

Character Stream Parsing of Mixed-lingual Text

Harald Romsdorfer and Beat Pfister

Text Analysis for TTS: Traditional

- "... Studio. They record the next record. Then ..."
⇒ [ðeɪ ? ðə nekst ?]

they		ðeɪ
record		?
the		ðə
next		nekst
record		?

Text Analysis for TTS: SVOX, ...

- "... Studio. They record the next record. Then ..."
⇒ [ðeɪ rɪ'kɔ:d ðə nekst 'rekɔ:d]

they		ðeɪ	personal pronoun
record		rɪ'kɔ:d	verb
the		ðə	determiner
next		nekst	adjective
record		'rekɔ:d	noun

⇒ Syntactic analysis is a must for TTS synthesis!

Text Analysis for TTS: Problems I & II

Contracted Word Forms & Ambiguous Punctuation Symbols

- "It's in St. Mary's St. Mary's at home."
⇒ [its in ? meəriz ? meəriz ət həʊm]

it's	its	?
in	in	preposition
st.	?	?
mary's	m'eəriz	?
st.	?	?
mary's	m'eəriz	?
at	ət	preposition
home	həʊm	?

Text Analysis for TTS: Problems I & II

Contracted Word Forms & Ambiguous Punctuation Symbols

- "It's in St. Mary's St. Mary's at home."
⇒ [ɪts ɪn sənt meəriz stri:t] [meəriz ət həʊm]

it	it	personal pronoun	mary	m'eəri	proper noun
's	s	auxiliary "be"	's	z	auxiliary "be"
in	ɪn	preposition	at	ət	preposition
st.	sənt	noun title "Saint"	home	həʊm	noun
mary's	m'eəriz	noun (possessive form)	.		full stop
st	stri:t	noun "street"			
.		full stop			

⇒ Use **syntactic words** instead of graphemic words as tokens!

Text Analysis for TTS: Problems III

Multi-word lexemes

- "He's in fine conditions in fine."
⇒ [hiz in **faj̩n** kəndɪʃ(ə)nz in **faj̩ni**]

he	hi	personal pronoun
's	z	auxiliary "be"
in	ɪn	preposition
fine	faj̩n	adjective
conditions	kəndɪʃ(ə)nz	noun
in fine	ɪn faj̩ni	adverb

Text Analysis for TTS: Problems IV

Languages without word delimiter: Chinese, Japanese, ...

研究生 亦般年闹大

他宅研究 生 命起源

Text Analysis for TTS: Problems IV

Languages without word delimiter: Chinese, Japanese, ...

研究生	亦般	年闹	大
yan2-jiu1-sheng1	yi4-ban1	nian2-ling2	da4
'Master student'	'generally'	'age'	'old'

他	宅	研究	生命起源
ta1	zhai	yan2-jiu1	sheng1-ming4-qi3-yuan2
'He'	'doing'	'research'	'the origin of life'

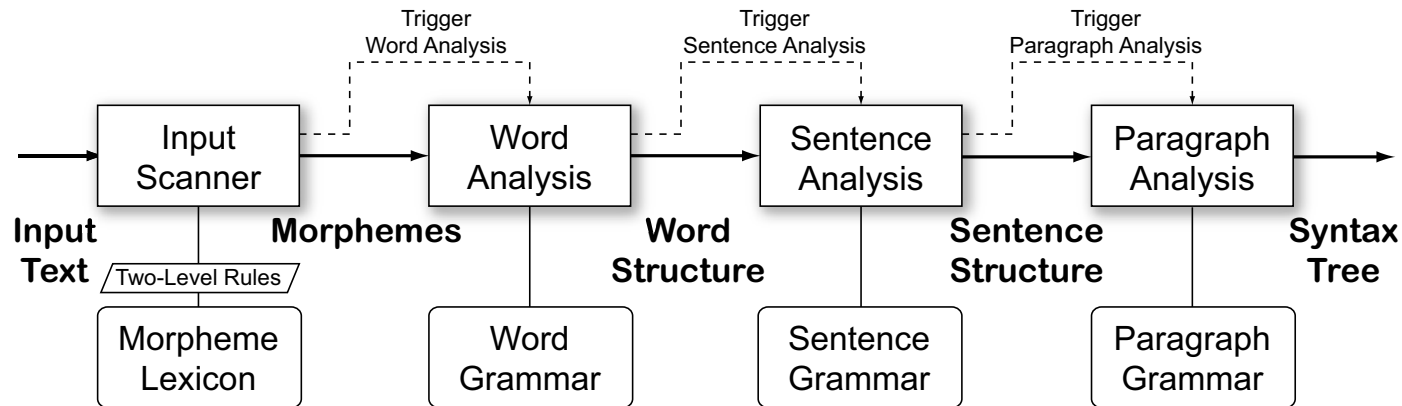
⇒ Boundary of syntactic words depends on **morphological and syntactic context!**

Basic Problem

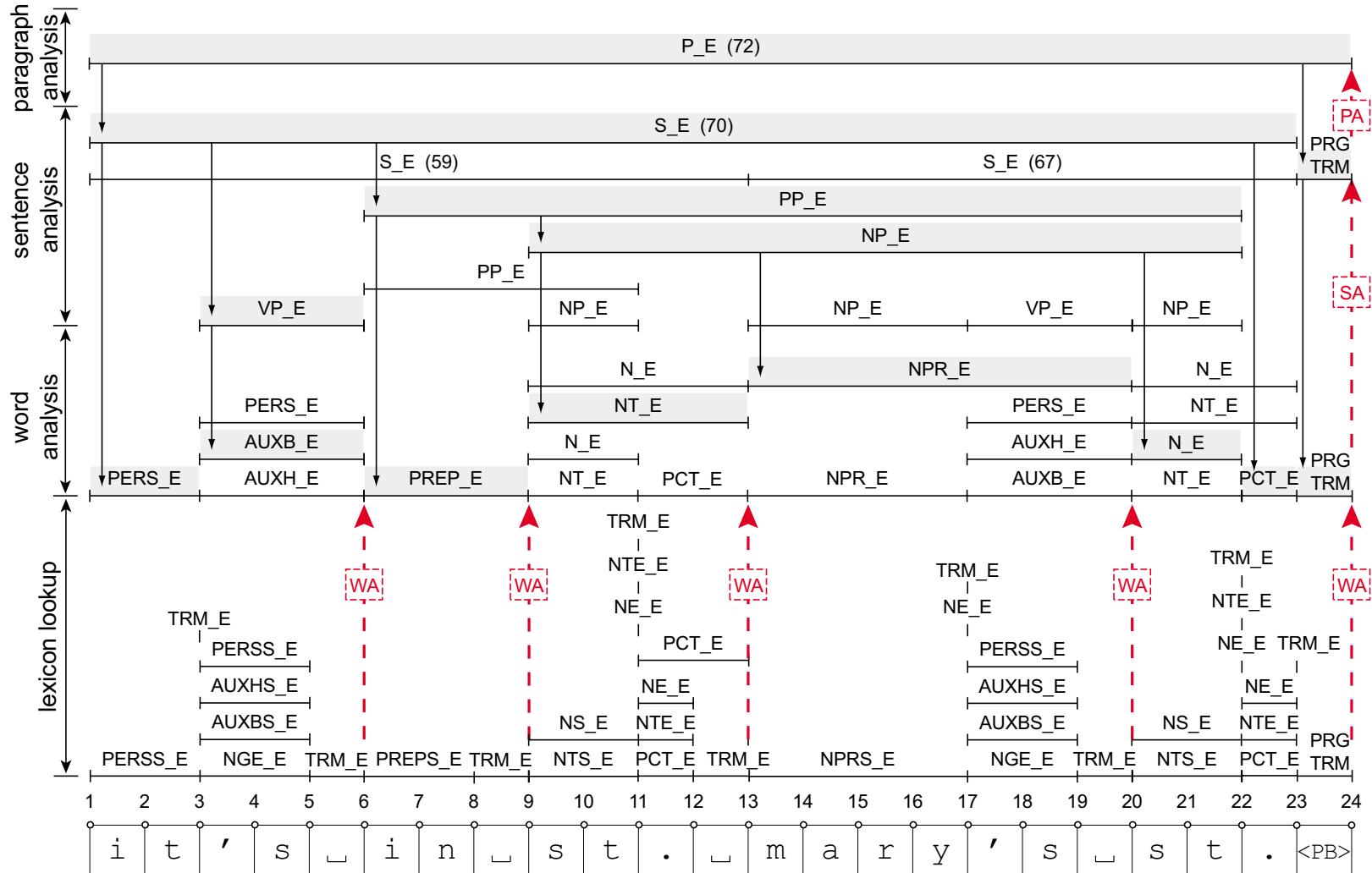
⇒ Syntactic words are the input to syntactic analysis.

⇒ But boundaries of syntactic words depend on result
of syntactic analysis!

Text Analysis of polySVOX



- ⇒ Scan input text character by character
- ⇒ Treat blanks, punctuation symbols, etc. like any other character
- ⇒ Define word boundaries to be blanks, “empty” symbols and certain punctuation symbols
- ⇒ Each module provides all alternatives for subsequent module
- ⇒ Each module triggers next module at unambiguous analysis positions



Text Analysis for TTS: Chinese/Japanese texts

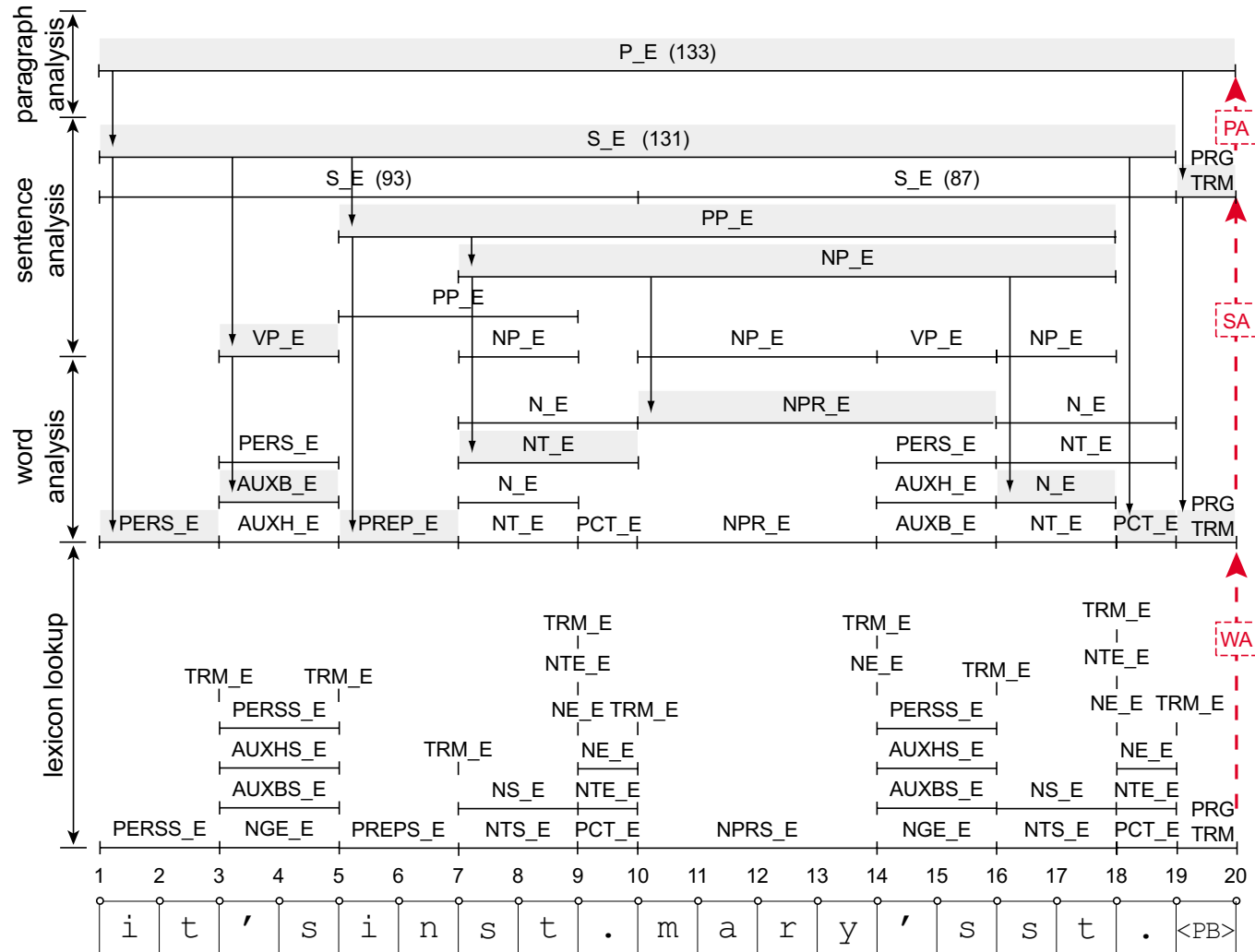
研究生	亦般	年闹	大
yan2-jiu1-sheng1	yi4-ban1	nian2-ling2	da4
'Master student'	'generally'	'age'	'old'

他	宅	研究	生命起源
ta1	zhai	yan2-jiu1	sheng1-ming4-qi3-yuan2
'He'	'doing'	'research'	'the origin of life'

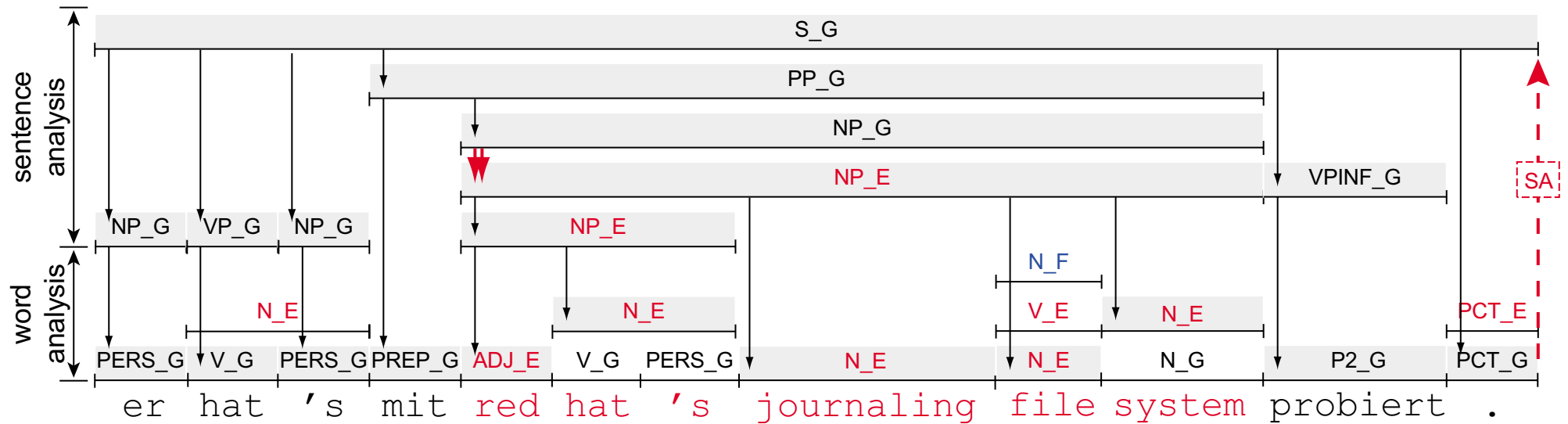
“Equivalent” example in English:

"it'sinst.mary'sst." ⇒ "It's in St. Mary's St."

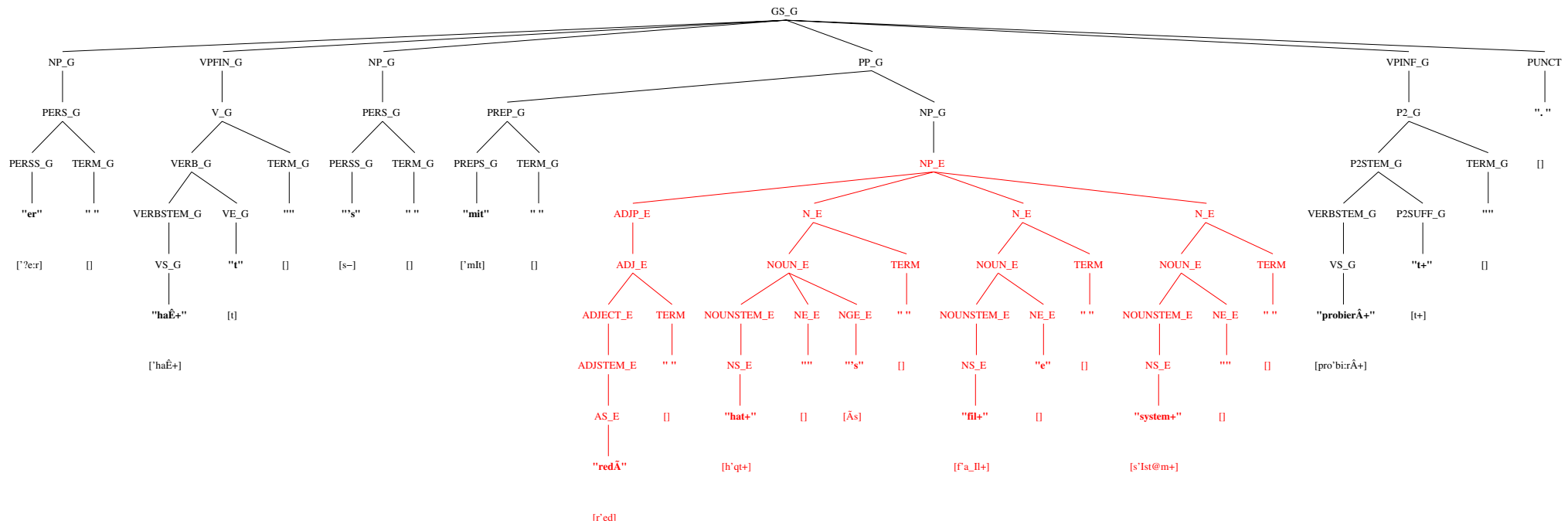
"itsinstwice." ⇒ "It sins twice."



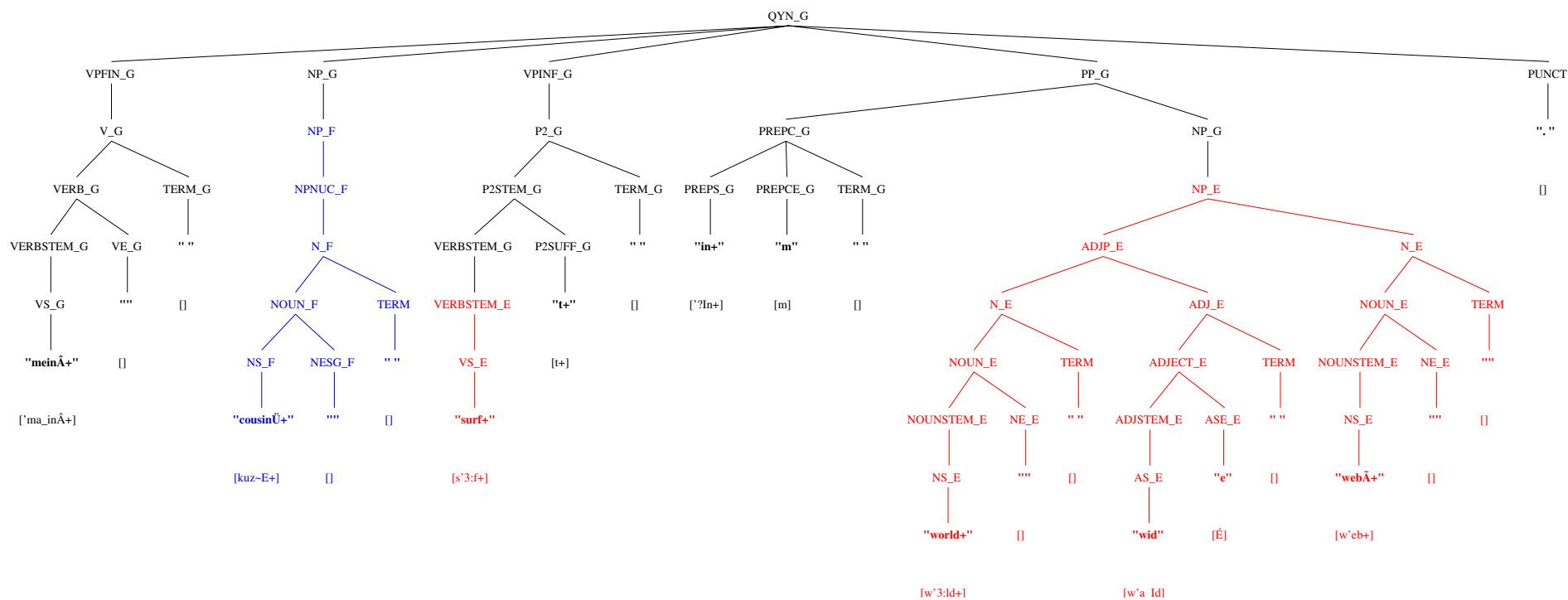
Mixed-lingual Sentences



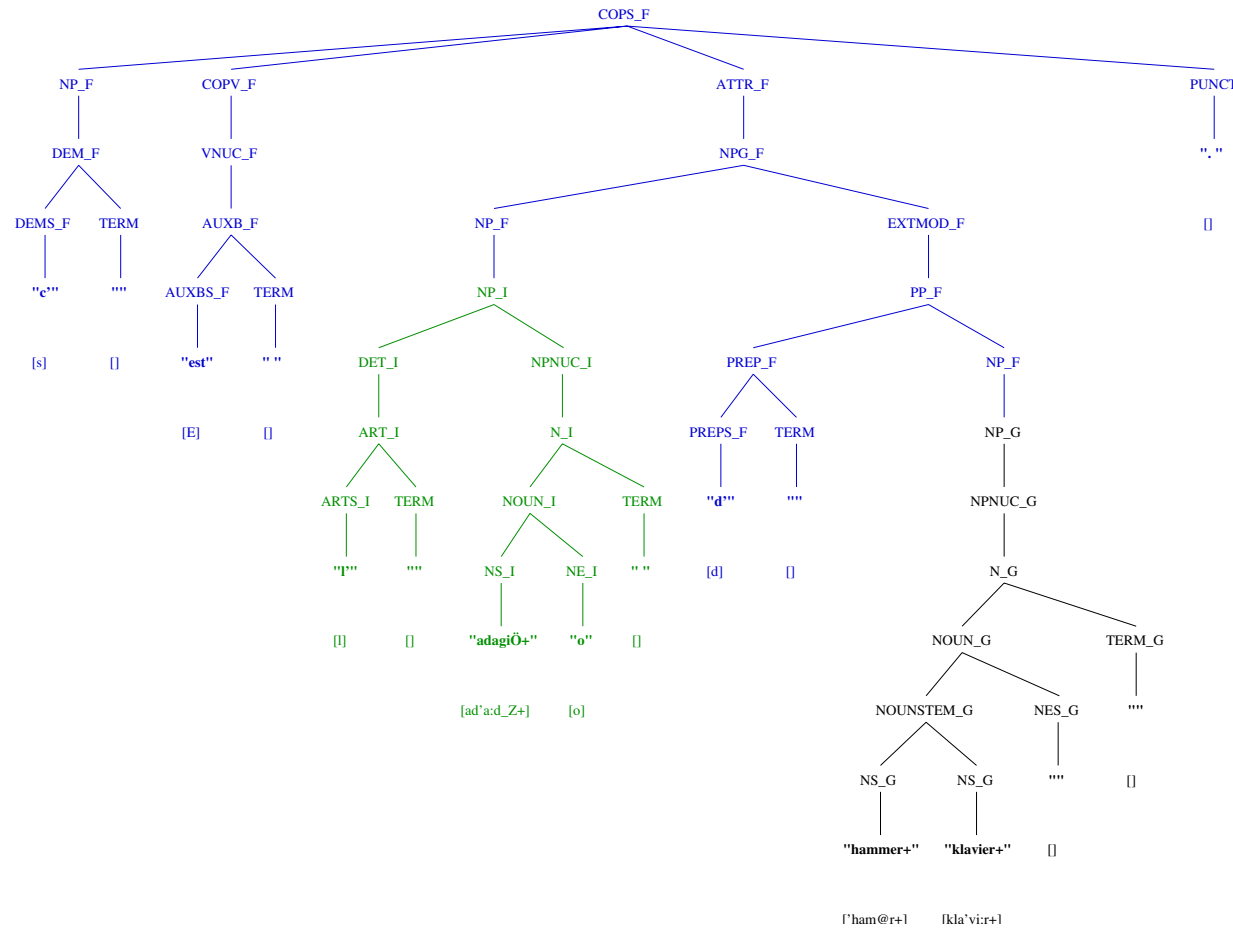
(He tried it with Red Hat's journaling file system.)



“Er hat’s mit **Red Hat’s File System** probiert.”



“Mein Cousin surft im World Wide Web.”



“C’est l’Adagio d’Hammerklavier.”

<<http://www.tik.ee.ethz.ch/~spr/SVOX/polysvoxdemo/>>

Thank you!