

Explaining and building trust in machine learning: the local and subgroup perspectives

Eliana Pastor

Seminar @ sqlRL/IDLab, University of Antwerp

Hello!

I'm Eliana!

I'm an assistant professor at Politecnico di Torino, Italy

I work on Trustworthy and Explainable AI

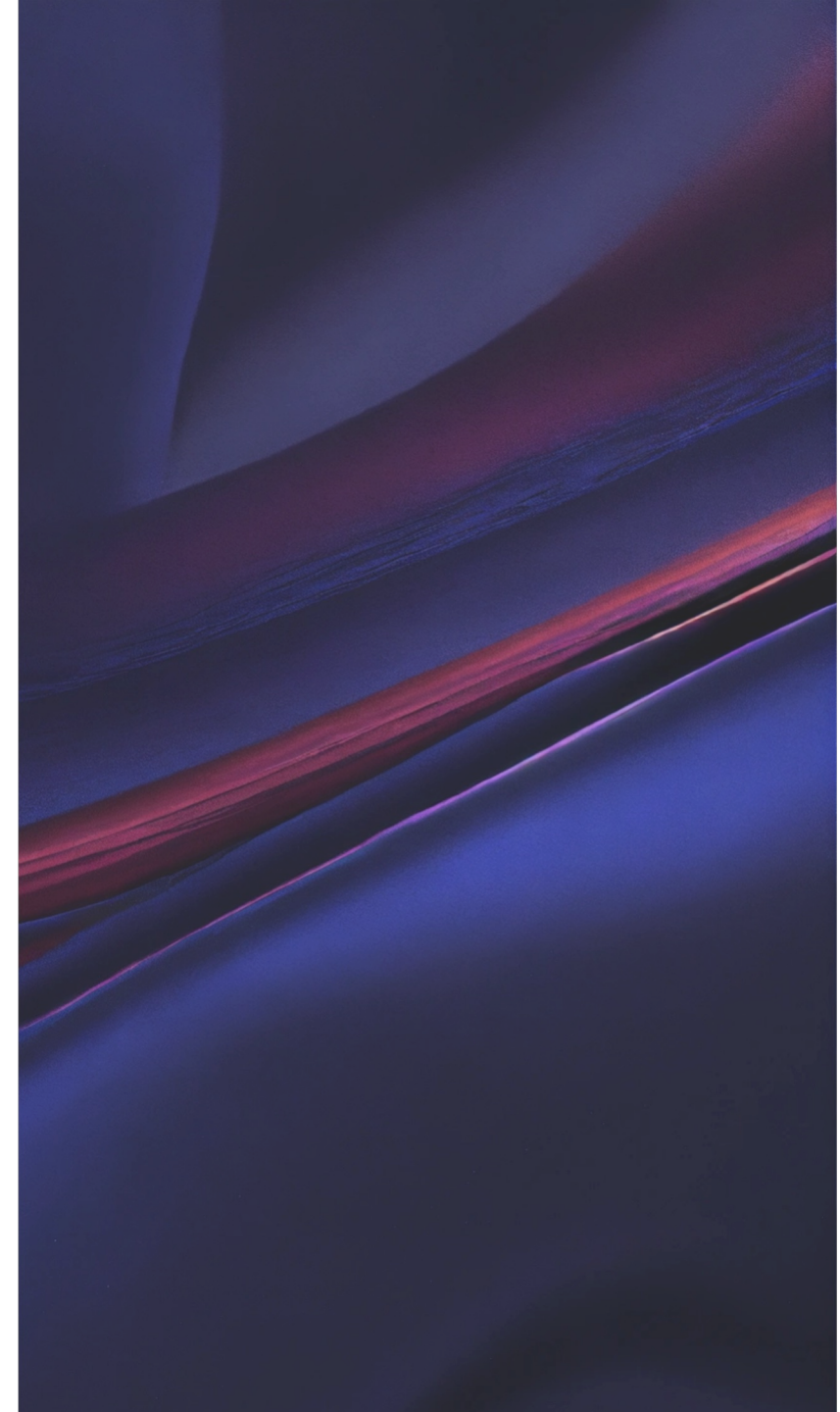


My interests.. in keywords

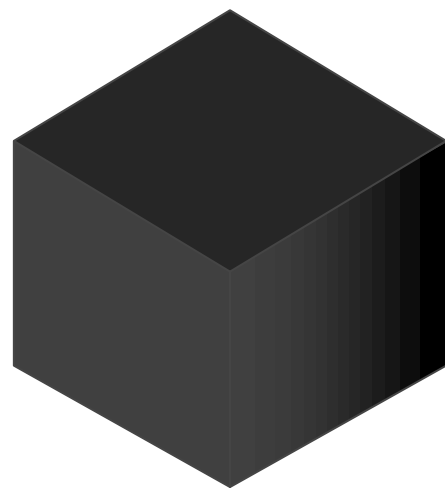
Trustworthy AI

Explainable AI, Fairness in AI, Robustness, Debugging

- Analysis of disparities in data subgroups
 - XAI for Speech & Sound
 - Post-hoc XAI for tabular & text data
 - KANs
 - Concept-based XAI
- & other stuff

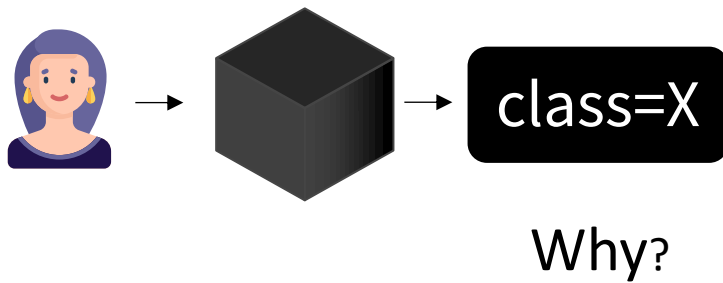


Our problem: open the box

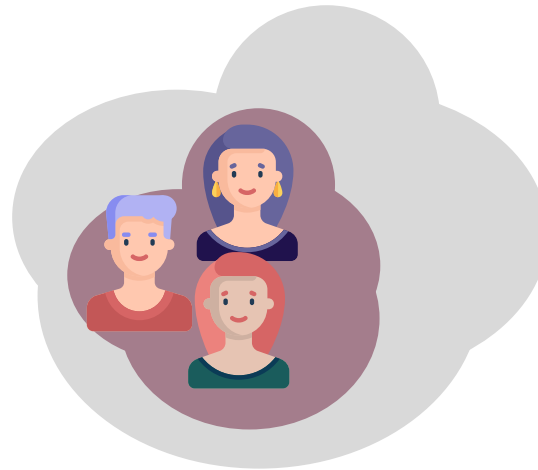


From which perspective we open the black box

Local perspective



Subgroup perspective



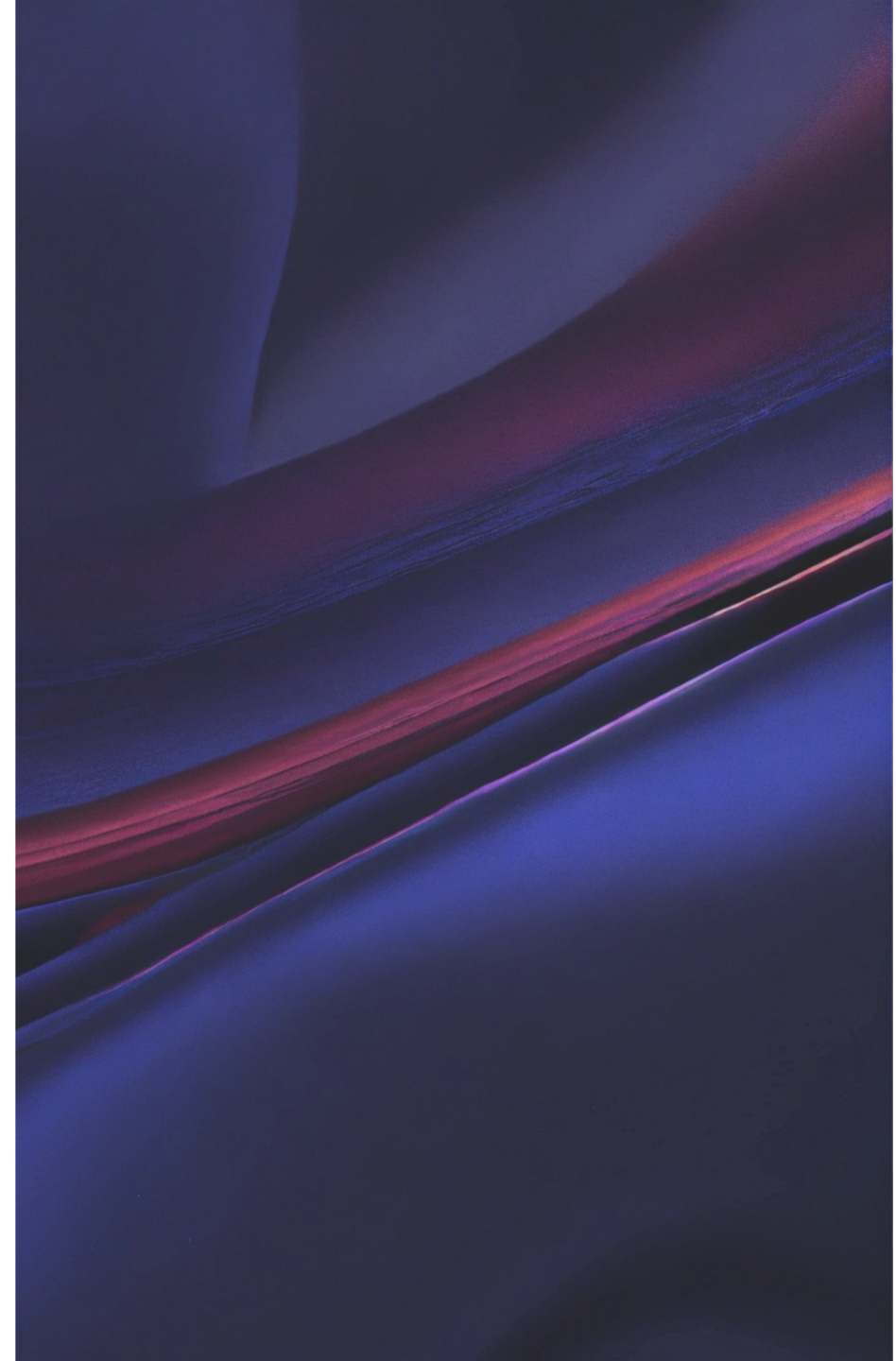
Global perspective



Outline

Subgroup perspective

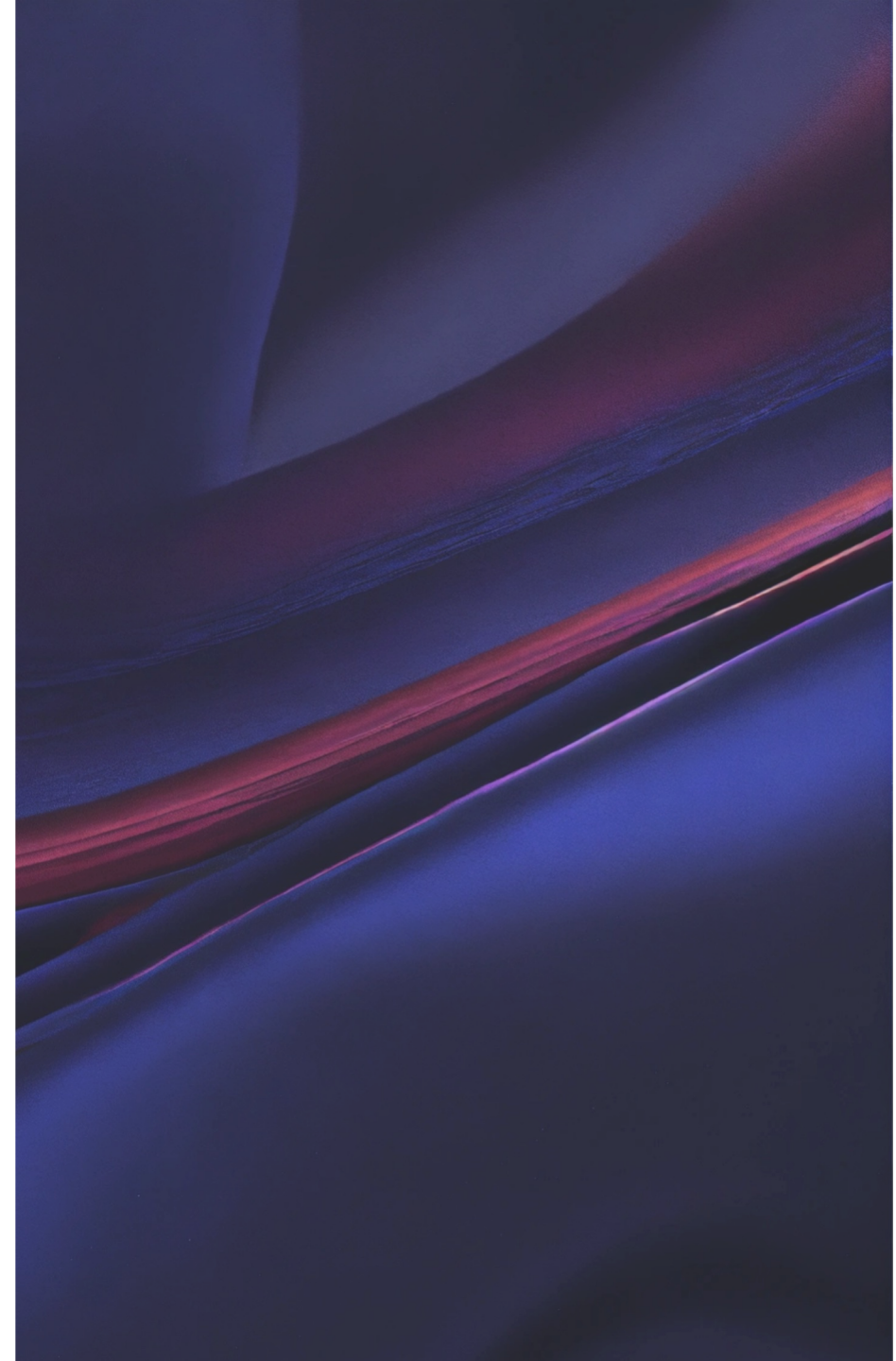
- Identification of subgroups with *divergent* classification behavior
- Divergent subgroup analysis in speech data
- Subgroup-based model comparison
- Mitigate subgroup disparities
- Interpretable subgroup drift detection



Outline

Local perspective

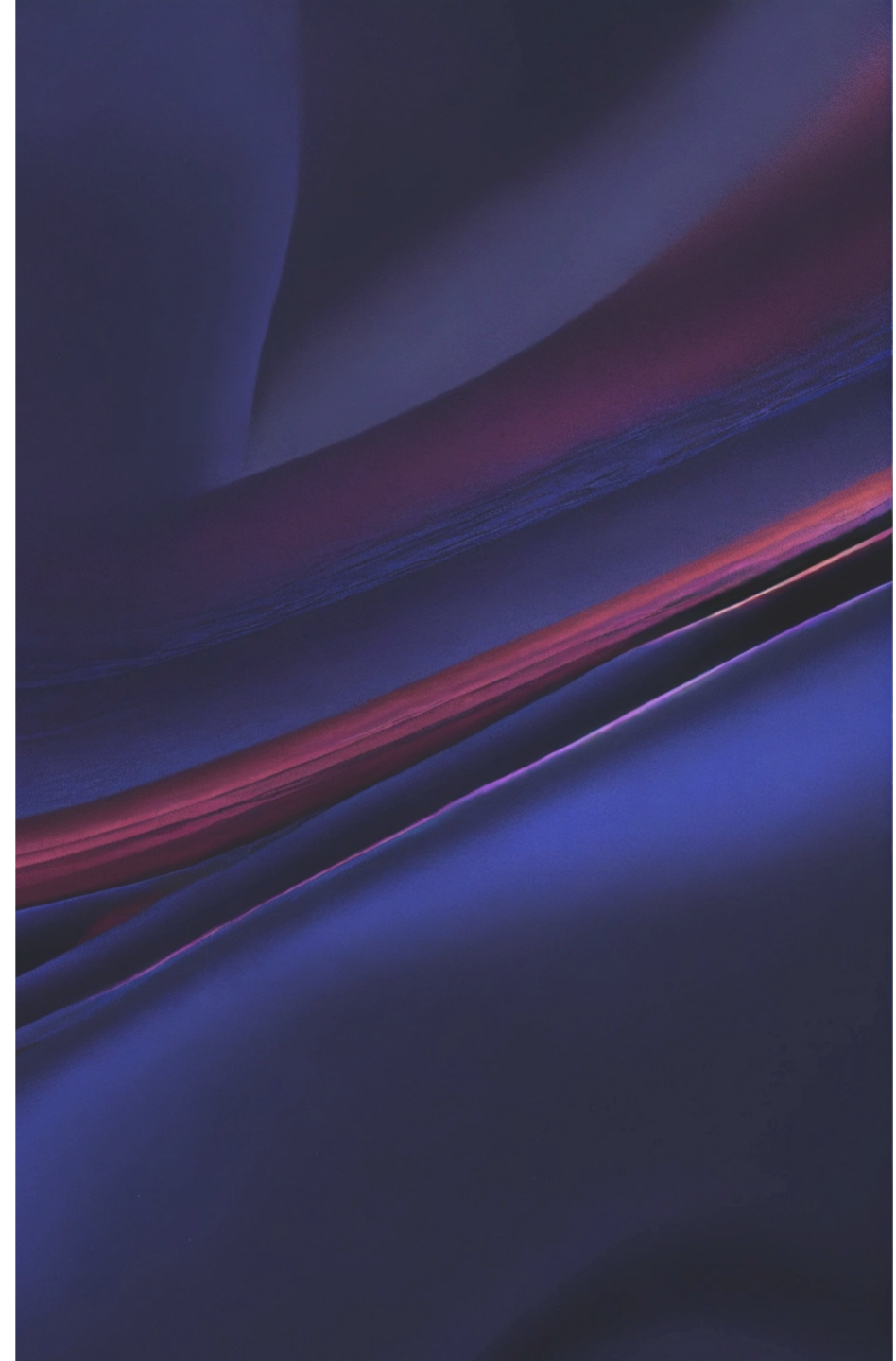
- Explaining prediction of speech models
- Assessing explainability methods for transformers models

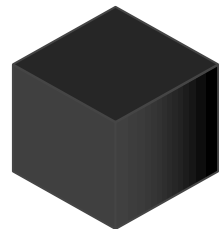


Outline

Subgroup perspective

- **Identification of subgroups with *divergent* classification behavior**
- Divergent subgroup analysis in speech data
- Subgroup-based model comparison
- Mitigate subgroup disparities
- Interpretable subgroup drift detection





PERFORMANCE

X%



HIGH ERROR RATE

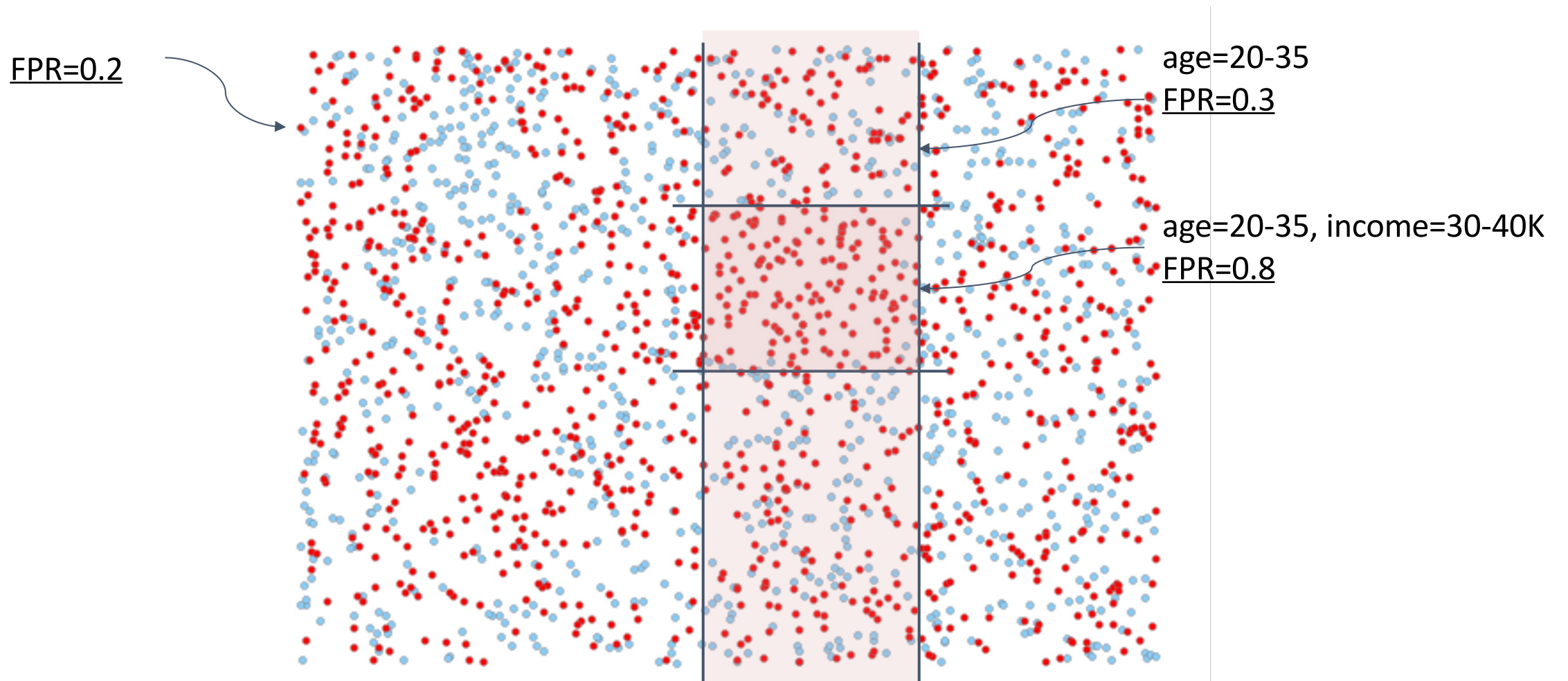


LOW ERROR RATE



AS OVERALL
MODEL BEHAVIOR

Motivation



Divergence of a subgroup

pattern, interpretable, e.g., {age=20-35, gender=female}

$$\Delta(S) = f(S) - f(D)$$

performance measure

all dataset

Generic & model agnostic

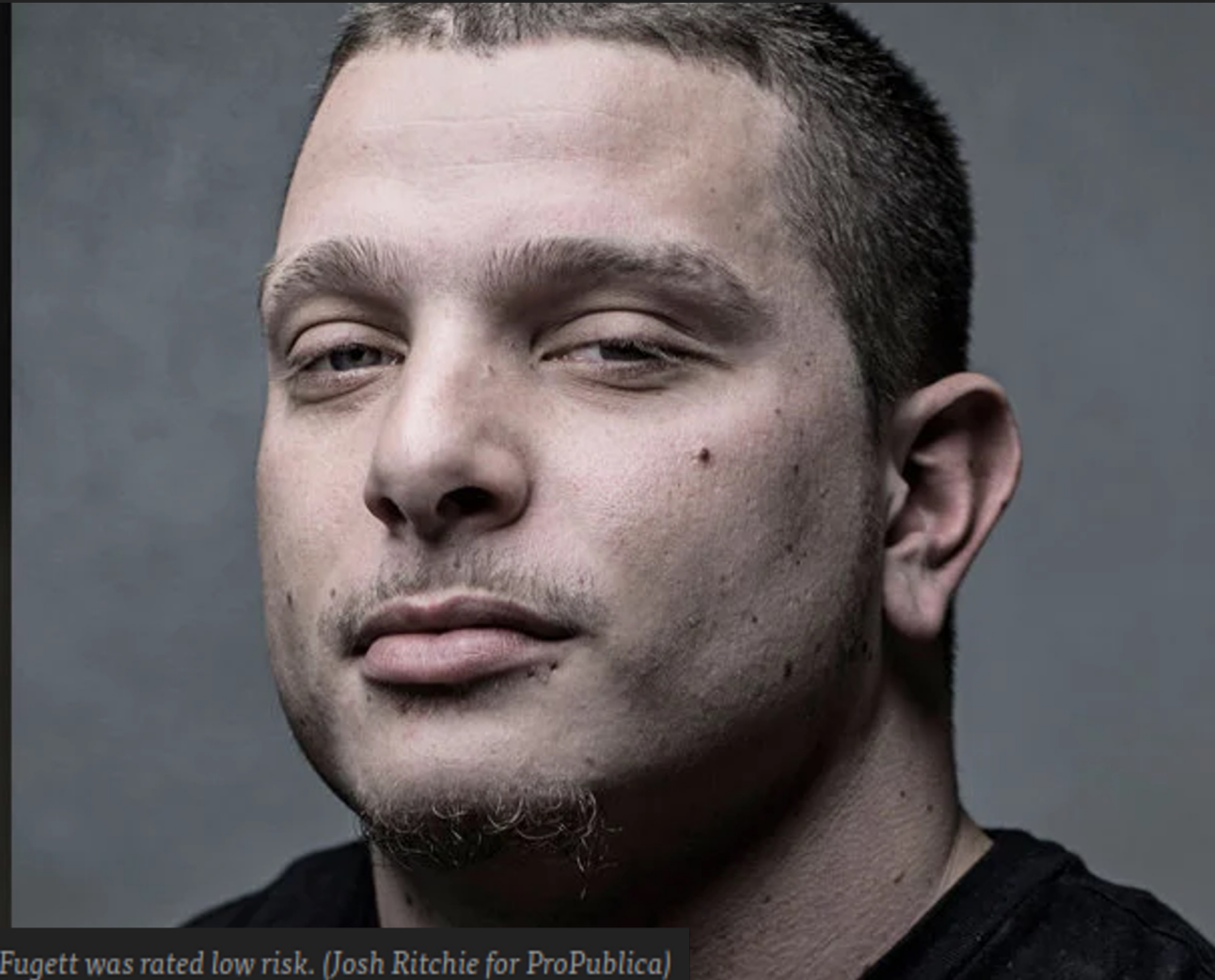
➤ **Automatic** identification of subgroups via **frequent pattern mining**

Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica

May 23, 2016



Bernard Parker, left, was rated high risk; Dylan Fugett was rated low risk. (Josh Ritchie for ProPublica)

Divergent subgroups - Example

COMPAS dataset. Recidivism predictions based on defendant information

		Divergence	Statistical significance	
	Itemset	Δ_{FPR}	support	t
Subgroup	age=25-45, #prior>3, race=Afr-Am, sex=Male	0.22	0.13	7.1
	age=25-45, #prior>3, race=Afr-Am	0.211	0.15	7.4
	age=25-45, charge=F, #prior>3, race=Afr-Am	0.202	0.11	6.2

Subgroup frequency

Contributions of items to divergence


Itemset	Δ_{FPR}
age=25-45, #prior>3, race=Afr-Am, sex=Male	0.22

What is the contribution
of each term?

Contributions of items to divergence



Team
&
Team score
↓
Subgroup I
&
Subgroup
divergence

 ?
→
Contribution
↓
Contribution
of an
attribute=value
(item) α

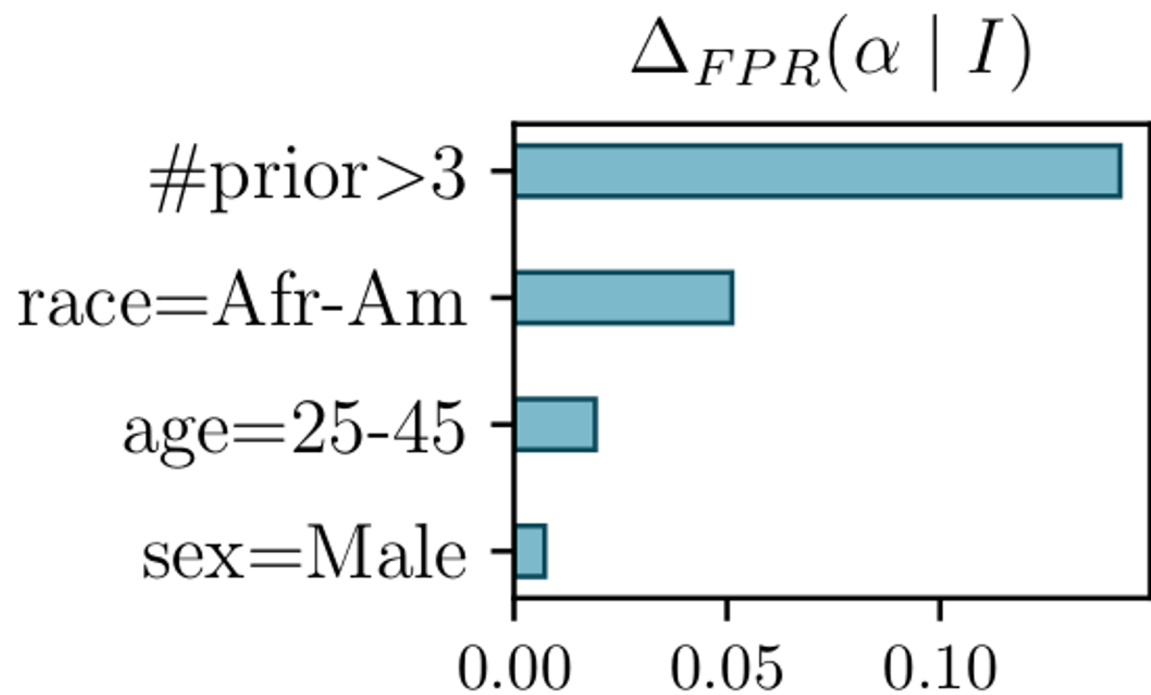
Shapley value
Given the score of subset of
players

Contribution of item α in I:

$$\Delta(\alpha | I) = \sum_{J \subseteq I \setminus \{\alpha\}} \frac{|J|!(|I| - |J| - 1)!}{|I|!} [\Delta(J \cup \alpha) - \Delta(J)]$$

Contributions of items to divergence

Itemset	Δ_{FPR}
age=25-45, #prior>3, race=Afr-Am, sex=Male	0.22



Divergent subgroups - Example

Itemset	Δ_{FPR}	support	t
age=25-45, #prior>3, race=Afr-Am, sex=Male	0.22	0.13	7.1
age=25-45, #prior>3, race=Afr-Am	0.211	0.15	7.4
age=25-45, charge=F, #prior>3, race=Afr-Am	0.202	0.11	6.2

Globally?

Global divergence

Global Shapley Value

A generalization of Shapley value that accounts for:

- Incompatible items (e.g. $\{age < 25, age > 45\}$)
- Minimum support threshold

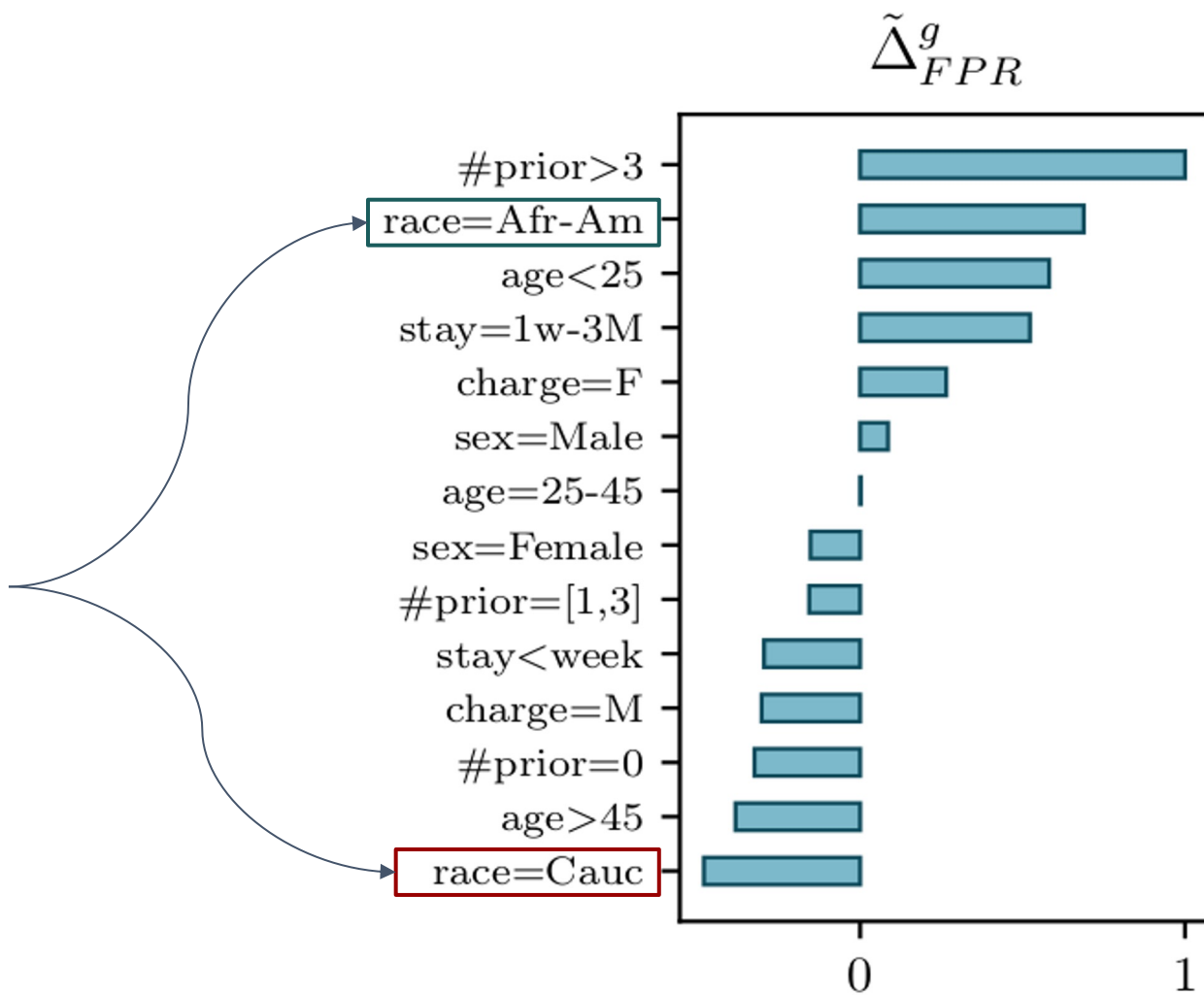
$$\tilde{\Delta}^g(I, s) = \sum_{B \subseteq A \setminus \text{attr}(I)} \frac{|B|! (|A| - |B| - |I|)!}{|A|! \prod_{b \in B \cup \text{attr}(I)} m_b} \sum_{J: J \cup I \in \mathcal{I}_{B \cup \text{attr}(I)}^*} [\Delta(J \cup I) - \Delta(J)]$$

normalization factor, where m_b
is # attribute values of b

set of frequent itemsets
with attributes $B \cup \text{attr}(I)$

Global divergence - COMPAS

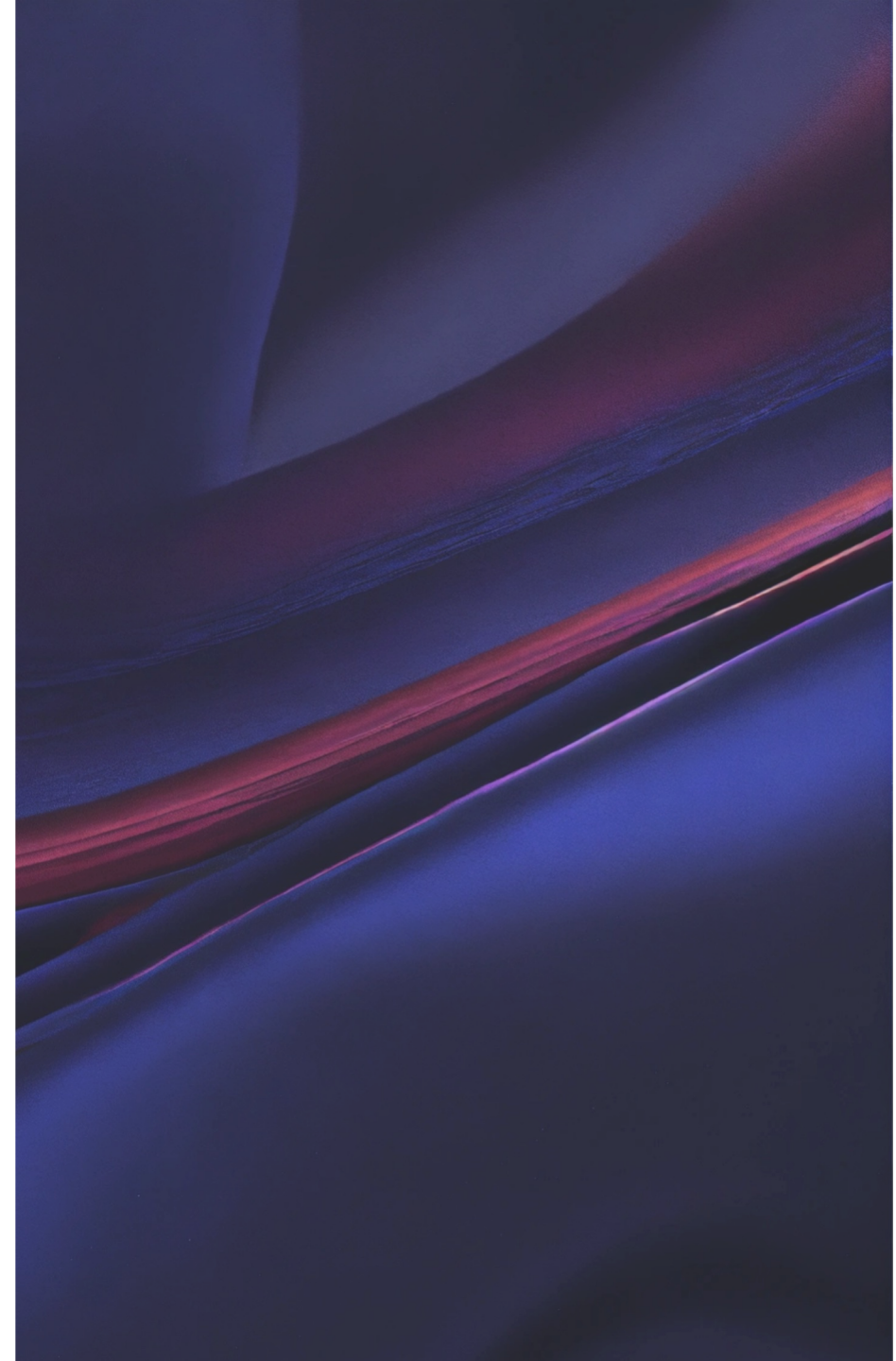
Role of ethnicity



Outline

Subgroup perspective

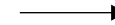
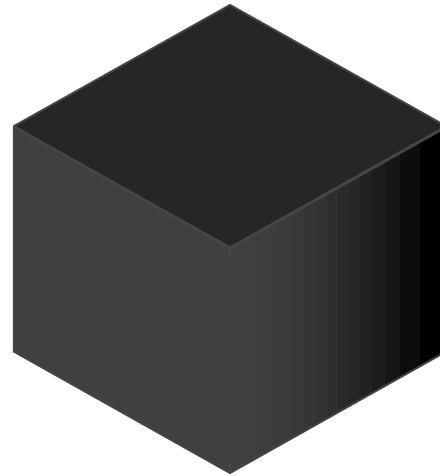
- Identification of subgroups with *divergent* classification behavior
- **Divergent subgroup analysis in speech data**
- Subgroup-based model comparison
- Mitigate subgroup disparities
- Interpretable subgroup drift detection



Our scenario



Turn on the kitchen lights



Action: activate

Object: lights

Location: kitchen

Desidered properties of a subgroup

Interpretable

- e.g., lower performance for *young women*

Adequately represented

- Statistically and operational significant

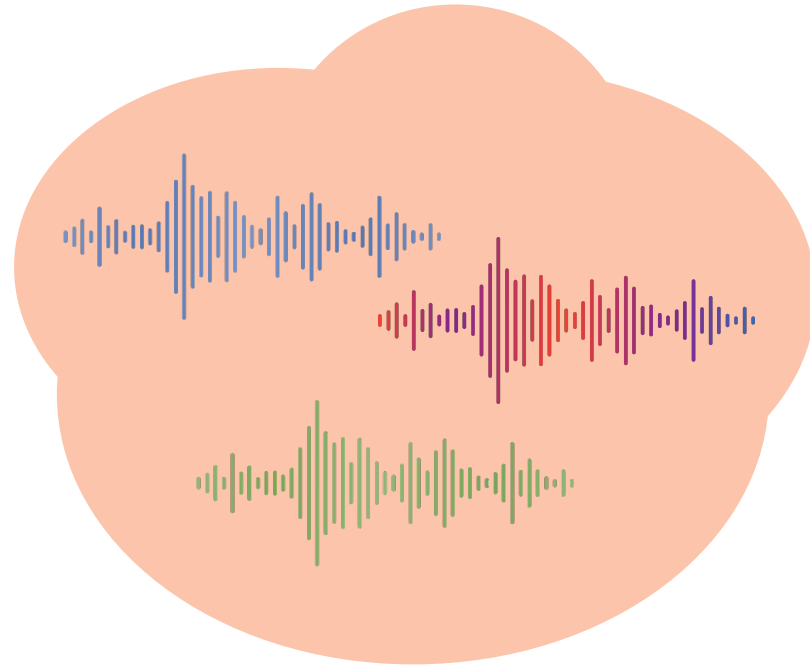
Highlighting peculiar behavior





How to make an
interpretable data
grouping?

Clustering?



But... clusters of utterances are not directly interpretable

Enhance utterance with interpretable metadata

**Speaker
demographics**

**Speaking and
recording conditions**

**Task- or dataset
specific features**



Metadata

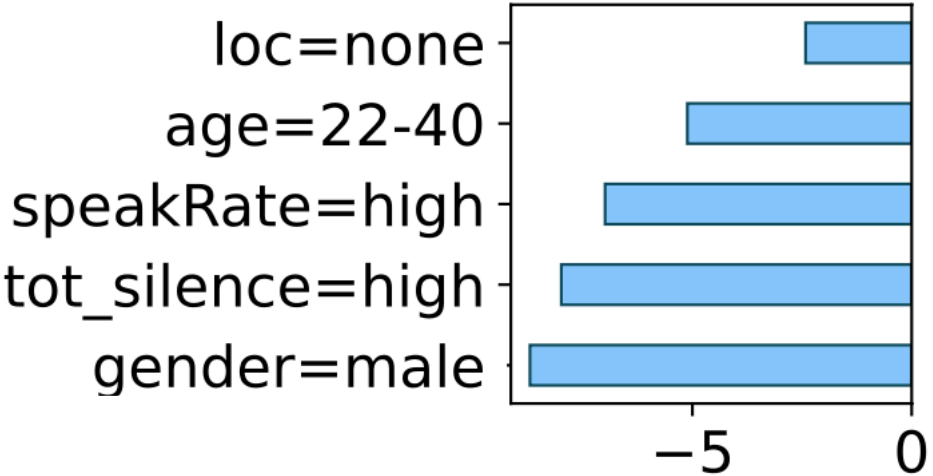
gender=female
country=Italian
noise-level=high
speaking rate=fast
...

Divergent subgroup

By 31.22 less accurate!

	<i>Subgroups</i>	<i>f</i>	Δ_f
I^-	{age=22-40, gender=male, location=none, speaking rate=high, tot silence=high}	60.50	-31.22
I^+	{age=22-40, location=washroom, speaking rate=low, trimmed duration=high}	100.0	8.28

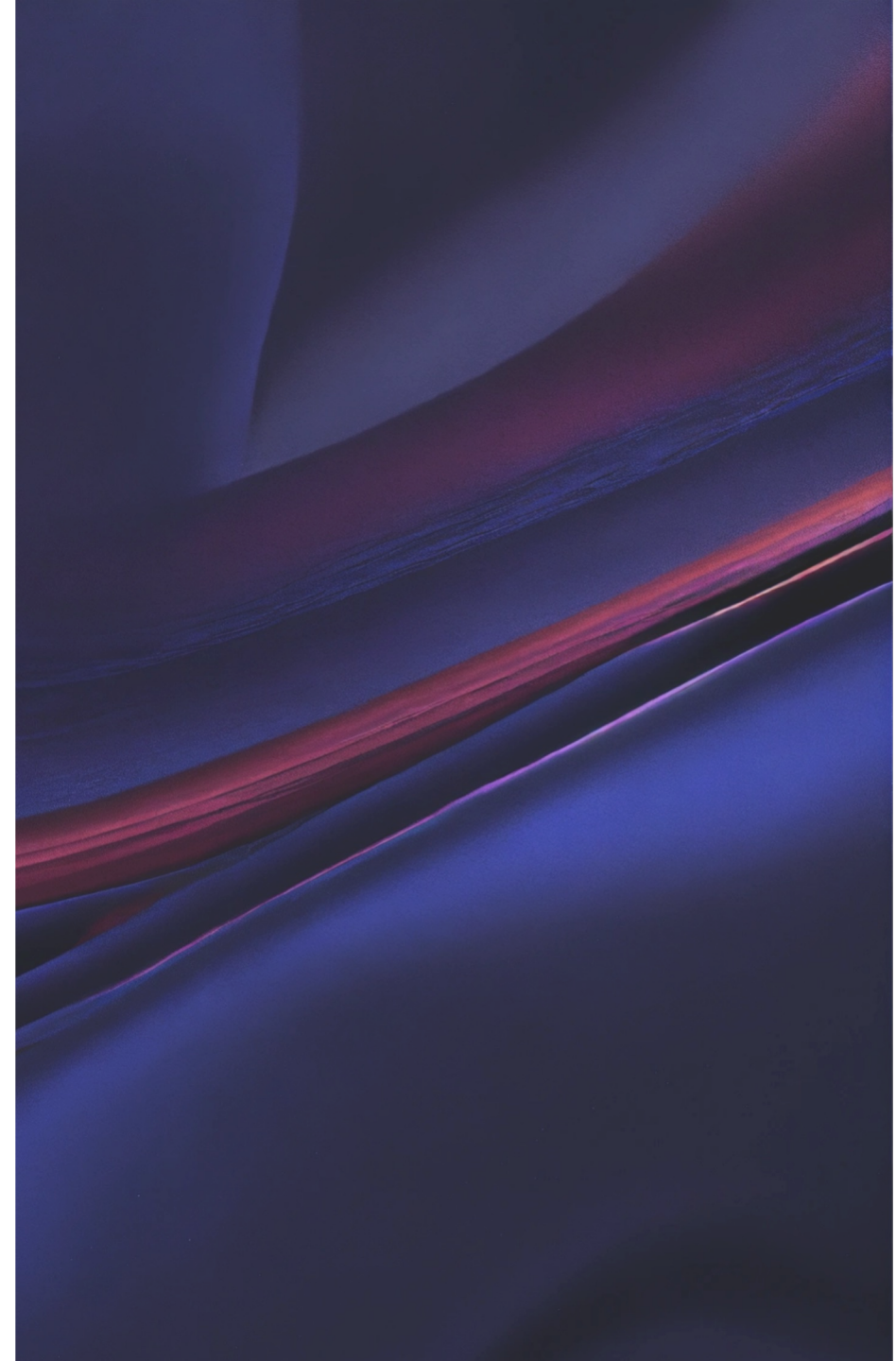
More accurate than average

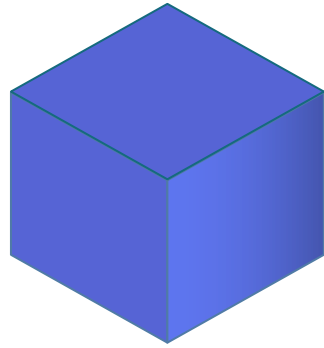


Outline

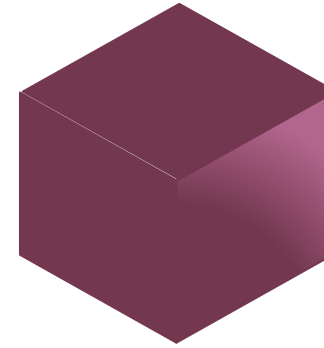
Subgroup perspective

- Identification of subgroups with *divergent* classification behavior
- Divergent subgroup analysis in speech data
- **Subgroup-based model comparison**
- Mitigate subgroup disparities
- Interpretable subgroup drift detection





Accuracy 91.72%



Accuracy 93.17%

Which model to choose?

.. most accurate..?

But on subgroups?

Inter-model performance gap

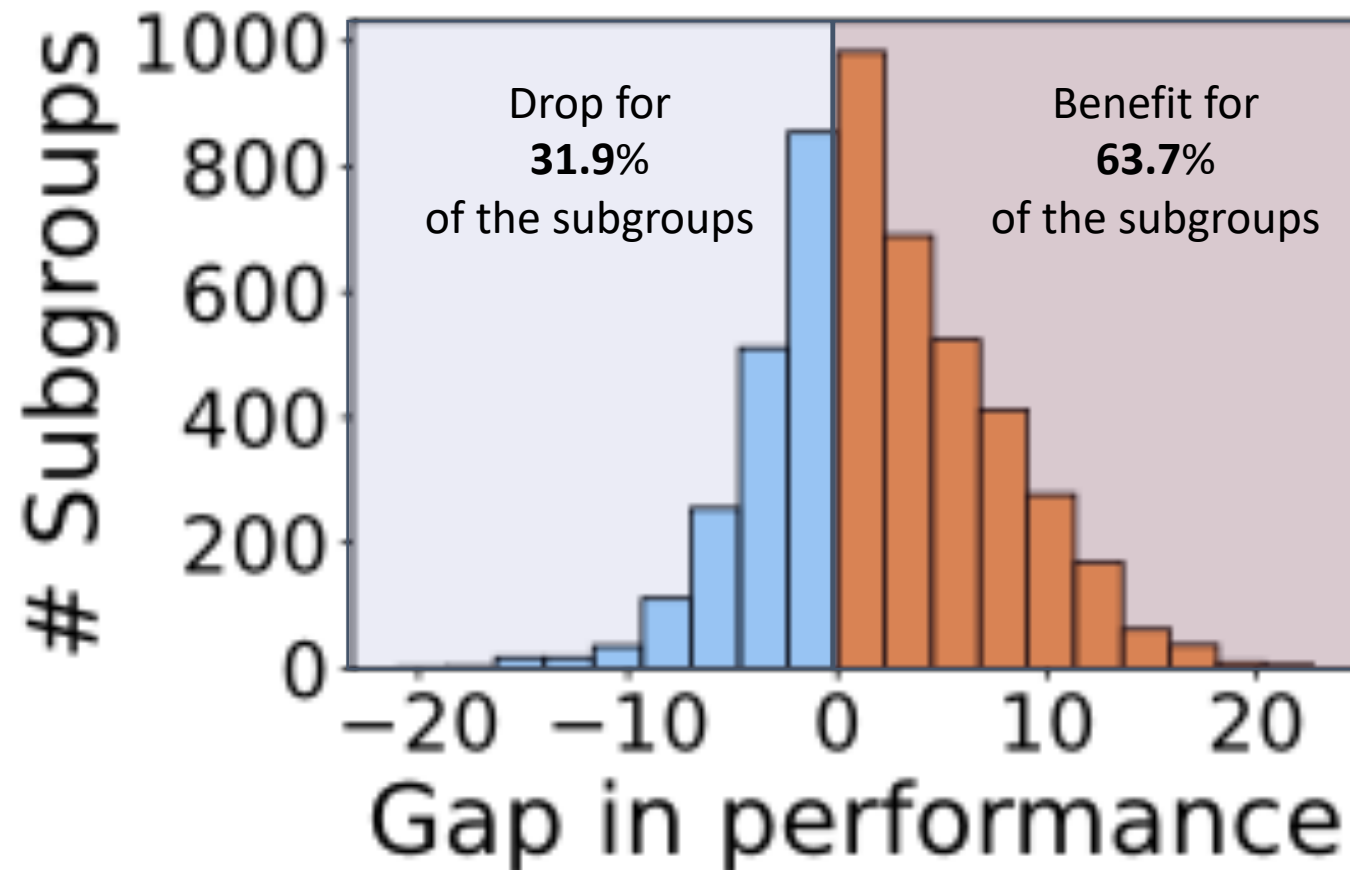
S = **pattern**, e.g., {age=20-35, gender=female}

$$gap_f(S, M_1, M_2) = f(S, M_2) - f(S, M_1)$$

performance on S of model M_2

performance on S of model M_1

Distribution of gain in performance



An example

<i>Subgroups</i>	<i>gap_f</i>	<i>f_{w2v2-b}</i>	<i>f_{w2v2-1}</i>
↑ { <i>action=increase, location=none, tot duration=low, trimmed speaking rate=low, trimmed duration=low</i> }	22.69	75.63	98.32
↓ { <i>action=activate, gender=male, speaking rate=low</i> }	-20.97	96.77	75.81

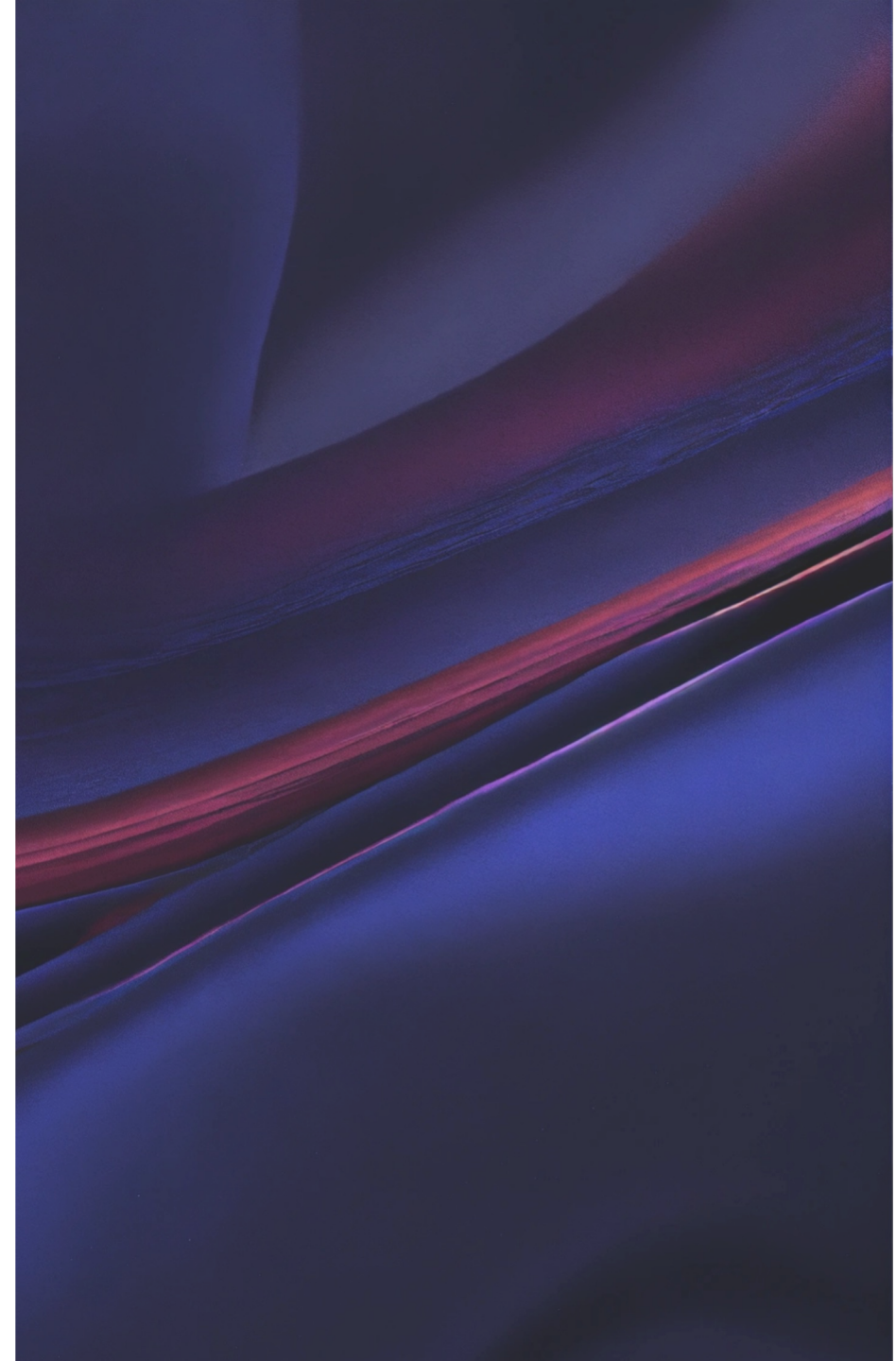
Increase in performance

Drop in performance

Outline

Subgroup perspective

- Identification of subgroups with *divergent* classification behavior
- Divergent subgroup analysis in speech data
- Subgroup-based model comparison
- **Mitigate subgroup disparities**
- Interpretable subgroup drift detection



From identification to mitigation

Once we identify **divergent patterns**.. actively operate on **mitigation**

Post-processing

- Subgroup-guided data acquisition

In-processing

- Divergence regularization
- Subgroup-based contrastive loss
- Targeted data augmentation

Post-processing Subgroup-guided data acquisition

Speaking rate=high, gender=male



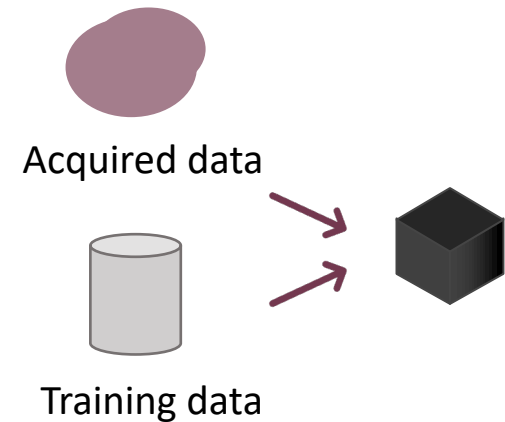
Step 1.

Identify the divergent patterns



Step 2.

Acquire data satisfying the patterns



Step 3.

Speech model re-training

In-processing Divergence regularization

Add a **divergence regularization term**

$$\mathcal{L}_{\Delta} = \sum_{x_i \in D} \max_{S \in \mathcal{S}(x_i)} |\Delta(S)| \mathcal{L}_{CE}(y_i, \hat{y}_i)$$

where $\mathcal{S}(x_i)$ is the set of subgroups satisfied by an instance x_i and \mathcal{L}_{CE} is the cross-entropy loss,

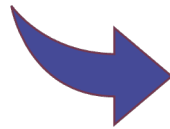
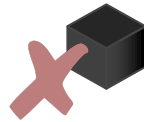
Higher weight for samples with high-divergence

In-processing Targeted data augmentation

Speaking rate=high, gender=male

Step 1.

Identify the divergent patterns



Step 2.

Data augmentation

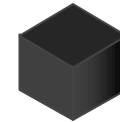
Time stretching,
background noise
injection, reverberation,
pitch shifting



Augmented data



Training data

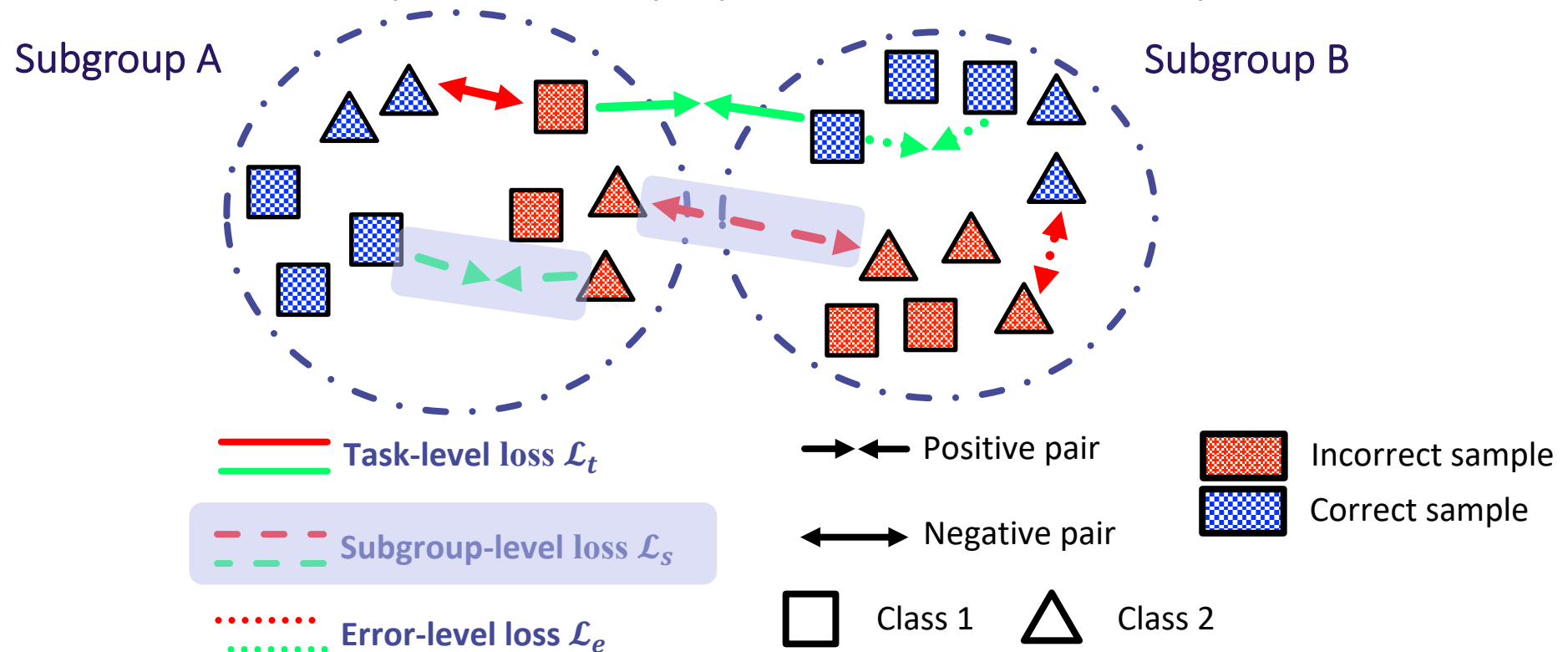


Step 3.

Speech model training

In-processing Subgroup-based contrastive training

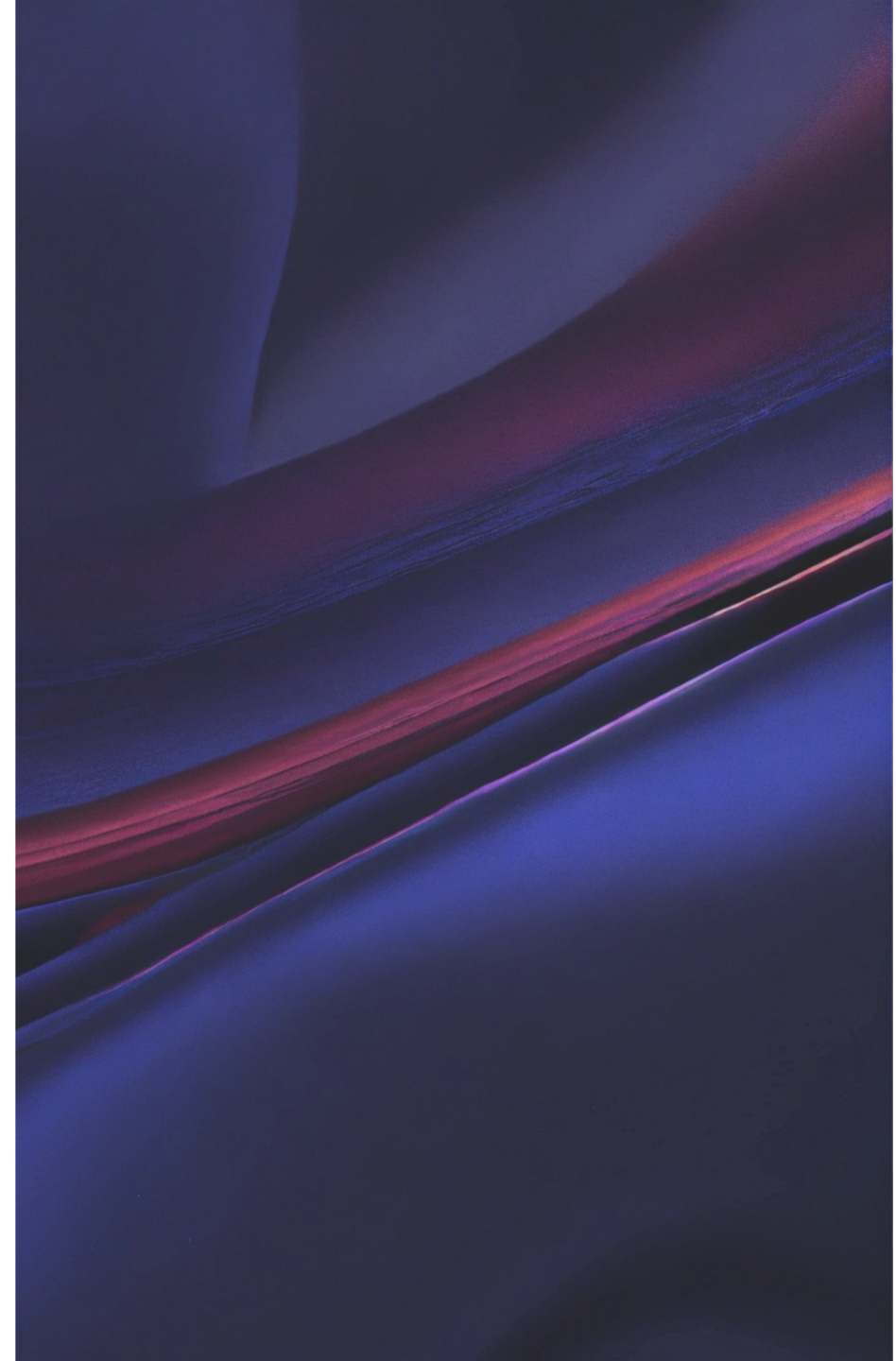
Three separate contrastive learning levels: **task**, **subgroup**, and **error**. At each level, we employ a multi-similarity (MS) loss to selectively contrast sample pairs based on their affinity



Outline

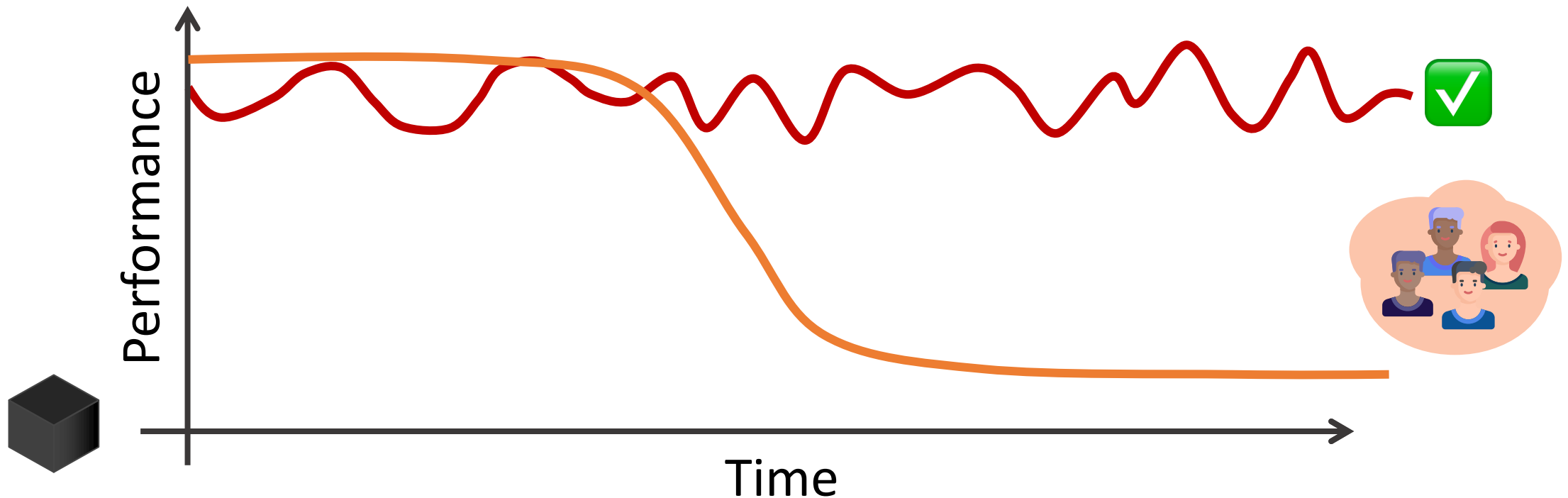
Subgroup perspective

- Identification of subgroups with *divergent* classification behavior
- Divergent subgroup analysis in speech data
- Subgroup-based model comparison
- Mitigate subgroup disparities
- **Interpretable subgroup drift detection**



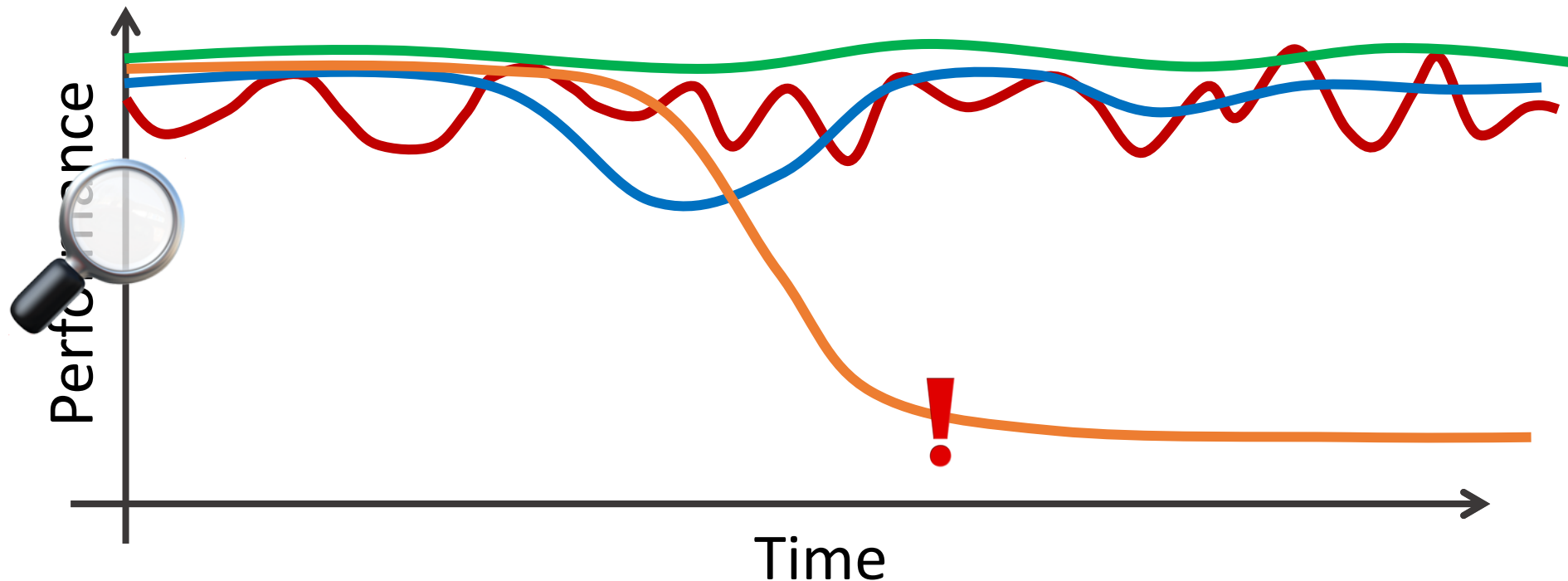
Subgroup based-drift detection

Typically, we monitor drift for overall performance



Subgroup based-drift detection

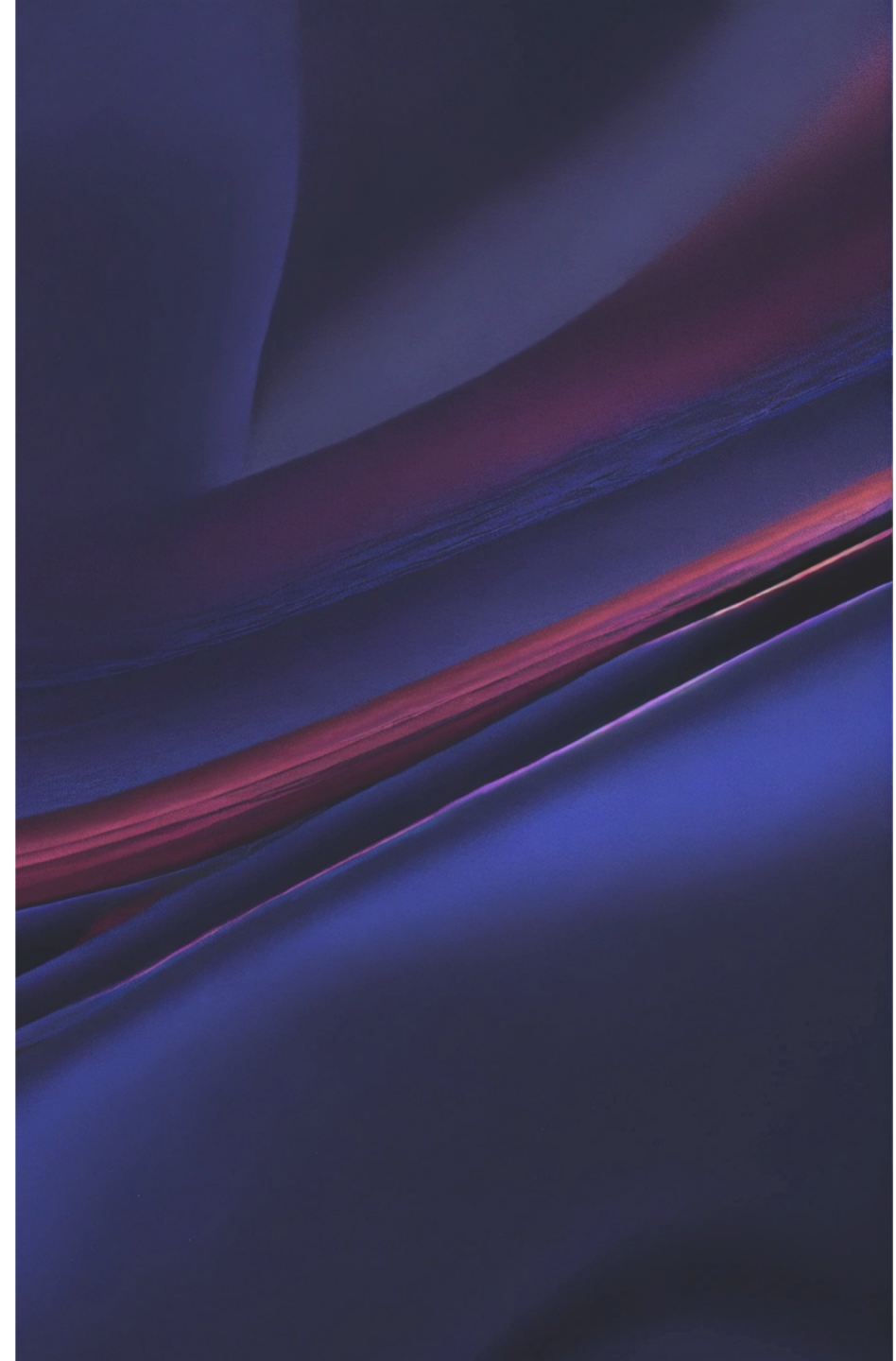
We propose an efficient algorithm to **monitorate subgroups** overtime and **detect subgroup drifts**



Outline

Subgroup perspective

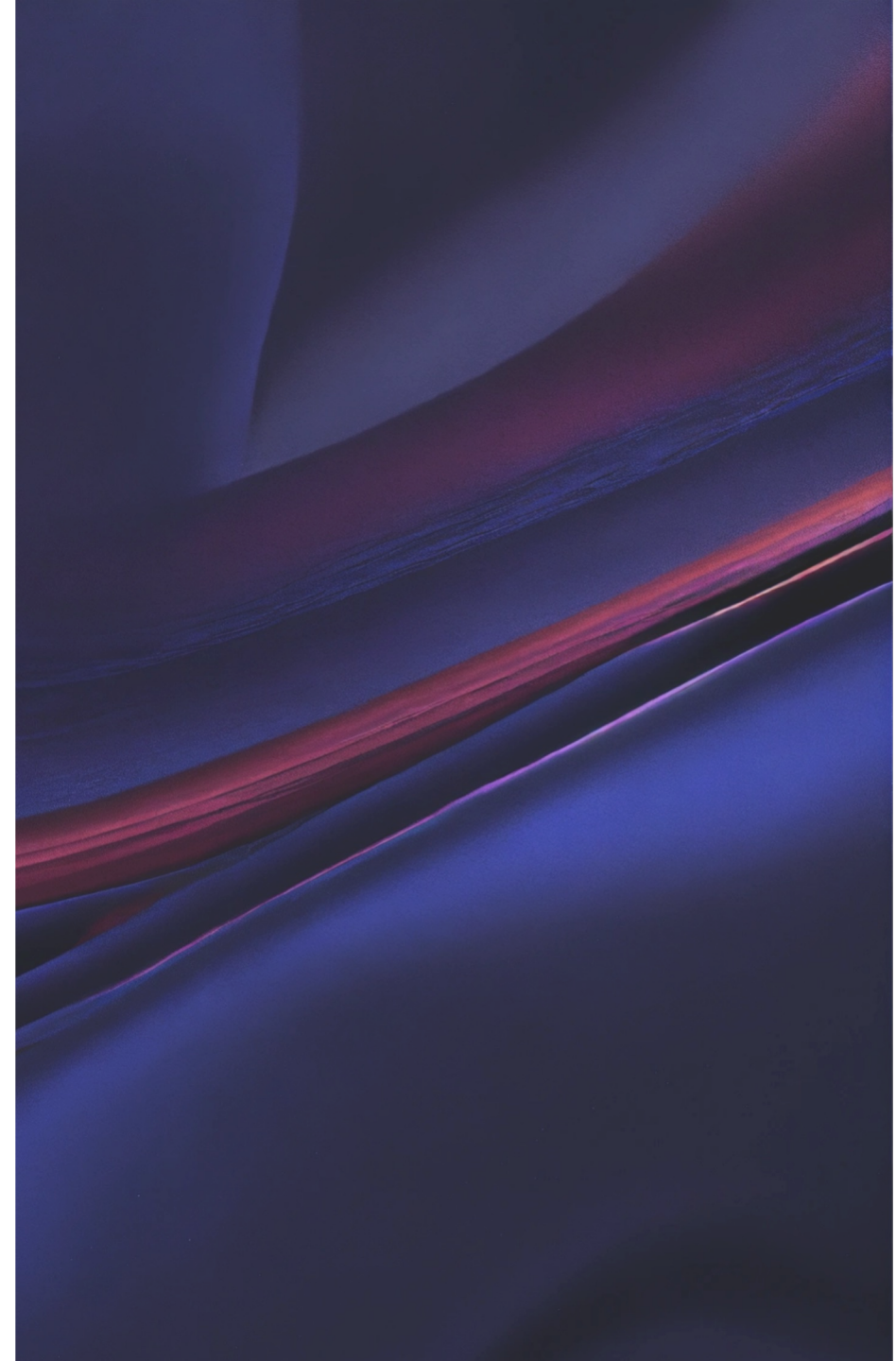
- Identification of subgroups with *divergent* classification behavior
- Divergent subgroup analysis in speech data
- Subgroup-based model comparison
- Mitigate subgroup disparities
- Interpretable subgroup drift detection



Outline

Local perspective

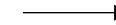
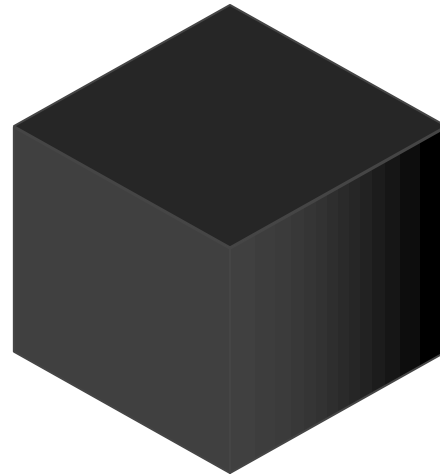
- **Explaining prediction of speech models**
- Assessing explainability methods for transformers models



Our scenario

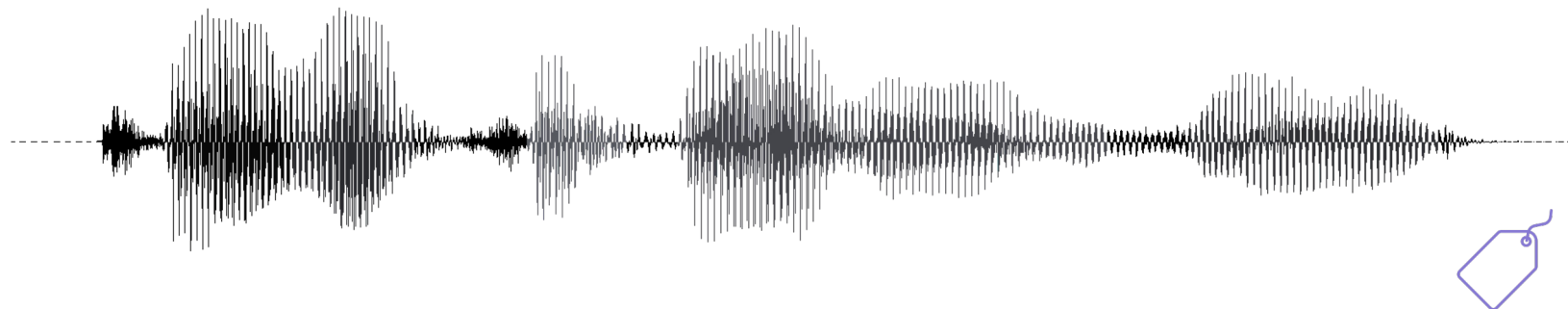


Turn on the kitchen lights



Action: activate
Object: lights
Location: kitchen

Why?

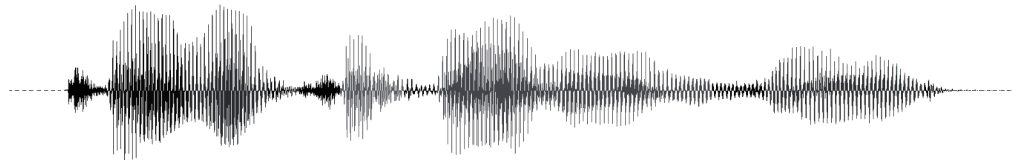


Explain the **interaction** between
utterance components and **predictions**
in a **human-understandable** manner

How do we define **interpretable representations** describing utterances?

Semantic

Spoken words



Turn up the bedroom heat

Paralinguistic

Prosody & external conditions



Pitch



Noise level



Speaking rate

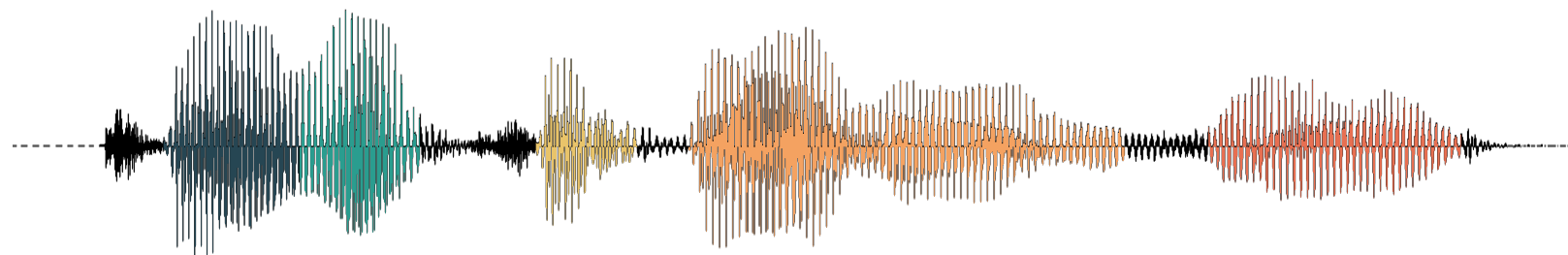
How do we **explain predictions** at the semantic and paralinguistic levels?

Perturbation-based approach

- Perturb the utterance based on an interpretable feature
- Measure the impact on predictions
- The greater the change, the more the model relies on this feature!

Semantic

Use a word-level
time alignment
model



Turn up

the

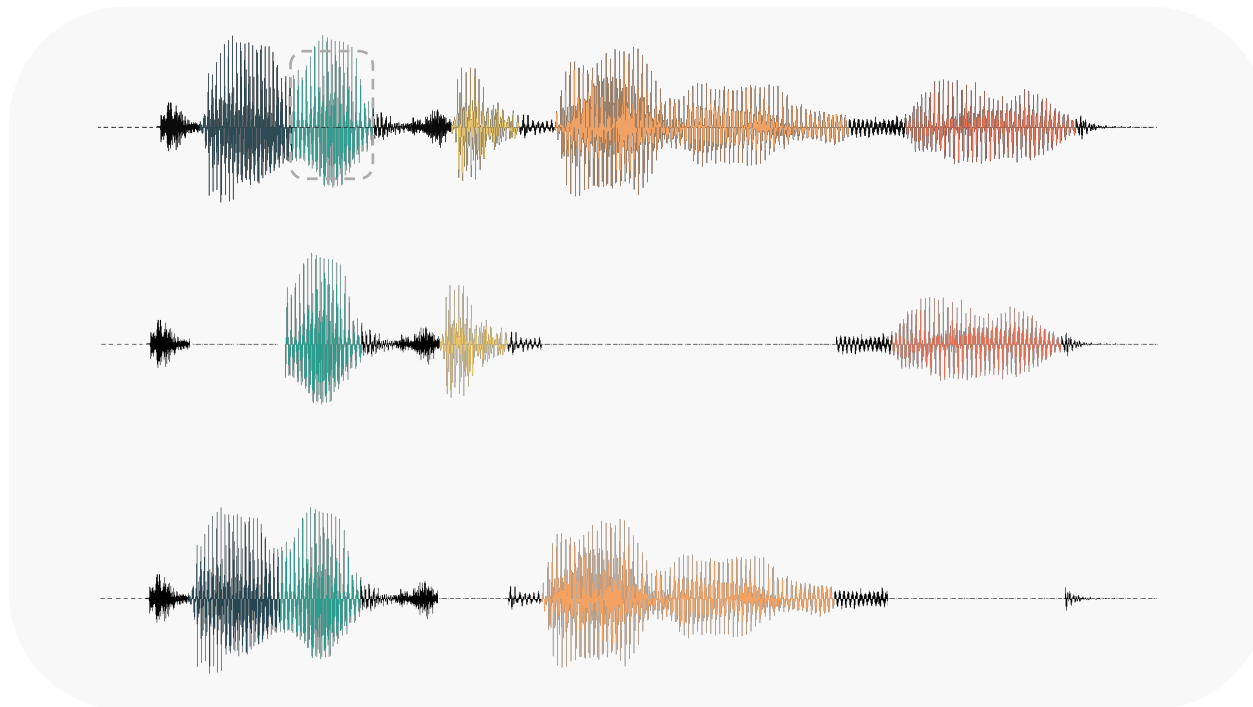
bedroom

heat



increase 98%
bedroom 85%

Mask audio segments



increase bedroom

25%

85%

36%

0%

98%

80%

Aggregate feature impact

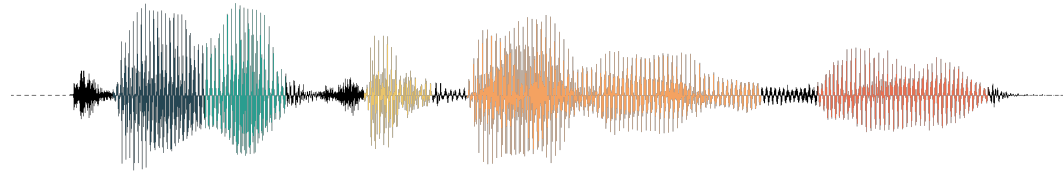
- Leave-one-out
- LIME

Turn up the bedroom heat



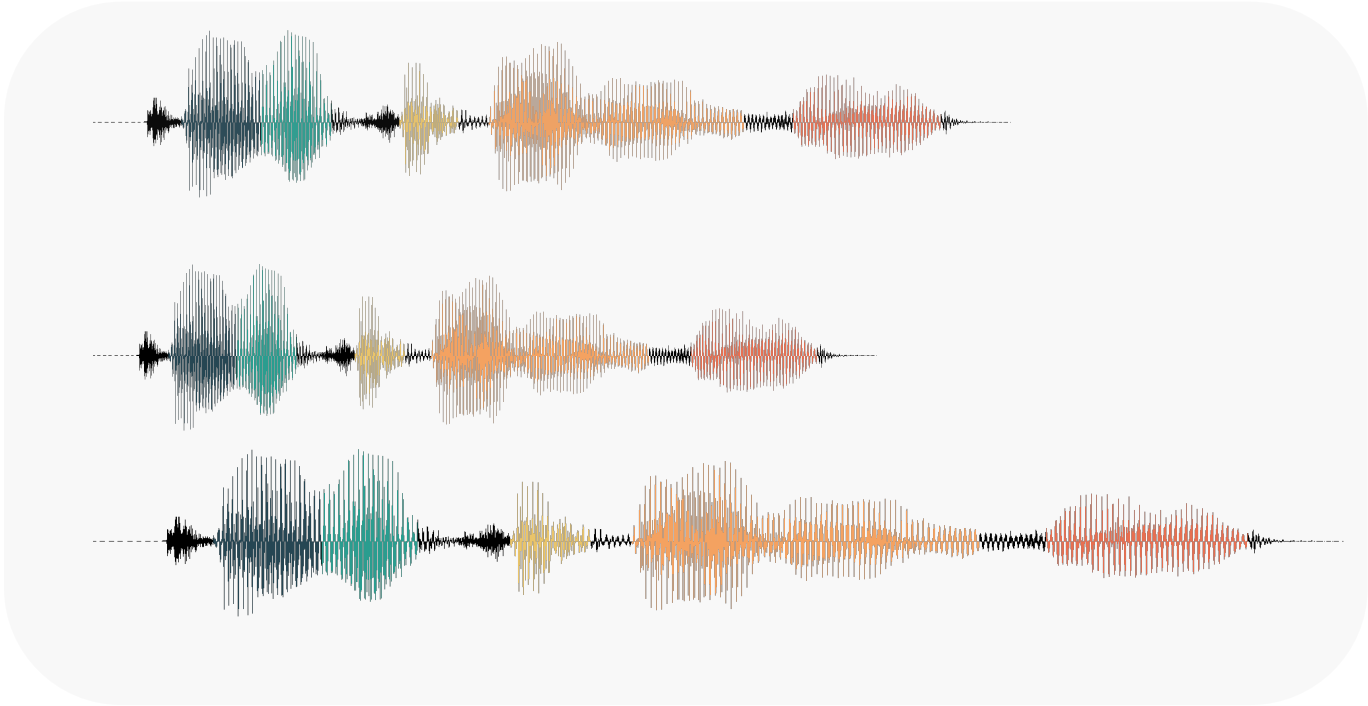
Word-level attributions

Paralinguistic

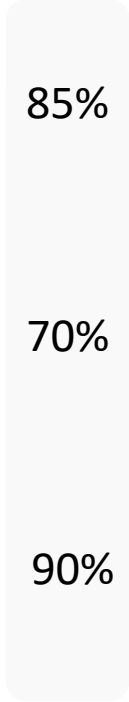
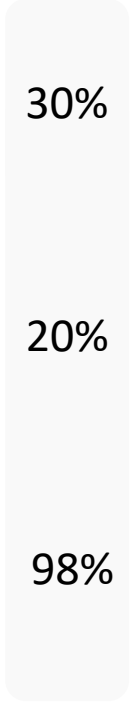


Paralinguistic

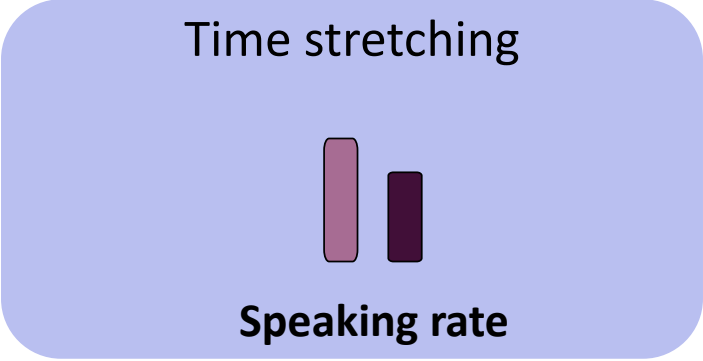
Perturb signal on paralinguistic



increase bedroom



Aggregate feature impact





Try it!

```
from speechxai import Benchmark
from transformers import Wav2Vec2ForSequenceClassification, Wav2Vec2FeatureExtractor

model = Wav2Vec2ForSequenceClassification.from_pretrained("superb/wav2vec2-base-superb-ic")
feature_extractor = Wav2Vec2FeatureExtractor.from_pretrained("superb/wav2vec2-base-superb-ic")

benchmark = Benchmark(model, feature_extractor)

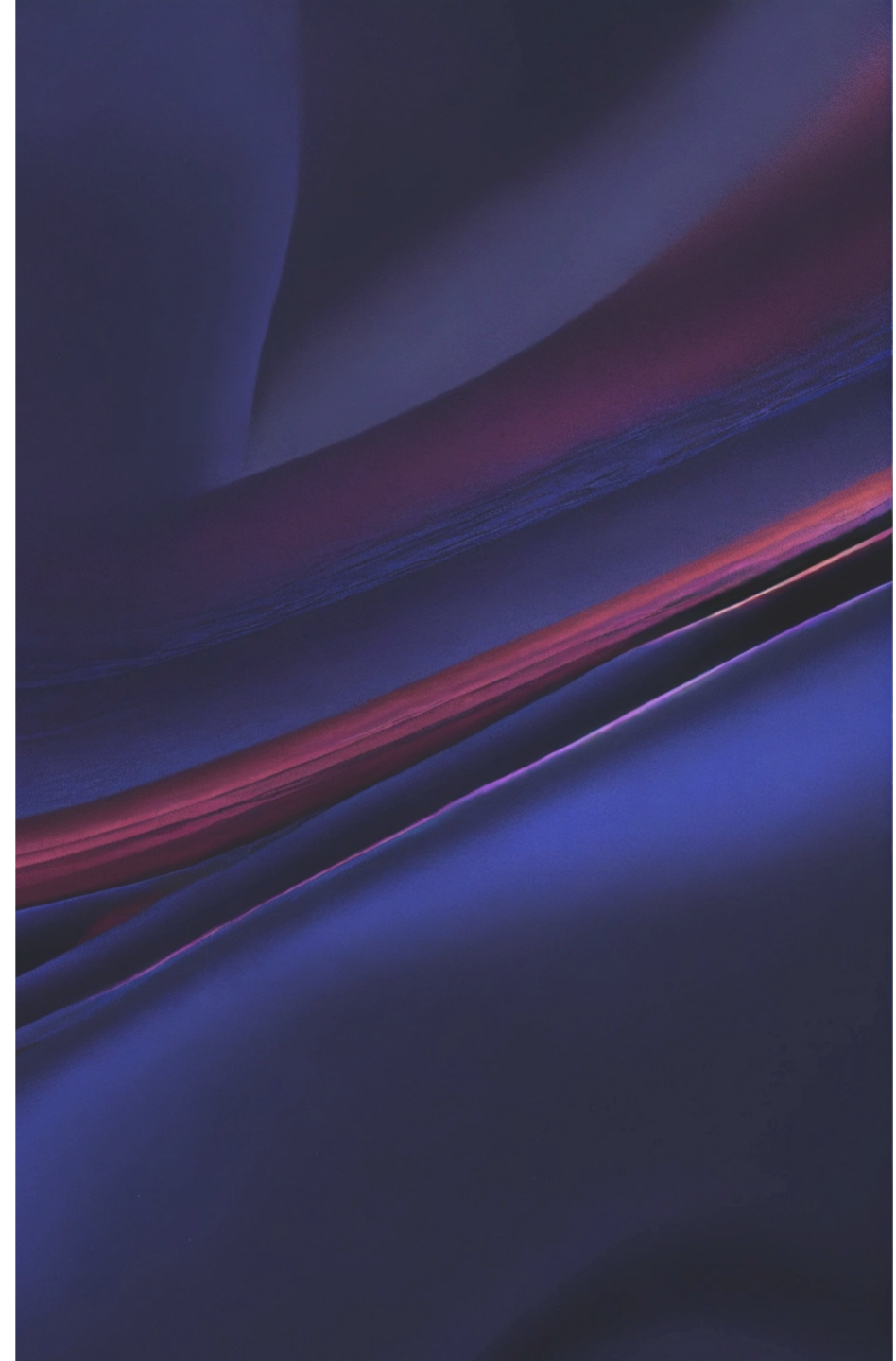
explanation = benchmark.explain(audio_path=audio_path, methodology="LIME")

benchmark.show_table(explanation)
```

Outline

Local perspective

- Explaining prediction of speech models
- **Assessing explainability methods for transformers models**



Library to explain and benchmark explainers

We proposed **ferret**, Python **library for benchmarking interpretability** techniques on **Transformers** for text and speech data



F E R R E T



What can you do with ferret?

Explaining individual prediction

Example. Sentiment classification – positive prediction

	Token	__Great	__movie	__for	__a	__great	__nap	!
Partition SHAP		0.35	0.12	0.05	0.06	0.35	-0.00	0.05
LIME		-0.07	-0.08	0.03	-0.01	-0.24	0.17	0.06
Gradient		0.12	0.17	0.06	0.04	0.14	0.23	0.05
Gradient (x Input)		-0.11	-0.09	-0.08	0.03	0.03	0.11	-0.05
Integrated Gradient		-0.09	0.10	0.11	-0.02	0.10	0.02	-0.03
Integrated Gradient (x Input)		-0.09	-0.15	-0.17	-0.15	-0.10	-0.24	-0.10



What can you do with ferret?

Evaluate explanations

Faithfulness

How accurately the explanation reflects the inner working of the model

Plausibility

How explanations are aligned with human reasoning



Faithfulness

	Token	__Great	__movie	__for	__a	__great	__nap	!
Partition SHAP		0.35	0.12	0.05	0.06	0.35	-0.00	0.05
LIME		-0.07	-0.08	0.03	-0.01	-0.24	0.17	0.06
Gradient		0.12	0.17	0.06	0.04	0.14	0.23	0.05
Gradient (x Input)		-0.11	-0.09	-0.08	0.03	0.03	0.11	-0.05
Integrated Gradient		-0.09	0.10	0.11	-0.02	0.10	0.02	-0.03
Integrated Gradient (x Input)		-0.09	-0.15	-0.17	-0.15	-0.10	-0.24	-0.10

	aopc_compr	aopc_suff	taucorr_loo
Partition SHAP	-0.13	-0.05	-0.33
LIME	-0.01	-0.20	0.24
Gradient	-0.16	-0.09	0.24
Gradient (x Input)	-0.00	-0.20	0.52
Integrated Gradient	-0.02	-0.17	0.33
Integrated Gradient (x Input)	0.00	1.00	-0.62



Plausibility

Human explanation

Token	__Great	__movie	__for	__a	__great	__nap	!
Partition SHAP	0.35	0.12	0.05	0.06	0.35	-0.00	0.05
LIME	-0.07	-0.08	0.03	-0.01	-0.24	0.17	0.06
Gradient	0.12	0.17	0.06	0.04	0.14	0.23	0.05
Gradient (x Input)	-0.11	-0.09	-0.08	0.03	0.03	0.11	-0.05
Integrated Gradient	-0.09	0.10	0.11	-0.02	0.10	0.02	-0.03
Integrated Gradient (x Input)	-0.09	-0.15	-0.17	-0.15	-0.10	-0.24	-0.10

auprc_plau **token_f1_plau** **token_iou_plau**

Partition SHAP	1.00	0.50	0.33
LIME	0.14	0.00	0.00
Gradient	0.29	0.44	0.29
Gradient (x Input)	0.24	0.40	0.25
Integrated Gradient	0.22	0.33	0.20
Integrated Gradient (x Input)	0.64	0.00	0.00



Try ferret!

ferret-xai 0.4.2

```
pip install ferret-xai
```



```
from transformers import AutoModelForSequenceClassification, AutoTokenizer
from ferret import Benchmark

name = "cardiffnlp/twitter-xlm-roberta-base-sentiment"
model = AutoModelForSequenceClassification.from_pretrained(name)
tokenizer = AutoTokenizer.from_pretrained(name)

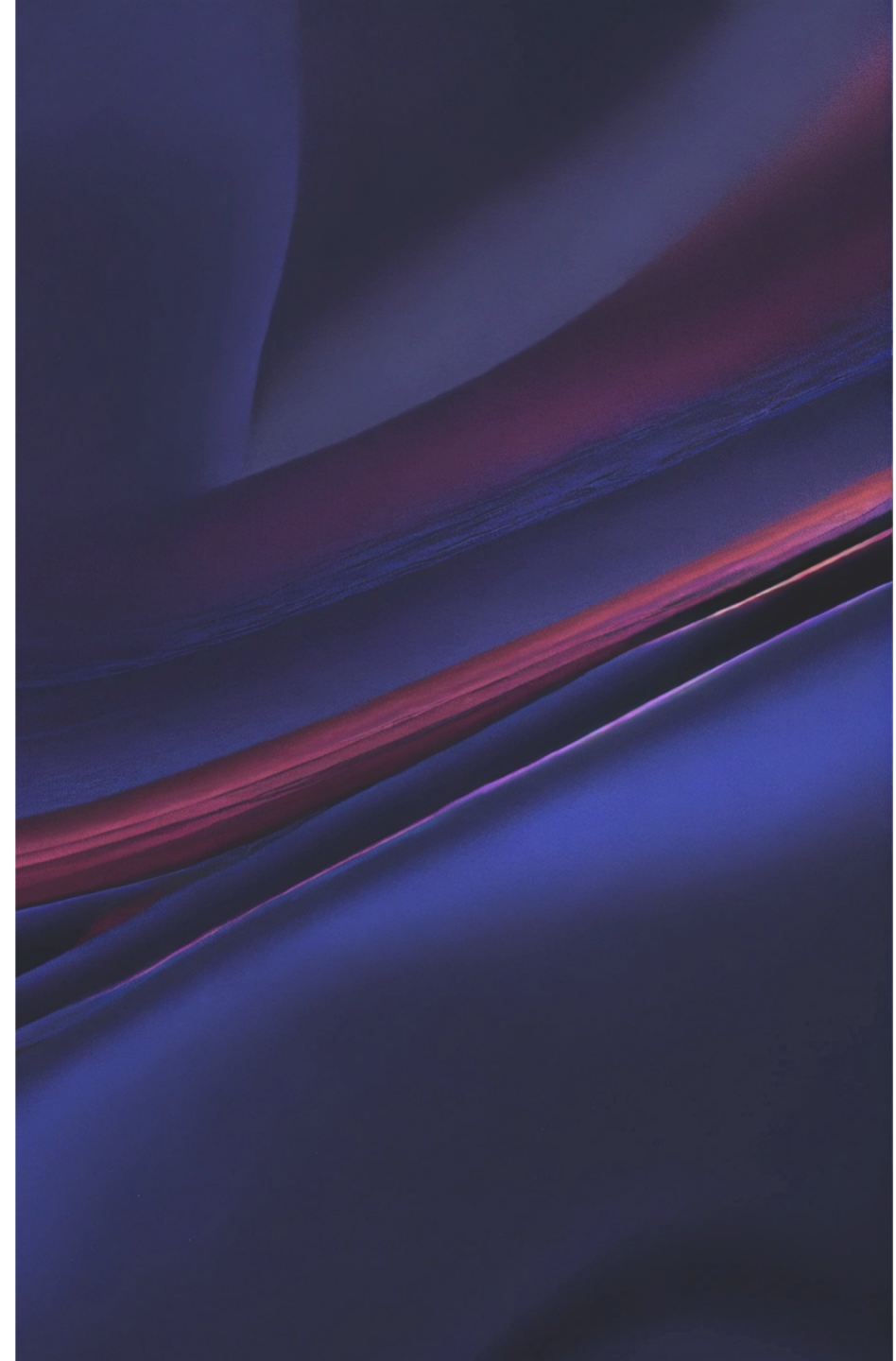
bench = Benchmark(model, tokenizer)
explanations = bench.explain("You look stunning!", target=1)
evaluations = bench.evaluate_explanations(explanations, target=1)

bench.show_evaluation_table(evaluations)
```

Outline

Local perspective

- Explaining prediction of speech models
- Assessing explainability methods for transformers models



Joint works with



Alkis Koudounas



Elena Baralis



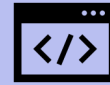
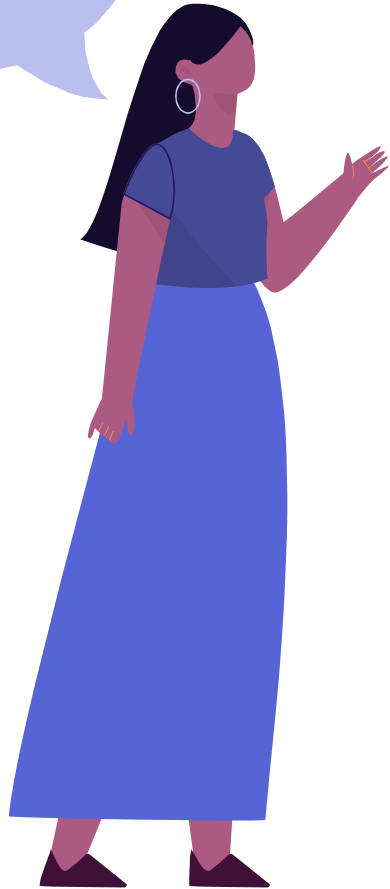
Luca de Alfaro



Flavio Giobergia

and all the other colleagues
& collaborators!

Thanks!



elianap.github.io/



eliana.pastor@polito.it



[eliana__pastor](https://twitter.com/eliana__pastor)



[@elianapastor.bsky.social](https://bsky.app/profile/@elianapastor.bsky.social)



[eliana-pastor](https://www.linkedin.com/in/eliana-pastor)