

LSA 125 – Psycholinguistics and Syntactic Corpora

**Today: Other ways to do corpus-
based psycholinguistics**

LSA Summer Institute 2009, UC Berkeley

Florian Jaeger

TAs: Judith Degen, Alex Fine, and Peter Graff



Today

- **Corpora vs. Experiments**
- **Other useful tools** and thoughts on ...
 - how to use *non*-syntactic corpora for research on sentence production.
 - how to use corpora for work on *comprehension*.
 - Some really cool resources



Common claims against corpus-based research to language production

- ... by linguists:
 - People make mistakes
 - Not clear that all speakers are native speakers
 - Rareness of relevant events → cannot argue from lack of occurrence
 - Grammaticality != judgments != occurrence in corpus
 - Corpora cannot be used to test usage-based gradience in data since they average over speakers



Common claims against corpus-based research to language production

- ... by psychologists:
 - Corpora provide only heterogeneous data which make reliable hypothesis testing impossible
 - Corpora cannot be used to learn about the underlying mechanisms of language production
 - Corpora only allow post-hoc testing of hypotheses



Trade-offs

- Ecological validity
- Heterogeneity (genre, styles, subject populations)
- Distributions



What happens when we present weird stimuli to our subjects?

- Recently on the streets:

“The plumber that the clerk saw in a pub bought a house at the river.”

“The boy that the girl kicked is tall.”

“Who do you think believes that people get used to unnatural stimuli?”

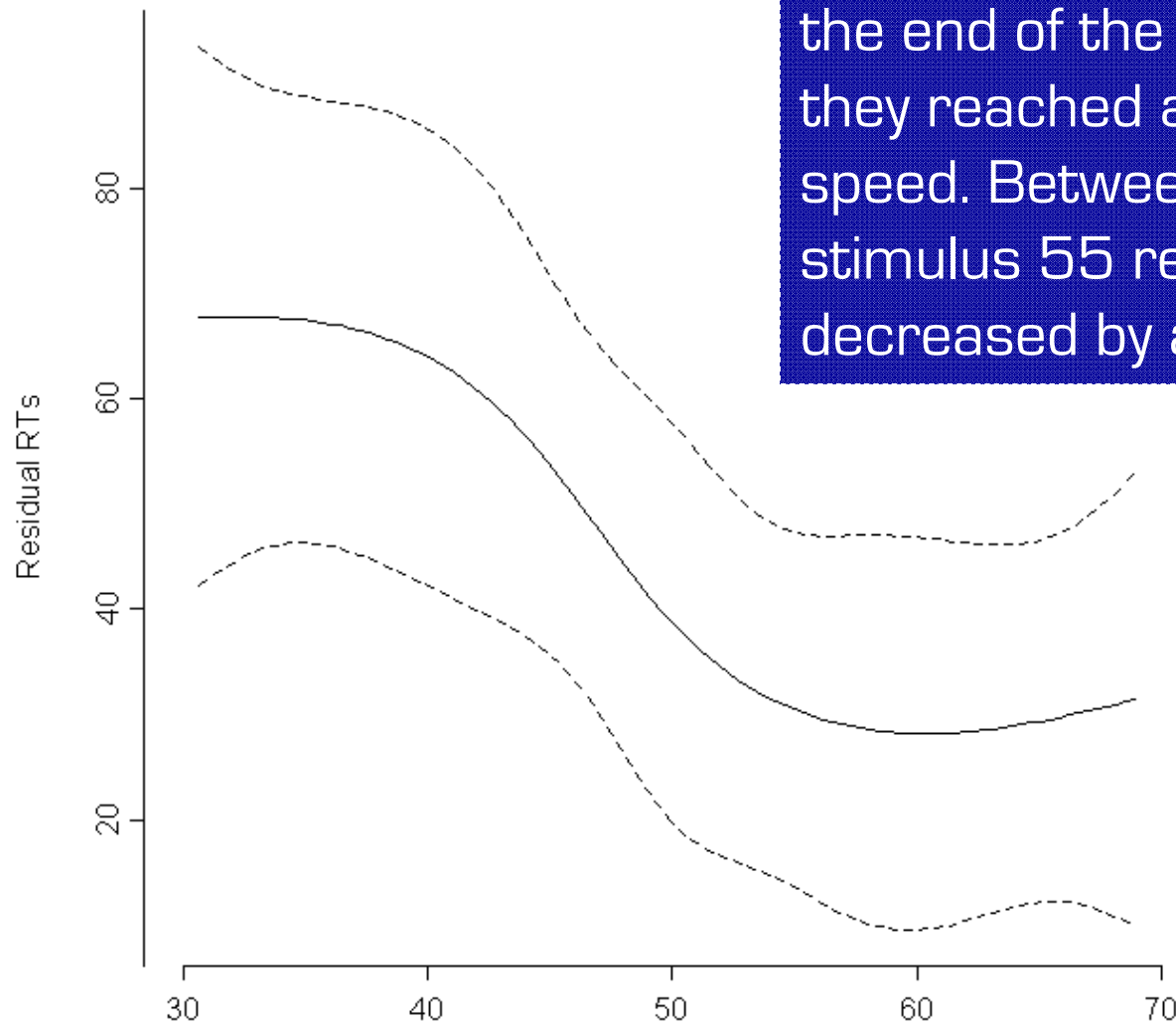


Processing Effects in Superiority-Violating *wh*-questions

- Self-paced reading study investigating the effect of *wh*-phrase type on the processing of binary superiority-violating *wh*-questions (in the target region)
- 2 (1st *wh*-phrase) x 2 (2nd *wh*-phrase) design:
 - BARE_BARE
Mary wondered *what* *who* read but later the teacher told her.
 - BARE_WHICH
Mary wondered *what* *which student* read but ...
 - WHICH_BARE
Mary wondered *which book* *who* read but ...
 - WHICH_WHICH
Mary wondered *which book* *which student* read but ...



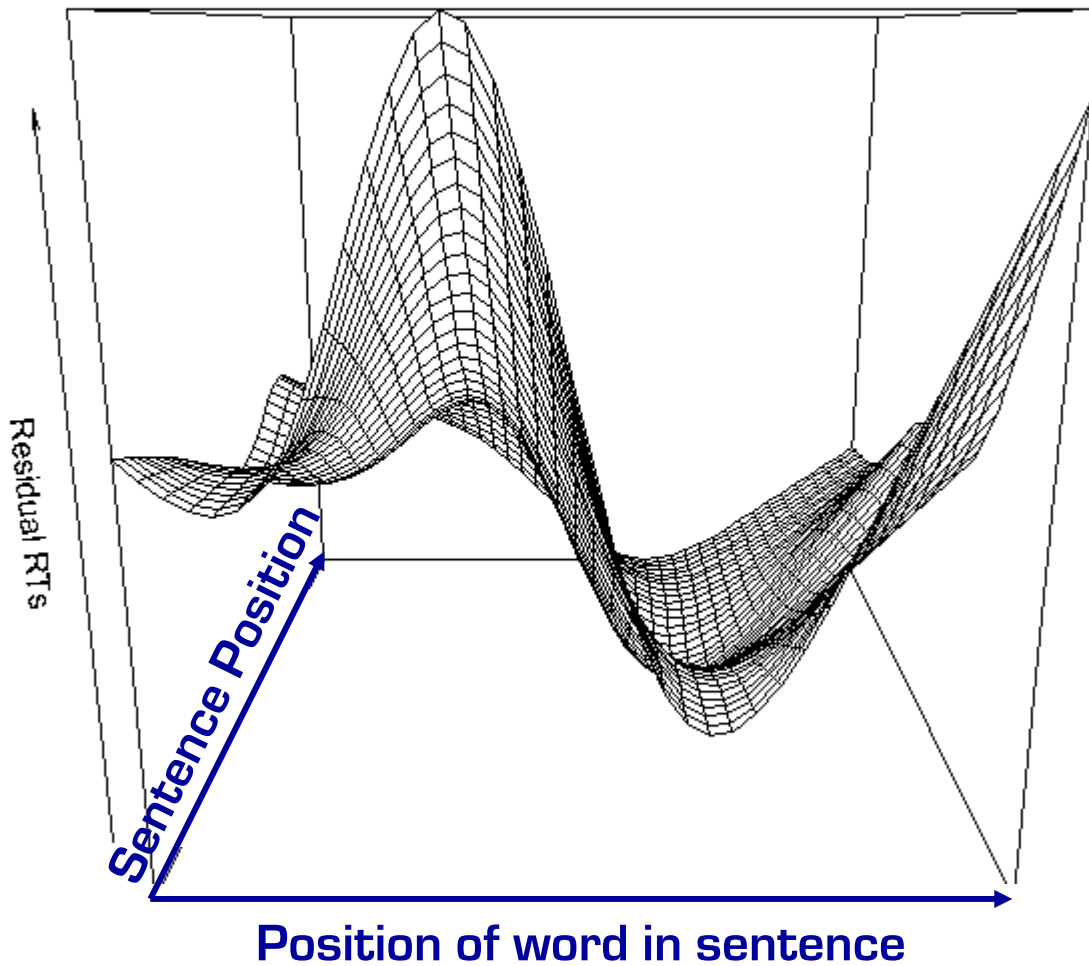
Participants read faster towards the end of the experiment until they reached a maximum reading speed. Between stimulus 40 and stimulus 55 residual reading times decreased by almost 50%.

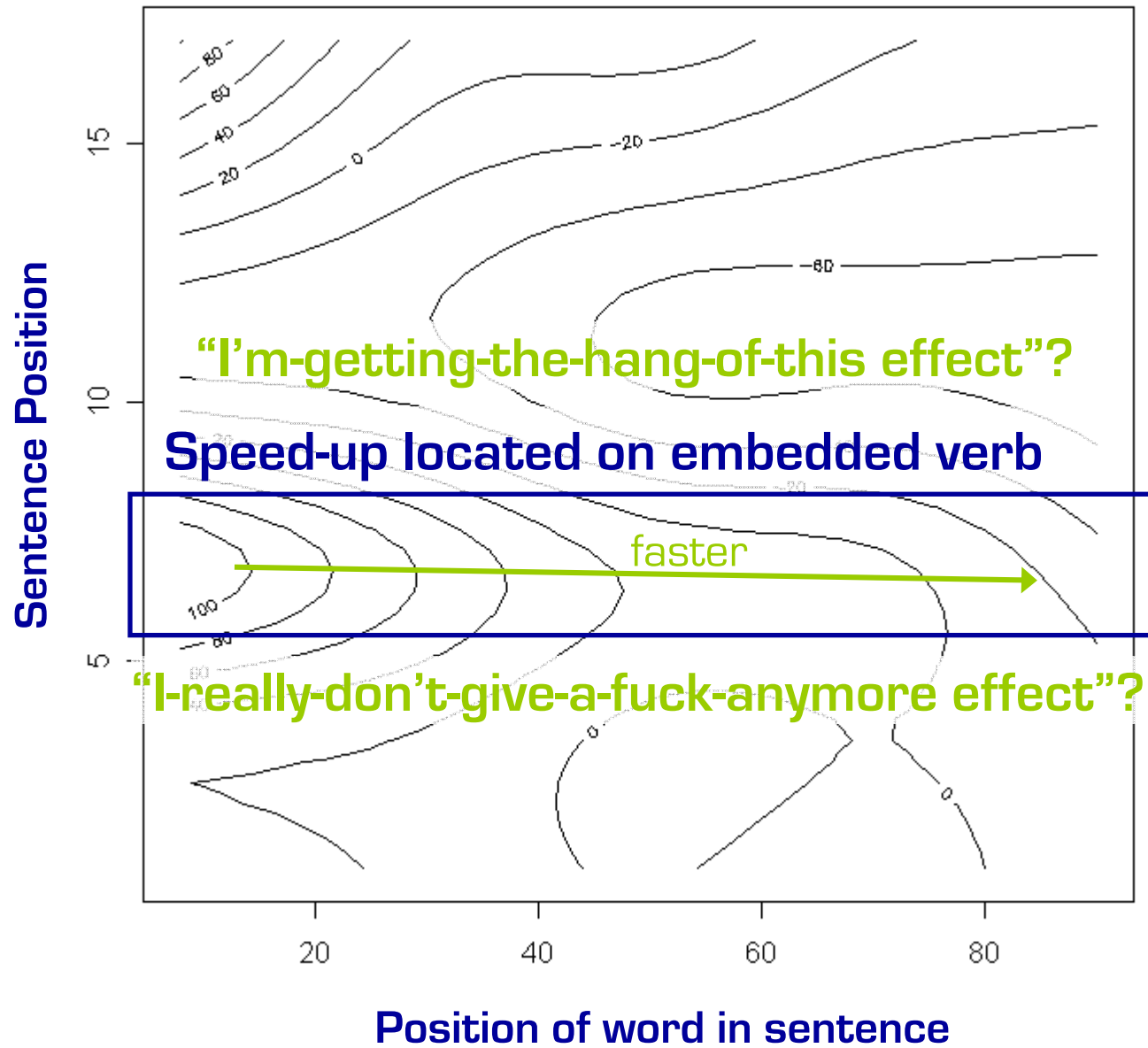


Position of sentence in list that the subject saw



But there's more to it than meets the eye. The position of a sentence in the list interact with position of the word in the sentence.





So what can be do?

- Why not sample from corpora when creating experiments, so that the stimuli are representative of natural language?
- The problem of unbalanced data can be addressed with modern statistical methods



UID's prediction for *that*-omission

- Speakers' decisions to pronounce or omit optional words should depend on processing complexity:
 - Optional words should be inserted, whenever processing would otherwise be difficult
 - Optional words should be omitted, whenever processing would be easy anyway

$$\log \frac{1}{p(RC \mid ctxt)}$$

**Phrase
onset**

$$\log \frac{1}{p(w_i \mid RC, ctxt)}$$

**Lexical
content**

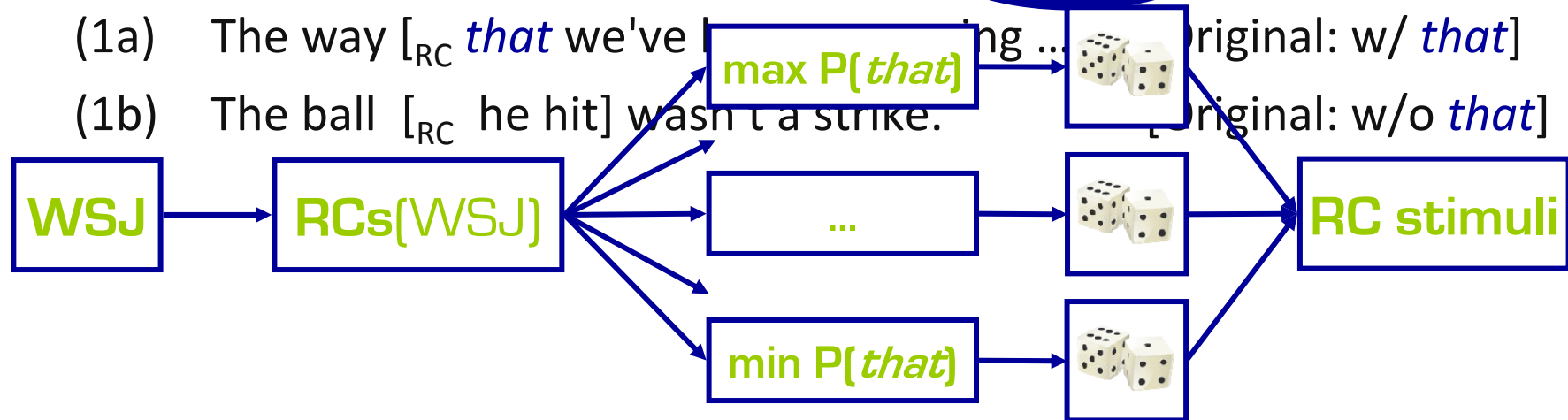
How big is the family { *that* *you* *cook for?* }
 { *you* *cook for?* }



Self-paced reading experiment :: Materials

- 24 RCs sampled from the parsed Wall Street Journal (Penn Treebank III release: Marcus et al. 1994) representative of the corpus in terms of RC frequency
 - 12 of with RC
 - 12 without RC

Used written corpus because this is a reading time study and because same patterns of *that*-omission have been observed in speech and writing (JaegerWasow05)



Self-paced reading experiment :: Materials

- (1a) The way [_{RC} *that* we've been managing ... [Original: w/ *that*]
(1b) The ball [_{RC} he hit] wasn't a strike. [Original: w/o *that*]

- For each stimulus, a matched stimulus was created that differed only in that-presence, (2a) for (1a) and (2b) for (1d):

- (2a) The way [_{RC} we've been managing ... [Original: w/ *that*]
(2b) The ball [_{RC} *that* he hit] wasn't a strike. [Original: w/o *that*]



Self-paced reading experiment :: Materials & Subjects

- 48 fillers (24 *wh*-question; 24 other types of stimuli resembling the experimental stimuli in complexity)
- 37 participants



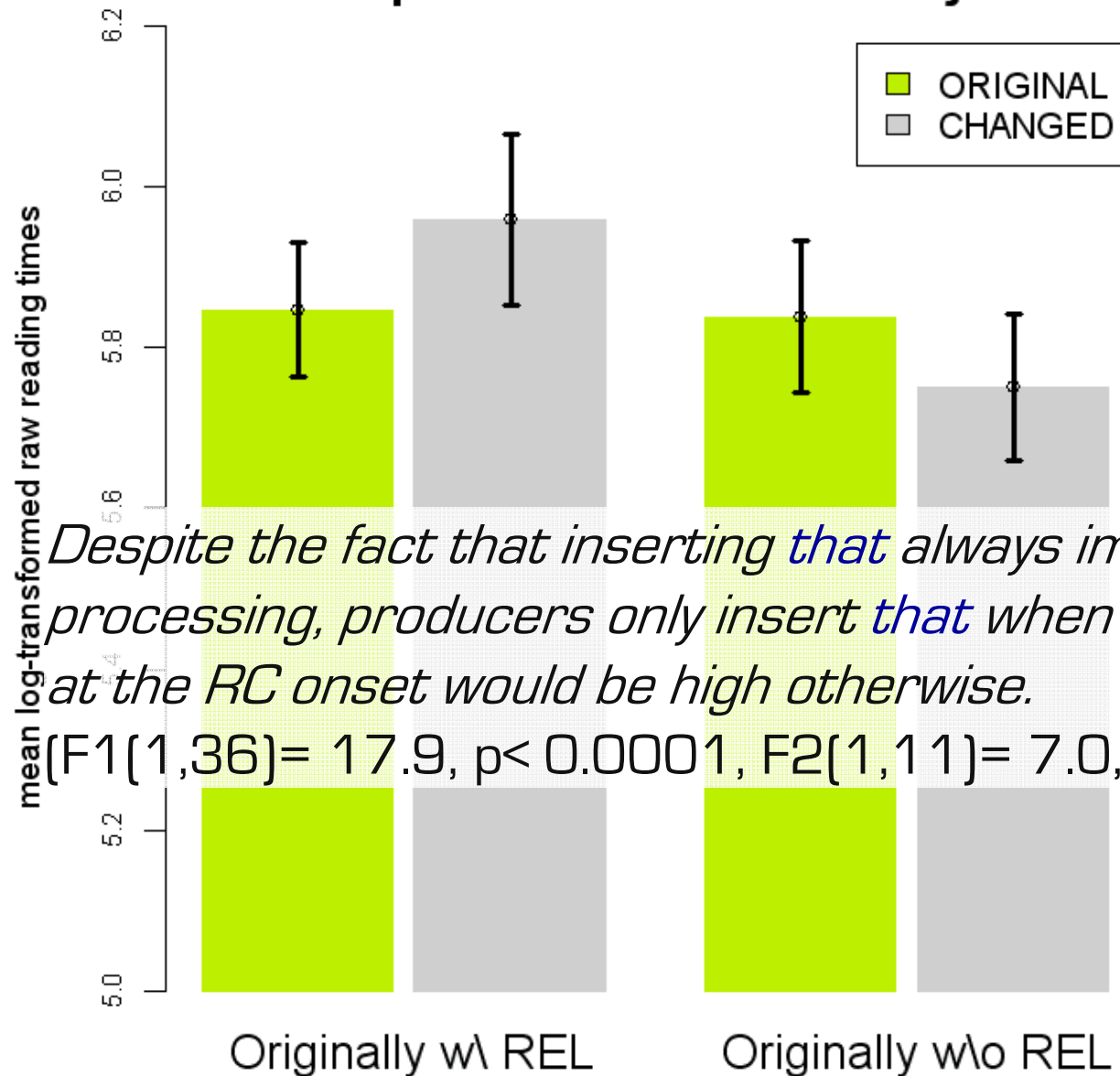
Self-paced reading experiment :: Procedure

- Stimuli presented word-by-word moving window

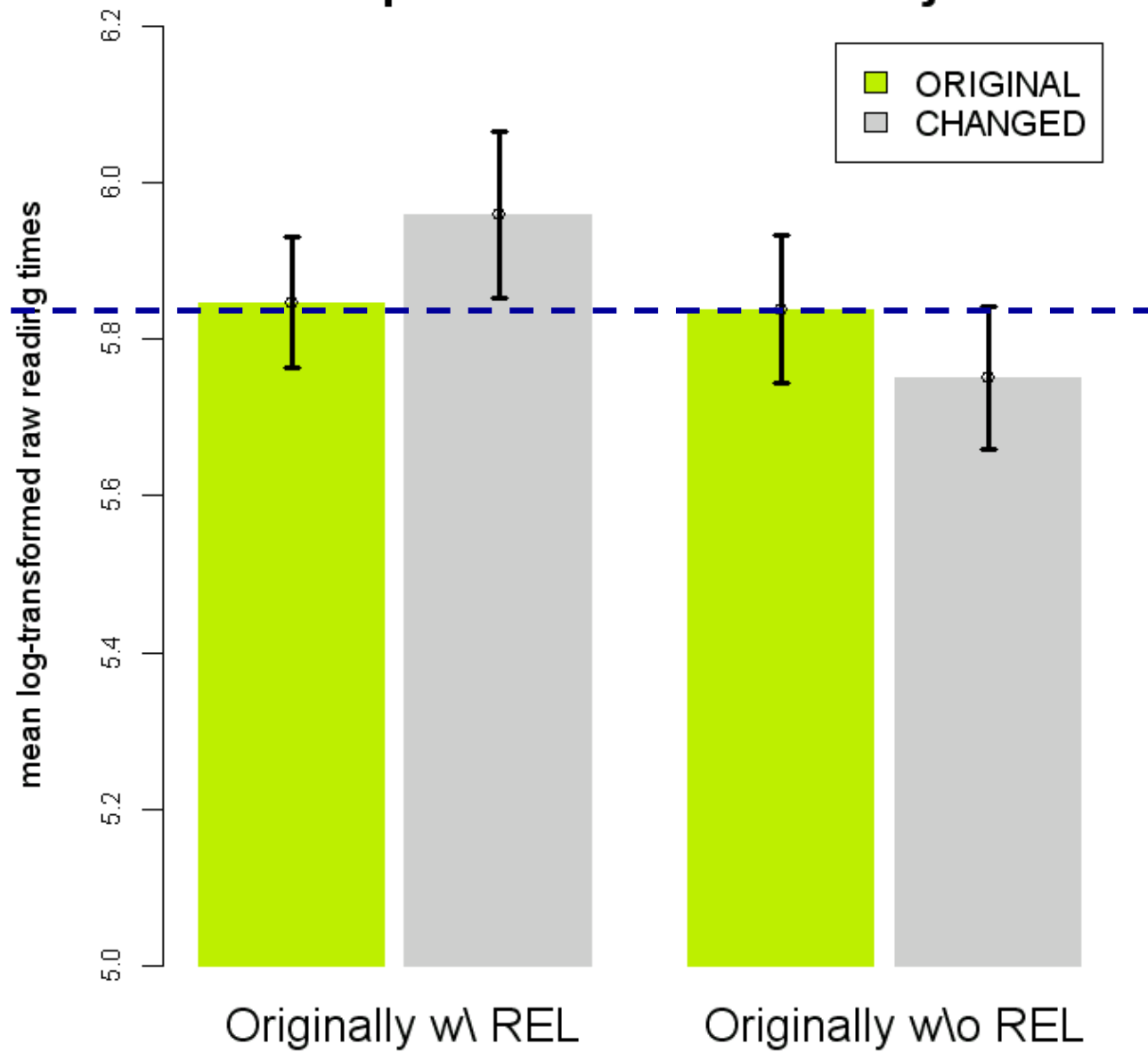
The way we've been managing



Comprehension of RC subject



Comprehension of RC subject



What else is out there?

- Treebanks – online and offline
- Automatically parsed corpora
- Non-syntactic online corpora with interesting search interfaces
- Other syntactic search tools: TigerSearch



Available Treebanks

- Offline

<http://www.ims.uni-stuttgart.de/projekte/TIGER/related/links.shtml>

or search LDC catalogue: <http://www ldc.upenn.edu/>

- American English (spoken & written)
- British English (spoken & written)
- German, Nordic Treebanks
- Mandarin Chinese, Korean, Japanese
- Czech, Bulgarian (planned), Russian
- Standard Arabic



Relative clauses in Mandarin Chinese :: Corpus

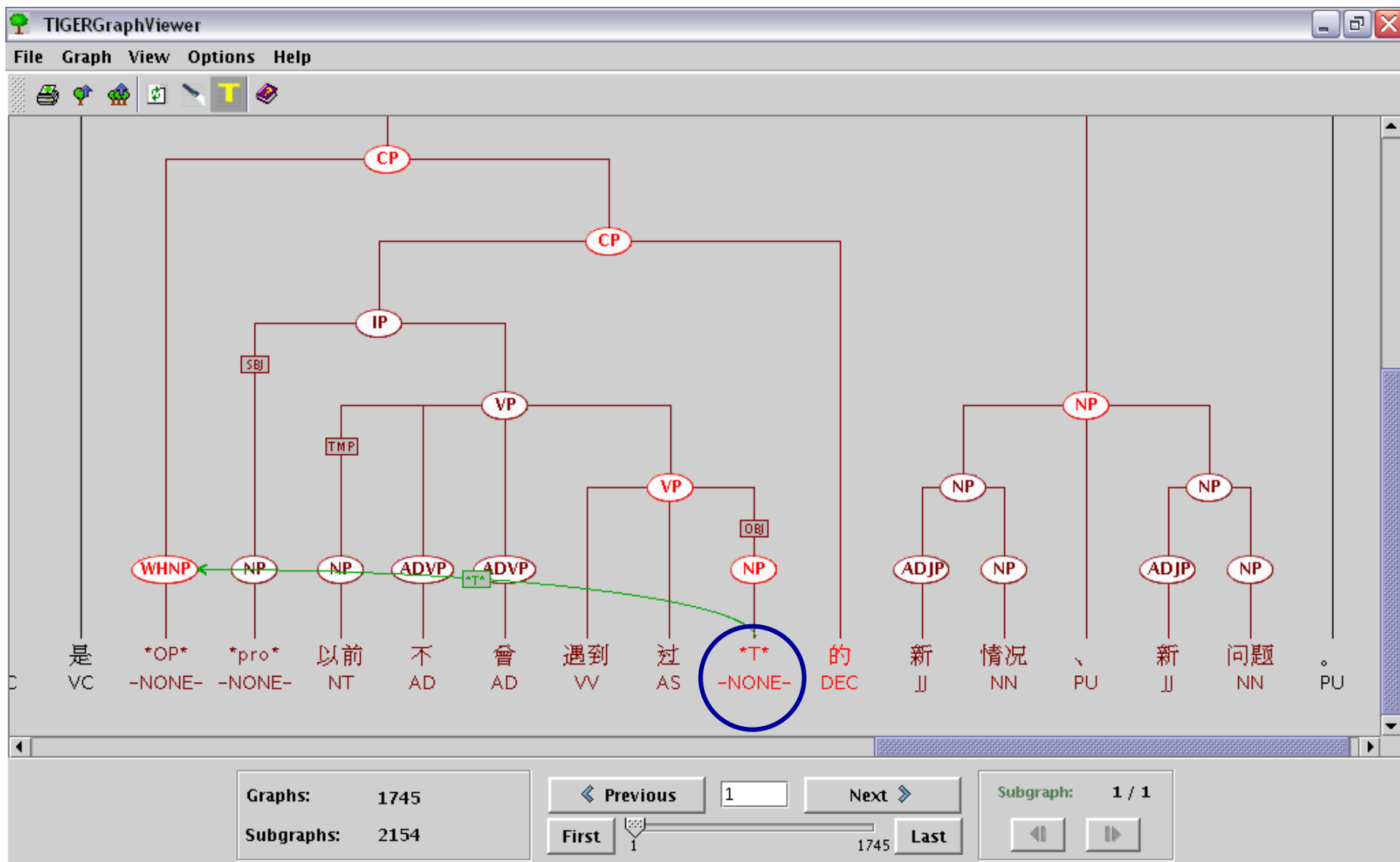
- Chinese Treebank (version 5.1), distributed by LDC
<http://www.cis.upenn.edu/~chinese/>
 - 18,525 graphs (\approx sentences)
 - 522,683 terminals (\approx words)
- Segmented
- Annotated for:
 - Part-of-speech
 - Syntactic bracketing (traces, fillers, PRO)
 - Functional structure (subject, topic, etc.)
 - In progress: *Semantic role labeling, discourse structure*



Mandarin *de*-omission

- Mandarin Chinese relative clauses are not always marked with a complementizer (*DE*).
 - $[_{NP} [_{RC} \dots de] \dots N]$
- Used **TigerSearch** to extract all 10,546 headed RCs (<http://www.ims.uni-stuttgart.de/projekte/TIGER/TIGERSearch/>)
 - 10,401 have *de*
 - 145 don't have *de*
 - *omission is rare* (1.4%; cf. 60-80% in Standard American English written corpora)





s4: 浦东开发开放是一项 *OP* *T* 振兴上海，建设现代化经济、贸易、金融中心的跨世纪工程，因此 *OP* *T* 大量出现的是 *OP* *pro* 以前不曾遇到过 *T* 的新情况、新问题。

Displaying matches (1745 matching corpus graphs, 2154 matching subgraphs).

Results :: *DE*-rate by type of extracted element

Extraction type	Total N	Percent <i>DE</i>
<i>Topic NP</i>	747	99.2%
<i>Subject NP</i>	8,002	98.9%
<i>Object NP</i>	2,002	97.8%
<i>PP</i>	263	93.2%

× 1.3

× 2.3

× 3.1



Available Treebanks (2)

- Smaller online corpora (e.g. <http://corp.hum.sdu.dk/>) also available for e.g.
 - Danish
 - Estonian
 - Portuguese
 - French
 - Esperanto



Automatically Parsed Corpora

- Beware of statistical biases, e.g. 100 million words of British English
 - Over exclusion
 - Over inclusion
 - NB: I would like to send you guys a paper some of the coming weeks that talks about this issue.
- Additional annotation may be needed



POS-annotated corpora

- CQP



Online Flat Text Corpora

- Web as corpus: <http://www.webcorp.org.uk/>
 - Limited regular expression syntax
- BYU collection of online searchable corpora (American and British English, Spanish, Portuguese) -- <http://view.byu.edu/>
 - Limited regular expression syntax
 - Part of speech
 - Lemmatizer
 - Synonym search
- Automated Google searches:
<http://www.linguistics.ucla.edu/people/hayes/QueryGoogle/qgapplet.html>



More links?

- David Lee's awesome list of links to corpora and tools:
<http://personal.cityu.edu.hk/~davidlee/devotedtocorpora/CBLLinks.htm>
- LDC pages: <http://www ldc.upenn.edu/>
- International corpus linguistics email list:
<http://gandalf.aksis.uib.no/corpora/>

