
UNSUPERVISED VS. TRANSFER LEARNING FOR MULTIMODAL ONE-SHOT MATCHING OF SPEECH AND IMAGES

INTERSPEECH 2020

authors

Leanne Nortje [nortjeleanne@gmail.com]

Herman Kamper [kamperh@sun.ac.za]

institute

E&E Engineering, Stellenbosch University, South Africa

» PROBLEM

» PROBLEM

- ▷ Vision/speech processing systems require large amounts of labelled data.

» PROBLEM

- ▷ Vision/speech processing systems require large amounts of labelled data.
- ▷ **Where to find a solution**

» PROBLEM

- ▷ Vision/speech processing systems require large amounts of labelled data.
- ▷ **Where to find a solution:** Children learn from few examples.

» PROBLEM

- ▷ Vision/speech processing systems require large amounts of labelled data.
- ▷ **Where to find a solution:** Children learn from few examples.

words \iff visual objects

» How do CHILDREN LEARN?

» How do CHILDREN LEARN?



“ice-cream”

» How do CHILDREN LEARN?



“ice-cream”



“cookie”

» How do CHILDREN LEARN?



“ice-cream”



“cookie”



“broccoli”

» How do CHILDREN LEARN?



“ice-cream”



“cookie”



“broccoli”

Which picture is that of an “ice-cream”?



» How do CHILDREN LEARN?



“ice-cream”



“cookie”



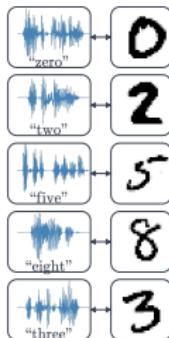
“broccoli”

Which picture is that of an “ice-cream”?



» MULTIMODAL ONE-SHOT MATCHING

Support set
 $\mathcal{S} = \{(\mathbf{x}_a^{(i)}, \mathbf{x}_v^{(i)})\}_{i=1}^L$

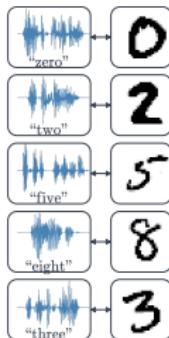


» MULTIMODAL ONE-SHOT MATCHING

Test question:

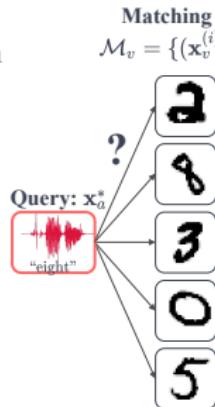
Support set

$$\mathcal{S} = \{(\mathbf{x}_a^{(i)}, \mathbf{x}_v^{(i)})\}_{i=1}^L$$



Matching set

$$\mathcal{M}_v = \{(\mathbf{x}_v^{(i)})\}_{i=1}^N$$



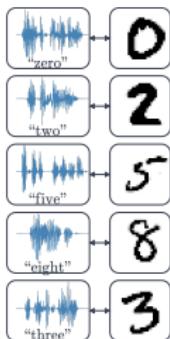
Which image matches the query "**eight**"?

» MULTIMODAL ONE-SHOT MATCHING

Test question:

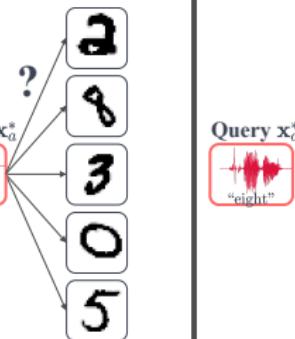
Support set

$$\mathcal{S} = \{(\mathbf{x}_a^{(i)}, \mathbf{x}_v^{(i)})\}_{i=1}^L$$



Matching set

$$\mathcal{M}_v = \{(\mathbf{x}_v^{(i)})\}_{i=1}^N$$



How?

Query \mathbf{x}_a^*



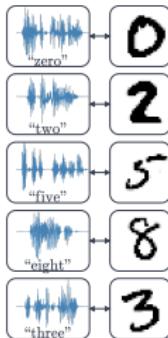
Which image matches the query "**eight**"?

» MULTIMODAL ONE-SHOT MATCHING

Test question:

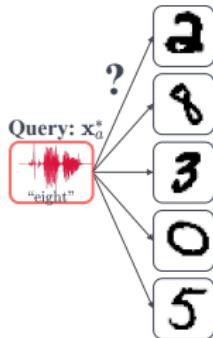
Support set

$$\mathcal{S} = \{(\mathbf{x}_a^{(i)}, \mathbf{x}_v^{(i)})\}_{i=1}^L$$



Matching set

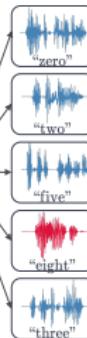
$$\mathcal{M}_v = \{(\mathbf{x}_v^{(i)})\}_{i=1}^N$$



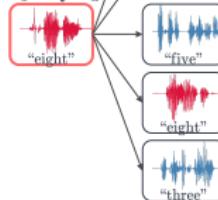
How?

Support set

$$\mathcal{S} = \{(\mathbf{x}_a^{(i)}, \mathbf{x}_v^{(i)})\}_{i=1}^L$$



Query \mathbf{x}_a^*



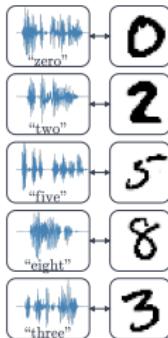
Which image matches the query **"eight"**?

» MULTIMODAL ONE-SHOT MATCHING

Test question:

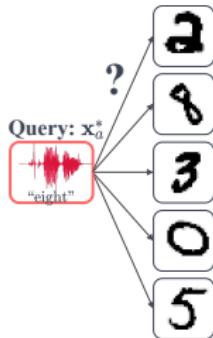
Support set

$$\mathcal{S} = \{(\mathbf{x}_a^{(i)}, \mathbf{x}_v^{(i)})\}_{i=1}^L$$



Matching set

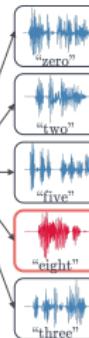
$$\mathcal{M}_v = \{(\mathbf{x}_v^{(i)})\}_{i=1}^N$$



How?

Support set

$$\mathcal{S} = \{(\mathbf{x}_a^{(i)}, \mathbf{x}_v^{(i)})\}_{i=1}^L$$



Query \mathbf{x}_a^*

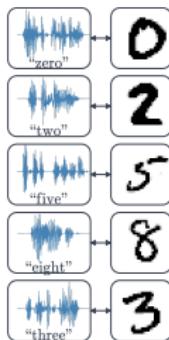
Which image matches the query **"eight"**?

» MULTIMODAL ONE-SHOT MATCHING

Test question:

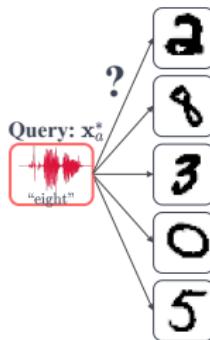
Support set

$$\mathcal{S} = \{(\mathbf{x}_a^{(i)}, \mathbf{x}_v^{(i)})\}_{i=1}^L$$



Matching set

$$\mathcal{M}_v = \{(\mathbf{x}_v^{(i)})\}_{i=1}^N$$



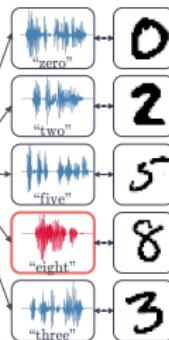
Query: \mathbf{x}_a^*

"eight"

How?

Support set

$$\mathcal{S} = \{(\mathbf{x}_a^{(i)}, \mathbf{x}_v^{(i)})\}_{i=1}^L$$



Query \mathbf{x}_a^*

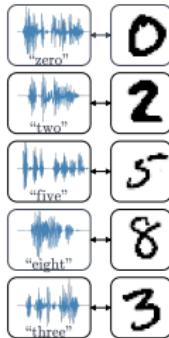
"eight"

Which image matches the query **"eight"**?

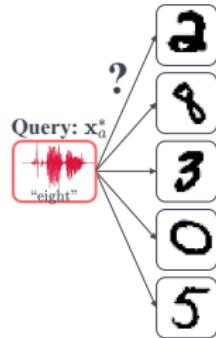
» MULTIMODAL ONE-SHOT MATCHING

Test question:

Support set
 $\mathcal{S} = \{(\mathbf{x}_a^{(i)}, \mathbf{x}_v^{(i)})\}_{i=1}^L$



Matching set
 $\mathcal{M}_v = \{(\mathbf{x}_v^{(i)})\}_{i=1}^N$

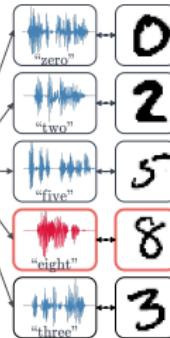


Query: \mathbf{x}_a^*

"eight"

How?

Support set
 $\mathcal{S} = \{(\mathbf{x}_a^{(i)}, \mathbf{x}_v^{(i)})\}_{i=1}^L$



Query \mathbf{x}_a^*

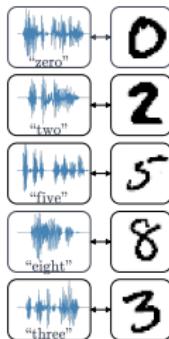
"eight"

Which image matches the query **"eight"**?

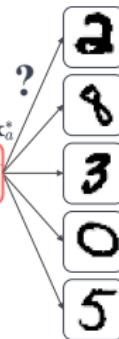
» MULTIMODAL ONE-SHOT MATCHING

Test question:

Support set
 $\mathcal{S} = \{(\mathbf{x}_a^{(i)}, \mathbf{x}_v^{(i)})\}_{i=1}^L$



Matching set
 $\mathcal{M}_v = \{(\mathbf{x}_v^{(i)})\}_{i=1}^N$

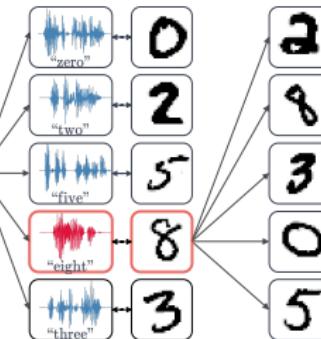


Query: \mathbf{x}_a^*

"eight"

How?

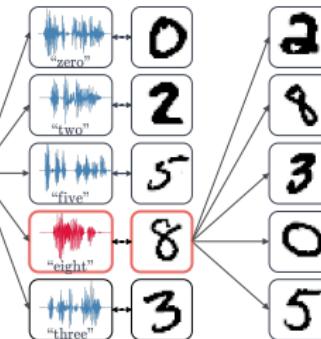
Support set
 $\mathcal{S} = \{(\mathbf{x}_a^{(i)}, \mathbf{x}_v^{(i)})\}_{i=1}^L$



Query \mathbf{x}_a^*

"eight"

Matching set
 $\mathcal{M}_v = \{(\mathbf{x}_v^{(i)})\}_{i=1}^N$

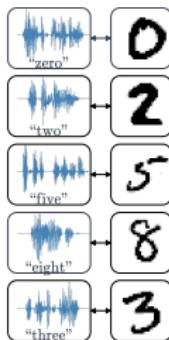


Which image matches the query "**eight**"?

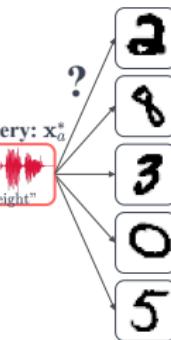
» MULTIMODAL ONE-SHOT MATCHING

Test question:

Support set
 $\mathcal{S} = \{(\mathbf{x}_a^{(i)}, \mathbf{x}_v^{(i)})\}_{i=1}^L$



Matching set
 $\mathcal{M}_v = \{(\mathbf{x}_v^{(i)})\}_{i=1}^N$



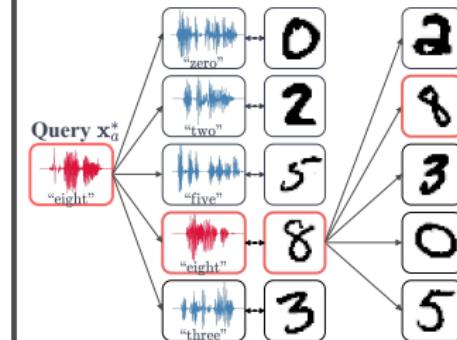
Query: \mathbf{x}_a^*

"eight"

How?

Support set
 $\mathcal{S} = \{(\mathbf{x}_a^{(i)}, \mathbf{x}_v^{(i)})\}_{i=1}^L$

Matching set
 $\mathcal{M}_v = \{(\mathbf{x}_v^{(i)})\}_{i=1}^N$



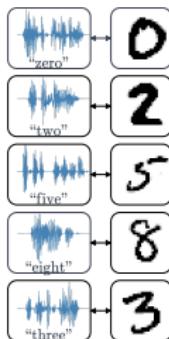
$$\mathcal{S} \rightarrow D_{\mathcal{S}}(\mathbf{x}_a, \mathbf{x}_v)$$

Which image matches the query **"eight"**?

» MULTIMODAL ONE-SHOT MATCHING

Test question:

Support set
 $\mathcal{S} = \{(\mathbf{x}_a^{(i)}, \mathbf{x}_v^{(i)})\}_{i=1}^L$



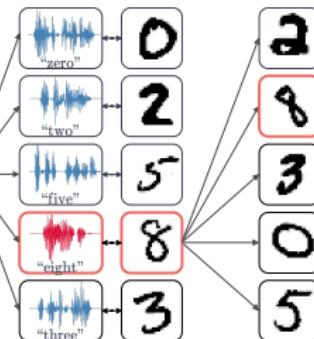
Matching set
 $\mathcal{M}_v = \{(\mathbf{x}_v^{(i)})\}_{i=1}^N$



How?

Support set
 $\mathcal{S} = \{(\mathbf{x}_a^{(i)}, \mathbf{x}_v^{(i)})\}_{i=1}^L$

Matching set
 $\mathcal{M}_v = \{(\mathbf{x}_v^{(i)})\}_{i=1}^N$



$$\mathcal{S} \rightarrow D_{\mathcal{S}}(\mathbf{x}_a, \mathbf{x}_v)$$

Multimodal one-shot learning: learn from one cross-modal paired example

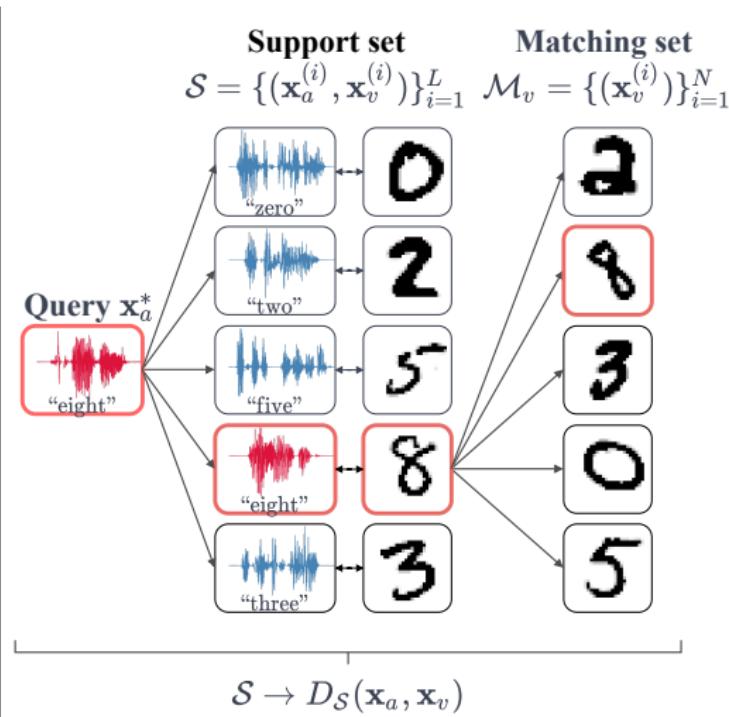
» MULTIMODAL ONE-SHOT MATCHING

Multimodal **one-shot** learning: learn from **one** cross-modal paired example.



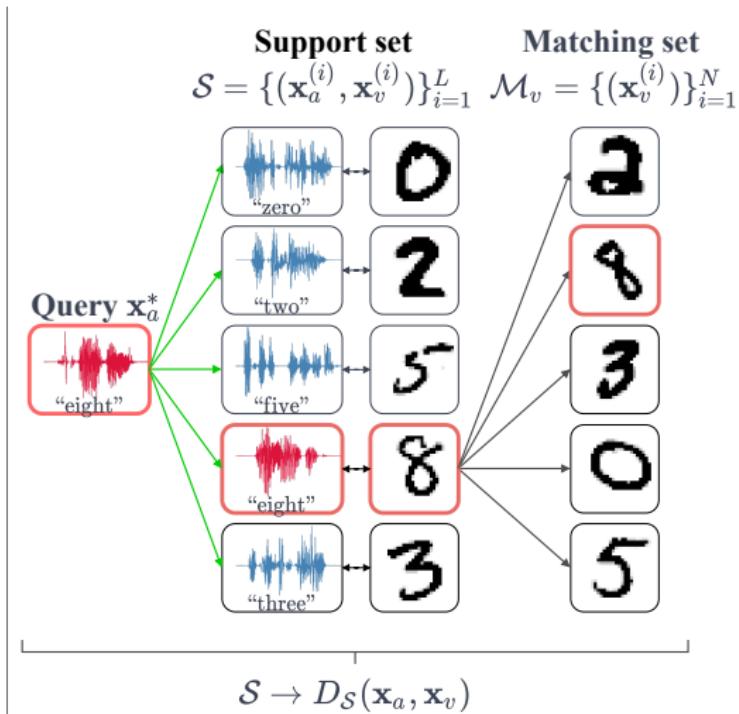
Multimodal **few-shot** learning: learn from a **few** cross-modal paired examples.

» MULTIMODAL ONE-SHOT MATCHING



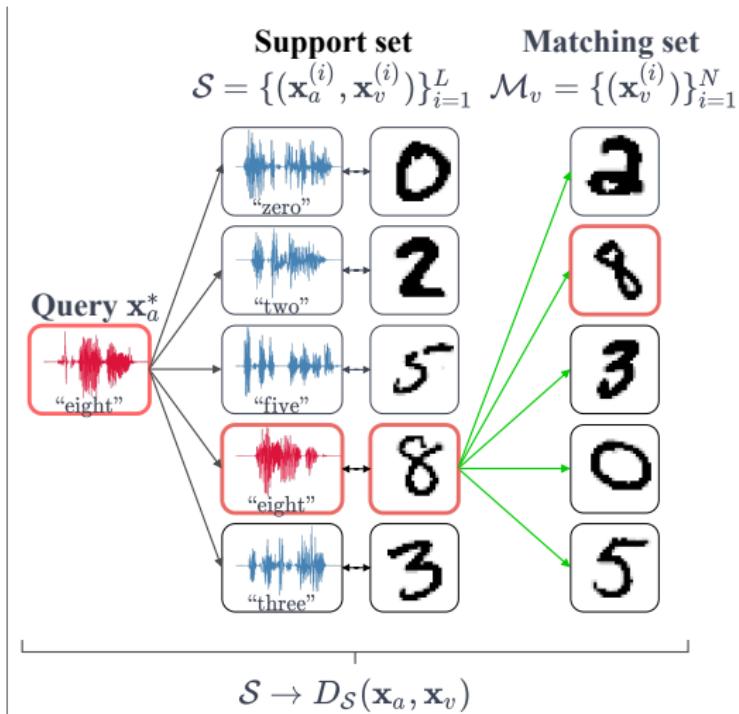
Our approach: a support set + two unimodal comparisons

» MULTIMODAL ONE-SHOT MATCHING



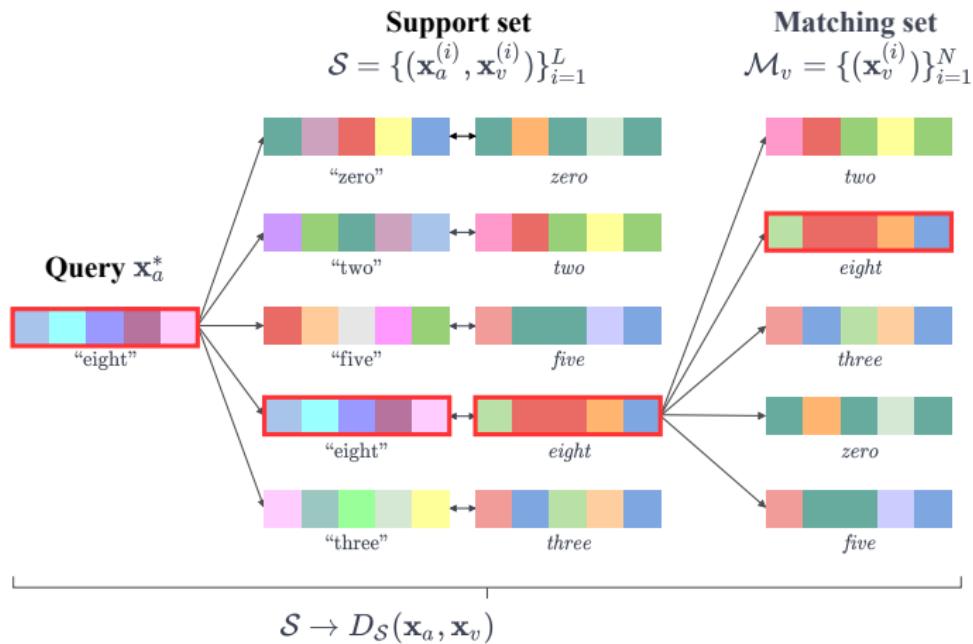
Our approach: a support set + two unimodal comparisons

» MULTIMODAL ONE-SHOT MATCHING



Our approach: a support set + two unimodal comparisons

» MULTIMODAL ONE-SHOT MATCHING



Speech comparison: cosine distance between two word representations
Image comparison: cosine distance between two image representations

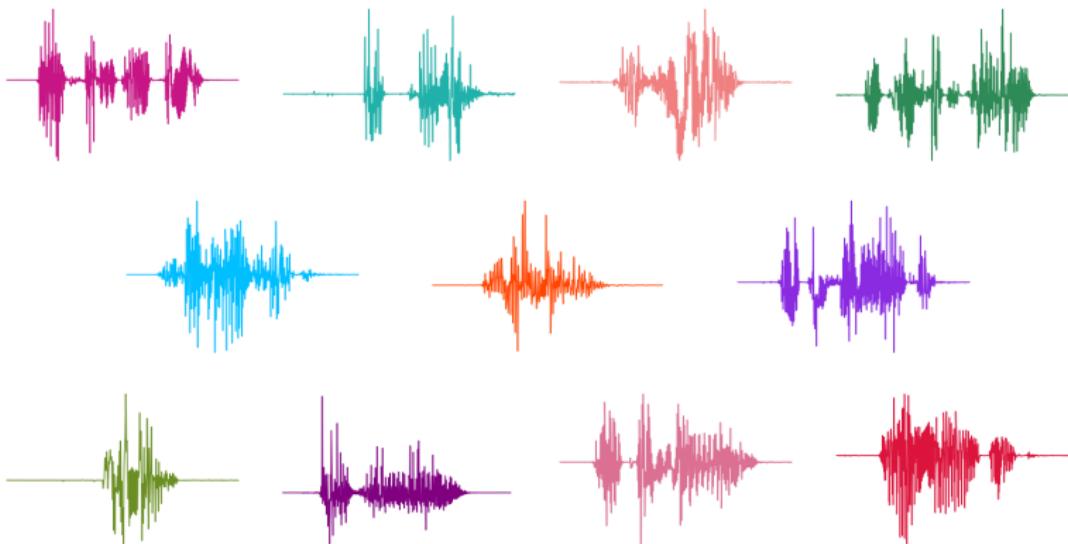
» UNSUPERVISED LEARNING

» UNSUPERVISED LEARNING

- ▷ Trained on unlabelled in-domain data.

» UNSUPERVISED LEARNING

- ▷ Trained on unlabelled in-domain data.
- ▷ Unlabelled in-domain data: **TIDigits**



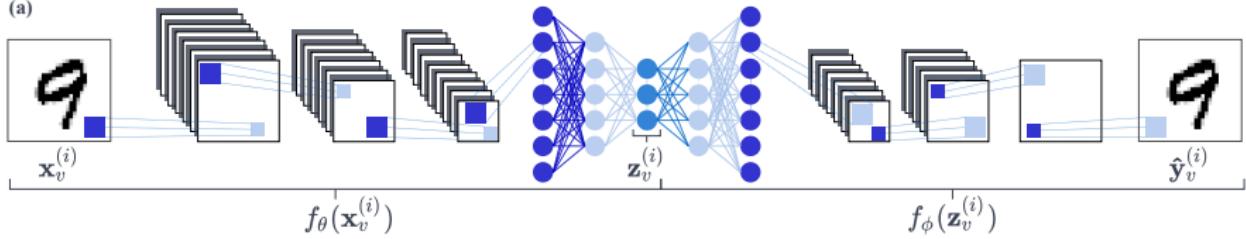
» UNSUPERVISED LEARNING

- ▷ Trained on unlabelled in-domain data.
- ▷ Unlabelled in-domain data: TIDigits and MNIST

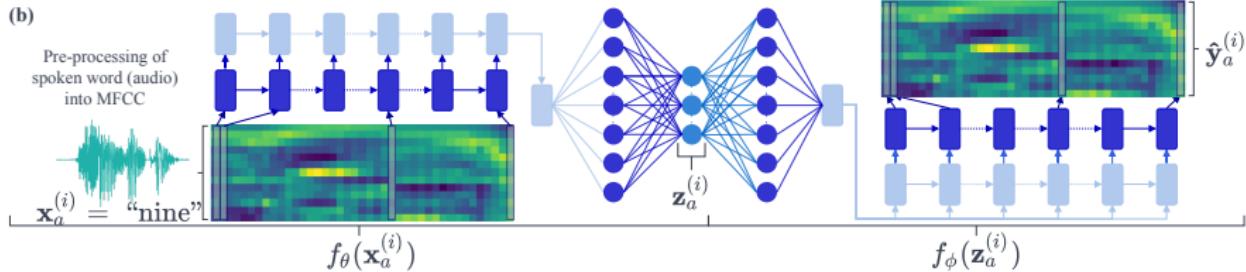
A 5x10 grid of handwritten digits, likely from the TIDigits dataset. The digits are arranged in five rows and ten columns. The digits are written in a cursive, black font on a white background. Some digits are clearly legible, while others are more stylized or partially obscured. The digits include 7, 4, 8, 0, 3, 1, 2, 9, 6, 0, 0, 1, 3, 5, 6, 0, 5, 4, 2, 9, 7, 5, 2, 7, 9, 3, 8, 0, 7, 8, 1, 0, 9, 6, 0, 5, 1, 6, 7, 4, 0, 2, 4, 4, 0, 9, 3, 1, 5, 0, 8, 9, 3, 6.

» UNSUPERVISED LEARNING

(a)

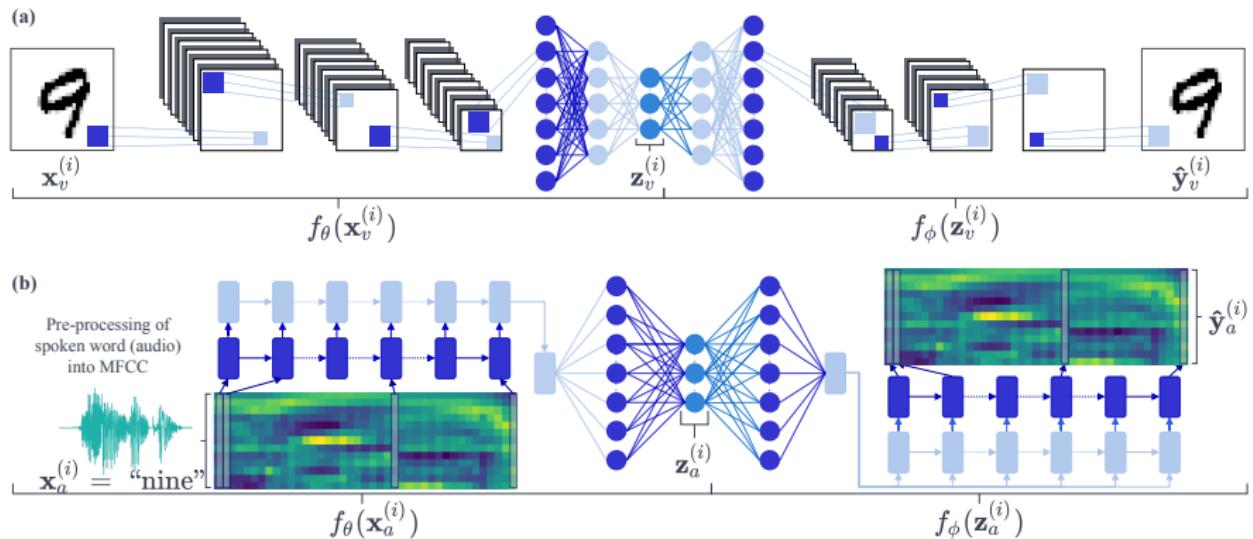


(b)



» UNSUPERVISED LEARNING

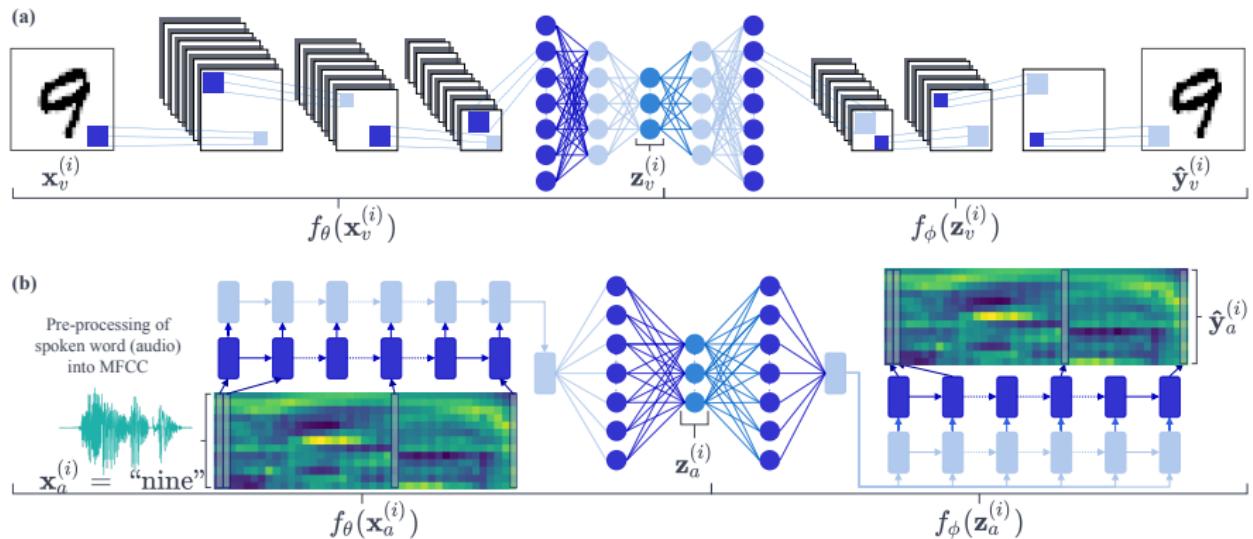
Autoencoder-like model architectures:



» UNSUPERVISED LEARNING

Autoencoder-like model architectures:

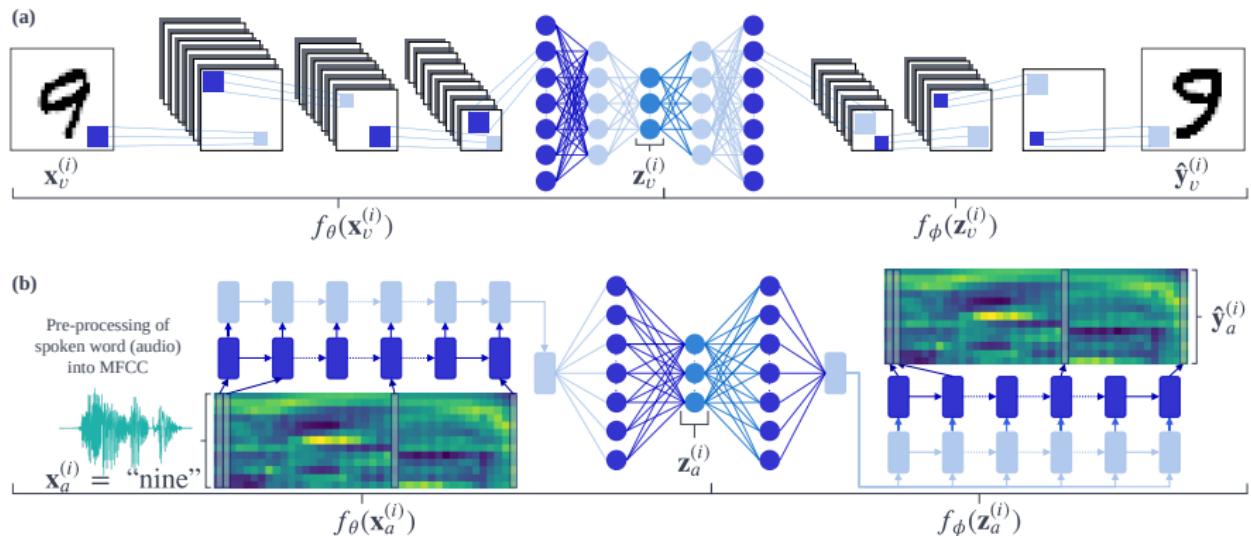
- ▷ Autoencoder (AE)



» UNSUPERVISED LEARNING

Autoencoder-like model architectures:

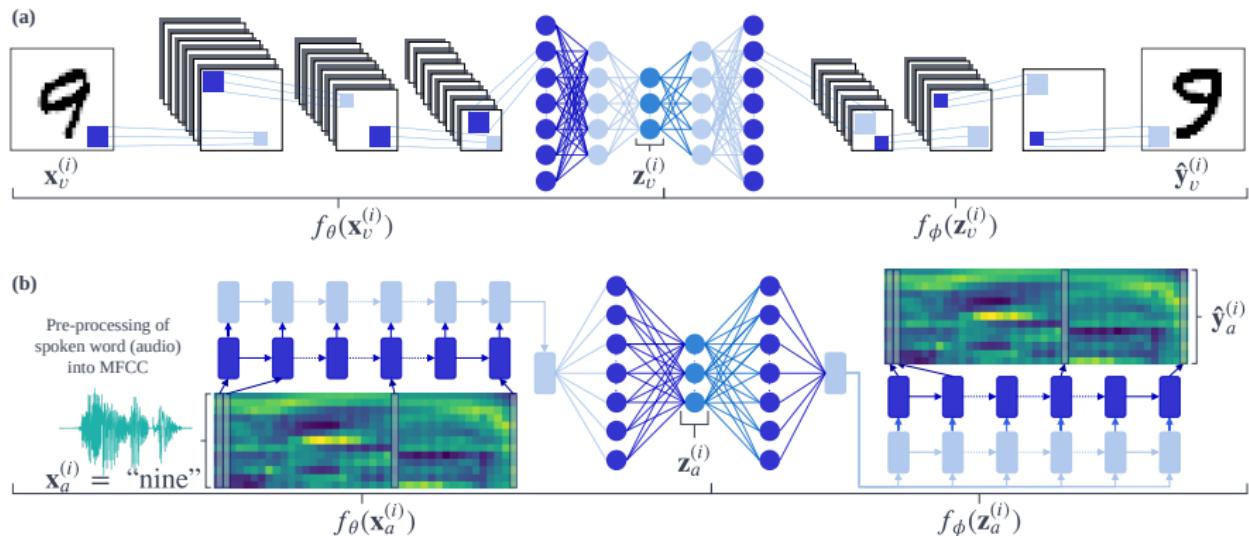
- ▷ Autoencoder (AE)
- ▷ Correspondence autoencoder (CAE)



» UNSUPERVISED LEARNING

Autoencoder-like model architectures:

- ▷ Autoencoder (AE)
- ▷ Correspondence autoencoder (CAE) (unsupervised within-modality pairs).



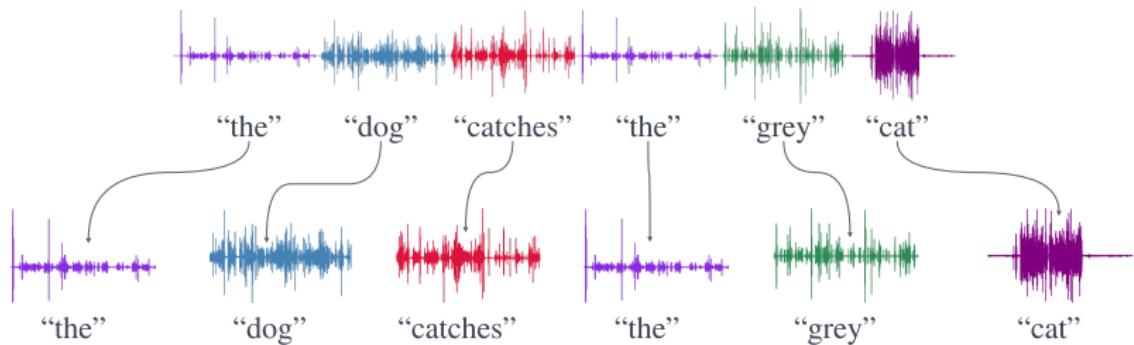
» TRANSFER LEARNING

» TRANSFER LEARNING

- ▷ Trained on labelled background data.

» TRANSFER LEARNING

- ▷ Trained on labelled background data.
- ▷ Labelled background data: **Buckeye**



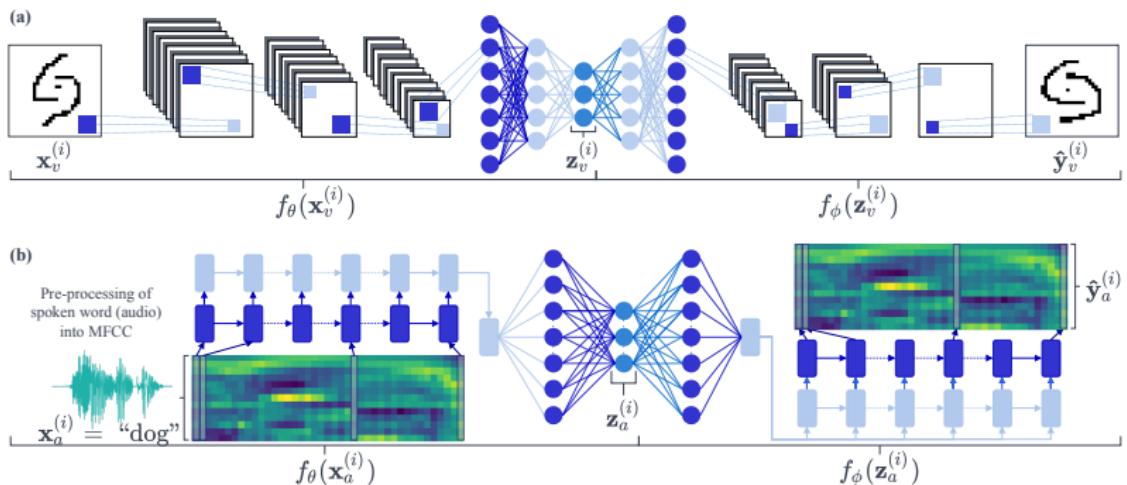
» TRANSFER LEARNING

- ▷ Trained on labelled background data.
- ▷ Labelled background data: Buckeye and Omniglot.

e ð l w g ſ ð q þ v ð œ ø ð
u ɹ n b h ɬ ɛ ɸ ɔ ɪ ɔ ʊ t
ɔ ʌ ɒ ɔ ʃ ɔ ɔ ɹ ʌ ɒ ɒ ɒ
ʌ ɻ ʌ ɔ ɹ ɹ ɹ ɹ ɹ ɹ ɹ ɹ
ð ɹ ɹ ɹ ɹ ɹ ɹ ɹ ɹ ɹ ɹ ɹ

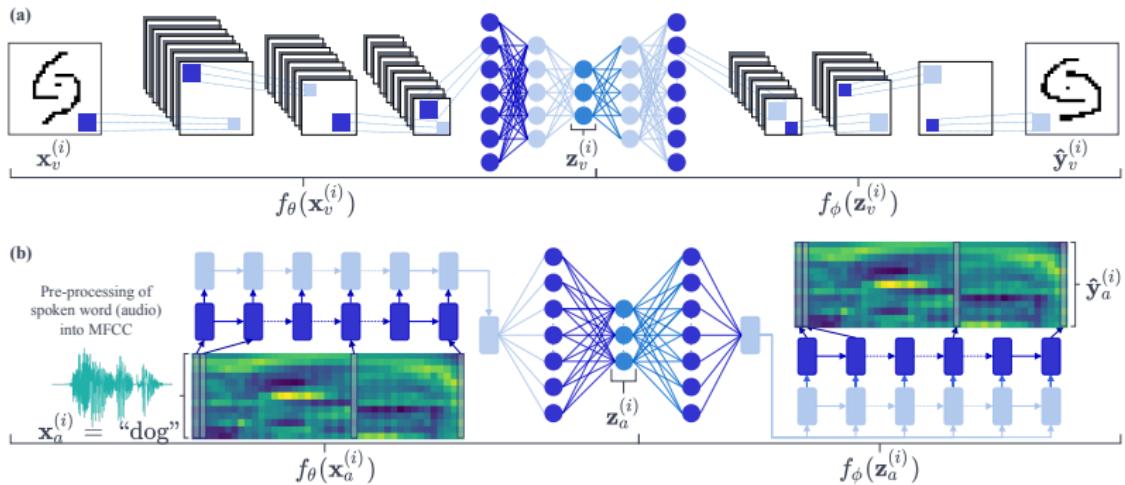
» TRANSFER LEARNING

- ▷ A transfer learned variant of the unsupervised CAE:



» TRANSFER LEARNING

- ▷ A transfer learned variant of the unsupervised CAE:
- ▷ Trained on **ground truth pairs**.



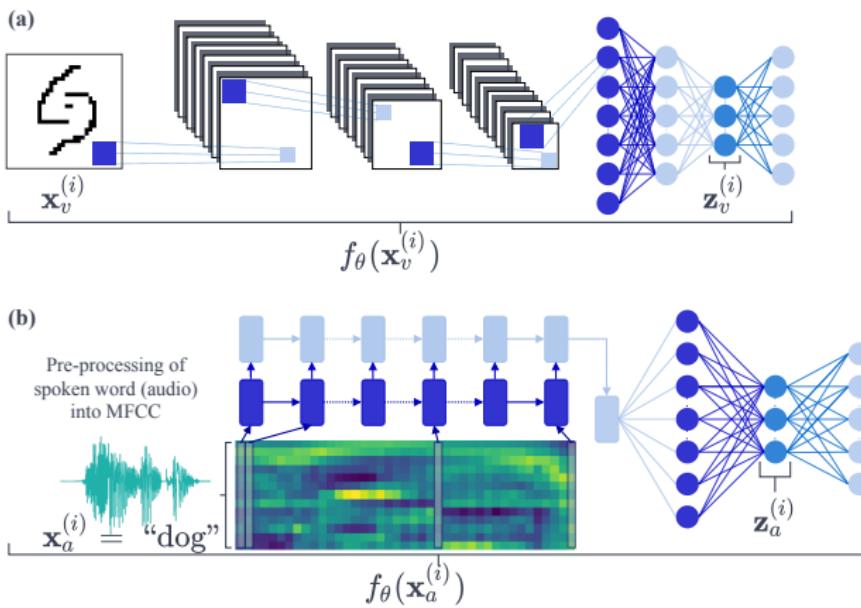
» TRANSFER LEARNING

Multimodal few-shot models from Elof et al. [1]:

» TRANSFER LEARNING

Multimodal few-shot models from Elof et al. [1]:

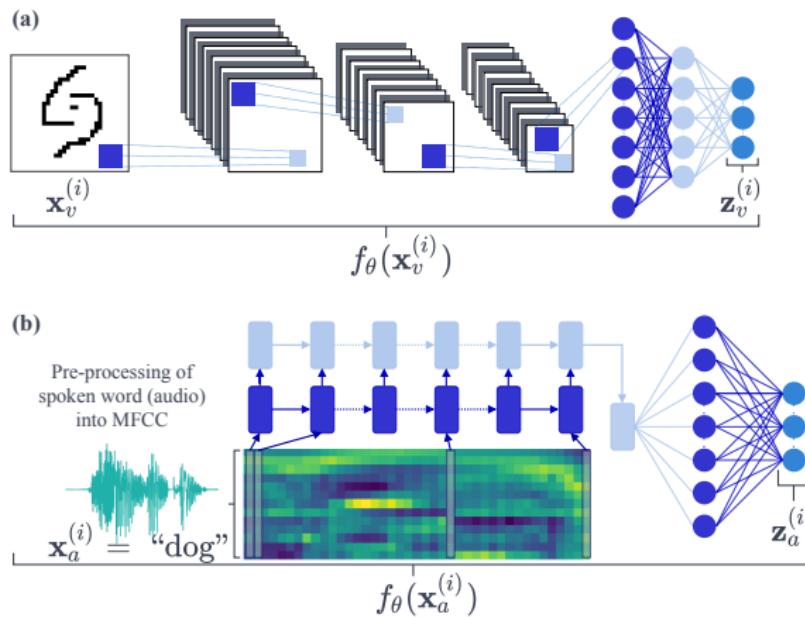
- ▷ Classifiers



» TRANSFER LEARNING

Multimodal few-shot models from Elof et al. [1]:

- ▷ Classifiers and
- ▷ Siamese Triplet networks.



» K-SHOT MULTIMODAL SPEECH AND IMAGE MATCHING

Model	DTW + Pixels	11-way accuracy (%)	
		one-shot	five-shot
Baseline	DTW + Pixels	31.80	41.88
Transfer learning models	Classifier	56.80 ± 1.19	59.67 ± 1.73
	Siamese	54.83 ± 1.80	59.25 ± 0.79
	CAE	46.60 ± 0.69	53.82 ± 1.07
Unsupervised models	AE	28.99 ± 0.84	38.68 ± 1.51
	CAE	42.75 ± 0.62	52.15 ± 0.69

» K-SHOT MULTIMODAL SPEECH AND IMAGE MATCHING

Model	DTW + Pixels	11-way accuracy (%)	
		one-shot	five-shot
Baseline	DTW + Pixels	31.80	41.88
Transfer learning models	Classifier	56.80 ± 1.19	59.67 ± 1.73
	Siamese	54.83 ± 1.80	59.25 ± 0.79
	CAE	46.60 ± 0.69	53.82 ± 1.07
Unsupervised models	AE	28.99 ± 0.84	38.68 ± 1.51
	CAE	42.75 ± 0.62	52.15 ± 0.69

» K-SHOT MULTIMODAL SPEECH AND IMAGE MATCHING

Model	DTW + Pixels	11-way accuracy (%)	
		one-shot	five-shot
Baseline	DTW + Pixels	31.80	41.88
Transfer learning models	Classifier	56.80 ± 1.19	59.67 ± 1.73
	Siamese	54.83 ± 1.80	59.25 ± 0.79
	CAE	46.60 ± 0.69	53.82 ± 1.07
Unsupervised models	AE	28.99 ± 0.84	38.68 ± 1.51
	CAE	42.75 ± 0.62	52.15 ± 0.69

» COULD UNSUPERVISED AND TRANSFER LEARNING BE COMPLEMENTARY?

» COULD UNSUPERVISED AND TRANSFER LEARNING BE COMPLEMENTARY?

- ▷ **Classifier pairs:** We find unsupervised training pairs using transfer learned classifiers.

» COULD UNSUPERVISED AND TRANSFER LEARNING BE COMPLEMENTARY?

- ▷ **Classifier pairs:** We find unsupervised training pairs using transfer learned classifiers.
- ▷ **CAE with classifier pairs:** Train an unsupervised CAE using these classifier pairs.

» COULD UNSUPERVISED AND TRANSFER LEARNING BE COMPLEMENTARY?

- ▷ **Classifier pairs:** We find unsupervised training pairs using transfer learned classifiers.
- ▷ **CAE with classifier pairs:** Train an unsupervised CAE using these classifier pairs.
- ▷ **Transfer learning + CAE fine-tuning:** Pretrain a CAE on ground truth background pairs and then train the CAE on these classifier pairs.

» TOWARDS COMBINED TRANSFER AND UNSUPERVISED LEARNING

Model	11-way accuracy (%)	
	one-shot	five-shot
Baseline: DTW + Pixels	31.80	41.88
Transfer learning: Classifier	56.80 ± 1.19	59.67 ± 1.73
CAE with cosine pairs	42.75 ± 0.62	52.15 ± 0.69
CAE with classifier pairs	48.66 ± 1.14	55.59 ± 0.71
Transfer learning + CAE fine-tuning	54.32 ± 2.19	59.37 ± 1.80
CAE with oracle pairs	89.19 ± 0.69	92.81 ± 0.47

» TOWARDS COMBINED TRANSFER AND UNSUPERVISED LEARNING

Model	11-way accuracy (%)	
	one-shot	five-shot
Baseline: DTW + Pixels	31.80	41.88
Transfer learning: Classifier	56.80 ± 1.19	59.67 ± 1.73
CAE with cosine pairs	42.75 ± 0.62	52.15 ± 0.69
CAE with classifier pairs	48.66 ± 1.14	55.59 ± 0.71
Transfer learning + CAE fine-tuning	54.32 ± 2.19	59.37 ± 1.80
CAE with oracle pairs	89.19 ± 0.69	92.81 ± 0.47

» TOWARDS COMBINED TRANSFER AND UNSUPERVISED LEARNING

Model	11-way accuracy (%)	
	one-shot	five-shot
Baseline: DTW + Pixels	31.80	41.88
Transfer learning: Classifier	56.80 ± 1.19	59.67 ± 1.73
CAE with cosine pairs	42.75 ± 0.62	52.15 ± 0.69
CAE with classifier pairs	48.66 ± 1.14	55.59 ± 0.71
Transfer learning + CAE fine-tuning	54.32 ± 2.19	59.37 ± 1.80
CAE with oracle pairs	89.19 ± 0.69	92.81 ± 0.47

» TOWARDS COMBINED TRANSFER AND UNSUPERVISED LEARNING

Model	11-way accuracy (%)	
	one-shot	five-shot
Baseline: DTW + Pixels	31.80	41.88
Transfer learning: Classifier	56.80 ± 1.19	59.67 ± 1.73
CAE with cosine pairs	42.75 ± 0.62	52.15 ± 0.69
CAE with classifier pairs	48.66 ± 1.14	55.59 ± 0.71
Transfer learning + CAE fine-tuning	54.32 ± 2.19	59.37 ± 1.80
CAE with oracle pairs	89.19 ± 0.69	92.81 ± 0.47

» CONCLUSIONS

» CONCLUSIONS

- ▷ Transfer learning outperforms unsupervised learning.

» CONCLUSIONS

- ▷ Transfer learning outperforms unsupervised learning.

- ▷ Unsupervised learning can be improved by using transfer learning.

» CONCLUSIONS

- ▷ Transfer learning outperforms unsupervised learning.
- ▷ Unsupervised learning can be improved by using transfer learning.
- ▷ Idealised experiments show the promise of unsupervised learning.

» REFERENCES

- [1] R. Eloff, H. A. Engelbrecht, and H. Kamper, "Multimodal one-shot learning of speech and images," in *Proc. ICCASP*, 2019.

» ACKNOWLEDGEMENTS

This work is supported in part by the National Research Foundation of South Africa (grant number: 120409), a Google Faculty Award for Herman Kamper, a DST CSIR scholarship for Leanne Nortje, and funding from Saigen.