# Assessing Speech Model Performance: A Subgroup Perspective

Alkis Koudounas, **Eliana Pastor**, Elena Baralis

Politecnico di Torino - Italy

SEBD 2024
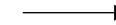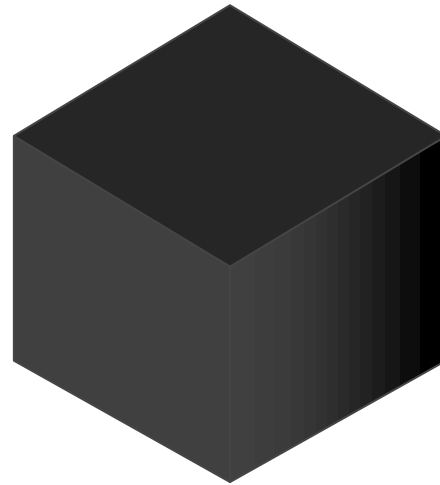
Politecnico di Torino    DBMG    ICSC Centro Nazionale di Ricerca in HPC, Big Data and Quantum Computing

In collaboration with amazon | science

# Our scenarios



**Automatic Speech Recognition**

Turn on the kitchen lights

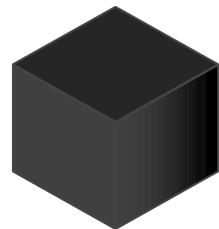**Intent classification**

Action:    activate
Object:    lights
Location: kitchen

**Emotion recognition**

Neutral

PERFORMANCE
X%

LOW ERROR RATE

HIGH ERROR RATE

AS OVERALL
MODEL
BEHAVIOR

# Outline

- **Identification of interpretable subgroups** with divergences in performance

- **Model comparison** from the subgroup perspective

- **Subgroup-guided acquisition** for model improvement

# Outline

- **Identification of interpretable subgroups** with divergences in performance

- **Model comparison** from the subgroup perspective

- **Subgroup-guided acquisition** for model improvement

How to make an interpretable data grouping?

# 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
...

# Subgroup identification

- **Automatic** identification of subgroups via **frequent pattern mining**
  - Slicing in the interpretable attibute space

- Compute subgroup **divergence**

performance measure

$$\Delta(S) = f(S) - f(D)$$
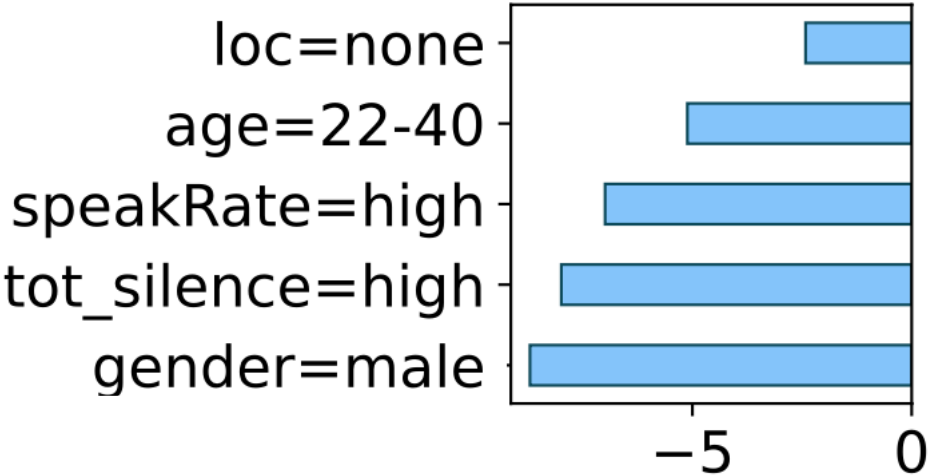
**pattern**, e.g., {age=20-35, gender=female}

**all dataset**

# Divergent subgroup

**By 31.22 less accurate!**

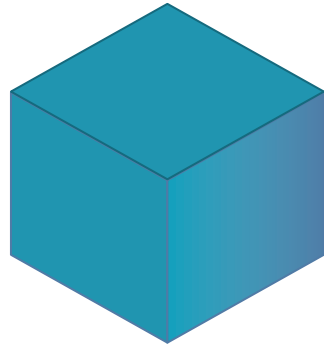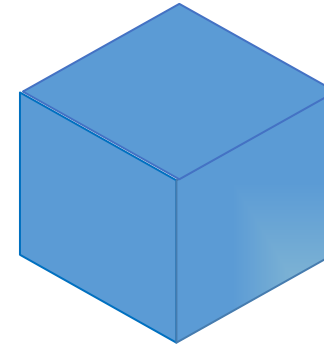| Subgroups | $f$ | $\Delta_f$ |
|---|---|---|
| $\mathrm{I}^-$ : {age=22-40, gender=male, location=none, speaking rate=high, tot silence=high} | 60.50 | -31.22 |
| $\mathrm{I}^+$ : {age=22-40, location=washroom, speaking rate=low, trimmed duration=high} | 100.0 | 8.28 |

**More accurate than average**

# Outline

- **Identification of interpretable subgroups** with divergences in performance

- **Model comparison** from the subgroup perspective

- **Subgroup-guided acquisition** for model improvement

Accuracy 91.72%

Accuracy 93.17%

# Which model to choose?

## .. most accurate..?

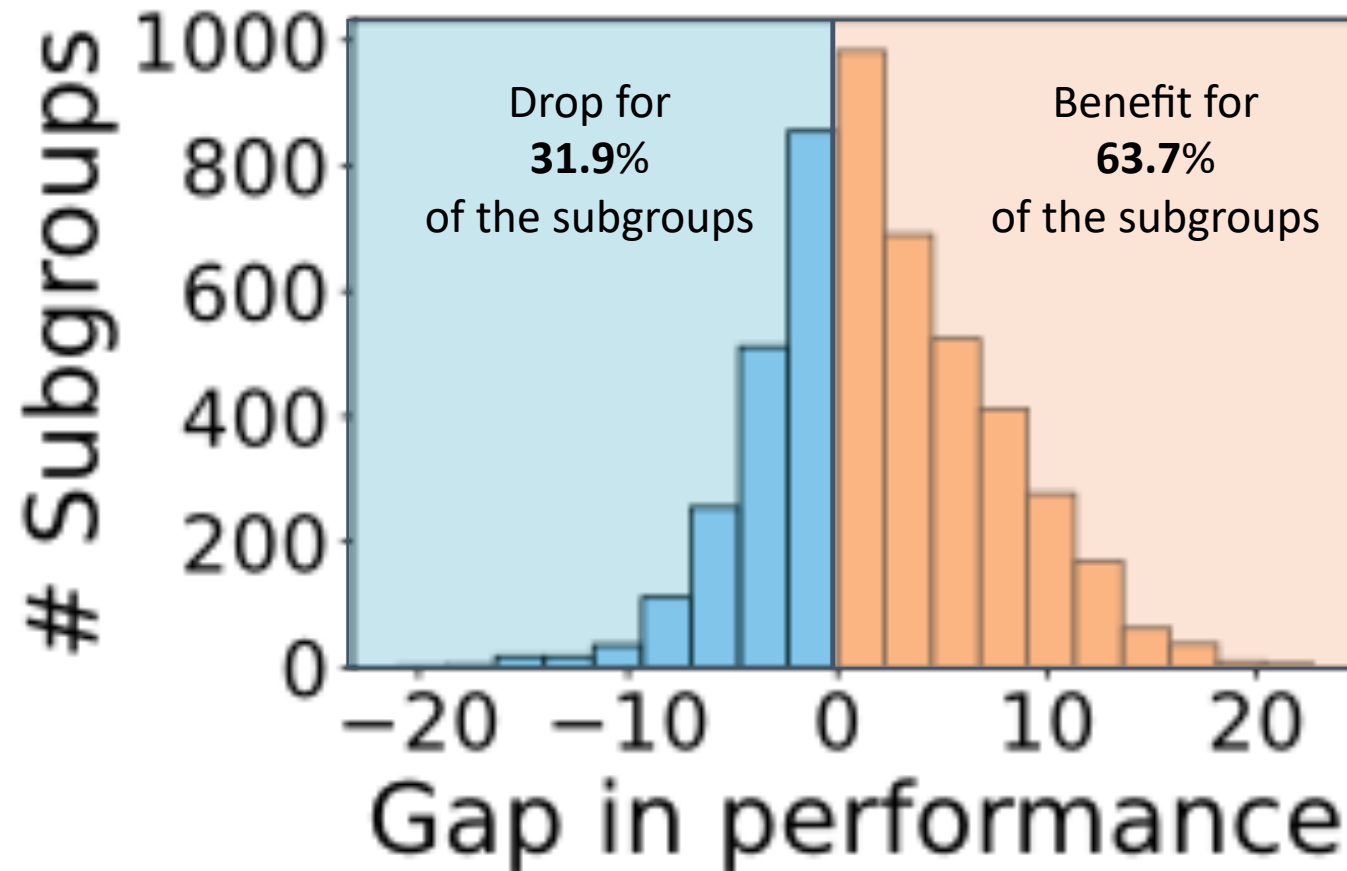## But on subgroups?

# 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



| Subgroups | $gap_f$ | $f_{\text{w2v2-b}}$ | $f_{\text{w2v2-l}}$ |
|---|---|---|---|
| ↑ {action=increase, location=none, tot duration=low, trimmed speaking rate=low, trimmed duration=low} | 22.69 | 75.63 | 98.32 |
| ↓ {action=activate, gender=male, speaking rate=low} | -20.97 | 96.77 | 75.81 |

**Increase in performance**

**Drop in performance**

# Outline

- **Identification of interpretable subgroups** with divergences in performance

- **Model comparison** from the subgroup perspective

- **Subgroup-guided acquisition** for model improvement

# Subgroup-guided data acquisition
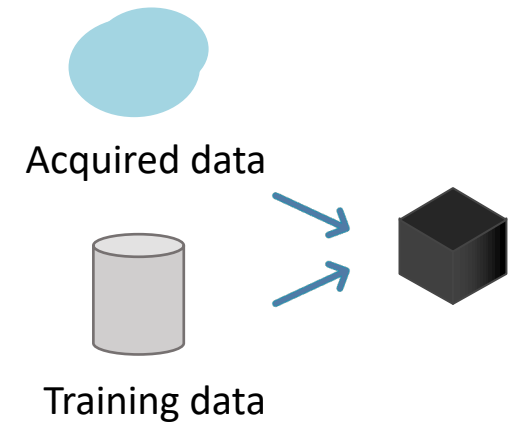
*Speaking rate=high, gender=male*

**Step 1.**
Identify the divergent patterns

**Step 2.**
Acquire data satisfying the patterns

Acquired data

Training data

**Step 3.**
Speech model re-training

# Results of subgroup-guided data acquisition

| Approach | #samples | Accuracy | F1 Macro | $\Delta^{-}_{max}$ | $\Delta^{-}_{avg-10}$ |
|---|---|---|---|---|---|
| original | 18506 | $91.58 \pm 0.08$ | $86.34 \pm 0.13$ | $-70.09 \pm 0.26$ | $-70.09 \pm 0.26$ |
| random | +226 | $92.56 \pm 0.44$ | $90.25 \pm 0.60$ | $-52.20 \pm 2.57$ | $-51.11 \pm 2.19$ |
| clustering | +226 | $89.77 \pm 0.88$ | $87.02 \pm 0.15$ | $-47.37 \pm 0.42$ | $-47.34 \pm 0.42$ |
| *ours* | +226 | $\mathbf{96.55 \pm 0.08}$ | $\mathbf{94.71 \pm 0.12}$ | $\mathbf{-40.60 \pm 0.35}$ | $\mathbf{-40.28 \pm 0.36}$ |
| all data | +4606 | $93.42 \pm 0.17$ | $93.11 \pm 0.17$ | $-53.18 \pm 0.15$ | $-50.89 \pm 0.09$ |

**Improvement compared to acquire all the data!**

**We improve overall performance!**

**We improve subgroup performance!**

Thanks!

✉ eliana.pastor@polito.it

✕ eliana__pastor

</> elianap.github.io/

or let's have a chat! ☀