# Addressing surprisal deficiencies in reading time models

Marten van Schijndel    William Schuler
December 11, 2016

Department of Linguistics, The Ohio State University

- Surprisal is a way to estimate text complexity

- Surprisal is a way to estimate text complexity
- Experienced complexity is reflected in reading speed

- Surprisal is a way to estimate text complexity
- Experienced complexity is reflected in reading speed

Claim:
Current surprisal models inadequately estimate reading complexity

- Surprisal is a way to estimate text complexity
- Experienced complexity is reflected in reading speed

Claim:
Current surprisal models inadequately estimate reading complexity

Consequence:
Other reading time predictors may get too much credit

The red apple that the girl ate …

The red apple that the $\boxed{\text{girl}}$ ate ...
$w_1$ $w_2$ $w_3$ $w_4$ $w_5$ $w_6$

Reading model of 'girl':
sentence position

4 chars

The red apple that the $w_6$ girl ate …

Reading model of 'girl':
sentence position, word length

The red apple that <u>the</u> girl ate …

Reading model of 'girl':
sentence position, word length, P(girl|the)

Reading model of 'girl':
sentence position, word length, P(girl|the)

The red apple that the girl ate ...

Reading model of 'girl':
sentence position, word length, P(girl|the)

Reading model of 'girl':
sentence position, word length, P(girl|the)

This study: *n*-gram and PCFG surprisal

This study: *n*-gram and PCFG surprisal

The red apple that <u>the</u> girl ate …

$N$-gram-surp(girl) $= -\log$ P(girl | the)

This study: *n*-gram and PCFG surprisal



$$\text{PCFG-surp(girl)} = -\sum_{T \in \text{Trees}} \log P(T_6 = \text{girl} \mid T_1 \dots T_5 = \text{The} \dots \text{the})$$

Cumulative *N*-gram Surprisal

$$\overset{1}{\text{The red apple that the}} \overset{2}{\text{girl ate}} \dots$$

Cumulative *N*-gram Surprisal

$$\text{The } \underset{1}{\underline{\text{red}}} \boxed{\text{apple}} \text{ that the } \overset{2}{\text{girl}} \text{ ate ...}$$

$$\text{cumu-}n\text{-gram}(w, f_{t-1}, f_t) = \sum_{i=f_{t-1}+1}^{f_t} -\log P(w_i \mid w_{i-n} \ldots w_{i-1})$$

Cumulative *N*-gram Surprisal

$$\text{The red} \overset{1}{\boxed{\text{apple}}} \overset{}{\boxed{\text{that}}} \text{the girl} \overset{2}{\text{ate}} \dots$$

$$\text{cumu-}n\text{-gram}(w, f_{t-1}, f_t) = \sum_{i=f_{t-1}+1}^{f_t} -\log P(w_i \mid w_{i-n} \dots w_{i-1})$$

Cumulative *N*-gram Surprisal

$$\text{The red } \overset{1}{\boxed{\text{apple}}} \boxed{\underline{\text{that}}} \boxed{\text{the}} \overset{2}{\text{girl ate ...}}$$

$$\text{cumu-}n\text{-gram}(w, f_{t-1}, f_t) = \sum_{i=f_{t-1}+1}^{f_t} -\log P(w_i \mid w_{i-n} \ldots w_{i-1})$$

Cumulative *N*-gram Surprisal

$$\text{The red } \overset{1}{\boxed{\text{apple}}} \boxed{\text{that}} \boxed{\underline{\text{the}}} \overset{2}{\boxed{\text{girl}}} \text{ ate } \ldots$$

$$\text{cumu-}n\text{-gram}(w, f_{t-1}, f_t) = \sum_{i=f_{t-1}+1}^{f_t} -\log P(w_i \mid w_{i-n} \ldots w_{i-1})$$

Cumulative PCFG Surprisal
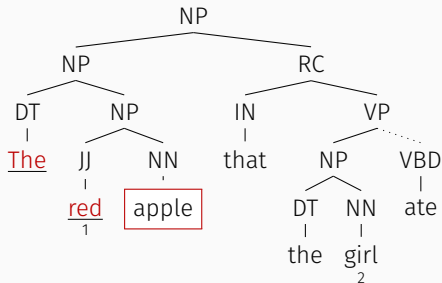


$$\text{Cumu-PCFG}(w, f_{t-1}, f_t) = \sum_{i=f_{t-1}}^{f_t} \sum_{T \in \text{Trees}} -\log P(T_i = w_i \mid T_1 \ldots T_{i-1} = w_1 \ldots w_{i-1})$$

Cumulative PCFG Surprisal



$$\text{Cumu-PCFG}(w, f_{t-1}, f_t) = \sum_{i=f_{t-1}}^{f_t} \sum_{T \in \text{Trees}} -\log P(T_i = w_i \mid T_1 \ldots T_{i-1} = w_1 \ldots w_{i-1})$$
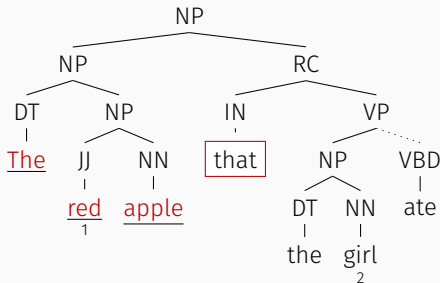
Cumulative PCFG Surprisal



$$\text{Cumu-PCFG}(w, f_{t-1}, f_t) = \sum_{i=f_{t-1}}^{f_t} \sum_{T \in \text{Trees}} -\log P(T_i = w_i \mid T_1 \dots T_{i-1} = w_1 \dots w_{i-1})$$
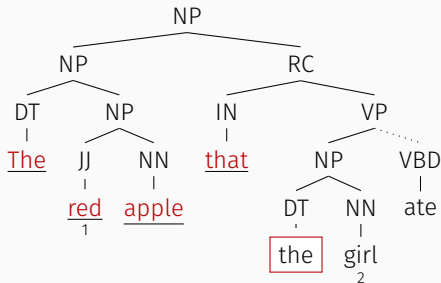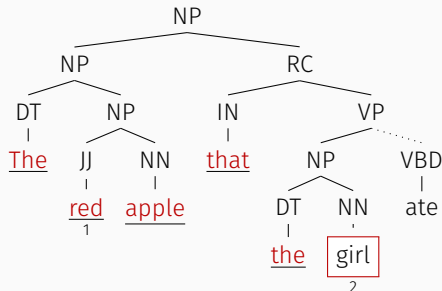
Cumulative PCFG Surprisal



$$\text{Cumu-PCFG}(w, f_{t-1}, f_t) = \sum_{i=f_{t-1}}^{f_t} \sum_{T \in \text{Trees}} -\log \mathsf{P}(T_i = w_i \mid T_1 \ldots T_{i-1} = w_1 \ldots w_{i-1})$$

# How well does this fix work?

*N*-gram surprisal

- 5-grams
- Trained on Gigaword 3.0 (Graff and Cieri, 2003)
- Computed with KenLM (Heafield et al., 2013)

## HOW WELL DOES THIS FIX WORK?

*N*-gram surprisal

- 5-grams
- Trained on Gigaword 3.0 (Graff and Cieri, 2003)
- Computed with KenLM (Heafield et al., 2013)

PCFG surprisal

- Trained on WSJ 02-21 (Marcus et al., 1993)
- Computed with van Schijndel et al., (2013) parser

University College London (UCL) Corpus (Frank et al., 2013)

- 43 subjects
- reading online novels
- frequent comprehension questions

Baseline mixed effects model

Fixed Factors

- sentence position
- word length
- region length
- whether the previous word was fixated

## HOW WELL DOES THIS FIX WORK?

Baseline mixed effects model

Fixed Factors

- sentence position
- word length
- region length
- whether the previous word was fixated

Random Factors

- All fixed factors as by-subject random slopes
- Item, subject and subject×sentence intercepts

Captured
reading time
variance

Cumu-N-grams

N-grams

10.61 ms          8.69 ms

Baseline

Cumu-N-grams

N-grams

Captured
reading time
variance

10.61 ms

6.69 ms

Baseline

After adding cumulative *n*-gram surprisal to model:

After adding cumulative *n*-gram surprisal to model:

- PCFG surprisal is not useful (p > 0.05)

## Accumulation does not help PCFG surprisal

After adding cumulative *n*-gram surprisal to model:

- PCFG surprisal is not useful (p > 0.05)
- Cumulative PCFG surprisal is not useful (p > 0.05)

Parafovial processing

1
The red apple that the girl ate …

Parafovial processing

Th(e red apple that t)he girl ate ...
[1]

Parafovial processing

$$\text{Th(e red apple that t)he girl ate ...}$$
<small>1</small> ... <small>2</small>

Prediction

The red apple that the girl ate …

Prediction

The re[1]d (apple that the girl) ate ...

Prediction

The re$\overset{1}{\text{d}}$ (apple that the gi$\overset{2}{\text{rl}}$) ate …

Subsequent regression

The red apple that the girl ate …

Subsequent regression

The red apple that the girl ate ...

Subsequent regression

$$\overset{1}{\text{The}} \overset{3}{\text{red}} \text{apple that the } \overset{2}{\text{girl}} \text{ate ...}$$

Subsequent regression

<div align="center">

1    3    4         2
The red apple that the girl ate …

</div>

Subsequent regression

<div align="center">

1    3     4       2   5

The red apple that the girl ate …

</div>

Cumulative PCFG surprisal only handles subsequent regression

Cumulative PCFG surprisal only handles subsequent regression

Parafovial: Th(e red apple that t)he girl ate …

Prediction: The red (apple that the girl) ate …

accumulated

Cumulative PCFG surprisal only handles subsequent regression

Parafovial: Th(e red apple that t)he girl ate …
<br>
$\overset{1}{}$ ... $\overset{2}{}$

Prediction: The red (apple that the girl) ate …
<br>
accumulated

Other accumulation mechanisms presuppose earlier accumulation

Upcoming material influences reading times

Upcoming material influences reading times

- Orthographic effects
  (Pynte, Kennedy, & Ducrot, 2004; Angele, Tran, & Rayner, 2013)

Upcoming material influences reading times

- Orthographic effects
  (Pynte, Kennedy, & Ducrot, 2004; Angele, Tran, & Rayner, 2013)
- Lexical effects
  (Kliegl et al., 2006; Li et al., 2014; Angele et al., 2015)

$$\overset{1}{\text{The red}} \text{ apple that the } \overset{2}{\text{girl}} \text{ ate } ...$$

$$\text{future-}n\text{-gram}(w, f_t, f_{t+1}) = \sum_{i=f_t}^{f_{t+1}} -\log \mathsf{P}(w_i \mid w_{i-n} \dots w_{i-1})$$

$$\text{The } \underset{1}{\underline{\text{red}}} \boxed{\text{apple}}^{2} \text{ that the girl ate ...}$$

$$\text{future-}n\text{-gram}(w, f_t, f_{t+1}) = \sum_{i=f_t}^{f_{t+1}} -\log \mathsf{P}(w_i \mid w_{i-n} \ldots w_{i-1})$$

The red $\boxed{\text{apple}}$ $\overset{2}{\boxed{\text{that}}}$ the girl ate …

(the word "red" is marked with a 1 above it)

$$\text{future-}n\text{-gram}(w, f_t, f_{t+1}) = \sum_{i=f_t}^{f_{t+1}} -\log P(w_i \mid w_{i-n} \ldots w_{i-1})$$

$$\text{The red} \underset{1}{\boxed{\text{apple}}}\ \boxed{\underline{\text{that}}}\ \underset{2}{\boxed{\text{the}}}\ \text{girl ate ...}$$

$$\text{future-}n\text{-gram}(w, f_t, f_{t+1}) = \sum_{i=f_t}^{f_{t+1}} -\log \mathsf{P}(w_i \mid w_{i-n} \ldots w_{i-1})$$

$$\text{The red} \overset{1}{\boxed{\text{apple}}} \boxed{\text{that}} \boxed{\underline{\text{the}}} \overset{2}{\boxed{\text{girl}}} \text{ate ...}$$

$$\text{future-}n\text{-gram}(w, f_t, f_{t+1}) = \sum_{i=f_t}^{f_{t+1}} -\log \mathsf{P}(w_i \mid w_{i-n} \ldots w_{i-1})$$

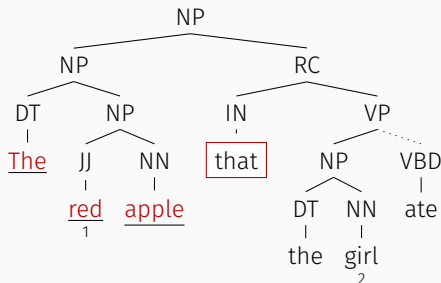$$\text{Future-PCFG}(w, f_t, f_{t+1}) = \sum_{i=f_t}^{f_{t+1}} \sum_{T \in \text{Trees}} -\log P(T_i = w_i \mid T_1 \dots T_{i-1} = w_1 \dots w_{i-1})$$

$$\text{Future-PCFG}(w, f_t, f_{t+1}) = \sum_{i=f_t}^{f_{t+1}} \sum_{T \in \text{Trees}} -\log P(T_i = w_i \mid T_1 \ldots T_{i-1} = w_1 \ldots w_{i-1})$$

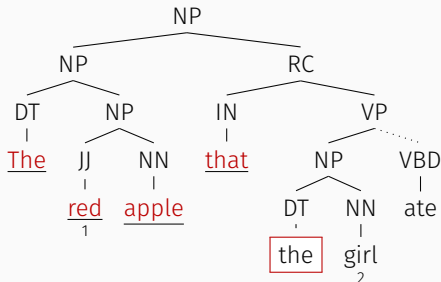$$\text{Future-PCFG}(w, f_t, f_{t+1}) = \sum_{i=f_t}^{f_{t+1}} \sum_{T \in \text{Trees}} -\log P(T_i = w_i \mid T_1 \ldots T_{i-1} = w_1 \ldots w_{i-1})$$
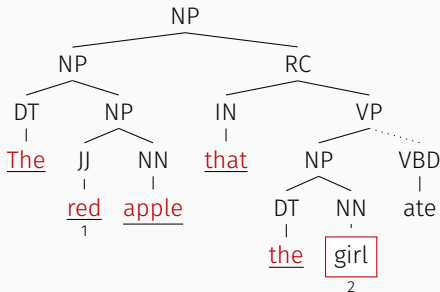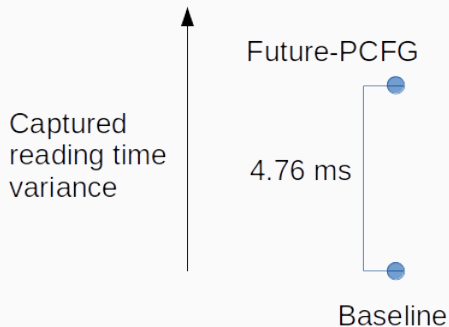
$$\text{Future-PCFG}(w, f_t, f_{t+1}) = \sum_{i=f_t}^{f_{t+1}} \sum_{T \in \text{Trees}} -\log P(T_i = w_i \mid T_1 \ldots T_{i-1} = w_1 \ldots w_{i-1})$$

Future-PCFG

Captured
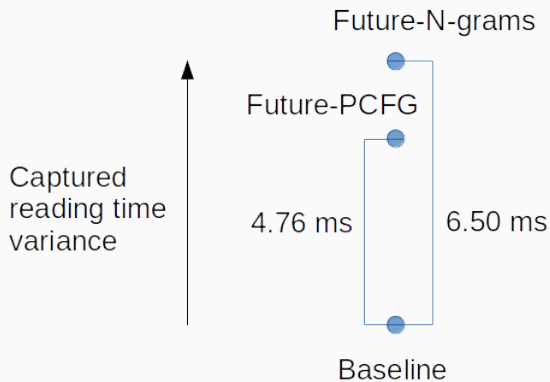reading time
variance

4.76 ms

Baseline

Future-N-grams

Future-PCFG

Captured reading time variance

4.76 ms

6.50 ms

Baseline

PCFG surprisal may require a richer grammar

Successor *n*-grams are most predictive for 2 future words (p < 0.001)

Successor *n*-grams are most predictive for 2 future words (p < 0.001)

6% of UCL saccades (n=3500) >2 words

- *N*-gram surprisal should be accumulated to predict reading times

- *N*-gram surprisal should be accumulated to predict reading times
- *N*-gram surprisal accumulates pre- and post-saccade

- *N*-gram surprisal should be accumulated to predict reading times
- *N*-gram surprisal accumulates pre- and post-saccade
    - Pre-saccade *n*-grams are limited

- *N*-gram surprisal should be accumulated to predict reading times
- *N*-gram surprisal accumulates pre- and post-saccade
    - Pre-saccade *n*-grams are limited
- PCFG surprisal does not accumulate

- *N*-gram surprisal should be accumulated to predict reading times
- *N*-gram surprisal accumulates pre- and post-saccade
    - Pre-saccade *n*-grams are limited
- PCFG surprisal does not accumulate
    - PCFG surprisal still predictive on Dundee

- *N*-gram surprisal should be accumulated to predict reading times
- *N*-gram surprisal accumulates pre- and post-saccade
    - Pre-saccade *n*-grams are limited
- PCFG surprisal does not accumulate
    - PCFG surprisal still predictive on Dundee
    - UCL corpus may not be syntactically complex enough

## Conclusion

- *N*-gram surprisal should be accumulated to predict reading times
- *N*-gram surprisal accumulates pre- and post-saccade
  - Pre-saccade *n*-grams are limited
- PCFG surprisal does not accumulate
  - PCFG surprisal still predictive on Dundee
  - UCL corpus may not be syntactically complex enough
  - PCFG surprisal may need a richer grammar

Thanks to:

- Stefan Frank
- National Science Foundation (DGE-1343012)

## Cumu-*N*-gram Results

| Model | *N*-gram vs Cumu-*N*-gram | | |
| | $\beta$ | Log-Likelihood | AIC |
| --- | --- | --- | --- |
| Baseline | | $-12702$ | 25476 |
| Base+Basic | 0.035 | $-12689^*$ | 25451 |
| Base+Cumulative | 0.055 | $-12683^*$ | 25440 |
| Base+Both | | $-12683^*$ | 25442 |

Base random: sentpos, wlen, rlen, prevfix, 5-gram, cumu-5-gram

Base fixed: sentpos, wlen, rlen, prevfix

Significance for the Base+Both model applies to improvement over the Base+Basic model.

| Model | Future-$N$-grams vs Future-PCFG | | |
|---|---|---|---|
| | $\beta$ | Log-Likelihood | AIC |
| Baseline | | $-12276$ | 24642 |
| Base+Future-$N$-grams | 0.034 | $-12259^*$ | 24610 |
| Base+Future-PCFG | 0.025 | $-12266^*$ | 24624 |
| Base+Both | | $-12259^*$ | 24612 |

Base random: sentpos, wlen, rlen, prevfix, cumu-5-gram, future-5-grams, future-PCFG
Base fixed: sentpos, wlen, rlen, prevfix, cumu-5-gram

Significance for the Base+Both model applies to improvement over the Base+Future-PCFG model.