

On the Formation of Phoneme Categories in DNN Acoustic Models

Tasha Nagamine

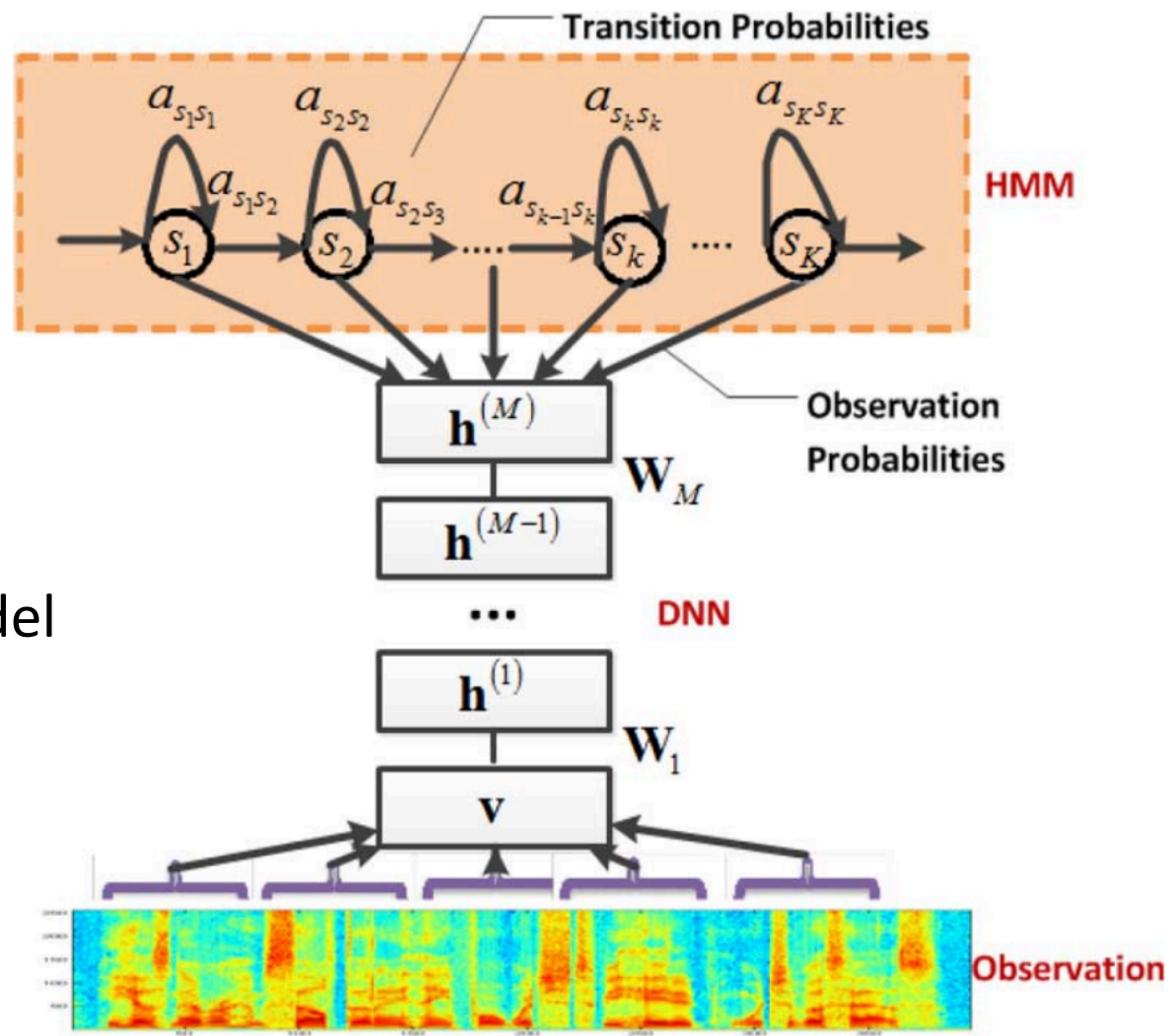
Department of Electrical Engineering, Columbia University

October 14, 2015

Motivation

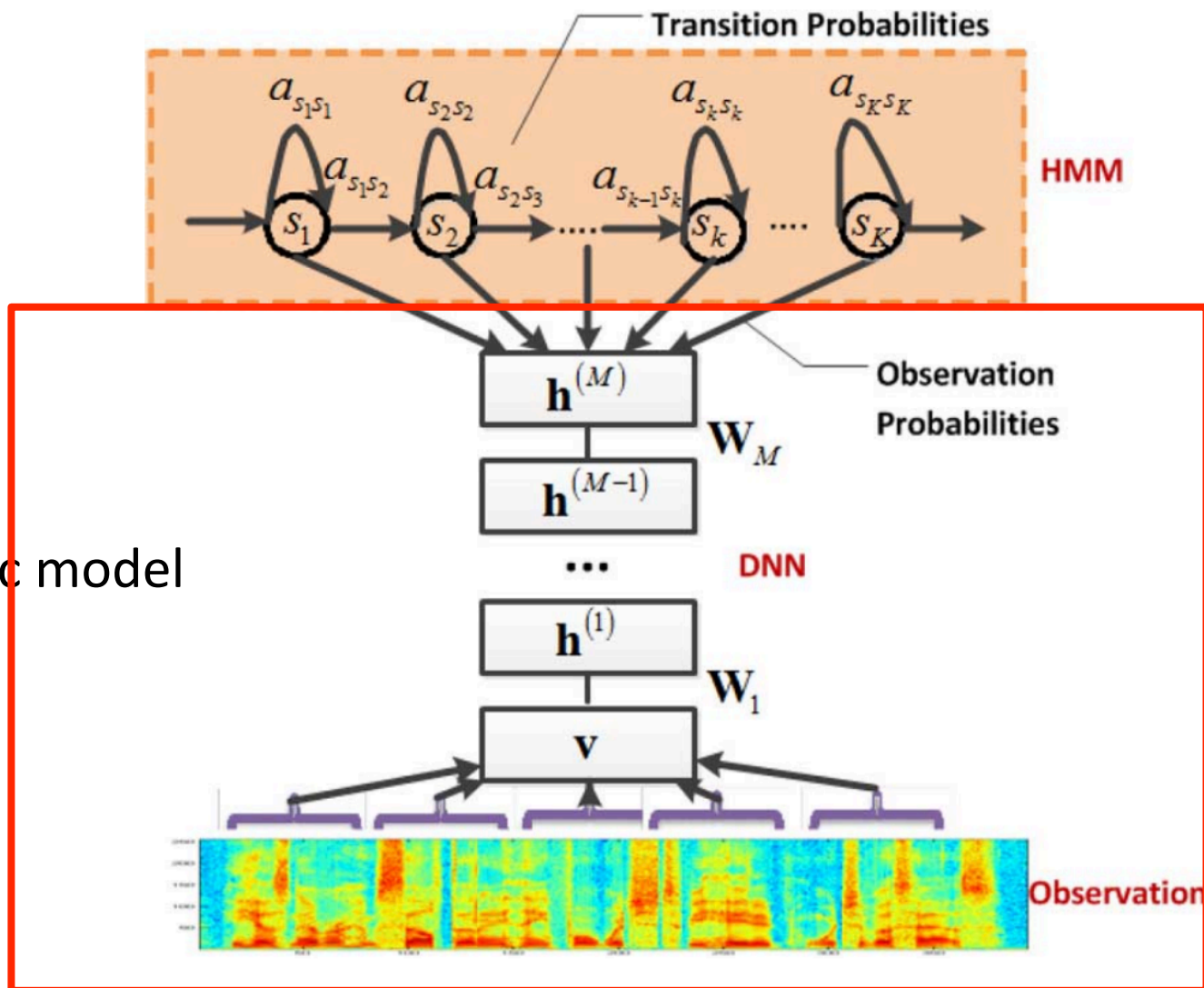
- Large performance gap between humans and state-of-the-art ASR systems
- Computational principles of DNNs remain elusive; they are analytically intractable
- Improving these models requires a better understanding of their transformations

Introduction to acoustic models



acoustic model

Introduction to acoustic models



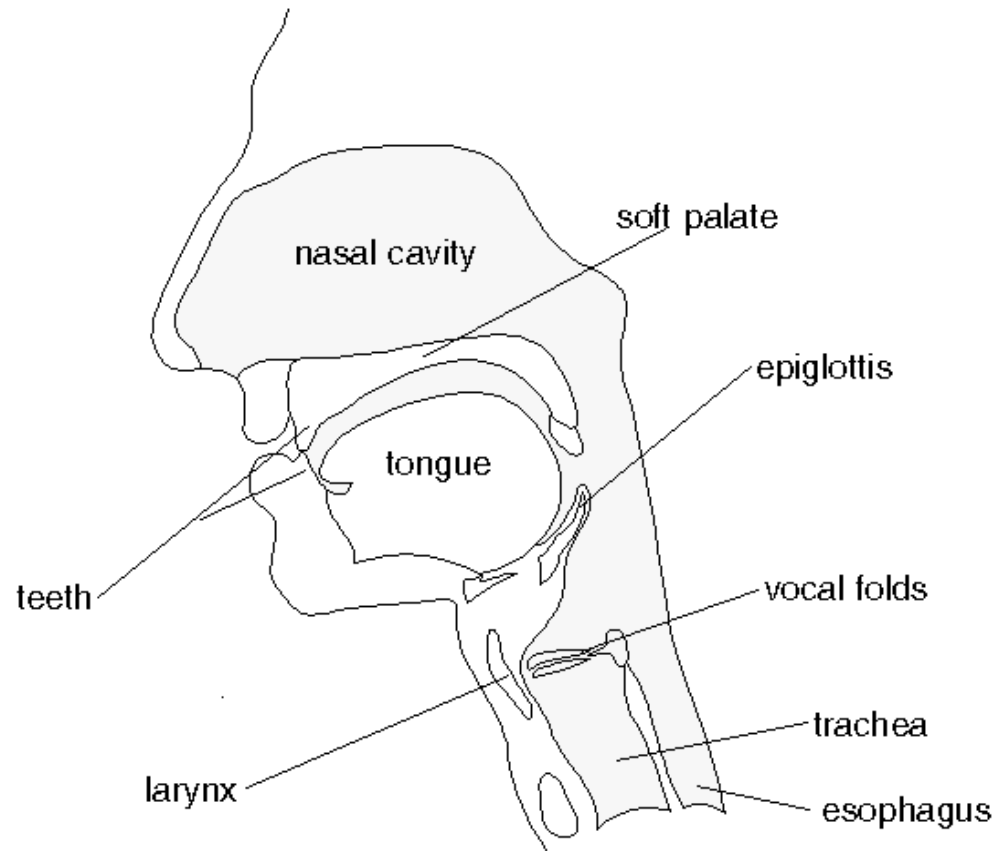
acoustic model

Phonemes

Smallest contrastive unit in language

- e.g., “k” vs. “b” in cat/bat
- ~40-60 in English

Output target in acoustic modeling



Phonetic Features

Manner of articulation

Place of articulation

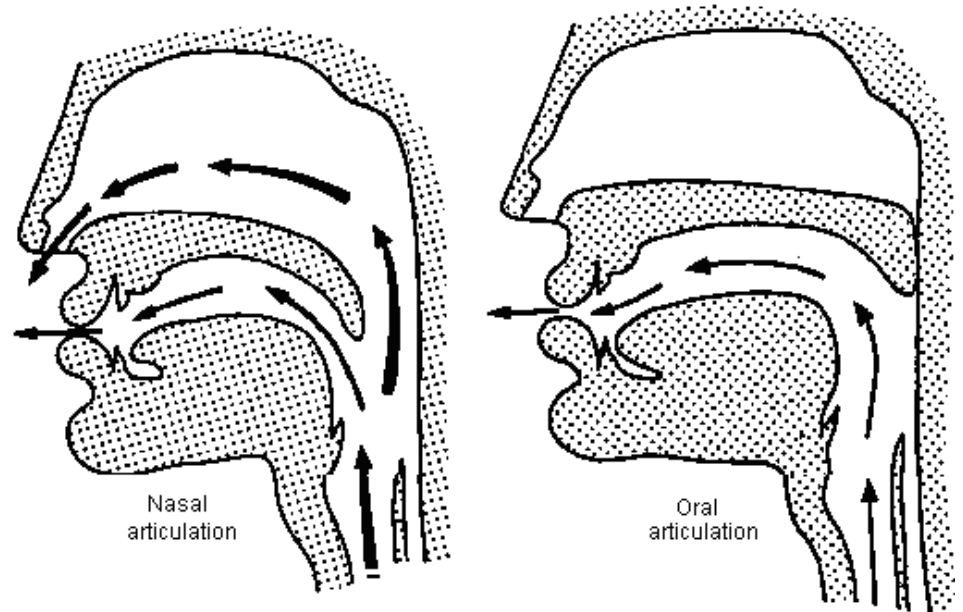
Voicing

Phonetic Features

Manner of articulation

Place of articulation

Voicing

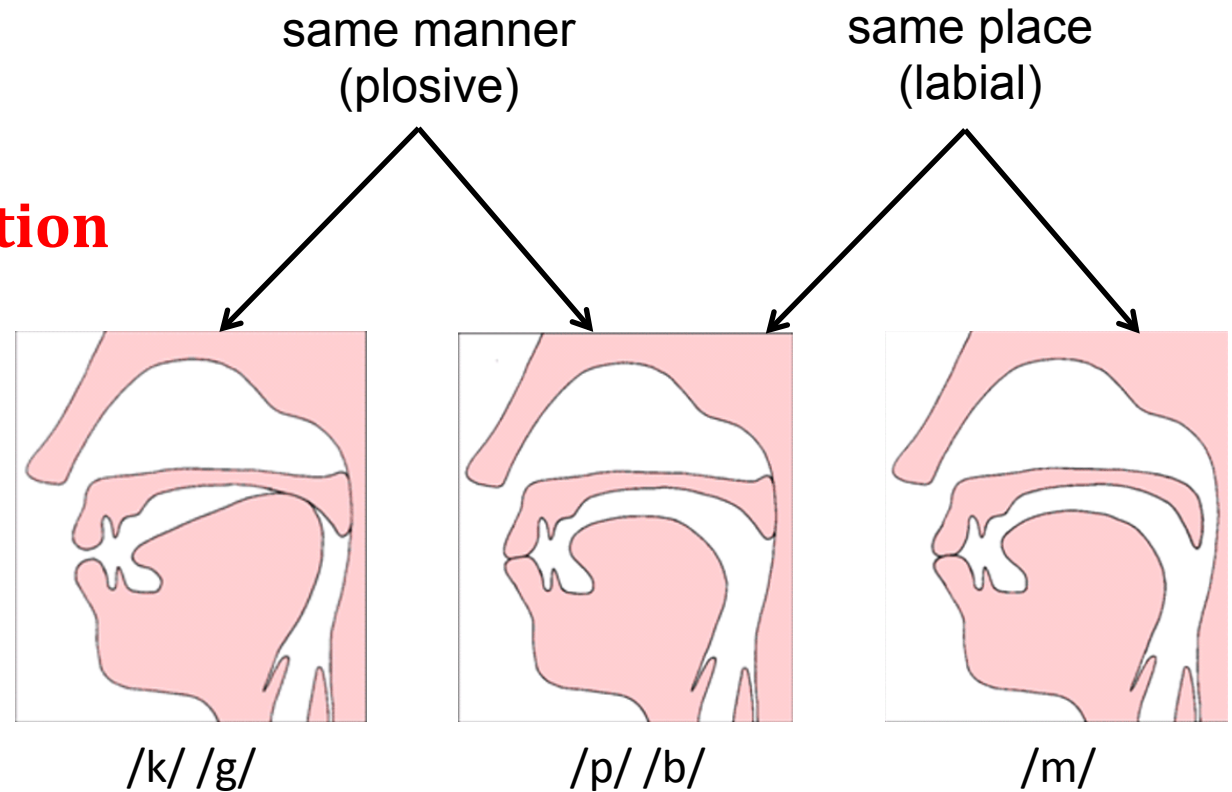


Phonetic Features

Manner of articulation

Place of articulation

Voicing

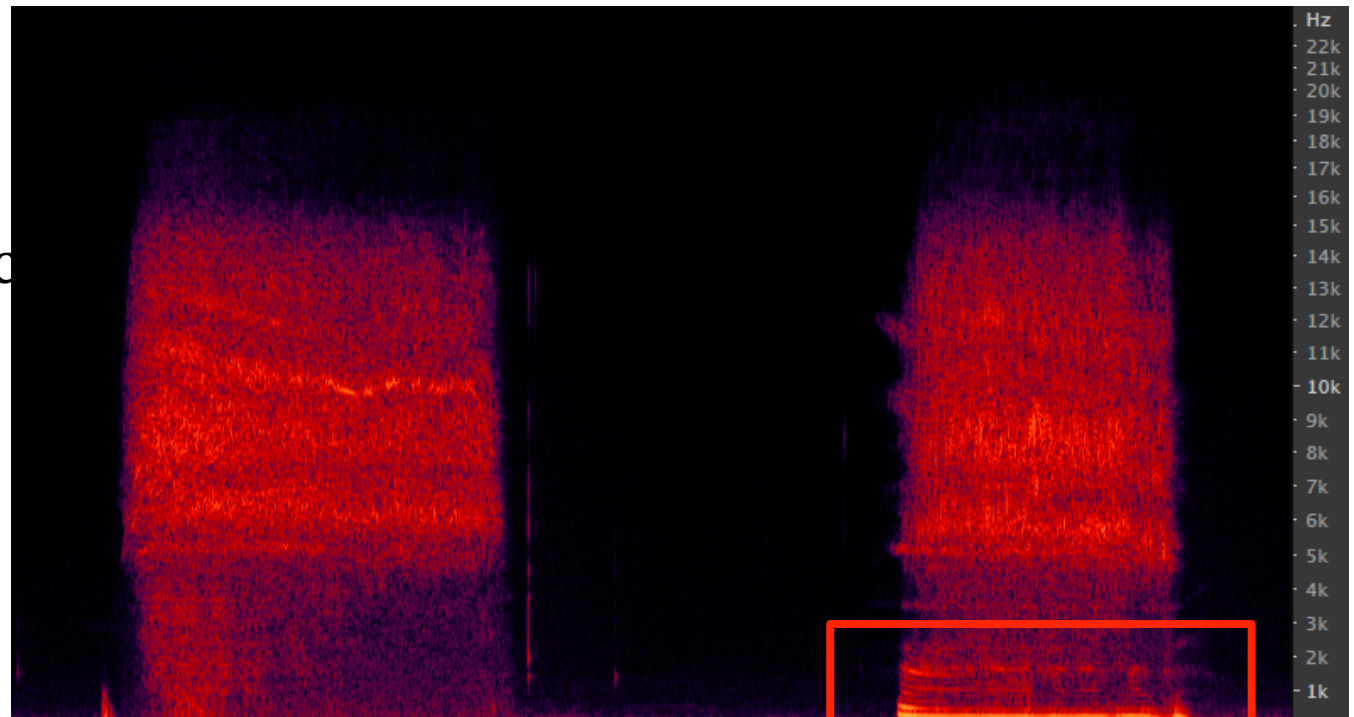


Phonetic Features

Manner of articulation

Place of articulation

Voicing

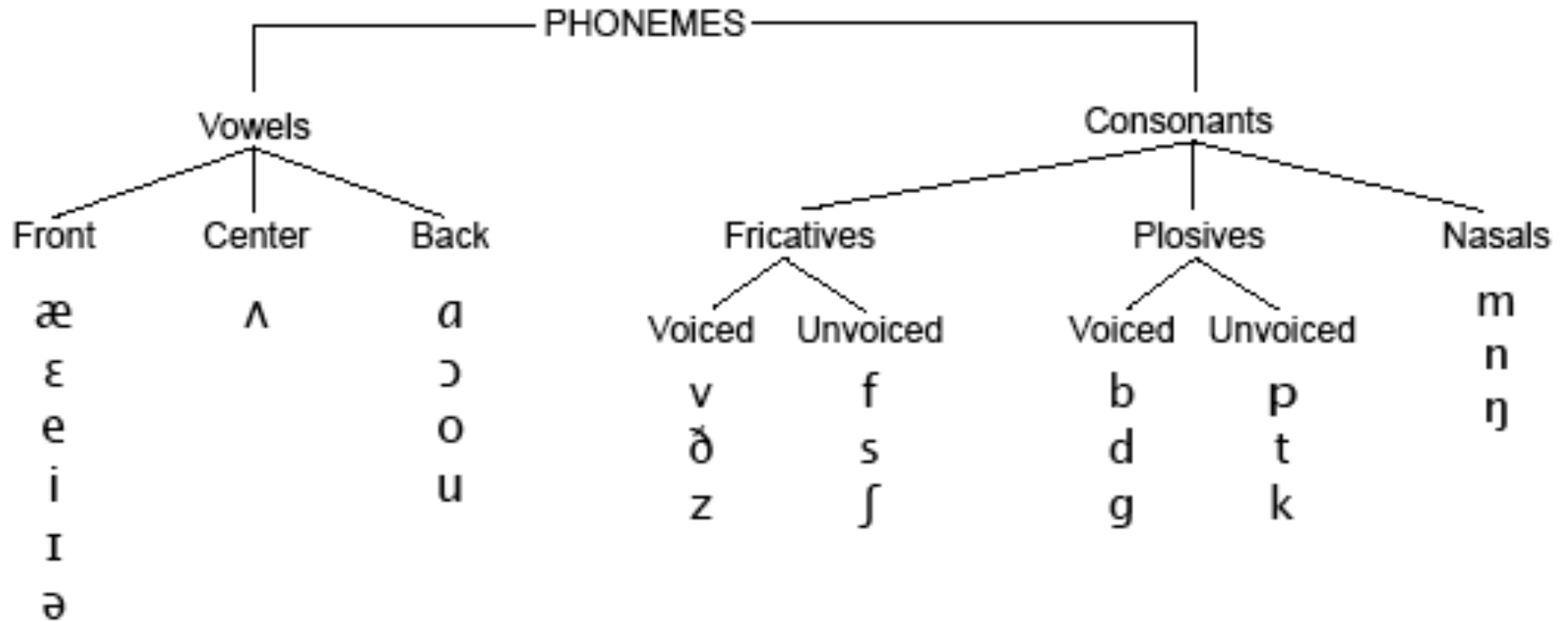


/s/ = unvoiced

/z/ = voiced

same manner (fricative)
same place (alveolar)

Phonetic Features



Distinctive Features
Chomsky, Halle, Stevens

Distinctive Features

Table 1. Distinctive Features of American English Consonants

	p	b	m	f	v	θ	ð	t	d	n	s	z	l	r	ʃ	ʒ	tʃ	dʒ	j	ɹ	k	g	ŋ	w	ʔ	h
Back	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	+	+	+	+	+	+
High	-	-	-	-	-	-	-	-	-	-	-	-	-	-	+	+	+	+	+	+	+	+	+	+	-	-
Coronal	-	-	-	-	-	+	+	+	+	+	+	+	+	+	+	+	+	+	-	-	-	-	-	-	-	-
Anterior	+	+	+	+	+	+	+	+	+	+	+	+	+	+	-	-	-	-	-	-	-	-	-	-	-	-
Labial	+	+	+	+	+	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	+	-	-
Continuant	-	-	-	+	+	+	+	-	-	-	+	+	+	-	+	+	-	-	+	+	-	-	-	+	-	+
Lateral	-	-	-	-	-	-	-	-	-	-	-	-	+	-	-	-	-	-	-	-	-	-	-	-	-	-
Nasal	-	-	+	-	-	-	-	-	-	+	-	-	-	-	-	-	-	-	-	-	-	-	+	-	-	-
Sonorant	-	-	+	-	-	-	-	-	-	+	-	-	+	+	-	-	-	-	+	+	-	-	+	+	-	-
Strident	-	-	-	+	+	-	-	-	-	-	+	+	-	-	+	+	+	+	-	-	-	-	-	-	-	-
Voiced	-	+	+	-	+	-	+	-	+	+	-	+	+	+	-	+	-	+	+	+	-	+	+	+	-	-

Table 2. Distinctive Features of American English Vowels

i	ɪ	e	ɛ	æ	u	ʊ	o	ɔ	a	ʌ	ə	
+	+	-	-	-	+	+	-	-	-	-	-	high
-	-	-	-	+	-	-	-	-	+	+	-	low
-	-	-	-	-	+	+	+	+	+	-	-	back
-	-	-	-	-	+	+	+	+	-	-	-	rounded
+	-	+	-	-	+	-	+	-	-	-	-	ATR

Phonemes and phones

Phoneme

Smallest contrastive unit in language.

Abstract idea.

Phone

Instances of phonemes in actual utterances.

Physical segments.

Example: “pat” vs. “bat”

4 phonemes

6 phones

Exploring How Deep Neural Networks Form Phonemic Categories

Tasha Nagamine¹, Michael L. Seltzer², Nima Mesgarani¹

¹Department of Electrical Engineering, Columbia University, New York, USA

²Microsoft Research, Redmond, USA

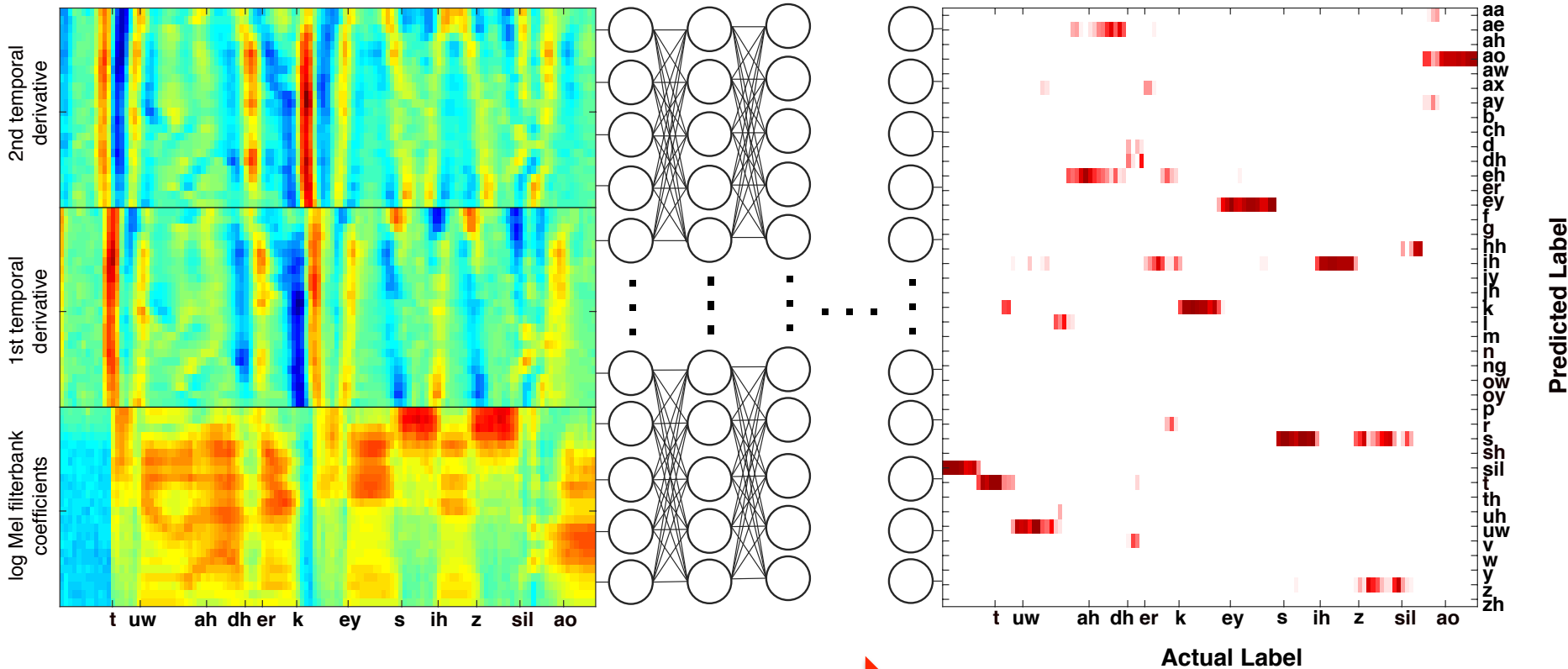
tasha.nagamine@columbia.edu, mseltzer@microsoft.com, nima@ee.columbia.edu



Input

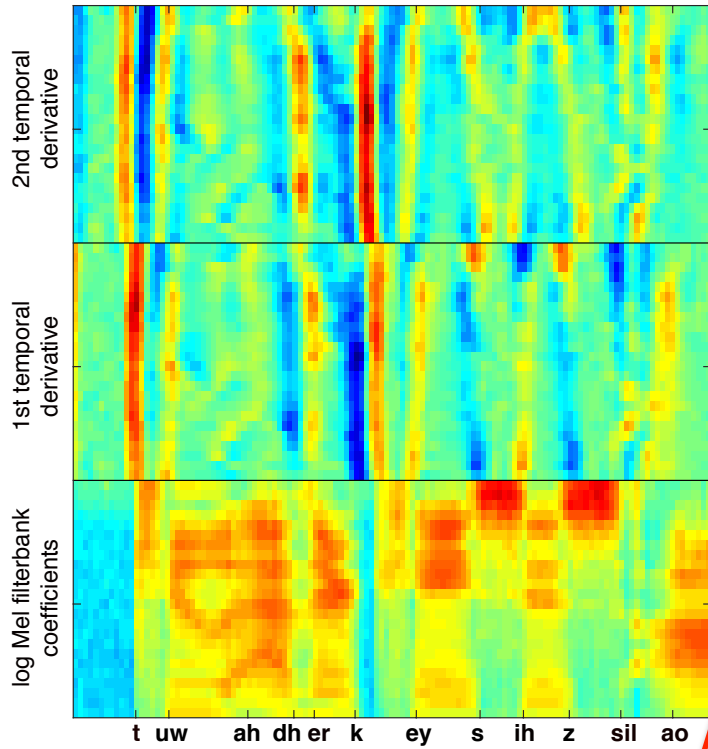
Hidden Layers

Output

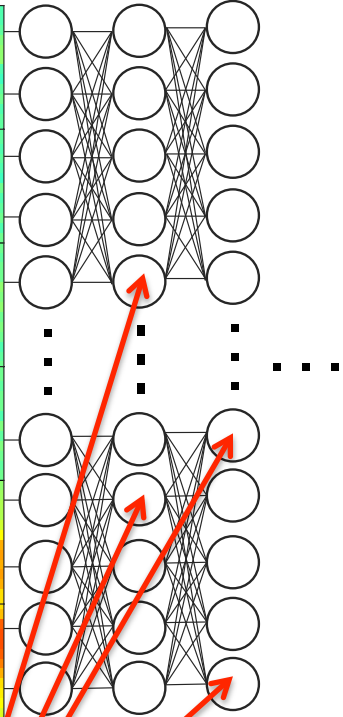


→
Feed-forward series of
nonlinear transformations...

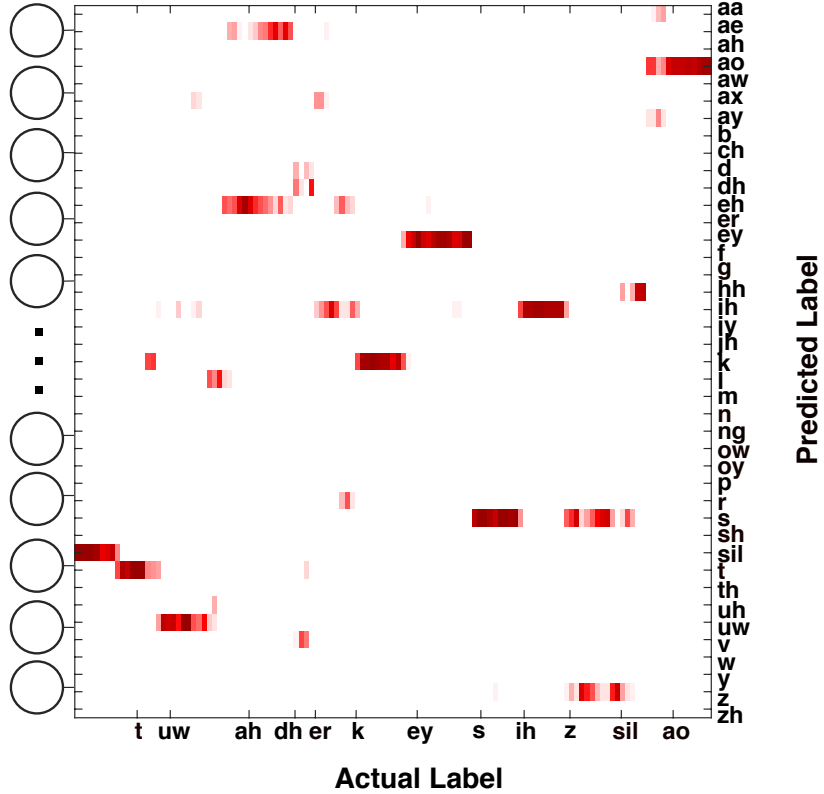
Input



Hidden Layers

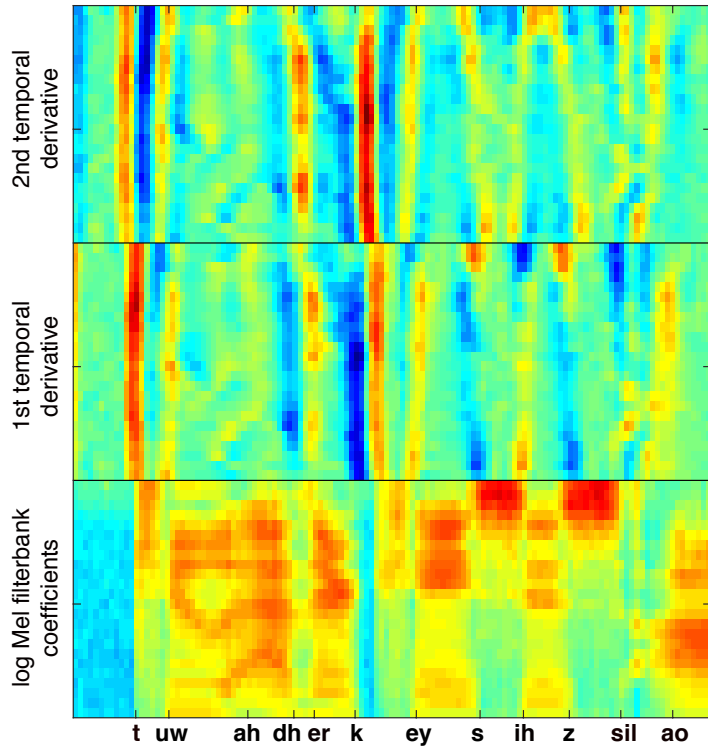


Output

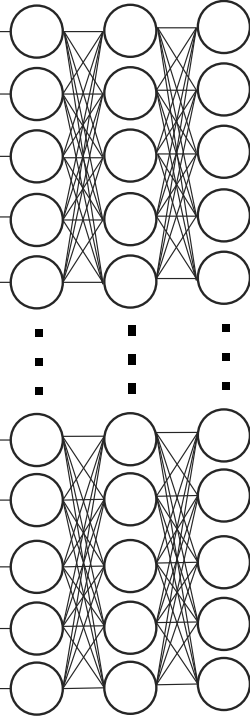


?

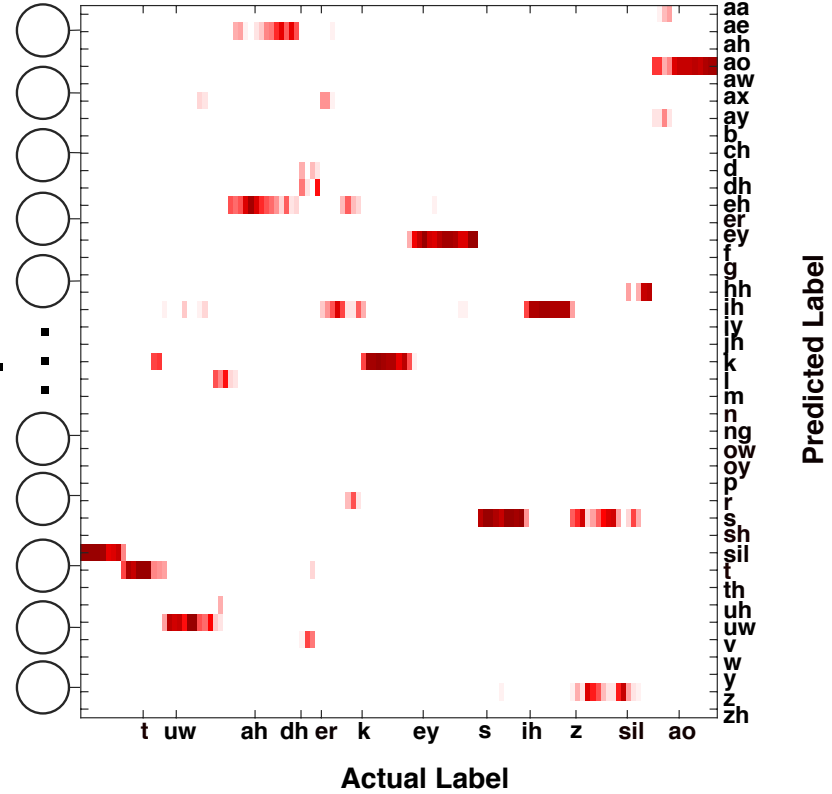
Input



Hidden Layers



Output



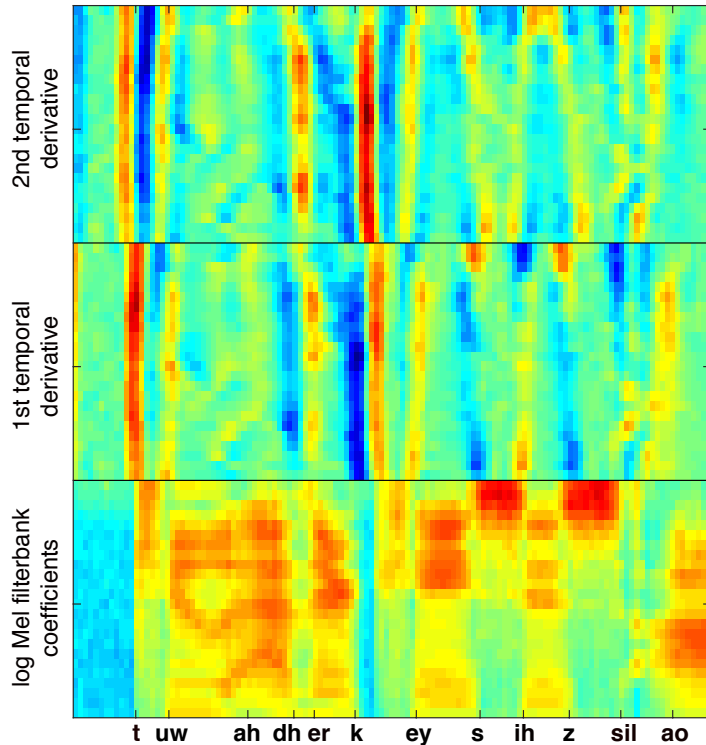
DNN Architecture

Input layer

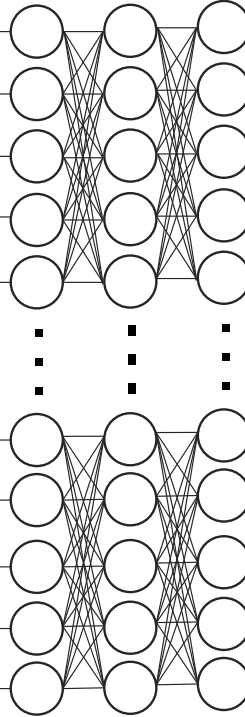
11 frames of 24-dimensional log Mel filter bank coefficients + deltas



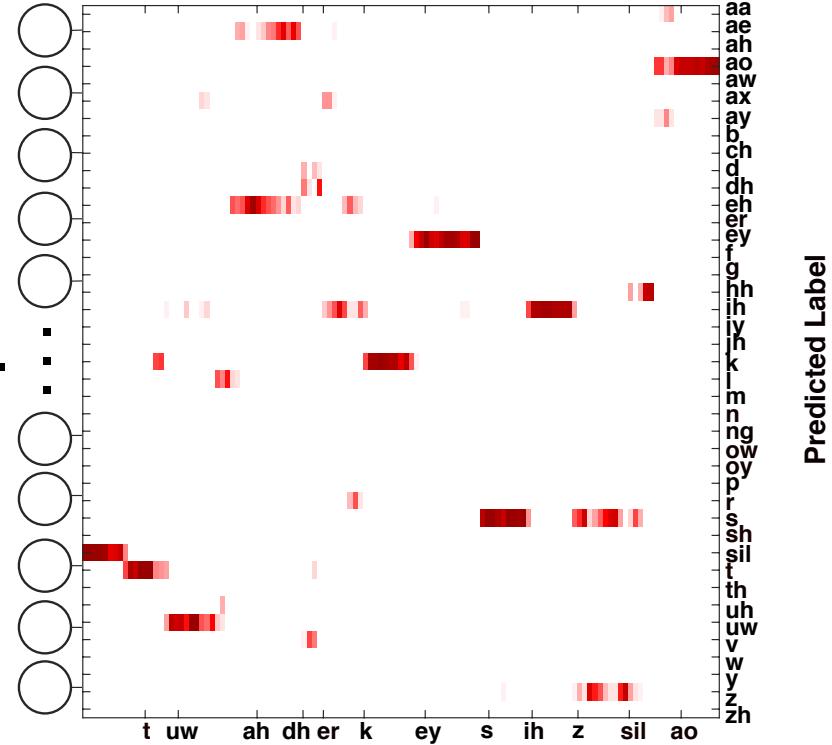
Input



Hidden Layers



Output



Actual Label

Predicted Label

DNN Architecture

Input layer

11 frames of 24-dimensional log Mel filter bank coefficients + deltas

5 sigmoid hidden layers

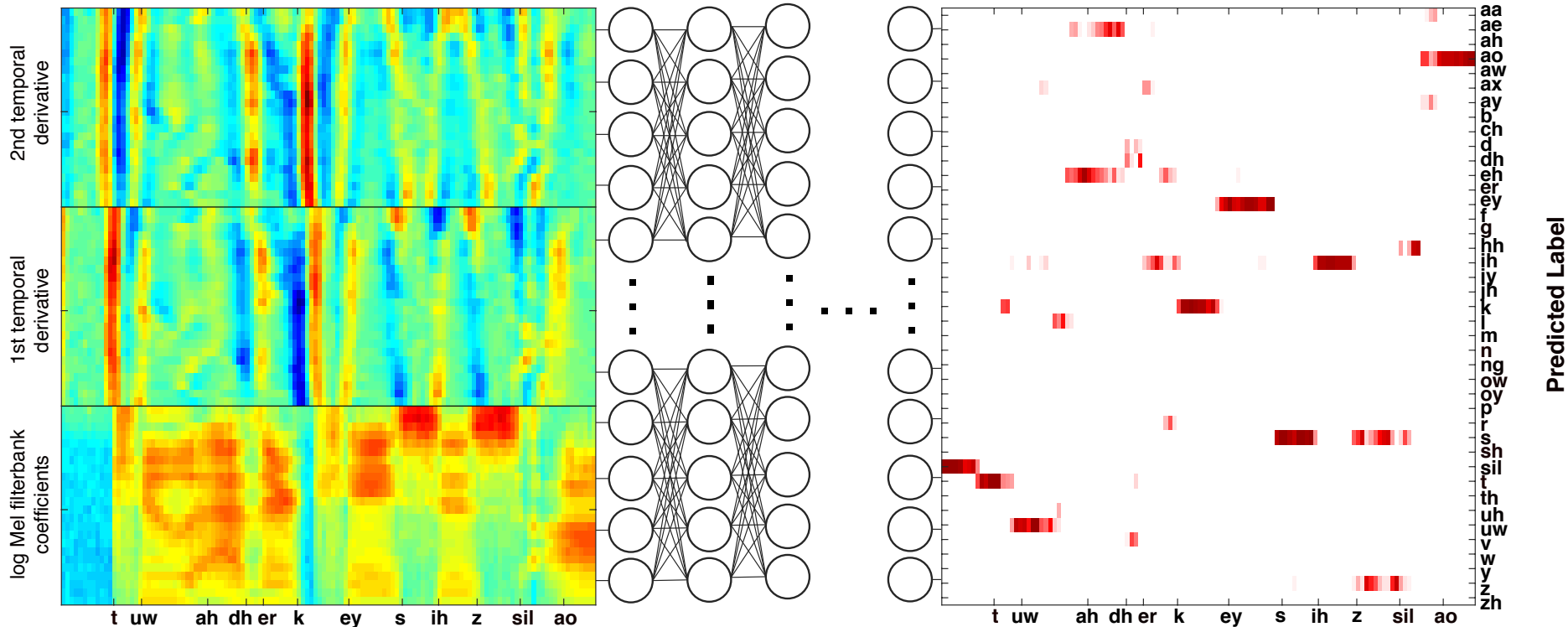
256 nodes each; fully connected feed-forward



Input

Hidden Layers

Output



Actual Label

Nagamine et al. Interspeech 2015

October 14, 2015

DNN Architecture

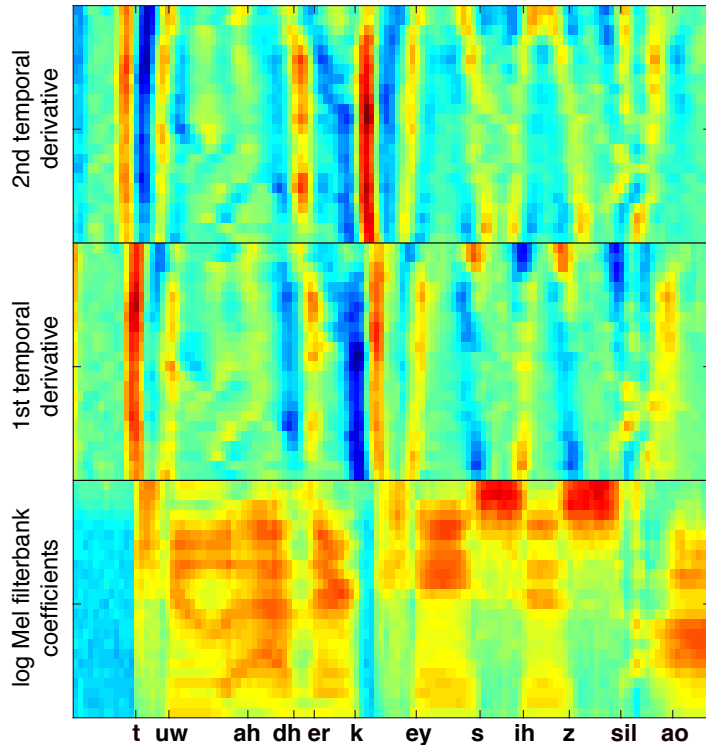
Input layer

11 frames of 24-dimensional log Mel filter bank coefficients + deltas

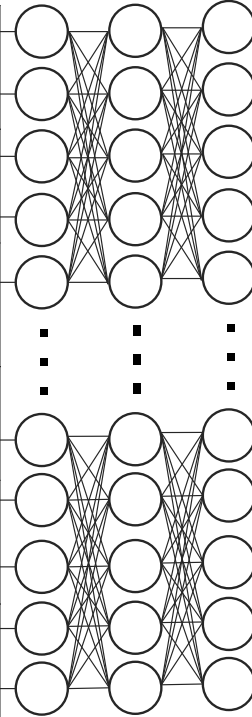
5 sigmoid hidden layers
256 nodes each;
fully connected feed-forward

Softmax output layer
41 nodes for 40 phonemes and silence; context independent

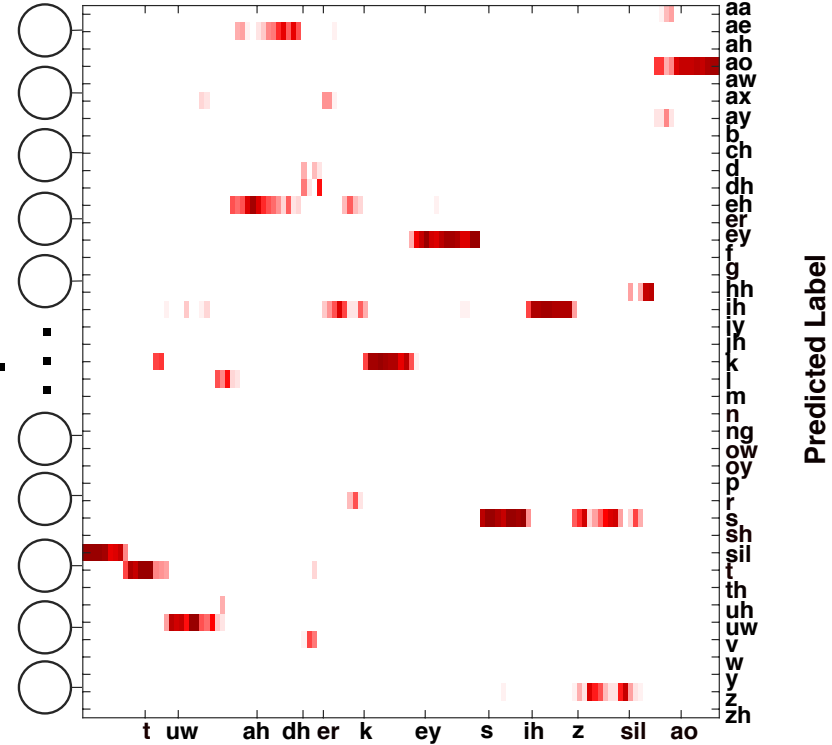
Input



Hidden Layers



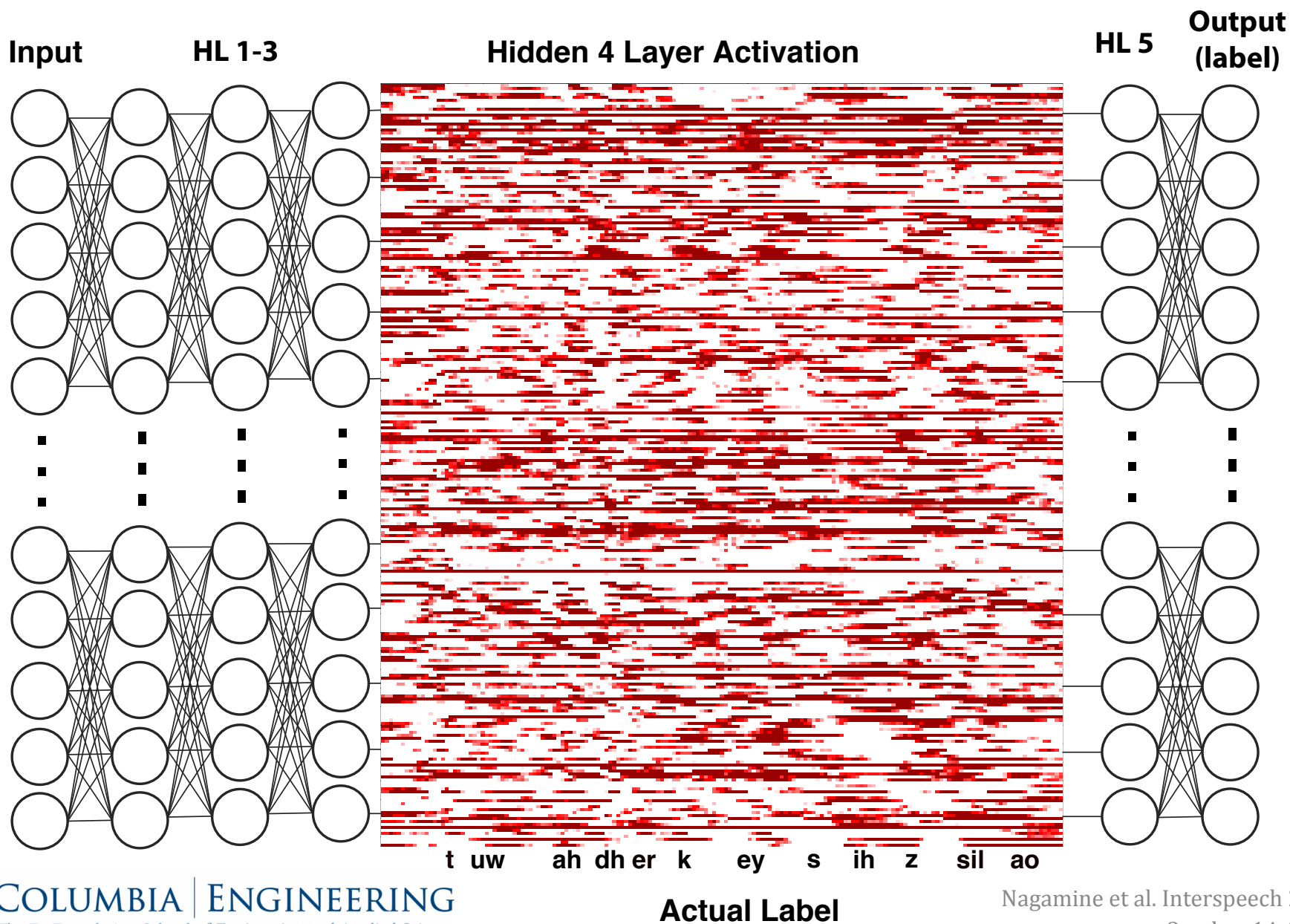
Output



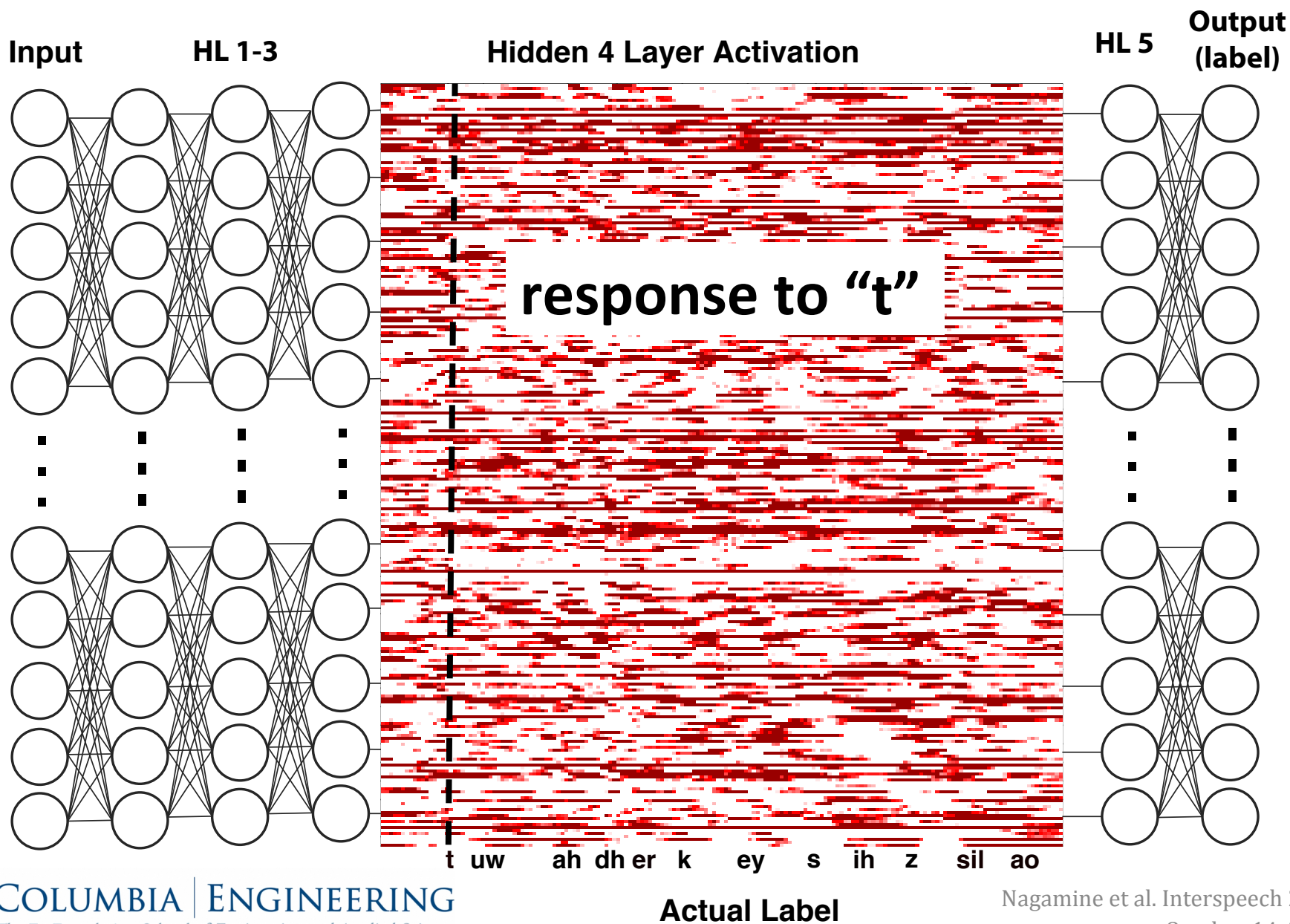
Actual Label

Predicted Label

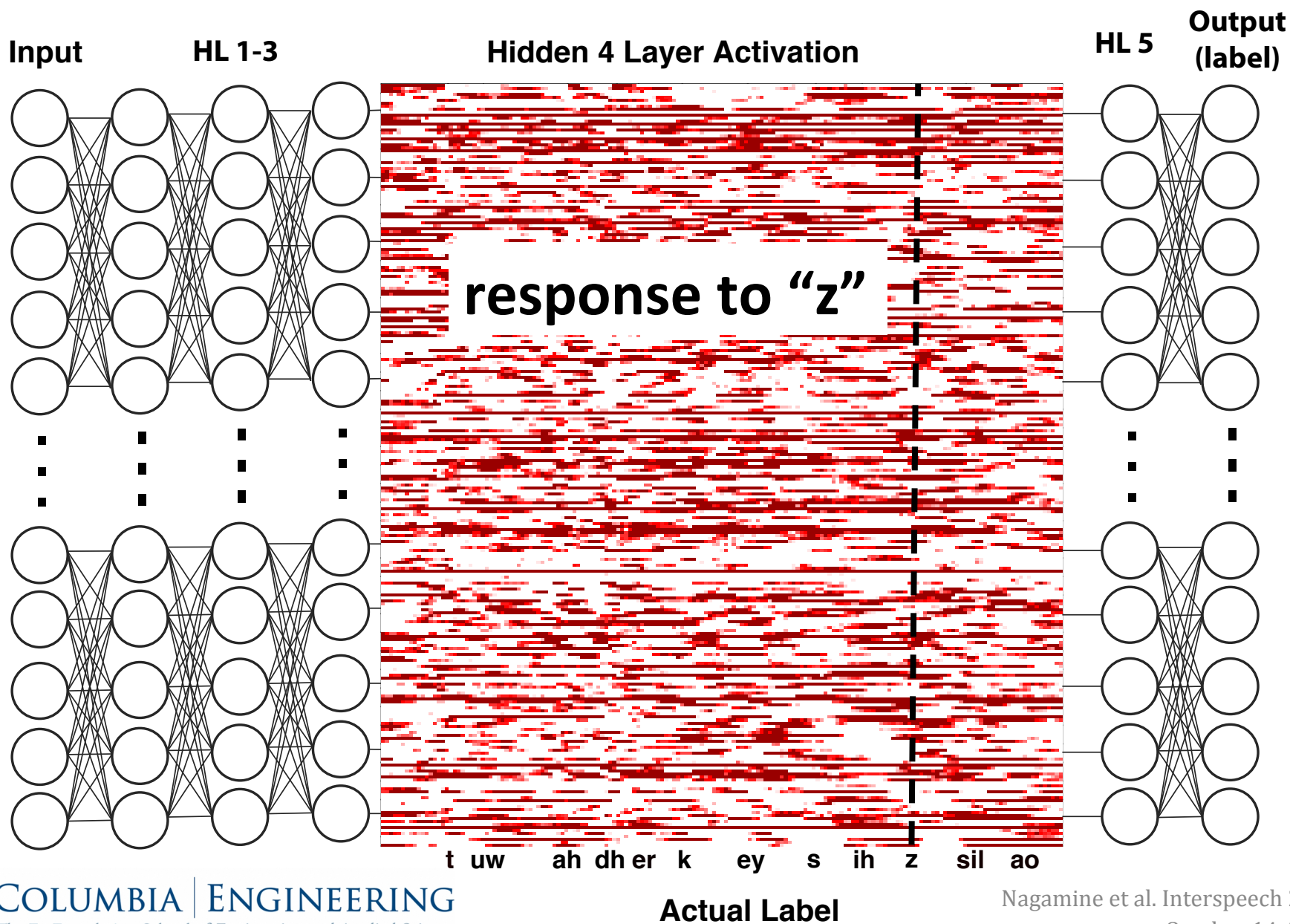
Speech stimuli & DNN activations



Speech stimuli & DNN activations



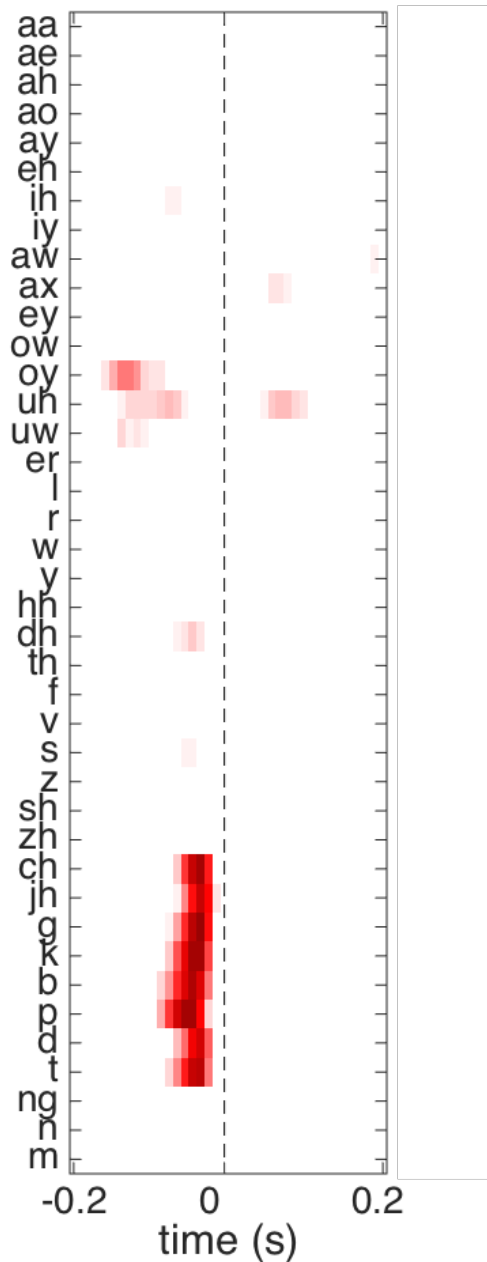
Speech stimuli & DNN activations



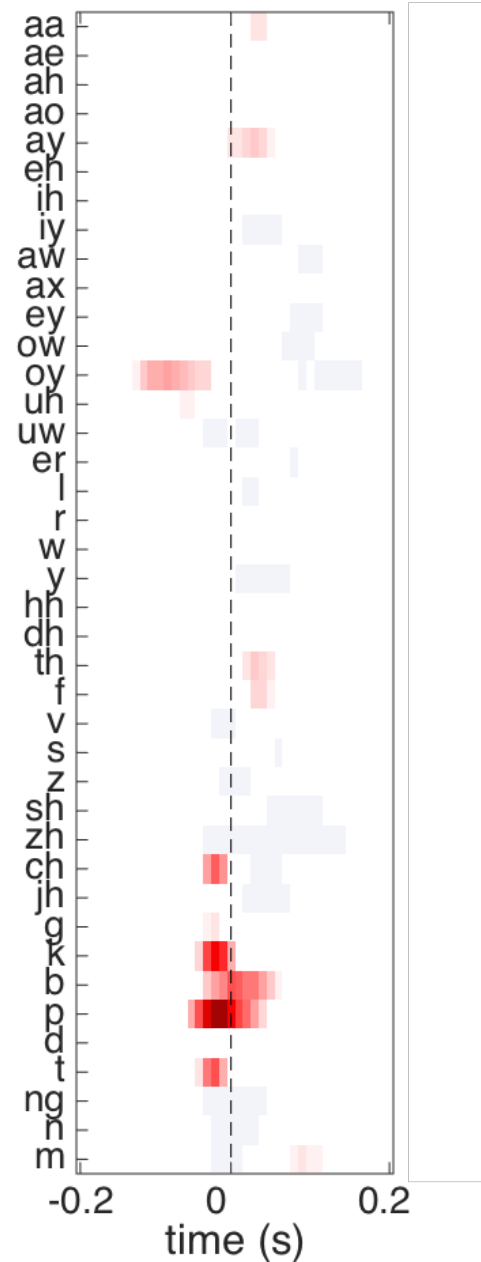
Summary of findings

1. Nodes are selective to phonetic features at the individual and population level

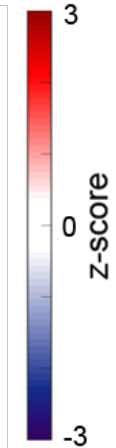
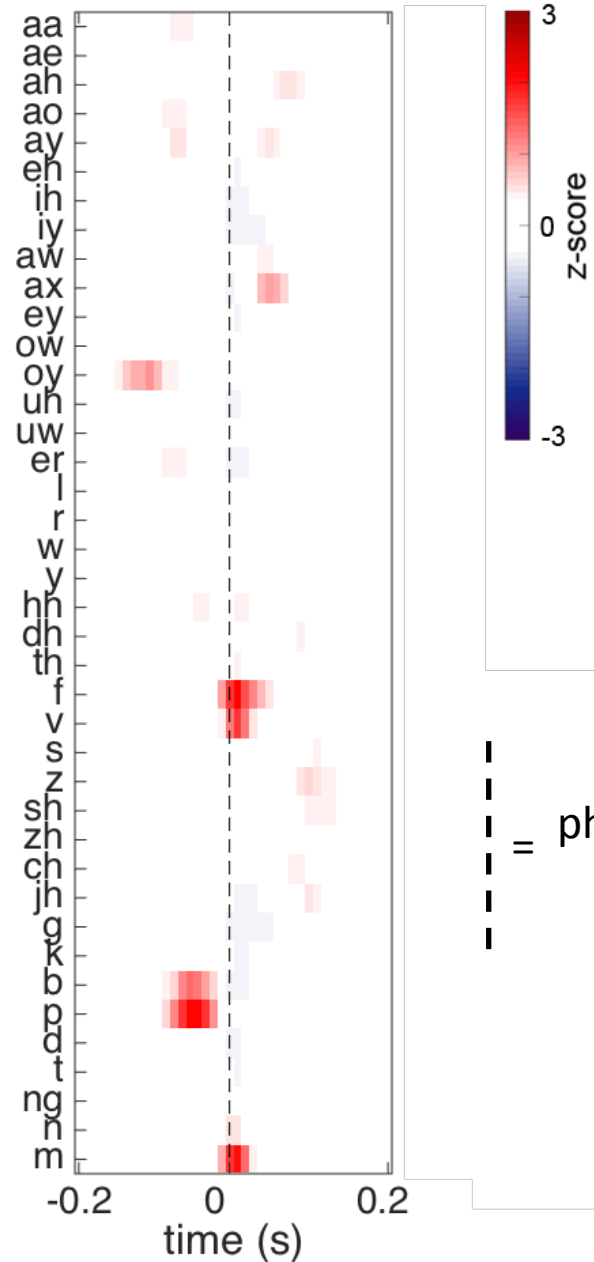
Node 16



Node 191

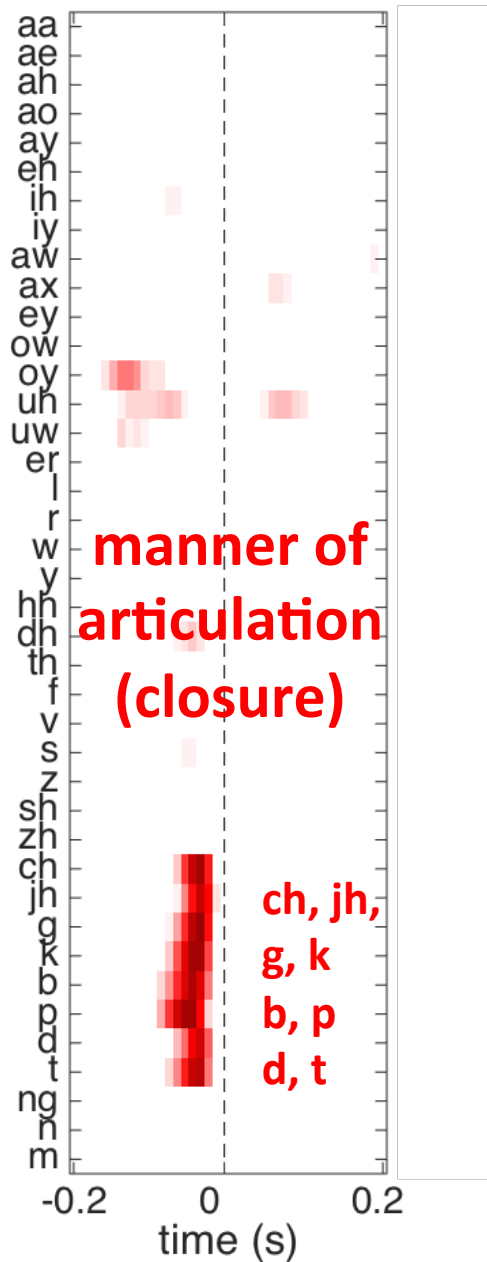


Node 165

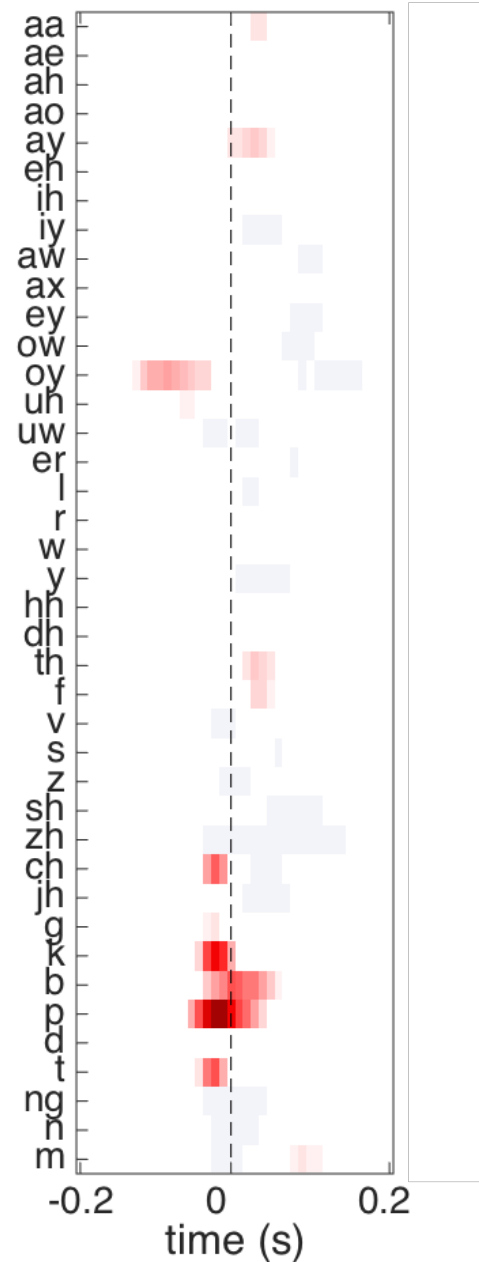


--- = phoneme onset

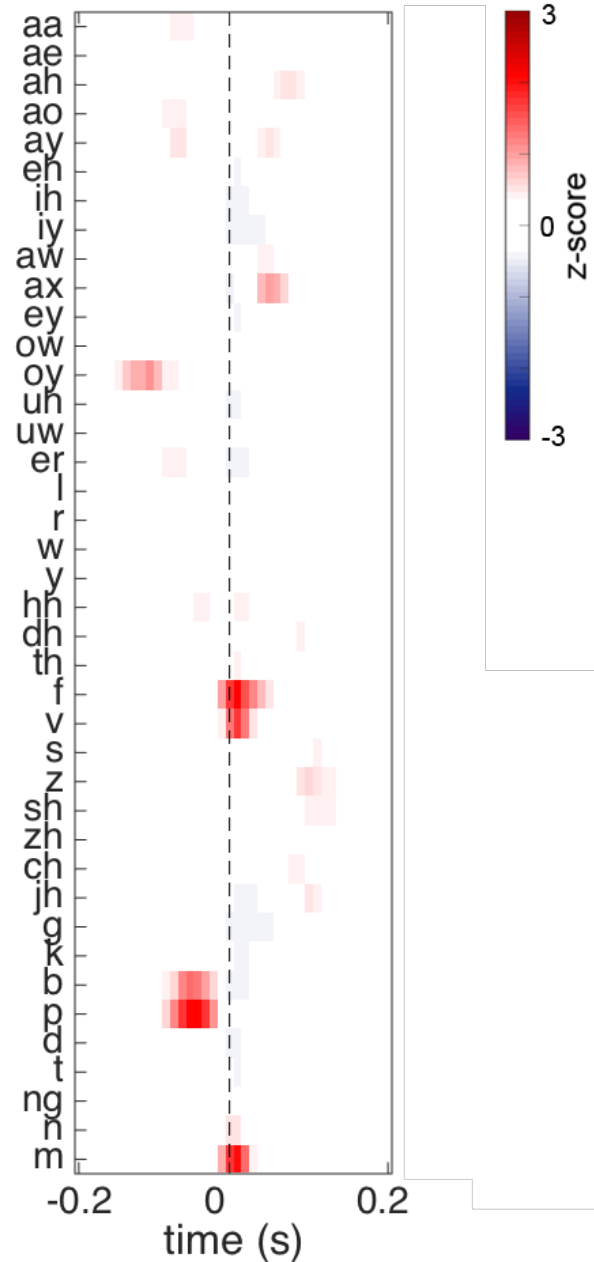
Node 16



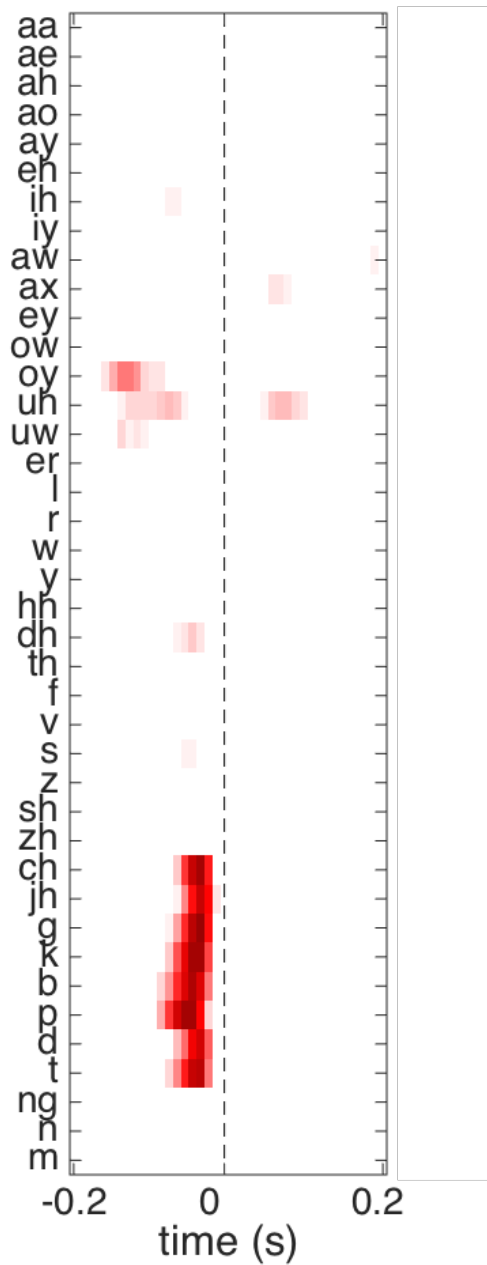
Node 191



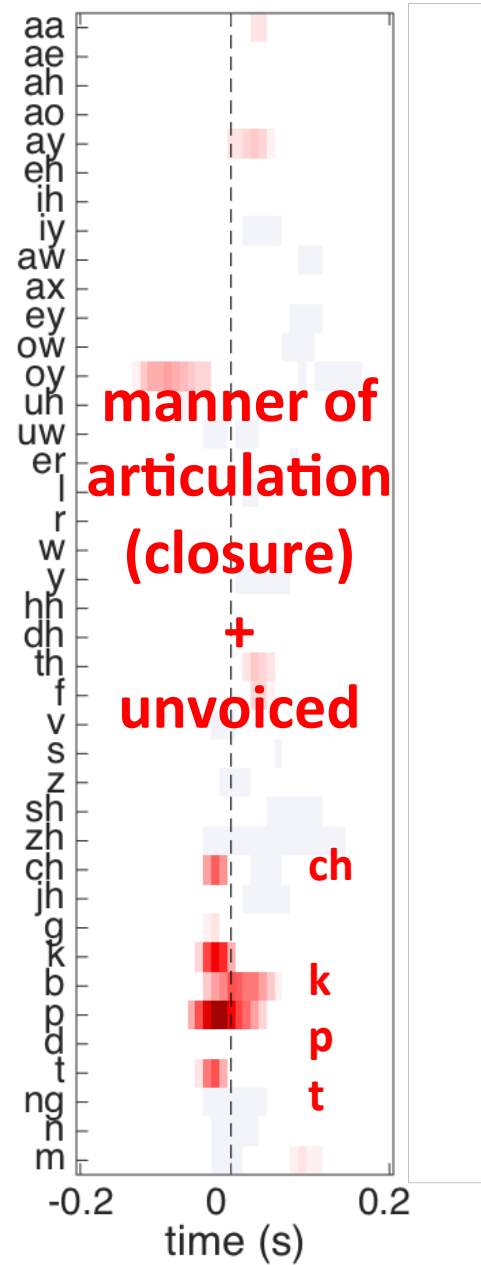
Node 165



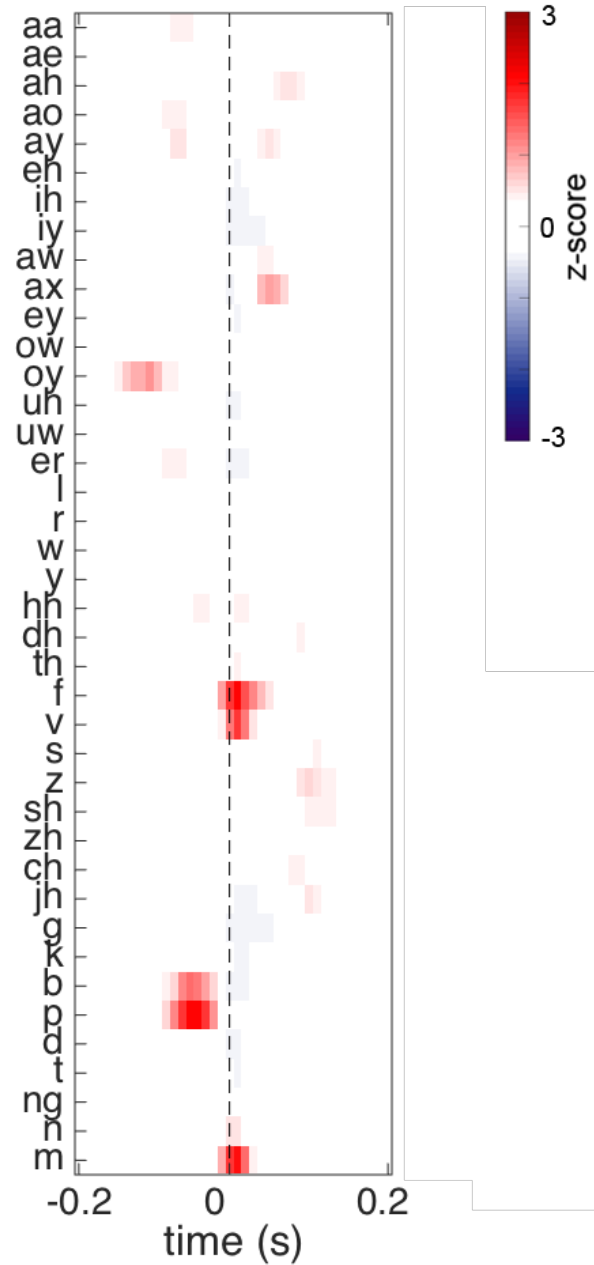
Node 16



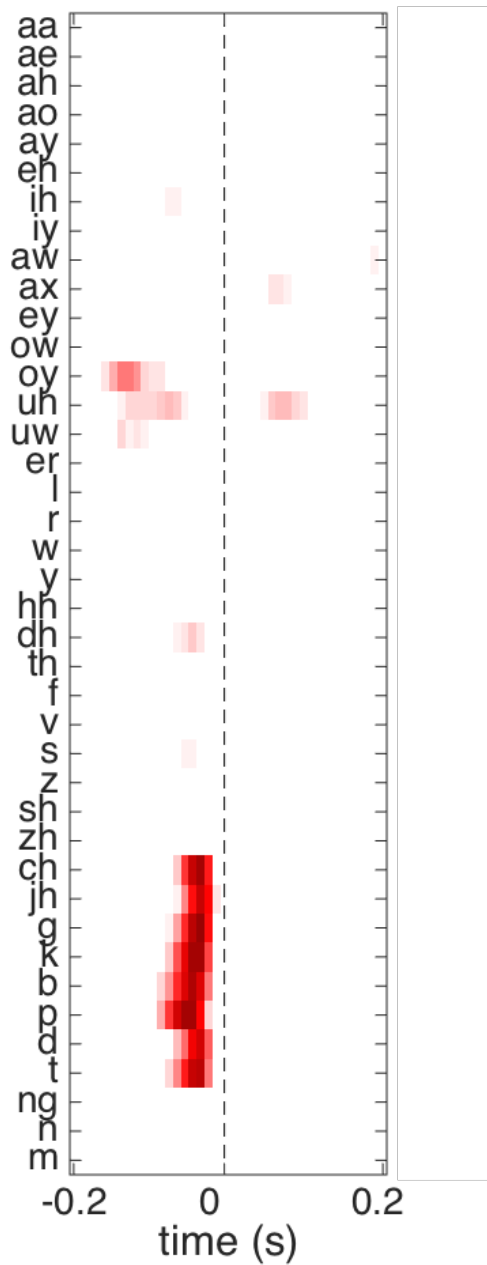
Node 191



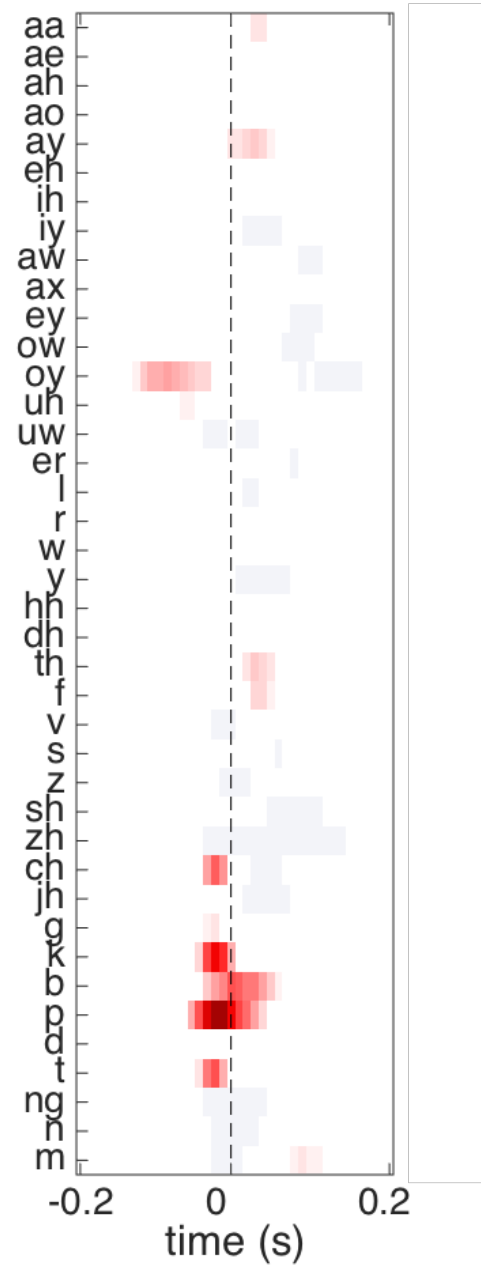
Node 165



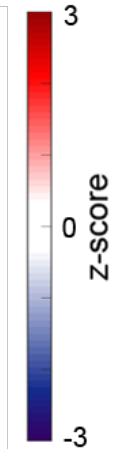
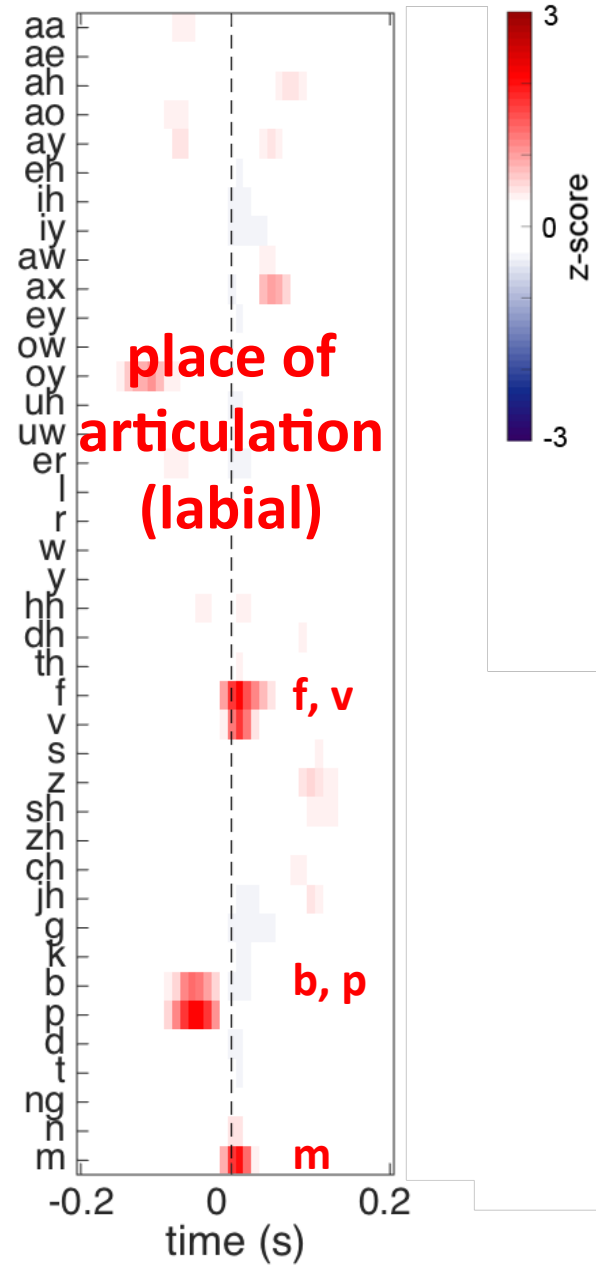
Node 16



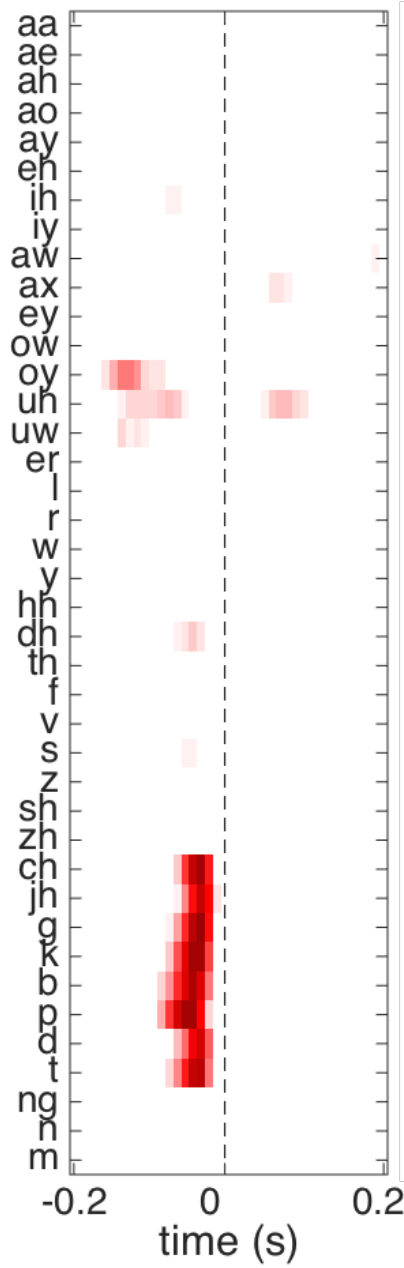
Node 191



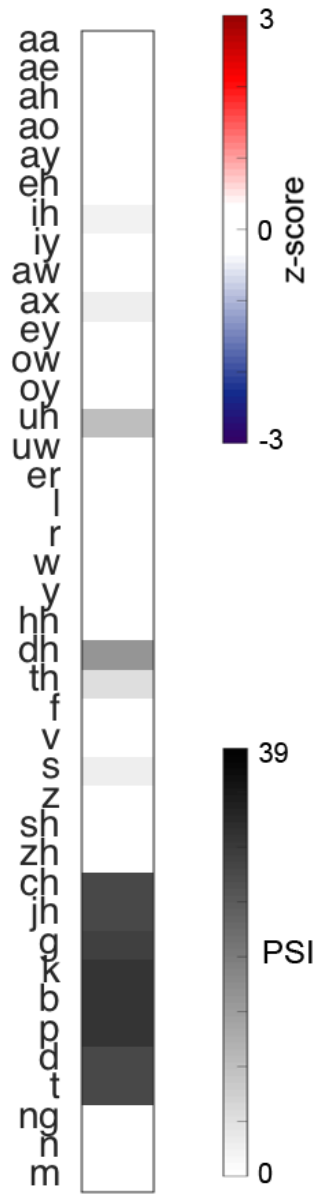
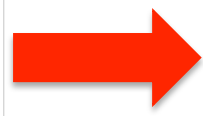
Node 165



Node 16

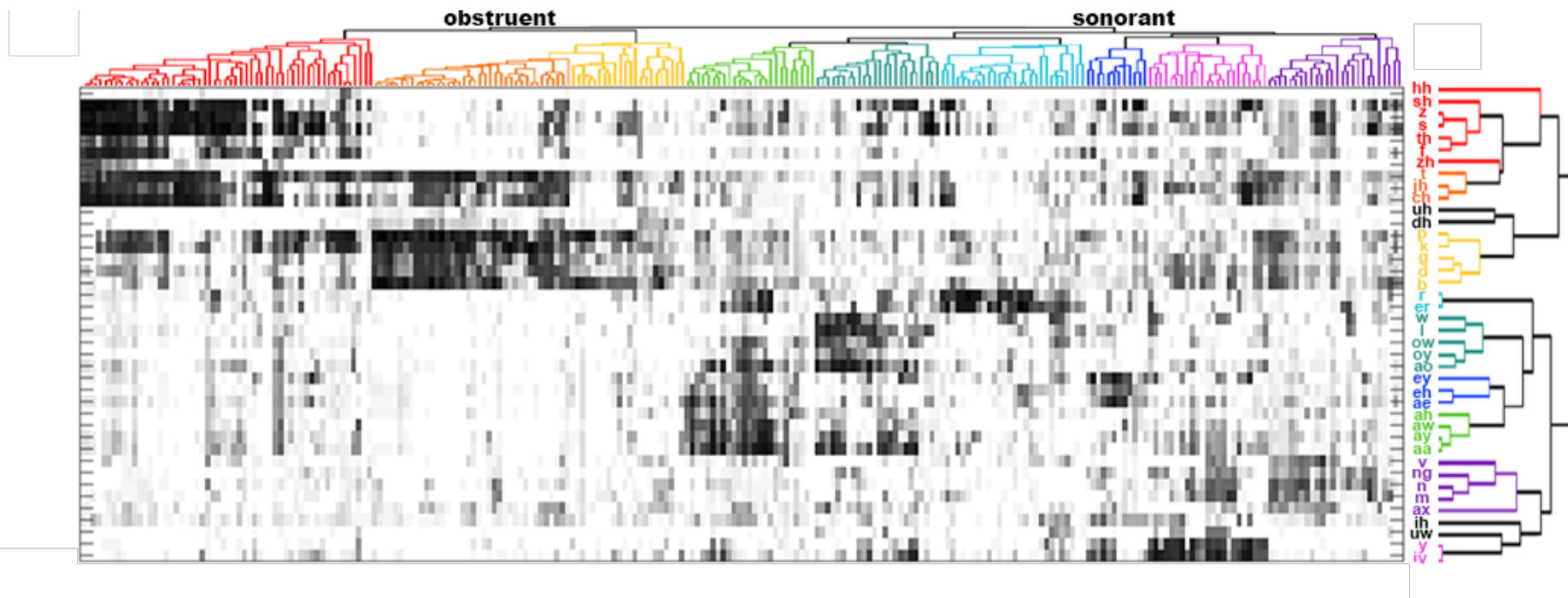


**Phoneme
Selectivity
Index
(PSI)**

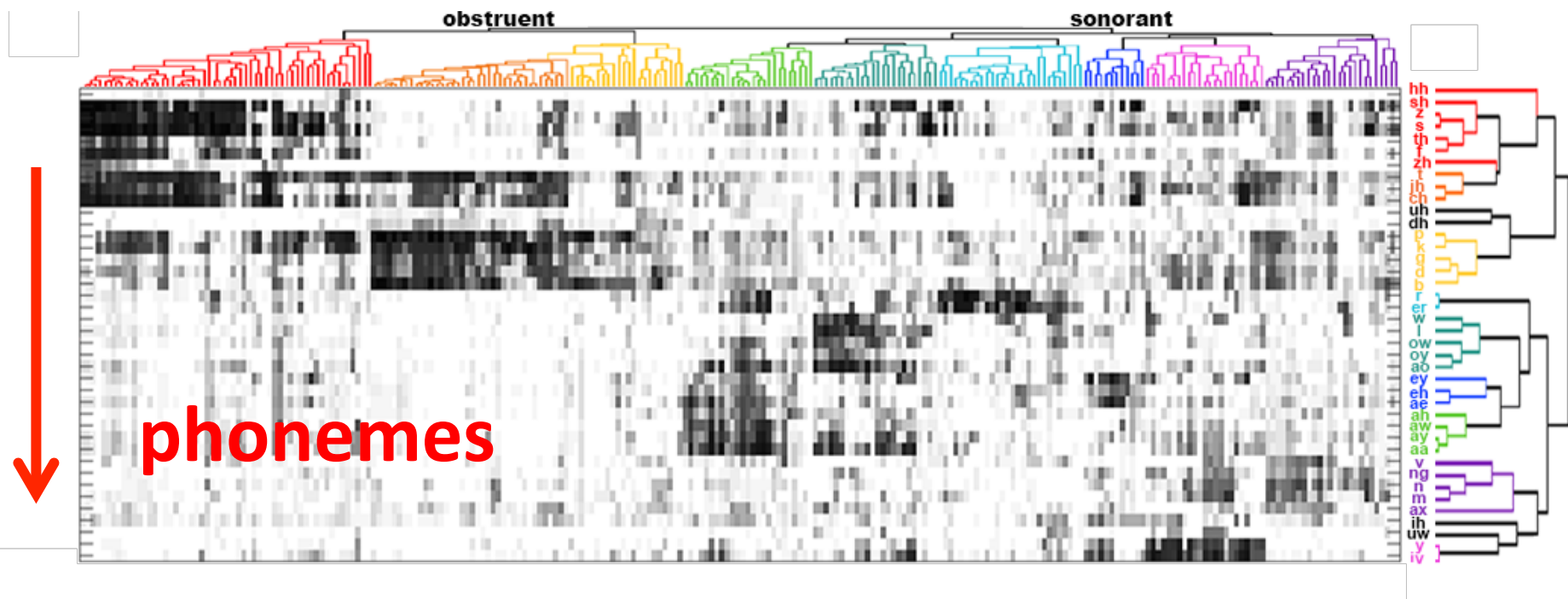


Hidden Layer 1

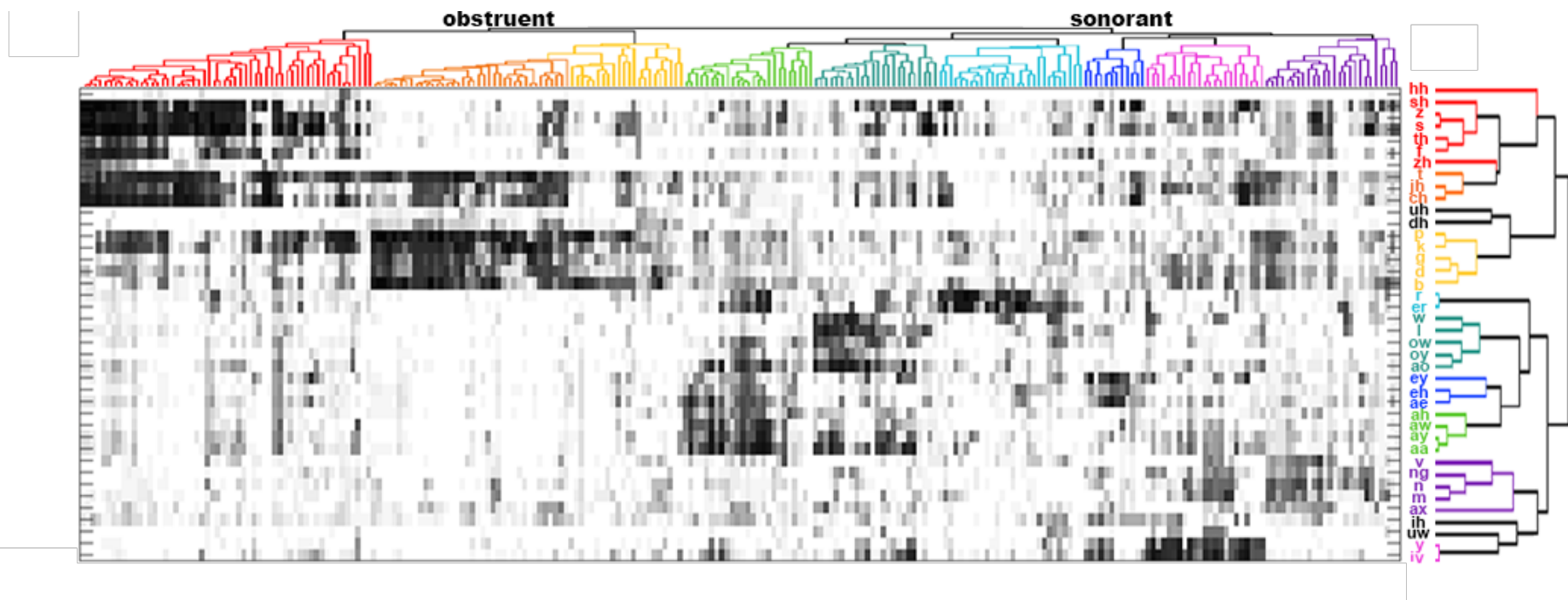
nodes



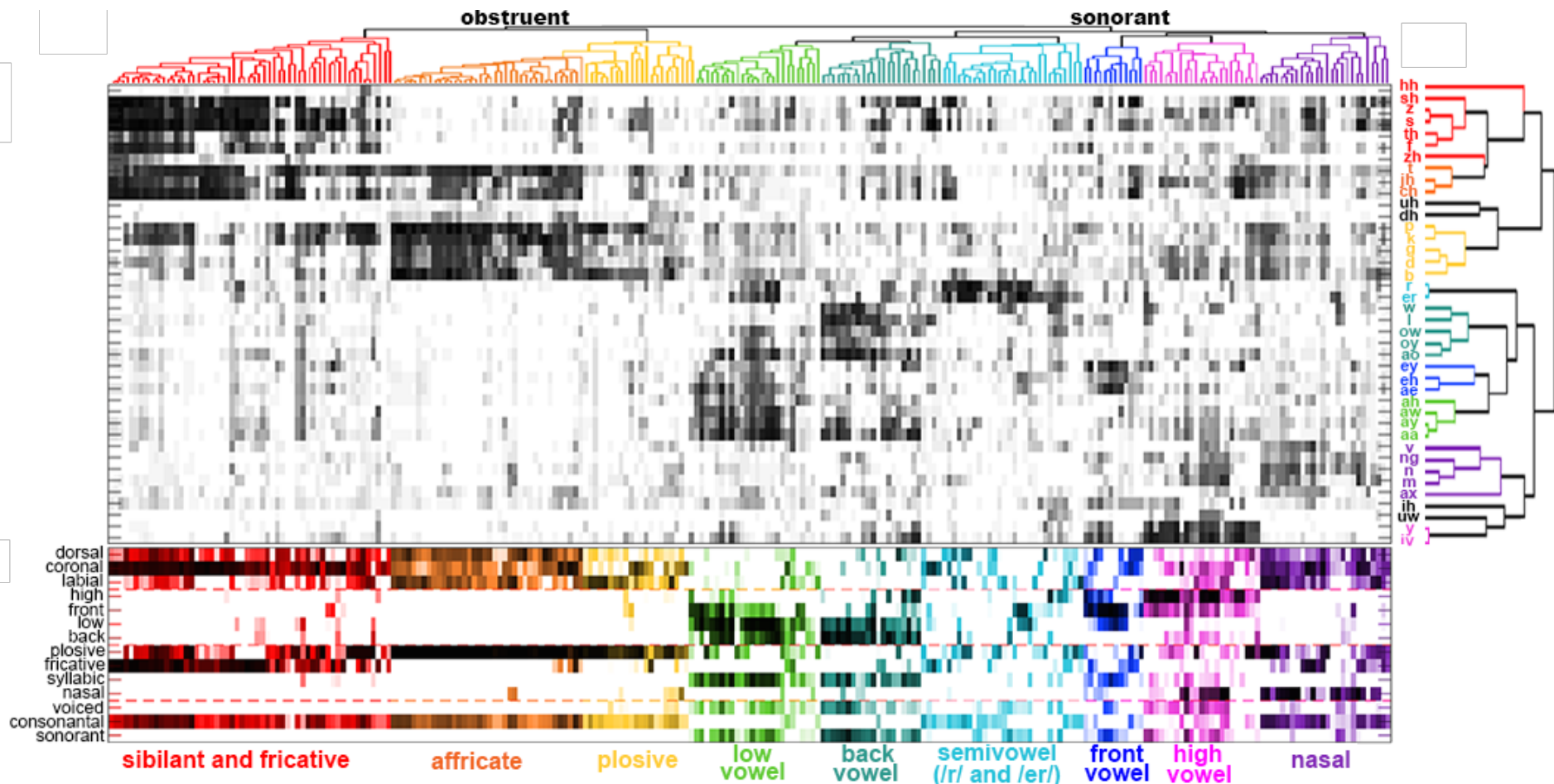
Hidden Layer 1



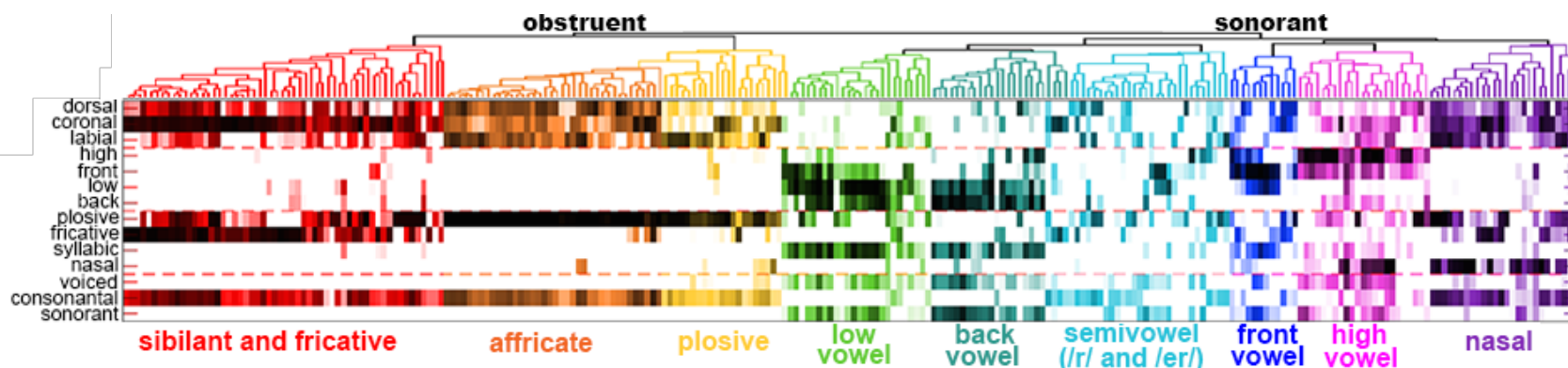
Hidden Layer 1



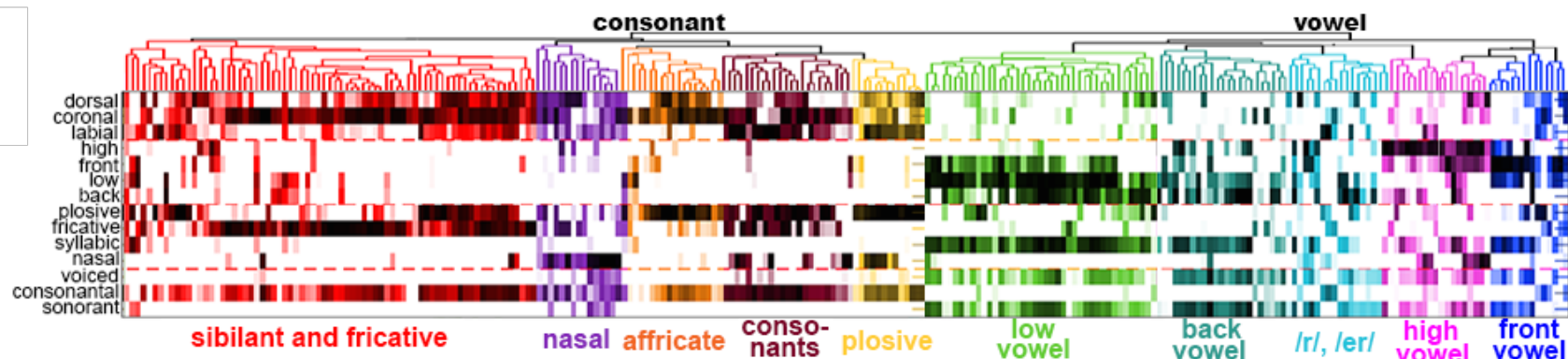
Hidden Layer 1



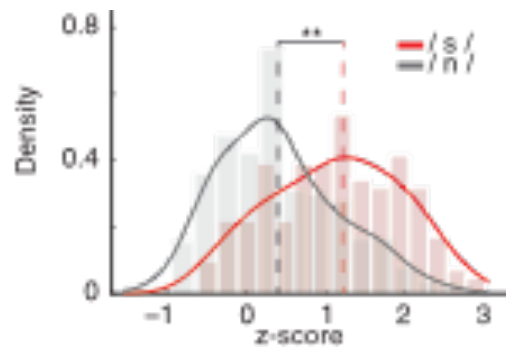
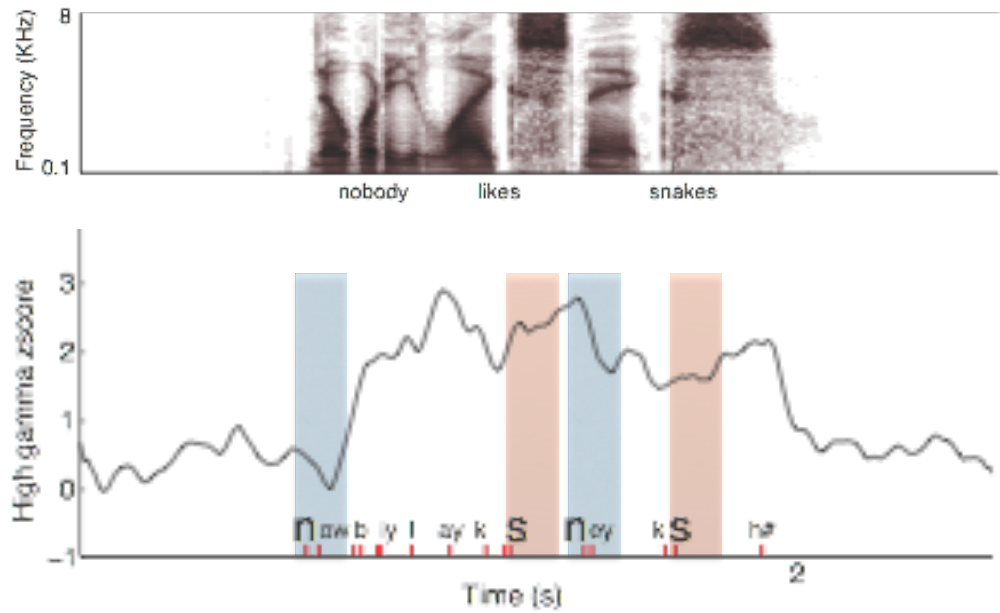
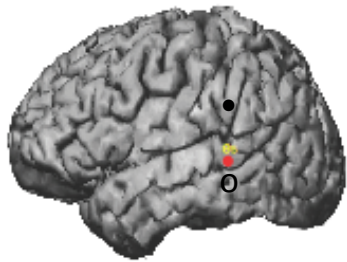
Hidden Layer 1



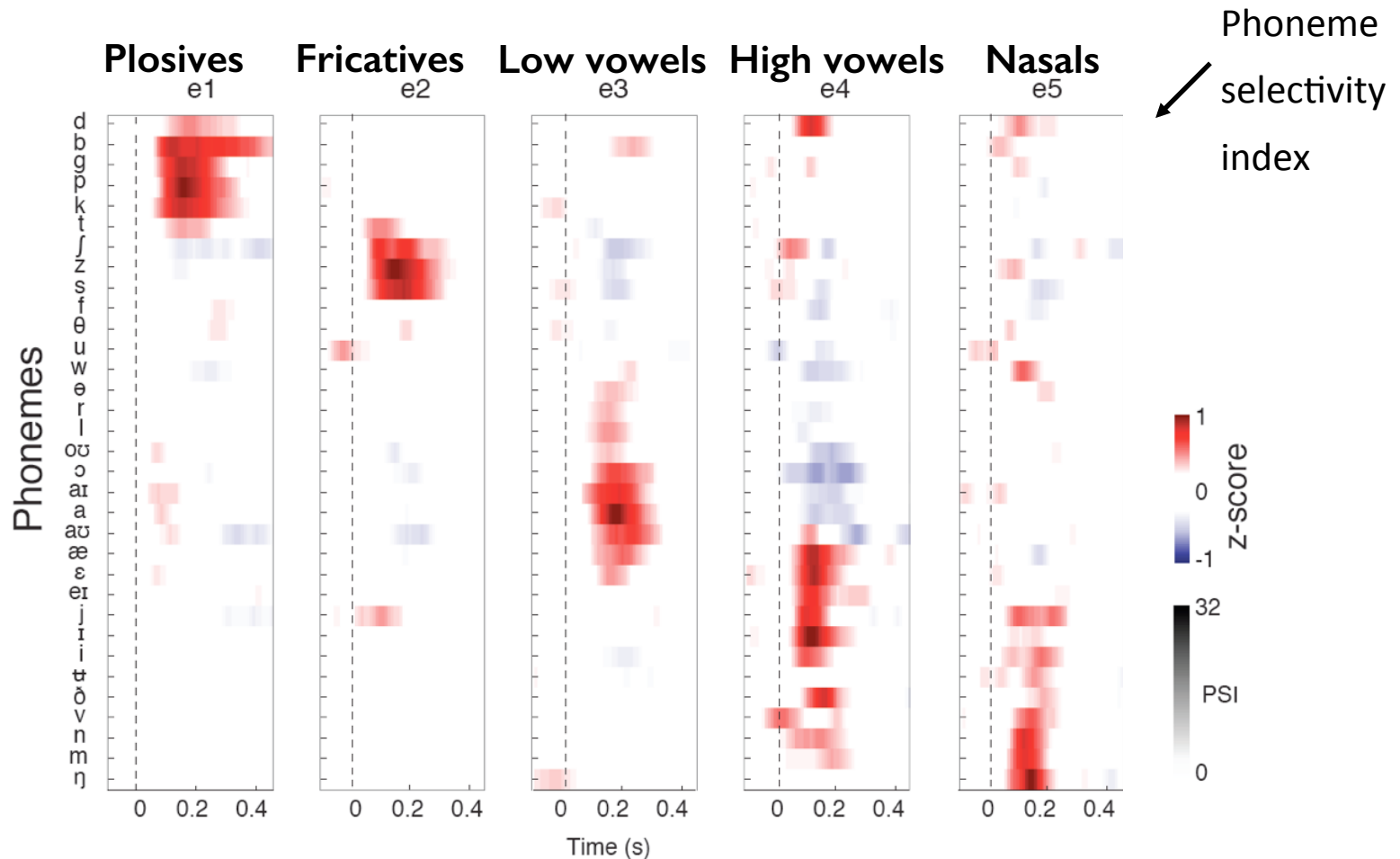
Hidden Layer 5



Neural responses to speech in human superior temporal gyrus (STG)

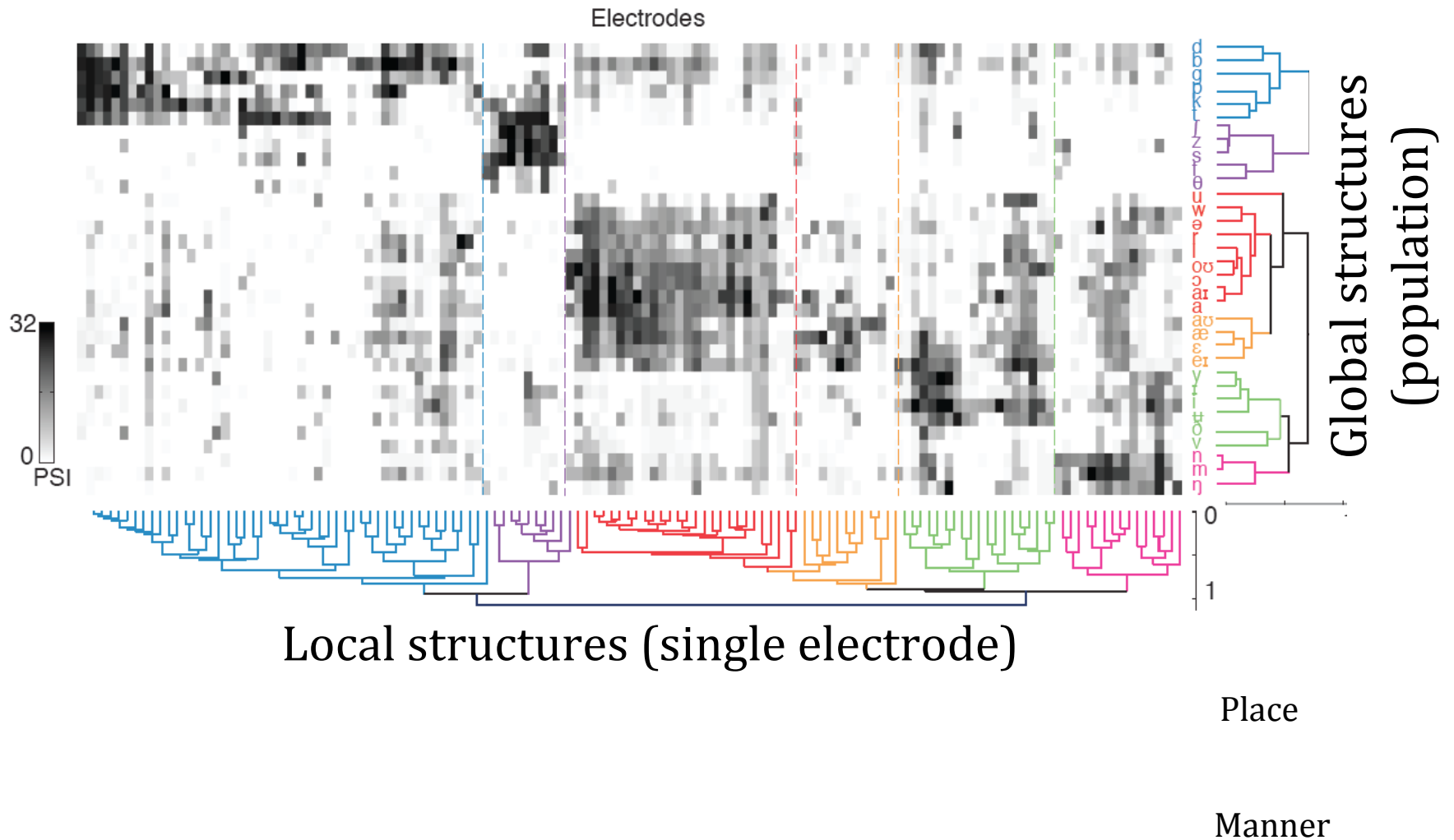


Examples of average phoneme responses in STG



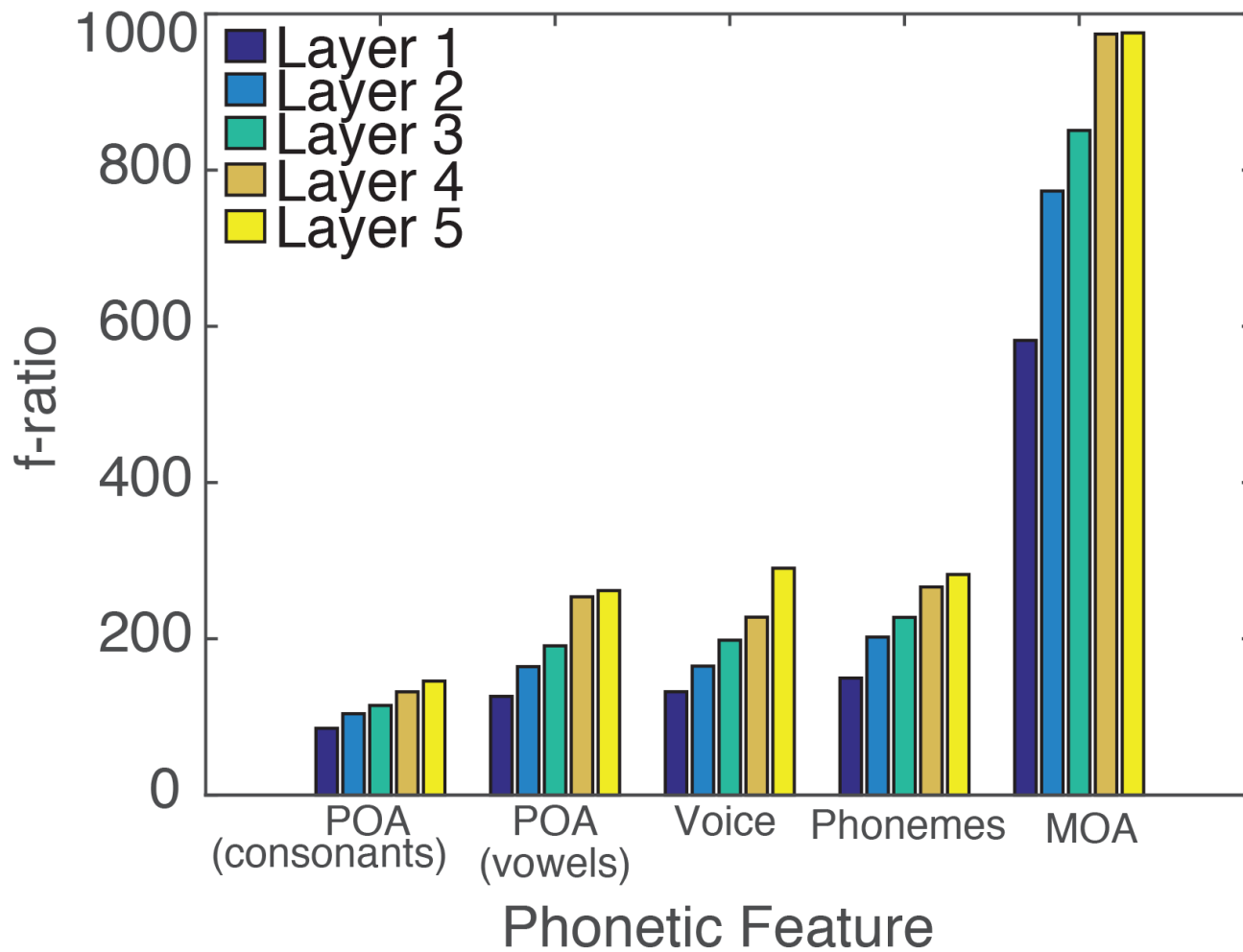
Diversity of responses: Strong preference at various STG sites to specific phoneme **groups** with shared attributes

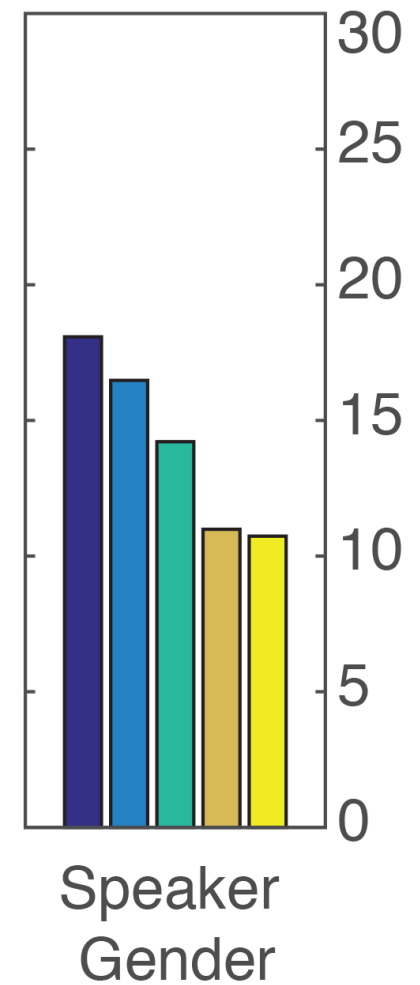
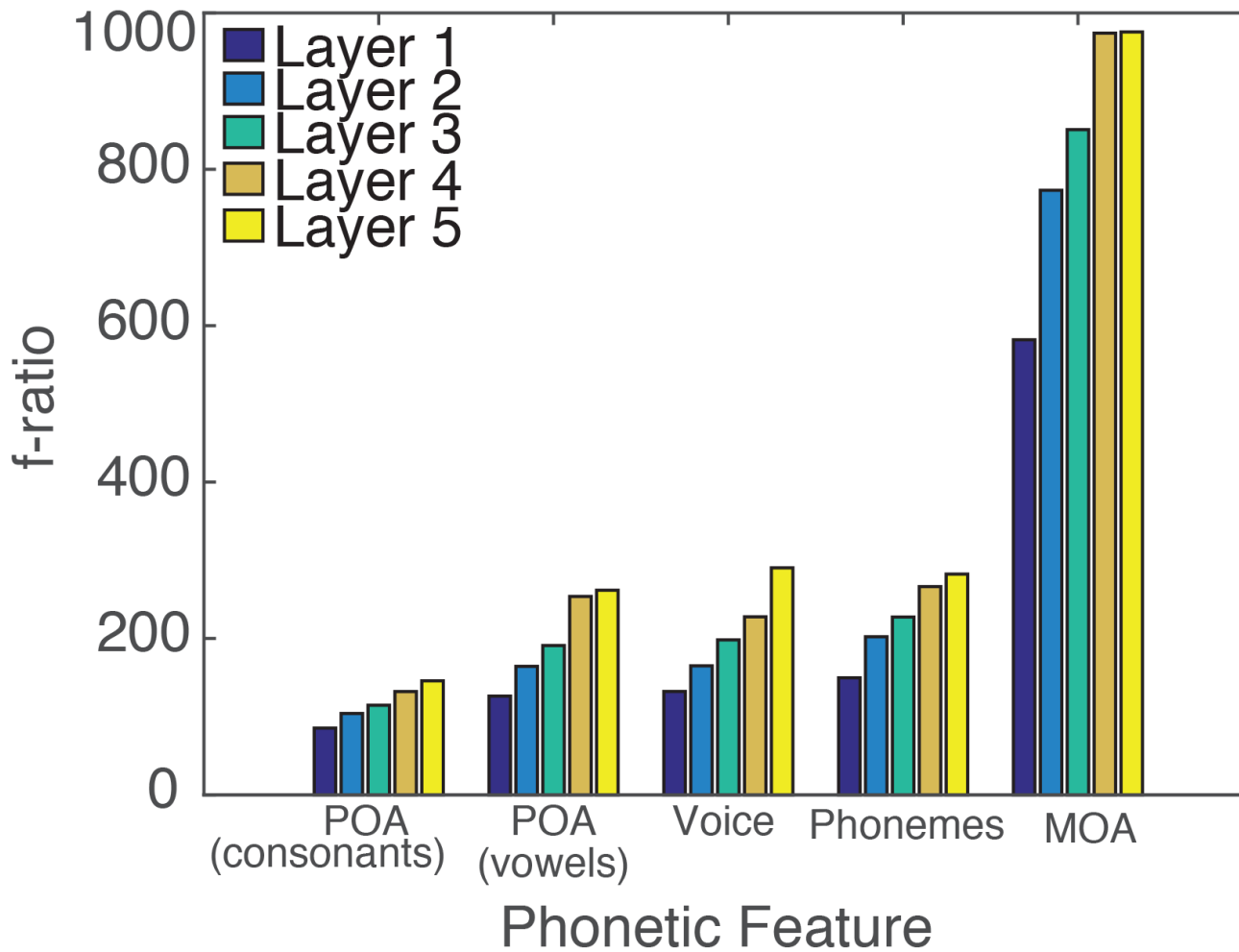
Clustering the PSI vectors



Summary of findings

1. Single nodes and populations of nodes in a layer are selective to phonetic features
2. Phonetic feature encoding becomes more explicit in deeper layers





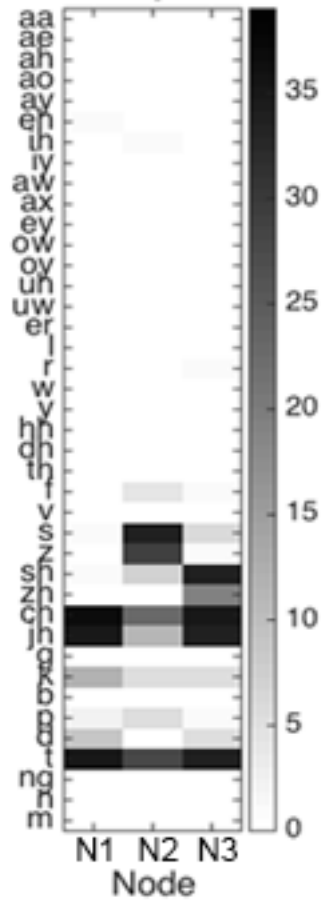
Summary of findings

1. Single nodes and populations of nodes in a layer are selective to phonetic features
2. Node selectivity to phonetic features becomes more explicit in deeper layers
3. Network invariance is learned through explicit representation of sources of variability

phoneme = "t"

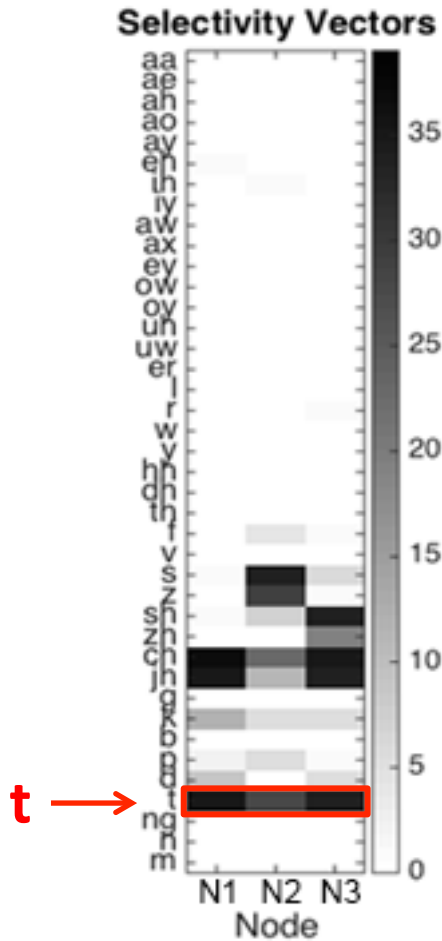
example selectivity for three nodes (N1, N2, N3)

Selectivity Vectors



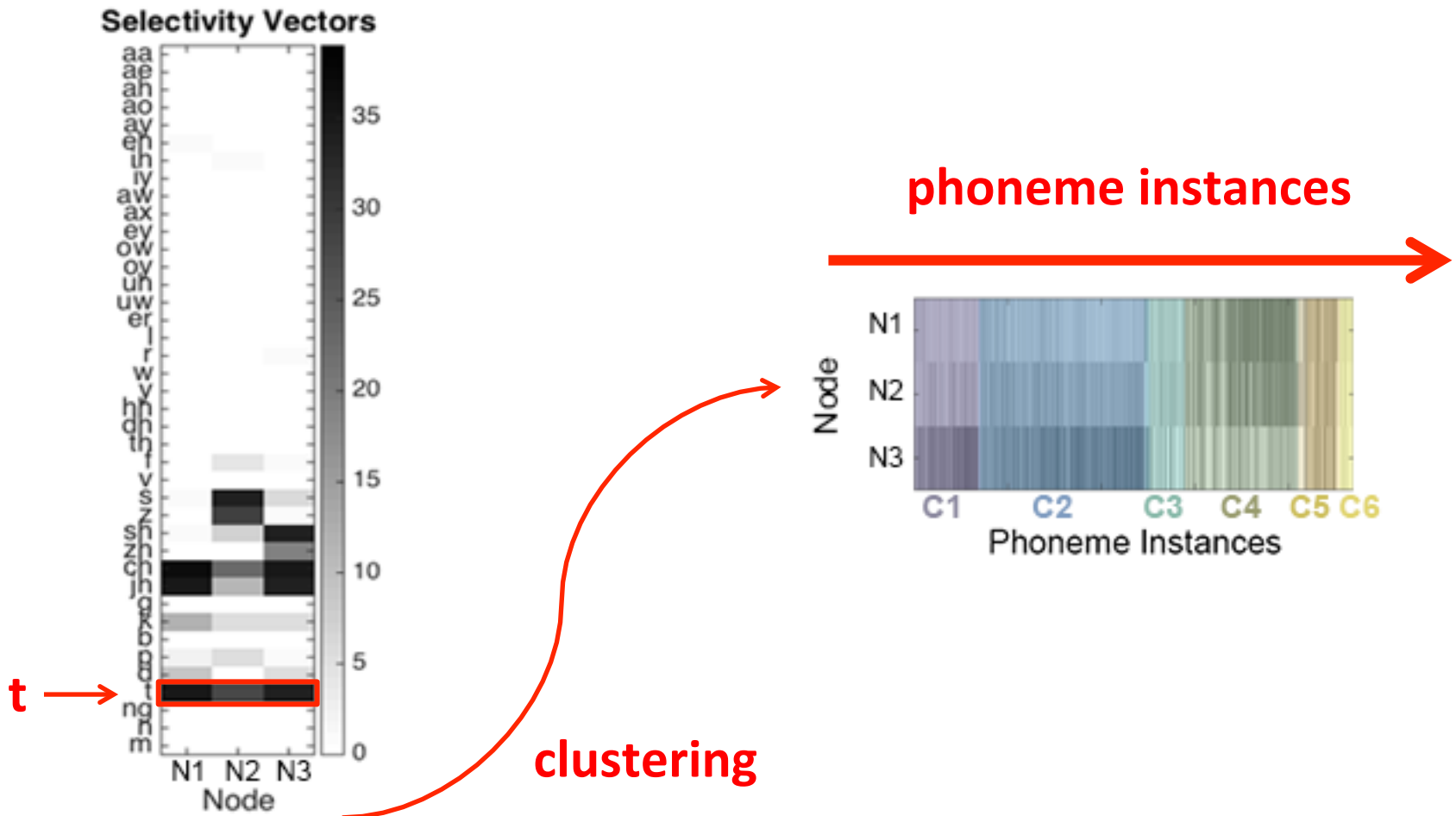
phoneme = "t"

example selectivity for three nodes (N1, N2, N3)



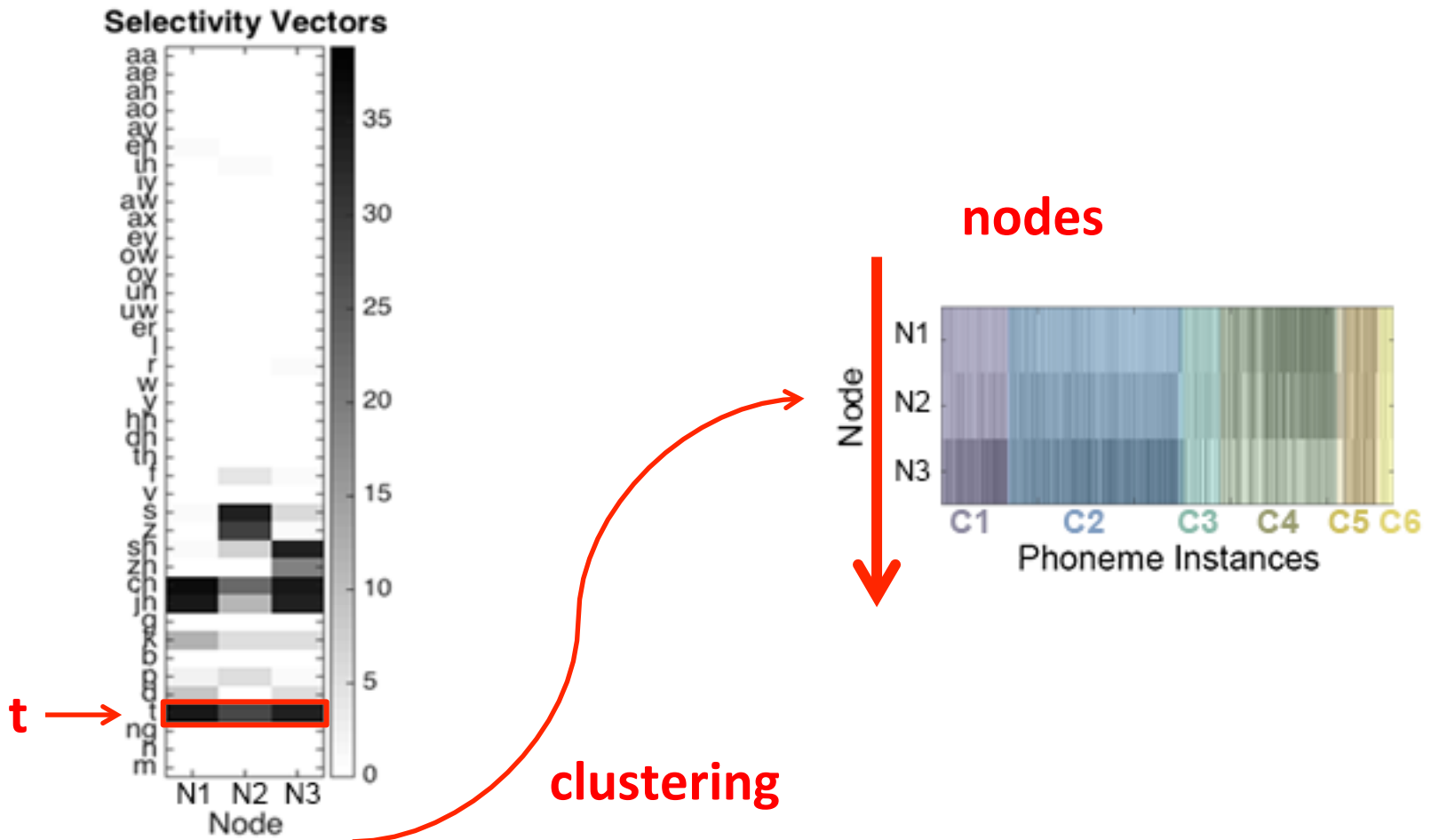
phoneme = "t"

example selectivity for three nodes (N1, N2, N3)

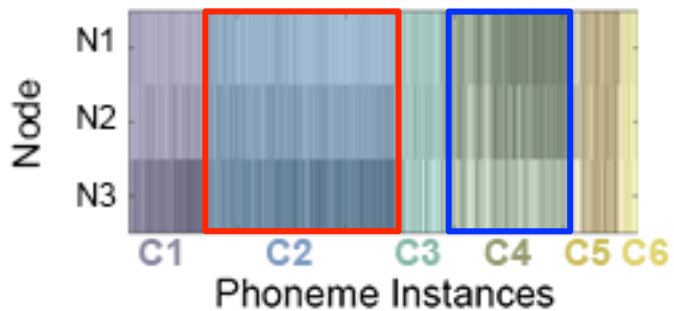
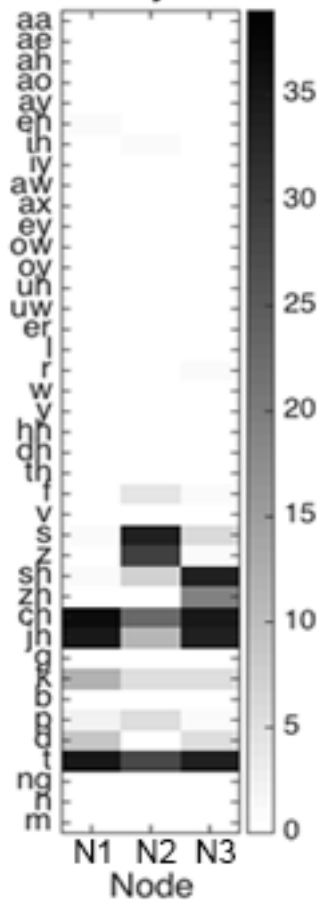


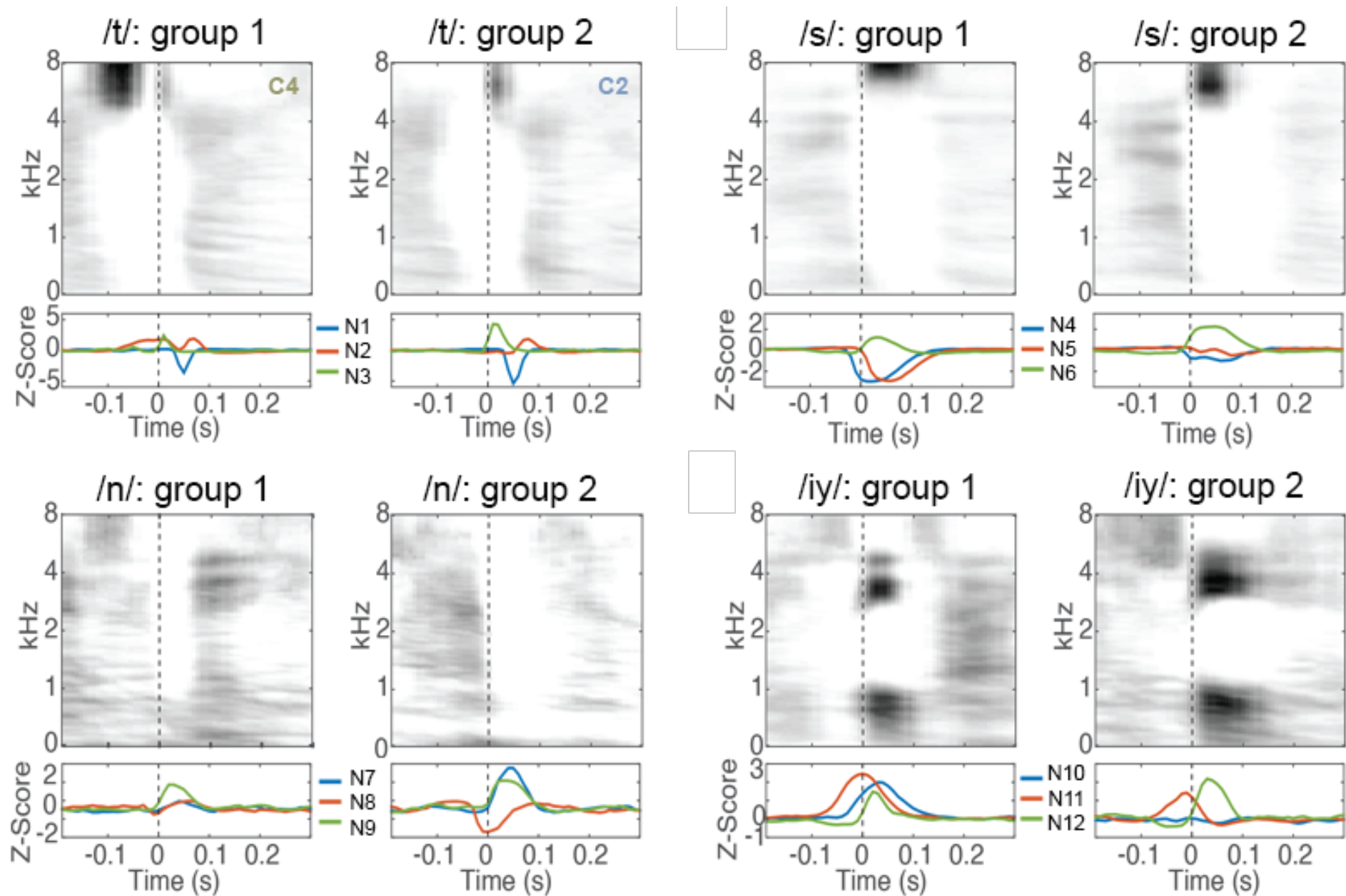
phoneme = "t"

example selectivity for three nodes (N1, N2, N3)



Selectivity Vectors





Summary of findings

1. Single nodes and populations of nodes in a layer are selective to phonetic features
2. Node selectivity to phonetic features becomes more explicit in deeper layers
3. Network invariance is learned through explicit representation of sources of variability

Questions?