

This is a repository copy of *Acoustic Indicators of Join Quality in Concatenated Video Game Commentary*.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/230212/>

Conference or Workshop Item:

Resti, Luca, Gully, Amelia orcid.org/0000-0002-8600-121X, McLoughlin, Michael Paul et al. (2 more authors) (2025) Acoustic Indicators of Join Quality in Concatenated Video Game Commentary. In: UK and Ireland Speech Workshop 2025, 16-17 Jun 2025.

Reuse

Items deposited in White Rose Research Online are protected by copyright, with all rights reserved unless indicated otherwise. They may be downloaded and/or printed for private study, or other acts as permitted by national copyright laws. The publisher or other rights holders may allow further reproduction and re-use of the full text version. This is indicated by the licence information on the White Rose Research Online record for the item.

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.

Acoustic Indicators of Join Quality in Concatenated Video Game Commentary

Luca Resti^{1,3}, Amelia J. Gully², Michael McLoughlin³, Gavin Kearney³, Alena Denisova¹

¹Department of Computer Science, University of York, United Kingdom

²Department of Language and Linguistic Science, University of York, United Kingdom

³AudioLab, School of Physics, Engineering and Technology, University of York, United Kingdom

{luca.resti, amelia.gully, michael.mcloughlin, gavin.kearney, alena.denisova}
@york.ac.uk

1. Introduction

In sports video games, commentary is sometimes generated by concatenating pre-recorded speech segments to dynamically narrate in-game events. While enabling reactive dialogue, this can introduce perceptual artifacts at join points, such as signal and prosodic discontinuities. Identifying these issues typically relies on manual quality control by audio designers, highlighting the need for data-driven approaches to support the review process. Our dataset consists of 18741 recordings of a professional sports commentator. Utterances are either names (teams or players) or actions (e.g., “passes the ball”). Full utterances are formed by randomly concatenating semantically compatible units on consonants, as per industrial practice. We analyzed 1,500 concatenated and 3,000 non-concatenated sports commentary samples using OpenSMILE’s *ComParE_2016* feature set (6,373 acoustic features) [1]. This analysis lays the foundation for subjective testing and the development of a predictive model to assist the quality assurance process by identifying perceptually problematic joins.

2. Methodology

After standardizing the extracted features, we applied PCA followed by t-SNE and UMAP (as shown in Figure 1). We analyzed feature importance by correlating individual acoustic features with embedding dimensions, performed HDBSCAN clustering [2] on the embeddings, and visualized feature distributions for concatenated and unconcatenated samples.

3. Results

3.1 Feature-Embedding Correlation

Correlation analysis between acoustic features and embedding dimensions shows that spectral and modulation-based features - particularly RASTA-filtered Linear Predictive Coding (LPC) and spectral energy coefficients - strongly influenced the spatial structure of both UMAP and t-SNE projections ($|r| > 0.7$ for top features, as shown in Table 1).

Feature (UMAP)	Corr
audspecRasta_lengthL1norm_sma_lpc0	0.726
audspecRasta_lengthL1norm_sma_lpc1	0.720
audSpec_Rfilt_sma[20]_lpc1	0.711

Table 1: Top features by mean absolute embedding correlation.

3.2 Clustering Analysis

The UMAP projection of acoustic features indicates a strong apparent separation between most concatenated and fluent utter-

This work was conducted in collaboration with Electronic Arts SEED and funded via the Sound Interactions in the Metaverse Centre for Doctoral Training and the EPSRC.

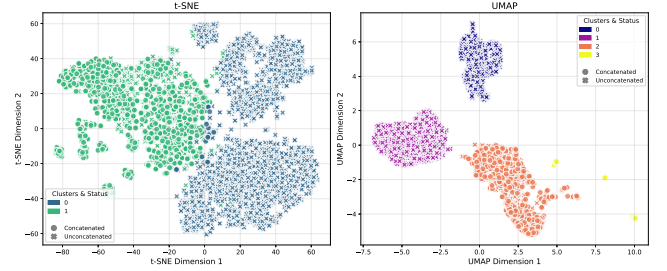


Figure 1: t-SNE and UMAP projections of OpenSMILE ComParE_2016 features.

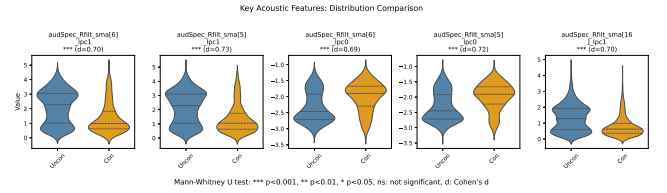


Figure 2: Top feature distribution across concatenated and unconcatenated speech.

ances. HDBSCAN clustering in UMAP space identified 4 clusters. Clusters 0 and 1 consist almost entirely of fluent speech, while cluster 3 is > 99% concatenated and exhibits high z-scores (relative to the dataset mean) in spectral variance and roll-off features. Cluster 2, which contains both fluent and concatenated samples, shows divergent trends in RASTA-filtered LPC-based features (e.g., +0.92 for $lpc1$ vs -0.91 for $lpc0$), indicating deviations in the spectral envelope.

3.3 Feature Distribution

We selected a subset of features with the highest correlation with UMAP dimensions and analyzed their distribution across samples and HDBSCAN clusters. As shown in Figure 2, top-ranked features display significant shifts between concatenated and fluent samples.

1. References

- [1] F. Eyben, M. Wöllmer, and B. Schuller, “openSMILE – the munich versatile and fast open-source audio feature extractor,” in *Proceedings of the 18th ACM International Conference on Multimedia (MM)*, 2010, pp. 1459–1462.
- [2] L. McInnes, J. Healy, and S. Astels, “hdbscan: Hierarchical density based clustering,” *The Journal of Open Source Software*, vol. 2, no. 11, p. 205, 2017.