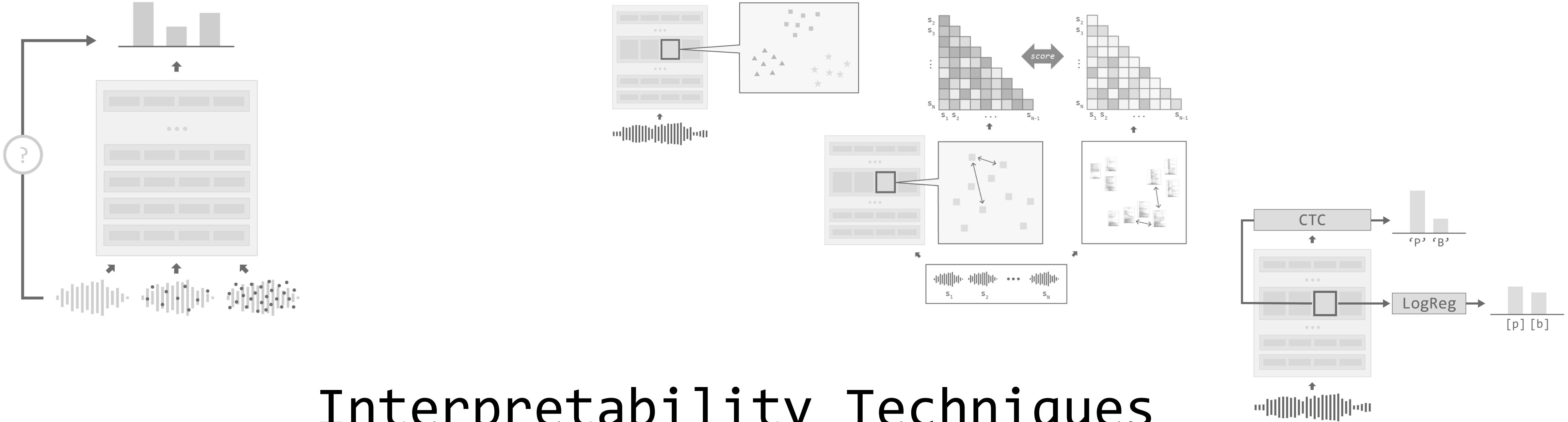


Interpretability Techniques for Speech Models

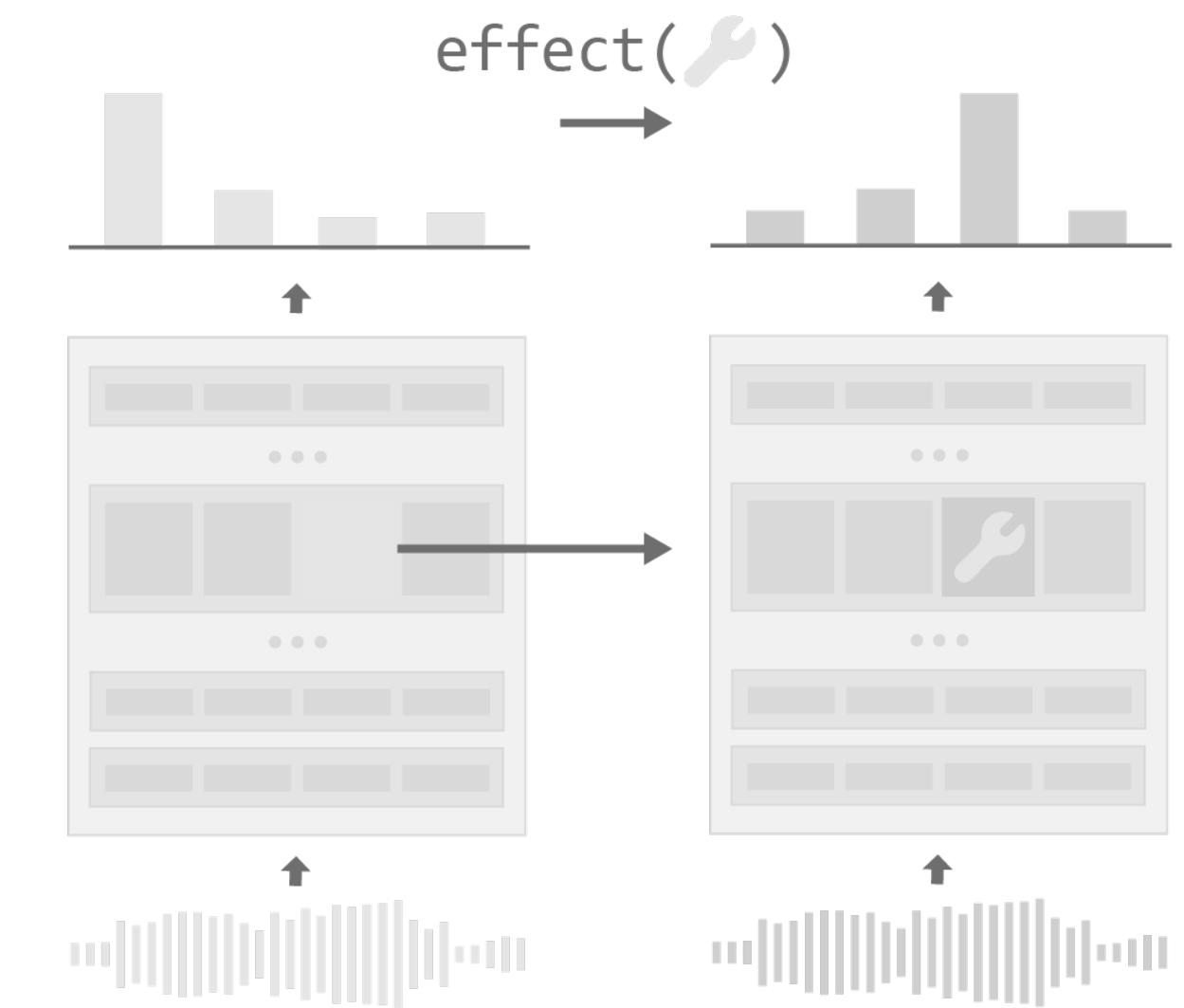
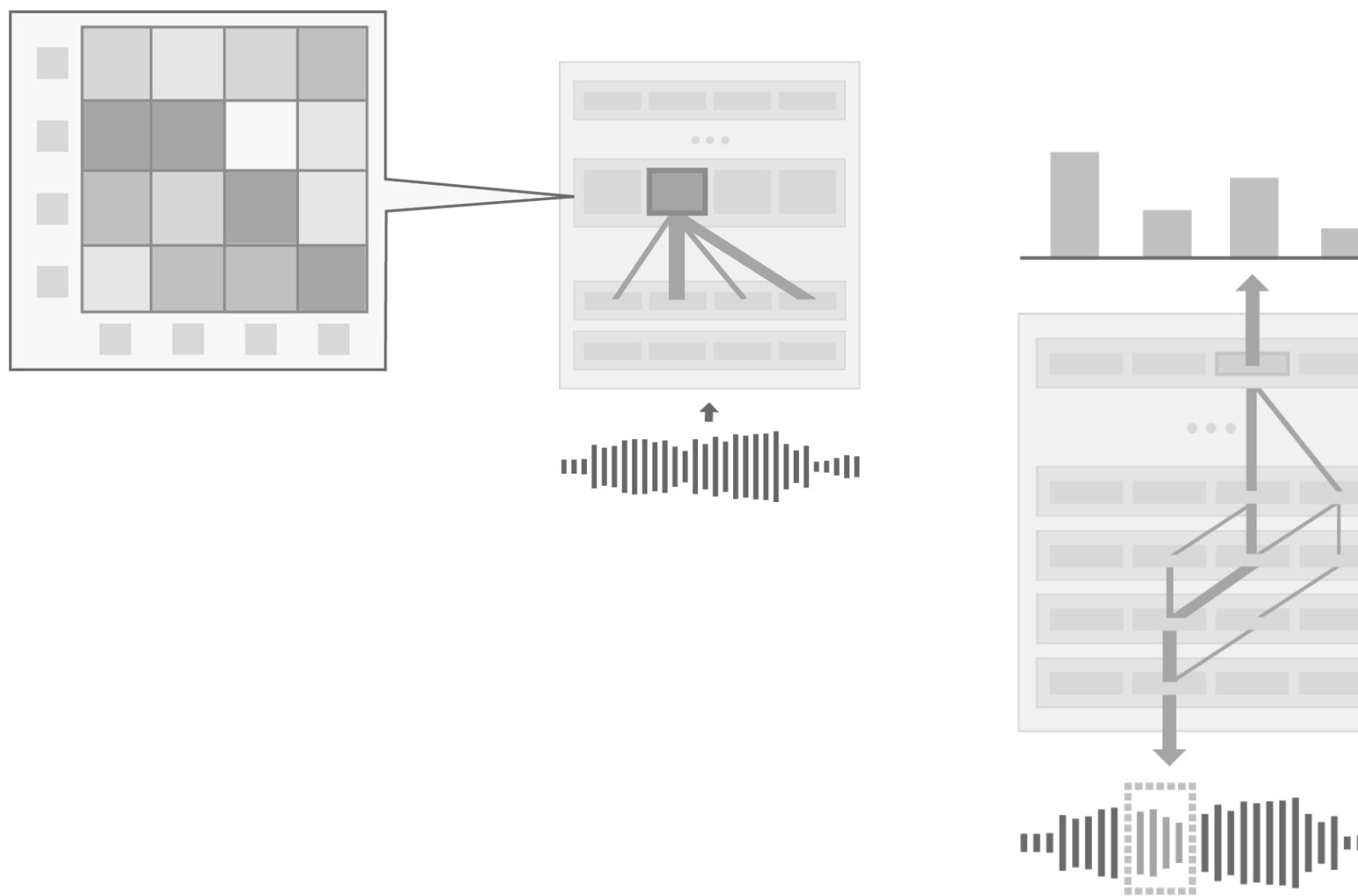
Conclusions & Outlook

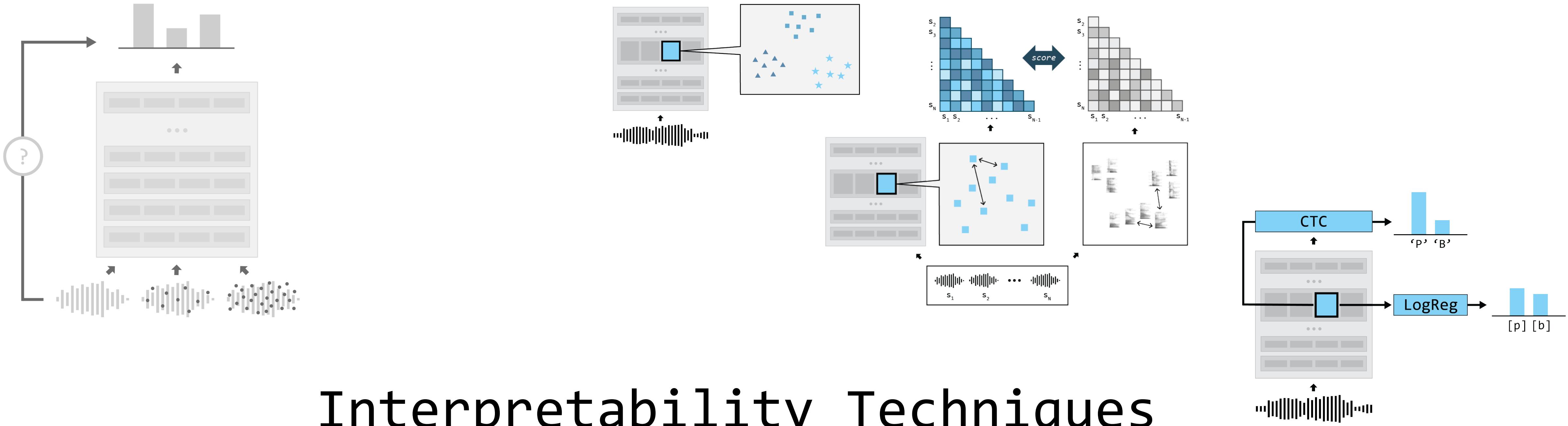


Marianne de Heer Kloots, 17-08-2025

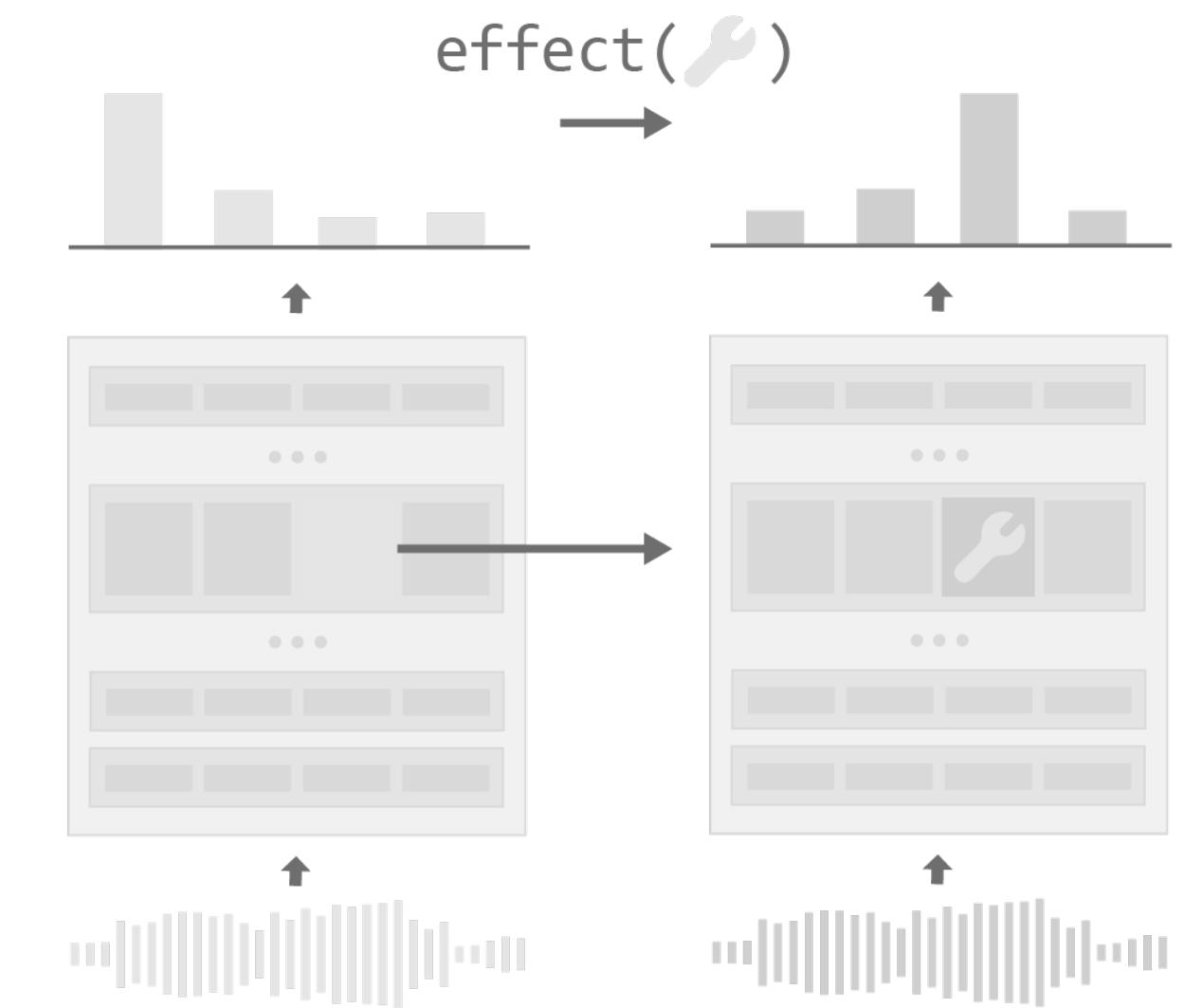
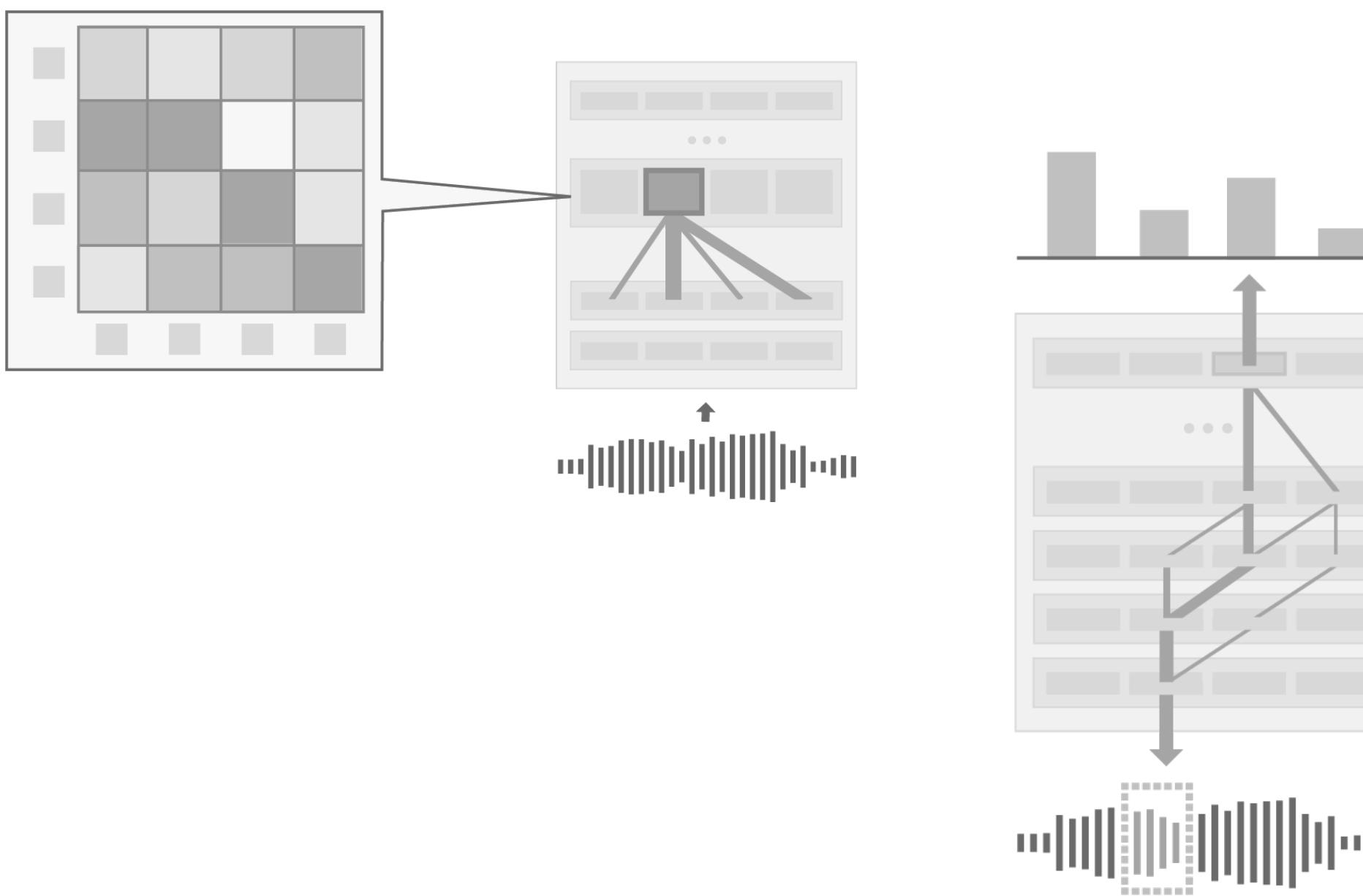


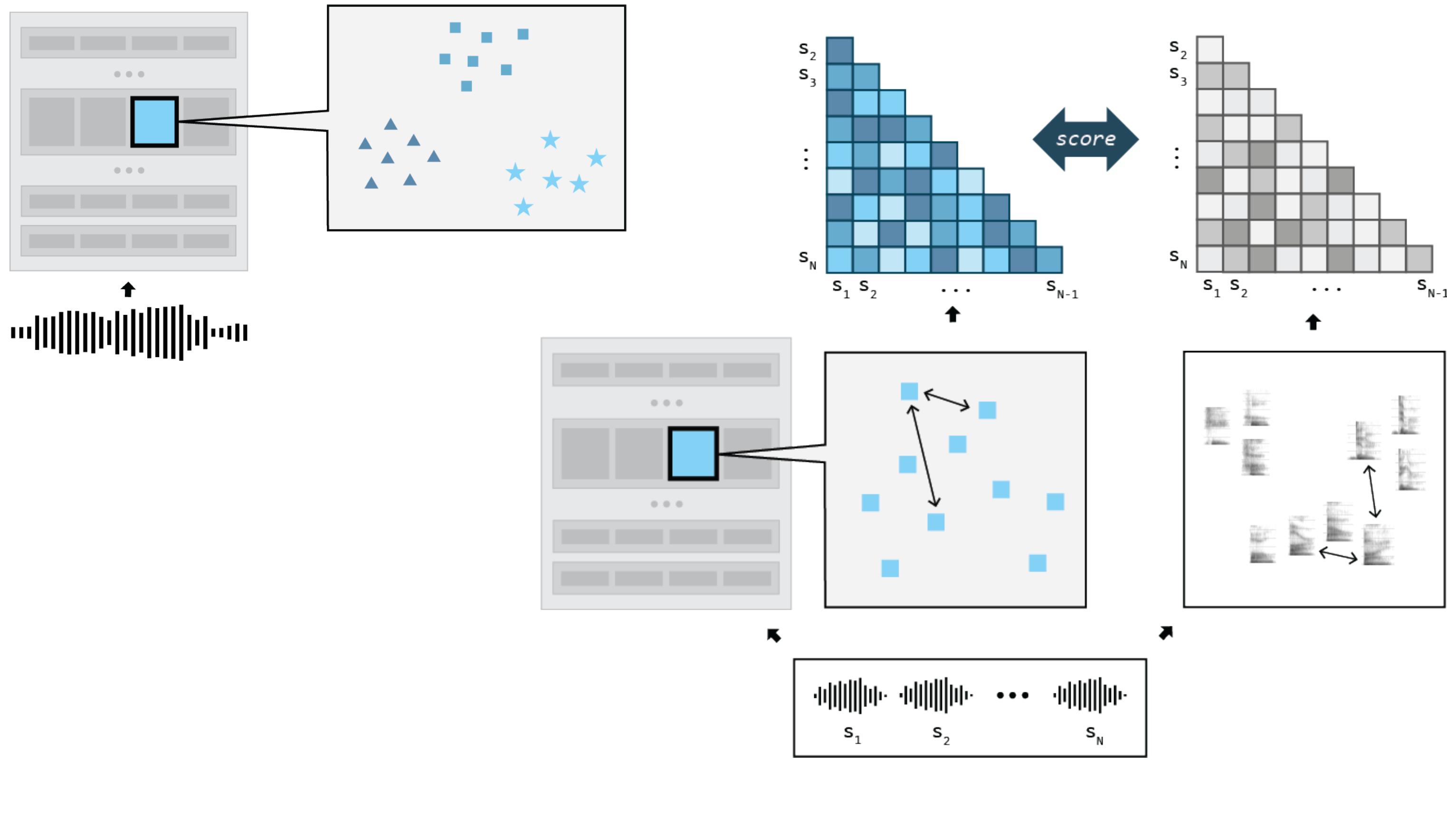
Interpretability Techniques for Speech Models





Interpretability Techniques for Speech Models

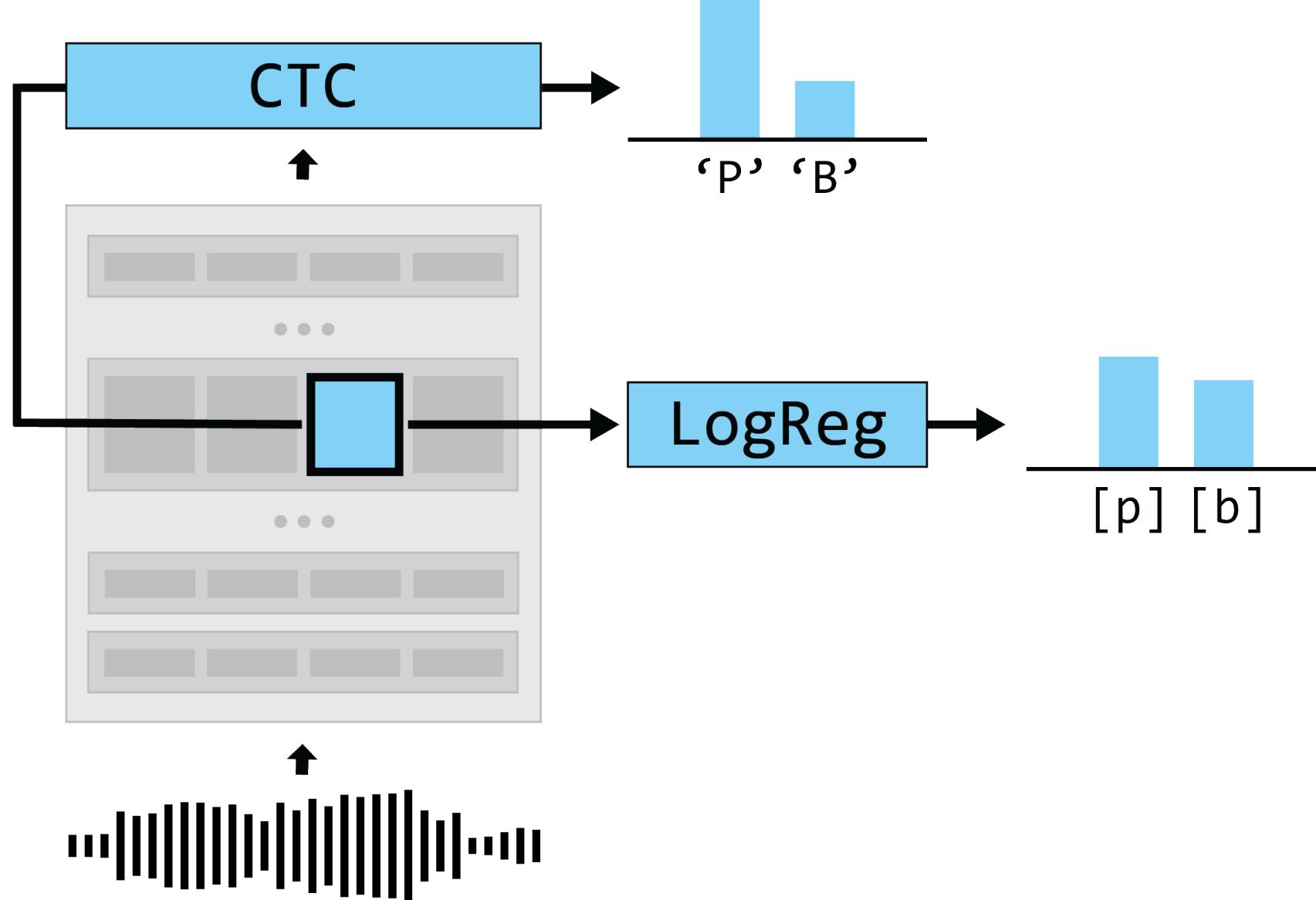




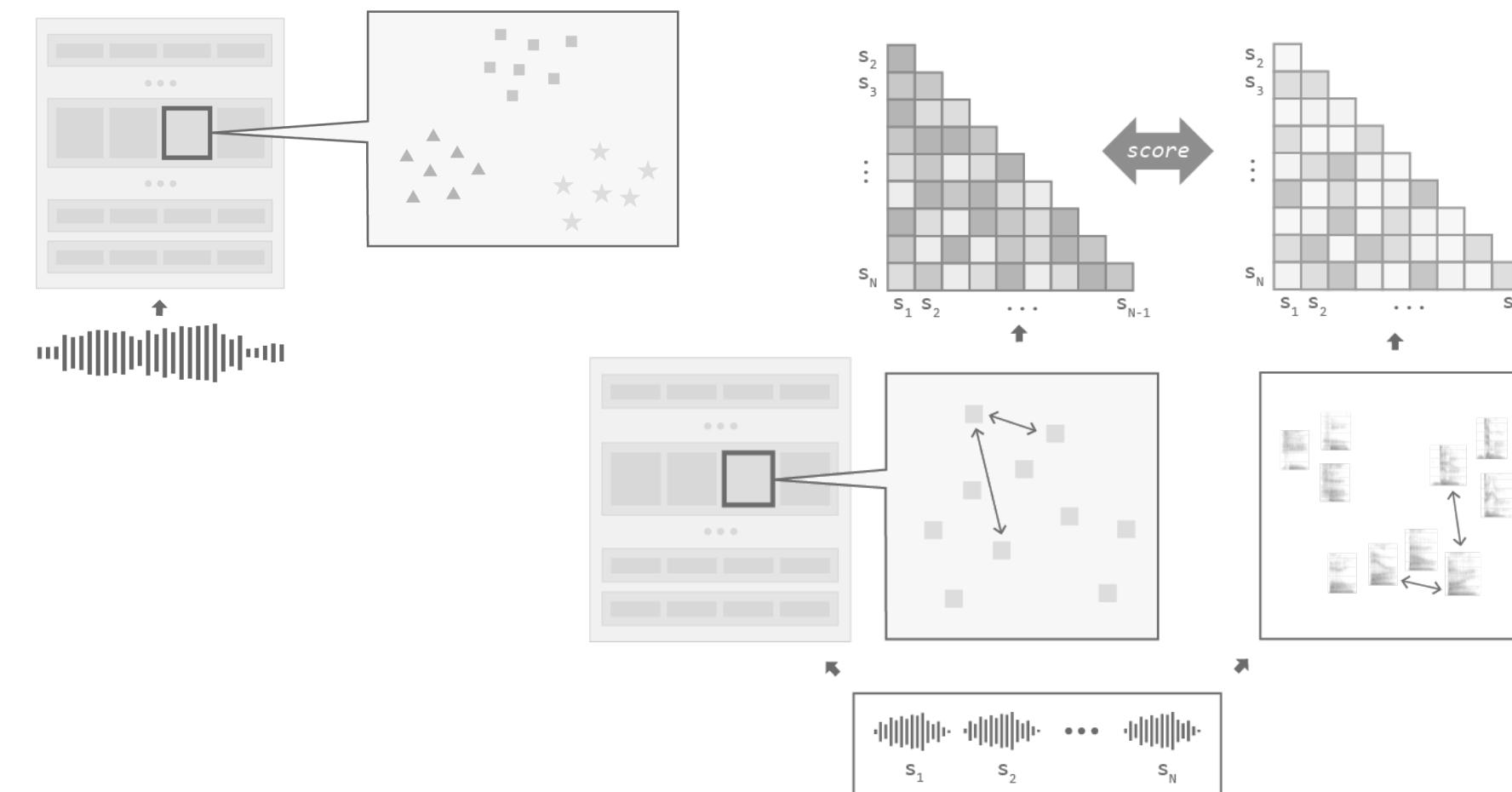
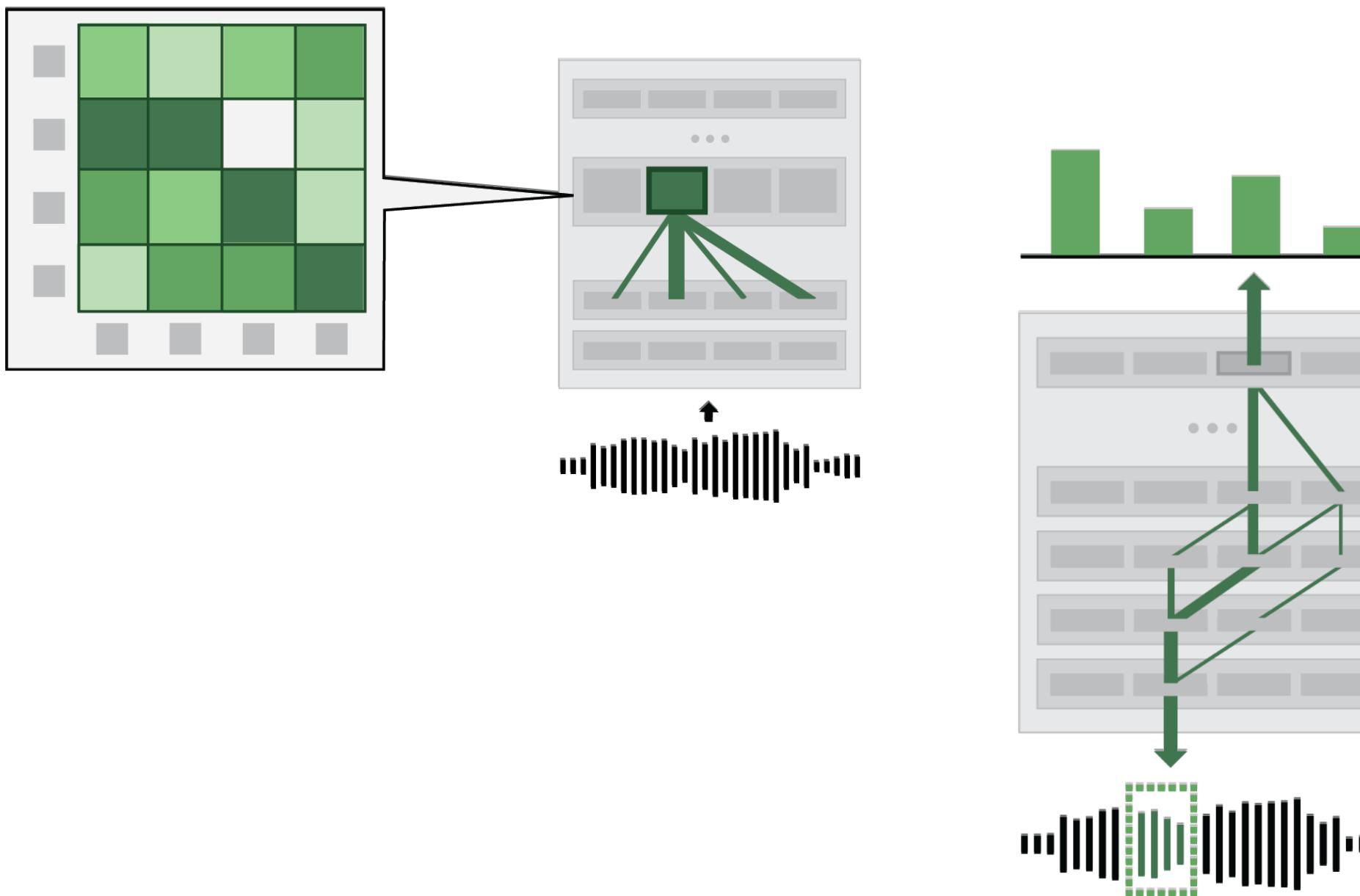
- Baselines and controls
- Zero-shot metrics vs. optimized diagnostic tools

Representational analyses

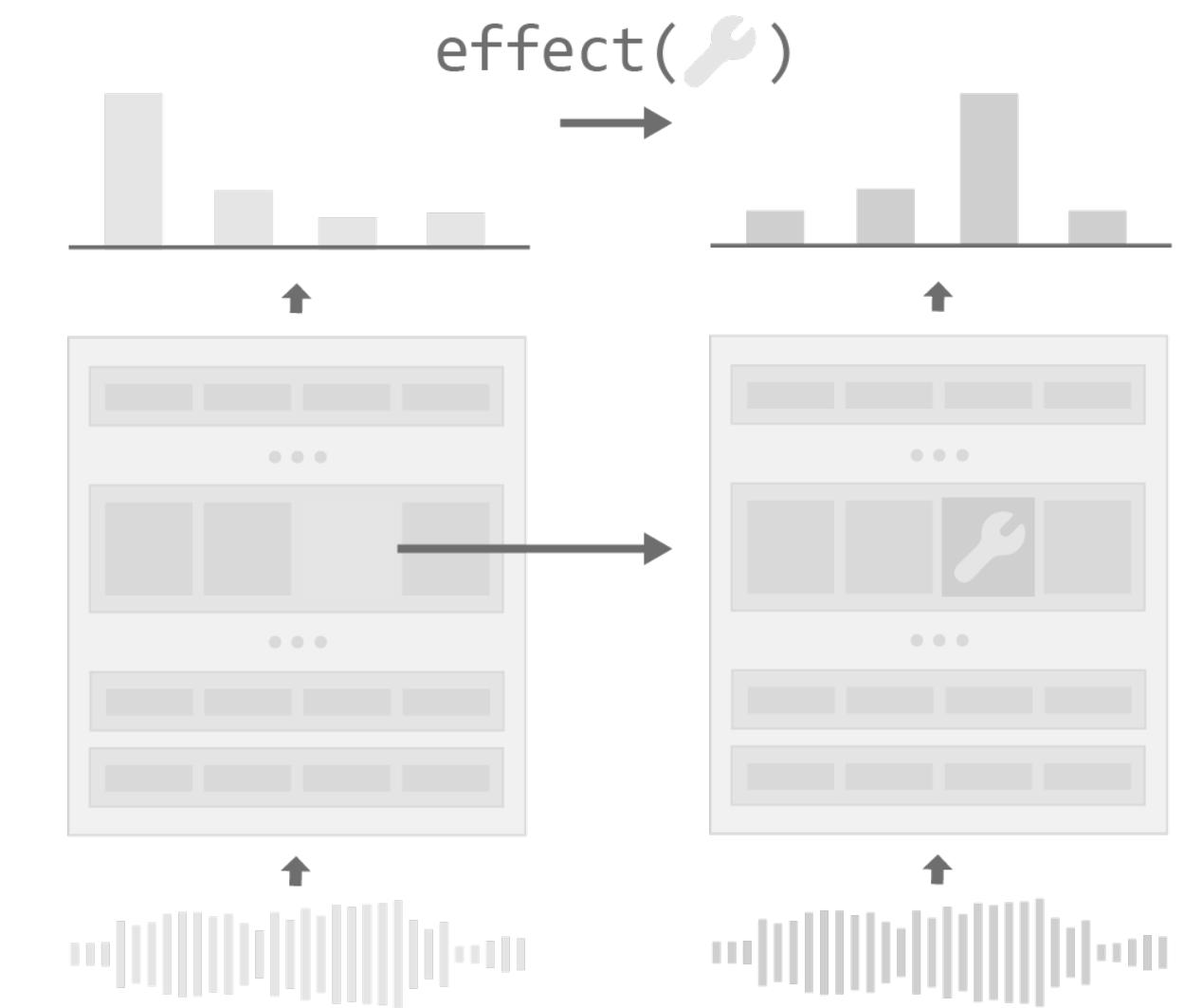
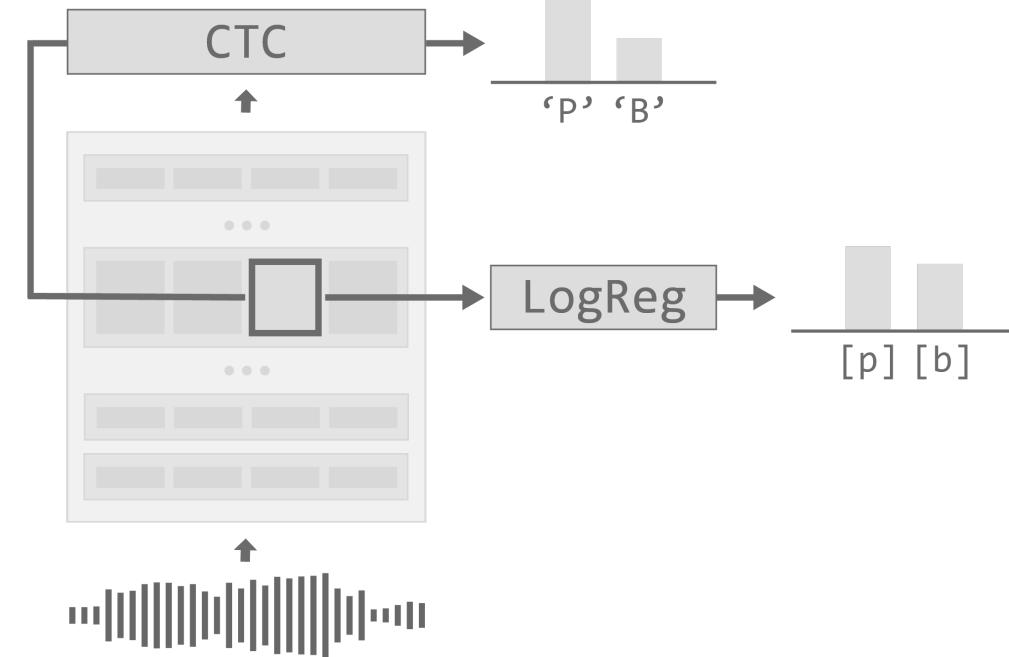
- Dimensionality reduction
- Probing classifiers
- Representation space comparisons

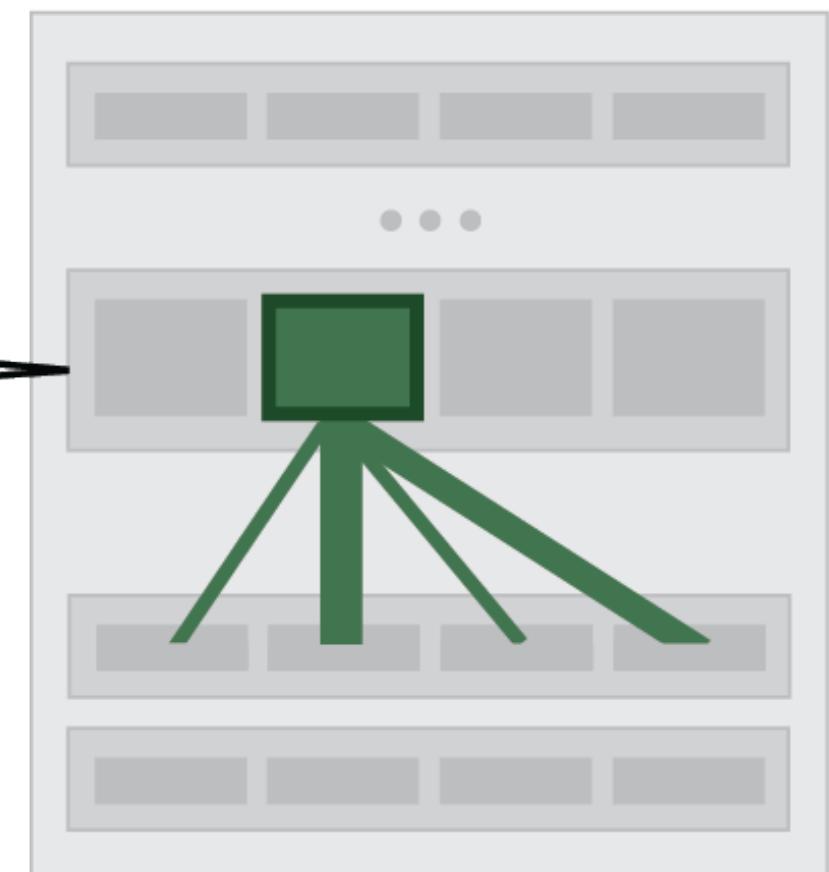
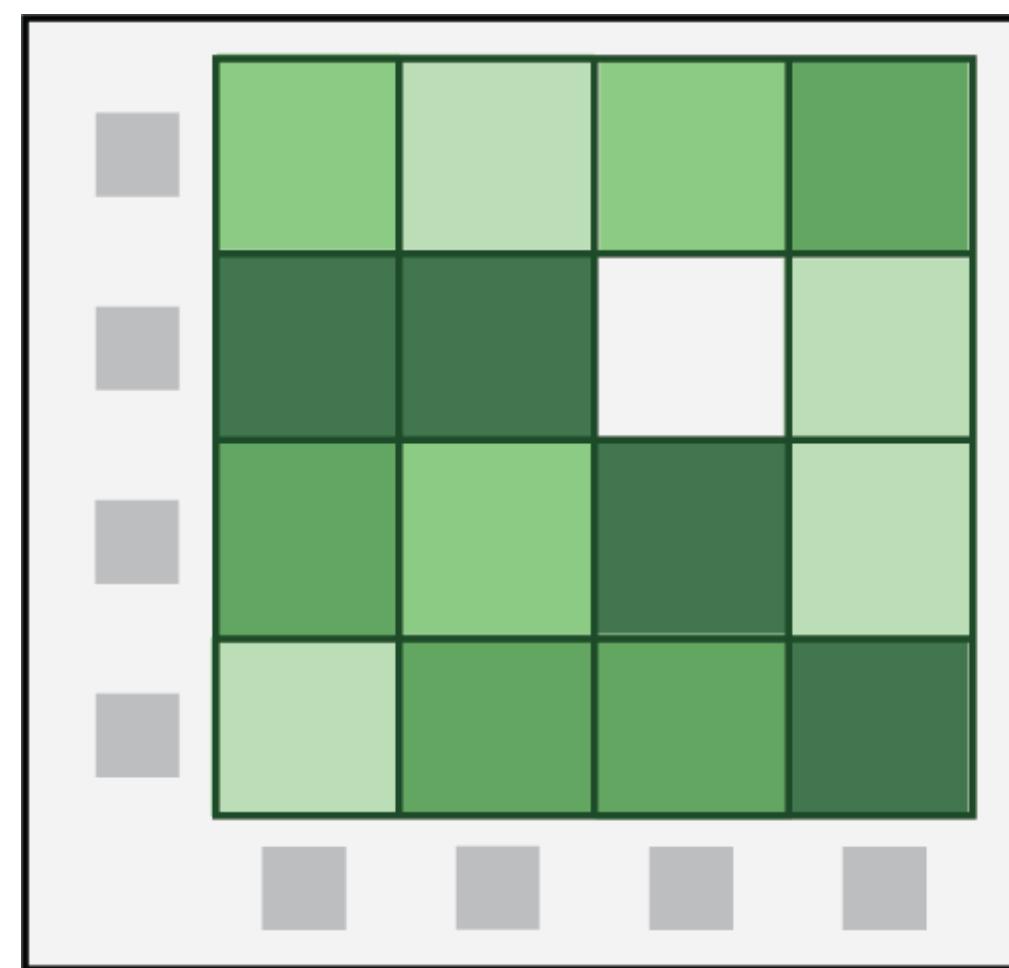


Interpretability Techniques for Speech Models



Representational analyses

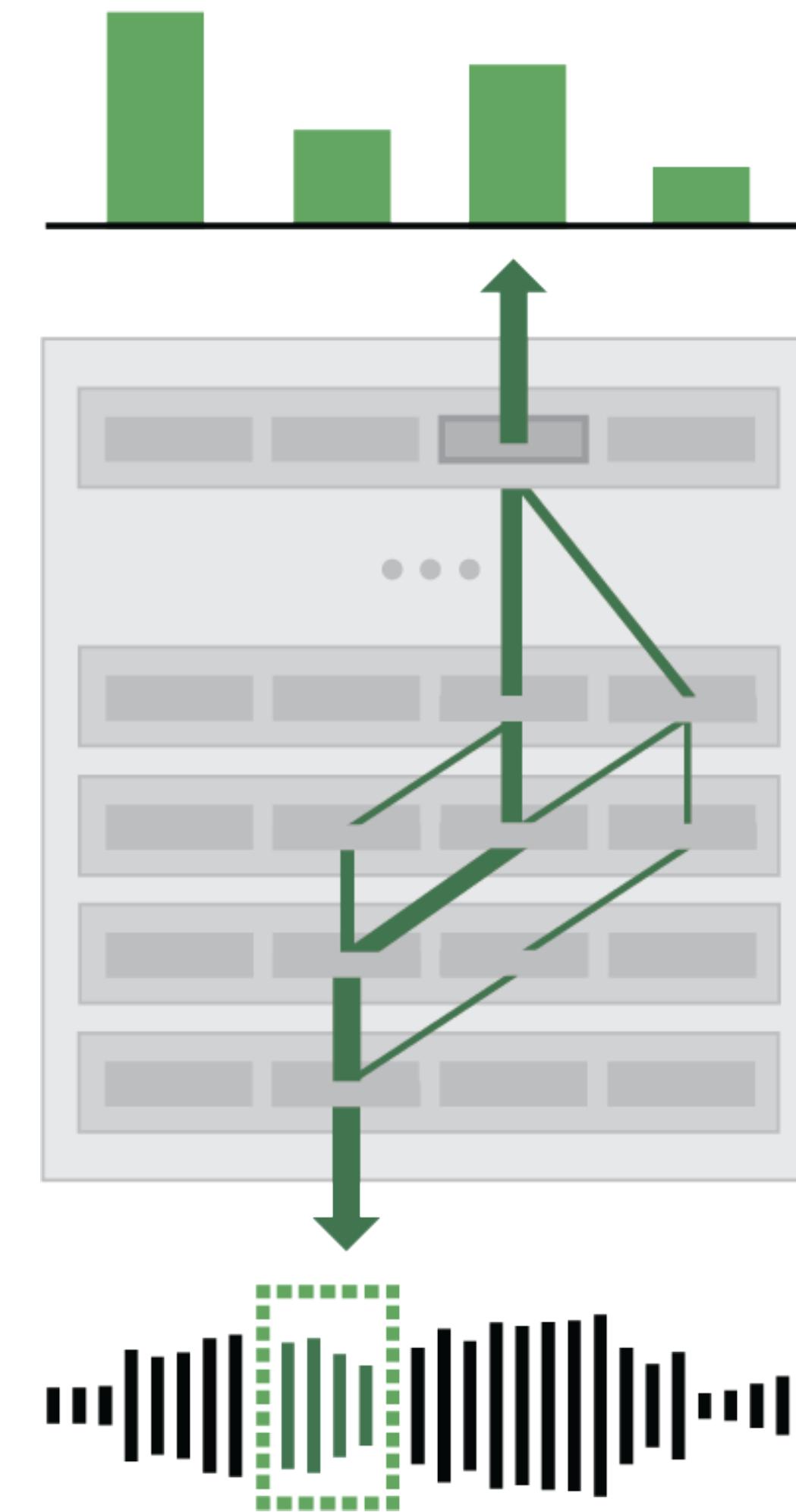




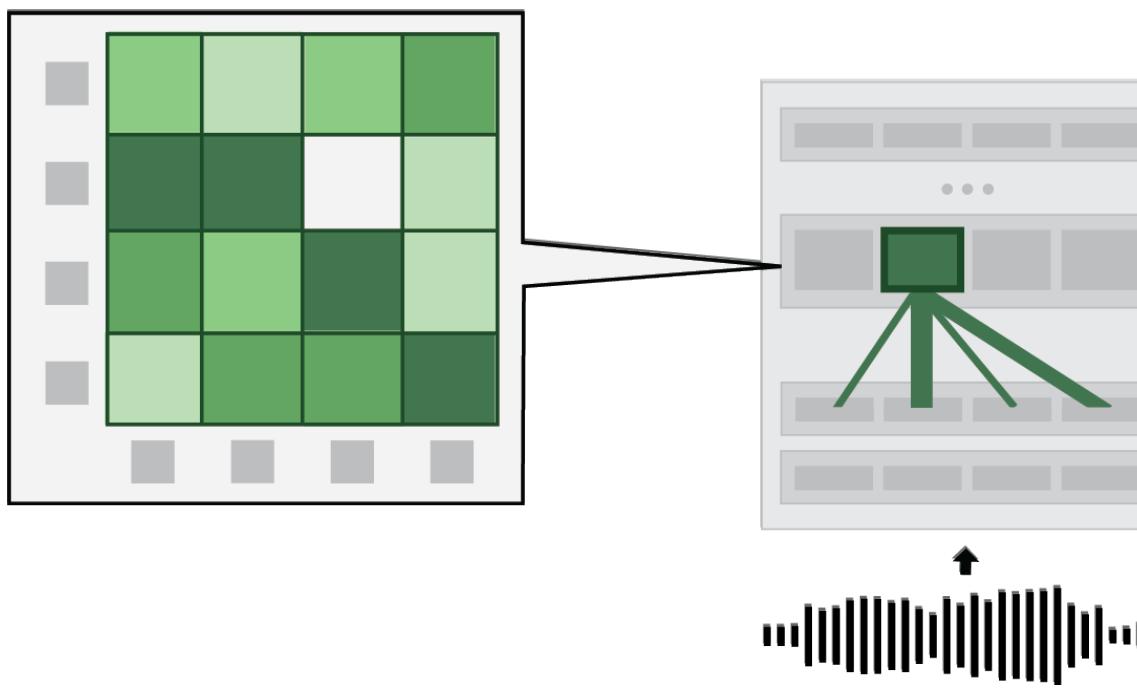
Feature Importance Scoring

- Context Mixing
- Feature Attribution

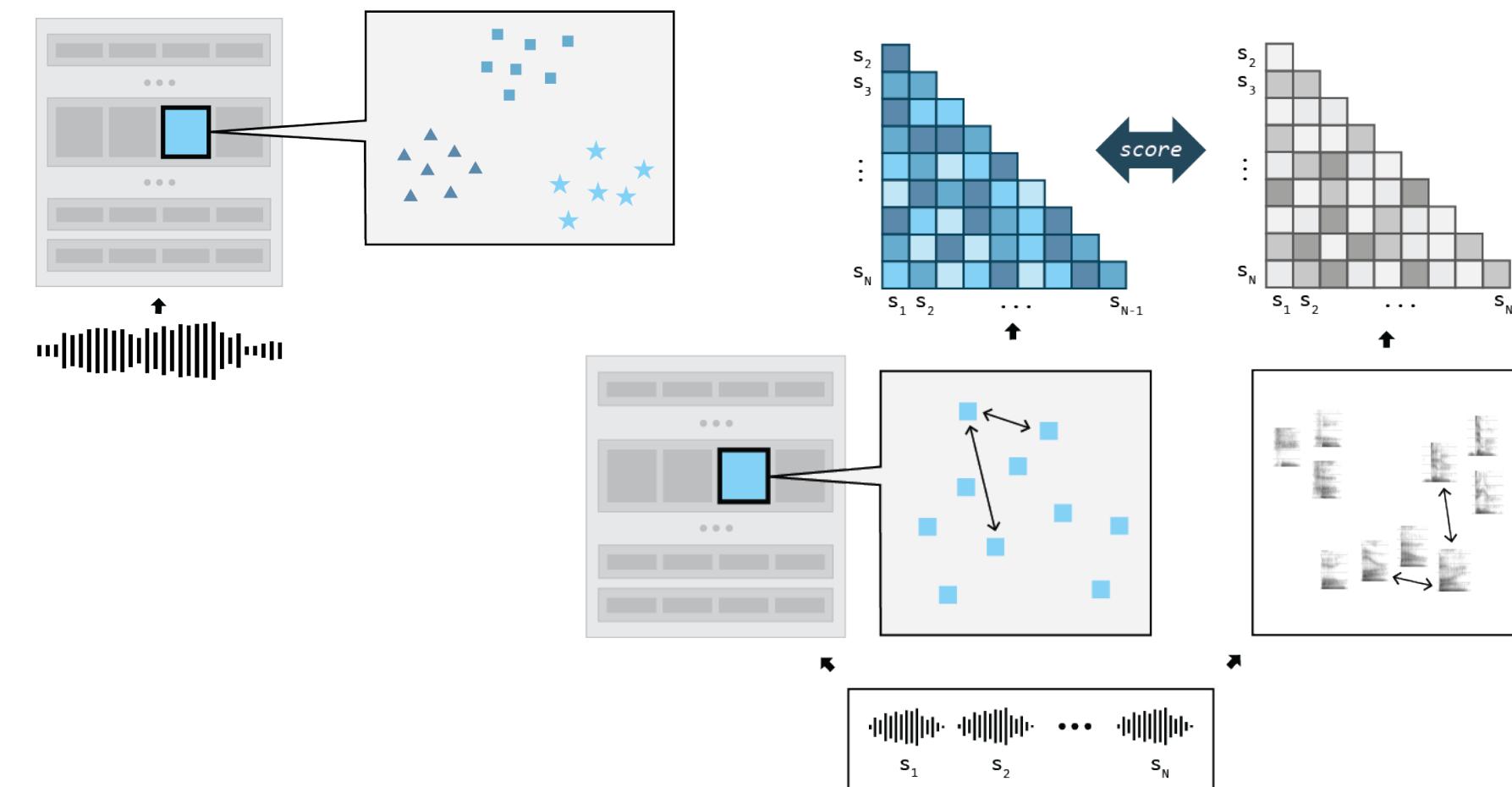
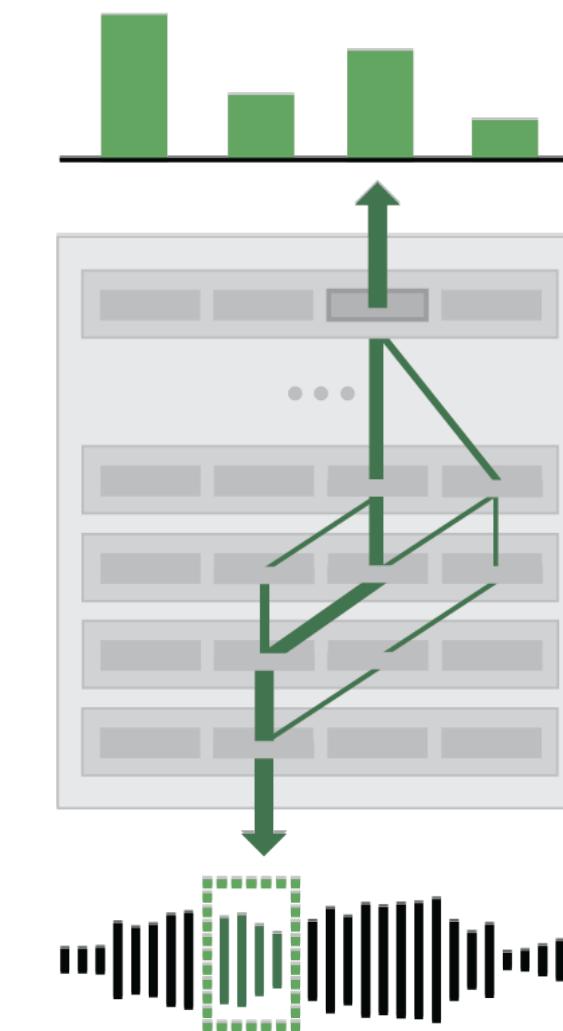
- Beyond attention weights
- Attribution reliability



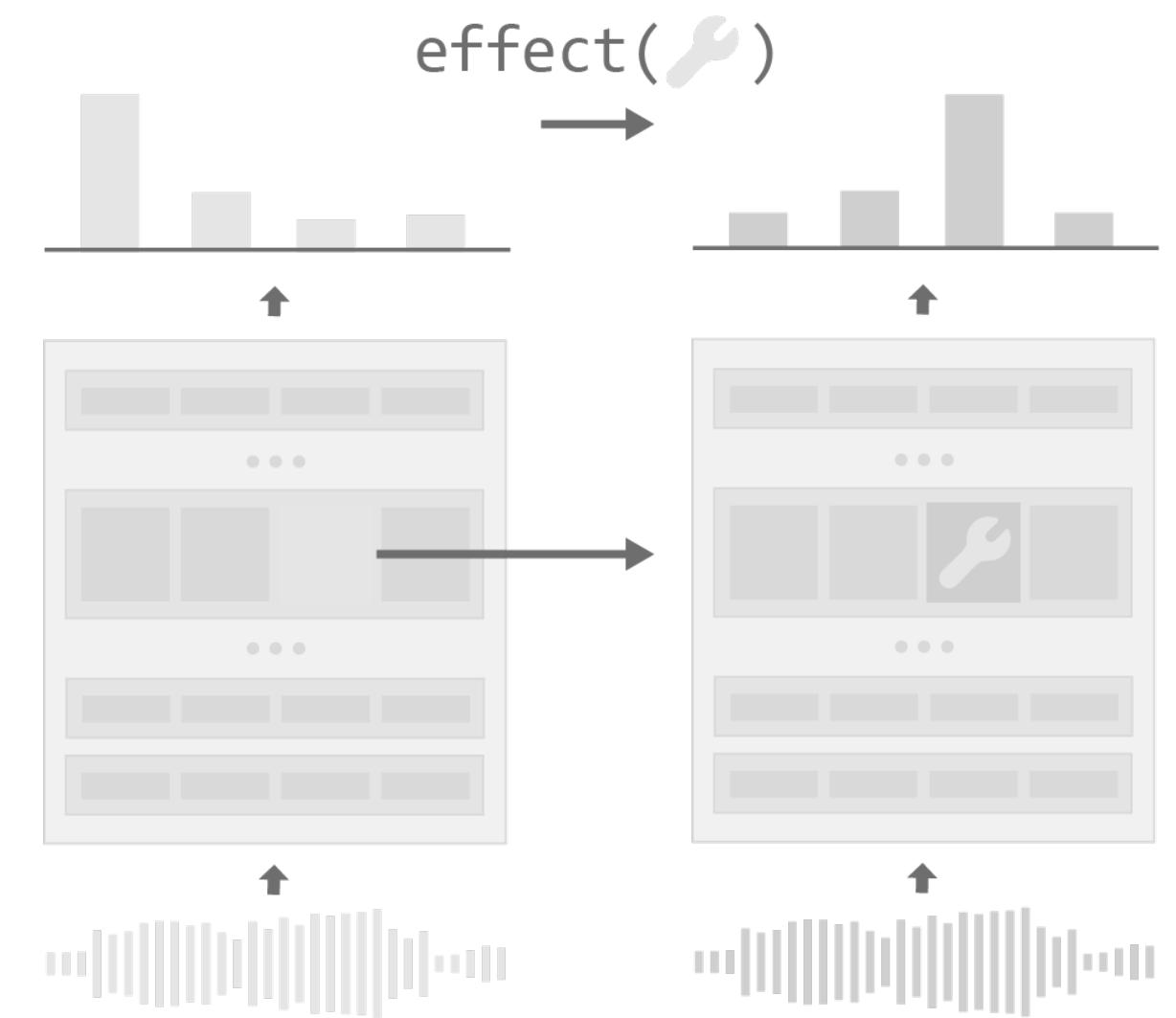
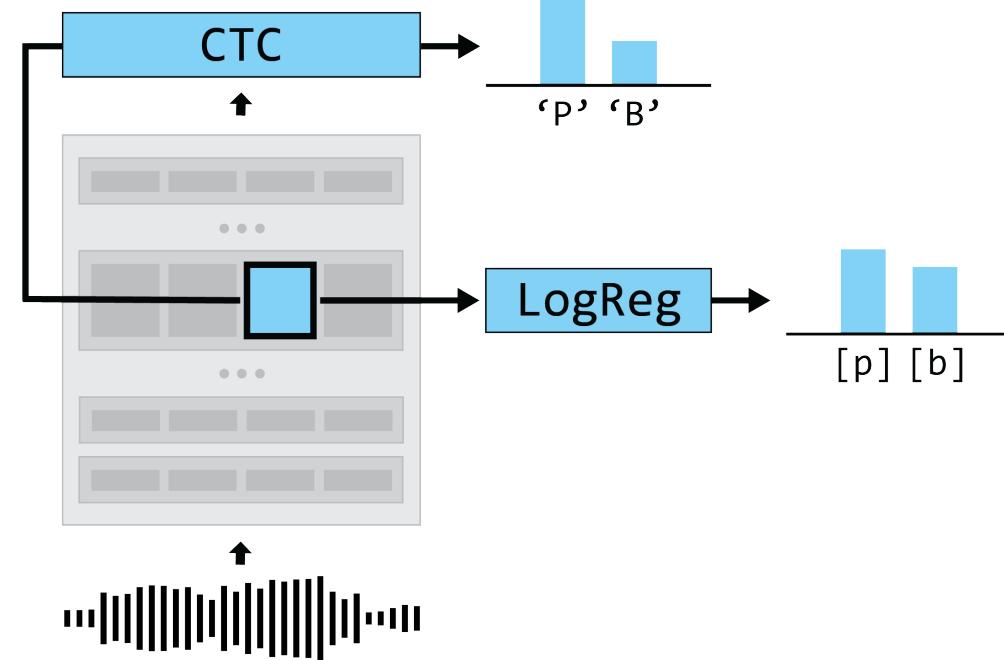
Interpretability Techniques for Speech Models

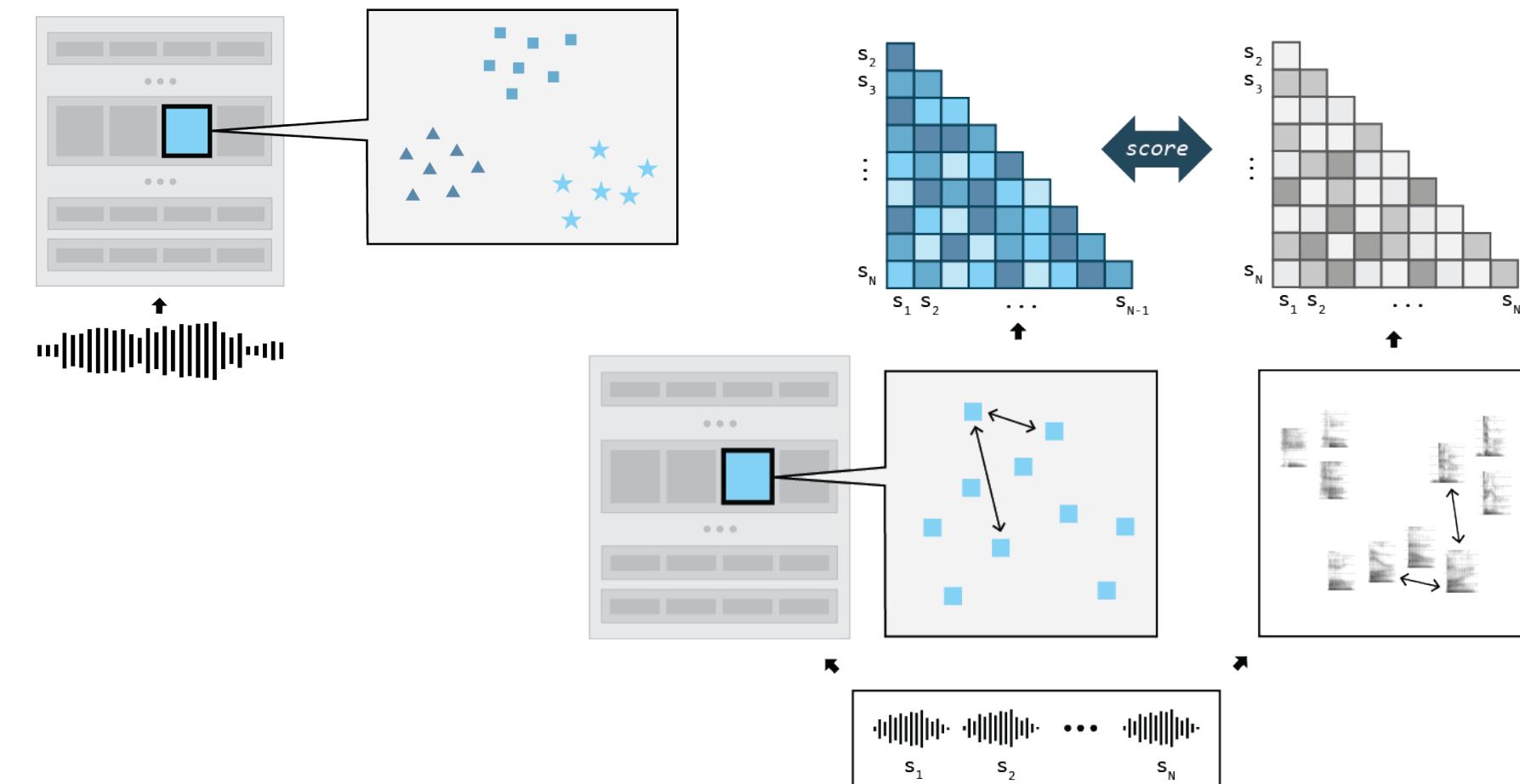
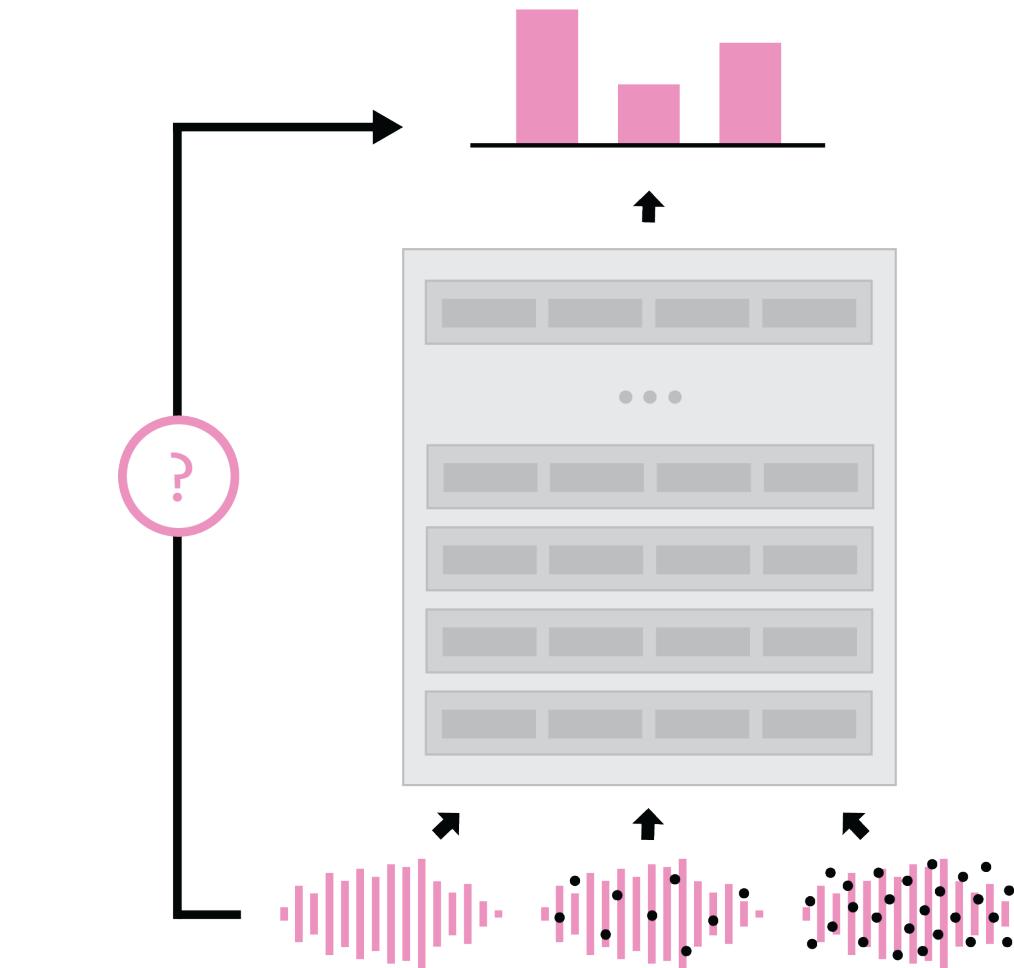


**Feature Importance
Scoring**

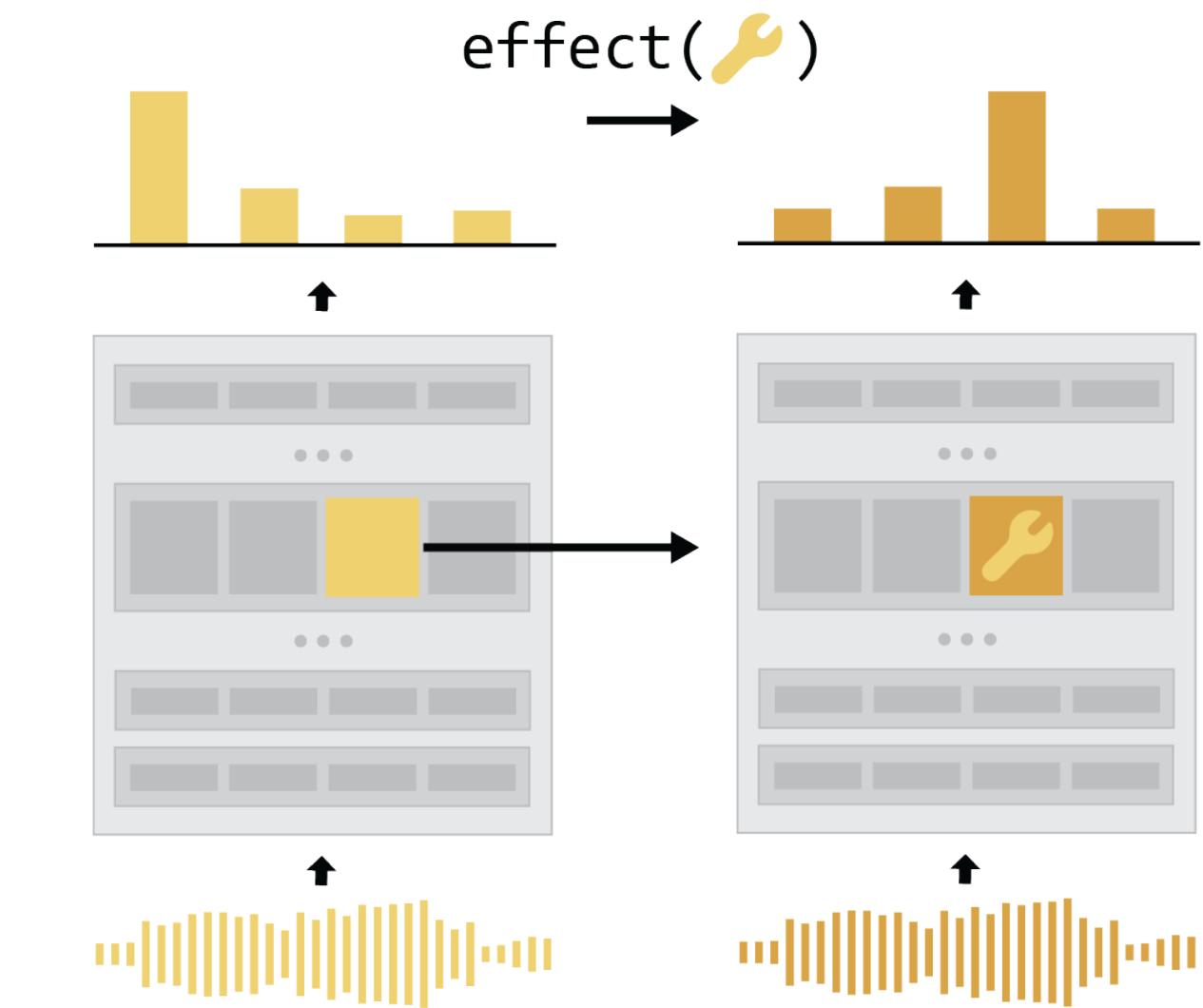
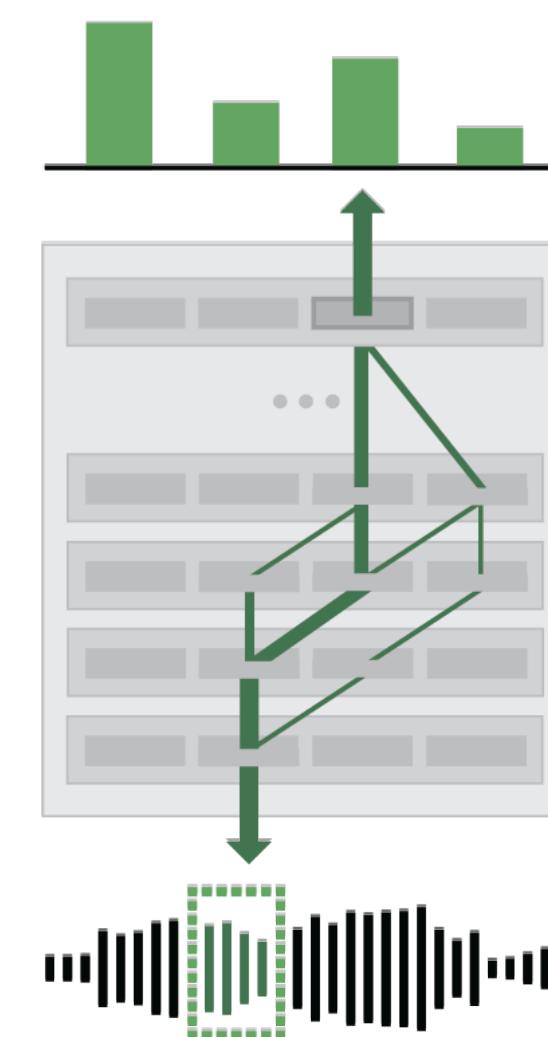
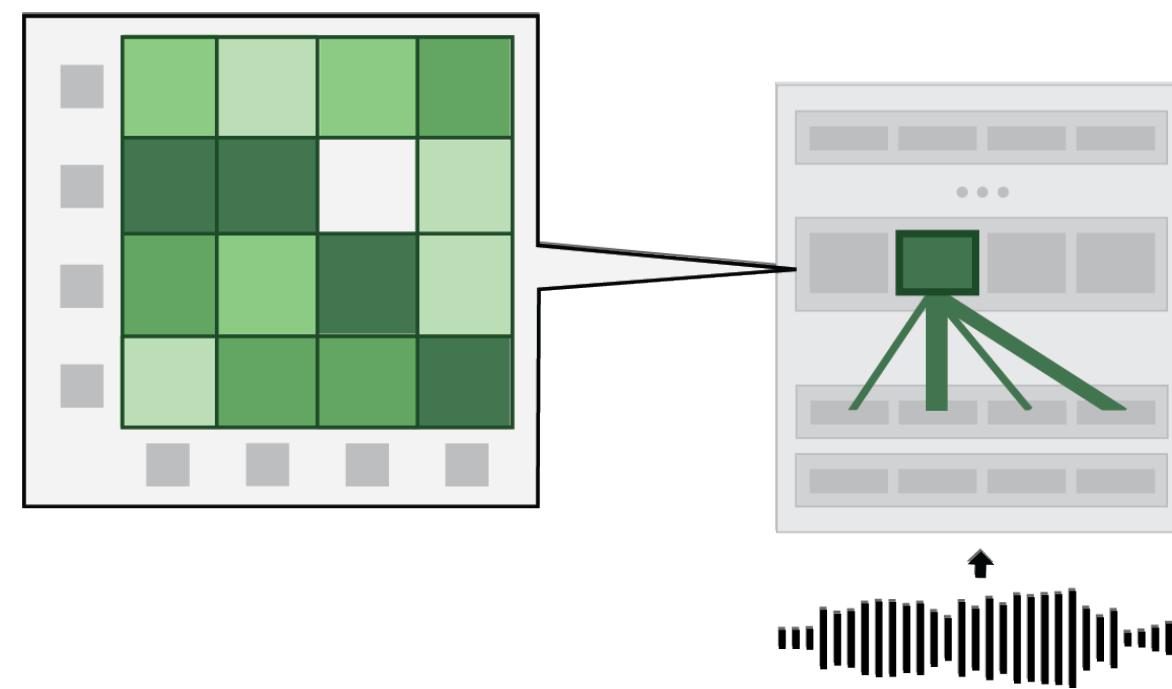


**Representational
analyses**





Interpretability Techniques for Speech Models



Outlook:

Can interpretability be useful?

Understanding model
training dynamics

Disentangling
representations &
attributions

Understanding links
between model
internals & output

Outlook:

Can interpretability be useful?

Understanding model
training dynamics

Disentangling
representations &
attributions

Understanding links
between model
internals & output

Understanding model training

Linguistic structure in speech models

Shen et al. (2024)

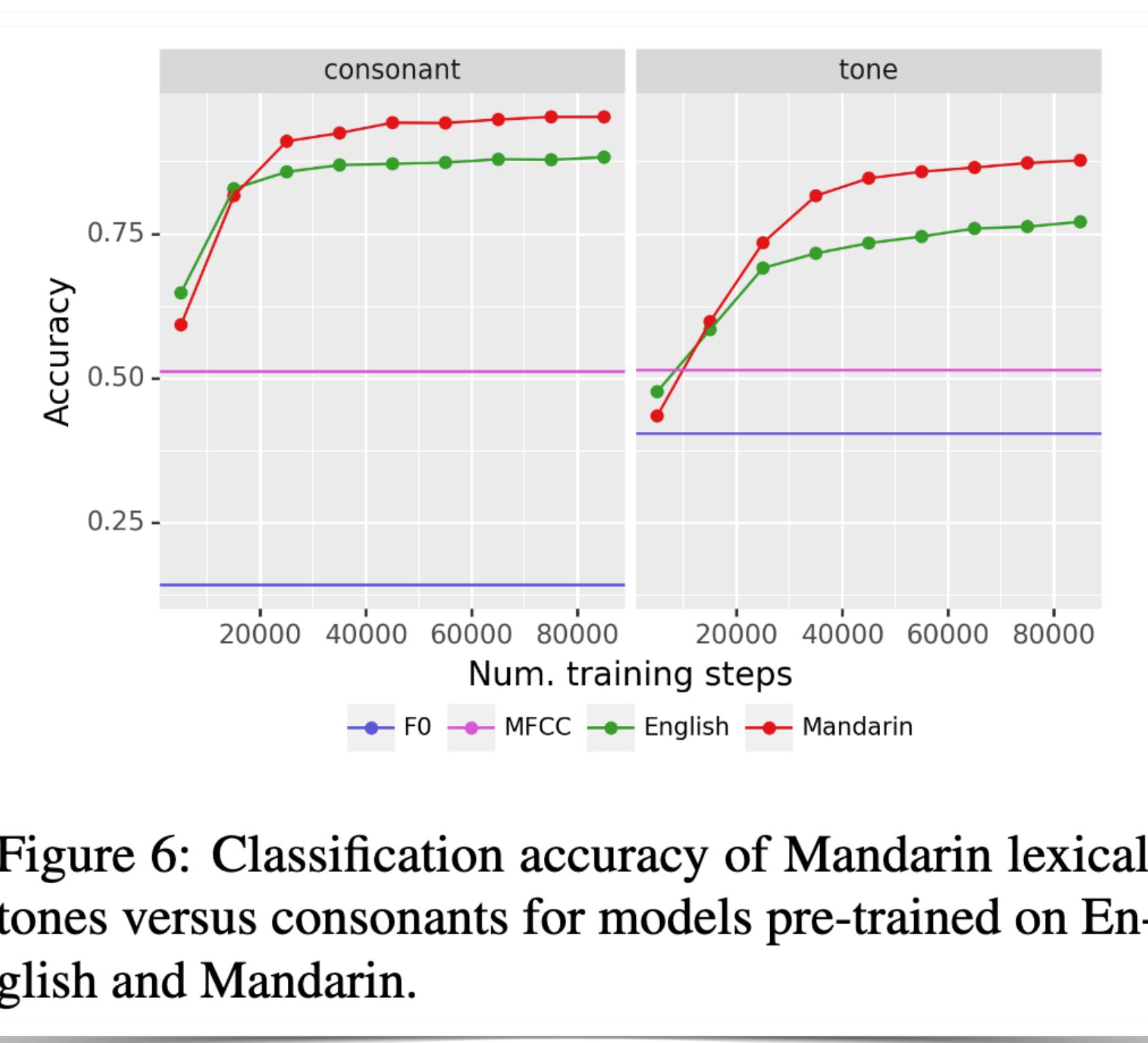


Figure 6: Classification accuracy of Mandarin lexical tones versus consonants for models pre-trained on English and Mandarin.

de Heer Kloots et al. (in prep)

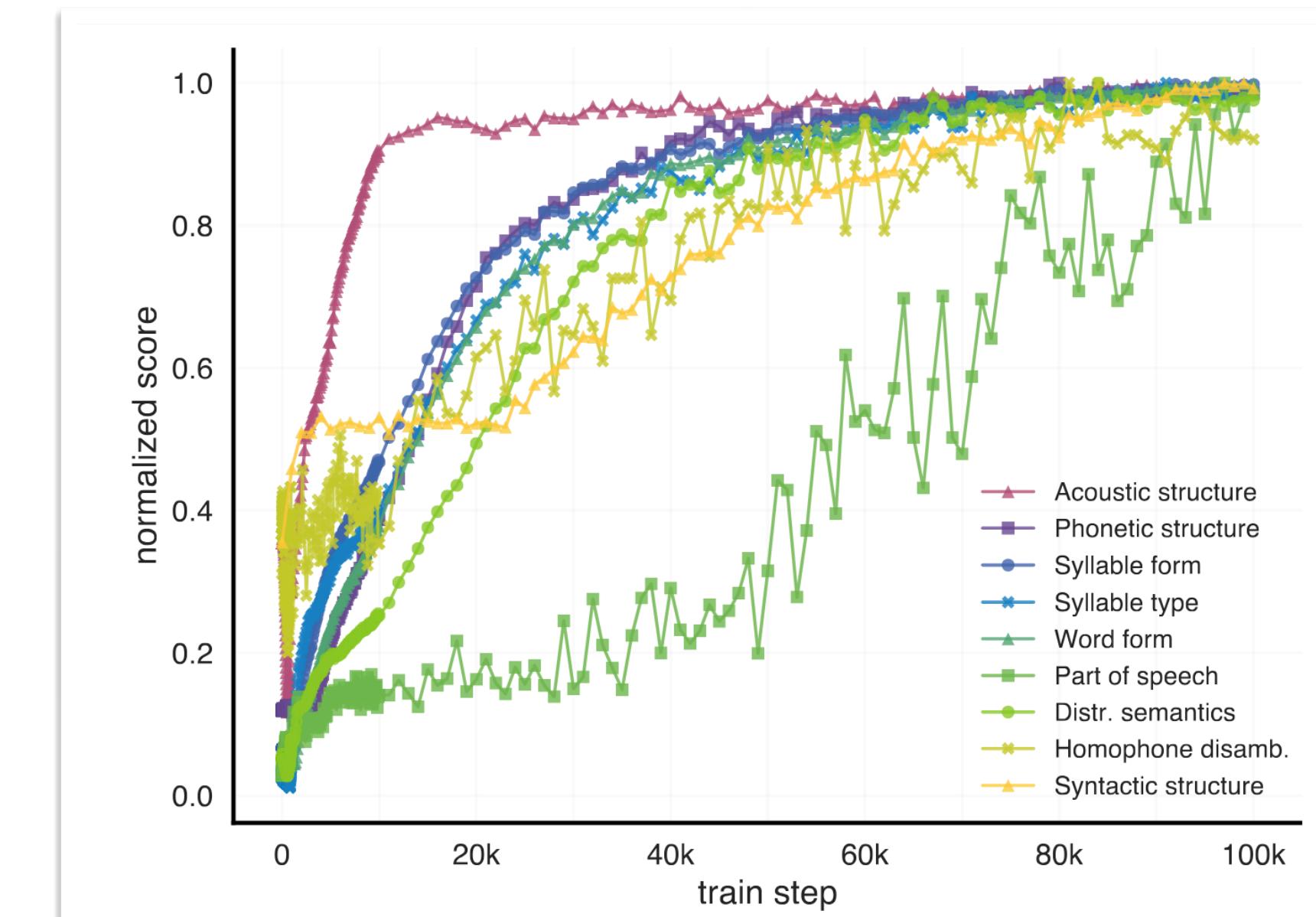


Figure 3: Development of linguistic encoding capacities across checkpoints in model training.

Understanding model training

Syntax acquisition in text models (Chen et al., 2025)

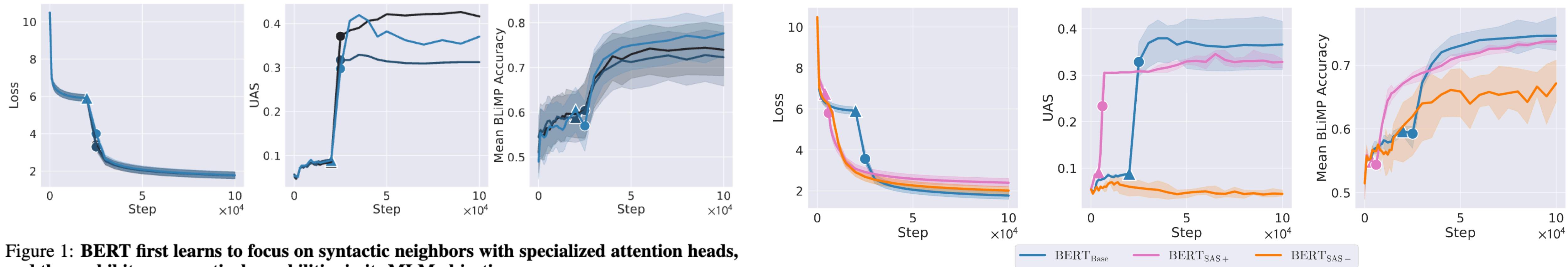


Figure 1: BERT first learns to focus on syntactic neighbors with specialized attention heads, and then exhibits grammatical capabilities in its MLM objective.

Outlook:

Can interpretability be useful?

Understanding model
training dynamics

Disentangling
representations &
attributions

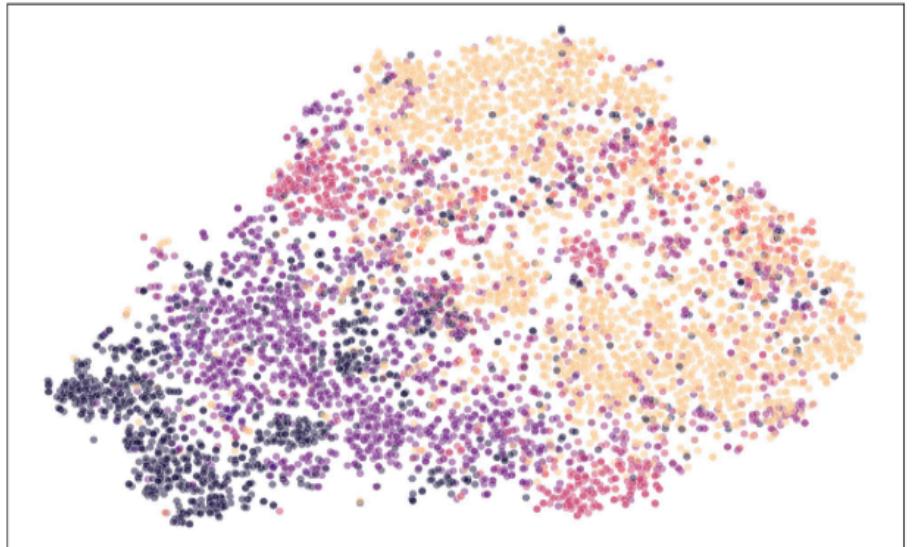
Understanding links
between model
internals & output

Disentangling representations & attribution

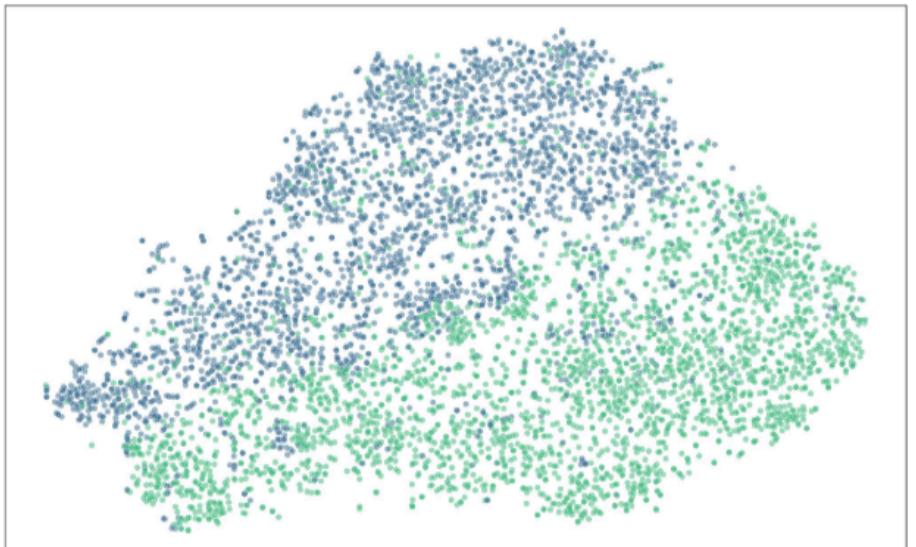
Speech model representations encode different kinds of interpretable information in an *entangled* way

- Can we disentangle such dimensions for specific tasks?
- Can we attribute model predictions to such interpretable dimensions?

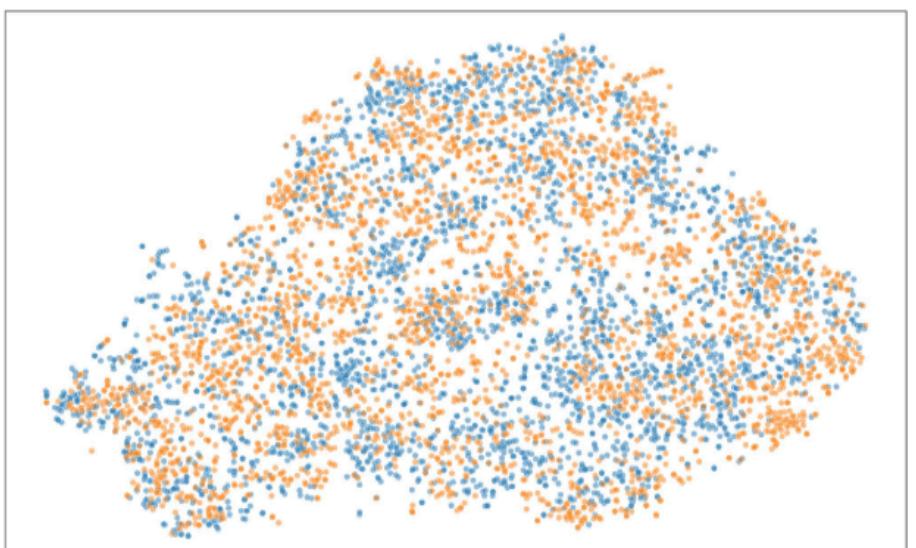
phone
class



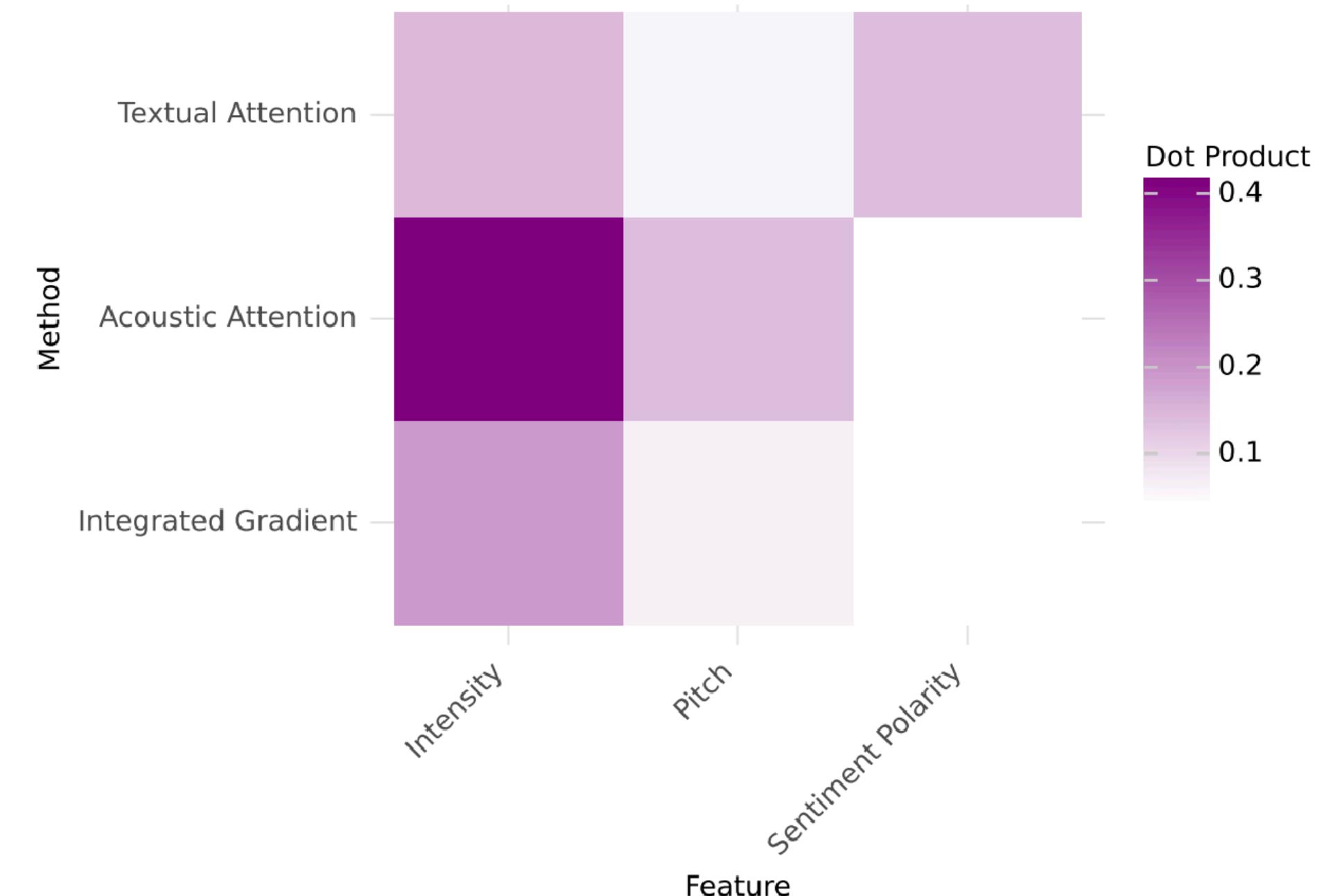
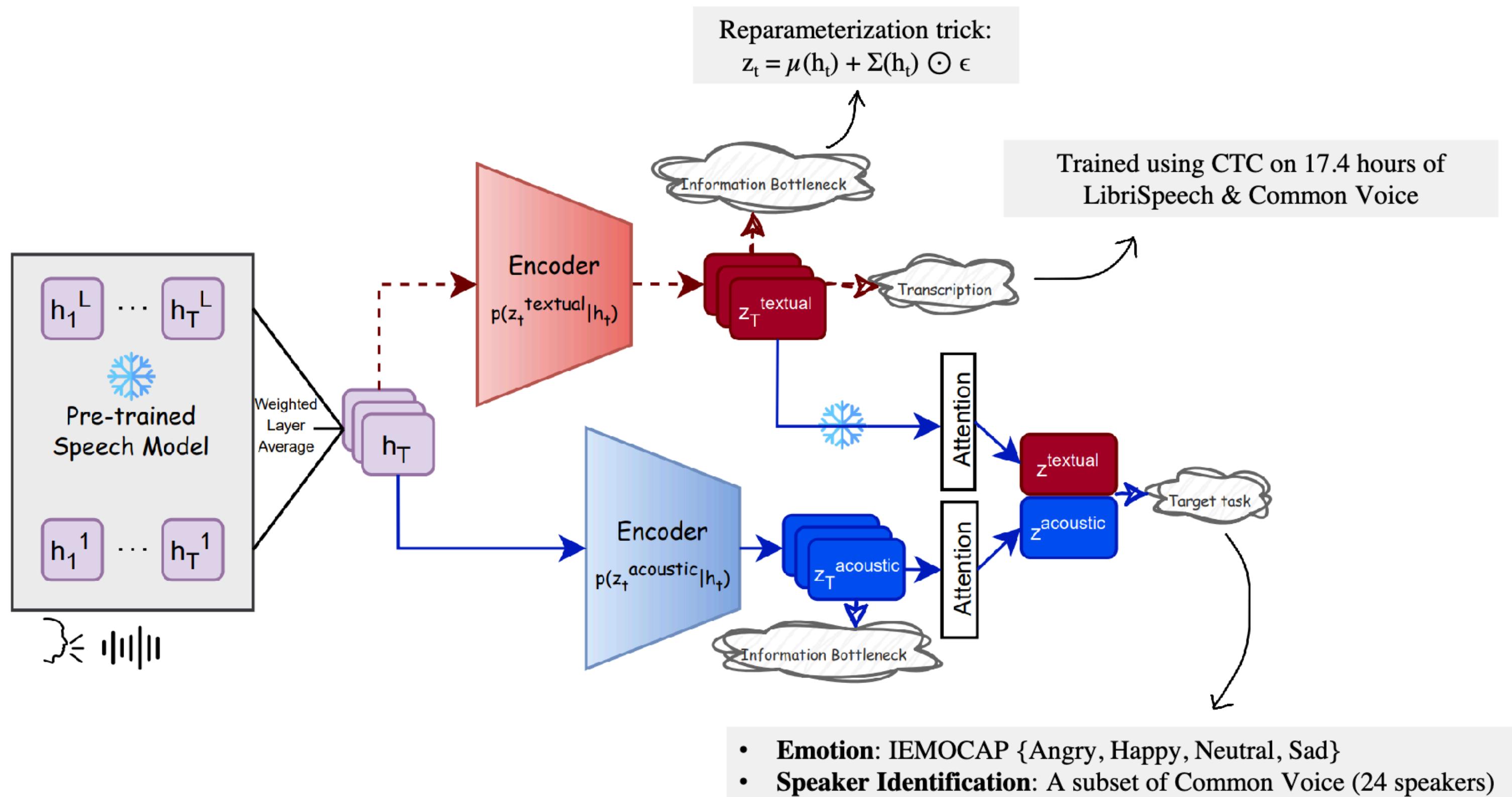
gender



language



Disentangling representations & attribution



Outlook:

Can interpretability be useful?

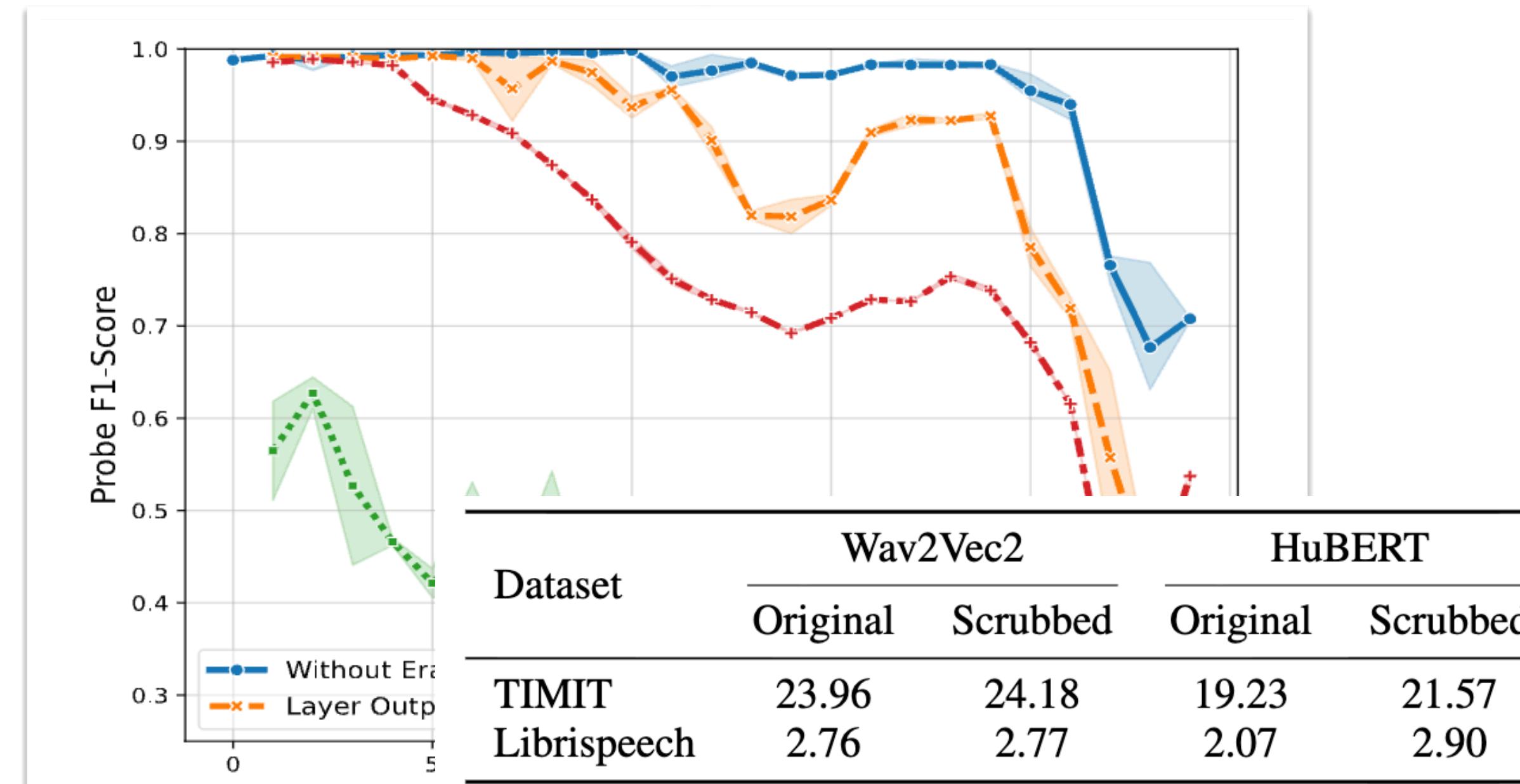
Understanding model
training dynamics

Disentangling
representations &
attributions

Understanding links
between model
internals & output

Understanding links between model internals and behaviour

Gender encoding & bias



(a)

Table 2: WER before and after gender scrubbing

Figure 1: Gender scrubbing for the Wav2Vec2 ASR model.

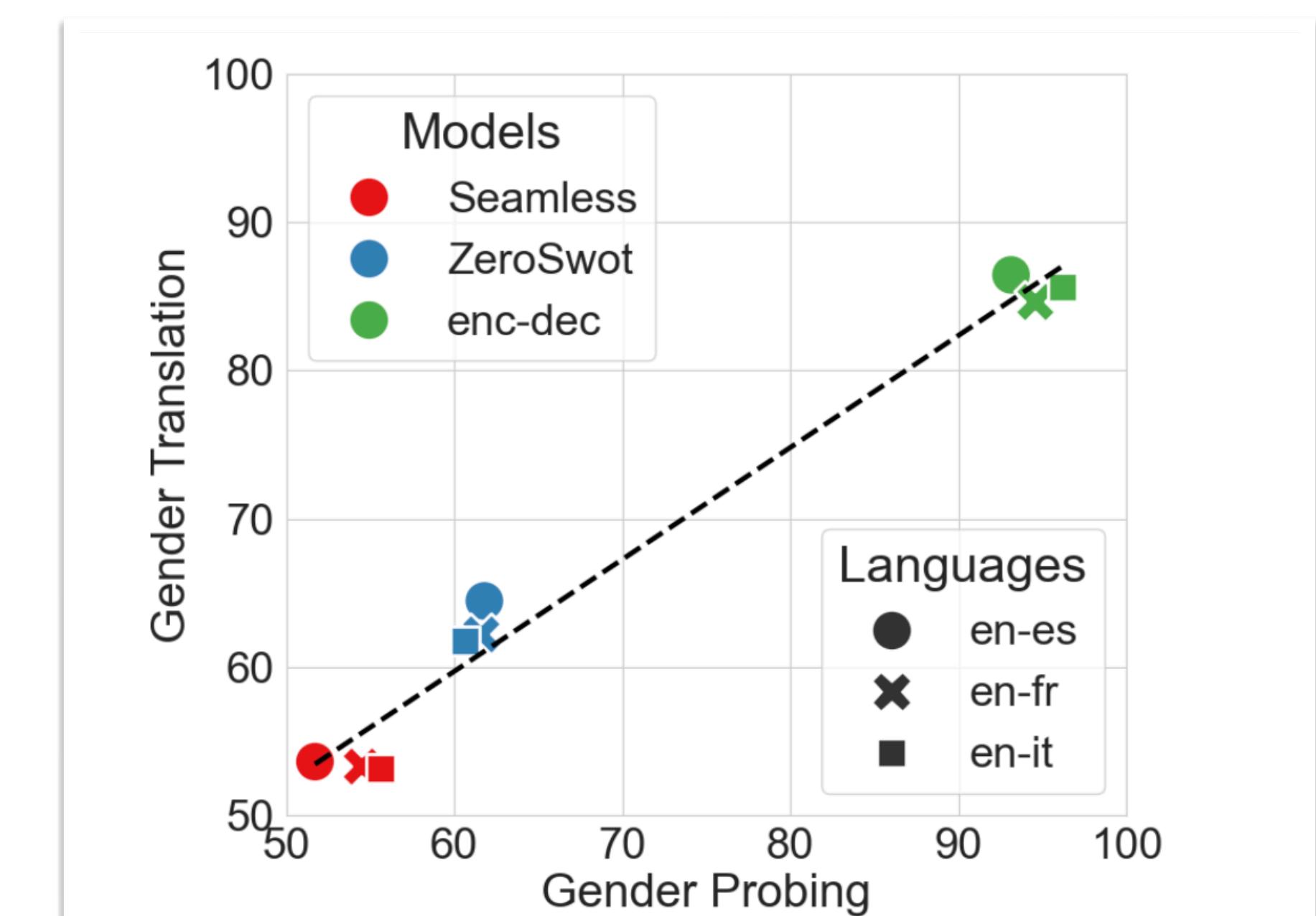


Figure 1: Correlation between overall gender probing performance (macro F1) and gender translation accuracy across models and languages on test-speaker.

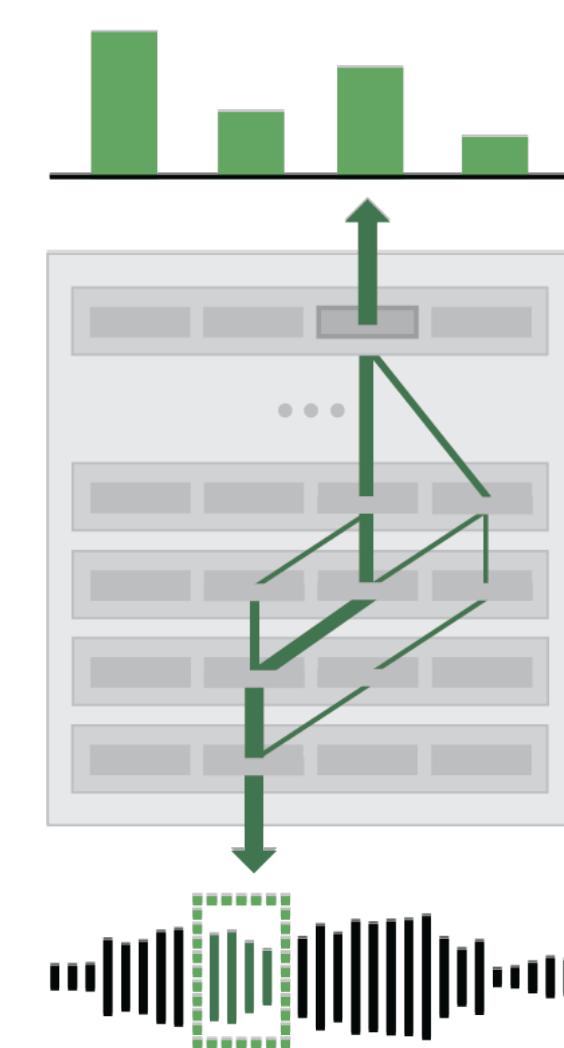
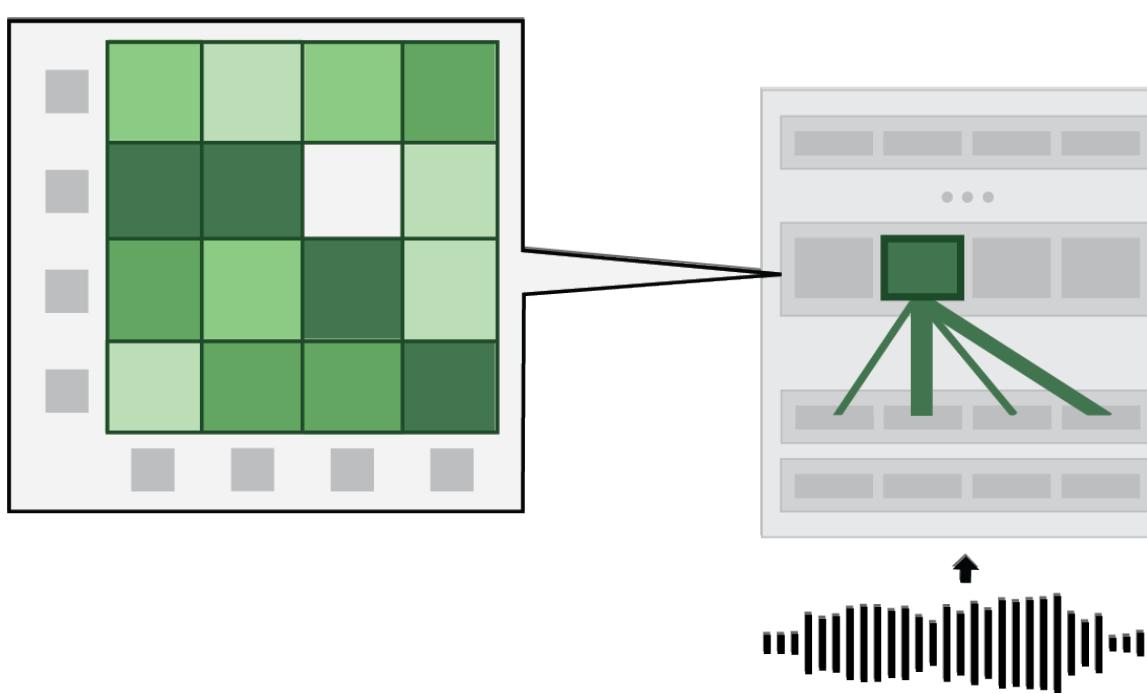
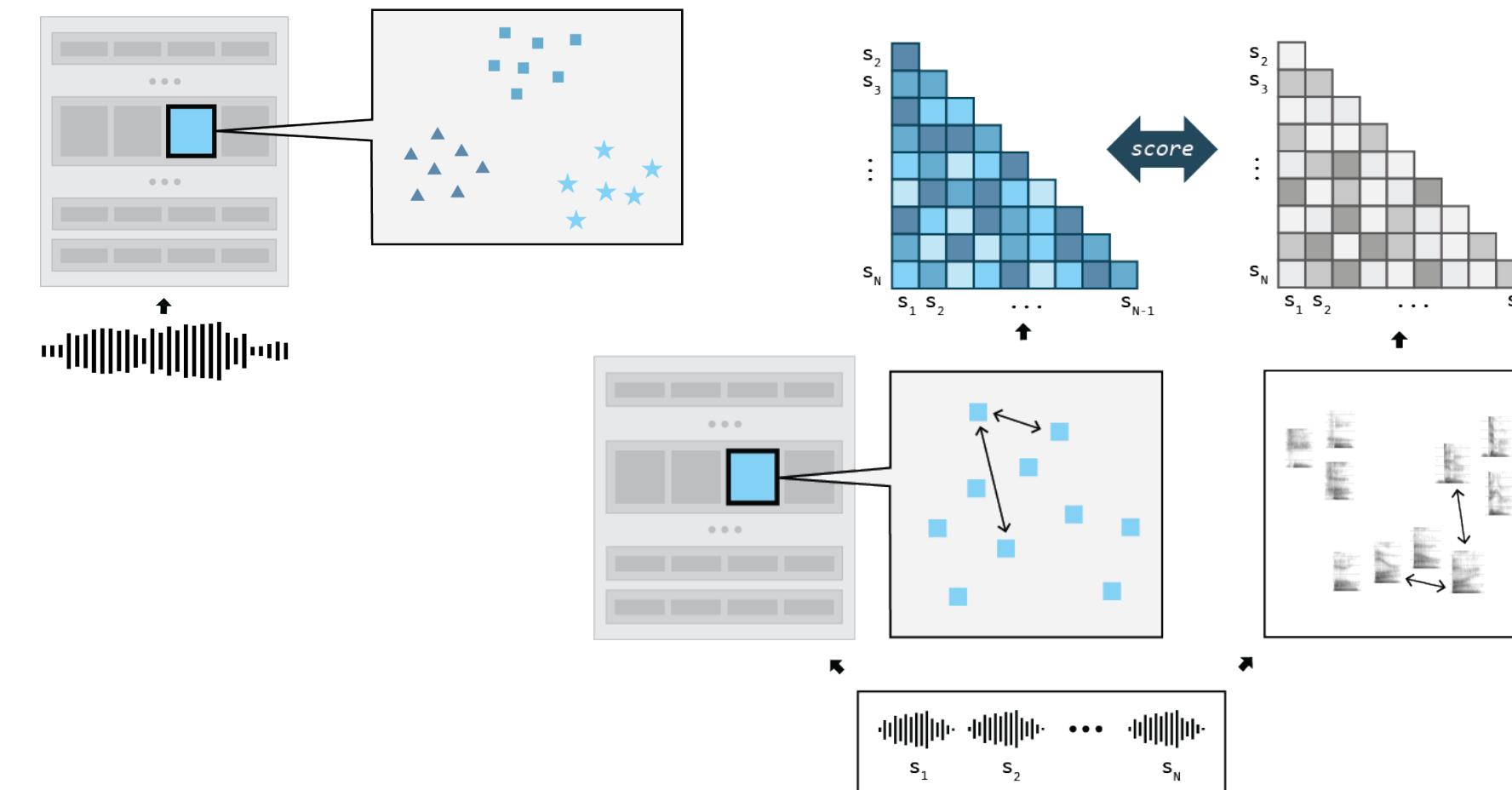
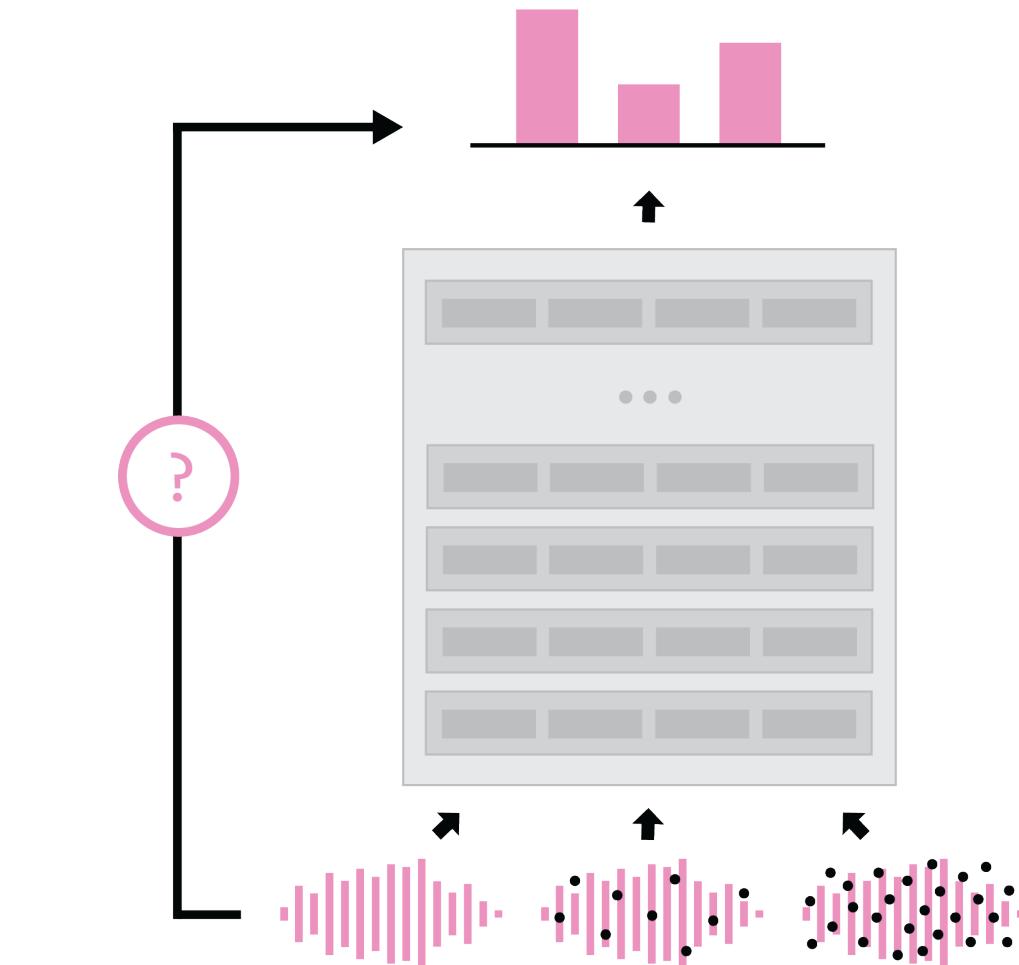
Outlook:

Can interpretability be useful?

Understanding model
training dynamics

Disentangling
representations &
attributions

Understanding links
between model
internals & output



Figures:
Marianne de Heer Kloots (2025)



Interpretability Techniques for Speech Models