



# TOWARD THE USE OF INFORMATION DENSITY BASED DESCRIPTIVE FEATURES IN HMM BASED SPEECH SYNTHESIS

Sébastien Le Maguer<sup>1</sup>, Bernd Möbius<sup>1</sup>, Ingmar Steiner<sup>1,2</sup>

<sup>1</sup> Computational Linguistics and Phonetics, Saarland University, Saarbrücken, Germany, <sup>2</sup>DFKI

## Introduction

### • Background

- Statistical TTS = huge effort assigned to acoustic modelling
- Descriptive feature set = almost the same for each system (the one presented in [1])

### • Problem

- How to enrich this descriptive feature set?

### • Proposition

- **New descriptive feature = unpredictability of an event**
- Based on information density & widely used in computational linguistics

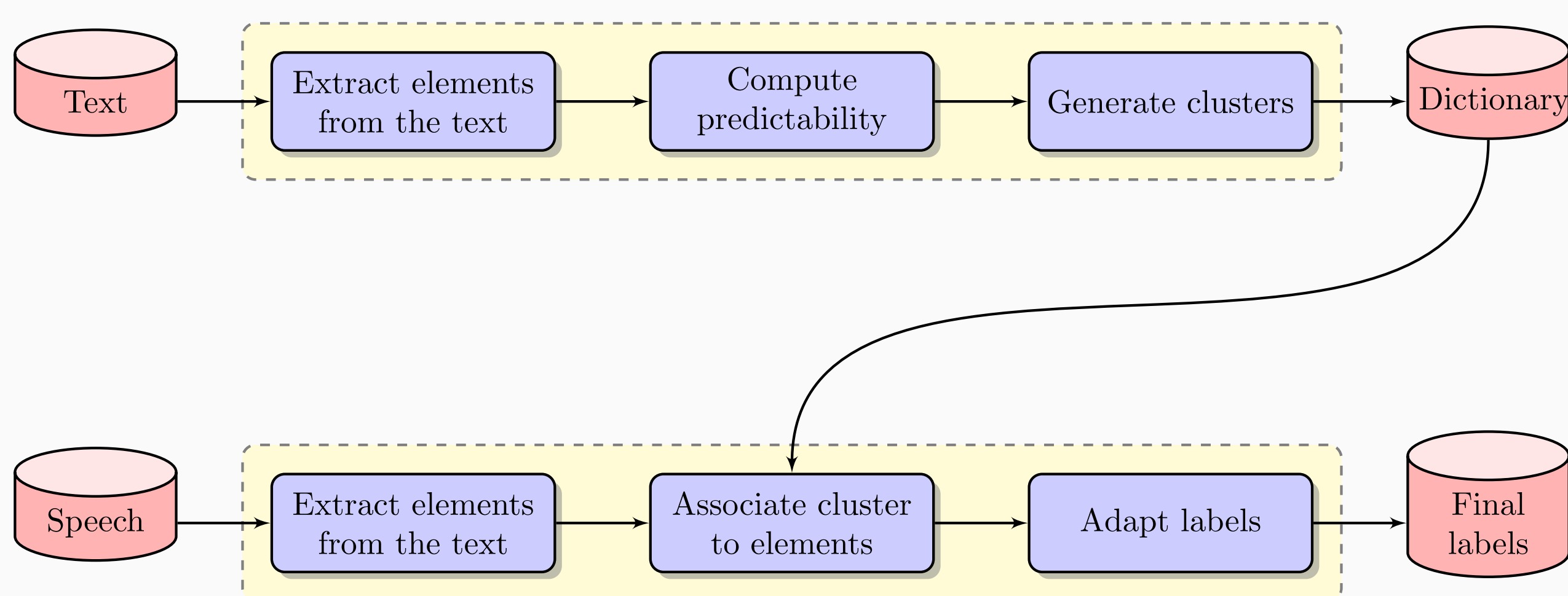
## Unpredictability (Surprisal)

### Once upon a ?

$$\text{Suprisal}(U_i) = -\log_2(P(U_i|U_{i-1}..U_{i-1-N})) \quad (1)$$

- Predictability of a word correlates with processing effort of pronouncing this word [2]
- Same correlation found at the syllable level [3].

## Process



## Feature Generation

### Problem: homogeneous representation (text/speech)

- **Syllable based**
  - IPA phoneme representation
- **Word based**
  - All **punctuation marks** are **discarded**;
  - A break mark is inserted at the **end of each paragraph**;
  - All words are converted to **lower case**

## Objective Evaluation

### Experimental Setup

#### • Speech corpus

- From “Black Beauty” (2013 Blizzard Challenge),
- 1 h (~470 utterances) = 13 522 syl., 7038 words,
- Segmented using EHMM + manually corrected

#### • Text corpus

- 2013 Blizzard Challenge – “Black Beauty”
- 82 books = 951 316 syl., 1 973 368 words

#### • System setup

- HTS 2.2 standard configuration
- Vocoder = STRAIGHT + MLSA filter
- MGC (50) + LF0 (1) + BAP (25) + Δ + ΔΔ

#### • 6 new descriptive features

- syl-unpredictability (P, C, N)
- word-unpredictability (P, C, N)

#### • 3 conditions

- *baseline*
- *unpred\_syllable*
- *unpred\_all*

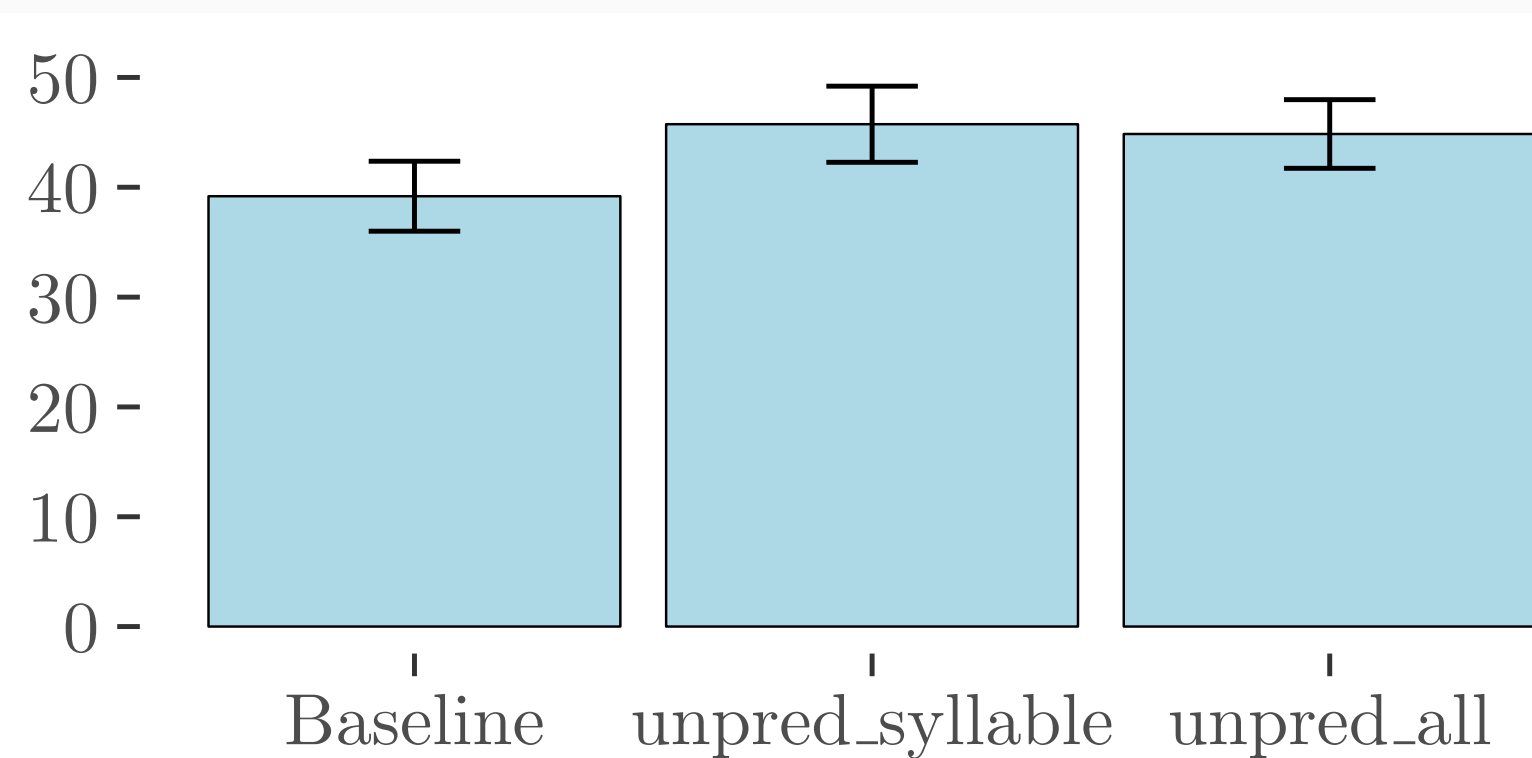
Condition	MCD	RMS-F0	VER	RMS-dur
<i>baseline</i>	6.45	475	15	11.1
<i>unpred_syllable</i>	6.33	463	14.6	10.6
<i>unpred_all</i>	6.33	467	14.8	10.4

### Tree Analysis

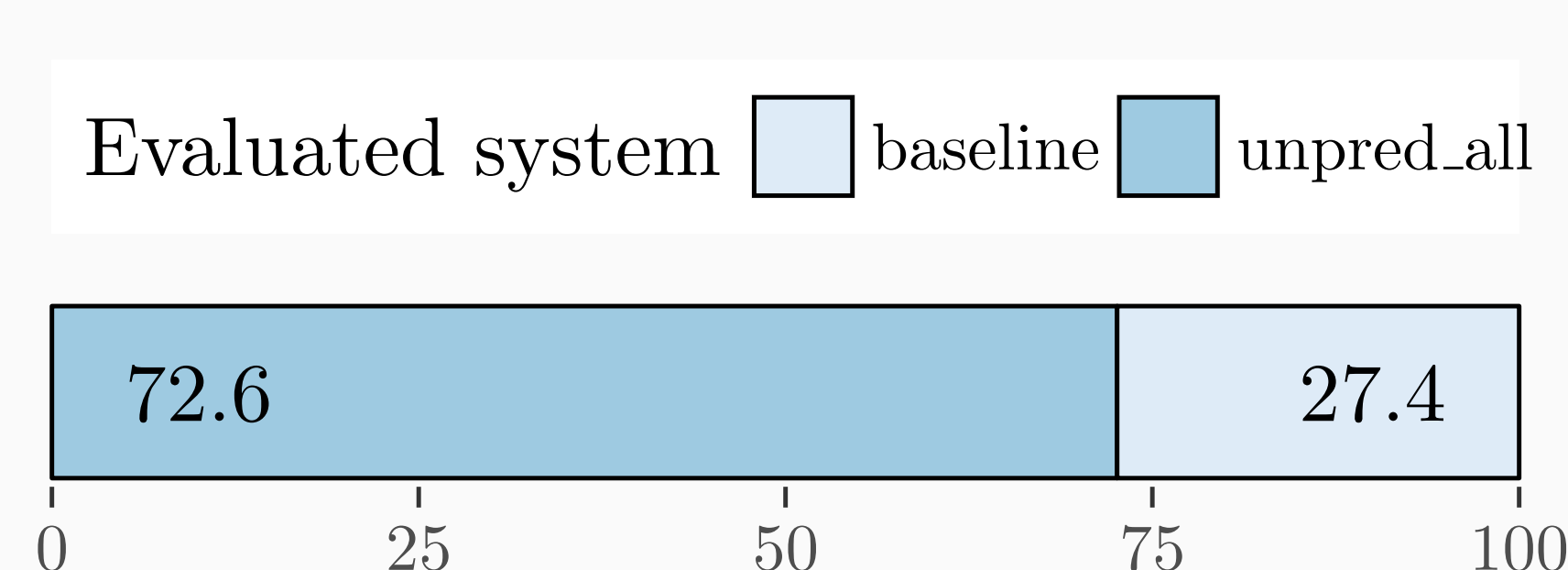
categories	baseline	unpred_syl.	unpred_all
<b>p1</b>	1648	1615	1594
<b>p3</b>	2174	2037	1175
<b>p5</b>	0	0	694
<b>syl-position</b>	<b>5799</b>	<b>5652</b>	<b>4056</b>
<b>syl-prosody</b>	98	128	183
<b>syl-unpredictability</b>	0	1657	<b>4163</b>
<b>word-position</b>	<b>7836</b>	<b>6787</b>	<b>4202</b>
<b>word-prosody</b>	2928	2817	2188
<b>word-unpredictability</b>	0	0	<b>7834</b>
<b>phrase-position</b>	1184	1573	802
<b>phrase-prosody</b>	<b>8723</b>	<b>8323</b>	<b>5892</b>
<b>utterance</b>	<b>7260</b>	<b>7429</b>	<b>6799</b>

## Subjective Evaluation

### MUSHRA Evaluation



### Preference Evaluation



### Analysis

- **AB**
  - Clear preference for the proposed system
- **MUSHRA**
  - Improvement ⇒ just a tendency
  - Evaluation: Spectrum vs. prosody?
- **Global**
  - Assumption = spectrum not impacted, prosody + natural

## Conclusion

- New descriptive feature: unpredictability (widely used in computational linguistics)
- Full process to compute and apply these features
- **Objective analysis**
  - Similarity not impacted
  - Model takes this feature into account
- **Subjective evaluation:** preference for our system ⇒ which dimension?

## Bibliography

- [1] K. Tokuda, H. Zen, and A. W. Black, “An HMM-based speech synthesis system applied to English,” in *Proceedings of the Speech Synthesis Workshop (SSW)*, 2002.
- [2] M. Kutas, K. A. DeLong, and N. J. Smith, “A look around at what lies ahead: Prediction and predictability in language processing,” in *Predictions in the Brain: Using Our Past to Generate a Future*, M. Bar, Ed. Oxford University Press, 2011, pp. 190–207.
- [3] T. F. Jaeger, “Redundancy and reduction: speakers manage syntactic information density,” *Cognitive Psychology*, vol. 61, pp. 23–62, 2010.