Lexical Sophistication and Zero-Shot Topic Modeling: Examining the Intersection of Word Choice and NLP Performance

Anonymous ACL submission

Abstract

Lexical choice-the selection of specific 2 words to convey meaning-plays a crucial 3 role in both human communication and 4 natural language processing (NLP). While 5 traditional topic modeling methods like 6 Latent Dirichlet Allocation (LDA) rely on 7 frequency and word co-occurrence 8 zero-shot topic patterns, modeling 9 leverages pre-trained language models to 10 classify unseen data without task-specific 11 training, making them inherently sensitive 12 to lexical choices. This study investigates 13 how variations in lexical sophistication 14 impact zero-shot topic modeling, focusing 15 on potential biases in topic classification. 16 Using the AG News dataset, original texts 17 were paired with paraphrased versions 18 generated by the PEGASUS model, and 19 lexical sophistication was measured 20 quantitatively. Analysis of RoBERTa's 21 predictions revealed topic moderate 22 sensitivity to lexical changes, with a 23 Lexical Bias Score (LBS) of 0.52. Instances 24 of topic shifts between original and 25 paraphrased texts further highlighted the 26 model's occasional misinterpretation of 27 context due to subtle lexical differences. 28 This study enhances our understanding of 29 how language models process lexical 30 sophistication, offering insights into 31 computational linguistics and 32 psycholinguistic theories. The findings 33 underscore the need for continuous 34 evaluation of pre-trained models to mitigate 35 biases and improve fairness in NLP 36 applications. Future research will explore 37 cross-linguistic analyses, model 38 comparisons, and the integration of human 39 judgments to deepen the study of lexical 40 sophistication in zero-shot learning 41 contexts. 42

43 1 Introduction

44 The rapid advancement of Natural Language 45 Processing (NLP) has led to the development of 46 sophisticated language models capable of 47 performing complex tasks with minimal 48 supervision. Zero-shot topic modeling, in 49 particular, has emerged as a powerful approach, 50 enabling models to classify and assign topics to 51 texts without task-specific training. This capability ⁵² is made possible by large-scale pre-trained models 53 like RoBERTa (Liu, 2019), which leverage 54 extensive corpora to generalize across diverse 55 linguistic contexts.

This study examines the sensitivity of zero-shot 56 57 models such as RoBERTa-large-mnli to lexical 58 sophistication—a dimension of lexical choice that 59 includes word complexity (Palfreyman & Karaki, 60 2019), difficulty (Vitta, Nicklin, & Albright, 2023), 61 and diversity (Baese-Berk, Drake, Foster, 62 Lee, Staggs, & Wright, 2021). Lexical 63 sophistication is integral to human language 64 comprehension (Liu & Dou, 2023), as humans 65 effortlessly process synonyms, paraphrases, and 66 nuanced linguistic variations. However, whether 67 state-of-the-art zero-shot models exhibit similar 68 adaptability remains an open question. This 69 research investigates whether changes in lexical 70 sophistication influence topic classification 71 outcomes, potentially revealing biases in model 72 predictions.

To address this, the study introduces the Lexical Bias Score (LBS) to quantify model sensitivity to Is lexical variation. By examining the relationship between lexical sophistication and topic classification accuracy, this research contributes to the broader discourse on fairness and robustness in NLP, underscoring the need to consider linguistic of diversity when developing and deploying language models in real-world applications.

82 1.1 Research Problem

129 83 Despite the widespread use of zero-shot topic 130 84 modeling in various applications, there is limited 131 ⁸⁵ understanding of how these models handle lexical 86 variation, particularly in terms of lexical 87 sophistication. Unlike classic topic modeling 133 88 methods that rely on word frequency within a 135 ⁸⁹ specific corpus, zero-shot models apply pre-trained 90 knowledge to unseen data, making them more 136 ⁹¹ susceptible to subtle lexical sophistication 92 differences. Preliminary findings suggest that 138 93 measurable changes in linguistic features can 139 94 significantly affect topic predictions (Liu & Guo, 141 95 2021).

This sensitivity raises concerns about the 142 96 143 97 reliability and fairness of zero-shot models, 98 especially in contexts where diverse linguistic 144 ⁹⁹ inputs are prevalent. For instance, in content ¹⁴⁵ 100 moderation, automated hiring systems, or 146 2 ¹⁰¹ multilingual NLP applications, inconsistent topic 102 predictions due to lexical sophistication biases 147 2.1 103 could lead to unfair outcomes and reduced model ¹⁰⁴ performance. Understanding the extent to which ¹⁴⁸ 149 105 lexical sophistication impacts zero-shot models is 106 crucial for ensuring robustness, fairness, and 150 151 107 accuracy in real-world NLP applications.

108 1.2 Objectives of the Study

The primary objective of this study is to evaluate
 the impact of lexical choices, particularly
 variations on lexical sophistication, on zero-shot
 topic modeling performance. Specifically, we aim
 to:

Analyze how zero-shot models like
 RoBERTa respond to variations in lexical
 sophistication.

Investigate the sensitivity of the zero-shot
 Investigate the sensitivity of the zero-shot
 topic modeling to lexical sophistication by
 introducing the Lexical Bias Score (LBS).

120 1.3 Significance of the Study

167 This research contributes to both computational 121 linguistics and psycholinguistics by providing a 122 deeper understanding of how language models 123 170 lexical choices and lexical process 124 171 sophistication. The integration of the Tool for 125 172 Analysis of the Automatic Lexical 126 173 Sophistication (TAALES)¹ offers a quantitative 127 174 lens through which lexical variation can be measured, enhancing our ability to evaluate model sensitivity and performance. The introduction of the Lexical Bias Score (LBS), alongside lexical sophistication metrics, provides a comprehensive toolset for assessing model robustness, fairness, and susceptibility to linguistic variation, with potential applications in model development, bias mitigation, and ethical AI practices.

From a psycholinguistic perspective, this study seeks to shed light on the parallels and divergences between human and machine language processing. By examining how zeroshot models handle lexical sophistication, we gain insights into their semantic flexibility and context sensitivity.

2 Prior Work

152

153

166

128

2.1 Zero-shot Learning

Zero-shot Learning (ZSL) has emerged as a transformative approach in Natural Language Processing (NLP), enabling models to perform tasks without task-specific training. Unlike traditional supervised methods that rely on labeled datasets, zero-shot models utilize pretrained knowledge to generalize across new, unseen tasks. This paradigm shift has been facilitated by large-scale language models which leverage extensive corpora and advanced training techniques to capture rich linguistic patterns and contextual nuances.

In zero-shot topic modeling, models classify texts into predefined categories without direct exposure to labeled examples for those categories. This approach has proven effective in various applications, including content moderation, document classification, and sentiment analysis. However, despite its strengths, zero-shot learning introduces unique challenges, especially regarding lexical choice and lexical sophistication. Unlike classic topic modeling approaches such as Latent Dirichlet Allocation (LDA), which rely on word frequency within a given corpus, zero-shot models depend on pre-trained embeddings, making them inherently more sensitive to subtle

¹

https://www.linguisticanalysistools.
org/taales.html

lexical changes. The absence of task-specific 223 175 fine-tuning amplifies the impact of lexical 224 176 variations, as models must rely solely on their 225 177 pre-existing linguistic knowledge. This 226 178 sensitivity raises questions about robustness, 227 179 particularly when faced with shifts in word 228 180 complexity, frequency, and diversity-elements 181 that are often overlooked in traditional models. 182 Existing studies on zero-shot learning have 230 183 largely focused on model architecture, training 231 184 efficiency, and performance across tasks (Wang, 232 185 Zheng, Yu, & Miao, 2019), but few have 233 186 addressed the impact of lexical sophistication on 234 187 zero-shot predictions (Lee, Cai, Meng, Wang, & 235 188 Wu, 2024). This gap highlights the need for 236 189 further investigation, particularly in 237 190 understanding how lexical variation influences 238 191 model classification accuracy. 192 239 193 240

194 2.2 RoBERTa for Zero-Shot Topic 241 195 Modeling 243 243 243 243

The RoBERTa-large-mnli model², developed by 244 196 Facebook AI and implemented through the 245 197 HuggingFace Transformers library, is a fine- 246 198 version of the RoBERTa tuned large 247 199 architecture, trained on the Multi-Genre Natural 248 200 Language Inference (MNLI) corpus. This 249 201 transformer-based model leverages masked 250 202 language modeling (MLM) during pretraining, 251 203 using a large and diverse corpus including 252 204 BookCorpus, Wikipedia, CC-News, 253 205 OpenWebText, and Stories, ensuring broad 254 206 linguistic exposure robustness and in 255 207 understanding contextual relationships. 256 208 The RoBERTa-large-mnli model excels in 257 209 zero-shot classification by reframing 258 210 classification tasks as Natural Language 259 211 Inference (NLI) problems. Given an input text 260 212 (premise) and a candidate label (hypothesis), it 261 213 evaluates the likelihood of the label being 262 214 applicable, making it highly adaptable to unseen 263 215 classification tasks. Its training on the MNLI 264 216 dataset, a benchmark for NLI tasks, allows it to 265 217 perform with high accuracy (90.2% on MNLI), 266 218 making it a reliable choice for zero-shot 267 219 applications. 220

This model's broad training data and dynamic 269 masking during pretraining make it particularly 270 sensitive to lexical choices—a crucial factor in this study's investigation of lexical sophistication. However, its reliance on unfiltered internet data may also introduce potential biases (Chae, & Davidson, 2023), particularly in handling diverse lexical inputs.

229 2.3 Lexical Choice and Lexical 230 Sophistication in NLP

Lexical choice, the selection of specific words and phrases to convey meaning, has been a critical area of study in both human language processing and machine learning applications. In psycholinguistics, lexical choices are influenced by context, audience, and cultural background, affecting how messages are interpreted (Kecskes & Cuenca, 2005). Within NLP, sensitivity to different lexical choices can significantly impact model performance, fairness, and robustness (Wang, Wang, & Yang, 2021).

In the realm of lexical analysis, several approaches have been used to understand aspects of lexical choice. WordNet via NLTK been employed to has explore word relationships and synonymy, shedding light on why one synonym might be chosen over another based on subtle semantic differences (Edmonds, & Hirst, 2002). Large corpora such as the Corpus of Contemporary American English (COCA) and the British National Corpus (BNC) have provided insights into word frequency and collocations, illustrating how the availability and familiarity of certain words influence the pool of lexical choices (Balota, & Chumbley, 1984). Additionally, LIWC (Linguistic Inquiry and Word Count) has offered frameworks for analyzing emotional and cognitive content in text, emphasizing how lexical choices are often driven by the degree of emotion or psychological intent a message aims to convey.

While these tools primarily address broad lexical variations, they underscore the complexity behind word selection processes. Building on these perspectives, lexical sophistication serves as a specific dimension of lexical choice. Metrics such as word frequency, lexical diversity, and rarity—captured by tools like TAALES (Tool for the Automatic Analysis

²

https://huggingface.co/FacebookAI/ro
berta-large-mnli

of Lexical Sophistication)—reflect 271 272 choices, providing a quantifiable approach to 318 quality paraphrased outputs. 273 evaluating lexical variation in terms of linguistic 319 274 sophistication. Studies using TAALES have 320 paraphrase generation process: 275 demonstrated lexical sophistication 321 that 276 277 278 279 Heintz, Choi, Batchelor, Karimi, & Malatinszky, 325 without truncation. 280 2023). 281

3 Methods 282

283 This study is grounded in the intersection of zero-284 shot learning and psycholinguistic theories of 285 lexical sophistication.

286 3.1 **Dataset Selection**

287 This study utilized the AG News dataset, a widely 288 recognized benchmark for topic classification, 289 accessed via the Hugging Face Datasets library. The dataset comprises news articles categorized 290 into four topics: World News, Sports, Business, and 291 Technology. Its established use in numerous NLP 292 ²⁹³ studies ensures that the dataset provides a reliable ²⁹⁴ foundation for evaluating zero-shot topic modeling 295 performance. The dataset's balanced structure and 342 3.3 ²⁹⁶ diverse content make it an appropriate choice for ³⁴³ 297 assessing the impact of lexical sophistication on 298 model predictions

299 3.2 **Paraphrase Generation**

300 The texts from the AG News dataset (7,600 texts with 4 categories and 1,900 texts per category) ³⁰² were paraphrased. Paraphrasing was implemented 303 to capture lexical sophistication variations while 304 maintaining the original semantic content. 305 PEGASUS (Pre-training with Extracted Gap-306 sentences for Abstractive Summarization) model, a 307 state-of-the-art transformer-based model designed 354 3.4 308 for text generation tasks such as summarization and 309 paraphrasing was selected for its superior 310 performance in maintaining semantic integrity 311 while introducing lexical variations, making it 312 well-suited for analyzing the impact of lexical 313 choices on zero-shot topic modeling. The 314 tuner007/pegasus paraphrase ³ model, 315 available via the Hugging Face Transformers

the 316 library, was employed. The model was fine-tuned complexity and deliberation behind word 317 specifically for paraphrasing tasks, ensuring high-

The following parameters were set during the

max length=300: This value aligns with the influences readability, comprehension, and even 322 average text length in the AG News dataset NLP model performance, particularly in text 323 (approximately 300 characters), ensuring that classification and sentiment analysis (Crossley, 324 paraphrased sentences retain essential information

> num return sequences=1: Each input 326 327 sentence was paraphrased once to avoid multiple 328 paraphrase options that could introduce additional variability. 329

> temperature=1.5: А relatively high 330 331 temperature value was chosen to encourage more ³³² diverse word choices while maintaining coherence. ³³³ The choice of PEGASUS was driven by its unique ³³⁴ pre-training objective, where sentences are masked 335 in a manner that simulates summarization, enabling 336 the model to learn contextual dependencies 337 effectively. This capability ensures that 338 paraphrased sentences are lexically diverse yet 339 semantically faithful to the original text, providing 340 a robust basis for investigating the sensitivity of ³⁴¹ zero-shot models to lexical sophistication.

Extracting Lexical Sophistication Measures

344 To quantify the lexical sophistication of both 345 original and paraphrased articles from the AG 346 News dataset, this study employed the Tool for the 347 Automatic Analysis of Lexical Sophistication 348 (Kyle, Crossley, & Berger, 2018). TAALES offers 349 a comprehensive set of over 400 indices related to 350 lexical sophistication, including word frequency, ³⁵¹ lexical diversity, word rarity, and psycholinguistic 352 features, making it an appropriate tool for assessing ³⁵³ variations in lexical sophistication in this research.

Zero-shot Topic Modeling

355 The zero-shot topic modeling in this study was 356 performed using the Hugging Face for 357 transformers pipeline zero-shot 358 classification:

from transformers import pipeline

360

https://huggingface.co/tuner007/pega sus paraphrase

```
classifier = pipeline("zero-shot-
361
362 classification",
                        model="roberta-
363 large-mnli")
```

The roberta-large-mnli assigns the most 364 365 probable topic to each input text based on a ³⁶⁶ provided set of candidate labels. For this study, the 367 four AG News categories—World, Sports, 368 Business, and Technology—were used as the 369 candidate labels. Texts from both the original AG 370 News dataset and paraphrased versions generated 371 using PEGASUS were classified using this 372 pipeline, and the assigned topics were analyzed to 373 explore the influence of lexical sophistication on 374 topic predictions.

375 3.5 **Computational Resources**

376 This study utilized the ROBERTa-large-mnli 377 model (355 million parameters; Liu, 2019) and the 378 PEGASUS paraphrase model (568 million 379 parameters; Zhang et al., 2020) for zero-shot topic 380 modeling and text paraphrasing, respectively. All 381 experiments were conducted on the Acer ³⁸² Supercomputer, equipped with 32 Intel® Xeon® 383 Silver 4208 CPUs @ 2.10GHz, with 220 GB of 384 RAM and no GPU acceleration.

385 was approximately 4 CPU hours, encompassing 412 consistent across all statistical measures, including 386 dataset preprocessing, paraphrasing, 388 sophistication analysis using TAALES 2.2, and 414 histogram in Figure 2 illustrates the distribution of 389 zero-shot classification using Hugging Face 415 text lengths for both sets, emphasizing the shorter Transformers in Python 3.9.12. 390

This high-performance computing infrastructure ⁴¹⁷ 392 facilitated the efficient execution of the study's 418 observed lexical sophistication, as shorter texts 393 experiments, despite constraints such as token 419 inherently limit the variety and complexity of ³⁹⁴ truncation during paraphrasing and the absence of ⁴²⁰ words used. However, this limitation was 395 GPU acceleration, which may have impacted 421 addressed by using TAALES, which extracts 396 processing time.

Results and Discussion 397

398 4.1 **Dataset Exploratory Analysis**

400 from the AG News dataset is visually represented 427 sophistication (Kyle, Crossley, & Berger, 2018) $_{402}$ equally represented (n=1,900).

403 404 followed showed that the text lengths between the 431 principal component (PC1) explaining 47% of the 405 original AG News dataset and the paraphrased 432 variance. This principal component was used as the 406 versions generated using the PEGASUS model 433 lexical sophistication measure for both datasets. 407 reveals a significant reduction in length post- 434 PC1, which depicts lexical richness and syntactic 408 paraphrasing. The mean length of original texts is 435 Complexity, encapsulates the use of elaborate while 409 approximately 237 characters,



Figure 1: Label distribution in AG News Dataset.



Figure 2: Text Length Relative Frequency Distribution of Original and Paraphrased Texts.

410 paraphrased texts have a significantly lower mean The total computational budget for this study 411 of 128 characters. This reduction in length is lexical 413 median, minimum, and maximum lengths. The 416 lengths of the paraphrased texts.

> This disparity in text length may impact the 422 lexical sophistication indices that go beyond 423 surface-level features like text length.

424 4.2 **Lexical Sophistication Measure**

425 To assess lexical sophistication, TAALES was used 399 The label distribution of the 7,600 news articles 426 to extract 484 granular measures of lexical in Figure 1, confirming that each category is 428 for both the original and paraphrased AG News 429 texts. A Principal Component Analysis (PCA) was The paraphrasing of original news articles that 430 performed to reduce dimensionality, with the first the 436 (academic) language and complex sentence

437 structures. The indices that loaded highly, i.e. 438 |factor loadings| > 0.50, on the principal 439 component highlight verbosity, syntactic depth, 440 and detailed clause structures, reinforcing its role 441 as a robust measure of lexical sophistication

Across all four categories, the results indicate Ata that paraphrased texts exhibit slightly lower principal component score (also referred to as *lexical sophistication measure*) to original texts, the though the differences exhibit only marginal significance. A paired t-test was conducted to ata evaluate the statistical significance of the difference in lexical sophistication between difference in lexical sophistication between to original and paraphrased texts, overall. The results decrease in lexical sophistication in paraphrased texts compared to the original texts, overall (t(7599) = 2.98, p = 0.053).

In the World News category, original texts had 455 456 a mean lexical sophistication measure of 0.130, 457 while paraphrased texts scored 0.108, (t(1900) = $_{458}$ 2.54, p = 0.052), reflecting a non-significant slight 459 reduction in lexical sophistication. For the Sports 460 category, original texts had a lexical sophistication ⁴⁶¹ measure of 0.120, and paraphrased texts had 0.097, $_{462}$ (t(1900) = 2.87, p = 0.050), indicating a marginally ⁴⁶³ significant decrease in lexical sophistication. In the 464 Business category, the mean lexical sophistication 465 scores were 0.148 for original texts and 0.125 for 466 paraphrased texts, (t(1900) = 2.74, p = 0.057), 467 showing a slight but not statistically robust 468 reduction in lexical sophistication. Lastly, the 469 Technology category demonstrated a non-470 significant decrease, with original texts scoring 506 471 0.122 and paraphrased texts scoring 0.099, (t(1900) $_{472} = 2.68$, p = 0.052). These findings suggest that 473 while paraphrased texts maintain overall semantic 474 integrity, they exhibit slightly lower lexical 475 sophistication across categories. It should be noted, 476 however, that the text length of the paraphrases 477 were shorter than most of the original texts which 478 could also be a factor that contributes to the 479 reduction in lexical sophistication. However upon 480 post hoc qualitative inspection after paraphrasing, 481 the reduction of lexical sophistication can be 482 readily observed across random samples.

483 4.3 Zero-shot Topic Modeling

⁴⁸⁴ The zero-shot topic modeling in this study was ⁴⁸⁵ performed using the ROBERTa-large-mnli ⁴⁸⁶ model, implemented through the Hugging Face ⁴⁸⁷ transformers library. This model, pre-trained on



Figure 3: Confusion Matrix of Original Text Topic Classifications.

⁴⁸⁸ large-scale Natural Language Inference (NLI)
⁴⁸⁹ tasks, is particularly suitable for zero-shot
⁴⁹⁰ classification due to its robust semantic
⁴⁹¹ representations and adaptability to unseen tasks.
⁴⁹² The classification was executed using the pipeline
⁴⁹³ API, with the following candidate labels provided
⁴⁹⁴ for each text input: World News, Sports, Business,
⁴⁹⁵ and Technology.

The topic (or classification) distributions for 497 both the original and paraphrased texts are 498 presented through the confusion matrices in 499 Figures 3 and 4. The original dataset displayed the 500 following distribution across categories: World 501 News (1330), Sports (1045), Business (570), and 502 Technology (4655). The paraphrased dataset 503 distribution is World News (633), Sports (1551), 504 Business (665), and Technology (4751).

505 4.3.1 Model Performance and Analytical 506 Focus

507 The zero-shot topic modeling yielded an accuracy 508 of 48.75% and 38.07% on the original texts and the ⁵⁰⁹ paraphrased texts respectively, with their confusion ⁵¹⁰ matrices presented in Figures 3 and 4. While these provide insight into the model's 511 metrics 512 classification performance, it is essential to 513 highlight that the focus of this study is not on the 514 overall accuracy or predictive performance of the 515 model. Instead, the primary objective is to 516 investigate how lexical sophistication and 517 paraphrasing influence topic classification, 518 particularly through observable topic shifts. Topic 519 shift is defined in this study as the change in the 520 predicted topic when a text is paraphrased, such 521 that the topic classification of the original text 522 differs from that of its paraphrased counterpart. 523 This phenomenon reflects how lexical alterations ⁵²⁴ influence the semantic interpretation of text by the



Figure 4: Confusion Matrix of Paraphrased Text Topic Classifications.

526 different topics to semantically related content. In 577 measured by the lexical sophistication measure. 527 addition, upon qualitative inspection of randomly 578 For the subset of predictions where topic labels 528 sampled instances, it was found that the observed 579 changed, the paraphrased texts consistently shifts-although 529 topic 530 misclassifications from а 531 perspective—were often semantically plausible. 582 futures and quarterly earnings reports was 532 For example, sentences initially labeled as World 583 classified as World News with a lexical 533 News that were paraphrased and subsequently 584 sophistication score of 0.142. Its paraphrased 534 classified as Business often reflected content that 585 version, which omitted references to external 535 straddled both domains. This suggests that the topic 586 economic factors, was classified as Business with 536 shifts captured by the model were not arbitrary but 587 a reduced score of 0.109. Another example is an 537 contextually coherent, underscoring 538 complexity and nuance of lexical choice in 589 initially labeled as World News with a score of ⁵³⁹ influencing semantic interpretation by language ⁵⁹⁰ 0.135, which was reclassified as Technology when 540 models. Therefore, while the reported accuracy 591 paraphrased, reflecting a decrease in lexical ⁵⁴¹ may seem low to moderate, it does not detract from ⁵⁹² sophistication to 0.102. 542 the central findings of this study, which emphasize 593 543 the role of lexical sophistication in shaping topic 594 lexical alterations can influence topic predictions, 544 classification decisions in zero-shot settings.

545 546 the sensitivity of the RoBERTa model to lexical 597 mean lexical sophistication score dropped from 547 changes introduced through paraphrasing. The 598 0.128 in original texts to 0.104 in paraphrased texts Technology category maintained a dominant ⁵⁹⁹ (t = 2.89, p = 0.004). This reinforces that lexical 549 presence in both distributions, suggesting that 600 sophistication significantly impacts zero-shot topic 550 lexical variations had minimal impact on the 601 modeling predictions, highlighting the sensitivity 551 model's classification for this category. However, 602 of models like ROBERTa-large-mnli to 552 notable shifts were observed in the Sports and 603 variations in lexical richness. 553 World News categories, with the former 604 554 experiencing an increase in paraphrased texts and 605 the paraphrased texts did not exhibit a decrease in 555 the latter a decrease.

These results suggest that paraphrasing not only 607 example below: 556 557 affects lexical sophistication but also influences the 608 Original Text predicted as Technology: 558 model's semantic interpretations, leading to 609 PeopleSoft's big bash See you next year in Las Vegas 559 variations in topic classification. The consistent 610, proclaimed a marquee at the PeopleSoft user 560 dominance of the Technology category could 611 conference in San Francisco in late September. It was ⁵⁶¹ indicate a more defined and stable lexical profile ⁶¹² one of many not-so-subtle attempts by the company to ⁵⁶² within this domain, while the fluctuations in the ⁶¹³ reassure its customers.

⁵⁶³ other categories underscore potential lexical biases ⁵⁶⁴ in zero-shot models when processing linguistically 565 varied inputs. In this study, stable refers to instances where there is no significant difference in lexical sophistication measure (PC1), 567 the 568 indicating that the lexical richness, syntactic 569 complexity, and overall linguistic depth remain 570 consistent between the original text and its 571 paraphrased counterpart, despite changes in 572 wording.

Investigating Topic Shifts 573 **4.4**

574 The topic shifts between original and paraphrased 575 texts were accompanied by statistically significant 525 language model, resulting in the assignment of 576 decreases in lexical sophistication scores, as technically 580 exhibited lower lexical sophistication. For ground-truth 581 instance, an original text discussing US stock the 588 original text on vehicle stability control systems,

These examples illustrate how even subtle ⁵⁹⁵ particularly when contextual richness is reduced. This divergence in topic distribution highlights 596 Across all analyzed pairs with topic shifts, the

> In some instances where topic shifts occurred, 606 lexical sophistication, as can be observed in the

614

615 Paraphrased Text predicted as Business:

616 617 September, a marquee proclaimed, "See you next 666 large-mnli, particularly when lexical variations 618 year in Las Vegas." It was one of many not-so- 667 are introduced through paraphrasing. 619 subtle attempts by the company to assure its 668 recalibration (topic shift) observed in certain 620 customers.

This suggests that the RoBERTa-large-mnli 670 significant change in lexical 621 622 model may recalibrate its classification based on 671 highlights the model's nuanced sensitivity to 623 subtle lexical changes that maintain complexity but 672 lexical semantics. This suggests that while lexical alter the focus. For example, paraphrased sentences 673 richness is a critical factor, the lexical semantics of 625 that retained intricate structures but shifted 674 word choices also influence model predictions, 626 thematic emphasis—such as removing specific 675 which is a known strength of transformer-based 627 company names or industry jargon-often resulted 676 topic modeling (Gruetzemacher, & Paradice, 628 in different topic predictions without lowering 677 2022). The stability of classifications within the 629 lexical sophistication scores.

630 4.5 **Investigating Lexical Bias**

631 To further investigate the influence of lexical 632 sophistication on zero-shot topic modeling, this 633 study introduces the Lexical Bias Score (LBS) as a 634 quantifiable metric. The LBS measures the 635 correlation between changes in lexical 636 sophistication and shifts in model predictions.

A comparison of original and paraphrased texts 637 was conducted, with lexical sophistication scores 638 639 (PC1). The LBS was computed as the Pearson 640 correlation coefficient between the differences in 641 lexical sophistication scores and the binary 642 indicator of prediction shifts (0 for consistent predictions, 1 for changes in predictions). 643

The analysis yielded an LBS of 0.52 (p < 0.05), 644 indicating a moderate positive correlation between sophistication variations 646 lexical and topic 647 prediction changes for the entire dataset. This 648 suggests that the ROBERTa-large-mnli model is sensitive to lexical richness, with more lexically 649 sophisticated texts being more likely to retain 650 consistent topic predictions. 651

Table 1 shows that the highest LBS values were 652 653 observed in the Sports and Business categories, 654 suggesting that lexical sophistication plays a 655 critical role in maintaining semantic integrity 656 within these domains. Conversely, the Technology 657 category exhibited the lowest LBS, aligning with 658 previous observations of its stable classification 659 despite paraphrasing.

660 5 Conclusion

661 The findings from this study underscore the 662 intricate relationship between lexical sophistication

and zero-shot topic modeling. The results reveal 664 that lexical sophistication plays a significant role in At the user conference in San Francisco in late 665 influencing the topic predictions of ROBERTA-The 669 paraphrased texts without a corresponding sophistication 678 Technology category, despite paraphrasing, points

Category	Mean LBS	t-value (p-values <0.05)
World News	0.57	2.34
	(sd=0.15)	
Sports	0.64	3.11
	(sd=0.12)	
Business	0.60	2.89
	(sd=0.10)	
Technology	0.39	1.74
	(sd=0.18)	

Table 1: Lexical Bias Score (LBS) Across Categories.

679 to a more consistent lexical profile in this domain, 680 whereas the variability observed in categories like 681 Sports and World News reflects potential lexical 682 biases.

These findings contribute to the growing body 684 of research on lexical sophistication in NLP, 685 emphasizing the need for further exploration of 686 lexical biases in language models (Navigli, Conia, 687 & Ross, 2023). Overall, this study highlights the 688 importance of lexical sophistication in zero-shot 689 learning, demonstrating that lexical choices-690 whether in original or paraphrased texts-can ⁶⁹¹ significantly influence model behavior. The ⁶⁹² insights gained from this research underscore the 693 need for more sophisticated handling of lexical 694 variation in language models to enhance their 695 fairness, robustness, and alignment with human 696 language processing (Bella, Helm, Koch, & 697 Giunchiglia, 2024), (Patil & Gudivada, 2024).

⁶⁹⁹ This study, while contributing to the understanding 700 of lexical sophistication in zero-shot topic 701 modeling, has several limitations that should be 750 acknowledged. 702 751

First, the analysis was conducted using only the 703 704 AG News dataset (7,600 samples) sourced from 705 Hugging Face. Although widely used in text 706 classification research, reliance on a single dataset 755 707 limits the generalizability of the findings to other 756 708 domains, genres, and languages. Future research 709 could explore more diverse datasets to validate and extend these results. 710 759

Second, the PEGASUS paraphrasing model 760 711 712 employed in this study introduces its own 761 Balota, D. A., & Chumbley, J. I. (1984). Are lexical 713 constraints. The model's token limit led to 762 714 truncation of longer texts, which may have affected 763 715 the lexical richness and overall text structure of the 764 716 paraphrased outputs. This truncation potentially 765 717 influenced both lexical sophistication measures 766 Chae, Y., & Davidson, T. (2023). Large language 718 and the zero-shot topic classification results, 767 719 introducing a source of bias that future work should 768 720 address by using paraphrasing tools capable of 721 handling longer text sequences. 770

Additionally, the study's findings rely on the 771 722 723 empirical Lexical Bias Score (LBS), a metric 724 introduced here to quantify the relationship 725 between lexical sophistication and topic 726 classification shifts. While LBS provides initial 727 insights, it is an empirical measure and may not 728 capture all dimensions of lexical bias. Further 729 validation of this metric across different models 778 730 and datasets is necessary to establish its robustness and utility. 731 780

this study focused on lexical Finally, 732 733 sophistication features using TAALES, without 734 incorporating deeper semantic or contextual 783 735 analyses. Future research could integrate more 736 advanced lexical measures, such as sentence and 737 document embeddings that capture broader 786 738 context, along with psycholinguistic frameworks, 787 739 such as Cohesion Network Analysis (McNamara, 788 740 Allen, Crossley, Dascalu, & Perret, 2017), to assess 741 whether topic shifts observed in paraphrased texts 790 742 result from disruptions in textual cohesion or 791 743 changes in lexical relationships. 792

744 Acknowledgments

745 Withheld for anonymity.

746 References

- 747 Baese-Berk, M. M., Drake, S., Foster, K., Lee, D. Y., Staggs, C., & Wright, J. M. (2021). Lexical diversity, lexical sophistication, and predictability for speech in multiple listening conditions. Frontiers in psychology, 12, 661415.
- 752 Bella, G., Helm, P., Koch, G., & Giunchiglia, F. (2024, June). Tackling Language Modelling Bias in Support of Linguistic Diversity. In The 2024 ACM Conference on Fairness, Accountability, and Transparency (pp. 562-572).
- 757 Crossley, S., Heintz, A., Choi, J. S., Batchelor, J., Karimi, M., & Malatinszky, A. (2023). A largescaled corpus for assessing text readability. Behavior Research Methods, 55(2), 491-507.
 - decisions a good measure of lexical access? The role of word frequency in the neglected decision stage. Journal of Experimental Psychology: Human perception and performance, 10(3), 340.
 - models for text classification: From zero-shot learning to fine-tuning. Open Science Foundation.
- 769 Edmonds, P., & Hirst, G. (2002). Near-synonymy and lexical choice. Computational linguistics, 28(2), 105-144.
- 772 Gruetzemacher, R., & Paradice, D. (2022). Deep transfer learning & beyond: Transformer language models in information systems research. ACM Computing Surveys (CSUR), 54(10s), 1-35.
- 776 Kecskes, I., & Cuenca, I. M. (2005). Lexical choice as a reflection of conceptual fluency. International Journal of Bilingualism, 9(1), 49-67.
- 779 Kyle, K., Crossley, S. A., & Berger, C. (2018). The tool for the analysis of lexical sophistication (TAALES): Version 2.0. Behavior Research Methods 50(3), pp. https://doi.org/10.3758/s13428-017-1030-1046. 0924-4
- 784 Lee, S., Cai, Y., Meng, D., Wang, Z., & Wu, Y. (2024, November). Unleashing Large Language Models' Proficiency in Zero-shot Essay Scoring. In Findings of the Association for Computational Linguistics: EMNLP 2024 (pp. 181-198).
- 789 Liu. T., & Guo, W. (2021, December). Topic classification on spoken documents using deep acoustic and linguistic features. In 2021 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU) (pp. 427-432). IEEE. 793
- 794 Liu, Y. (2019). Roberta: A robustly optimized bert pretraining approach. arXiv preprint 795 arXiv:1907.11692, 364. 796

797 Liu, Z., & Dou, J. (2023). Lexical density, lexical diversity, and lexical sophistication in 798 799 simultaneously interpreted texts: a cognitive perspective. Frontiers in Psychology, 14, 1276705. 800

801 McNamara, D. S., Allen, L. K., Crossley, S. A., Dascalu, M., & Perret, C. A. (2017). Natural 802 Language Processing and Learning Analytics. 803

Grantee Submission. 804

805 Navigli, R., Conia, S., & Ross, B. (2023). Biases in large language models: origins, inventory, and 806 discussion. ACM Journal of Data and Information 807

Quality, 15(2), 1-21. 808

alfreyman, D. M., & Karaki, S. (2019). Lexical F 809 sophistication across languages: a preliminary study 810

of undergraduate writing in Arabic (L1) and English

811 (L2). International Journal of Bilingual Education

812 and Bilingualism, 22(8), 992-1015.

813

Patil, R., & Gudivada, V. (2024). A review of current 814 trends, techniques, and challenges in large language 815 models (llms). Applied Sciences, 14(5), 2074. 816

Vitta, J. P., Nicklin, C., & Albright, S. W. (2023). 817

Academic word difficulty and multidimensional 818

English - for lexical sophistication: 819 An

academic - purposes - focused conceptual 820

replication of Hashimoto and Egbert (2019). The 821

Modern Language Journal, 107(1), 373-397. 822

Wang, X., Wang, H., & Yang, D. (2021). Measure and 823

improve robustness in NLP models: A survey. arXiv 824 preprint arXiv:2112.08313. 825

Wang, W., Zheng, V. W., Yu, H., & Miao, C. (2019). A 826

survey of zero-shot learning: Settings, methods, and 827

applications. ACM Transactions on Intelligent 828

Systems and Technology (TIST), 10(2), 1-37. 829