

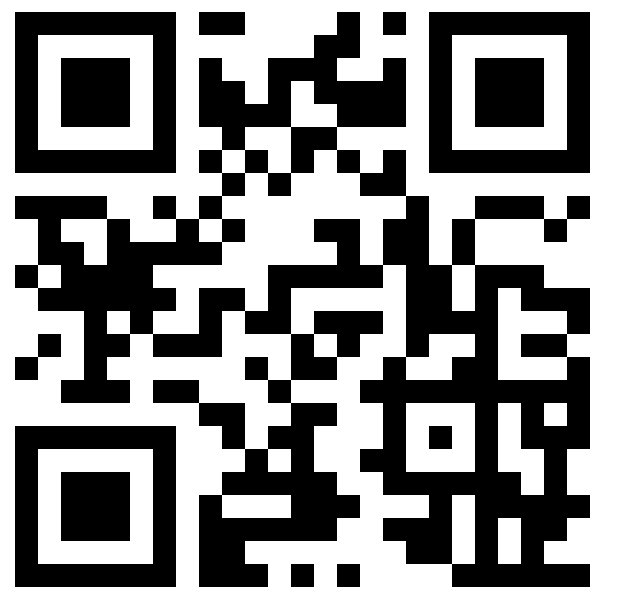
Introduction

- ▶ Benchmark data are an important tool for developing theories and evaluating model predictions.
- ▶ The majority of benchmark data in sentence processing are limited to English (e.g., [3, 8]).

Our Work (in Progress)

- ▶ We collect eye-tracking benchmark data for a battery of postulated effects in German (10 phenomena).
- ▶ In parallel, we also collect self-paced reading (SPR) data on the same materials.
- ▶ So far (July 2025):
 - ▶ 155 in-lab participants have been tested with eye tracking. 1 was excluded due to low accuracy on comprehension questions.
 - ▶ 736 Prolific participants have been tested with SPR. 47 were excluded due to low accuracy on comprehension questions.
- ▶ We show the results so far, next to predictions based on qualitative theories, GPT-2 surprisal [4, 11, 16, 18], and lossy-context GPT-2 surprisal [1, 2, 6, 16].
- ▶ Model comparisons using Pareto-smoothed importance sampling [20] assess out-of-sample predictiveness.

Pre-Registration Protocol (SPR)



osf.io/wpra9

Predictions, Results, and Model Comparison

Experimental Designs

GPSD (2×2): Garden Paths From Subject-vs.-Direct-Object Ambiguity
Ambiguous/Unambiguous × S-O/O-S — closely replicating [13]

GPSI (2×2): Garden Paths From Subject-vs.-Indirect-Object Ambiguity
Ambiguous/Unambiguous × Active/Passive — loosely replicating [14]

GPCA (2×2): Garden Paths From Coordination Ambiguity
NP-/VP-Coordination × AP-/PP-Modifier — closely replicating [10]

GPMI (2×2): Garden Paths From Modifier-vs.-Indirect-Object Ambiguity
Modifier/No-Modifier × Ambiguous/Unambiguous — closely replicating [9]

AGAT (2×2): Agreement Attraction in Grammatical Sentences
Singular-/Plural-Controller × Match/Mismatch — closely replicating [5]

LOCO (2×2): Local Coherence
Coherent/Incoherent × Intervener/No-Intervener — closely replicating [15]

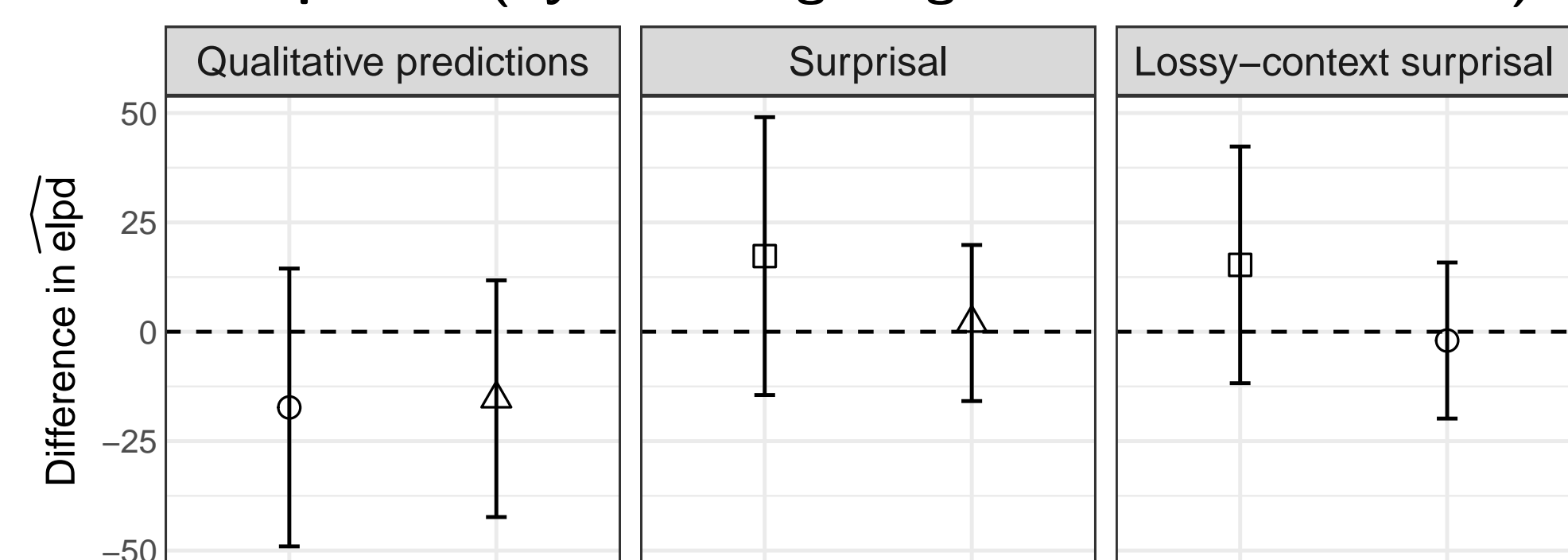
SBIN (2×2): Similarity-Based Interference
Subject-Cue [Yes/No] × Animacy-Cue [Yes/No] — closely replicating [17]

RCSO (2×2): Subject vs. Object Relative Clauses
Subject/Object × Double-/Single-Embedding — German adaptation of [7]

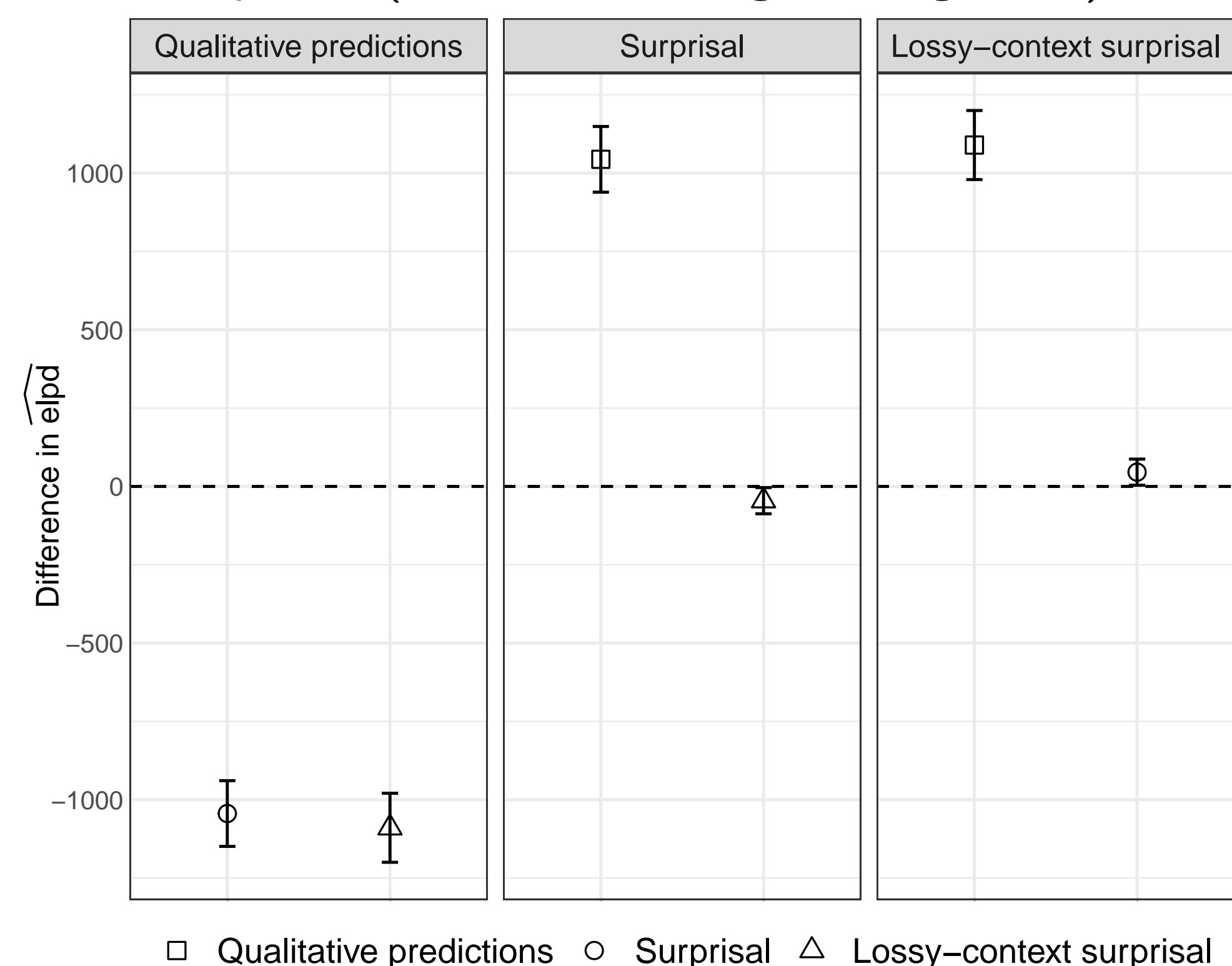
SYAA (3×1): Syntax-Based Attachment Ambiguity
High-/Low-/Ambiguous-Attachment — closely replicating [12]

SEAA (3×1): Semantics-Based Attachment Ambiguity
High-/Low-/Ambiguous-Attachment — German adaptation of [19]

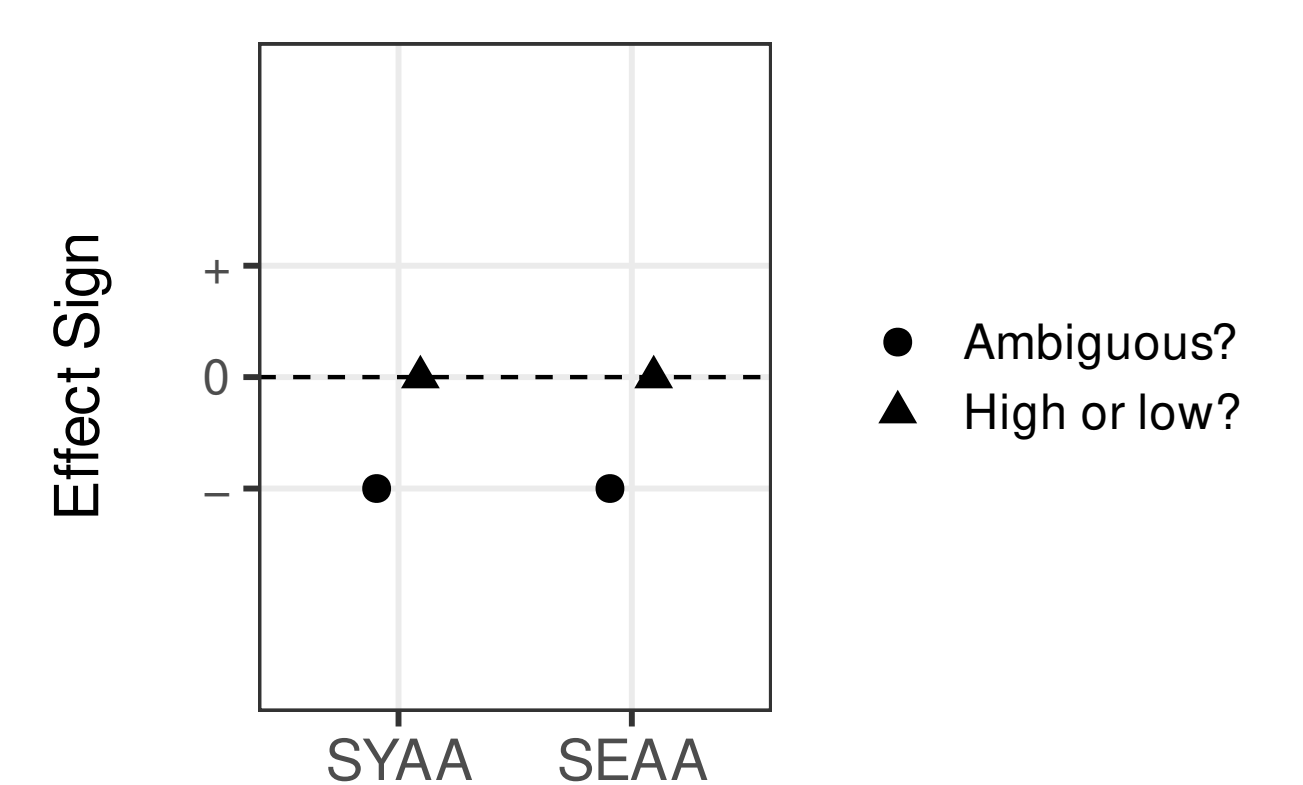
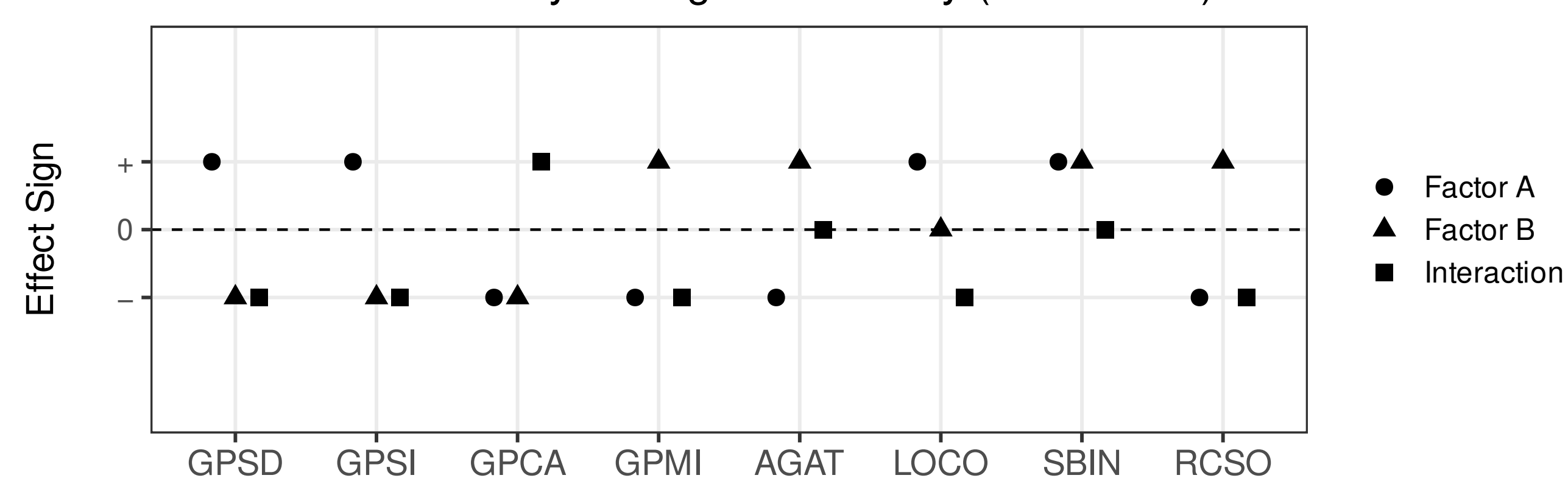
Model Comparison (Eye Tracking: Regression Path Durations)



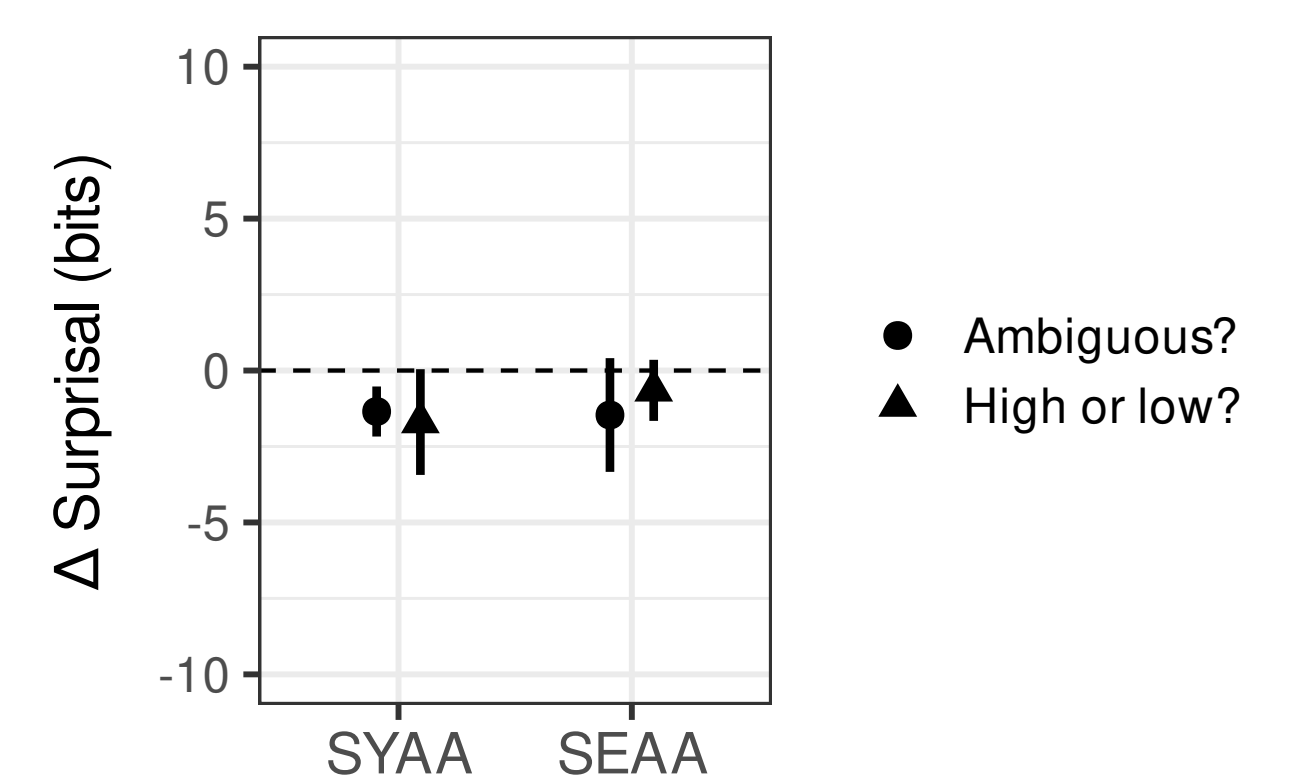
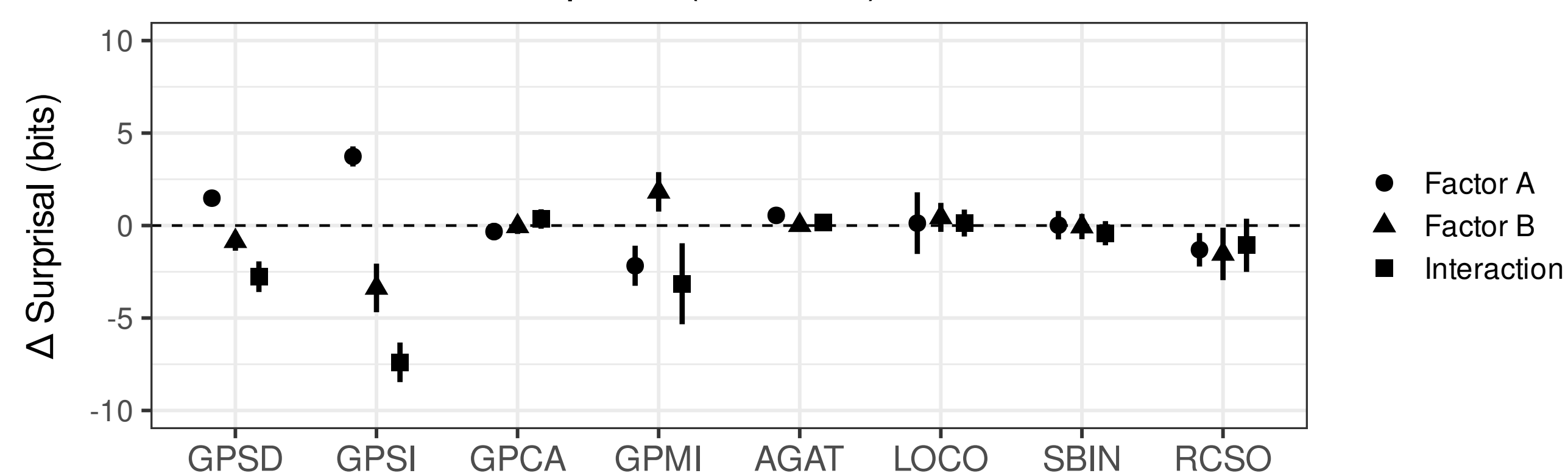
Model Comparison (Self-Paced Reading: Reading Times)



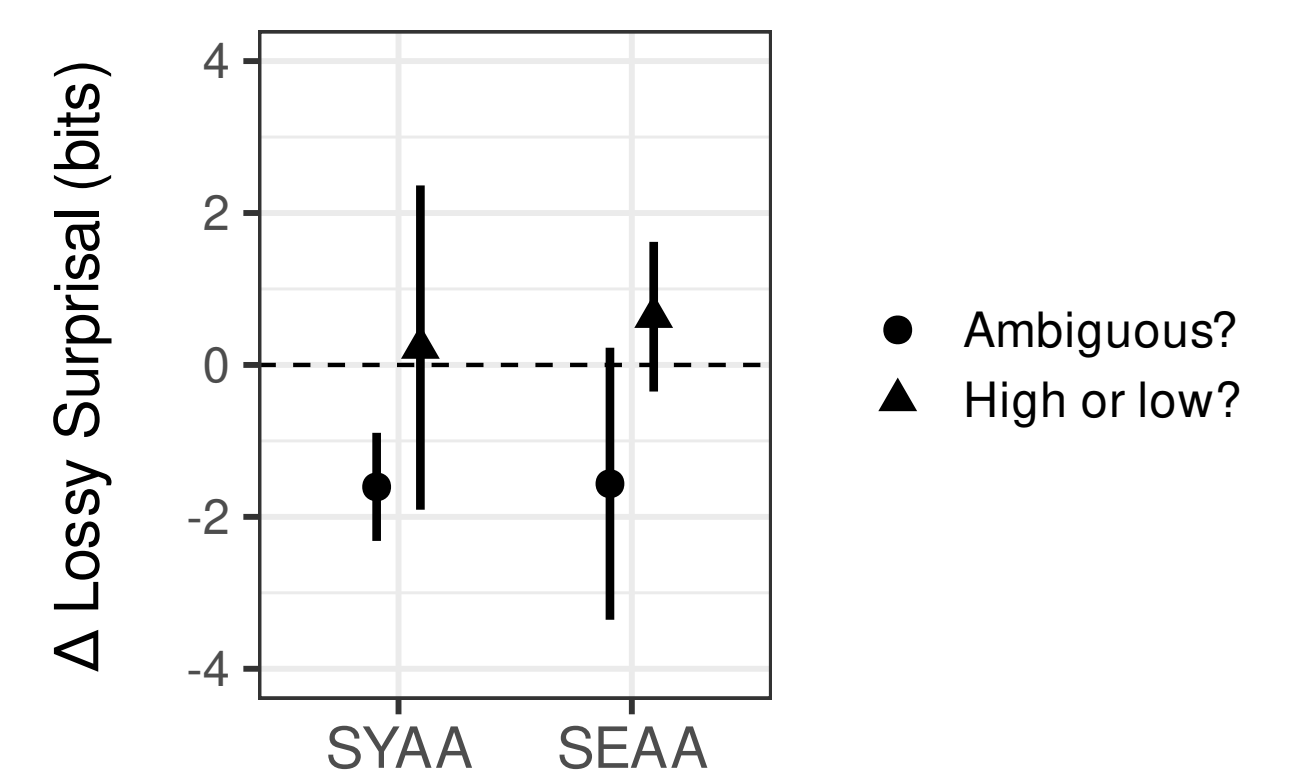
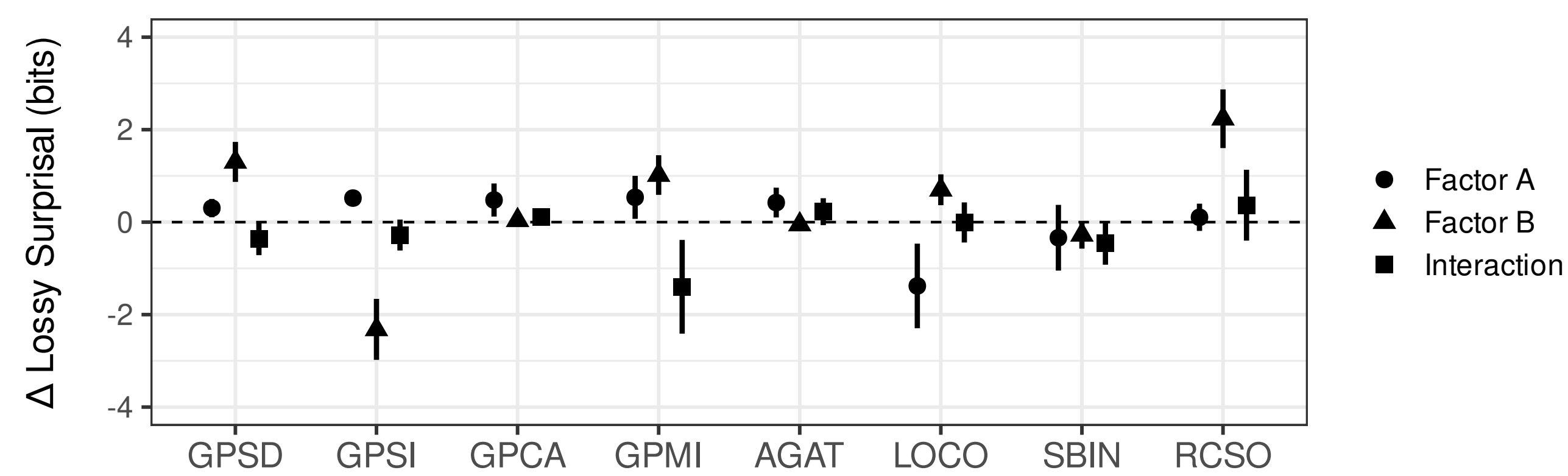
Predictions From Psycholinguistic Theory (Qualitative)



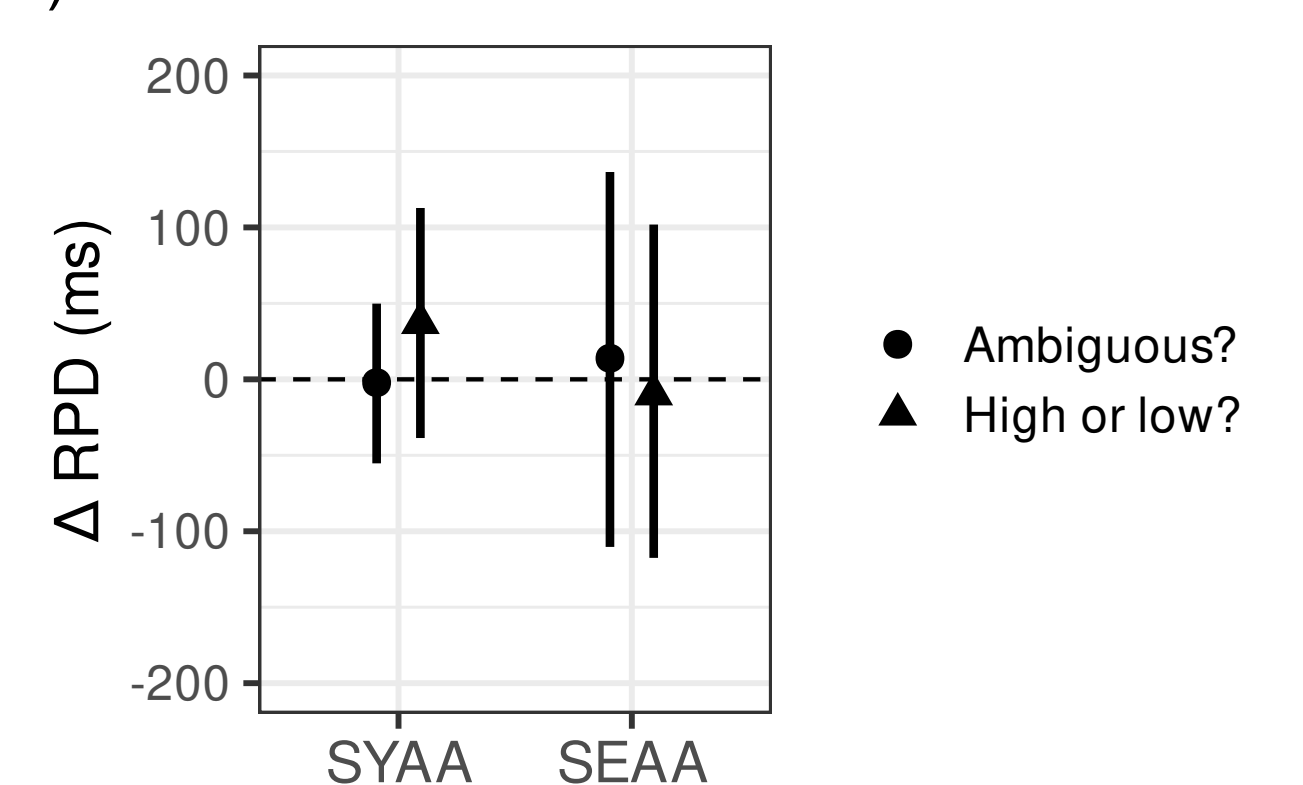
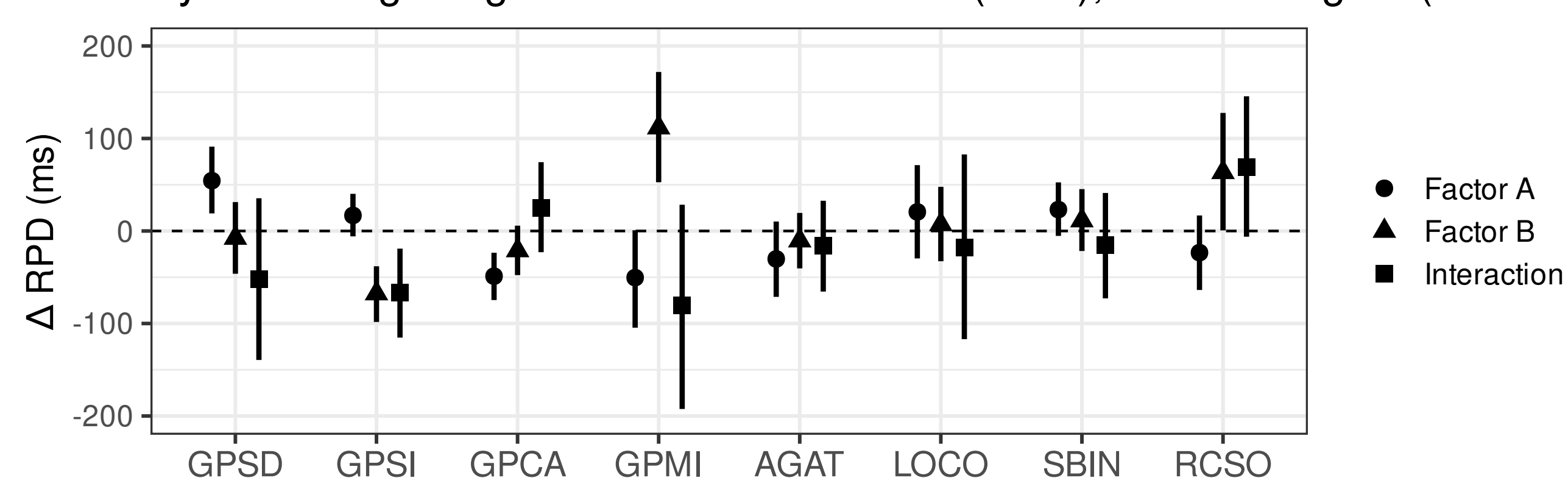
Predictions From Surprisal (95% CIs)



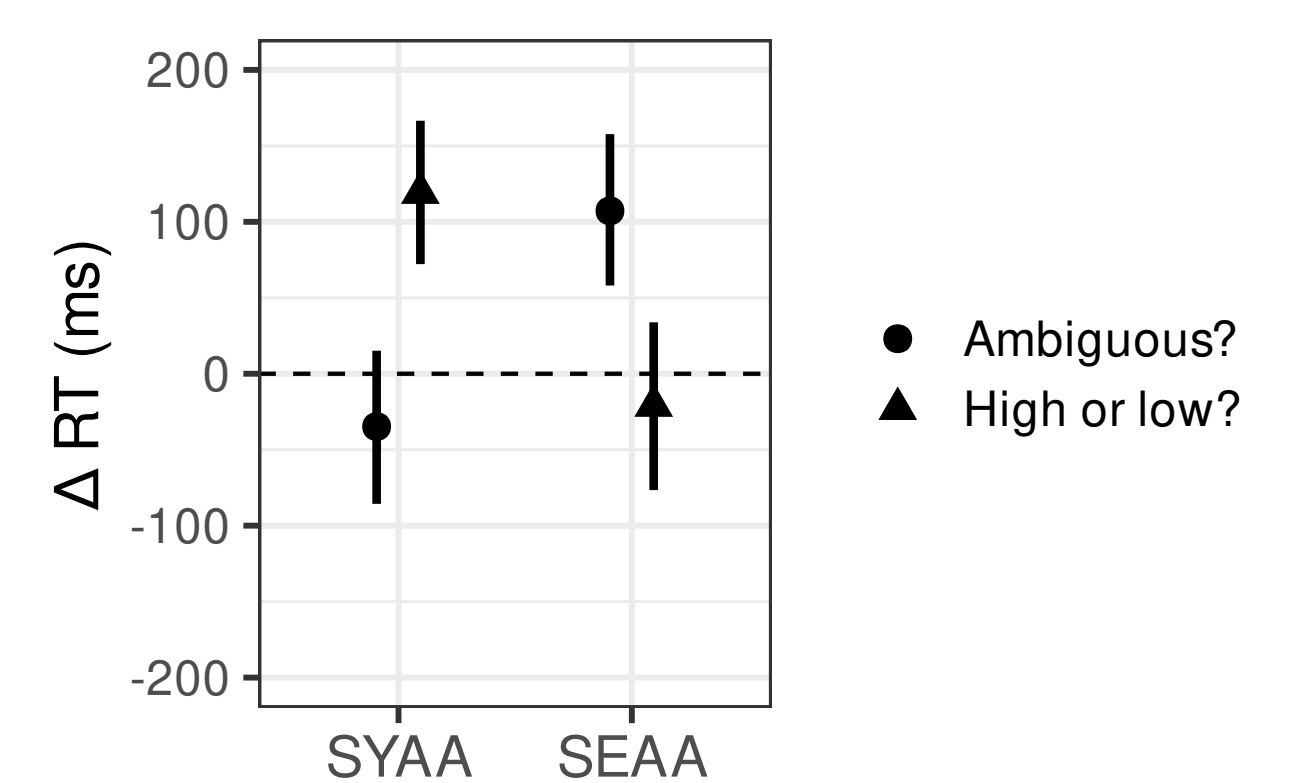
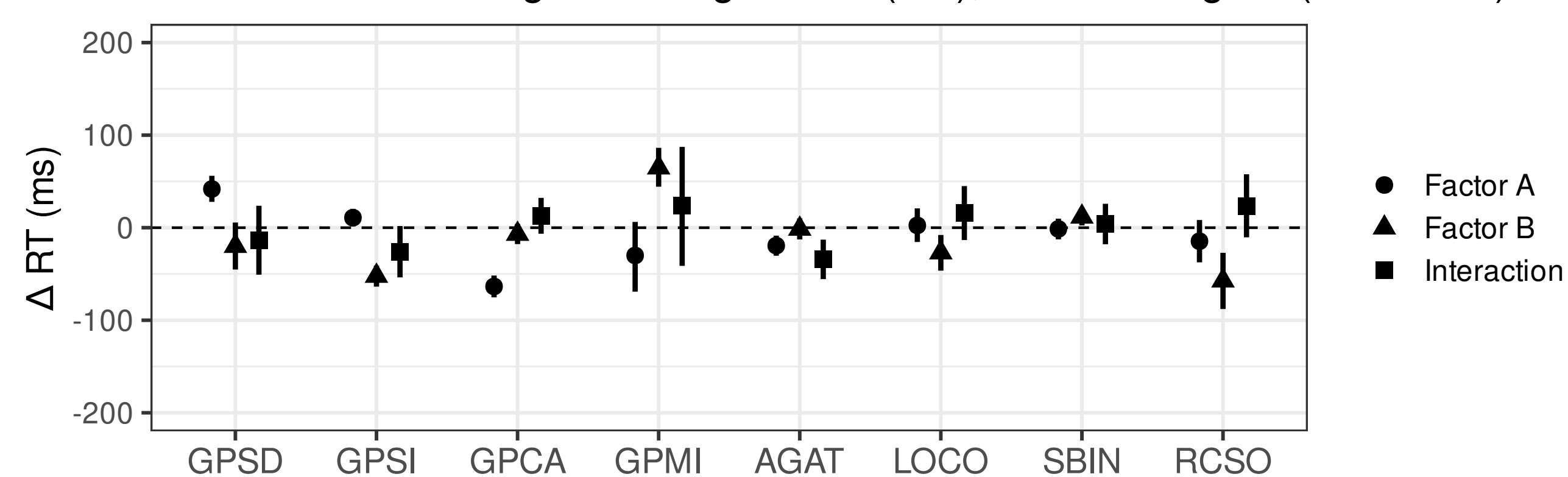
Predictions From Lossy-Context Surprisal (95% CIs)



Eye Tracking: Regression Path Durations (RPD), Critical Region (95% Crls)



Self-Paced Reading: Reading Times (RT), Critical Region (95% Crls)



References

[1] J. Devlin et al. In: *Proceedings of the 2019 Conference of the NAACL: Human Language Technologies, Volume 1 (long and short papers)*. 2019, pp. 4171–4186. [2] R. Futrell, E. Gibson, and R. P. Levy. In: *Cognitive Science* 44.3 (2021), e12814. [3] R. Futrell et al. In: *Language Resources and Evaluation* 55 (2021), pp. 63–77. [4] J. T. Hale. In: *Proceedings of the Second Meeting of the NAACL*. Pittsburgh, PA, 2001. [5] J. Häussler. PhD thesis. University of Konstanz, 2009. [6] J. Hennert et al. Unpublished manuscript. 2025. [7] F. Hsiao and E. Gibson. In: *Cognition* 90.1 (2003), pp. 3–27. [8] K.-J. Huang et al. In: *Journal of Memory and Language* 137 (2024), p. 104510. [9] A. van Kampen. PhD thesis. Free University of Berlin, 2001. [10] L. Konieczny, B. Hemforth, and C. Scheepers. In: *German Sentence Processing*. Ed. by B. Hemforth and L. Konieczny. Springer, 2000, pp. 247–278. [11] R. Levy. In: *Cognition* 106.3 (2008), pp. 1126–1177. [12] P. Logačev. In: *Journal of Experimental Psychology: Learning, Memory, and Cognition* 49.9 (2023), p. 1471. [13] M. Meng and M. Bader. In: *Language and Speech* 43.1 (2000), pp. 43–74. [14] M. Meng and M. Bader. In: *Language and Cognitive Processes* 15.6 (2000), pp. 615–666. [15] D. Paape and S. Vasishth. In: *Language and Speech* 59.3 (2016), pp. 387–403. [16] A. Radford et al. In: *OpenAI Blog* 1.8 (2019), p. 9. [17] P. Schoknecht, H. Yadav, and S. Vasishth. In: *Journal of Memory and Language* 141 (2025), p. 104599. [18] B. Staatsbibliothek. <https://huggingface.co/dbmdz/german-gpt2>. 2020. [19] M. J. Traxler, M. J. Pickering, and C. Clifton Jr. In: *Journal of Memory and Language* 39.4 (1998), pp. 558–592. [20] A. Vehtari, A. Gelman, and J. Gabry. In: *Statistics and Computing* 27 (2017), pp. 1413–1432.