

ML4Seismic Partners Meeting 2023

A New Seismic Fault Label Uncertainty Dataset: Insights from Expertise, Certainty, and Consistency

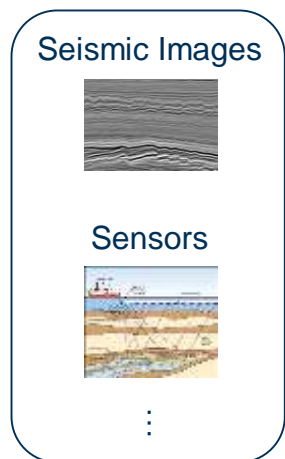
Jorge Quesada, Chen Zhou, Mohit Prabhushankar, and Ghassan AlRegib

Introduction

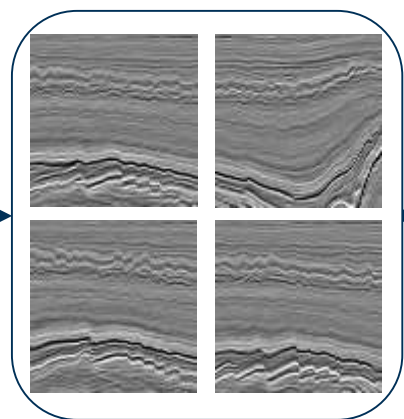
Where does annotator disagreement stand in the interpretation pipeline?

Proposed Uncertainty Framework

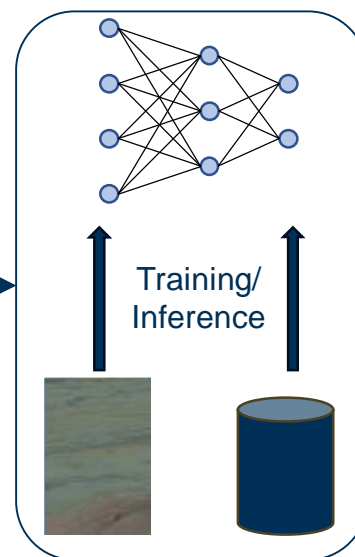
Data Parameters



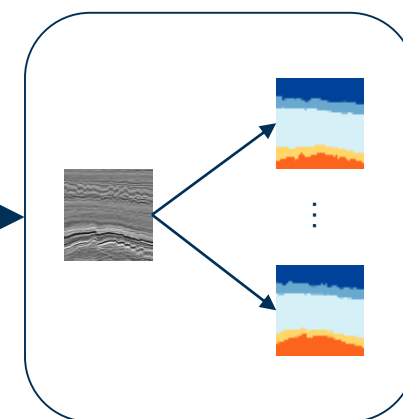
Seismic Data



Algorithm



Labels



Data
 ξ

$$p(X|\xi)$$

Images
 X

$$p(\alpha|X, W)$$

Interpretations
 $\alpha = \text{"Labels"}$
Interpretational
Uncertainty

Data
Uncertainty

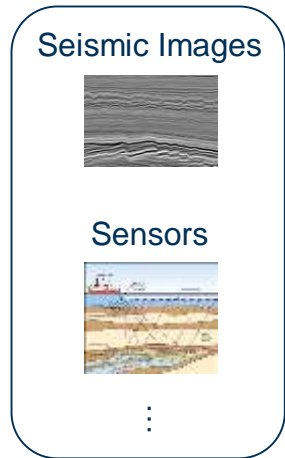
Model
Uncertainty

Introduction

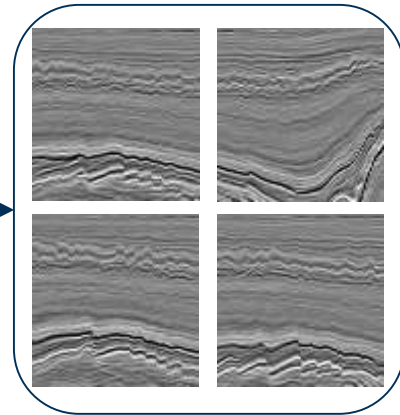
Where does annotator disagreement stand in the interpretation pipeline?

Proposed Uncertainty Framework

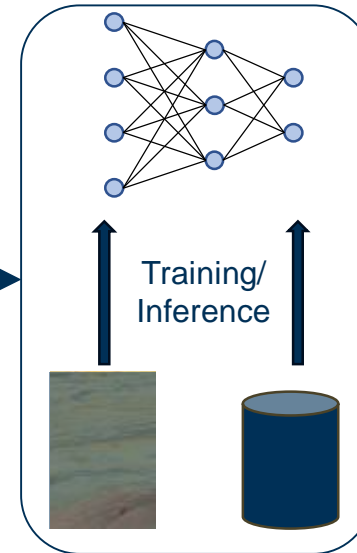
Data Parameters



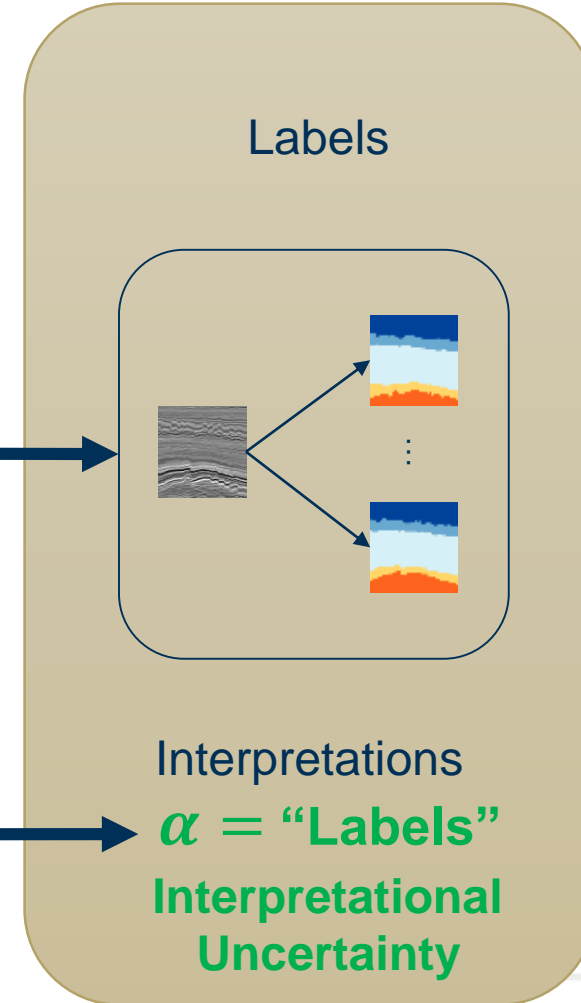
Seismic Data



Algorithm



Labels



Data
 ξ

$$p(X|\xi)$$

Images

X

$$p(\alpha|X, W)$$

Interpretations
 $\alpha = \text{“Labels”}$
Interpretational
Uncertainty

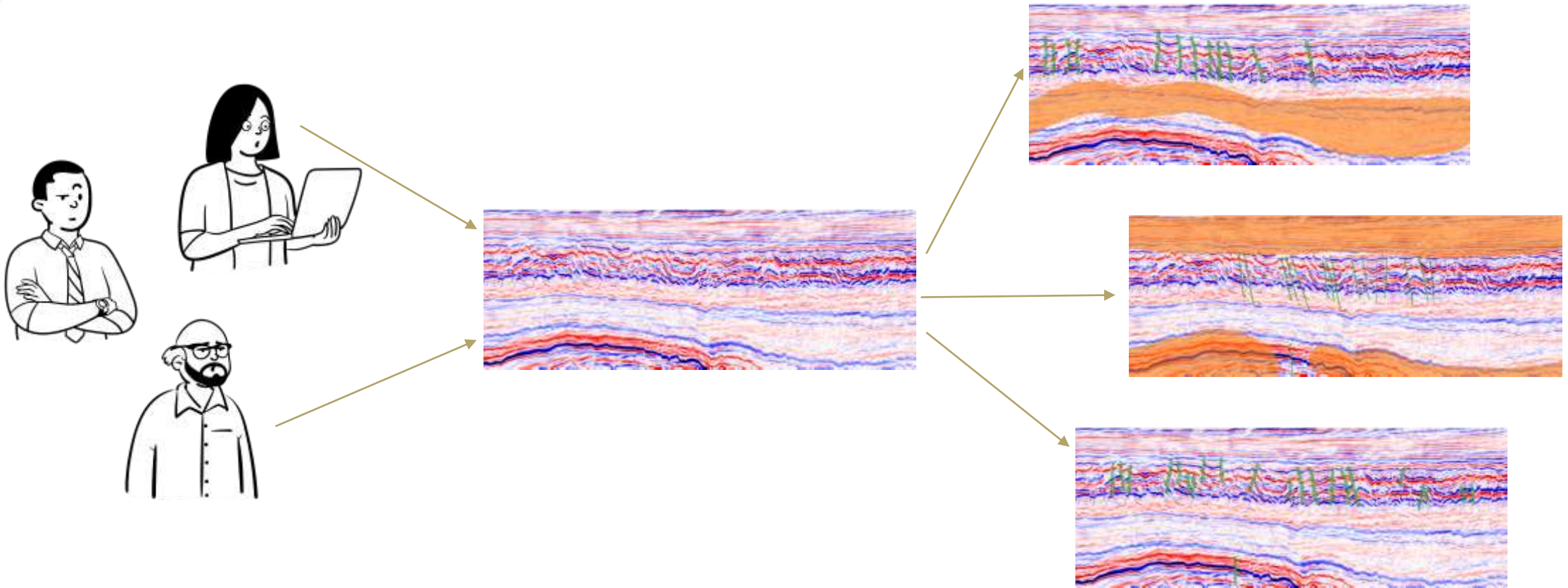
Data
Uncertainty

Model
Uncertainty

Introduction

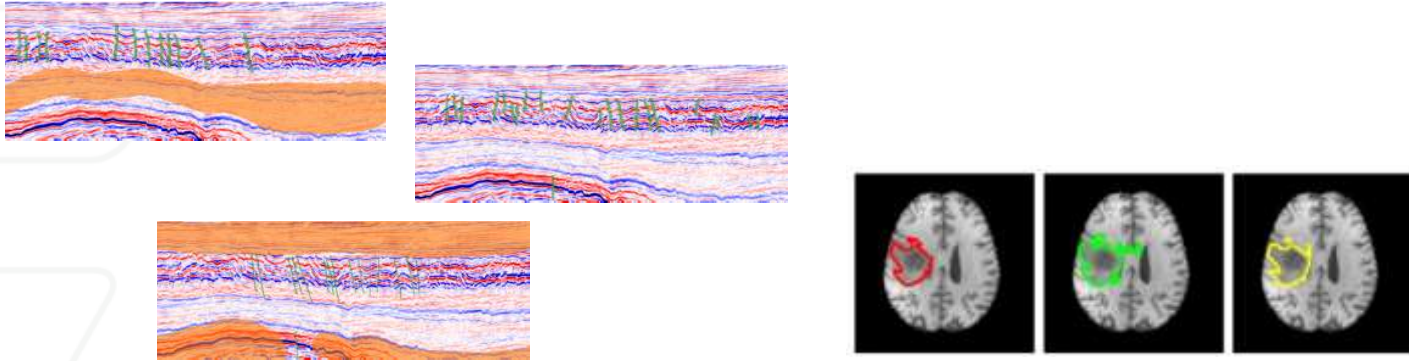
Multiple annotators can disagree during data labeling

An image can yield different labels due to annotator confidence, expertise, etc

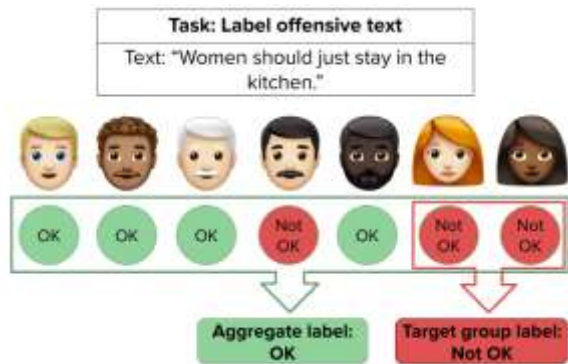


Introduction

When not addressed, annotator disagreement hinders downstream performance



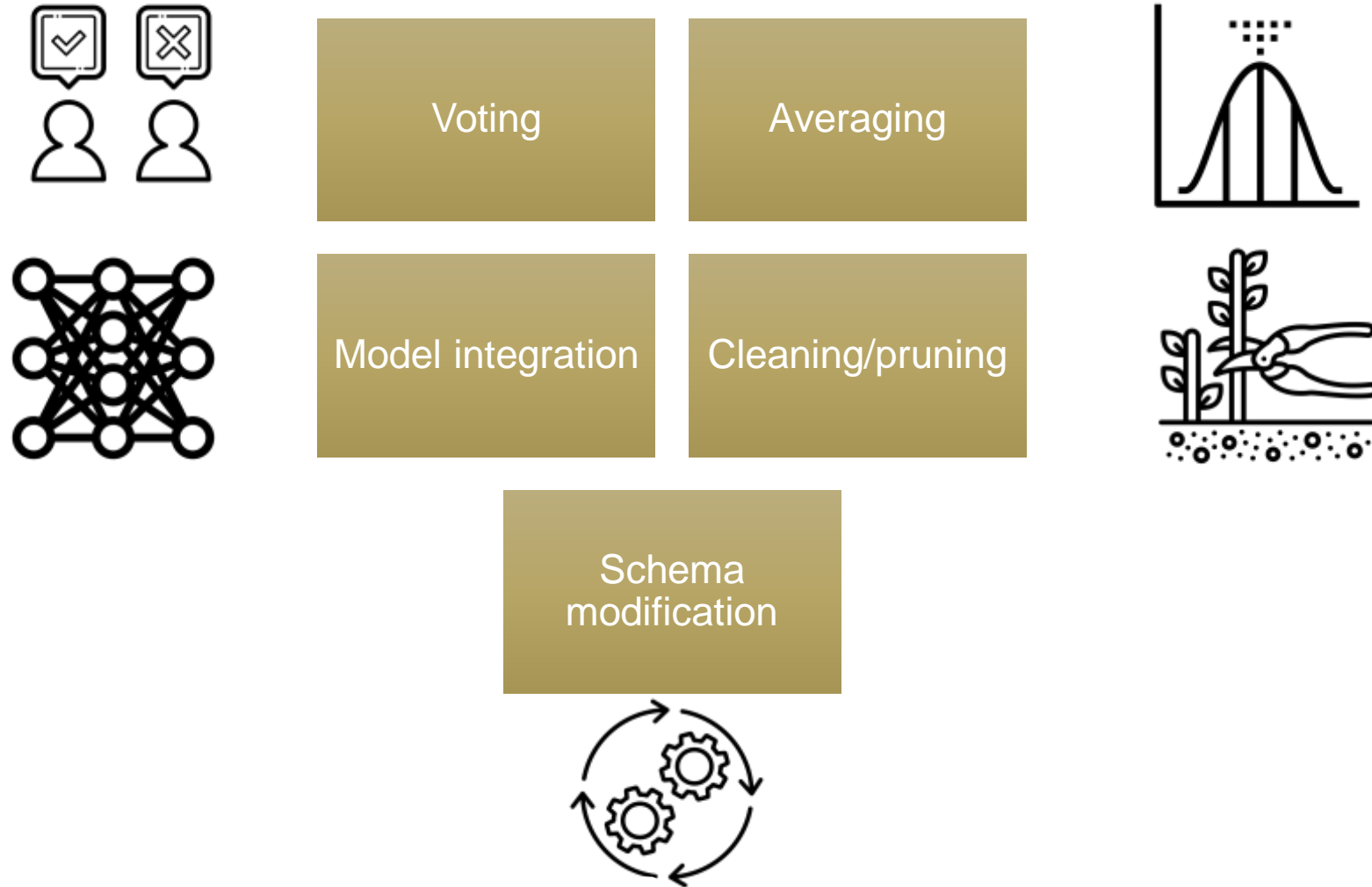
Degraded model/task performance



Noisy labels

Introduction

How is annotator disagreement generally dealt with?

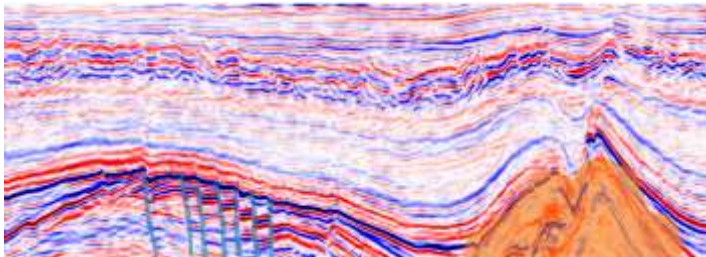
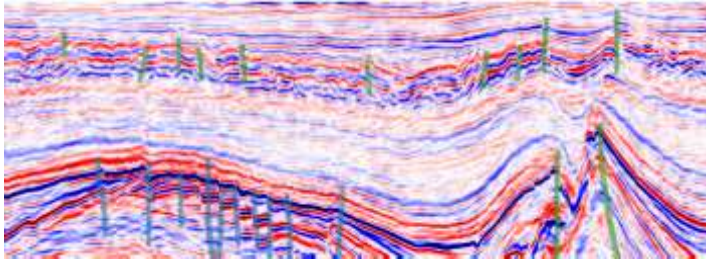


Introduction

Where does annotator disagreement stem from?

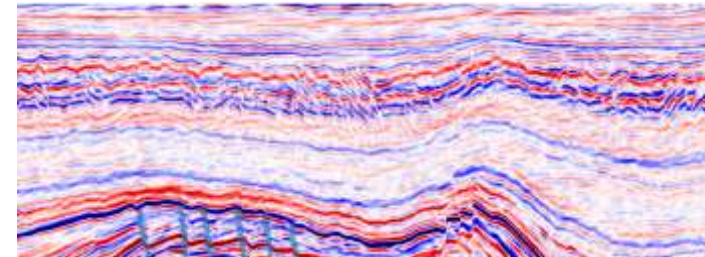
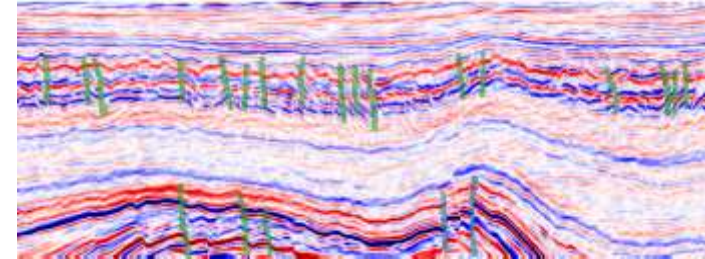
Expertise

How much of a role does domain knowledge play?



Confidence

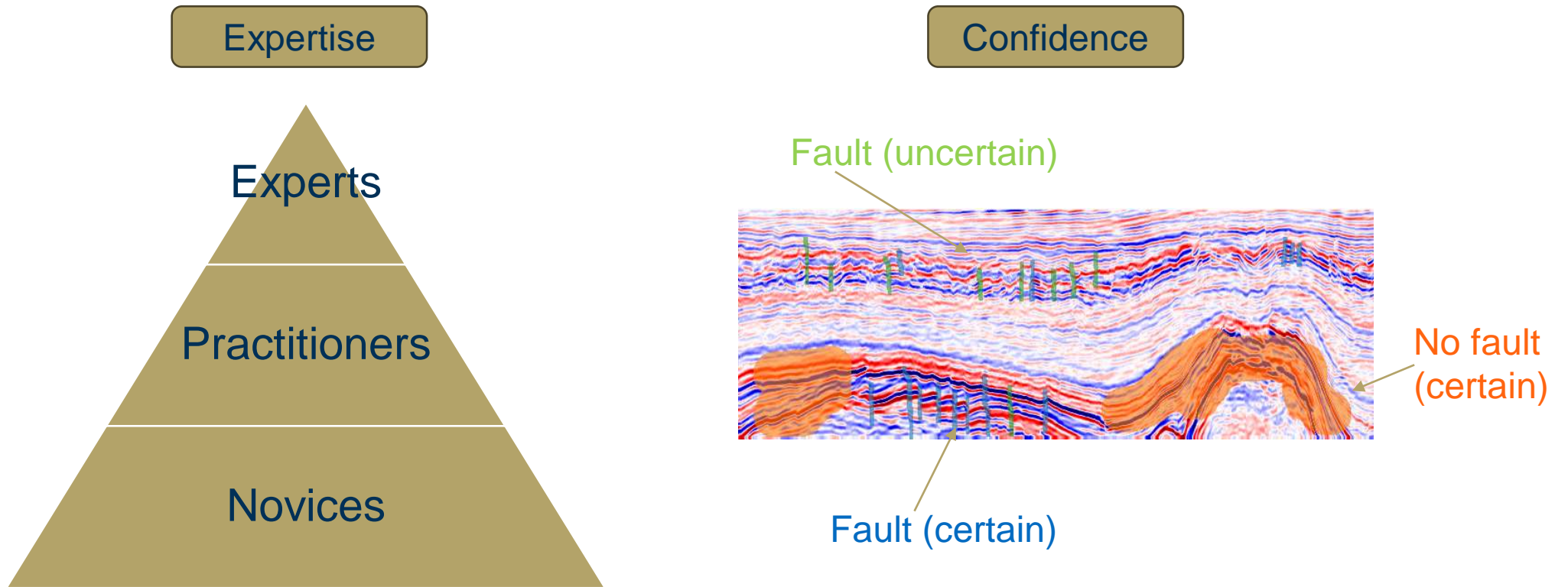
How certain is each labeler of their own annotations?



Introduction

Novelty of Our Work

A dataset comprising labels across multiple levels of **expertise** and **confidence**



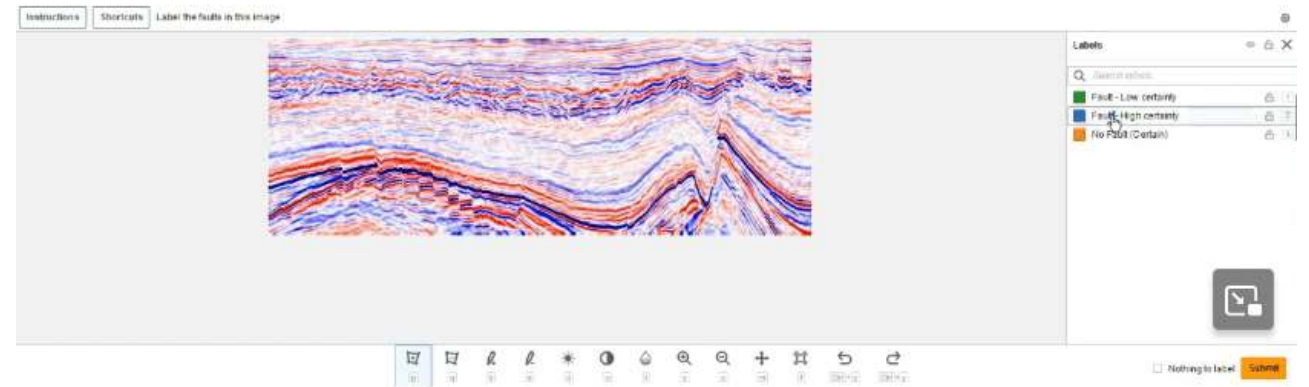
Summary of our dataset

We leverage Amazon Mechanical Turk for the labeling process

- 400 images, divided into 20 batches
- For each batch, 2 images are repeated 3 times for quality assessment
- 2 bonuses:
 - Number of images, promotes full dataset completion
 - Consistency, promotes thorough labeling

Current contribution:

- 1 expert, 8 practitioners



Experimental design - Platform

What is Amazon Mechanical Turk (MTurk)?

MTurk allows for a distributed outsourcing of virtual tasks

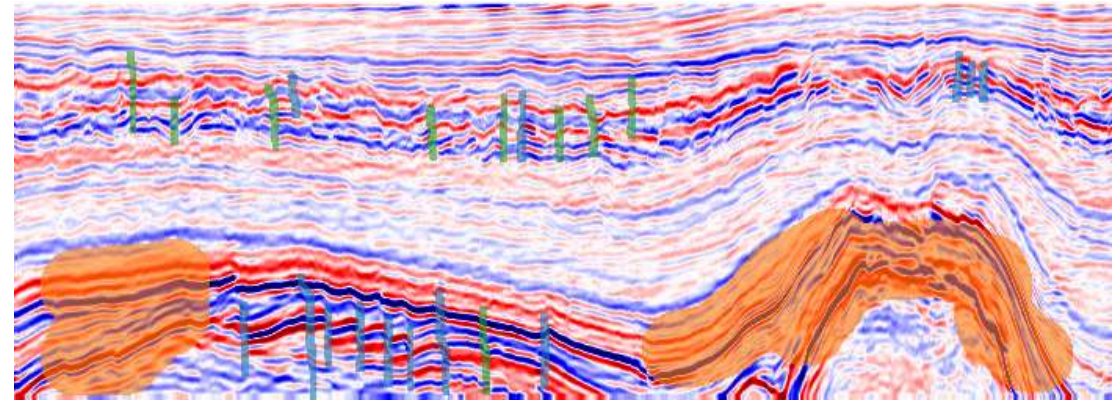
- Crowdsourcing marketplace
- Distributed workforce can perform tasks virtually
- Data validation and research, survey participation, content moderation, etc



Experimental design - Labels

Three label categories enable flexible annotations

- We present 1 image per HIT, and consider 3 fault label categories for different certainty levels:
 - Fault (certain)
 - Fault (Uncertain)
 - No fault (Certain)
- Other schemes considered:
 - Multiple imgs per HIT
 - Sliders and text input boxes for certainty



Experimental design - Batches

Task is divided in 20 batches for easy navigation

- We divide the dataset into 20 batches, with 20 unique images each
- Each batch also contains 3 copies of 2 redundant images for quality assessment, totaling 6 QA images per batch
- Total batch size --> 24 imgs

HIT Groups (1-20 of 685)

Requester	Title	Hits	Amount	Created	Status	Actions
anonat	RYD/WR Feedback	256	\$3.00	1h ago	Pending	Accept & Work
DataScience	SCAuch_2023102408281028101974719_0107waly-1cd1-400a-c28ac2f03ba	1	\$3.00	2h ago	Pending	Accept & Work
Amazon Requester Inc. - RYD Team	Internal testing, please ignore. This is a cycle to test some things we are working on	5	\$3.01	2h ago	Pending	Accept & Work
CS	DriftNetEyes-Test0	1	\$3.01	2h ago	Pending	Accept & Work
Amazon Requester Inc. - Tensile	AFRON Pattern	238	\$3.00	13h ago	Pending	Accept & Work
we-intel	Customer feedback Oct-04-2023	10	\$3.00	25h ago	Pending	Accept & Work
we-intel	US Privacy Oct-04-2023 - Tennessee Department Remediation Request	29	\$3.00	25h ago	Pending	Accept & Work
Amazon Requester Inc. - Tensile	MiniTrainingPositiveTests-2023-10-04-102220	3	\$3.00	27h ago	Pending	Accept & Work
Amazon Requester Inc. - Tensile	MiniTrainingPositiveValidationTests-2023-10-04-155344	3	\$3.00	28h ago	Pending	Accept & Work
altesa.gesuch	Label the geological faults in seismic images 800	21	\$3.00	28h ago	Pending	Accept & Work
altesa.gesuch	Label the geological faults in seismic images 800	21	\$3.00	27h ago	Pending	Accept & Work
CS	DriftNetEyes-Test0	1	\$3.01	31h ago	Pending	Accept & Work
Explorer Heidelberg University	TEST TITLE 0	1	\$1.00	34h ago	Pending	Accept & Work
altesa.gesuch	Label the geological faults in seismic images 700	21	\$3.00	38h ago	Pending	Accept & Work
Explorer Heidelberg University	TEST TITLE 0	4	\$1.00	38h ago	Pending	Accept & Work
altesa.gesuch	Label the geological faults in seismic images 600	21	\$3.00	39h ago	Pending	Accept & Work
altesa.gesuch	Label the geological faults in seismic images 500	21	\$3.00	41h ago	Pending	Accept & Work
altesa.gesuch	Label the geological faults in seismic images 400	21	\$3.00	41h ago	Pending	Accept & Work
ProductPulse	Perform a single human check to unlock ProductPulse HITs: Regions a Mi	1	\$3.10	42h ago	Pending	Accept & Work

Experimental design - Instructions

Concise guidelines aim to allow novices to annotate effectively

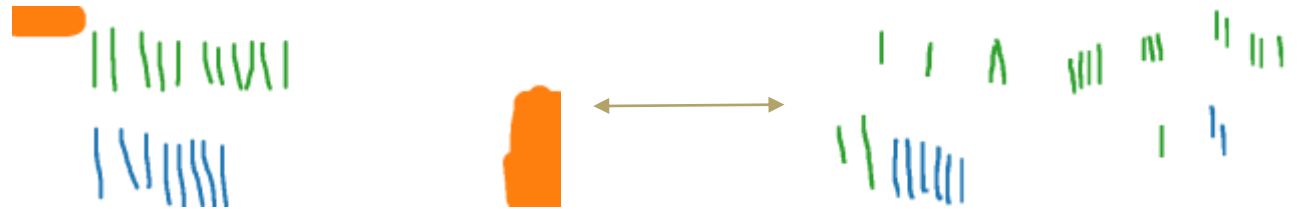
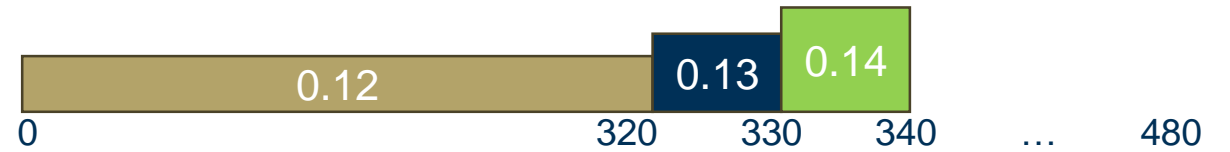
- We provide an instructional video in the task website and inside the layout, with the following details:
 - Fault definition
 - Sample image and label meaning
 - Platform usage
 - Payment scheme



Experimental design - Payment

Concise guidelines aim to allow novices to annotate effectively

- Base pay:
 - Reviewed existing mturk tasks and literature to arrive at initial pay rate
 - Adjusted pay rate for average time after internal task completion
- Number of image bonus:
 - Prorated bonus applied to final section of dataset to motivate completion
- Consistency bonus
 - Internal self-agreement metric for quality assessment

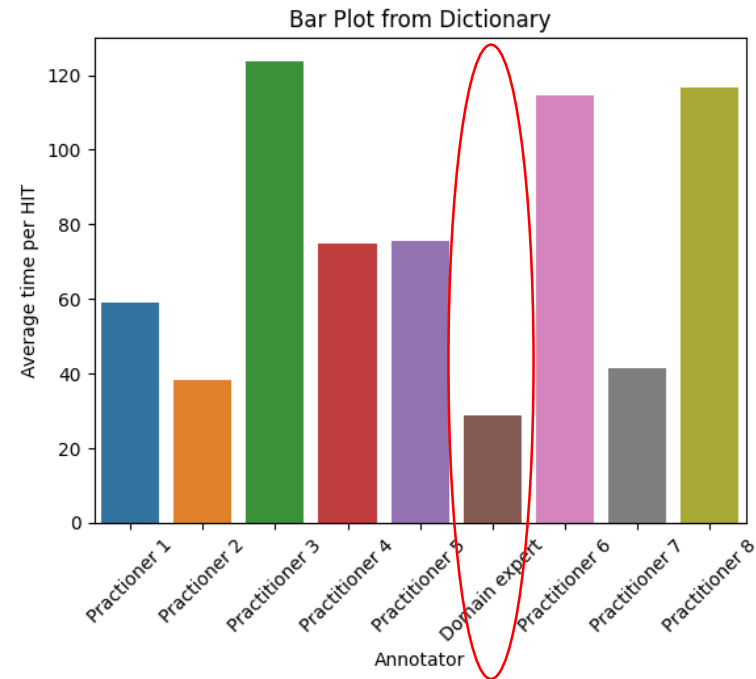


Insights from current annotators - Speed

Domain expert is generally faster than intermediate users

Expert is significantly faster than intermediate users

- Average labeling time: 75s
- Expert average time: 30s

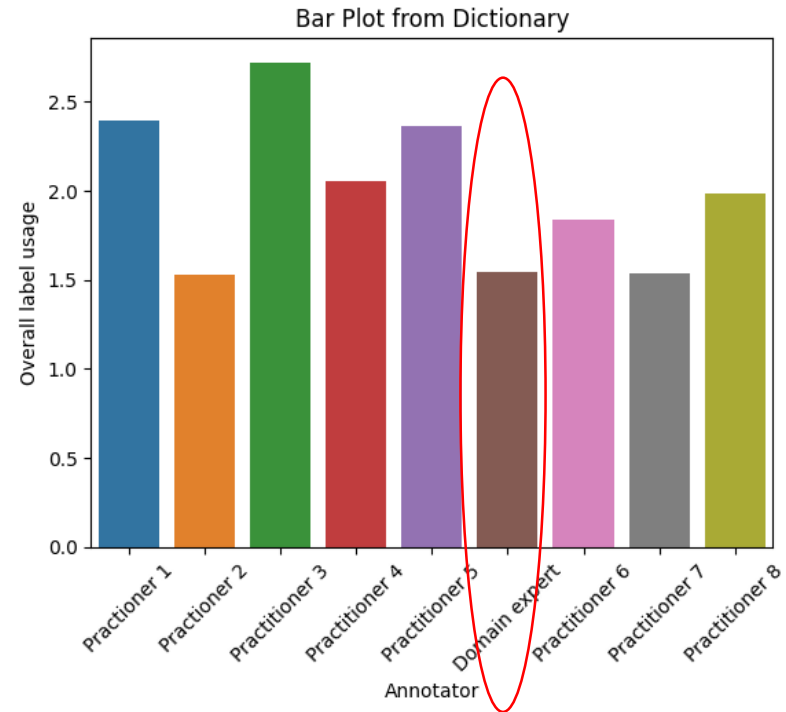


Insights from current annotators - Label usage

Annotators use less labels over time

Most people use only two labels or less (on average)

- Most people use only two labels on average
- Expert uses generally less labels than practitioners

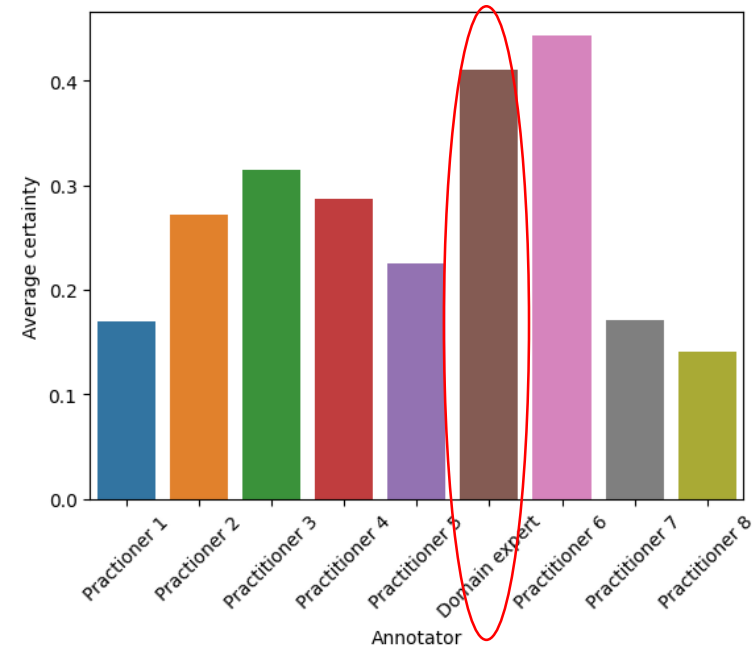


Insights from current annotators - Certainty

Confidence oscillates throughout task

Annotators with the most exposure make more confident labeling

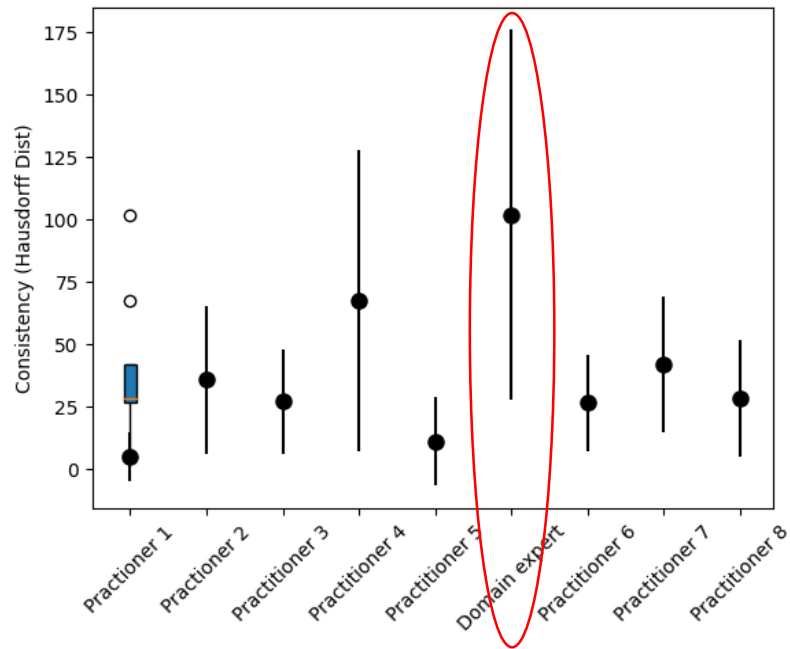
- Confidence taken only for the 2 fault labels
- Most experienced annotators label more confidently



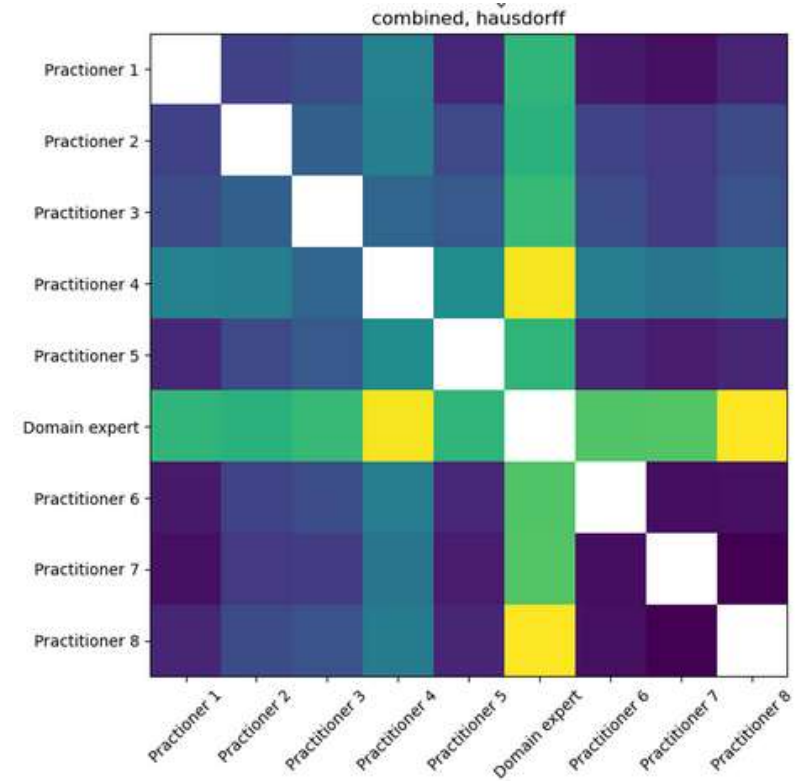
Agreement measurements

Intra and inter annotator agreement varies significantly

Expert is very distanced from other annotators as well as himself



Intra-annotator disagreement



Inter-annotator disagreement

Conclusion

Cross-expertise and multi-confidence dataset

- Uncertainty insights can help machine learning community build better models and methodologies to account for annotator disagreement
- A better understanding of the expertise gap can lead to more efficient fault labeling pipelines, reduced expert workload, and better fault detection models
- Our experimental design also constitutes a valuable resource for the seismic community to harness crowdsourcing platforms for efficient data labeling and annotation

For more OLIVES content,
please visit:

GitHub



Publications

