# Can entropy explain successor surprisal effects in reading?

Marten van Schijndel and Tal Linzen

vansky@jhu.edu

Department of Cognitive Science, Johns Hopkins University

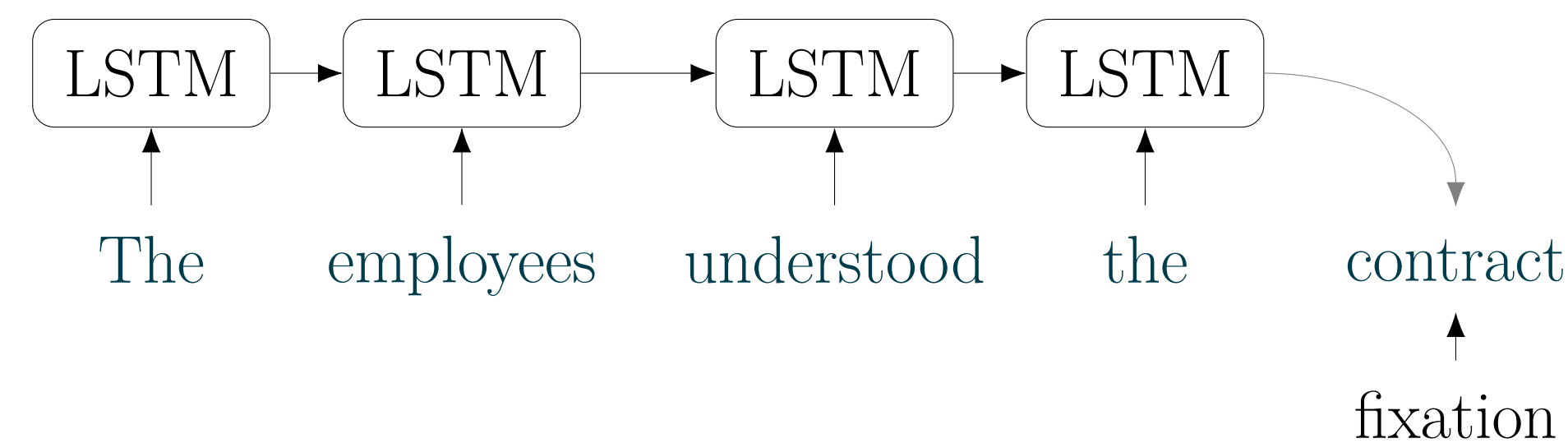**JOHNS HOPKINS**
U N I V E R S I T Y

## Abstract

It is well-established that reading times are influenced by word probabilities [7], but strangely this holds true even for words that have *not been viewed yet* and which are *not visible* to the reader [1, 8]. Angele et al. hypothesize that this effect may be driven by entropy, but previous studies have relied on multiple separate models to compute the relevant measures, which muddies interpretation of their results. We test their hypothesis using a single neural language model to estimate the relevant computational measures.

## Model

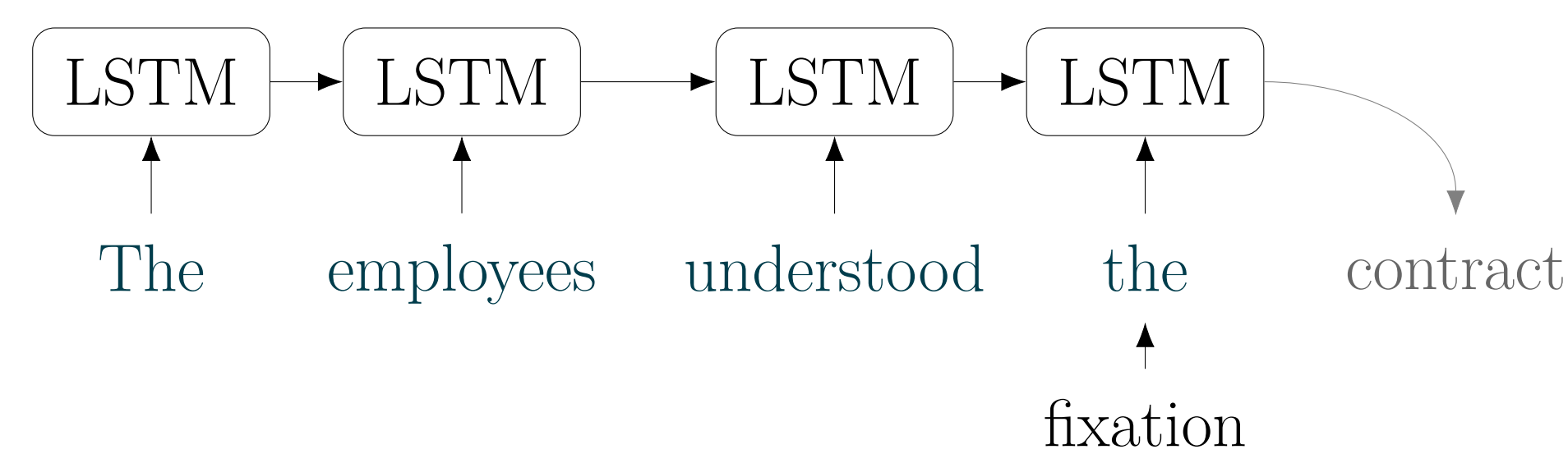LSTM LM trained on 90M words of English Wikipedia [3]

## Measures

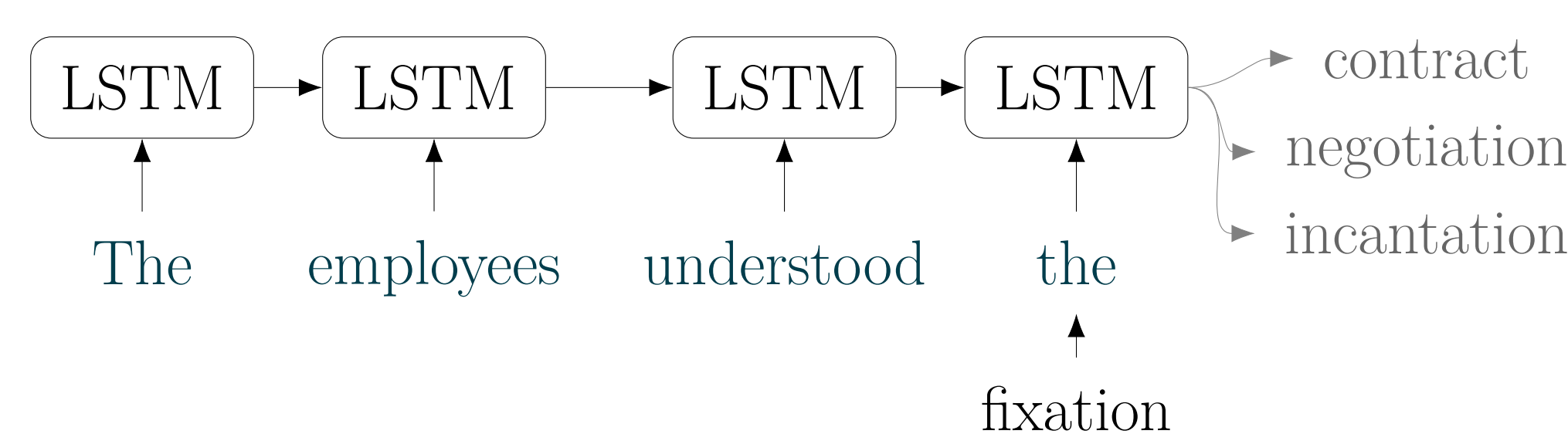**Surprisal** [6, 4] estimates the amount of new information:

| LSTM | → | LSTM | → | LSTM | → | LSTM |

The     employees     understood     the     contract

fixation

$$\text{surprisal}(w_t) = -\log \text{P}(w_t \mid w_{1\ldots t-1}) \qquad (1)$$

**Successor surprisal** [5] estimates upcoming information:

| LSTM | → | LSTM | → | LSTM | → | LSTM |

The     employees     understood     the     contract

fixation

$$\text{successor surprisal}(w_t) = -\log \text{P}(w_{t+1} \mid w_{1\ldots t}) \qquad (2)$$
$$= \text{surprisal}(w_{t+1}) \qquad (3)$$

**Entropy** [6] estimates the amount of uncertainty:

| LSTM | → | LSTM | → | LSTM | → | LSTM |

contract
negotiation
incantation

The     employees     understood     the

fixation

$$H(w_t) = -\sum_{w_{t+1} \in V} \text{P}(w_{t+1} \mid w_{1\ldots t}) \log \text{P}(w_{t+1} \mid w_{1\ldots t}) \qquad (4)$$
$$= E[\text{surprisal}(w_{t+1})] \qquad (5)$$
$$= E[\text{successor surprisal}(w_t)] \qquad (6)$$

## Data

Natural Stories Corpus [2]

- 10 texts (485 sentences)
- Self-paced reading times
- 181 participants
- We omit multi-token words (e.g., *boar·!·'*)
- We partition the sentences:
  1/3 exploration : 2/3 confirmation

## Successor Surprisal as Entropy Estimator

In practice, with a finite set of observations $T$ which are regressed simultaneously, successor surprisal should provide a Monte Carlo estimator of entropy in that corpus:

$$\hat{H}(T) \approx -\sum_{t=1}^{|T|} \frac{1}{|T|} \log \text{P}(w_{t+1} \mid w_{1\ldots t}) \qquad (7)$$
$$= \sum_{t=1}^{|T|} \frac{1}{|T|} \text{successor surprisal}(w_{t+1}) \qquad (8)$$
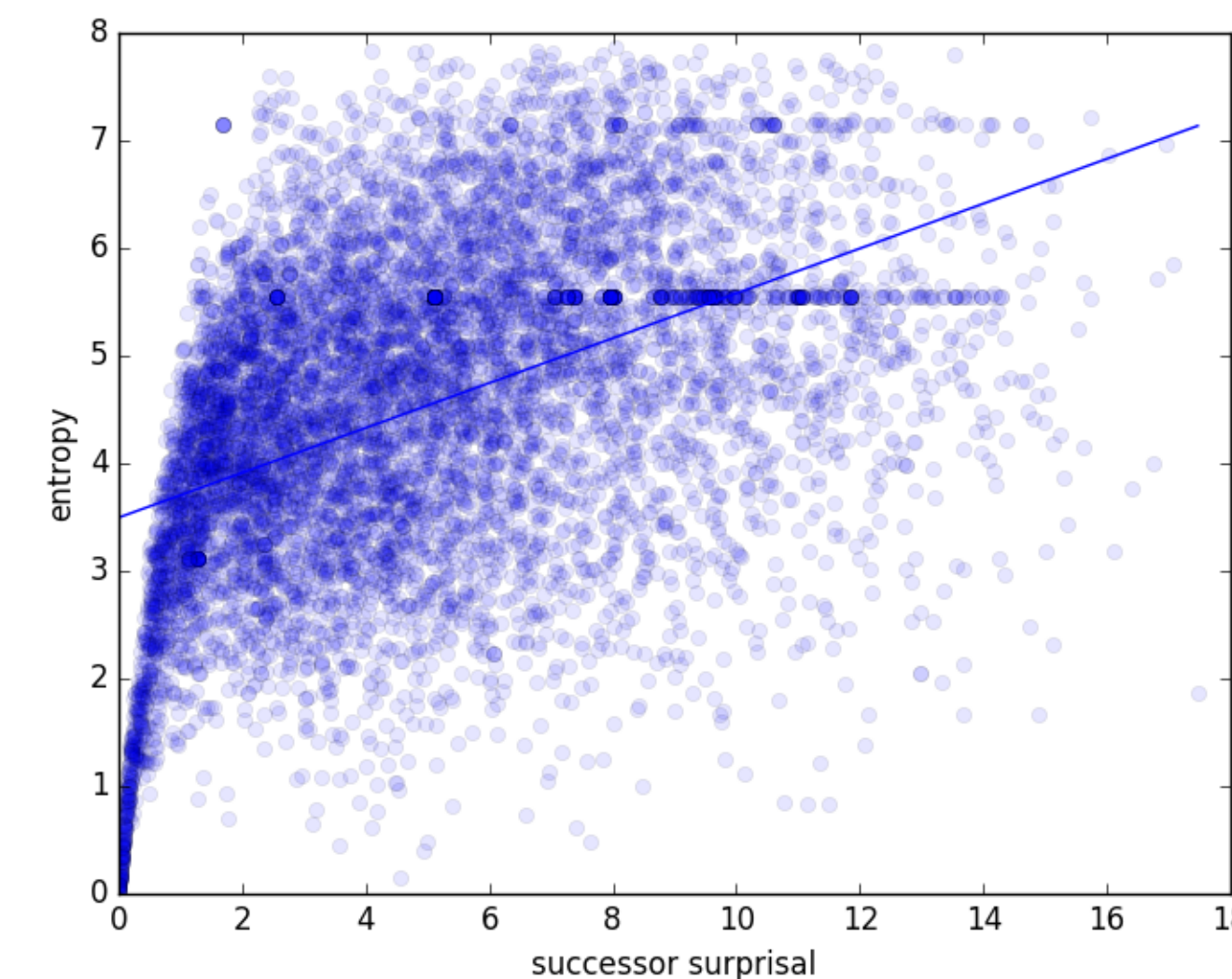


**Figure 1:** Successor surprisal plotted against entropy for each word in the Natural Stories Corpus.

The Pearson correlation is 0.45, providing empirical validation that the limit-case relation between the measures applies even in a relatively small corpus setting.

## Reading Time Predictions

|  | $\hat{\beta}$ | $\hat{\sigma}$ | $t$ |
|---|---|---|---|
| (Intercept) | 332.50 | 6.30 | 52.76 |
| Sentence position | 0.82 | 0.52 | 1.59 |
| Word length | 5.13 | 1.01 | 5.10 |
| Surprisal | 5.77 | 0.57 | 10.06 |
| Successor surprisal | 3.40 | 0.40 | 8.53 |
| Entropy | 3.21 | 0.55 | 5.81 |

**Table 1:** Fixed effect coefficients from fitting self-paced reading times. Since predictors were z-transformed, the $\hat{\beta}$ coefficients indicate the change in ms per standard deviation of each predictor.

## Change Number of Possible Continuations

By changing the number of possible continuations considered by the model, we can probe the rough number of continuations readers are sensitive to. Further, if people consider a small number of continuations, that could account for successor surprisal's continued influence.

$$H(w_t) = -\sum_{w_{t+1} \in V} \text{P}(w_{t+1} \mid w_{1\ldots t}) \log \text{P}(w_{t+1} \mid w_{1\ldots t}) \qquad (9)$$
$$\approx -\sum_{w_{t+1} \in K} \text{P}(w_{t+1} \mid w_{1\ldots t}) \log \text{P}(w_{t+1} \mid w_{1\ldots t}) \qquad (10)$$

| $K$ | Successor surprisal | Total entropy |
|---|---|---|
| 5 | 0.212 | 0.541 |
| 50 | 0.335 | 0.820 |
| 500 | 0.397 | 0.947 |
| 5000 | 0.434 | 0.992 |
| 50000 | 0.454 | 1 |

**Figure 2:** Correlation between successor surprisal and entropy when entropy is computed over the most probable $K$ continuations.

| $K$ | $\hat{\beta}_H$ | $\hat{\sigma}_H$ | $\hat{\beta}_s$ | $\hat{\sigma}_s$ |
|---|---|---|---|---|
| 5 | 3.19 | 0.69 | 3.96 | 0.53 |
| 50 | 3.43 | 0.70 | 3.85 | 0.54 |
| 500 | 4.11 | 0.69 | 3.66 | 0.54 |
| 5000 | 4.67 | 0.70 | 3.52 | 0.54 |
| 50000 | 4.87 | 0.70 | 3.47 | 0.54 |

**Figure 3:** Entropy ($H$) and successor surprisal ($s$) coefficients in a regression model for the exploratory data partition, when $H$ is calculated over the $K$ most probable continuations.

## Conclusion

- Findings support Angele et al. hypothesis that uncertainty drives the successor surprisal effect in reading times.
- Entropy is unable to account for full successor effect; some other driver likely present.
- Readers are sensitive to a large number of possible continuations.

Mixed model formula:
RT ~ word_length + sentence_position + surprisal + successor_surprisal + entropy + (1 | item) +
(0 + word_length + sentence_position + surprisal + successor_surprisal + entropy | subject)

## References

[1] Bernhard Angele, Elizabeth R. Schotter, Timothy J. Slattery, Tara L. Tenenbaum, Klinton Bicknell, and Keith Rayner.
Do successor effects in reading reflect lexical parafoveal processing? evidence from corpus-based and experimental eye movement data.
*Journal of Memory and Language*, 79–80:76–96, 2015.

[2] Richard Futrell, Edward Gibson, Hal Tily, Anastasia Vishnevetsky, Steve Piantadosi, and Evelina Fedorenko.
The natural stories corpus.
In *Language Resources and Evaluation Conference*, pages 76–82, 2018.

[3] Kristina Gulordava, Piotr Bojanowski, Edouard Grave, Tal Linzen, and Marco Baroni.
Colorless green recurrent networks dream hierarchically.
In *Proceedings of the North American Chapter of the Association for Computational Linguistics*, 2018.

[4] John Hale.
A probabilistic early parser as a psycholinguistic model.
In *Proceedings of NAACL*, 2001.

[5] R. Kliegl, A. Nuthmann, and R. Engbert.
Tracking the mind during reading: the influence of past, present, and future words on fixation durations.
*Journal of Experimental Psychology: General*, 135:12–35, 2006.

[6] Claude Shannon.
A mathematical theory of communication.
*Bell System Technical Journal*, 27:379–423, 623–656, 1948.

[7] Nathaniel J. Smith and Roger Levy.
The effect of word predictability on reading time is logarithmic.
*Cognition*, 128:302–319, 2013.

[8] Marten van Schijndel and William Schuler.
Addressing surprisal deficiencies in reading time models.
In *Proceedings of the Computational Linguistics for Linguistic Complexity Workshop*. Association for Computational Linguistics, 2016.