| cases | doc_1 | | doc_2 | | decision | id |
|---|---|---|---|---|---|---|
| | **authors** | • Ryant, N.<br>• Liberman, M.<br>• Yuan, J. | **authors** | • Neville Ryant<br>• M. Liberman<br>• Jiahong Yuan | DUPLICATES | 188 |
| | **title** | Speech activity detection on youtube using deep neural networks | **title** | Speech activity detection on youtube using deep neural networks | | |
| | **publication_date** | 2013-08-25 00:00:00 | **publication_date** | None | | |
| | **source** | SupportedSources.CROSSREF | **source** | SupportedSources.SEMANTIC_SCHOLAR | | |
| | **journal** | | **journal** | | | |
| | **volume** | | **volume** | | | |
| | **doi** | 10.21437/interspeech.2013-203 | **doi** | | | |
| | **urls** | • http://dx.doi.org/10.21437/interspeech.2013-203 | **urls** | • https://www.semanticscholar.org/paper/e651c1ec20460ae0f0fcd95f705370090a3bb335 | | |
| | **id** | id4589494756250490056 | **id** | id-8190058057386089388 | | |
| | **abstract** | | **abstract** | Speech activity detection (SAD) is an important first step in speech processing. Commonly used methods (e.g., frame-level classification using gaussian mixture models (GMMs)) work well under stationary noise conditions, but do not generalize well to domains such as YouTube, where videos may exhibit a diverse range of environmental conditions. One solution is to augment the conventional cepstral features with additional, hand-engineered features (e.g., spectral flux, spectral centroid, multiband spectral entropies) which are robust to changes in environment and recording condition. An alternative approach, explored here, is to learn robust features during the course of training using an appropriate architecture such as deep neural networks (DNNs). In this paper we demonstrate that a DNN with input consisting of multiple frames of mel frequency cepstral coefficients (MFCCs) yields drastically lower frame-wise error rates (19.6%) on YouTube videos compared to a conventional GMM based system (40%). | | |
| | **versions** | | **versions** | | | |