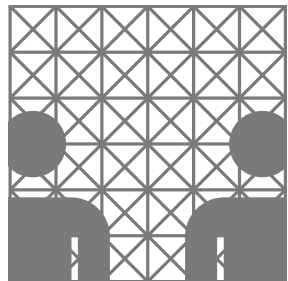


Decision Tree Usage for Incremental Parametric Speech Synthesis

Timo Baumann

baumann@informatik.uni-hamburg.de



Decision Tree Usage for
**Incremental Parametric
Speech Synthesis**

Decision Tree Usage for
**Incremental Parametric
Speech Synthesis**

Decision Tree Usage for **Incremental Parametric Speech Synthesis**

- ability to change ongoing speech output
- needed to cope with unexpected events in interactive use-cases

Speech Output in Typical Systems

current point in time

There's an appointment today at 8:30 titled: 'presentation' with the note: 'do not miss'.

- conventionally:
generate, synthesize and deliver utterance as a whole

Speech Output in Typical Systems

current point in time



There's an appointment today at 8:30 titled: 'presentation' with the note: 'do not miss'.

- potentially slow, as all processing is utterance-initial
→ reason for canned speech in deployed dialogue systems

Speech Output in Typical Systems

current point in time

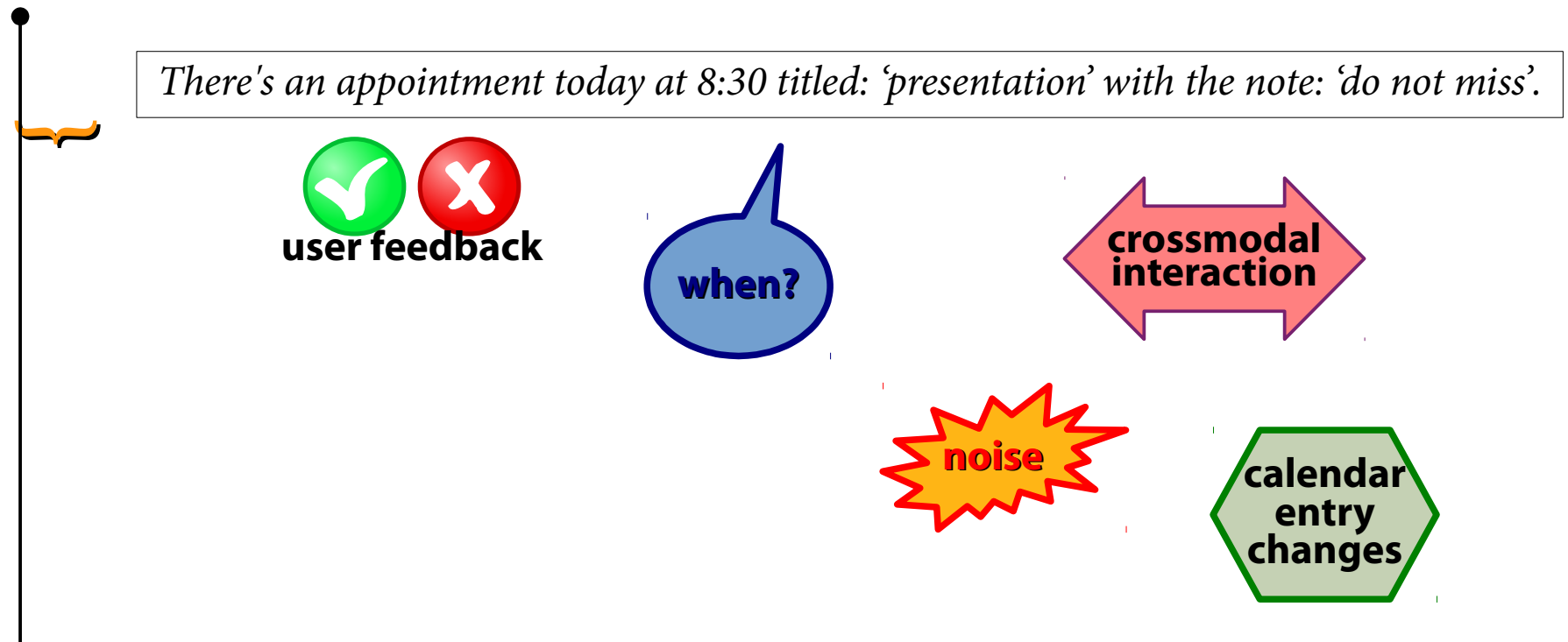


There's an appointment today at 8:30 titled: 'presentation' with the note: 'do not miss'.

- inflexible: unable to change the ongoing utterance (neither the content nor the delivery parameters)
 - no way to react to the listener or the environment

Speech Output in Typical Systems

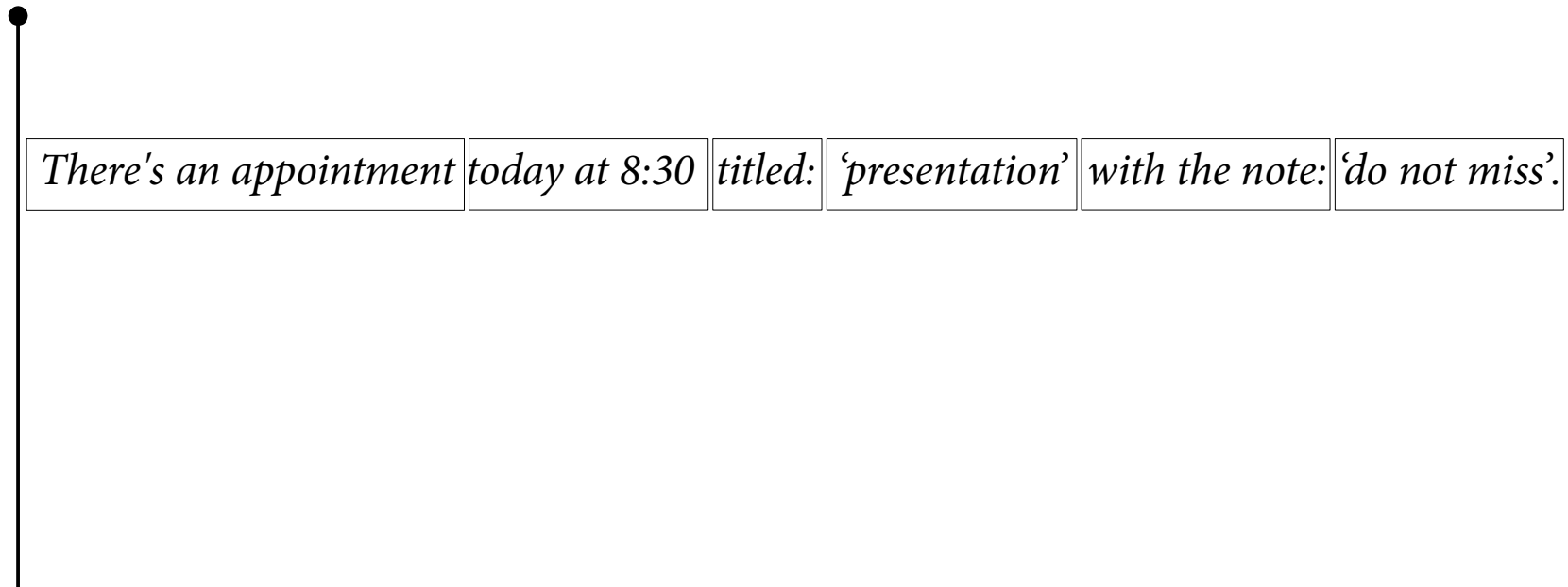
current point in time



- inflexible: unable to change the ongoing utterance (neither the content nor the delivery parameters)
 - no way to react to the listener or the environment

Potentially Better: Incremental Speech Output

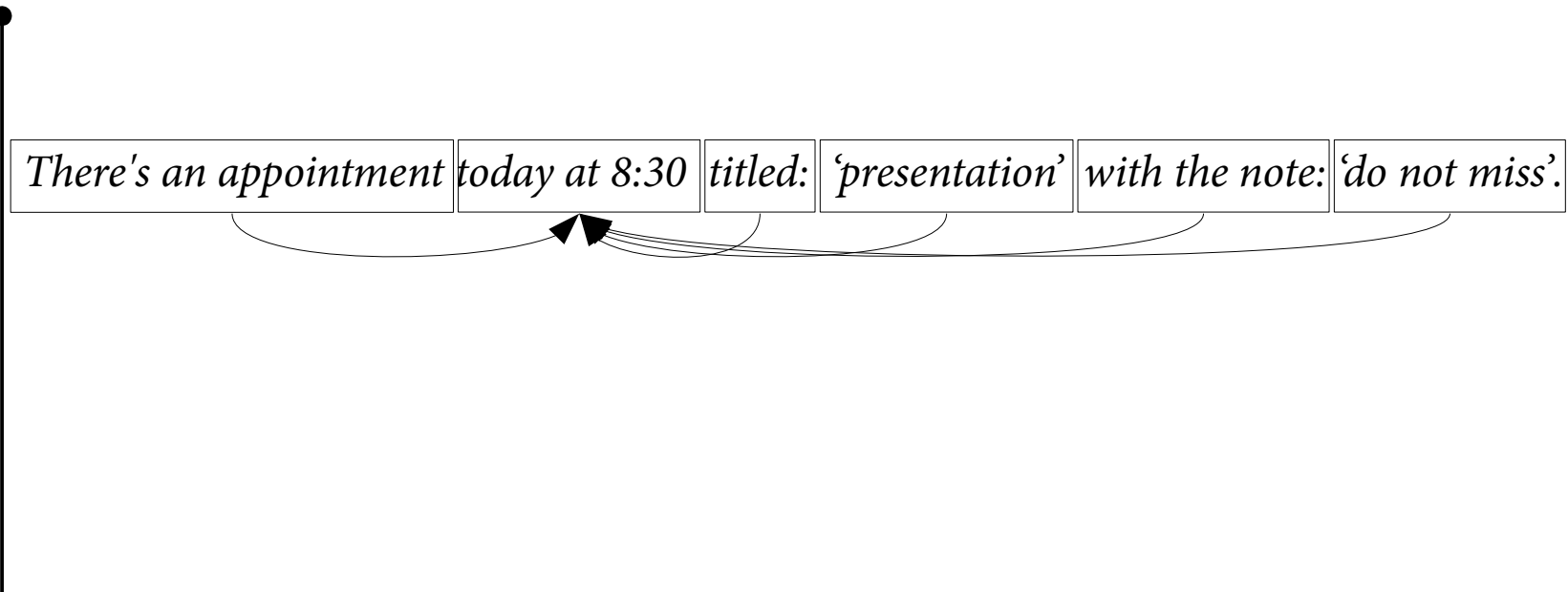
current point in time



- generate, synthesize and deliver the utterance in small *chunks*
 - smaller chunks, higher flexibility
 - but (re)compute with as much context as is available or needed

Potentially Better: Incremental Speech Output

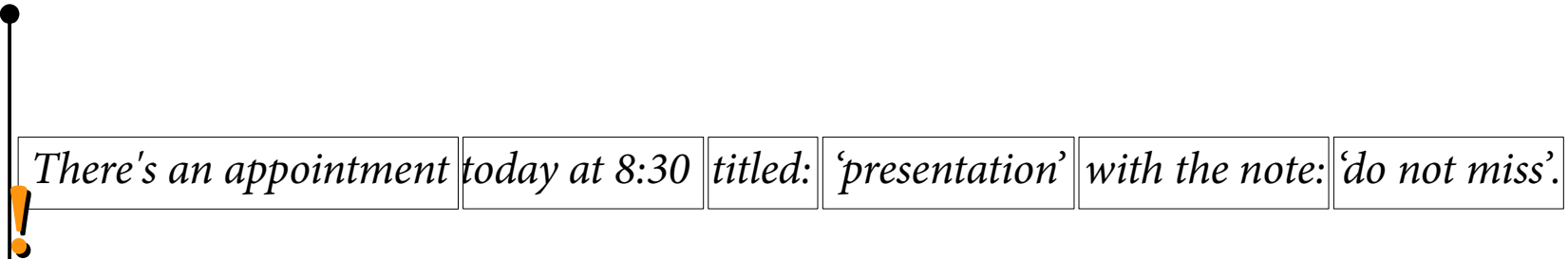
current point in time



- generate, synthesize and deliver the utterance in small *chunks*
 - smaller chunks, higher flexibility
 - but (re)compute with as much context as is available or needed

Potentially Better: Incremental Speech Output

current point in time



- less utterance-initial processing → faster onset

Potentially Better: Incremental Speech Output

current point in time

There's an appointment today at 8:30 titled: 'presentation' with the note: 'do not miss'.



at 8:30, titled: 'presentation' ...

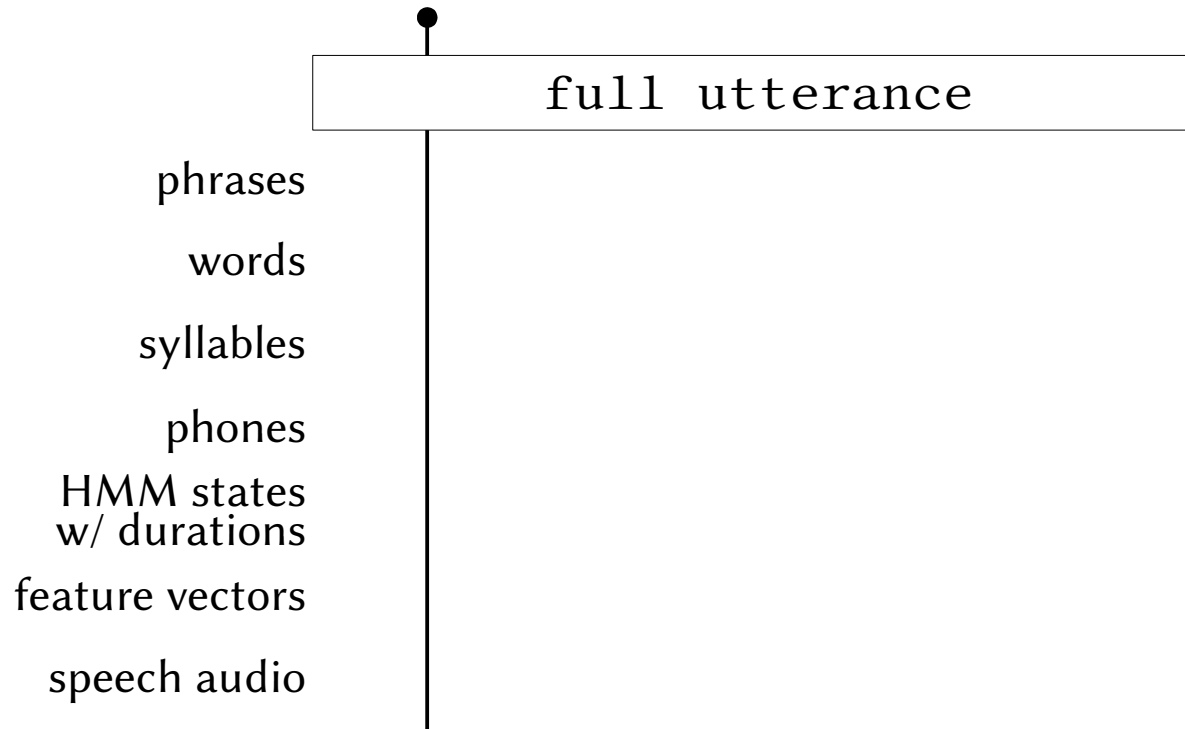
- incremental output may take *changes* into account
- react and adapt to user feedback / requests / noise

what I'm trying to say is:

incremental speech synthesis
is a **requirement** for
highly responsive behaviour
in interactive systems

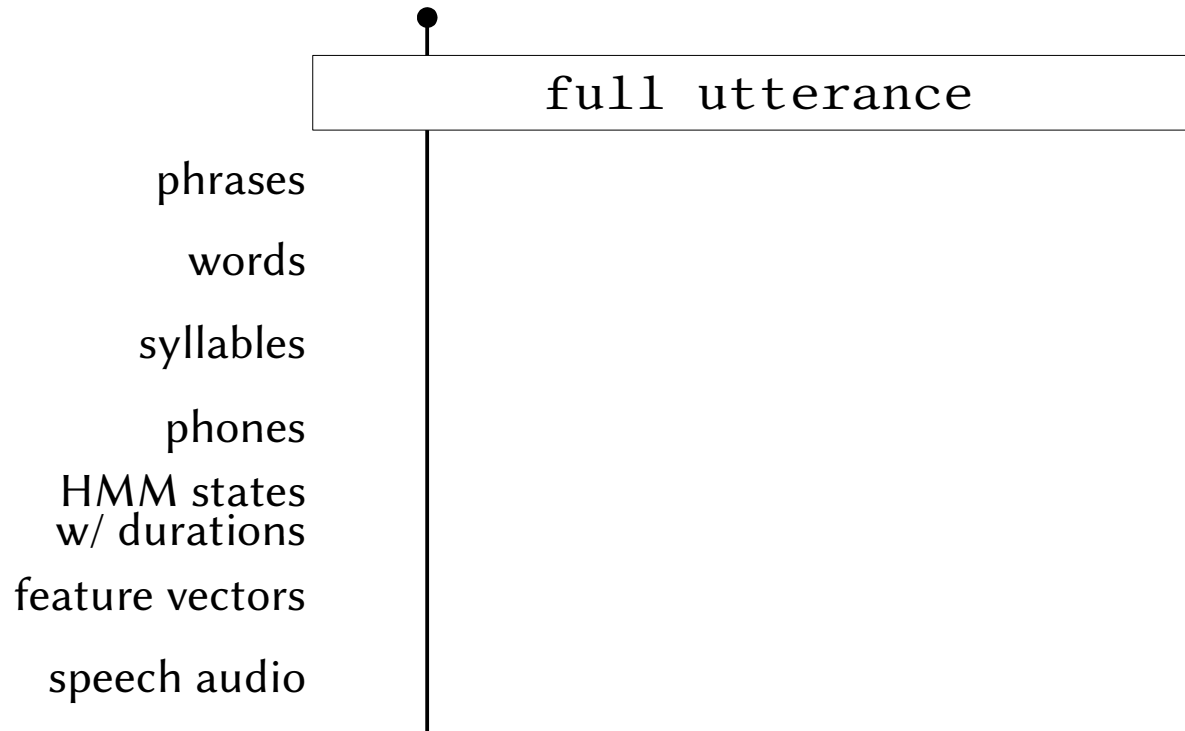
„Just-in-Time“ Incremental Speech Synthesis

„Your flight in May, to Florence, has been confirmed by the airline.“
not?



„Just-in-Time“ Incremental Speech Synthesis

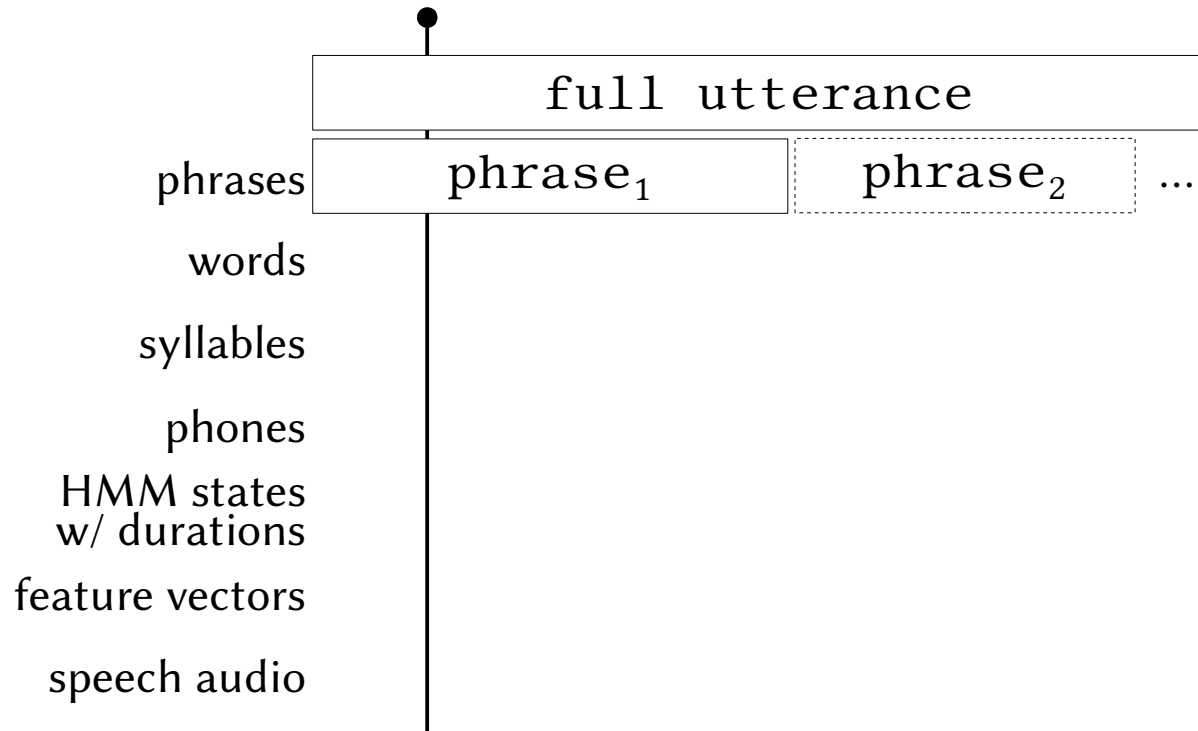
„Your flight in May, to Florence, has been confirmed by the airline.“



as implemented in InproTK (Baumann&Schlangen SDCTD 2012)

„Just-in-Time“ Incremental Speech Synthesis

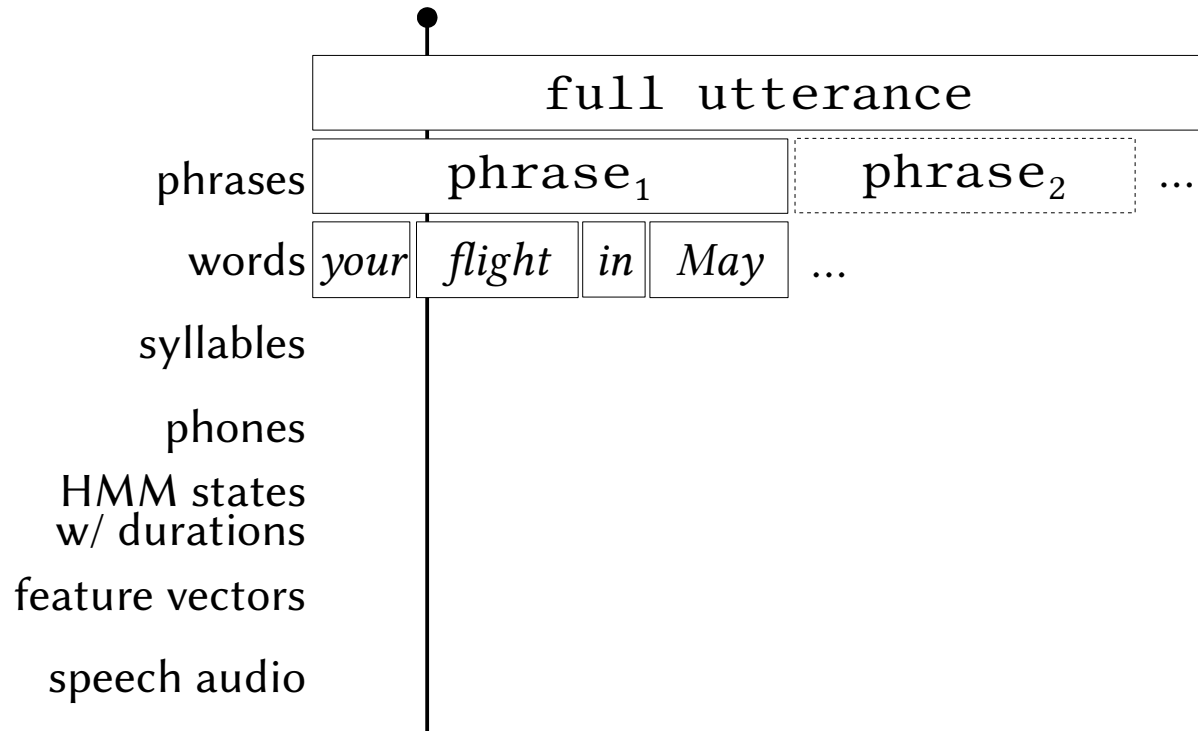
„Your flight in May, to Florence, has been confirmed by the airline.“



as implemented in InproTK (Baumann&Schlangen SDCTD 2012)

„Just-in-Time“ Incremental Speech Synthesis

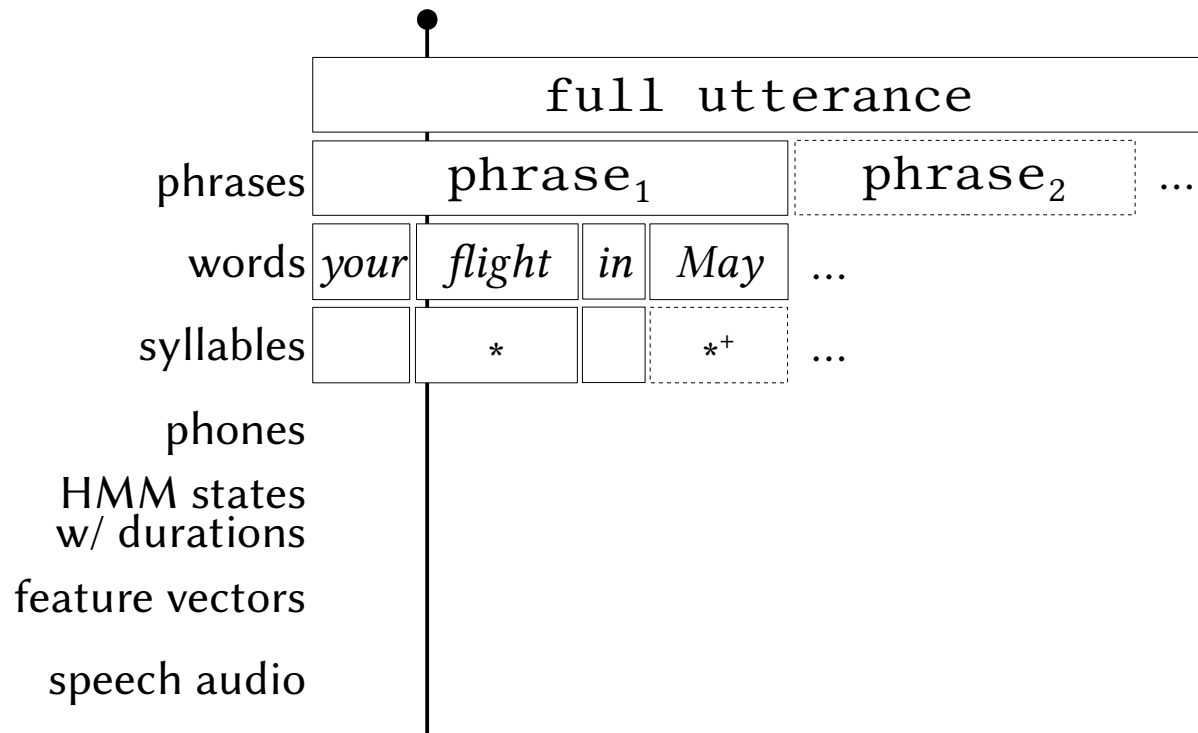
„Your flight in May, to Florence, has been confirmed by the airline.“



as implemented in InproTK (Baumann&Schlangen SDCTD 2012)

„Just-in-Time“ Incremental Speech Synthesis

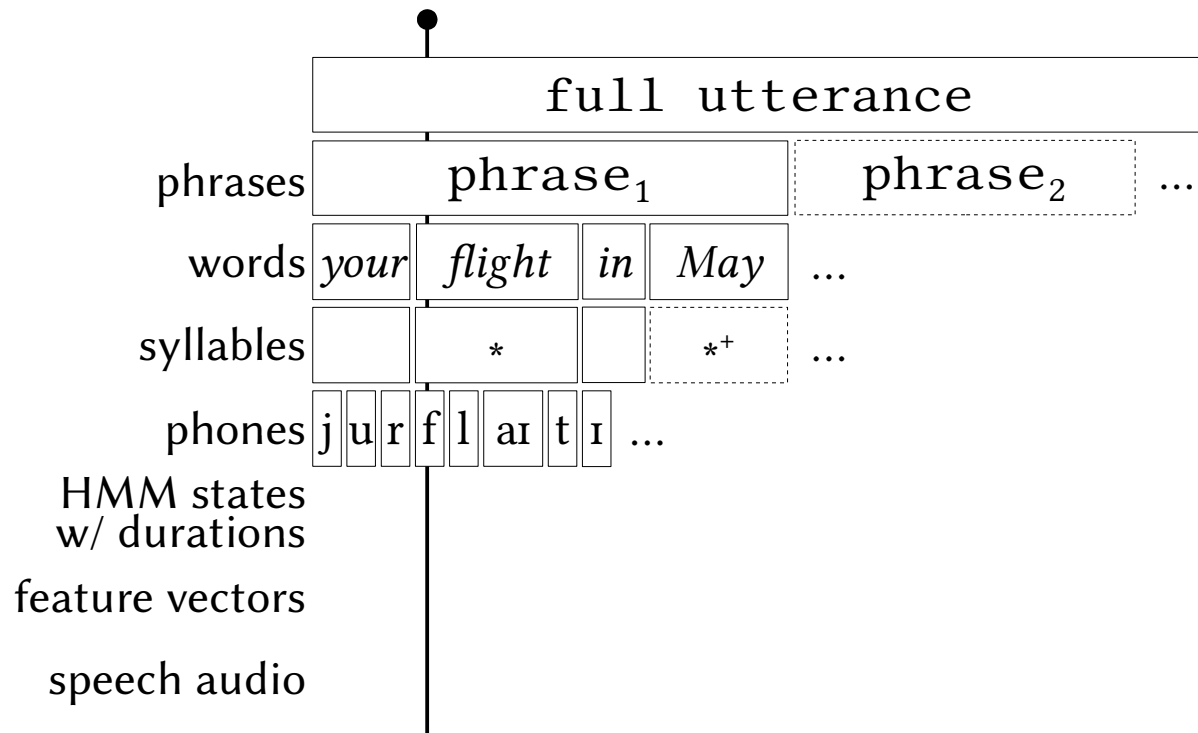
„Your flight in May, to Florence, has been confirmed by the airline.“



as implemented in InproTK (Baumann&Schlangen SDCTD 2012)

„Just-in-Time“ Incremental Speech Synthesis

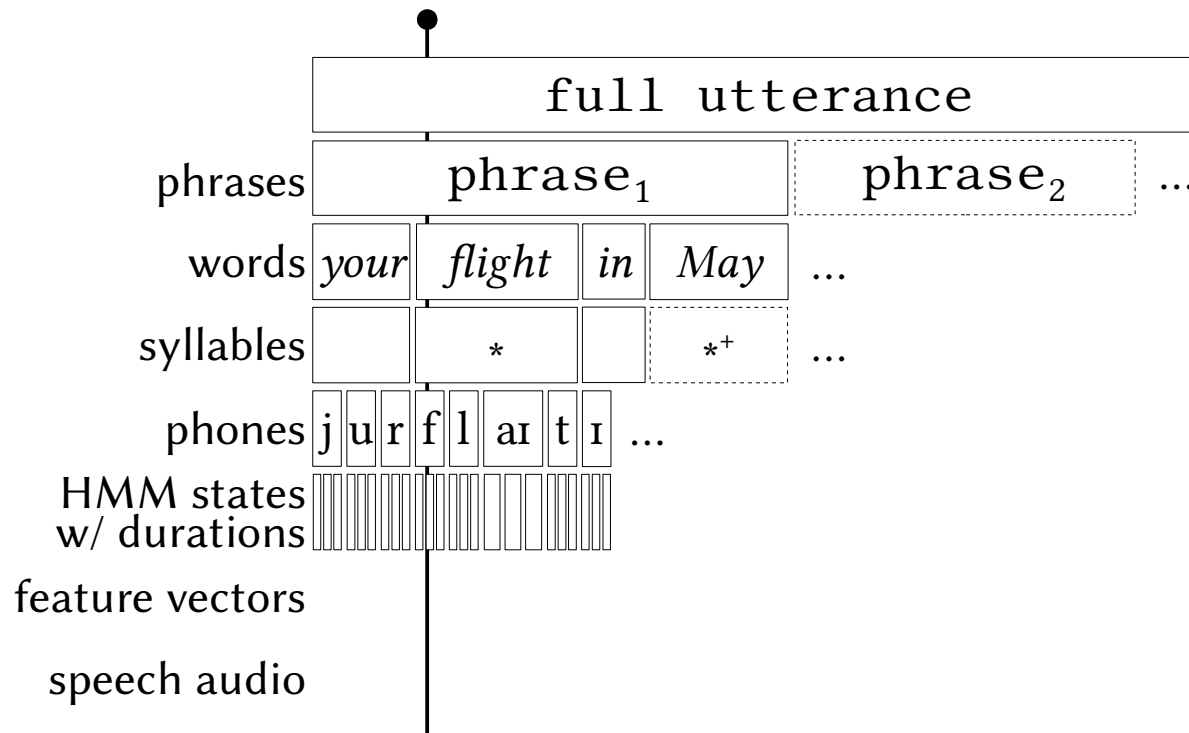
„Your flight in May, to Florence, has been confirmed by the airline.“



as implemented in InproTK (Baumann&Schlangen SDCTD 2012)

„Just-in-Time“ Incremental Speech Synthesis

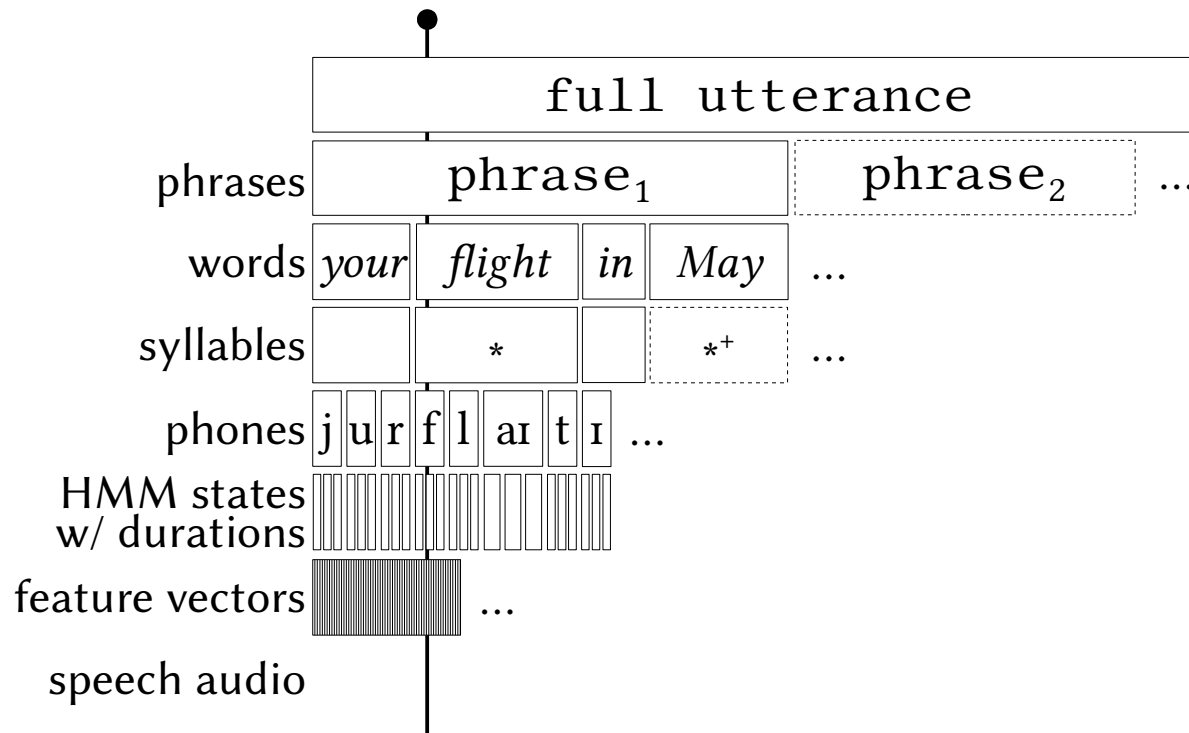
„Your flight in May, to Florence, has been confirmed by the airline.“



as implemented in InproTK (Baumann&Schlangen SDCTD 2012)

„Just-in-Time“ Incremental Speech Synthesis

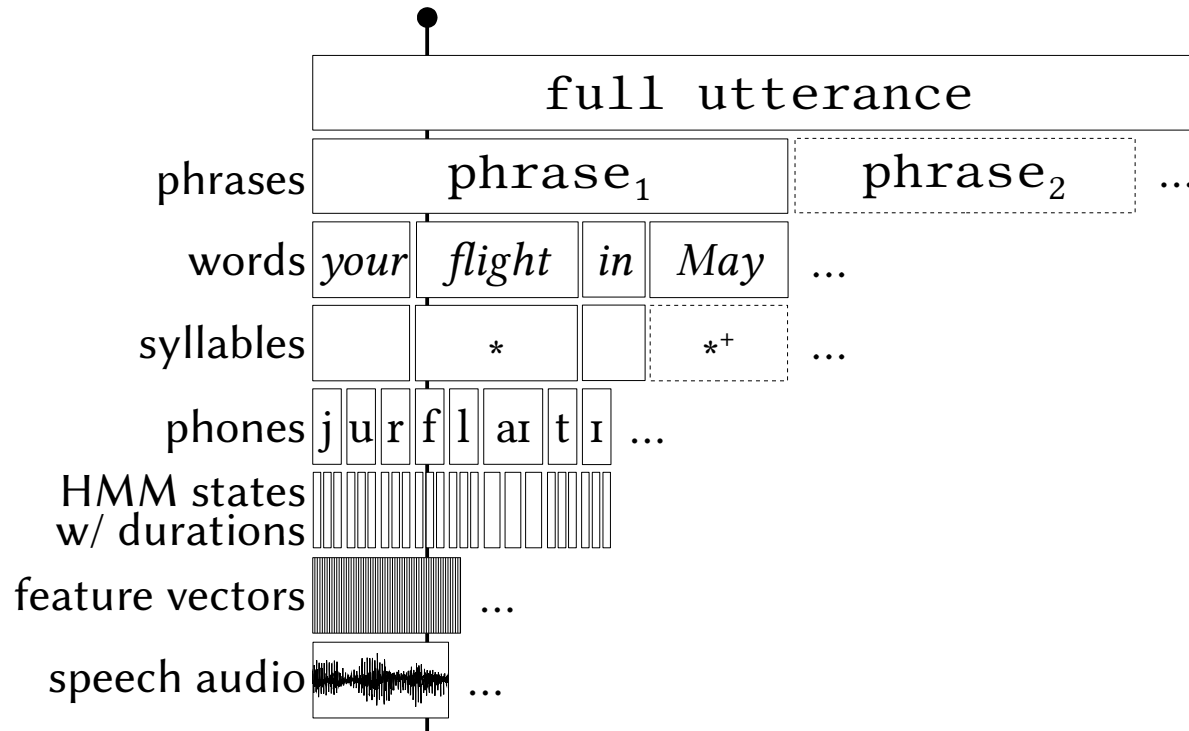
„Your flight in May, to Florence, has been confirmed by the airline.“



as implemented in InproTK (Baumann&Schlangen SDCTD 2012)

„Just-in-Time“ Incremental Speech Synthesis

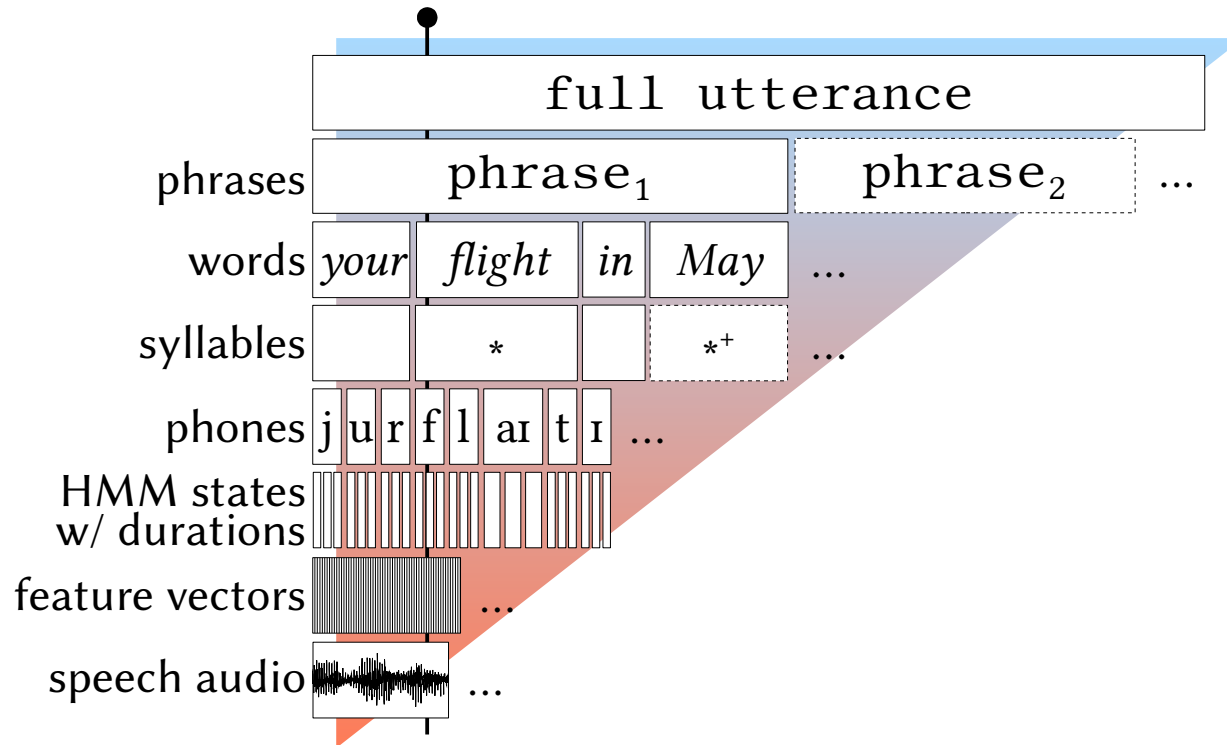
„Your flight in May, to Florence, has been confirmed by the airline.“



as implemented in InproTK (Baumann&Schlangen SDCTD 2012)

„Just-in-Time“ Incremental Speech Synthesis

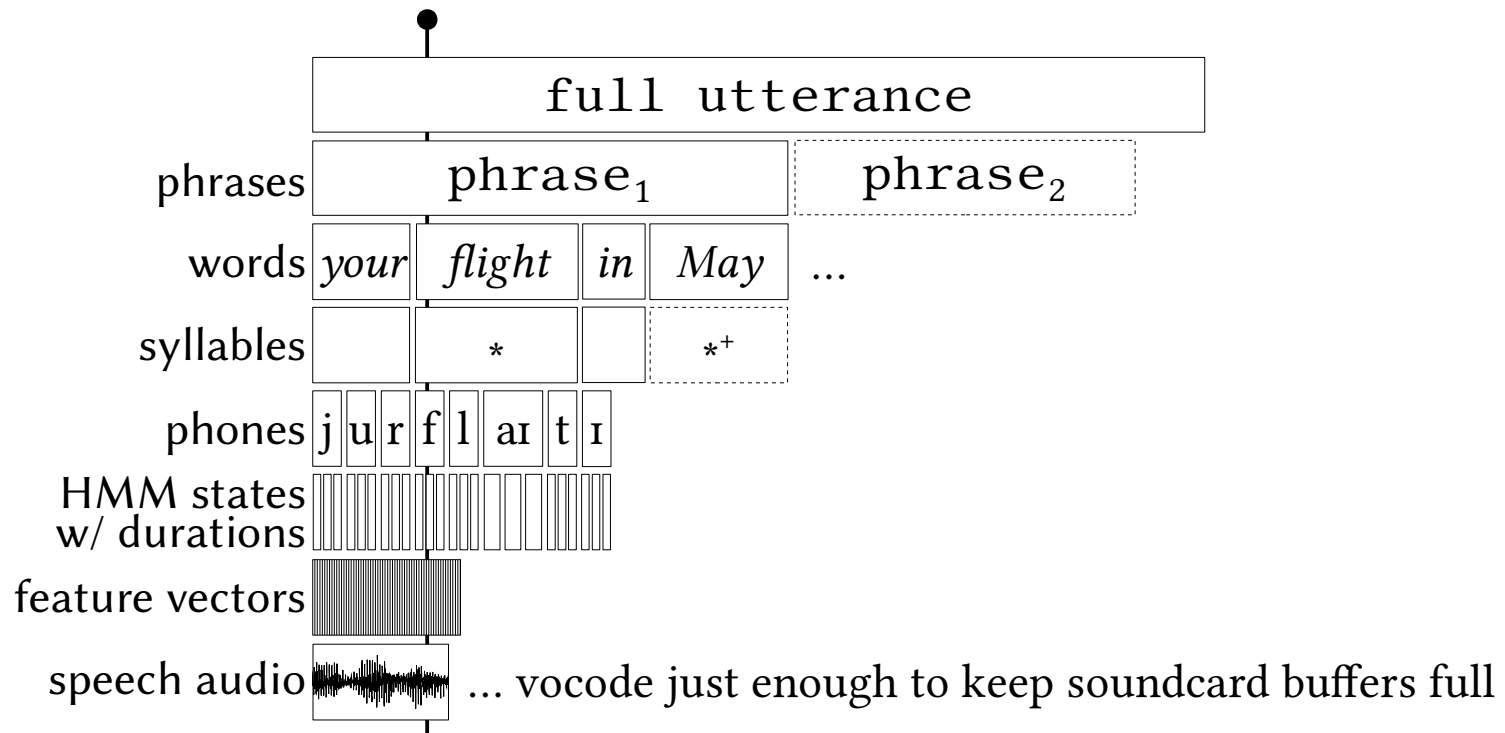
„Your flight in May, to Florence, has been confirmed by the airline.“



as implemented in InproTK (Baumann&Schlangen SDCTD 2012)

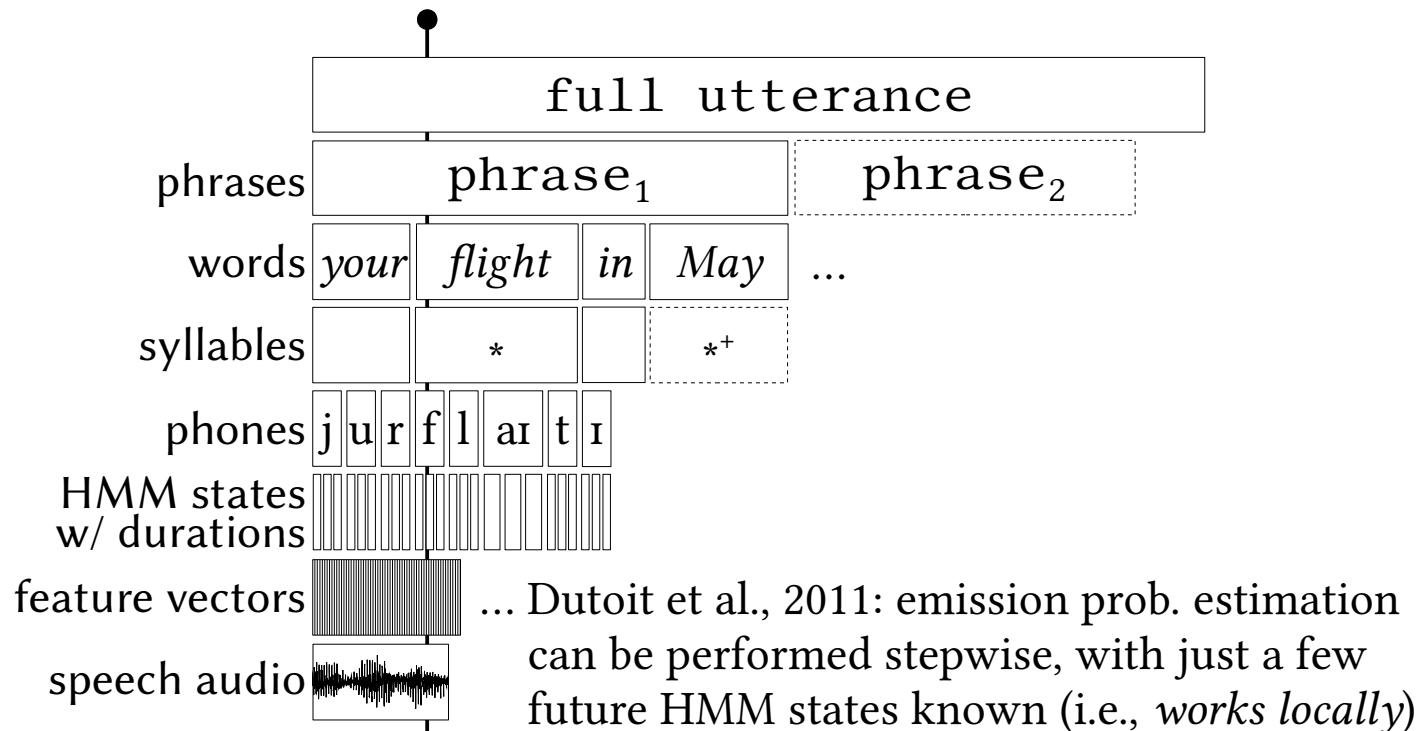
„Just-in-Time“ Incremental Speech Synthesis

„Your flight in May, to Florence, has been confirmed by the airline.“



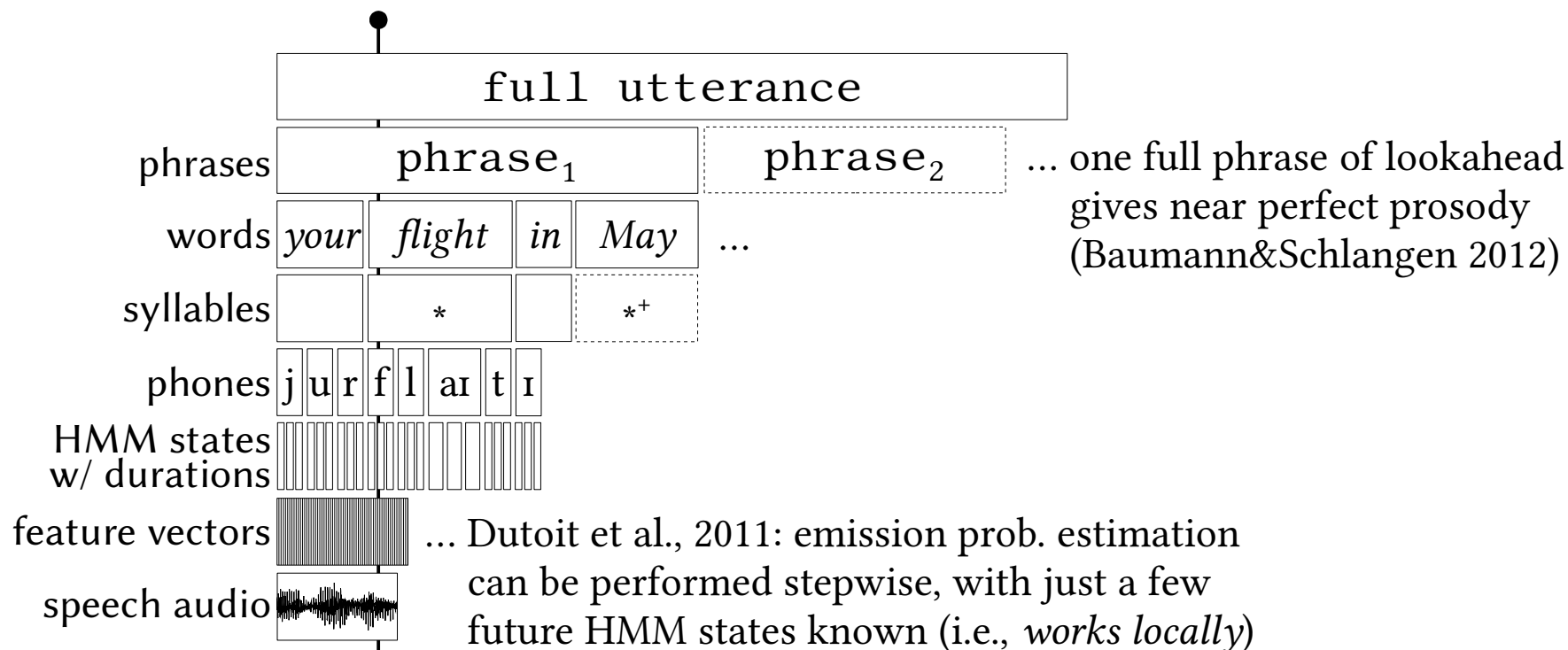
„Just-in-Time“ Incremental Speech Synthesis

„Your flight in May, to Florence, has been confirmed by the airline.“



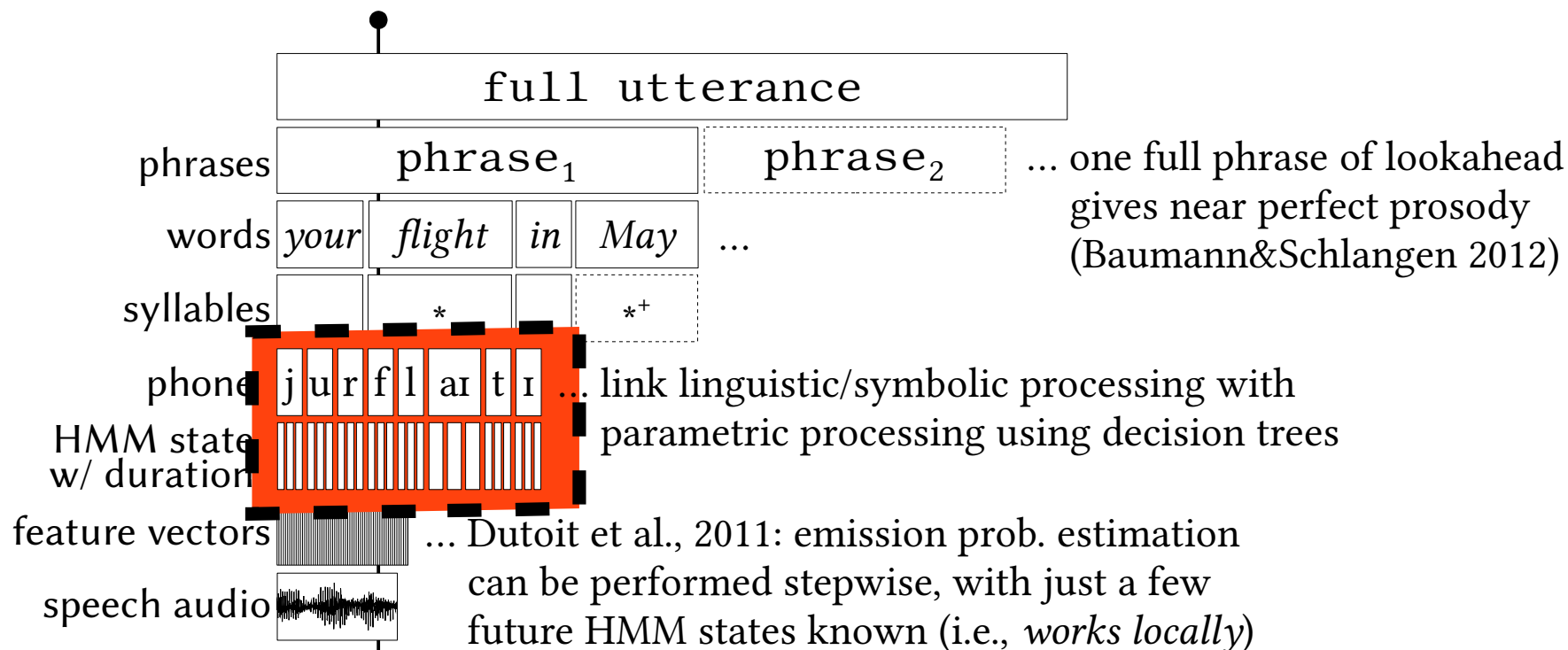
„Just-in-Time“ Incremental Speech Synthesis

„Your flight in May, to Florence, has been confirmed by the airline.“



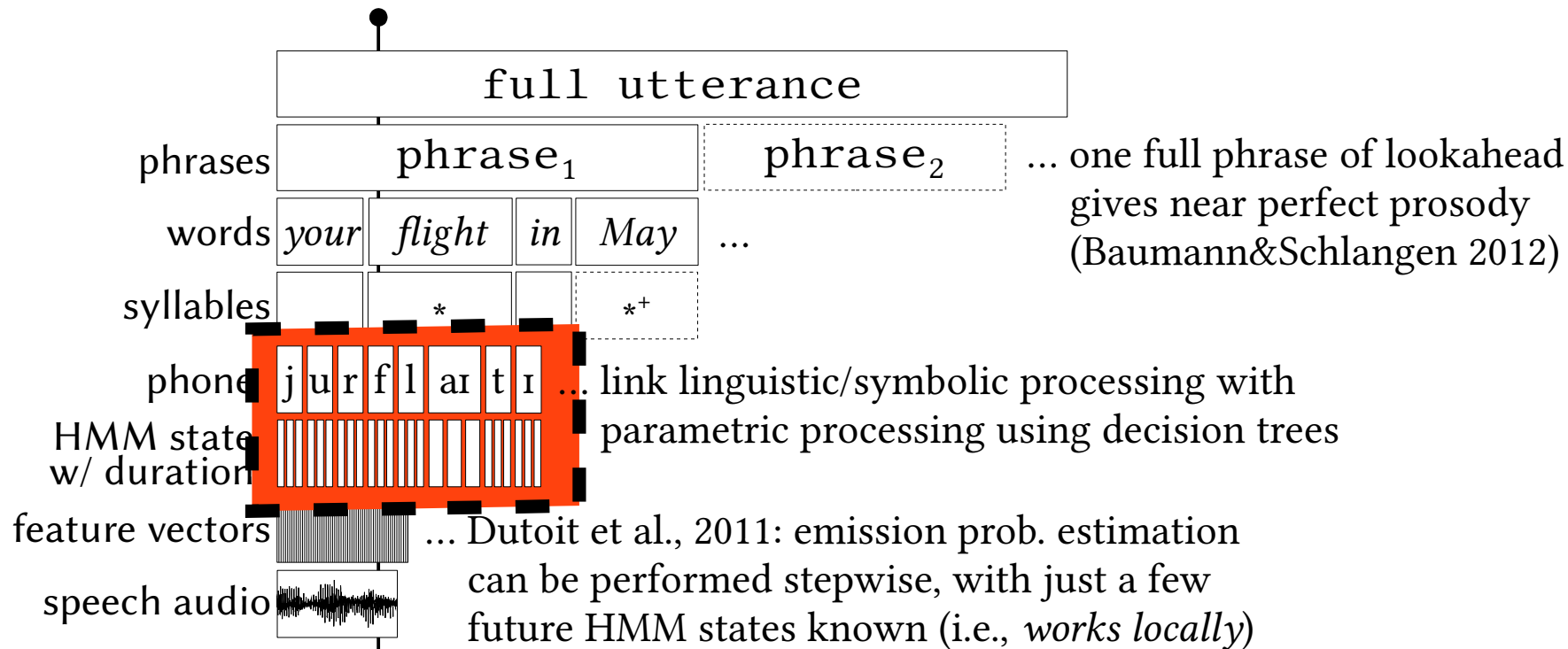
„Just-in-Time“ Incremental Speech Synthesis

„Your flight in May, to Florence, has been confirmed by the airline.“



„Just-in-Time“ Incremental Speech Synthesis

„Your flight in May, to Florence, has been confirmed by the airline.“



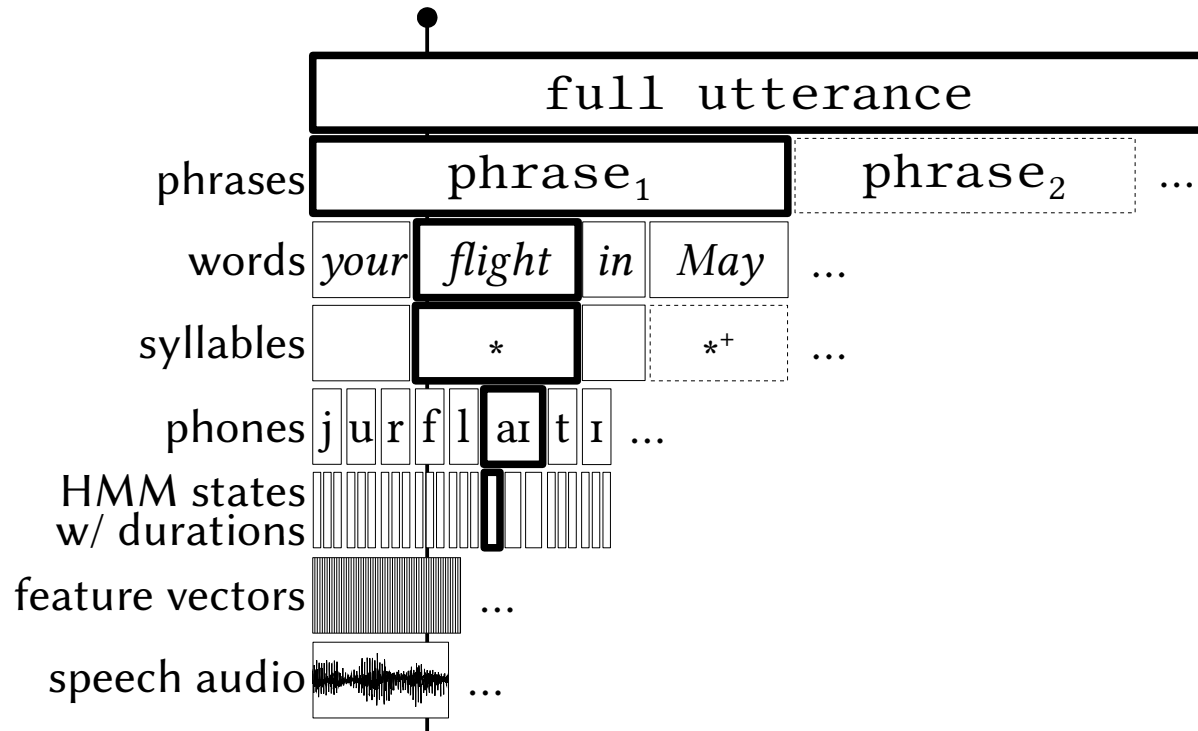
- goal: better incremental HMM state selection
 - without re-training synthesis voices
 - without sacrificing non-incremental performance

HMM State Selection

- most TTS systems (such as MaryTTS) use decision trees
 - separate trees for MCEP, STR, f0 streams, and state duration
- feature sets using various types of information
 - MaryTTS: roughly 100 features
- many features are *non-local*
 - such as „how many phonemes until end of utterance?“
- non-local features
are not available in incremental processing

Classification of features among two dimensions

„Your flight in May, to Florence, has been confirmed by the airline.“

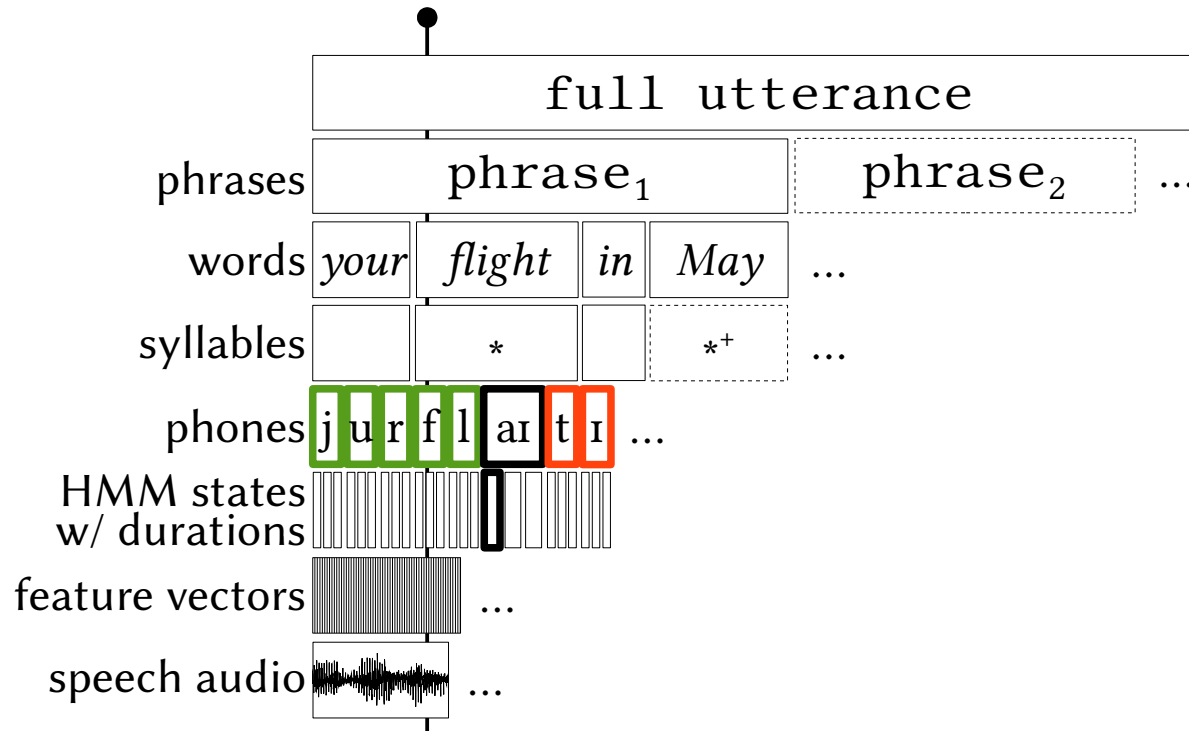


- level of linguistic abstraction
 - higher-level information spans longer time frames

as implemented in InproTK (Baumann&Schlangen SDCTD 2012)

Classification of features among two dimensions

„Your flight in May, to Florence, has been confirmed by the airline.“



- temporal direction: **past**, present, **future**
 - past is generally available, future requires lookahead

as implemented in InproTK (Baumann&Schlangen SDCTD 2012)

A Classification of Features: MaryTTS feature counts for German

	past	current	future
full sentence			
phrase/accenuation			
word			
syllable			
phone			

A Classification of Features: MaryTTS feature counts for German

	past	current	future
full sentence	—	5	—
phrase/accenuation	11	10	10
word	2	7	3
syllable	3	8	2
phone	20	10	19

A Classification of Features: MaryTTS feature counts for German

	past	current	future
full sentence	—	5	—
phrase/accenuation	11	10	10
word	2	7	3
syllable	3	8	2
phone	20	10	19

- generalize features into classes that represent lookahead requirements in an incremental system

A Classification of Features: MaryTTS feature counts for German

	past	current	future
full sentence	—	5	—
phrase/accenuation	11	10	10
word	2	7	3
syllable	3	8	2
phone	20	10	19

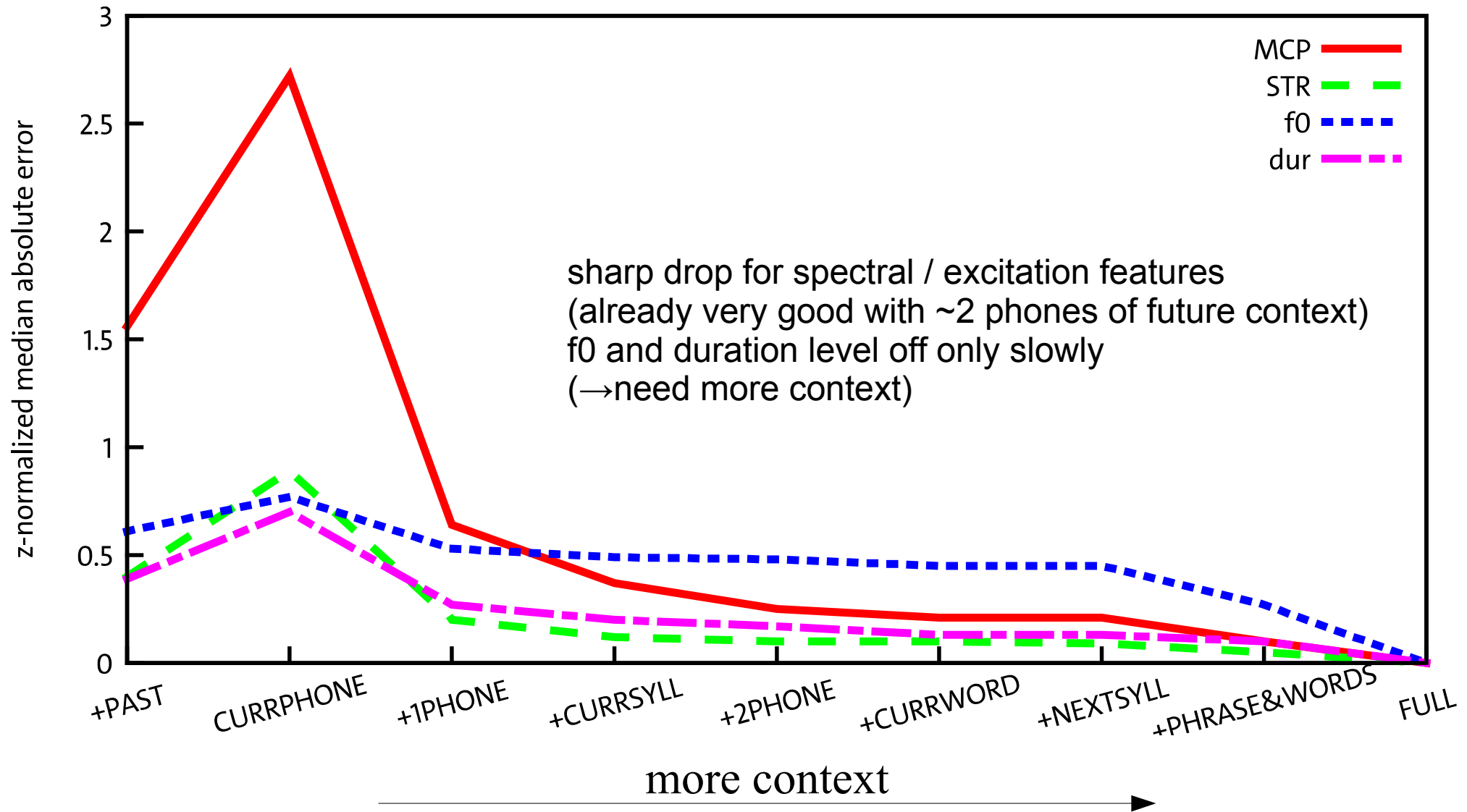
- generalize features into classes that represent lookahead requirements in an incremental system

Experiment:

What if a feature class is not available?

- substitute a default feature value
- training:
 - synthesized 600 utterances and recorded all feature usages (and their values) in the decision trees
 - determine default feature values
 - most common value for categorical features
 - mean value for numeric features
- test:
 - re-process, substituting features of a class by their defaults
- measure (numeric) deterioration of resulting HMM states
(z-normalized mean absolute error)

Results per Feature Class



Conclusion

- decision tree features can be missing during incremental processing
 - substitute with default values
- classified features into classes which are meant to correspond to context/lookahead requirements
- the more context, the better the results
 - relatively small lookahead (1 syllable/2 phones) enough for voice quality (MCEP and STR)
 - prosody (duration and f_0) in contrast, requires a large lookahead, or more advanced methods
- there's no simple good or bad, but a continuous improvement the more context is available

Thank you.

baumann@informatik.uni-hamburg.de,
get the code at inprotk.sf.net.

Funded by a Daimler and Benz Foundation PostDoc grant.

Thanks to Sven Mutzl, who performed some initial analyses while on an internship with Prof. Kai Yu at Yiao Tong University, Shanghai, China.

Raum für Notizen

Default feature values vs.
„properly“ dealing with missing features