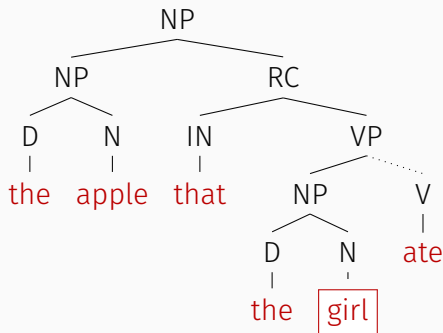


HIERARCHIC SYNTAX IMPROVES READING TIME PREDICTION

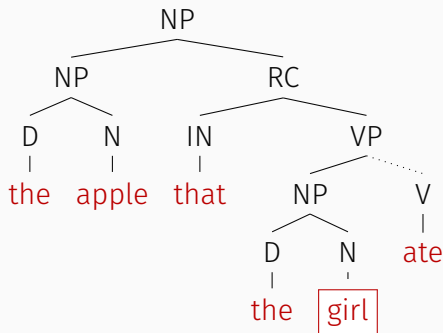
Marten van Schijndel and William Schuler
Department of Linguistics
The Ohio State University
June 3, 2015

Previous studies have debated whether humans use hierarchic syntax

Previous studies have debated whether humans use hierarchic syntax



Previous studies have debated whether humans use hierarchic syntax



But standard baseline predictors may be deficient

This work shows that:

This work shows that:

Baselines can be greatly improved (accumulation)

This work shows that:

Baselines can be greatly improved (accumulation)

Hierarchic syntax is still predictive over stronger baseline

This work shows that:

Baselines can be greatly improved (accumulation)

Hierarchical syntax is still predictive over stronger baseline

Hierarchical syntax not improved by accumulation

This work shows that:

Baselines can be greatly improved (accumulation)

Hierarchic syntax is still predictive over stronger baseline

Hierarchic syntax not improved by accumulation

Long distance dependencies independently improve model

HIERARCHIC SYNTAX IN READING?

The red apple that the ¹girl ²ate ...

FRANK & BOD (2011)

The red apple that the ¹girl² ate ...

FRANK & BOD (2011)

Baseline:

- Sentence Position
- Word length
- N-grams (Unigram, bigram)

The red apple that the girl ate ...
 w_1 w_2 w_3 w_4 w_5 w_6

FRANK & BOD (2011)

Baseline:

- Sentence Position
- Word length
- N-grams (Unigram, bigram)

The red apple that the girl ate ...

4 chars
W₆

FRANK & BOD (2011)

Baseline:

- Sentence Position
- Word length
- N-grams (Unigram, bigram)

The red apple that the girl ate ...

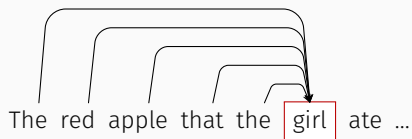
4 chars
W₆

FRANK & BOD (2011)

Baseline:

- Sentence Position
- Word length
- N-grams (Unigram, bigram)

HIERARCHIC SYNTAX IN READING?



FRANK & BOD (2011)

Baseline:

- Sentence Position
- Word length
- N-grams (Unigram, bigram)

Test POS Predictors:

- Echo State Network (ESN)
- Phrase Structure Grammar (PSG)

HIERARCHIC SYNTAX IN READING?



FRANK & BOD (2011)

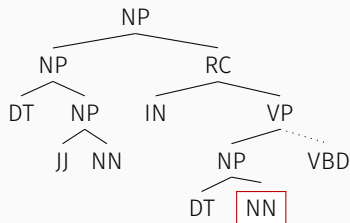
Baseline:

- Sentence Position
- Word length
- *N*-grams (Unigram, bigram)

Test POS Predictors:

- Echo State Network (ESN)
- Phrase Structure Grammar (PSG)

HIERARCHIC SYNTAX IN READING?



FRANK & BOD (2011)

Baseline:

- Sentence Position
- Word length
- *N*-grams (Unigram, bigram)

Test POS Predictors:

- Echo State Network (ESN)
- Phrase Structure Grammar (PSG)

FRANK & BOD (2011)

Baseline:

- Sentence Position
- Word length
- *N*-grams (Unigram, bigram)

Test POS Predictors:

- Echo State Network (ESN)
- Phrase Structure Grammar (PSG)

FRANK & BOD (2011)

Baseline:

- Sentence Position
- Word length
- *N*-grams (Unigram, bigram)

Test POS Predictors:

- Echo State Network (ESN)
- Phrase Structure Grammar (PSG)

Outcome:

$PSG < ESN + PSG$

$ESN = ESN + PSG$

FRANK & BOD (2011)

Baseline:

- Sentence Position
- Word length
- *N*-grams (Unigram, bigram)

Test POS Predictors:

- Echo State Network (ESN)
- Phrase Structure Grammar (PSG)

Outcome:

$PSG < ESN + PSG$ Sequential helps over hierarchic

$ESN = ESN + PSG$

FRANK & BOD (2011)

Baseline:

- Sentence Position
- Word length
- *N*-grams (Unigram, bigram)

Test POS Predictors:

- Echo State Network (ESN)
- Phrase Structure Grammar (PSG)

Outcome:

$PSG < ESN + PSG$

$ESN = ESN + PSG$ Hierarchic doesn't help over sequential

FOSSUM & LEVY (2012)

Replicated Frank & Bod (2011):

PSG < ESN + PSG

ESN = ESN + PSG

FOSSUM & LEVY (2012)

Replicated Frank & Bod (2011):

PSG < ESN + PSG

ESN = ESN + PSG

Better *n*-gram baseline (more data) changes result:

PSG ESN + PSG

ESN = ESN + PSG

FOSSUM & LEVY (2012)

Replicated Frank & Bod (2011):

PSG < ESN + PSG

ESN = ESN + PSG

Better *n*-gram baseline (more data) changes result:

PSG ESN + PSG Sequential doesn't help over hierarchic

ESN = ESN + PSG

FOSSUM & LEVY (2012)

Replicated Frank & Bod (2011):

PSG < ESN + PSG

ESN = ESN + PSG

Better *n*-gram baseline (more data) changes result:

PSG \equiv ESN + PSG Sequential doesn't help over hierarchic

ESN = ESN + PSG

Also: lexicalized syntax improves PSG fit

Previous reading time studies:

- Unigrams/Bigrams/Trigrams
Trained on WSJ, Dundee, BNC


Previous reading time studies:

- Unigrams/Bigrams/Trigrams
Trained on WSJ, Dundee, BNC
- Only from region boundaries

BIGRAM EXAMPLE

Reading time of *girl* after *red*

The ¹red apple that the ²girl ate ...




region

X: bigram target X: bigram condition

BIGRAM EXAMPLE

Reading time of *girl* after *red*

The ¹red apple that the ²girl ate ...



region


X: bigram target X: bigram condition

- Fails to capture entire sequence;

BIGRAM EXAMPLE

Reading time of *girl* after *red*

The ¹red apple that the ²girl ate ...



region


X: bigram target X: bigram condition

- Fails to capture entire sequence;
- Conditions never generated;

BIGRAM EXAMPLE

Reading time of *girl* after *red*

The ¹red apple that the ²girl ate ...



region

X: bigram target X: bigram condition

- Fails to capture entire sequence;
- Conditions never generated;
- Probability of sequence is deficient

CUMULATIVE BIGRAM EXAMPLE

Reading time of *girl* after *red*:

The red¹ apple that the girl² ate ...

X: bigram targets X: bigram conditions

CUMULATIVE BIGRAM EXAMPLE

Reading time of *girl* after *red*:

The ¹red apple that the ²girl ate ...

X: bigram targets X: bigram conditions

- Captures entire sequence;
- Well-formed sequence probability;
- Reflects processing that must be done by humans

Previous reading time studies:

- Unigrams/Bigrams/Trigrams
- Trained on WSJ, Dundee, BNC
- Only from region boundaries

Previous reading time studies:

- Unigrams/Bigrams/Trigrams
- Trained on WSJ, Dundee, BNC
- Only from region boundaries

This study:

- 5-grams (w/ backoff)
- Trained on Gigaword 4.0
- Cumulative and Non-cumulative

Dundee Corpus (Kennedy et al., 2003)

- 10 subjects
- 2,388 sentences
- 58,439 words
- 194,882 first pass durations
- 193,709 go-past durations

Exclusions:

- Unknown words (5 tokens)
- First and last of a line
- Regions larger than 4 words (track loss)

Baseline:

Fixed Effects

- Sentence Position
- Word length
- Region Length
- Preceding word fixated?

Random Effects

- Item/Subject Intercepts
- By Subject Slopes:
 - All Fixed Effects
 - N -grams (5-grams)
 - N -grams (Cumulative-5-grams)

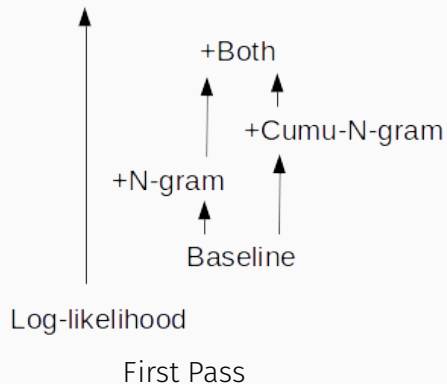
Baseline:

Fixed Effects

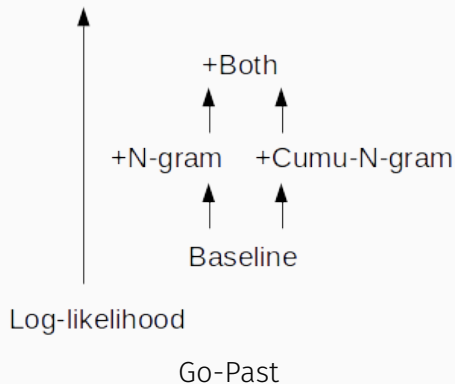
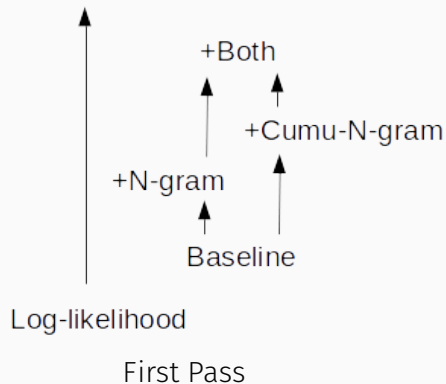
- Sentence Position
- Word length
- Region Length
- Preceding word fixated?

Random Effects

- Item/Subject Intercepts
- By Subject Slopes:
 - All Fixed Effects
 - N -grams (5-grams) ←
 - N -grams (Cumulative-5-grams) ←



CUMU- n -GRAMS PREDICT READING TIMES



- Is hierarchic surprisal useful over the better baseline?

- Is hierarchic surprisal useful over the better baseline?
- If so, can it be similarly improved through accumulation?

- Is hierarchic surprisal useful over the better baseline?
- If so, can it be similarly improved through accumulation?
van Schijndel & Schuler (2013) found it could over weaker baselines

Grammar:

Berkeley parser, WSJ, 5 split-merge cycles (Petrov & Klein 2007)

Baseline:

Fixed Effects

- Same as before
- *N*-grams (5-grams)
- *N*-grams (Cumulative-5-grams)

Baseline:

Fixed Effects

- Same as before
- *N*-grams (5-grams)
- *N*-grams (Cumulative-5-grams)

Random Effects

- Same as before
- By Subject Slopes:
 - Hierarchic surprisal
 - Cumulative-Hierarchic surprisal

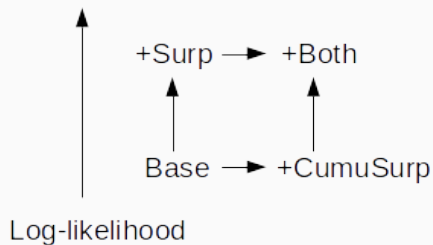
Baseline:

Fixed Effects

- Same as before
- *N*-grams (5-grams)
- *N*-grams (Cumulative-5-grams)

Random Effects

- Same as before
- By Subject Slopes:
 - Hierarchic surprisal ←
 - Cumulative-Hierarchic surprisal ←



First Pass and Go-Past

- Suggests previous findings were due to weaker n -gram baseline

- Suggests previous findings were due to weaker n -gram baseline
- Suggests only local PCFG surprisal affects reading times

- Suggests previous findings were due to weaker n -gram baseline
- Suggests only local PCFG surprisal affects reading times

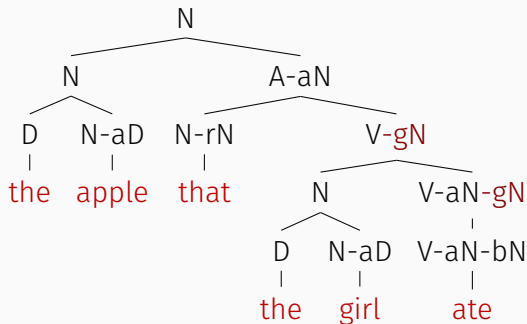
But... long-distance dependencies should affect reading times!

- Suggests previous findings were due to weaker n -gram baseline
- Suggests only local PCFG surprisal affects reading times

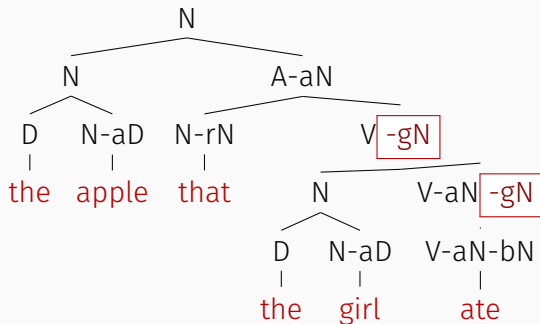
But... long-distance dependencies should affect reading times!

Let's try a PCFG that tracks long-distance deps

Nguyen et al. (2012)



Nguyen et al. (2012)



Baseline:

Fixed Effects

- Same as before

Random Effects

- Same as before
- By Subject Slopes:
 - Hierarchic PTB surprisal
 - Hierarchic GCG surprisal

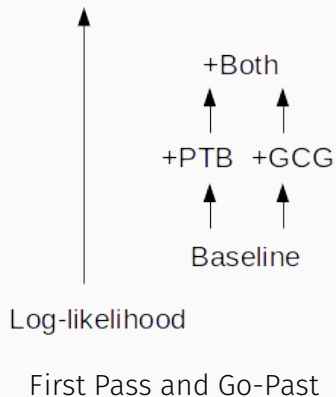
Baseline:

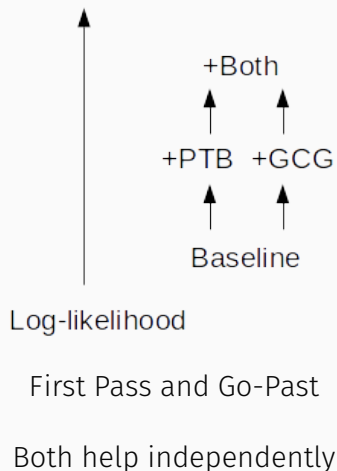
Fixed Effects

- Same as before

Random Effects

- Same as before
- By Subject Slopes:
 - Hierarchic PTB surprisal ←
 - Hierarchic GCG surprisal ←





Hierarchic syntax predicts reading times over strong linear baseline

Hierarchic syntax predicts reading times over strong linear baseline

Long-distance dependencies *do* affect reading times

Hierarchical syntax predicts reading times over strong linear baseline

Long-distance dependencies *do* affect reading times

Studies should use cumu-*n*-grams in their baselines

Compare to Echo State Networks

Compare to Echo State Networks

Test anticipatory accumulation

Thanks to:

- Stefan Frank
- Attendees of CUNY 2015
- National Science Foundation (DGE-1343012)

First Pass Evaluation (Log-Likelihood):

Base -1212399	
Base+N-gram -1212396 ($p < 0.05$)	Base+Cumu-n-gram -1212392 ($p < 0.01$)
Base+Both -1212387 ($p < 0.01$)	Base+Both -1212387 ($p < 0.01$)

Comparable with go-past durations

Go-Past Evaluation (Log-Likelihood):

Base -1261582	
Base+N-gram -1261577 ($p < 0.01$)	Base+Cumu-n-gram -1261576 ($p < 0.01$)
Base+Both -1261570 ($p < 0.01$)	Base+Both -1261570 ($p < 0.01$)

First Pass Evaluation (Log-Likelihood):

Base -1212260	
Base+Surp -1212253 ($p < 0.01$)	Base+CumuSurp -1212259
Base+Both -1212253	Base+Both -1212253 ($p < 0.01$)

Comparable with go-past durations

Go-Past Evaluation (Log-Likelihood):

Base -1261488	
Base+Surp -1261481 ($p < 0.01$)	Base+CumuSurp -1261487
Base+Both -1261481	Base+Both -1261481 ($p < 0.01$)

First Pass Evaluation (Log-Likelihood):

Base	
-1212242	
Base+PTB	Base+GCG
-1212239 ($p < 0.01$)	-1212239 ($p < 0.05$)
Base+Both	Base+Both
-1212235 ($p < 0.05$)	-1212235 ($p < 0.01$)

Both help independently

PCFG surprisal helps more with go-past durations

Go-Past Evaluation (Log-Likelihood):

Base −1261474	
Base+PTB −1261468 ($p < 0.01$)	Base+GCG −1261470 ($p < 0.01$)
Base+Both −1261465 ($p < 0.01$)	Base+Both −1261465 ($p < 0.01$)

Again, both help independently.

FIXED EFFECT COEFFICIENTS FOR BASE+PTB+GCG

Predictor	First Pass		Go-Past	
	coef	t value	coef	t value
sentpos	-2.47	-3.59	-2.82	-3.38
wlen	25.90	8.67	28.98	9.97
prevfix	-30.16	-7.81	-37.42	-11.49
<i>n</i> -gram	-2.39	-1.81	-6.70	-3.36
cumu- <i>n</i> -gram	-14.69	-7.36	-11.68	-5.01
rln	-5.67	-1.31	-12.51	-2.59
surp-GCG	4.97	2.87	5.74	2.73
surp-PTB	4.20	3.23	4.85	3.29