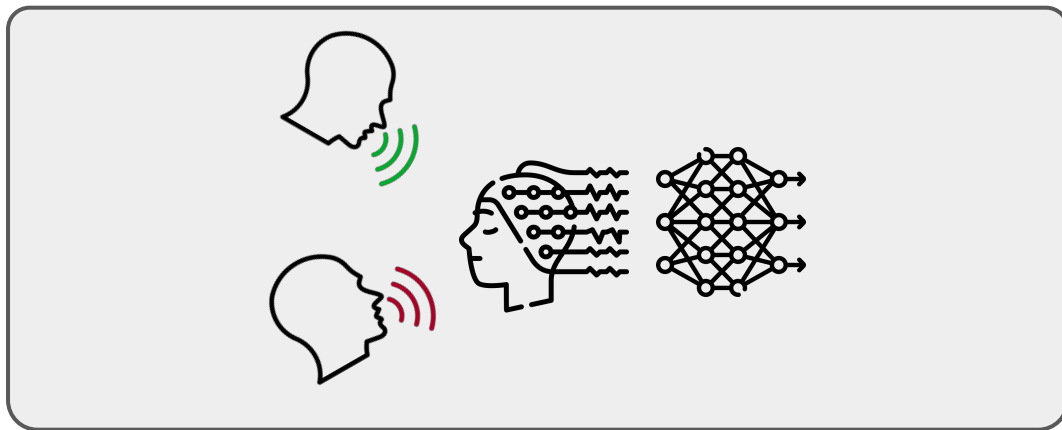# Exploring Foundation Models for Auditory Attention Decoding

Rasmus Steen Mikkelsen (s204135)

Victor Tolsager Olesen (s204141)
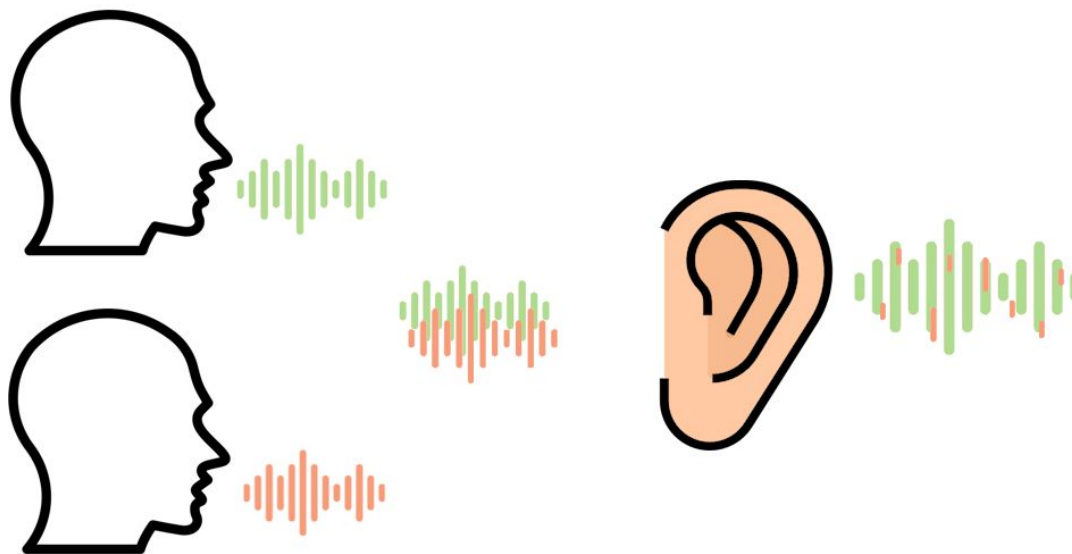
# Introduction

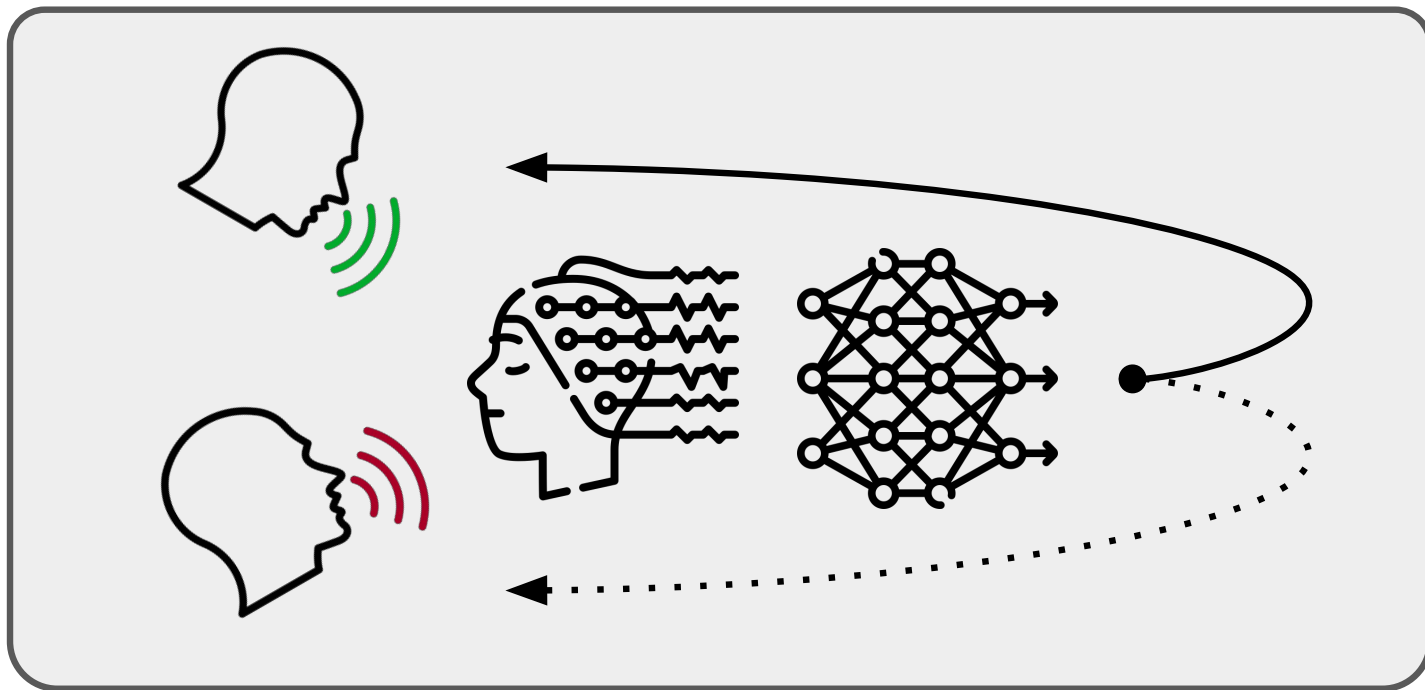# Introduction

- Cocktail party effect
- Hearing aid users

- AAD: Audio+EEG → Attention
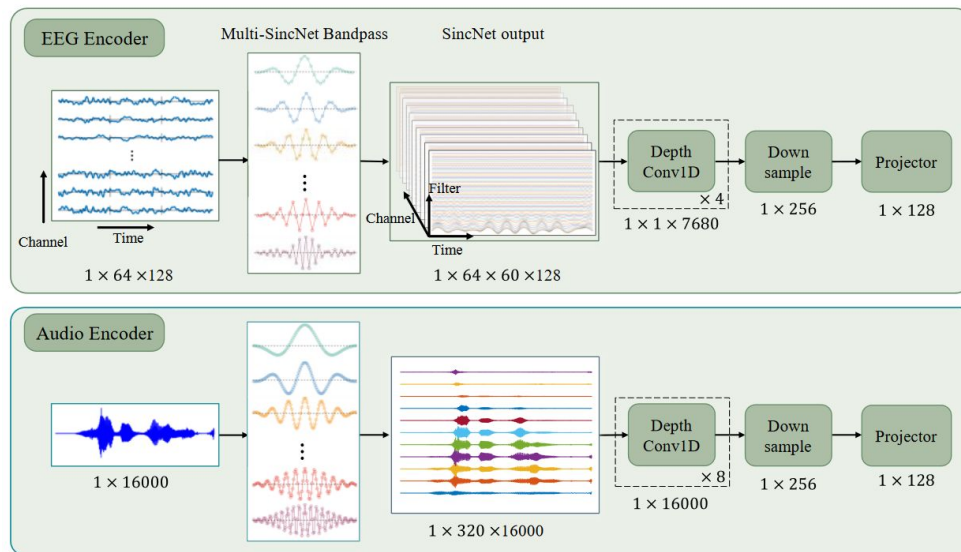- Decision window: Time segment used to predict

# Introduction

## Foundation models

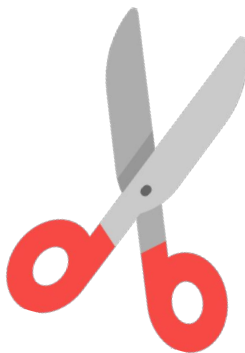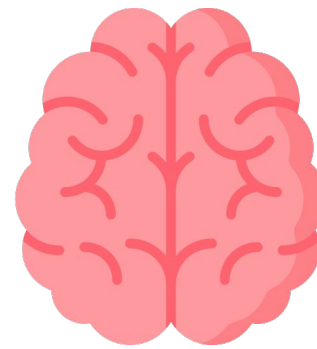- Foundation Models
- SOTA AAD Models



**NLP**
**BERT**

**Vision+Text**
**CLIP**

**Audio+Text**
**CLAP**

**EEG**
**LaBraM**

**RQ1:** How do CLAP and LaBraM perform as pretrained feature extractors for auditory attention decoding?

**RQ2:** How does contrastive learning compare to supervised classification for training robust AAD models using CLAP and LaBraM?

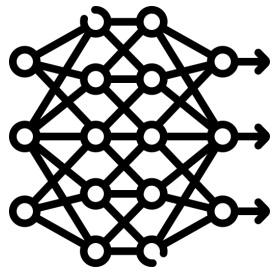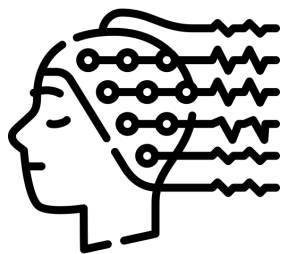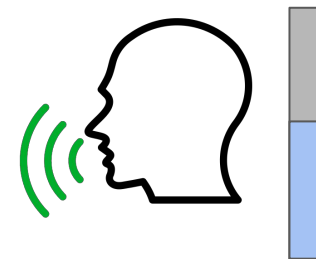**RQ3:** How does the length of decision windows affect performance?
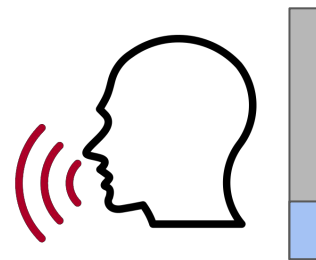
# Literature Review

**Backwards Approach**

**Correlation**

Feature Fusion

Audio Encoder

EEG Encoder

0.9

**Auditory Spatial Attention Decoding**

## Why Direct Classification?

> [..] the process of stimulus reconstruction [..] is not optimized to effectively detect attention. [...] the compression of multichannel EEG signals into a single waveform through stimulus reconstruction reduces the available information for analysis[1]

> [...] correlation between the reconstructed and the attended speech envelopes is generally weak[2]

[1]: Siqi Cai et al. "EEG-based Auditory Attention Detection in Cocktail Party Environment."
[2]: Enze Su et al. "STAnet: A Spatiotemporal Attention Network for Decoding Auditory Spatial Attention From EEG."

- Larger models
- Our model: LAION-CLAP

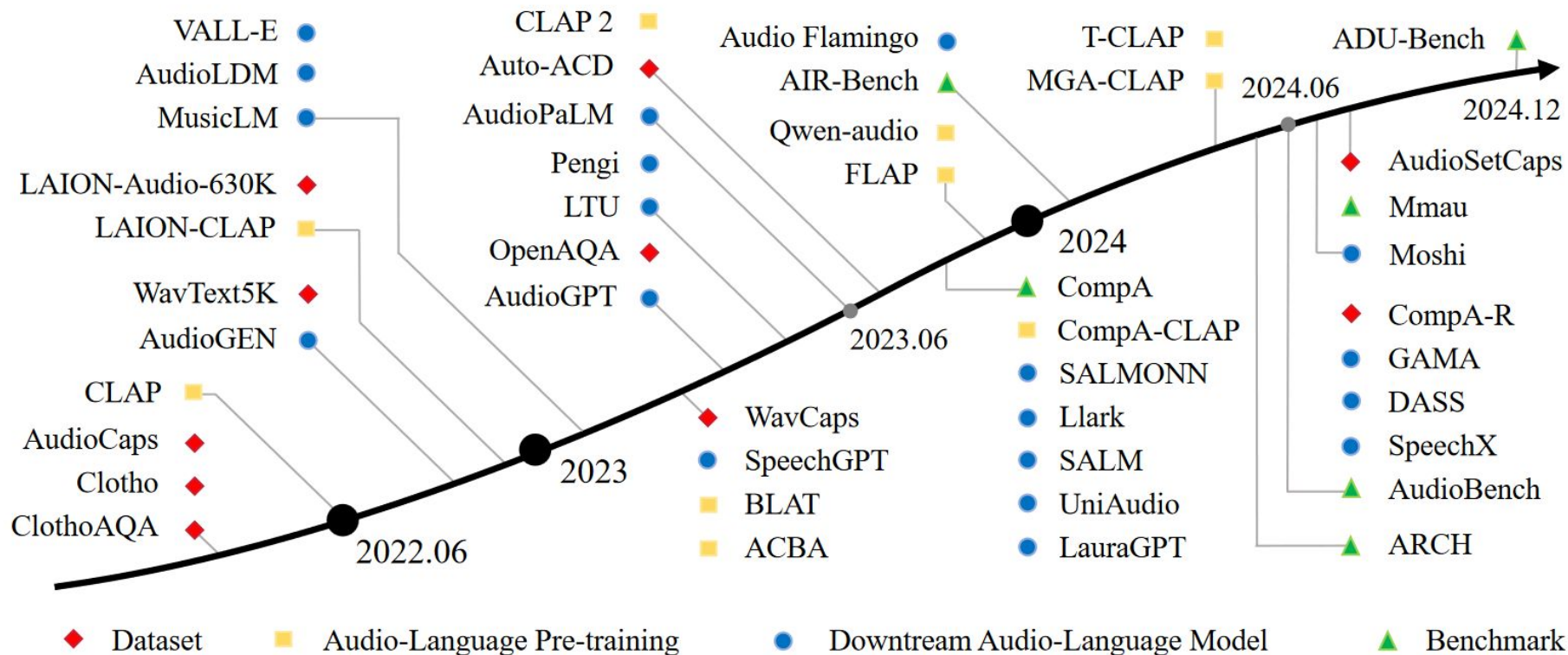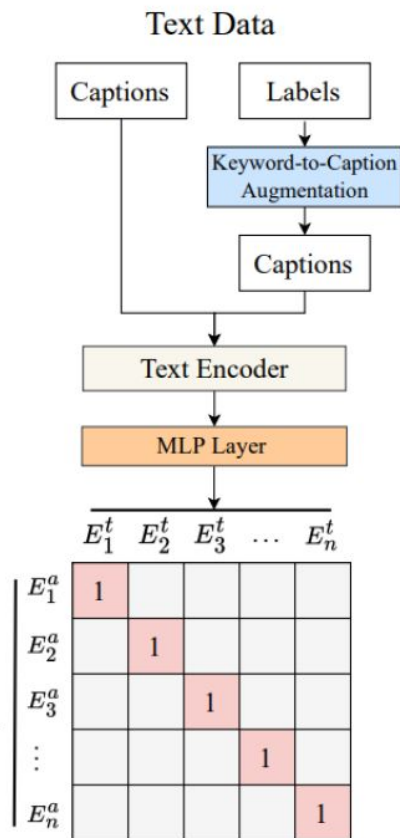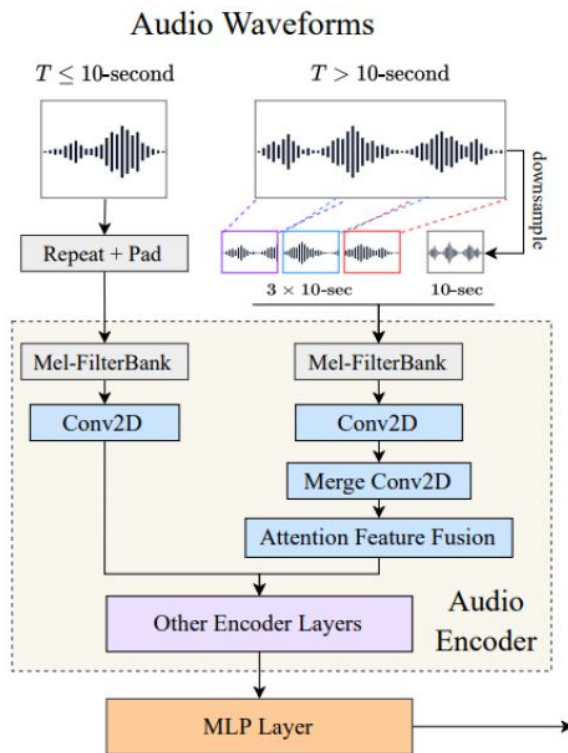## LAION-CLAP

- Contrastive Language Audio Pretraining (CLAP)
- Trained on multiple datasets

🎵 Traffic_Light.wav 🎵

*A group of people standing on the street near a busy freeway.*



Audio Waveforms

$T \leq 10\text{-second}$     $T > 10\text{-second}$

downsample

Repeat + Pad

$3 \times 10\text{-sec}$    $10\text{-sec}$

Mel-FilterBank     Mel-FilterBank

Conv2D     Conv2D

Merge Conv2D

Attention Feature Fusion

Other Encoder Layers    Audio Encoder

MLP Layer

Text Data

Captions     Labels

Keyword-to-Caption Augmentation

Captions

Text Encoder

MLP Layer

$E_1^t$   $E_2^t$   $E_3^t$   $\cdots$   $E_n^t$

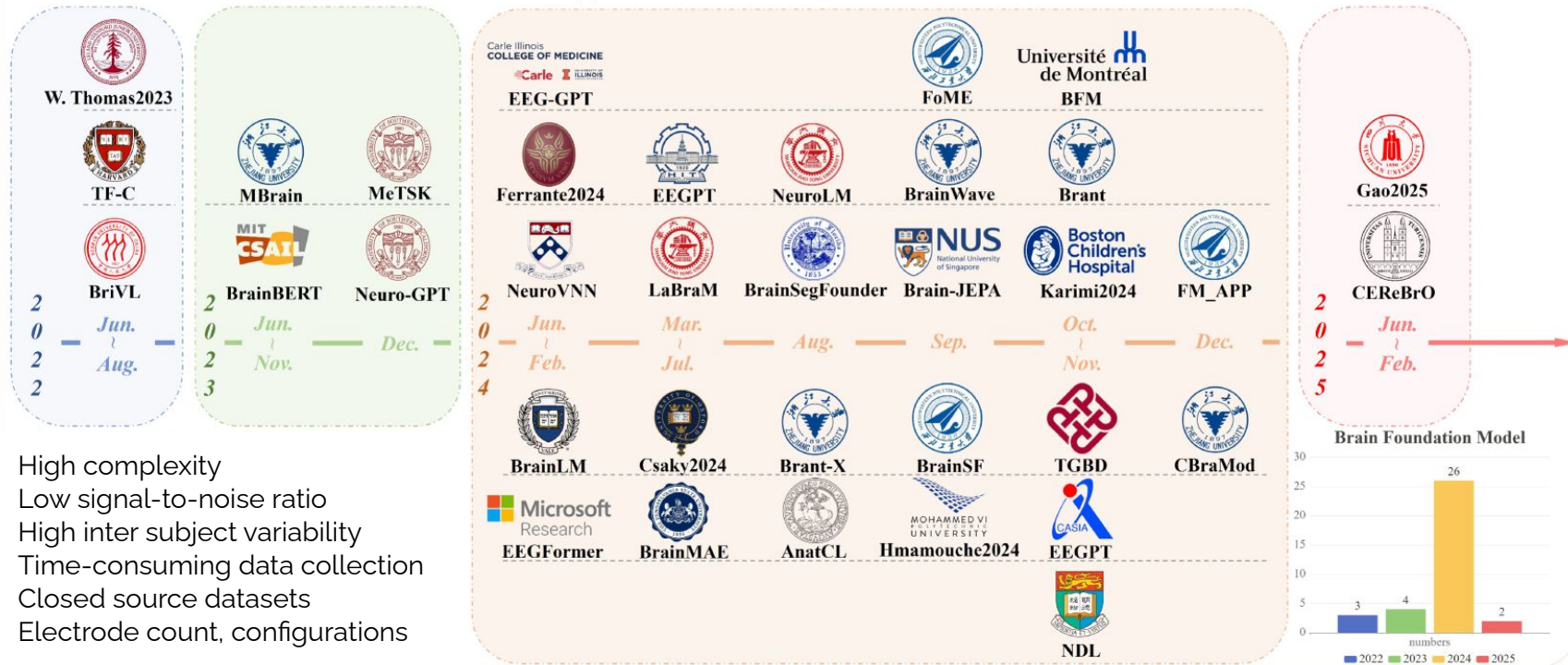| | $E_1^t$ | $E_2^t$ | $E_3^t$ | $\cdots$ | $E_n^t$ |
|---|---|---|---|---|---|
| $E_1^a$ | 1 | | | | |
| $E_2^a$ | | 1 | | | |
| $E_3^a$ | | | 1 | | |
| $\vdots$ | | | | 1 | |
| $E_n^a$ | | | | | 1 |

Brain Foundation Models



(d) Brain Foundation Models *Until February 7, 2025*

- High complexity
- Low signal-to-noise ratio
- High inter subject variability
- Time-consuming data collection
- Closed source datasets
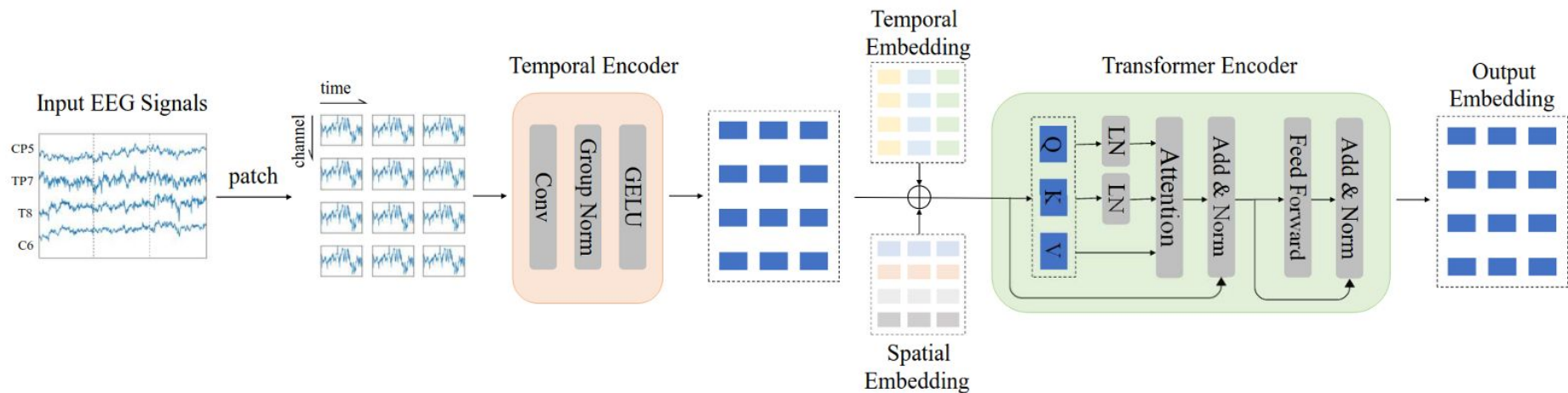- Electrode count, configurations

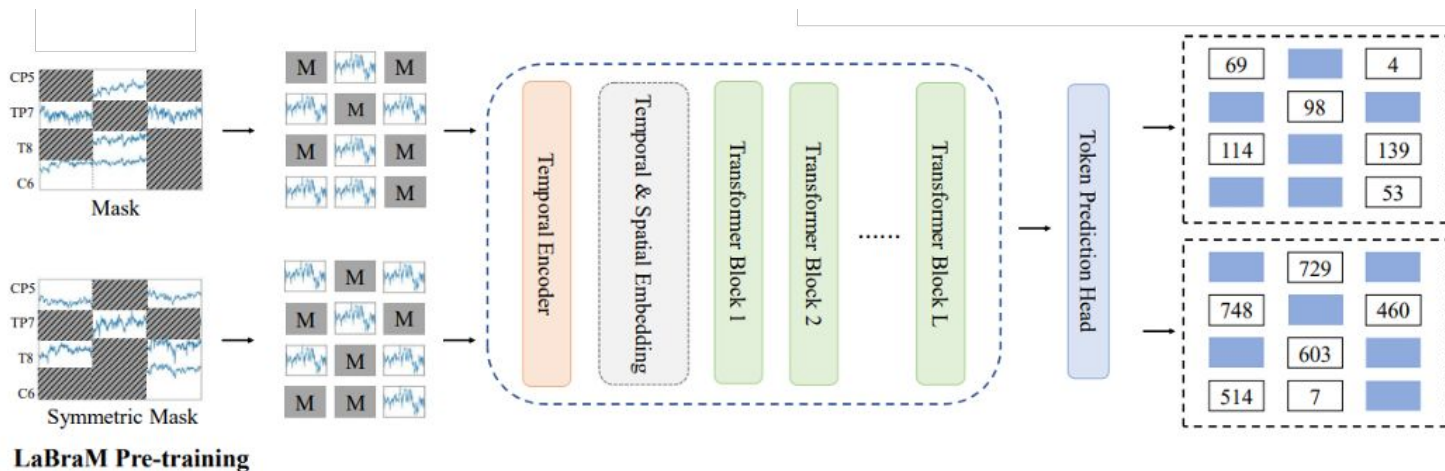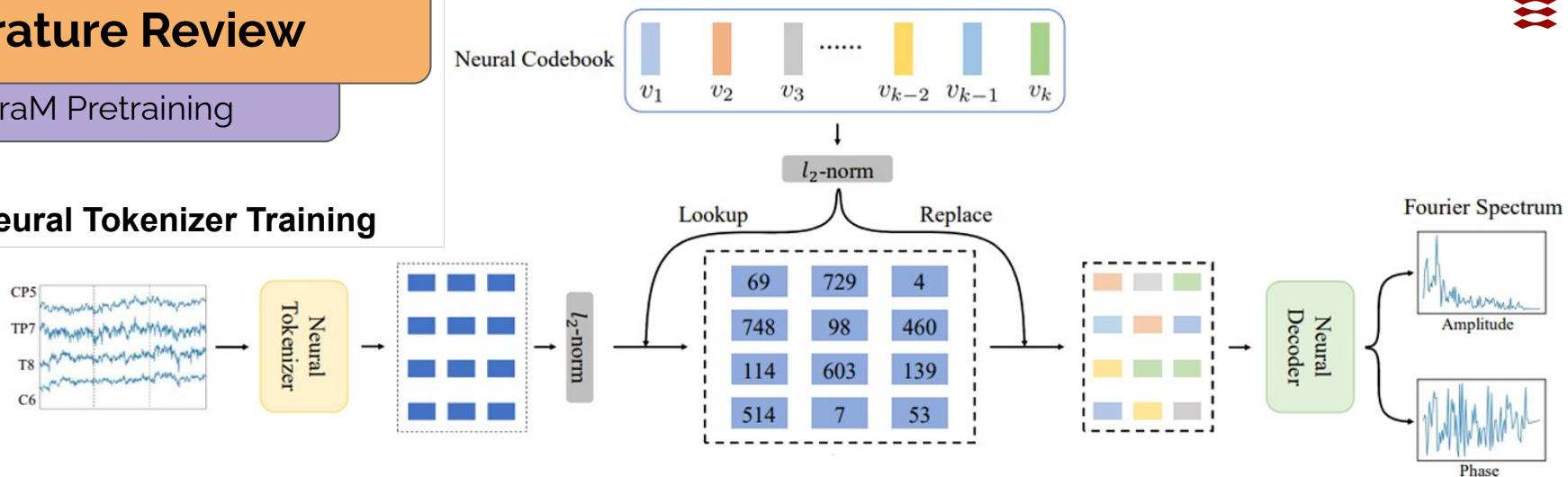## LaBraM

- Large Brain Model (LaBraM)

## Neural Transformer

**Neural Tokenizer Training**



**LaBraM Pre-training**
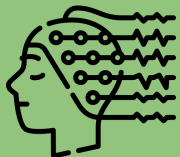
# Data
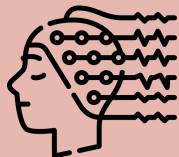
- 26 subjects
- Five conditions
- Male audio clips: 200, Female audio clips: 165
- Trial length: 1 minute

**Freefield**

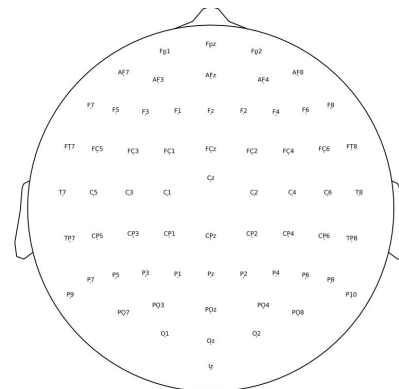**Insertphones**
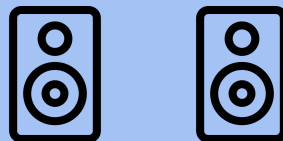
**-1,-4,-7dB**

- 3 subjects missing, left with 23 subjects
- 3364 trials

| Subject | 1 | 2 | 4 | 5 | 8 | 14 | 15 | 16 | 23 | 25 |
|---------|---|---|---|---|---|----|----|----|----|----|
| Insert | ✗ | ✗ | ✓ | ✓ | ✗ | ✓ | ✗ | ✗ | ✓ | ✓ |
| Free | ✗ | ✗ | ✓ | ✗ | ✓ | ✓ | ✗ | ✗ | ✓ | ✗ |
| -1dB | ✗ | ✗ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ |
| -4dB | ✗ | ✗ | ✗ | ✓ | ✓ | ✓ | ✗ | ✓ | ✓ | ✓ |
| -7dB | ✗ | ✗ | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ | ✗ | ✓ |

| Subject | Condition | # Missing Trials |
|---------|-----------|------------------|
| 10 | Insert | 16 |
| 20 | -7dB | 11 |
| 26 | Insert | 16 |
| 26 | -4dB | 15 |

# Data

2 yes/no questions per trial



Accuracy per Participant

Accuracy per Condition

Preprocessing

1. EEG was bandpass filtered between 0.5-30Hz

2. Independent Component Analysis (ICA) to remove EEG artifacts

3. EEG downsampled from 8192Hz → 200Hz
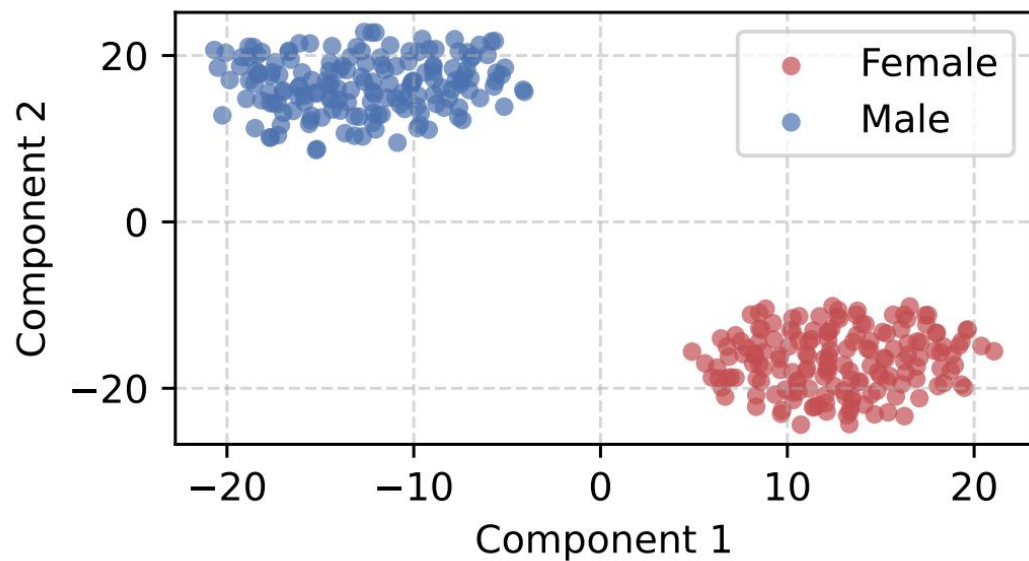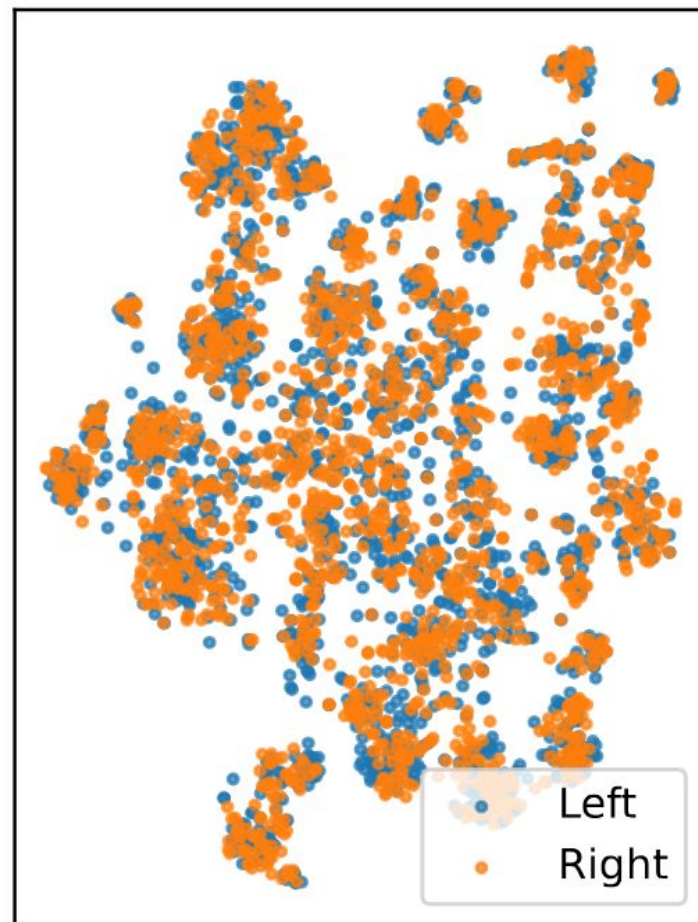
4. Audio upsampled from 44100Hz → 48000Hz

Direction as label

t-SNE of CLAP Embeddings by Gender

# Methodology

# Methodology

| Male Audio 1 | | |
|---|---|---|
| Trial 168 | Trial 588 | Trial 1065 |
| Trial 1237 | Trial 2215 | Trial 2716 |

↓

Training split

| Male Audio 4 | | |
|---|---|---|
| Trial 47 | Trial 216 | Trial 321 |
| Trial 473 | Trial 907 | Trial 1364 |

↓

Validation split

# Methodology

- Randomized trial segments
- Fixed validation segments
- Three augmentations:
  - Channel dropout
  - FT Surrogate
  - Time Reverse

```
1  # eeg_embed - EEG model embedding [n, d]
2  # audio_embed - Audio model embedding [n, d]
3  # target_ids - ids of audio segments [n]
4  # b, t_prime - learnable bias and temperature
5  # n - mini-batch size
6  eeg_embed_z = l2_normalize(eeg_embed)
7  audio_embed_z = l2_normalize(audio_embed)
8  t = exp(t_prime)
9  # ~ is used as a short hand for adding a new axis to an array to allow array
     broadcasting
10 labels = 2 * (target_ids[:, ~] == target_ids[~, :]) - ones(n,n)
11 logits = dot(eeg_embed_z, audio_embed_z.T) * t + b
12 loss = -sum(log_sigmoid(labels * logits)) / n
```
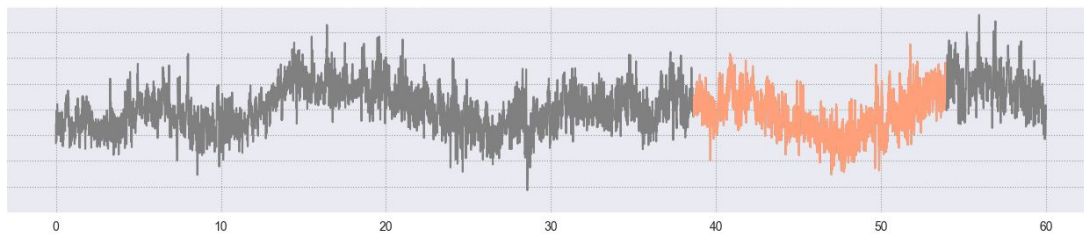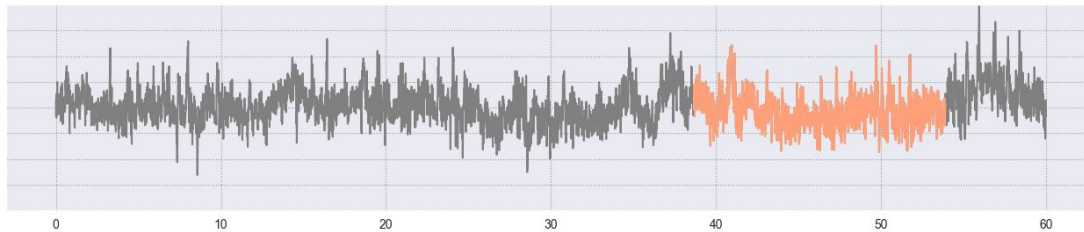
# Methodology

Contrastive learning

BX200  BX200  BX128

**LaBraM**  GELU  DROPOUT  Linear  GELU  DROPOUT  Linear

Dropout: 0.08
LR: 5e-4
Scheduler: OneCycle
Batch size: 32

SigLIP

BX512  BX200  BX128

**CLAP**  DROPOUT  Linear  GELU  DROPOUT  Linear

# Results & Discussion

# Results & Discussion

Baseline

- Each experiment used a 15 second decision window
- Only ran experiments with a single seed
- Backwards model

| # Conditions | Validation accuracy | Test accuracy |
|---|---|---|
| Two conditions | 0.643 | 0.604 |
| Five conditions | 0.564 | 0.568 |

Condition classification

# Results & Discussion

Contrastive learning

- Overfitting
- Memorization

## Temporal Split

## Contrastive learning

# Results & Discussion

Contrastive learning

## Temporal Split with mismatched EEG

## Noise-free conditions

Contrastive split accuracy



## All conditions

Contrastive split accuracy

## All conditions

- Better than random guessing
- High response accuracy + no missing data-> high model accuracy (9, 12, 24)

## Subject accuracy on all conditions

# Data

Response accuracy

2 yes/no questions per trial

**Accuracy per Participant**

Contrastive learning



Accuracy by Decision Window and Split

Augmentation results

- TR: Time Reversal
- DR: Channel Dropout
- FTS: Fourier Transform Surrogate

|  | Train | Val | Test |
|---|---|---|---|
| No aug | **0.989** | **0.752** | 0.702 |
| TR | 0.987 | 0.748 | **0.725** |
| DR | 0.946 | 0.711 | 0.667 |
| FTS | 0.905 | 0.714 | 0.694 |

ASAD



Direction as label — Gender as label. Accuracy bar charts. Direction as label: Train 0.874, Validation 0.665, Test 0.639. Gender as label: Train 0.926, Validation 0.695, Test 0.675.

## Direct classification

|  | Train | Validation | Test |
|---|---|---|---|
| Linear probe | 0.572 | 0.521 | 0.522 |
| LaBraM finetuning | **0.984** | **0.707** | **0.676** |
| Full finetuning | 0.722 | 0.523 | 0.492 |

# Conclusion

**RQ1:** How do CLAP and LaBraM perform as pretrained feature extractors for auditory attention decoding?

|  | Train | Validation | Test |
|---|---|---|---|
| Linear probe | 0.572 | 0.521 | 0.522 |

**RQ1:** How do CLAP and LaBraM perform as pretrained feature extractors for auditory attention decoding?



t-SNE of LaBraM EEG embeddings



t-SNE of CLAP Embeddings by Gender

**RQ2:** How does contrastive learning compare to supervised classification for training robust AAD models using CLAP and LaBraM?



Supervised split accuracy



Contrastive split accuracy

RQ3

**RQ3:** How does the length of decision windows affect performance?



Accuracy by Decision Window and Split

# Thank you for your Attention

# Appendix

## Baseline

- Each experiment used a 15 second decision window
- Only ran experiments with a single seed
- Backwards TRF model

### Two condition performance

| Split | Validation accuracy | Test accuracy |
|---|---|---|
| Temporal | 0.588 | 0.633 |
| Audio-disjoint | 0.643 | 0.604 |

### Five condition performance

| Split | Validation accuracy | Test accuracy |
|---|---|---|
| Temporal | 0.593 | 0.599 |
| Audio-disjoint | 0.564 | 0.568 |

# Literature Review

*[..] the process of stimulus reconstruction [..] is not optimized to effectively detect attention. [...] the compression of multichannel EEG signals into a single waveform through stimulus reconstruction reduces the available information for analysis[1]*

*[The neural network] outperforms the baseline linear stimulus reconstruction method, improving decoding accuracy [...] from 59% to 87%[2]*

*[...] correlation between the reconstructed and the attended speech envelopes is generally weak[3]*

[1]: Siqi Cai et al. "EEG-based Auditory Attention Detection in Cocktail Party Environment."
[2]: Gregory Ciccarelli et al. "Comparison of Two-Talker Attention Decoding from EEG with Nonlinear Neural Networks and Linear Methods."
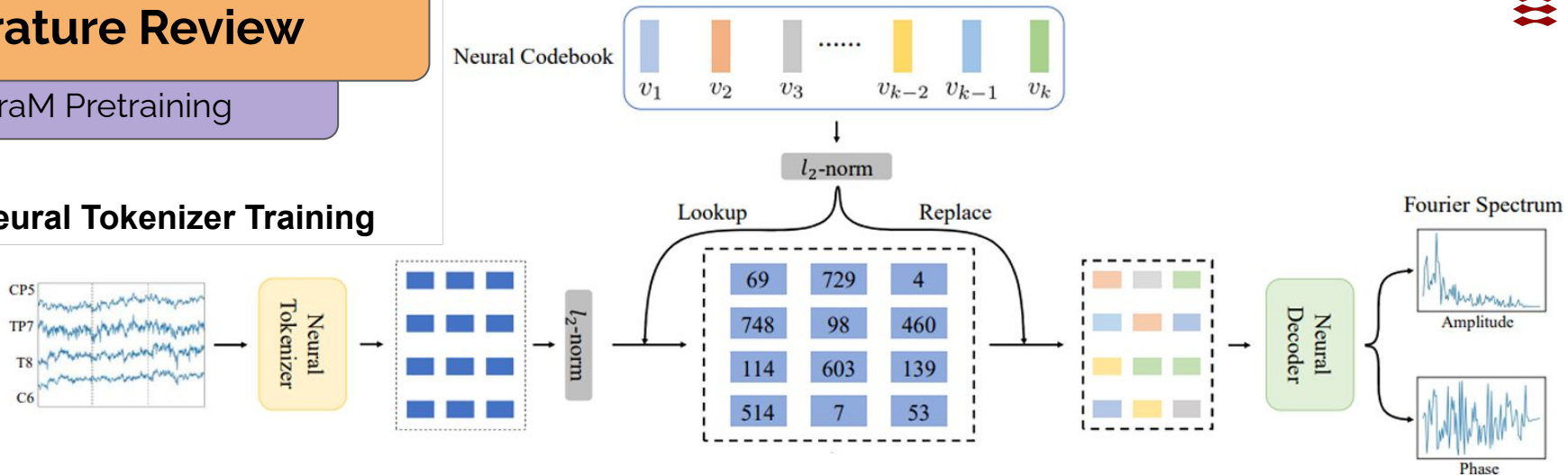[3]: Enze Su et al. "STAnet: A Spatiotemporal Attention Network for Decoding Auditory Spatial Attention From EEG."

LaBraM Pretraining

**Neural Tokenizer Training**



$$\mathcal{L}_T = \sum_{x \in \mathcal{D}} \sum_{i=1}^{N} \left\| o_i^A - A_i \right\|_2^2 + \left\| o_i^\phi - \phi_i \right\|_2^2 + \left\| \text{sg}(\ell_2(p_i)) - \ell_2(v_{z_i}) \right\|_2^2 + \left\| \ell_2(p_i) - \text{sg}(\ell_2(v_{z_i})) \right\|_2^2$$

Predicted amplitude   Predicted phase
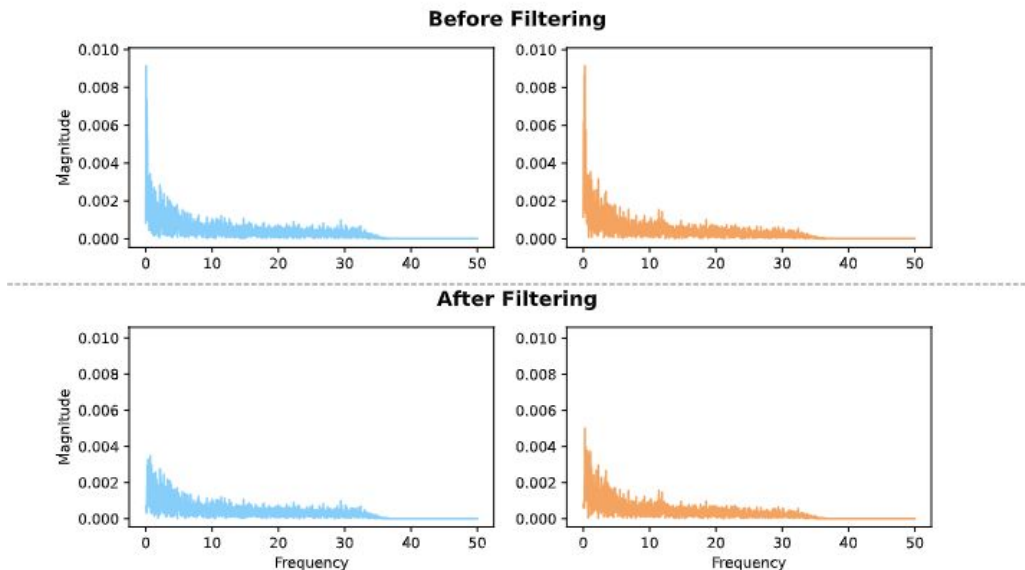
Tokenizer Vector

Actual amplitude   Actual phase
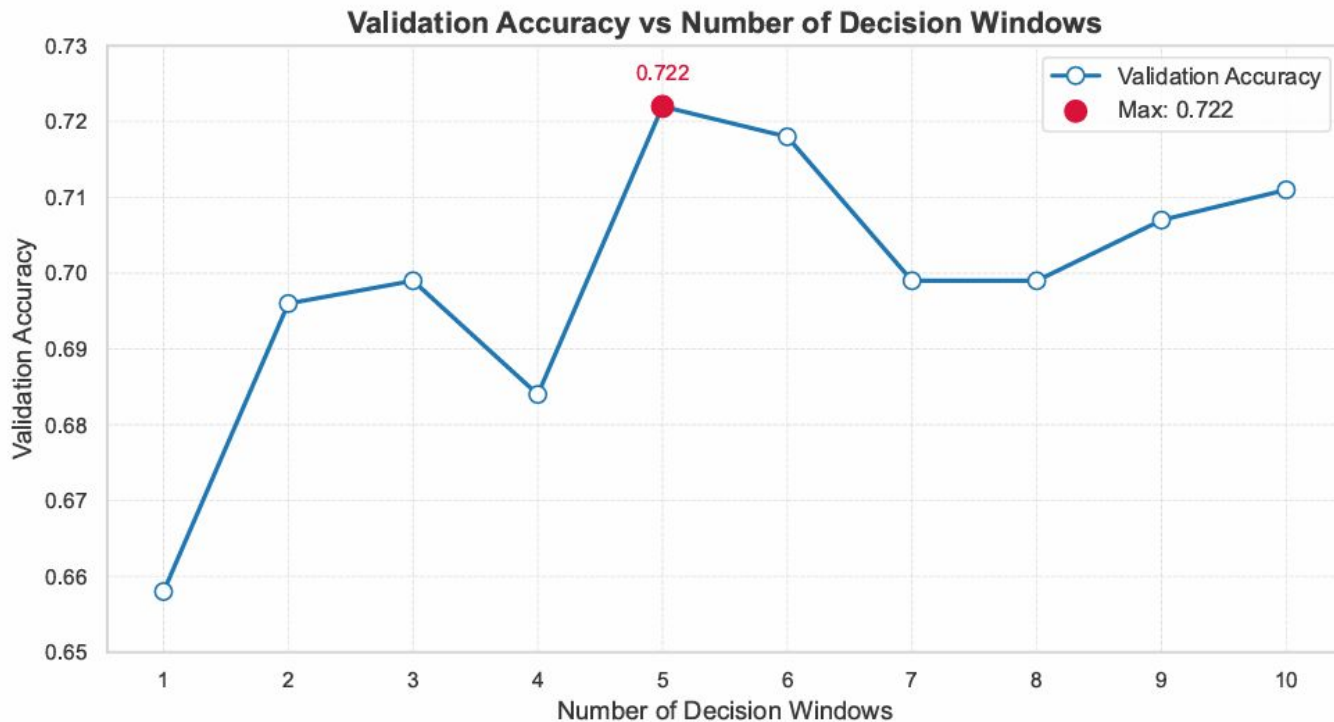
Codebook Vector

# Data

Preprocessing

- EEG was bandpass filtered between 0.5-30Hz
- ICA to remove EEG artifacts
- EEG downsampled from 8192Hz → 200Hz
- Audio upsamples from 44100Hz → 48000Hz



Before Filtering

After Filtering

# Results & Discussion

## Comparisons

**Lund Contrastive**
- Hearing impaired subjects
- Unspecified background noise
- CNN + attention
- Subject specific architecture

**Lund DCCA**
- No added background noise
- Whisper + Deep Canonical-correlation analysis

| | Lund Contrastive[1] | Lund DCCA[2] | Our Model |
|---|---|---|---|
| **Accuracy** | 71.5% | 67.9% | 67.0% |

(5 second decision window)

[1] Gautam Sridhar et al. "Improving auditory attention decoding i noisy environments for listeners with hearing impairment through contrastive learning"
[2] Alessandro Celoria et al. "An ASR-based Hybrid Approach for Auditory Attention Decoding"

Leave-one-out condition AAD performance