

# ADDRESSING SURPRISAL DEFICIENCIES IN READING TIME MODELS

---

Marten van Schijndel   William Schuler

December 11, 2016

Department of Linguistics, The Ohio State University

- Surprisal (PCFG,  $N$ -gram) is a way to estimate text complexity

- Surprisal (PCFG,  $N$ -gram) is a way to estimate text complexity
- Experienced complexity is reflected in reading speed

- Surprisal (PCFG,  $N$ -gram) is a way to estimate text complexity
- Experienced complexity is reflected in reading speed

Claim:

Current surprisal models inadequately estimate reading complexity

- Surprisal (PCFG,  $N$ -gram) is a way to estimate text complexity
- Experienced complexity is reflected in reading speed

Claim:

Current surprisal models inadequately estimate reading complexity

This work:

A simple tweak to fix the surprisal measures

# READING COMPLEXITY IS ESTIMATED BASED ON REGION ENDING

The red apple that the <sup>1</sup>girl<sup>2</sup> ate ...

# READING COMPLEXITY IS ESTIMATED BASED ON REGION ENDING

The red apple that the girl ate ...

$w_1$   $w_2$   $w_3$   $w_4$   $w_5$   $w_6$

Reading model of 'girl':  
sentence position

## READING COMPLEXITY IS ESTIMATED BASED ON REGION ENDING

The red apple that the girl ate ...

4 chars  
 $w_6$

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



The red apple that the girl ate ...

4 chars  
w<sub>6</sub>

Reading model of 'girl':  
sentence position, word length,  $P(\text{girl}|\text{the})$

The red apple that the <sup>1</sup>girl<sup>2</sup> ate ...  
↑  
*important*

Reading model of 'girl':  
sentence position, word length,  $P(\text{girl}|\text{the})$

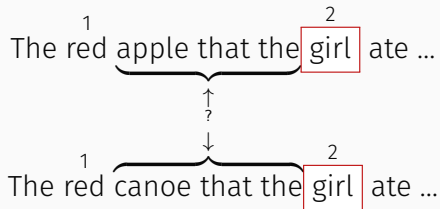
# READING COMPLEXITY IS ESTIMATED BASED ON REGION ENDING

The red <sup>1</sup> apple that the <sup>2</sup> girl ate ...

The diagram shows the sentence "The red apple that the girl ate ...". A bracket is drawn under the words "apple that the", with an arrow pointing down to a question mark. The word "girl" is enclosed in a red rectangular box, and a small number "2" is positioned above it. A small number "1" is positioned above the word "apple".

Reading model of 'girl':  
sentence position, word length,  $P(\text{girl}|\text{the})$

## READING COMPLEXITY IS ESTIMATED BASED ON REGION ENDING



Reading model of 'girl':

sentence position, word length,  $P(\text{girl}|\text{the})$

This study:  $n$ -gram and PCFG surprisal

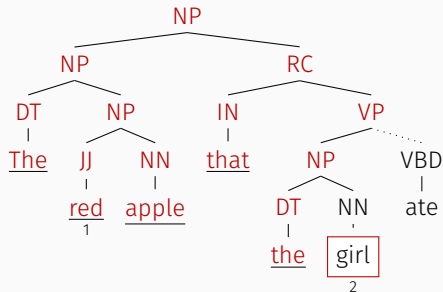
This study:  $n$ -gram and PCFG surprisal

The red apple that the girl ate ...

$$N\text{-gram-surp}(\text{girl}) = -\log P(\text{girl} \mid \text{the})$$

# SURPRISAL: PROBABILITY OF OBSERVATION GIVEN CONTEXT

This study:  $n$ -gram and PCFG surprisal



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

## Cumulative $N$ -gram Surprisal

The red<sup>1</sup> apple that the girl<sup>2</sup> ate ...



Cumulative  $N$ -gram Surprisal

The <sup>1</sup>red apple that the <sup>2</sup>girl ate ...

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

Cumulative  $N$ -gram Surprisal

The red <sup>1</sup> apple that the girl <sup>2</sup> ate ...

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

Cumulative  $N$ -gram Surprisal

The red<sup>1</sup> apple that the girl<sup>2</sup> ate ...

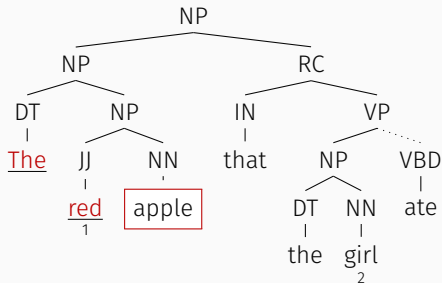
$$\text{cumu-}n\text{-gram}(w, f_{t-1}, f_t) = \sum_{i=f_{t-1}+1}^{f_t} -\log P(w_i \mid w_{i-n} \dots w_{i-1})$$

Cumulative  $N$ -gram Surprisal

The red <sup>1</sup> apple that the <sup>2</sup> girl ate ...

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

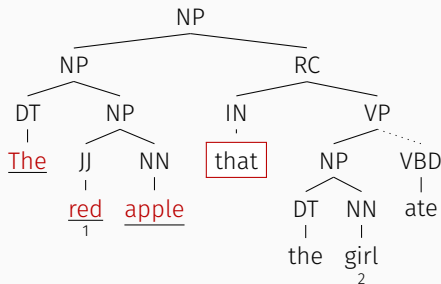
## Cumulative PCFG Surprisal



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

# ACCUMULATED SURPRISAL FIXES THE THEORETICAL PROBLEM

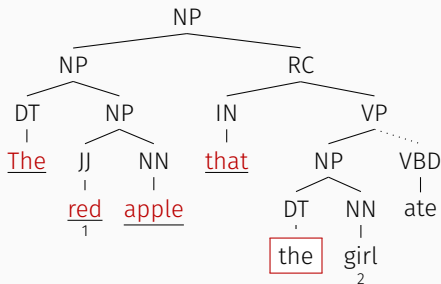
## Cumulative PCFG Surprisal



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

# ACCUMULATED SURPRISEL FIXES THE THEORETICAL PROBLEM

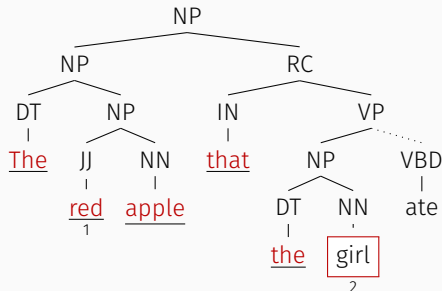
## Cumulative PCFG Surprisal



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

# ACCUMULATED SURPRISAL FIXES THE THEORETICAL PROBLEM

## Cumulative PCFG Surprisal



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



## *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 short sentences from online novels
- frequent comprehension questions

# HOW WELL DOES THIS FIX WORK?

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

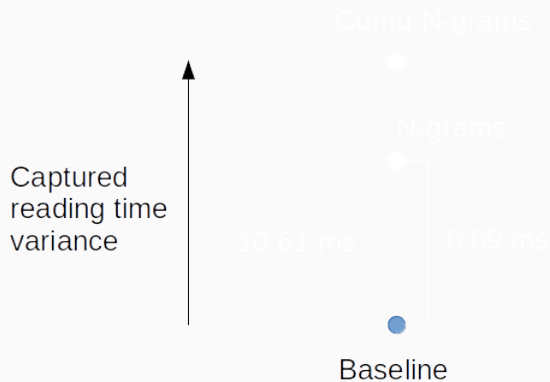
Fixed Factors

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

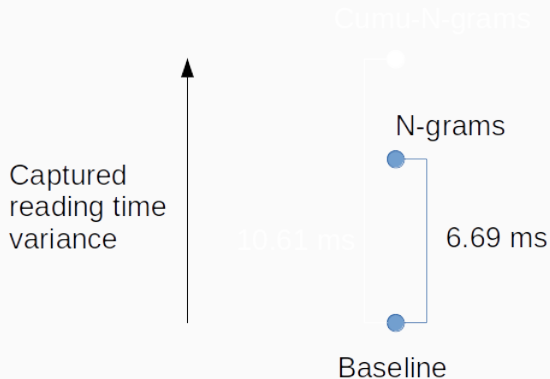
Random Factors

- All fixed factors as by-subject random slopes
- Item, subject and subject $\times$ sentence intercepts

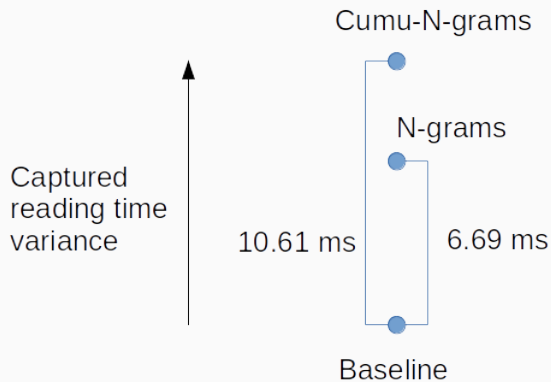
# ACCUMULATION IMPROVES N-GRAM SURPRISAL



# ACCUMULATION IMPROVES N-GRAM SURPRISAL



# ACCUMULATION IMPROVES N-GRAM SURPRISAL





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

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

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

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

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

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

- PCFG surprisal is not useful ( $p > 0.05$ )
- Cumulative PCFG surprisal is not useful ( $p > 0.05$ )
- †Cumulative PCFG is useful with richer grammar ( $p < 0.001$ )

What does accumulation model?

Subsequent regression

The red<sup>1</sup> apple that the girl ate ...

Subsequent regression

The red<sup>1</sup> apple that the girl<sup>2</sup> ate ...

Subsequent regression

The red<sup>1</sup> apple<sup>3</sup> that the girl<sup>2</sup> ate ...



Subsequent regression

<sup>1</sup>   <sup>3</sup>   <sup>4</sup>   <sup>2</sup>  
The red apple that the girl ate ...

Subsequent regression

<sup>1</sup>   <sup>3</sup>   <sup>4</sup>   <sup>2</sup>   <sup>5</sup>  
The red apple that the girl ate ...

Parafoveal processing

<sup>1</sup>  
The red apple that the girl ate ...

Parafoveal processing

Th(e<sup>1</sup> red apple that t)he girl ate ...

Parafoveal processing

Th(e<sup>1</sup> red apple that t)he<sup>2</sup> girl ate ...

Prediction (entropy)

The red<sup>1</sup> apple that the girl ate ...

Prediction (entropy)

The red <sup>1</sup> (apple that the girl) ate ...

Prediction (entropy)

The red<sup>1</sup> (apple that the girl<sup>2</sup>) ate ...



Cumulative surprisal only handles subsequent regression



Cumulative surprisal only handles subsequent regression

Parafoveal: Th(e red <sup>1</sup> apple that t)he <sup>2</sup> girl ate ...

Prediction: The red <sup>1</sup> (apple that the <sup>2</sup> girl) ate ...  
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)

The red <sup>1</sup>apple that the girl <sup>2</sup>ate ...

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

The <sup>1</sup>red apple that the girl <sup>2</sup>ate ...

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



The red <sup>1</sup> apple that the girl <sup>2</sup> ate ...

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

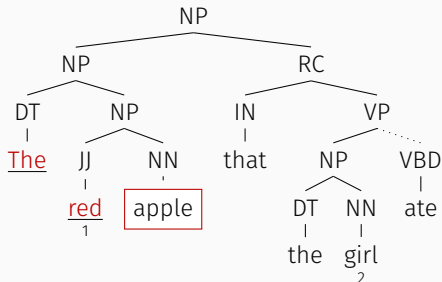
The red <sup>1</sup> apple that the <sup>2</sup> girl ate ...

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

The red <sup>1</sup> apple that the <sup>2</sup> girl ate ...

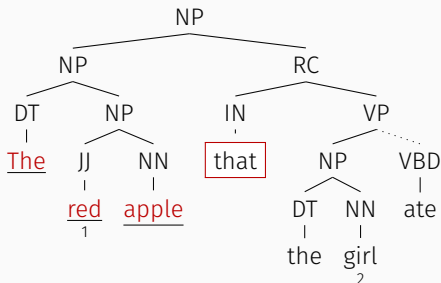
$$\text{future-}n\text{-gram}(w, f_t, f_{t+1}) = \sum_{i=f_t}^{f_{t+1}} -\log P(w_i \mid w_{i-n} \dots w_{i-1})$$

# SUCCESSOR PCFG SURPRISAL



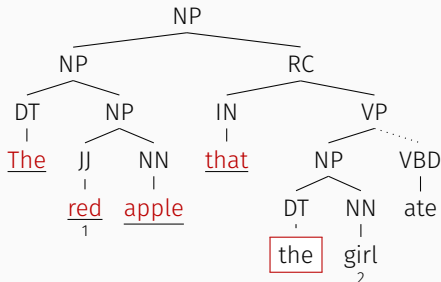
$$\text{Future-PCFG}(w, f_t, f_{t+1}) = \sum_{i=f_t}^{f_{t+1}} -\log P(T_i = w_i \mid T_1 \dots T_{i-1} = w_1 \dots w_{i-1})$$

# SUCCESSOR PCFG SURPRISAL



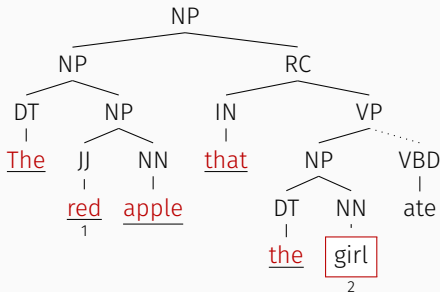
$$\text{Future-PCFG}(w, f_t, f_{t+1}) = \sum_{i=f_t}^{f_{t+1}} -\log P(T_i = w_i \mid T_1 \dots T_{i-1} = w_1 \dots w_{i-1})$$

# SUCCESSOR PCFG SURPRISAL

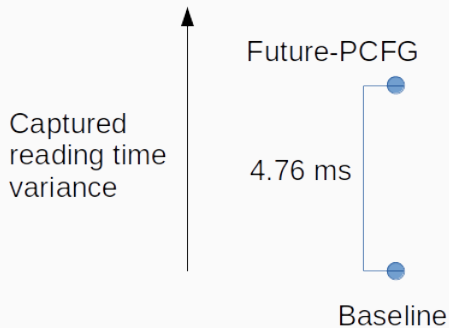


$$\text{Future-PCFG}(w, f_t, f_{t+1}) = \sum_{i=f_t}^{f_{t+1}} -\log P(T_i = w_i \mid T_1 \dots T_{i-1} = w_1 \dots w_{i-1})$$

# SUCCESSOR PCFG SURPRISAL

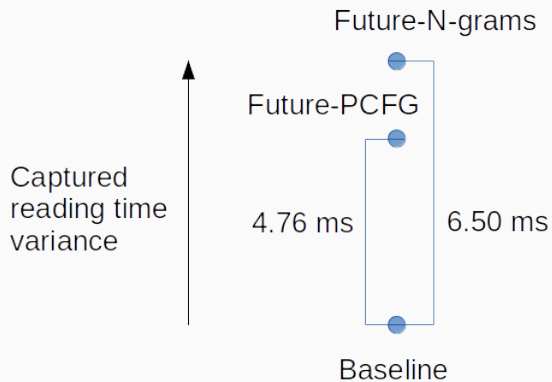


$$\text{Future-PCFG}(w, f_t, f_{t+1}) = \sum_{i=f_t}^{f_{t+1}} -\log P(T_i = w_i \mid T_1 \dots T_{i-1} = w_1 \dots w_{i-1})$$

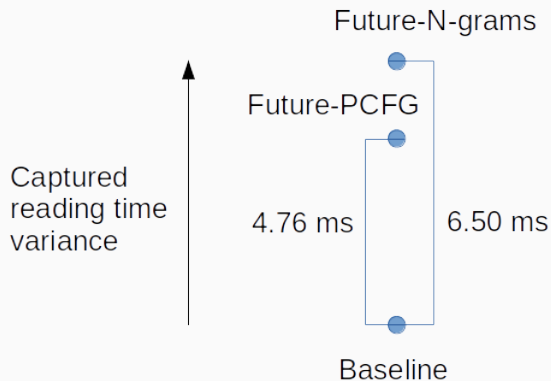




# SUCCESSOR N-GRAMS WORK BETTER



# SUCCESSOR N-GRAMS WORK BETTER



PCFG surprisal may require a richer grammar

# SUCCESSOR *N*-GRAMS HAVE LIMITED INFLUENCE

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

## CONCLUSION: ACCUMULATE SURPRISAL!

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

## CONCLUSION: ACCUMULATE SURPRISAL!

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

## CONCLUSION: ACCUMULATE SURPRISAL!

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



## CONCLUSION: ACCUMULATE SURPRISAL!

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

## CONCLUSION: ACCUMULATE SURPRISAL!

- $N$ -gram surprisal should be accumulated to predict reading times
- $N$ -gram surprisal accumulates pre- and post-saccade
  - Pre-saccade  $n$ -grams are limited
- PTB PCFG surprisal does not accumulate
- †Richer grammars may accumulate better

Thanks to:

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

## UCL EFFECT SIZE REFERENCE

| Model                            | Effect Size (ms) |
|----------------------------------|------------------|
| Future <i>N</i> -grams           | 6.5*             |
| <i>N</i> -grams                  | 6.69             |
| Cumulative GCG-PCFG <sup>†</sup> | 8.25*            |
| Cumulative <i>N</i> -grams       | 10.61*           |

\* $p < 0.001$

*N*-gram model has the given effect size before adding cumu-*n*-grams.

# CUMU-*N*-GRAM RESULTS

| Model           | N-gram vs Cumu- <i>N</i> -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, prefix, 5-gram, cumu-5-gram

Base fixed: sentpos, wlen, rlen, prefix

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

## FUTURE SURPRISAL RESULTS

| 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, prefix, cumu-5-gram,  
future-5-grams, future-PCFG

Base fixed: sentpos, wlen, rlen, prefix, cumu-5-gram

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