If Attention Serves as a Cognitive Model of Human Memory Retrieval, What is the Plausible Memory Representation?

Anonymous ACL submission

Abstract

Recent work in computational psycholinguistics has revealed intriguing parallels between attention mechanisms and human memory retrieval, focusing primarily on Transformer architectures that operate on token-level representations. However, computational psycholinguistic research has also established that syntactic 800 structures provide compelling explanations for human sentence processing that word-level factors alone cannot fully account for. In this study, we investigate whether the attention mechanism of Transformer Grammar (TG), which uniquely operates on syntactic structures as representa-013 tional units, can serve as a cognitive model of human memory retrieval, using Normalized Attention Entropy (NAE) as a linking hypothesis 017 between model behavior and human processing difficulty. Our experiments demonstrate that TG's attention achieves superior predictive power for self-paced reading times compared to vanilla Transformer's, with further analyses revealing independent contributions from both models. These findings suggest that human sentence processing involves dual memory representations-one based on syntactic structures and another on token sequenceswith attention serving as the general retrieval algorithm, while highlighting the importance of incorporating syntactic structures as representational units.

1 Introduction

032Whether language models (LMs) developed in nat-
ural language processing (NLP) are plausible as
cognitive models of human sentence processing
is a central question in computational psycholin-
guistics. Over the past two decades, this ques-
tion has been primarily addressed from the per-
spective of *expectation-based theories*—one of the
two major classes of human sentence processing
theory—examining whether LMs' next-word pre-
diction can serve as a model of human predictive

processing (Hale, 2001; Levy, 2008; Wilcox et al., 2020; Merkx and Frank, 2021; *inter alia*).

The recent success of Transformer architectures (Vaswani et al., 2017) in NLP has unexpectedly opened a new avenue of investigation from the perspective of *memory-based theories*, the other major class of sentence processing theory. Researchers have proposed that the attention mechanism, despite its engineering origins, can implement a human memory retrieval theory known as cue-based retrieval (Van Dyke and Lewis, 2003). Recent studies have revealed intriguing parallels between the weighted reference patterns exhibited by the attention mechanism and the elements that humans may retrieve during online sentence comprehension (Ryu and Lewis, 2021; Oh and Schuler, 2022; Timkey and Linzen, 2023).

Computational psycholinguistics has also established that human sentence processing cannot be fully explained by word-level factors alone; rather, *syntactic structures* have provided compelling explanations for it. For instance, next-word prediction from LMs that explicitly incorporate syntactic structure building demonstrates superior performance in accounting for human brain activity compared to vanilla RNNs and Transformers (Hale et al., 2018; Wolfman et al., 2024); the number of syntactic nodes hypothesized to be constructed per word correlates significantly with both reading times (Kajikawa et al., 2024) and neural activity patterns (Brennan et al., 2012).

Given these findings, if attention can serve as a general algorithm for memory retrieval in human sentence processing, human memory retrieval should be captured by the attention mechanism operating on syntactic structures as well as that operating on token sequences. In this study, we investigate whether the attention mechanism of Transformer Grammar (TG; Sartran et al., 2022), which uniquely operates on syntactic structures as representational units, can serve as a cognitive model of

081

042

043

human memory retrieval, using Normalized Attention Entropy (NAE; Oh and Schuler, 2022) as the 084 linking hypothesis between model behavior and human processing difficulty. Our experiments demonstrate that TG's attention achieves superior predictive power for self-paced reading times compared to vanilla Transformer's, with further analyses revealing independent contributions from both models. These findings suggest that human sentence 091 processing involves dual memory representationsone based on syntactic structures and another on token sequences-with attention serving as the general retrieval algorithm, while highlighting the importance of incorporating syntactic structures as representational units.

2 Background

101

102

103

105

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

127

128

129

131

2.1 Normalized Attention Entropy (NAE)

Many psycholinguistic studies assume that human sentence processing involves memory retrieval, where based on the various cues provided by the current input word (e.g., verbs), elements (e.g., their arguments) are retrieved from working memory. For example, in (1), which is taken from Van Dyke (2002), when the verb *was complaining* is input, its subject *the resident* must be retrieved from working memory.

- a. The worker was surprised that the resident[subj,anim] [who was living near the dangerous warehouse] was complaining about the investigation.
 - b. The worker was surprised that the **resident**_[subj,anim] [who said that the warehouse_[subj] was dangerous] *was complaining* about the investigation.

According to the cue-based retrieval theory (Van Dyke and Lewis, 2003), such retrieval becomes more difficult when similar elements exist in the sentence because the cues are overloaded; for example, only in (1b), *warehouse* may interfere with *resident* since they both have the feature [subj] as a retrieval cue. Van Dyke (2002) showed that humans read *was complaining* more slowly in (1b) than in (1a), providing empirical support for the cue-based retrieval theory.

In recent computational psycholinguistics, attempts have been made to interpret the attention mechanism—a weighted reference of preceding tokens based on Query and Key vectors—as a computational implementation of cue-based retrieval. Notably, Ryu and Lewis (2021) proposed Attention Entropy (AE) as a linking hypothesis, where the diffuseness of attention weights is assumed to quantify the degree of retrieval interference. While AE was initially proposed for modeling interference effects in specific constructions, Oh and Schuler (2022) extended it to naturally occurring text by introducing two normalizations: (i) division by the maximum entropy achievable given the number of preceding tokens, and (ii) sum-to-1 renormalization of attention weights over 1-(i - 1)-th tokens (**Normalized AE, NAE**).¹

$$NAE_{l,h,i} =$$
 144

132

133

134

135

136

137

138

139

140

141

142

143

145

146

147

148

149

150

151

152

153

154

155

156

157

159

160

161

162

164

165

$$\frac{\mathbf{a}_{l,h,i[1:i-1]}^{\top}}{\log_2(i-1)\mathbf{1}^{\top}\mathbf{a}_{l,h,i[1:i-1]}} (\log_2 \frac{\mathbf{a}_{l,h,i[1:i-1]}}{\mathbf{1}^{\top}\mathbf{a}_{l,h,i[1:i-1]}}),$$
(1)

where $a_{l,h,i}$ represents the attention weight vector when using the *i*-th token as the query in the *h*-th head of layer l^2 . In this paper, we employ this NAE as a linking hypothesis between attention mechanisms and human memory retrieval.³

2.2 Transformer Grammar (TG)

Transformer Grammar (TG; Sartran et al., 2022) is a type of syntactic LM, a generative model that jointly generates token sequences x and their corresponding syntactic structures y. TG formulates the generation of (x, y) as modeling a sequence of actions, a (e.g., (S (NP The blue bird NP) (VP sings VP) S)), constructing both token sequences and their syntactic structures in a top-down, leftto-right manner. The action sequence a comprises three types of operations:

- (X: Generate a non-terminal symbol (X, where X represents a phrasal tag such as NP;
- w: Generate a terminal symbol w, where w represents a token such as bird;

¹Oh and Schuler (2022) showed that regression models for predicting reading times fail to converge with vanilla AE.

²Oh and Schuler (2022) explored NAE calculation using various attention weight formulations, but in this study, we adopt the norm-based attention weight formulation (Kobayashi et al., 2020), which achieved the highest predictive power on the self-paced reading time corpus.

³While Oh and Schuler (2022) also proposed other metrics based on distances between attention weights at consecutive time steps, we exclusively adopt NAE because (i) in TG, the number of preceding elements varies with time, making distance definition non-trivial, and (ii) Oh and Schuler (2022) demonstrated that NAE's predictive power subsumes that of distance-based metrics in the self-paced reading time corpus.



Figure 1: TG's attention mask with COMPOSE/STACK attention mechanisms, adapted from Sartran et al. (2022). COMPOSE generates a vector representation of the closed phrase, while subsequent STACK operations reference this single vector as the representation of the closed phrase. Red boxes indicate the attention weights used to calculate NAE for each word.

• X): Generate X) to close the most recent opened non-terminal symbol, where X matches the phrasal tag of the targeted nonterminal symbol.

166

167

169

170

171

172

173

174

176

177

178

179

180

181

182

183

187

189

190

191

193

The probability of action sequence a (a_1, a_2, \cdots, a_n) is decomposed using the chain rule. Formally, TG is defined as:

$$p(\boldsymbol{x}, \boldsymbol{y}) = p(\boldsymbol{a}) = \sum_{t=1}^{n} p(a_t | a_{< t}).$$
 (2)

TG's key innovation lies in its handling of closed phrases: immediately after generating X), it computes a vector representation of the closed phrase, which subsequent next-action predictions use as the representation for that phrase. Technically, this operation is realized via two components: X) action duplication and a specialized attention mask. The duplication process transforms a into a' by duplicating all X) actions (e.g., (S (NP The blue bird NP) NP) (VP sings VP) VP) S) S)), while preserving the modeling space p(a) by preventing predictions for duplicated positions. The atten-185 tion mask implements two distinct attention mechanisms: COMPOSE and STACK (Figure 1). COMPOSE operates exclusively at the first occurrence of each X) to generate the phrasal representation by attending only to vectors between the corresponding (X and X) (without making predictions). STACK operates at all other positions to compute representations for next-action prediction, with attention

restricted to positions on the *stack* (comprising unclosed non-terminals, not-composed terminals, and closed phrases).

194

195

196

197

198

199

200

201

202

203

204

205

207

208

209

210

211

212

213

214

215

216

217

218

219

220

222

223

224

225

226

227

228

229

230

232

233

234

235

Previous research has demonstrated that TG's probability estimates align more closely with human offline grammaticality judgments (Sartran et al., 2022) and online brain activity than vanilla Transformers (Wolfman et al., 2024). This study investigates whether the attention mechanism of TG, which uniquely operates on syntactic structures as representational units, can serve as a cognitive model of human memory retrieval.

3 Methods

3.1 NAE calculation with TG

The calculation of NAE with TG requires assumptions regarding two key perspectives:

- 1. What syntactic structures should be assumed for a given token sequence?
- 2. How should the cognitive load of attention from non-lexical tokens (i.e., (X and X)) be attributed to words?

In response to these considerations, we make the following assumptions:

- 1-A. We assume only the globally correct syntactic structure (i.e., "perfect oracle"; Brennan, 2016).
- 2-A. We consider only attention from words, excluding attention from non-lexical tokens.

The adoption of 1-A. is motivated by two factors. First, the self-paced reading time corpus we utilized here provides gold-standard syntactic structures for each sentence, and previous studies have developed predictors based on these annotations (Shain et al., 2020; Isono, 2024). Using the same structural assumptions enables fair comparison with these established predictors, considering the possibility of parsing errors. Second, TG's current implementation lacks beam search procedure (Stern et al., 2017; Crabbé et al., 2019), an inference technique commonly used in cognitive modeling to handle local ambiguities through parallel parsing (Hale et al., 2018; Sugimoto et al., 2024).⁴

⁴As a proof of concept, we also conducted experiments using multiple syntactic structures generated by wordsynchronous beam search with Recurrent Neural Network Grammar (Dyer et al., 2016; Kuncoro et al., 2017; Noji and Oseki, 2021), obtaining similar results (Appendix D).

Regarding 2-A., given the multiple possible approaches to attributing processing load from nonlexical tokens to words, we adopt the most straightforward and theoretically neutral approach. Figure 1 denotes the attention weights used to calculate NAE for each word, with red boxes.

3.2 Settings

236

237

240

241

242

243

245

246

247

249

252

254

259

261

262

264

266

269

271

275

276

281

Language models We used 16-layer, 8-head TG and Transformer (252M parameters).⁵ All hyperparameters followed the default settings described in Sartran et al. (2022) (see Appendix A). Following Oh and Schuler (2022), we computed NAE separately for each attention head at the topmost layers and then summed the values across heads.

Training data We used BLLIP-LG, a dataset containing 42M tokens (1.8M sentences) from the Brown Laboratory for Linguistic Information Processing (BLLIP) 1987–89 WSJ Corpus Release 1 (Charniak et al., 2000).⁶ The corpus was re-parsed using a state-of-the-art constituency parser (Kitaev and Klein, 2018) and split into trainval-test sets by Hu et al. (2020). BLLIP-LG has been widely used for training syntactic LMs, including TG. Following Sartran et al. (2022), we trained a 32K SentencePiece tokenizer (Kudo and Richardson, 2018) on the training set and segmented each sentence into subword units.

Both TG and Transformer were trained at the sentence level: TG maximized the joint probability p(x, y) on action sequences, while Transformer maximized the probability p(x) on terminal subword sequences. For training hyperparameters, we largely followed the default settings in Sartran et al. (2022) but adjusted the batch size to fit within the memory constraints of our hardware (NVIDIA A100, 40GB). Accordingly, we tuned other hyperparameters (e.g., learning rate) to maintain training stability. We trained three models with different random seeds and selected the checkpoint with the lowest validation loss for each run.

Reading time data We used the Natural Stories corpus (Futrell et al., 2018)⁷ consisting of 10 stories (485 sentences, 10,245 words) with self-paced reading times collected from 181 anonymized native English speakers. Following Futrell et al.'s preprocessing, data points were removed if (i) a

⁶https://catalog.ldc.upenn.edu/LDC2000T43

participant scored less than 5/6 on comprehension questions for a story or (ii) individual reading times were less than 100 ms or greater than 3,000 ms. Following Oh and Schuler (2022), we also excluded sentence-initial and sentence-final data points. We further removed sentence-second data points, as they lack a log trigram frequency of the previous token required for our baseline regression model. After preprocessing, 724,883 data points from 180 participants remained for statistical analysis, out of the original 848,747 data points. 282

284

287

288

289

290

291

292

293

294

296

297

298

299

301

302

303

304

305

306

307

309

310

311

312

313

314

315

316

317

318

319

321

322

323

324

Statistical analysis We evaluated each LM's NAE contribution to reading time prediction by measuring improvements in regression model fit when adding NAE as predictors. For each model (TG/Transformer), we included both the current word's NAE (tg_nae/tf_nae) and the previous word's NAE (tg_nae_so/tf_nae_so) to account for spillover effects.⁸⁹ Model improvement was quantified as the increase in log-likelihood (Δ LogLik). This evaluation was conducted for each random seed, and we report the mean Δ LogLik with standard deviation.

The baseline regression model included standard predictors from the Natural Stories corpus:

- zone and position: word position in the story and sentence;
- wordlen: number of characters in the word;
- unigram, bigram, and trigram: logtransformed n-gram frequencies.

We additionally included the following predictors:

- tg_surp and tf_surp: surprisal from TG and Transformer;
- stack_count: number of elements in the *stack* (comprising unclosed non-terminals, not-composed terminals, and closed phrases).

Following Oh and Schuler (2022), we included surprisal to test NAE's significance in the presence of surprisal predictors from the same LMs.¹⁰ Stack count was included to isolate the cost of holding elements (Joshi, 1990; Abney and Johnson, 1991; Resnik, 1992) from their interference effects, which TG's NAE was designed to capture.

⁵https://github.com/google-deepmind/ transformer_grammars

⁷https://github.com/languageMIT/naturalstories

⁸_so indicates spillover.

⁹Following Oh and Schuler (2022), we summed the subword NAE values for each word.

¹⁰For an experiment on the predictive power of surprisal itself, see Appendix E.

Model	$\Delta LogLik(\uparrow)$	Predictor	Effect size [ms]	p-value range	Significant seeds
TG	96.3 (±3.0)	tg_nae tg_nae_so	$\begin{array}{c} 2.91 \ (\pm \ 0.1) \\ 0.655 \ (\pm \ 0.1) \end{array}$	<0.001 <0.05	3/3 3/3
Transformer	35.2 (±7.0)	tf_nae tf_nae_so	$\begin{array}{c} 1.77\ (\pm\ 0.2)\\ 0.293\ (\pm\ 0.2)\end{array}$	< 0.001 0.04–0.38	3/3 1/3

Table 1: TG's and Transformer's NAE contribution to reading time prediction (Δ LogLik). The effect size per standard deviation is shown for each model-derived predictor, along with the *p*-value range across random seeds and the number of seeds showing significant contributions. Standard deviations across seeds for Δ LogLik and effect sizes are shown in parentheses. The mean reading time in the analysis is 335 ms.

All predictors were *z*-transformed, and we also included the previous word's values as predictors to model spillover, except for the positional information. The baseline regression model was a linear mixed-effects model (Baayen et al., 2008) with these fixed effects and by-subjects and by-story random intercepts:

32	$\log(extsf{RT}) \sim extsf{zone} + extsf{position} + extsf{wordlen} +$
33	${\tt unigram+bigram+trigram+}$
34	${\tt tf_surp} + {\tt tg_surp} + \\$
35	${\tt stack_count+wordlen_so+}$
36	$\verb"unigram_so+bigram_so+$
37	$\tt trigram_so + tf_surp_so +$
38	$\verb"tg_surp_so+stack_count_so+$
39	(1 participant) + (1 story) (3)

To assess each LM's independent contribution to reading time prediction, we also conducted likelihood ratio tests by extending Equation 3 in two ways: adding both LMs' NAE versus adding only one LM's NAE. Note that a larger Δ LogLik from one LM does not necessarily indicate that it contributes above and beyond the other LM, nor does a smaller Δ LogLik indicate no unique contribution. Following Aurnhammer and Frank (2019), we used NAE and surprisal values averaged across random seeds for these nested model comparisons.

4 Results

326

327

328

329

330

331

3

3

3

341

347

352

4.1 Does TG's NAE have predictive power for reading times?

354Table 1 presents the contributions of TG's and355Transformer's NAE to reading time prediction.356First, Transformer's NAE exhibited significant pre-357dictive power for reading times, independent of358baseline predictors such as surprisal. The effect359size was in the expected positive direction (higher

NAE values corresponding to longer reading times), primarily showing the immediate effect. This corroborates the arguments of Ryu and Lewis and Oh and Schuler that the attention mechanism—the weighted reference of preceding tokens—functions as a cognitive model of human memory retrieval, despite its engineering-oriented origins. 360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

383

384

385

387

389

390

391

392

393

394

395

396

397

398

Second, TG's NAE exhibited robust predictive power, demonstrating significant positive effects in both immediate and spillover contexts. This finding not only provides additional evidence for incremental construction of syntactic structures in human sentence processing (e.g., Fossum and Levy, 2012), but also suggests that TG's attention mechanism effectively models memory retrieval from these constructed syntactic representations.

Finally, TG's NAE made a substantially stronger contribution to reading time prediction $(\Delta LogLik=96.3)$ compared to Transformer's NAE $(\Delta LogLik=35.2)$. This finding suggests that retrieval from syntactic memory representations plays a more dominant role in human sentence processing than retrieval from lexical memory representations. This underscores the importance of incorporating syntactic structures as a unit of memory representation, which we implemented through the integration of TG and NAE here.

4.2 Do TG's and Transformer's NAE have independent contributions?

Figure 2 presents the results of likelihood ratio tests examining the independence of TG's and Transformer's NAE contributions. The regression model incorporating NAE from both LMs ('TG & Transformer') demonstrated significantly higher predictive power than the models containing NAE from either LM alone ('TG' or 'Transformer'). This reveals that TG's NAE certainly captures variance in reading times that Transformer's NAE cannot explain, while Transformer's NAE, despite its lower



Model Δ LogLik Predictor *p*-value 96.1 *_nae (3/3)TG (± 15.9) ** (3/3)*_nae_so *** 86.3 (3/3)*_nae TG_{-comp} (± 31.8) n.s. (0/3) *_nae_so

Table 2: TG's and TG_{-comp}'s contribution to reading time prediction. The rightmost column shows the *p*value range across random seeds (*** p < 0.001, ** p < 0.01, and *n.s.* not significant), along with the number of seeds showing significant contributions. Due to the potential multicollinearity between the Transformer's NAE and TG/TG_{-comp}'s NAE, the column of the effect size is omitted.

Figure 2: Likelihood ratio test results examining the independence of NAE's predictive power

overall predictive power, accounts for unique variance not captured by TG's NAE. This finding aligns with psycholinguistic literature, where cognitive models of memory retrieval encompass both syntax-based approaches (e.g., verb-argument relationships; Lewis and Vasishth, 2005) and semanticbased approaches (e.g., bag-of-words-like similarity; Brouwer et al., 2012), suggesting that the attention mechanisms of TG and Transformer serve as complementary cognitive models, each capturing distinct aspects of human memory retrieval.

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

4.3 What aspect of memory retrieval do TG's and Transformer's NAE capture?

To investigate the aspects of human memory retrieval captured by TG's and Transformer's NAE, we analyzed differences in prediction improvement across part-of-speech (POS) tags annotated in the Natural Stories corpus. Our analysis followed three steps: (i) selecting POS tags with more than 1,000 occurrences, (ii) for each POS tag, testing the significance of improvement from the baseline regression model (measured in Δ Root Mean Squared Error, $\Delta RMSE$) when adding NAE of the current and previous word as fixed effects,¹¹ and (iii) examining the significance of *differences* in Δ RMSE between TG and Transformer for POS tags where either model showed significant improvement. We assessed significance using Wilcoxon signed-rank tests with Bonferroni correction (p < 0.05).

Figure 3 presents the differences in prediction improvement across POS tags.¹² Consistent with the larger Δ LogLik value, TG's NAE demonstrated advantages over Transformer's NAE across a broader range of POS tags. Notably, TG's NAE exhibited superior improvement across verbs (VB, VBD, VBG, VBP), while Transformer's NAE excelled only for possessive pronouns (\$PRP). These findings indicate that different types of retrieval operations—verb-triggered retrieval (e.g., subject retrieval) and possessive pronoun-triggered retrieval (e.g., antecedent retrieval)—are better modeled by distinct cognitive mechanisms: attention with syntactic and token memory representations, respectively. This pattern supports our earlier argument regarding the complementary nature of these models (Section 4.2). 431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

5 Follow-up analysis

5.1 Do TG's advantages stem from the COMPOSE attention?

As described in Section 2.2, TG's key feature is the COMPOSE attention, which explicitly generates single vector representations for closed phrases. Here, we investigate whether TG's predictive power derives from merely considering syntactic structures or from explicitly treating closed phrases as single representations (see Hale et al., 2018; Brennan et al., 2020). To address this question, we developed TG_{-comp} , a TG variant that processes each action in the action sequence a as an individual token (i.e., Choe and Charniak's approach), without the COMPOSE attention. We trained TG_{-comp} with identical hyperparameters as TG. The baseline regression model (Equation 3) was augmented with (i) TG_{-comp}'s surprisal and (ii) Transformer's NAE to (i) ensure a fair comparison between TG and TG_{-comp} and (ii) distinguish

¹¹We used the same regression models as in Section 4.2, where surprisal and NAE values were averaged across seeds.

¹²For a complete list of POS tags in the Natural Stories corpus, see Appendix C.



Figure 3: Differences in reading time prediction improvement (Δ RMSE) between TG and Transformer across POS tags (TG - Transformer). The y-axis shows the mean differences per word, with the error bars representing standard errors. Only POS tags showing significant improvement in either model and significant differences between models are displayed. Statistical significance after Bonferroni correction: ** p < 0.01, *** p < 0.001.



Figure 4: Differences in $\triangle RMSE$ between TG and TG_{-comp} across POS tags (TG - TG_{-comp})

between the effects of direct terminal token access and syntactic structure consideration in TG_{-comp} .

Table 2 presents the Δ LogLik values obtained when incorporating either TG's or TG_{-comp}'s NAE as fixed effects into the baseline regression model. Note that due to the potential multicollinearity between Transformer's NAE and TG/TG_{-comp}'s NAE, we focus on the Δ LogLik values and significance of the contribution rather than individual effect sizes. Our analysis reveals two key findings. First, TG_{-comp}'s NAE demonstrates significant predictive power for reading times, even in the presence of Transformer's NAE, implying that consideration of syntactic structures alone captures certain memory retrievals based on syntactic information. Second, although we should note that the standard deviation across seeds is relatively large, TG's NAE outperforms TG_{-comp}'s, suggesting that the attention mechanism operating on syntactic memory representations more effectively captures variance in syntax-based memory retrieval. Likelihood ratio tests further revealed that TG's NAE captured reading time patterns unexplainable by TG_{-comp} ('TG & TG_{-comp}'> 'TG_{-comp}', p < 0.001). However, we also found that TG_{-comp}, despite its lower overall predictive power, also explained unique variance ('TG & TG_{-comp}'> 'TG', p < 0.001).

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

504

506

507

508

509

510

511

512

To investigate the underlying mechanisms of this complementary relationship, we analyzed the Δ RMSE differences across POS tags (Figure 4). The results showed that TG's NAE demonstrated advantages over TG_{-comp}'s NAE across most POS tags, with notable exceptions for VBG, VBP, and WP-POS tags that typically involve immediate modification of the previously composed phrase (e.g., (NP (NP ... NP) (SBAR (WHNP (WP who). This pattern leads us to hypothesize that memory retrieval based on syntactic cues typically operates on composed representations for efficiency, but in contexts where immediate modification of a composed phrase occurs (e.g., relative clauses), humans may strategically retrieve individual lexical elements within the composed phrases.

5.2 Does TG's NAE capture interference effects?

Psycholinguistic research has identified two primary types of memory retrieval costs: *interference* effects, which NAE aims to capture, and *decay*

effects-the cognitive load associated with access-513 ing elements at greater linear distances (e.g., Gib-514 son, 1998, 2000). Here, we examine whether 515 TG's NAE genuinely captures interference effects 516 by testing its independence from variables that model memory decay effects. For modeling decay 518 effects, we employed Category Locality Theory 519 (CLT; Isono, 2024),¹³ which treats phrases in syntactic structure¹⁴ as representational units of mem-521 ory and quantifies decay effects using the distance 522 (measured in content words) between an input and 524 the phrases to be composed with it.

526

527

528

530

531

539

540

543

544

545

547

551

552

553

554

To assess independence, we tested whether TG's NAE and CLT maintain their contributions when simultaneously included in the baseline regression model (Equation 3), and examined their independence through likelihood ratio tests.¹⁵ The results (Table 3) show that TG's NAE exhibited significant effects in both immediate and spillover conditions, and CLT demonstrated a significant immediate effect (with a marginally significant spillover effect). A nested model comparison confirmed that these effects were independent ('TG & CLT'>'CLT', p < 0.001; 'TG & CLT'>'TG', p < 0.05).

These results provide empirical evidence that NAE quantifies interference rather than decay in memory retrieval—extending beyond previous studies on NAE (Ryu and Lewis, 2021; Oh and Schuler, 2022). This finding is significant because, as far as we are aware, while psycholinguistics has developed various implementations of memory decay effects, it has lacked broad-coverage implementations of interference effects applicable to naturally occurring texts. Our results suggest that NAE represents a promising approach for quantifying interference effects in a broad-coverage manner.

6 Level of description

In cognitive modeling studies based on surprisal theory, explanations typically follow the form "if these LMs were models of human prediction, the difficulty of next-word disambiguation that humans solve would be approximated as follows." Such ex-

Model	Predictor	Effect size [ms]	
TG & CLT	tg_nae tg_nae_so	2.93*** 0.69**	
	clt clt_so	0.32^{*} 0.25^{\cdot}	

Table 3: Effect sizes per standard deviation are shown for TG's NAE and CLT predictors. Significance levels: *** p < 0.001, ** p < 0.01, * p < 0.05, p < 0.1.

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

planations typically operate at the most abstract of Marr's three levels of description-the computational level. Recently, Futrell et al. (2020) proposed lossy-context surprisal to integrate memory representation perspectives into surprisal theory. However, as they explicitly stated, this theory remains at the computational level, relaxing assumptions about memory representations in human predictive processing. In contrast, cognitive models of human memory, such as cue-based retrieval, generally provide explanations about mechanisms that deal with specific mental representations. These explanations typically move down one level to the algorithmic level of description. While not explicitly stated by the authors, we argue that the work of Ryu and Lewis and Oh and Schuler-conceptualizing attention mechanisms as implementations of cue-based retrieval—also operates at the algorithmic level, similar to cue-based retrieval itself. Our research can be characterized as an investigation into the nature of memory representations, addressing fundamental questions at this level (see Hale, 2014).

7 Conclusion

In this paper, we have demonstrated that attention can serve as the general algorithm for modeling human memory retrieval from two representational systems. Furthermore, we have shown that among the LMs examined in this study (TG, TG_{-comp} , and Transformer), TG—whose attention mechanism uniquely operates on syntactic structures as representational units—best captures dominant factors in human sentence processing. Our results suggest that the integration of attention mechanisms (developed in NLP) with syntactic structures (theorized in linguistics) constitutes a broad-coverage candidate implementation for human memory retrieval. We hope these findings will foster greater collaboration between these two fields.

¹³Although Dependency Locality Theory (DLT; Gibson, 1998, 2000) is widely recognized as one of the most prominent models for capturing decay effects, we opted for CLT in this study, following Isono's finding that DLT-based predictors fail to achieve statistical significance in explaining reading times in the Natural Stories corpus.

¹⁴CLT assumes syntactic structure based on Combinatory Categorial Grammar (Steedman, 2000).

¹⁵As in other likelihood ratio tests, we used surprisal and NAE values averaged across random seeds.

593

594

595

596

598

602

610

611

612

614

615

617

618

619

622

623

625

631

639

640

643

Limitations

Our NAE calculation comprised three steps: (i) computing NAE for each attention head in the topmost layers, (ii) adding the values across heads, and (iii) summing subword-level values into word level. While this procedure strictly adhered to Oh and Schuler (2022), alternative approaches to handling layers, attention heads (Ryu and Lewis, 2021), and subword tokens (Oh and Schuler, 2024; Giulianelli et al., 2024) warrant investigation.

While our study provides an in-depth investigation using the Natural Stories corpus-an English self-paced reading time dataset-the breadth of our analysis has certain limitations. The generalizability of our findings to different languages (e.g., Japanese self-paced reading time corpus from Asahara, 2022) and other cognitive load (e.g., gaze duration from Kennedy et al., 2003 or EEG and fMRI from Bhattasali et al., 2020) remains to be investigated.

As discussed in Section 3.1, we employed "perfect oracles" as syntactic structures behind token sequences. This idealization leaves the resolution of local ambiguities, which humans encounter during actual sentence processing, outside the scope of our study (for a conceptual case study, see Appendix D). By incorporating these kinds of entirely new factors, more detailed models could emerge.

Following prior work (Wolfman et al., 2024), we adopted the default TG implementation of a top-down parsing strategy. However, psycholinguistic literature has suggested that a left-corner parsing strategy might be more plausible for human sentence processing (Abney and Johnson, 1991; Resnik, 1992). While previous studies have primarily evaluated the plausibility of parsing strategies from a memory capacity perspective (cf. stack_count), TG's NAE might offer a new opportunity to revisit the question of cognitively plausible parsing strategies from the perspective of memory interference.

Finally, while this paper focused on investigating the attention mechanism of TG through the lens of memory-based theory, exploring TG as an integrated implementation for expectation-based theory (via surprisal) and memory-based theory (via NAE) represents a promising future direction (Michaelov et al., 2021; Ryu and Lewis, 2022). Specifically, future work could investigate the attention mechanism of TG as the underlying driver of surprisal's predictive power (Appendix E), ana-

lyzing the relationship between surprisal and NAE.	644
Ethical considerations	645
We employed AI-based tools (Claude, ChatGPT,	646
GitHub Copilot, and Grammarly) for writing and	647
coding assistance. These tools were used in compli-	648
ance with the ACL Policy on the Use of AI Writing	649
Assistance.	650
References	651
Steven P. Abney and Mark Johnson. 1991. Mem-	652
ory requirements and local ambiguities of parsing	653
strategies. <i>Journal of Psycholinguistic Research</i> ,	654
20(3):233–250.	655
Masayuki Asahara. 2022. Reading Time and Vo-	656
cabulary Rating in the Japanese Language: Large-	657
Scale Japanese Reading Time Data Collection Us-	658
ing Crowdsourcing. In <i>Proceedings of the Thir-</i>	659
<i>teenth Language Resources and Evaluation Confer-</i>	660
<i>ence</i> , pages 5178–5187, Marseille, France. European	661
Language Resources Association.	662
Christoph Aurnhammer and Stefan L. Frank. 2019.	663
Comparing Gated and Simple Recurrent Neural Net-	664
work Architectures as Modelsof Human Sentence	665
Processing. <i>Proceedings of the Annual Meeting of</i>	666
<i>the Cognitive Science Society</i> , 41(0).	667
R. H. Baayen, D. J. Davidson, and D. M. Bates. 2008.	668
Mixed-effects modeling with crossed random effects	669
for subjects and items. <i>Journal of Memory and Lan-</i>	670
<i>guage</i> , 59(4):390–412.	671
Douglas Bates, Martin Mächler, Ben Bolker, and Steve	672
Walker. 2015. Fitting linear mixed-effects models	673
using lme4. <i>Journal of Statistical Software</i> , 67(1):1–	674
48.	675
Shohini Bhattasali, Jonathan Brennan, Wen-Ming Luh,	676
Berta Franzluebbers, and John Hale. 2020. The Al-	677
ice Datasets: fMRI & EEG Observations of Natural	678
Language Comprehension. In <i>Proceedings of the</i>	679
<i>Twelfth Language Resources and Evaluation Confer-</i>	680
<i>ence</i> , pages 120–125, Marseille, France. European	681
Language Resources Association.	682
Jonathan Brennan. 2016. Naturalistic Sentence Com-	683
prehension in the Brain. <i>Language and Linguistics</i>	684
<i>Compass</i> , 10(7):299–313.	685
Jonathan Brennan, Yuval Nir, Uri Hasson, Rafael	686
Malach, David J. Heeger, and Liina Pylkkänen. 2012.	687
Syntactic structure building in the anterior tempo-	688
ral lobe during natural story listening. <i>Brain and</i>	689
<i>Language</i> , 120(2):163–173.	690
Jonathan R. Brennan, Chris Dyer, Adhiguna Kuncoro,	691
and John T. Hale. 2020. Localizing syntactic pre-	692
dictions using recurrent neural network grammars.	693

694

Neuropsychologia, 146:107479.

- 702 707 710 711 712 713 714 715 716 717 718 721 722 723 724 725 726 727 729 734 735 736 737 738 739 740 741 742 743 744

- 745
- 747 748
- 750

- Harm Brouwer, Hartmut Fitz, and John Hoeks. 2012. Getting real about Semantic Illusions: Rethinking the functional role of the P600 in language comprehension. Brain Research, 1446:127–143.
- Eugene Charniak, Don Blaheta, Niyu Ge, Keith Hall, John Hale, and Mark Johnson. 2000. BLLIP 1987-89 WSJ Corpus Release 1.
- Do Kook Choe and Eugene Charniak. 2016. Parsing as Language Modeling. In Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, pages 2331–2336, Austin, Texas. Association for Computational Linguistics.
- Benoit Crabbé, Murielle Fabre, and Christophe Pallier. 2019. Variable beam search for generative neural parsing and its relevance for the analysis of neuroimaging signal. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 1150–1160, Hong Kong, China. Association for Computational Linguistics.
- Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc Le, and Ruslan Salakhutdinov. 2019. Transformer-XL: Attentive Language Models beyond a Fixed-Length Context. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 2978–2988, Florence, Italy. Association for Computational Linguistics.
 - Chris Dyer, Adhiguna Kuncoro, Miguel Ballesteros, and Noah A. Smith. 2016. Recurrent Neural Network Grammars. In Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 199–209, San Diego, California. Association for Computational Linguistics.
- Victoria Fossum and Roger Levy. 2012. Sequential vs. Hierarchical Syntactic Models of Human Incremental Sentence Processing. In Proceedings of the 3rd Workshop on Cognitive Modeling and Computational Linguistics (CMCL 2012), pages 61-69, Montréal, Canada. Association for Computational Linguistics.
- Richard Futrell, Edward Gibson, and Roger P. Levy. 2020. Lossy-Context Surprisal: An Information-Theoretic Model of Memory Effects in Sentence Processing. Cognitive Science, 44(3):e12814.
- Richard Futrell, Edward Gibson, Harry J. Tily, Idan Blank, Anastasia Vishnevetsky, Steven Piantadosi, and Evelina Fedorenko. 2018. The Natural Stories Corpus. In Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018), Miyazaki, Japan. European Language Resources Association (ELRA).
- Edward Gibson. 1998. Linguistic complexity: Locality of syntactic dependencies. Cognition, 68(1):1-76.
- Edward Gibson. 2000. The dependency locality theory: A distance-based theory of linguistic complexity. In

Image, Language, Brain: Papers from the First Mind Articulation Project Symposium, pages 94–126. The MIT Press, Cambridge, MA, US.

751

752

753

754

755

758

759

760

761

762

763

764

765

766

767

768

769

770

771

774

775

776

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

- Mario Giulianelli, Luca Malagutti, Juan Luis Gastaldi, Brian DuSell, Tim Vieira, and Ryan Cotterell. 2024. On the Proper Treatment of Tokenization in Psycholinguistics. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 18556–18572, Miami, Florida, USA. Association for Computational Linguistics.
- John Hale. 2001. A Probabilistic Earley Parser as a Psycholinguistic Model. In Second Meeting of the North American Chapter of the Association for Computational Linguistics.
- John Hale, Chris Dyer, Adhiguna Kuncoro, and Jonathan Brennan. 2018. Finding syntax in human encephalography with beam search. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 2727-2736, Melbourne, Australia. Association for Computational Linguistics.
- John T. Hale. 2014. Automaton Theories of Human Sentence Comprehension. CSLI Studies in Computational Linguistics. CSLI Publications, Stanford, CA.
- Jennifer Hu, Jon Gauthier, Peng Oian, Ethan Wilcox, and Roger Levy. 2020. A Systematic Assessment of Syntactic Generalization in Neural Language Models. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 1725–1744, Online. Association for Computational Linguistics.
- Shinnosuke Isono. 2024. Category Locality Theory: A unified account of locality effects in sentence comprehension. Cognition, 247:105766.
- Aravind K. Joshi. 1990. Processing crossed and nested dependencies: An automation perspective on the psycholinguistic results. Language and Cognitive Pro*cesses*, 5(1):1–27.
- Kohei Kajikawa, Ryo Yoshida, and Yohei Oseki. 2024. Dissociating Syntactic Operations via Composition Count. Proceedings of the Annual Meeting of the Cognitive Science Society, 46(0).
- Alan Kennedy, Robin Hill, and Joël Pynte. 2003. The dundee corpus. In Proceedings of the 12th European Conference on Eye Movement.
- Diederik P. Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings.
- Nikita Kitaev and Dan Klein. 2018. Constituency Parsing with a Self-Attentive Encoder. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 2676–2686, Melbourne, Australia. Association for Computational Linguistics.

910

911

912

913

914

915

916

917

862

863

- 811
- 814
- 815
- 817

- 823
- 825
- 828

- 833

835 836

837

- 841
- 844 845
- 847
- 850
- 851
- 853
- 855 856

857 858

- 861

Goro Kobayashi, Tatsuki Kuribayashi, Sho Yokoi, and Kentaro Inui. 2020. Attention is Not Only a Weight: Analyzing Transformers with Vector Norms. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 7057–7075, Online. Association for Computational Linguistics.

- Taku Kudo and John Richardson. 2018. SentencePiece: A simple and language independent subword tokenizer and detokenizer for Neural Text Processing. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, pages 66-71, Brussels, Belgium. Association for Computational Linguistics.
- Adhiguna Kuncoro, Miguel Ballesteros, Lingpeng Kong, Chris Dyer, Graham Neubig, and Noah A. Smith. 2017. What Do Recurrent Neural Network Grammars Learn About Syntax? In Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers, pages 1249-1258, Valencia, Spain. Association for Computational Linguistics.
- Alexandra Kuznetsova, Per B. Brockhoff, and Rune H. B. Christensen. 2017. ImerTest package: Tests in linear mixed effects models. Journal of Statistical Software, 82(13):1–26.
- Roger Levy. 2008. Expectation-based syntactic comprehension. Cognition, 106(3):1126-1177.
- Richard L. Lewis and Shravan Vasishth. 2005. An activation-based model of sentence processing as skilled memory retrieval. Cognitive Science, 29(3):375-419.
- David Marr. 1982. Vision: A Computational Investigation into the Human Representation and Processing of Visual Information. W. H. Freeman and Company, San Francisco.
- Danny Merkx and Stefan L. Frank. 2021. Human Sentence Processing: Recurrence or Attention? In Proceedings of the Workshop on Cognitive Modeling and Computational Linguistics, pages 12-22, Online. Association for Computational Linguistics.
- James A. Michaelov, Megan D. Bardolph, Seana Coulson, and Benjamin Bergen. 2021. Different kinds of cognitive plausibility: Why are transformers better than RNNs at predicting N400 amplitude? Proceedings of the Annual Meeting of the Cognitive Science Society, 43(43).
- Hiroshi Noji and Yohei Oseki. 2021. Effective Batching for Recurrent Neural Network Grammars. In Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021, pages 4340-4352, Online. Association for Computational Linguistics.
- Byung-Doh Oh and William Schuler. 2022. Entropyand Distance-Based Predictors From GPT-2 Attention Patterns Predict Reading Times Over and Above

GPT-2 Surprisal. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 9324–9334, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

- Byung-Doh Oh and William Schuler. 2024. Leading Whitespaces of Language Models' Subword Vocabulary Pose a Confound for Calculating Word Probabilities. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 3464-3472, Miami, Florida, USA. Association for Computational Linguistics.
- R Core Team. 2024. R: A Language and Environment for Statistical Computing. Vienna, Austria.
- Philip Resnik. 1992. Left-Corner Parsing and Psychological Plausibility. In COLING 1992 Volume 1: The 14th International Conference on Computational Linguistics.
- Soo Hyun Ryu and Richard Lewis. 2021. Accounting for Agreement Phenomena in Sentence Comprehension with Transformer Language Models: Effects of Similarity-based Interference on Surprisal and Attention. In Proceedings of the Workshop on Cognitive Modeling and Computational Linguistics, pages 61-71, Online. Association for Computational Linguistics.
- Soo Hyun Ryu and Richard L. Lewis. 2022. Using Transformer language model to integrate surprisal, entropy, and working memory retrieval accounts of sentence processing. In Proceedings of the 35th Annual Conference on Human Sentence Processing, Santa Cruz, CA, USA.
- Laurent Sartran, Samuel Barrett, Adhiguna Kuncoro, Miloš Stanojević, Phil Blunsom, and Chris Dyer. 2022. Transformer Grammars: Augmenting Transformer Language Models with Syntactic Inductive Biases at Scale. Transactions of the Association for Computational Linguistics, 10:1423–1439.
- Cory Shain, Idan Asher Blank, Marten van Schijndel, William Schuler, and Evelina Fedorenko. 2020. fMRI reveals language-specific predictive coding during naturalistic sentence comprehension. Neuropsychologia, 138:107307.
- Mark Steedman. 2000. The Syntactic Process. MIT Press, Cambridge, MA, USA.
- Mitchell Stern, Daniel Fried, and Dan Klein. 2017. Effective Inference for Generative Neural Parsing. In Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, pages 1695–1700, Copenhagen, Denmark. Association for Computational Linguistics.
- Yushi Sugimoto, Ryo Yoshida, Hyeonjeong Jeong, Masatoshi Koizumi, Jonathan R. Brennan, and Yohei Oseki. 2024. Localizing Syntactic Composition with Left-Corner Recurrent Neural Network Grammars. Neurobiology of Language, 5(1):201–224.

918

- 9 9
- 924 925
- 926
- 928
- 929 930
- 931
- 932 933
- 934 935
- 936
- 937 938
- 939 940
- 941 942

943 944

94 94

947 948

949

950 951

952

954

- 955
- 950
- 958

960

961

962

963

964

William Timkey and Tal Linzen. 2023. A Language Model with Limited Memory Capacity Captures Interference in Human Sentence Processing. In Findings of the Association for Computational Linguistics: EMNLP 2023, pages 8705–8720, Singapore. Association for Computational Linguistics.

- Julie A Van Dyke and Richard L Lewis. 2003. Distinguishing effects of structure and decay on attachment and repair: A cue-based parsing account of recovery from misanalyzed ambiguities. *Journal of Memory and Language*, 49(3):285–316.
- Julie Ann Van Dyke. 2002. *Retrieval Effects in Sentence Parsing and Interpretation*. University of Pittsburgh ETD, University of Pittsburgh.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In Advances in Neural Information Processing Systems, volume 30. Curran Associates, Inc.
- Ethan G. Wilcox, Jon Gauthier, Jennifer Hu, Peng Qian, and Roger P. Levy. 2020. On the Predictive Power of Neural Language Models for Human Real-TimeComprehension Behavior. In *Proceedings of the Annual Meeting of the Cognitive Science Society*, volume 42.
- Michael Wolfman, Donald Dunagan, Jonathan Brennan, and John Hale. 2024. Hierarchical syntactic structure in human-like language models. In *Proceedings of the Workshop on Cognitive Modeling and Computational Linguistics*, pages 72–80, Bangkok, Thailand. Association for Computational Linguistics.

A Hyperparameters

The hyperparameters are shown in Table 4. All model hyperparameters follow Sartran et al. (2022); Wolfman et al. (2024), while training hyperparameters were adjusted to accommodate the batch size suitable for our computational resources (NVIDIA A100, 40GB). The total computational cost required for all experiments was approximately 225 GPU hours.

B Correlations between predictors

The correlations between predictors in our statistical analysis are shown in Table 5. While the NAE from different LMs shows a very high correlation with each other, their predictive power for the self-paced reading times remains independent (see Section 4.2 and 5.1).¹⁶

C Part-of-speech tags

Table 5 presents the complete list of part-of-speech (POS) and symbol tags in the Natural Stories corpus. As reading times are annotated for each whitespace-delimited region, for data points containing symbol tags (e.g., NNP.), we used the stripped version (e.g., NNP) in our analysis. Additionally, we excluded from our analysis any data points containing multiple POS tags (e.g., NNP POS).

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1004

D Parallel parsing experiment

As a conceptual case study for the local ambiguity resolution in syntactic structures behind token sequences, we implemented TG's NAE calculation using 10 syntactic structures obtained through word-synchronous beam search (Stern et al., 2017) with Recurrent Neural Network Grammar (RNNG; Dyer et al., 2016; Kuncoro et al., 2017; Noji and Oseki, 2021).¹⁷¹⁸ NAE was computed individually for each syntactic structure and then aggregated as a weighted average:

$$\text{NAE}_{\text{TG}_{l,h,i}} \coloneqq \frac{\sum_{t \in \text{Beam}_i} p(t) \cdot \text{NAE}_{l,h,i}^t}{\sum_{t \in \text{Beam}_i} p(t)}, \quad (4)$$

where Beam represents the set of syntactic structures synchronized at the *i*-th word (|Beam| = 10).¹⁹

The analysis revealed patterns consistent with those observed when considering only the globally correct syntactic structure: both LMs' NAE demonstrated significant predictive power for reading times, with TG's NAE showing stronger contributions compared to Transformer's (Table 6).²⁰ The likelihood ratio test further confirmed independent contributions from both LMs (p < 0.001 for both comparisons: 'TG & Transformer'> 'Transformer' and 'TG & Transformer'> 'TG').

E Surprisal experiment

We analyzed each LM's surprisal contribution to reading time prediction using a baseline regression model that excluded both LMs' surprisal from Equation 3 but included their NAE (Table 7). While

²⁰_bs indicates beam search.

¹⁶_mcomp indicates -comp.

¹⁷https://github.com/aistairc/rnng-pytorch

¹⁸RNNG was trained on BLLIP-LG using default hyperparameters. For inference, action beam size and fast track were set to 100 and 1, respectively.

¹⁹stack_count was similarly calculated as the weighted average across syntactic structures in Beam.

Model architecture	Transformer-XL (Dai et al., 2019)
Vocabulary size	32,768
Model dimension	1,024
Feed-forward dimension	4,096
Number of layers	16
Number of heads	8
Segment length	256
Memory length	256
Optimizer	Adam ($\beta_1 = 0.9, \beta_2 = 0.999$) (Kingma and Ba, 2015)
Batch size	16
Number of training steps	400,000
Learning rate scheduler	Linear warm-up & cosine annealing
Number of warm-up steps	32,000
Initial learning rate	2.5×10^{-8}
Maximum learning rate	$3.75 imes 10^{-5}$
Final learning rate	7.5×10^{-8}
Dropout rate	0.1

Table 4: Model and training hyperparameters



Figure 5: Correlations between predictors in our statistical analysis

CC	Coordinating conjunction	PRP\$	Possessive pronoun
CD	Cardinal number	RB	Adverb
DT	Determiner	RBR	Adverb, comparative
EX	Existential there	RBS	Adverb, superlative
FW	Foreign word	RP	Particle
IN	Preposition or subordinating conjunction	TO	to
JJ	Adjective	UH	Interjection
JJR	Adjective, comparative	VB	Verb, base form
JJS	Adjective, superlative	VBD	Verb, past tense
MD	Modal	VBG	Verb, gerund or present participle
NN	Noun, singular or mass	VBN	Verb, past participle
NNS	Noun, plural	VBP	Verb, non-3rd person singular present
NNP	Proper noun, singuler	VBZ	Verb, 3rd person singular present
NNPS	Proper noun, plural	WDT	Wh-determiner
PDT	Predeterminer	WP	Wh-pronoun
POS	Possessive ending	WP\$	Possessive wh-pronoun
PRP	Personal pronoun	WRB	Wh-adverb
-LRB-	Left round bracket	,	Comma
-RRB-	Right round bracket		Period
"	Open double quotes	:	Colon
"	Closing double quotes		

Table 5: POS and symbol tags in the Natural Stories corpus

Model	$\Delta LogLik(\uparrow)$	Predictor	Effect size [ms]	<i>p</i> -value range	Significant seeds
TG	76.6 (±6.6)	tg_bs_nae tg_bs_nae_so	$\begin{array}{c} 2.53 \ (\pm \ 0.2) \\ 0.849 \ (\pm \ 0.1) \end{array}$	<0.001 <0.001	3/3 3/3
Transformer	28.0 (±6.9)	tf_nae tf_nae_so	$\begin{array}{c} 1.58 \ (\pm \ 0.2) \\ 0.457 \ (\pm \ 0.2) \end{array}$	< 0.001 0.006–0.13	3/3 1/3

Table 6: TG's and Transformer's NAE contribution to reading time prediction, where TG's NAE was calculated with multiple syntactic structures generated by word-synchronous beam search with RNNG

Model	$\Delta LogLik(\uparrow)$	Predictor	Effect size [ms]	<i>p</i> -value range	Significant seeds
TG	265 (±11)	tg_surp tg_surp_so	$\begin{array}{l} 4.63 \ (\pm \ 0.1) \\ 1.64 \ (\pm \ 0.1) \end{array}$	<0.001 <0.001	3/3 3/3
Transformer	299 (±30)	tf_surp tf_surp_so	$5.31 (\pm 0.2) \\ 1.68 (\pm 0.2)$	<0.001 <0.001	3/3 3/3

Table 7: TG's and Transformer's surprisal contribution to reading time prediction

both LMs' surprisal demonstrated significant pre-1005 dictive power for reading times, Transformer's sur-1006 prisal exhibited a stronger contribution compared to 1007 TG's. Additionally, our likelihood ratio test using 1008 the averaged surprisal revealed that the regression 1009 model incorporating both LMs' surprisal showed 1010 significantly higher predictive power compared to 1011 models with only one LM (p < 0.001 for both 1012 comparisons: 'TG & Transformer'>'Transformer' 1013 and 'TG & Transformer'>'TG'). These findings 1014 suggest that (i) unlike attention mechanisms, next-1015 word prediction based solely on token sequences 1016 more effectively captures dominant factors of hu-1017 man prediction processing, but (ii) similar to 1018 attention mechanisms, both types of next-word 1019 prediction-those based on token sequences alone and those leveraging both syntactic structures and 1021 token sequences-may coexist as models that cap-1022 ture distinct aspects of human predictive process-1023 1024 ing.

F Comparision between TG and TG_{-comp} with a weak baseline regression model

1025

1026

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1041

1042

1043

1044

1046

1047

1048

1049

1050

1051

1052

1054

In Section 5.1, we explored the advantage of treating closed phrases as single representations beyond explicit syntactic structure consideration. Our analysis incorporated Transformer's NAE in the baseline regression model to distinguish between two effects in TG_{-comp} : syntactic structure consideration and direct terminal token access.

To evaluate which model—TG or TG_{-comp} better captures more dominant factors in human sentence processing as a single model, we assessed their predictive power without Transformer's NAE in the baseline regression model (Table 8). The analysis revealed TG's superior predictive power (Δ LogLik=97.3) compared to TG_{-comp} ($\Delta LogLik=92.2$). Additionally, TG demonstrated both immediate and spillover effects, while TG_{-comp} primarily showed an immediate effect. These results highlight that TG, which explicitly treats closed phrases as single representations, outperforms TG_{-comp}, even when considering TG_{-comp}'s advantage in direct terminal token access. Consistent with findings in Section 5.1, likelihood tests confirmed TG's independent predictive power from TG_{-comp} ('TG & TG_{-comp} '>'T G_{-comp} ', p <0.001); TG_{-comp}, despite its lower overall predictive power, accounted for unique variance ('TG & TG_{-comp}'>'TG', p < 0.01).

G License

Table 9 summarizes the licenses of the data and
tools employed in this paper. All data and tools1056were used under their respective license terms.1057

Model	Δ LogLik (†)	Predictor	Effect size [ms]	<i>p</i> -value range	Significant seeds
TG	97.3 (±4.0)	tg_nae tg_nae_so	$\begin{array}{c} 2.93 \ (\pm \ 0.1) \\ 0.657 \ (\pm \ 0.1) \end{array}$	<0.001 <0.01	3/3 3/3
$TG_{-\mathrm{comp}}$	92.2 (±10)	tg_mcomp_nae tg_mcomp_nae_so	$\begin{array}{c} 2.88\ (\pm\ 0.2)\\ 0.416\ (\pm\ 0.1) \end{array}$	< 0.001 0.01–0.17	3/3 1/3

Table 8: TG's and TG $_{\rm comp}$'s NAE contribution to reading time prediction with Transformer's NAE excluded from the regression baseline model

Dataset/Tool	License
BLLIP (Charniak et al., 2000) Natural Stories corpus (Futrell et al., 2018)	BLLIP 1987–89 WSJ Corpus Release 1 CC BY-NC-SA 4.0
transformer_grammar (Sartran et al., 2022) rnng-pytorch (Noji and Oseki, 2021) SentencePiece (Kudo and Richardson, 2018) R (version 4.4.2) (R Core Team, 2024)	Apache 2.0 MIT License Apache 2.0 GNU GPL ≥ 2
Ime4 (version 1.1.34) (Bates et al., 2015)ImerTest (version 3.1.3) (Kuznetsova et al., 2017)	$\begin{array}{l} \text{GNU GPL} \geq 2 \\ \text{GNU GPL} \geq 2 \end{array}$

Table 9: Licenses of datasets and tools