# Don't Blame the Data, Blame the Model:
## Understanding Noise and Bias
## When Learning from Subjective Annotations

Abhishek Anand, Negar Mokhberian, Prathyusha Naresh Kumar, Anweasha Saha, Zihao He, Ashwin Rao, Fred Morstatter, Kristina Lerman

# Motivation



Prompt: Do you see an airplane in the given picture?

# Motivation

# Motivation

Prompt: Do you see an airplane in the given picture?

Prompt: Do you think the given text is offensive?

Majority Vote Aggregation: Minority votes are discarded as Noise

Prompt: Do you see an airplane in the given picture?

Prompt: Do you think the given text is offensive?

**Objective**

Majority Vote Aggregation: Minority votes are discarded as Noise

Prompt: Do you see an airplane in the given picture?

Prompt: Do you think the given text is offensive?

**Objective**

**Subjective**

Majority Vote Aggregation: Minority votes are discarded as Noise

Prompt: Do you see an airplane in the given picture?

Prompt: Do you think the given text is offensive?

**Objective**

**Subjective**

Majority Vote Aggregation: Minority votes are discarded as Noise

Bias in Annotation

Prompt: Do you see an airplane in the given picture?

Prompt: Do you think the given text is offensive?

**Objective**

**Subjective**

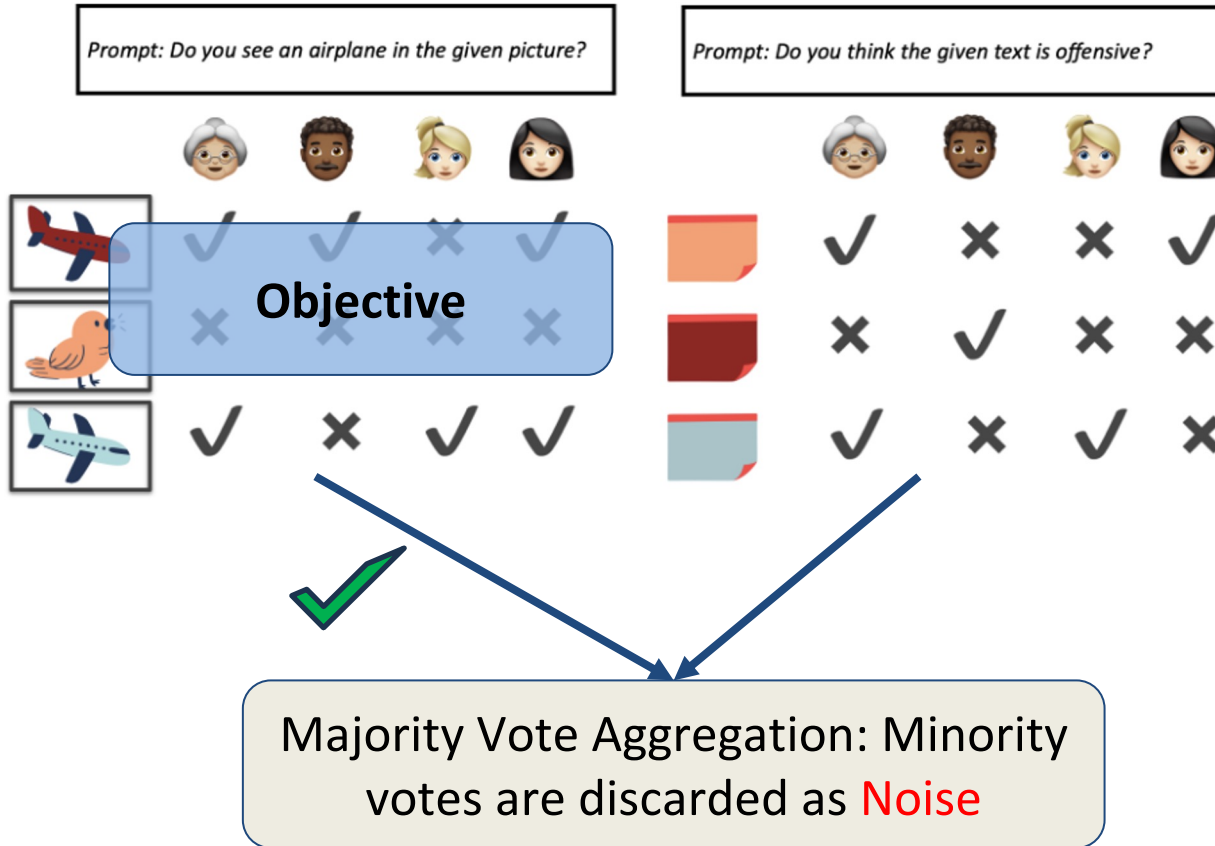Different perspectives essential for subjective tasks?

Majority Vote Aggregation: Minority votes are discarded as Noise

Bias in Annotation

# Motivation

Prompt: Do you see an airplane in the given picture?

Prompt: Do you think the given text is offensive?

**Objective**

**Subjective**

Different perspectives essential for subjective tasks?

Majority Vote Aggregation: Minority votes are discarded as Noise

Bias in Annotation

# Noise vs Bias?

*Information Sciences Institute*

# Research Questions?

Assumption – Single correct label exists

# Research Questions?

Assumption – Single correct label exists

Correlation between human disagreement on instances and model's uncertainty in prediction when using majority labels?

# Research Questions?

Assumption – Single correct label exists

Correlation between human disagreement on instances and model's uncertainty in prediction when using majority labels?

All annotations available

# Research Questions?

Assumption – Single correct label exists

Correlation between human disagreement on instances and model's uncertainty in prediction when using majority labels?

All annotations available

Does learning from raw annotations enhance the model's confidence?
Are perspectivist classification models effective?

Model uncertainty

# Model uncertainty

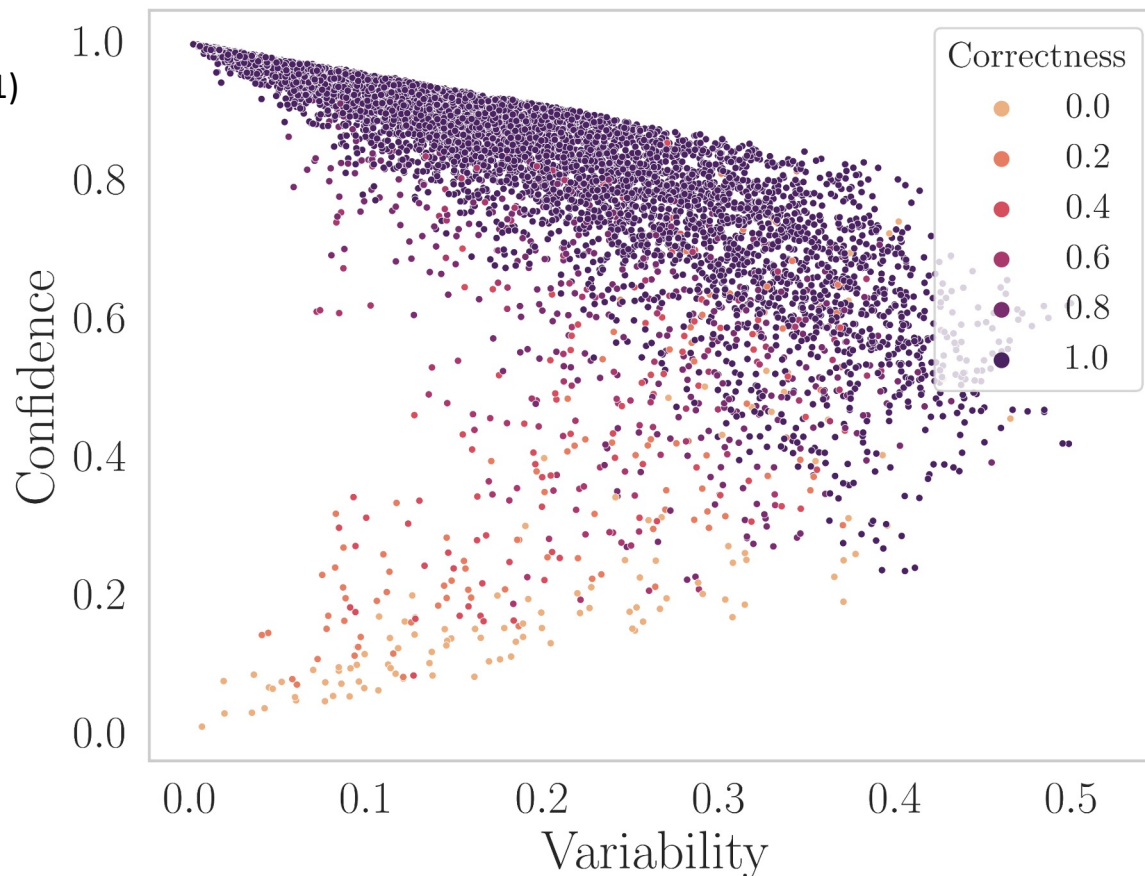Assumption – Noisy samples lead to uncertainty in modelling.

# Model uncertainty

Assumption – Noisy samples lead to uncertainty in modelling.

Quantifying uncertainty in modelling -  Data Maps (Swayamdipta et al., 2020)
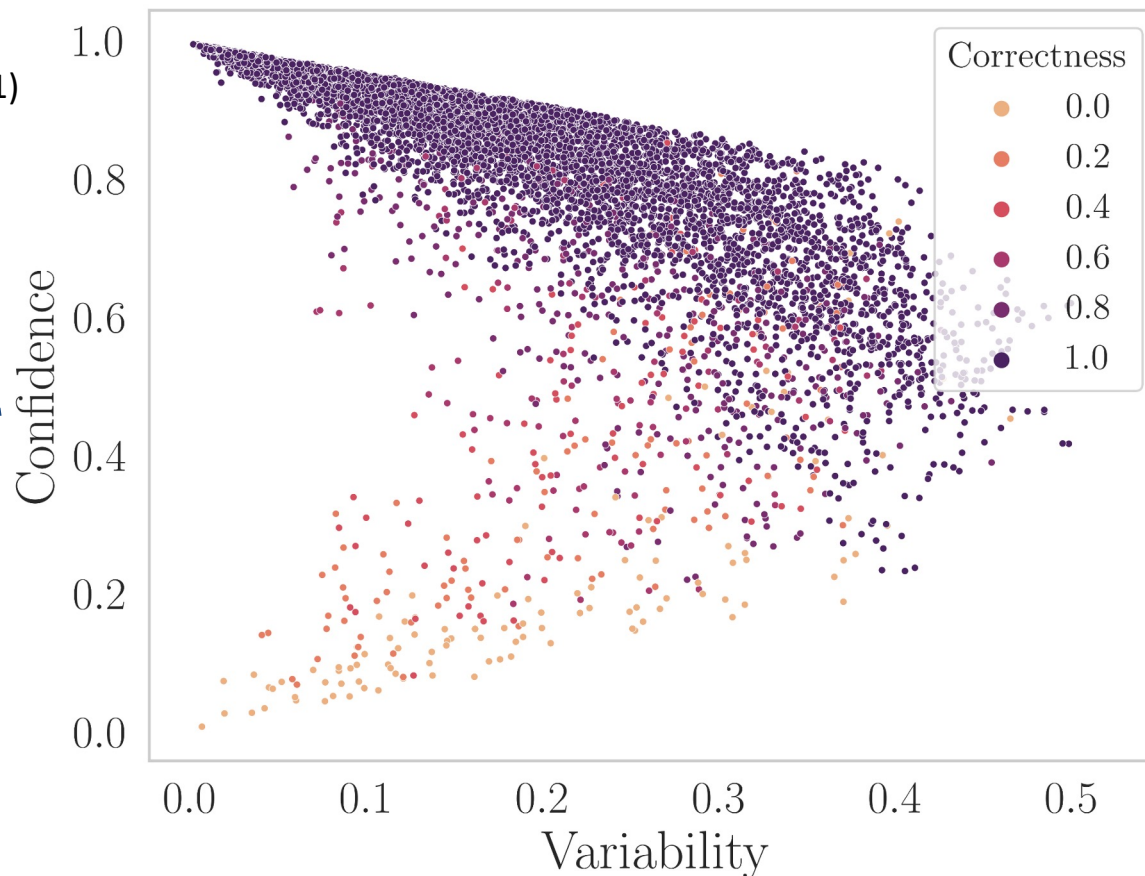
# Model uncertainty

Assumption – Noisy samples lead to uncertainty in modelling.

Quantifying uncertainty in modelling -  Data Maps (Swayamdipta et al., 2020)

**MDA** (Leonardelli et al., EMNLP 2021)
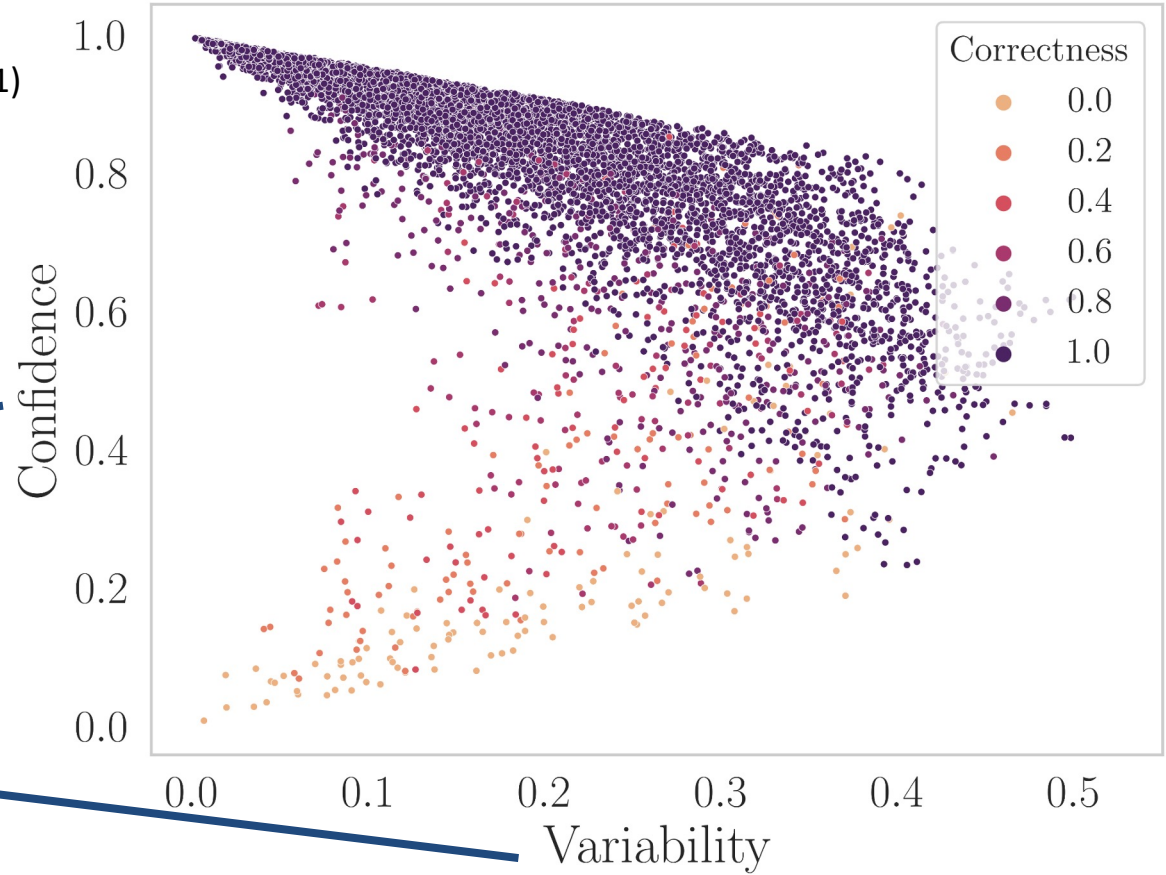
# Model uncertainty

Assumption – Noisy samples lead to uncertainty in modelling.

Quantifying uncertainty in modelling -  Data Maps (Swayamdipta et al., 2020)

**MDA** (Leonardelli et al., EMNLP 2021)

Mean of probabilities for gold label across epochs.

# Model uncertainty

Assumption – Noisy samples lead to uncertainty in modelling.

Quantifying uncertainty in modelling - Data Maps (Swayamdipta et al., 2020)

**MDA** (Leonardelli et al., EMNLP 2021)

Mean of probabilities for gold label across epochs.

Standard Deviation of probabilities for gold label across epochs.

Correlation between human disagreement on instances and model's uncertainty when using majority labels?

**Annotator Agreement Level ($a_m$) :**
fraction of annotations that align with majority vote for a text sample

Correlation between **human disagreement** on instances and model's uncertainty when using majority labels?

A measure to quantify disagreement between annotators on a label for a given sample.

**Annotator Agreement Level ($a_m$) :**
fraction of annotations that align with majority vote for a text sample
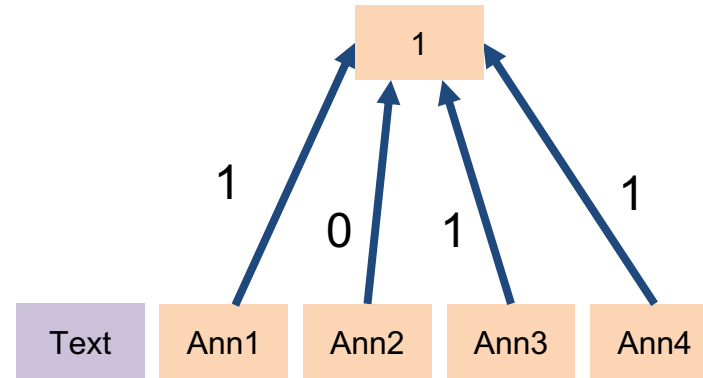
Correlation between **human disagreement** on instances and model's uncertainty when using majority labels?
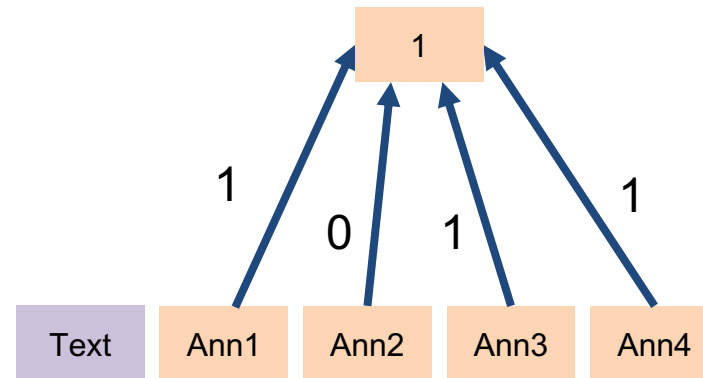
A measure to quantify disagreement between annotators on a label for a given sample.

**Annotator Agreement Level ($a_m$) :**
fraction of annotations that align with majority vote for a text sample

| Text | Ann1 | Ann2 | Ann3 | Ann4 |
|------|------|------|------|------|

Correlation between **human disagreement** on instances and model's uncertainty when using majority labels?

A measure to quantify disagreement between annotators on a label for a given sample.

**Annotator Agreement Level ($a_m$) :** fraction of annotations that align with majority vote for a text sample

Correlation between **human disagreement** on instances and model's uncertainty when using majority labels?

| | 1 | 0 | 1 | 1 |
| --- | --- | --- | --- | --- |
| Text | Ann1 | Ann2 | Ann3 | Ann4 |

A measure to quantify disagreement between annotators on a label for a given sample.

**Annotator Agreement Level**

Majority Vote = 1 ($a_m$=0.75)

**Annotator Agreement Level ($a_m$) :** fraction of annotations that align with majority vote for a text sample

1

1   0   1   1

Text   Ann1   Ann2   Ann3   Ann4

Correlation between **human disagreement** on instances and model's uncertainty when using majority labels?

Correlation between human disagreement on instances and model's uncertainty when using majority labels?

Data

| | |
|---|---|
| Text1 | Majority Label1 |
| Text2 | Majority Label2 |

Correlation between human disagreement on instances and model's uncertainty when using majority labels?

Single-GT
model

**Single Ground Truth Model**:
Majority vote label is considered as
ground truth. We fine-tune RoBERTa
for our study

Modelling

Data

| Text1 | Majority Label1 |
| Text2 | Majority Label2 |

Correlation between human disagreement on instances
and model's uncertainty when using majority labels?

# Single-GT model

**Single Ground Truth Model**:
Majority vote label is considered as ground truth. We fine-tune **RoBERTa** for our study

Modelling

## Data

| Text1 | Majority Label1 |
| Text2 | Majority Label2 |

Prediction for Text

Text Classification Model (RoBERTa)

Input Text

Correlation between human disagreement on instances and model's uncertainty when using majority labels?

# Datasets

| | Toxicity or Hate speech | | |
|---|---|---|---|
| | **MDA** (Leonardelli et al., EMNLP 2021) | **SBIC** (Sap et al., ACL 2020) | **MHS** (Kennedy et al., 2020) |
| **# Annotators** | 819 | 307 | 7,912 |
| **# Annotations per annotator** | 63.7±139 | 479.3±829.6 | 17.1±3.8 |
| **# Unique texts** | 10,440 | 45,318 | 39,565 |
| **# Annotations per text** | 5 | 3.2±1.2 | 2.3±1.0 |
| **# Number of labels** | 2 | 2 | 3 |

# Datasets

| | Toxicity or Hate speech | | |
|---|---|---|---|
| | **MDA** (Leonardelli et al., EMNLP 2021) | **SBIC** (Sap et al., ACL 2020) | **MHS** (Kennedy et al., 2020) |
| **# Annotators** | 819 | 307 | 7,912 |
| **# Annotations per annotator** | 63.7±139 | 479.3±829.6 | 17.1±3.8 |
| **# Unique texts** | 10,440 | 45,318 | 39,565 |
| **# Annotations per text** | 5 | 3.2±1.2 | 2.3±1.0 |
| **# Number of labels** | 2 | 2 | 3 |

# Datasets

| | Toxicity or Hate speech | | |
|---|---|---|---|
| | **MDA** (Leonardelli et al., EMNLP 2021) | **SBIC** (Sap et al., ACL 2020) | **MHS** (Kennedy et al., 2020) |
| **# Annotators** | 819 | 307 | 7,912 |
| **# Annotations per annotator** | 63.7±139 | 479.3±829.6 | 17.1±3.8 |
| **# Unique texts** | 10,440 | 45,318 | 39,565 |
| **# Annotations per text** | 5 | 3.2±1.2 | 2.3±1.0 |
| **# Number of labels** | 2 | 2 | 3 |

# Confidence vs Annotator Agreement Level

# Confidence vs Annotator Agreement Level

MDA

# Confidence vs Annotator Agreement Level

MDA

# Confidence vs Annotator Agreement Level

MDA

# Confidence vs Annotator Agreement Level

# Confidence vs Annotator Agreement Level

SBIC

# Confidence vs Annotator Agreement Level
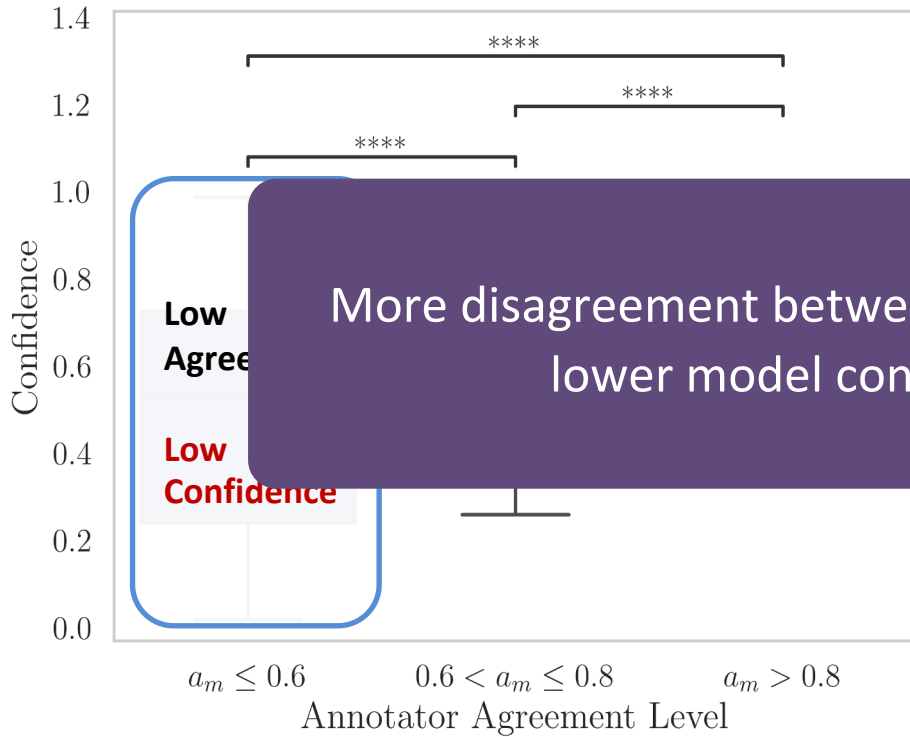


SBIC

MHS

# Confidence vs Annotator Agreement Level



SBIC

MHS

# Confidence vs Annotator Agreement Level



SBIC

MHS

# Confidence vs Annotator Agreement Level



SBIC

MHS

More disagreement between annotators correlates with lower model confidence for sample

# Confidence vs Annotator Agreement Level

SBIC

MHS



Low Agree...

**Low Confidence**

**Confidence**

More disagreement between annotators correlates with lower model confidence for sample

Does learning from raw annotations enhance the model's confidence for the high disagreement instances?

Does learning from raw annotations enhance the model's confidence for the high disagreement instances?
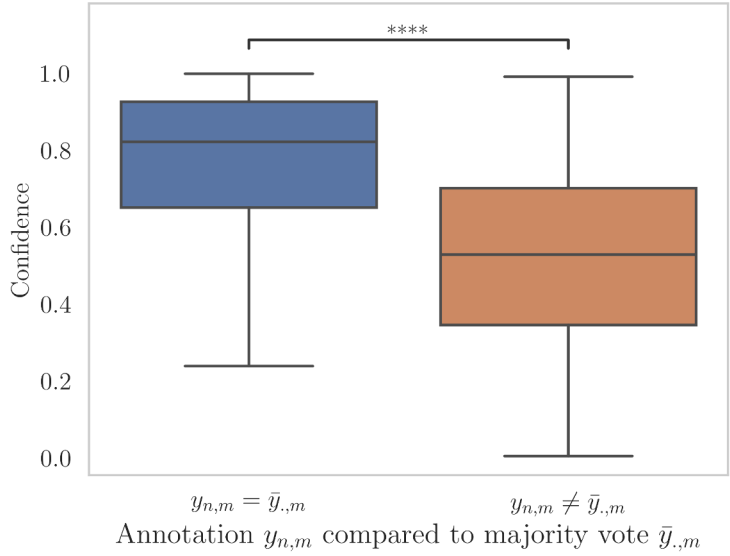
Data

| Text1 | Ann1 | Ann2 | Ann3 | Ann4 |
| Text2 | Ann1 | Ann2 | Ann3 | Ann4 |

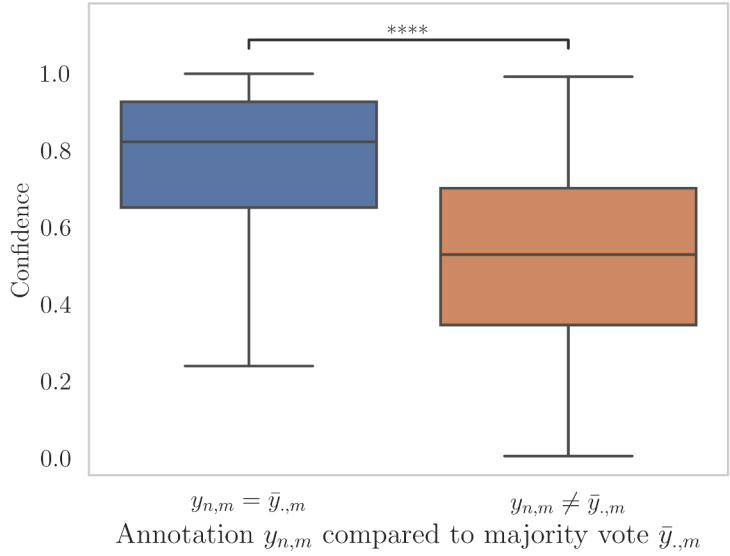Does learning from **raw annotations** enhance the model's confidence for the high disagreement instances?
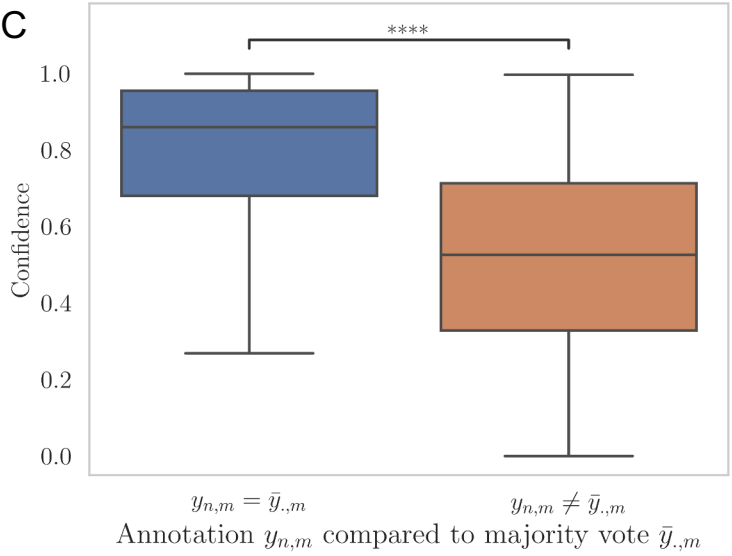
**Multi-GT model**

**Multiple Ground Truth Model**: Each annotation by an annotator is considered a ground truth. We consider DISCO (Weerasooriya et al., 2023) for our study

Requirement: Models that take raw annotations as input

**Data**

| Text1 | Ann1 | Ann2 | Ann3 | Ann4 |
| Text2 | Ann1 | Ann2 | Ann3 | Ann4 |

Does learning from **raw annotations** enhance the model's confidence for the high disagreement instances?

# Multi-GT model

**Multiple Ground Truth Model**: Each annotation by an annotator is considered a ground truth. We consider **DISCO** (Weerasooriya et al., 2023) for our study

Prediction for Text, Annotator Id

Requirement: Models that take raw annotations as input

# Data

| Text1 | Ann1 | Ann2 | Ann3 | Ann4 |
| Text2 | Ann1 | Ann2 | Ann3 | Ann4 |

Input Text    Annotator Id

Does learning from **raw annotations** enhance the model's confidence for the high disagreement samples?

USC Viterbi
School of Engineering

# Multi-GT model

MDA



*Information Sciences Institute*

# Multi-GT model

*Information Sciences Institute*

# Multi-GT model
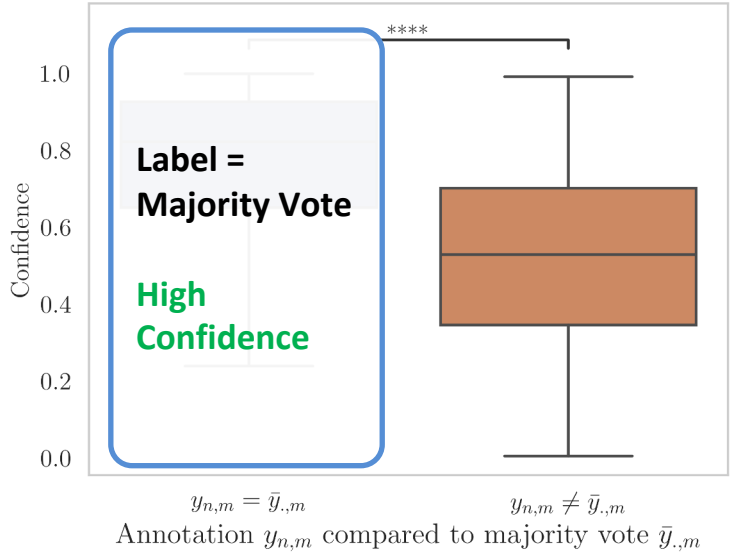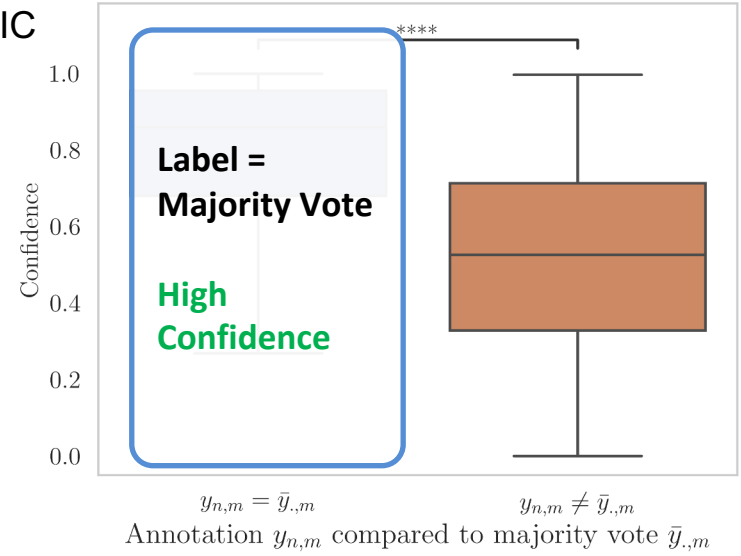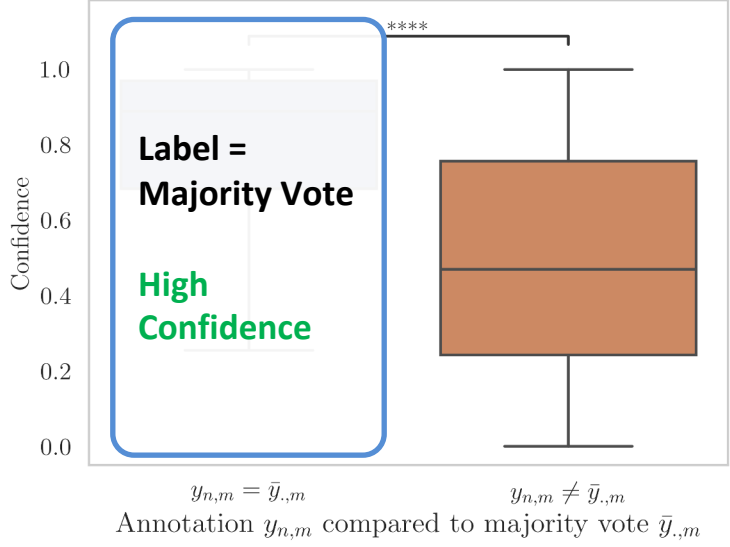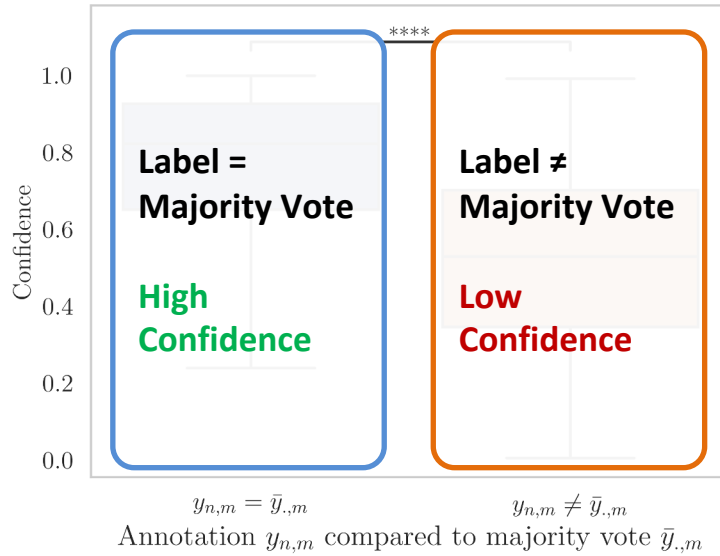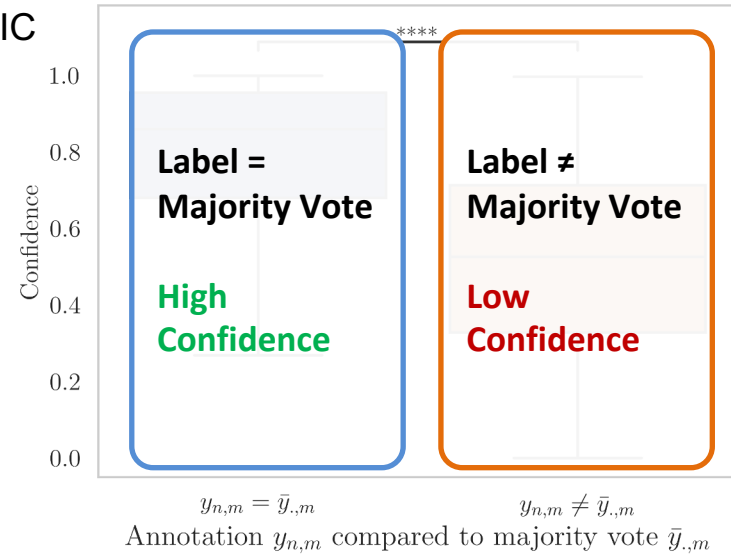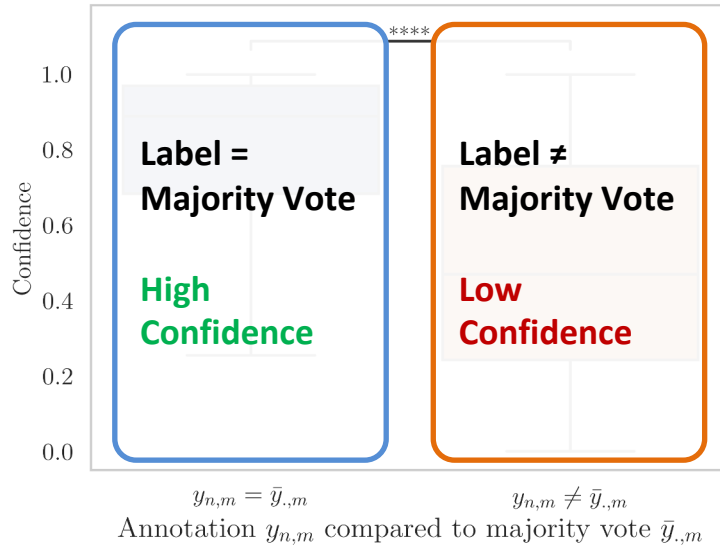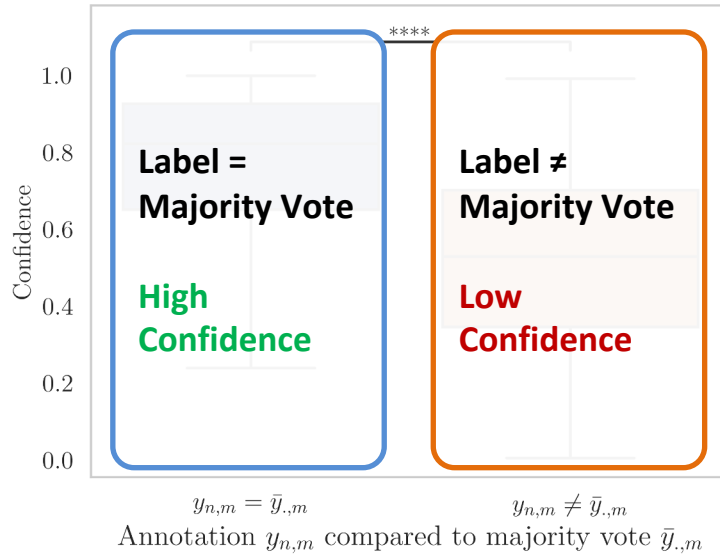
## Model Confidence vs Agreement with Majority vote

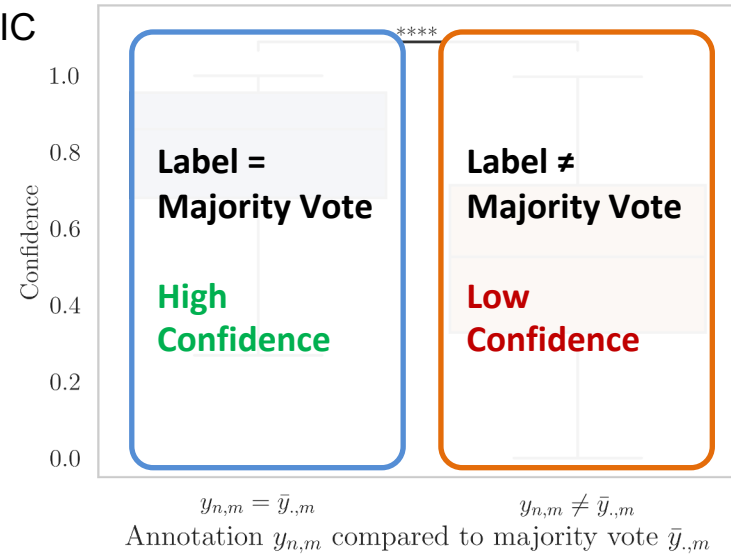# Multi-GT model

MDA



**Label = Majority Vote**

**High Confidence**

****

$y_{n,m} = \bar{y}_{.,m}$    $y_{n,m} \neq \bar{y}_{.,m}$

Annotation $y_{n,m}$ compared to majority vote $\bar{y}_{.,m}$

SBIC

**Label = Majority Vote**

**High Confidence**

****

$y_{n,m} = \bar{y}_{.,m}$    $y_{n,m} \neq \bar{y}_{.,m}$

Annotation $y_{n,m}$ compared to majority vote $\bar{y}_{.,m}$

MHS

**Label = Majority Vote**

**High Confidence**

****

$y_{n,m} = \bar{y}_{.,m}$    $y_{n,m} \neq \bar{y}_{.,m}$

Annotation $y_{n,m}$ compared to majority vote $\bar{y}_{.,m}$

# Multi-GT model

## Model Confidence vs Agreement with Majority vote

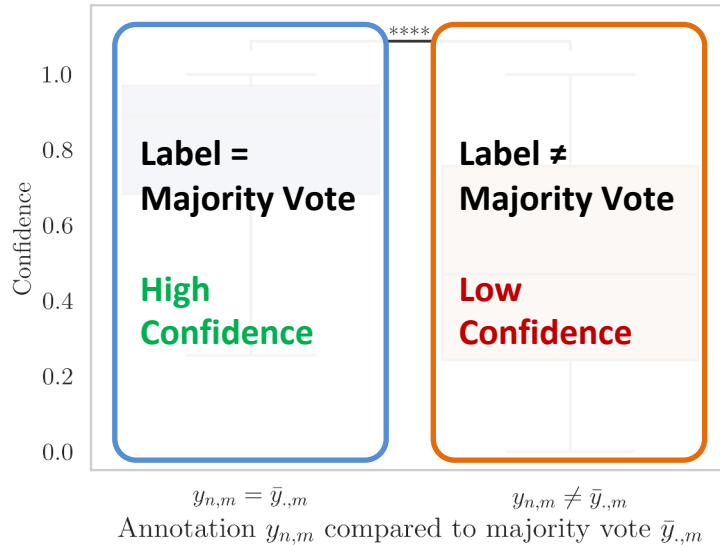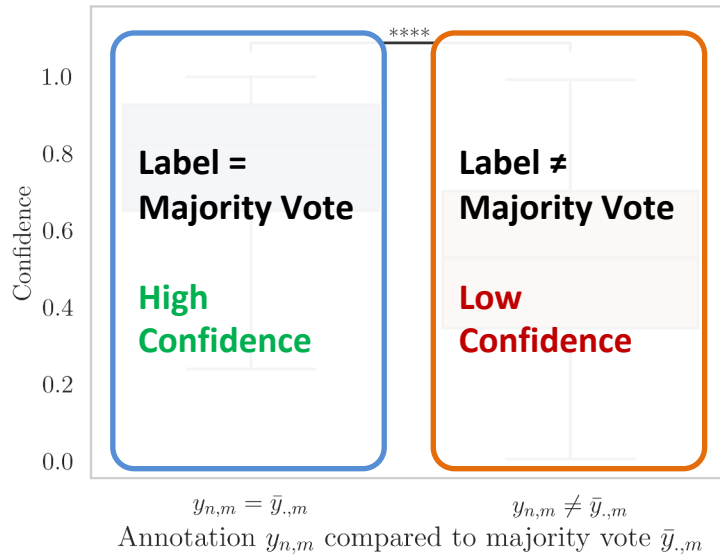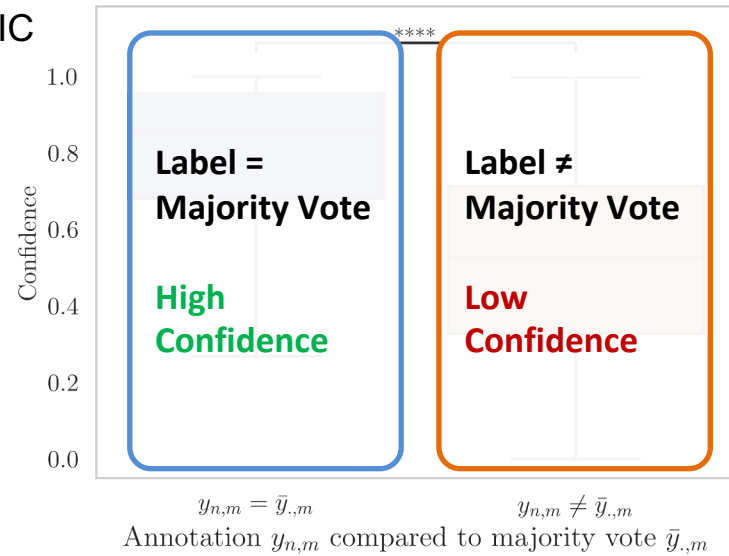**MDA**



Label =
Majority Vote

**High
Confidence**

Label ≠
Majority Vote

**Low
Confidence**

****

$y_{n,m} = \bar{y}_{.,m}$   $y_{n,m} \neq \bar{y}_{.,m}$

Annotation $y_{n,m}$ compared to majority vote $\bar{y}_{.,m}$

**SBIC**



Label =
Majority Vote

**High
Confidence**
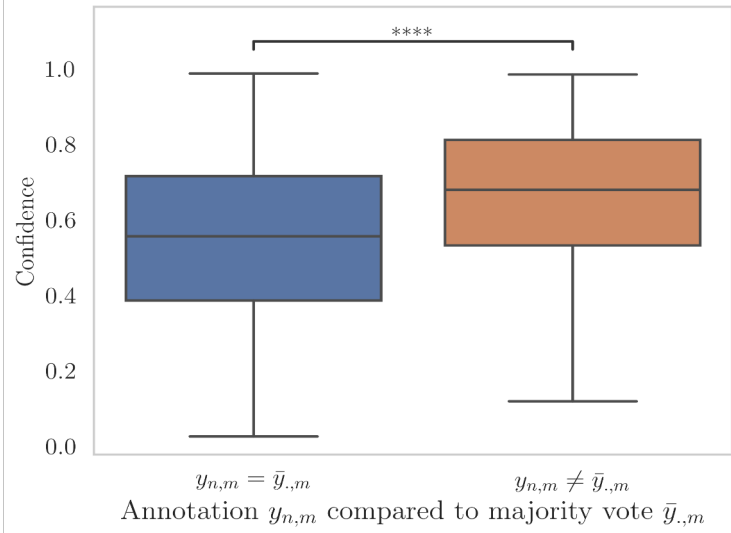
Label ≠
Majority Vote

**Low
Confidence**

****

$y_{n,m} = \bar{y}_{.,m}$   $y_{n,m} \neq \bar{y}_{.,m}$

Annotation $y_{n,m}$ compared to majority vote $\bar{y}_{.,m}$

**MHS**



Label =
Majority Vote

**High
Confidence**

Label ≠
Majority Vote

**Low
Confidence**

****

$y_{n,m} = \bar{y}_{.,m}$   $y_{n,m} \neq \bar{y}_{.,m}$

Annotation $y_{n,m}$ compared to majority vote $\bar{y}_{.,m}$

*Information Sciences Institute*

USC Viterbi
School of Engineering

# Multi-GT model

## MDA



Label = Majority Vote

**High Confidence**

Label ≠ Majority Vote

**Low Confidence**

****

$y_{n,m} = \bar{y}_{.,m}$     $y_{n,m} \neq \bar{y}_{.,m}$

Annotation $y_{n,m}$ compared to majority vote $\bar{y}_{.,m}$

Confidence

## SBIC



Label = Majority Vote

**High Confidence**

Label ≠ Majority Vote

**Low Confidence**

****

$y_{n,m} = \bar{y}_{.,m}$     $y_{n,m} \neq \bar{y}_{.,m}$

Annotation $y_{n,m}$ compared to majority vote $\bar{y}_{.,m}$

Confidence

## MHS



Label = Majority Vote

**High Confidence**

Label ≠ Majority Vote

**Low Confidence**

****

$y_{n,m} = \bar{y}_{.,m}$     $y_{n,m} \neq \bar{y}_{.,m}$

Annotation $y_{n,m}$ compared to majority vote $\bar{y}_{.,m}$

Confidence

Consistent with observed trends in Single-GT model, higher agreement correlates with higher confidence

*Information Sciences Institute*

USC Viterbi
School of Engineering

# Multi-GT model

MDA

Label = Majority Vote

High Confidence

Label ≠ Majority Vote

Low Confidence

****

$y_{n,m} = \bar{y}_{.,m}$

$y_{n,m} \neq \bar{y}_{.,m}$

Annotation $y_{n,m}$ compared to majority vote $\bar{y}_{.,m}$

SBIC

Label = Majority Vote

High Confidence

Label ≠ Majority Vote

Low Confidence

****

$y_{n,m} = \bar{y}_{.,m}$

$y_{n,m} \neq \bar{y}_{.,m}$

Annotation $y_{n,m}$ compared to majority vote $\bar{y}_{.,m}$

MHS

Label = Majority Vote

High Confidence

Label ≠ Majority Vote

Low Confidence

****

$y_{n,m} = \bar{y}_{.,m}$

$y_{n,m} \neq \bar{y}_{.,m}$

Annotation $y_{n,m}$ compared to majority vote $\bar{y}_{.,m}$

Consistent with observed trends in Single-GT model, higher agreement correlates with higher confidence

Can Multi-GT improve on high disagreement samples from Single-GT?

USC Viterbi
School of Engineering

# Multi-GT model    (Low confidence (< 0.5) samples in Single-GT)

MDA

# Multi-GT model   (Low confidence (< 0.5) samples in Single-GT)

MDA



SBIC

# Multi-GT model    (Low confidence (< 0.5) samples in Single-GT)



*Information Sciences Institute*

# Multi-GT model    (Low confidence (< 0.5) samples in Single-GT)

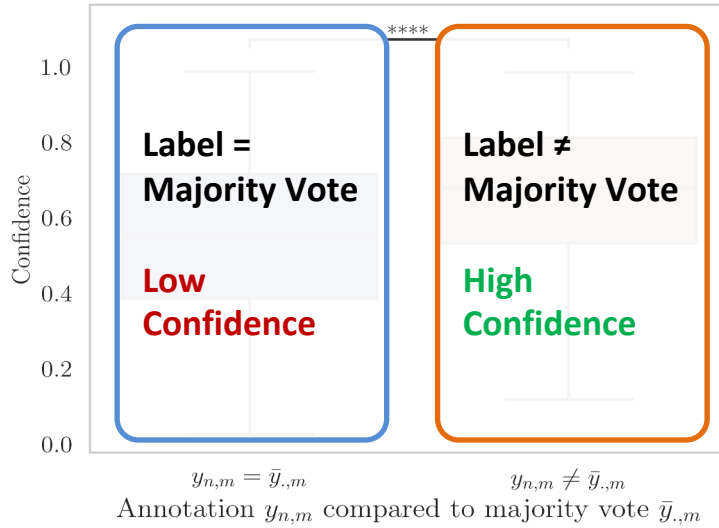MDA



SBIC



MHS

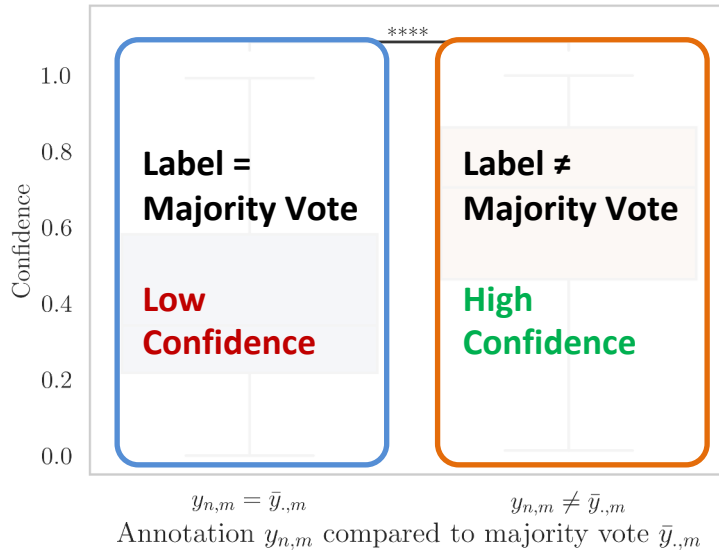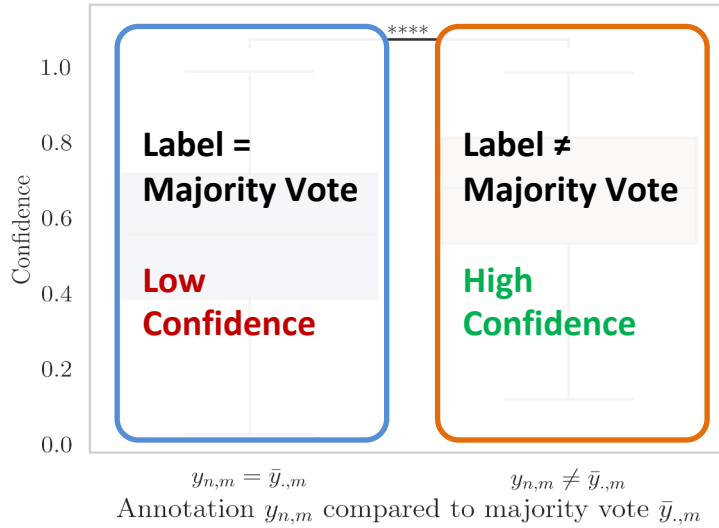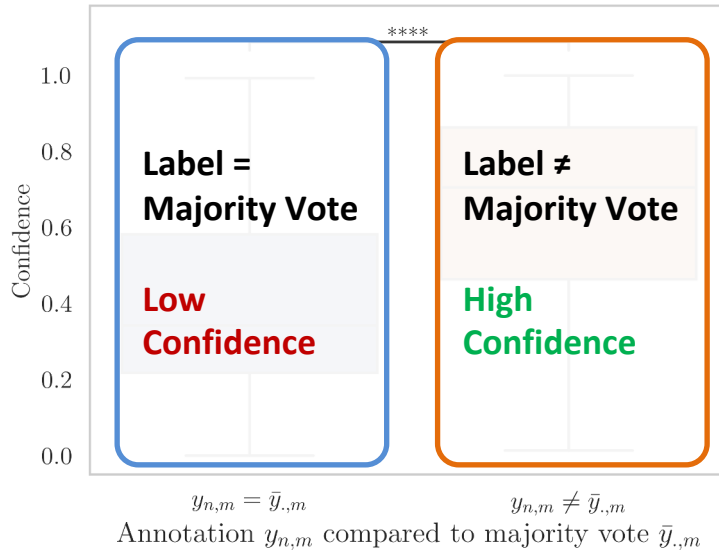# Multi-GT model  (Low confidence (< 0.5) samples in Single-GT)



MDA

Confidence

1.0
0.8
0.6
0.4
0.2
0.0

**Label =
Majority Vote**

**Low
Confidence**

****

**Label ≠
Majority Vote**

**High
Confidence**

$y_{n,m} = \bar{y}_{.,m}$  $y_{n,m} \neq \bar{y}_{.,m}$
Annotation $y_{n,m}$ compared to majority vote $\bar{y}_{.,m}$

SBIC

Confidence

1.0
0.8
0.6
0.4
0.2
0.0

**Label =
Majority Vote**

**Low
Confidence**

****

**Label ≠
Majority Vote**

**High
Confidence**

$y_{n,m} = \bar{y}_{.,m}$  $y_{n,m} \neq \bar{y}_{.,m}$
Annotation $y_{n,m}$ compared to majority vote $\bar{y}_{.,m}$

MHS

Confidence

1.0
0.8
0.6
0.4
0.2
0.0

**Label =
Majority Vote**

**Low
Confidence**

****

**Label ≠
Majority Vote**

**High
Confidence**

$y_{n,m} = \bar{y}_{.,m}$  $y_{n,m} \neq \bar{y}_{.,m}$
Annotation $y_{n,m}$ compared to majority vote $\bar{y}_{.,m}$

USC Viterbi
School of Engineering

# Multi-GT model (Low confidence (< 0.5) samples in Single-GT)

MDA



SBIC



MHS



Provided the annotations that were discarded as noise, DISCO learns valuable signals boosting confidence on these samples

# Multi-GT model



Prompt: Do you think the given text is offensive?

Prompt: Do you think the given text is offensive?

**Subjective**

# Multi-GT model

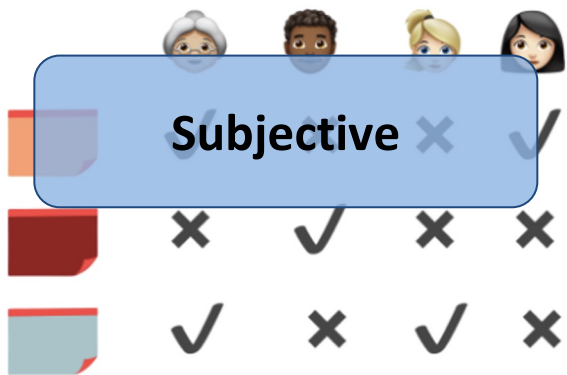Prompt: Do you think the given text is offensive?

**Subjective**

Different perspectives essential for subjective tasks

Prompt: Do you think the given text is offensive?

**Subjective**

Different perspectives essential for subjective tasks

Can the model learn multiple annotators' perspectives?

# Multi-GT model

Prompt: Do you think the given text is offensive?

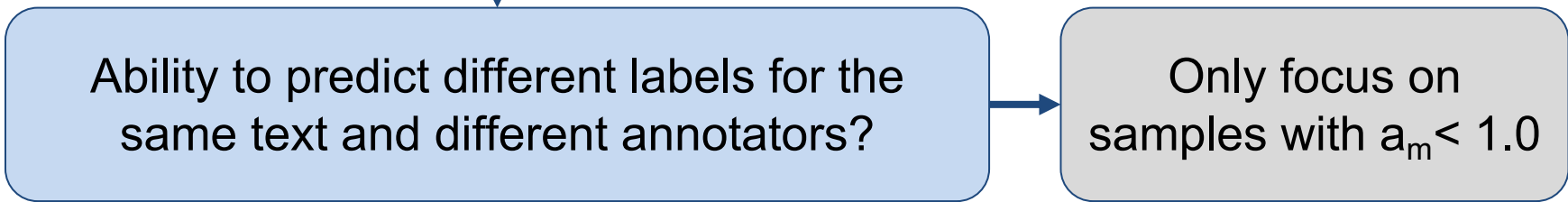**Subjective**

Different perspectives essential for subjective tasks

Can the model learn multiple annotators' perspectives?

Ability to predict different labels for the same text and different annotators?

# Multi-GT model



Prompt: Do you think the given text is offensive?

**Subjective**

Different perspectives essential for subjective tasks

Can the model learn multiple annotators' perspectives?

Ability to predict different labels for the same text and different annotators?
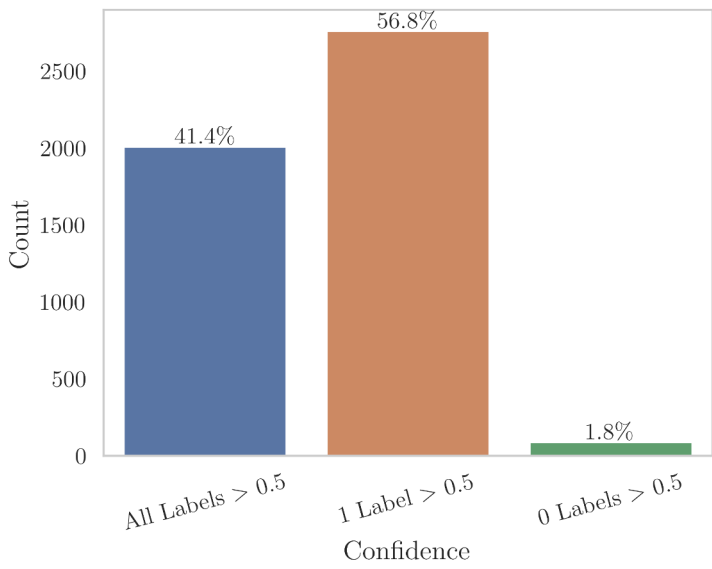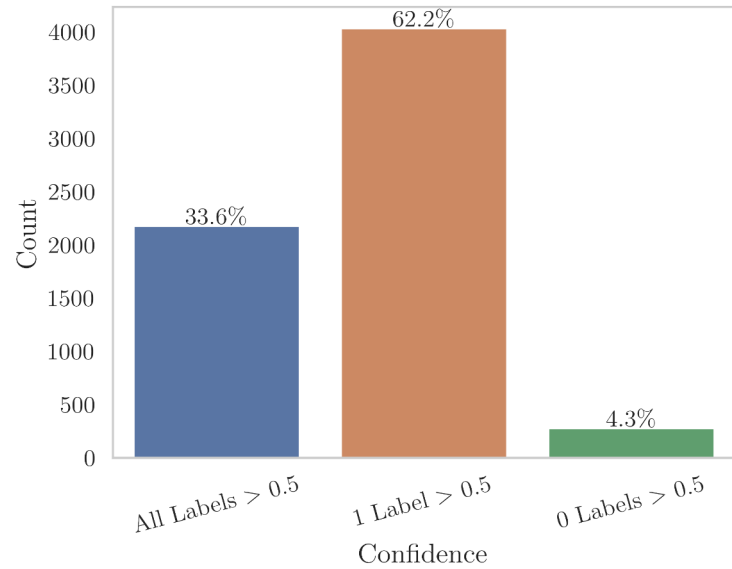
Only focus on samples with $a_m < 1.0$

# Multi-GT model

Can it learn multiple annotators' perspectives?

# Multi-GT model

## Can it learn multiple annotators' perspectives?

MDA

# Multi-GT model

## Can it learn multiple annotators' perspectives?



MDA

SBIC

# Multi-GT model
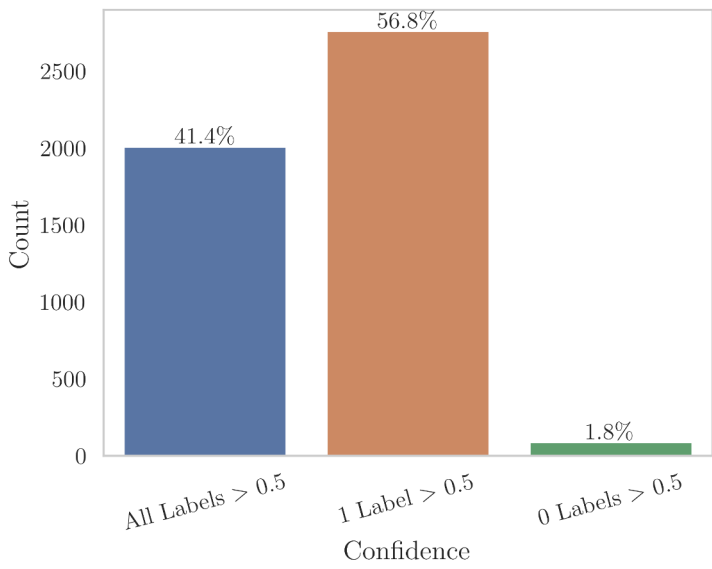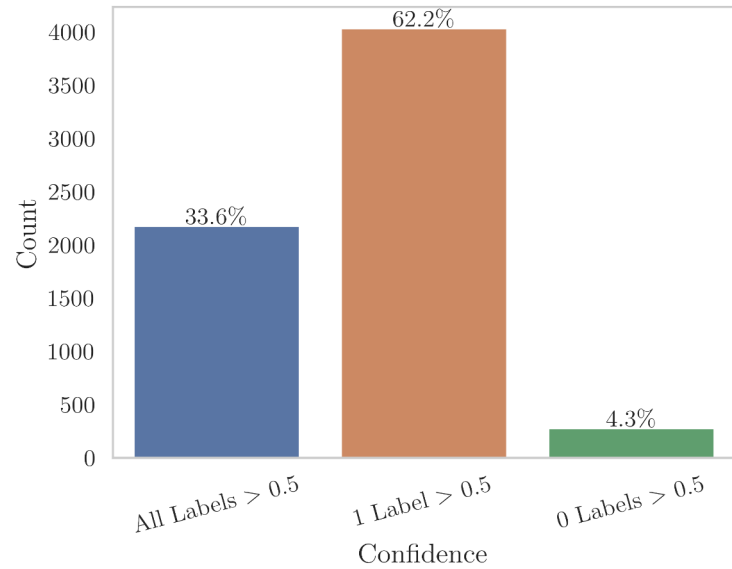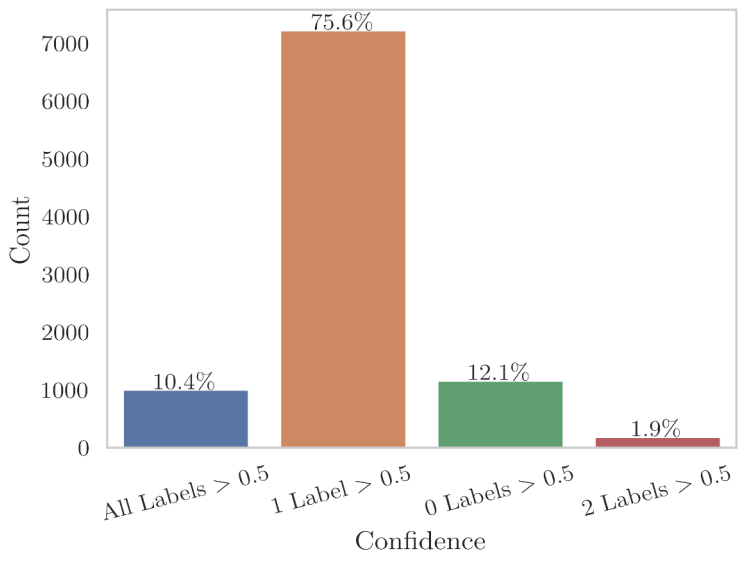
## Can it learn multiple annotators' perspectives?

# Multi-GT model

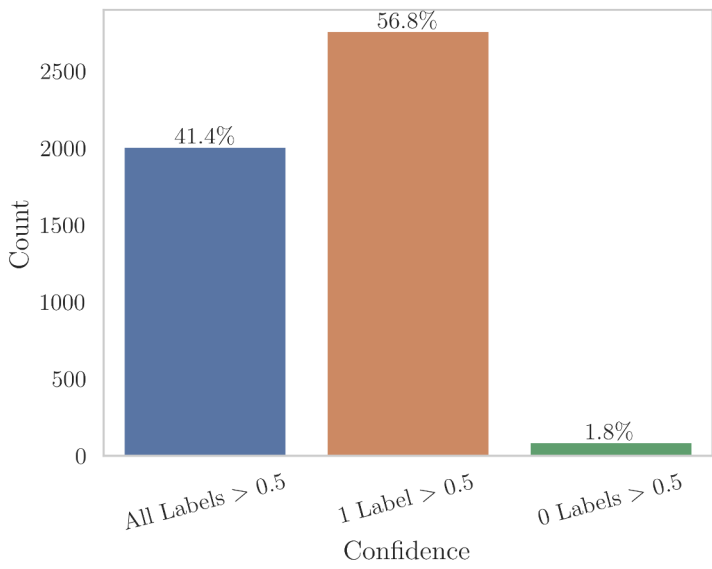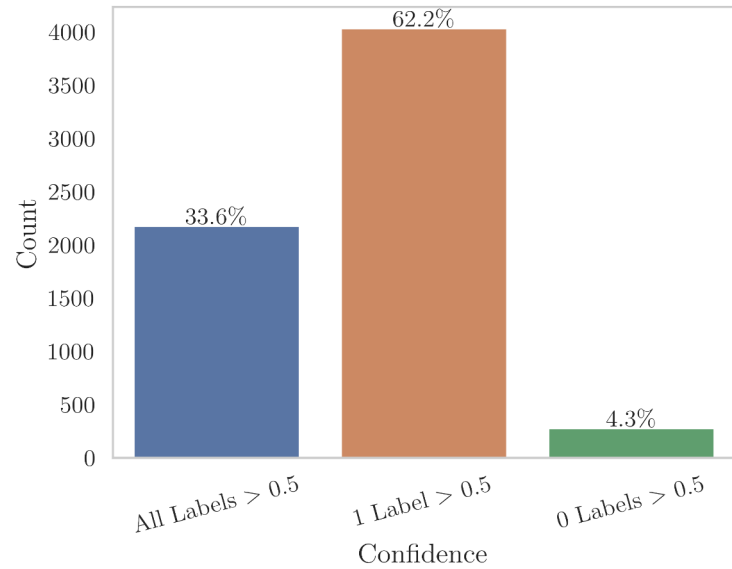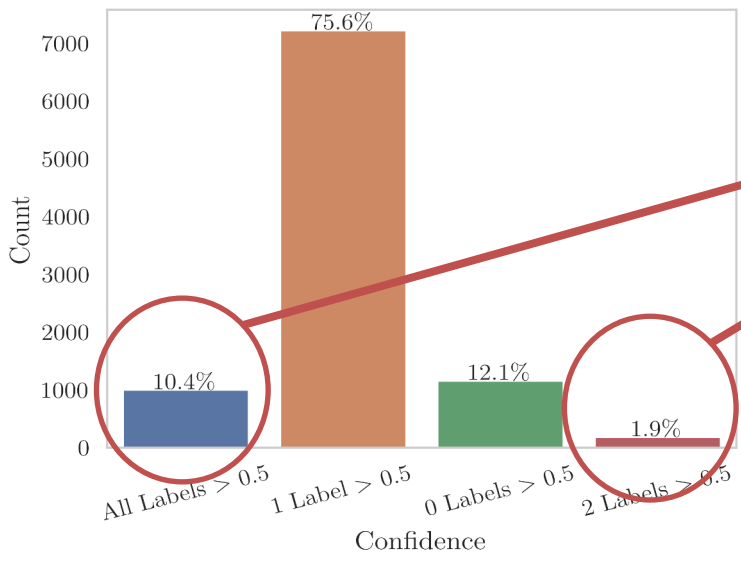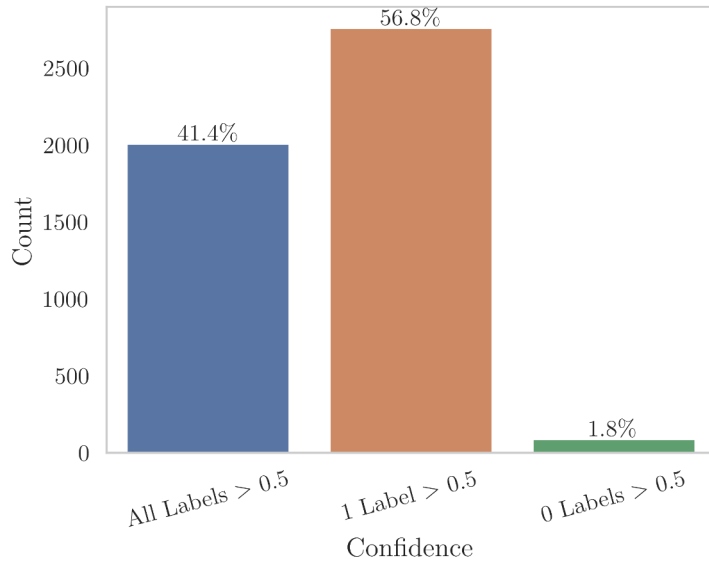## Can it learn multiple annotators' perspectives?



MDA

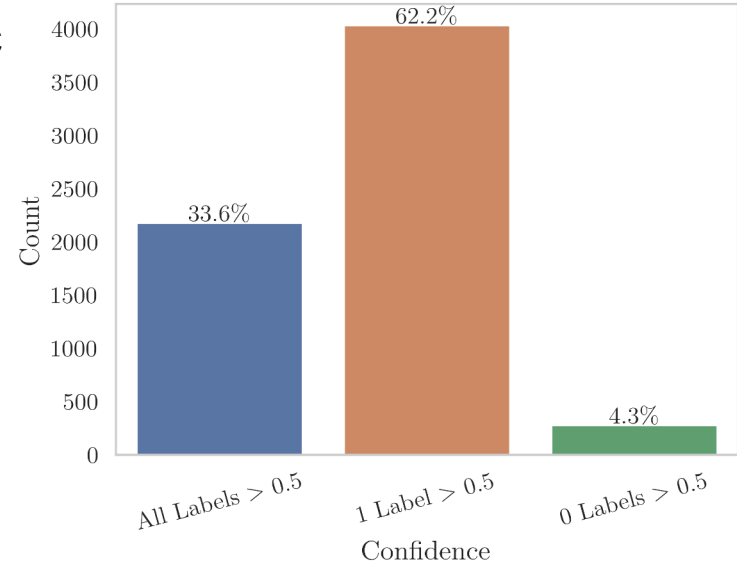SBIC

MHS

Finds it difficult to learn multiple labels for MHS

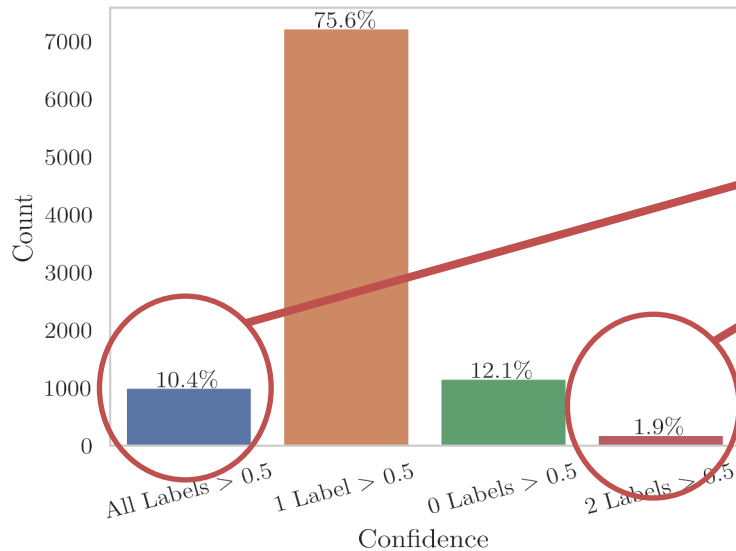# Multi-GT model

## Can it learn multiple annotators' perspectives?



Finds it difficult to learn multiple labels for MHS

This dataset is extra challenging because avg number of annotations per annotator is ~ 17!

# Takeaways

## **Noise** vs **Bias**?

> More disagreement between annotators correlates with low model confidence

# Takeaways

## **Noise** vs **Bias**?

More disagreement between annotators correlates with low model confidence

Majority vote captures just a single perspective - insufficient for subjective tasks

# Takeaways

## **Noise** vs **Bias**?

> More disagreement between annotators correlates with low model confidence

> Majority vote captures just a single perspective - insufficient for subjective tasks

> Multi-Gt model effectively utilizes minority vote annotations that are usually discarded as noise

*Information Sciences Institute*

# Takeaways

## **Noise** vs **Bias**?

More disagreement between annotators correlates with low model confidence

Majority vote captures just a single perspective - insufficient for subjective tasks

Multi-Gt model effectively utilizes minority vote annotations that are usually discarded as noise

Number of annotations per annotator important in modelling their perspective

*Information Sciences Institute*