

# Multi-level representations in speech processing in brain and machine: Evidence from EMEG and RSA

✉ cw417@cam.ac.uk

Cai Wingfield<sup>1,3</sup>, Li Su<sup>2</sup>, Barry Devereux<sup>1</sup>, Xunying Liu<sup>3,4</sup>, Chao Zhang<sup>3</sup>, Phil Woodland<sup>3</sup>, Elisabeth Fonteneau<sup>1</sup>, Andrew Thwaites<sup>1</sup>, William Marslen-Wilson<sup>1,5</sup>

<sup>1</sup>Department of Psychology, <sup>2</sup>Department of Psychiatry, <sup>3</sup>Department of Engineering; University of Cambridge

<sup>4</sup>Department of Systems Engineering and Engineering Management; Chinese University of Hong Kong

<sup>5</sup>MRC Cognition & Brain Sciences Unit; Cambridge



## Introduction

### Human speech recognition

- For humans, speech recognition feels effortless and automatic.
- There is only limited neurocomputational understanding of how this is achieved.
- Recent evidence suggests that responses to speech may be represented in a low-dimensional space of articulatory features<sup>1,2</sup>.

### Machine speech recognition

- Automatic speech recognition (ASR) systems using deep neural network (DNN) acoustic models approach human levels of performance, with word identification rates well over 90%<sup>3</sup>.
- They provide a computationally specific model of how successful speech recognition can be achieved.
- DNNs have been successful models for brain responses in other domains, e.g. vision<sup>4</sup>.

### The questions we ask

- Can we model human neural responses to speech with an ASR-derived DNN acoustic model?
- Can we characterise hidden-layer representations?
- Do activations in different layers of such a DNN differentially explain cortical speech responses through space and time?

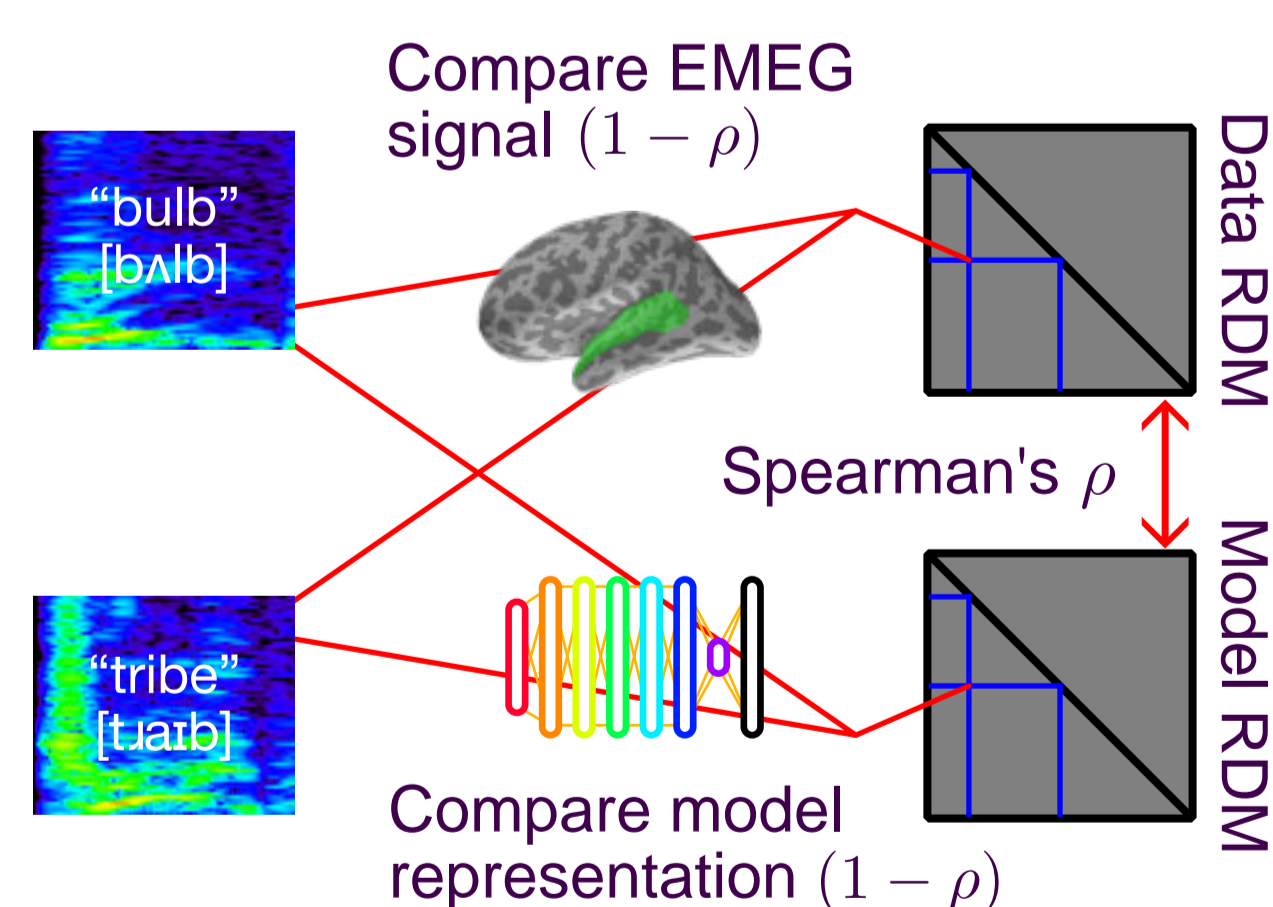
- Here we use a DNN-based ASR system to model brain responses to speech using multivariate searchlight techniques<sup>5</sup>.

## EMEG experiment

- Whole-brain simultaneous electro- and magneto-encephalography (EMEG) data was recorded from 16 participants while they listened to recordings of 400 British English words.
- EMEG data was source-localised and warped to a subject-average cortical mesh.
- The mesh was restricted to an auditory cortex (AC) mask comprising bilateral superior temporal cortices (STC) and Heschl's gyrus (HG).
- The mask matches locations previously found to exhibit phonetic feature sensitivity to speech<sup>1,2</sup>.

## ssRSA

- Representational similarity analysis (RSA<sup>6</sup>) compares condition space representations using a representational distance matrix (RDM) with pairwise condition distances as entries.
- Data RDMs are computed from brain responses; model RDMs from model representations or predicted distances.



- RSA abstracts away from specific responses, allowing comparison of representations in very different formats<sup>4,6</sup>.
- Spatiotemporal searchlight RSA (ssRSA<sup>7</sup>) computes data RDMs from a continuously moving regular spatiotemporal searchlight patch.

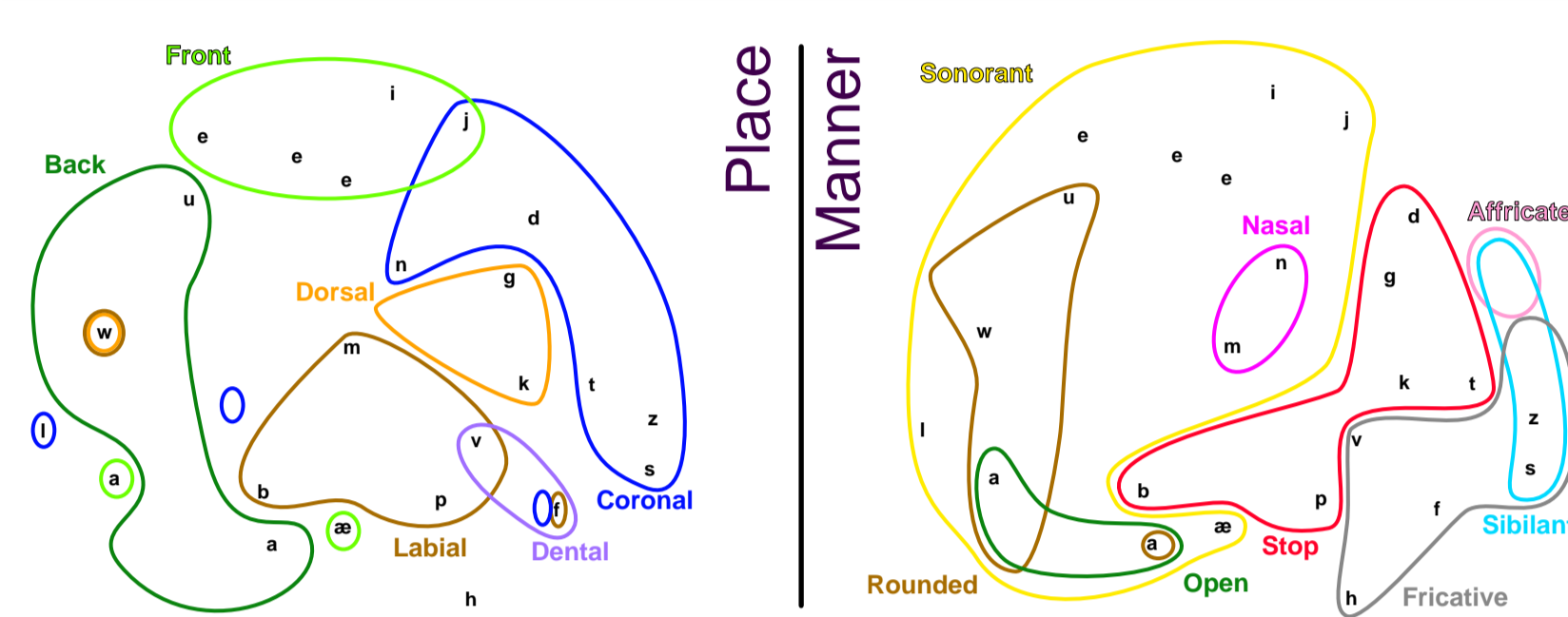
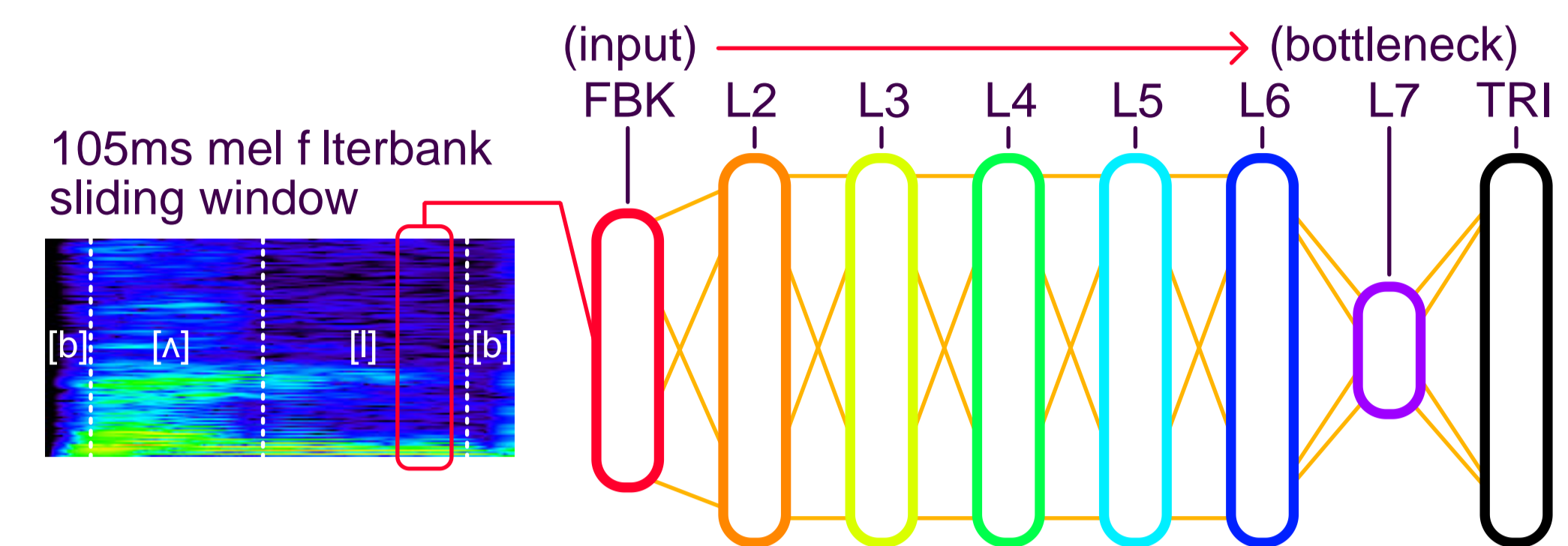
## References

1. Wingfield C, et al. (submitted) *PLOS Comp Biol.*
2. Mesgarani N, et al. (2014) *Science.*
3. Young S, et al. (2015) *The HTK Book.*
4. Khaligh-Razavi S, et al. (2014) *PLOS Comp Biol.*
5. Kriegeskorte N, et al. (2006) *PNAS.*
6. Kriegeskorte N, et al. (2008) *Front Syst Neurosci.*
7. Su L, et al. (2012) *PRNI, IEEE.*
8. Smith S, et al. (2009) *NeuroImage.*

## Computational modelling and analysis

### Deep neural network speech recogniser

- To model the speech recognition process, we used HTK<sup>3</sup>, an ASR system.
- HTK uses a fully-connected, feed-forward DNN as an acoustic model, mapping speech sounds to phonetic targets through time.
- DNN trained on ~700 hours of subtitled TV audio.
- Input is speech audio with a mel-frequency filterbank applied (FBK).
- Between input (FBK) and output (TRI) layers, information passes through five high-dimensional (1000-node) hidden layers (L2–L6) and a low-dimensional (26-node) "bottleneck" layer (L7).
- Output layer gives posteriors for ~6000 phonetic targets to a set of hidden Markov models (HMMs) associated with phonemes used in the ASR acoustic model.



### Visualising hidden-layer representations

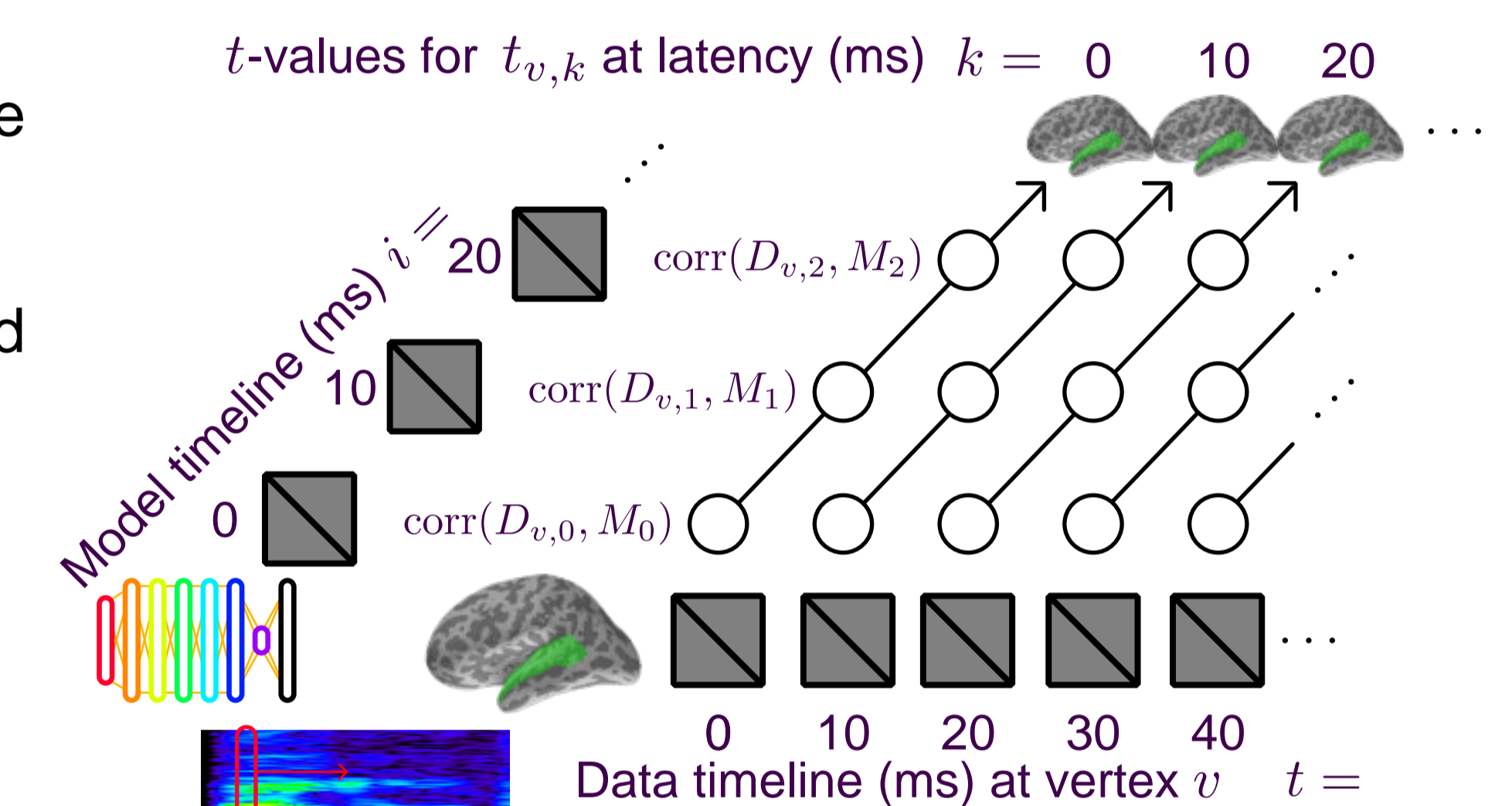
- Average L7 response to different phones can be visualised by Sammon nonlinear multidimensional scaling.
- Bottleneck responses clustered according to place and manner of articulation of consonants, position of vowels, and broad category distinctions (e.g. sonorant–obstruent).

### Computational mapping

- Model RDMs were produced from activations in each layer of the DNN as HTK listened to the same recorded words as the participants.
- We fitted the model RDM time course to the data RDM time course systematically at different processing lags (0–250ms).
- This gave time-resolved maps of the fit of each model to the brain data in individual subjects at each lag, which we converted to *t*-maps across subjects.
- Threshold-free cluster enhancement (TFCE) was performed on the resultant *t*-maps<sup>8</sup>.

$$\text{TFCE}(t_{v,k}) = \int_{h=0}^{t_{v,k}} \sqrt{e(h)} \cdot h^2 dh$$

- Random-effects permutation across subjects was used to determine the significance of the TFCE values.
- Resultant maps were thresholded at  $p < 0.001$ .

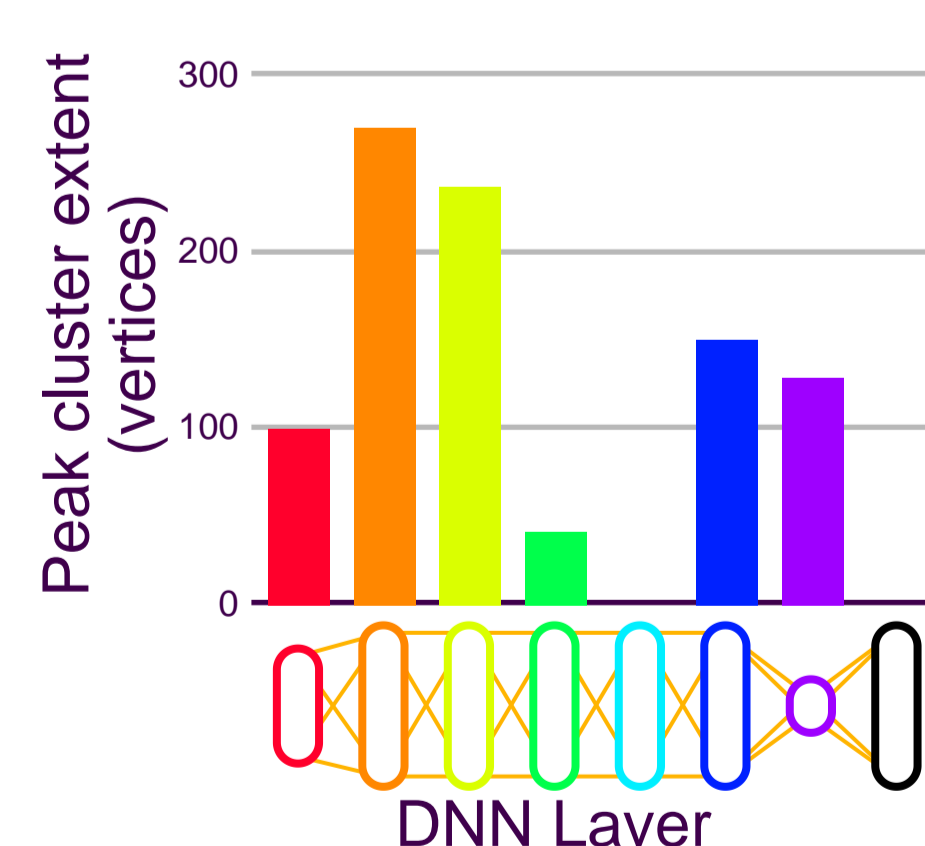


## Results

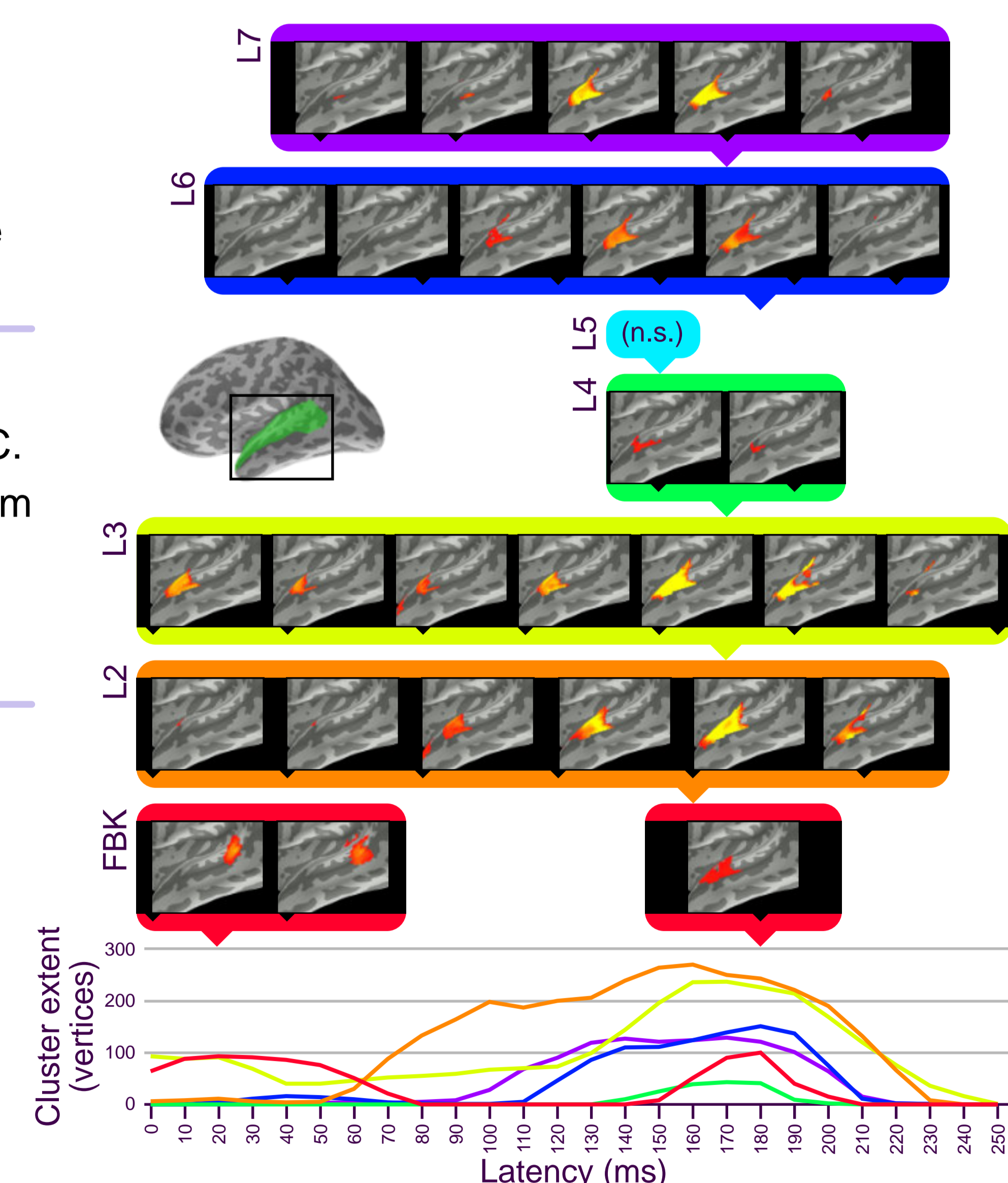
- Brain renderings show left-hemisphere supra-threshold vertices of TFCE *t*-maps ( $p < 0.01$ ).
- Line graph shows the time-course of each layer as it attains its maximum cluster size in left STC.
- Bar graph shows peak left-hemisphere cluster size over the epoch.

### Differential fits of DNN layers in left STC

- Input layer FBK peaked early (0–70ms) in left posterior STC.
- Layers L2–L4 and L6–L7 peaked later, achieving a maximum at ~170ms.
- Layers L5 and TRI showed no significant fit in the left hemisphere at this threshold.



- Overall, fit improved between layers FBK–L3, diminished for L4–L5, and re-emerged for L6–L7 (see bar graph).
- Right STC showed an early peak (~0–120ms) which did not distinguish between layers.



## Conclusions

- We related different layers of a DNN-based computational model of speech recognition to EMEG data recorded from humans listening to speech.
- Several hidden-layer models fit the brain data better than the input and output layer models.
- Earlier hidden-layer models L2 and L3 (closer to acoustic representation) fit brain data well.
- Higher hidden layers (L4 and L5) failed to fit brain data at the same level, perhaps indicating a divergence of human and machine representations.
- However at L6 and the compressed code of L7, where activations seem to cluster by articulatory features, the hidden-layer models again correlated with brain representations.
- Posterior–anterior STC showed changes in model fit over lags and layers, which may reflect progressively more specific representations of speech, with anterior STC representations relating to phonetic rather than acoustic features.