

Speech systems that emulate human language acquisition

ILCC, University of Edinburgh, Mar. 2024

Herman Kamper

E&E Engineering, Stellenbosch University, South Africa

<http://www.kamperh.com/>





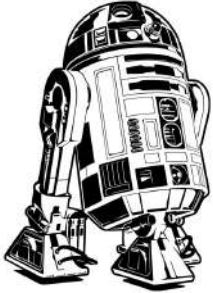
Supervised speech recognition and synthesis

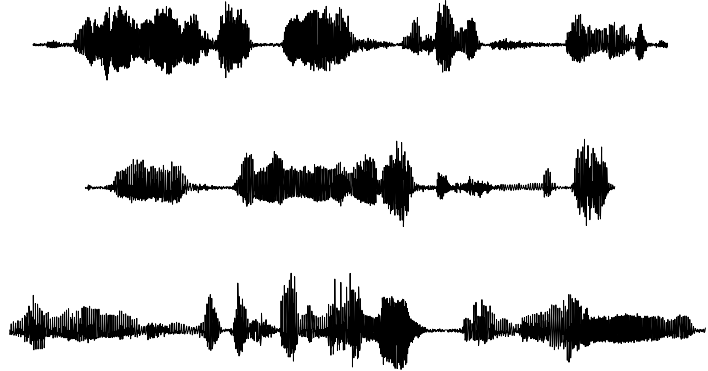
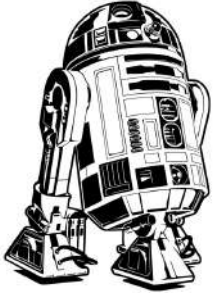


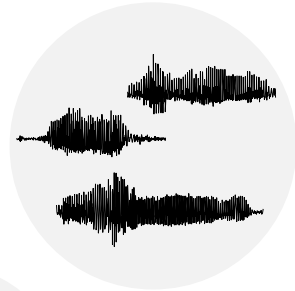
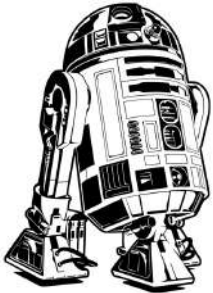
i had to think of some example speech

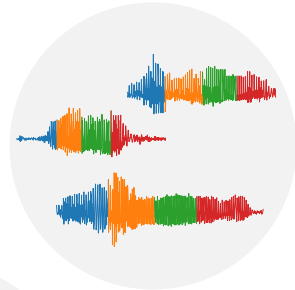
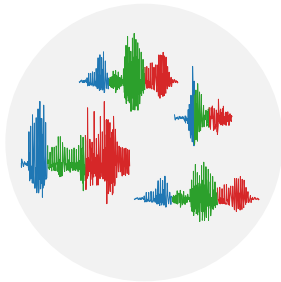
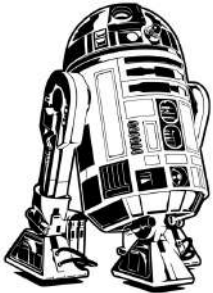


since speech recognition is really cool











Why attempt to emulate language acquisition?



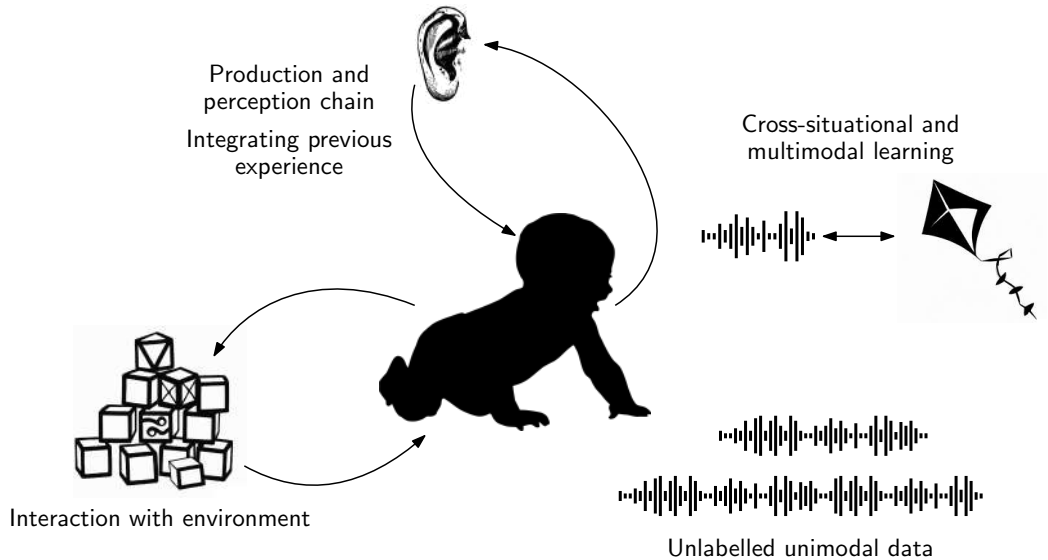
Improvements in speech technology

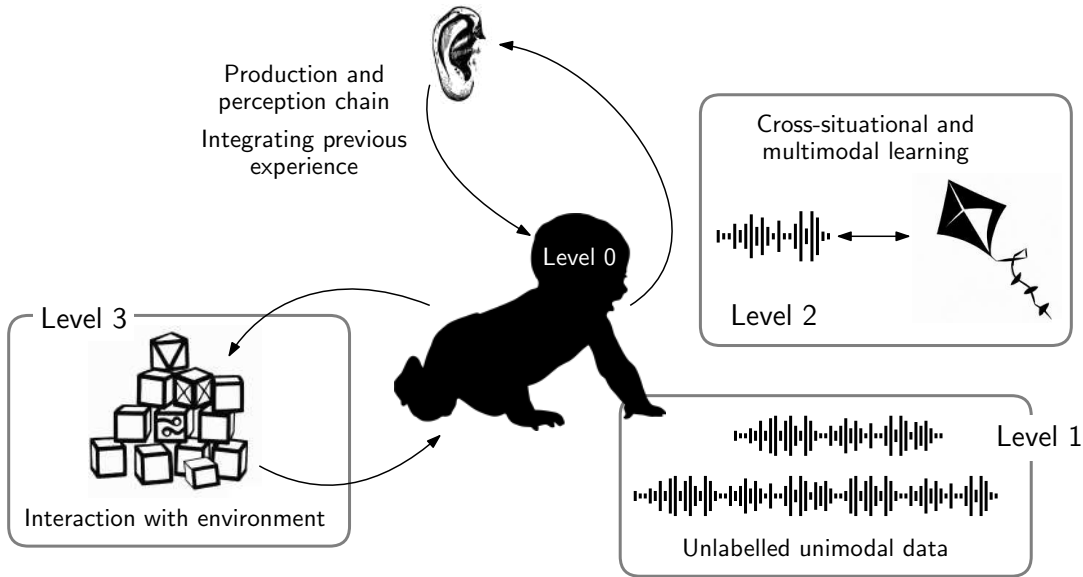


New insights and approaches for machines that learn



New insights into human learning





1. Mutual exclusivity in visually grounded speech models



Leanne
Nortje



Dan
Oneață



Yevgen
Matusevych

Nortje et al., "Visually grounded few-shot word acquisition with fewer shots," in *Interspeech*, 2023.

Nortje et al., "Visually grounded speech models have a mutual exclusivity bias," *Accepted*, 2024.

Children's Use of Mutual Exclusivity to Constrain the Meanings of Words

ELLEN M. MARKMAN

AND

GWYN F. WACHTEL

Stanford University

For children to acquire vocabulary as rapidly as they do, they must be able to eliminate many potential meanings of words. One way children may do this is to assume category terms are mutually exclusive. Thus, if a child already knows a label for an object, a new label for that object should be rejected.

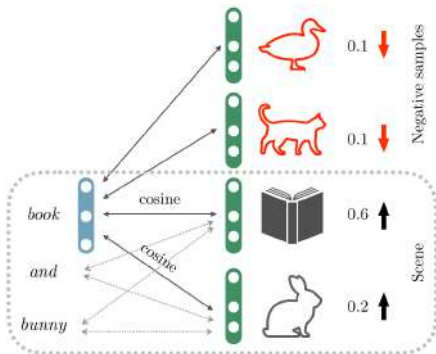


?

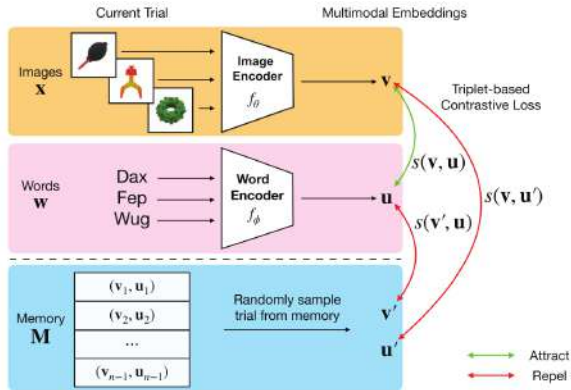


Previous computational studies

Previous computational studies

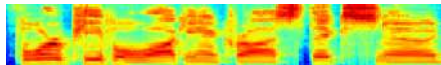
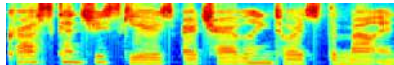
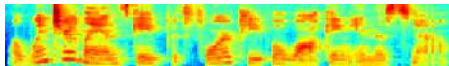


(Gulordava et al., 2020)

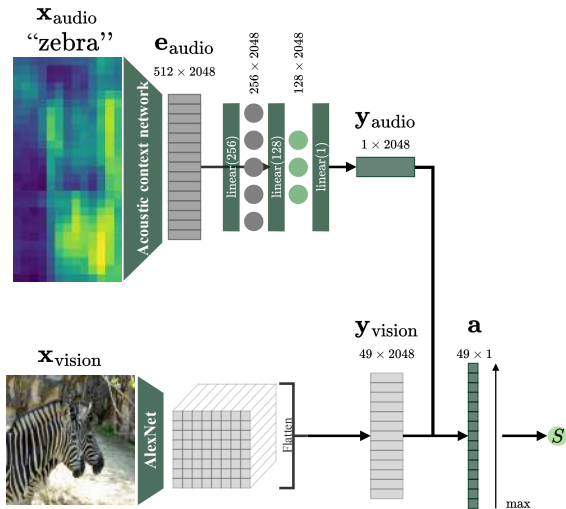


(Vong and Lake, 2022)

Visually grounded speech models



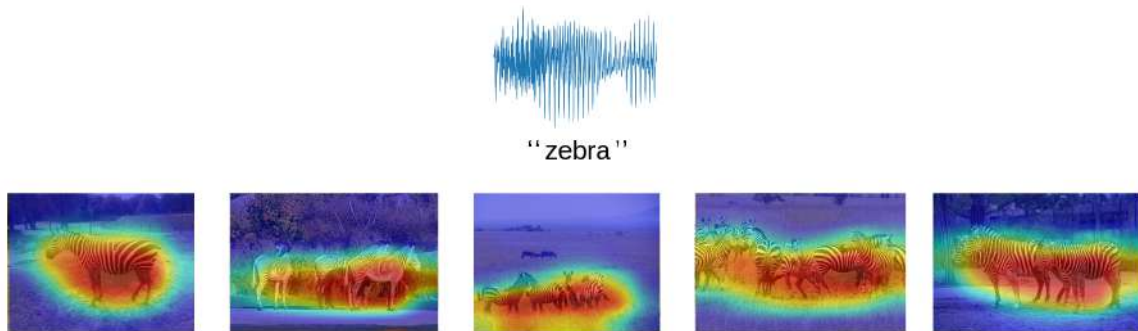
Multimodal attention network (MattNet)



The acoustic context network is initialised with a CPC model trained on Places and LibriSpeech (level 1).

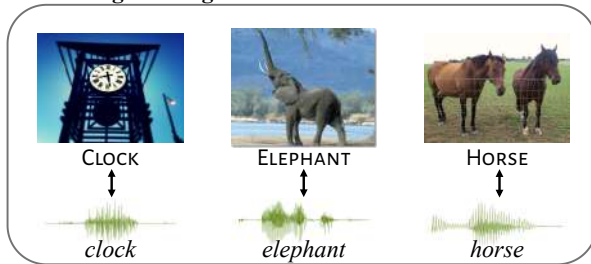
The vision branch is initialised with a self-supervised variant of AlexNet (level 1).

Attention visualisation

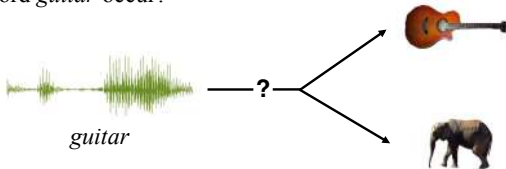


Testing a visually grounded speech model for the ME bias

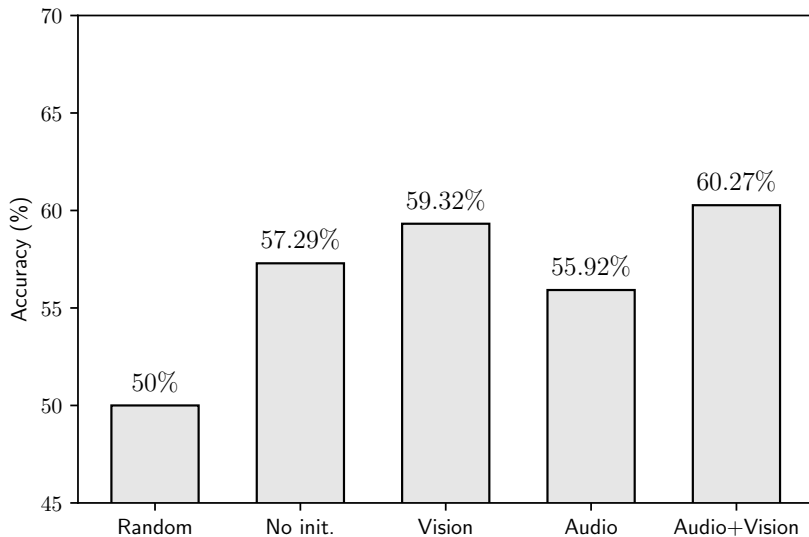
Given during training:

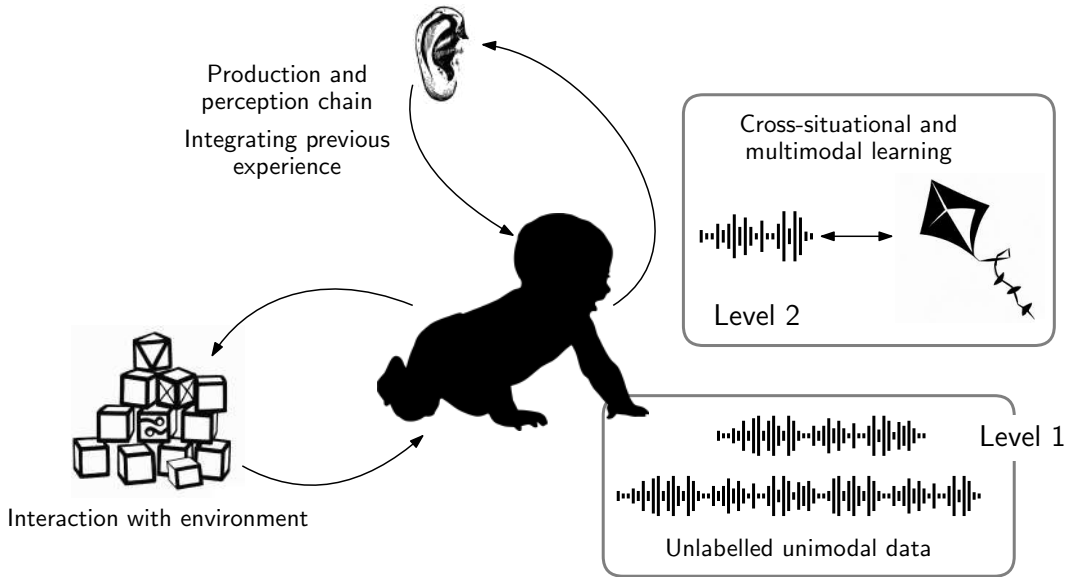


Test-time question: In which picture does the novel spoken keyword *guitar* occur?

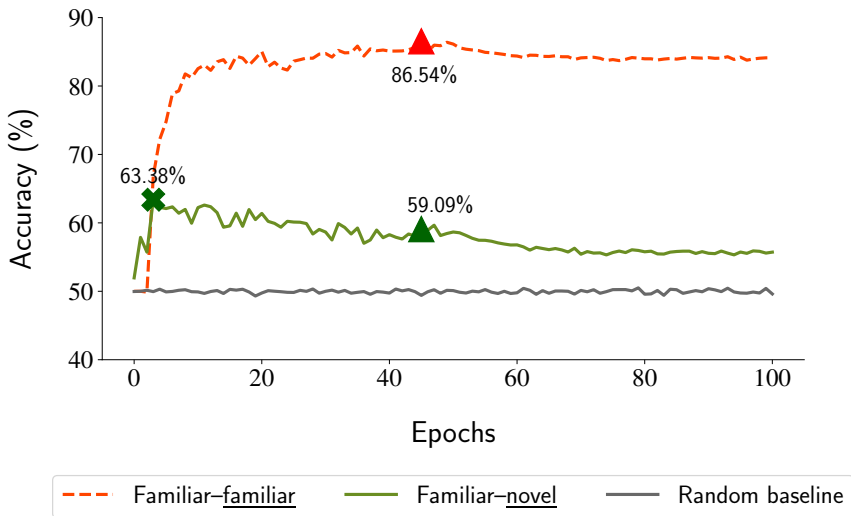


Mutual exclusivity results



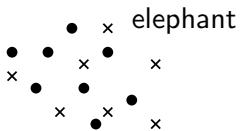


Mutual exclusivity results

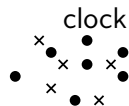
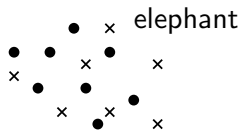


How is the representation space organised?

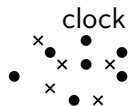
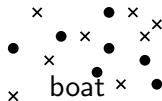
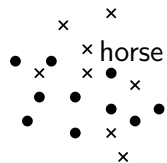
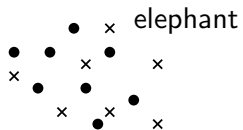
How is the representation space organised?



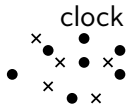
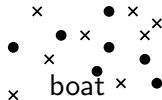
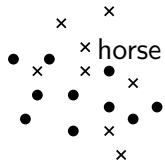
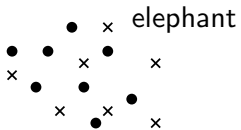
How is the representation space organised?



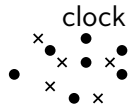
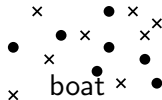
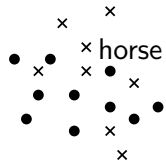
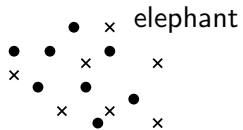
How is the representation space organised?



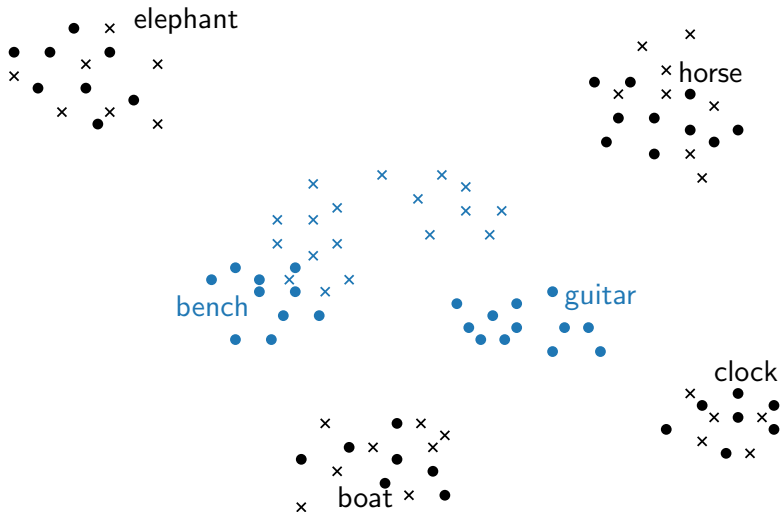
How is the representation space organised?



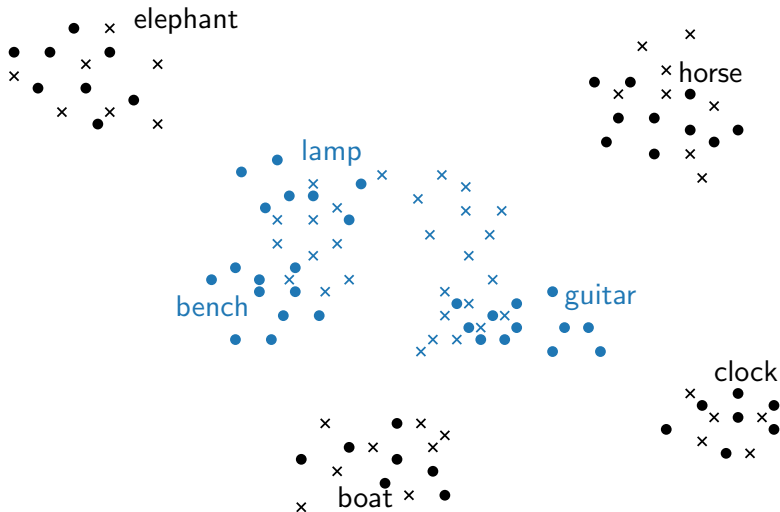
How is the representation space organised?



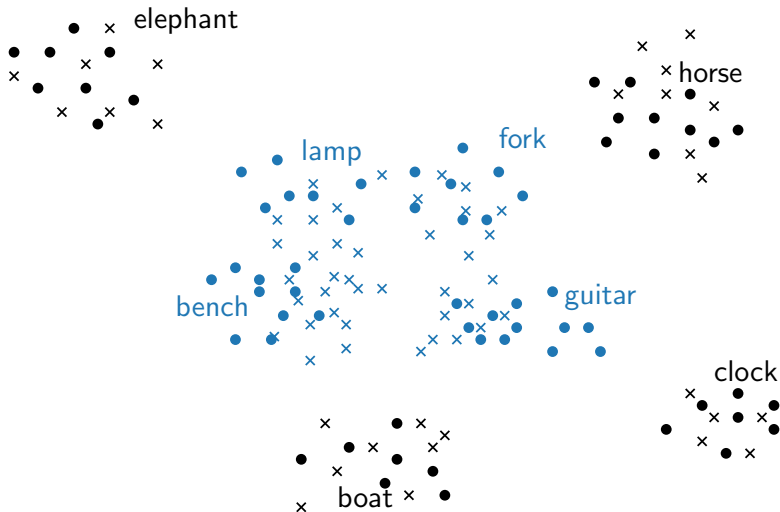
How is the representation space organised?



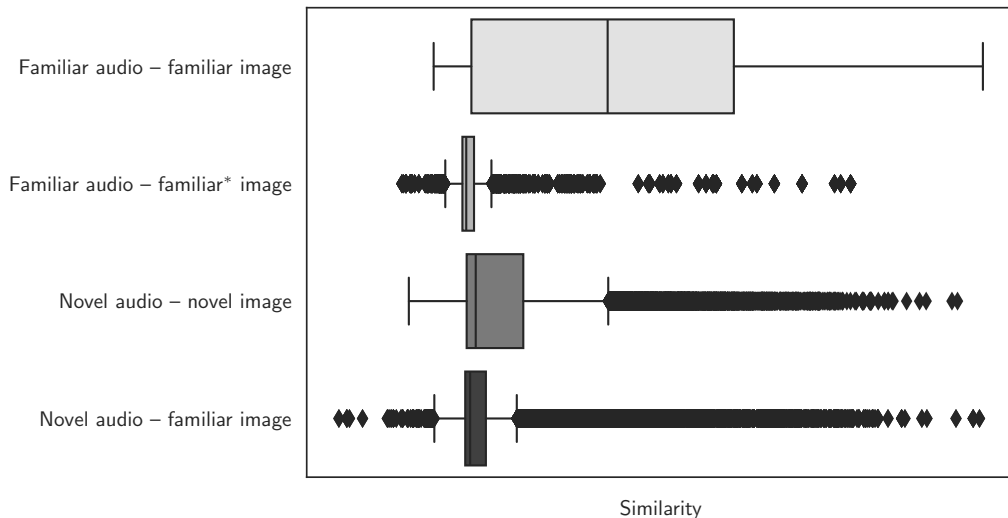
How is the representation space organised?



How is the representation space organised?



How is the representation space organised?

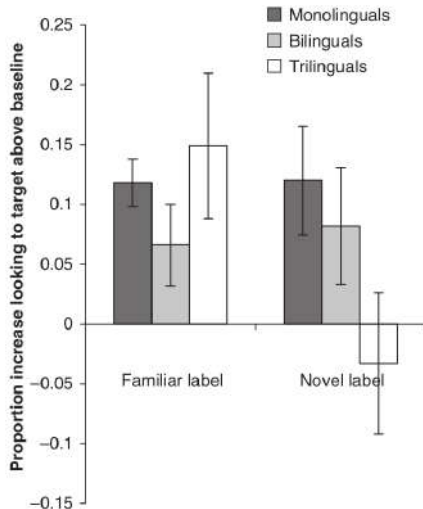


Conclusions and future work

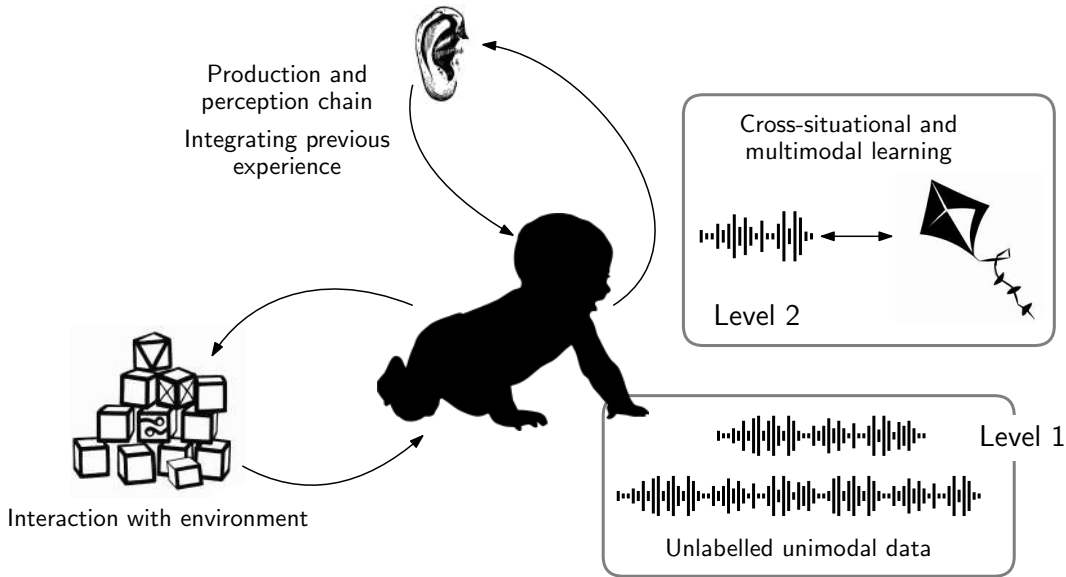
- Showed an example of how we can compare an artificial learner to human infants
- Use speech and not written words
- Adds weight that visually grounded speech model could be studied as a cognitive proxy

Conclusions and future work

- Showed an example of how we can compare an artificial learner to human infants
- Use speech and not written words
- Adds weight that visually grounded speech model could be studied as a cognitive proxy
- Future work: Mutual exclusivity in multilingual models



(Byers-Heinlein and Werker, 2009)



2. Probing self-supervised speech models by listening



Benjamin
van Niekerc



Matthew
Baas



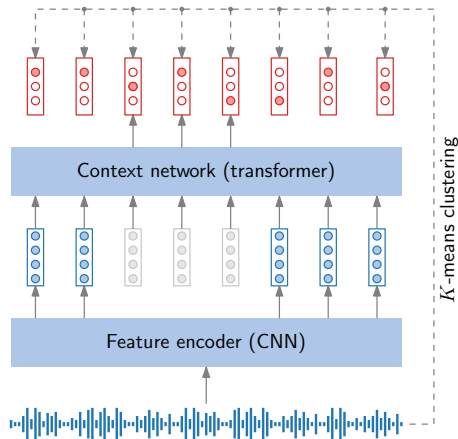
Marc-André
Carbonneau

Baas et al., "Voice conversion with just nearest neighbors," in *Interspeech*, 2023.

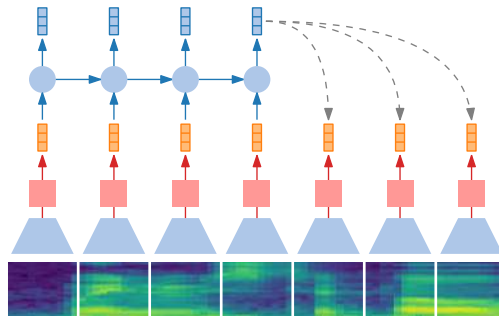
van Niekerc et al., "Rhythm modeling for voice conversion," *IEEE SPL*, 2023.

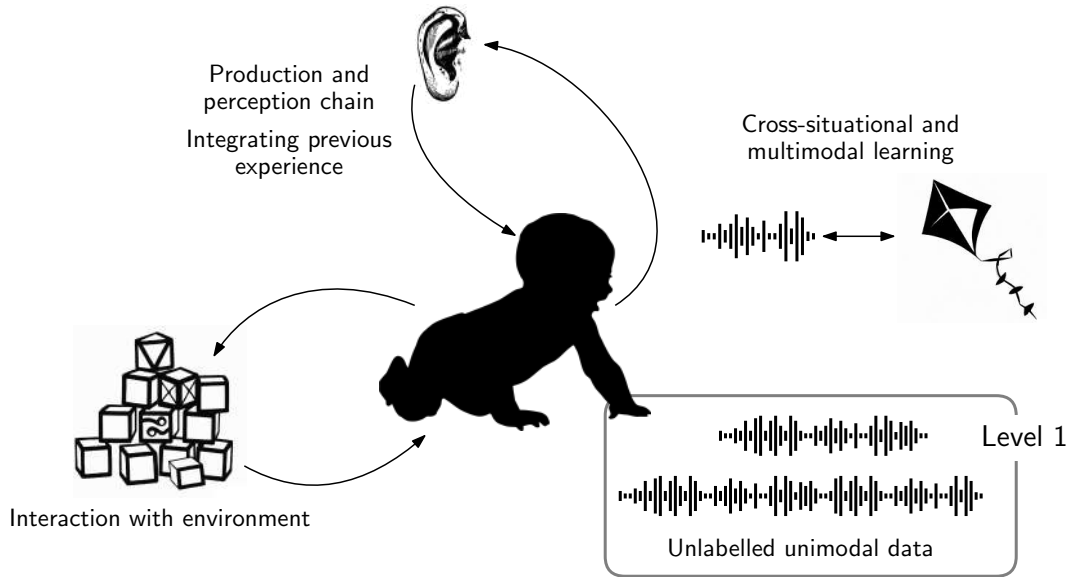
Self-supervised spoken language models

HuBERT / WavLM:



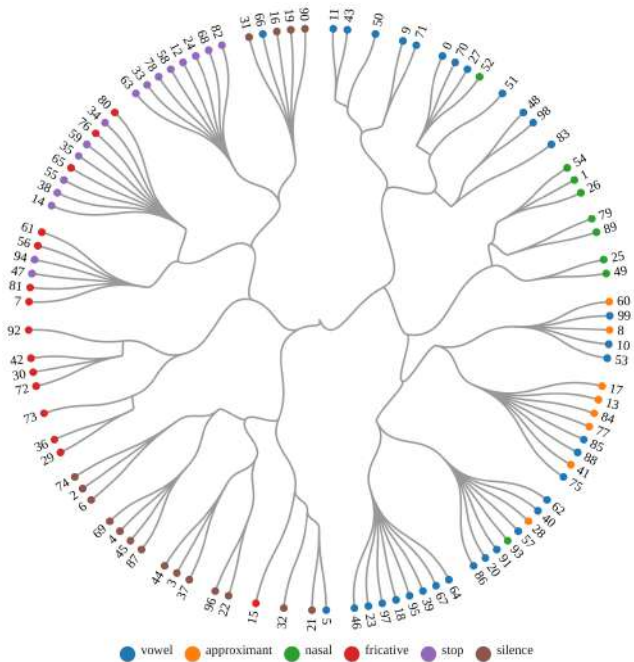
Contrastive predictive coding (CPC):

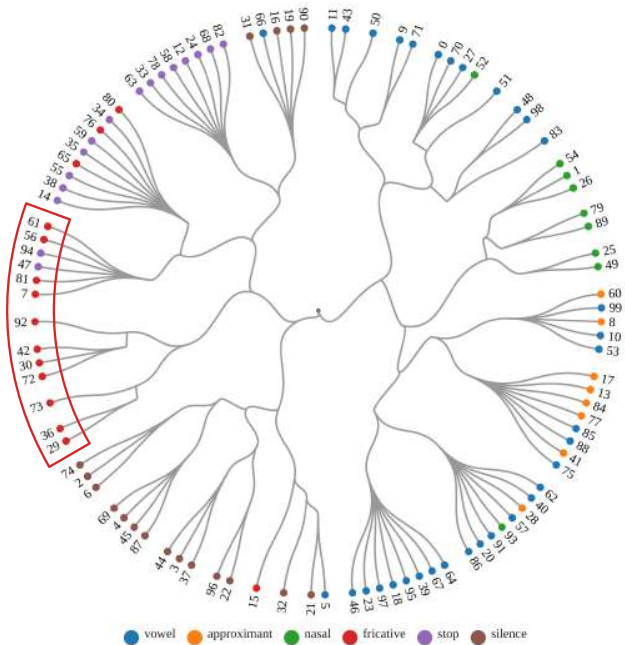




We use voice alteration and voice conversion as a probe to show you how phonetic content and speaker are captured.

(But it's really just an excuse . . .)



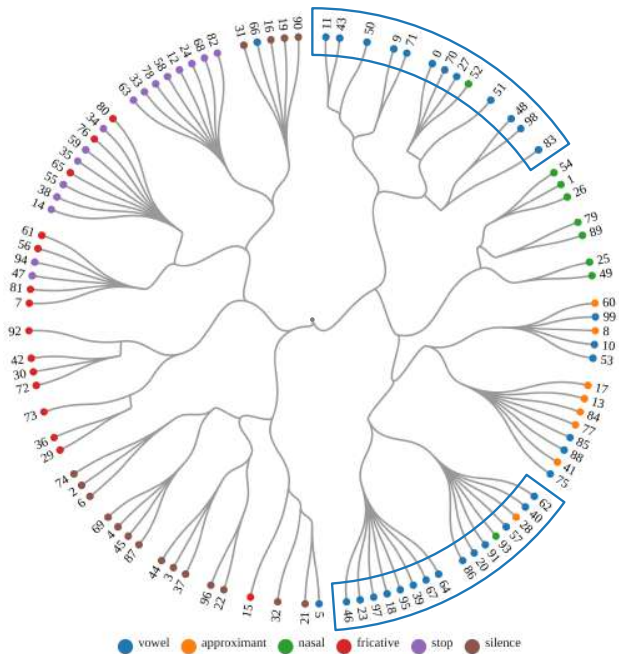


No modification:

Play

Fricatives:

Play

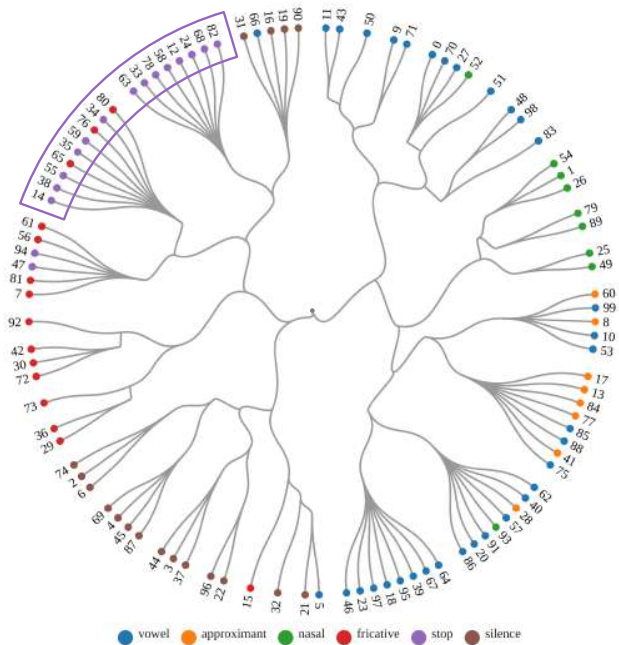


No modification:

Play

Vowels:

Play

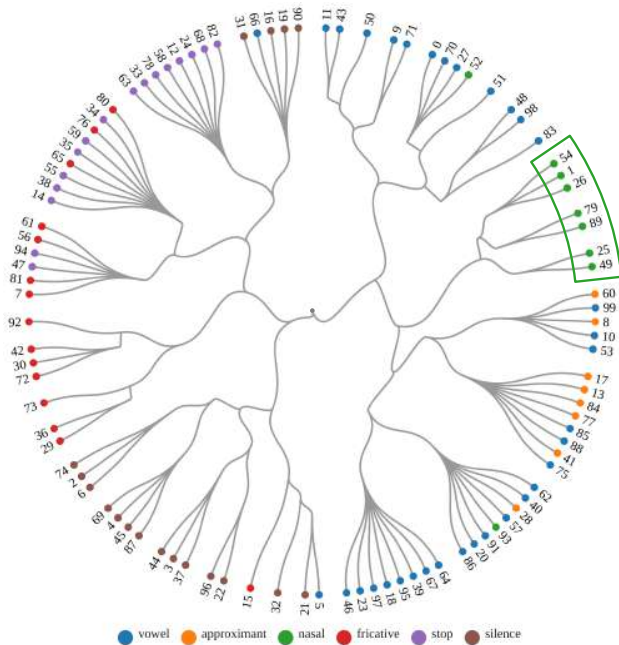


No modification:

Play

Stops:

Play



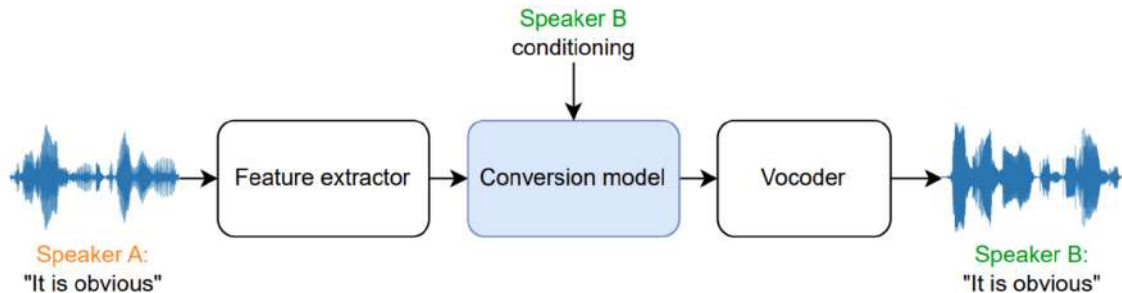
No modification:

Play

Nasals:

Play

Voice conversion

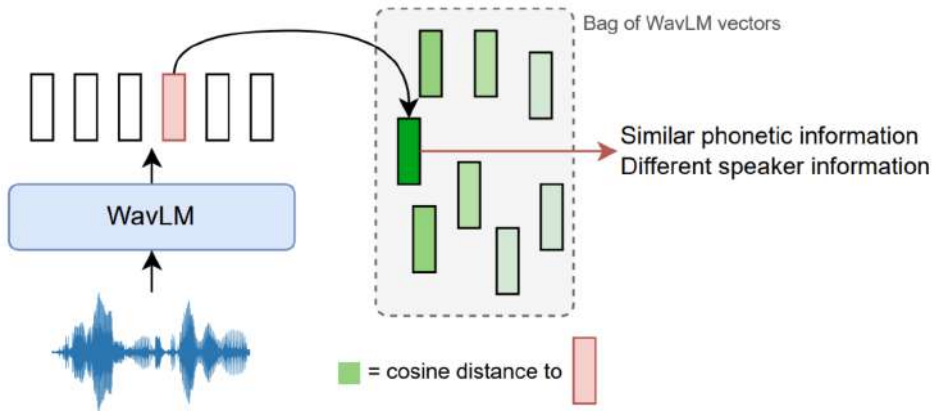


Source: [Play](#)

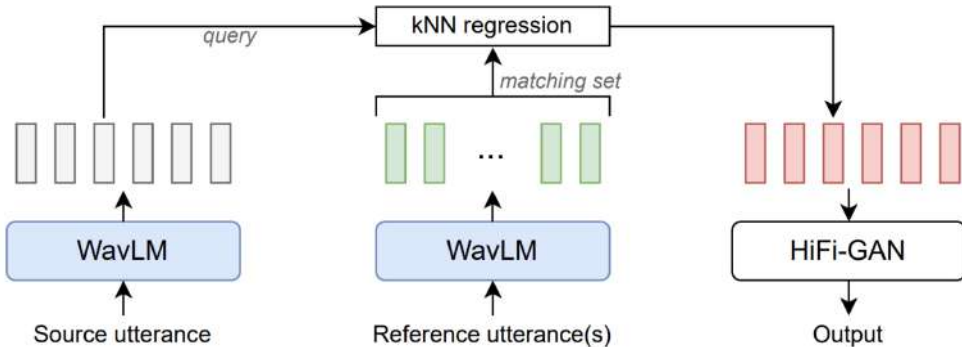
Reference: [Play](#)

Output: [Play](#)

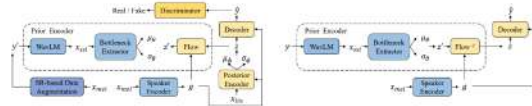
Our key idea



k -nearest neighbours voice conversion (kNN-VC)



Existing voice conversion systems



FreeVC [2022]

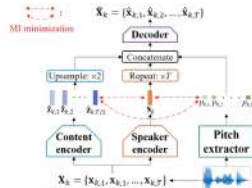
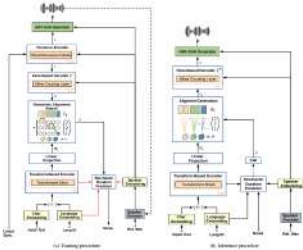


Figure 1: Diagram of the proposed VQMIVC system.

VQMIVC [2021]



YourTTS [2023]

Voice conversion results


Model	WER ↓	EER ↑	MOS ↑	SIM ↑
<i>Testset topline</i>	5.96	–	4.24	3.19
VQMIVC (Wang et al., 2021)	59.46	2.22	2.70	2.09
YourTTS (Casanova et al., 2022)	11.93	25.32	3.53	2.57
FreeVC (Li et al., 2022)	7.61	8.97	4.07	2.38
kNN-VC	7.36	37.15	4.03	2.91

Fun samples

Cross-lingual conversion:

Source: 


Reference: 

Output: 

Whispered music conversion:

Source: 


Reference: 

Output: 

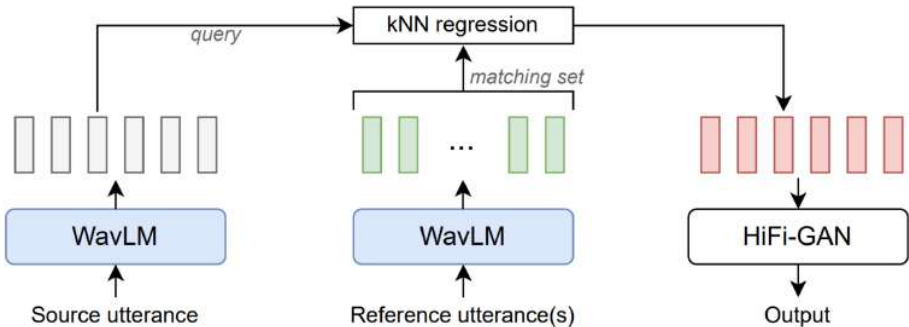
Human-to-animal conversion:

Source: 

Reference: 

Output: 

Voice conversion with stuttered reference speech



Source: [Play](#)

Reference: [Play](#)

Output: [Play](#)

Baseline: [Play](#) (TTS)

Source: [Play](#)

Reference: [Play](#)

Output: [Play](#)

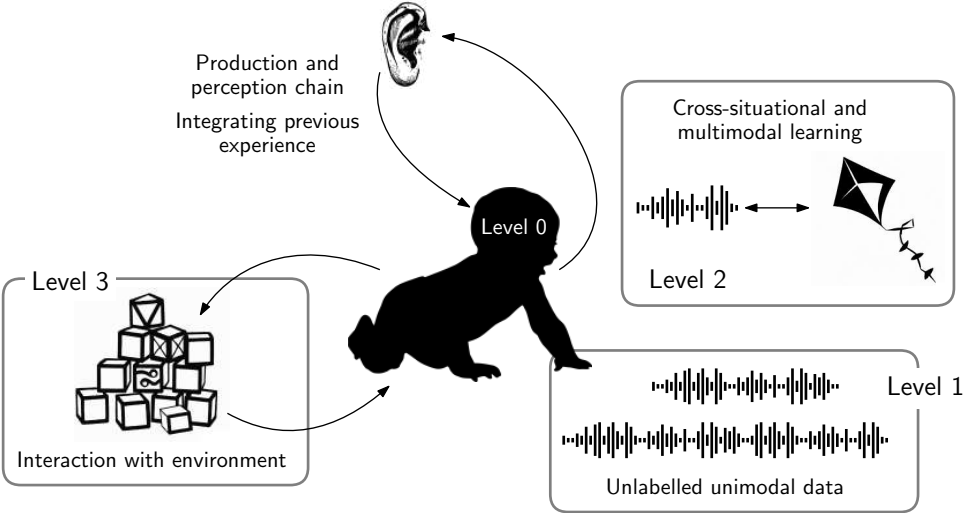
Baseline: [Play](#) (manual)

What does this tell us about self-supervised speech models?

- Broader phonetic categories are captured in hierarchy
- Phonetic content is matched through cosine distance
- But speaker characteristics are also still strongly captured

All of this is kind of expected, but it is still cool to be able to hear it!

Conclusion



<https://bshall.github.io/knn-vc/>

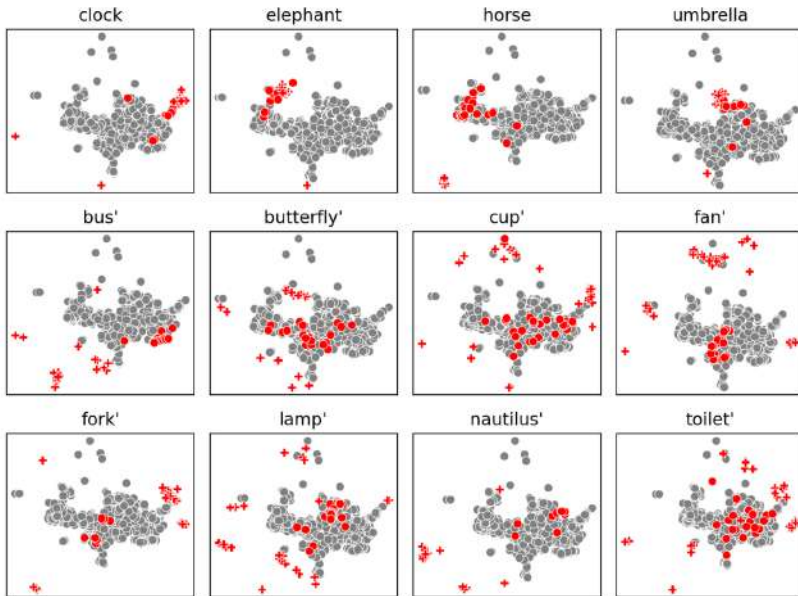
<https://www.kamperh.com/>

Please email me if you want to chat:

kamperh@gmail.com

Mutual exclusivity results

		Model initialisation		Accuracy (%)	
		Audio (CPC)	Vision (AlexNet)	Familiar- <u>familiar</u>	Familiar- <u>novel</u>
1	Random baseline	N/A	N/A	50.19	49.92
2	MattNet	✗	✗	72.86	57.29
3		✗	✓	85.89	59.32
4		✓	✗	75.78	55.92
5		✓	✓	83.20	60.27



Attention visualisation



“fire hydrant”

