

TTS-System “Papageno”

Overview

- i **Scaleable speech synthesis engine**
- i **Unified kernel for embedded and server solutions**
- i **Covers applications from 5 MB to 100 MB footprint**
- i **Multi-lingual system due to separation of engine and speaker/language dependent knowledge bases**
- i **3 main parts:**
 - i **Preprocessing: tokeniser, tagger, grapheme-to-phoneme**
 - i **Prosody: symbolic and acoustic part**
 - i **Acoustic: unit selection and concatenation**



**Professional
Speech
Processing**

Knowledge bases – Tokeniser

- i **Language dependent**
- i **Different types of characters:**
 - i **Separators: space, tabulator, new line ...**
 - i **Numbers: digits, floating point (comma), ordinal point (th, nd, rd, th), dates (colon, point), relations (slash, colon) ...**
 - i **Words: graphemes, dashes, points (abbreviations)**
 - i **Punctuation marks: brackets, points, comma, slashes ...**
- i **Training material for handling of numbers (date, time, ordinal – cardinal, relations)**
- i **Handling of abbreviations, acronyms**



Knowledge bases – Tagger

- i **Language dependent**
- i **Part-of-speech (POS) tagging**
- i **Additional features (gender, number, case, definiteness), esp. for German ordinal numbers “der 1. (erste) Versuch” – “ein 1. (erster) Versuch”**
- i **Requirements: lexicon and tagged text**
- i **Features as detailed as possible**
- i **Training of neural network or n-gram**



**Professional
Speech
Processing**

Knowledge bases – G2P

- i **Language dependent**
- i **Pronunciation of words, abbreviations, acronyms and numbers**
- i **Requirements: pronunciation dictionary**
- i **Features:**
 - i **Huge dictionary**
 - i **Special handling of foreign words (different phoneme sets, mapping between the sets)**
 - i **List of proper names (or special label in the dictionary)**
 - i **Distinguishing features for homographs (POS, semantic information)**
- i **List of syllable cores and stressable phonemes**
- i **G2P of words by neural networks (first : phoneme sequence incl. syllable boundaries; second: word stress)**
- i **Pronunciation of numbers by graphs trained on examples**
- i **Abbreviations and acronyms with dictionaries**



Knowledge bases – symbolic prosody

- i **Speaker dependent (but may be used for other speakers too)**
- i **Position of phrase accents and breaks**
- i **Requirements: labelled text**
- i **Features: one break (B3) and one accent level (PA)**
- i **Labels manually by ONE expert, given the acoustic representation of the text**
- i **inter- and intra-labeller variance too big**
- i **Training of neural networks (first: position of breaks; second: position of accent within one phrase)**



Knowledge bases – text

- i **Should contain every type of sentences (declaration, question, exclamation)**
- i **Long and short sentences**
- i **Consistency between written and spoken text:**
 - i **Abbreviations**
 - i **Cardinal numbers (100 can be “hundert” or “einhundert”)**
 - i **ordinal numbers in graphemic form (“erstens” instead of “1.”)**
 - i **For homographs different entries in the dictionary (“modern_adj” for modern and “modern_verb” for molder)**
- i **Consistent file names and extensions**

