



OPEN ACCESS

EDITED BY
Yuqi Han,
Tsinghua University, China

REVIEWED BY
Yuxuan Zhao,
Institute of Automation (CAS), China
Alec Marantz,
New York University, United States

*CORRESPONDENCE
Li Su
l.su@sheffield.ac.uk

†These authors have contributed
equally to this work and share first
authorship

RECEIVED 29 September 2022
ACCEPTED 29 November 2022
PUBLISHED 21 December 2022

CITATION
Wingfield C, Zhang C, Devereux B,
Fonteneau E, Thwaites A, Liu X,
Woodland P, Marslen-Wilson W and
Su L (2022) On the similarities of
representations in artificial and brain
neural networks for speech
recognition.
Front. Comput. Neurosci. 16:1057439.
doi: 10.3389/fncom.2022.1057439

COPYRIGHT
© 2022 Wingfield, Zhang, Devereux,
Fonteneau, Thwaites, Liu, Woodland,
Marslen-Wilson and Su. This is an
open-access article distributed under
the terms of the [Creative Commons
Attribution License \(CC BY\)](https://creativecommons.org/licenses/by/4.0/). The use,
distribution or reproduction in other
forums is permitted, provided the
original author(s) and the copyright
owner(s) are credited and that the
original publication in this journal is
cited, in accordance with accepted
academic practice. No use, distribution
or reproduction is permitted which
does not comply with these terms.

On the similarities of representations in artificial and brain neural networks for speech recognition

Cai Wingfield^{1†}, Chao Zhang^{2†}, Barry Devereux³,
Elisabeth Fonteneau⁴, Andrew Thwaites⁵, Xunying Liu⁶,
Phil Woodland², William Marslen-Wilson⁵ and Li Su^{7,8*}

¹Department of Psychology, Lancaster University, Lancaster, United Kingdom, ²Department of Engineering, University of Cambridge, Cambridge, United Kingdom, ³School of Electronics, Electrical Engineering and Computer Science, Queens University Belfast, Belfast, United Kingdom, ⁴Department of Psychology, University Paul Valéry Montpellier, Montpellier, France, ⁵Department of Psychology, University of Cambridge, Cambridge, United Kingdom, ⁶Department of Systems Engineering and Engineering Management, The Chinese University of Hong Kong, Shatin, Hong Kong SAR, China, ⁷Department of Neuroscience, Neuroscience Institute, Insigneo Institute for in silico Medicine, University of Sheffield, Sheffield, United Kingdom, ⁸Department of Psychiatry, University of Cambridge, Cambridge, United Kingdom

Introduction: In recent years, machines powered by deep learning have achieved near-human levels of performance in speech recognition. The fields of artificial intelligence and cognitive neuroscience have finally reached a similar level of performance, despite their huge differences in implementation, and so deep learning models can—in principle—serve as candidates for mechanistic models of the human auditory system.

Methods: Utilizing high-performance automatic speech recognition systems, and advanced non-invasive human neuroimaging technology such as magnetoencephalography and multivariate pattern-information analysis, the current study aimed to relate machine-learned representations of speech to recorded human brain representations of the same speech.

Results: In one direction, we found a quasi-hierarchical functional organization in human auditory cortex qualitatively matched with the hidden layers of deep artificial neural networks trained as part of an automatic speech recognizer. In the reverse direction, we modified the hidden layer organization of the artificial neural network based on neural activation patterns in human brains. The result was a substantial improvement in word recognition accuracy and learned speech representations.

Discussion: We have demonstrated that artificial and brain neural networks can be mutually informative in the domain of speech recognition.

KEYWORDS

automatic speech recognition, deep neural network, representational similarity analysis, auditory cortex, speech recognition

1. Introduction

Speech comprehension—the ability to accurately identify words and meaning in a continuous auditory stream—is a cornerstone of the human communicative faculty. Nonetheless, there is still limited understanding of the neurocomputational representations and processes in the human brain which underpin it. In this paper we approach a fundamental component of speech comprehension—namely the recognition of word identities from the sound of speech—in reverse: to find artificial systems which can accomplish the task, and use them to model and probe the brain's solution. In the domain of engineering, automatic speech recognition (ASR) systems are designed to identify words from recorded speech audio. In this way, ASR systems provide a computationally explicit account of how speech recognition *can* be achieved, so correspondences between the human and machine systems are of particular interest; specifically, the question of whether the learned representations in an ASR can be linked to those found in human brains. Modern advances in high-resolution neuroimaging and multivariate pattern-information analysis have made this investigation feasible.

In the present research, we took a bidirectional approach, relating machine-learned representations of speech to recorded brain representations of the same speech. First, we used the representations learned by an ASR system with deep neural network (DNN) acoustic models (Hinton et al., 2012) to probe the representations of heard speech in the brains of human participants undergoing continuous brain imaging. This provided a mechanistic model of speech recognition, and evidence of it matching responses in human auditory cortex. Then, in the opposite direction, we used the architectural patterns of neural activation we found in the brains to refine the DNN architecture and demonstrated that this improves ASR performance. This bidirectional approach was made possible by recently developed multivariate pattern analysis methods capable of comparing learned speech representations in living brain tissue and computational models.

ASR encompasses a family of computationally specified processes which perform the task of converting recorded speech sounds to the underlying word identities. Modern ASR systems employing DNN acoustic and language models now approach human levels of word recognition accuracy on specific tasks. For instance, regarding English, the word error rate (WER) of transcribing careful reading speech with no background noise can be lower than 2% (Luscher et al., 2019; Park et al., 2019), and the WER of transcribing spontaneous conversational telephone speech can be lower than 6% (Saon et al., 2017; Xiong et al., 2018).

For the present study, our ASR system was constructed based on a set of hidden Markov models (HMMs). For each, a designated context-dependent phonetic unit handled the transitions between the hidden states. A DNN model was used

to provide the observation probability of a speech feature vector given each HMM state. This framework is often called a “hybrid system” in the ASR literature (Bourlard and Morgan, 1993; Hinton et al., 2012). The Hidden Markov Model Toolkit (HTK; Young et al., 2015; Zhang and Woodland, 2015a) represents a historical state-of-the-art ASR system, and is still among the most widely used. We used HTK to train the DNN-HMMs and construct the overall ASR pipeline of audio to text. A version of this model comprised a key part of the first-place winner of the multi-genre broadcast (MGB) challenge of the IEEE Automatic Speech Recognition and Understanding Workshop 2015 (Bell et al., 2015; Woodland et al., 2015). In this paper, all ASR systems were built in HTK using 200 h of training data from the MGB challenge. We designed the experimental setup carefully to use only British English speech and reduce the channel difference caused by different recording devices.

Of particular importance for the present study is the inclusion of a low-dimensional *bottleneck* layer in the DNN structure of our initial model. Each of the first five hidden layers contains 1,000 nodes, while the sixth hidden layer has just 26 nodes. Our choice to include six hidden layers in the DNN is not arbitrary. The performance of different DNN structures in the MGB challenge has previously been studied. Empirically, having a fewer hidden layers result in worse WERs, while more hidden layers result in unstable training performance due to the increased difficulty when optimizing deeper models. Similar structures were often adopted on different datasets and by different groups (e.g., Karafiát et al., 2013; Doddipatla et al., 2014; Yu et al., 2014; Liu et al., 2015). Since the layers in our DNN are feed-forward and fully connected, each node in each layer is connected only with the nodes from its immediately preceding layer, and as such the acoustic feature representations of the input speech are forced to pass through each layer in turn to derive the final output probabilities of the context-dependent phonetic units. The bottleneck layer representations are highly compressed and discriminative, and are therefore widely used as an alternative type of input features to acoustic models in ASR literature¹ (Grézl et al., 2007; Tüske et al., 2014; Woodland et al., 2015). In addition, the inclusion of this bottleneck layer greatly reduces the number of DNN parameters without significantly diminishing the accuracy of word recognition (Woodland et al., 2015), since it can prevent the model from over-fitting to the training data (Bishop, 2006). Thus, the bottleneck layer representation provides a learned, low-dimensional representation of speech which is both parsimonious and sufficient for high-performance speech recognition. This is especially interesting for the present study, given the inherently low-dimensional parameterization of speech that is given by articulatory features, which are a

¹ Bottleneck layers which are trained alongside the other layers in a model have been shown to be superior to other methods of lowering dimensions, such as simple PCA (Grézl et al., 2007).

candidate characterization of responses to speech in human auditory cortex.

Recent electrocorticography (ECoG: Mesgarani et al., 2008, 2014; Chang et al., 2010; Di Liberto et al., 2015; Moses et al., 2016, 2018) and functional magnetic resonance imaging (fMRI: Arsenault and Buchsbaum, 2015; Correia et al., 2015) studies in humans show differential responses to speech sounds exhibiting different articulatory features in superior temporal speech areas. Heschl's gyrus (HG) and surrounding areas of the bilateral superior temporal cortices (STC) have also shown selective sensitivity to perceptual features of speech sounds earlier in the recognition process (Chan et al., 2014; Moerel et al., 2014; Saenz and Langers, 2014; Su et al., 2014; Thwaites et al., 2016). Building on our previous work investigating phonetic feature sensitivity in human auditory cortex (Wingfield et al., 2017), we focus our present analysis within language-related brain regions: STC and HG.

The neuroimaging data used in this study comes from electroencephalography and magnetoencephalography (MEG) recordings of participants listening to spoken words in a magnetoencephalography (MEG) brain scanner. High-resolution magnetic resonance imaging (MRI) was acquired using a 3T MRI scanner for better source localization. As in our previous studies (Fonteneau et al., 2014; Su et al., 2014; Wingfield et al., 2017), the data (MEG and MRI) has been combined to generate a source-space reconstruction of the electrophysiological activity which gave rise to the measurements at the electroencephalography (EEG) and MEG sensors. Using standard minimum-norm estimation (MNE) procedures guided by anatomical constraints from structural MRIs of the participants (Hämäläinen and Ilmoniemi, 1994; Gramfort et al., 2014), sources were localized to a cortical mesh at the gray-matter–white-matter boundary. Working with source-space activity allows us to retain the high temporal resolution of MEG, while gaining access to resolved spatial pattern information. It also provides the opportunity to restrict the analysis to specific regions of interest on the cortex, where an effect of interest is most likely to be found.

Recent developments in multivariate neuroimaging pattern analysis methods have made it possible to probe the representational content of recorded brain activity patterns. Among these, representational similarity analysis (RSA: Kriegeskorte et al., 2008a) provides a flexible approach which is well-suited to complex computational models of rich stimulus sets. The fundamental principle of our RSA procedures was the computation of the similarity structures of the brain's response to experimental stimuli, and comparing the similarity structures with those derived from computational models. In a typical RSA study, this similarity structure is captured in a representational dissimilarity matrix (RDM), a square symmetric matrix whose rows and columns are indexed by the experimental stimuli, and whose entries give values for the dissimilarity of two conditions, as given by their correlation distance in the response space.

A key strength of RSA is that RDMs abstract away from the specific implementation of the DNN model or measured neural response, allowing direct comparisons between artificial and human speech recognition systems; the so-called “dissimilarity trick” (Kriegeskorte and Kievit, 2013). The comparison between RDMs computed from the ASR model and RDMs from human brains take the form of a Spearman's rank correlation ρ between the two (Nili et al., 2014).

RSA has been extended using the fMRI searchlight-mapping framework (Kriegeskorte et al., 2006; Nili et al., 2014) so that representations can be mapped through image volumes. Subsequently, searchlight RSA has been further extended into the temporal dimension afforded by MEG data: spatiotemporal searchlight RSA (ssRSA: Su et al., 2012, 2014). Here, as in other studies using computational cognitive models (e.g., Khaligh-Razavi and Kriegeskorte, 2014; Mack et al., 2016), ssRSA facilitates the comparison to a machine representation of the stimulus space which may otherwise be incommensurable with a distributed brain response.

In the machine-to-human direction, using ssRSA and the ASR system as a reference, we found that the early layers of the DNN corresponded to early neural activation in primary auditory cortex, i.e., bilateral Heschl's gyrus, while the later layers of the DNN corresponded to late activation in higher level auditory brain regions surrounding the primary sensory cortex. This finding reveals that the neural network located within HG is likely to have a similar functional role as early layers of the DNN model, extracting basic acoustic features (though see Hamilton et al., 2021 for a recent contrasting study). The neurocomputational function of superior temporal gyrus regions is akin to later layers of the DNN, computing complex auditory features such as articulation and phonemic information.

In the reverse human-to-machine direction, using the pattern of results in the brain-image analysis, we improved the architecture of the DNN. The spatial extent of neural activation explained by the hidden-layer representations progressively reduced for higher layers, before expanding again for the bottleneck layer. This pattern, which mirrored the structure of the DNN itself, and (assuming an efficient and parsimonious processing stream in the brain) suggests that some pre-bottleneck layers might be superfluous in preparing the low-dimensional bottleneck compression. We restructured the DNN model with the bottleneck layer moved to more closely resemble the pattern of activation observed in the brain, hypothesizing that this would lead to a better transformation. With this simple, brain-inspired modification, we significantly improved the performance of the ASR system. It is notable that similar DNN structures have been developed independently elsewhere in order to optimize the low-dimensional speech feature representations from the DNN bottleneck layer. However, “reverse-engineer” human learning systems implemented in brain tissue in such a bidirectional fashion provides a

complementary approach in developing and refining DNN learning algorithms.

2. Study 1: Investigating ASR DNN representations

2.1. Materials and methods

2.1.1. Building DNN-HMM acoustic models for ASR

Here we construct a DNN which can each be included as a component in the hybrid DNN-HMM set-up of HTK. This is a widely used speech recognition set-up in both academic and industrial communities (Hinton et al., 2012), whose architecture is illustrated in Figure 1. Each network comprises an input layer, six hidden layers, and an output layer, which are all fully-connected feed-forward layers.

The DNN acoustic model was trained to classify each input frame into one of the triphone units at each time step. We used it as the acoustic model of our DNN-HMM ASR system to estimate the triphone unit likelihoods corresponding to each frame. The log-Mel filter bank (FBK) acoustic features were used throughout the paper, which were extracted with a 25 ms duration and 10 ms frame shift. The first order differentials of the FBK features were also included to extend the acoustic feature vectors. Each of these windows was transformed into a 40-dimensional FBK feature vector representing a speech frame with an offset of 10 ms. When being fed into the DNN input layer, the 40-dimensional feature vectors were augmented with their first-order time derivatives (also termed as *delta features* in the speech-recognition literature) to form an 80-dimensional vector o_t for the t -th frame. The final DNN input feature vector, x_t , was formed by stacking nine consecutive acoustic vectors around t , i.e., $x_t = \{o_{t-4}, o_{t-3}, \dots, o_{t+4}\}$. Therefore, the DNN input layer (denoted as the FBK layer from Figure 2 to Figure 1) has 720 nodes and covers a 125 ms long input window starting at $(10 \times t - 50)$ ms and ending at $(10 \times t + 75)$ ms. Where this wider context window extended beyond the limits of the recording (i.e., at the beginning and end of the recording), boundary frames were duplicated to make up the nine consecutive frames.

Following the input layer FBK, there are five 1,000-node hidden layers (L2–L6), a 26-node “bottleneck” layer (L7), and the output layer (TRI). This network is therefore denoted as DNN-BN₇ since the bottleneck layer is the seventh layer (L7). All hidden nodes use a sigmoid activation function and the output layer uses a softmax activation function to estimate pseudo posterior probabilities for 6,027 output units. There are 6,026 such units corresponding to the tied triphone HMM states which are obtained by the decision tree clustering algorithm (Young et al., 1994). The last output unit is relevant to the non-speech HMM states. The DNN was trained on a corpus consisting of 200 h of British English speech selected from 7

weeks of TV broadcast shows by the BBC covering all genres. Using such a training set with a reasonably large amount of realistic speech samples guarantees our DNN model to be properly trained and close to the models used in real-world speech recognition applications. The DNN model was trained to classify each of the speech frames in the training set into one of the output units based on the cross-entropy loss function. All DNN-BN models were trained with the same configuration. The training was conducted using a modified NewBob learning rate scheduler (Zhang and Woodland, 2015a), with each minibatch having 800 frames, and with an initial learning rate of 2.0×10^{-3} and a momentum factor of 0.5. A layer-by-layer pre-training approach was adopted, which started by training a shallow artificial neural network with only one hidden layer for one epoch, and gradually adding in more hidden layers as the penultimate layer, one layer per epoch until the final DNN structure is achieved (Hinton et al., 2012). Afterwards the entire DNN model is jointly fine-tuned for 20 epochs. More details about the training configuration and data processing procedure can be found in (Woodland et al., 2015; Zhang and Woodland, 2015a).

When performing speech recognition at test-time, the posterior probabilities, $P(s_k | x_t)$, were converted to log-likelihoods to use as the observation density probabilities of the triphone HMM states. Specifically, the conversion was performed by

$$\ln p(x_t | s_k) = \ln P(s_k | x_t) + \ln p(x_t) - \ln P(s_k), \quad (1)$$

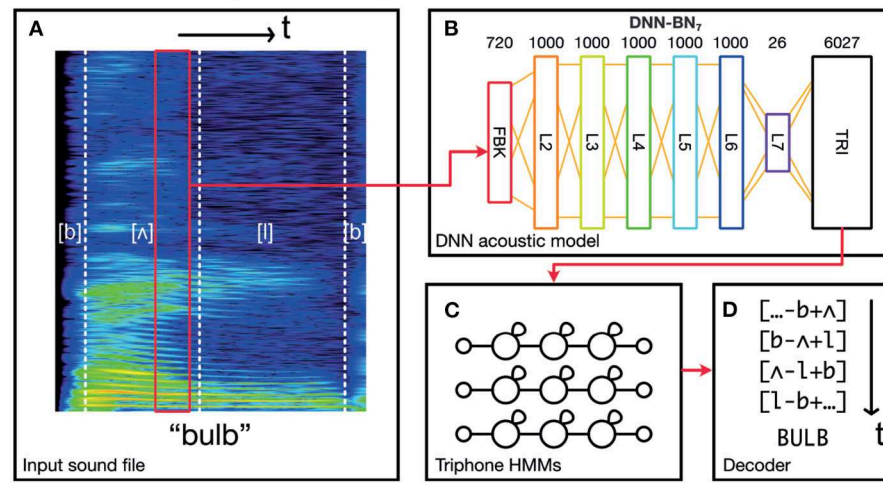
where s_k is a DNN output for target k , and $P(s_k)$ is the frequency of frames corresponding to the units associated with target k in the frame-to-HMM-state alignments of the training set (Hinton et al., 2012).

2.1.2. Recorded speech stimuli

This study used speech stimulus recordings from Fonteneau et al. (2014), which consists of 400 English words spoken by a native British English female speaker. The set of words consists of nouns and verbs (e.g., *talk*, *claim*), some of which were past-tense inflected (e.g., *arrived*, *jumped*). We assume that the words' linguistic properties are independent of the acoustic-phonetic properties presently under investigation. We also assume that this sample of recorded speech provides a reasonable representation of naturally occurring phonetic variants of British English, with the caveat that the sampled utterances are restricted to isolated words and a single speaker.

Audio stimuli, which were originally recorded and presented to subjects with a 22.1 kHz sampling rate, were down-sampled to 16 kHz before building models, as the DNN was trained on a 16 kHz audio training set. After the DNN was first trained on the data from BBC TV programs, it was further adapted to fit the characteristics of the speaker and the recording channel

HTK automatic speech recognizer



RSA Procedure

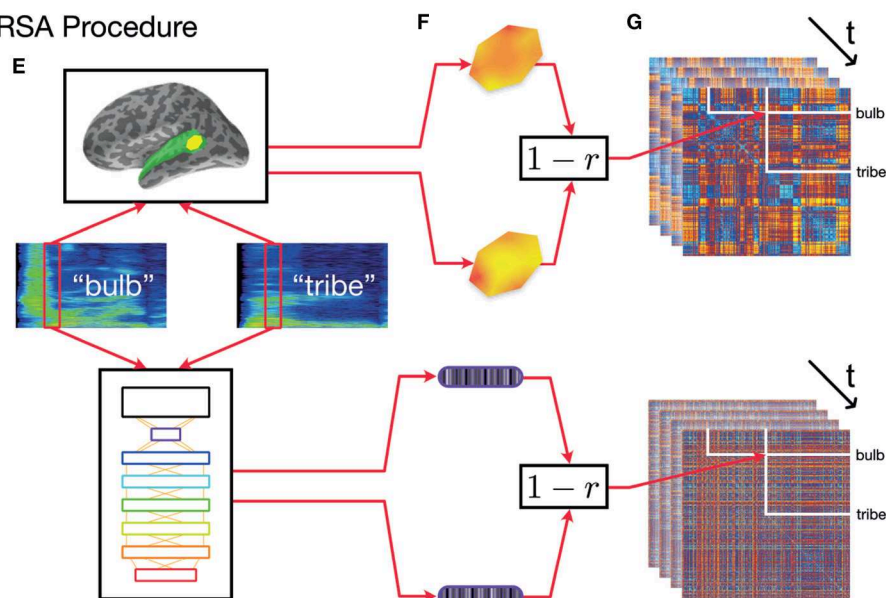
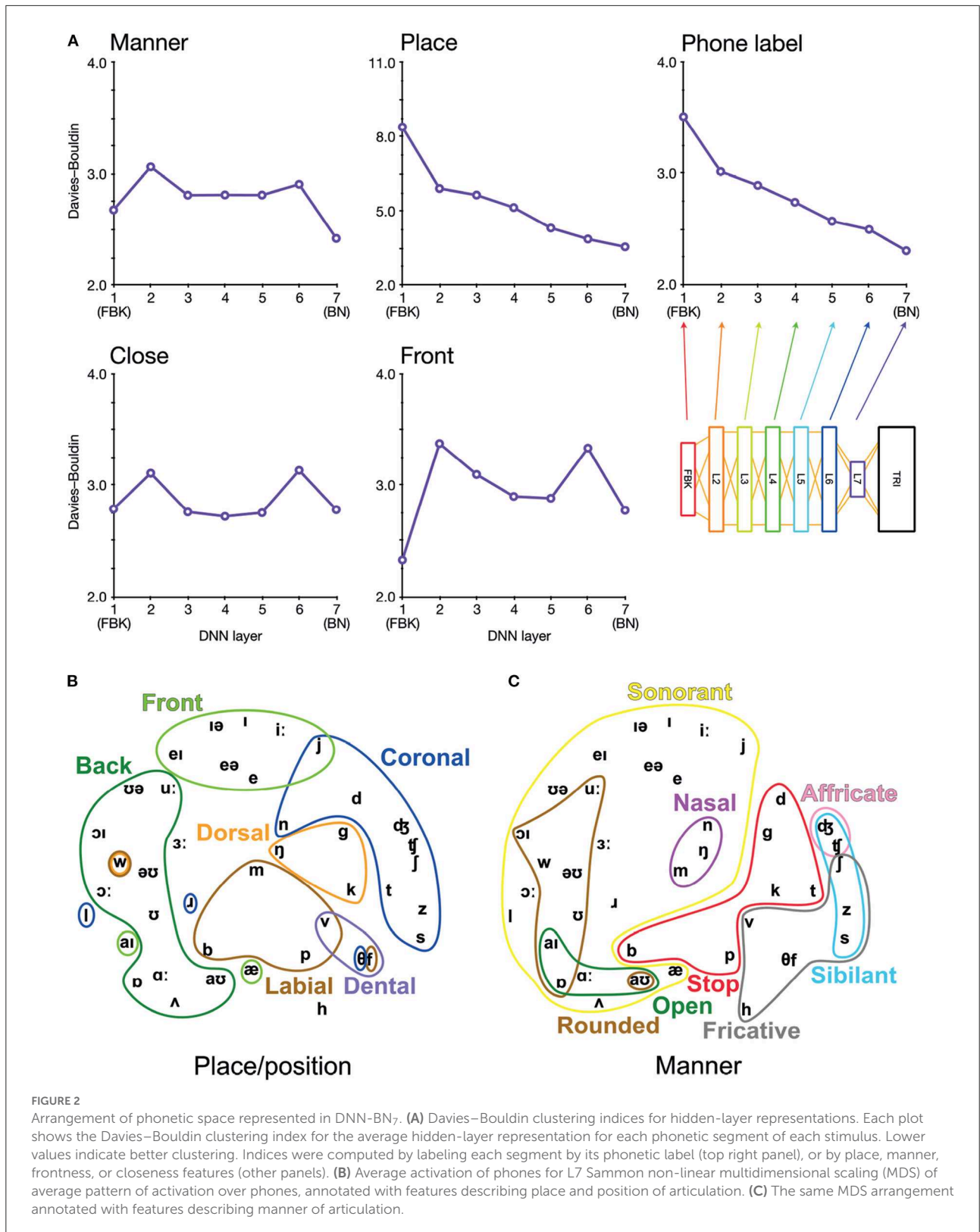


FIGURE 1

Schematic of the overall procedure. (A–D) Schematic representation of our automatic speech recognition system. Our ASR model is a hybrid DNN–HMM system built with HTK (Young et al., 2015; Zhang and Woodland, 2015a). (A) An acoustic vector is built from a window of recorded speech. (B) This is used as an input for a DNN acoustic model which estimates posterior probabilities of triphonic units. Numbers above the figure indicate the size of each layer. Hidden layer L7 is the bottleneck layer for DNN-BN₇. (C) The triphone posteriors (TRI) are converted into log likelihoods, and used in a set of phonetic HMMs. (D) A decoder computes word identities from the HMM states. (E–G) Computing dynamic RDMs. (E) A pair of stimuli is presented to each subject, and the subjects’ brain responses are recorded over time. The same stimuli are processed using HTK, and the hidden-layer activations recorded over time. (F) The spatiotemporal response pattern within a patch of each subject’s cortex is compared using correlation distance. The same comparison is made between hidden-layer activation vectors. (G) This is repeated for each pair of stimuli, and distances entered into a pairwise comparison matrix called a representational dissimilarity matrix (RDM). As both brain response and DNN response evolve over time, additional frames of the dynamic RDM are computed.

of the stimuli data using an extra adaptation stage with 976 isolated words (see Zhang and Woodland, 2015b for details of the approach). This is to avoid any potential bias to our experimental results caused by the differences between the DNN model training set and the stimuli set, without requiring the

collection of a large amount of speech samples in the same setting as the stimuli set to build a DNN model from scratch. There are no overlapping speech samples (words) between the adaptation and stimuli sets. This guarantees that the model RDM obtained using our stimuli set is not over-fitted into the



seen data, and guarantees our results and conclusions to be as general as possible.

2.1.3. Evaluating clustered representations

To investigate how the assignment of phonetic and featural labels to each segment of the stimuli could explain hidden-layer representations in DNN-BN₇, we computed Davies–Bouldin clustering indices for representational spaces at each layer.

Davies–Bouldin indices (Davies and Bouldin, 1979) are defined as the average ratio of within- and between-cluster distances for each cluster with its closest neighboring cluster. They indicate the suitability of category label assignment to clusters in high-dimensional data, with lower values indicating better suitability and with 0 the minimum possible value (obtained only if labels are shared only between identical points). This in turn serves as an indication of how suitably phonetic and feature labels might be assigned to hidden-layer representations. To compute Davies–Bouldin indices, we recorded the vector of hidden-layer activations elicited by each input time window of the stimuli for each layer in each DNN. There was a high level of correlation between many activation vectors resulting from overlapping adjacent input vectors. To minimize the effect of this, we used average vectors from each hidden layer over each contiguous phonetic segment. For example, in the word “bulb”, the hidden-layer representations associated with each frame corresponding to the acoustic implementation of the first [b] were combined, and separately the representations for the final [b] were combined. Then, to each combined vector, we assigned a label under five separate labeling schemes: closeness features, frontness features, place features, manner features, and phonetic label. For place and manner features, we considered only phones which exhibited a place or manner feature (i.e., obstruents). For frontness and closeness features, we likewise considered only phones which exhibited frontness or closeness features (i.e., syllabic vowels). Where a phone had more than one appropriate feature assignment, we used the most appropriate feature. The full assignment of feature labels for phones used in the clustering analysis is given in [Supplementary Figure 1](#).

We computed *p*-values for each Davies–Bouldin index calculation using a permutation procedure in which phone labels were randomized after averaging activation vectors for each segment of input (5,000 permutations). *p*-values were computed by randomizing the labels and recomputing Davies–Bouldin indices 5,000 times, building a distribution of Davies–Bouldin indices under the null hypothesis that phone and feature labels did not systematically explain differences in hidden-layer activations. In all cases, the observed Davies–Bouldin index was lower than the minimum value in the null distribution, yielding an estimated *p*-value of exactly 0.0002. Since the precision of this value is limited by the number of permutations performed, we report it as $p < 0.001$. All Davies–Bouldin index values reported were significant at the $p < 0.001$ level.

2.2. Results and discussion

Davies–Bouldin indices for each layer and categorization scheme are shown in [Figure 2A](#). Of particular interest is the improvement of feature-based clustering in bottleneck layer L7 of DNN-BN₇, which shows that it is, in some sense, reconstructing the featural *articulatory* dimensions of the speaker. That is, though this was not included in the teaching signal, when forced to parsimoniously pass comprehension-relevant information through the bottleneck, DNN-BN₇ finds a representation of the input space which maps well onto the constraints on speech sounds inherent in the mechanics of the speaker. L7 showed the best clustering indices out of all layers for manner and place features and phone labels, and the second-best for frontness features. For closeness alone, L7 was not the best, but was still better than its adjacent layer L6. The general trend was that clustering improved for successively higher layers. Layers prior to the bottleneck tended to have larger clustering indices, indicating that their activations were not as well accounted for by phonetic or featural descriptions.

To further illustrate and visualize the representational space for L7, we used the phonetic partitioning of our stimuli provided by HTK, and averaged the activation across hidden nodes in L7 for each window of our 400 stimulus words which was eventually labeled with each phone. This gave us an average L7 response vector for each phone. We visualized this response space using the Sammon non-linear multidimensional scaling (MDS) technique in which true high-dimensional distances between points are compressed into two dimensions so as to minimize distortion (Sammon, 1969). Place/position features are highlighted in [Figure 2B](#), and manner features are highlighted in [Figure 2C](#).

To be clear, the presence of these feature clusters does not imply that there are individual nodes in L7 which track specific articulatory features. However, using the reasoning of RSA, we can see that articulatory features are descriptive of the overall arrangement of phones in the L7 response space. This ability to characterize and model an overall pattern ensemble in a way abstracted from the specific response format and distributed neural representations is one of the strengths of the RSA technique.

3. Study 2: Representational similarity mapping of auditory cortex with DNN representations

3.1. Materials and methods

3.1.1. Computing model RDMs from incremental machine states

To encapsulate the representational space of each of the DNN’s hidden layer representations through time, we

computed model RDMs from the activation of each layer using the following procedure, illustrated in [Figure 1](#). RSA computations were performed in Matlab using the RSA toolbox ([Nili et al., 2014](#)).

As described previously, the input layer of the DNN had access to 125 ms of audio input at each time step, to estimate the triphone-HMM-state likelihoods. Since we can only compute model RDMs where the DNN has activations for every word in the stimuli set, only the activations corresponding to the frames whose ending time is smaller than 285 ms (the duration of the shortest word) are used in our experiments. Since each frame has a 25 ms duration and a 10 ms shift, only the activations of the first 27 frames of each word are reserved to construct our model RDMs (as the frame index t is required to satisfy $10 \times t + 25 \leq 285$).

For each fixed position of the sliding time window on each pair of our 400 stimulus words, we obtained the pattern of activation over the nodes in a particular layer of the DNN. By computing Pearson's correlation distance ($1 - r$) between activation pattern for each pair of words, we built a 400×400 model RDM whose rows and columns were indexed by the stimulus words. Then, by moving the sliding time window in 10 ms increments and recomputing model RDM frames in this way, we produced a series of model RDMs which varied throughout the first 260 ms of the stimuli. We repeated this procedure for each hidden layer L2–L7, as well as the input and output layers FBK and TRI, producing in total eight series of model RDMs, or 216 individual model RDM frames. When building a model RDM frame from the input layer FBK, we used only the 40 log-mel filterbank values within the central 25 ms window (and did not include the first derivatives or overlapping context windows).

3.1.2. EMEG data collection

Sixteen right-handed native speakers of British English (six male, aged 19–35 years, self-reported normal hearing) participated in the study. For each participant, recordings of 400 English words, as spoken by a female native British English speaker were presented binaurally. Each word was repeated once. The study was approved by the Peterborough and Fenland Ethical Committee (UK). Continuous MEG data were recorded using a 306 channels VectorView system (Elektra-Neuromag, Helsinki, Finland). EEG was recorded simultaneously from 70 Ag-AgCl electrodes placed within an elastic cap (EASYCAP GmbH, Herrsching-Breitbrunn, Germany) according to the extended 10/20 system and using a nose electrode as the recording reference. All data [Fonteneau et al. \(2014\)](#).

3.1.3. EMEG source estimation

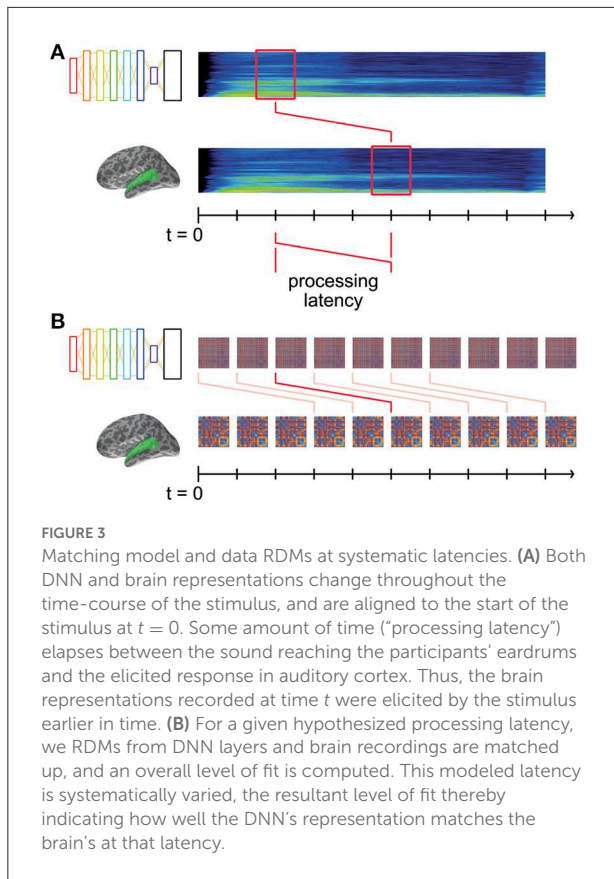
In order to track the cortical locations of brain–model correspondence, we estimated the location of cortical sources using the anatomically constrained MNE ([Hämäläinen and Ilmoniemi, 1994](#)) with identical parameters to those used in our previous work ([Fonteneau et al., 2014](#); [Su et al., 2014](#); [Wingfield et al., 2017](#)). MR structural images for each participant were obtained using a GRAPPA 3D MPRAGE sequence (TR = 2250 ms; TE = 2.99 ms; flip-angle = 9 deg; acceleration factor = 2) on a 3 T Trio (Siemens, Erlangen, Germany) with 1 mm isotropic voxels. From the MRI data, a representation of each participant's cerebral cortex was constructed using FreeSurfer software (<https://surfer.nmr.mgh.harvard.edu/>). The forward model was calculated with a three-layer boundary element model using the outer surface of the scalp as well as the outer and inner surfaces of the skull identified in the anatomical MRI. This combination of MRI, MEG, and EEG data provides better source localization than MEG or EEG alone ([Molins et al., 2008](#)).

The constructed cortical surface was decimated to yield approximately 12,000 vertices that were used as the locations of the dipoles. This was further restricted to the bilateral superior temporal mask as discussed previously. After applying the bilateral region of interest mask, 661 vertices remained in the left hemisphere and 613 in the right. To perform group analysis, the cortical surfaces of individual subjects were inflated and aligned using a spherical morphing technique implemented by MNE ([Gramfort et al., 2014](#)). Sensitivity to neural sources was improved by calculating a noise covariance matrix based on the 100 ms pre-stimulus period. The activations at each location of the cortical surface were estimated over 1 ms windows.

This source-reconstructed representation of the electrophysiological activity of the brain as the listeners heard the target set of 400 words was used to compute brain RDMs.

3.1.4. Computing brain RDMs in a spatiotemporal searchlight

To match the similarity structures computed from each layer of the DNN to those found in human participants, in the ssRSA procedure, RDMs were calculated from the EMEG data contained within a regular spatial searchlight patch and fixed-width sliding temporal window. We used a patch of vertices of radius 20 mm, and a 25 ms sliding window to match the 25 ms frames used in ASR. The searchlight patch was moved to center on each vertex in the masked source mesh, while the sliding window is moved throughout the epoch in fixed time-steps of 10 ms. From within each searchlight patch, we extracted the spatiotemporal response pattern from each subject's EMEG data. We computed word-by-word RDMs using Pearson's correlation distance ($1 - r$) on the resulting response vectors. These RDMs were averaged across subjects, resulting in one brain RDM for



each within-mask vertex. Our 25 ms ssRSA sliding window moved in increments of 10 ms throughout an EMEG epoch of [0, 540] ms, giving us a series of RDMs at each vertex for sliding windows $[t, t + 25]$ ms for each value of $t = 0, 10, \dots, 510$. In total, this resulted in a total of 66,300 brain RDM frames. By using the ssRSA framework, we make this vast number of comparisons tractable by systematizing the comparison.

3.1.5. Systematic brain–model RDM comparisons

The model RDMs computed from the DNN layer activations describe the changing representational dissimilarity space of each layer throughout the duration of the stimulus words. We can think of this as a dynamic model timeline for each layer; a collection of RDMs indexed by time throughout the stimulus. Similarly, the brain data-derived RDMs computed from brain recordings describe the changing representational dissimilarity space of the brain responses at each searchlight location throughout the epoch, which we can think of as a dynamic data timeline. It takes non-zero time for vibrations at the eardrum to elicit responses in auditory cortex (Figure 3A). Therefore, it does not make sense to only compare the DNN RDM from a given time window to the precisely corresponding

brain RDM for the same window of stimulus: to do so would be to hypothesize instantaneous auditory processing in auditory nerves and in the brain.

Instead, we offset the brain RDM's timeline by a fixed latency, k ms (Figure 3B). Then, matching corresponding DNN and brain RDMs at latency k tests the hypothesis that the DNN's representations explain those in auditory cortex k ms later. By systematically varying k , we are able to find the time at which the brain's representations are best explained by those in the DNN layers.

Thus, for each such potential processing latency, we obtain a spatial map describing the degree to which a DNN layer explains the brain's representations at that latency (i.e., mean Spearman's rank correlation coefficient between DNN and brain RDMs at that latency). Varying the latency then adds a temporal dimension to the maps of fit.

This process is repeated for each subject, and data combined by a t -test of the ρ values across subjects at each vertex within the mask and each latency. This resulted in one spatiotemporal t -map for each layer of the DNN. For this analysis, we used latencies ranging from 0 to 250 ms, in 10 ms increments.

3.1.6. Threshold-free cluster enhancement

We applied threshold-free cluster² enhancement (TFCE: Smith and Nichols, 2009) to the t -maps from each layer of the DNN. TFCE is an image-enhancement technique which enables the use of cluster-sensitive statistical methods without the requirement to make an arbitrary choice of initial cluster-forming threshold and is used as the standard statistical method by the FSL software package (Jenkinson et al., 2012).

TFCE transforms a statistical image in such a way that the value at each point becomes a weighted sum of local supporting clustered signal. Importantly, the shape of isocontours, and hence locations of local maxima, are unchanged by the TFCE transformation. For a t -map comprised of values $t_{v,k}$ for vertices v and latencies k , the TFCE transformation is given by

$$\text{TFCE}(t_{v,k}) = \int_0^{t_{v,k}} h^2 \sqrt{e(h)} dh \quad (2)$$

where $e(h)$ is the cluster extent of the connected component of (v, k) at threshold h . We approximated (2) with the sum

$$\sum_{i\Delta h \leq t_{v,k} < (i+1)\Delta h} (i\Delta h)^2 \sqrt{e(i\Delta h)} \quad (3)$$

² The term *cluster* here refers to spatiotemporally contiguous sets of datapoints in statistical maps of activation or model fit. This is a different term to *cluster* as used in the previous section to refer to sets of points located close-by in a high-dimensional abstract space. It is unfortunate that both of these concepts have the same name, but we hope their distinct meanings will be clear from the context.

where Δh was set to 0.1. The choice of Δh affects the accuracy of the approximation (3) but should not substantially bias the results.

All t -maps presented for the remainder of this paper have TFCE applied.

3.1.7. Group statistics and correction for multiple comparisons

To assess the statistical significance of the t -maps, we converted the t -values to p -values using a random-effects randomization method over subjects, under which p -values are corrected for multiple spatiotemporal comparisons (Nichols and Holmes, 2002; Smith and Nichols, 2009; Su et al., 2012). In the random-effects test, a null-distribution of t -values is simulated under the null hypothesis that Spearman's rank correlation values ρ are symmetrically distributed about 0 (i.e., no effect). By randomly flipping the sign of each individual subject's ρ -maps before computing the t -tests across subjects and applying the TFCE transformation, we simulate t -maps under the null hypothesis that experimental conditions are not differentially represented in EMEG responses. From each such simulated map, we record the map-maximum t -value, and collect these into a null distribution over all permutations. For this analysis we repeated the randomization 1,000 times, and collected separate null distributions for each hemisphere. To assess the statistical significance of a true t -value, we see in which quantile it lies in the simulated null distribution of map-maximum randomization t -values.

We performed this procedure separately for the models derived from each layer of the DNN, allowing us to obtain t -maps which could be easily thresholded at a fixed, corrected p -value.

3.2. Results

We used the dynamic representations from each layer of DNN-BN₇ to model spatiotemporal representations in the auditory cortices of human participants in an EMEG study by applying ssRSA. Areas of auditory cortex (Figure 4A) were defined using the Desikan–Killiany Atlas (STC and HG).

Figure 4 shows the left hemisphere results of this analysis. The brain maps in Figure 4B show threshold-free-cluster-enhanced t -maps (Smith and Nichols, 2009) computed from the model RDMs of each hidden layer, thresholded at $p < 0.01$. Model RDMs computed from all DNN layers except L5 showed significant fit in left STC and HG. Input layer FBK peaked early in left posterior STC at 0–70 ms, and later in left anterior STC and HG at 140–210 ms. Hidden-layer models L2–L4 and L6–L7 peaked later than FBK, achieving

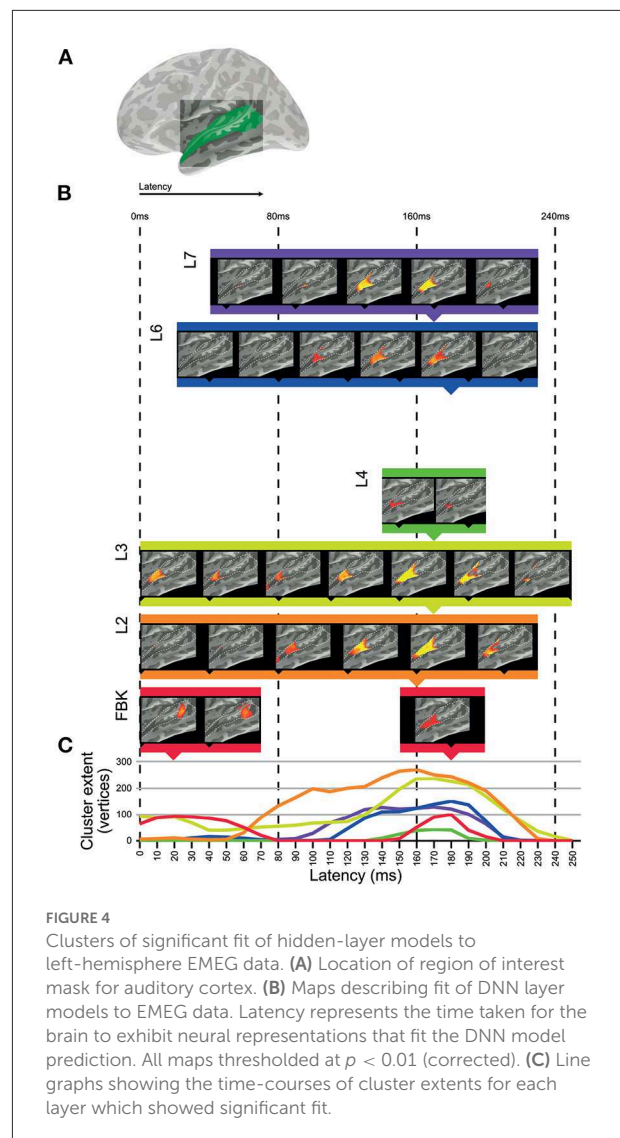


FIGURE 4
Clusters of significant fit of hidden-layer models to left-hemisphere EMEG data. (A) Location of region of interest mask for auditory cortex. (B) Maps describing fit of DNN layer models to EMEG data. Latency represents the time taken for the brain to exhibit neural representations that fit the DNN model prediction. All maps thresholded at $p < 0.01$ (corrected). (C) Line graphs showing the time-courses of cluster extents for each layer which showed significant fit.

maximum cluster size at approximately 170 ms. Layers L5 and TRI showed no significant fit in the regions of interest. Overall, significant cluster size increased between layers FBK–L3, diminished for L4 and L5, and re-emerged for L6 and L7.

The line graphs in Figure 4C show the time-courses of each layer as they attain their maximum cluster extent. In general, there appeared to be two distinct clusters across the superior temporal region: an early cluster peaked in left posterior STC for the DNN input layer FBK, and another late cluster peaked in left anterior STC for DNN layers L1–L4 and L6–L7, throughout the whole epoch, but attaining a maximum cluster size at approx 170 ms. Details of timings for each layer are shown in Supplementary Table 1. Right hemisphere results are included in Supplementary Figure 2.

3.3. Discussion

The input layer FBK representing purely acoustic information (i.e., not a learned or task-relevant representation) showed a later and smaller effect (cluster in human posterior STC) than that of higher layers L2 and L3. The strongest cluster for FBK was early, and the later cluster appears to be a weaker version of those for higher hidden-layer models. The late cluster for FBK indicates that there is some involvement of both low-level acoustic features and higher-level phonetic information in the later neural processes at around 170 ms. However, since there is an intrinsic correlation between acoustic information and phonetic information, it is hard to completely dissociate them. Another explanation for the mixture of high and low levels of speech representations in a single brain region at the same time is the existence of feedback connections in human perceptual systems (However, the ASR systems used in this paper can achieve high degree of accuracy without the top-down feedback loop from higher to lower hidden layers.). It should be noted that while the FBK, L2 and L4 clusters all register as significant at a latency of 0 ms, timings correspond to a 25 ms window of EMEG data being matched against model state computed for the central 25 ms of 125 ms windows of audio, so only approximates the actual latency.

Moving up to hidden layers L2 and L3, we saw later clusters which fit the brain data more strongly than FBK in the left hemisphere. All hidden layers including L2 and L3 activate according to learned parameters. Progressively higher layers L4 and L5 fit with smaller clusters in human STC, with L5 showing no significant vertices at any time point ($p > 0.01$) in the left hemisphere but a very small cluster in the right hemisphere. However, the highest hidden layers L6 and L7 once again showed string fit with activations in left anterior STC.

Of particular interest is this re-emergence of fit in anterior STC to the representations in the bottleneck layer L7. In this layer of the DNN, the 1,000-node representation of L6 is substantially constrained by the reduced size of the 26-node L7. In particular, the fact that ASR accuracy is not greatly reduced by the inclusion of this bottleneck layer indicates that, for the machine solution, 26 nodes provide sufficient degrees of freedom to describe a phonetic space for purposes of word recognition. This, in conjunction with the re-emergence of fit for L7 to STC representations makes the representations of this layer of particular interest. The hidden layers in the DNN learn to sequentially transform acoustic information into phonetic probabilities in a way which generalizes across speakers and background acoustic conditions. There is no guarantee that the features the DNN learns to identify for recognition are comparable to those learned by the brain, so the fact that significant matches in the RDMS were found between machine and human solutions of the same problem is worthy of further consideration.

4. Study 3: Improving DNN design

4.1. Materials and methods

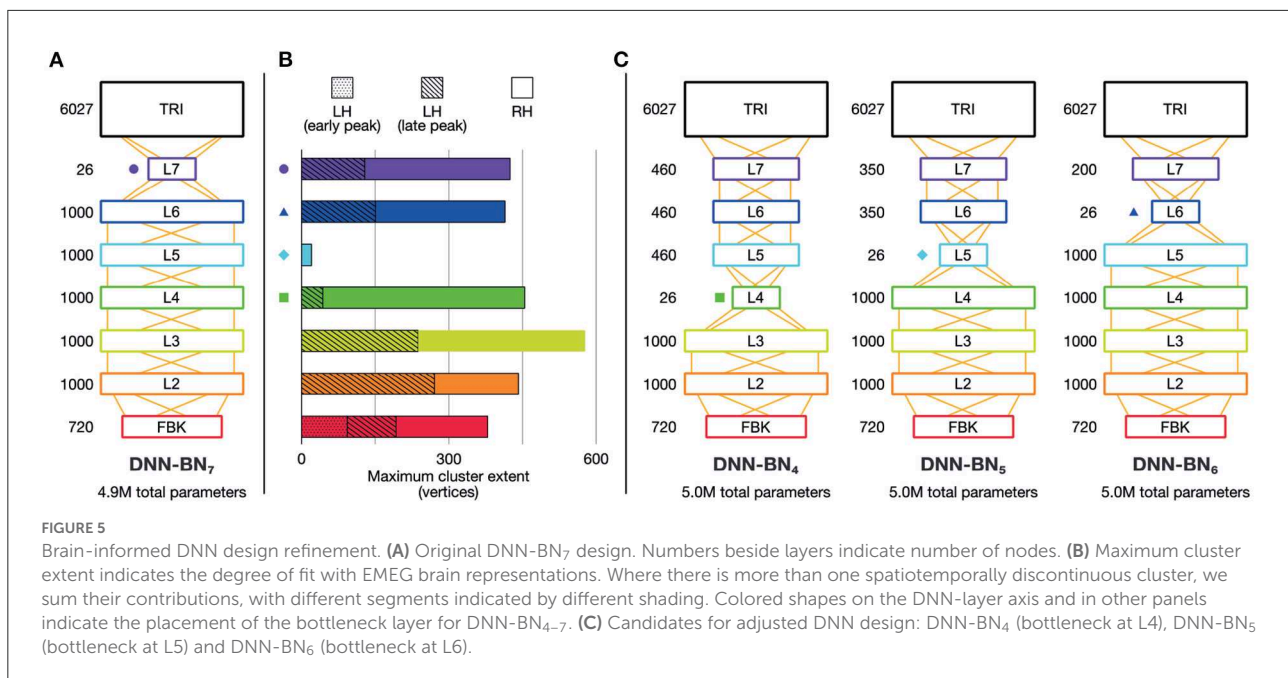
From the maximum cluster extents of the DNN layers shown in [Figure 4](#), the activations of the DNN acoustic model significantly correspond to the activity in the left-hemisphere of human brain when listening to the same speech samples. This suggests that the DNN and human brain rely on similar mechanisms and internal representations for speech recognition.

Human speech recognition still has superior performance and robustness in comparison to even the most advanced ASR systems, so we reasoned that it could be possible to improve the DNN model structure based on the evidence recorded from the brain.

The overall minimal spatiotemporal clusters for L5 of DNN-BN₇ suggested that while early layers (L2–L3) were performing analogous transformations to early auditory cortex, and that the bottleneck (L7) was representing speech audio with a similarly parsimonious basis as left auditory cortex, there was a divergence of representation at intermediate layers (L4–L6). This indicates the possibility that the calculations in DNN layer L5 are less important for recognizing the speech accurately since brain does not appear to use such representations in the recognition process. On the other hand, although a bottleneck layer is positioned at L7, its strong correspondence to the brain reveals the importance of the calculations performed in that layer. Thus, it is natural to assume that more parameters and calculations in important layers can improve speech recognition performance, while fewer calculations can reduce the complexity of the model DNN structure without sacrificing the performance too much. With the supposition that the arrangement of auditory cortex would be adapted specifically to speech processing, we hypothesized that by moving the bottleneck layer into the positions occupied by divergent layers in DNN-BN₇, the network might learn representations that closer resemble those of human cortex, and thus improve the performance of the model.

To this end, we built and studied another DNN model, DNN-BN₅, which has the same number of parameters as DNN-BN₇ but has the bottleneck layer moved from L7 to L5. The details of the new DNN structures are shown in [Figure 5C](#). For purposes of comparison, and following the same naming convention, we expanded our investigation with another two DNN models, DNN-BN₄ and DNN-BN₆ were also built for DNNs whose bottleneck layers are L4 and L6 respectively. In all models the number of parameters was kept to 5.0 million, matching the 4.9 million parameters of DNN-BN₇.

It may appear to the reader as if an alternative modification would be to re-locate the bottleneck layer relative to the input layer as we have done so, but attach it directly to the TRI layer (as in DNN-BN₇) without intermediate levels. However



we chose to fix the number of DNN layers and simply move the position of the bottleneck in order to keep the total number of parameters fixed at 5 million, since number of trainable parameters is a strong determiner of performance ceiling. We could have retained 5 million parameters by inflating the size of the hidden layers between the input and the bottleneck, but this would have forced upstream representations to change between models, making DNN-BN₇ harder to compare to DNN-BN₄₋₆. Additionally, early DNN studies demonstrated that, for a fixed number of parameters, deeper, thinner models (i.e., those with more layers containing fewer units) performed significantly better than shallower, wider models, and this is now a standard practice (Morgan, 2011; Hinton et al., 2012). Alternative DNN design choices may have different effects, and we hope to investigate this in future work.

We tested the derived DNN models with different bottleneck layer positions using two tasks: general large-vocabulary continuous speech recognition with recordings from BBC TV programs, and in-domain isolated-word recognition using the stimuli set. The MGB Dev set was derived as a subset of the official development set of the MGB speech recognition challenge (Bell et al., 2015), which includes 5.5 h of speech. Since the MGB testing set involves sufficient samples (8,713 utterances and 1.98M frames) from 285 speakers and 12 shows with diversified genres, and the related WER results are reliable metrics to evaluate the general performance of the DNN models for speech recognition. In contrast, the WERs on the stimuli set are much more noisier since it only consists of 400 isolated words from a single female speaker. However, the stimuli set WERs are still important metrics since the same 400 words are

TABLE 1 The performance of DNN-HMM systems with different bottleneck layer positions.

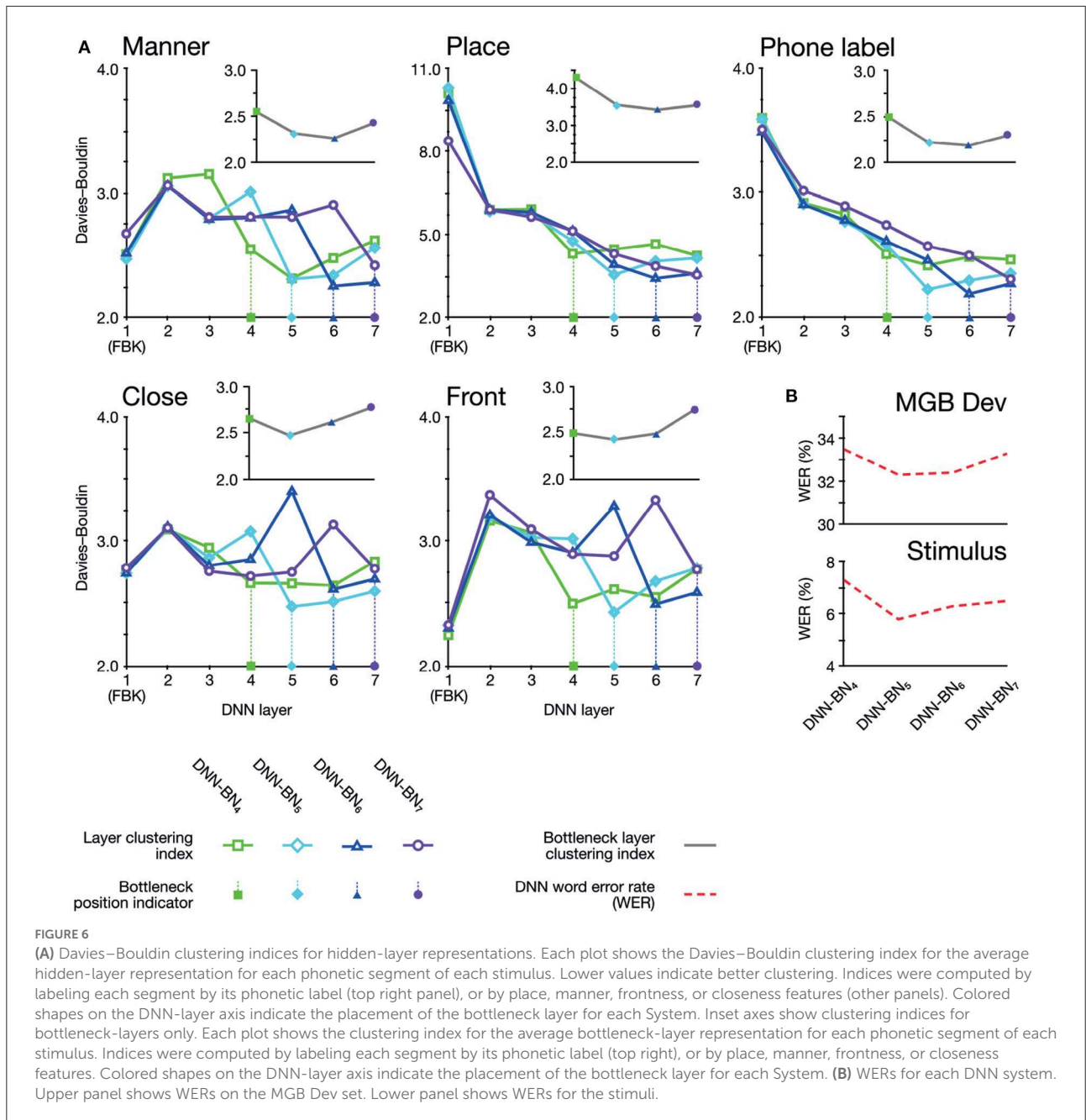
System	Bottleneck layer	Accuracy%		WER%	
		Train	HV	MGB Dev	Stimuli
DNN-BN ₇	L7	44.0	41.5	33.3	6.5
DNN-BN ₆	L6	44.6	42.3	32.4	6.3
DNN-BN ₅	L5	44.2	42.3	32.3	5.8
DNN-BN ₄	L4	42.6	41.1	33.5	7.3

The WERs (the lower the better) were given on both the MGB challenge official development subset (MGB Dev), which is a general purpose large vocabulary continuous speech recognition testing set, as well as the 400 isolated words used as the stimuli in our listening experiments to derive the RDM (Stimuli). The MGB Dev WERs are reliable indicators for the general performance of the systems in realistic ASR tasks. The Stimuli WERs are the most direct indicators of the model performance on the data used in our brain-machine comparison experiments. The classification accuracy values (the higher the better) were obtained by classifying each frame into one of the 6,027 triphonic DNN output units were obtained on both the training and held-out validation (HV) sets. For fair comparisons, DNN structures of all systems were constrained to have the same amount of model parameters (about 5M for each model, as shown in Figure 5). Accuracy can be considered as an auxiliary performance metric, which indicates that DNN-BN₆ suffered more from over-fitting compared to DNN-BN₅, since DNN-BN₆ is better in the training accuracy but not in the HV accuracy.

used to build the RDMs used in the key experiments. These results are presented in Table 1 and Figure 6C.

4.2. Results

As shown in Table 1 and Figure 6C, adjusting the design of the DNN structure to better fit with the representations



exhibited in the human subjects led to improved DNN performance in terms of WER in DNN-BN₅ and DNN-BN₆. The MGB Dev set contains sufficient testing samples with diversified speaker and genre variations. When testing on MGB Dev, a 4-g language model (Woodland et al., 2015) was used to provide word-level contexts by rescoring each hypothesis in decoding as in general large vocabulary continuous speech recognition applications. The 1.0% absolute WER reduction (relatively 3.3%) obtained by comparing DNN-BN₇ with DNN-BN₅ is substantial (Bell et al., 2015; Woodland et al., 2015).

Such an improvement was achieved without increasing the number of parameters, and hence demonstrates the superiority of the structure of DNN-BN₅. DNN-BN₆ also performed 0.9% (absolute WER) better than DNN-BN₇, but 0.1% WER worse than DNN-BN₅. This can also be observed from the frame classification accuracy values, as DNN-BN₆ has the same HV accuracy but better train accuracy compared with DNN-BN₅, indicating that placing the bottleneck layer at L6 results in overfitting. Regarding the stimulus set, no language model was used since each stimulus utterance has only one word and the

recognition requires no word-level context. Still, the changes in WERs are consistent with those on the MGB Dev set. Comparing Table 1 to Figure 5, the WERs and the maximum cluster extent values of these DNN models are also consistent on the Stimulus test set.

As well as altering the position of the bottleneck layer, we also trained and tested a DNN without a bottleneck layer, but using the same 5.0M parameters. This DNN achieved 44.0% train accuracy and 42.3% HV accuracy, and 32.3% MGB Dev WER and 5.8% WER on the stimuli. In other words, close-to, but just falling short of (albeit insignificantly), the overall best model including a bottleneck layer: DNN-BN₅. The inclusion of a bottleneck layer was included in DNN-BN₅ was motivated both for machine-learning and computational-modeling reasons, as we have described. It is therefore notable that even though DNN-BN₅ contains a bottleneck layer, and thus forces a compression of the speech representation from 1,000 down to 26 dimensions, it was still able to achieve the overall best performance.

What is not immediately clear, however, is whether this improvement in performance arises from a corresponding improvement in the model's ability to extract a feature-based representation. In other words, if the bottleneck layer learns a representation akin to articulatory features, by moving the layer to improve performance does this enhance this learned representation? To answer this question, we investigated how the assignment of phonetic and featural labels to each segment of the stimuli could explain their hidden-layer representations. As before, we probed the organization of the representational space of each hidden layer according to phones and features using Davies–Bouldin clustering indices.

The clustering results exhibited two overall patterns of note. First, clustering (i.e., suitability of assignment of phonetic and featural labels to hidden layer representations) was improved on the DNNs whose design had been inspired by the human brains. Second, the optimum clustering level was often found in the bottleneck layer itself (highlighted on the graphs in Figure 6A). The clustering index at the bottleneck layers alone are separately graphed in inset panels in Figure 6A, and show that bottleneck layer clustering was also improved in DNN-BN₅ and DNN-BN₆.

In other words, the placement of the bottleneck layer in position 5 and 6 yielded, as predicted, the best clustering results both overall and in the bottleneck layer itself. Moving the bottleneck layer too far back (DNN-BN₄) yielded worse clustering results generally and in the bottleneck layer—indicated by the characteristic U-shaped curves in Figure 6B.

4.3. Discussion

Artificial Intelligence (AI) and machine learning have already been extensively applied in neuroscience primarily in analyzing and decoding large and complex neuroimaging

or cell recording data sets. Here, DNN-based ASR systems were used as a model for developing and testing hypothesis and neuroscientific theories about how human brains perform speech recognition. This type of mechanistic or generative model—where the computational model can perform the behavioral task with realistic data (in this case, spoken word recognition)—can serve as a comprehensive framework for testing claims about neurocognitive functional organization (Kriegeskorte and Douglas, 2018). Moreover, insights can flow both ways; the neuroimaging data can also guide the exploration of the model space and lead to improvements in model performance, as we have seen.

While our use of neurological data only indirectly informed the improvements to ASR architecture, the present work can be seen as an initial step toward extracting system-level designs for neuromorphic computing from human auditory systems. This goal in itself is not new (see e.g., Toneva and Wehbe, 2019), however the key novel element of our approach is the ability to relate the machine and human solutions in complementary directions. The power of RSA, and in particular ssRSA, to relate the different forms of representations in these systems is key in this work. In summary, the methodology illustrated here paves the way for future integration of neuroscience and AI with the two fields driving each other forwards.

5. General discussion

We have used a DNN-based ASR system and spatiotemporal imaging data of human auditory cortex in a mutually informative study. In the machine-to-human direction, we have used a computational model of speech processing to examine representations of speech throughout space and time in human auditory cortex measured as source-localized EMEG data. In so doing, we have produced a functional map in human subjects for each part of the multi-stage computational model. We were able to relate dynamic states in the operating machine speech recognizer to dynamic brain states in human participants by using ssRSA, extended to account for a dynamically changing model. In a complementary analysis, we have improved the performance of the DNN-based ASR model by adapting the layered network architecture inspired by the staged neural activation patterns observed in human auditory cortex.

5.1. Relating dynamic brain and machine states: Comparing and contrasting computational models in vision and audition

There has been some recent successes in comparing machine models of perception to human neuroimaging data. This has primarily been in the domain of visual object perception (e.g.,

Kriegeskorte et al., 2008b; Cadieu et al., 2014; Clarke et al., 2014; Khaligh-Razavi and Kriegeskorte, 2014; Güçlü and van Gerven, 2015; Kriegeskorte, 2015; Cichy et al., 2016; Kheradpisheh et al., 2016; Devereux et al., 2018), with less progress made in speech perception (though see our previous work: Su et al., 2014; Wingfield et al., 2017).

The visual systems of humans and other primates are highly related, both in their architecture and in accounts of the neurocomputational processes they facilitate. There is evidence of a hierarchical organization of cortical regions in the early visual systems of human and non-human primates. There are also detailed accounts of process sequencing from early visual cortex through higher perceptual and semantic representation which exist for visual object perception in several primate models (e.g., Van Essen et al., 2001; Tootell et al., 2003; Denys et al., 2004; Orban et al., 2004; Kriegeskorte et al., 2008b). This is not so the case for speech processing and audition to the same degree.

In parallel, machine models for vision have often been designed based on theories of primate cortical processing hierarchies. This extends to recent work employing deep convolutional neural networks (CNN) for visual object processing, in particular those featuring layers of convolution and pooling. Furthermore, the convolutional layers in CNNs appear to learn features resembling those in the receptive fields of early visual cortex, and higher layers' representational spaces also match those found in higher visual cortex, and other regions in the visual object perception networks (Khaligh-Razavi and Kriegeskorte, 2014; Güçlü and van Gerven, 2015; Wen et al., 2018). Importantly, this means that the internal structures of machine vision systems are potentially informative and relevant to our understanding of the neurocomputational architecture of the natural system (and vice versa), and not just whether they generate equivalent outputs (for example in object classification tasks). To date, these common features are not well-established for DNNs or other type of acoustic models widely used for ASR systems.

Certain aspects of the human auditory processing system have resemblances to those in other primate models (Rauschecker and Scott, 2009; Baumann et al., 2013). However, no non-human primate supports anything like human speech communication, where intricately modulated sequences of speech sounds map onto hundreds of thousands of learned linguistic elements (words and morphemes), each with its own combination of acoustic-phonetic identifiers.

Perhaps due to this lack of neurocomputationally explicit models of spoken word recognition, the design of ASR systems has typically not been guided by existing biological models. Rather, by optimizing for engineering-relevant properties such as statistical learning efficiency, they have nonetheless achieved impressive accuracy and robustness.

It is striking, therefore, that we have been able to show that the regularities that successful ASR systems encode in the

mapping between speech input and word-level phonetic labeling can indeed be related to the regularities extracted by the human system. In addition, like animal visual systems have inspired the field of computer vision, we have demonstrated that human auditory cortex can improve ASR systems using ssRSA.

6. Conclusion and future work

We have shown that our deep artificial neural network model of speech processing bears resemblance to patterns of activation in the human auditory cortex using the combination of ssRSA with multimodal neuroimaging data. The results also showed that the low-dimensional bottleneck layer in the DNN could learn representations that characterize articulatory features of human speech. In ASR research, although the development of systems based around the extraction of articulatory features has a long history (e.g., Deng and Sun, 1994), except for a small number of exemplars (e.g., Zhang et al., 2011; Mitra et al., 2013), recent studies mostly rely on written-form-based word piece units (Schuster and Nakajima, 2012; Wu et al., 2016) that are not directly associated with phonetic units. Our findings imply that developing appropriate intermediate representations for articulatory features may be central to speech recognition in both human and machine solutions. In human neuroscience studies, this account is consistent with previous findings of articulatory feature representation in the human auditory cortex (Mesgarani et al., 2014; Correia et al., 2015; Wingfield et al., 2017), but awaits further investigation and exploitation in machine solutions for speech recognition. In particular, previous work by Hamilton et al. (2021) has shown that—unlike our DNN architecture—the organization of early speech areas in the brain are not purely hierarchical, suggesting new potential avenues of model architectures including layer-bypassing connections.

The results we have presented here prompt further questions regarding how modifications to the design and training of DNN-based ASR models affects their representations, how to most effectively tailor a model to match the representational organization of the human brain, and which of these modifications lead to improved performance at the task. We hope to continue similar investigations to other types of artificial neural network models in our future work, such as different hidden activation functions, time-delay neural networks (Waibel et al., 1989; Peddinti et al., 2015), CNNs (LeCun et al., 1998; Krizhevsky et al., 2012), and recurrent neural networks (Rumelhart et al., 1986; Hochreiter and Schmidhuber, 1997), etc.

There is a difference between speech recognition (i.e., the extraction of word identities from speech audio) and speech comprehension (i.e., understanding and the elicitation of meaning). In this paper we have tackled only recognition. The HTK model we used is established and highly used in

the literature, and while it is able to incorporate context *via* the sliding window and hidden Markov language model, we certainly would not claim that it understands or comprehends speech as humans can. Recently, large deep artificial neural network models pre-trained on a massive amount of unlabeled waveform features (e.g., Baevski et al., 2020; Hsu et al., 2021; Chen et al., 2022), have demonstrated strong generalization abilities to ASR and many para-linguistic speech tasks (Mohamed et al., 2022). While we would not claim that these larger models are capable of true understanding, it would nonetheless be interesting to apply the methods used in this paper to study similar types of models and tasks. This may contribute to understanding the functional organization of human auditory cortex and improve such large scale speech-based computational models.

Data availability statement

Publicly available datasets were analyzed in this study. This data can be found here: The 200 h Multi-genre Broadcast (MGB) dataset used to train the ASR model was only available to the 2015 MGB-1 Challenge (<http://www.mgb-challenge.org/>) participants with copyright restrictions from BBC. The 26-dimensional hidden layer representations extracted from the L7 layer of the DNN-BN₇ model can be found in (<http://mi.eng.cam.ac.uk/~cz277/stimuli>). Masked, preprocessed human neuroimaging data used for this analysis is available from figshare (<https://doi.org/10.6084/m9.figshare.5313484.v1>). The DNN-based ASR system was created using an open-source toolkit, the HTK toolkit version 3.5 (<https://htk.eng.cam.ac.uk/>). The RSA procedure for this paper was performed using the open-source RSA toolbox (https://github.com/rsagroup/rsatoolbox_matlab), with the addition of specific extensions for ssRSA for EMEG (<https://github.com/lisulab/rsatoolbox> and <https://github.com/lisulab/rsa-dnn-mapping>). RDMS were computed from DNN layer representations using publicly available scripts (<https://github.com/lisulab/htk-postprocessing>).

Ethics statement

The studies involving human participants were reviewed and approved by Peterborough and Fenland Ethical Committee (UK). The patients/participants provided their written informed consent to participate in this study.

Author contributions

CW: conceptualization, formal analysis, software, methodology, writing, and editing. CZ: conceptualization, formal analysis, software, methodology, writing, editing, and data curation. BD: methodology and editing. EF: data

acquisition, data curation, and editing. AT: software, data curation, and editing. XL: software, methodology, and editing. PW and WM-W: conceptualization, funding acquisition, supervision, and editing. LS: conceptualization, software, methodology, funding acquisition, supervision, writing, and editing. All authors contributed to the article and approved the submitted version.

Funding

This research was supported financially by a Senior Research Fellowship to LS from Alzheimer's Research UK (ARUK-SRF2017B-1), an Advanced Investigator grant to WM-W from the European Research Council (AdG 230570 NEUROLEX), by MRC Cognition and Brain Sciences Unit (CBSU) funding to WM-W (U.1055.04.002.00001.01), and by a European Research Council Advanced Investigator grant under the European Community's Horizon 2020 Research and Innovation Programme (2014-2020 ERC Grant agreement no 669820) to Lorraine K. Tyler.

Acknowledgments

The authors thank Anastasia Klimovich-Smith, Hun Choi, Lorraine Tyler, Andreas Marouchos, and Geoffrey Hinton for thoughtful comments and discussions. RSA computation was done in the RSA toolbox for Matlab (Nili et al., 2014) using custom EMEG and ssRSA extensions, to which Isma Zulfqar, Fawad Jamshed, and Jana Klimová also contributed.

Conflict of interest

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

Supplementary material

The Supplementary Material for this article can be found online at: <https://www.frontiersin.org/articles/10.3389/fncom.2022.1057439/full#supplementary-material>

References

- Arsenault, J. S., and Buchsbaum, B. R. (2015). Distributed neural representations of phonological features during speech perception. *J. Neurosci.* 35, 634–642. doi: 10.1523/JNEUROSCI.2454-14.2015
- Baevski, A., Zhou, Y., Mohamed, A., and Auli, M. (2020). “wav2vec 2.0: a framework for self-supervised learning of speech representations,” in *Proceedings of the 34th International Conference on Neural Information Processing Systems: NIPS'20, Vol. 33* (Vancouver, BC), 12449–12460.
- Baumann, S., Petkov, C. I., and Griffiths, T. D. (2013). A unified framework for the organization of the primate auditory cortex. *Front. Syst. Neurosci.* 7, 11. doi: 10.3389/fnsys.2013.00011
- Bell, P., Gales, M., Hain, T., Kilgour, J., Lanchantin, P., Liu, X., et al. (2015). “The MGB challenge: evaluating multi-genre broadcast media transcription,” in *Proc. ASRU* (Scottsdale, AZ), 687–693. doi: 10.1109/ASRU.2015.7404863
- Bishop, C. (2006). *Pattern Recognition and Machine Learning*. New York, NY: Springer.
- Boullard, H., and Morgan, N. (1993). *Connectionist Speech Recognition: A Hybrid Approach*. Norwell, MA: Kluwer Academic Publishers. doi: 10.1007/978-1-4615-3210-1
- Cadiou, C. F., Hong, H., Yamins, D. L., Pinto, N., Ardila, D., Solomon, E. A., et al. (2014). Deep neural networks rival the representation of primate IT cortex for core visual object recognition. *PLOS Comput. Biol.* 10, e1003963. doi: 10.1371/journal.pcbi.1003963
- Chan, A. M., Dykstra, A. R., Jayaram, V., Leonard, M. K., Travis, K. E., Gygi, B., et al. (2014). Speech-specific tuning of neurons in human superior temporal gyrus. *Cereb. Cortex* 24, 2679–2693. doi: 10.1093/cercor/bht127
- Chang, E. F., Rieger, J. W., Johnson, K., Berger, M. S., Barbaro, N. M., and Knight, R. T. (2010). Categorical speech representation in human superior temporal gyrus. *Nat. Neurosci.* 13, 1428–1432. doi: 10.1038/nn.2641
- Chen, S., Wang, C., Chen, Z., Wu, Y., Liu, S., Chen, Z., et al. (2022). WavLM: large-scale self-supervised pre-training for full stack speech processing. *arXiv preprint arXiv:2110.13900*. doi: 10.1109/JSTSP.2022.3188113
- Cichy, R. M., Khosla, A., Pantazis, D., and Oliva, A. (2016). Dynamics of scene representations in the human brain revealed by magnetoencephalography and deep neural networks. *Neuroimage*. 153, 346–358. doi: 10.1016/j.neuroimage.2016.03.063
- Clarke, A., Devereux, B. J., Randall, B., and Tyler, L. K. (2014). Predicting the time course of individual objects with MEG. *Cereb. Cortex* 25, 3602–3612. doi: 10.1093/cercor/bhu203
- Correia, J. M., Jansma, B. M., and Bonte, M. (2015). Decoding articulatory features from fMRI responses in dorsal speech regions. *J. Neurosci.* 35, 15015–15025. doi: 10.1523/JNEUROSCI.0977-15.2015
- Davies, D. L., and Bouldin, D. W. (1979). A cluster separation measure. *IEEE Trans. Pattern Anal. Mach. Intell.* PAMI-1, 224–227. doi: 10.1109/TPAMI.1979.4766909
- Deng, L., and Sun, D. X. (1994). A statistical approach to automatic speech recognition using the atomic speech units constructed from overlapping articulatory features. *J. Acoust. Soc. Am.* 95, 2702–2719. doi: 10.1121/1.409839
- Denys, K., Vanduffel, W., Fize, D., Nelissen, K., Peuskens, H., Van Essen, D., et al. (2004). The processing of visual shape in the cerebral cortex of human and nonhuman primates: a functional magnetic resonance imaging study. *J. Neurosci.* 24, 2551–2565. doi: 10.1523/JNEUROSCI.3569-03.2004
- Devereux, B. J., Clarke, A., and Tyler, L. K. (2018). Integrated deep visual and semantic attractor neural networks predict fMRI pattern-information along the ventral object processing pathway. *Nat. Sci. Rep.* 8, 10636. doi: 10.1038/s41598-018-28865-1
- Di Liberto, G. M., O’Sullivan, J. A., and Lalor, E. C. (2015). Low-frequency cortical entrainment to speech reflects phoneme-level processing. *Curr. Biol.* 25, 2457–2465. doi: 10.1016/j.cub.2015.08.030
- Doddipatla, R., Hasan, M., and Hain, T. (2014). “Speaker dependent bottleneck layer training for speaker adaptation in automatic speech recognition,” in *Proc. Interspeech*, 2199–2203. doi: 10.21437/Interspeech.2014-492
- Fonteneau, E., Bozic, M., and Marslen-Wilson, W. D. (2014). Brain network connectivity during language comprehension: interacting linguistic and perceptual subsystems. *Cereb. Cortex* 25, 3962–3976. doi: 10.1093/cercor/bhu283
- Gramfort, A., Luessi, M., Larson, E., Engemann, D. A., Strohmeier, D., Brodbeck, C., et al. (2014). MNE software for processing MEG and EEG data. *Neuroimage* 86, 446–460. doi: 10.1016/j.neuroimage.2013.10.027
- Grézl, F., Karafiát, M., Kontár, S., and Černocký, J. (2007). “Probabilistic and bottle-neck features for LVCSR of meetings,” in *Proc. ICASSP* (Honolulu, HI), 757–760.
- Güçlü, U., and van Gerven, M. A. (2015). Deep neural networks reveal a gradient in the complexity of neural representations across the ventral stream. *J. Neurosci.* 35, 10005–10014. doi: 10.1523/JNEUROSCI.5023-14.2015
- Hämäläinen, M. S., and Ilmoniemi, R. (1994). Interpreting magnetic fields of the brain: minimum norm estimates. *Med. Biol. Eng. Comput.* 32, 35–42. doi: 10.1007/BF02512476
- Hamilton, L. S., Oganian, Y., Hall, J., and Chang, E. F. (2021). Parallel and distributed encoding of speech across human auditory cortex. *Cell* 12, 4626.e13–4639.e13. doi: 10.1016/j.cell.2021.07.019
- Hinton, G., Deng, L., Yu, D., Dahl, G. E., Mohamed, A.-R., Jaitly, N., et al. (2012). Deep neural networks for acoustic modeling in speech recognition: the shared views of four research groups. *IEEE Signal Process. Mag.* 29, 82–97. doi: 10.1109/MSP.2012.2205597
- Hochreiter, S., and Schmidhuber, J. (1997). Long short-term memory. *Neural Comput.* 9, 1735–1780. doi: 10.1162/neco.1997.9.8.1735
- Hsu, W.-N., Bolte, B., Tsai, Y.-H. H., K., L., Salakhutdinov, R., et al. (2021). HuBERT: self-supervised speech representation learning by masked prediction of hidden units. *IEEE/ACM Trans. Audio Speech Lang. Process.* 29, 3451–3460. doi: 10.1109/TASLP.2021.3122291
- Jenkinson, M., Beckmann, C. F., Behrens, T. E., Woolrich, M. W., and Smith, S. M. (2012). FSL. *Neuroimage* 62, 782–790. doi: 10.1016/j.neuroimage.2011.09.015
- Karafiát, M., Grézl, F., Hannemann, M., Veselý, K., and Černocký, J. H. (2013). “BUT BABEL system for spontaneous Cantonese,” in *Proc. Interspeech* (Lyon), 2589–2593. doi: 10.21437/Interspeech.2013-582
- Khaligh-Razavi, S.-M., and Kriegeskorte, N. (2014). Deep supervised, but not unsupervised, models may explain IT cortical representation. *PLoS Comput. Biol.* 10, e1003915. doi: 10.1371/journal.pcbi.1003915
- Kheradpisheh, S. R., Ghodrati, M., Ganjtabesh, M., and Masquelier, T. (2016). Deep networks can resemble human feed-forward vision in invariant object recognition. *Nat. Sci. Rep.* 6, 32672. doi: 10.1038/srep32672
- Kriegeskorte, N. (2015). Deep neural networks: a new framework for modeling biological vision and brain information processing. *Annu. Rev. Vision Sci.* 1, 417–446. doi: 10.1146/annurev-vision-082114-035447
- Kriegeskorte, N., and Douglas, P. K. (2018). Cognitive computational neuroscience. *Nat. Neurosci.* 21, 1148–1160. doi: 10.1038/s41593-018-0210-5
- Kriegeskorte, N., Goebel, R., and Bandettini, P. (2006). Information-based functional brain mapping. *Proc. Natl. Acad. Sci. U.S.A.* 103, 3863–3868. doi: 10.1073/pnas.0600244103
- Kriegeskorte, N., and Kievit, R. A. (2013). Representational geometry: Integrating cognition, computation, and the brain. *Trends Cogn. Sci.* 17, 401–412. doi: 10.1016/j.tics.2013.06.007
- Kriegeskorte, N., Mur, M., and Bandettini, P. (2008a). Representational similarity analysis-connecting the branches of systems neuroscience. *Front. Syst. Neurosci.* 2, 4. doi: 10.3389/neuro.06.004.2008
- Kriegeskorte, N., Mur, M., Ruff, D. A., Kiani, R., Bodurka, J., Esteky, H., et al. (2008b). Matching categorical object representations in inferior temporal cortex of man and monkey. *Neuron* 60, 1126–1141. doi: 10.1016/j.neuron.2008.10.043
- Krizhevsky, A., Sutskever, I., and Hinton, G. (2012). “Imagenet classification with deep convolutional neural networks,” in *Proc. NIPS* (New York, NY).
- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proc. IEEE* 86, 2278–2324. doi: 10.1109/5.726791
- Liu, X., Flego, F., Wang, L., Zhang, C., Gales, M., and Woodland, P. (2015). “The Cambridge university 2014 BOLT conversational telephone Mandarin Chinese LVCSR system for speech translation,” in *Proc. Interspeech* (Dresden), 3145–3149. doi: 10.21437/Interspeech.2015-633
- Luscher, C., Beck, E., Irie, K., Kitza, M., Michel, W., Zeyer, A., et al. (2019). “RWTH ASR systems for LibriSpeech: hybrid vs attention,” in *Proc. Interspeech* (Graz), 231–235. doi: 10.21437/Interspeech.2019-1780
- Mack, M. L., Love, B. C., and Preston, A. R. (2016). Dynamic updating of hippocampal object representations reflects new conceptual knowledge. *Proc. Natl. Acad. Sci. U.S.A.* 201614048, 113, 13203–13208. doi: 10.1073/pnas.1614048113

- Mesgarani, N., Cheung, C., Johnson, K., and Chang, E. F. (2014). Phonetic feature encoding in human superior temporal gyrus. *Science* 343, 1006–1010. doi: 10.1126/science.1245994
- Mesgarani, N., David, S. V., Fritz, J. B., and Shamma, S. A. (2008). Phoneme representation and classification in primary auditory cortex. *J. Acoust. Soc. Am.* 123, 899–909. doi: 10.1121/1.2816572
- Mitra, V., Wang, W., Stolcke, A., Nam, H., Richey, C., Yuan, J., et al. (2013). “Articulatory trajectories for large-vocabulary speech recognition,” in *Proc. ICASSP* (Vancouver, BC: IEEE), 7145–7149. doi: 10.1109/ICASSP.2013.6639049
- Moerel, M., De Martino, F., and Formisano, E. (2014). An anatomical and functional topography of human auditory cortical areas. *Front. Neurosci.* 8, 225. doi: 10.3389/fnins.2014.00225
- Mohamed, A., Lee, H.-Y., Borgholt, L., Havtorn, J., Edin, J., Igel, C., et al. (2022). Self-supervised speech representation learning: a review. *arXiv preprint arXiv:2205.10643*. doi: 10.1109/JSTSP.2022.3207050
- Molins, A., Stufflebeam, S. M., Brown, E. N., and Hämläinen, M. S. (2008). Quantification of the benefit from integrating MEG and EEG data in minimum ℓ_2 -norm estimation. *Neuroimage* 42, 1069–1077. doi: 10.1016/j.neuroimage.2008.05.064
- Morgan, N. (2011). Deep and wide: multiple layers in automatic speech recognition. *IEEE Trans. Audio Speech Lang. Process.* 20, 7–13. doi: 10.1109/TASL.2011.2116010
- Moses, D. A., Leonard, M. K., and Chang, E. F. (2018). Real-time classification of auditory sentences using evoked cortical activity in humans. *J. Neural Eng.* 15, 036005. doi: 10.1088/1741-2552/aaab6f
- Moses, D. A., Mesgarani, N., Leonard, M. K., and Chang, E. F. (2016). Neural speech recognition: Continuous phoneme decoding using spatiotemporal representations of human cortical activity. *J. Neural Eng.* 13, 056004. doi: 10.1088/1741-2560/13/5/056004
- Nichols, T. E., and Holmes, A. P. (2002). Nonparametric permutation tests for functional neuroimaging: a primer with examples. *Hum. Brain Mapp.* 15, 1–25. doi: 10.1002/hbm.1058
- Nili, H., Wingfield, C., Walther, A., Su, L., Marslen-Wilson, W., and Kriegeskorte, N. (2014). A toolbox for representational similarity analysis. *PLoS Comput. Biol.* 10, e1003553. doi: 10.1371/journal.pcbi.1003553
- Orban, G. A., Van Essen, D., and Vanduffel, W. (2004). Comparative mapping of higher visual areas in monkeys and humans. *Trends Cogn. Sci.* 8, 315–324. doi: 10.1016/j.tics.2004.05.009
- Park, D., Chan, W., Zhang, Y., Chiu, C.-C., Zoph, B., Cubuk, E., et al. (2019). “SpecAugment: a simple data augmentation method for automatic speech recognition,” in *Proc. Interspeech* (Graz), 2613–2617. doi: 10.21437/Interspeech.2019-2680
- Peddinti, V., Povey, D., and Khudanpur, S. (2015). “A time delay neural network architecture for efficient modeling of long temporal contexts,” in *Proc. Interspeech* (Dresden), 3214–3218. doi: 10.21437/Interspeech.2015-647
- Rauschecker, J. P., and Scott, S. K. (2009). Maps and streams in the auditory cortex: nonhuman primates illuminate human speech processing. *Nat. Neurosci.* 12, 718–724. doi: 10.1038/nn.2331
- Rumelhart, D., McClelland, J., and PDP Research Group (1986). *Parallel Distributed Processing: Explorations in the Microstructure of Cognition, Volume 1: Foundations*. Cambridge, MA: MIT Press. doi: 10.7551/mitpress/5236.001.0001
- Saenz, M., and Langers, D. R. (2014). Tonotopic mapping of human auditory cortex. *Hear. Res.* 307, 42–52. doi: 10.1016/j.heares.2013.07.016
- Sammon, J. W. (1969). A nonlinear mapping for data structure analysis. *IEEE Trans. Comput.* 18, 401–409. doi: 10.1109/T-C.1969.222678
- Saon, G., Kurata, G., Sercu, T., Audhkhasi, K., Thomas, S., Dimitriadis, D., et al. (2017). “English conversational telephone speech recognition by humans and machines,” in *Proc. Interspeech* (Stockholm), 132–136. doi: 10.21437/Interspeech.2017-405
- Schuster, M., and Nakajima, K. (2012). “Japanese and Korean voice search,” in *Proc. ICASSP* (Kyoto), 5149–5152. doi: 10.1109/ICASSP.2012.6289079
- Smith, S. M., and Nichols, T. E. (2009). Threshold-free cluster enhancement: addressing problems of smoothing, threshold dependence and localisation in cluster inference. *Neuroimage* 44, 83–98. doi: 10.1016/j.neuroimage.2008.03.061
- Su, L., Fonteneau, E., Marslen-Wilson, W., and Kriegeskorte, N. (2012). “Spatiotemporal searchlight representational similarity analysis in EMEG source space,” in *Proc. PRNI* (London), 97–100. doi: 10.1109/PRNI.2012.26
- Su, L., Zulfikar, I., Jamshed, F., Fonteneau, E., and Marslen-Wilson, W. (2014). Mapping tonotopic organization in human temporal cortex: representational similarity analysis in EMEG source space. *Front. Neurosci.* 8, 368. doi: 10.3389/fnins.2014.00368
- Thwaites, A., Glasberg, B. R., Nimmo-Smith, I., Marslen-Wilson, W. D., and Moore, B. C. (2016). Representation of instantaneous and short-term loudness in the human cortex. *Front. Neurosci.* 10, 183. doi: 10.3389/fnins.2016.00183
- Toneva, M., and Wehbe, L. (2019). “Interpreting and improving natural-language processing (in machines) with natural language-processing (in the brain),” in *33rd Conference on Neural Information Processing Systems (NeurIPS 2019)*, eds H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox and R. Garnett (Vancouver).
- Tootell, R. B., Tsao, D., and Vanduffel, W. (2003). Neuroimaging weighs in: humans meet macaques in “primate” visual cortex. *J. Neurosci.* 23, 3981–3989. doi: 10.1523/JNEUROSCI.23-10-03981.2003
- Tüske, Z., Golik, P., Schlüter, R., and Ney, H. (2014). “Acoustic modeling with deep neural networks using raw time signal for LVCSR,” in *Proc. Interspeech* (Singapore), 890–894. doi: 10.21437/Interspeech.2014-223
- Van Essen, D. C., Lewis, J. W., Drury, H. A., Hadjikhani, N., Tootell, R. B., Bakircioglu, M., et al. (2001). Mapping visual cortex in monkeys and humans using surface-based atlases. *Vision Res.* 41, 1359–1378. doi: 10.1016/S0042-6989(01)00045-1
- Waibel, A., Hanazawa, T., Hinton, G., Shikano, K., and Lang, K. (1989). Phoneme recognition using time-delay neural networks. *IEEE Trans. Acoust. Speech Signal Process.* 37, 328–339. doi: 10.1109/29.21701
- Wen, H., Shi, J., Zhang, Y., Lu, K. -H., Cao, J., and Liu, Z. (2018). Neural encoding and decoding with deep learning for dynamic natural vision. *Cereb. Cortex* 28, 4136–4160. doi: 10.1093/cercor/bhx268
- Wingfield, C., Su, L., Liu, X., Zhang, C., Woodland, P., Thwaites, A., et al. (2017). Relating dynamic brain states to dynamic machine states: human and machine solutions to the speech recognition problem. *PLoS Comput. Biol.* 13, e1005617. doi: 10.1371/journal.pcbi.1005617
- Woodland, P., Liu, X., Qian, Y., Zhang, C., Gales, M., Karanasou, P., et al. (2015). “Cambridge University transcription systems for the Multi-genre Broadcast Challenge,” in *Proc. ASRU* (Scottsdale, AZ), 639–646. doi: 10.1109/ASRU.2015.7404856
- Wu, Y., Schuster, M., Chen, Z., Le, Q., Norouzi, M., Macherey, W., et al. (2016). Google’s neural machine translation system: bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*.
- Xiong, W., Wu, L., Droppo, J., Huang, X., and Stolcke, A. (2018). “The Microsoft 2016 conversational speech recognition system,” in *Proc. ICASSP* (New Orleans), 5255–5259. doi: 10.1109/ICASSP.2017.7953159
- Young, S., Evermann, G., Gales, M., Hain, T., Kershaw, D., Liu, X., et al. (2015). *The HTK Book (for HTK version 3.5)*. Cambridge: Cambridge University Engineering Department.
- Young, S. J., Odell, J. J., and Woodland, P. C. (1994). “Tree-based state tying for high accuracy acoustic modelling,” in *Proc. HLT* (Stroudsburg, PA: Association for Computational Linguistics), 307–312. doi: 10.3115/1075812.1075885
- Yu, Z., Chuangsuwanich, E., and Glass, J. (2014). “Extracting deep neural network bottleneck features using low-rank matrix factorization,” in *Proc. ICASSP* (Florence), 185–189. doi: 10.1109/ICASSP.2014.6853583
- Zhang, C., Liu, Y., and Lee, C.-H. (2011). “Detection-based accented speech recognition using articulatory features,” in *Proc. ASRU* (Waikoloa), 500–505. doi: 10.1109/ASRU.2011.6163982
- Zhang, C., and Woodland, P. C. (2015a). “A general artificial neural network extension for HTK,” in *Proc. Interspeech* (Dresden), 3581–3585. doi: 10.21437/Interspeech.2015-710
- Zhang, C., and Woodland, P. C. (2015b). “Parameterised sigmoid and ReLU hidden activation functions for DNN acoustic modelling,” in *Proc. Interspeech* (Dresden), 3224–3228. doi: 10.21437/Interspeech.2015-649