| | | doc_1 | | doc_2 | decision | id |
|---|---|---|---|---|---|---|
| cases | authors | <ul><li>Mayu Otani</li><li>Yuta Nakashima</li><li>Esa Rahtu</li><li>J. HeikkilÃ¤</li><li>N. Yokoya</li></ul> | authors | <ul><li>Mayu Otani</li><li>Yuta Nakashima</li><li>Esa Rahtu</li><li>Janne HeikkilÃ¤</li><li>Naokazu Yokoya</li></ul> | DUPLICATES | 373 |
| | title | Learning Joint Representations of Videos and Sentences with Web Image Search | title | Learning Joint Representations of Videos and Sentences with Web Image Search | | |
| | publication_date | 2016-08-08 00:00:00 | publication_date | 2016-08-08 09:54:15+00:00 | | |
| | source | SupportedSources.SEMANTIC_SCHOLAR | source | SupportedSources.ARXIV | | |
| | journal | | journal | None | | |
| | volume | | volume | | | |
| | doi | 10.1007/978-3-319-46604-0_46 | doi | | | |
| | urls | <ul><li>https://www.semanticscholar.org/paper/aa76f655c2ad655080593a191c4b479ab9f18117</li></ul> | urls | <ul><li>http://arxiv.org/pdf/1608.02367v1</li><li>http://arxiv.org/abs/1608.02367v1</li><li>http://arxiv.org/pdf/1608.02367v1</li></ul> | | |
| | id | id-7750524676858768777 | id | id2121869158855415783 | | |
| | abstract | None | abstract | Our objective is video retrieval based on natural language queries. In addition, we consider the analogous problem of retrieving sentences or generating descriptions given an input video. Recent work has addressed the problem by embedding visual and textual inputs into a common space where semantic similarities correlate to distances. We also adopt the embedding approach, and make the following contributions: First, we utilize web image search in sentence embedding process to disambiguate fine-grained visual concepts. Second, we propose embedding models for sentence, image, and video inputs whose parameters are learned simultaneously. Finally, we show how the proposed model can be applied to description generation. Overall, we observe a clear improvement over the state-of-the-art methods in the video and sentence retrieval tasks. In description generation, the performance level is comparable to the current state-of-the-art, although our embeddings were trained for the retrieval tasks. | | |
| | versions | | versions | | | |