

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



Overview

- We participated in the **slide-group SQ-SCR**.
- General idea:
 - Augment text-retrieval methods with prosodic features: pitch (F0), loudness, and duration.
 - Compute an **acoustic score** for each term.
 - Promote the rank of segments containing acoustically prominent terms.

Motivation

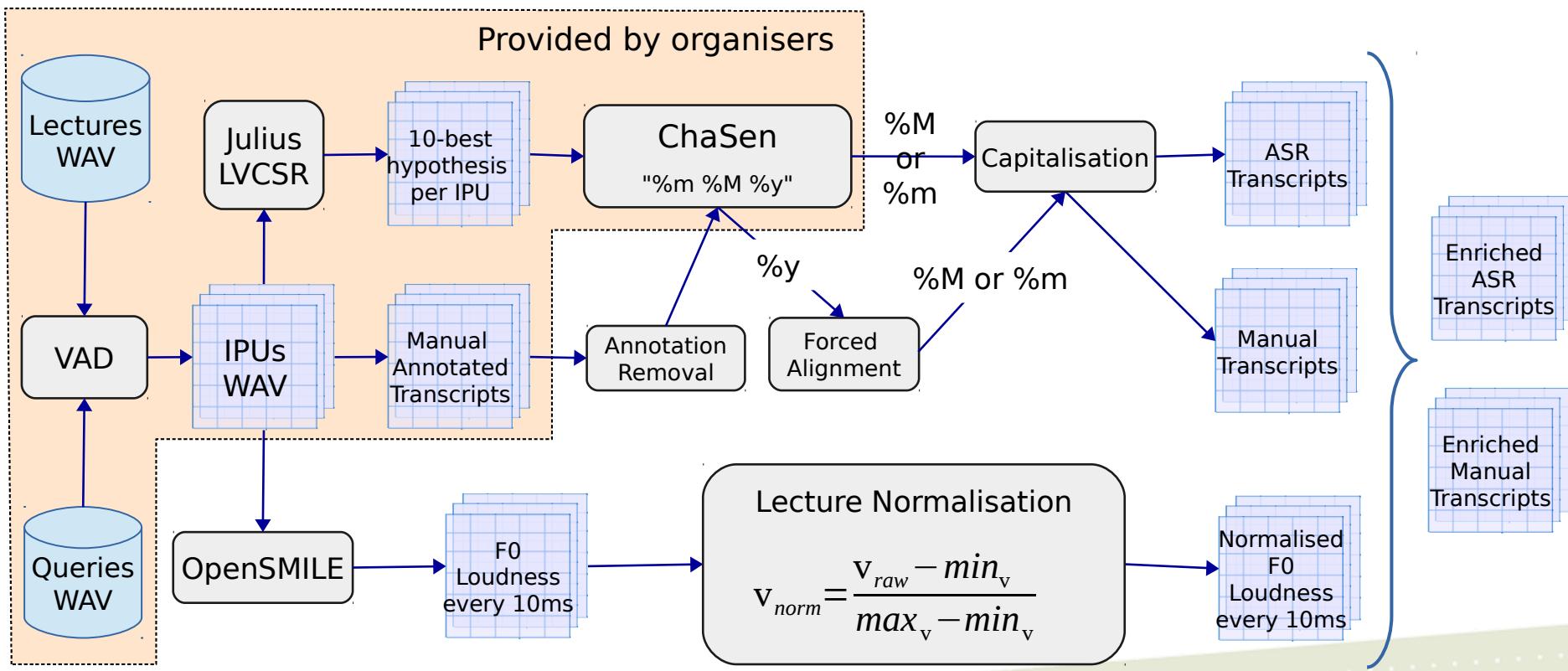
- Prosody:
 - Rhythm, stress, intonation, duration, loudness.
- Shown useful in many speech processing tasks:
 - Emotions, discourse structure, speech acts, speaker ID, topic segmentation.
- Prominent speech units stand-out from their context.
- Information status: old vs new information.

Related Work

- **Crestani [1]**: possible correlation between acoustic stress and TF-IDF scores (English).
- **Chen et al [2]**: signal amplitude and duration in a spoken document retrieval (SDR) task (Mandarin).
- **Guinaudeau [3]**: F0 and RMS energy in a topic tracking task (French).
- **Racca et al [4]**: F0, loudness, and duration in SCR (English).

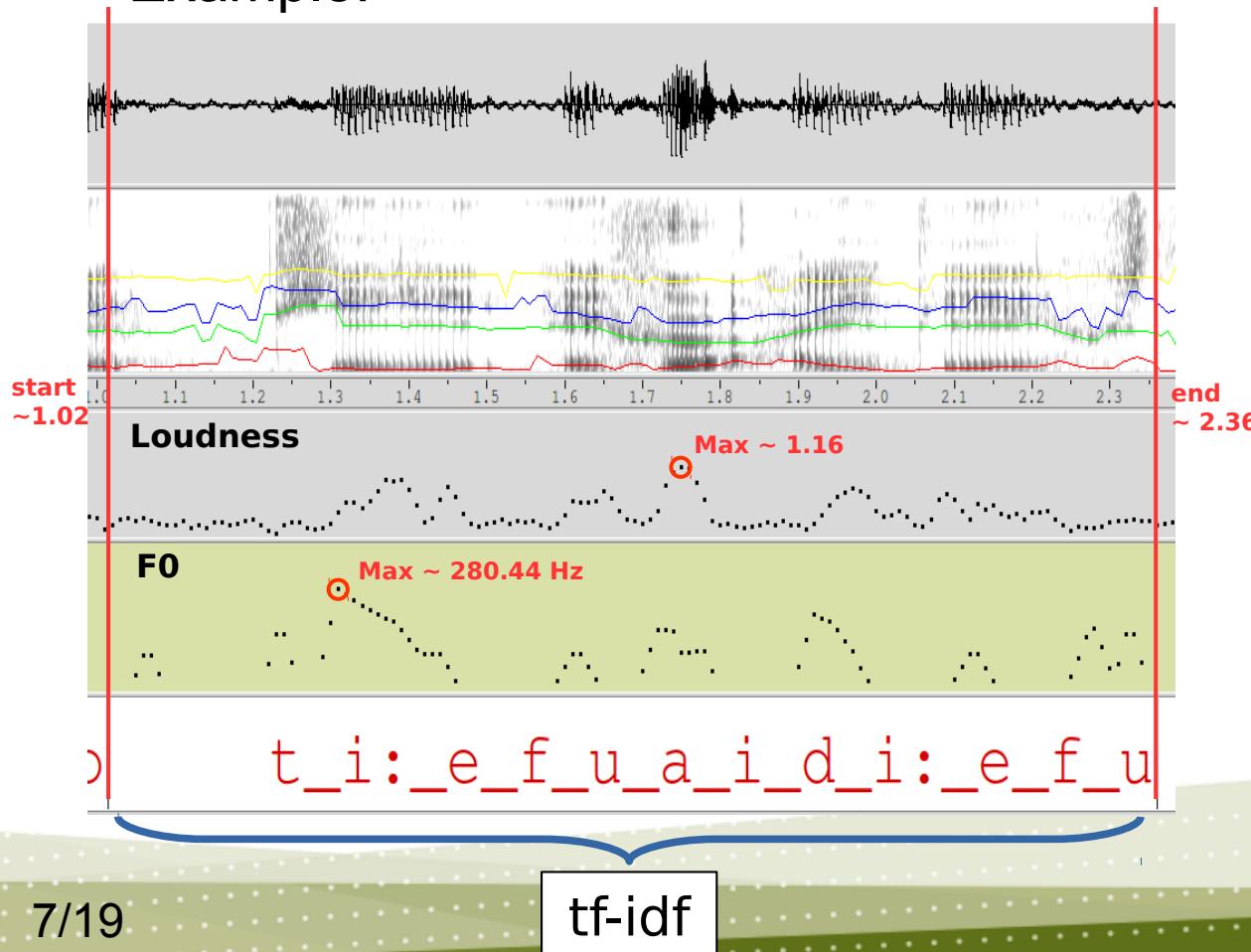
Data Pre-processing

- 1-best WORD *match*, *unmatchAMLM*, and manual transcripts.



Prosodic Features

- Raw duration, lecture-normalised F0 and loudness.
- Example:



Duration

$$d = 2.36 s - 1.02 s = 1.34 s$$

Lecture Normalisation

$$v_{norm} = \frac{v_{raw} - min_v}{max_v - min_v}$$

Loudness

$$\max(\mathbf{l}_{i,j}^k) = 1.16 \quad \text{Raw}$$

$$\max(\mathbf{l}_{i,j}^k) = 0.37 \quad \text{Normalised}$$

Pitch (F0)

$$\max(\mathbf{f0}_{i,j}^k) = 280.44 \text{ Hz} \quad \text{Raw}$$

$$\max(\mathbf{f0}_{i,j}^k) = 0.58 \quad \text{Normalised}$$

Prosodic Features

—F0, loudness, and duration for the term “*i*” term in segment “*j*”.

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

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

$$d(i, j) = \max_k \left\{ d_{i, j}^k \right\}$$

$$f0_{range}(i, j) = \max_k \left\{ \max \left(\mathbf{f0}_{i, j}^k \right) \right\} - \min_k \left\{ \min \left(\mathbf{f0}_{i, j}^k \right) \right\}$$

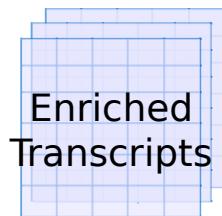
Acoustic Score

- We experimented with six definitions for the acoustic score of term “*i*” in segment “*j*”.

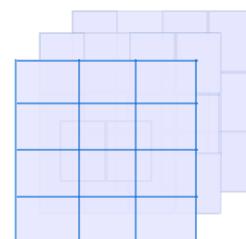
$$ac(i, j) = \begin{cases} f0(i, j) & \textbf{Pitch [P]} \\ l(i, j) & \textbf{Loudness [L]} \\ d(i, j) & \textbf{Duration [Dur]} \\ f0_{\text{range}}(i, j) & \textbf{Pitch Range [Pr]} \\ l(i, j) \cdot f0(i, j) & \textbf{[LP]} \\ l(i, j) \cdot f0_{\text{range}}(i, j) & \textbf{[LPr]} \end{cases}$$

Indexing

IPUs with Prosody

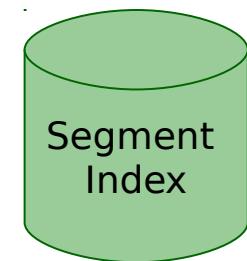
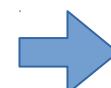


Slide-group segments with Prosody



IPU Grouping

Terrier Indexing



Segment Index

- Slide-group segments indexed using Terrier IR Framework.
- Index stores F0, loudness and duration for each term occurrence along with text statistics.

Retrieval

- Probabilistic model with BM25 weighting:

$$\text{rel}(\mathbf{q}, \mathbf{s}_j) = \sum_i^M w(i, j)$$

- Three definitions for $w(i, j)$ were explored:

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

$$\text{idf}(i, C) = \log \left(\frac{N}{n_i} + 1 \right)$$

$$\text{tf}(i, j) = \frac{k_1 \cdot \text{tf}_{i,j}}{\text{tf}_{i,j} + k_1 \left(1 - b + b \frac{\text{dl}_j}{\text{avdl}} \right)}$$

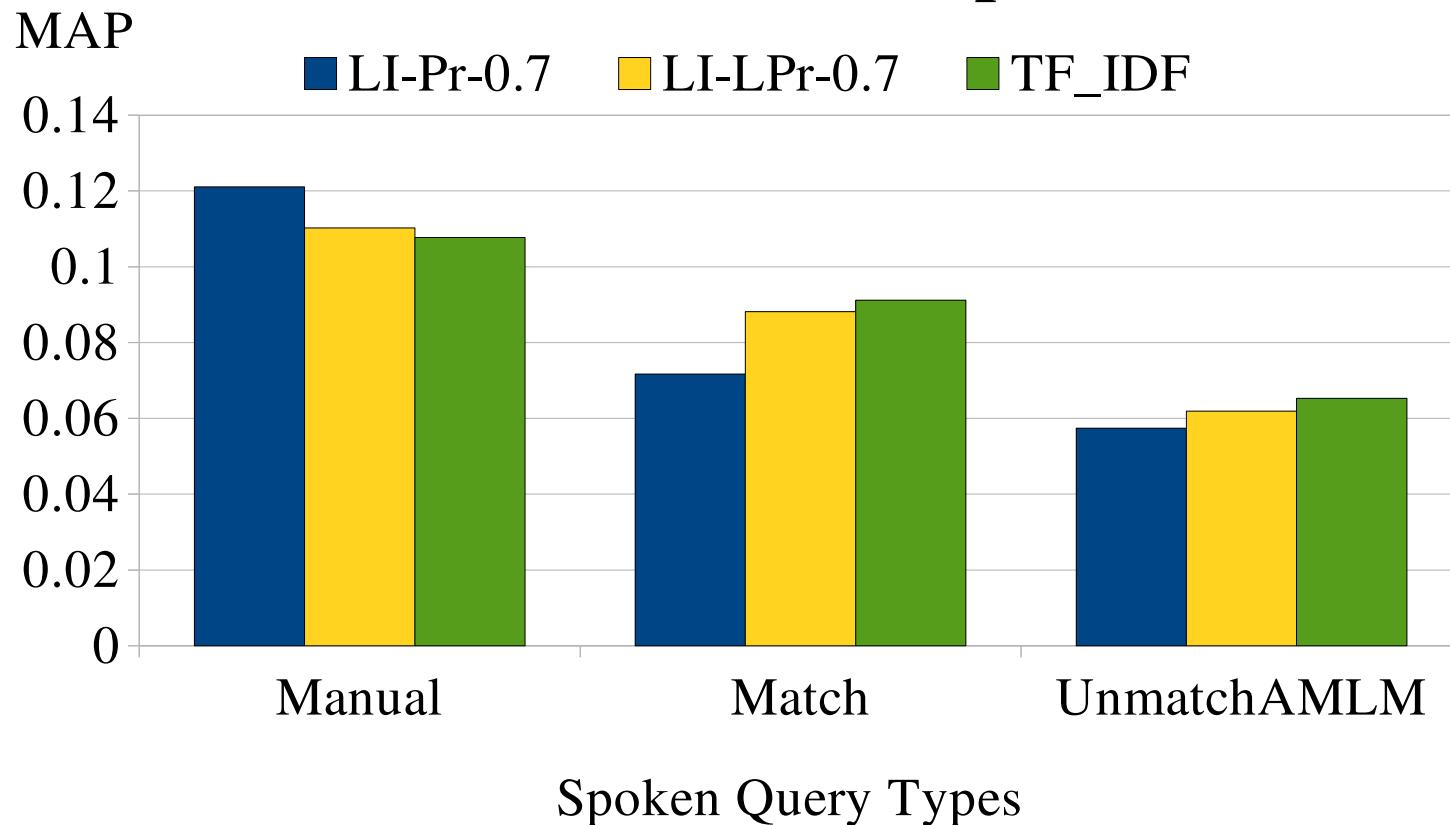
Parameter Tuning

— SpokenDoc-2 passage retrieval: 120 text queries

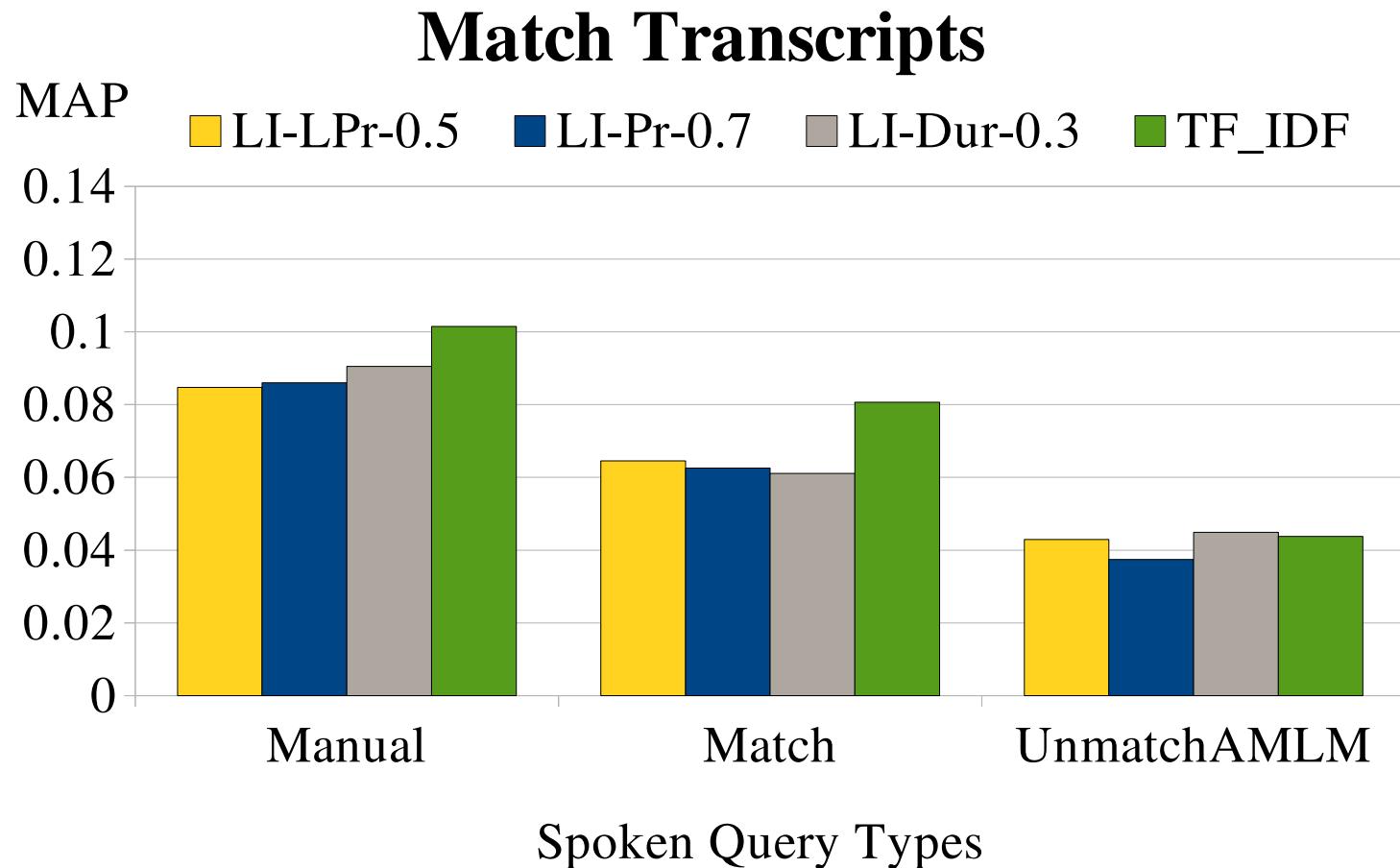
Lecture Transcript	w(i, j)	ac(i, j)	α	θ_{ir}	θ_{ac}	uMAP	pwMAP	fMAP
Manual	LI	LPr	0.7			.1369	.0976	.1005
	LI	Pr	0.7			.1369	.0951	.0995
	G	LP		1	1	.1326	.0960	.0989
	TF-IDF					.1270	.0950	.0972
Match	LI	LPr	0.5			.0842	.0508	.0524
	LI	Dur	0.3			.0819	.0498	.0521
	G	Pr		1	1	.0786	.0473	.0499
	LI	Pr	0.7			.0778	.0490	.0501
	TF-IDF					.0682	.0477	.0486
UnmatchAMLM	G	P		3	1	.0288	.0208	.0131
	LI	LP	0.5			.0278	.0210	.0135
	LI	LPr	0.2			.0271	.0205	.0132
	LI	P	0.9			.0227	.0206	.0129
	TF-IDF					.0222	.0203	.0128

Results: SpokenQuery&Doc

Manual Transcripts

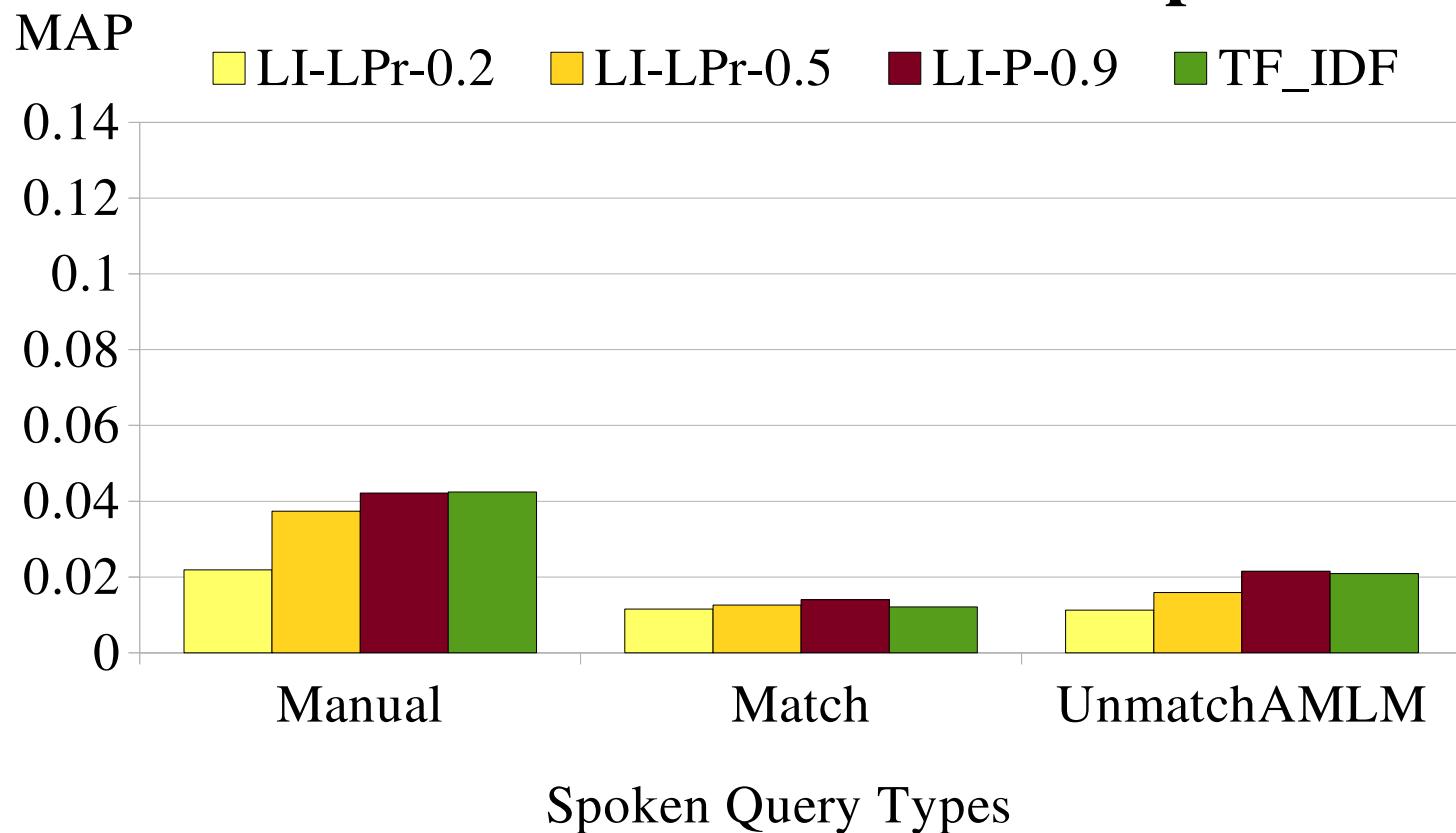


Results: SpokenQuery&Doc



Results: SpokenQuery&Doc

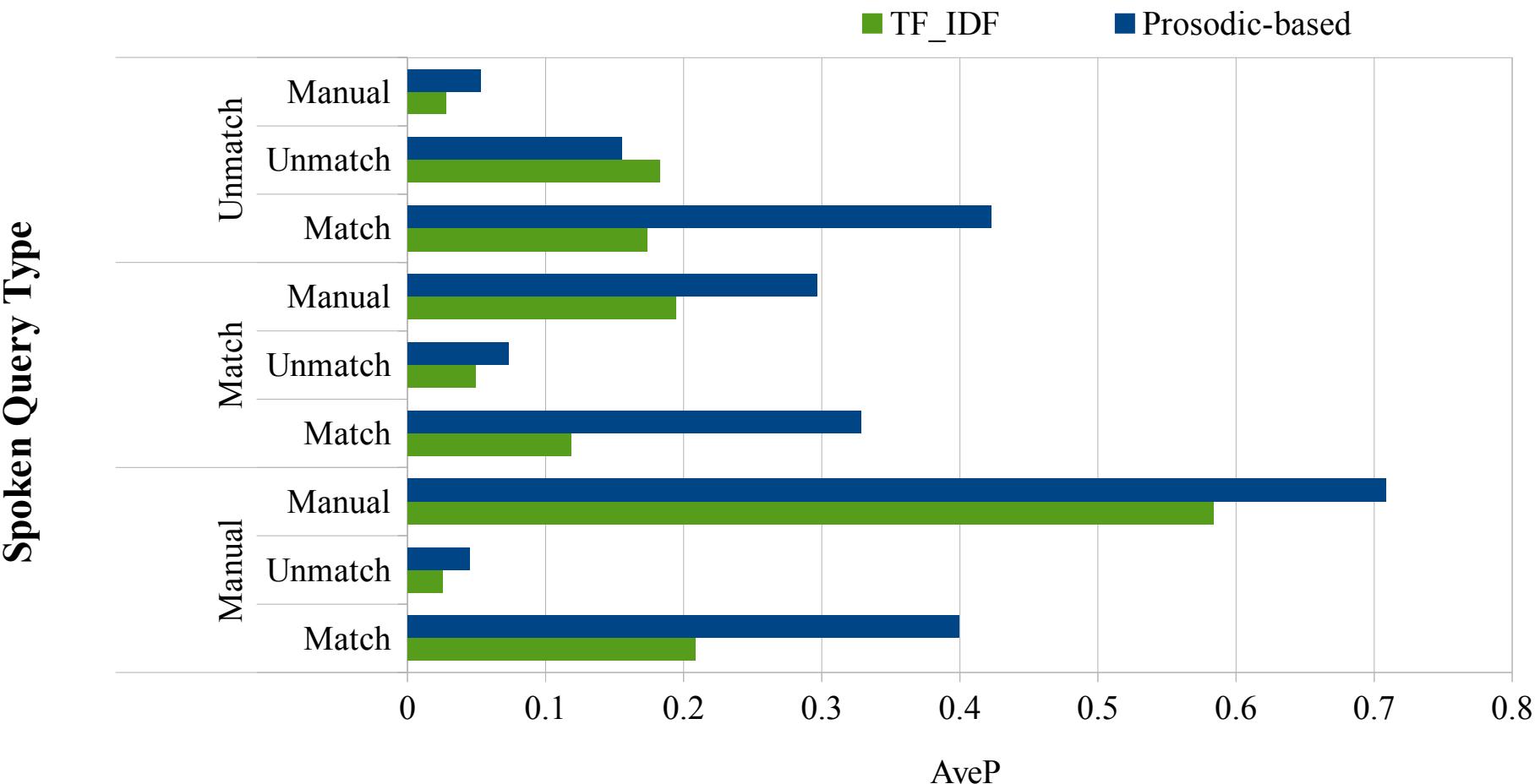
UnmatchAMLM Transcripts



Results: SpokenQuery&Doc

2 relevant segments

Query 1: Prosodic-based vs TF_IDF



Conclusions & Further Work



- Continued exploring if prosodic prominence can be used to improve retrieval effectiveness.
- No significant differences between prosodic and text based runs (t student's test $\sim 95\%$ conf. level).
- Transcript quality affects retrieval effectiveness.
- Prosodic-based models may be useful for some queries/target segments:
 - Future work: predict when this happens.

References

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

Questions?