

# DCU at the NTCIR-11 SpokenQuery&Doc Task

David N. Racca

Gareth J.F. Jones

CNGL Centre for Global Intelligent Content, School of Computing, Dublin City University, Dublin, Ireland  
 {dracca, gjones}@computing.dcu.ie

## Introduction

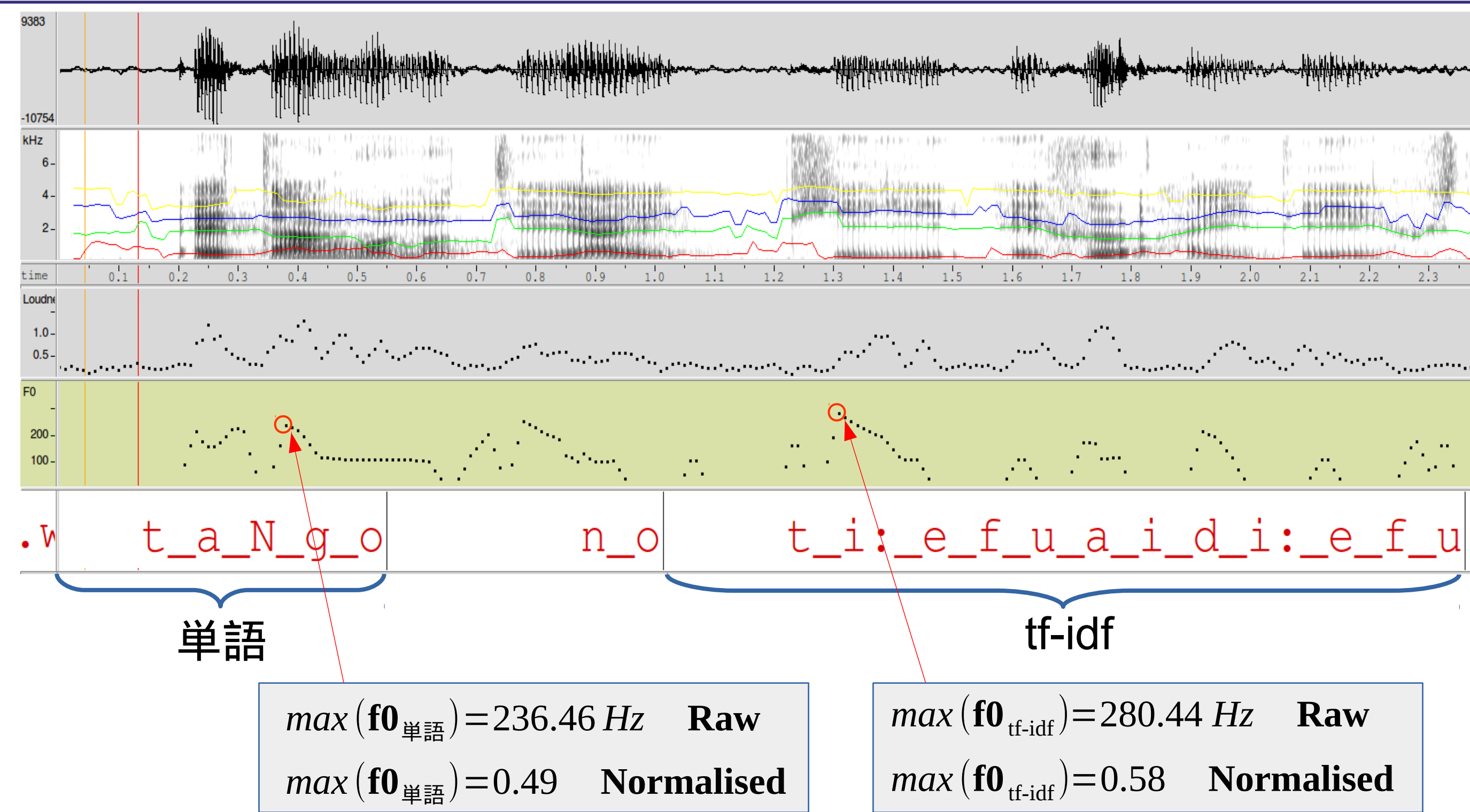
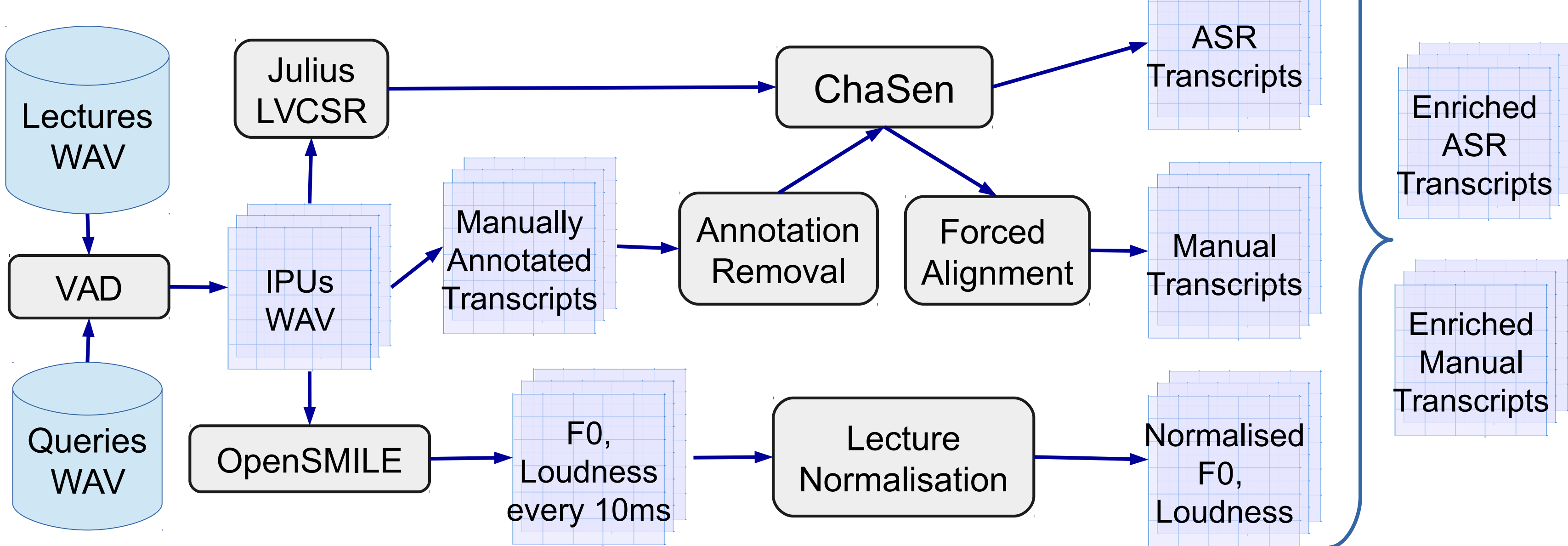
- Speech is more than a simple sequence of words.
- Prosodic variation encode rich information about:
  - emotions, discourse structure, dialogue acts, focus, emphasis, contrast, topic shifting, etc.
- We examined the potential of prosodic prominence in the NTCIR-11 SpokenQuery&Doc Task.

## Background and Previous Work

Prosody may be useful in speech search:

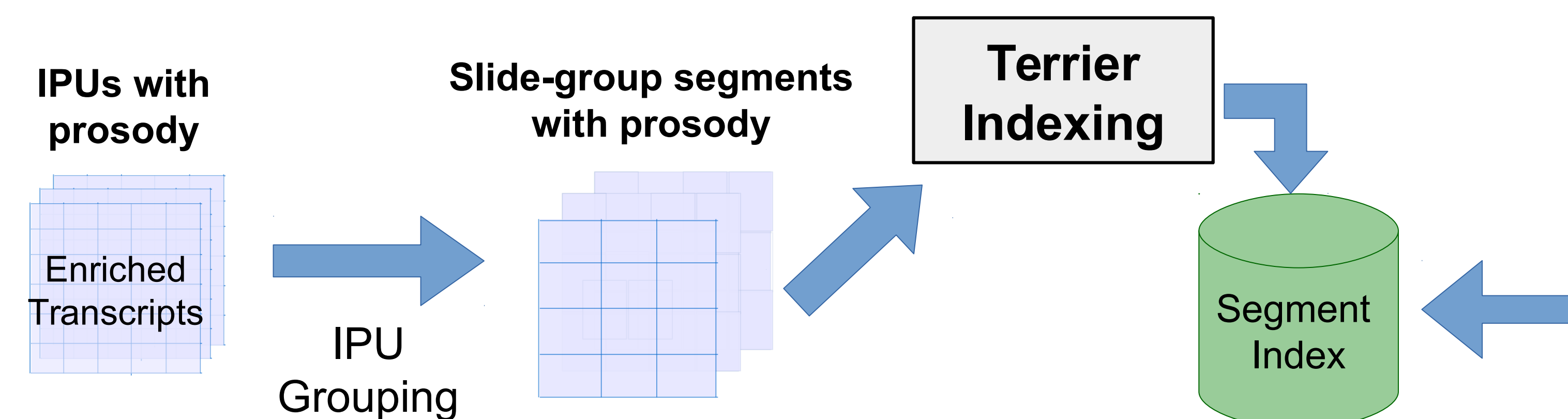
- Relationship between stress and TF-IDF scores [1].
- SDR exploiting amplitude and duration [2].
- Topic tracking exploiting energy and pitch [3].
- SCR exploiting pitch, loudness, and duration [4].

## Data Pre-processing



## Indexing

- Stores normalised prosodic features for each term.



## Retrieval

- Increases weights of prominent terms.

### Terrier Matching

$$\text{rel}(q, s_j) = \sum_i w(i, j)$$

$$w(i, j) = \begin{cases} \text{idf}(i, C) [\alpha \cdot \text{tf}(i, j) + (1 - \alpha) \text{ac}(i, j)] & \text{LI} \\ \frac{\theta_{ir} \cdot \text{tf}(i, j) \cdot \text{idf}(i, C) + \theta_{ac} \cdot \text{ac}(i, j)}{\theta_{ir} + \theta_{ac}} & \text{G} \\ \text{tf}(i, j) \cdot \text{idf}(i, C) & \text{TF\_IDF} \end{cases}$$

$$\text{ac}(i, j) = \begin{cases} \text{f0}(i, j) & \text{Pitch [P]} \\ \text{l}(i, j) & \text{Loudness [L]} \\ \text{d}(i, j) & \text{Duration [Dur]} \\ \text{f0}_{\text{range}}(i, j) & \text{Pitch Range [Pr]} \\ \text{l}(i, j) \cdot \text{f0}(i, j) & \text{[LP]} \\ \text{l}(i, j) \cdot \text{f0}_{\text{range}}(i, j) & \text{[LPr]} \end{cases}$$

$$\text{f0}(i, j) = \max_k \{ \max(\mathbf{f0}_{i,j}^k) \}$$

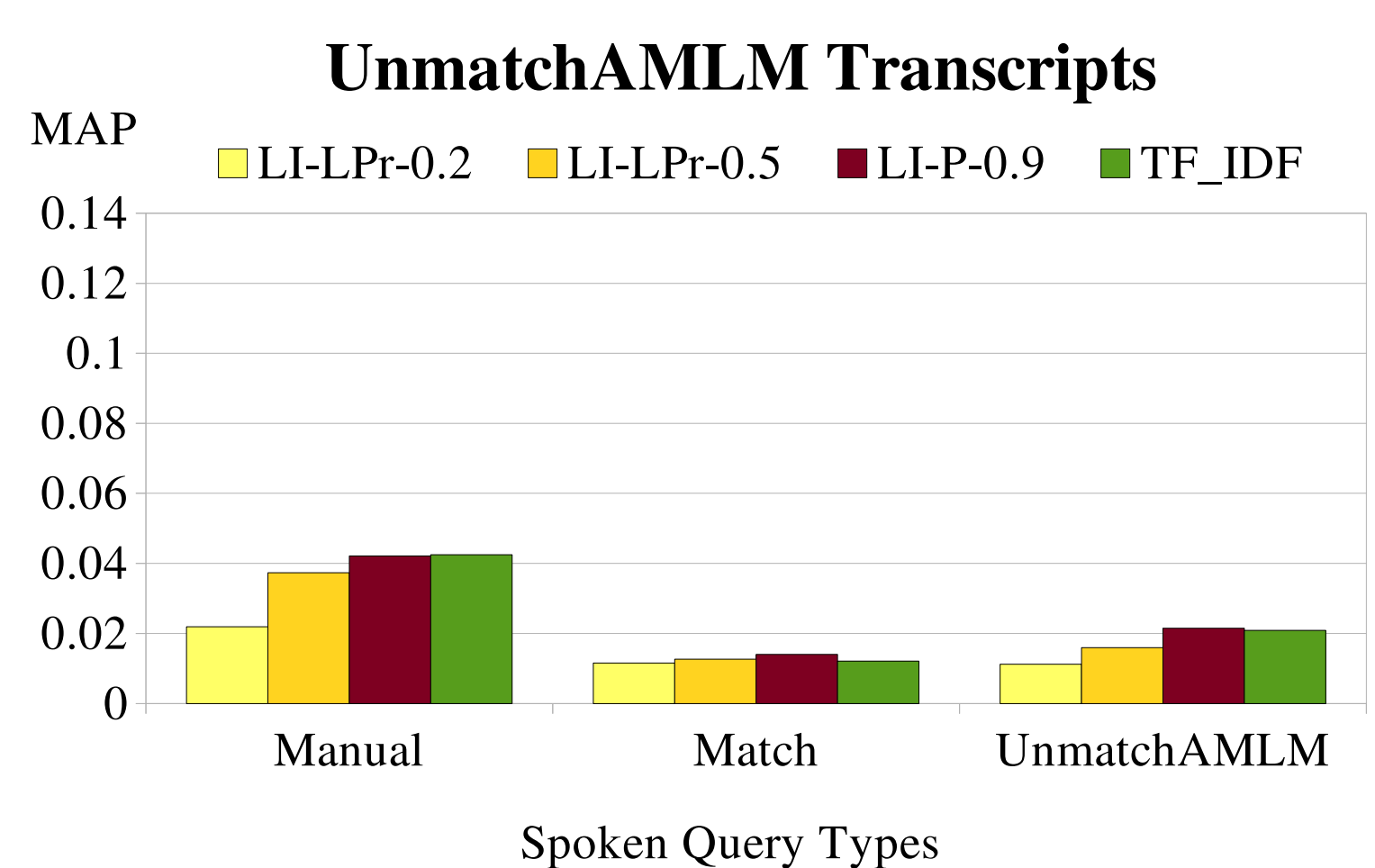
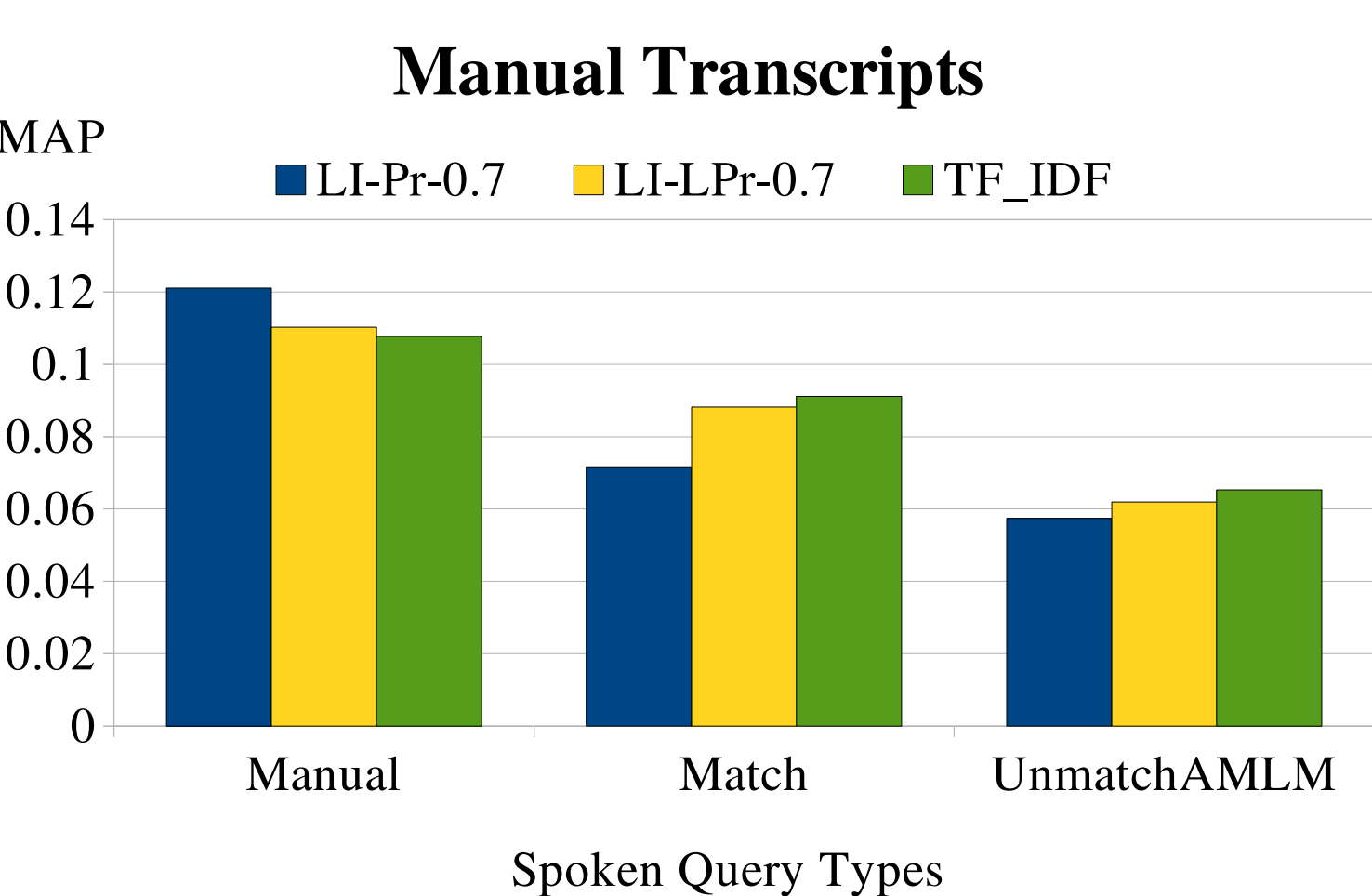
$$\text{l}(i, j) = \max_k \{ \max(\mathbf{l}_{i,j}^k) \}$$

$$\text{d}(i, j) = \max_k \{ \mathbf{d}_{i,j}^k \}$$

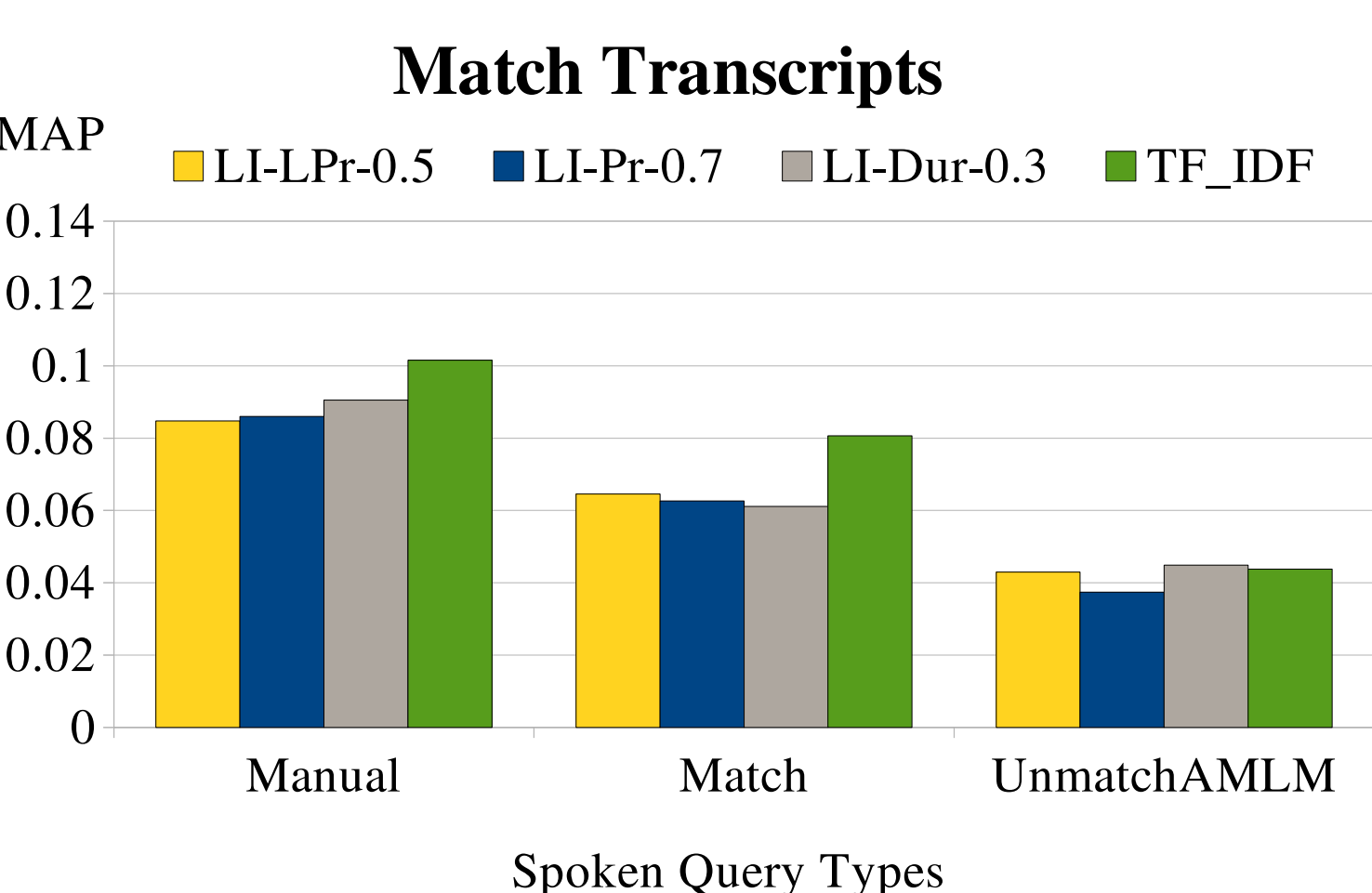
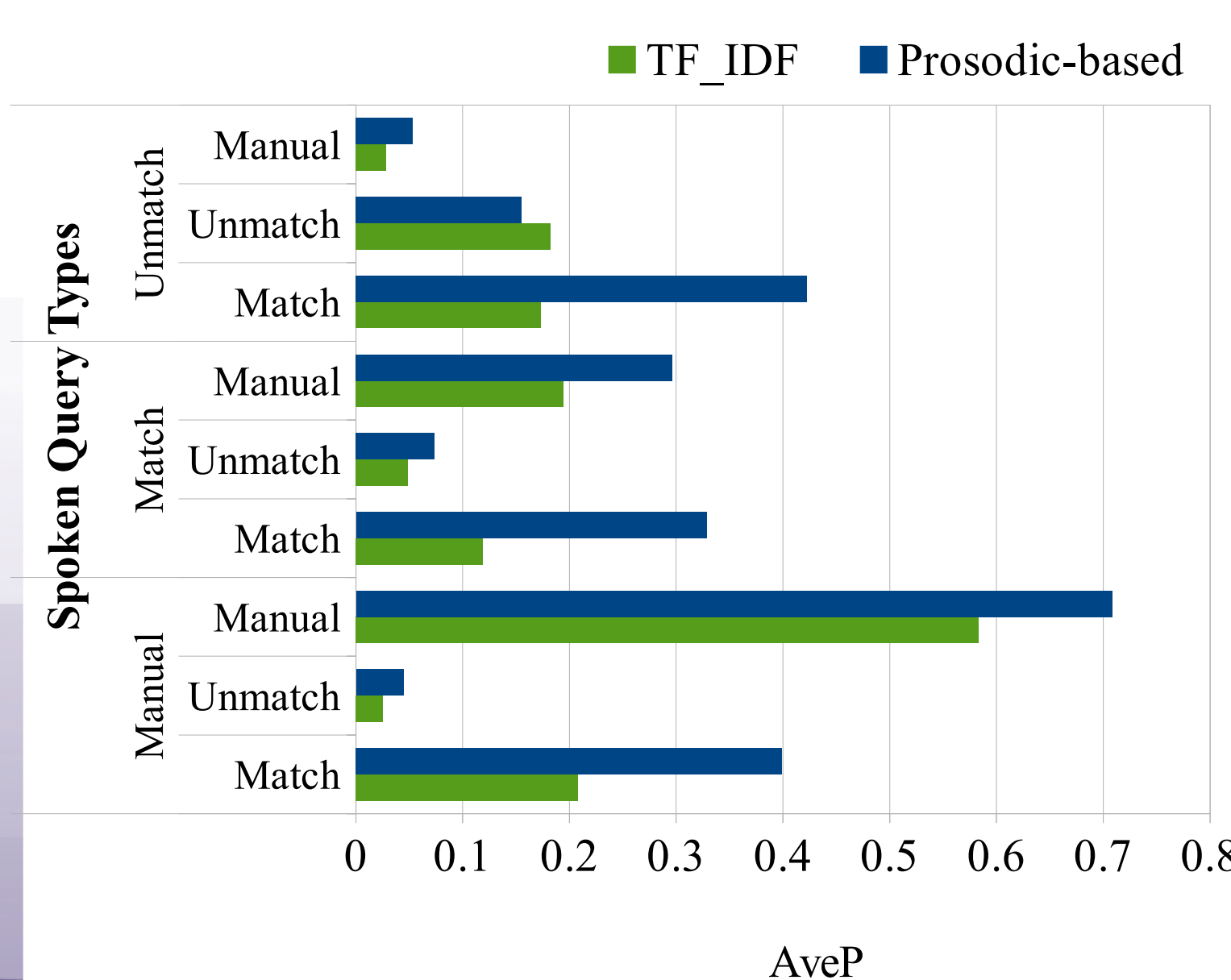
$$\text{f0}_{\text{range}}(i, j) = \max_k \{ \max(\mathbf{f0}_{i,j}^k) \} - \min_k \{ \min(\mathbf{f0}_{i,j}^k) \}$$

$$\text{tf}(i, j) = \frac{k_1 \text{tf}_{i,j}}{\text{tf}_{i,j} + k_1 \left( 1 - b + b \frac{\text{dl}_j}{\text{avdl}} \right)} \quad \text{idf}(i, C) = \log \left( \frac{N}{n_i} + 1 \right)$$

## Results



Query 1: Prosodic-based vs TF\_IDF



## Conclusions

- No significant differences between prosodic and text based runs.
- Transcript quality affects retrieval effectiveness.
- Prosodic-based models may be useful for certain queries.

## References

[1] F. Crestani. *Towards the use of prosodic information for spoken document retrieval*. SIGIR'01, 2001.  
 [2] B. Chen et al. *Improved spoken document retrieval by exploring extra acoustic and linguistic cues*. INTERSPEECH'01, 2001.  
 [3] C. Guinaudeau et al. *Accounting for prosodic information to improve ASR-based topic tracking for TV broadcast news*. INTERSPEECH'11, 2011.  
 [4] D.N. Racca et al. *DCU search runs at MediaEval 2014 Search and Hyperlinking*. MediaEval 2014 Multimedia Benchmark Workshop, 2014.