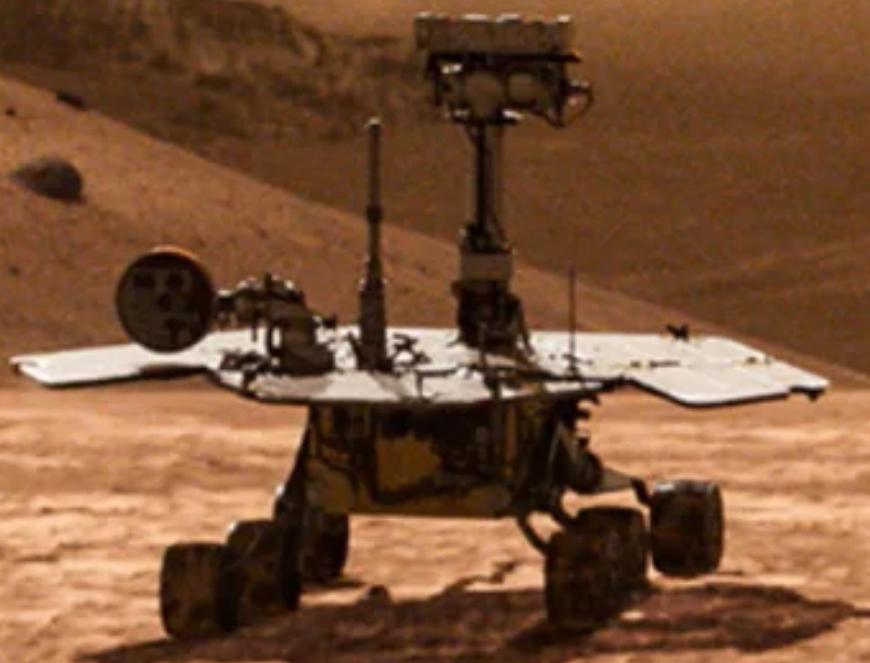
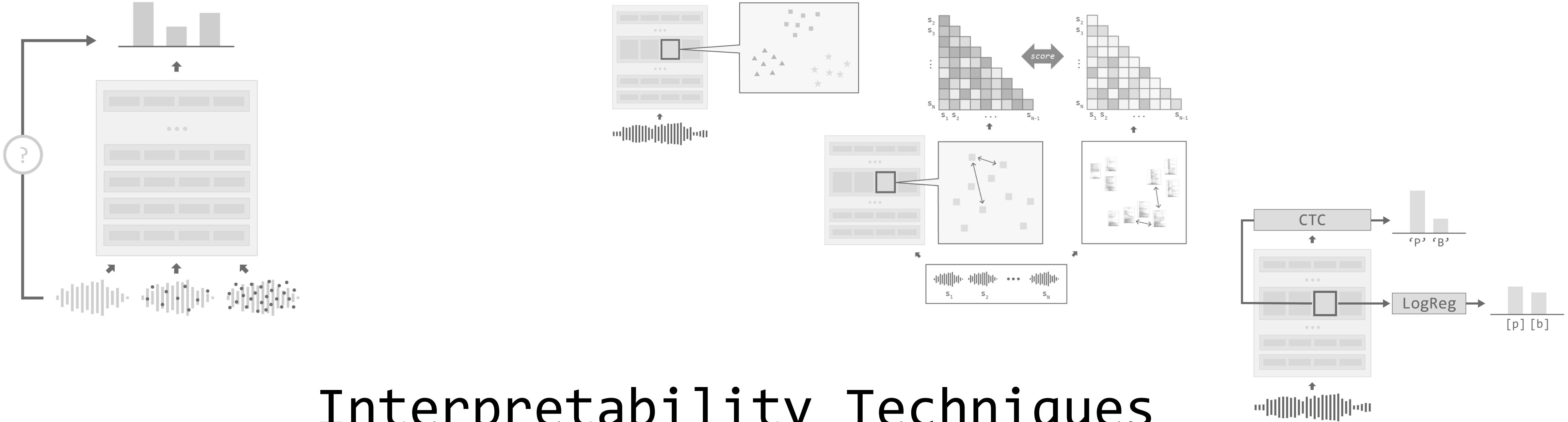


Interpretability Techniques for Speech Models

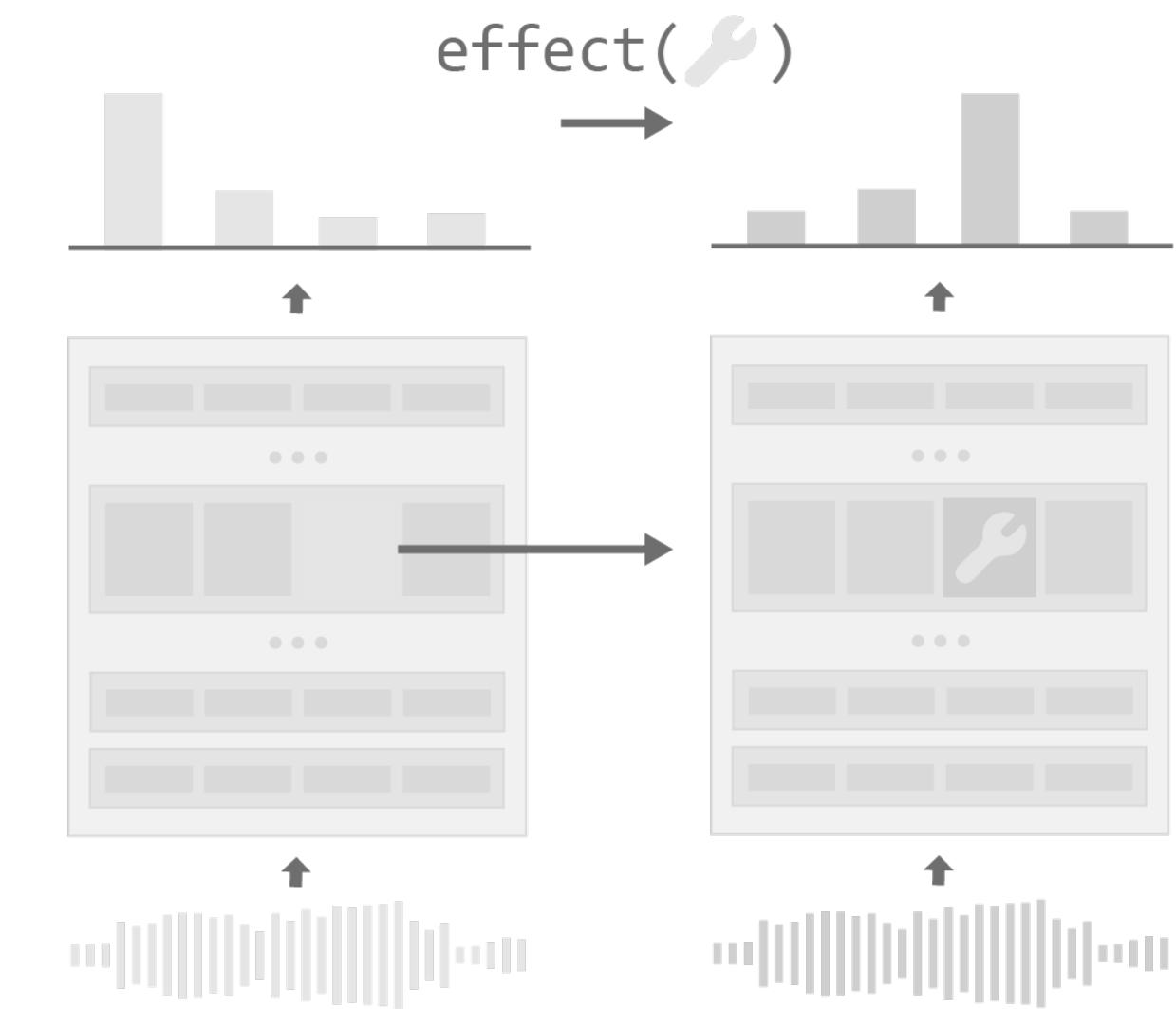
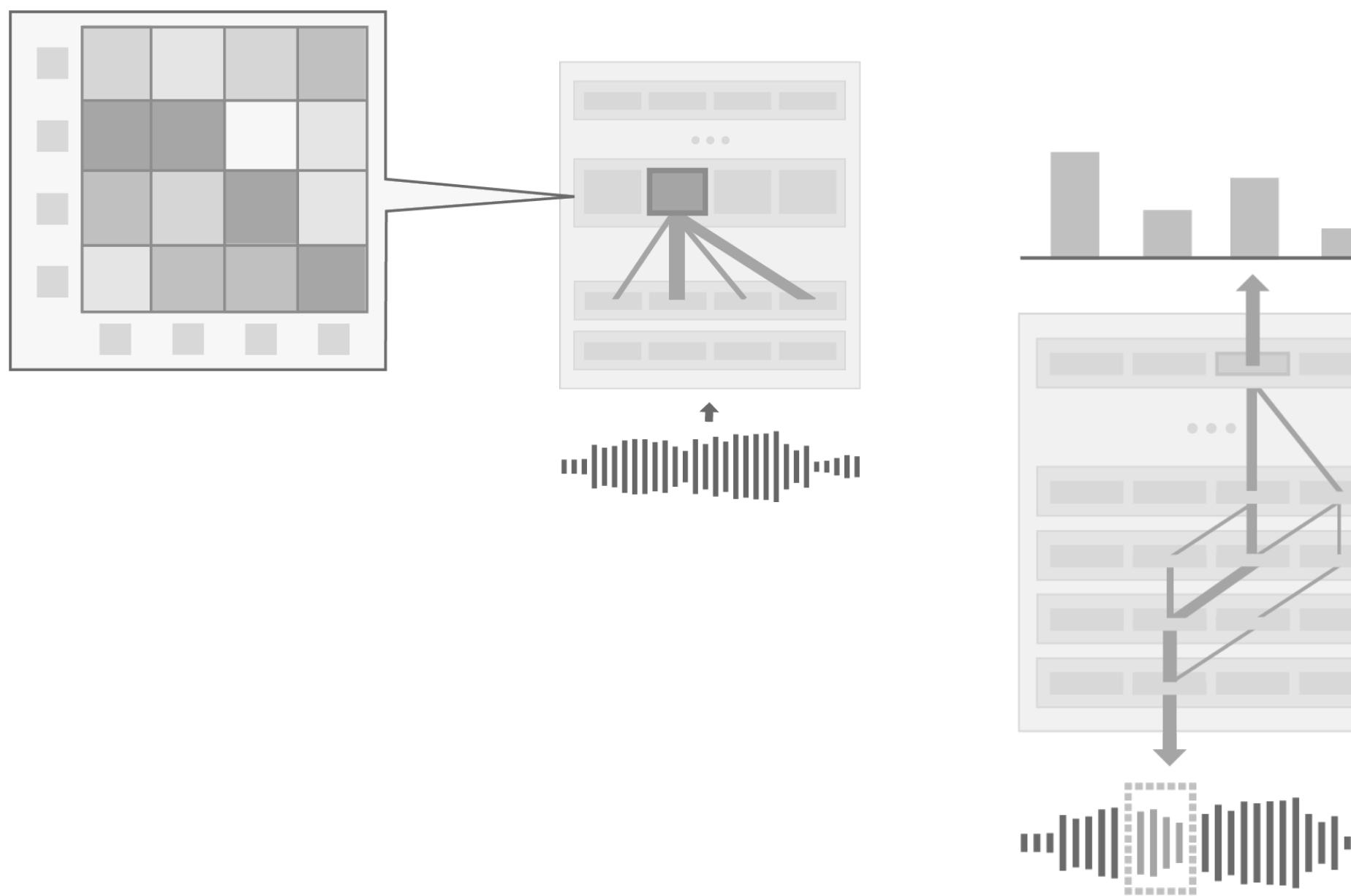
Conclusions & Outlook

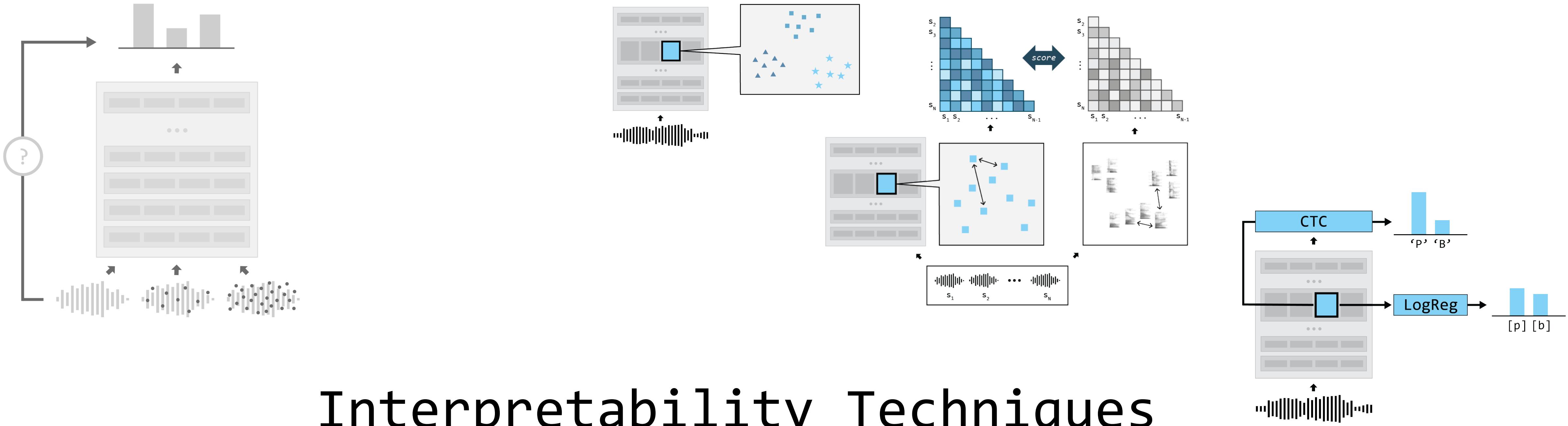


Marianne de Heer Kloots, 17-08-2025

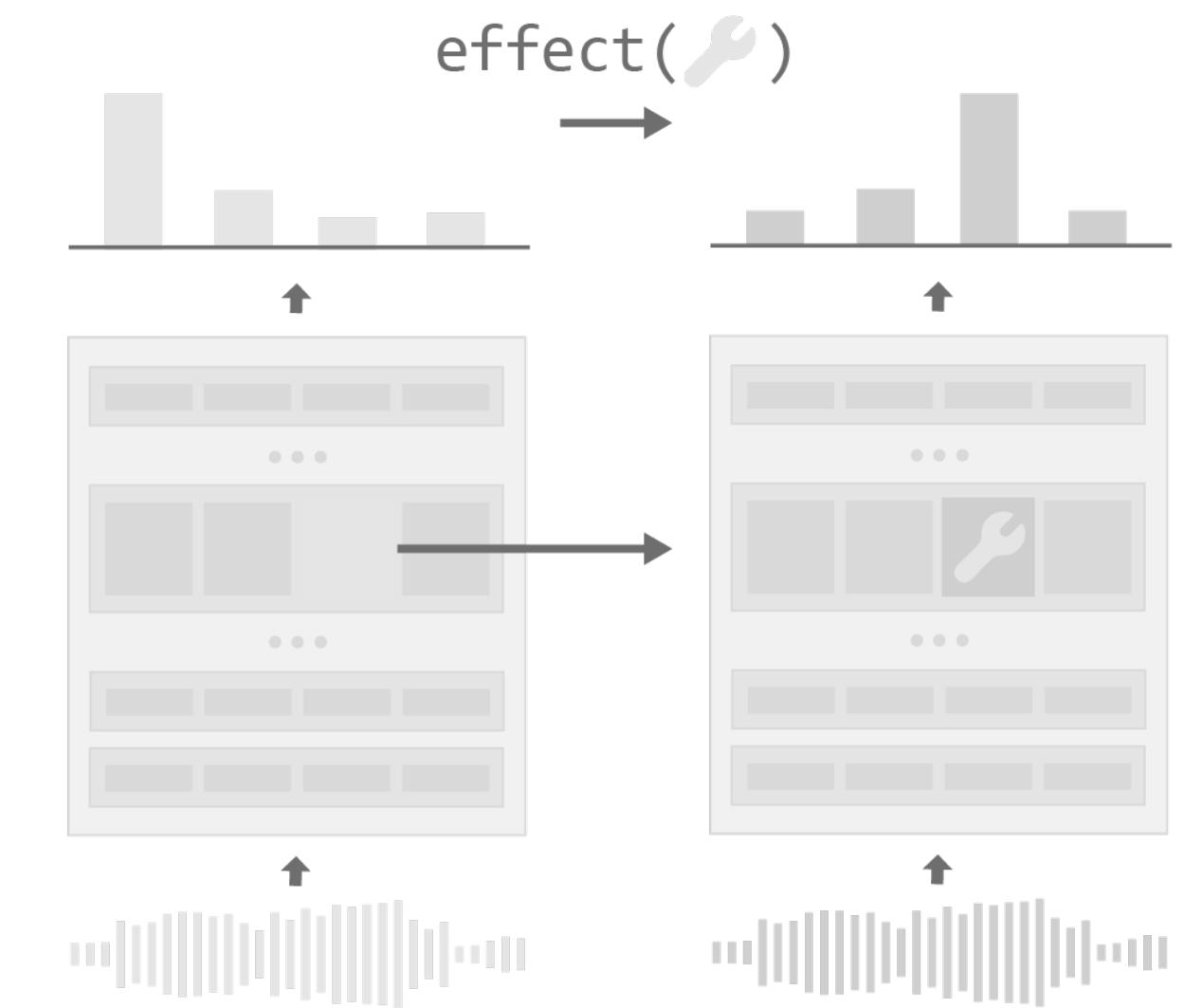
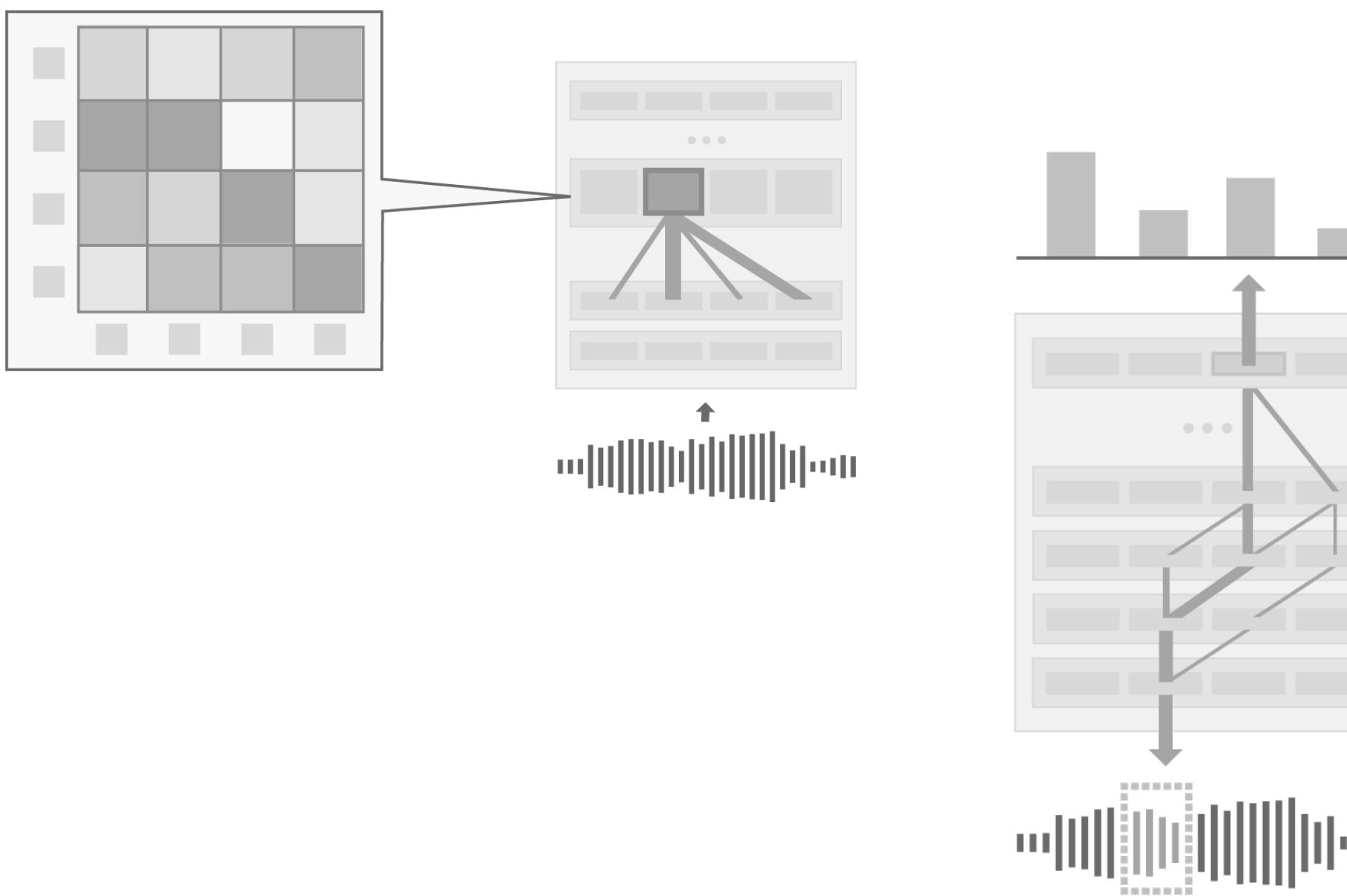


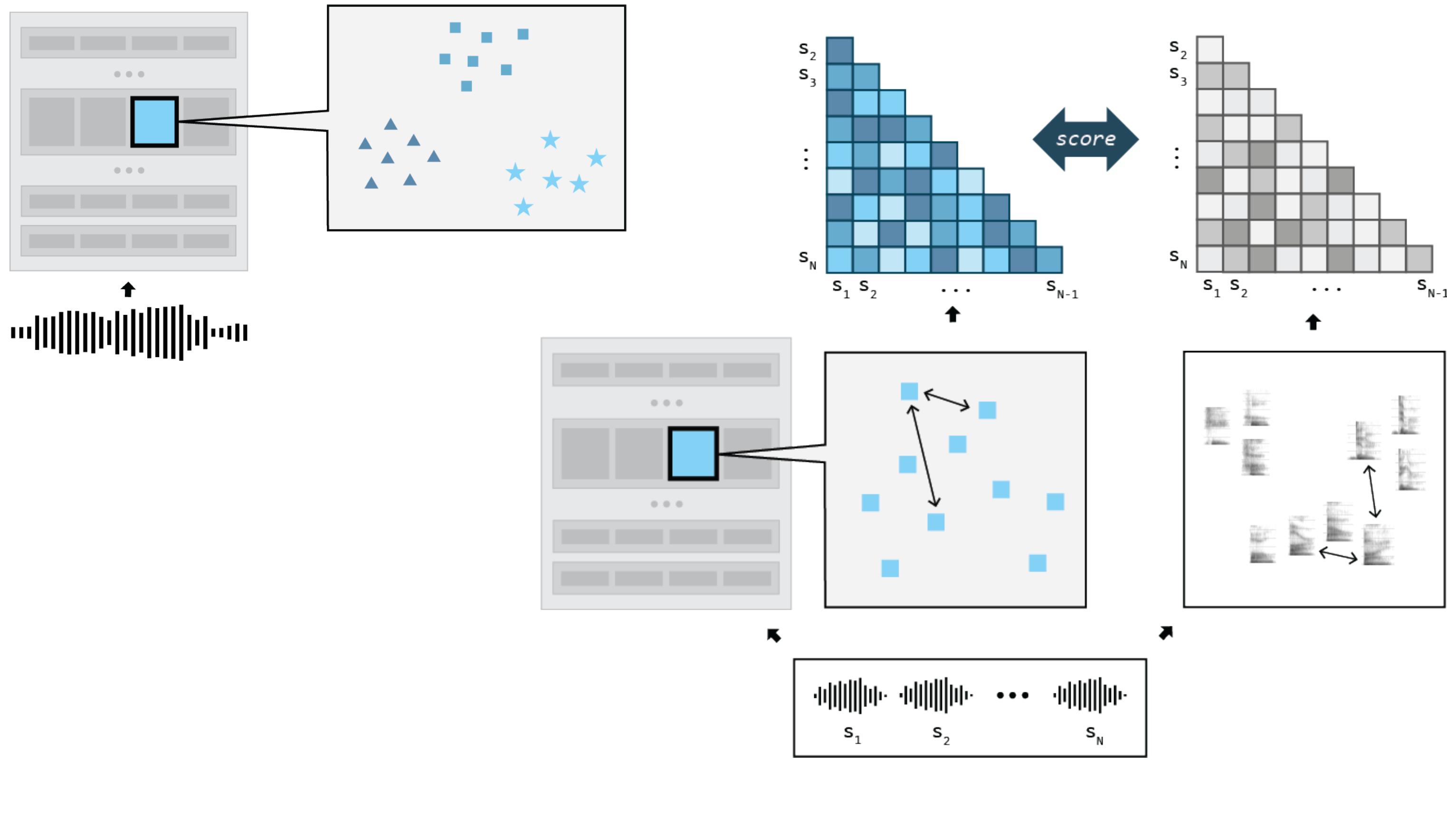
Interpretability Techniques for Speech Models





Interpretability Techniques for Speech Models

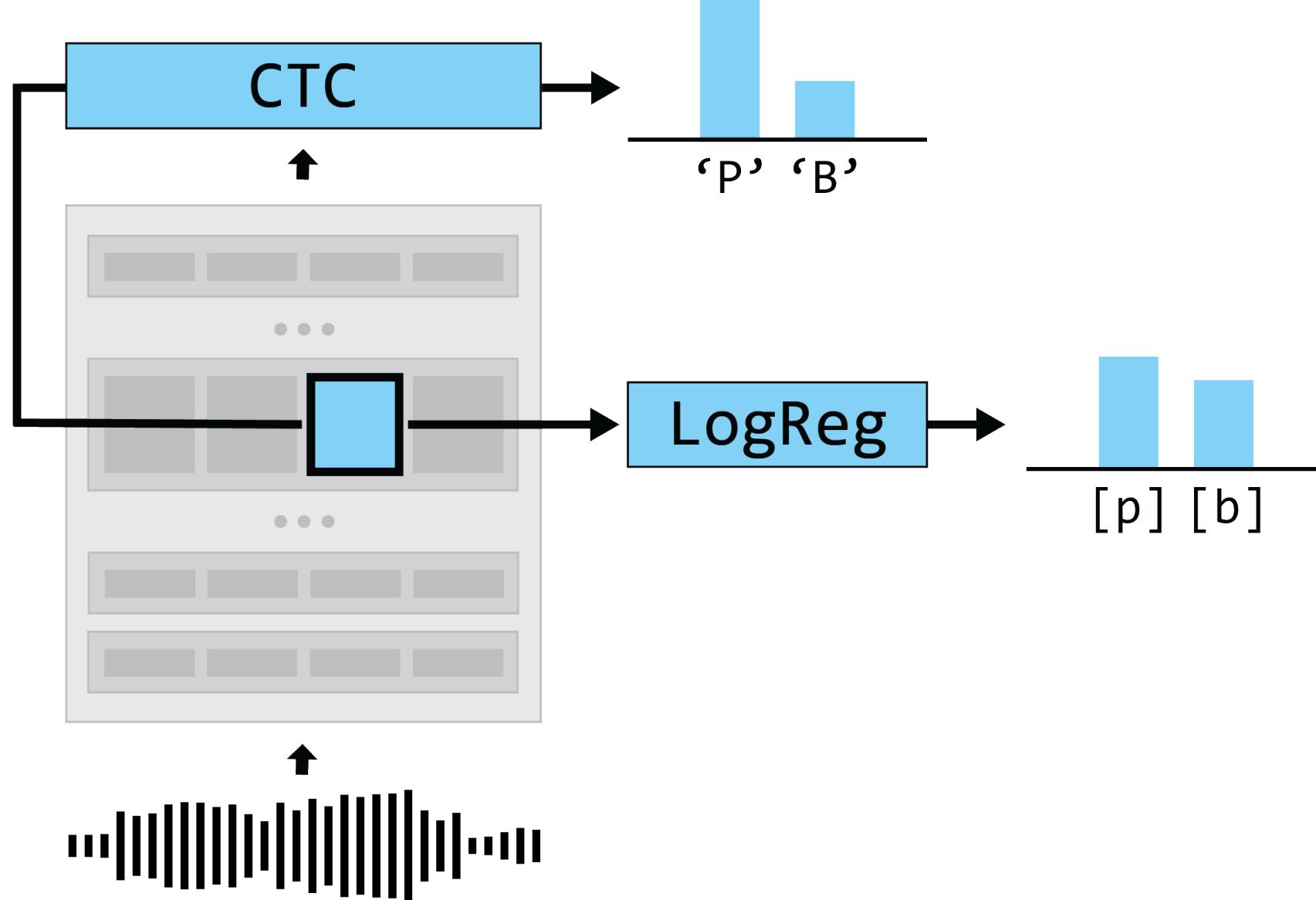




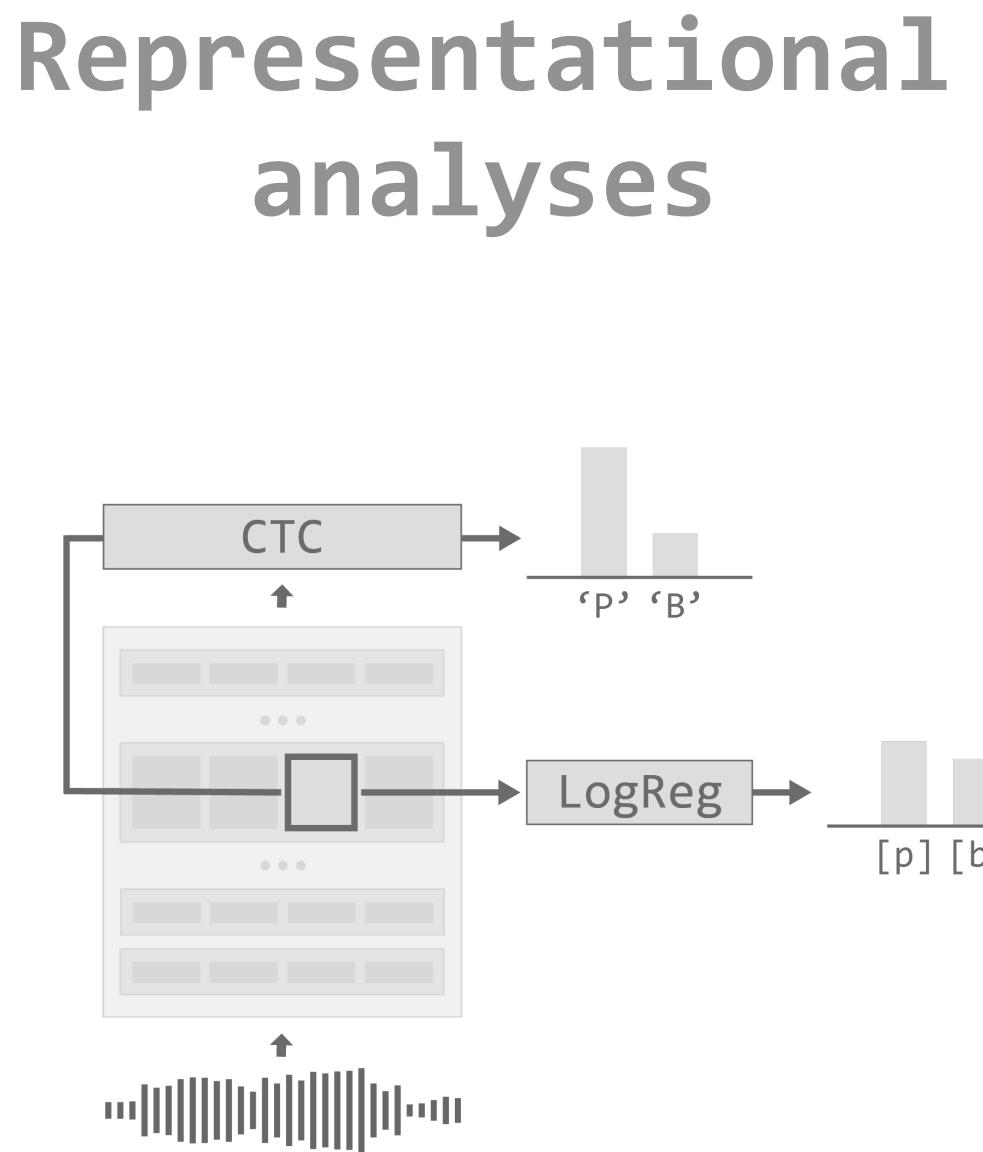
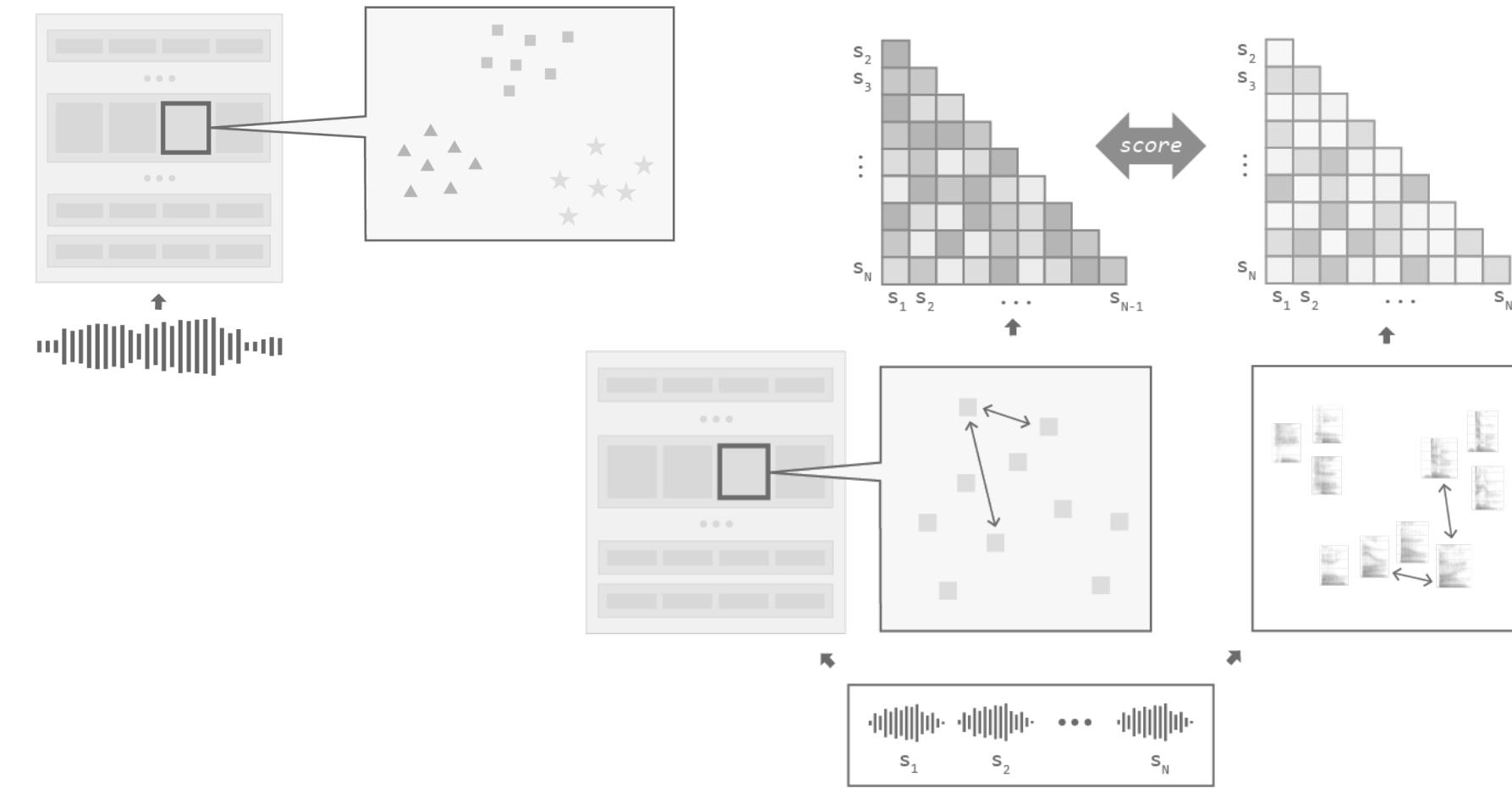
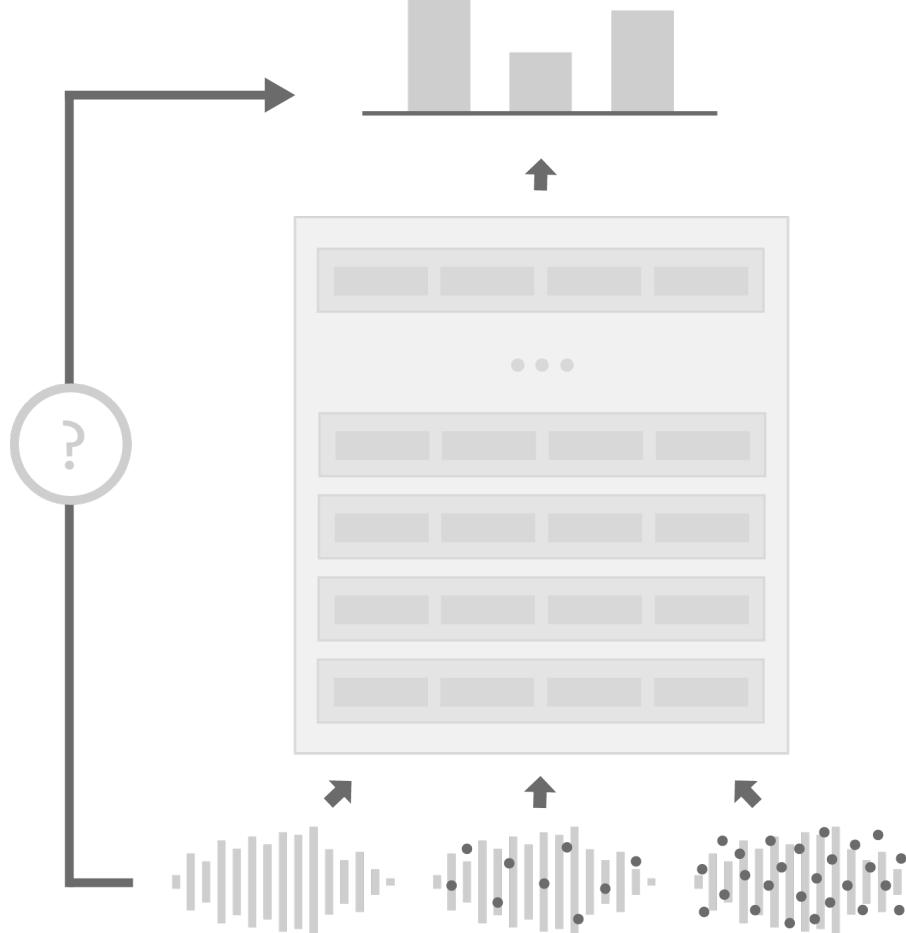
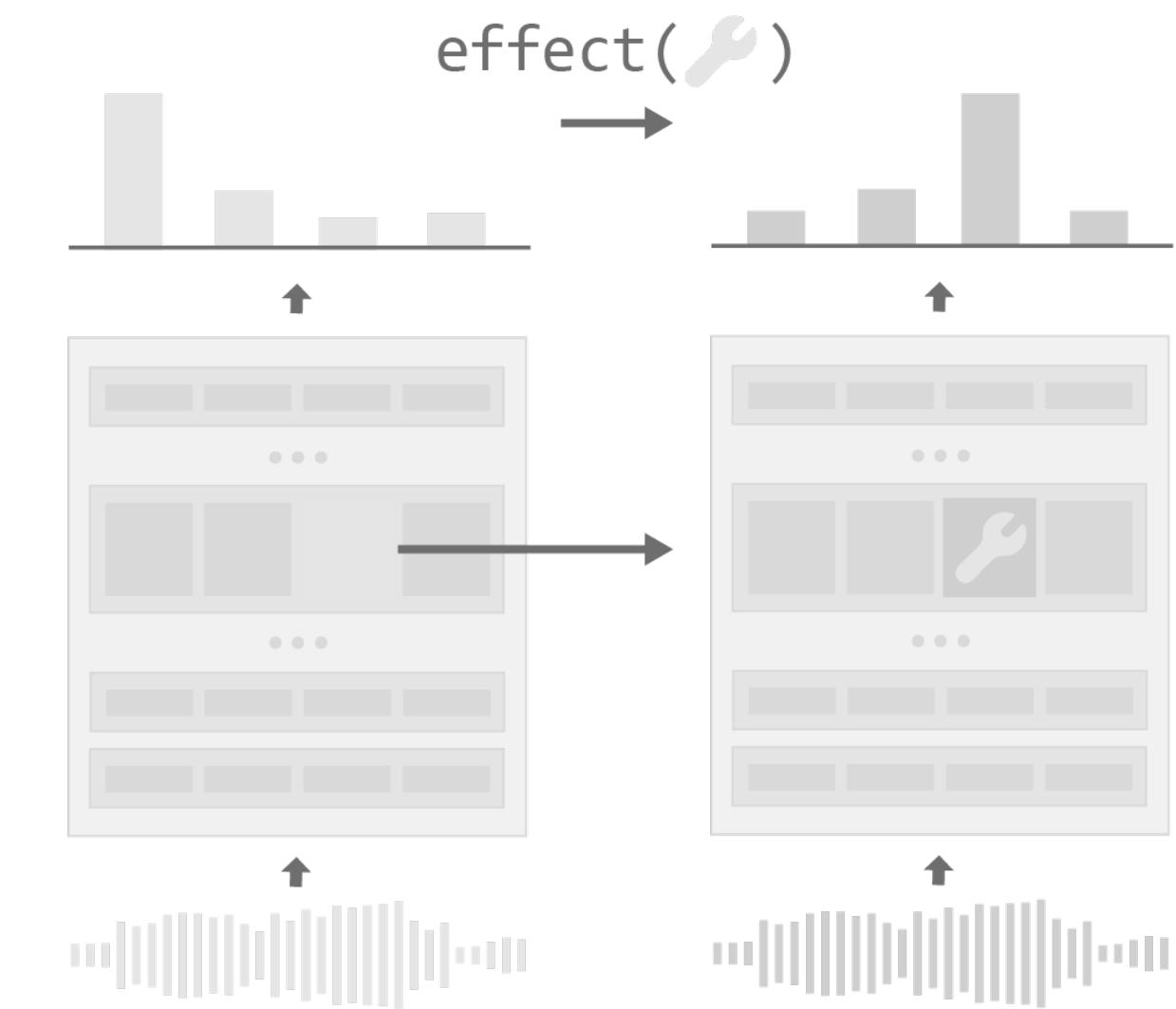
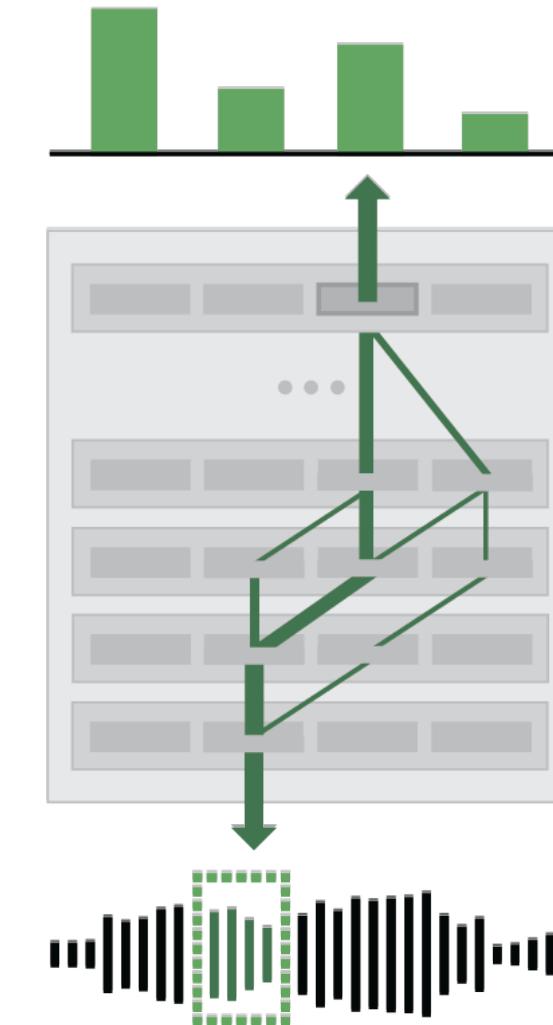
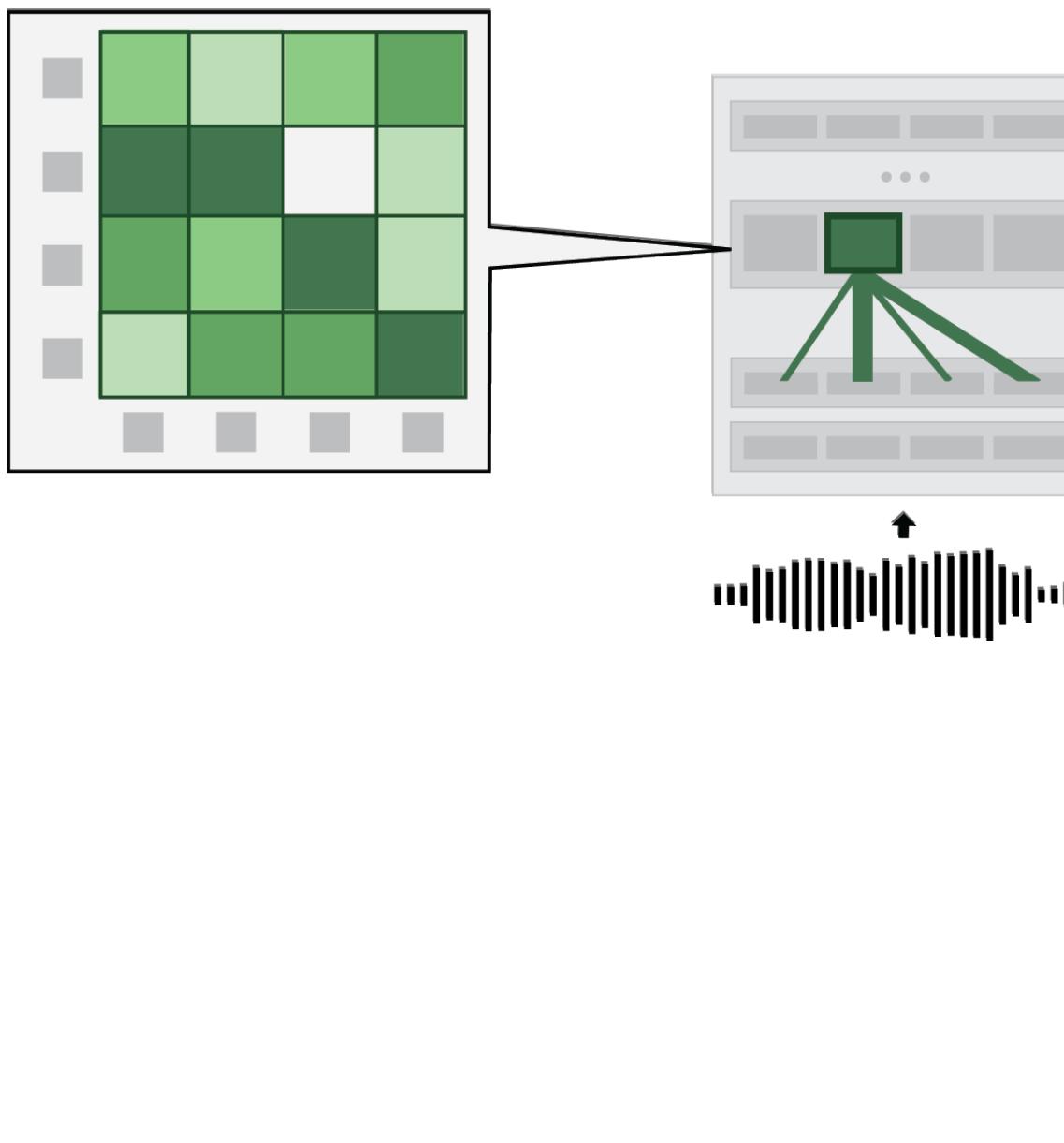
- Baselines and controls
- Zero-shot metrics vs. optimized diagnostic tools

Representational analyses

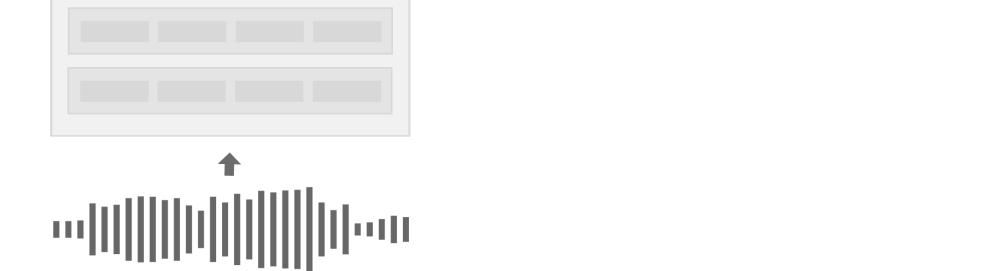
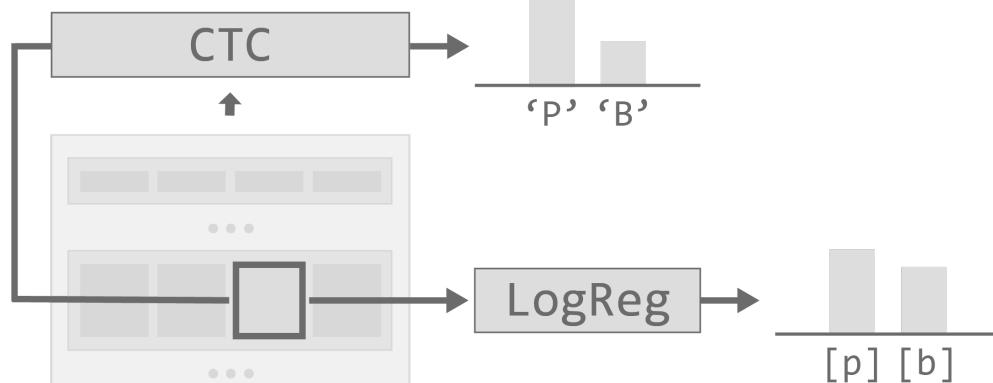
- Dimensionality reduction
- Probing classifiers
- Representation space comparisons

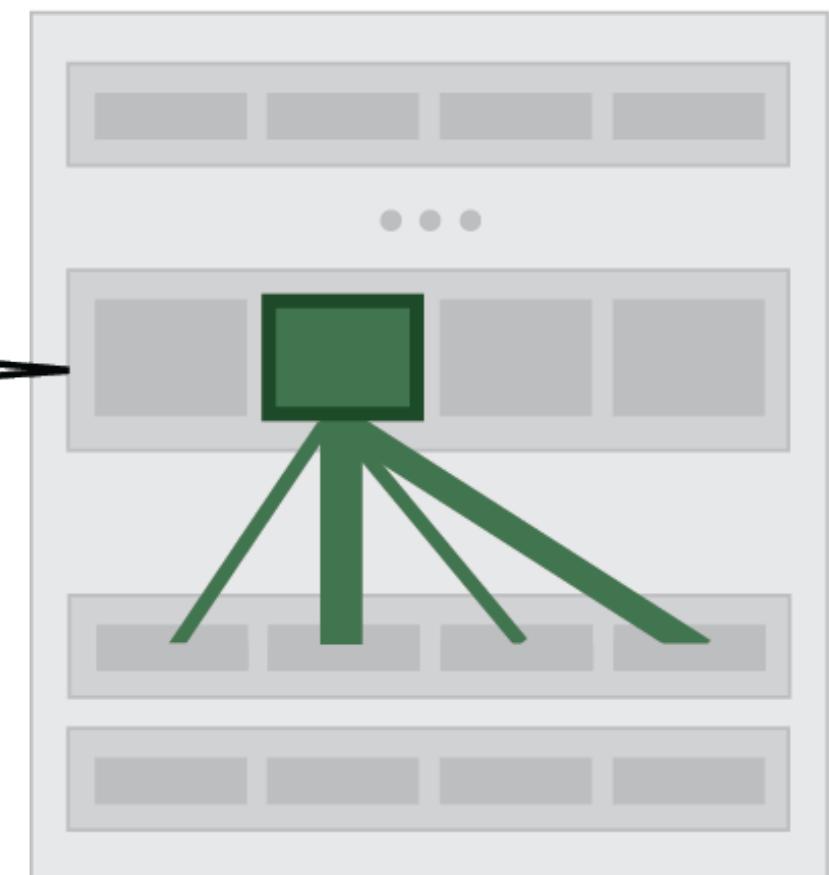
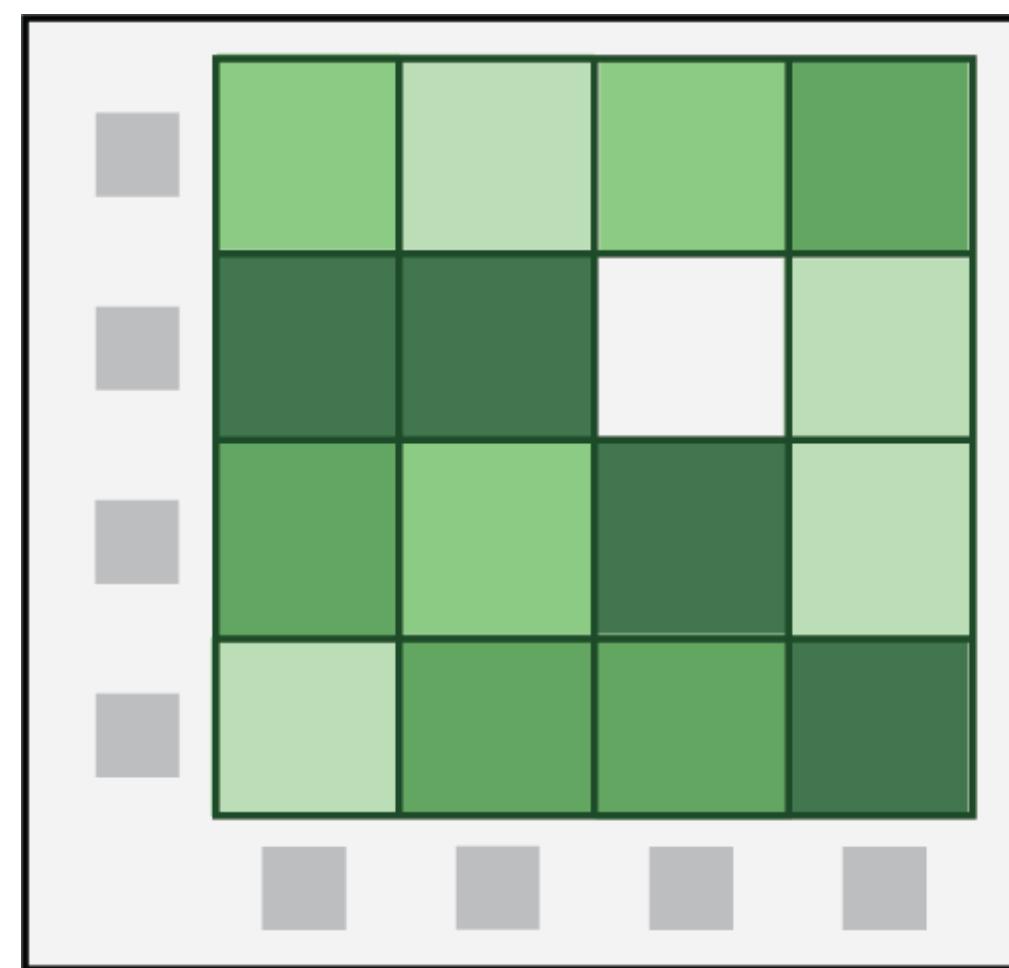


Interpretability Techniques for Speech Models



Representational
analyses

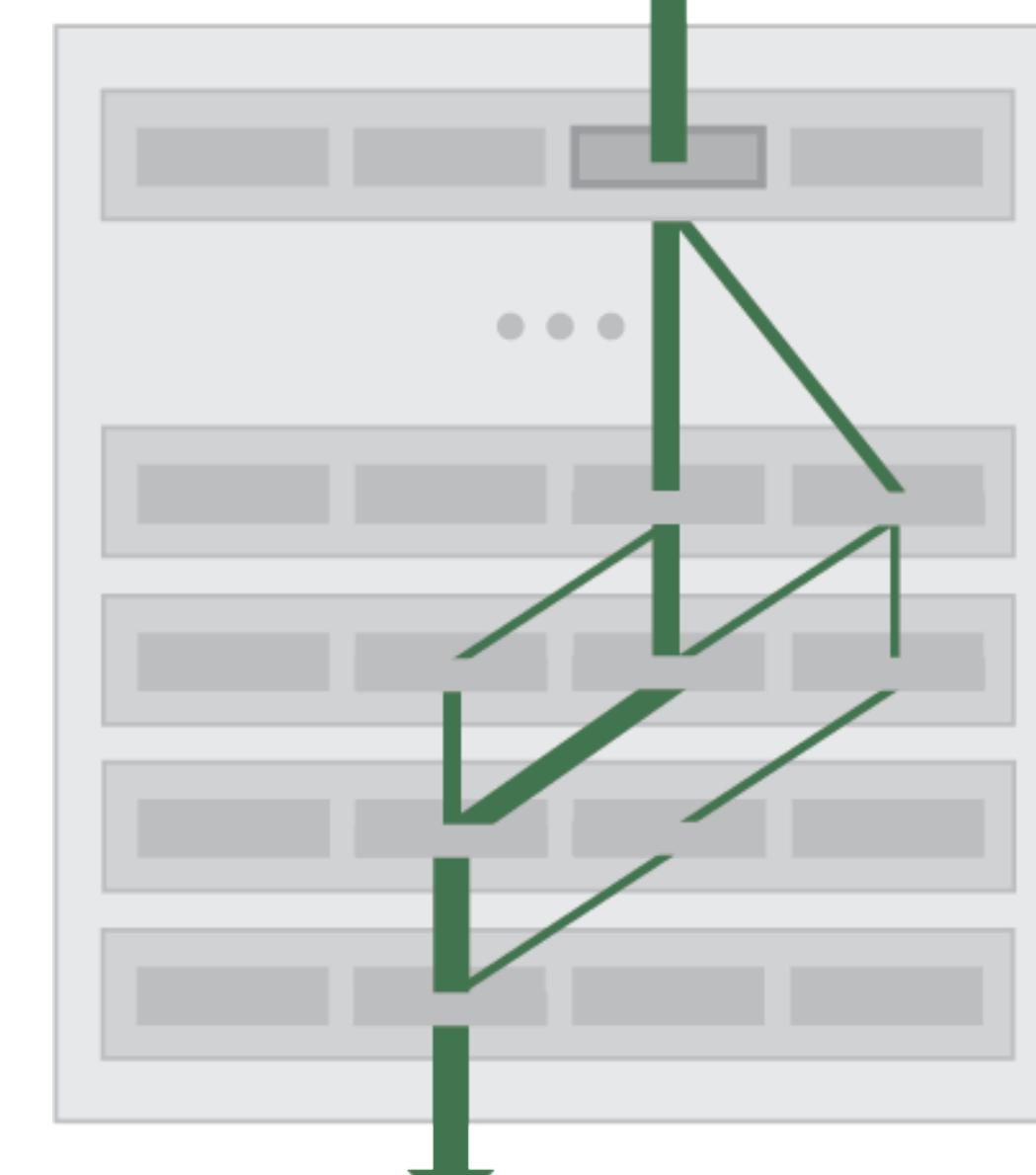
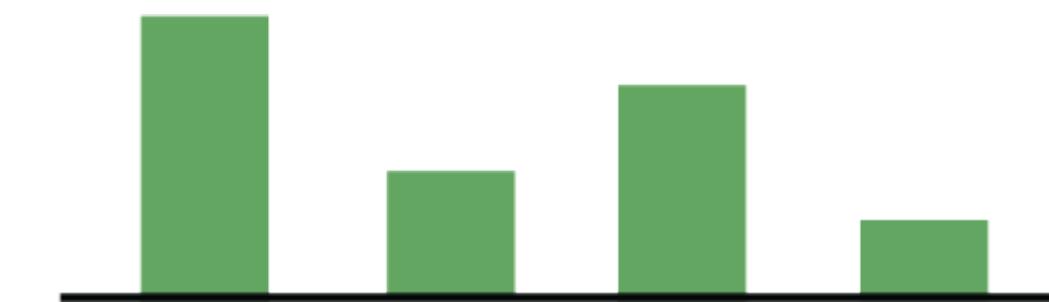




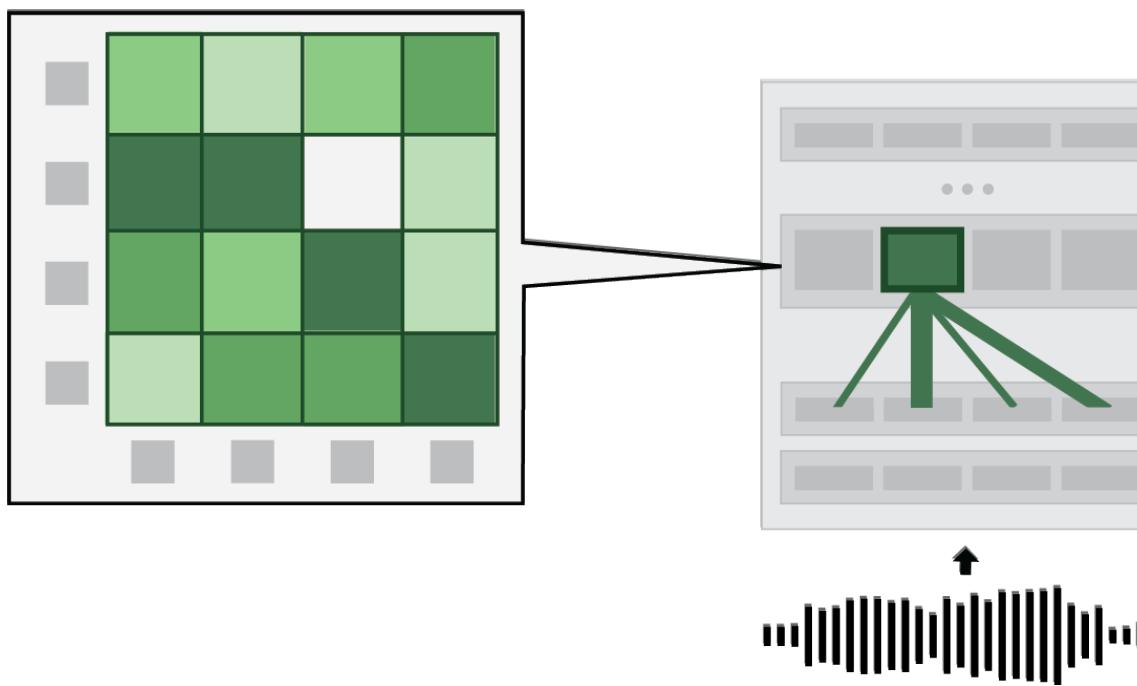
Feature Importance Scoring

- Context Mixing
- Feature Attribution

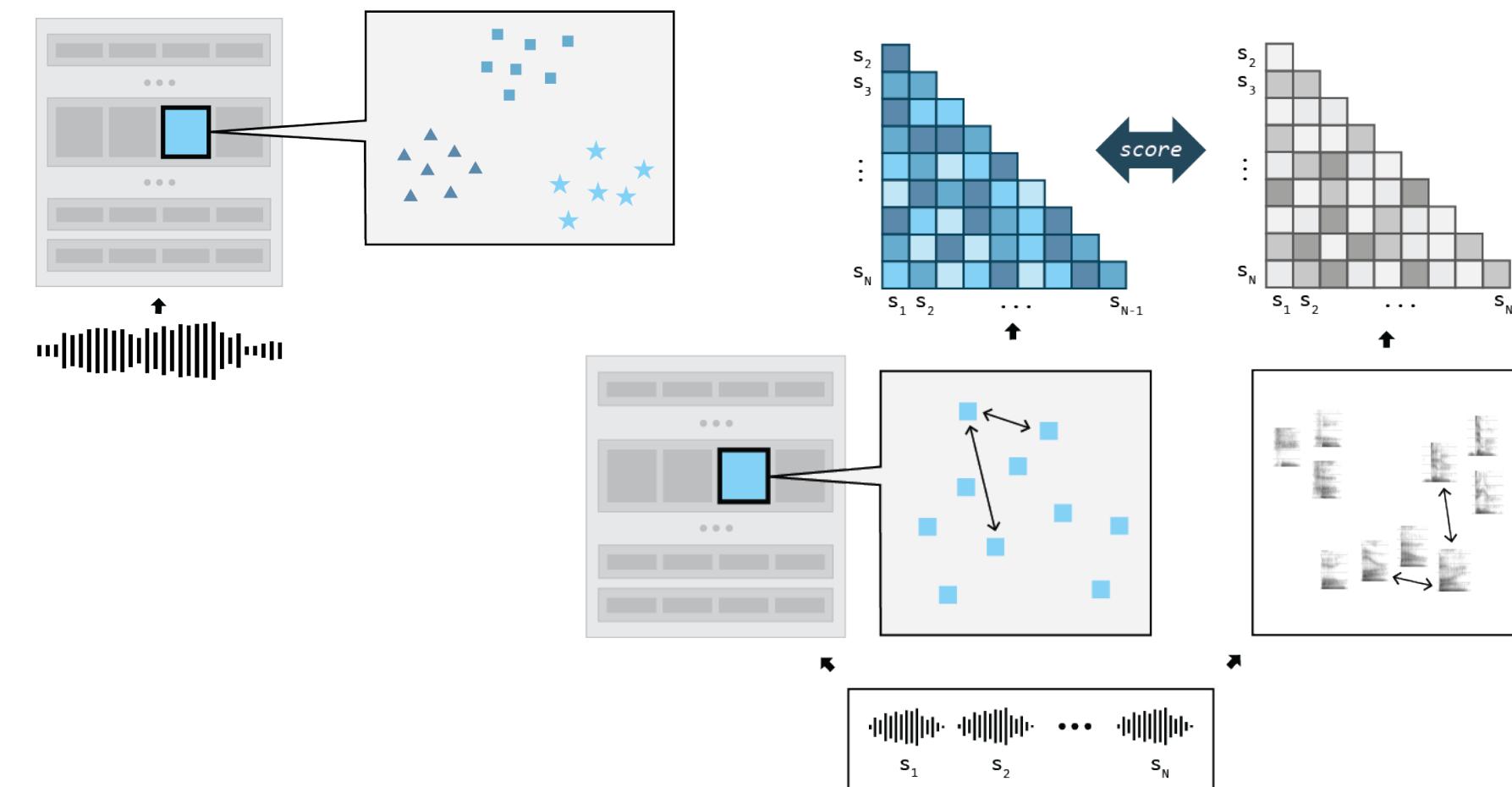
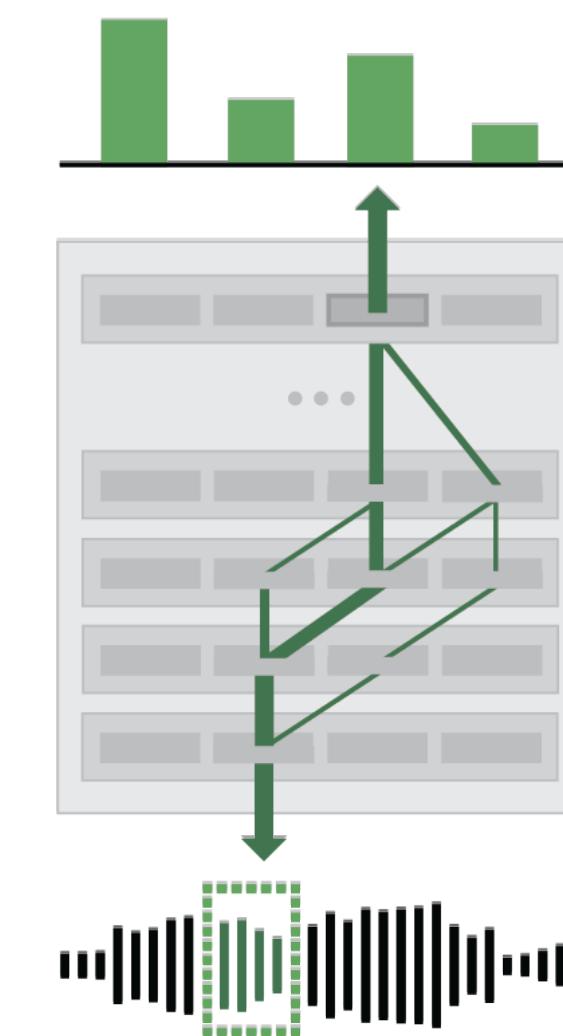
- Beyond attention weights
- Attribution reliability



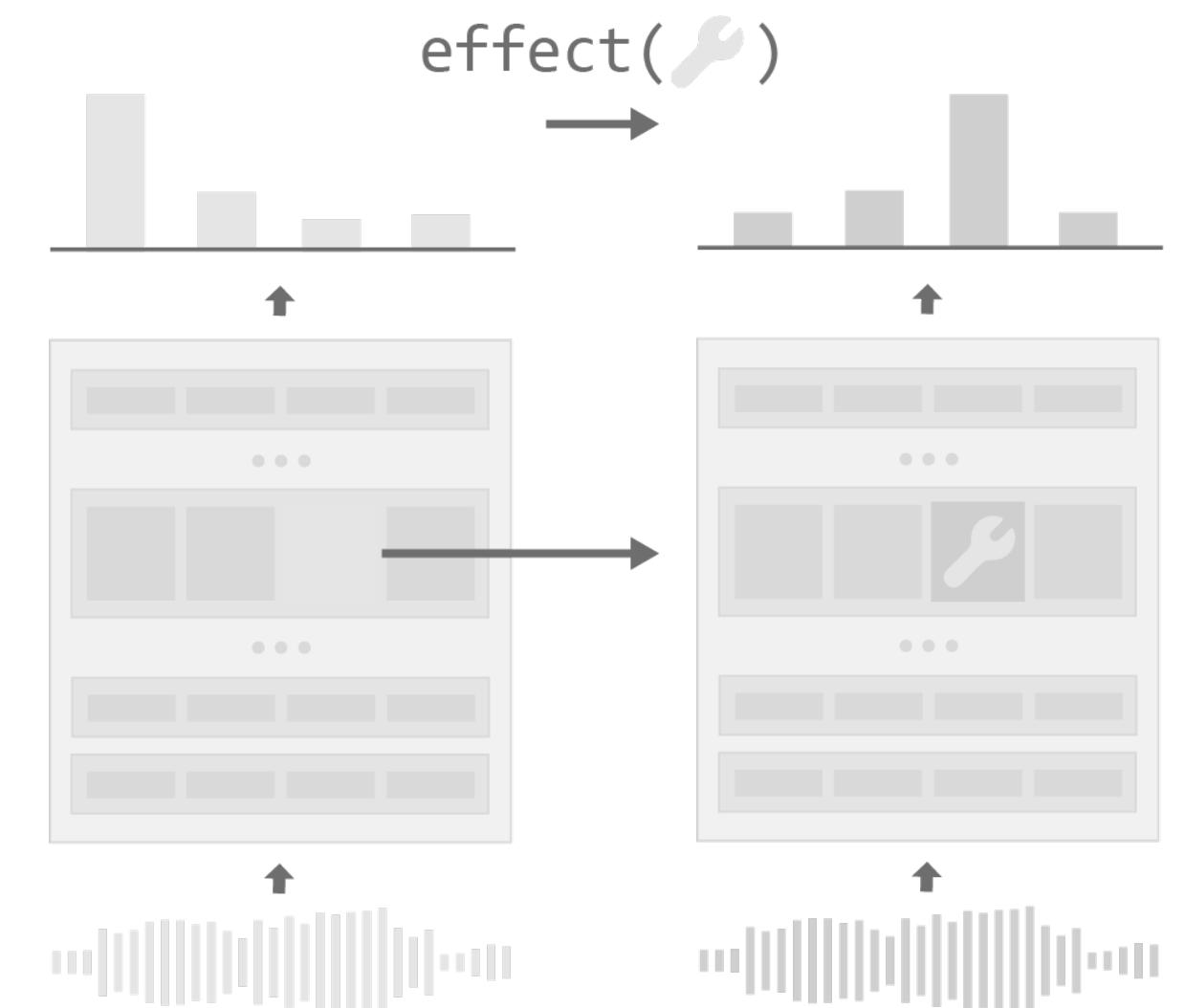
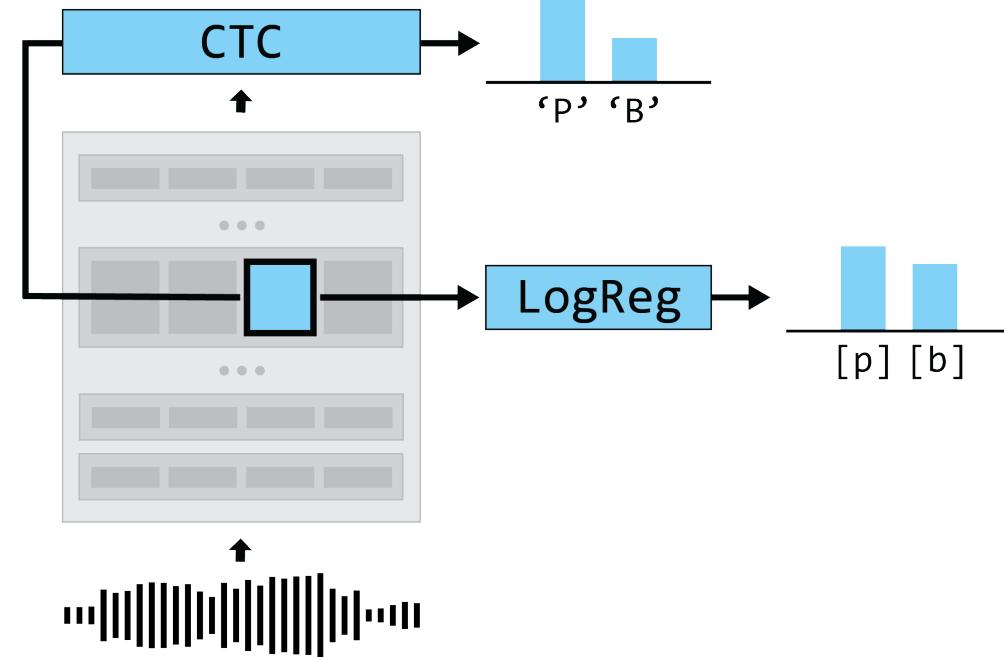
Interpretability Techniques for Speech Models



**Feature Importance
Scoring**



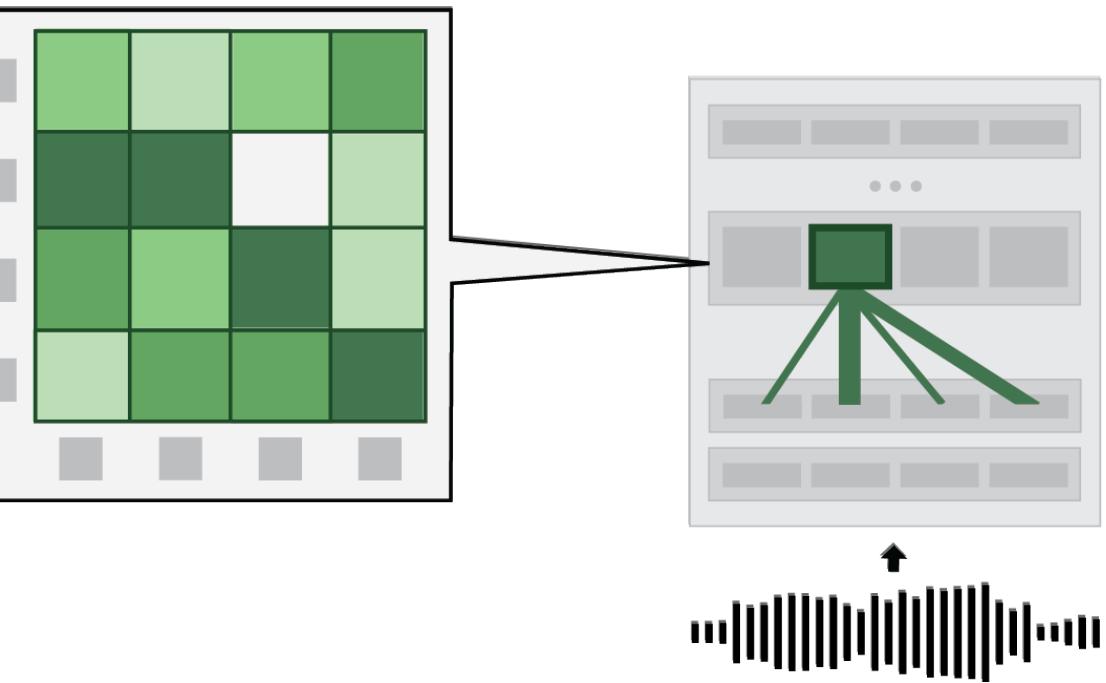
**Representational
analyses**



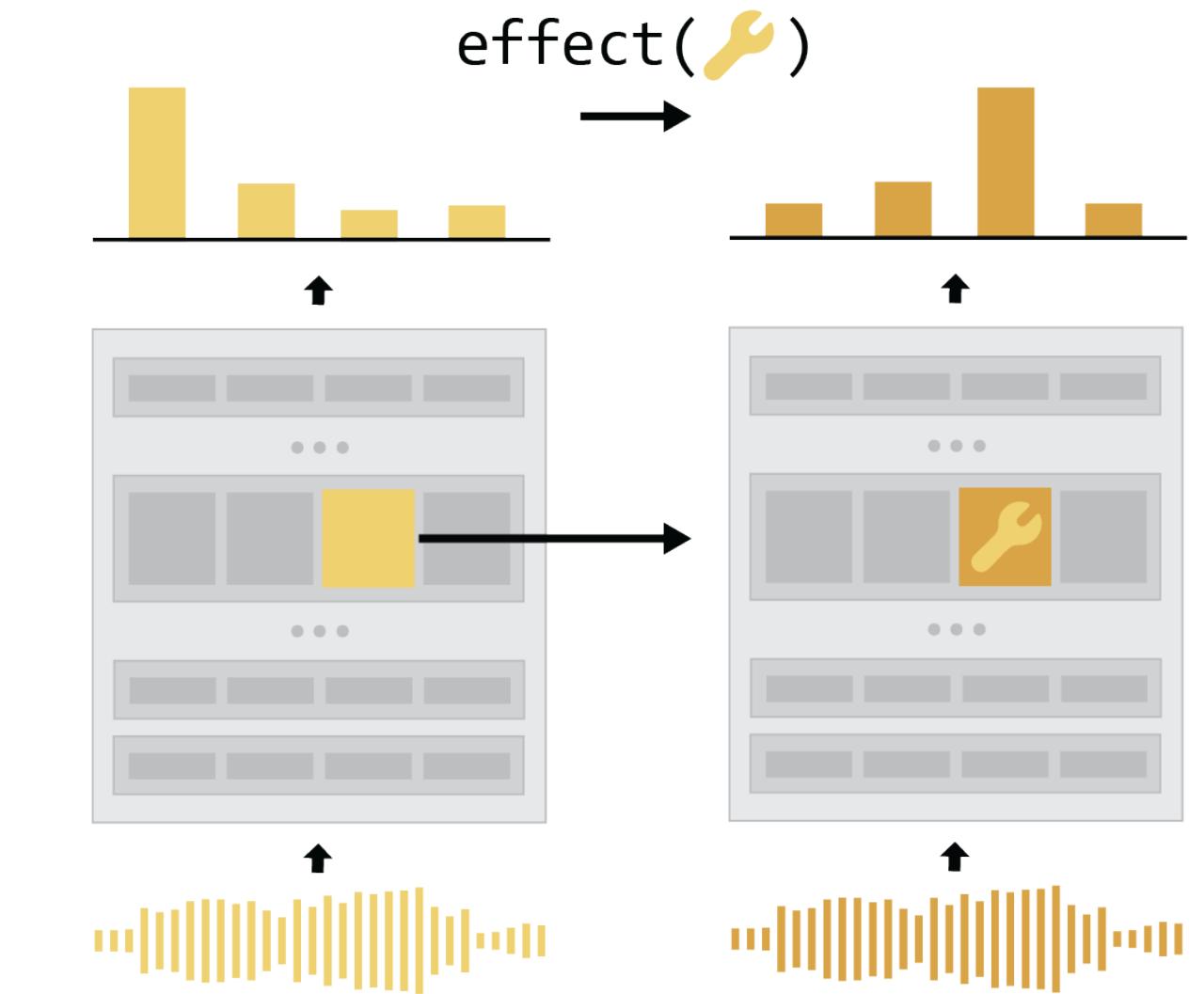
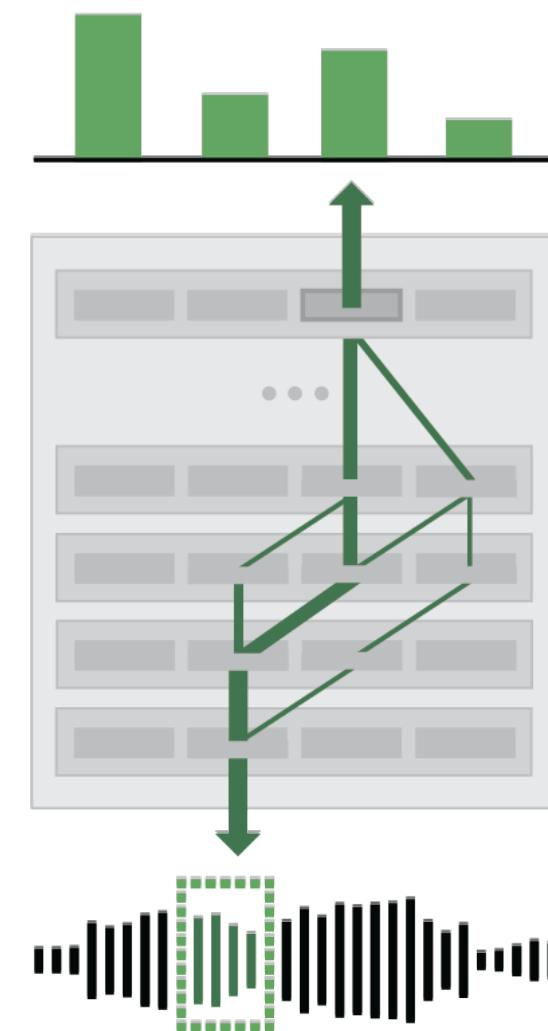
Behavioural analyses



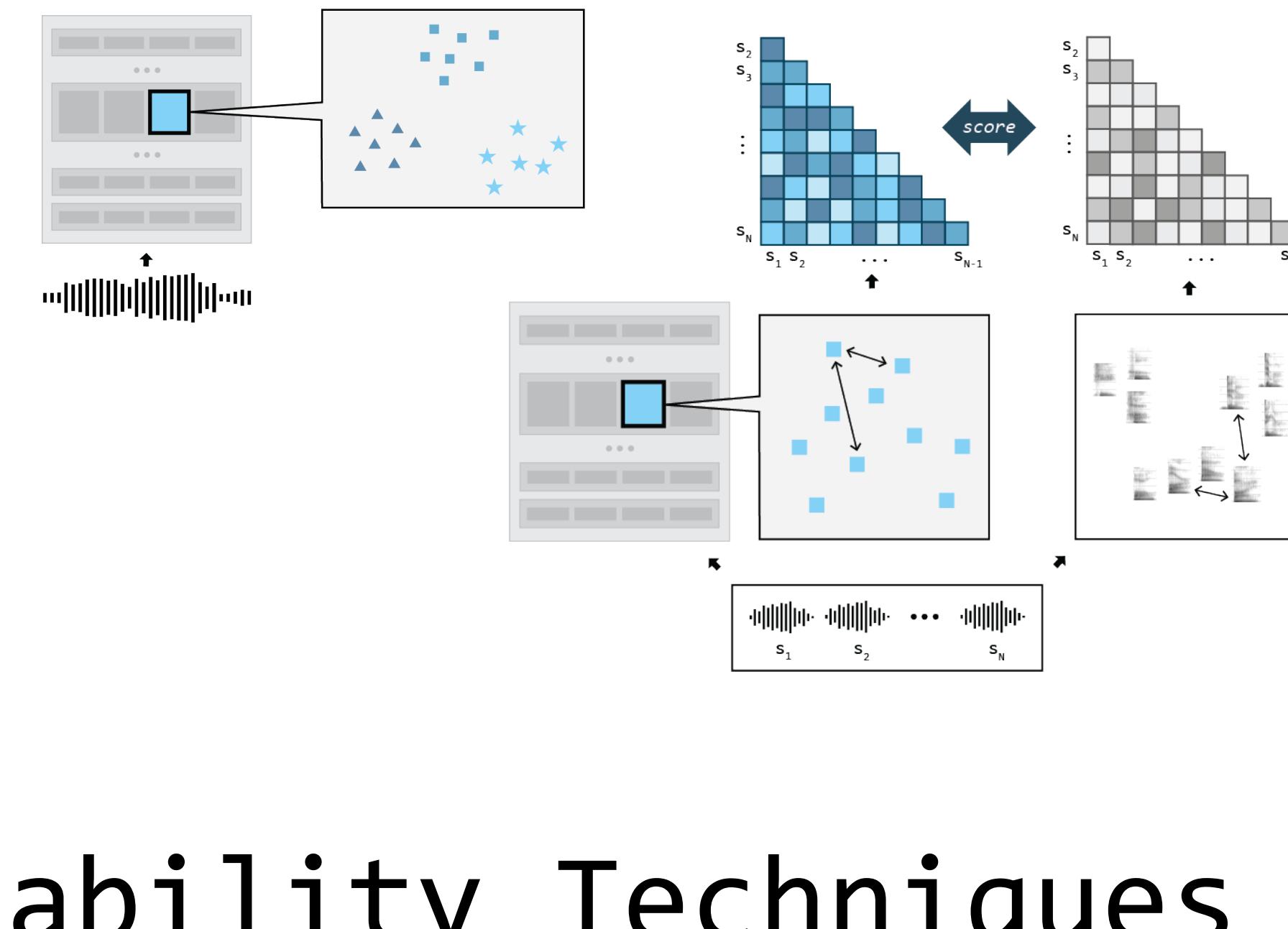
Interpretability Techniques for Speech Models



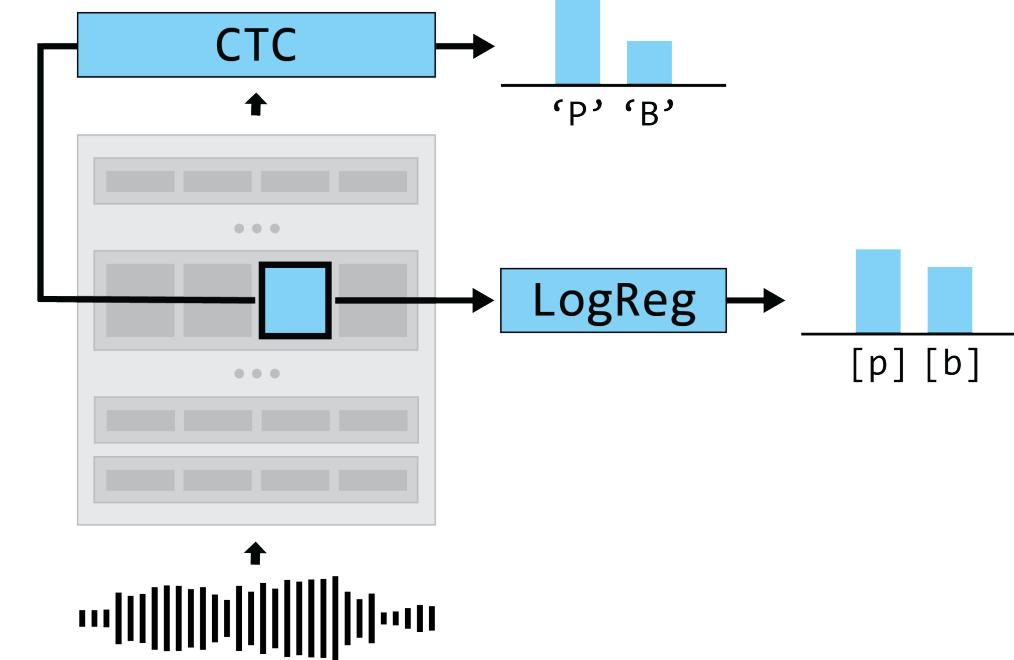
Feature Importance Scoring



Causal analyses



Representational analyses



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)

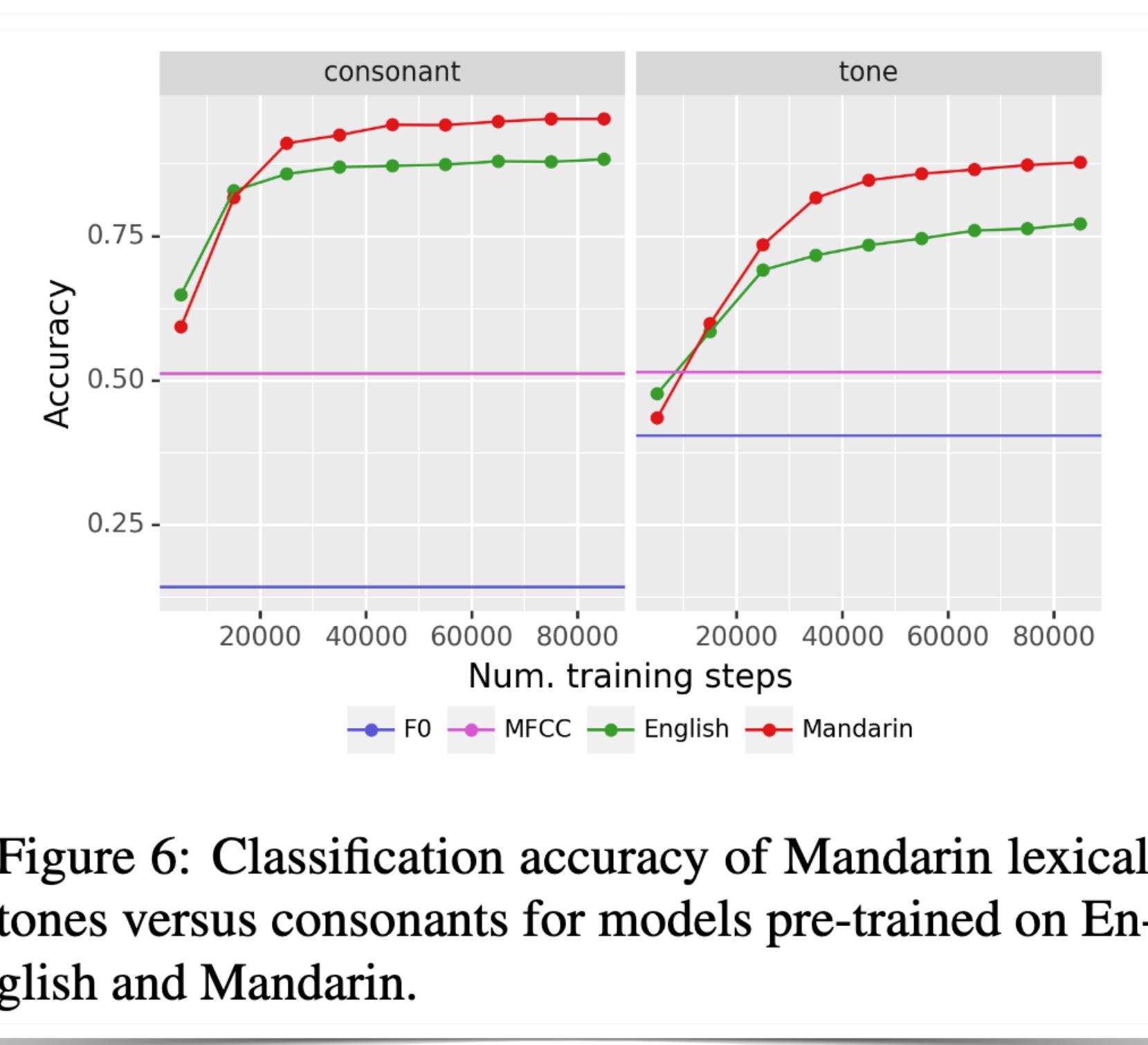


Figure 6: Classification accuracy of Mandarin lexical tones versus consonants for models pre-trained on English and Mandarin.

de Heer Kloots et al. (in prep)

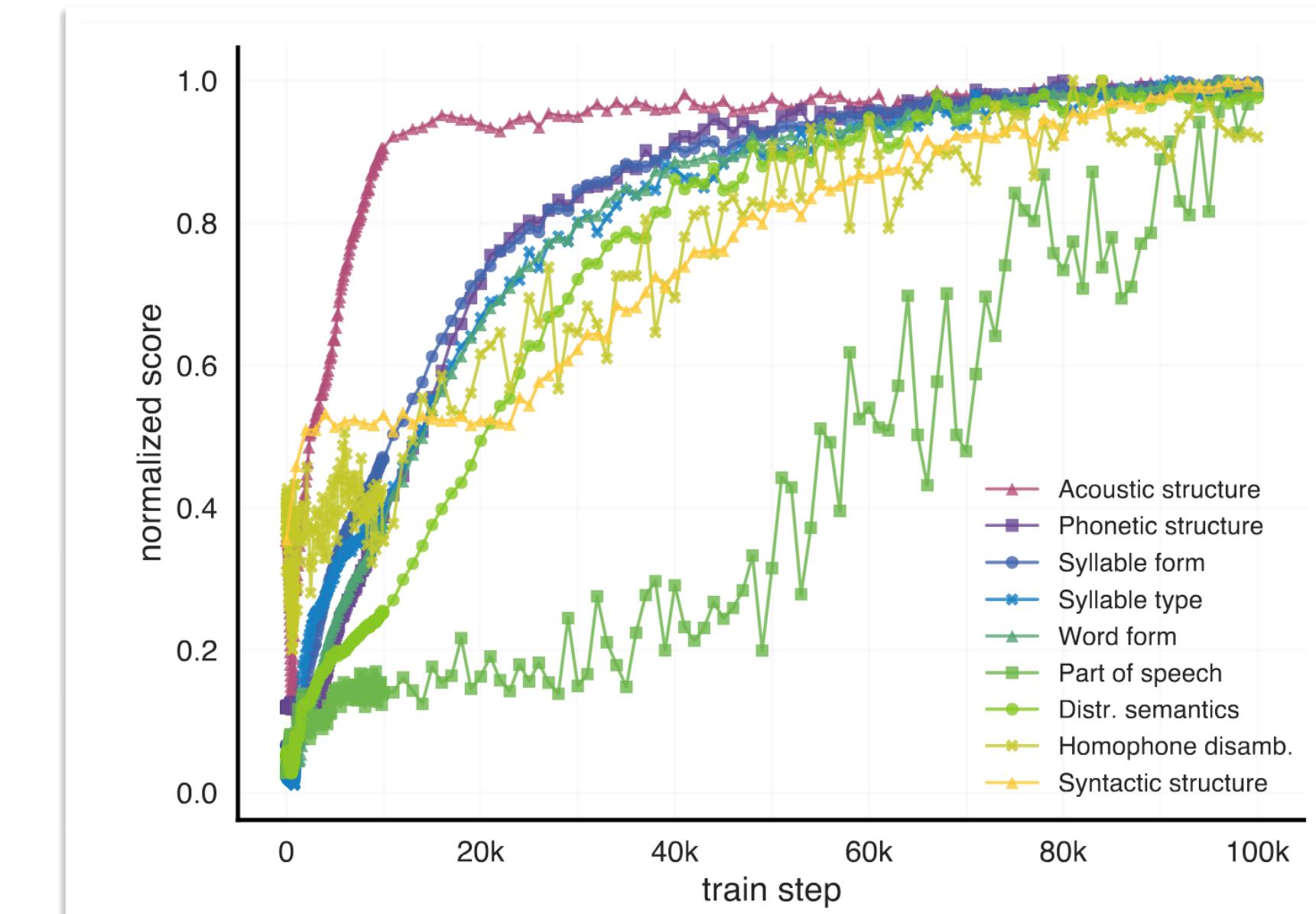


Figure 3: Development of linguistic encoding capacities across checkpoints in model training.

Understanding model training

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

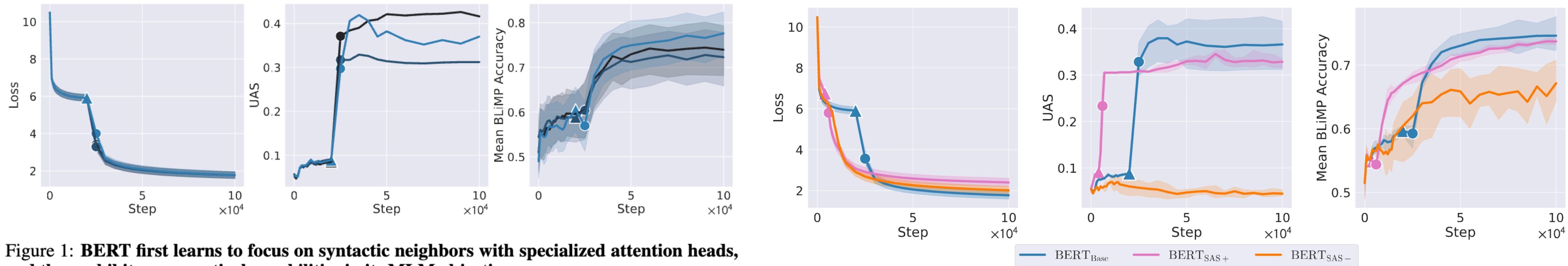


Figure 1: **BERT first learns to focus on syntactic neighbors with specialized attention heads, and then exhibits grammatical capabilities in its MLM objective.**

Outlook:

Can interpretability be useful?

Understanding model
training dynamics

Disentangling
representations &
attributions

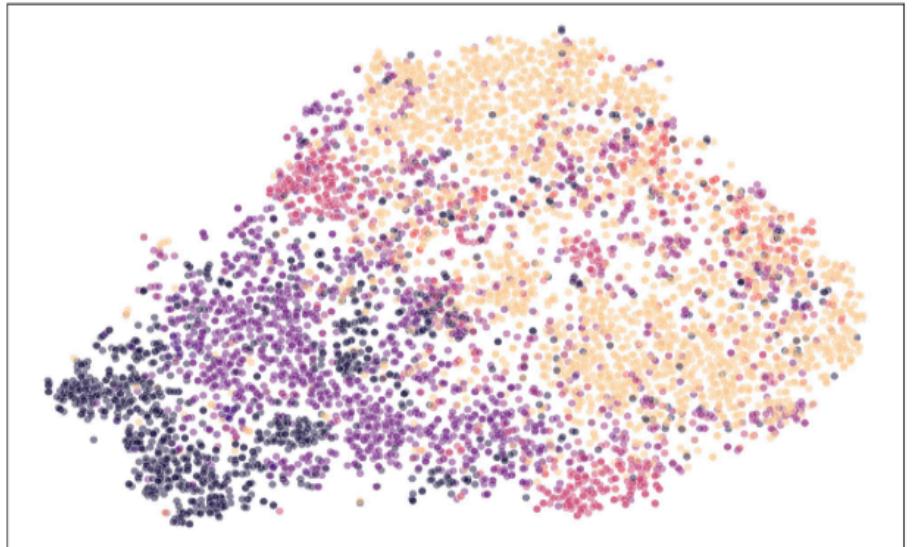
Understanding links
between model
internals & output

Disentangling representations & attribution

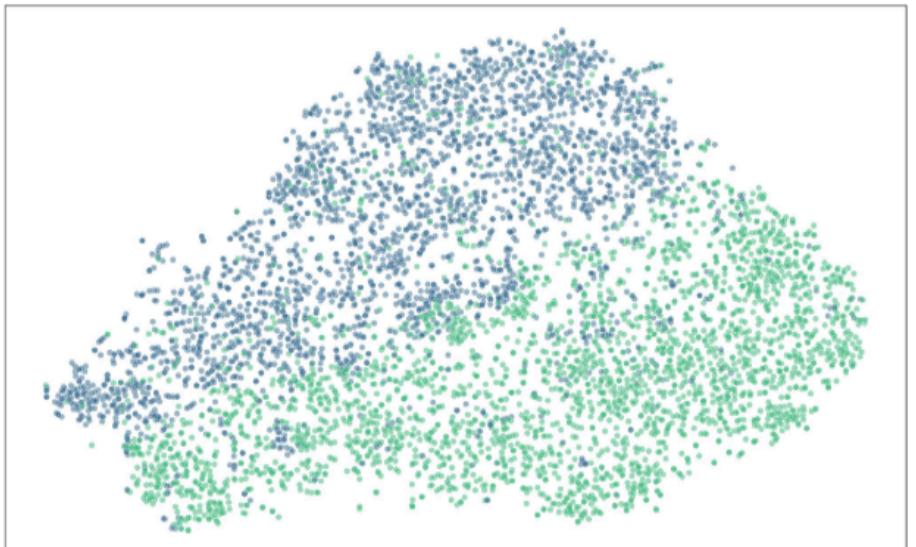
Speech model representations encode different kinds of interpretable information in an *entangled* way

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

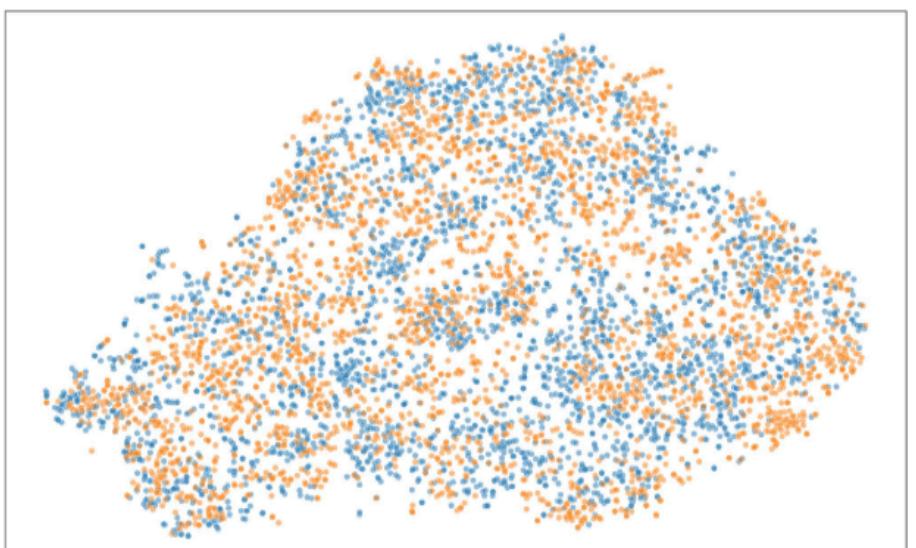
phone
class



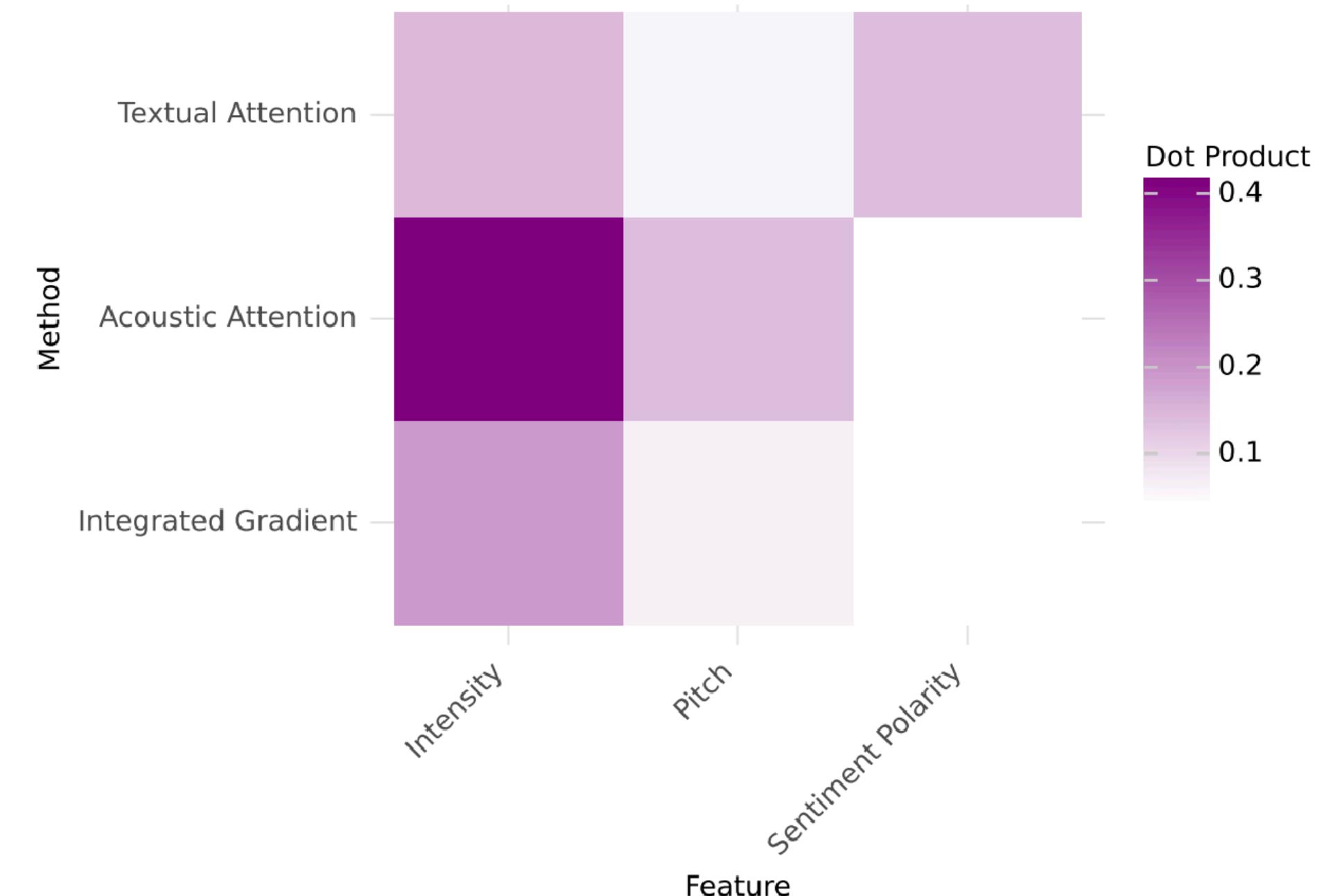
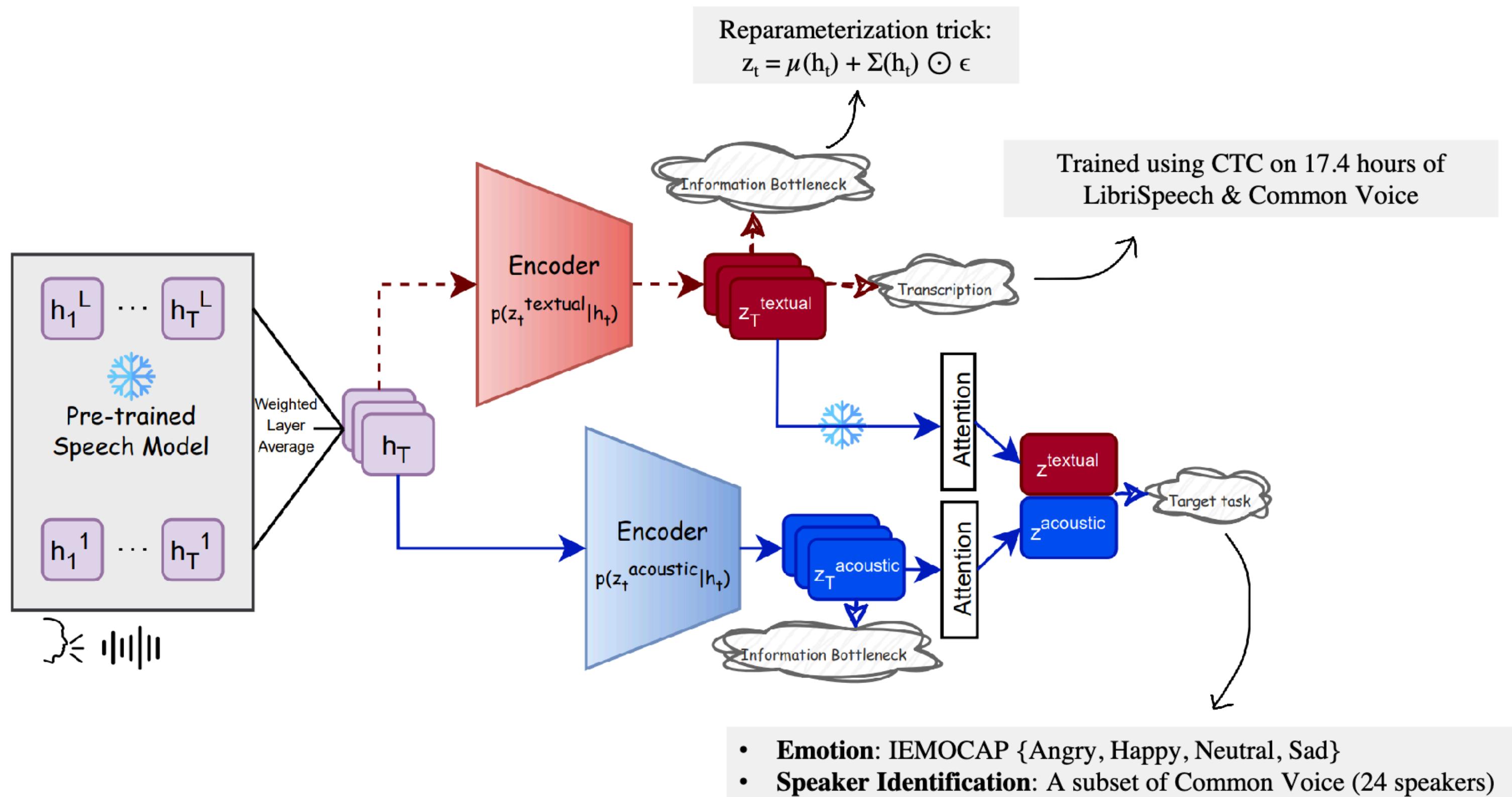
gender



language



Disentangling representations & attribution



Outlook:

Can interpretability be useful?

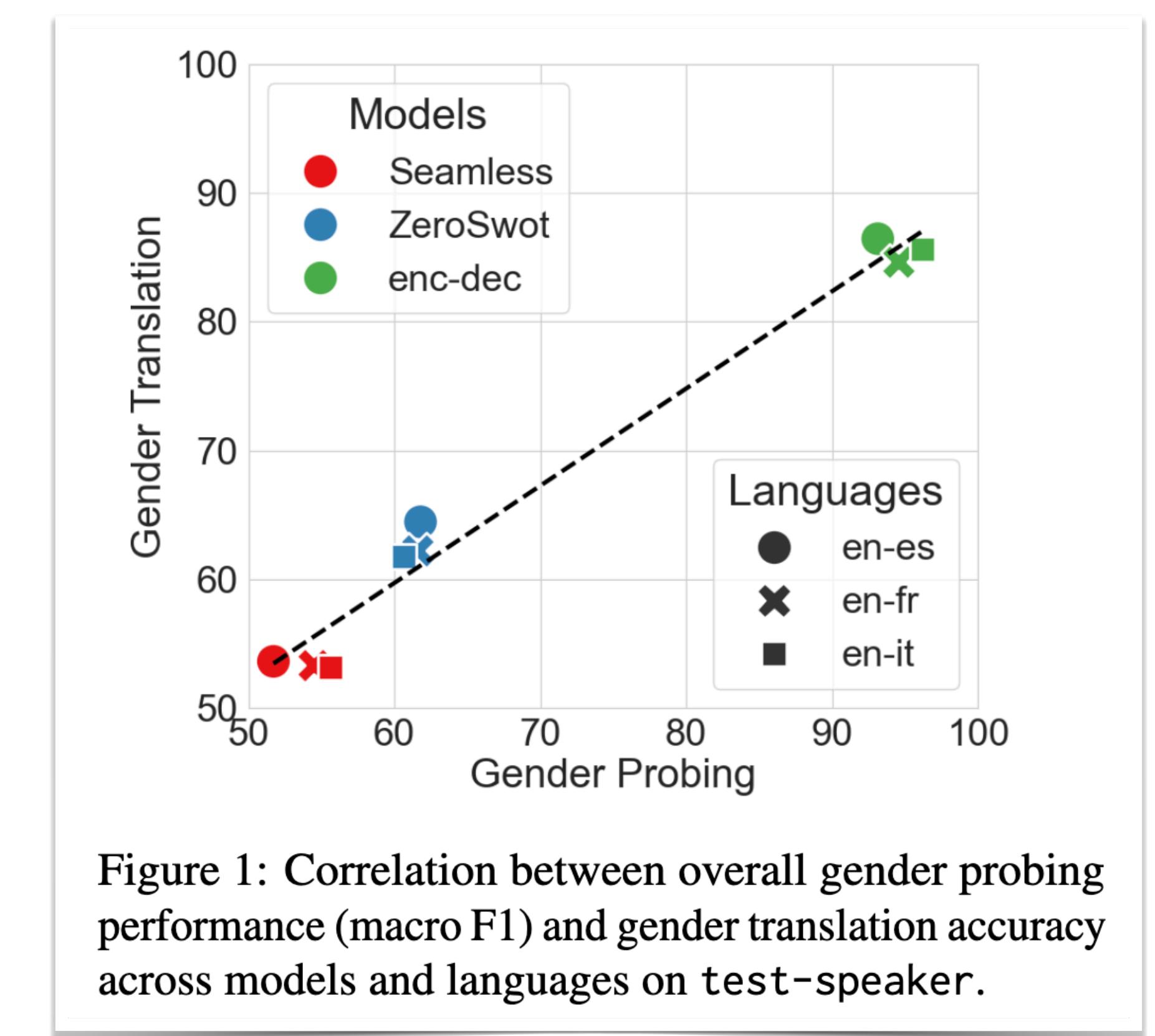
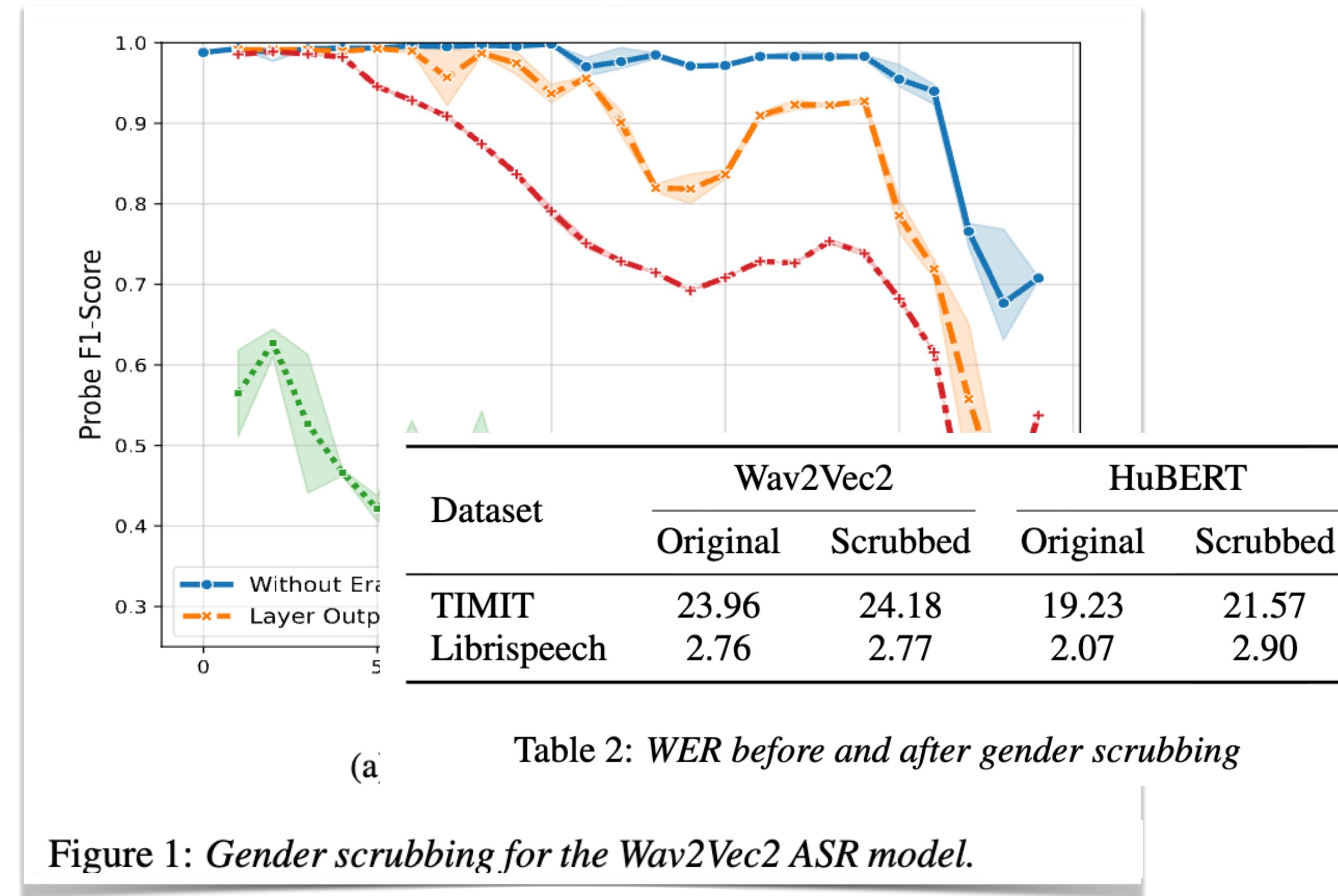
Understanding model
training dynamics

Disentangling
representations &
attributions

Understanding links
between model
internals & output

Understanding links between model internals and behaviour

Gender encoding & bias



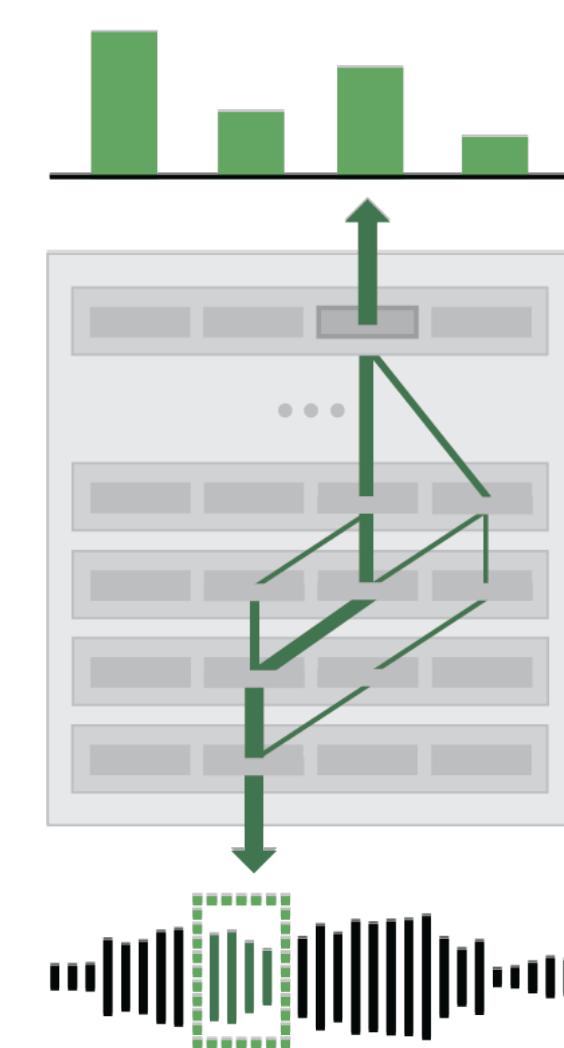
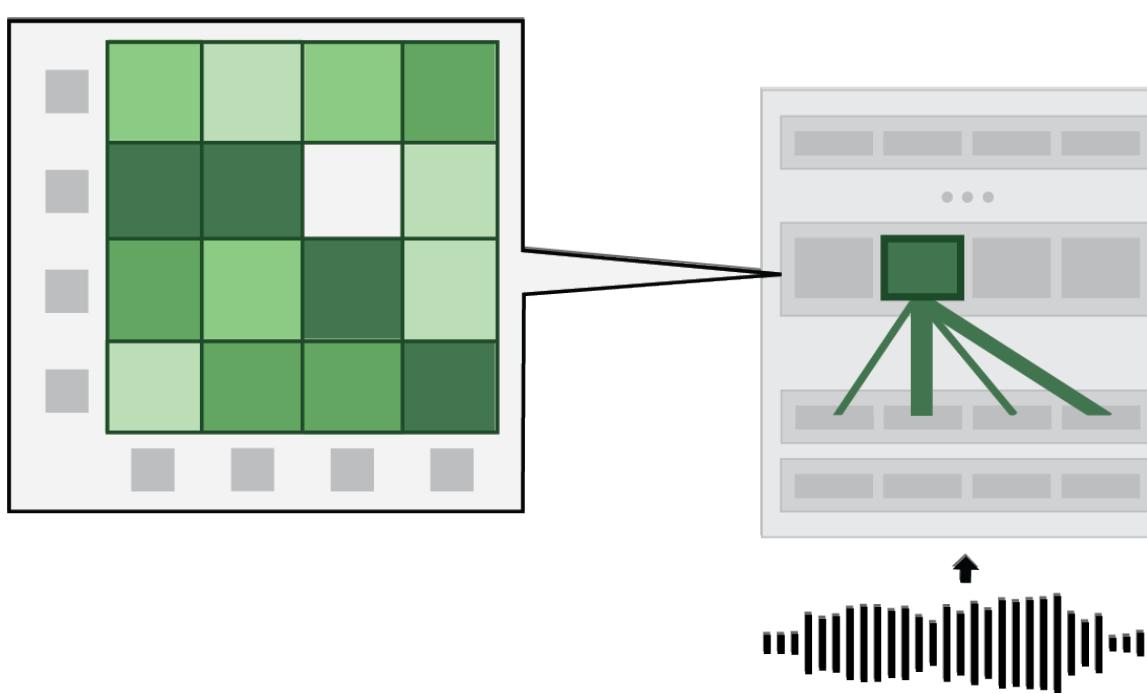
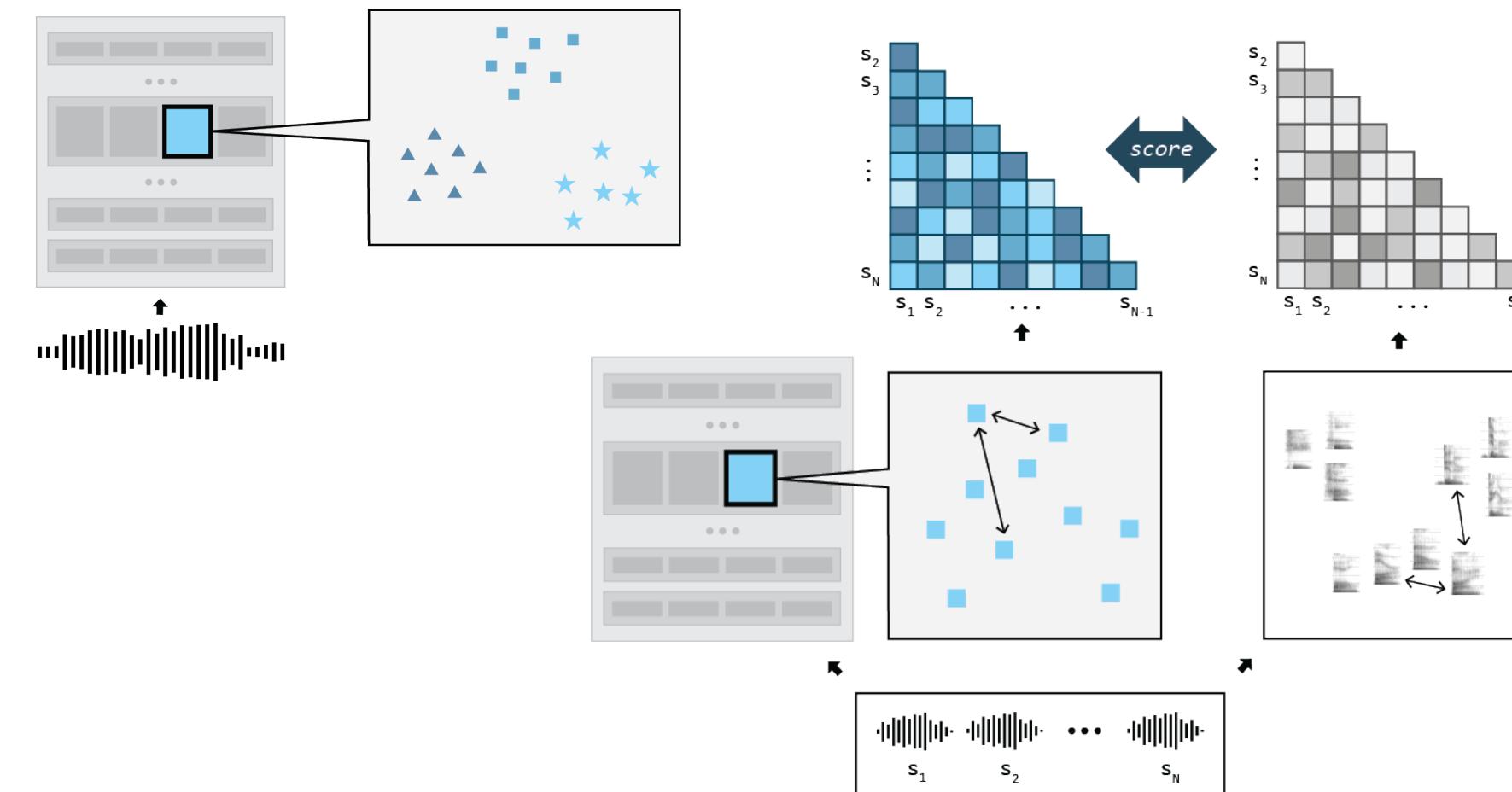
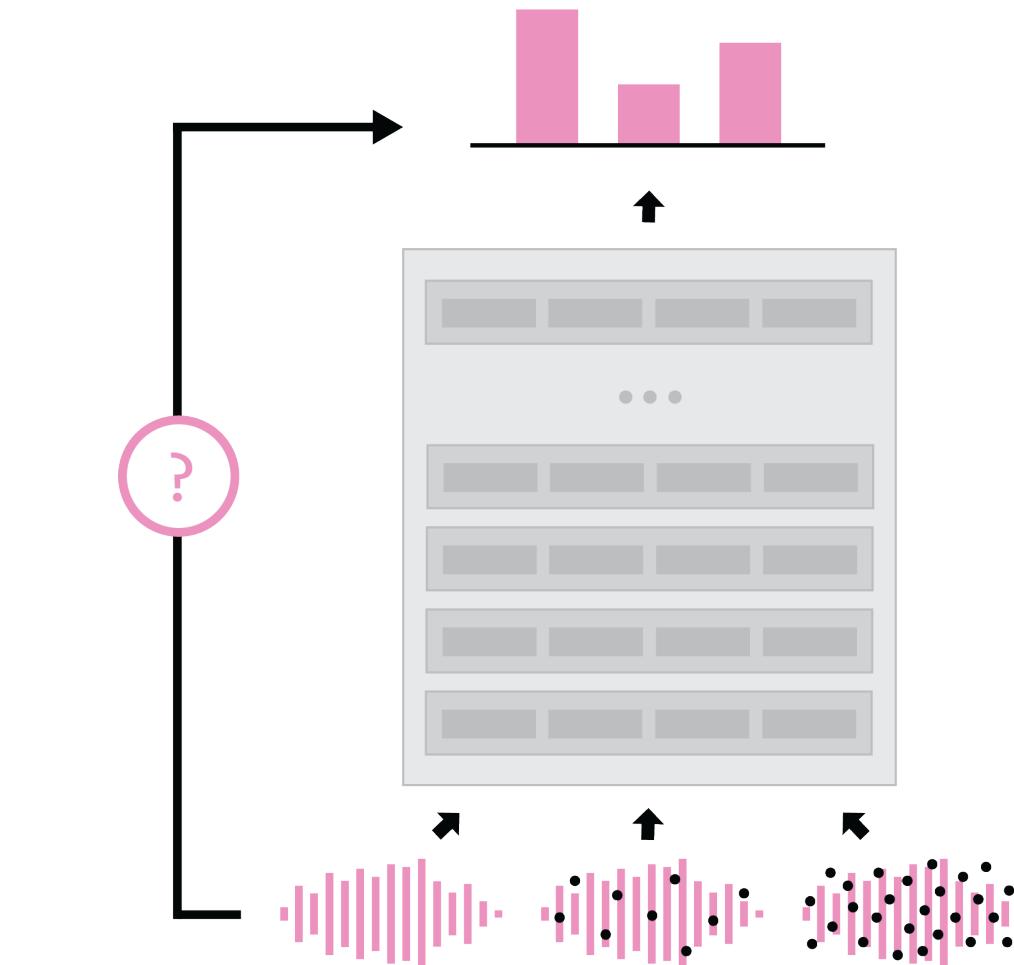
Outlook:

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

The Bitter Lesson

“The biggest lesson that can be read from 70 years of AI research is that general methods that leverage computation are ultimately the most effective, and by a large margin...”

the human-knowledge approach tends to complicate methods in ways that make them less suited to taking advantage of general methods leveraging computation.”

1.

Speech technology will (also) in the future
be dominated by machine-learned,
blackbox models. Posthoc Interpretability
will be an essential aspect of its
responsible deployment.

Interpretability is challenging

End-to-end deep learning is so successful, because it allows models to develop ‘meso-scale’ representations, involving non-linear transforms and interactions between input features, from training on gigantic datasets.

Scale, nonlinearities, feature interactions are exactly what makes posthoc interpretability challenging

Success depends on isolating the nonlinear feature interactions, but the best way to do this will not trivially become apparent by bottom-up techniques

2.

No Success Without Engaging in Earnest
with Theory and without Elaborate Reasoning

2.

No Success Without Engaging in Earnest
with Theory and without Elaborate Reasoning

The No-SWEETER lesson!

3.

The rich tradition of theories and descriptive concepts in perception research, phonetics and phonology offer great opportunities for posthoc interpretability of speech technology

1. Speech technology will (also) in the future be dominated by machine-learned, blackbox models. Posthoc Interpretability will be an essential aspect of its responsible deployment.
2. No Success Without Engaging in Earnest with Theory and without Elaborate Reasoning
3. The rich tradition of theories and descriptive concepts in perception research, phonetics and phonology offer great opportunities for posthoc interpretability of speech technology