Temperature-scaling surprisal estimates improve fit to human reading times – but does it do so for the "right reasons"?

Anonymous ACL submission

Abstract

A wide body of evidence shows that human language processing difficulty is predicted by the information-theoretic measure surprisal, a word's negative log probability in context. 005 However, it is still unclear how to best estimate these probabilities needed for predicting human processing difficulty - while a long-standing 007 belief held that models with lower perplexity would provide more accurate estimates of word predictability, and therefore lead to better reading time predictions, recent work has shown 011 that for very large models, psycholinguistic predictive power decreases. One reason could be that language models might be more confident of their predictions than humans, because they have had exposure to several magnitudes more data. In this paper, we test what effect 017 temperature-scaling of large language model (LLM) predictions has on surprisal estimates 019 and their predictive power of reading times of English texts. Firstly, we show that calibration of large language models typically improves with model size, i.e. poorer calibration cannot account for poorer fit to reading times. Secondly, we find that temperature-scaling probabilities lead to a systematically better fit to reading times (up to 89% improvement in delta 027 log likelihood), across several reading time corpora. Finally, we show that this improvement in fit is chiefly driven by words that are composed of multiple subword tokens.¹

1 Introduction

032

041

In psycholinguistics, a key finding is that words with higher surprisal (= negative log probability of the word in context) require more time for processing (Hale, 2001; Levy, 2008). Numerous studies provided experimental evidence supporting this theory, demonstrating that surprisal is a powerful predictive measure of processing complexity (e.g., Demberg and Keller, 2008; Wilcox et al., 2020, 2023; Shain et al., 2022), and that the relationship between surprisal and reading times (RTs) indeed seems to be linear (Smith and Levy, 2013; Wilcox et al., 2020; Shain et al., 2022).

042

043

044

047

049

051

054

055

058

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

079

081

However, prior work implicitly made the assumption that human predictability estimates would be similar to the actual probability of a word occurring in a given context, and that therefore, surprisal values estimated from models that achieve lower perplexities should also approximate human processing difficulty better (Goodkind and Bicknell, 2018; Merkx and Frank, 2021).

Recent research has however found that this is not true - surprisal values from very large LLMs provide in fact a very poor fit to reading times. Oh and Schuler (2023b) hypothesize that this might be due to LLMs being "too confident" in their estimates of rare named entities compared to humans, thanks to their manifold larger exposure to data and greater memory capacity compared to humans. Furthermore, work on NLP applications like question answering has reported that probability estimates from pretrained language models are often overconfident, i.e. they are higher than the ground truth probability (Si et al., 2022; Kumar, 2022). These findings hence beg the question whether current LLMs are well-calibrated with respect to "objective" word occurrence probabilities. Relatedly, we ask whether LLM probability estimates are overconfident compared to human estimates (as observed in reading times).

One approach to address calibration problems is to use *temperature scaling*, as done e.g., in vision tasks (Guo et al., 2017; Hendrycks et al., 2019). Temperature-scaling with a temperature T > 1has the effect that the probability distribution is flattened such that it becomes more similar to a uniform distribution. Temperature-scaling hence incorporates uncertainty into the probability estimates from LLMs.

We note that the idea to work with flattened distributions instead of the original probability dis-

¹Code in this paper will be released upon paper acceptance.

tributions from LLMs is also related to contextual Rényi Entropy as discussed by Pimentel et al. (2023), as well as the super/sub-linear surprisal effect by Shain et al. (2022). However, rather than merely adjust the power of surprisal in super/sublogarithmic patterns or the power of probability in Rényi entropy, our work represents a distinct branch of study (i.e., probability calibration) in machine learning: shaping the probability distribution itself through shaping the logits before softmax. We also discuss the motivation for why a slightly flattened distribution may be more suitable, and whether this change in distribution is applied when calculating surprisal vs. when calculating entropy.

084

091

097

099

100

101 102

103

104

105

106

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

127

128

129

130

131

132

Our experimental results show that scaling probabilities can largely improve the fit to reading times in all 12 settings (3 corpora \times 4 neural LMs). Our contributions are summarized as follows: (1) We propose temperature-scaled surprisal, where surprisal is calculated from temperaturescaled probabilities. (2) We demonstrate that temperature-scaling with temperature T \approx 2.5 improves predictability of human reading times of English texts compared to T=1. (3) We identify linguistic phenomena that correlate with the benefit of temperature-scaled surprisal by analyzing residual errors from regression models.

2 Predictive Power for Reading Times

In psycholinguistics, RTs on a word are believed to correlate with its processing difficulty. RTs can be gathered using different paradigms, including eyetracking while reading text on a screen (Rayner, 1998), self-paced reading (Aaronson and Scarborough, 1976; Mitchell and Green, 1978) and the Maze task (Forster et al., 2009).

The most common procedure for predicting words' RT is first to select a set of predictor variables thought to impact RTs $\mathbf{v} = [v^{(1)}, ..., v^{(d)}]^{\top} \in \mathbb{R}^d$, which include, e.g., the length of a word w_t , $|w_t|$, the frequency of a word freq (w_t) . Let $f_{\phi} : \mathbb{R}^d \to \mathbb{R}$ be a regression model parametrized by ϕ used to fit these predictors for the prediction of human RTs rt: $rt(w_t | \mathbf{w}_{< t}) \sim f_{\phi}(\mathbf{v})$, given the previous context $\mathbf{w}_{< t}$. The performance of f_{ϕ} is quantified by its log-likelihood, with a higher loglikelihood indicating a better psychometric predictive power for human RTs (Frank and Bod, 2011; Fossum and Levy, 2012).

Besides the word length $|w_t|$ and word frequency freq (w_t) , a word's surprisal (i.e., its negative log-

probability in context) (Hale, 2001; Levy, 2008) has been shown to be predictive of RTs (Demberg and Keller, 2008; Goodkind and Bicknell, 2018; Wilcox et al., 2020; Shain et al., 2022).

133

134

135

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

3 Methods

In this section, we delve into key aspects of information-theoretic measures in language comprehension. We start with surprisal, a method connecting processing difficulty to word predictability. As word predictability is empirically estimated by LLMs, we introduce the notion of calibration errors, metrics quantifying how good the estimation of word predictability is. Further, we lay out temperature-scaled surprisal, and the relation between varying temperature vs. varying α in contextual Rényi entropy.

3.1 Surprisal

Starting from Shannon (1948), the information conveyed by a word w_t has been quantified as the negative log probability of the word w_t given its previous context $w_{<t}$. In Surprisal Theory (Hale, 2001; Levy, 2008), this quantity is called surprisal $s(w_t)$ and proposed to be predictive of the word's processing difficulty, typically quantified as its RT. Surprisal values are typically estimated from language models $\hat{p}(w_t|w_{<t})$.

$$s(w_t) = -\log_2 p(w_t | \boldsymbol{w}_{< t}), \qquad (1)$$

3.2 Calibration error

Definitions Let $\mathcal{D} = \{(x_i, y_i)\}_i^N$ be a data set where $x_i \in \mathcal{X}$ is an sample (i.e., context) and $y_i \in \mathcal{K} = [K]$ is a category label. Let g_{θ} and $\hat{z}_i = g_{\theta}(x_i)$ denote a language model parametrized by θ and the output logit vector of sample *i*, respectively. The predicted class label \hat{y}_i for sample *i* is given by $\hat{y}_i = \arg \max_{k \in \mathcal{K}} g(x_i)_k$ and confidence for sample *i* is given by $\hat{p}_i = \max_{k \in \mathcal{K}} g(x_i)_k$. A model is perfectly calibrated when the confidence \hat{p} is equal to the frequency of correctness, i.e., $\mathbb{P}(\hat{y}_i = y_i | \hat{p}_i = p) = p$ holding for all $p \in [0, 1]$ and any sample *i*. Any difference between the left and right sides of the above equation indicates there exists a *calibration error*.

Expected calibration error (ECE) (Guo et al., 2017) ECE is the most popular calibration metric, which empirically approximates the calibration error by discretizing the probability interval into a fixed number of bins (B_m with $m \in \{1, 2, ..., M\}$),

181

182

- 183

- 186
- 188

190

192

193

194

195

196

197

198

199

201

204

209

215

 $= \frac{1}{NK} \sum_{k=1}^{K} \sum_{m=1}^{M} |\sum_{i \in B_m} \hat{p}_{i,k} - \sum_{i \in B_m} \mathbb{1}[k = y_i]|,$

and measures the gaps of averaged confidence and

 $\text{ECE} = \frac{1}{N} \sum_{m=1}^{M} |\sum_{i \in B_m} \hat{p}_i - \sum_{i \in B_m} \mathbb{1}[\hat{y}_i = y_i]|, \ (2)$

where $\mathbb{1}$ is the indicator function. However, it does

not necessarily measure the actual-word probabil-

ity, which is the probability required for calculating

surprisal in Eq. 1. It focuses only on the top-label

Classwise-ECE (CECE) (Kumar et al., 2019;

Kull et al., 2019) In comparison, CECE mea-

sures probabilities of all classes. For each bin and

every class k, it assesses the difference between the

average confidence of samples for class k and the

actual proportion of class k. If assuming all classes

averaged accuracy in each bin B_m .

probability for a given sample.

weigh equally, we have:

CECE

where $\hat{p}_{i,k}$ is the predicted probability of sample ifor class k.

Human-likeness calibration error (HCE) We define the HCE as the Kullback-Leibler divergence (KL divergence) between predicted probability \hat{p} from a neural LM and actual probability p^* of *hu*man language model.

$$HCE = D_{KL}(\hat{\boldsymbol{p}}||\boldsymbol{p}^*). \tag{4}$$

Empirically, since p^* is not directly observable, we approximate it by the estimates of a temperaturescaled model that best fits human reading times (as discussed later). We denote the approximated HCE using such a method as HCE_{TS}.

3.3 Temperature-scaled surprisal

Temperature scaling (Guo et al., 2017) is a widelyused method to improve model calibration. Given 211 the output logit vector $\hat{\mathbf{z}}_i$ for sample *i*, a single 212 scalar T > 0 is applied to rescale $\hat{\mathbf{z}}_i$ before the 213 softmax activation: 214

$$\hat{q}_i = \max_k \sigma_{SM} (\frac{\hat{\mathbf{z}}_i}{T})^{(k)}, \tag{5}$$

where \hat{q}_i is the calibrated confidence for sample 216 *i*, and σ_{SM} is the softmax function. Scaling by 217

a scalar T does not alter the ranking; hence, the predicted label \hat{y}_i remains unchanged. As T > 1, it "softens" the probability distribution (i.e., makes the distribution more uniform), increasing uncertainty and entropy of the probability distribution, while T < 1 peaks the distribution. The parameter T in research on calibration is optimized by minimizing the negative log-likelihood on the validation set. In our experiments of fit to human RTs, we manually tune this temperature with T > 1.

218

219

221

222

223

224

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

267

Temperature scaling has been successfully applied in several applications: In knowledge distillation (Hinton et al., 2015), temperature scaling (with T > 1) is used to "soften" the knowledge (i.e., probability distribution) provided by the teacher model; in text generation, temperature is used to shape the probability distribution to ease certain aspects of the problems of top-k sampling (e.g., choosing an appropriate k value across varying contexts) (Ficler and Goldberg, 2017; Fan et al., 2018). Temperature tuning inherently shifts the model's output in the generation's quality/diversity spectrum (Caccia et al., 2018), with higher temperature decreasing the quality of generation while improving its diversity. This also aligns with our consideration of a possibility that human probability distributions might be flatter than the ones learned by language models and thus increasing the predictive diversity of surprisal provided by LLMs could potentially yield more human-like distributions.

Given Eq. 5, temperature-scaled surprisal is:

$$s_T(w_t, T) = -\log_2(\sigma_{SM}(\hat{\mathbf{z}}_{w_t}/T)^{(k^*)}),$$
 (6)

where $\hat{\mathbf{z}}_{w_t}$ and $k^* = y_{w_t}$ denote the logit vector and the actual word w_t class, respectively. For given $t \in (0,\infty)$, we simply denote $s_T(w_t, T = t)$ as $s_T|_{T=t}$. A temperature T with its best performance of final fit to RTs is denoted as T^* .

The extent to which a word's surprisal is affected by temperature scaling depends on the distribution and thus correlates with the entropy at word w_t . Consider an example of two five-class probability distributions $p_i = [0.8, 0.05, 0.05, 0.05, 0.05]$ and $p_i = [0.8, 0.2, 0, 0, 0]$, for which the word indicated by the first position in the probability vector has identical surprisal in both p_i and p_i . Notably, p_i is more uniform and p_i is more peaked, resulting in distinct entropy characteristics: $H(w_i|\boldsymbol{w}_{< i}) > H(w_i|\boldsymbol{w}_{< i})$ where the entropy defined as the expectation



Figure 1: Temperature-scaled surprisal $s_T(w_t, T)$ with corresponding $T \in [1, 2.5]$ for two random five-class probability distributions: $p_i =$ [0.8, 0.05, 0.05, 0.05, 0.05] and $p_j = [0.8, 0.2, 0, 0, 0]$. Dashed lines show Shannon entropy (H₁). Loosely dashed lines show Rényi entropy with $\alpha = 1/2$ (H_{1/2}).

of surprisal of current word w_t over vocabulary, $H(w_t|\boldsymbol{w}_{< t}) = \mathbb{E}_{w' \sim p(\cdot|\boldsymbol{w}_{< t})}[s(w')] = -\sum_{w' \in \overline{W}} p(w'|\boldsymbol{w}_{< t}) \log_2 p(w'|\boldsymbol{w}_{< t})$, where $\overline{W} = W \cup \{EOS\}$ denotes the set of vocabulary W with EOS token. Fig. 1 illustrates **a greater increase in surprisal** for a word with a more uniform distribution than with a more peaked distribution.

This figure also an ecdotally shows that the effect of applying temperature scaling with T > 1 is similar to the effect of setting $\alpha < 1$ in Rényi entropy. We will discuss the relationship between these parameters in more detail in Appendix A.

4 Experimental setup

4.1 Datasets

268

269

271

272

276

277

278

281

282

284

296

We conduct analyses on two self-paced reading corpora, the Natural Stories Corpus (Futrell et al., 2018) and the Brown Corpus (Smith and Levy, 2013), as well as on the Dundee Corpus (Kennedy et al., 2003), which contains the eye-movement record; our analyses in this paper focus on firstpass times² from the Dundee corpus. We follow previous work with respect to the preprocessing steps for each corpus (Kuribayashi et al., 2022; Shain et al., 2022). Appendix C includes details about the preprocessing steps of each corpus.

4.2 Language Models

Recent observations showed that surprisal provided by LLMs with more parameters and lower perplexity is less predictive of self-paced reading times 297 and eye-gaze durations (Shain et al., 2022; Oh 298 and Schuler, 2023b); across different experiments, 299 GPT-2 (Radford et al., 2019) surprisals were found 300 to predict human RTs the best. Therefore, we take 301 four variants of pretrained GPT-2 (small, medium, 302 large, xl) as our language models in all experiments. 303 Following prior work, we obtain the surprisal for 304 words composed of more than one subword by sum-305 ming up the surprisal estimates of the subwords. 306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

327

329

330

331

332

333

334

335

337

338

339

340

341

4.3 Metrics and evaluation

We measure the predictive power of surprisal estimates from different language models, which is denoted as the log-likelihood difference between a linear mixed-effects (LME) regression model using lme4 package (Bates et al., 2015) with a predictor of surprisal estimates (target model) and a model without surprisal (base model), following Goodkind and Bicknell (2018); Wilcox et al. (2020). More specifically, the metric of delta log-likelihood is defined as:

$$\Delta_{\text{llh}} = \text{llh}(f_{\phi}(\mathbf{v}^{tgt})) - \text{llh}(f_{\phi}(\mathbf{v}^{base})), \quad (7)$$

where \mathbf{v}^{tgt} is target predictor variables that include baseline predictor variables as well as predictor variables of our interest, such as surprisal or temperature-scaled surprisal. \mathbf{v}^{base} is base predictor variables only including baseline predictor variables. The greater the value of Δ_{llh} , the more valuable the additional surprisal estimates are for predicting human reading times.

For the calibration error evaluation, we set the number of bins M to 15 for both ECE and CECE, aligning with prior literature, such as works by Guo et al. (2017); Kumar et al. (2019); Rahimi et al. (2020b), to ensure consistency in addressing problems where comparable probability ranges are relevant. The calibration metrics (ECE and CECE) are evaluated separately on each of the reading time corpus D. For simplicity, our calibration evaluation is conducted at the token level. Given that many words have extremely low probabilities and thus are often grouped into a single bin, we also evaluate the calibration error *under the log probability binning scheme*. For further descriptions regarding the metrics and evaluation, see Appendix D.

²First pass times are calculated as the sum of all fixation durations from first entering to first leaving the word during the first pass, i.e., only those cases are counted where no words further advanced in the text have been fixated.

5 Results

342

343

345

347

351

353

354

361

365

367

370

372

373

377

379

384

385

387

389

5.1 Calibration of LLMs

Table 1 shows ECE and CECE in log binning scheme for GPT-2 models of different sizes. LLMs are in general well calibrated on language modeling. Besides, LLM calibration improves with scale. Larger LMs are better calibrated. This conclusion is consistent with calibration investigation evaluated in BIG-bench multiple-choice tasks in Srivastava et al. (2023) as well as in several tasks including language modelling in Zhu et al. (2023).

5.2 Main result: temperature-scaled surprisal improves human reading time prediction

We evaluate the predictive power of temperaturescaled surprisal. We scale T in the range of [1, 10]and measure Δ_{llh} , see Fig. 2. First, a confirmatory observation regarding the relationship between model size and predictive power: At T = 1, GPT-2 small exhibits the best predictive performance, and as the model size increases, Δ_{llh} declines, which is consistent with previous studies (Shain et al., 2022; Oh et al., 2022; Oh and Schuler, 2023b). Secondly, scaling the surprisal with T > 1 can significantly improve the predictive power across all corpora and LLMs. With optimal T^* , on Dundee, Natural Stories, and Brown, the Δ_{llh} improvement is 23-43%, 60-89%, and 14-24%, respectively. We access statistical significance of GPT-2 small in Appendix H, where we report a result of p < 0.001on three corpora. We also observe a consistent pattern: when increasing T, Δ_{llh} first rises then declines; the optimal value T^* falls within the range of (2, 3) (around 2.5) across all models and corpora in our setting. At T^* , even though the impact of model size on final performance is not fully recovered, the disparity diminishes. Smaller models continue to outperform, but the extent of model sizes influencing performance is reduced.

> Finally, **larger LMs typically have a larger human-likeness calibration error**, shown in Table 1. Larger LMs also require a higher value of T to reach their best performance and have a greater increase by temperature-scaled surprisal.

5.3 Calibration error vs. RT prediction error

Table 2 shows ECE and CECE in both equallyspaced and log binning schemes when T equals 1 and T^* on three corpora. Probability distribution shaped by an optimal T^* learnt for fit to human

		T^*	$\Delta_{llh} +$	$\text{HCE}_{\text{TS}}\downarrow$	$ECE_{log} \downarrow$	$CECE_{log} \downarrow$
Dundee	s	2.75	22.5	3.11	1.59	4.07E-03
	m	3.0	42.0	3.61	1.74	4.13E-03
	1	3.0	39.9	3.82	1.55	3.99E-03
	xl	3.25	43.2	4.13	1.29	3.84E-03
	S	2.5	60.3	3.31	1.91	1.53E-02
NC	m	2.5	63.0	3.50	1.80	1.50E-02
183	1	2.5	82.6	3.97	1.70	1.40E-02
	xl	2.5	89.0	4.07	1.56	1.35E-02
	s	2.5	13.7	3.10	1.69	1.53E-02
Daorra	m	2.5	16.2	3.29	2.27	1.51E-02
Brown	1	2.75	21.8	4.18	1.58	1.44E-02
	\mathbf{xl}	2.75	24.4	4.29	1.56	1.38E-02

Table 1: Optimal T^* , $\Delta_{\rm llh}$ improvement (%) ($\Delta_{\rm llh} + = (\Delta_{\rm llh}(T = T^*) - \Delta_{\rm llh}(T = 1))/\Delta_{\rm llh}(T = 1))$, and calibration errors (HCE_{TS}, % ECE and % CECE) for GPT2s on Dundee, Natural Stories (NS) and Brown. $\Delta_{\rm llh}$ values are multiplied by 1000. ECE and CECE are evaluated on log binning scheme.

RTs drastically hurts the model calibration regarding these two metrics. ECE and CECE with T^* are more than 10 times worse than values with T = 1. This discrepancy can be attributed to the different minima of deviations in LM human RT prediction and expected calibration error. The former is minimized towards words where LMs surprisal significantly deviates from human processing difficulty, while the latter is typically minimized with respect to the negative log-likelihood on a hold-out dataset (Guo et al., 2017; Rahimi et al., 2020a). 390

391

392

393

394

395

396

397

399

400

401

402

403

404

405

406

407

408

409

6 Linguistic analysis

Next we want to gain insight into what words benefit the most from temperature scaling. To this end, we analyze residuals from fitting LME regression models, identifying data points where scaling the temperature parameter notably enhances the fit of human RTs. Specifically, we quantify the improvement in fit by comparing the mean squared error (MSE) before and after adjusting the temperature

		ECE↓	$\text{ECE}_{\log} \downarrow$	CECE↓	$CECE_{log} \downarrow$
Dundee	T = 1	1.43	1.59	4.05E-03	4.07E-03
	$T=T^*$	28.68	28.68	7.30E-03	9.88E-03
NC	T = 1	2.48	1.91	1.83E-02	1.53E-02
143	$T = T^*$	35.85	35.85	3.16E-02	3.97E-02
D	T = 1	1.82	1.69	1.67E-02	1.53E-02
DIOWII	$T=T^*$	33.16	33.16	2.75E-02	3.34E-02

Table 2: Expected calibration errors (% ECE and % CECE) for GPT-2 small on Dundee, Natural Stories (NS) and Brown. Results are all evaluated on the equally-spaced binning scheme and log binning scheme.



Figure 2: Relationship between Δ_{llh} of GPT-2 models and corresponding temperature. T is scaled from 1.0 to 10.



Figure 3: Relationship between $\Delta_{\rm MSE}$ and negative log actual-word probability (surprisal). We take the number of bins to 20. Black dashed lines denote $\Delta_{\rm MSE} = 0$. Subsets containing less than 1% of data are ignored for each corpus.

410 to its optimal value as follows:

411

412

413

414

415

416

417

418

419

420

421

$$\Delta_{\text{MSE}}(F) = \text{MSE}_{T=1}(x_F) - \text{MSE}_{T=T^*}(x_F),$$
(8)

where $\text{MSE}_{T=T'}(x_F)$ is the MSE calculated by all the data x_F under the linguistic factor F. The difference $\Delta_{\text{MSE}}(F)$ thus quantifies the impact of scaling relative to the linguistic factor F. A higher $\Delta_{\text{MSE}}(F)$ signifies a greater influence of temperature-scaled surprisal of factor F. To ensure sufficient data in each subset, we only consider subsets including more than 1% of the data in each corpus.

6.1 Influence of low probability words

Given that temperature scaling enhances human 422 likeness by shaping the probability distribution, 423 it is natural to think about investigating whether 424 there exists an inherent relationship between the 425 distribution of probability and Δ_{MSE} . Specifically, 426 one might ask questions like if samples with low 427 probability gain more from temperature scaling or 428 the other way around. We find that high surprisal 429 words benefit more from temperature scaling than 430 low surprisal words, across all corpora, see Fig. 3. 431

6.2 Influence of word types

We investigate the effects of word-level properties, which include:

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

Named entities. Research has substantiated that named entities (NEs) require increased reading time for humans since during the processing of such words (Damasio et al., 2004; Wang et al., 2013). Oh and Schuler (2023b) showed that NEs are among the top two significant factors contributing to the discrepancies of large and small LMs across all corpus-by-LM combinations. Therefore, we were wondering whether the effect of temperature-scaling might be driven by NE. To test this, we automatically tagged NEs using a BERT base model (Devlin et al., 2019) fined-tuned for NER³.

Part-of-speech tags. Similarly, previous research has argued that the poor fit of large LMs is primarily due to assigning too low surprisal estimates to open-class words like nouns and adjectives (Oh and Schuler, 2023b). We POS-tagged the corpora using the NLTK toolkit (Bird et al., 2009) with the default Penn Treebank Tag set. In the following, we mainly focus on the four classes of open-class tags, as well as a subset of the whole closed-class tags (CC).

				d entities	POS tags					
	GPT2	Avg.	NE	non-NE	NN	ADJ	VERB	ADV	CC	
Dundee	s	26.3	87.0	23.4	33.8	100.5	-2.0	2.6	10.4	
	m	41.7	152.3	36.4	57.0	123.3	7.8	27.6	16.4	
	1	40.1	158.2	34.5	56.3	126.5	4.8	19.2	14.0	
	xl	41.4	168.2	35.4	60.0	125.5	6.9	19.7	13.5	
	s	105.7	186.8	104.6	148.7	152.5	122.0	49.0	77.1	
NC	m	108.5	155.9	107.9	145.3	152.0	130.1	60.8	80.8	
183	1	127.7	151.6	127.3	175.6	158.6	152.9	74.8	94.3	
	xl	123.3	141.8	123.1	163.6	145.4	<u>161.2</u>	81.5	89.0	
	s	37.2	266.0	28.1	54.3	-65.2	138.1	32.1	5.9	
Destro	m	41.4	257.6	32.8	71.4	-60.6	137.5	38.6	3.5	
DIOWII	1	42.6	265.3	51.1	69.9	-110.3	160.8	17.2	24.7	
	xl	54.8	282.3	45.8	90.5	-90.2	151.3	32.2	20.0	

Table 3: Δ_{MSE} measurement on word-level properties of GPT-2 models on Dundee, Natural Stories (NS) and Brown. Top-3 on each corpus-by-LM are underlined.

³Link: https://huggingface.co/dslim/bert-base-NER

Results. The result, as shown in Table 3, shows primary factors responsible for the benefit of using $s_T(w_t, T)$ for each corpus-by-LM combination. The top three influential subsets for each corpus are underlined. Among all datasets and models, **named entities perform to be the most beneficial word-level attribute.** In contrast, **closed-class words profit the least from temperature scaling.** Performance trends are consistent across different model variants on the same corpus.

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

We also measured empirically how often temperature scaling increased vs. decreased the surprisal estimate of a word. Our results show that for ca. 90% of words, surprisal estimates are increased through temperature scaling across all word classes. For the subset of named entities, a slightly smaller percentage exhibits increased surprisal estimates. For a full analysis across different corpora and models, see Table 5 in Appendix B.

We further investigate the benefit of temperaturescaled surprisal (quantified by Δ_{MSE}) given the subset of words whose probability decreases (or increases). The results are in Table 4. On Dundee, the main gain arises from the reduction of large probabilities via temperature scaling. Conversely, for Natural Stories, the primary benefit comes more strongly from words with originally very low probability, which become more probable. For Brown, the effects are evenly split. These findings align with our theoretical intuition that temperature scaling enhances the fit performance by making probabilities more smooth, which means not only making high probabilities lower but also making very low probabilities higher and close to 1/K, since a very low probability also means the model is confident in the incorrectness of certain classes.

Considering effects on named entities more specifically, we find that on Natural Stories and Brown, the benefit of temperature scaling can mostly be attributed to reducing the probability estimates of highly predictable entities, while on Dundee the beneficial effect mostly arises from increasing probabilities of named entities. We speculate that this could be due to the types of most frequent named entities that occur in the different text sorts, and present a more detailed analysis of this aspect in Appendix B.

6.3 Influence of multiple-token words

A fact that is often ignored (but see Nair and Resnik, 2023) is that modern LLMs use subword tokeniza-

				Named entities				
		Av	Avg.		ЛЕ	non-NE		
Corpus	GPT2	$p_{w_t}\downarrow$	p_{w_t}	$p_{w_t}\downarrow$	$p_{w_t}\uparrow^*$	$p_{w_t}\downarrow$	$p_{w_t}\uparrow$	
	s	27.4	18.2	81.3	107.2	25.1	10.1	
Dundee	m	41.9	39.8	139.1	205.6	37.8	23.9	
	1	41.0	31.3	156.1	166.6	36.2	18.0	
	xl	42.5	29.8	170.2	158.8	37.0	16.9	
	S	94.5	275.6	218.5	3.0	92.9	284.9	
NC	m	105.7	158.3	179.3	-34.9	104.7	163.9	
113	1	125.0	166.1	197.5	-224.8	124	175.4	
	xl	121.8	140.7	197.3	-272.6	120.8	149.5	
-	s	37.6	32.6	329.7	-170.6	26.6	45.5	
Drown	m	39.1	72.3	276.0	143.6	30.5	66.3	
Brown	1	52.7	28.1	325.8	-205.9	42.5	44.4	
	xl	50.9	111.5	298.2	168.2	41.7	107.1	

Table 4: Given words whose probability decreases (and increases), the corresponding $\Delta_{MSE}(p_{wt}\downarrow)$ (and $\Delta_{MSE}(p_{wt}\uparrow)$) measurement for GPT-2 models on Dundee, Natural Stories (NS) and Brown. A higher Δ_{MSE} is displayed in bold in the average across all word types (Avg.), named entities (NE), and non-named entities (non-NE) columns, respectively, for each corpusby-LM combination. The column with * indicates insufficient (less than 1%) data.

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

tion. This means that long words may consist of several tokens. In this case, the probability of the complete word is calculated by multiplying the probabilities of the subword tokens (and the word's surprisal is correspondingly calculated by adding the surprisals of the subwords). While this may often not matter, whether a word is tokenized into a single subword or several subwords can make a remarkable difference when applying temperature scaling: imagine a long / difficult word which has a low probability (and correspondingly a high surprisal). If this word were to be represented as a single subword token, temperature scaling might have the effect that the probability of this word gets *increased* during temperature scaling, and its surprisal estimate is hence decreased at T > 1.

If, on the other hand, the same word were to be composed of two subword tokens, one or both of the subword tokens can be expected to have a higher probability (than a hypothetical single subword token), and it is possible that during temperature scaling, the probabilities of the subword tokens would each be *decreased* at T > 1, such that the sum of the surprisals of the subword tokens would be much higher, compared to the word's surprisal estimate at T = 1.

To summarize, whether the surprisal of a certain word would increase or decrease after temperature scaling could depend on whether that word happens to be included in the subword token vocabulary or



Figure 4: Relationship between Δ_{llh} of GPT-2 s on three corpora and corresponding temperature T.

not.⁴ In order to quantify to what extent subword tokenization affects surprisal estimates, we conducted several analyses.

Fig. 4 shows Δ_{llh} under various conditions: scaling all words (consistent with experiments in Section 5.2) vs. taking into the analysis only the subset of single-token words and multiple-token words. The comparison between the full, dotted, and dashed lines highlights that **the benefit of temperature-scaled surprisal comes primarily** from the scaling of multiple-token words.

Next, it is interesting to consider for what percentage of multiple-token words temperaturescaling *increases* the surprisal. We find that the surprisal of more than 90% of multiple-token words increases, and the **ratio is higher than across singletoken words** by ca. 6% on Dundee and Brown, see Table 12 in Appendix L for more details.

7 Discussion

Our experiments demonstrate that choosing a temperature around 2.5 improves the fit to human reading times. Furthermore, we find that this effect is chiefly driven by an improved fit for words which consist of several subword tokens.⁵ Named entities and other open class words tend to have a larger tendency to contain several subword tokens, which can explain why temperature scaling is particularly effective for these words.

So what does all of this mean for surprisal estimates from LLMs and reading time prediction? Firstly, following the argumentation of Oh and Schuler (2023b), it is possible that indeed the effect is driven by humans failing to accurately estimate the probability of rare words (rare words being the ones that are split up into several subwords), because they do not reach sufficient language experience or because human language models do not track these probabilities well. In this case, temperature-scaling rare words to which the LLM assigns a too high probability (and hence a low surprisal) would be a good strategy to counteract the discrepancy between humans and LLMs. From LLMs' perspective, recalling the observation from Section 5.3 that larger LLMs that yield poorer fits to RTs are actually better calibrated, hence the massive training dataset might be at the cause of driving these models away from the human-like predictive processing, aligning with Oh and Schuler (2023a). 574

575

576

577

578

579

580

582

583

584

585

586

587

588

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

610

611

612

613

614

615

616

617

618

619

620

621

622

623

Secondly, it is likely that the beneficial effect of temperature scaling is an artifact of subword tokenization, and that this effect would disappear if all words were composed of only a single subword token (cf. our explanation in Section 6.3). That is, temperature scaling would not be beneficial because of the reasons that motivated this research originally, but only because it is a way of assigning higher surprisal to words consisting of several subword tokens. In order to test this hypothesis, one would have to re-train a GPT-2 model using a vocabulary that at least includes all words that are contained in the reading time corpora, and then rerunning the analysis to check whether a beneficial effect of temperature scaling can still be found.

Finally, it is also possible that the splitting of a word into subwords coincides with the reader fixating a word several times, and that these added fixations lead to an overestimate in RTs compared to the actual surprisal experienced by a human reader. Future work could investigate this hypothesis by analysing RTs on subwords instead of aggregated words (with the caveat that subword tokens may not be cognitively plausible units).

8 Conclusion

This paper studies the prediction of human RTs from the perspective of probability distribution. We make the following contributions: (1) We demonstrate that the prediction of RTs can be significantly improved via temperature scaling of LLM probability estimates. (2) We demonstrate that the primary benefit of temperature-scaled surprisal is driven by words composed of several subword tokens. These words also tend to be rarer / long open-class words. Future work should investigate the interaction of subword tokenization and temperature scaling, as well as the issue of tokenization in the analysis of eye-tracking data.

572

573

538

⁴Distributions of surprisal for single vs. multiple token words before and after temperature scaling are provided in Fig. 8 in Appendix L.

⁵Appendix K shows that subword tokenization has larger explanatory power than word class.

Limitations

624

650

663

664

666

667

669

670

672

In this work, the identification of the optimal T for temperature-scaled surprisal is manually tuned. Future research could develop an automated method 627 to determine this optimal value, e.g., from specific characteristics of LLMs or corpora. Additionally, a question may be asked whether the possible nonlinear relationship between surprisal and reading times (Shain et al., 2022; Hoover et al., 2023) could influence the temperature-scaled surprisal's superiority over original surprisal. Investigating the effectiveness of temperature-scaled surprisal using generalized additive models, a branch of models that assume less about the linearity than linear mixed effect models employed here, would be an extension. Finally, exploring effects of temperature-scaled sur-639 prisal on different measures of fixation duration could be considered in future work.

Ethical Considerations

The datasets and packages we used are all publiclyavailable and have no privacy issues.

References

- Doris Aaronson and Hollis S Scarborough. 1976. Performance theories for sentence coding: Some quantitative evidence. *Journal of Experimental Psychology: Human perception and performance*, 2(1):56.
- Bernhard Angele, Elizabeth R Schotter, Timothy J Slattery, Tara L Tenenbaum, Klinton Bicknell, and Keith Rayner. 2015. Do successor effects in reading reflect lexical parafoveal processing? evidence from corpusbased and experimental eye movement data. *Journal* of Memory and Language, 79:76–96.
- Christoph Aurnhammer and Stefan L Frank. 2019. Evaluating information-theoretic measures of word prediction in naturalistic sentence reading. *Neuropsychologia*, 134:107198.
- Douglas Bates, Martin Mächler, Ben Bolker, and Steve Walker. 2015. Fitting linear mixed-effects models using lme4. *Journal of Statistical Software*, 67:1–48.
- Steven Bird, Ewan Klein, and Edward Loper. 2009. Natural language processing with Python: analyzing text with the natural language toolkit. " O'Reilly Media, Inc.".
- Massimo Caccia, Lucas Caccia, William Fedus, Hugo Larochelle, Joelle Pineau, and Laurent Charlin.
 2018. Language gans falling short. arXiv preprint arXiv:1811.02549.
- Hanna Damasio, Daniel Tranel, Thomas Grabowski, Ralph Adolphs, and Antonio Damasio. 2004. Neural

systems behind word and concept retrieval. *Cognition*, 92(1-2):179–229. 673

674

675

676

677

678

679

680

681

682

683

684

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

- Vera Demberg and Frank Keller. 2008. Data from eyetracking corpora as evidence for theories of syntactic processing complexity. *Cognition*, 109(2):193–210.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- 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.
- Jessica Ficler and Yoav Goldberg. 2017. Controlling linguistic style aspects in neural language generation. In *Proceedings of the Workshop on Stylistic Variation*, pages 94–104, Copenhagen, Denmark. Association for Computational Linguistics.
- Kenneth I Forster, Christine Guerrera, and Lisa Elliot. 2009. The maze task: Measuring forced incremental sentence processing time. *Behavior research methods*, 41:163–171.
- 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.
- Stefan L Frank and Rens Bod. 2011. Insensitivity of the human sentence-processing system to hierarchical structure. *Psychological science*, 22(6):829–834.
- 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).
- 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.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. 2017. On calibration of modern neural networks. In *International conference on machine learning*, pages 1321–1330. PMLR.

827

828

829

830

831

832

779

780

728

729

- 770 773 774 775
- 778

- 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. 2003. The information conveyed by words in sentences. Journal of psycholinguistic research, 32:101-123.
- John Hale. 2006. Uncertainty about the rest of the sentence. Cognitive science, 30(4):643-672.
- Dan Hendrycks, Kimin Lee, and Mantas Mazeika. 2019. Using pre-training can improve model robustness and uncertainty. In International conference on machine learning, pages 2712–2721. PMLR.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. arXiv preprint arXiv:1503.02531.
- Jacob Louis Hoover, Morgan Sonderegger, Steven T Piantadosi, and Timothy J O'Donnell. 2023. The plausibility of sampling as an algorithmic theory of sentence processing. Open Mind, 7:350-391.
- Alan Kennedy, Robin Hill, and Joël Pynte. 2003. The dundee corpus. In Proceedings of the 12th European Conference on Eye Movement.
- Meelis Kull, Miquel Perello Nieto, Markus Kängsepp, Telmo Silva Filho, Hao Song, and Peter Flach. 2019. Beyond temperature scaling: Obtaining wellcalibrated multi-class probabilities with dirichlet calibration. Advances in neural information processing systems, 32.
- Ananya Kumar, Percy S Liang, and Tengyu Ma. 2019. Verified uncertainty calibration. Advances in Neural Information Processing Systems, 32.
- Sawan Kumar. 2022. Answer-level calibration for freeform multiple choice question answering. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 665-679, Dublin, Ireland. Association for Computational Linguistics.
- Tatsuki Kuribayashi, Yohei Oseki, Ana Brassard, and Kentaro Inui. 2022. Context limitations make neural language models more human-like. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 10421–10436, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Roger Levy. 2008. Expectation-based syntactic comprehension. Cognition, 106(3):1126-1177.
- Tal Linzen and T Florian Jaeger. 2014. Investigating the role of entropy in sentence processing. In Proceedings of the fifth workshop on cognitive modeling and computational linguistics, pages 10-18.

- 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.
- Don C Mitchell and David W Green. 1978. The effects of context and content on immediate processing in reading. The quarterly journal of experimental psychology, 30(4):609-636.
- Sathvik Nair and Philip Resnik. 2023. Words, subwords, and morphemes: What really matters in the surprisalreading time relationship? In Findings of the Association for Computational Linguistics: EMNLP 2023, pages 11251–11260, Singapore. Association for Computational Linguistics.
- Byung-Doh Oh, Christian Clark, and William Schuler. 2022. Comparison of structural parsers and neural language models as surprisal estimators. Frontiers in Artificial Intelligence, 5:777963.
- Byung-Doh Oh and William Schuler. 2023a. Transformer-based language model surprisal predicts human reading times best with about two billion training tokens. In Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023, pages 1915–1921. Association for Computational Linguistics.
- Byung-Doh Oh and William Schuler, 2023b. Why does surprisal from larger transformer-based language models provide a poorer fit to human reading times? Transactions of the Association for Computational Linguistics, 11:336–350.
- Tiago Pimentel, Clara Meister, Ethan G. Wilcox, Roger Levy, and Ryan Cotterell. 2023. On the effect of anticipation on reading times. Transactions of the Association for Computational Linguistics.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. OpenAI blog, 1(8):9.
- Amir Rahimi, Kartik Gupta, Thalaiyasingam Ajanthan, Thomas Mensink, Cristian Sminchisescu, and Richard Hartley. 2020a. Post-hoc calibration of neural networks. arXiv preprint arXiv:2006.12807, 2.
- Amir Rahimi, Amirreza Shaban, Ching-An Cheng, Richard Hartley, and Byron Boots. 2020b. Intra order-preserving functions for calibration of multiclass neural networks. Advances in Neural Information Processing Systems, 33:13456–13467.
- Keith Rayner. 1998. Eye movements in reading and information processing: 20 years of research. Psychological bulletin, 124(3):372.
- David Reeb and Michael M Wolf. 2015. Tight bound on relative entropy by entropy difference. IEEE Transactions on Information Theory, 61(3):1458–1473.

- 852 853
- 854 855 857 858
- 865
- 868

871

872

874

875

876

877

878

885

844

842

833

834

836

839

845

847

849

850

Association for Computational Linguistics.

Press.

ing time.

27(3):379-423.

Nathaniel J Smith and Roger Levy. 2013. The effect of word predictability on reading time is logarithmic. Cognition, 128(3):302-319.

Alfréd Rényi. 1961. On measures of entropy and infor-

mation. In Proceedings of the Fourth Berkeley Sym-

posium on Mathematical Statistics and Probability,

Volume 1: Contributions to the Theory of Statistics,

volume 4, pages 547–562. University of California

Cory Shain, Clara Meister, Tiago Pimentel, Ryan Cot-

Claude Elwood Shannon. 1948. A mathematical theory

Chenglei Si, Chen Zhao, Sewon Min, and Jordan Boyd-

Graber. 2022. Re-examining calibration: The case

of question answering. In Findings of the Associa-

tion for Computational Linguistics: EMNLP 2022,

pages 2814–2829, Abu Dhabi, United Arab Emirates.

of communication. The Bell system technical journal,

terell, and Roger Levy. 2022. Large-scale evidence

for logarithmic effects of word predictability on read-

Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. 2023. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. Transactions on Machine Learning Research.

Marten van Schijndel and Tal Linzen. 2019. Can entropy explain successor surprisal effects in reading? In Proceedings of the Society for Computation in Linguistics (SCiL) 2019, pages 1-7.

Lin Wang, Zude Zhu, Marcel Bastiaansen, Peter Hagoort, and Yufang Yang. 2013. Recognizing the emotional valence of names: An erp study. Brain and Language, 125(1):118-127.

Ethan Gotlieb Wilcox, Jon Gauthier, Jennifer Hu, Peng Qian, and Roger P. Levy. 2020. On the predictive power of neural language models for human realtime comprehension behavior. In Proceedings of the 42nd Annual Meeting of the Cognitive Science Society, page 1707–1713.

Ethan Gotlieb Wilcox, Tiago Pimentel, Clara Meister, Ryan Cotterell, and Roger P Levy. 2023. Testing the predictions of surprisal theory in 11 languages. Transactions of the Association for Computational Linguistics.

Chiwei Zhu, Benfeng Xu, Quan Wang, Yongdong Zhang, and Zhendong Mao. 2023. On the calibration of large language models and alignment. In Findings of the Association for Computational Linguistics: EMNLP 2023, pages 9778–9795, Singapore. Association for Computational Linguistics.

Α **Connection to Contextual Rényi** Entropy

While a lot of work has investigated the effect of next word entropy on reading times (Hale, 2003, 2006; Linzen and Jaeger, 2014; Angele et al., 2015; van Schijndel and Linzen, 2019; Aurnhammer and Frank, 2019; Pimentel et al., 2023), we will here focus on contextual Rényi entropy (the entropy of the probability distribution at the current time stamp, which is parameterized by α), as proposed in Pimentel et al. (2023) to represent human anticipatory reading process. Pimentel et al. (2023) find that Rényi entropy with an optimal α^* in the range of (0,1) (around 1/2) obtains the best performance in reading time prediction (compared to Shannon Entropy ($\alpha = 1$) or compared to unscaled surprisal estimates).

Mathematically, Contextual Rényi entropy (Rényi, 1961) is defined as:

 H_{α}

$$(w_t \mid \boldsymbol{w}_{< t}) = \lim_{\beta \to \alpha} \frac{1}{1 - \beta} \log_2 \sum_{w \in \overline{\mathcal{W}}} (p(w \mid \boldsymbol{w}_{< t}))^{\beta}.$$
(9)

For given $\alpha' \in (0, \infty)$, we simply denote $H_{\alpha}(w_t \mid$ $|w_{< t}\rangle|_{\alpha = \alpha'}$ as $H_{\alpha}|_{\alpha = \alpha'}$.

Theorem 1 (Monotonicity of $s_T(w_t, T)$ and $H_{\alpha}(w_t \mid \boldsymbol{w}_{\leq t}))$. Given any probability distribution **p** with actual-word probability $p_{w_t} > 1/K$, where K is the number of classes, temperature-scaled surprisal $s_T(w_t, T)$ is strictly monotonically increasing in $\Delta_T \in [1, \infty]$, Rényi entropy $H_{\alpha}(w_t \mid \boldsymbol{w}_{\leq t})$ is strictly monotonically decreasing in $\Delta_{\alpha} \in [0, 1]$, especially,

$$s_T|_{T=1} < s_T|_{T=T^*} < \lim_{T \to \infty} s_T(w_t, T)$$
 (10)

$$H_{\alpha}|_{\alpha=1} < H_{\alpha}|_{\alpha=1/2} < H_{\alpha}|_{\alpha=0}, \quad (11)$$

where T^* is the optimal T of fit to RTs in the range of Δ_T .

Proof. Eq. (10) can be easily verified by considering the monotonicity of temperature-scaled softmax output $\sigma_{SM}(\hat{\boldsymbol{z}}_{w_t}/T)$. The second part of Eq. (11) can be rewritten as:

$$H_{\alpha}|_{\alpha=1/2} = 2\log_2 \sum_{w \in \overline{\mathcal{W}}} \sqrt{p(w|\boldsymbol{w}_{< t})}$$
(12)

$$< 2\log_2 \sqrt{K\sum_{w\in\overline{W}} p(w|\boldsymbol{w}_{< t})}$$
 (13) 920

$$= -\log_2(1/K) = H_{\alpha}|_{\alpha=0},$$
 (14) 92

906

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

907 908

909 910

911 912

913

914 915

916

917

918

919

920

921

922

923

928where for the step from Eq. (12) to Eq. (13) we929use AM-QM inequality and K is the number of930classes in tokenizer. The first part of Eq. (11) can931be rewritten as:

932
$$H_{\alpha}|_{\alpha=1/2} = 2\log_2 \sum_{w \in \overline{\mathcal{W}}} \sqrt{p(w|\boldsymbol{w}_{< t})}$$
(15)

933
$$> 2\log_2 \sqrt{\prod_{w \in \overline{\mathcal{W}}} (\frac{1}{p(w|\boldsymbol{w}_{< t})})^{p(w|\boldsymbol{w}_{< t})}}$$
(1

934

949

953

$$= \sum_{w \in \overline{\mathcal{W}}} p(w|\boldsymbol{w}_{< t}) \log_2 p(w|\boldsymbol{w}_{< t}) = \mathbf{H}_{\alpha}|_{\alpha}$$
(17)

(16)

where from Eq. (15) to Eq. (16) we use AM-GM inequality.

Theorem 2 *Rényi entropy with* $\alpha = 0$ *is equivalent to temperature-scaled surprisal with* $T \rightarrow \infty$.

$$\mathcal{H}_{\alpha}(w_t \mid \boldsymbol{w}_{< t})|_{\alpha = 0} = \lim_{T \to \infty} s_T(w_t, T). \quad (18)$$

940*Proof.* By plugging in $\alpha = 0$, Contextual Rényi941entropy recovers to be the entropy that readers942concentrate on the count of potential words with943nonzero probabilities, which is defined in Eq. (5)944in Pimentel et al. (2023). As $T \to \infty$, temperature-945scaled surprisal converges to the surprisal induced946by random guessing. Given the assumtion that947 $p(w|w_{<t}) > 0$ for each word $w \in \overline{W}$, LHS be-948comes:

$$LHS = -\log_2(1/K), \tag{19}$$

950 where K is the number of classes. As $T \to \infty$, 951 RHS becomes:

952
$$RHS = -\lim_{T \to \infty} \log_2 \frac{e^{z_{w_t}/T}}{\sum_{w \in \overline{W}} e^{z_w/T}}$$
(20)

$$= -\log_2(1/K) \tag{21}$$

Theorem 3 For $K \ge 2$, the expectation of the L1 norm between Rényi entropy with $\alpha = 1$ and temperature-scaled surprisal with T = 1 has an upper bound.

958
$$\mathbb{E}[|s_T|_{T=1} - \mathcal{H}_{\alpha}|_{\alpha=1}|] < \sqrt{\frac{1}{4}\log^2(K-1) + 1}$$
(22)

Proof. With Jensen's inequality, we have:

$$\mathbb{E}[|s_T|_{T=1} - H_{\alpha}|_{\alpha=1}|]$$
(23) 960

$$\leq \sqrt{\mathbb{E}[(s_T|_{T=1} - H_{\alpha}|_{\alpha=1})^2]}$$
 (24) 961

$$= \sqrt{\mathbb{E}[(-\log_2 p_{w_t} - \sum_{w \in \overline{\mathcal{W}}} p(w)(-\log_2 p(w)))^2]}$$
 96

959

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

$$=\sqrt{\operatorname{Var}[s_T|_{T=1}]} \tag{26}$$

$$<\sqrt{\frac{1}{4}\log^2(K-1)}+1,$$
 (27) 9

where $Var[\cdot]$ denotes the variance. The last inequality is shown by Lemma 4, completing the proof of this theorem.

Lemma 4 (Maximum variance of the surprisal). (See Theorem 8 and Lemma 15 in (Reeb and Wolf, 2015)). Let $\rho = \text{diag}(p_1, p_2, ..., p_d)$ be a state on a *d*-dimensional system. Let $-\log p_i$ be the surprisal of the output *i* in this system. Define N_d to be:

$$N_d := \frac{1}{4} \log^2(d-1) + 1.$$
 (28)

For $d \ge 2$, the variance of surprisal has a tight upper bound:

$$\operatorname{var}_{\rho}(-\log \rho) < N_d \tag{29}$$

Theorem 2 claims the equivalence of temperaturescaled surprisal $s_T(w_t, T)$ and Rényi entropy H_{α} when $T \to \infty$ and $\alpha = 0$. Theorem 3, on the other side, gives an upper bound when T = 1and $\alpha = 1$. Intuitively, when $T \in (1, \infty), s_T$ can be considered as a softened version of $s_T|_{T=1}$. Similarly, when $\alpha \in (0,1)$, H_{α} can be considered as a softened version of $H_{\alpha}|_{\alpha=1}$. Mathematically, Theorem 1 provides the monotonicity of both functions within their respective domains. Hypothetically, given the above conditions, when tuning both functions with the aim of a better fit to RTs, $s_T|_{T=T^*}$ and $H_{\alpha}|_{\alpha=1/2}$ might be close. Empirically, Fig. 5 illustrates the relationship between averaged Rényi entropy $H_{\alpha}|_{\alpha=\{0,1/2,1\}}$ and $\overline{s}_T|_{T=\{1,T^*,\infty\}}$ on probabilities on three corpora. Notably, $\overline{\mathrm{H}}_{\alpha}|_{\alpha=1/2}$ and $\overline{s}_{T}|_{T=T^{*}}$ are closely aligned, especially when compared with other entropy and surprisal data points. This empirical evidence partly verifies Theorem 2, Theorem 3 and our hypothesis.



Figure 5: A comparison of averaged temperaturescaled surprisal $\overline{s}_T|_{T=\{1,T^*,\infty\}}$ and Rényi entropy $\overline{\mathrm{H}}_{\alpha}|_{\alpha=\{0,1/2,1\}}$.

B Further analysis in Section 6.2

998

999

1000

1002

1003

1005

1006

1007

1008

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1023

1024

1025

1027

1028

1029

1030

1031

1032

1033

1034

We observe that larger LMs exhibit an increased $\Delta_{\rm MSE}$ by utilizing temperature-scaled surprisal, as shown in the average column (Avg.) of Table 3. Specifically, on Dundee, the top 2 models achieving the largest improvement through temperature scaling are GPT-2 medium and xl, while GPT-2 large and xl have the most benefit on Natural Stories and Brown. This result is consistent with previously observed Δ_{llh} improvement (Δ_{llh} +) across the corpus-by-LM reported in Table 1, suggesting a correlation between model likelihood and MSEs of the regression models. We do not observe a mismatch between them, as posited by Oh and Schuler (2023b) that LME models achieve similar MSEs irrespective of obvious differences in model likelihood.

Regarding the effect of the change (increase or decrease) of actual-word probability on the final fit to RTs, we first analyzed the ratio of probabilities decreasing (or increasing) for all words, as well as for subsets with specific word-level properties, choosing named entities as the representative, as shown in Table 5. We observed that probabilities of the majority of words (around 80-90%) decrease by temperature scaling. Compared with the average across all word types (as indicated in the 'Avg.' column), named entities exhibit a lower ratio of probability reduction. Larger LMs tend to have a higher ratio, especially the ratio for named entities, likely because smaller models may lack the specific knowledge of less common terms, such as named entities.

Recalling one of the results in Section 6.2 that the main advantage of temperature-scaled surprisal arises from reduction of large probabilities on Dundee and the amplification of small probabilities on Natural Stories. However, for named entities, 1035 the story is converse on Dundee vs. on Natural Sto-1036 ries and Brown, where for the latter two corpora, 1037 the advantage is primarily due to reducing the prob-1038 abilities of highly predictable entities. We shed 1039 light to the possible reason of such a discrepancy in 1040 Fig. 6, which displays the top 15 frequent words for 1041 GPT-2 small on three corpora. Notably, Natural 1042 Stories and Brown show a marked lack of words 1043 with increased probabilities (blue bins) compared 1044 to Dundee. This lack weakens the overall impact 1045 of rising probabilities (quantified by $\Delta_{MSE}(p_{w_t}\uparrow)$). 1046 Specifically, on Brown, only 4 out of 15 top fre-1047 quent words have the part of increased probabilities 1048 (blue bins), correlating with the largest discrep-1049 ancy in Δ_{MSE} between probabilities that decrease 1050 (329.7) and those that increase (-170.6) in Table 4. 1051

		А	vg.	Name	d entities
Corpus	GPT2	$p_{w_t}\downarrow$	$ res \downarrow$	$p_{w_t}\downarrow$	$ res \downarrow$
Dundee	S	88.0	51.8	78.1	52.3
	m	89.6	52.5	80.1	54.1
	1	90.2	52.3	80.1	53.5
	xl	91.4	52.4	82.7	54.3
	S	93.8	55.0	85.3	51.8
Natural	m	94.7	55.2	89.1	53.2
Stories	1	93.5	55.7	89.1	53.4
	xl	92.1	55.5	88.2	52.8
	S	91.8	51.5	87.3	50.9
Drown	m	93.2	51.5	86.1	50.9
DIOWII	1	93.3	51.8	88.6	52.1
	xl	93.5	51.7	87.8	53.3

Table 5: The ratio of probability of predicted word p_{w_t} getting smaller and the absolute value of residuals |res| getting smaller for GPT-2 models on three corpora.

C Preprocessing steps

On Dundee ET corpus (Kennedy et al., 2003), 1053 we use the first-pass gaze duration. Following 1054 prior work (Kuribayashi et al., 2022), we remove 1055 words containing numbers or punctuation, words 1056 that are either the first or the last one in a line, 1057 as well as words whose previous words contain numbers or punctuation. On Natural Stories SPR 1059 corpus (Futrell et al., 2018), following Shain et al. (2022), we remove words if the RT is less than 1061 100ms or greater than 3,000ms, if the words are in 1062 the first or last position of each story, if participants 1063 answered less than 5 out of 8 comprehension questions correctly, if words contain numbers or punc-1065



Figure 6: Top 15 frequent named entities for GPT-2 small on Dundee, Natural Stories and Brown. \uparrow and \downarrow denote probability being higher and smaller, respectively. \bigcirc and \times denote unbeneficial words (absolute residual error increases) and beneficial words (absolute residual error decreases) by temperature scaling, respectively.

tuation, and if words whose previous words containing numbers or punctuation. On Brown SPR corpus (Smith and Levy, 2013), following Shain et al. (2022), we remove words if the RT is less than 100ms or greater than 3,000ms and if words contain numbers or punctuation.

1066

1067

1068

1069

1070

1071

1074

1075

1077

1078

1079

1081

1082

1083

1084

1085

1086

1087

1088

1089

1090

1091

1094

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104

1105

D Further descriptions on metrics and evaluation

We evaluate calibration error (% ECE and % CECE) in both equally-spaced and log binning schemes. In equally-spaced binning scheme, the samples are grouped into $M \in \mathbb{N}$ equally-spaced interval bins based on their confidences \hat{p}_i . Conversely, the log binning scheme operates under an *empirical upper limit* for $-\log_2 \hat{p}_i$, denoted as $\max(-\log_2 \hat{p})$. Table 6 shows ranges of \hat{p} and $-\log_2 \hat{p}$ for GPT2s on three corpora. For this scheme, we establish $M \in \mathbb{N}$ log-equally-spaced interval bins within the range of $(0, \max(-\log_2 \hat{p})]$.

We investigate scaling $T \in [1, 10]$, considering both densely and sparsely distributed points. The values examined are detailed as follows: [1.0, 1.1, ..., 1.9] for dense intervals, [2.0, 2.25, ..., 3.25] for moderately spaced intervals, and [3.5, 4.0, ..., 10.0] for sparse intervals.

Following Kuribayashi et al. (2022), reading times of a base model are modelled by the following formula:

$$rt \sim \text{freq} * \text{length} + \text{freq}_{rev_1} * \text{length}_{rev_1} + (1|\text{article}) + (1|\text{subj}_{id})$$

(30)

A target model additionally includes surprisal estimates of current words and previous words:

$$rt \sim surprisal + surprisal_prev_1 + surprisal_prev_1$$

+ freq * length + freq_prev_1 * length_prev_1 + (1|article) + (1|subj id).

(31)

On Dundee corpus, both models also include features of [screenN, lineN, segmentN]. We also perform experiments with both models without interactions among predictors in Appendix I.

E Exploring further effectiveness of temperature-scaled surprisal over basic predictors

1106In this section, we explore the question of whether1107the benefit of temperature-scaled surprisal holds

only for regression models already containing other 1108 predictors such as length and frequency. We con-1109 duct experiments similar to those detailed in Sec-1110 tion 5.2 while setting base predictor variables \mathbf{v}^{base} 1111 to 0 and target predictor variables \mathbf{v}^{tgt} to only 1112 temperature-scaled surprisal $s_T(w_t, T)$ in Eq. 7. 1113 Fig. 7 shows that excluding base predictors de-1114 crease but not totally impact the effectiveness of 1115 temperature-scaled surprisal. 1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

F Calibration error for single-token and multiple-token words

In Table 7, we demonstrate the calibration error (%ECE) for single-token and multiple-token words for GPT-2 small. Calibration evaluation is conducted at the token level as before. Results indicate that **multiple-token words show larger calibration errors than single-token words**.

G Probability distribution before and after temperature scaling

Fig. 8 shows actual-word probability distribution before and after temperature scaling for GPT-2 small on three corpora. **Multiple-token words tend to have smaller probabilities than singletoken words**, both before and after temperature scaling.

H Significant test of temperature-scaled surprisal

We report the statistical significance based on selecting the most representative model, GPT2s, on three corpora in Table 8. Models with temperaturescaled surprisal lead to statistically significant posi-2 tive Δ_{llh} (p < 0.001).

I Analysis on correlations among predictors

We investigate the question of whether the benefit 1142 of temperature-scaled surprisal is primarily due to 1143 the interactions and correlations among predictors. 1144 We first run experiments with the original target 1145 LME model as in Eq. 31 (denoted as model 1), a 1146 model that has no interactions between frequency 1147 and length as in Eq. 32 (denoted as model 2) and a 1148 third model that has no interactions and addition-1149 ally includes random slopes for subject as in Eq. 33 1150 (denoted as model 3). 1151

	$\hat{m{p}}$	$-\log_2 \hat{\boldsymbol{p}}$
Dundee	[4.99e-03, 1)	(0, 7.65]
Natural Stories	[8.567e-03, 1)	(0, 6.87]
Brown	[8.15e-03, 1)	(0, 6.94]

Table 6: Ranges of \hat{p} and $-\log_2 \hat{p}$ for GPT2s on Dundee, Natural Stories and Brown.



Figure 7: Relationship between Δ_{llh} of GPT-2 small and corresponding temperature. T is scaled from 1.0 to 10. Base predictor variables \mathbf{v}^{base} and target predictor variables are 0 and temperature-scaled surprisal $s_T(w_t, T)$, respectively.

		ECE _{single}	$ECE_{multiple}$
Dundaa	T = 1	1.98	2.05
Dulluee	T = T *	25.58	36.10
Natural Starias	T = 1	2.20	3.78
Inatural Stories	T = T *	32.38	47.02
Drown	T = 1	1.69	3.86
DIOWII	T = T *	28.70	42.99

Table 7: Expected calibration errors of tokens in single-token (% ECE_{single}) and multiple-token words (% $ECE_{multiple}$) before and after temperature scaling for GPT-2 small on Dundee, Natural Stories and Brown. Results are all evaluated on the equally-spaced binning scheme.

Corpora	Models	р
Dundee	target vs. base	< 0.001
NS	target vs. base	< 0.001
Brown	target vs. base	< 0.001

Table 8: Significance of temperature-scaled surprisal for GPT2 small on three corpora with $T = T^*$.



Figure 8: Distribution of negative log actual-word probability (surprisal) before (left side of figure) and after (right side of figure) temperature scaling for single-token and multiple-token words for GPT-2 small on three corpora. Values of surprisal with probability of 0.1, 0.01 and 1/K (random guessing) are displayed using dash lines.

1153

1154

1155

- 1156 1157 1158 1159 1160 1161 1162 1163 1164
- 1165
- 1166 1167
- 1168
- 1169
- 1170

1172 1173

1174 1175

1175

1177

1178

1179

$\begin{aligned} rt &\sim \text{surprisal} + \text{surprisal_prev_1} + \text{surprisal_prev_2} \\ &+ \text{freq} + \text{length} + \text{freq_prev_1} + \text{length_prev_1} \\ &+ (1|\text{article}) + (1|\text{subj_id}). \end{aligned}$ the ratio of that category words among the total words: $\bar{\Delta}_{\text{MSE}}(F) = \Delta_{\text{MSE}(F)} \cdot \text{ratio}(F)$. Table 11 shows that **multiple-token words drive the much** stronger averaged benefit of temperature-scaled stronger averaged s

(32)

 $rt \sim surprisal + surprisal_prev_1 + surprisal_prev_2L$ + freq + length + freq_prev_1 + length_prev_1 + (1|article) + (surprisal|subj_id).

(33)

The results are in Table 9. **Removing the interactions among predictors or additionally including random slopes does not influence the effectiveness of temperature-scaled surprisal.**

Furthermore, we also investigated the correlations among predictors by examining the correlation matrix for GPT2 small on three corpora (model 1). Table 9, 10 and 11 indicate that **temperaturescaled surprisal does not exhibit a stronger correlation with the other predictors in comparison to the original surprisal**, as shown in the surprisal column ('surp'), which excludes the concern that the primary benefits are simply due to correlations between the baseline predictor and temperaturescaled surprisal.

J Influence of multiple-token words vs. model size

Table 10 shows the increase of Δ_{llh} of temperaturescaled surprisal by only taking into the analysis the subset of multiple-token words. The benefit of temperature-scaled surprisal being primarily from the scaling of multiple-token words still holds for larger LLMs. For larger LLMs, the influence of multiple-token words is also larger.

K Influence of word-level attributes vs. influence of multiple-token words

We explore which of these two factors has a 1180 stronger effect on the benefit of temperature-scaled 1181 surprisal, word-level attributes in Section 6.2 or 1182 multiple-token words in Section 6.3. For word 1183 types, we select named entities as the representa-1184 tive attribute since they perform to be the most 1185 beneficial ones as discussed in Section 6.2. For 1186 multiple-token words, we select all multiple-token 1187 words with more-than-one tokens. In order to fairly 1188 compare the influence, we normalize Δ_{MSE} of 1189

each category under the linguistic factor F with1190the ratio of that category words among the total1191words: $\overline{\Delta}_{MSE}(F) = \Delta_{MSE(F)} \cdot ratio(F)$. Table 111192shows that multiple-token words drive the much1193stronger averaged benefit of temperature-scaled1194surprisal, compared with the averaged benefit of1195named entities.1196

L Other results in Section 6

Corpora	Models	T^*	$\Delta_{\rm llh}(T=1)$	$\Delta_{\rm llh}(T=T^*)$	$\Delta_{\text{llh}}+$
Dundee	model1	2.75	6.90	8.45	22.5
Dundee	model2	2.75	6.79	8.12	19.6
Dundee	model3	2.75	7.81	9.12	16.8
Natural Stories	model1	2.5	4.36	6.99	60.3
Natural Stories	model2	2.5	4.35	6.99	60.7
Natural Stories	model3	*	*	*	*
Brown	model1	2.5	6.62	7.53	13.7
Brown	model2	2.25	6.62	7.30	10.3
Brown	model3	*	*	*	*

Table 9: Optimal T^* , $\Delta_{\text{llh}}(T = 1)$, $\Delta_{\text{llh}}(T = T^*)$, and Δ_{llh} + for three models for GPT2 small on three corpora. * indicates regression models not converged.

	(Intr)	surp	surp_1	surp_2	log_frq	length	log_frq_1	length_1	log_frq_2
surp	0.004								
surp_1	0.000	-0.147							
surp_2	-0.001	-0.057	-0.101						
log_frq	0.0200	0.238	0.002	-0.03					
length	0.019	-0.272	0.027	0.04	0.602				
log_frq_1	0.022	-0.085	0.332	-0.048	0.034	-0.021			
length_1	0.028	0.034	-0.200	0.031	0.003	-0.025	0.650		
log_frq_2	0.032	-0.081	0.002	0.000	0.374	0.626	-0.009	0.014	
length_2	0.038	-0.013	-0.033	0.003	-0.003	0.043	0.509	0.578	0.014
(a) $T = 1$									
	(Intr)	surp	surp_1	surp_2	log_frq	length	log_frq_1	length_1	log_frq_2
surp	0.006								
surp_1	0.005	-0.145							
surp_2	-0.003	-0.074	-0.154						
log_frq	0.020	-0.055	0.050	-0.013					
length	0.017	-0.395	0.042	0.044	0.676				
log_frq_1	0.024	-0.058	0.063	0.011	0.051	-0.018			
length_1	0.025	0.060	-0.353	0.075	-0.016	-0.035	0.702		
log_frq_2	0.031	-0.156	0.004	0.004	0.409	0.634	-0.005	0.011	
length_2	0.037	0.001	-0.088	-0.006	-0.003	0.038	0.542	0.574	0.014
				(b)	$T = T^*$				

Figure 9: Correlation matrix for GPT2s on Dundee with (a) T = 1 and (b) $T = T^*$.

	(Intr)	surp	surp_1	surp_2	log_frq	length	log_frq_1	length_1	log_frq_2
surp	0.002								
surp_1	0.002	-0.009							
surp_2	0.001	0.003	-0.019						
log_frq	0.017	0.237	0.013	-0.016					
length	0.022	-0.181	0.018	0.011	0.692				
log_frq_1	0.018	0.019	0.238	-0.051	0.067	-0.015			
length_1	0.022	0.013	-0.183	0.029	-0.01	0.011	0.672		
log_frq_2	0.030	0.005	0.030	0.010	0.472	0.586	0.008	0.017	
length_2	0.030	0.010	0.011	0.018	-0.005	0.02	0.468	0.589	0.023
(a) $T = 1$									
	(Intr)	surp	surp_1	surp_2	log_frq	length	log_frq_1	length_1	log_frq_2
surp	0.013								
surp_1	0.011	-0.108							
surp_2									
	-0.002	-0.034	-0.080						
log_frq	-0.002 0.020	-0.034 0.200	-0.080 0.009	-0.021					
log_frq length	-0.002 0.020 0.020	-0.034 0.200 -0.194	-0.080 0.009 0.014	-0.021 0.010	0.700				
log_frq length log_frq_1	-0.002 0.020 0.020 0.019	-0.034 0.200 -0.194 -0.09	-0.080 0.009 0.014 0.231	-0.021 0.010 -0.026	0.700 0.048	0.001			
log_frq length log_frq_1 length_1	-0.002 0.020 0.020 0.019 0.020	-0.034 0.200 -0.194 -0.09 0.016	-0.080 0.009 0.014 0.231 -0.203	-0.021 0.010 -0.026 0.045	0.700 0.048 -0.013	0.001 0.014	0.667		
log_frq length log_frq_1 length_1 log_frq_2	-0.002 0.020 0.020 0.019 0.020 0.031	-0.034 0.200 -0.194 -0.09 0.016 0.035	-0.080 0.009 0.014 0.231 -0.203 0.004	-0.021 0.010 -0.026 0.045 -0.007	0.700 0.048 -0.013 0.482	0.001 0.014 0.578	0.667 0.000	0.020	
log_frq length log_frq_1 length_1 log_frq_2 length_2	-0.002 0.020 0.020 0.019 0.020 0.031 0.031	-0.034 0.200 -0.194 -0.09 0.016 0.035 0.015	-0.080 0.009 0.014 0.231 -0.203 0.004 0.038	-0.021 0.010 -0.026 0.045 -0.007 -0.036	0.700 0.048 -0.013 0.482 -0.003	0.001 0.014 0.578 0.019	0.667 0.000 0.474	0.020 0.579	0.023

Figure 10: Correlation matrix for GPT2s on Natural Stories with (a) T = 1 and (b) $T = T^*$.

	(Intr)	surp	surp_1	surp_2	log_frq	length	log_frq_1	length_1	log_frq_2
surp	0.003								
surp_1	-0.003	-0.058							
surp_2	-0.001	-0.021	-0.039						
log_frq	0.032	0.269	0.007	-0.053					
length	0.036	-0.206	0.005	-0.007	0.691				
log_frq_1	0.007	-0.068	0.212	-0.044	0.084	0.021			
length_1	0.012	-0.018	-0.379	0.036	0.022	0.060	0.484		
log_frq_2	0.045	-0.003	0.000	-0.009	0.539	0.593	-0.013	0.016	
length_2	0.028	-0.019	-0.09	0.018	0.020	0.054	0.247	0.347	-0.012
(a) $T = 1$									
	(Intr)	surp	surp_1	surp_2	log_frq	length	log_frq_1	length_1	log_frq_2
surp	0.019								
surp_1	-0.010	-0.114							
surp_2	-0.002	-0.046	-0.096						
log_frq	0.035	0.165	0.010	-0.049					
length	0.032	-0.241	0.027	-0.010	0.719				
log_frq_1	0.008	-0.103	-0.124	0.011	0.078	0.034			
length_1	0.015	0.019	-0.572	0.079	0.018	0.043	0.580		
log_frq_2	0.045	0.015	-0.024	-0.023	0.554	0.584	-0.012	0.026	
length_2	0.029	0.009	-0.263	0.008	0.022	0.046	0.295	0.418	-0.005
$\frac{10020}{1000} -0.205 -0.205 -0.000 -0.022 -0.040 -0.275 -0.410 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005 -0.005$									

Figure 11: Correlation matrix for GPT2s on Brown with (a) T = 1 and (b) $T = T^*$.

	GPT2	$\Delta_{\mathbf{llh}} + (\mathrm{multiple})$
	s	23.6
Dundaa	m	36.4
Dunuee	1	38.0
	xl	42.9
	S	45.2
NC	m	50.1
142	1	62.0
	xl	67.8
	s	9.2
Drown	m	13.4
Brown	1	17.9
	xl	5.49

Table 10: $\Delta_{\rm llh}$ improvement by only scaling tokens in multiple-token words (%) ($\Delta_{\rm llh}$ + (multiple) = $(\Delta_{\rm llh}(T = T^*, \text{multiple}) - \Delta_{\rm llh}(T = 1))/\Delta_{\rm llh}(T = 1))$ for GPT2s on Dundee, Natural Stories (NS) and Brown.

	GPT2	NE	#>1
	s	3.9	17.0
Dundaa	m	6.9	26.7
Dunuee	1	7.2	27.0
	xl	7.6	27.9
	s	2.6	35.9
NS	m	2.2	38.4
143	1	2.1	43.3
	xl	2.0	40.6
	S	10.2	27.0
Drown	m	9.8	28.9
DIOMI	1	10.1	30.7
	xl	10.8	36.0

Table 11: $\overline{\Delta}_{MSE}$ measurement on named entites (NE) and multiple-token words (#>1) for GPT-2 models on Dundee, Natural Stories (NS) and Brown.

	ratio of $p_{w_t} \downarrow$				ratio of named entities			
	#=1	#>1	#=2	#=3	#=1	#>1	#=2	#=3
Dundee	87.6	93.7	90.6	98.3	3.7	16.3	16.6	17.4
Natural Stories	92.1	93.0	92.2	97.2*	1.3	3.5	3.3	4.7*
Brown	93.0	98.1	97.6	35.2*	3.3	12.3	10.9	17.0*

Table 12: This table displays the ratio of words with decreasing probability $(p_{w_t}\downarrow)$ and the ratio of named entities on subsets for both single-token words (#=1) and multiple-token words (#>1) for GPT-2 small on three corpora. Numbers marked with * indicate subsets with insufficient (less than 1%) data.

	#=1		#>1		#=2		#=3	
	$p_{w_t} \downarrow$	p_{w_t}	$p_{w_t} \downarrow$	p_{w_t}	$p_{w_t}\downarrow$	$p_{w_t}\uparrow$	$p_{w_t}\downarrow$	$p_{w_t}\uparrow$
Dundee	8.0	19.6	269.5	-20.3*	50.5	26.6*	497.4	125.4**
NS	117.3	142.3	242.5	93.0*	312.6	95.8*	-123.9*	50.6**
Brown	35.2	-61.0	327.3	5290.2**	17.3	5290.2**	655.0*	0.0**

Table 13: Given words with decreasing (and increasing) probability, the corresponding $\Delta_{MSE}(p_{w_t}\downarrow)$ (and $\Delta_{MSE}(p_{w_t}\uparrow)$) measurement for both single-token words (#=1) and multiple-token words (#>1) for GPT-2 small on three corpora. Numbers marked with * indicate subsets with insufficient (less than 1%) data. Numbers marked with ** indicate subsets with super insufficient (around or less than 0.1%) data.