

Entropy- and Distance-Based Predictors From GPT-2 Attention Patterns

Predict Reading Times Over and Above GPT-2 Surprisal

PAPER



oh.531@osu.edu
github.com/byungdoh/attn_dist

Byung-Doh Oh William Schuler

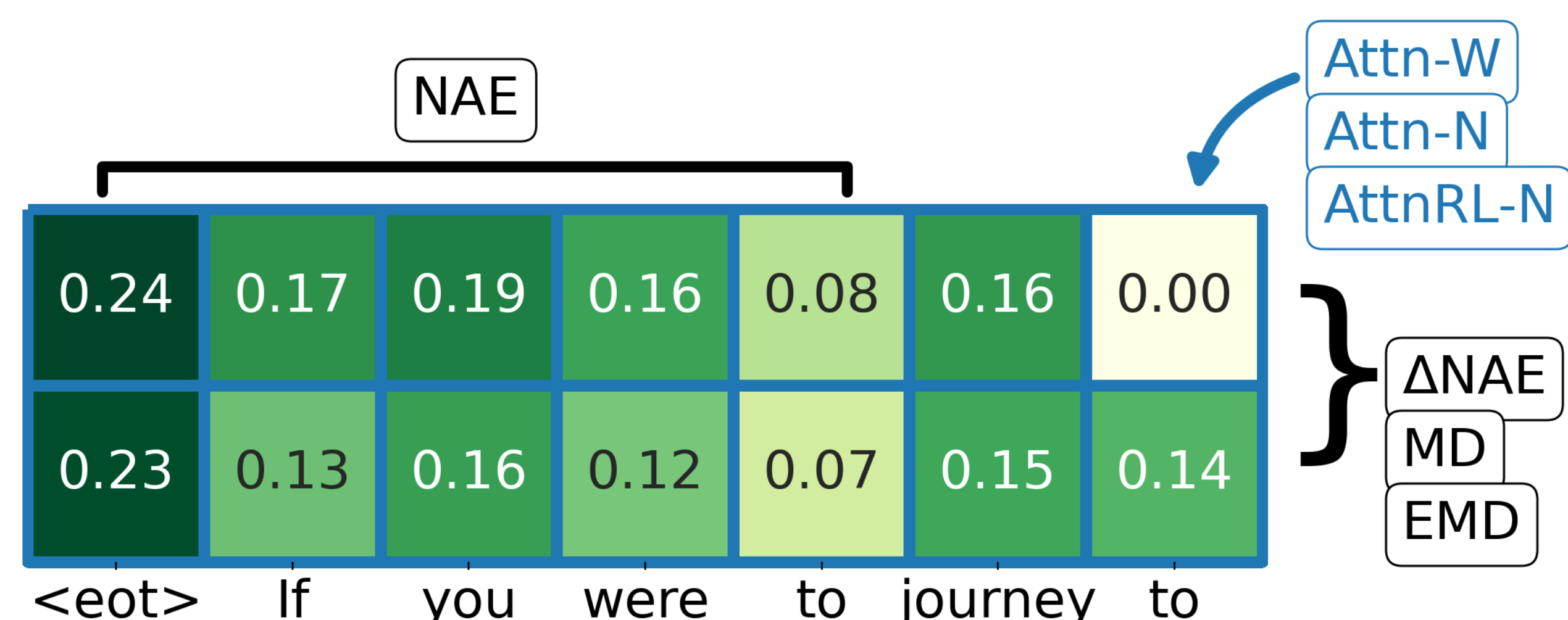
The Ohio State University

Introduction

- There are empirical shortcomings of LM surprisal as expectation-based predictors of comprehension difficulty, such as underprediction of garden-path effects [9]
- As such, there are recent efforts to identify memory-based effects from LM representations
- For example, a connection has been made between Transformer self-attention weights and cue-based retrieval [7], but their entropy was not predictive over surprisal [8]
- Self-attention weights proper do not accurately reflect the importance of each token in context [2, 4]

Entropy- and Distance-Based Predictors

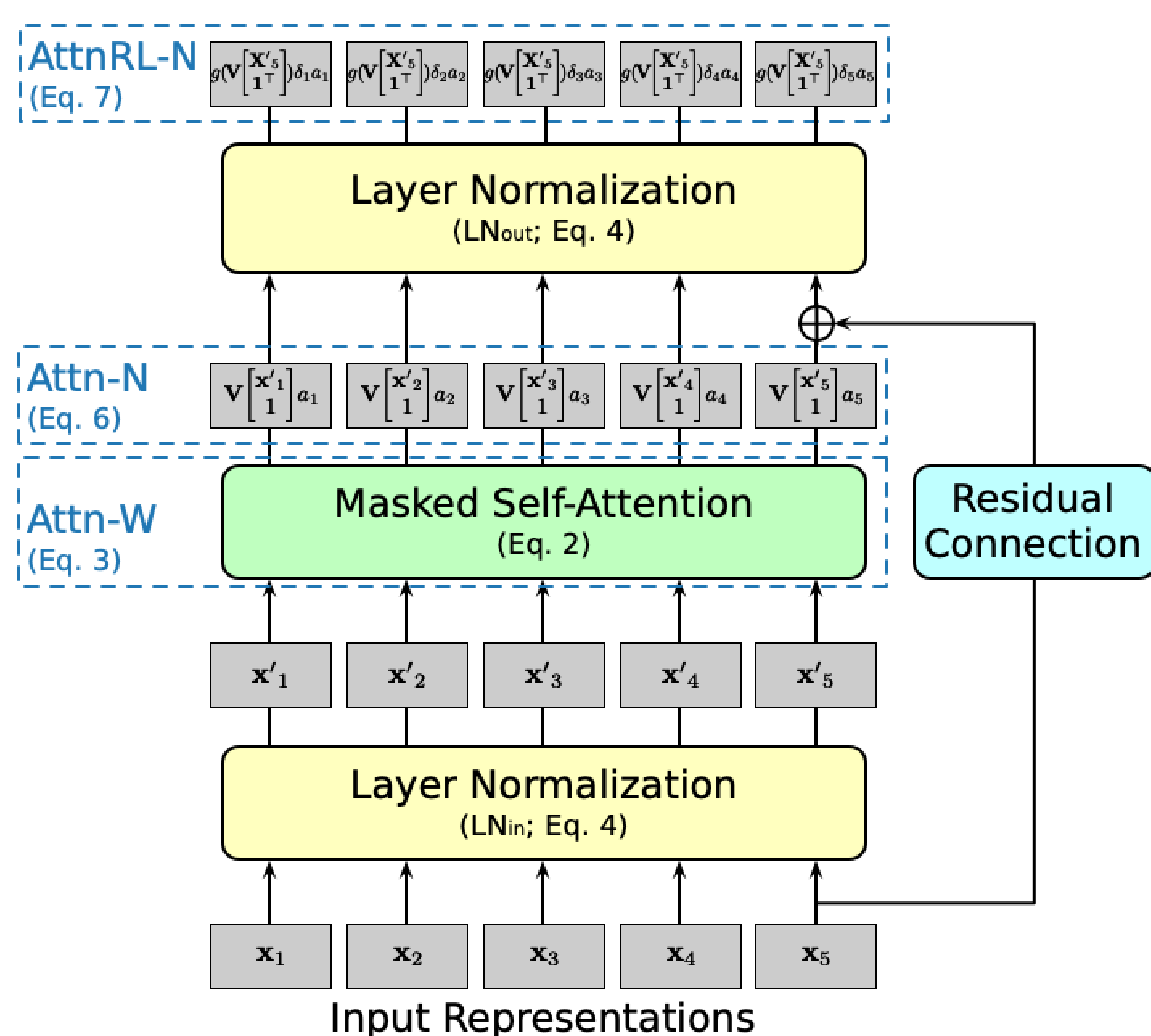
This work defines **entropy- and distance-based predictors** of comprehension difficulty under **different formulations of attention patterns**:



- 1 Normalized attention entropy (NAE): Entropy of normalized weights over $w_{1..i-1}$ divided by maximum entropy
- 2 Δ Normalized attention entropy (Δ NAE): Absolute value of change in NAE across consecutive timesteps
- 3 Manhattan distance (MD): 1-norm of difference in attention weight vectors across consecutive timesteps
- 4 Earth Mover's Distance (EMD): Minimum amount of 'work' necessary to transform the current attention weight vector to the next

Formulations of GPT-2 [6] Attention Patterns

- Linear nature of the computations in a self-attention block allows the aggregation of representations to be deferred [4, 5]
- Vector norms are normalized to yield weights (ATTN-N, ATTNRL-N) that are comparable to self-attention weights (ATTN-W)



Evaluation on Human Reading Times

- Evaluation on the Natural Stories Corpus [1] and the Dundee Corpus [3]
- Baseline: low-level predictors, unigram surprisal, and GPT-2 surprisal
- Predictors of interest calculated from topmost attention heads of GPT-2

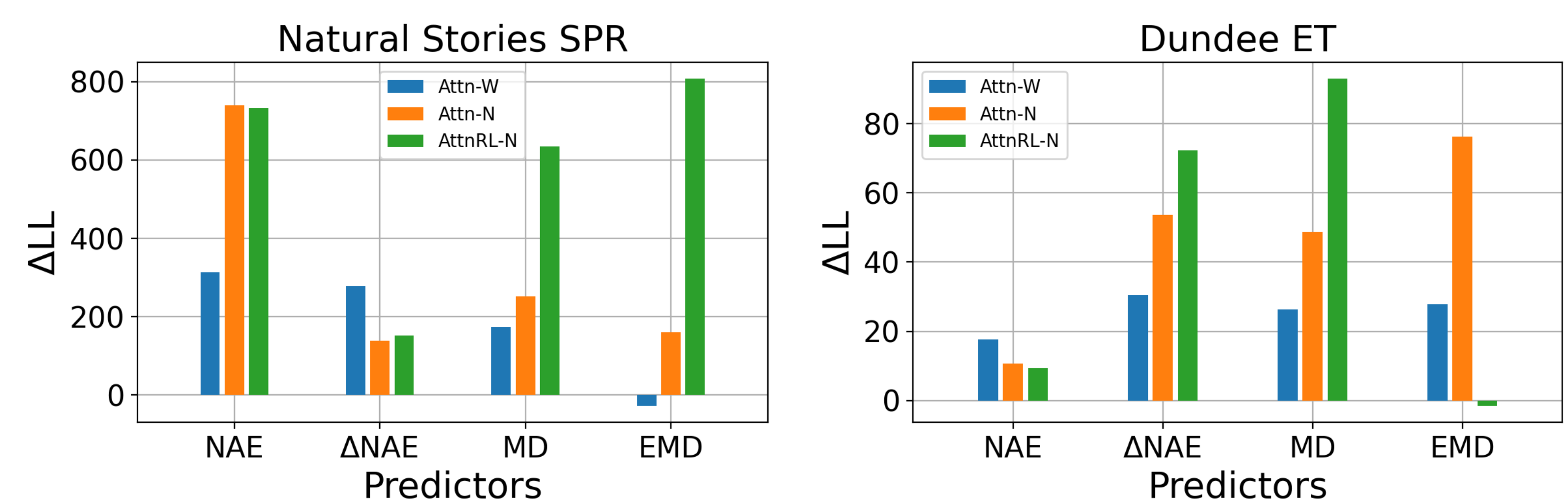


Figure 1: Improvements in regression model log-likelihood from including each predictor on the exploratory (dev) partition.

| Corpus | Predictor | Effect Size (p -value) |
|-----------------|-------------|---------------------------|
| Natural Stories | ATTN-N+NAE | 6.87 ms ($p < 0.001$) |
| | GPT2SURP | 2.56 ms |
| | ATTNRL-N+MD | 6.59 ms ($p < 0.001$) |
| Dundee | GPT2SURP | 2.82 ms |
| | ATTN-N+NAE | N/A (n.s.) |
| | ATTNRL-N+MD | 1.05 ms ($p < 0.001$) |
| | GPT2SURP | 3.81 ms |

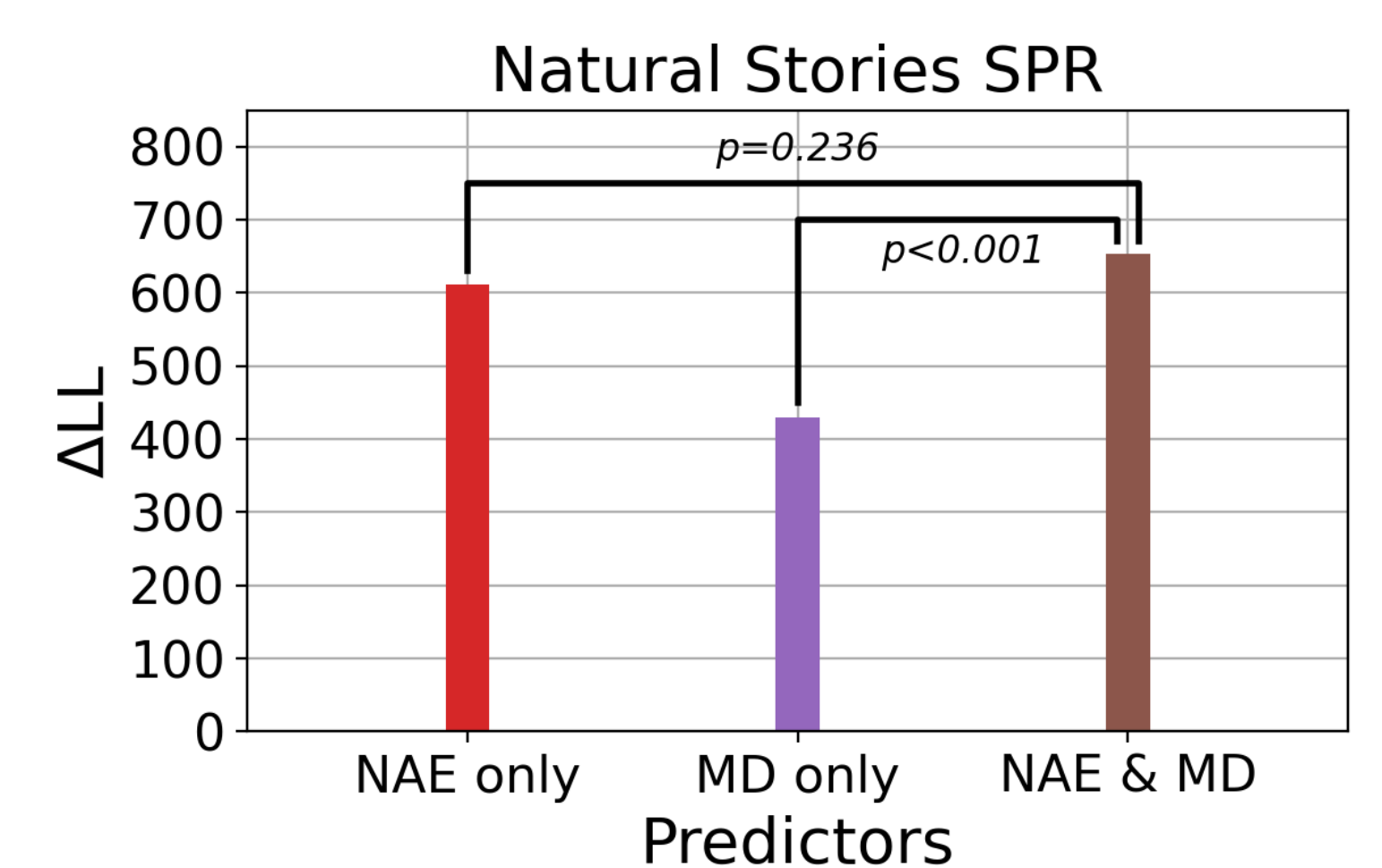


Table 1: Effect sizes per standard deviation on the held-out (test) partition.

Figure 2: Improvements in log-likelihood on the held-out (test) partition.

| Predictor | Natural Stories | | | | | | | | | | | | | |
|--------------|-----------------|--------------|------|------|------|--------------|------|------|------|--------------|------|------|-------|---------|
| | NAE | Δ NAE | MD | EMD | NAE | Δ NAE | MD | EMD | NAE | Δ NAE | MD | EMD | GPT-2 | Unigram |
| Attn-W | | 0.50 | 0.76 | 0.53 | 0.92 | 0.47 | 0.75 | 0.45 | 0.92 | 0.48 | 0.77 | 0.48 | 0.22 | 0.46 |
| Δ NAE | 0.55 | | 0.84 | 0.95 | 0.67 | 0.88 | 0.86 | 0.92 | 0.67 | 0.88 | 0.86 | 0.75 | 0.29 | 0.52 |
| MD | 0.78 | 0.87 | | 0.88 | 0.79 | 0.74 | 0.97 | 0.76 | 0.79 | 0.74 | 0.90 | 0.71 | 0.22 | 0.55 |
| EMD | 0.58 | 0.94 | 0.91 | | 0.67 | 0.81 | 0.87 | 0.91 | 0.67 | 0.82 | 0.85 | 0.74 | 0.28 | 0.53 |
| Attn-N | | 0.92 | 0.71 | 0.81 | 0.70 | | 0.63 | 0.83 | 0.63 | 1.00 | 0.64 | 0.90 | 0.61 | 0.33 |
| Δ NAE | 0.48 | 0.87 | 0.71 | 0.79 | 0.66 | | 0.82 | 0.93 | 0.63 | 1.00 | 0.84 | 0.74 | 0.27 | 0.48 |
| MD | 0.76 | 0.89 | 0.96 | 0.89 | 0.86 | 0.82 | | 0.83 | 0.83 | 0.83 | 0.96 | 0.75 | 0.26 | 0.58 |
| EMD | 0.48 | 0.91 | 0.76 | 0.89 | 0.66 | 0.93 | 0.85 | | 0.63 | 0.93 | 0.85 | 0.80 | 0.30 | 0.50 |
| AttnRL-N | | 0.92 | 0.71 | 0.81 | 0.70 | 1.00 | 0.66 | 0.86 | 0.67 | | 0.64 | 0.90 | 0.61 | 0.34 |
| Δ NAE | 0.49 | 0.88 | 0.72 | 0.79 | 0.66 | 1.00 | 0.83 | 0.93 | 0.67 | | 0.85 | 0.74 | 0.27 | 0.48 |
| MD | 0.76 | 0.87 | 0.88 | 0.85 | 0.91 | 0.86 | 0.96 | 0.87 | 0.91 | 0.86 | | 0.79 | 0.32 | 0.60 |
| EMD | 0.50 | 0.74 | 0.69 | 0.72 | 0.62 | 0.73 | 0.75 | 0.78 | 0.62 | 0.73 | 0.77 | | 0.26 | 0.47 |
| Surprisal | | 0.30 | 0.42 | 0.36 | 0.40 | 0.43 | 0.40 | 0.40 | 0.42 | 0.44 | 0.40 | 0.44 | 0.36 | |
| Uni. | 0.42 | 0.52 | 0.52 | 0.51 | 0.52 | 0.49 | 0.56 | 0.50 | 0.52 | 0.49 | 0.57 | 0.45 | 0.60 | |
| Dundee | | 0.50 | 0.76 | 0.53 | 0.92 | 0.47 | 0.75 | 0.45 | 0.92 | 0.48 | 0.77 | 0.48 | 0.22 | 0.46 |

Figure 3: Pearson correlation coefficients between predictors.

Conclusion

Results show robust effects of Transformer attention-based predictors in predicting reading times of broad-coverage naturalistic data

References

- [1] R. Futrell, E. Gibson, H. J. Tily, I. Blank, A. Vishnevetsky, S. Piantadosi, and E. Fedorenko. The Natural Stories corpus: A reading-time corpus of English texts containing rare syntactic constructions. *Language Resources and Evaluation*, 2021.
- [2] S. Jain and B. C. Wallace. Attention is not explanation. In *Proc. NAACL*, 2019.
- [3] A. Kennedy, R. Hill, and J. Pynte. The Dundee Corpus. In *Proc. ECEM*, 2003.
- [4] G. Kobayashi, T. Kuribayashi, S. Yokoi, and K. Inui. Attention is not only a weight: Analyzing transformers with vector norms. In *Proc. EMNLP*, 2020.
- [5] G. Kobayashi, T. Kuribayashi, S. Yokoi, and K. Inui. Incorporating residual and normalization layers into analysis of masked language models. In *Proc. EMNLP*, 2021.
- [6] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever. Language models are unsupervised multitask learners. *OpenAI Technical Report*, 2019.
- [7] S. H. Ryu and R. L. Lewis. Accounting for agreement phenomena in sentence comprehension with transformer language models: Effects of similarity-based interference on surprisal and attention. In *Proc. CMCL*, 2021.
- [8] S. H. Ryu and R. L. Lewis. Using Transformer language model to integrate surprisal, entropy, and working memory retrieval accounts of sentence processing. In *35th Annual Conference on Human Sentence Processing*, 2022.
- [9] M. van Schijndel and T. Linzen. Single-stage prediction models do not explain the magnitude of syntactic disambiguation difficulty. *Cognitive Science*, 2021.