The representation of speech in the human and artificial brain

Dr. Odette Scharenborg

Associate Professor & Delft Technology Fellow Delft University of Technology The Netherlands o.e.scharenborg@tudelft.nl



The representation of speech in the human and deep neural networks

Dr. Odette Scharenborg

Associate Professor & Delft Technology Fellow Delft University of Technology The Netherlands o.e.scharenborg@tudelft.nl



# Acknowledgements



- Polina Drozdova
- Roeland van Hout



- Sebastian Tiesmeyer
- Nikki van der Gouw
- Martha Larson
- Najim Dehak
- Junrui Ni
- Mark Hasegawa-Johnson





# **Ultimate goal**

Automatic speech recognition (ASR) for all the world's languages & all types of speech



# Problem

Transcription 1



Transcription 2



Transcription 3

#### signal.

nary code with which the present ls may take various forms, all of e property that the symbol (or 'epresenting each number (or sign differs from the ones represent: er and the next higher number litude) in only one digit (or puls Because this code in its primar; built up from the conventional a sort of reflection process and I rms may in turn be built up freform in signal.

which hnary code with which the present nated in 1s may take various forms, all of s the "refe property that the symbol (or a received presenting each number (or sign

s the "refe property that the symbol (or a receiver epresenting each number (or sigr differs from the ones represent: er and the next higher number ( litude) in only one digit (or puls Because this code in its primar built up from the conventional a sort of reflection process and I rms may in turn be built up fr form in similar fashion, the c , which has as yet no recognized nated in this specification and s the "reflected binary code."

a receiver station, reflected binar



# **Possible solution**

Create a flexible ASR that is trained on language/speech type X and mapped to language/speech type Y



# What we need

- 1. Invariant units of speech which transfer easily and accurately to new languages & types of speech
- 2. ASR system that can flexibly adapt to new languages & types of speech
- 3. ASR system that can *decide* when to create a new sound category
- 4. ASR system that can *create* a new sound category



# What we need

- 1. Invariant units of speech which transfer easily and accurately to new languages & types of speech
- 2. ASR system that can flexibly adapt to new languages & types of speech
- 3. ASR system that can *decide* when to create a new sound category
- 4. ASR system that can *create* a new sound category



# Human listeners

1. Adapt to all types of pronunciations/speakers

2. Create new sound categories when learning a new language



# **Comparing humans and machines**

Different hardware:

Same task: recognition of words from speech

⇒ Comparing humans and machines [Scharenborg, 2007]:

- provide insights into human speech processing (computational modelling)
- improve ASR technology (MFCCs, PLPs, templatebased ASR)



# **Ultimate question**

 What is the optimal representation of speech for a DNN-based ASR to be able to map one language/type of speech to another?



# This talk

- How do humans adapt? → Perceptual learning in humans
- Perceptual learning in DNNs
- Perceptual learning in DNNs over time
- Representation of speech in DNNs



# **Perceptual learning**

"[...] relatively long-lasting changes to an organism's perceptual system that improve its ability to respond to its environment and are caused by this environment" [Goldstone, 1998, p. 586]

# **Perceptual learning**

"[...] relatively **long-lasting changes** to an organism's perceptual system that improve its ability to **respond to its environment** and are caused by this environment" [Goldstone, 1998, p. 586]





UDEITT [Drozdova, van Hout & Scharenborg, 2015, 2016; Norris, McQueen & Cutler, 2003]<sup>15</sup>

# **General procedure**



# Phonetic categorisation results



**TUDelft** [Scharenborg, Mitterer & McQueen, Interspeech, 2011]

# Lexically-guided perceptual learning

- Causes a temporary change in phonetic category boundaries
- Needs lexical or phonotactic knowledge
- Generalises to words that have not been presented earlier
- Speaker dependent
- Fast: only 10-15 ambiguous items needed
- Long-lasting effect

**JDelft** [Clarke-Davidson et al., 2008; Eisner & McQueen, 2005, 2006; Kraljic & Samuel, 2005, 2007; McQueen et al., 2006; Scharenborg & Janse, 2013; Drozdova, van Hout, & Scharenborg, 2016]

# **Perceptual learning in DNNs**



... but can they flexibly adapt to different speech types like human listeners can?

**TUDelft** [Scharenborg, Tiesmeyer, Hasegawa-Johnson, Dehak, 2018]

If so, **visualize** the mechanism by which the DNNs perform this adaptation  $\rightarrow$ 

intermediate representations that might have correlates in human perceptual adaptation





# Lexical retuning

- 1. Train baseline DNN on Dutch read speech ~64h:
- 2. Retrain the baseline model using the acoustic stimuli of the human experiment [Scharenborg & Janse, 2013]:

Baseline	AmbR	AmbL
<ul> <li>120 natural words</li> <li>40 [J]-final words</li> <li>40 [I]-final words</li> </ul>	<ul> <li>120 natural words</li> <li>40 [I]-final words</li> <li>40 [J]-final words: <ul> <li>[J] replaced by</li> <li>[I/J]</li> </ul> </li> </ul>	<ul> <li>120 natural words</li> <li>40 [J]-final words</li> <li>40 [I]-final words: <ul> <li>[I] replaced by</li> <li>[I/J]</li> </ul> </li> </ul>



More [J] responses when exposed to wekke[l/J] → Learn to interpret [l/J] as [J]

More [I] responses when exposed to appe[l/」]

 $\rightarrow$  Learn to interpret [l/J] as [l]

# **Expectations**



Differences particularly between AmbL and AmbR models



# Results

- Lexical retuning results
- Inter-segment distances
- Visualisations of the clusters in the hidden layers



# Lexical retuning results

#### Test set = train set

	Sound	Sound(s) classified (%)	
	Baseline, retrained model		
	[1]	<b>l(97.6)</b> , m(3.4)	
	[J]	<b>J(95.0)</b> , sil(5.0)	
	[L/I]	<b>l(46.9)</b> , sil(23.5), ə(19.8), <b>ɹ(8.6)</b> , εi (1.2)	
	AmbL model		
	[1]	o(78.0), ɔ(10.4), sil(2.4), e(2.4), ɛi (2.4), øː(2.4)	
	[4]	<b>J(87.5)</b> , e(7.5), ∧u(2.5), εi(2.5)	
$\rightarrow$	[L/l]	<mark>l(81.5)</mark> , ə(12.3), ə(2.5), εi(2.5), t(1.2)	
	AmbR model		
	[1]	<b>I(97.6)</b> , sil(3.4)	
	[L]	<b>J(72.5)</b> , ə(15.0), sil(10.0), t(2.5)	
<b>G</b>	[./ا]	<b>J(88.9)</b> , sil(6.2), ə(4.9)	
<b>ÍU</b> Delft			

# Inter-segment distances for each layer

For higher layers, distance

- → [I]-[I/J] decreases for AmbL model
- $\rightarrow$  []-[L] decreases for AmbR model



# Visualisations of the learned clusters in the hidden layers



# PCA visualisations of the 4th hidden layers

#### **Baseline model**



**ŤU**Delft

# PCA visualisations of the 4th hidden layers

#### AmbR model





[L/J]

# PCA visualisations of the 4th hidden layers

AmbL model





# So ...

- DNNs trained with ambiguous sounds show perceptual learning
- Not only at the output layer but also at intermediate layers



# Perceptual learning in DNNs over time

 Human listeners only need 10-15 ambiguous items

• How about DNNs?



# Same experimental set-up

- Retraining in 10 bins of 4 ambiguous items
- Test on unseen data from next bin





# Proportion of [J] responses over time

#### **Baseline model**

ŤI



# Proportion of [J] responses over time



# Proportion of [J] responses over time

#### AmbL model



# So ...

#### DNNs

- Need only a few ambiguous items
- Show a similar step-like function as humans



# Representation of speech in DNNs

• Free reign to the DNN

What speech representations are learned when faced with the large variability in speech?

# **Research question**

Can a generic DNN-based ASR trained to distinguish high-level speech sounds learn the underlying structures used by human listeners to understand speech?



# An example from vision





# Methodology

- Naïve, generic feed-forward deep neural network
   3 hidden layers with 1024 nodes each
- Corpus Spoken Dutch, 64h read speech
- Consonsant/vowel classification task
- Visualize the activations of the speech representations at the hidden layers
  - Using different linguistic labels, to check for clustering



# Linguistic labels

- <u>Vowels/consonants</u> differ in the way they are produced: absence/presence of a constriction in the vocal tract
- <u>Phoneme</u>: the smallest unit that changes the meaning of a word
   e.g. *ball* and *wall*
- <u>Articulatory feature</u>: acoustic correlate of the articulatory properties of speech sounds
  - *Manner of articulation:* type of constriction (e.g., full closure, narrowing)
  - Place of articulation: location of the constriction (consonants only)
  - *Tongue position:* location of the bunch of the tongue (vowels only)

# **ŤU**Delft

# C/V classification results

- Frame level accuracy: 85.5%
  - Averaged over 5 training runs
  - Consonants: 85.19% / Vowels: 86.69% correct



# Visualisations: V/C classification task



**ŤU**Delft

Layer 2 of the network

Layer 3 of the network

#### Manner of articulation, 3rd layer Consonants Vowels



- •: approximant
- •: nasal

**TU**Delft

a.

- •: fricative
  - •: plosive



- •: diphthong •: short vowel
- •: long vowel



- •: alveolar consonant
- : glottal consonant
- •: palatal consonant

# **ŤU**Delft

- •: bilabial consonant
- •: labiodental consonant •: front vowel
- •: velar consonant

- •: central vowel •: back vowel

# Summary

### DNNs

- Progressive abstraction in subsequent hidden layers
- Capture the structure in speech by clustering the speech signal, without explicit training
- Adapt to nonstandard speech on the basis of only a few labelled examples, showing a step-like function



# Our needs

- Invariant units of speech which transfer easily and accurately to other languages
- → DNN phone categories are flexible
- > DNNs automatically create linguistic categories
- ASR system that can flexibly adapt to new types of speech
- DNNs are able to do that

Next: ASR system that can *decide* when to create a new sound category, and do so

