Genre Matters: How Text Types Interact with Decoding Strategies and Lexical Predictors in Shaping Reading Behavior

Anonymous ACL submission

Abstract

The type of a text profoundly shapes reading behavior, yet little is known about how different text types interact with word-level features and the properties of machine-generated texts to influence how readers process language. In this study, we investigate how different text types affect eye movements during reading, how decoding strategies used to generate texts interact with text type, and how text types modulate the influence of word-level psycholinguistic fea-011 tures such as surprisal, word length, and lexical frequency. Leveraging EMTeC (Bolliger et al., 2025), the first eye-tracking corpus of LLM-generated texts across six text types and multiple decoding algorithms, we show that text type strongly modulates cognitive effort 018 during reading, that psycholinguistic effects 019 induced by word-level features vary systematically across genres, and that decoding strategies interact with text types to shape reading behavior. These findings offer insights into genre-specific cognitive processing and have implications for the human-centric design of AI-generated texts.

1 Introduction

027

042

The type or genre of a text influences the cognitive effort we expend at different stages of language processing (Blohm et al., 2022). A proxy for this cognitive load in language processing consists in the way we move our eyes during reading: not only do eye movements contain information about the properties and structure of the text being read, but they also provide insights into the cognitive mechanisms underlying language processing, as different words require a different amount of cognitive effort to be processed (Rayner, 1998; Rayner and Clifton, 2009). Given these qualities, eye movements have been leveraged to investigate readers' interactions with different text types, observing, for instance, that poetry leads to more regressions (Corcoran et al., 2023) or that fiction is read faster than nonfiction (Brysbaert, 2019). However, most of these studies have examined different genres in isolation and not directly pitted them against each other under the same experimental conditions, which would be crucial to make direct comparisons. 043

045

047

049

051

054

055

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

079

083

Moreover, while these studies do look at reading behavior in different text types, they do so in a coarse-grained manner by, for instance, considering overall reading time at the text level, thereby not accounting for word-level features which prompt reading patterns. These word-level features constitute psycholinguistic phenomena whose effects have long been established and include the word length effect — longer words take more time to read than shorter ones (Rayner, 2009; Hyönä and Olson, 1995; Just and Carpenter, 1980; Kliegl et al., 2004) -, the lexical frequency effect — frequent words are processed faster than infrequent ones (Forster and Chambers, 1973; Inhoff and Rayner, 1986) —, and the surprisal effect — high-surprisal words take longer to process than low-surprisal ones (Hale, 2001; Levy, 2008; Gruteke Klein et al., 2024; Xu et al., 2023b). That these effects exist in different text genres has been corroborated extensively (Pimentel et al., 2023; Frank and Aumeistere, 2024; Kuperman et al., 2024; Torres et al., 2021, *i.a.*) but mainly in isolation. Examining how word-level features play out across text types, however, can reveal interactions between these features and text type properties and contribute insights to cognitive science by showing that the influence of certain psycholinguistic effects might be genre-dependent, such as that a reader's sensitivity to predictability in processing is a function of text type.

Recently, a growing body of research has examined the relationship between textual outputs by language models (LMs) and humans and whether there is similarity in language production or language understanding processes between the two (Venkatraman et al., 2023; Giulianelli et al., 2023). An

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

integral aspect of these textual outputs by the LMs is the decoding strategy used to generate the text and its alignment with cognitive processing strate-086 gies. So far, only one study (Bolliger et al., 2024) investigated how humans read texts generated by large language models (LLMs), focusing on how different models and decoding algorithms affect 090 cognitive processing during reading and suggesting that decoding strategies can affect its readability. However, this line of work has not yet considered how these effects may interact with the type of text being generated. Investigating this interaction can highlight whether certain decoding methods are better suited, in terms of processing ease, for particular genres and can help ensure that AI systems generate texts in a way that aligns with our genre-specific processing strategies. The interplay 100 between decoding method and genre-specific prop-101 erties is thus an important but underexplored area. 102

This study investigates the effect of text type on reading behavior and its interaction with psycholinguistic phenomena as well as with neural decoding algorithms by tackling the following questions:

103

104

107

109

110

111

112

113

114

115

116

117

118

119

121

122

123

124

125

126

127

128

130

131

132

134

- RQ₁ Do different text types elicit different reading patterns, reflecting different cognitive demands during reading?
- RQ₂ Do well-established word-level predictors of reading behavior, such as surprisal, word length, and lexical frequency, interact with text type in shaping how people read?
- RQ₃ Do the neural decoding strategies used to generate texts of different text types interact with those text types in shaping reading behavior?

To this end, we leverage the Eye Movements on Machine-Generated **Texts Corpus** (Bolliger et al., 2025, EMTeC), the first dataset containing eyetracking data on LLM-generated texts across six different text types, generated using a variety of decoding algorithms. This dataset does not only allow for a direct comparison of reading behavior across different text types and how psycholinguistic effects vary between them, but also how they interact with decoding algorithms.

Our findings suggest that text type exerts a strong influence on cognitive effort during reading, that the magnitude of the psycholinguistic effects elicited by lexical features is modulated by text type, and that the decoding strategies used by language models interact with text types to shape the ease of processing generated texts.

2 Related Work

Text type or genre has long been recognized as a key factor in shaping reading behavior. Poetry, for example, induces longer fixations and more regressions due to its atypical syntax, ambiguity, and foregrounded language (Blohm et al., 2022; Corcoran et al., 2023), and readers' eye movements differ even when identical content is presented in poetic versus prosaic layout (Fechino et al., 2020). In contrast, narrative fiction elicits more linear reading patterns, attributed to its predictability (Graesser et al., 2003). Studies comparing fiction and nonfiction suggest that fiction is read more quickly, a difference largely driven by word length and lexical complexity (Brysbaert, 2019; Corcoran et al., 2023). While these studies demonstrate genrespecific reading patterns, they typically examine one genre at a time, under differing experimental conditions, thereby limiting comparability. Our work fills this gap by comparing six genres directly within a controlled, unified dataset.

In parallel, a large body of work has investigated psycholinguistic predictors of reading difficulty, such as surprisal (Hale, 2001; Levy, 2008; Gruteke Klein et al., 2024; Xu et al., 2023b; Shain et al., 2024), word length (Rayner, 1998, 2009; Hyönä and Olson, 1995; Just and Carpenter, 1980; Kliegl et al., 2004; Gerth and Festman, 2021; Kuperman et al., 2024), and lexical frequency (Forster and Chambers, 1973; Inhoff and Rayner, 1986; Chen and Ko, 2011; Torres et al., 2021). These effects have been consistently observed across a wide range of genres, including narrative (Luke and Christianson, 2016, 2018; Cop et al., 2017; Salicchi et al., 2023; Frank and Aumeistere, 2024), expository (Kennedy et al., 2003; Xu et al., 2023b; Goodkind and Bicknell, 2018), and scientific texts (Klein et al., 2025; Jakobi et al., 2025). Even stylistic deviations such as foregrounding in literary texts modulate these effects (Van den Hoven et al., 2016). Although these findings highlight the robustness of psycholinguistic predictors, few studies have investigated whether their magnitude or nature differs across text types. Our study addresses this by systematically analyzing interactions between genre and psycholinguistic effects within the same experimental setup.

Finally, recent research has begun examining how texts generated by large language models are processed by human readers. Bolliger et al. (2024) showed that decoding strategies, such as top-p sam-

pling or greedy decoding, can influence reading 186 behavior, although no single strategy consistently 187 outperformed others across measures or models in terms of processing ease. Other studies have explored the structure and information distribution of LLM outputs from the perspective of predictabil-191 ity or information density (Giulianelli et al., 2023; 192 Venkatraman et al., 2023), but these analyses were 193 conducted at the sentence level and did not incorpo-194 rate eye-tracking data or account for text type. To 195 date, no study has examined whether and how the impact of decoding strategies interacts with word-197 level features. Our study fills this gap by leveraging 198 EMTeC (Bolliger et al., 2025), which combines 199 multiple genres, multiple decoding strategies, and 200 human eye-tracking data.

Experiments¹ 3

3.1 Data

207

208

209

210

211

212

213

214

215

216

218

219

220

224

227

231

EMTeC We employ reading data from the Eye Movements on Machine-Generated Texts Corpus (Bolliger et al., 2025, EMTeC), an English eyetracking-while reading corpus whose stimuli were created with three different large language models (LLMs) — Phi-2 (Javaheripi et al., 2023), Mistral 7B Instruct (Jiang et al., 2023), and WizardLM (Xu et al., 2023a) — using five decoding algorithms - greedy search, beam search, ancestral sampling, top-k sampling (Fan et al., 2018), and top-p sampling (Holtzman et al., 2020). The generated stimuli belong to six different types of text categories: Non-fiction, where the models were prompted to either write a description or an argumentation; Fiction, where the LLMs were instructed to write a short story or a dialogue between two characters; Poetry, where the LLMs were prompted to write a poem; Summarization, where they were asked to summarize an input text; Article, where they ought to craft a news article out of an article synopsis; and Key-word text, where the LLMs had to create texts based on a range of input key words.

Reading Measures We consider the binary reading measures (RMs) fixated (Fix; whether or not a word was fixated) and first-pass regression (FPReg; whether or not a regression was initiated in the firstpass reading of the word) and the continuous RMs total fixation time (TFT; the sum of all fixations on a word), first-pass reading time (FPRT; the sum of the durations of all first-pass fixations on a word), 233

re-reading time (RRT; the sum of the durations of all fixations on a word that do not belong to the first pass), and regression path duration (RPD; the sum of all fixation durations starting from the first first-pass fixation on a word until fixating a word to the right of this word). While TFT and Fix indicate global language processing, FPRT and FPReg indicate early and RRT and RPD late stages of processing.

3.2 Predictors

Word-level features. We include word-level predictors, namely surprisal, lexical frequency, and word length, whose impact on eye movement behavior in reading is well-established and key to psycholinguistic theories of reading and, more broadly, language comprehension (Reichle et al., 2003; Engbert et al., 2005; Veldre et al., 2020; Rabe et al., 2024). Surprisal quantifies the predictability of a word. It is based on surprisal theory (Hale, 2001; Levy, 2008), which operationalizes the relationship between cognitive processing effort and word predictability and posits that the cognitive effort needed to process a word is a function of that word's predictability. More specifically, surprisal is the negative log-probability of a word conditioned on its preceding (linguistic and extralinguistic) context. This quantity is approximated by neural language models which only take the preceding linguistic context into account. As such, given a vocabulary Σ and an augmented vocabulary $\overline{\Sigma} = \Sigma \cup \{EOS\}$ that contains a special EOS (end-of-sentence) token, the surprisal s of a word $w \in \Sigma$ at position t is defined as

$$s(w_t) := -\log_2 p_\phi(w_t \mid \boldsymbol{w}_{< t}), \qquad (1)$$

where $p_{\phi}(\cdot \mid \boldsymbol{w}_{< t})$ is the language model's approximate distribution of the true distribution $p(\cdot \mid \boldsymbol{w}_{< t})$ over words $w \in \overline{\Sigma}$ in context $w_{\leq t}$.² In the following, surprisal is estimated with GPT-2 base (Radford et al., 2019), which has been shown to have the highest predictive power on reading times among LMs (Shain et al., 2024). As the reading measures are computed on the level of white-space separated words but LMs use tokenizers that separate words into sub-word tokens (Sennrich et al., 2016; Song et al., 2021), we aggregate surprisal to the word level by summing up the surprisal values of the individual sub-word tokens.³ The *lexical frequency*

267

268

269

270

271

272

273

274

275

276

277

278

279

234

235

236

237

238

239

240

241

242

243

244

245

246

²*I.e.*, surprisal is computed across sentence boundaries.

³For elaborations on the pooling of sub-word token surprisal values, refer to Appendix A.

281of a word is the Zipf frequency obtained from the282wordfreq library,4 which presents the frequency283of a word on a logarithmic scale5 and is the word's284base-10 logarithm of the number of times it ap-285pears in a billion words. Word length refers to the286number of characters of a word, including adjacent287punctuation.

Contrast Coding of Text Type and Decoding Strategy Both the factor text type, consisting of the levels non-fiction, fiction, poetry, summarization, article, and key-word text, as well as the factor decoding strategy, consisting of the levels beam search, ancestral sampling, top-k sampling, top*p* sampling, and greedy search, are sum-contrast coded. Sum-contrast coding compares the grand mean of the dependent variable — the reading measure — for each but one level of the factor to the grand mean across all levels. I.e., for a factor with k levels, it generates k - 1 contrast variables. The levels key-word text and greedy search serve as the reference level and are only implicitly represented in the grand mean intercept. The comparisons are factor level minus grand mean. The factor levels are coded as 1, the grand mean as -1. The contrast matrices are depicted in Appendix B.

3.3 Methods

291

292

295

296

297

301

305

306

307

308

311

312

313

314

315

316

317

322

324

For the analyses, we utilize linear mixed-effects models: linear regressions for continuous variables, and logistic regressions with a logistic linking function for binary variables. Let $f_{\theta} : v_i \mapsto y_{ij}$ be a linear mixed model parametrized by θ , mapping from the predictors v_i of word *i* to the log-transformed reading measure y_{ij} of word *i* read by subject *j*, following a log-normal distribution. The predictors v_i include surprisal s_i , the z-score standardized lexical Zipf frequency f_i and word length l_i , and the sum-contrast coded factors text type tt_i and decoding strategy dec_i . All models include a by-subject random intercept θ_{0j} .⁶ We fit all models using the R library 1me4 (Bates et al., 2015).

3.4 RQ₁: The Effect of Text Types

To examine whether the text type influences reading behavior overall, *i.e.*, across all decoding strategies, we fit a regression model f_{θ} defined as

$$\frac{f_{\theta}: y_{ij} \sim \theta_0 + \theta_{0j} + \theta_1 l_i + \theta_2 f_i + \theta_3 s_i + \theta_4 t t_i. (2)}{{}^4 \text{https://pypi.org/project/wordfreg/}}$$

⁵There exists a linear relationship between log-frequency and reading times.

Results Figure 1 depicts the effect estimates of the sum-contrast coded text types on the prediction of the different reading measures. The reading pattern elicited by the different text types is mostly consistent across the different RMs and the effects are mostly significant, even when controlling for the psycholinguistic covariates surprisal, word length, and lexical frequency. Poetry exhibits the strongest positive effects: readers spend more time overall reading words in poems; they have higher FPRTs and RRTs, and poetry induces more FPRegs as well as number of fixations on words. Fiction and non-fiction, on the other hand, show the strongest negative effects: they cause significantly fewer fixations and FPRegs and lower reading times at any stage of processing (TFT, FPRT, RRT). Summarization and articles are both close to average, although summarizations cause slightly more-than-average fixations and FPRegs, while articles cause slightly less.

326

327

328

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

347

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

369

370

371

372

373

374

3.5 RQ₂: The Interaction between Word-Level Features and Text Types

In order to investigate how the psycholinguistic predictors surprisal, word length, and lexical frequency interact with text type to influence reading behavior as measured in a variety of reading measures and to assess whether the strength of these linguistic effects changes depending on the text type, we fit a regression model f_{θ} defined as

$$\begin{aligned} f_{\theta} : y_{ij} \sim \theta_0 + \theta_{0j} + \theta_1 l_i + \theta_2 f_i + \theta_3 s_i + \theta_4 t t_i + \\ \theta_5 (l_i \times t t_i) + \beta_6 (f_i \times t t_i) + \beta_7 (s_i \times t t_i), \end{aligned}$$
(3)

where $(l_i \times tt_i), (f_i \times tt_i), \text{and}(s_i \times tt_i)$ are the interactions between the three psycholinguistic predictors and the text types.

Results The fixed effects of the psycholinguistic predictors serve as a sanity check: they are as expected and are plotted in Appendix C.

Figure 2 depicts the interaction effects between sum-contrast coded text types and the psycholinguistic predictors and reveals nuanced patterns. In summarization texts, surprisal effects are stronger than average for early binary measures (Fix and FPReg) but weaker for early and late reading times (FPRT and RPD), while lexical frequency effects were generally smaller. In poetry, surprisal exerted a smaller-than-average effect on FPRTs and TFTs but a greater-than-average effect on RPDs. Lexical frequency effects were amplified during FPRTs and RRTs in poetry, and word length exerted stronger effects on both fixation probability

⁶We do not include random effects for items, as the number of unique items is too low.



Figure 1: Effect estimates (mean and 95% CI) of sum-contrast coded text types on the prediction of different reading measures. Filled dots indicate that the effect is significantly different from the grand mean.

and reading durations. For non-fiction, surprisal had a stronger effect on fixation probability and FPRT, while lexical frequency and word length effects were weaker. Fiction texts amplified lexical frequency effects across almost all reading measures, with high-frequency words particularly facilitating faster reading, and exhibited reduced word length effects except for FPRegs. Finally, article texts showed stronger surprisal effects on TFTs and RRTs, stronger word length effects, and mixed frequency effects.

375

376

377 378

384

386

391

3.6 RQ₃: The Interaction Between Decoding Strategies and Text Types

In order to assess how the different decoding strategies used to generate the texts and the text types that the LLMs were prompted to generate interact in influencing human reading behavior, we fit a regression model f_{θ} defined as

$$f_{\boldsymbol{\theta}}: y_{ij} \sim \theta_0 + \theta_{0j} + \theta_1 l_i + \theta_2 f_i + \theta_3 s_i + \\ \theta_4 t t_i + \theta_5 dec_i + \theta_6 (t t_i \times dec_i),$$
(4)

where $(tt_i \times dec_i)$ is the interaction between between the sum-contrast coded factors text type and decoding strategy.

397 **Results** The fixed effects of the psycholinguistic398 predictors are plotted in Appendix D as a sanity

check and are as expected for the psycholinguistic predictors and the text types. The main effects of the decoding strategies are mostly not significantly different from the grand mean.

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

Figure 3 shows the interaction effects between text type and decoding strategy. For poetry, texts generated with ancestral sampling and top-k sampling exhibited shorter FPRTs, shorter RPDs, and lower TFTs compared to the grand mean, while texts generated with top-p decoding exhibited longer RPDs and higher TFTs. For fiction, beam search was associated with fewer fixations and reduced RRTs, whereas top-p decoding increased fixation probability and sampling increased RRTs. In non-fiction texts, top-p decoding was associated with fewer FPRegs, fewer fixations, and shorter TFTs, while top-k decoding was associated with longer TFTs. In summarization texts, top-k decoding was associated with fewer FPRegs, shorter RPDs, lower RRTs, and reduced TFTs, whereas beam search and sampling were associated with increased RRTs and TFTs. For articles, top-k decoding was associated with increased fixation probability, longer RRTs, and higher TFTs, while top-pdecoding was associated with shorter TFTs.



Figure 2: Interaction effects (mean and 95% CI) between text types and psycholinguistic predictors. A filled dot indicates that the interaction is significant.

4 Discussion

The experimental results presented in this study contribute to the understanding of how *text types* influence reading behavior and how they *interact* with an *LLM's decoding strategy* and wellestablished *psycholinguistic phenomena* such as a word's predicability. The findings for RQ₁ clearly demonstrate genre-driven divergences in reading behavior. Poetry emerged as the genre associated with the highest cognitive load across all stages of reading, which aligns with psycholinguistic theories that poetry's unconventional syntax and dense metaphoric context demand deeper interpretative processing and frequent re-analysis (Blohm et al., 2022; Corcoran et al., 2023; Fechino et al., 2020). Conversely, fiction and non-fiction texts were associated with significantly reduced cognitive demands, which suggests that narrative and expository prose align with readers' genre expectations and facilitate fluent reading (Graesser et al., 2003). These findings confirm that the properties of different genres profoundly shape real-time cognitive processing during reading. They also underscore that poetry remains cognitively unique among genres — a pattern that persisted even though the stimuli were machine-generated, highlighting the robustness of genre-specific processing strategies. 438

439

440

441

442

443

444

445

446

447

448

449

450

451

The genre-specificity in reading behavior is fur-

435

436

437



Figure 3: Interaction effects (mean and 95% CI) between sum-contrast coded decoding strategy and text type. A filled dot indicates that the interaction is statistically significant, meaning that the effect of a combination of text type and decoding strategy significantly differs from the effect predicted by the additive fixed effects alone.

ther corroborated and expanded upon in the results of the experiment answering RQ₂. In poetry, surprisal had a weaker-than-average effect on early reading measures (FPRT), but a stronger-thanaverage effect on regression paths. This implies that readers tolerate local unpredictability in poetry during initial reading, but experience delayed integration difficulties requiring re-reading and reevaluation. Fictional texts amplified the influence of lexical frequency across nearly all reading measures: high-frequency words in fiction were read especially quickly. This suggests that in familiar narrative structures, readers rely more heavily on

452

453

454

455

456

457

458

459

460

461

462

463

464

lexical familiarity to facilitate fluent reading. Word length effects were diminished, except for first-pass regressions, indicating that in fiction, processing difficulties are less driven by orthographic length and more by broader discourse-level factors. In non-fiction, surprisal effects on fixation probability and FPRTs were heightened, while lexical frequency and word length effects were weaker: readers seem to engage more heavily with predictive mechanisms during informational text reading, possibly due to the structured, factual nature of the content. These findings underline that while classic psycholinguistic predictors like surprisal, lexical 465

466

467

468

469

470

471

472

473

474

475

476

486

487

488

489

490

491

492

493

494

495

496

497

498

500

501

502

503

504

505

506

508

510

511 512

513

514

516

517

518

519

521

523

525

529

478

frequency, and word length remain robust across genres, the magnitude and timing of their effects vary systematically with text type. Readers dynamically adapt their cognitive strategies depending on genre-specific expectations and structures.

We further found that while the main effects of the decoding strategies utilized to generate the texts were minimal, their interactions with genre revealed meaningful patterns. In poetry, texts generated with sampling-based strategies were easier to process — yielding shorter FPRTs, shorter RPDs, and lower TFTs - compared to those generated with top-p sampling, which paradoxically increased cognitive effort. This suggests that moderate stochasticity benefits poetry by fostering the unpredictability and variability that readers expect, whereas the specific distribution of probabilities under top-p sampling may have introduced irregularities detrimental to coherent interpretation. In fiction, deterministic decoding via beam search facilitated the reading experience, reducing fixation probability and re-reading times, whereas stochastic decoding strategies (sampling, top-p) introduced mild disruptions. This aligns with the expectation that narratives benefit from high predictability and coherence. In non-fiction, moderate randomness introduced by top-p decoding surprisingly facilitated reading — reducing regressions, fixations, and TFTs — while top-k decoding complicated it. This finding suggests that informational texts may benefit from slight variability, which might enhance engagement without compromising clarity. In summarization texts, top-k decoding led to the easiest reading (fewer regressions, shorter reading times), while both beam search and sampling complicated processing. This is intriguing because one might expect beam search to yield clear, coherent summaries — highlighting that stochastic decoding may sometimes better balance informativeness and readability. For articles, top-k decoding increased cognitive load, while top-p decoding decreased it, again emphasizing that subtle differences in decoding randomness can have substantial cognitive effects depending on genre.

In sum, these results demonstrate that no single decoding strategy universally optimizes readability. Rather, the ideal decoding method is crucially dependent on the genre and its associated cognitive demands as well as genre-specific expectations. This has direct implications for the design of human-centric LLM applications: depending on the desired use case or target population, generation systems may adapt decoding strategies to optimize user comprehension by facilitating reading ease, thereby matching the desired properties of different text types. 530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

Overall, our findings have important implications for both cognitive science and AI research. From a cognitive perspective, the study reinforces the view that genre deeply shapes cognitive processing strategies during reading. Not only does it affect the baseline ease or difficulty of reading, but it also modulates the impact of core psycholinguistic variables like surprisal, lexical frequency, and word length. These results imply that cognitive models of reading must account for genre as a systematic source of variance, not merely as a surface-level property. From an AI and NLP perspective, our results highlight that how a text is generated matters just as much as what genre it is intended to emulate. Different decoding strategies differentially align with text types in terms of ease of processing, affecting the cognitive accessibility of LLM-generated texts. As LLMs increasingly generate content for educational, journalistic, and entertainment purposes, understanding and optimizing for genre-appropriate readability will be crucial. Finally, studying AI-generated texts provides a new lens through which to test cognitive theories: by controlling genre and text structure via generation parameters, we can probe the flexibility and robustness of human reading strategies in a way that complements traditional studies on human-written texts.

5 Conclusion

This study shows that text type significantly shapes reading behavior, modulating not only overall cognitive demands but also the strength and manifestation of core psycholinguistic effects. Genres like poetry induce higher effort, while fiction and nonfiction support easier processing. We further find that decoding strategies interact with genre in nontrivial ways, indicating that optimizing readability in machine-generated texts requires genre-sensitive approaches. These results highlight the need for adaptive generation systems that align with genrespecific cognitive norms.

575 Limitations

Several limitations must be acknowledged. First, 576 while EMTeC provides a unique opportunity to 577 study eye movements across machine-generated 578 texts of different types, it does not include human-579 written baselines, which limits direct comparisons between human and machine text processing. Sec-581 ond, the texts were generated using only three 582 LLMs and five decoding strategies, which may not 583 capture the full diversity of possible outputs or de-584 coding configurations. Third, the study focuses on adult readers and English texts; results may not 586 generalize to different age groups, languages, or literacy backgrounds. Finally, while we account for core psycholinguistic predictors, other linguistic variables such as syntactic complexity or discourse 590 coherence were not directly controlled and could 591 influence reading behavior.

References

594

598

600

604

610

611

614

617

618

619

621

626

- Douglas Bates, Martin Mächler, Ben Bolker, and Steve Walker. 2015. Fitting linear mixed-effects models using lme4. *Journal of Statistical Software*, 67(1):1– 48.
- Stefan Blohm, Stefano Versace, Sanja Methner, Valentin Wagner, Matthias Schlesewsky, and Winfried Menninghaus and. 2022. Reading poetry and prose: Eye movements and acoustic evidence. *Discourse Processes*, 59(3):159–183.
- Lena S. Bolliger, Patrick Haller, Isabelle C. R. Cretton, David R. Reich, Tannon Kew, and Lena A. Jäger. 2025. EMTeC: A corpus of eye movements on machine-generated texts. *Behavior Research Methods*. In press. https://github.com/DiLi-Lab/ EMTeC/.
- Lena Sophia Bolliger, Patrick Haller, and Lena Ann Jäger. 2024. On the alignment of LM language generation and human language comprehension. In *Proceedings of the 7th BlackboxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP*, pages 217–231, Miami, Florida, US. Association for Computational Linguistics.
- Marc Brysbaert. 2019. How many words do we read per minute? a review and meta-analysis of reading rate. *Journal of Memory and Language*, 109:104047.
- Minglei Chen and Hwawei Ko. 2011. Exploring the eyemovement patterns as chinese children read texts: a developmental perspective. *Journal of Research in Reading*, 34(2):232–246.
- Uschi Cop, Nicolas Dirix, Denis Drieghe, and Wouter Duyck. 2017. Presenting geco: An eyetracking corpus of monolingual and bilingual sentence reading. *Behavior research methods*, 49:602–615.

Rhiannon Corcoran, Christophe de Bezenac, and Philip Davis. 2023. 'looking before and after': Can simple eye tracking patterns distinguish poetic from prosaic texts? *Frontiers in Psychology*, 14:1066303. 627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

- Ralf Engbert, Antje Nuthmann, Eike M. Richter, and Reinhold Kliegl. 2005. SWIFT: A dynamical model of saccade generation during reading. *Psychological Review*, 112(4):777.
- Angela Fan, Mike Lewis, and Yann Dauphin. 2018. Hierarchical neural story generation. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 889–898, Melbourne, Australia. Association for Computational Linguistics.
- Marion Fechino, Arthur M Jacobs, and Jana Lüdtke. 2020. Following in jakobson and lévi-strauss' footsteps: A neurocognitive poetics investigation of eye movements during the reading of baudelaire's 'les chats'. *Journal of Eye Movement Research*, 13(3):10– 16910.
- Kenneth I Forster and Susan M Chambers. 1973. Lexical access and naming time. *Journal of verbal learning and verbal behavior*, 12(6):627–635.
- Stefan L. Frank and Anna Aumeistere. 2024. An eyetracking-with-eeg coregistration corpus of narrative sentences. *Language Resources and Evaluation*, 58(2):641–657.
- Sabrina Gerth and Julia Festman. 2021. Reading development, word length and frequency effects: An eye-tracking study with slow and fast readers. *Frontiers in Communication*, 6:743113.
- Mario Giulianelli, Sarenne Wallbridge, and Raquel Fernández. 2023. Information value: Measuring utterance predictability as distance from plausible alternatives. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5633–5653, Singapore. Association for Computational Linguistics.
- Adam Goodkind and Klinton Bicknell. 2018. Predictive power of word surprisal for reading times is a linear function of language model quality. In *Proceedings of the 8th Workshop on Cognitive Modeling and Computational Linguistics (CMCL 2018)*, pages 10–18, Salt Lake City, Utah. Association for Computational Linguistics.
- Arthur C Graesser, Danielle S McNamara, and Max M Louwerse. 2003. What do readers need to learn in order to process coherence relations in narrative and expository text. *Rethinking Reading Comprehension*, 82:98.
- Keren Gruteke Klein, Yoav Meiri, Omer Shubi, and Yevgeni Berzak. 2024. The effect of surprisal on reading times in information seeking and repeated reading. In *Proceedings of the 28th Conference on Computational Natural Language Learning*, pages 219–230, Miami, FL, USA. Association for Computational Linguistics.

2024. Word length and frequency effects on text cholinguistic model. In Second meeting of the North 739 American Chapter of the Association for Computareading are highly similar in 12 alphabetic languages. 740 tional Linguistics. Journal of Memory and Language, 135:104497. 741 Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Roger Levy. 2008. Expectation-based syntactic compre-742 Yejin Choi. 2020. The curious case of neural text hension. Cognition, 106(3):1126–1177. degeneration. Preprint, arXiv:1904.09751. Steven G Luke and Kiel Christianson. 2016. Limits on 744 Jukka Hyönä and Richard K Olson. 1995. Eye fixlexical prediction during reading. Cognitive psychol-745 ation patterns among dyslexic and normal readers: ogy, 88:22-60. 746 effects of word length and word frequency. Journal of Experimental Psychology: Learning, Memory, and Steven G Luke and Kiel Christianson. 2018. The Provo 747 Cognition, 21(6):1430. Corpus: A large eye-tracking corpus with predictability norms. Behavior Research Methods, 50:826–833. Albrecht Werner Inhoff and Keith Rayner. 1986. Tiago Pimentel, Clara Meister, Ethan G Wilcox, Roger P Parafoveal word processing during eye fixations in 750 reading: Effects of word frequency. Perception & Levy, and Ryan Cotterell. 2023. On the effect of psychophysics, 40(6):431-439. anticipation on reading times. Transactions of the 752 Association for Computational Linguistics, 11:1624-Deborah N Jakobi, Thomas Kern, David R Reich, 1642. 754 Patrick Haller, and Lena A Jäger. 2025. PoTeC: Maximilian M. Rabe, Dario Paape, Daniela Mertzen, A German naturalistic eye-tracking-while-reading 755 corpus. Behavior Research Methods. In press. Shravan Vasishth, and Ralf Engbert. 2024. Seam: 756 github.com/DiLi-Lab/PoTeC. An integrated activation-coupled model of sentence 757 processing and eye movements in reading. Journal 758 Mojan Javaheripi, Sébastien Bubeck, Marah Abdin, of Memory and Language, 135:104496. 759 Jyoti Aneja, Sebastien Bubeck, Caio César Teodoro Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Mendes, Weizhu Chen, Allie Del Giorno, Ronen 760 Eldan, Sivakanth Gopi, Suriya Gunasekar, Mojan Dario Amodei, Ilya Sutskever, and 1 others. 2019. 761 Javaheripi, Piero Kauffmann, Yin Tat Lee, Yuanzhi Language models are unsupervised multitask learn-762 Li, Anh Nguyen, Gustavo de Rosa, Olli Saarikivi, ers. 763 Adil Salim, and 9 others. 2023. Phi-2: The surprising power of small language models. https: Keith Rayner. 1998. Eye movements in reading and 764 //www.microsoft.com/en-us/research/blog/ information processing: 20 years of research. Psy-765 phi-2-the-surprising-power-of-small-language-modelagical Bulletin, 124(3):372. 766 Albert Q. Jiang, Alexandre Sablayrolles, Arthur Men-Keith Rayner. 2009. The 35th sir frederick bartlett lec-767 sch, Chris Bamford, Devendra Singh Chaplot, Diego ture: Eye movements and attention in reading, scene 768 de las Casas, Florian Bressand, Gianna Lengyel, Guilperception, and visual search. Quarterly Journal of 769 laume Lample, Lucile Saulnier, Lélio Renard Lavaud, Experimental Psychology, 62(8):1457–1506. 770 Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, Keith Rayner and Charles Clifton. 2009. Language pro-771 and William El Sayed. 2023. Mistral 7B. Preprint, cessing in reading and speech perception is fast and arXiv:2310.06825. incremental: Implications for event-related potential 773 research. Biological Psychology, 80(1):4-9. Before 774 Marcel A Just and Patricia A Carpenter. 1980. A thethe N400: Early Latency Language ERPs. 775 ory of reading: from eye fixations to comprehension. Psychological Review, 87(4):329. Erik D. Reichle, Keith Rayner, and Alexander Pollatsek. 776 2003. The E-Z reader model of eye-movement con-Alan Kennedy, Robin Hill, and Joël Pynte. 2003. The trol in reading: Comparisons to other models. The 778 Behavioral and Brain Sciences, 26:445-526. dundee corpus. In Proceedings of the 12th European 779 Conference on Eye Movement. Lavinia Salicchi, Emmanuele Chersoni, and Alessandro 780 Keren Gruteke Klein, Shachar Frenkel, Omer Shubi, Lenci. 2023. A study on surprisal and semantic relat-781 and Yevgeni Berzak. 2025. Surprisal takes it all: Eye edness for eye-tracking data prediction. Frontiers in 782 tracking based cognitive evaluation of text readability Psychology, 14:1112365. 783 measures. arXiv preprint arXiv:2502.11150. Rico Sennrich, Barry Haddow, and Alexandra Birch. 784 Reinhold Kliegl, Ellen Grabner, Martin Rolfs, and Ralf 2016. Neural machine translation of rare words with 785 subword units. In Proceedings of the 54th Annual Engbert. 2004. Length, frequency, and predictabil-786 Meeting of the Association for Computational Linity effects of words on eye movements in reading. 787 European Journal of Cognitive Psychology, 16(1guistics, pages 1715–1725, Berlin, Germany. Associ-2):262-284. ation for Computational Linguistics. 789 10

Victor Kuperman, Sascha Schroeder, and Daniil Gnetov.

738

John Hale. 2001. A probabilistic earley parser as a psy-

707

710

711

712

713

714

715

716

717

718

721 722

723

724

726

727

728

729

730

731

732

733

734

Cory Shain, Clara Meister, Tiago Pimentel, Ryan Cotterell, and Roger Levy. 2024. Large-scale evidence for logarithmic effects of word predictability on reading time. *Proceedings of the National Academy of Sciences*, 121(10):e2307876121.

790

791

793

795

796

804

805 806

807

810

811

812 813

814

815

816

817

818 819

820

822

824

826

- Xinying Song, Alex Salcianu, Yang Song, Dave Dopson, and Denny Zhou. 2021. Fast WordPiece tokenization. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 2089–2103, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Débora Torres, Wagner R Sena, Humberto A Carmona, André A Moreira, Hernán A Makse, and José S Andrade Jr. 2021. Eye-tracking as a proxy for coherence and complexity of texts. *PloS one*, 16(12):e0260236.
- Emiel Van den Hoven, Franziska Hartung, Michael Burke, and Roel M Willems. 2016. Individual differences in sensitivity to style during literary reading: Insights from eye-tracking. *Collabra*, 2(1):25.
- Aaron Veldre, Lili Yu, Sally Andrews, and Erik D Reichle. 2020. Towards a complete model of reading: Simulating lexical decision, word naming, and sentence reading with Über-Reader. In *Proceedings of* the 42nd annual conference of the cognitive science society. Cognitive Science Society.
- Saranya Venkatraman, He He, and David Reitter. 2023. How do decoding algorithms distribute information in dialogue responses? In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 953–962, Dubrovnik, Croatia. Association for Computational Linguistics.
- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023a. WizardLM: Empowering large language models to follow complex instructions. *Preprint*, arXiv:2304.12244.
- Weijie Xu, Jason Chon, Tianran Liu, and Richard Futrell. 2023b. The linearity of the effect of surprisal on reading times across languages. In *Findings of the Association for Computational Linguistics: EMNLP* 2023, pages 15711–15721, Singapore. Association for Computational Linguistics.

Appendix for Genre Matters: How Text Types Interact with Decoding Strategies and Lexical Predictors in Shaping Reading Behavior

A Pooling of Surprisal

832

833

837

841

843

846

The word-level surprisal values utilized in this study are already contained within EMTeC (Bolliger et al., 2025), where surprisal has been estimated with a range of language models, including GPT-2 *base* (Radford et al., 2019). Since language models employ tokenizers that separate words into sub-word tokens (Sennrich et al., 2016; Song et al., 2021) but the reading measure data is on word-level, the surprisal values must be pooled to word level.

Since the sum of two logarithms is equal to the logarithm of the product of their arguments, *i.e.*, $\log a + \log b = \log [a \cdot b]$, surprisal is pooled to word-level as follows: given k sub-word tokens $w_n, w_{n+1}, \ldots, w_{n+k}$ that belong to the same word token w, the word-level surprisal of w is computed as

$$s(w_n, w_{n+1}, \dots, w_{n+k}) = -\log p(w_n, w_{n+1}, \dots, w_{n+k} \mid \boldsymbol{w}_{< n})$$

= $-\log \left[p(w_n \mid \boldsymbol{w}_{< n}) \cdot p(w_{n+1} \mid \boldsymbol{w}_{< n+1}) \cdot \dots \cdot p(w_{n+k} \mid \boldsymbol{v}_{< n+k}) \right]$
= $-\log p(w_n \mid \boldsymbol{w}_{< n}) - \log p(w_{n+1} \mid \boldsymbol{w}_{< n+1})$
 $- \dots - \log p(w_{n+k} \mid \boldsymbol{v}_{< n+k}).$

This shows that summing up sub-word level surprisal values is equivalent to computing the surprisal of the joint probability distribution of the sub-word tokens.

B Contrast Matrices

Below the contrast matrices used in the experiments are depicted. Table 1 shows the sum-contrast coded factor text type, and Table 2 shows the sum-contrast coded factor decoding strategy.

Factor Level	Article synopsis vs grand-mean	Summarization vs grand-mean	Non-fiction vs grand-mean	Fiction vs grand-mean	Poetry vs grand-mean
Article synopsis	1	0	0	0	0
Summarization	0	1	0	0	0
Non-fiction	0	0	1	0	0
Fiction	0	0	0	1	0
Poetry	0	0	0	0	1
Key-word text	-1	-1	-1	-1	-1

Table	1:	Sum	contrast	matrix	for the	factor t	text	type.

$C RQ_2$

The fixed effects of the psycholinguistic predictors are plotted in Figure 4 as a sanity check. Across all predictors and reading measures, the direction of the effect is as expected: the effects of lexical frequency are significantly negative (high-frequency words cause lower reading times), the effects of surprisal are significantly positive (high-surprisal words cause longer reading times), as are the effects of word length (longer words cause longer reading times). The only exception is surprisal as a predictor for the binary variable first-pass regression.

Factor Level	Beam search vs grand mean	Sampling vs grand mean	Top- <i>k</i> vs grand mean	Top- <i>p</i> vs grand mean
Doom coorch	1	0	0	0
Dealli search	1	0	0	0
Sampling	0	1	0	0
Top-k	0	0	1	0
Top-p	0	0	0	1
Greedy search	-1	-1	-1	-1

Table 2: Sum contrast matrix for the factor decoding strategy.

D RQ₃

Figure 5 depicts the estimates of the fixed effects of the psycholinguistic predictors and the sum-contrast coded predictors text type and decoding strategy. This serves as a sanity check to corroborate that the effects of the psycholinguistic predictors are as would be expected: the effects of lexical frequency are negative (frequent words cause lower reading times), the effects of word length are positive (longer words cause longer reading times), as are the effects of surprisal (high-surprisal words cause longer reading times). Moreover, the main effects of the text types exhibit the same pattern as in the results for RQ_1 (see § 3.4). The main effects of the different decoding strategies, on the other hand, are mostly not significantly different from the grand mean with the exception of beam search, indicating that texts generated with this decoding strategy elicit longer-than-average re-reading time.



Figure 4: Estimates (mean and 95% CI) of the fixed effects of the psycholinguistic predictors lexical frequency, word length, and surprisal. All effects are significantly different from zero.



Figure 5: Estimates (mean and 95% CI) of the fixed effects of the psycholinguistic predictors lexical frequency, word length, and surprisal, and of the sum-contrast coded factors text type and decoding strategies. A filled dot indicates that the effect is significantly different from zero for the continuous psycholinguistic predictors, or significantly different from the grand mean for the sum-contrast coded text type and decoding strategy.