

Overview

- **Background:** Visual grounding is a commonly used source of weak supervision for tasks involving untranscribed spoken data (e.g. [1,2]).
- **Open question:** Does visual grounding still help if we have text annotations during training?
- **Our setting:** A low-resource setting where a fraction of the spoken training corpus is transcribed.
- **Our work:** Explores how to best combine the two modalities for *semantic speech retrieval*

Query word	Retrieved utterance
kids	a group of young boys playing soccer
beach	a dog retrieves a branch from a beach

Multi-Task Learning (MTL) Approach

What are the multiple tasks?

Visually supervised task (MTL-visSup)

Trained on *image-speech pairs*¹. An external *image tagger*² provides weak labels as ground truth.

Textually supervised task (MTL-textSup)

Trained on *speech-text pairs*³. Each transcript provides a multi-hot bag-of-words vector as ground truth.

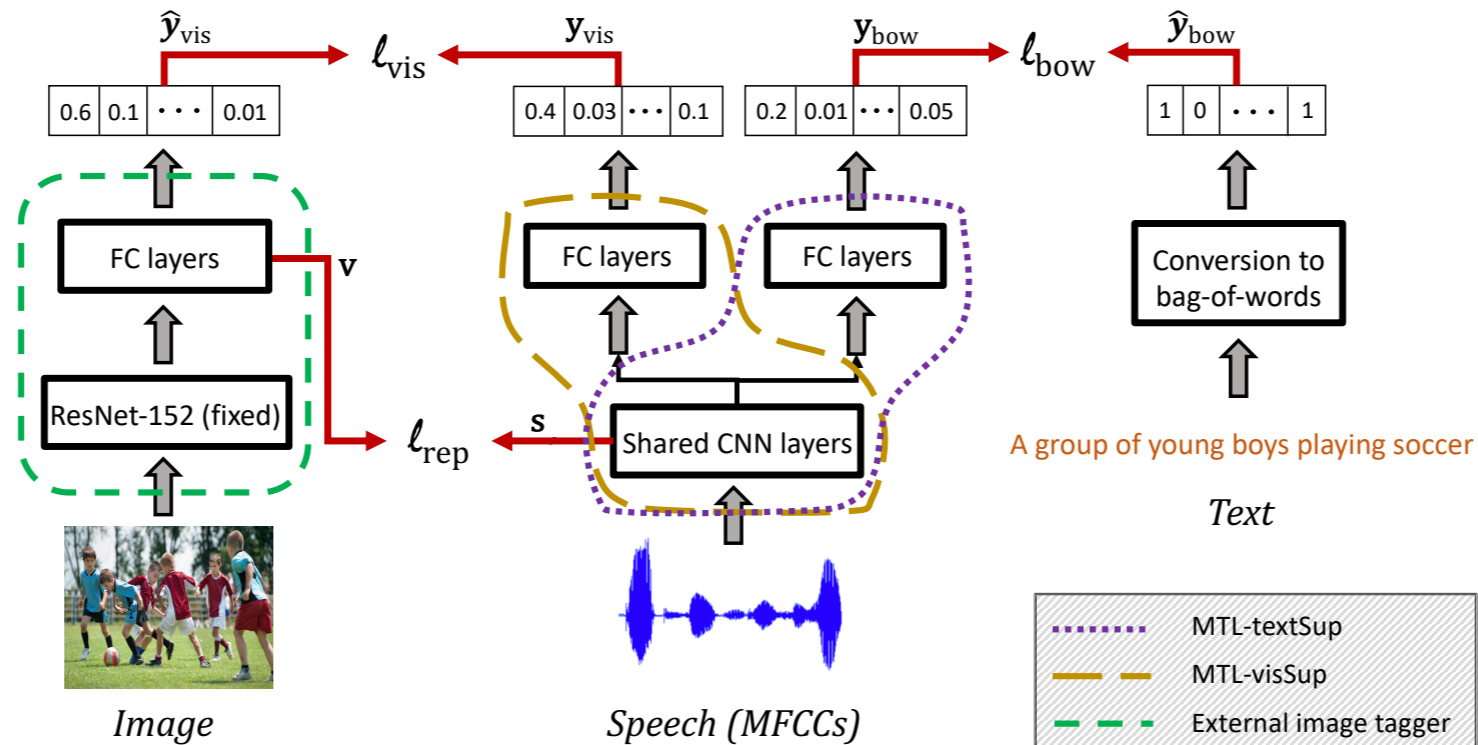
Unsupervised representation learning

Trained using the intermediate visual and speech representations. The former is fixed; the latter is updated during training.

What are the loss functions?

Supervised task losses ($sup \in \{vis, bow\}$): summed cross entropy between the predicted and ground truth vectors.

$$\ell_{sup} = - \sum_{w=1}^{N_{sup}} \left\{ \hat{y}_{sup,w} \log y_{sup,w} + (1 - \hat{y}_{sup,w}) \log [1 - y_{sup,w}] \right\}$$



Representation loss: margin-based contrastive loss with margin m , positive pair $\{v, s\}$, negative pairs $\{v', s\}$ and $\{v, s'\}$, and cosine distance.

$$\ell_{rep} = \left\{ \frac{1}{|V|} \sum_{v' \in V} \max[0, m + d_{\cos}(v, s) - d_{\cos}(v', s)] + \frac{1}{|S|} \sum_{s' \in S} \max[0, m + d_{\cos}(v, s) - d_{\cos}(v, s')] \right\}$$

Total loss: weighted sum of the three losses.

$$\ell = \alpha \cdot \ell_{vis} + \beta \cdot \ell_{bow} + (1 - \alpha - \beta) \cdot \ell_{rep}$$

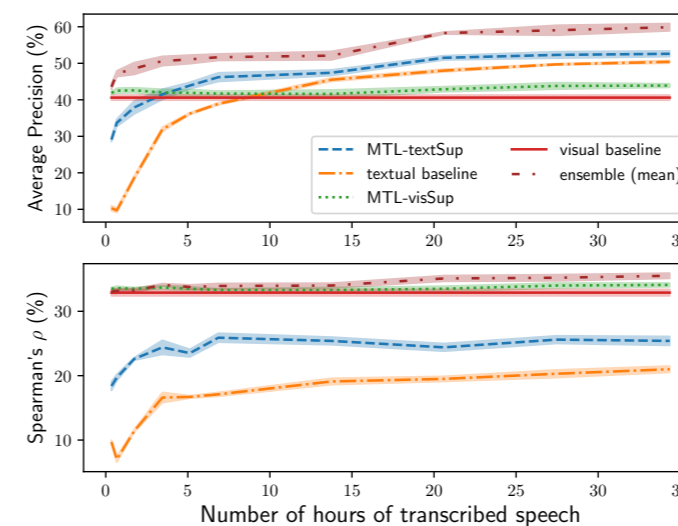
How is the inference done?

Input: Spoken utterances

Output: Scores from either MTL-visSup (y_{vis}) or MTL-textSup (y_{bow}) or a combination of the two

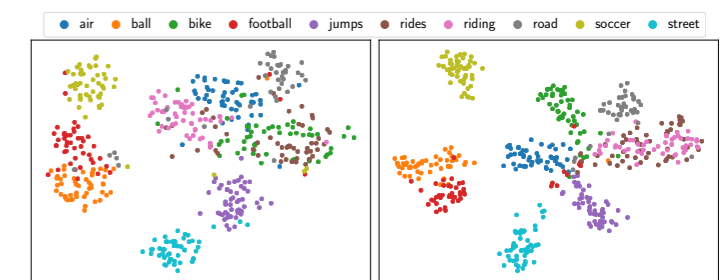
Main Results

The models are evaluated on a corpus of *semantic relevance judgements*⁴ (collected by Kamper et al. [1]).



Additional Observations

- Adding representation loss gives a gain of 7-15% on average precision.
- Higher output dimensionality acts as a regularizer in lower supervision conditions.
- The proposed model outperforms pre-training and hierarchical MTL.
- t-SNE visualization of the learned representations in the text baseline (left) and MTL-textSup (right)



Conclusion

- Visual grounding helps even in the presence of textual supervision.
- Proposed MTL approach significantly improves performance at all levels of supervision.
- Joint training with representation loss helps.

Future Work

Domain extension: Does our visually grounded model perform well on speech not describing visual scenes?

Modify text encoder: Can we explicitly encode semantics in the textual supervision?

References

- [1] D. Harwath et al. Unsupervised Learning of Spoken Language with Visual Context, *NIPS 2016*.
- [2] H. Kamper et al. Semantic Speech Retrieval with a Visually Grounded Model of Untranscribed Speech, *IEEE TASLP, vol. 27, pp. 89-98, 2019*.

¹The Flickr8k Audio Captions Corpus consisting of ~8k images paired with 5 spoken captions each amounting to a total of ~46 hours of speech data (~34 hours training, ~6 hours dev, and ~6 hours test).

²ImageNet pre-trained fixed ResNet followed by fully connected layers trained on the union of MSCOCO and Flickr30k, with ~149k images (~107k training, ~42k dev).

³Written transcripts of the Flickr8k Audio Captions. We use subsets of these transcripts with varying sizes: from just ~21 minutes to the complete ~34 hours of labelled speech.

⁴~1k utterances from the Flickr8k Audio Captions Corpus with their semantic relevance for each of 67 query words. Each (utterance, keyword) pair was labeled by 5 annotators. Majority vote of the annotators ("hard labels") and the actual number of votes ("soft labels") for evaluation.