
All-for-One: Combining Small Convolutional Models for MEG-Based Speech Detection

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Decoding speech from non-invasive brain recordings remains a major challenge for
2 brain–computer interfaces. The LibriBrain Competition (NeurIPS 2025) (Landau
3 et al., 2025) addresses this question, with the aim of classifying moments of speech
4 and silence from magnetoencephalography (MEG) recordings. To accomplish this
5 task, we adapted a convolutional neural network (CNN) previously developed to
6 model a similar task (Défossez et al., 2023). Instead of training a single large model,
7 we took an ensemble modeling approach and trained multiple small convolutional
8 networks and combined their outputs using a lightweight transformer. Our results
9 show that in this task, small convolutional architectures, when aggregated, can
10 achieve strong performance.

11 1 Introduction

12 The 2025 LibriBrain competition aims to catalyze progress in non-invasive speech decoding by
13 providing a large-scale MEG dataset and standardized benchmarks (Landau et al., 2025). The
14 competition consisted of a classification task that classified speech versus silence from the MEG
15 signals of participants listening to continuous speech. We based our model on BrainMagick (Défossez
16 et al., 2023), which demonstrated the potential of CNNs for decoding speech from MEG signals.
17 One specificity of our team was that we worked with limited computational resources (two 4GB
18 Quadro T1000 GPUs). This prevented us from training large end-to-end architectures. Despite
19 these computational constraints, we show that small convolutional models trained independently can
20 achieve strong individual performance. Second, we demonstrate that combining several such models
21 through a transformer yields synergistic improvements, outperforming the best individual network.

22 2 Material and Methods

23 2.1 Dataset and Task

24 We used the LibriBrain dataset (Özdogan et al., 2025), comprising over 50 hours of preprocessed
25 MEG recordings (306 channels, 250 Hz) from a single participant listening to LibriVox audiobooks.
26 The input data x is a MEG epoch of $306 \text{ sensor} \times T \text{ samples}$ tensor. For each input data x we have
27 a corresponding label y , where $y = 0$ if silence, $y = 1$ if speech, $y = 2$ if speech onset and $y = 3$
28 if speech offset. "Speech onsets" corresponded to the first 25 samples, or 100 ms, of each speech
29 utterance (Hamilton et al., 2018). Similarly, "speech offsets" corresponded to the first 25 samples, or
30 100 ms, of silence after each speech utterance. Input data was sampled from the full MEG signals
31 with a stride of 1 sample.

32 Training data consisted in the standard LibriBrain training dataset (51.57 hours) minus the 8th chapter
33 of the 4th book that contained annotations errors ("0,8,*Sherlock4,1*"). Validation and test data were
34 the LibriBrain validation (0.36 hours) and test (0.38 hours) dataset respectively.

35 **2.2 Model Architecture**

36 We used a two-stage architecture in which (1) multiple small CNNs generated class logits that were
37 (2) subsequently combined by a transformer for final prediction.

38 The CNN architecture was based on BrainMagick(Défossez et al., 2023) and adapted for small-scale
39 GPU setups. We used the same spatial attention layer and the same temporal convolutional structure.
40 We decreased the width of the convolutional blocks from 320 to 64 channels. We increased the
41 number of convolutional blocks with residual connections from 5 to 7. Finally, we added a classifier
42 head at the end. The final architecture consisted in a spatial attention pooling layer with 128 channels,
43 a layer of temporal convolution with 64 channels and kernel size 1, a layer of temporal convolution
44 with 32 channels and kernel size 1, 7 blocks of temporal convolution with kernel size 3 and dilation
45 from 1 to 64 with residual skip connections, a classification head of 3 fully connected layers with
46 256, 32 and 2 channels respectively.

47 We trained four first-stage CNNs varying in window length and label set: 3 s, 5 s, and 7 s windows
48 for binary speech/silence classification, and a fourth 5 s model distinguishing speech, silence, speech
49 onset, and speech offset (Hamilton et al., 2018) (see Table 1).

50 The transformer received as input the concatenated logits from all CNNs. It used positional encoding
51 and operated on sequences of length 200 with a batch size of 32, a model dimension of 64, four
52 attention heads, two layers, and a dropout rate of 0.2.

53 **2.3 Regularization**

54 The final prediction of the transformer model was further regularized to account for the global
55 statistics of the training dataset: silence segments of less than 100 ms were replaced by speech and
56 speech segments of less than 300 ms were replaced by silence.

57 **2.4 Training details**

58 Due to limited computational resources, each model was trained only once without hyperparameter
59 optimization. For the CNNs, we used a cross-entropy loss and the AdamW optimizer (learning rate =
60 3×10^{-4} , batch size = 128). Validation loss was computed every 1k steps. Training was stopped
61 after 10k steps without improvement in the validation loss. For the transformer, we used a binary
62 cross-entropy with logits loss and the AdamW optimizer (learning rate = 3×10^{-4} , batch size = 32),
63 training until the full training dataset had been traversed once.

64 **3 Results**

65 All four convolutional models achieved strong individual performance, with F1-macro scores con-
66 sistently above 85%, largely beating the competition benchmark (F1-macro = 68%) on the test
67 set. Combining their logits through the transformer led to a consistent improvement, reaching an

Table 1: Speech detection performance (F1-macro, %) of individual convolutional models, their transformer ensemble, and the final regularized system on the LibriBrain dataset.

Model	# labels	Window size [s]	F1-macro (test)	F1-macro (held-out)
Individual CNN 1	2	2	86.9	–
Individual CNN 2	2	5	88.8	–
Individual CNN 3	2	7	87.2	–
Individual CNN 4	4	5	86.6	–
+ transformer ensemble	–	–	89.3	–
+ regularization	–	–	89.5	90.3

68 F1-macro of 89.3%. Applying the corpus-level regularization further increased performance to 89.5%
69 on the test data and 90.3% on the held-out evaluation set.

70 **4 Discussion**

71 Small convolutional networks remain efficient baselines for MEG-based speech decoding, even under
72 computational constraints. When combined through a lightweight transformer, these models form
73 an ensemble that combines the complementary representations learned by each CNN, resulting in
74 improved performance. Our suggests that ensemble-based strategies allow to enhance performance
75 without requiring large models, providing a simple approach to non-invasive speech decoding.

76 **References**

- 77 Landau, G.; Özdogan, M.; Elvers, G.; Mantegna, F.; Somaia, P.; Jayalath, D.; Kurth, L.; Kwon, T.;
78 Shillingford, B.; Farquhar, G.; others The 2025 PNPL competition: Speech detection and phoneme
79 classification in the LibriBrain dataset. *arXiv preprint arXiv:2506.10165* **2025**,
- 80 Défossez, A.; Caucheteux, C.; Rapin, J.; Kabeli, O.; King, J.-R. Decoding speech perception from
81 non-invasive brain recordings. *Nature Machine Intelligence* **2023**, 5, 1097–1107.
- 82 Özdogan, M.; Landau, G.; Elvers, G.; Jayalath, D.; Somaia, P.; Mantegna, F.; Woolrich, M.;
83 Jones, O. P. LibriBrain: Over 50 Hours of Within-Subject MEG to Improve Speech Decoding
84 Methods at Scale. *arXiv preprint arXiv:2506.02098* **2025**,
- 85 Hamilton, L. S.; Edwards, E.; Chang, E. F. A spatial map of onset and sustained responses to speech
86 in the human superior temporal gyrus. *Current Biology* **2018**, 28, 1860–1871.