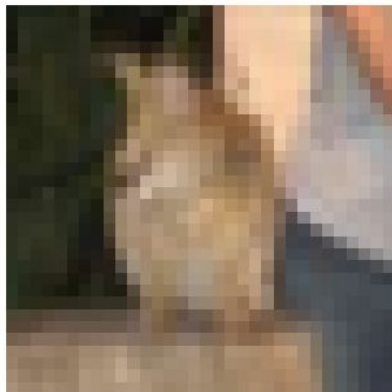Observation

AI prediction

**?**



Correct/good prediction?

**"bird"**

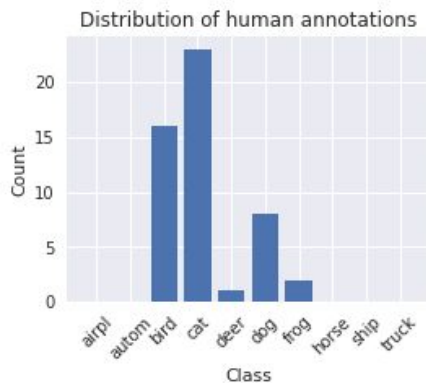# Standard evaluation of supervised models

Unknown
true label

**?**

Observation

Annotations

AI prediction

?

**"bird"**
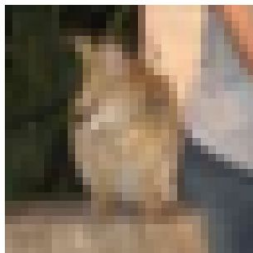
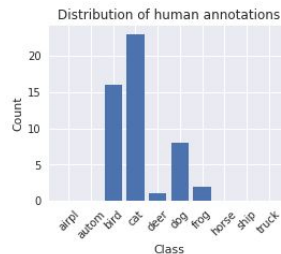Annotators disagree!

# Standard evaluation of supervised models

Unknown
true label

**?**

Observation

Annotations

Majority vote

AI prediction



Distribution of human annotations

Incorrect!

**"cat"** ⟷ **"bird"**

Ignores
disagreement!

# 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

Y. Liu et al. A deep learning system for differential diagnosis of skin diseases. Nature Medicine, 2020.

# Standard evaluation of supervised models

Unknown true label

Observation

Annotations

AI prediction

**?**

$b^1$: {*Pyogenic granuloma* (Low)} {***Hemangioma*** (Med)} {*Melanoma* (High)}

$b^2$ {*Angiokeratoma of skin* (Low)} {*Atypical Nevus* (Med)}

$b^3$: {***Hemangioma*** (Med)} {*Melanocytic Nevus* (Low), *Melanoma* (High), *O/E – ecchymoses present* (Low)}

$b^4$: {***Hemangioma*** (Med), *Melanoma* (High), *Skin Tag* (Low)}

$b^5$: {*Melanoma* (High)}

$b^6$: {***Hemangioma*** (Med)} {*Melanoma* (High)} {*Melanocytic Nevus* (Low)}

**?**

**"Hemangioma"**

Majority voting is non-trivial

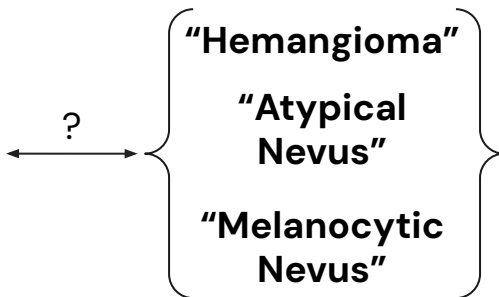# Standard evaluation of supervised models

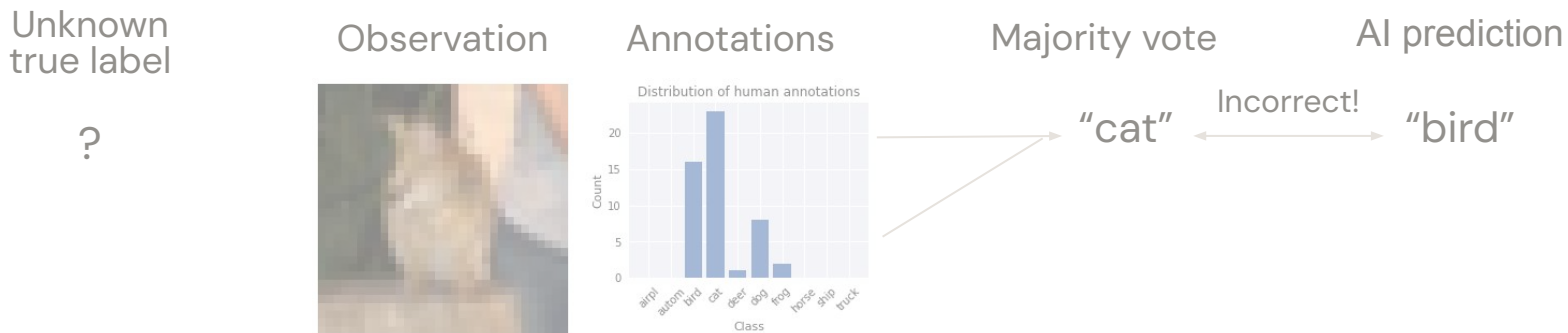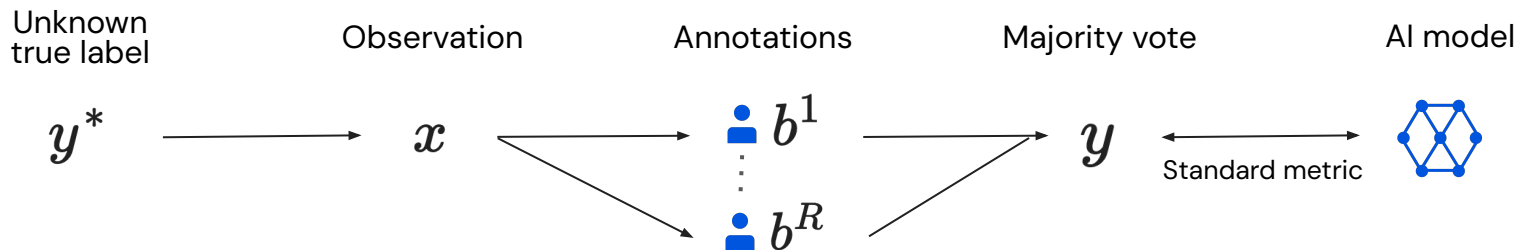Unknown true label

**?**

Observation

Annotations

$b^1$: {*Pyogenic granuloma* (Low)} {*Hemangioma* (Med)} {*Melanoma* (High)}

$b^2$ {*Angiokeratoma of skin* (Low)} {*Atypical Nevus* (Med)}

$b^3$: {*Hemangioma* (Med)} {*Melanocytic Nevus* (Low), *Melanoma* (High), *O/E – ecchymoses present* (Low)}

$b^4$: {*Hemangioma* (Med), *Melanoma* (High), *Skin Tag* (Low)}

$b^5$: {*Melanoma* (High)}

$b^6$: {*Hemangioma* (Med)} {*Melanoma* (High)} {*Melanocytic Nevus* (Low)}
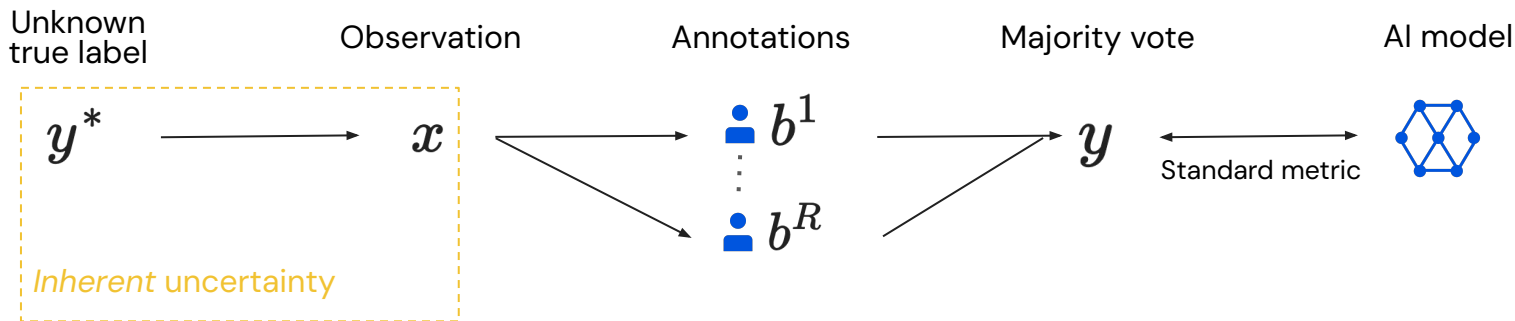
?

AI prediction *set*

**"Hemangioma"**

**"Atypical Nevus"**

**"Melanocytic Nevus"**

# Standard evaluation of supervised models

# *Inherent* uncertainty



| Unknown true label | Observation | Annotations | Majority vote | AI model |

$y^*$ → $x$ → $b^1$ ... $b^R$ → $y$ ← Standard metric → AI model
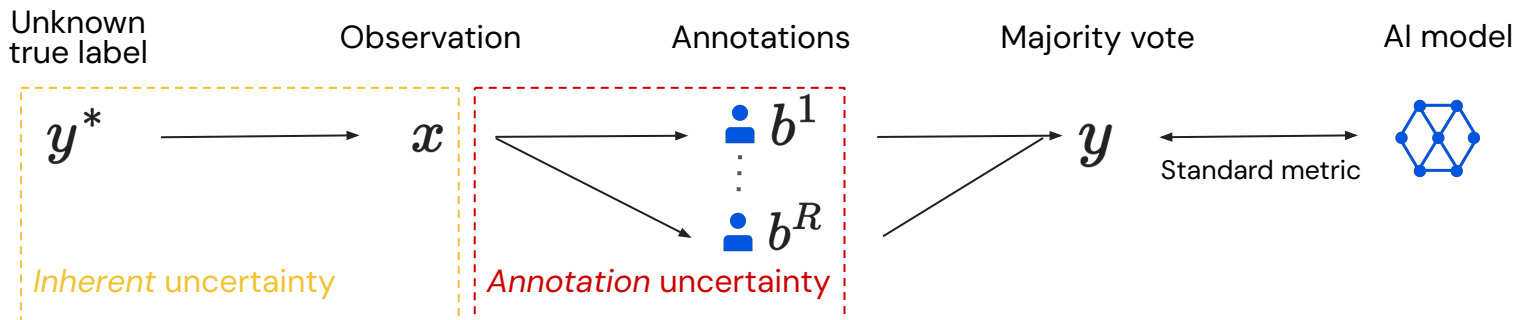
*Inherent* uncertainty

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

Unknown true label      Observation      Annotations      Majority vote      AI model

$$y^*$$

$$x$$

$$b^1$$

$$b^R$$

$$y$$

Standard metric
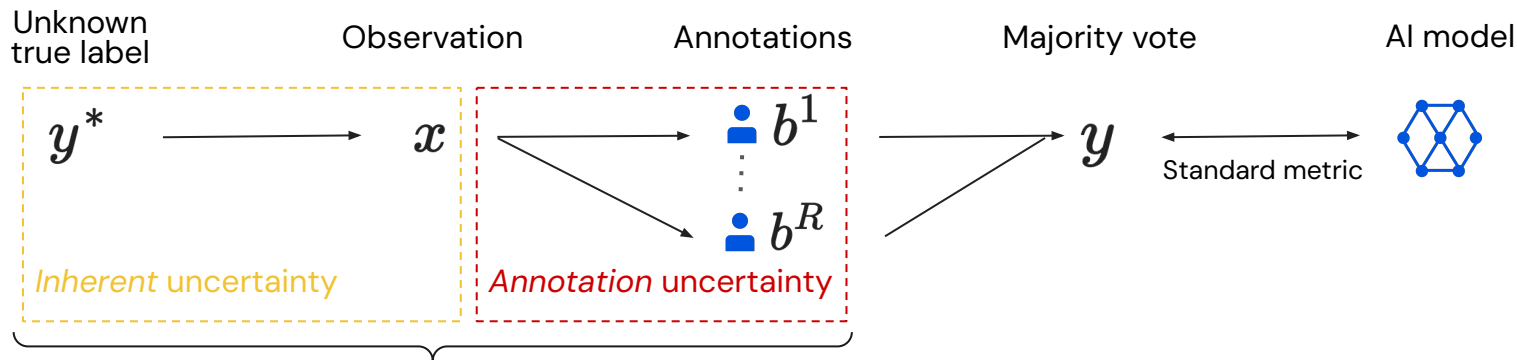
*Inherent* uncertainty

*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

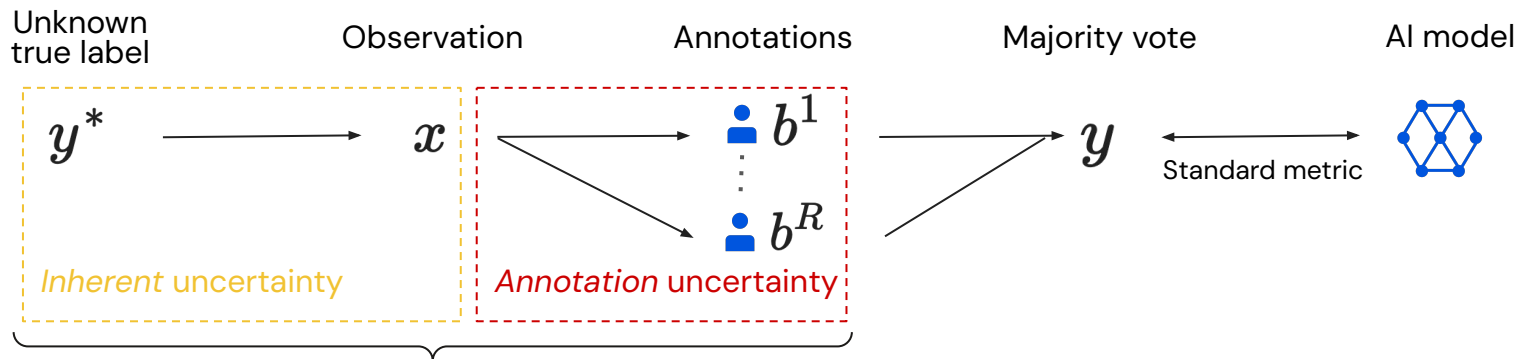TL;DR: annotation is difficult and we have to trust experts.

# Ground truth uncertainty

Unknown true label    Observation    Annotations    Majority vote    AI model

$$y^*$$

$$x$$

$$b^1$$

$$b^R$$

$$y$$

Standard metric

*Inherent* uncertainty

*Annotation* uncertainty

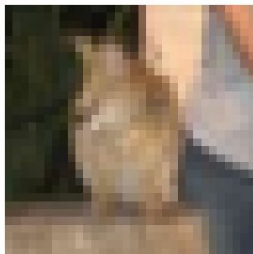Ground truth uncertainty = inherent + annotation uncertainty

- We observe both through annotator **disagreement**
- Often impossible to disentangle inherent and annotation uncertainty

# Ground truth uncertainty

Unknown true label      Observation      Annotations      Majority vote      AI model

$y^*$     →     $x$      $b^1$ ... $b^R$      $y$

Standard metric

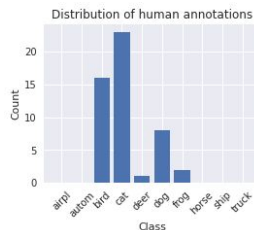*Inherent* uncertainty      *Annotation* uncertainty

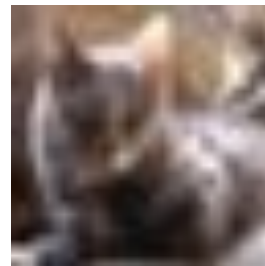Ground truth uncertainty = inherent + annotation uncertainty

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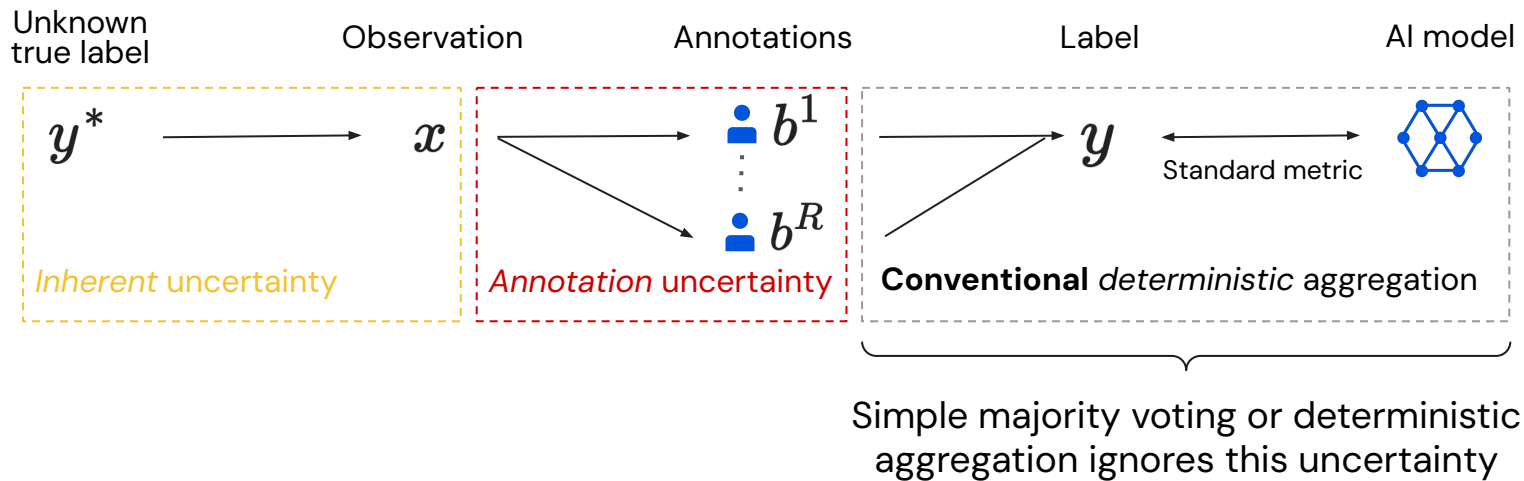- Usually we cannot disentangle between inherent and annotation uncertainty
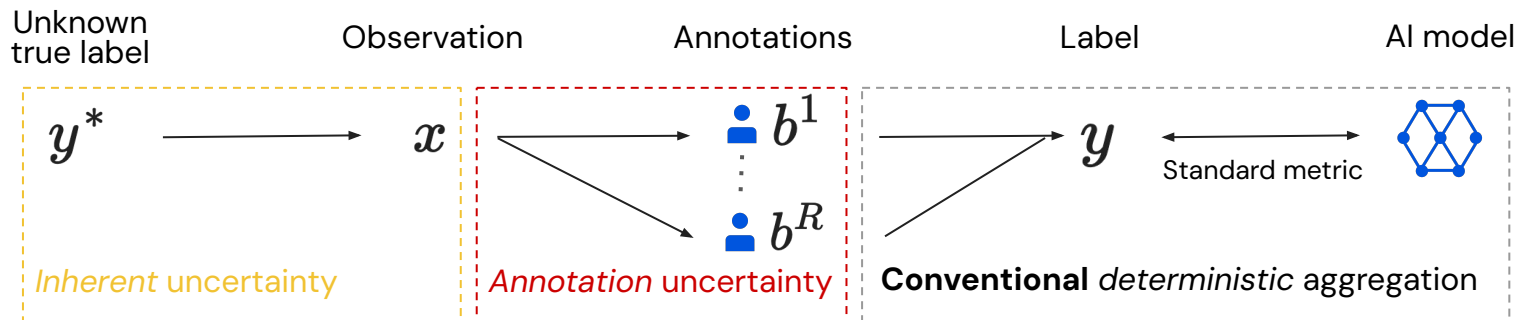


Inherent uncertainty      Distribution of human annotations      Annotation uncertainty

# Deterministic aggregation ignores uncertainty

Unknown true label → Observation → Annotations → Label → AI model

$y^*$ → $x$ → $b^1$ ... $b^R$ → $y$ ← Standard metric → (AI model)

*Inherent* uncertainty

*Annotation* uncertainty

**Conventional** *deterministic* aggregation

Simple majority voting or deterministic aggregation ignores this uncertainty

# Deterministic aggregation ignores uncertainty

| Unknown true label | Observation | Annotations | Label | AI model |

$y^*$ → $x$

$b^1$

$b^R$

$y$

Standard metric

*Inherent* uncertainty

*Annotation* uncertainty

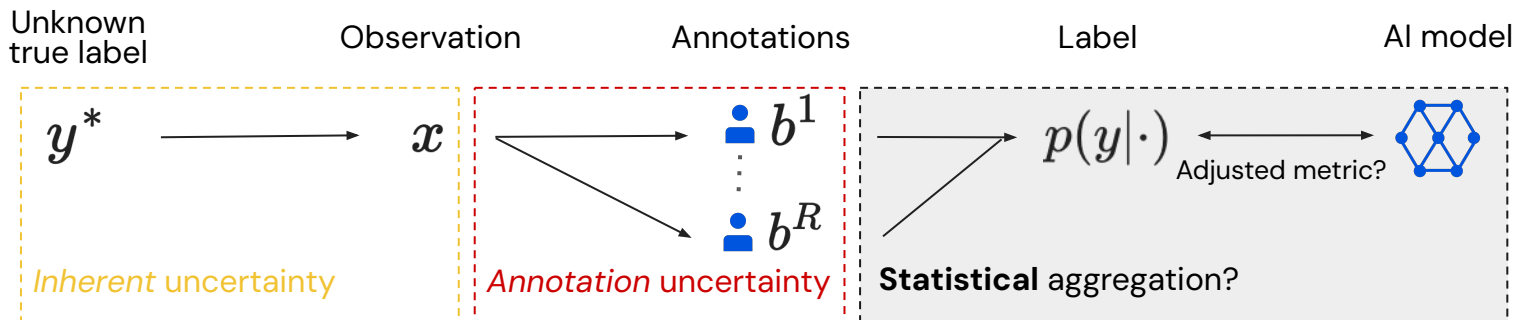**Conventional** *deterministic* aggregation

Deterministic aggregation ignores ground truth uncertainty:

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

# Deterministic aggregation ignores uncertainty

Unknown true label | Observation | Annotations | Label | AI model

$y^*$ → $x$ → $b^1$ ... $b^R$ → $p(y|\cdot)$ ↔ (AI model)

Adjusted metric?

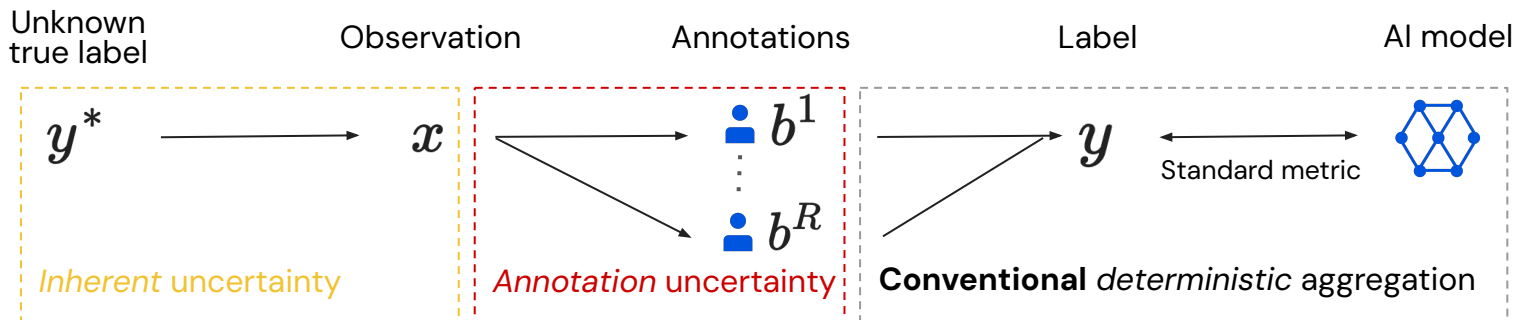*Inherent* uncertainty | *Annotation* uncertainty | **Statistical** aggregation?

Can we use a statistical aggregation model to account for uncertainty?

- Statistical aggregation of annotations
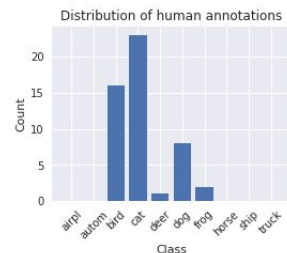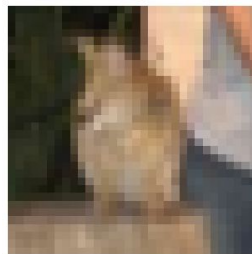- Adjust evaluation metrics by uncertainty

# Deterministic aggregation ignores uncertainty

Unknown true label | Observation | Annotations | Label | AI model



$y^*$ → $x$ → $b^1$ ⋮ $b^R$ → $y$ ← Standard metric →

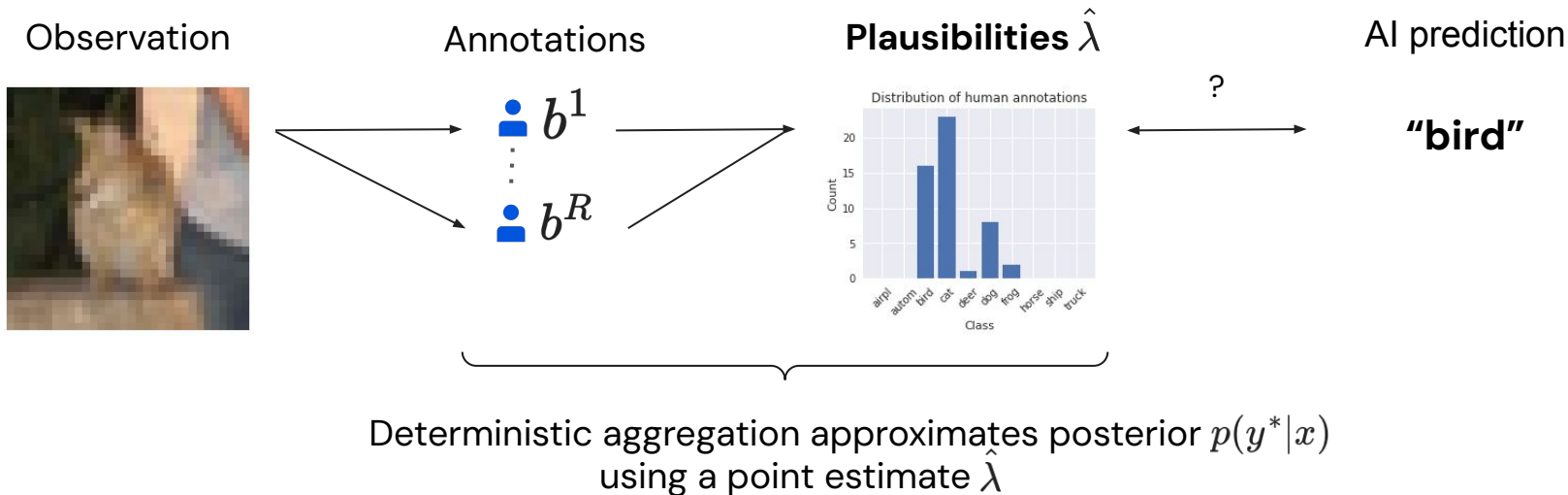*Inherent* uncertainty | *Annotation* uncertainty | **Conventional** *deterministic* aggregation

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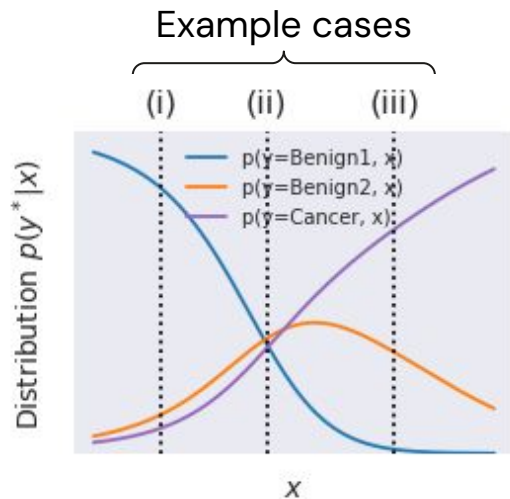➔ Does not quantify uncertainty on top of metrics



Distribution of human annotations

# Introducing *plausibilities*

| Observation | Annotations | Plausibilities $\hat{\lambda}$ | AI prediction |
|---|---|---|---|



$b^1$

$\vdots$

$b^R$

?

**"bird"**

Distribution of human annotations

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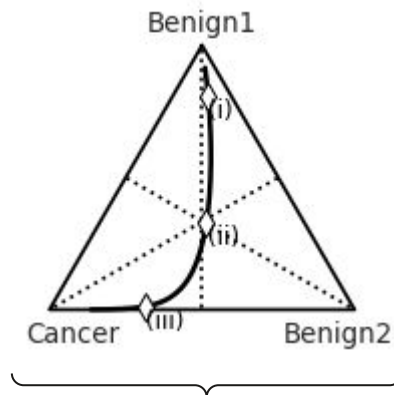using a point estimate $\hat{\lambda}$

- "Plausibilities" = how *plausible* is a label given the annotations
- In this talk: categorical distributions over classes
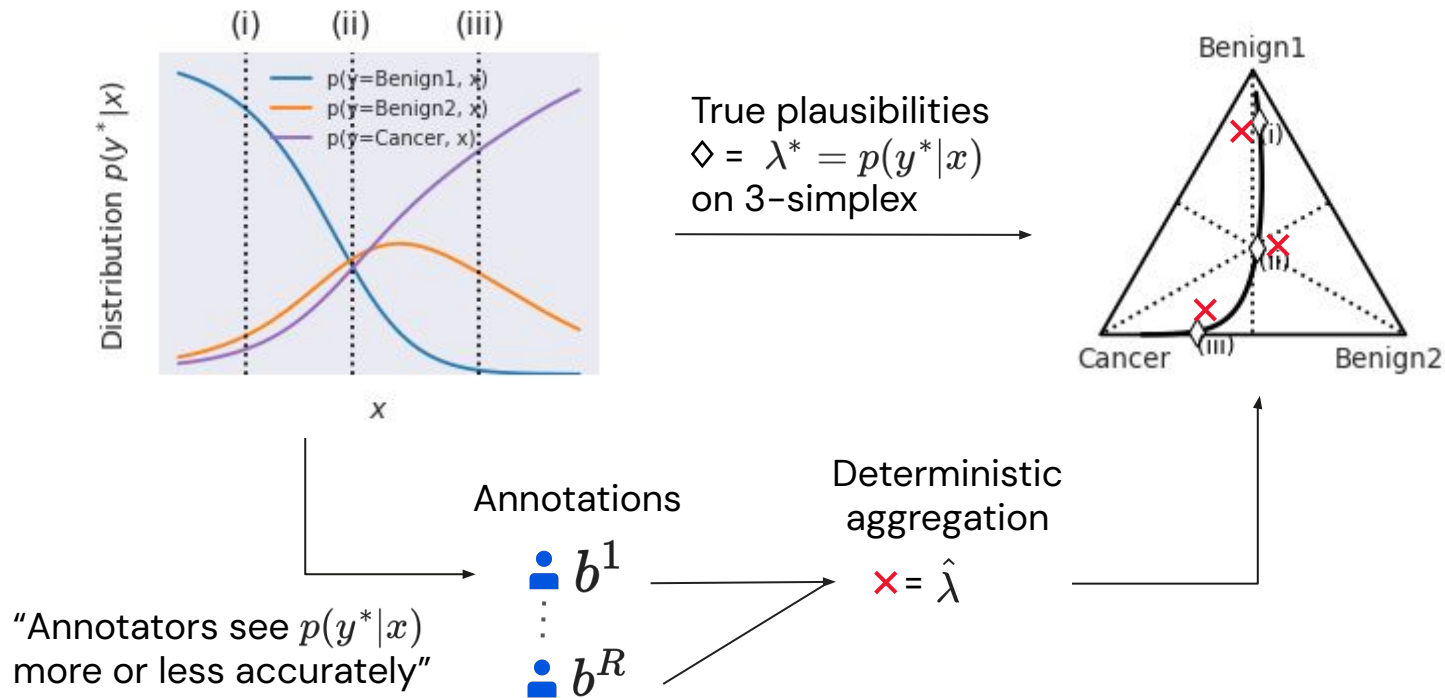
# Plausibilities on one-dimensional toy example

Example cases



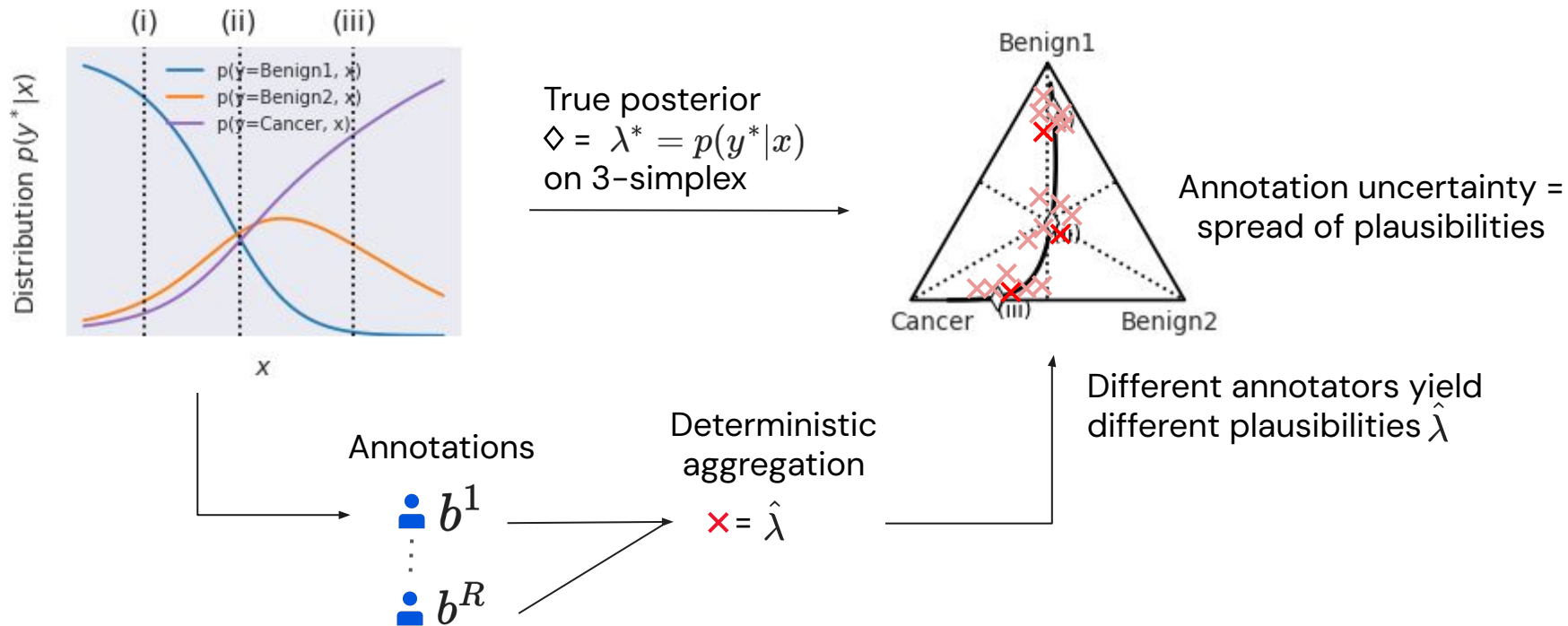True plausibilities
$\Diamond = \lambda^* = p(y^*|x)$
on 3-simplex
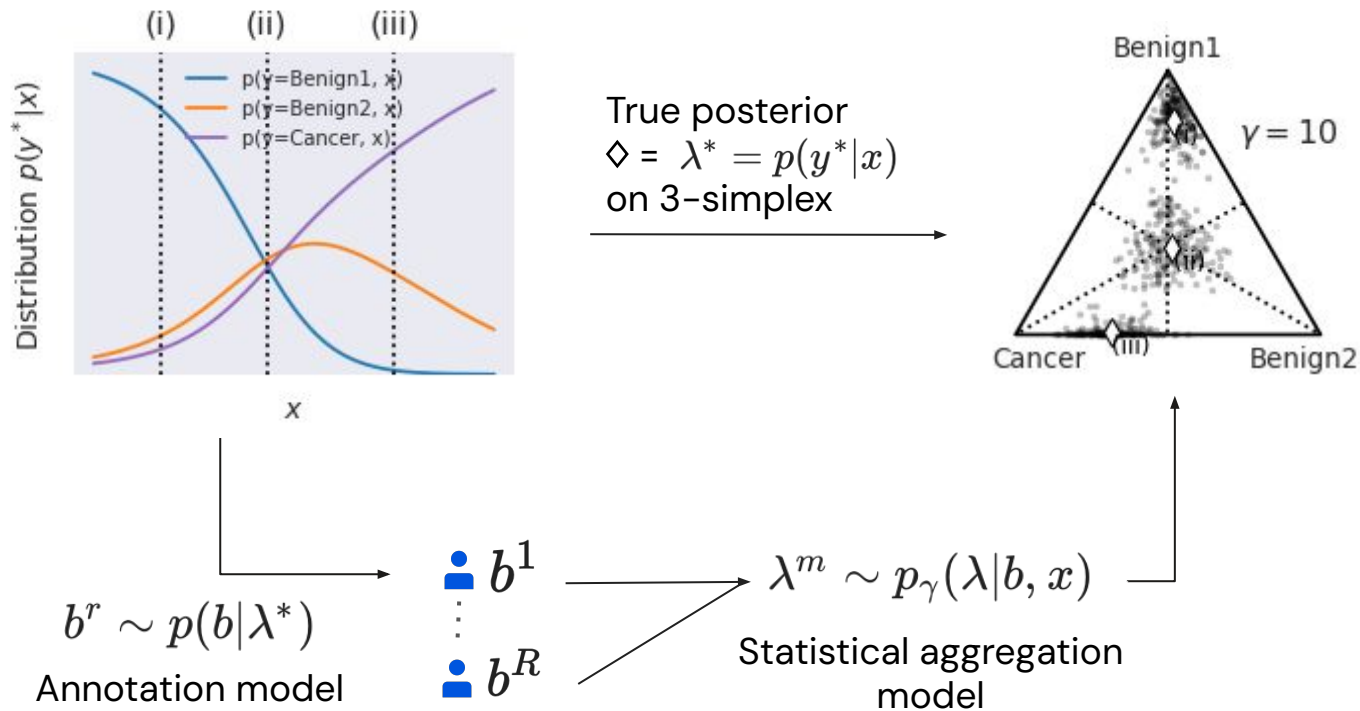


Inherent uncertainty =
location on simplex

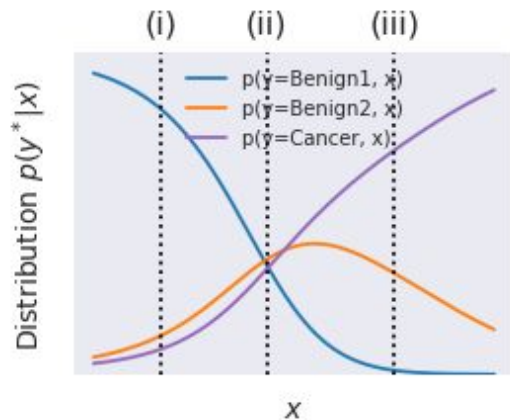# Point estimates from deterministic aggregation

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

Annotations

$\bullet \ b^1$

$\vdots$

$\bullet \ b^R$

Deterministic
aggregation

$\times = \hat{\lambda}$

"Annotators see $p(y^*|x)$
more or less accurately"

# Variation in plausibilities through re-annotating

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

Annotation uncertainty = spread of plausibilities

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

Annotations
$b^1$
⋮
$b^R$

Deterministic aggregation
✗ = $\hat{\lambda}$

# *Statistical* aggregation



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

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

Annotation model

$b^1$
⋮
$b^R$

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

Statistical aggregation
model

# Annotator *reliability* in statistical aggregation

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

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

$b^1$
$\vdots$
$b^R$

$\lambda^m \sim p_\gamma(\lambda|b, x)$
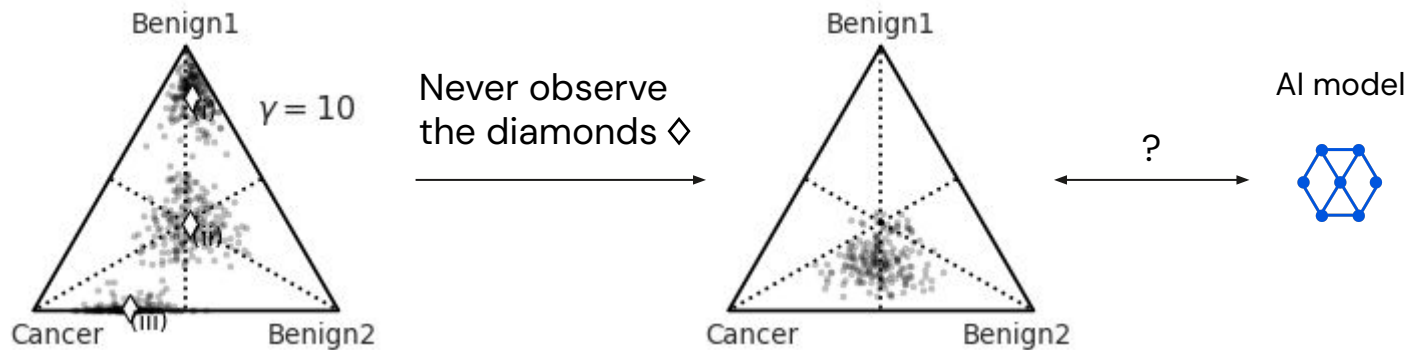
Reliability $\gamma$ = lower or higher
prior trust in annotators

# Conclusion: plausibilities on toy example

Ground truth uncertainty on the simplex:

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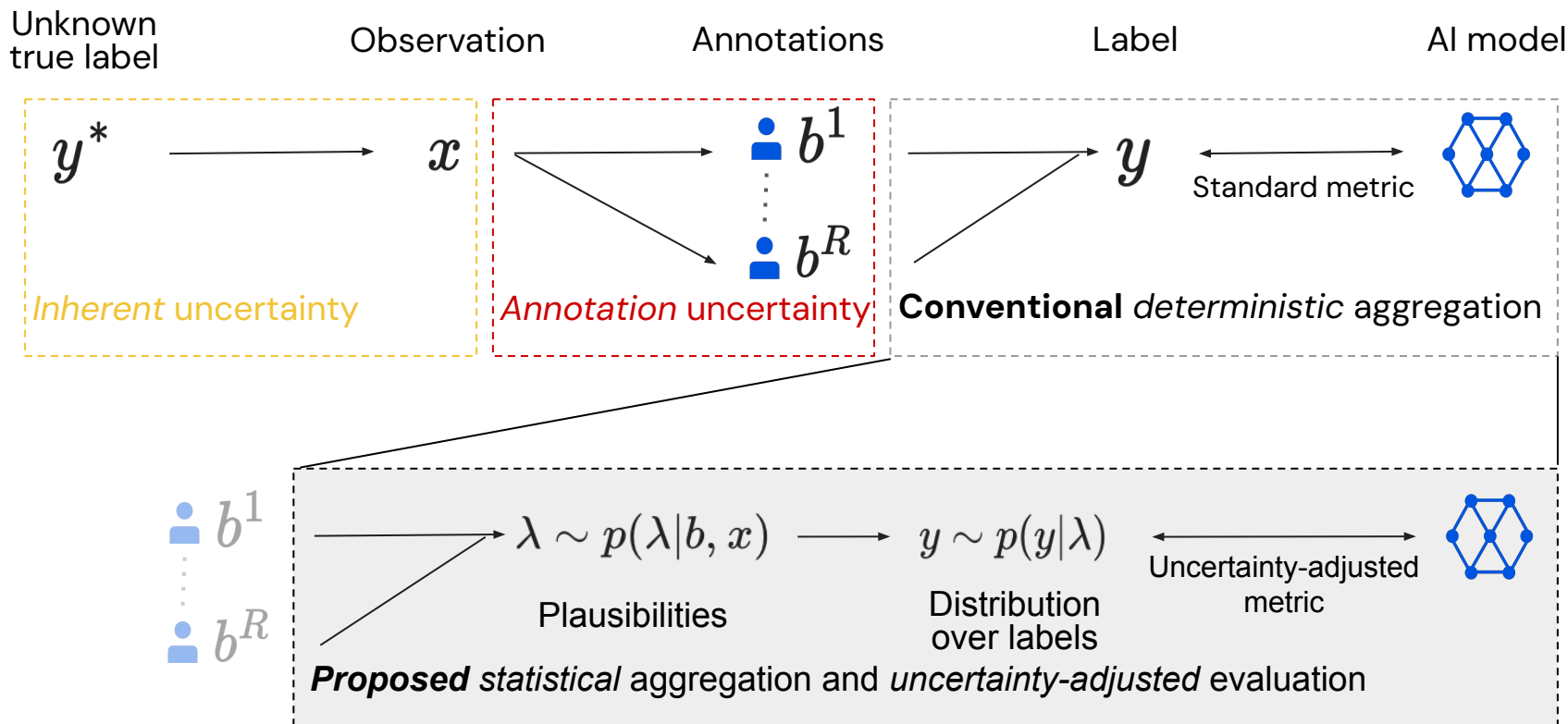- Spread of plausibilities = annotation uncertainty

# Conclusion: plausibilities on toy example

Statistically modeling aggregation:

- Allows to disentangle inherent and annotation uncertainty to some extent (subject to modeling assumptions, depending on reliability)
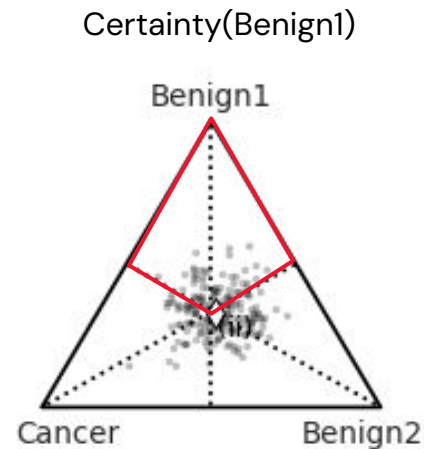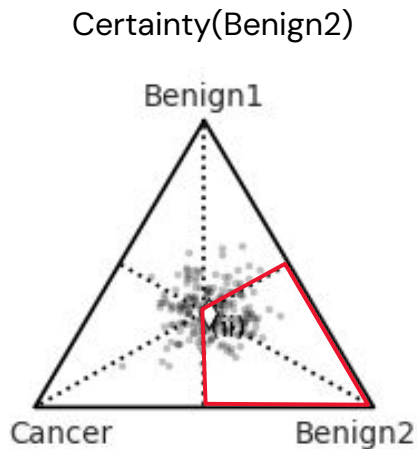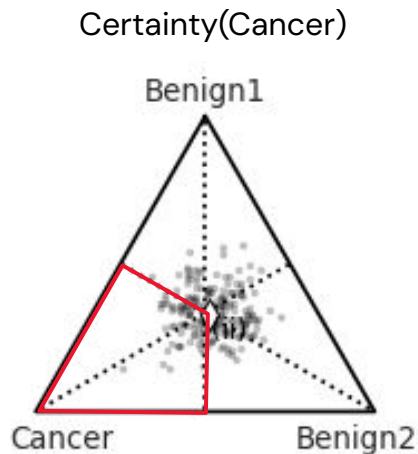- Avoids expensive re-annotation to get uncertainty estimates
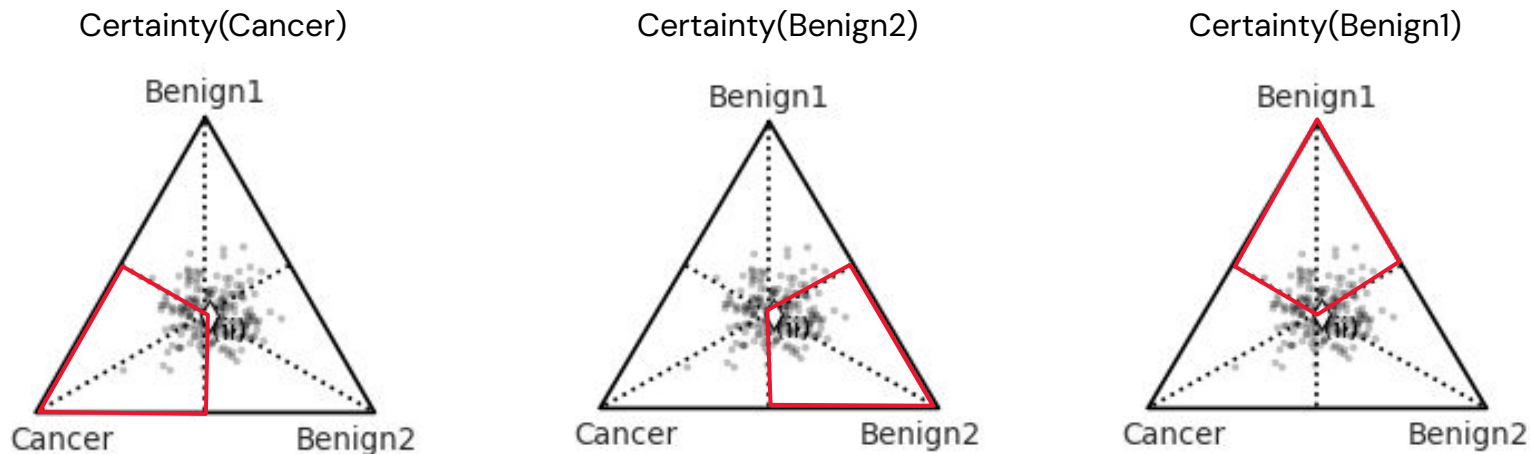
# Summary: proposed statistical framework

# Measuring annotation uncertainty

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

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

Certainty(Cancer)

Certainty(Benign2)

Certainty(Benign1)

# 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?

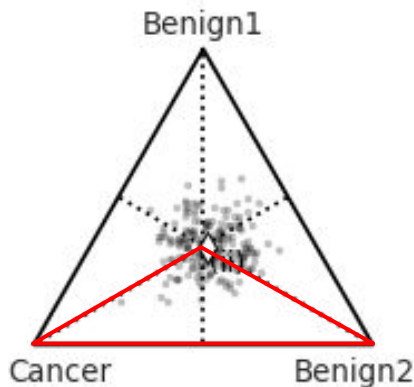$$\text{AnnotationCertainty}(b, x) = \max_y \text{Certainty}(y; b, x)$$

Certainty(Cancer)     Certainty(Benign2)     Certainty(Benign1)
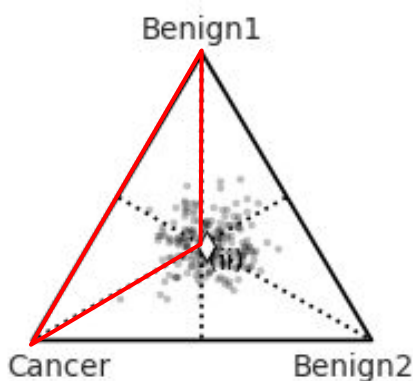
# Measuring annotation uncertainty

- Can also quantify certainty of label sets $Y$:

$$\text{Certainty}(Y; b, x) = \mathbb{E}_{p(\lambda|b,x)} \left[ \delta[Y = \text{top\_k}(\lambda)] \right]$$
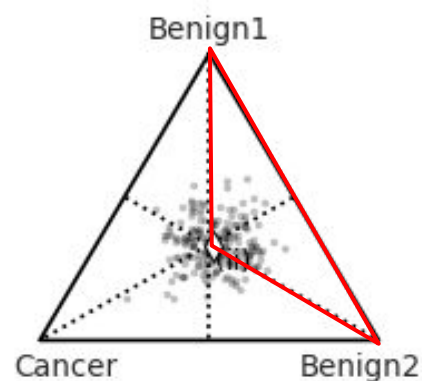
Certainty({Cancer, Benign2})

Certainty({Cancer, Benign2})

Certainty({Benign1, Benign2})

# Measuring annotation uncertainty

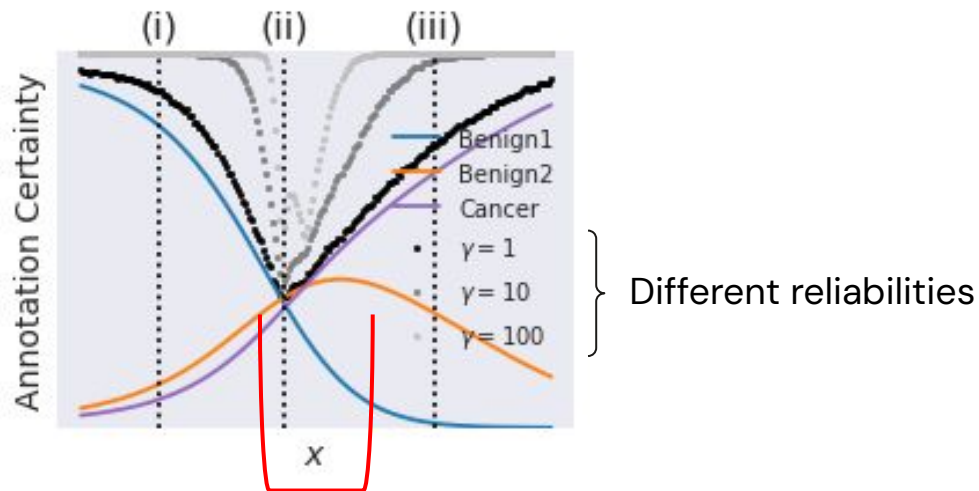- How *certain* is it that y is the top–1 label?

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

- What is the highest certainty across labels?

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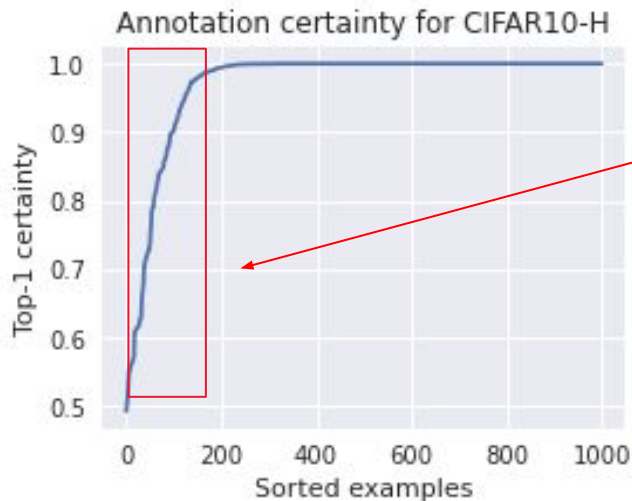
# Measuring annotation uncertainty

- Annotation certainty on toy example for different reliabilities $\gamma$ :



Different reliabilities

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

# Measuring annotation uncertainty

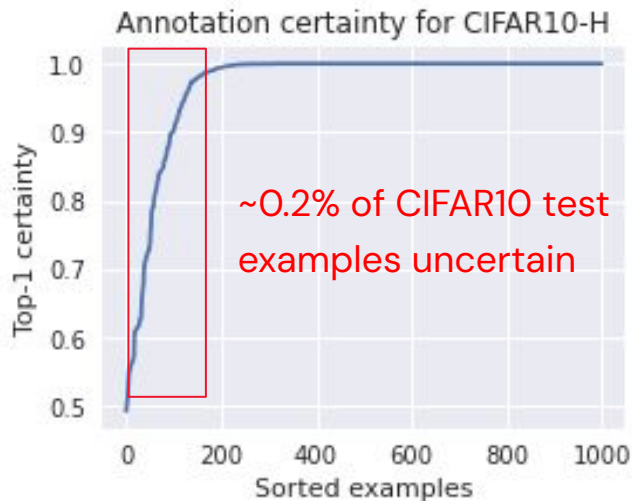- Annotation certainty on CIFAR10 using annotations from CIFAR10–H:



Annotation certainty for CIFAR10-H

- 178 examples with annotation certainty < 99%
- This is ~0.2% of all CIFAR10 test examples

# Measuring annotation uncertainty

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



Annotation certainty for CIFAR10-H

~0.2% of CIFAR10 test examples uncertain

Papers with Code leaderboard:

| µ2Net | 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%
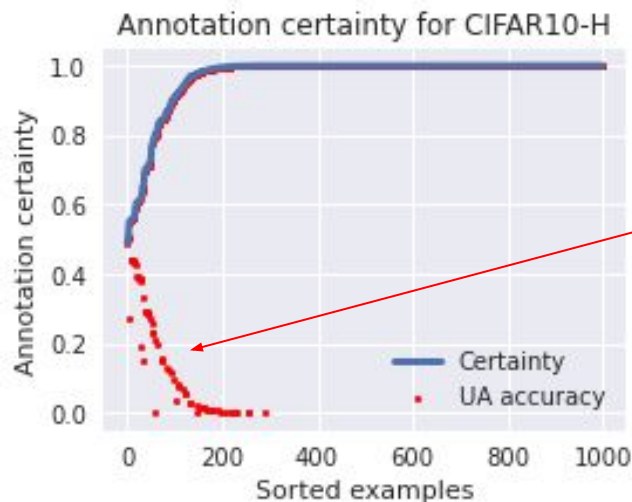
# *Uncertainty–adjusted* (top-k) accuracy

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 $C_{\text{top-}k}(x)$:

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



Annotation certainty for CIFAR10-H

- $C_{\text{top-}k}(x)$ = <u>original CIFAR10 labels</u> (k = 1)
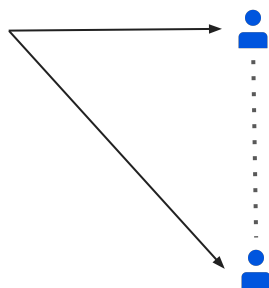- ➔ Even CIFAR10 labels perform poorly on uncertain examples!

# Case study: dermatology

Observation

Annotations



$b^1$: {*Pyogenic granuloma* (Low)} {*Hemangioma* (Med)} {*Melanoma* (High)}
$b^2$ {*Angiokeratoma of skin* (Low)} {*Atypical Nevus* (Med)}
$b^3$: {*Hemangioma* (Med)} {*Melanocytic Nevus* (Low), *Melanoma* (High), *O/E – ecchymoses present* (Low)}
$b^4$: {*Hemangioma* (Med), *Melanoma* (High), *Skin Tag* (Low)}
$b^5$: {*Melanoma* (High)}
$b^6$: {*Hemangioma* (Med)} {*Melanoma* (High)} {*Melanocytic Nevus* (Low)}

Partial rankings to model differential diagnoses

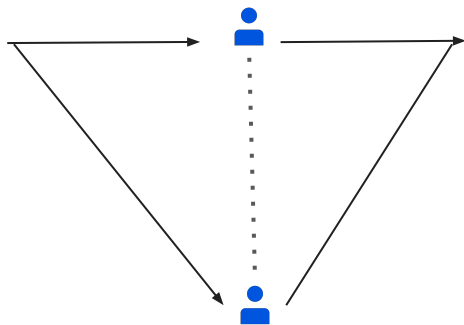# Case study: deterministic aggregation using IRN

Task: predict dermatological conditions from images.

- Inverse rank normalization (IRN) to aggregate annotators' differential diagnoses.
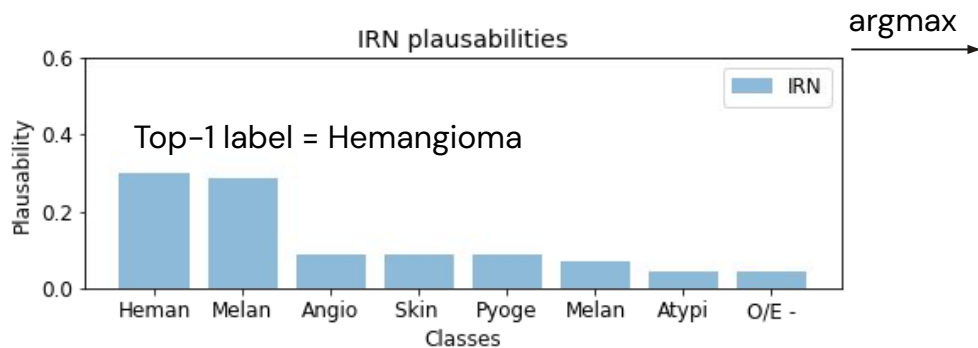
Observation          Annotations          Plausibilities $\hat{\lambda}$



argmax

Top-1 label = Hemangioma

# Case study: statistical aggregation using PrIRN

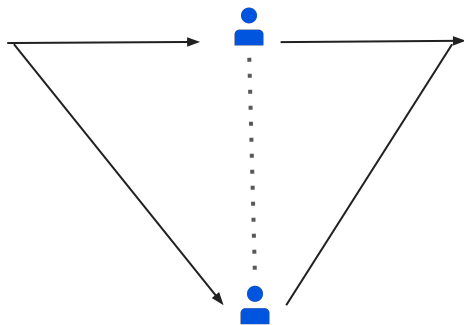Task: predict dermatological conditions from images.

- Plackett–Luce or probabilistic IRN (PrIRN) to model $p(\lambda|b)$
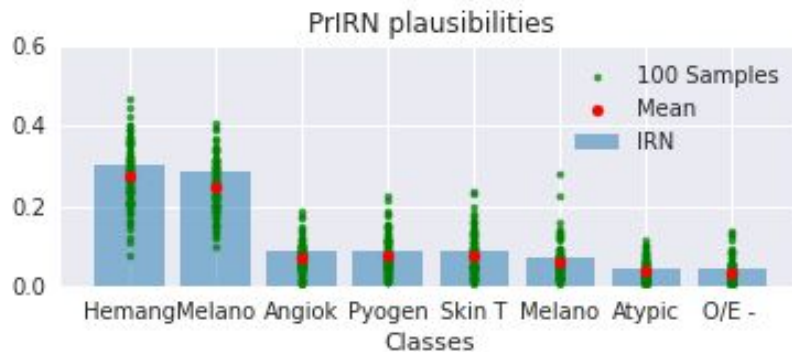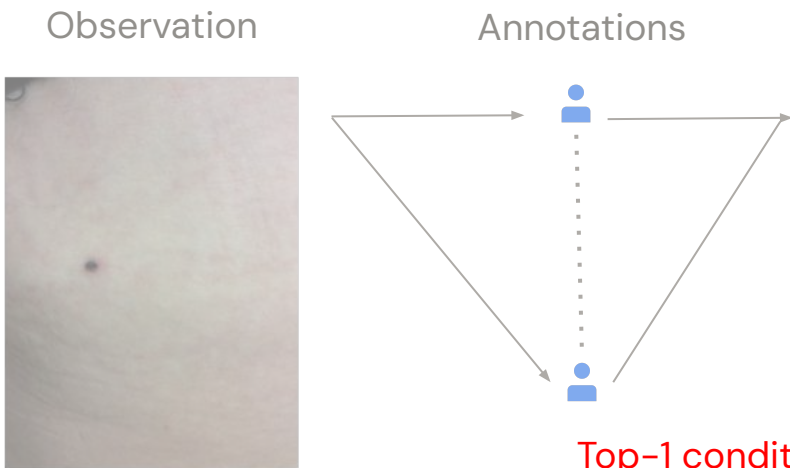
Observation                    Annotations              Plausibilities $\lambda^m \sim p_\gamma(\lambda|b)$
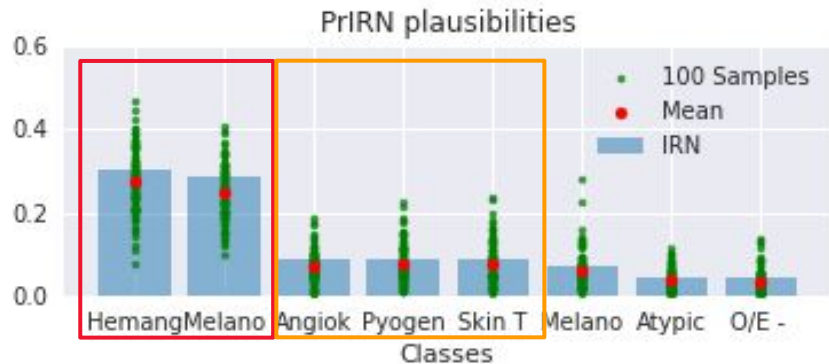
# Case study: statistical aggregation using PrIRN

Task: predict dermatological conditions from images.

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

**Observation**

**Annotations**

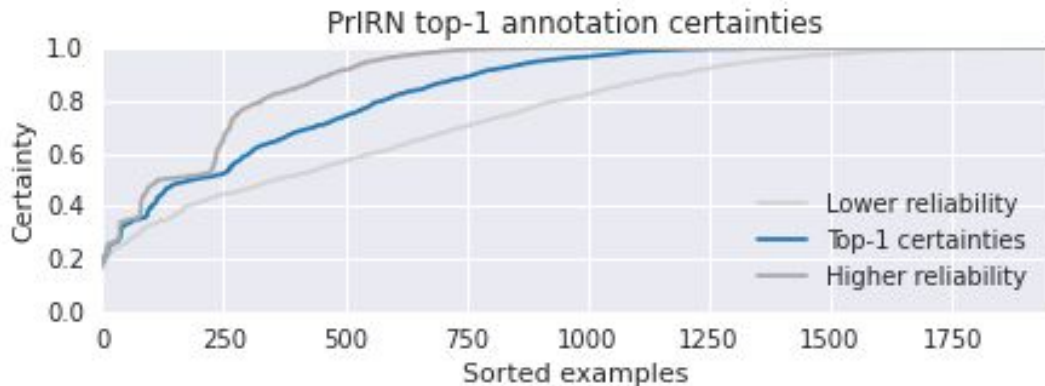Plausibilities $\lambda^m \sim p_\gamma(\lambda|b)$



Top-1 condition changes easily
= low annotation certainty

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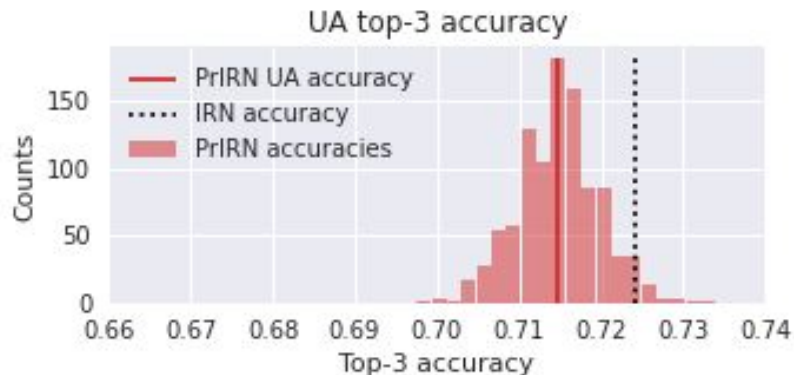change easily

# High annotation uncertainty

- Significant portions of cases with high annotation uncertainty:



PrIRN top-1 annotation certainties

- Lower reliability
- Top-1 certainties
- Higher reliability

➔ In discussions with dermatologists often attributed to inherent uncertainty

# Uncertainty-adjusted top-3 accuracy

- Across cases / per plausibility:



UA top-3 accuracy

- PrIRN UA accuracy
- IRN accuracy
- PrIRN accuracies

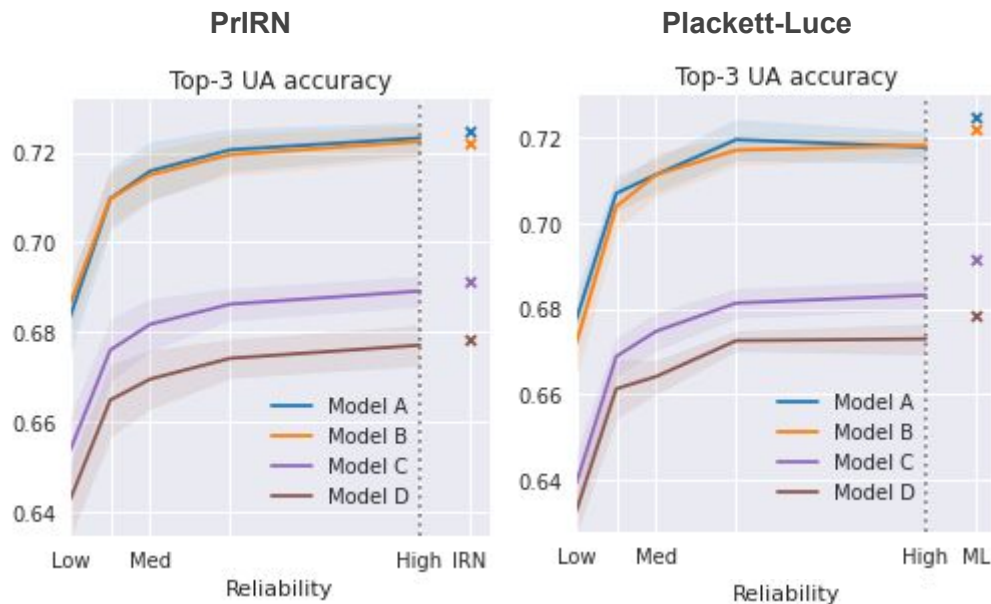➔ Significant variation in top-3 accuracy

# 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:



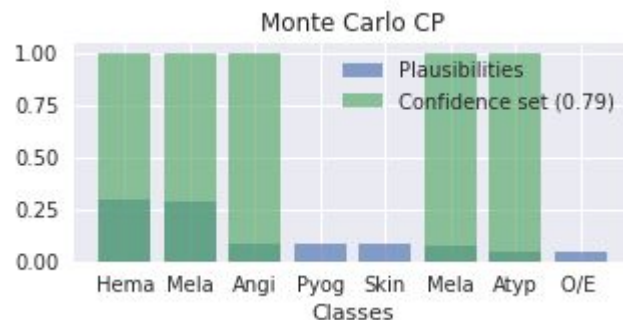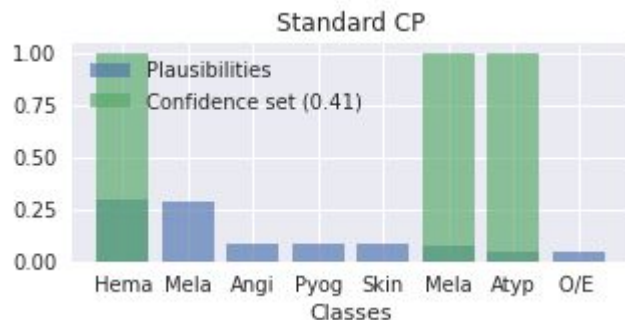➔ Aggregation is a mode choice usually not made explicit!

# 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 $\tau$
- Threshold used to predict confidence sets of classes at test time instead of the top–k:
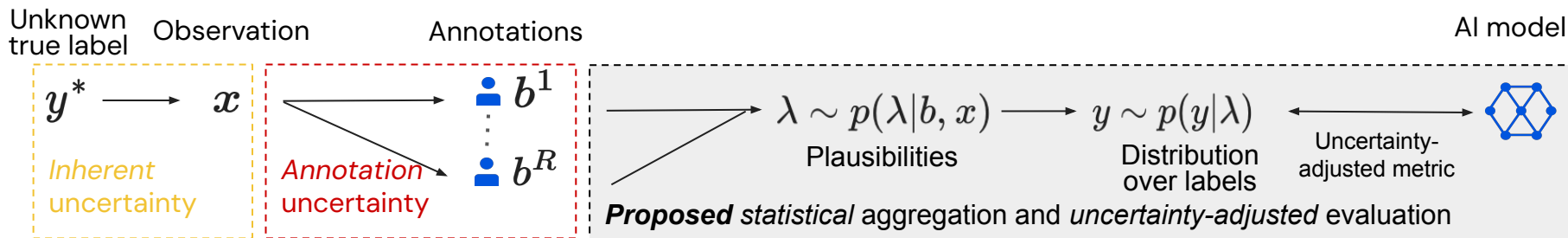
$$C_{\text{top-}k}(x) \quad\longrightarrow\quad C_{\text{CP}}(x) := \{k \in [K] : k - \text{th softmax} \geq \tau\}$$

- We propose *Monte Carlo* conformal prediction to calibrate directly against the annotations
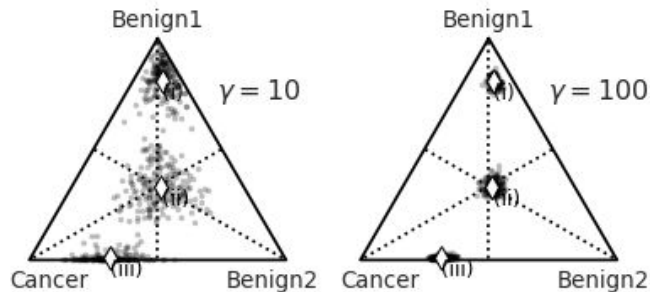
# Conclusion

Proposed a statistical framework for dealing with ground truth uncertainty:



Unknown true label — $y^*$ — *Inherent* uncertainty

Observation — $x$

Annotations — $b^1$ ... $b^R$ — *Annotation* uncertainty

$\lambda \sim p(\lambda|b, x)$ — Plausibilities

$y \sim p(y|\lambda)$ — Distribution over labels

Uncertainty-adjusted metric

AI model

***Proposed** statistical* aggregation and *uncertainty-adjusted* evaluation

➔ 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: arxiv.org/abs/2307.02191 | arxiv.org/abs/2307.09302 | davidstutz.de | dstutz@google.com