

# Leveraging LEXICAL and GRAMMATICAL Errors: Extending ASR Error Measurements through NLP

Anonymous ACL submission

## Abstract

This paper addresses the limitations of current Automatic Speech Recognition (ASR) evaluation metrics by highlighting the inadequacies of overall error rates, particularly Word Error Rate. While this offers a broad assessment, it lacks the granularity needed to discern specific linguistic categories affected by errors. We offer an NLP-driven metric based on parts of speech and grammatical categories, to provide a more in-depth analysis of ASR errors. Using the Whisper ASR system on English, Japanese, and Spanish, within the CommonVoice 15 dataset, we analyze GRAMMATICAL and LEXICAL error rates. Results show that GRAMMATICAL words trigger less errors than LEXICAL words across all languages, and Proper Nouns in Japanese combined with case markers are related to higher accuracy. By categorizing errors based on these linguistic attributes, our methodology aims to enhance the explanatory power of error analysis in ASR, contributing to a more precise evaluation of system performance based on NLP approaches.

## 1 Introduction

Automatic Speech Recognition (ASR) technologies have undergone significant advancements (O’Shaughnessy, 2023; Reitmaier et al., 2022) and the widespread adoption of ASR systems in various industries (e.g., Healthcare, Defence and Automotive) highlight the critical role of accurate evaluation to ensure their effectiveness, reliability and user satisfaction.

Word Error Rate (WER) is a crucial metric used to evaluate the performance of ASR systems. It measures the operational accuracy of an ASR system by calculating the ratio of the total number

of errors – comprising substitutions, deletions, and insertions in the transcription output – to the number of words in the audio signal input to the ASR system (Kumalija and Nakamoto, 2022).

While WER is used widely as a standard metric (Ali and Renals, 2004; NithyaKalyani and Jothilakshmi, 2019; Park et al., 2023), it has been reported to have some critical limitations (He et al., 2011). The primary limitation of WER lies in treating all errors equally, regardless of their impact on the meaning of the transcribed text. For instance, misrecognizing a key word might change the meaning of a sentence significantly than other non-key words, but WER weighs all the errors the same.

Additionally, WER cannot gauge the relative importance of specific words in the ground truth transcription, prompting the proposal of alternative metrics that account for semantics (Kafle and Huenerfauth, 2017), entity recognition (Garofolo et al., 1998), and parts of speech (Roux et al., 2022).

Prior studies indicate that WER does not consistently correlate with human judgment on ASR system performance (Morris et al., 2004; Whetten and Kennigton, 2023). These findings highlight the necessity for refined linguistic metrics that enable a more granular analysis of errors.

Another critical limitation of the existing metric is its inability to unveil the specific characteristics of errors. Although ASR systems may exhibit similar overall error rates as measured by WER, this metric fails to distinguish between errors affecting different linguistic categories. For example, two ASR systems with equivalent WERs might impact linguistic accuracy differently: one may disproportionately affect GRAMMATICAL words, while the other could affect more frequently LEXICAL or content words. Consequently, a deeper examination of error complexities within linguistic categories is essential to identify and specify areas

81 of vulnerability within ASR systems (Adegbegha  
82 et al., 2024; Errattahi et al., 2018; Kheddar et al.,  
83 2023; Lee et al., 2011; Li et al., 2023).

84 Recognizing the limitations of current  
85 methodologies, we propose the integration of  
86 linguistic metrics into the evaluation of ASR  
87 systems. An in-depth analysis based on linguistic  
88 categories, including parts of speech and  
89 grammatical classifications, enriches our  
90 understanding of error complexities. By  
91 categorizing errors based on linguistic attributes,  
92 we gain valuable insights into the nature of errors  
93 and how they behave within the context of these  
94 systems. This approach not only clarifies the types  
95 of errors but also enhances the explanatory power  
96 of error analysis, providing a more comprehensive  
97 understanding of ASR system performance.

98 Our proposed methodology presents an  
99 approach to analyze and report errors in ASR  
100 outputs. Adopting a multilingual perspective, we  
101 examine errors in English, Japanese, and Spanish,  
102 leveraging the Whisper ASR system (Radford et  
103 al., 2023) on the CommonVoice 15 dataset (Ardila  
104 et al., 2020). Utilizing Parts of Speech tagging  
105 (POS), we differentiate errors into two specific  
106 categories: those affecting grammatical or function  
107 words (referred as GRAMMATICAL error), and those  
108 impacting lexical or content words (referred as  
109 LEXICAL error).

110 The distinction between GRAMMATICAL and  
111 LEXICAL categories facilitates a layered  
112 comparison of ASR errors, examining not only  
113 their aggregate impact but also their specific  
114 manifestation across different linguistic types. This  
115 dual-level analysis enhances our understanding of  
116 ASR errors, significantly addressing the gaps  
117 identified in literature. This approach provides a  
118 more specific and informative perspective of ASR  
119 performance, addressing to the need for detailed  
120 error analysis in the advancing field of speech  
121 recognition technologies.

122 This work significantly advances the reporting  
123 of detailed linguistic layers in ASR systems,  
124 establishing a more consistent methodology that  
125 extends beyond the limited scope of previous  
126 research, which often confined analyses to specific  
127 databases/languages. We propose a systematic  
128 approach applicable across all languages supported  
129 by the ASR system with available universal  
130 dependencies. By developing metrics within a  
131 single widely used ASR system, we enable refined  
132 comparison between the reference text (REF) and

133 the generated hypothesized text (HYP) across  
134 languages with varying typological characteristics.

## 135 **2 Related Work**

### 136 **2.1 Current Progress on Metrics**

137 ASR evaluation methodologies have undergone  
138 some refinements, incorporating diverse error  
139 metrics that surpass mere word counts, including  
140 *word embeddings* (Devlin et al., 2019), *sentence*  
141 *embeddings* (Reimers and Gurevych, 2019), and  
142 *semantic proximity* (Zhang et al., 2020).

143 Taking inspiration from machine translation,  
144 wherein linguistic metrics significantly enhance  
145 translation accuracy, this paper adapts similar  
146 methodologies to ASR. Popović and Ney (2007)  
147 effectively incorporated linguistic attributes, like  
148 parts of speech, into translation evaluation and  
149 introduced the *Position Independent Error Rate*  
150 (PER) to measure the impact of each POS class on  
151 overall word error rates. While their study analyzed  
152 POS-based in two languages, English and Spanish,  
153 and compared them with human assessments, it  
154 remains to be determined whether these findings  
155 can be generalized to other languages with different  
156 typological characteristics.

### 157 **2.2 Main Gaps in Previous Work**

158 Although previous studies have covered relevant  
159 aspects of the error assessment and description in  
160 ASR system outputs, there are yet three significant  
161 gaps remaining. Firstly, these studies have relied on  
162 custom-built or less standardized ASR systems,  
163 limiting the generalizability and reproducibility of  
164 their findings. Our research counters this limitation  
165 by employing the widely recognized and  
166 standardized Whisper ASR system; we ensure that  
167 our findings are more applicable to a broader range  
168 of applications and that our methodology can be  
169 more easily replicated by other researchers in the  
170 field.

171 Secondly, earlier studies have often focused on  
172 analyzing ASR errors in a single language or  
173 closely related languages, which limits the  
174 understanding of ASR performance across diverse  
175 linguistic parameters. Our study expands this scope  
176 by examining ASR errors in three linguistically  
177 differentiated languages – English, Japanese and  
178 Spanish – thus broadening the evaluation of our  
179 proposed error metrics and enhancing  
180 understanding of ASR systems across varied  
181 language families.

182 Lastly, although previous studies have measured 233  
 183 ASR errors at various linguistic levels, they have 234  
 184 not being consistent in proposing a direct 235  
 185 integration of these detailed measurements into the 236  
 186 overall reporting of ASR outputs. Our study  
 187 addresses this by not only detailing these  
 188 measurements but also integrating them with the  
 189 general reporting of ASR errors. This integrated  
 190 approach provides a more comprehensive and  
 191 informative analysis of ASR performance.

### 192 3 Methodology

193 In our methodology, we build upon the linguistic-  
 194 based error metrics found in Cao et al., (2023) and  
 195 Roux et al. (2022), which provide a finer-grained  
 196 analysis of errors and discrepancies. Our approach  
 197 enhances this framework through two distinct  
 198 strategies. First, we conduct a comparative analysis  
 199 of three linguistically diverse languages: English,  
 200 Japanese and Spanish. Each language exhibits  
 201 different levels of linguistic inflections, such as  
 202 changes in word form to mark distinctions such as  
 203 tense, person, and number. For example, verb  
 204 conjugations are a type of inflections and regular  
 205 plurals in English. This comparative study allows  
 206 us to assess how inflectional complexity impacts  
 207 ASR accuracy across different linguistic systems.

208 In the second aspect of our methodology, we  
 209 differentiate between GRAMMATICAL and LEXICAL  
 210 categories, grouping POS categories into these  
 211 classes since their errors have different impacts on  
 212 the **HYP** text. This categorization is crucial  
 213 because LEXICAL errors directly lead to the  
 214 misunderstanding of the intended message and the  
 215 incorrect interpretation of the text (Hemchua and  
 216 Schmitt, 2006). In contrast, while GRAMMATICAL  
 217 errors can also cause misunderstandings, their  
 218 effect on the overall comprehension of the text is  
 219 generally less disruptive than that of LEXICAL  
 220 errors. Table 1 shows examples of GRAMMATICAL  
 221 and LEXICAL errors.

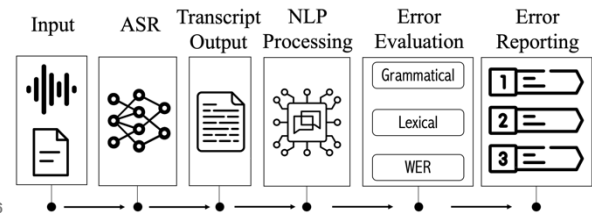
222 In developing performance metrics for ASR  
 223 systems, we adhered to four essential criteria  
 224 (Morris et al., 2004; McCowan et al., 2004). First,  
 225 it should reflect some level of human judgment,  
 226 aiding in the identification of how much  
 227 information is effectively communicated and how  
 228 much is lost. Second, it must be straightforward to  
 229 apply, facilitating comparisons across various ASR  
 230 systems. Third, it should be language-independent  
 231 to ensure unbiased error comparisons across  
 232 languages from different typological

233 classifications. Finally, the metric should be easy to  
 234 interpret from the outputs. These principles ensure  
 235 that our metrics are both practical and applicable in  
 236 real-world settings.

	Words					
REF text	The	cat	sat	on	the	mat
LEXICAL error	The	dog	sat	on	the	mat
GRAM. error	The	cat	sat	on	that	mat

237 Table 1: Examples of Differences between  
 238 GRAMMATICAL and LEXICAL errors.

239 Figure 1 below summarizes the main six stages  
 240 followed in this paper, and these are expanded in  
 241 the following sections. The first three correspond  
 242 to the ASR processing, four and five correspond to  
 243 the NLP processing and categorization of errors,  
 244 and the final stage corresponds to the error  
 245 reporting.



246 Figure 1: Data Processing and Analysis Stages.

#### 248 3.1 Languages Chosen

249 The selection of languages was driven by both data  
 250 availability and the authors’ expertise, resulting in  
 251 the choice of English, Japanese, and Spanish.  
 252 These languages serve as robust testing grounds  
 253 due to their shared characteristics and notable  
 254 differences. Both English and Spanish belong to  
 255 the Indo-European language family, and Japanese  
 256 belongs to the Japonic language family  
 257 (Ethnologue, 1999). They also exhibit divergences  
 258 in their levels of inflection, a factor relevant to ASR  
 259 system errors.

260 Some research has found that word classes with  
 261 higher inflection are more prone to errors  
 262 compared to those with less or no inflection (Berg  
 263 et al., 2024; Smith-Lock, 1991). For instance, the  
 264 English article *the* remains uninflected, while its  
 265 Spanish counterparts carry gender and number  
 266 inflections (feminine singular: “la”, masculine  
 267 singular: “el”, feminine plural: “las,” masculine  
 268 plural: “los”). Additionally, variations in inflection  
 269 levels are evident in verb paradigms. While English  
 270 may have six main isolated forms (base, infinitive,  
 271 past simple, past participle, gerund, and third

person singular) (Lee and Seneff, 2008), Japanese has 12 inflections (Hisamitsu and Nitta, 1994), and Spanish can have 52 distinct forms reflecting person, number, tense, aspect, and mood (Centeno and Obler, 2001).

Another major difference in Japanese, unlike English and Spanish, is that it does not generally use white spaces to separate words. To identify morphemes and words in Japanese, two main approaches have been taken. The first approach is to define base unit words by first identifying syntactic words, which is done in Universal Dependencies (UD) by using *Short Unit Words* (SUW). The second approach uses *Long Unit Words* (LUW) as the base units in Japanese. Although similar results were achieved from both approaches, Omura et al. (2021) argue that lemmatization of LUW is more complex for a morphologically rich language. The pretrained model used in this paper utilized the SUW approach.

The choice of these three languages allows for typological comparisons, highlighting characteristics that are shared by all three languages, by two of them, or individually. The comparison is shown in Table 2. This summarizes the main linguistic classifications across the three languages. We present the major morphological features relevant for the current study.

	EN	JA	SP
<b>Deriv. Morphology</b>	Prefixes and Suffixes		
<b>Morphology</b>	Analytic	Synthetic	
<b>Gender</b>	No		Nouns
<b>Word Order</b>	SVO	SOV	SVO
<b>Word Formation</b>	Mostly analytic	Agglutination	Inflections
<b>Inflection</b>	Limited	Verbs and Adjectives	Rich
<b>Case Marking</b>	Pronouns	Extensive	Nominative Accusative Dative

Table 2: Main Morphological Descriptions for each Language.

These linguistic differences in inflection levels contribute to the richness of errors observed in ASR systems.

## 3.2 Speech Datasets

We utilized the Common Voice 15 dataset, a publicly available collection of multilingual and open voice data provided by the *Mozilla Common Voice Project* (Ardila et al., 2020). Designed for training and validating automatic speech recognition systems, the dataset encompasses a diverse range of voices and linguistic contexts. The data does not contain personally identifying information. Table 3 below displays the characteristics of the datasets per language.

	Descriptors	EN	JA	SP
audio	Number of Files	16,386	4,978	15,796
	Duration (hr)	26.9	6.6	26.8
	Speech Dur. (hr)	22.7	5.4	23.2
	Mean Dur. (sec)	5.9	4.8	6.11
	Mean Speech Dur. (sec)	4.9	3.9	5.3
text	Total Characters	890K	105K	960K
	Total Words	153K	55K	156K
	Unique Words	21K	8K	23K
	Characters p/Text	54	21	61
	Words p/Text	10	10	10

Table 3: Dataset Descriptions for each Language.

The dataset encompasses contributions from a substantial number of speakers, providing a rich variety of linguistic and acoustic characteristics. In our analysis, we focused on a subset consisting of recordings from the test sets for the three languages. The dataset comprises over 16,000 sentences for English, approximately 5,000 for Japanese, and more than 15,000 sentences for Spanish. This offers a comprehensive sample of spoken language for evaluating ASR systems. The inclusion of a broad range of sentences and speakers enhances the robustness and generalizability of our findings, contributing to a more comprehensive understanding of the performance of the ASR system in diverse linguistic contexts. This includes variations in syntactic, semantic, and phonetic-phonological contexts.

## 3.3 ASR System

All our experiments were conducted using *OpenAI Whisper* (Radford et al., 2023). Whisper comprises multilingual multitask models trained on 680,000 hours of labelled and curated speech data from diverse internet sources. In this experiment, we employed *Whisper-Tiny (T)*, *Whisper-Medium (M)*, *Whisper-Large-v2 (LV2)* and *Whisper Large-*

v3 (LV3). Comparing these four model sizes allows us to examine whether there are relevant accuracy gains across all ASR models.

### 3.4 Analysis

**Word Error Rate:** WER is computed by comparing the reference transcript (ground truth) with the output generated by the ASR system. The formula for WER is given by:

$$WER = (S+D+I) / N$$

Where,  $S$  represents the number of substitutions,  $D$  represents the number of deletions,  $I$  represents the number of insertions, and  $N$  is the total number of words in the reference transcript.

The analysis was conducted in R (R Core Team, 2023) using the outputs of Whisper. Our focus lies in ASR errors when comparing the reference text (REF) to the hypothesis text (HYP). SCLITE was employed for error calculation, identifying substitutions, insertions, and deletions per sentence. SCLITE, part of the NIST SCT<sup>1</sup> Scoring Toolkit, is a tool for scoring and evaluating speech recognition system output. It compares the HYP to the correct REF. Post-comparison, statistics are gathered, and various reports can be generated to summarize recognition system performance. To assess the performance of the ASR system, we utilized the WER metric, a widely accepted measure for transcription accuracy assessment.

**Parts of Speech and Lexical Items:** Linguistic tagging was conducted using the UDPIPE (Straka and Straková, 2017) library (Wijffels et al., 2023) in R to enhance the textual analysis of transcribed speech data. UDPIPE, a state-of-the-art Natural Language Processing (NLP) library, incorporates pre-trained models for various linguistic tasks, which are based on Universal Dependencies (UD). Specifically, we employed UDPIPE’s pipeline for POS tagging. The tagging process consisted of three main steps.

Firstly, in text preprocessing, raw transcripts underwent preprocessing to eliminate artifacts or noise that might impact tagging accuracy. Secondly, during tokenization, preprocessed transcripts were tokenized into individual words or sub-word units using UDPIPE’s tokenization module. The third step involved POS Tagging, where the POS tagging module assigned

grammatical categories – such as nouns, verbs, adjectives – to each token in the transcripts. This information was crucial for understanding the syntactic structure of the spoken content. Careful consideration of punctuation, case sensitivity, and text normalization procedures was carried out to ensure accurate comparisons between REF and ASR-generated transcripts. UDPIPE outputs have been reported to exhibit varying levels of performance. Straka and Straková (2017) reported that the automatic identification of POS has an accuracy of 93.50% for English, 88.19% for Japanese, and 98.14% for Spanish. We assess our outputs based on these reported accuracy levels.

**Linguistic Metric Analysis:** We propose a metric that categorizes errors based on whether they occur in any of the two categories within a *Word Class*: GRAMMATICAL and LEXICAL. From the UDPIPE output, each POS was grouped into either the GRAMMATICAL group (ADP, AUX, CCONJ, DET, PART, PRON, SCONJ) or the LEXICAL Group (ADJ, ADV, NOUN, NUM, PROPN, VERB). From this, we calculated errors at the POS tagging in the REF and HYP texts, defined as  $POS_{er}$ , and we also calculated differences at the LEXICAL and GRAMMATICAL levels, following the formulae below:

$$POS_{er} = (S_{POS} + D_{POS} + I_{POS}) / N_{POS}$$

$$LEX_{er} = (S_{LEX} + D_{LEX} + I_{LEX}) / N_{LEX}$$

$$GRAM_{er} = (S_{GRAM} + D_{GRAM} + I_{GRAM}) / N_{GRAM}$$

*Word Class* errors are then calculated for each group across the entire dataset per language and ASR model size: LEXICAL errors ( $LEX_{er}$ ) and GRAMMATICAL errors ( $GRAM_{er}$ ).

	Total	Size	Subset	Diff.
EN	16386	T	14459	12%
		M	14855	9%
		LV2	14886	9%
		LV3	14896	9%
JA	4978	T	573	88%
		M	1544	69%
		LV2	1634	67%
		LV3	1724	65%
SP	15796	T	12569	20%
		M	14833	6%
		LV2	14910	6%
		LV3	14967	5%

Table 4: Number of Sentences included in the Analyses.

<sup>1</sup> <https://github.com/usnistgov/SCT>

426 To ensure fair comparisons, the analysis was  
 427 conducted on sentences with matching number of  
 428 words, i.e., when **REF** and **HYP** have the same  
 429 number of words, avoiding penalization for  
 430 incorrect pairs due to deletions and insertions.  
 431 Table 4 above shows the differences in the number  
 432 of sentences in the original transcriptions and the  
 433 ones filtered in the analyses.

## 434 4 Experiment and Results

### 435 4.1 ASR Model Size Comparison

436 Table 5 provides a summary of the experiment  
 437 results, highlighting significant differences in  
 438 performance across the **T**, **M**, **LV2** and **LV3**  
 439 models for all evaluated languages.

	Category	ASR Model (%)			
		LV3	LV2	M	T
EN	WER	8.3	8.9	9.9	23.7
	POS_er	5.2	5.5	6.1	17
	LEX_er	11	11.4	11.9	22.4
	GRAM_er	4.6	4.8	5.3	15.1
JA	WER	5.7	6.4	7.5	24.6
	POS_er	1.7	2	2.4	12.7
	LEX_er	3.6	3.9	5.4	22.9
	GRAM_er	1.9	2.7	3.2	13.7
SP	WER	4	4.9	5.8	23.5
	POS_er	1.5	1.9	2.2	10.9
	LEX_er	4.2	4.7	5	12.4
	GRAM_er	1.6	1.3	3.4	8.5

440 Table 5: Breakdown of Error Rates Results. Values are  
 441 shown in Percentage of Errors (%).

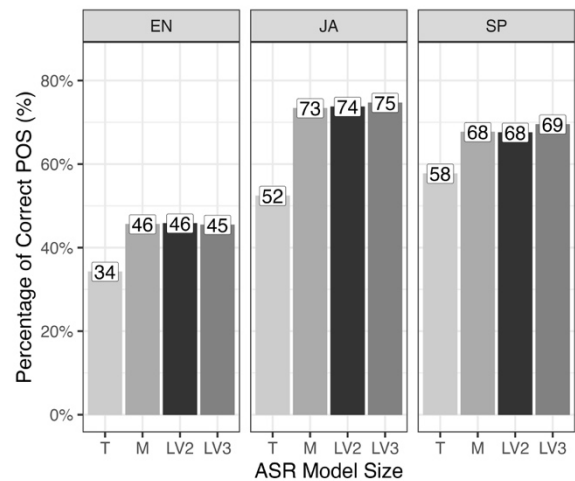
442 The **T** models consistently show the highest error  
 443 rates (English = 23.7%; Japanese = 24.6%; Spanish  
 444 = 23.5%), while the other models (**M**, **LV2** and  
 445 **LV3**) demonstrate notably lower and more uniform  
 446 WERs across all languages. Among these, the **LV3**  
 447 model yields the most accurate results (English =  
 448 8.3%; Japanese = 5.7%; Spanish = 4%). It is  
 449 evident that the **T** models show comparable WERs  
 450 for the three languages, whereas the larger models  
 451 exhibit higher accuracy, with Spanish being the  
 452 most accurate and English the least accurate.

### 453 4.2 Parts of Speech Comparison

454 When delving into the other metrics, our results  
 455 discover a more nuanced understanding, shedding  
 456 light on the categories to which ASR systems are

457 more susceptible for errors. **POS\_er** results  
 458 demonstrate lower error rates in comparison to  
 459 WER. This is notably more distinctive for Japanese  
 460 and Spanish (English = 8.3%WER vs 5.2%  
 461 **POS\_er**; Japanese = 5.7% vs 1.7%; Spanish = 4%  
 462 vs 1.5%).

463 These results indicate that errors are more  
 464 generalizable at the POS level, as compared to the  
 465 word level. As such, this can help better our  
 466 understanding of what types of errors can be  
 467 consistently expected from ASR outputs, and in  
 468 what morphological contexts. A more in-depth  
 469 analysis looked at those cases where the word form  
 470 was incorrect (which counts to more WER) but it  
 471 still had the same POS (which did not count as error  
 472 for the **POS\_er**).



473

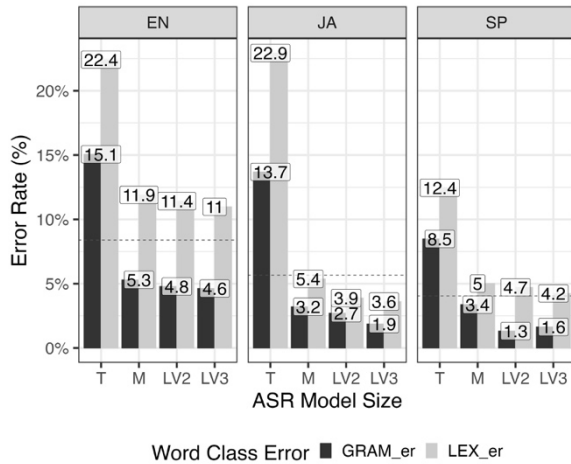
474 Figure 2: Percentage of Cases (and Counts) of  
 475 Wrong Words but with Correct POS.

476 In the case of highly inflectional languages, this  
 477 difference can be observed when the **HYP** text has  
 478 a singular form of a noun (e.g., *cat*), but the **REF**  
 479 text was the same word in the plural form (e.g.,  
 480 *cats*). This observation underscores the limitations  
 481 of relying solely on WER, as it fails to capture  
 482 subtle linguistic nuances retained in **POS\_er**. Figure  
 483 2 shows a breakdown by language and model size  
 484 for this experiment. It shows that Japanese and  
 485 Spanish have more cases where errors are  
 486 explained by inflectional differences between  
 487 words (i.e., words are different but not their POS),  
 488 as compared to English.

### 489 4.3 Word Class Comparison

490 The third layer of analysis distinguishes between  
 491 LEXICAL error rate (**LEX\_er**) and GRAMMATICAL  
 492 error rate (**GRAM\_er**), revealing patterns not

493 captured by the previous two layers (WER and  
 494 POS\_er). Figure 3 presents the error rates broken  
 495 down by language, model size, and word class  
 496 (GRAMMATICAL or LEXICAL) with a horizontal  
 497 dotted line indicating the overall POS\_er as  
 498 reference.



499

500 Figure 3: Error Rates across all Languages and  
 501 Model Sizes split by Word Class Errors.

502 Among the languages examined, Spanish  
 503 consistently shows the lowest overall error rates,  
 504 while English presents the highest. In the LV3  
 505 model analysis, for LEX\_er, Japanese records  
 506 slightly lower rates than Spanish, while English  
 507 exhibits the highest error rates (English = 11%;  
 508 Japanese = 3.6%; Spanish = 4.2%). This variation  
 509 can be explained linguistically by the fact that  
 510 LEXICAL categories in Japanese and Spanish have  
 511 higher inflections than in English, and these  
 512 inflections are presented as affixes in both  
 513 languages, helping the ASR system to understand  
 514 the patterns of occurrence, useful to identify and  
 515 predict the word form and its function in the  
 516 language. This indicates that correct inflectional  
 517 words significantly enhance predicting LEXICAL  
 518 words. Although this finding is in contrast with  
 519 Berg et al. (2024) and Smith-Lock. (1991), our  
 520 results show that higher inflections are related to  
 521 higher accuracy.

522 Further examination explored the extent to  
 523 which predictable inflections helped in correctly  
 524 identifying words for the ASR system. For this, we  
 525 chose PROPER NOUNS (PROPN), a subclass of the  
 526 LEXICAL words (See Appendices for reference).  
 527 Our results show that Japanese is the language with  
 528 least error rates, and English with the most errors  
 529 (English = 39.6%; Japanese = 5.7%; Spanish =

530 18.6%). This is attributed to the use of case markers  
 531 for Proper Nouns in Japanese, feature that is absent  
 532 in English and Spanish, facilitating more accurate  
 533 identification and prediction of Proper Nouns.

534 The analysis revealed that the top six occurring  
 535 words after PROPN were the case markers *さん*  
 536 (3.4% – honorific particle), *の* (3.6% – possessive),  
 537 *に* (2.4% – place), *と* (1.8% – joining nouns), *は*  
 538 (1.8% – topic marking particle), and *が* (1.5% –  
 539 grammatical subject), all accounting for  
 540 approximately 15% of all words after PROPN in the  
 541 Japanese dataset.

542 GRAM\_er results show that Spanish had the  
 543 lowest error rates compared to Japanese (slight  
 544 difference of 0.3%) and, more significantly, to  
 545 English (English = 4.6%; Japanese = 1.9%;  
 546 Spanish = 1.6%). An in-depth analysis highlighted  
 547 that the primary errors in English were associated  
 548 with subordinating conjunctions (e.g., *if*, *that*,  
 549 *while*) whereas the coordinating conjunctions were  
 550 the ones driving more errors in Japanese (e.g., *と*  
 551 *and*; *も* *also*) and Spanish (e.g., *y* *and*; *o* *or*). This  
 552 indicates that a combination of grammatical  
 553 assessment and linguistic function helps in a deeper  
 554 understanding of how languages use specific words  
 555 and the impact it has on the ASR accuracy. This  
 556 approach is not necessarily language-dependent,  
 557 but rather relies more on the typological function a  
 558 word class has across multiple languages.

#### 559 4.4 Assessment across all Comparisons

560 When assessing all error rate metrics, we find that  
 561 POS\_er, LEX\_er and GRAM\_er contribute to a more  
 562 robust understanding of ASR errors. The examples  
 563 below are used to analyze the different layers. The  
 564 first line is the transcription and below each, the  
 565 POS tagging is given for each word.

##### 567 (a) English

568 REF: *what is the matter with a thousand dollars*  
 569 PRON V DET NOUN ADP DET NUM N  
 570 HYP: *what is the matter with the thousand dollar*  
 571 PRON V DET NOUN ADP DET NUM N

##### 573 (b) Spanish

574 REF: *las batallas se libraron primero en los territorios*  
 575 DET N PRON V ADJ ADP DET N  
 576 HYP: *las batallas el vibrado primero en los territorios*  
 577 DET N DET N ADJ ADP DET N

579 Sentence (a) has a WER of 25% (two wrong words  
 580 over eight words in total). The errors are found in  
 581 the words *a* and *dollars* in the REF text. The word  
 582 *a*, the indefinite article, changed to definite article

583 *the* in the **HYP** text. The change happened in a  
584 GRAMMATICAL word. The second word, *dollars*,  
585 changed to the singular form in the **HYP** text,  
586 *dollar*. In terms of classification, it is a LEXICAL  
587 word, however, the change occurred in a  
588 morpheme that conveys plurality in English.  
589 Since these errors did not change the word class,  
590 POS\_er, LEX\_er, and GRAM\_er have 0% error  
591 rate. This also shows that the NOUN *dollar* was not  
592 changed to another NOUN, but just its plurality,  
593 which did not compromise the meaning as  
594 compared to being changed to another word, like  
595 *scholar*, for example.

596 Sentence (b) has the same WER as (a), 25%.  
597 However, the error patterns are different. The  
598 errors are found in the words “se” (reflexive  
599 PRON) and “libraron” (simple past on V *to fight*)  
600 in the **REF** text. Both words were substituted with  
601 different words and with different POS categories,  
602 however, they belong to the same word class  
603 (PRON > DET; V > N). POS\_er is 25%, LEX\_er is  
604 12.5%, and GRAM\_er is 12.5%. Compared to  
605 sentence (a), only the WER is the same, but the  
606 other error metrics are all different. These two  
607 sentences exemplify the complementary nature of  
608 the application of all metrics, rather than being  
609 competing measures.

## 610 **5 Discussion**

611 This study examines two NLP-driven error metrics  
612 across English, Japanese, and Spanish to assess the  
613 accuracy of the Whisper ASR system. Given the  
614 higher inflectional complexity, Japanese and  
615 Spanish provide a valuable context for analysis,  
616 particularly in GRAMMATICAL words. Our findings  
617 reveal that relying solely on WER can obscure  
618 nuanced aspect of ASR performance. For instance,  
619 while WER figures maybe comparable across  
620 models, such as in the **T** model results, a closer  
621 examination at the GRAMMATICAL vs LEXICAL  
622 level unveils distinct accuracies. Conversely, even  
623 with differing WERs, such as in the Large models  
624 (**LV2**, **LV3**) where Japanese and Spanish  
625 outperform English, a detailed analysis exposes the  
626 ASR system’s consistent performance on LEXICAL  
627 words but divergence in handling GRAMMATICAL  
628 words, notably with English struggling more with  
629 subordinating conjunctions.

630 These metrics bridge the gap between general  
631 WER assessment and more granular POS error  
632 analyses. While reporting each POS category  
633 individually could complicate comparisons across  
634 languages, our metrics offer a balanced approach.

635 This approach allows for a detailed identification  
636 of strengths and weaknesses in ASR systems at  
637 crucial linguistic levels, enhancing both the  
638 interpretability and practical applicability of ASR  
639 performance evaluations.

640 For the areas of improvement pertinent to ASR  
641 systems, we can make suggestions based on the  
642 observed patterns and the datasets used. First, the  
643 training of the ASR systems should include more  
644 accurate weighting of words based on whether they  
645 are LEXICAL vs GRAMMATICAL in function. This  
646 adjustment will enable a better-informed decision  
647 driven by word usage in context. Second,  
648 performance can also be improved if more  
649 spontaneous speech is used in the training of ASR  
650 systems. This can lead to the observation of more  
651 LEXICAL words, such as low-frequency words or  
652 those more commonly found in spoken language,  
653 in contexts that are less typical of controlled  
654 speech. Finally, the errors observed strongly  
655 suggest that errors follow specific linguistic  
656 patterns (e.g., LEXICAL vs GRAMMATICAL, or  
657 PROPON vs NOUN). In this sense, they go beyond  
658 language-dependent patterns and can be better  
659 understood under linguistic typologies.

## 660 **6 Conclusions**

661 The automatic processing and annotation of natural  
662 speech are complex tasks influenced by both the  
663 systems themselves and, most importantly, by the  
664 inherent characteristics of languages and their  
665 typological differences. Current systems have  
666 made significant progress in addressing these  
667 complexities. One notable advancement is the  
668 ability to perform automatic grammatical error  
669 comparisons across languages with different  
670 typological classifications. This advancement  
671 necessitates a cautious approach to understanding  
672 intrinsic language differences and variations based  
673 on the ASR system or the data used for training.

674 Linguistically informed metrics play a crucial  
675 role in interpreting performance. This, combined  
676 with robust NLP approaches, prove to be efficient  
677 for this task. Our metric and implementation  
678 developed for assessing ASR performance help  
679 identifying areas for improvement and linguistic  
680 aspects that pose specific challenges. Additionally,  
681 these metrics assist those working on ASR systems  
682 and datasets in developing more efficient  
683 algorithms and infrastructures.



## 684 7 Limitations

685 In this study, we investigated three languages with  
686 different typological characteristics, highlighting  
687 both shared and unique features. However, the  
688 scope of our research did not extend to  
689 polysynthetic languages, presenting a limitation in  
690 the diversity of language types analyzed. Future  
691 work should include a broader range of languages  
692 to determine if the observed patterns are replicated  
693 or consistent across major language families.

694 Our analysis was confined to the Whisper ASR  
695 system, chosen for its broad usage and  
696 accessibility. This focus presents a limitation as it  
697 does not address a comparative evaluation across  
698 different ASR systems. Future research should  
699 examine a variety of ASR systems to ascertain  
700 whether the observed errors are influenced by  
701 specific systems or linguistic characteristics  
702 themselves.

703 Finally, we did not address fine-tuning of ASR  
704 models. This can help identify whether the errors  
705 are specific to the ASR model or the data used for  
706 training, and to what extent we can generalize these  
707 errors. Future research should also investigate the  
708 impact of linguistically-based fine-tuning on the  
709 performance of ASR systems.

## 710 Acknowledgments

711 To be added in final version.

## 712 References

713 Ahmed Ali and Steve Renals. 2018. Word Error Rate  
714 Estimation for Speech Recognition: e-WER. In  
715 *Proceedings of the 56th Annual Meeting of the*  
716 *Association for Computational Linguistics*, Volume  
717 2, pp. 20-24, Melbourne, Australia. Association for  
718 Computational Linguistics.

719 Rosana Ardila, Megan Branson, Kelly Davis, Michael  
720 Henretty, Michael Kohler, Josh Meyer, Reuben  
721 Morais, Lindsay Saunders, Francis M. Tyers, and  
722 Gregor Weber. 2020. Common Voice: A Massively-  
723 Multilingual Speech Corpus. In *Proceedings of the*  
724 *12th Conference on Language Resources and*  
725 *Evaluation (LREC 2020)*, pp. 421-4215.

726 Kristina Berg, Stefan Hartmann, and Daniel Claeser.  
727 2024. Are some morphological units more prone to  
728 spelling variation than others? A case study using  
729 spontaneous handwritten data. *Morphology*, Volume  
730 34: pp. 173-188.

731 Jie Cao, Ananya Ganesh, Jon Cai, Rosy Southwell, E.  
732 Margaret Perkoff, Michael Regan, Katharina Kann,  
733 James H. Martin, Martha Palmer, and Sidney

734 D'Mello. 2023. A Comparative Analysis of  
735 Automatic Speech Recognition Errors in Small  
736 Group Classroom Discourse. In *UMAP '23:*  
737 *Proceedings of the 31st ACM Conference on User*  
738 *Modeling, Adaptation and Personalization (UMAP*  
739 *'23)*, June 26--29, 2023, Limassol, Cyprus. ACM,  
740 New York, NY, USA 13 Pages.

741 Jose G. Centeno and Loraine K. Obler. 2001.  
742 Agrammatic verb errors in Spanish speakers and  
743 their normal discourse correlates. *Journal of*  
744 *Neurolinguistics*, Volume 14, pp. 349-363.

745 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and  
746 Kristina Toutanova. 2019. Bert: Pretraining of deep  
747 bidirectional transformers for language  
748 understanding. *North American Chapter of the*  
749 *Association for Computational Linguistics*.

750 Yetunde E. Adegbegeha, Aarav Minocha, and Renu  
751 Balyan. 2024. Analyzing Multilingual Automatic  
752 Speech Recognition Systems Performance. In F.  
753 Zhao, and D. Miao, (eds.), *Communications in*  
754 *Computer and Information Science, Volume 1946*.  
755 Springer, Singapore.

756 Rahhal Errattahi, Asmaa E. Hannani, and Hassan  
757 Ouahmane. 2018. Automatic Speech Recognition  
758 Errors Detection and Correction: A Review. In  
759 *International Conference on Natural Language and*  
760 *Speech Processing, Volume 128*, pp. 32-37.

761 Ethnologue: languages of the world. 1999.  
762 Dallas:Texas, SIL International.

763 John S. Garofolo, Ellen M. Voorhees, C G. Auzanne,  
764 Vincent M. Stanford, and B A. Lund. 1998. 1998  
765 TREC-7 spoken document retrieval track overview  
766 and results. In *Proceedings of the 7th Text REtrieval*  
767 *Conference*. NIST, pp. 79-89.

768 Xiaodong He, Li Deng, and Alex Acero. 2011. Why  
769 word error rate is not a good metric for speech  
770 recognizer training for the speech translation task?.  
771 In *IEEE International Conference on Acoustics,*  
772 *Speech and Signal Processing (ICASSP)*, pp. 5632-  
773 5635.

774 Saengchan Hemchua and Norbert Schmitt. 2006. An  
775 Analysis of Lexical Errors in The English  
776 Compositions of Thai Learners. *Prospect: An*  
777 *Australian Journal of TESOL*, 21(3):3-25.

778 Toru Hisamitsu and Yoshihiko Nitta. 1994. An  
779 Efficient Treatment of Japanese Verb Inflection for  
780 Morphological Analysis. In *International*  
781 *Conference on Computational Linguistics*, pp. 194-  
782 200.

783 Sushant Kafle and Matt Huenerfauth. 2017. Evaluating  
784 the usability of automatically generated captions for  
785 people who are deaf or hard of hearing. In  
786 *Proceedings of the 19th International ACM*

- 877 *SIGACCESS Conference on Computers and*  
878 *Accessibility*, pp. 165-174.
- 879 Hamza Kheddar, Yassine Himeur, Somaya Al-  
890 Maadeed, Abbas Amira, and Faycal Bensaali. 2023.  
891 Deep transfer learning for automatic speech  
892 recognition: Towards better generalization.  
893 *Knowledge-Based Systems, Volume 277*.
- 894 Elhard Kumalija and Yukikazu Nakamoto. 2022.  
895 Performance evaluation of automatic speech  
896 recognition systems on integrated noise-network  
897 distorted speech. *Frontiers in Signal Processing*,  
898 2:1-10.
- 899 John Lee and Stephanie Seneff. 2008. Correcting  
900 Misuse of Verb Forms. In *Proceedings of ACL-08:*  
901 *HLT, Columbus, Ohio, USA, June 2008, Association*  
902 *for Computational Linguistics*, pp. 174-182.
- 903 Sungjin Lee, Hyungjong Noh, Kyusong Lee, and Gary  
904 G. Lee. 2011. Grammatical error detection for  
905 corrective feedback provision in oral conversations.  
906 In *Proceedings of the Twenty-Fifth AAAI*  
907 *Conference on Artificial Intelligence (AAAI'11)*,  
908 AAAI Press, pp. 797-802.
- 909 Chengzhang Li, Ming Zhang, Xuejun Zhang, and  
910 Yonghong Yan. 2023. MCRSpell: A metric learning  
911 of correct representation for Chinese spelling  
912 correction. *Expert Systems with Applications*,  
913 *Volume 237, Part B*.
- 914 Iain A. McCowan, Darren Moore, John Dines, Daniel  
915 Gatica-Perez, Mike Flynn, Pierre Wellner, and  
916 Herve Bourlard. 2004. *On the use of information*  
917 *retrieval measures for speech recognition*  
918 *evaluation*. Technical Report. IDIAP.
- 919 Andrew C. Morris, Viktoria Maier, and Phil D. Green.  
920 2004. From WER and RIL to MER and WIL:  
921 Improved evaluation measures for connected speech  
922 recognition. In *Proceedings INTERSPEECH 2004 –*  
923 *8th Annual Conference of the International Speech*  
924 *Communication Association*, Jeju Island, Korea,  
925 October 2004, pp. 2765-2768.
- 926 A. NithyaKalyani and S. Jothilakshmi. 2019. Speech  
927 Summarization for Tamil Language. In Ed. Nilanjan  
928 Dey, *Intelligent Speech Signal Processing*,  
929 Academic Press, pp. 113-138.
- 930 Douglas O'Shaughnessy. 2023. Trends and  
931 developments in automatic speech recognition  
932 research. *Computer Speech & Language, Volume*  
933 *83*.
- 934 Mai Omura, Aya Wasaka, and Masayuki Asahara.  
935 2021. Word Delimitation Issues in UD Japanese. In  
936 *Proceedings of the Fifth Workshop on Universal*  
937 *Dependencies (UDW, SyntaxFest 2021)*, Sofia,  
938 Bulgaria, Association for Computational  
939 Linguistics, pp. 142-150.
- 940 Chanho Park et al. 2023. Fast word error rate  
941 estimation using self-supervised representations for  
942 speech and text. *arXiv preprint arXiv:2310.08225*.
- 943 Maja Popović and Hermann Ney. 2007. Word Error  
944 Rates: Decomposition over POS classes and  
945 Applications for Error Analysis. In *Proceedings of*  
946 *the Second Workshop on Statistical Machine*  
947 *Translation, Association for Computational*  
948 *Linguistics*, Prague, Czech Republic, pp. 48-55.
- 949 R Core Team. 2023. R: A language and environment  
950 for statistical computing. *R Foundation for*  
951 *Statistical Computing*, Vienna:Austria, URL  
952 <https://www.R-project.org/>.
- 953 Alec Radford, Jong W. Kim, Tao Xu, Greg Brockman,  
954 Christine McLeavey, and Ilya Sutskever. 2023.  
955 Robust speech recognition via large-scale weak  
956 supervision. In *Proceedings of the 40th*  
957 *International Conference on Machine Learning*,  
958 July 2023, pp. 28492-28518.
- 959 Nils Reimers and Iryna Gurevych. 2019. Sentence-  
960 bert: Sentence embeddings using siamese bert-  
961 networks. In *Proceedings of the 2019 Conference on*  
962 *Empirical Methods in Natural Language*  
963 *Processing and the 9th International Joint*  
964 *Conference on Natural Language Processing*,  
965 *Association for Computational Linguistics*, pp.  
966 3973-3983.
- 967 Thomas Reitmaier, Electra Wallington, Dani K. Raju,  
968 Ondrej Klejch, Jennifer Pearson, Matt Jones, Peter  
969 Bell, and Simon Robinson. 2022. Opportunities and  
970 Challenges of Automatic Speech Recognition  
971 Systems for Low-Resource Language Speakers. In  
972 *Proceedings of the 2022 CHI Conference on Human*  
973 *Factors in Computing Systems (CHI '22)*, New  
974 York: USA, 299, pp. 1-17.
- 975 Thibault B. Roux, Mickael Rouvier, Jane Wottawa, and  
976 Richard Dufour. 2022. Qualitative evaluation of  
977 language model rescoring in automatic speech  
978 recognition. In *Proceedings INTERSPEECH 2022 –*  
979 *23rd Annual Conference of the International Speech*  
980 *Communication Association*, Incheon, Korea, Sep.  
981 2022, pp. 3968-3972.
- 982 Karen M. Smith-Lock. 1991. Errors of Inflection in the  
983 Writing of Normal and Poor Readers. *Language and*  
984 *Speech*, 34(4), pp. 341-350.
- 985 Milan Straka and Jana Straková. 2017. Tokenizing,  
986 POS Tagging, Lemmatizing and Parsing UD 2.0  
987 with UDPipe. In *Proceedings of the CoNLL 2017*  
988 *Shared Task: Multilingual Parsing from Raw Text to*  
989 *Universal Dependencies*, pp. 88-99, Vancouver,  
990 Canada. Association for Computational Linguistics.
- 991 Ryan Whetten and Casey Kennigton. 2023. Evaluating  
992 and Improving Automatic Speech Recognition  
993 using Severity. In *The 22nd Workshop on*

894 *Biomedical Natural Language Processing and* 925 **Japanese Errors Breakdown**  
 895 *BioNLP Shared Tasks*, July 13, pp. 79-91.

896 Jan Wijffels, Milan Straka, and Jana Straková. 2023.  
 897 *Udpipe: Tokenization, Parts of Speech Tagging,*  
 898 *Lemmatization and Dependency Parsing with the*  
 899 *Udpipe Nlp Toolkit*. [https://CRAN.R-](https://CRAN.R-project.org/package=udpipe)  
 900 [project.org/package=udpipe](https://CRAN.R-project.org/package=udpipe).

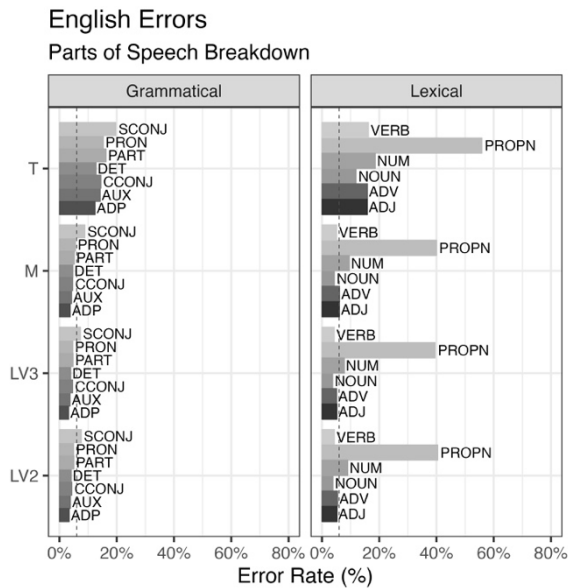
901 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q.  
 902 Weinberger, and Yoav Artzi. 2020. *Bertscore:*  
 903 *Evaluating text generation with bert*. In *8th*  
 904 *International Conference on Learning*  
 905 *Representations, ICLR 2020*, Addis Ababa,  
 906 Ethiopia, April 26–30.

## 907 A Appendices

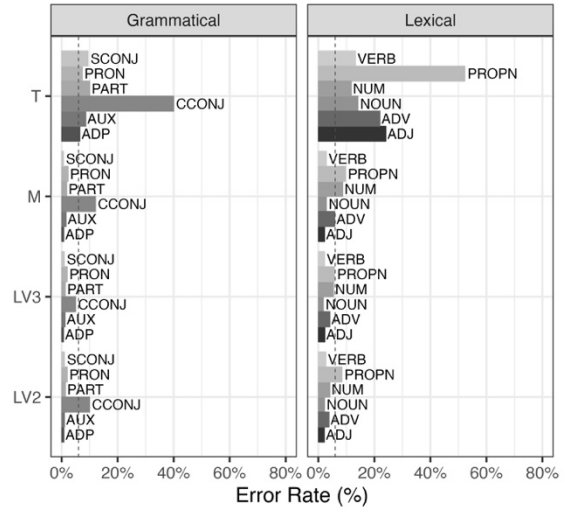
908 This section includes the breakdown of all errors  
 909 for Parts of Speech across all languages. The  
 910 horizontal dotted lines indicate the overall POS\_er 926  
 911 as reference.

912

### 913 English Errors Breakdown



### Japanese Errors Parts of Speech Breakdown



### 927 Spanish Errors Breakdown

#### Spanish Errors Parts of Speech Breakdown

