

cases	doc_1		doc_2				decision	id
			<div>authors<div><ul style="list-style-type: none"><li>Ankita Pasad</li><li>Bowen Shi</li><li>Herman Kamper</li><li>Karen Livescu</li></ul></div></div> <div>title<div>On the Contributions of Visual and Textual Supervision in Low-Resource Semantic Speech Retrieval</div></div> <div>publication_date<div>2019-04-24 17:44:06+00:00</div></div> <div>source<div>SupportedSources.ARXIV</div></div> <div>journal<div>None</div></div> <div>volume<div></div></div> <div>doi<div></div></div> <div>urls<div><ul style="list-style-type: none"><li>http://arxiv.org/pdf/1904.10947v2</li><li>http://arxiv.org/abs/1904.10947v2</li><li>http://arxiv.org/pdf/1904.10947v2</li></ul></div></div> <div>id<div>id3492340363340870377</div></div> <div>abstract<div>Recent work has shown that speech paired with images can be used to learn semantically meaningful speech representations even without any textual supervision. In real-world low-resource settings, however, we often have access to some transcribed speech. We study whether and how visual grounding is useful in the presence of varying amounts of textual supervision. In particular, we consider the task of semantic speech retrieval in a low-resource setting. We use a previously studied data set and task, where models are trained on images with spoken captions and evaluated on human judgments of semantic relevance. We propose a multitask learning approach to leverage both visual and textual modalities, with visual supervision in the form of keyword probabilities from an external tagger. We find that visual grounding is helpful even in the presence of textual supervision, and we analyze this effect over a range of sizes of transcribed data sets. With ~5 hours of transcribed speech, we obtain 23% higher average precision when also using visual supervision.</div></div> <div>versions<div></div></div>	DUPLICATES	345			
	authors	<ul style="list-style-type: none"><li>Pasad, A.</li><li>Shi, B.</li><li>Kamper, H.</li><li>Livescu, K.</li></ul>						
	title	On the Contributions of Visual and Textual Supervision in Low-Resource Semantic Speech Retrieval						
	publication_date	2019-09-15 00:00:00						
	source	SupportedSources.CROSSREF						
	journal							
	volume							
	doi	10.21437/interspeech.2019-3051						
	urls	<ul style="list-style-type: none"><li>http://dx.doi.org/10.21437/interspeech.2019-3051</li></ul>						
	id	id-1606990110734108228						
	abstract							
	versions							