

# Towards Localisation of Keywords in Speech using Weak Supervision

Kayode Olaleye   Benjamin van Niekerk  
Herman Kamper



UNIVERSITEIT  
iYUNIVESITHI  
STELLENBOSCH  
UNIVERSITY



Department of Electrical and Electronic Engineering  
Stellenbosch University

December 11, 2020

## Detection

Query: *snow*

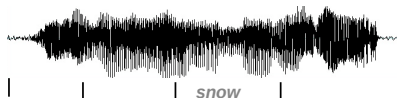


Utterance

Does it occur?

## Localisation

Query: *snow*



Utterance

Where does it occur?

## Visual supervision



Image



Utterance

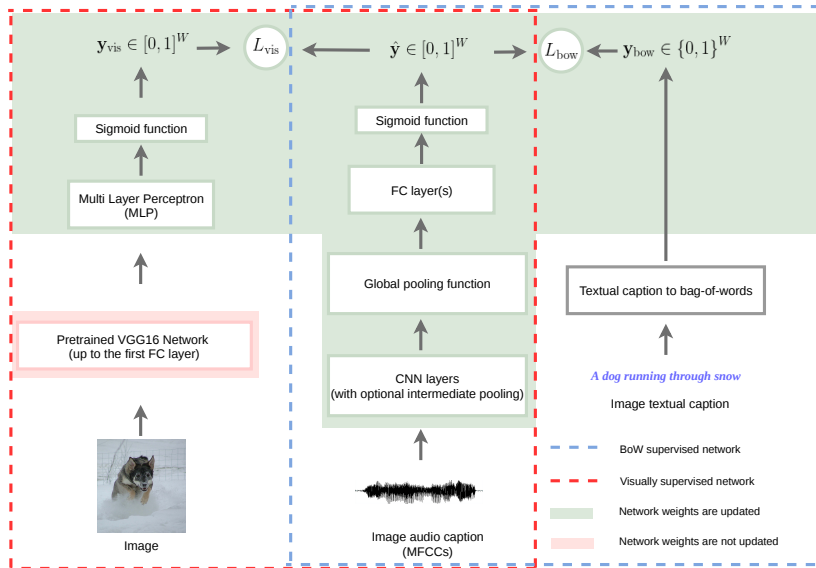
## Bag-of-Word (BoW) supervision

*snow, running, through, dog*

BoW



Utterance

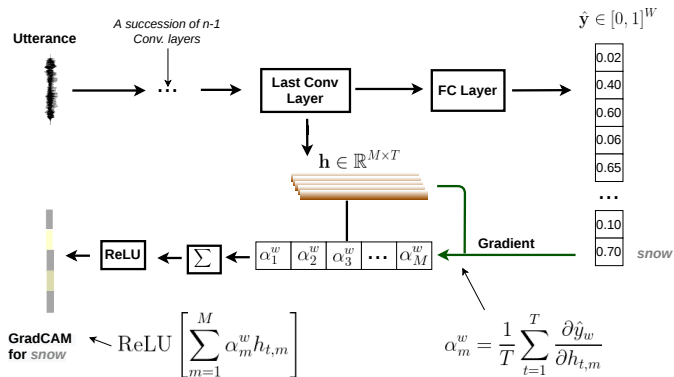


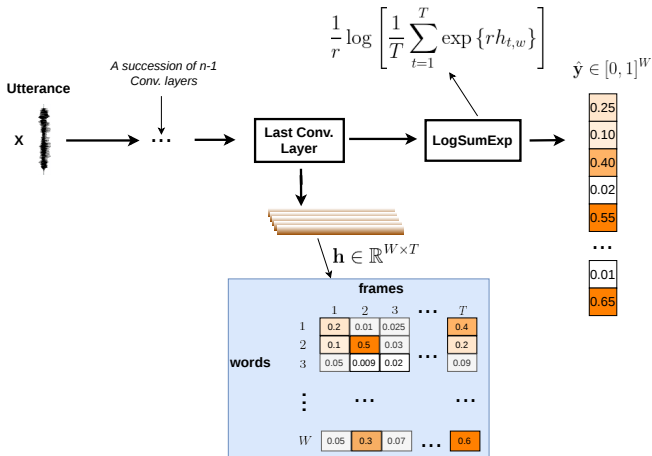
## ► GradCAM

- Introduced in the vision domain to localise an object in an image.
- Works with any **trained** CNN architecture.
- Determines the portion of an input that contributes to a decision of interest using gradient information.

## ► Palaz, Synnaeve, and Collobert (PSC)

- Designed to simultaneously perform detection and localisation of keywords in speech utterance.
- The CNN architecture is restricted in some ways (*No intermediate max-pooling; no fully-connected layers; LogSumExp function as the global pooling function*).





Two forms of weak supervision:

- ▶ BoW
- ▶ Visual

Two localisation methods:

- ▶ PSC
- ▶ GradCAM



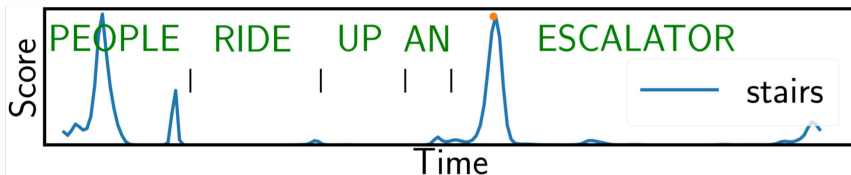
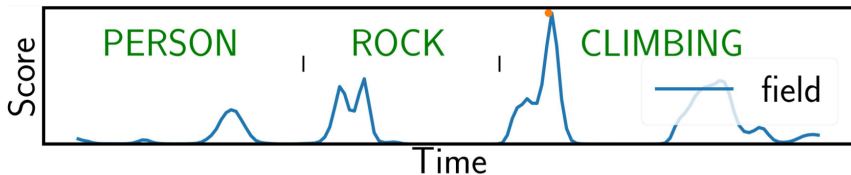
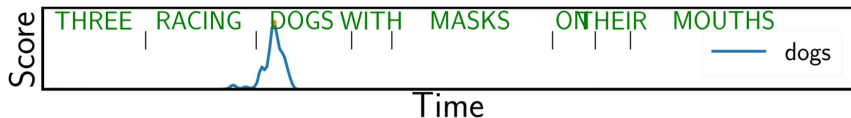
## Oracle accuracy

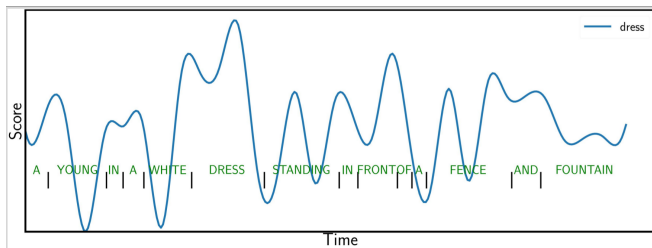
Mechanism	Supervision method	
	BoW	Visual
<b>PSC</b>	63.6	19.1
<b>GradCAM</b>	17.8	16.0

## Actual accuracy

Mechanism	BoW				Visual			
	<i>P</i>	<i>R</i>	<i>F1</i>	Accuracy	<i>P</i>	<i>R</i>	<i>F1</i>	Accuracy
<b>PSC</b>	75.2	53.0	62.2	50.4	28.6	8.0	12.5	7.6
<b>GradCAM</b>	17.7	24.5	20.5	13.2	5.0	5.7	5.3	4.4

# Examples from PSC model (with visual supervision)





- Method is optimised for multi-class classification: high probability for a particular class implies low probabilities for others.
- We use it here for a multi-label classification. Hence, during backward pass, Gradcam puts peaks over all the words in the utterance while locating "*dress*"

- ▶ Our question: Is keyword localisation in speech possible with two forms of weak supervision where location information is not provided?
- ▶ Compared BoW supervision versus visual supervision, and PSC versus GradCAM.
- ▶ BoW-trained model outperformed visually-trained model. PSC outperformed GradCAM on the localisation task.
- ▶ Visual supervision provides potential for high precision localisation.
- ▶ Mismatch between GradGAM and multi-label classification loss: poor performance.
- ▶ Should investigate better localisation methods.

Thank you for listening!

Keyword detection scores (without considering localisation) with threshold  $\alpha$ .

Model	$\alpha = 0.4$			$\alpha = 0.6$		
	$P$	$R$	$F1$	$P$	$R$	$F1$
<i>Visual supervision:</i>						
<b>PSC</b>	44.5	9.8	16.1	74.7	4.3	8.1
<b>GradCAM</b>	29.3	22.0	25.1	42.7	12.7	19.6
<i>BoW supervision:</i>						
<b>PSC</b>	82.2	49.0	61.4	87.8	46.1	60.4
<b>GradCAM</b>	79.3	52.6	63.2	82.5	50.9	63.0

1. Palaz, Dimitri, Gabriel Synnaeve, and Ronan Collobert. "Jointly Learning to Locate and Classify Words Using Convolutional Networks." INTERSPEECH, 2016.
2. Chrupała, Grzegorz, Lieke Gelderloos, and Afra Alishahi. "Representations of language in a model of visually grounded speech signal." arXiv preprint arXiv:1702.01991, 2017.
3. Kamper, Herman, Aristotelis Anastassiou, and Karen Livescu. "Semantic query-by-example speech search using visual grounding." ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE, 2019.
4. Abdel-Hamid, Ossama, et al. "Convolutional neural networks for speech recognition." IEEE/ACM Transactions on audio, speech, and language processing, 2014.
5. Doersch, Carl, and Andrew Zisserman. "Multi-task self-supervised visual learning." Proceedings of the IEEE International Conference on Computer Vision, 2017.

6. Aytaç, Yusuf, Carl Vondrick, and Antonio Torralba. "Soundnet: Learning sound representations from unlabeled video." Advances in neural information processing systems, 2016.
7. Harwath, David, and James Glass. "Deep multimodal semantic embeddings for speech and images." 2015 IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU). IEEE, 2015.
8. Synnaeve, Gabriel, Maarten Versteegh, and Emmanuel Dupoux. "Learning words from images and speech." NIPS Workshop Learn. Semantics, 2014.
9. Kamper, Herman, and Michael Roth. "Visually grounded cross-lingual keyword spotting in speech." In Proc. SLTU, 2018.
10. Selvaraju, Ramprasaath R., et al. "Grad-cam: Visual explanations from deep networks via gradient-based localization." Proceedings of the IEEE international conference on computer vision, 2017.
11. Harwath, David, et al. "Jointly discovering visual objects and spoken words from raw sensory input." Proceedings of the European conference on computer vision (ECCV), 2018.



12. Harwath, David, and James R. Glass. "Learning word-like units from joint audio-visual analysis." In Proc. ACL, 2017.
13. Settle, Shane, et al. "Query-by-example search with discriminative neural acoustic word embeddings." In Proc. Interspeech, 2017.
14. Bansal, Sameer, et al. "Towards speech-to-text translation without speech recognition." In Proc. EACL, 2017.
15. Kamper, Herman, Gregory Shakhnarovich, and Karen Livescu. "Semantic speech retrieval with a visually grounded model of untranscribed speech." IEEE/ACM Trans. Acoust., Speech, Signal Process, 2018.
16. Kamper, Herman, Aristotelis Anastassiou, and Karen Livescu. "Semantic query-by-example speech search using visual grounding." In Proc. ICASSP, 2019.
17. Pasad, Ankita, et al. "On the contributions of visual and textual supervision in low-resource semantic speech retrieval." In Proc. Interspeech, 2019.