

Tracing L1 Interference in English Learner Writing: A Longitudinal Corpus with Error Annotations

Anonymous ACL submission

Abstract

The availability of suitable learner corpora is paramount for the study of second language acquisition (SLA) and language transfer. However, curating learner corpora is a challenging endeavor as high quality learner data is rarely publicly available. This results in only a few such corpora, such as ICLE and TOEFL-11, available to the community. To address this important gap, in this paper we present ANONYMOUS,¹ a novel English learner corpus with longitudinal data. ANONYMOUS contains texts written by adult learners taking English as a second language courses in the USA with the goal of either preparing for university admission or improving their language proficiency while starting their university degrees. ANONYMOUS contains 687 instances written by speakers of 15 different L1s. Unlike most learner corpora, this corpus contains longitudinal data which enables researchers to investigate language learning over time. We present two case studies using ANONYMOUS at the intersection of SLA and Computational Linguistics: (1) Native Language Identification (NLI); and (2) a quantitative and qualitative study using LLMs on linguistic features influenced by L1.²

1 Introduction

A language learner’s native language (L1, or first language) often influences fluency, grammatical patterns, and vocabulary usage in their second language (L2). This influence can result in L2 production containing distinctive linguistic features that may be unfamiliar or questionable to a native speaker.

These features have been the focus of research in Second Language Acquisition (SLA) and Computational Linguistics, particularly through the use of

learner corpora. These corpora enable the systematic analysis of language learner production and support tasks such as Native Language Identification (NLI), which involves automatically identifying a learner’s L1 based on linguistic patterns in their L2 writing or speech.

Early NLI research primarily concentrated on developing machine learning (ML) models for spoken language, analyzing features such as pronunciation, stress, and prosodic patterns (Krishna et al., 2019). More recent work in text-based NLI, however, leverages written features — including word choice, syntax, and spelling — to make predictions about an individual’s native language.

Although text-based NLI has a range of applications — including author profiling, forensics, spam and phishing detection, and various educational uses (Malmasi et al., 2017) — it has received comparatively less research attention. This paper investigates the use of NLI to identify authors’ native languages in student essays. Each student essay is analyzed for various error types to explore potential correlations between the student’s L1 and the types and frequencies of errors produced. We demonstrate how NLI can be used not only to automatically identify an author’s L1, but also to contribute to research in Second Language Acquisition (SLA).

The contributions of our work are the following:

1. We introduce ANONYMOUS, a novel corpus of L2 writing with longitudinal data. The corpus can be used for a variety of purposes in Computational Linguistics and SLA.
2. We describe the first linguistically-informed LLM-based study of features of L1 to L2 transfer on longitudinal data.
3. We present various NLI experiments using this corpus. We evaluated the performance of various models, from traditional classifiers like SVMs, to state-of-the-art LLMs such as GPT-4.

¹Anonymized to ensure double-blind review.

²All data and code will be made publicly available upon acceptance of this manuscript.

2 Related Work

Second Language Acquisition and Error Taxonomies A substantial body of SLA research has documented systematic learner errors often attributed to L1 interference (Richards, 1971; Odlin, 1989). While large learner corpora, such as the Cambridge Learner Corpus (Nicholls, 2003) and NUCLE (Dahlmeier et al., 2013), include valuable metadata about each writer’s L1, they do not typically annotate individual errors for cross-linguistic influence. Instead, error frameworks tend to focus on the nature of the error by classifying its locus (lexis, syntax, morphology) and the type of surface modification required (e.g., omission, addition, substitution), rather than investigating *why* it emerged (Díaz-Negrillo and Fernández-Domínguez, 2006).

Grammatical Error Correction and LLMs Recent advancements in grammatical error correction (GEC) have been driven by the emergence of large language models (LLMs) like GPT-3 and GPT-4, which have been evaluated for their performance on GEC tasks (Song et al., 2024; Kobayashi et al., 2024). For instance, studies have investigated the effectiveness of LLMs in GEC evaluation by employing prompts designed to incorporate various evaluation criteria (Loem et al., 2023; Fang et al., 2023). However, these models primarily focus on correcting errors rather than explaining them in the context of a learner’s L1.

Recent research has attempted to go beyond grammatical error correction by considering L1 influences in academic writing. Zomer and Frankenberg-Garcia 2021 proposed a pre-trained encoder-decoder model designed to improve research writing by adapting corrections to the writer’s L1 background. Their approach recognizes that L1 influences writing style and errors, offering targeted corrections based on linguistic transfer effects. However, the study primarily focuses on enhancing research writing rather than systematically analyzing or categorizing L1 interference at a linguistic level, and the model does not explicitly attribute errors to specific sources of transfer, such as phonological, orthographic, or syntactic influences from the L1.

Our Contribution In contrast to these approaches, our work is, to our knowledge, the first to use LLMs paired with human oversight for explicit L1 interference analysis. We require the model to identify whether an error stems from L1 interfer-

ence and at what level (e.g., syntax, morphology) and justify the label with concrete linguistic features from the learner’s native language. By integrating SLA insights, we generate fine-grained annotations that capture L1 influence. This structured, L1-aware output goes *beyond* standard GEC tasks, helping to bridge the gap between automatic correction and the emphasis on deeper linguistic analysis in SLA research.

Native Language Identification NLI operates on the assumption that a learner’s native language shapes the acquisition and production of a second language, a phenomenon referred to as cross-linguistic influence or language transfer (Krashen, 1981; Ellis, 2015). Language transfer results in L1 features manifesting in L2 production, allowing computational models to recognize patterns shared by speakers of the same L1 when communicating in a given L2. Text-based NLI has a number of important applications, such as serving as a corpus-driven approach for SLA (Jarvis and Crossley, 2012) and enabling the development of effective L2 teaching materials and computer-aided language learning (CALL) software. Additionally, NLI has been shown to improve NLP systems when dealing with texts from non-native speakers, contributing to tasks like author profiling, forensics, spam and phishing detection (Malmasi et al., 2017).

As evidenced by a recent survey (Goswami et al., 2024), traditional statistical models such as Support Vector Machines (SVMs) trained on n -grams as features have historically delivered the best performance for text-based NLI. A few recent studies (Lotfi et al., 2020; Uluslu and Schneider, 2022; Zhang and Salle, 2023; Ng and Markov, 2025), however, have shown that fine-tuned LLMs such as GPT-4 deliver state-of-the-art performance for English NLI. In this paper, we test multiple approaches on this corpus, capturing the full breadth of the available toolkit from including SVM ensembles all the way to the recently released GPT-4o.

3 The ANONYMOUS Corpus

Collection context ANONYMOUS was gathered from an introductory academic-writing course taken by international post-graduate students at a U.S. R1 university between 2022 and 2024.³ Learners produced three assignment types (short answers, long essays, and group reflections) based on

³Ethics and data-governance details appear in §7.

Dataset	L1 Languages	Size	L1 Information		Annotations	
			L1 Metadata	L1-Annotated Errors	Fine-Grained Errors	Longitudinal
Cambridge Learner Corpus (CLC)	80+	~2.9M words	✓	✗	✓	✗
NUCLE (CoNLL-2014)	—	~1.9M words	✗	✗	✓	✗
FCE Corpus	16 (EU/Asia)	~1,200 essays	✗	✗	✓	✗
ICNALE	10 (E/SE Asia)	~3.8M words	✓	?	✓	✗
TOEFL11	11	~12,000 essays	✓	✗	✓	✗
EFCAMDAT	9	~100K learners	?	✗	✓	✓
BEA-2019 (W&I+LOCNESS)	—	334 learners	✗	✗	✓	✗
ANONYMOUS	Arabic, Chinese, Vietnamese (+12 others)	57 (+26) learners	✓	✓	✓	✓

Table 1: Comparison of ANONYMOUS with widely used SLA learner corpora. Sources: CLC (Nicholls, 2003), NUCLE (Dahlmeier et al., 2013), FCE (Yannakoudakis et al., 2011), ICNALE (Ishikawa, 2023), TOEFL11 (Blanchard et al., 2013), EFCAMDAT (Geertzen et al., 2014), BEA-2019 (Bryant et al., 2019). ✓ = present; ✗ = absent; ? = partial/indirect.

their university and U.S. acculturation experiences, submitting work electronically via the learning-management system.

Learner profile Students self-reported their L1 and country of origin when enrolling; no participant listed multiple L1s. All had demonstrated advanced English proficiency (IELTS ≥ 7) and were enrolled in Master’s programs. The full corpus contains 687 essays from 15 L1s.⁴ Because many L1s are represented by only one or two learners, all analyses in this paper focus on the three languages with at least three writers: Arabic, Chinese, and Vietnamese.

Corpus Composition and Per-L1 Breakdown ANONYMOUS includes a subset comprising texts from Arabic, Chinese, and Vietnamese learners. Table 2 presents the overall size of the subset analyzed in this paper, including the per-L1 breakdown, which details the number of learners, documents, tokens, median document length, mean submissions per learner, and the median number of weeks between first and last submission (a proxy for longitudinal depth).

Note that the three cohorts differ in how often they submitted short versus long tasks (Table 2).

Positioning among existing learner corpora Table 1 contrasts ANONYMOUS with the most frequently used English-learner resources. While several corpora include either longitudinal data (EFCAMDAT) or L1 metadata (ICNALE, TOEFL11), ANONYMOUS is, to our knowledge, the first to combine *all four* of the following in a single resource: fine-grained error labels, explicit L1 transfer annotations, detailed L1 metadata for each learner, and multiple submissions per learner over time.

⁴The full list of L1s is in Appendix Table 10.

This unique combination enables research questions that have been difficult to pursue with previous datasets, such as modeling the trajectory of cross-linguistic influence throughout an academic term (§4) or leveraging L1-aware error signals for few-shot native-language identification (§6).

Examples and splits Table 4 shows anonymized excerpts of each assignment type, while Table 3 gives the training, development, and test split used in our experiments.

4 Error Annotations

SLA-Grounded Annotation We draw on established SLA research to develop an annotation framework for learner errors. The categories—phonetic misrepresentations, morphological overgeneralizations, and L1-based orthographic interference—reflect well-documented SLA phenomena, such as Spanish speakers inserting an “e” before /s/ clusters or the over extension of regular morphological rules, e.g., “bayed” for “bought” (Richards and Schmidt, 2011; Freeman et al., 2016). Grounding our schema in SLA principles ensures theoretical and pedagogical relevance.

Using LLMs for L1-Based Annotation Our key methodological contribution is leveraging LLMs to generate SLA-informed annotations at scale, significantly reducing the labor-intensive nature of traditional error annotation.

Conventional annotation processes require thousands of expert-annotator hours to construct large corpora, with estimates suggesting that annotating one million words could take 2000-5000 hours⁵.

⁵For context, manually annotating a corpus of this scale—similar to NUCLE (Dahlmeier et al., 2013)—at an estimated rate of 500 words per hour would require extensive expert labor. This estimate accounts for multiple annotation passes, as is standard in error correction corpora, and is derived from previous annotation efforts (Dahlmeier et al., 2013;

L1	Learners	Docs	Tokens	Median tok/doc	Entries per learner	Span(wks)	Count		Proportion	
							Long	Short	Long	Short
Arabic	35	345	63090	79	9.86	10	158	187	0.47	0.53
Chinese	18	133	28835	88	7.39	4	69	64	0.50	0.50
Vietnamese	4	47	12471	199	11.75	11.9	35	12	0.70	0.30
Total	57	525	104396	-	-	-	-	-	-	-

Table 2: Corpus subset composition and per-L1 breakdown, including the total number of documents, tokens, learners, and document types analyzed in this paper, a full breakdown of the corpus composition can be found in Appendix D, Table 10.

L1	Train	Dev	Test	Total
Arabic	275	35	35	345
Chinese	107	13	13	133
Vietnamese	37	5	5	47
Total	419	53	53	525

Table 3: Document counts in the train/dev/test split.

In contrast, our approach harnesses a prompt-driven LLM to systematically classify errors, integrating SLA insights to provide structured, L1-aware annotations at scale. The prompt (see Appendix A) guides the model to:

- Identify each error’s subcategory (orthographic, morphological, lexical, grammatical, etc.).
- Flag L1 interference when observed, referencing specific native-language forms (e.g., a Spanish “e+s” cluster or Arabic morphological patterns).

We then extract the exact error span. Figure 1a shows examples of Chinese and Arabic L1 interference, verified by native bilingual speakers.

4.1 Modeling Error Rate Differences Across Assignment Types

To account for repeated submissions by the same learners, we fit a Poisson Generalized Estimating Equations (GEE) model (log link, token count as offset) with robust standard errors, clustering by writer. This approach accounts for within-writer correlation without modeling random effects directly, allowing for population-averaged estimates of assignment type effects. Short answers exhibit approximately 3.47 times the error rate of essays ($\beta=1.24$, $p<0.001$), while group assignments show a 45% increase compared to individual assignments ($\beta=0.37$, $p=0.004$). This finding suggests that LLM-based error detection may systematically under-report errors in longer submissions, highlighting a limitation when comparing error rates across texts of varying lengths.

Ng et al., 2014).

```
{
  "incorrect": "in the learning aspect",
  "correct": "in terms of learning",
  "type": {
    "L1InterferenceSubcategory.SYNTACTIC_INTERFERENCE": 1
  },
  "l1_interference_reason": "Chinese syntax often uses phrases like '在...方面' which translates directly to 'in the... aspect', leading to syntactic interference.",
  "span_start": 3136,
  "span_end": 3158
}
```

(a) Annotated learner errors illustrating syntactic interference from Chinese L1, where direct translations of native constructions result in non-standard English expressions.

```
{
  "incorrect": "attande",
  "correct": "attend",
  "type": {
    "L1InterferenceSubcategory.ORTHOGRAPHIC_INTERFERENCE": 0.7,
    "OrthographySubcategory.PHONETIC": 0.3
  },
  "l1_interference_reason": "Arabic speakers might add extra vowels or alter consonant sounds due to the absence of certain English phonemes in Arabic, leading to 'attande' instead of 'attend'."
}
```

(b) Annotated learner errors illustrating orthographic interference from Arabic L1, where phonetic spelling errors arise from the lack of vowel marking in Arabic.

Figure 1: Each entry contains the incorrect phrase, its span, the corrected form, and an explanation of the interference type.

4.2 Detecting keyboard typos

To determine whether the LLM occasionally assigns high-stakes labels to errors that are really just keyboard slips, we compared the QWERTY keyboard-distance distribution of the Typo category with every other sub-category using Welch’s *t*-test (Table 8).

The keyboard-distance analysis suggests that the LLM is generally well-calibrated for high-level categories (e.g., Grammatical, Lexical, L1-Interference). However, it tends to over-label several low-level orthographic phenomena.

In particular, consonant-doubling, consonant-substitution, morphological, and phonetic errors often resemble typos in terms of key proximity. Statistical analysis revealed no significant differ-

Assignment Type	Question	Student Answer
Short	Why are we asking you about the “type of learning” that is happening at UNIVERSITY?	To know about what I get benefit from it.
Long	Dissertation Paper – Write about your experience at UNIVERSITY.	After few hours fly, two plant transfer finely I got to the destination...
Group	Describe what you have learned from the group project.	The first, take away is that I can talk with me from the language activity is that most people have a perfect specking skill...

Table 4: Anonymized examples of the three assignment types.

ence in mean distances for these error types when compared to genuine typos ($p > .05$). This indicates that many such tokens could be re-classified as benign slips rather than systematic errors. Conversely, errors with significantly larger mean distances than typos ($|t| \geq 2.08$, $p < .05$) include grammatical, lexical, L1-interference, hyphenation/spacing, silent-letter/irregular, and vowel-substitution/omission. These categories typically involve changes that go beyond adjacent-key slips, suggesting a more substantive error rather than a mere typo. Interestingly, capitalization/punctuation and the broader punctuation class showed smaller average distances compared to typos ($t = 2.49$ and 3.11 ; $p = .017$ and $.004$). This pattern is consistent with same-key mistakes, such as missed shift keys, rather than cross-key substitutions.

These findings motivate two main adjustments: (i) implementing a post-processing rule to downgrade low-distance instances within borderline sub-categories, and (ii) refining prompt engineering to explicitly consider keyboard proximity when distinguishing between typos and more substantial errors.

However, we do not remove labels for errors that resemble typos solely based on keyboard proximity. The fact that some morphological or phonetic errors have similar distances to genuine typos does not imply they are typographical mistakes; such errors may still arise from systematic L1 interference or language processing challenges. Therefore, we interpret the similarity as a potential confounding factor rather than grounds for exclusion.

4.3 Human Verification of GPT-4 Annotations

To assess the reliability of our automatically-generated labels, we employed a two-tier human-in-the-loop verification process. This approach combines document-level recall checks with native-speaker scrutiny of L1 interference claims, providing a principled estimate of annotation quality.

Verification Process All essays and error snippets were presented in a web interface that allowed span-level confirmation or correction; corrections were stored as an additional layer in the corpus. Disagreements were discussed in weekly meetings to ensure consistent annotation practices.

The verification process involved two stages:

- **L1-Specific Check:** Two native-speaker linguists (Arabic and Mandarin) independently evaluated 10 randomly-selected errors per language flagged as L1 interference by GPT-4. They answered the following questions:

- **Q1 (Plausibility):** Is this a plausible case of L1 interference? (Yes / No + rationale).
- **Q2 (Explanation):** Is GPT-4’s explanation of the interference accurate? (Yes / No + rationale).

Native-speaker acceptance rates were 100% for both Arabic and Mandarin.

- **Document-Level Audit:** A third linguist, experienced in corpus annotation, audited 13% of the essays (stratified by L1 and assignment type). The linguist evaluated whether:

- GPT-4 correctly identified errors or missed any errors.
- Identified errors were correctly typed (orthographic, morphological, grammatical, etc.).
- For errors labeled as L1 interference, both the attribution and the explanation were accurate.

4.4 Evaluation

Table 5 presents the precision, recall, and F1 score for each evaluation aspect.

Interpretation The results indicate that GPT-4o is highly effective at detecting learner errors with high precision (92%), meaning that when the model flags an issue, it is usually correct. This strong performance is particularly evident in detecting surface-level issues such as typos and lexical inter-

Metric	Precision	Recall	F1 Score
Error Detection	0.916	0.107	0.191
Correction Agreement	0.697	0.083	0.149
Type Agreement	0.613	0.074	0.132
L1 Reason Agreement	0.837	0.038	0.072

Table 5: Overall performance metrics for LLM annotations compared to human annotations.

ference (see Table 9), where the model consistently provides accurate and relevant corrections (70% precision).

However, the model demonstrates lower recall (11%), indicating that it often misses less obvious or contextually embedded errors, particularly in longer texts where errors are more dispersed. This pattern suggests that while GPT-4o is reliable when identifying clear, surface-level errors, there is room for improvement in capturing more nuanced linguistic issues, such as morphological over-generalization and affix errors.

The L1-specific verification further supports the model’s strengths: native-speaker linguists unanimously confirmed (100%) the plausibility of L1 interference in both Arabic and Mandarin cases flagged by the model. This finding highlights the model’s ability to accurately identify L1-related errors when they are detected.

Overall, the results demonstrate that GPT-4o is an effective tool for high-precision error detection, especially in tasks where surface-level accuracy is critical. Future improvements could focus on enhancing recall, particularly for longer or more syntactically complex texts, in order to maximize the model’s utility for comprehensive learner error analysis.

Annotator	LLM	
	Error	NotError
Error	113	912
NotError	7	0

Table 6: Confusion matrix (correct, wrong) for error detection between LLM and human annotator.

5 Data Analysis

5.1 Tracking Student Errors Over Time

As timestamped writing submissions enable longitudinal analysis at both individual and cohort levels, we track student error patterns over time to analyze student development and learning trajectories.

To ensure comparability across time periods, we normalize error counts against text length and assignment counts. This allows us to assess whether certain error types diminish with proficiency gains or persist, indicating deeper linguistic challenges. Of course, the expectation for an English proficiency course is that learner errors diminish over time.

None of the observed fluctuations (e.g., rising error counts in certain months, subsequent declines) reach statistical significance (see Appendix C). However, the fine-grained L1-based labels reveal that certain patterns persist—such as Arabic speakers’ difficulties with vowel representation or literal syntactic translations from Chinese—suggesting that some cross-linguistic influences remain stable over time rather than disappearing with increased exposure to English (Odlin, 1989).

Our results seem to contradict our hypothesis that error frequencies should reduce – for the 2022 cohort, for instance, error frequencies largely increase from one assignment to the other until the last assignment. For 2024, the story is somewhat reversed. We plan to explore several possible explanations for these observations. For example, it might be the case that students do become better L2 speakers, but their assignments also become harder, leading to more errors. Or, perhaps, it could be the case that the first assignment was, by design, an easy one, leading to fewer errors, and, if we discard it, for the 2023 and 2024 cohorts we might actually confirm our hypothesis that learner error frequencies reduce over time. We plan to explore these explanations more deeply in future work, engaging with the instructors of the class as well as with the students themselves.

5.2 Lexical Development

Beyond tracking general error trends, we also explore lexical development in relation to Romance and Germanic vocabulary acquisition. Previous studies have documented that Germanic and Romance L1 speakers tend to overuse cognates from their respective L1s in English at lower proficiency levels, with this reliance decreasing as proficiency increases (Nativ et al., 2024). However, our focus dataset consists of Arabic, Chinese, and Vietnamese L1 speakers, for whom English lacks a strong lexical overlap with their native languages. Analyzing how these learners acquire vocabulary from different etymological sources represents a novel contribution to SLA research.

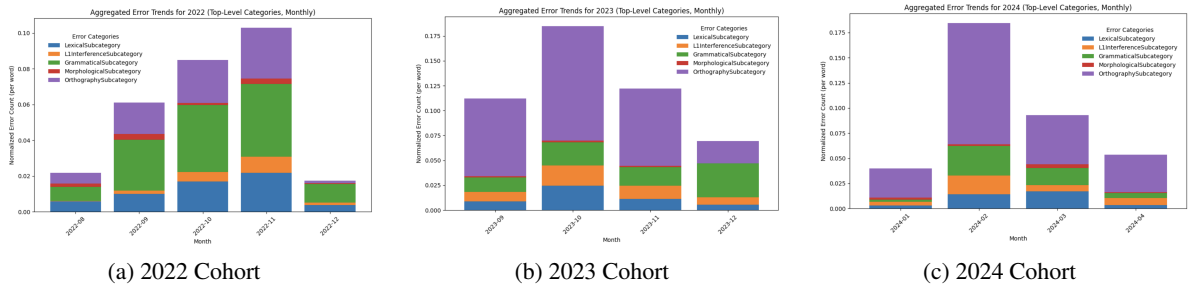


Figure 2: Aggregated Error Trends for Different Cohorts (Top-Level Categories, Monthly). The 2022 cohort shows a gradual increase in errors, peaking in November. The 2023 cohort exhibits higher orthographic errors throughout, while the 2024 cohort displays a sharp peak in February before declining.

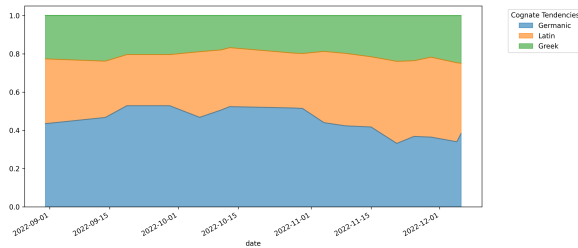


Figure 3: Proportion of Germanic, Latin, and Greek-derived vocabulary in learner writing over time (2022 cohort). The increase in Latin-based words suggests a shift toward academic vocabulary, while Germanic words remain dominant.

In theory, we expect to see an increasing tendency toward Romance-derived vocabulary as students advance in proficiency, given that academic and formal English draws heavily from Latin and French (Hernandez et al., 2021). Our analysis partially supports this: the 2022 cohort (see Figure 3) shows a statistically significant rise in Latin-based vocabulary over time ($p = 0.0199$). However, this trend vanishes in the 2023 and 2024 cohorts, raising questions about how learners from non-Indo-European backgrounds acquire academic vocabulary. Differences in instructional input, cognitive processing, or exposure to academic vocabulary may contribute to these variations. The observed increase in the 2022 cohort suggests that under certain conditions, learners do shift toward more Latin-derived vocabulary as they progress, highlighting the need for further research into the factors that influence this shift. Future studies should examine whether these trends persist across larger datasets and explore pedagogical interventions that could facilitate the acquisition of academic English vocabulary for learners from diverse linguistic backgrounds.

5.3 Further Syntactic Pattern Analysis

Syntactic analysis in NLP and SLA research has traditionally relied on head-dependent relations within dependency trees (Constant et al., 2017). However, these relations often fail to capture multi-word syntactic units that function as a single structural unit. This is also the issue with analyses that focus on common Part-of-Speech n -grams.

Here, we propose to use *syntactic catenae* as the unit of analysis to remedy these issues. Osborne et al. (2012) introduced *catenae* as a more flexible syntactic representation, defining them as *any sequence of words that maintains a continuous dominance relationship in a dependency tree*. This definition allows catenae to include non-constituent structures and discontinuous elements that are crucial for syntactic analysis.

Catenae have been used in syntactic theory to describe verb complexes, idiomatic expressions, and discontinuous dependencies (Osborne et al., 2012; Imrényi, 2013). However, their application in corpus-based computational linguistics, particularly in L2 syntactic variation analysis, remains unexplored. We investigate whether catenae distributions exhibit L1-specific patterns in learner writing, exploring whether different L1 groups favor certain syntactic constructions when producing English.

We additionally conduct a supplementary investigation using POS bigrams, which capture short-range syntactic dependencies (De Gregorio et al., 2024). While less structurally expressive than catenae, POS bigrams offer a more conventional means of detecting syntactic variation across L1 groups.

Methodology Using Stanza (Qi et al., 2020), we extract catenae from dependency-parsed texts, representing them as sequences of (*dependency relation*, *POS tag*) pairs (e.g., `det-DT | comp:obj-NN | mod-JJ`). This allows for a structural analysis

independent of lexical choice. For interpretability, we also retain corresponding lexical sequences.

To supplement the catenae analysis, we also extract POS bigrams from learner texts, identifying adjacent POS sequences (e.g., DT NN, NN VBZ) as a proxy for syntactic tendencies across L1 groups.

Cross-L1 Comparison For both catenae and POS bigrams, we compute relative frequencies and apply TF-IDF weighting to identify structures that were more prominent in one L1 group relative to others. Across both analyses, we do not observe *strong* L1-specific syntactic patterns. Frequent catenae were largely **shared across L1 groups**, with no consistent L1-driven structural tendencies. That said, we do observe some interesting differences across different L1s. For example, compound noun constructions feature more prominently in Vietnamese L1 speakers and much less common in Chinese ones, even though one might expect the opposite due to the extensive compounding in Chinese.

We should note that the large space of possible catenae combinations and our rather sparse corpus limit our ability to detect robust differences. The relatively small number of speakers per L1 further constrained cross-L1 generalizability. We maintain, though, that catenae are the appropriate unit of analysis for uncovering L1-influenced syntactic patterns, and leave such a larger scale analysis encompassing more corpora for future work.

6 Native Language Identification

As a further showcase of the utility of our dataset for other downstream tasks, we carry out multiple NLI experiments with results presented in Table 7. We report results in terms of accuracy and macro F1 score following the literature in this task (Goswami et al., 2024).

Models We train multiple SVM systems using various features such as POS n -grams of $n \in [1, 4]$ and word n -grams of $n \in [1, 2]$. We then combine them in a majority voting ensemble (Malmasi and Dras, 2017) and we refer to this model as SVM Ensemble in the table. We also fine-tune multiple BERT-based models on ANONYMOUS namely BERT, mBERT (Devlin et al., 2019), and RoBERTa (Liu et al., 2019). For these, we use a learning rate of $1e-5$ for all models and early stopping on our development set. Last, we benchmark three LLMs on ANONYMOUS, namely FLAN-T5 (Chung et al.,

Approach	Models	Acc.	F1
Statistical			
	SVM Ensemble	0.75	0.73
BERT-based			
	roBERTa	0.79	0.75
	BERT	0.77	0.72
	mBERT	0.70	0.68
LLM Zero-shot			
	GPT 4o	0.66	0.66
	LLaMa3.1	0.41	0.43
	FLAN T5	0.32	0.37
LLM Fine-tuning			
	GPT 4o	0.97	0.96
	LLaMa3.1	0.87	0.84
	FLAN T5	0.66	0.53

Table 7: Results of different models on the ANONYMOUS dataset. LLMs require fine-tuning to outperform BERT-based and simple statistical approaches.

2024), GPT-4o (Achiam et al., 2023), and the 70B parameter LLaMa 3.1 (Touvron et al., 2023). We benchmark the three models using both zero-shot prompting as well as task-specific fine-tuning on the training set.

NLI Takeaways Corroborating the results reported in recent studies using popular NLI datasets like TOEFL 11 (Ng and Markov, 2025), we observe that fine-tuned models achieve the highest performance on ANONYMOUS. All three LLMs obtain significant performance improvement from zero-shot prompting to task fine-tuning. The performance of LLMs using zero-shot prompting is, in turn, inferior to the performance of both SVM ensemble and the three BERT models. This indicates that off-the-shelf LLMs do not fare particularly well in identifying L1s without any specific task fine-tuning.

7 Conclusion and Future Work

We present ANONYMOUS, a first-of-its-kind dataset of learner English, which stands apart from others due to encompassing longitudinal data and fine-grained L1 interference annotations. We showcase interesting analysis on three L1s, introduce new syntactic analysis units, and perform NLI experiments on a subset of our dataset.

Importantly, ANONYMOUS will continue expanding every year with each incoming student cohort. As a result, ANONYMOUS will facilitate promising research directions in Second Language Acquisition research, while also presenting opportunities for challenging setups in the development of language learning applications.

Limitations

Our approach likely performs best for high-resource languages, as LLMs are trained predominantly on well-documented linguistic data. For low-resource languages with limited digital presence or sparse learner corpora, the model’s ability to identify and explain L1 interference may be weaker, leading to noisier or less reliable annotations. This perhaps limits the generalizability of our approach, but we believe this limitation is mitigated by the fact that most second language learners opt to learn high-resource languages.

Additionally, while we conduct careful manual verification of a subset of model-generated annotations for the three L1s that we study in this paper, a more extensive validation process is likely needed to ensure consistency and reliability across diverse L1s.

A major challenge for the reproducibility of our work is the rapid evolution of LLMs (e.g., GPT-3.5, GPT-4), as results can depend on a specific model version that later might become unavailable. We chose to rely on the best currently available model to ensure higher quality annotations for our dataset, but future work could reproduce this effort with open-sourced/open-weight models to explore robustness to model variation. In addition, future work should evaluate performance across a broader range of linguistic backgrounds and explore strategies for maintaining reproducibility despite ongoing model updates.

Ethical Considerations & Data Governance

The dataset collected for this research is undergoing Institutional Review Board (IRB) approval at the authors university. IRB is a committee at (research) institutions that reviews research proposals that involve human subjects with the goal of ensuring that all research follows ethical guidelines and regulations.

We explicitly address ethical considerations related to the collection, processing, and sharing of our dataset.

Data Collection Consent Our dataset contains writing samples from non-native English speakers at ANONYMOUS. Permission for dataset sharing was obtained from the appropriate university departments. All data was collected with ethical guidelines for linguistic research.

Data Retention & Access The anonymized dataset will be made available to researchers for non-commercial purposes. Although highly unlikely given the steps described below, access requires agreement to terms that prohibit attempts to re-identify individuals or use the data for purposes beyond research. While we maintain language origin information (L1) for linguistic analysis purposes, all identifying information has been removed to protect student privacy.

Anonymization Process The data contained in ANONYMOUS is non-sensitive in nature as it consists of responses to exercise prompts. Nevertheless, to protect students’ privacy and anonymity, we implemented a rigorous anonymization pipeline. We first eliminate all meta-data that could potentially review identifiable information keeping only non-identifiable meta-data such as the writers’ L1, course taken, and the exercise prompts for each instance. Secondly, in line with best practices in the field, we replace any potentially identifiable information within the text instances with a placeholder token (Megyesi et al., 2018). This includes place or people names. This process was done semi-automatically using Python scripts with the aid of a researcher working in the project who manually checked instances for any remaining information that could potentially reveal a students’ identity.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Daniel Blanchard, Joel Tetreault, Derrick Higgins, Aoife Cahill, and Martin Chodorow. 2013. *Toefl11: A corpus of non-native english*. *ETS Research Report Series*, 2013(2):i–15.
- Christopher Bryant, Mariano Felice, Øistein E. Andersen, and Ted Briscoe. 2019. *The BEA-2019 shared task on grammatical error correction*. In *Proceedings of the Fourteenth Workshop on Innovative Use of NLP for Building Educational Applications*, pages 52–75, Florence, Italy. Association for Computational Linguistics.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2024. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25:1–53.

698	Mathieu Constant, Gülşen Eryiğit, Johanna Monti, Lonneke van der Plas, Carlos Ramisch, Michael Rosner, and Amalia Todirascu. 2017. Multiword expression processing: A survey . <i>Computational Linguistics</i> , 43(4):837–892.	753
699		754
700		755
701		
702		
703	Daniel Dahlmeier, Hwee Tou Ng, and Siew Mei Wu. 2013. Building a large annotated corpus of learner English: The NUS corpus of learner English . In <i>Proceedings of the Eighth Workshop on Innovative Use of NLP for Building Educational Applications</i> , pages 22–31, Atlanta, Georgia. Association for Computational Linguistics.	756
704		757
705		758
706		759
707		
708	J. De Gregorio, R. Toral, and D. Sánchez. 2024. Exploring language relations through syntactic distances and geographic proximity . <i>EPJ Data Science</i> , 13:61.	760
709		761
710		762
711		763
712		764
713	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In <i>Proceedings of NAACL</i> .	765
714		766
715		
716		
717	Ana Díaz-Negrillo and Jesús Fernández-Domínguez. 2006. Error Tagging Systems for Learner Corpora. <i>Revista española de lingüística aplicada, ISSN 0213-2028, Vol. 19, 2006, pags. 83-102</i> , 19.	767
718		768
719		
720		
721	Rod Ellis. 2015. <i>Understanding second language acquisition 2nd edition</i> . Oxford university press.	769
722		770
723		771
724	Tao Fang, Shu Yang, Kaixin Lan, Derek F. Wong, Jinpeng Hu, Lidia S. Chao, and Yue Zhang. 2023. Is chatgpt a highly fluent grammatical error correction system? a comprehensive evaluation . <i>Preprint</i> , arXiv:2304.01746.	772
725		
726		
727		
728	Max R. Freeman, Henrike K. Blumenfeld, and Viorica Marian. 2016. Phonotactic constraints are activated across languages in bilinguals . <i>Frontiers in Psychology</i> , 7:702.	773
729		774
730		775
731		776
732		777
733	Jeroen Geertzen, Theodora Alexopoulou, and Anna Korhonen. 2014. Efcamdat: The ef-cambridge open language database. In <i>Selected Proceedings of the 2012 Second Language Research Forum</i> , pages 240–254, Somerville, MA. Cascadilla Proceedings Project.	778
734		779
735		780
736		781
737		782
738	Dhiman Goswami, Sharanya Thilagan, Kai North, Shervin Malmasi, and Marcos Zampieri. 2024. Native language identification in texts: A survey. In <i>Proceedings of NAACL</i> .	783
739		784
740		785
741		
742	Arturo E. Hernandez, Juliana Ronderos, Jean Philippe Bodet, Hannah Claussenius-Kalman, My V. H. Nguyen, and Ferenc Bunta. 2021. German in childhood and latin in adolescence: On the bidialectal nature of lexical access in english . <i>Humanities and Social Sciences Communications</i> , 8(1):162.	786
743		787
744		788
745		
746		
747	András Imrényi. 2013. The syntax of Hungarian auxiliaries: A dependency grammar account . In <i>Proceedings of the Second International Conference on Dependency Linguistics (DepLing 2013)</i> , pages 118–127, Prague, Czech Republic. Charles University in Prague, Matfyzpress, Prague, Czech Republic.	789
748		790
749		791
750		
751		
752		
	Shin’ichiro Ishikawa. 2023. <i>The ICNALE Guide: An Introduction to a Learner Corpus Study on Asian Learners’ L2 English</i> . Routledge, London.	792
		793
		794
		795
	Scott Jarvis and Scott A Crossley. 2012. <i>Approaching Language Transfer Through Text Classification: Explorations in the Detectionbased Approach</i> . Multilingual Matters.	796
	Masamune Kobayashi, Masato Mita, and Mamoru Komachi. 2024. Large language models are state-of-the-art evaluator for grammatical error correction . In <i>Proceedings of the 19th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2024)</i> , pages 68–77, Mexico City, Mexico. Association for Computational Linguistics.	797
		798
		799
		800
		801
		802
		803
		804
	Stephen Krashen. 1981. Second language acquisition. <i>Second Language Learning</i> .	805
		806
	G Radha Krishna, R Krishnan, and VK Mittal. 2019. An automated system for regional nativity identification of indian speakers from english speech. In <i>Proceedings of IEEE INDICON</i> .	807
		808
	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized BERT pretraining approach. <i>CoRR</i> , abs/1907.11692.	
	Mengsay Loem, Masahiro Kaneko, Sho Takase, and Naoaki Okazaki. 2023. Exploring effectiveness of GPT-3 in grammatical error correction: A study on performance and controllability in prompt-based methods . In <i>Proceedings of the 18th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2023)</i> , pages 205–219, Toronto, Canada. Association for Computational Linguistics.	
	Ehsan Lotfi, Ilia Markov, and Walter Daelemans. 2020. A deep generative approach to native language identification. In <i>Proceedings of COLING</i> .	
	Shervin Malmasi and Mark Dras. 2017. Multilingual native language identification. <i>Natural Language Engineering</i> .	
	Shervin Malmasi, Keelan Evanini, Aoife Cahill, Joel Tetreault, Robert Pugh, Christopher Hamill, Diane Napolitano, and Yao Qian. 2017. A report on the 2017 native language identification shared task. In <i>Proceedings of BEA</i> .	
	Beata Megyesi, Lena Granstedt, Sofia Johansson, Julia Prentice, Daniel Rosén, Carl-Johan Schenström, and Elena Volodina. 2018. Learner corpus anonymization in the age of gdpr: Insights from the creation of a learner corpus of swedish . In <i>Proceedings of the 7th Workshop on NLP for Computer Assisted Language Learning at SLTC 2018 (NLP4CALL 2018)</i> , pages 47–56.	
	Liat Nativ, Yuval Nov, Noam Ordan, Shuly Wintner, and Anat Prior. 2024. Do more proficient writers use fewer cognates in l2? a computational approach . <i>Bilingualism: Language and Cognition</i> , 27(1):84–94.	

- Hwee Tou Ng, Siew Mei Wu, Ted Briscoe, Christian Hadiwinoto, Raymond Hendy Susanto, and Christopher Bryant. 2014. [The CoNLL-2014 shared task on grammatical error correction](#). In *Proceedings of the Eighteenth Conference on Computational Natural Language Learning: Shared Task*, pages 1–14, Baltimore, Maryland. Association for Computational Linguistics.
- Yee Man Ng and Ilia Markov. 2025. [Leveraging open-source large language models for native language identification](#). In *Proceedings of the 12th Workshop on NLP for Similar Languages, Varieties and Dialects*, pages 20–28, Abu Dhabi, UAE. Association for Computational Linguistics.
- Diane Nicholls. 2003. [The cambridge learner corpus: Error coding and analysis for lexicography and elt](#). In *Proceedings of the Corpus Linguistics 2003 conference*, volume 16, page 572–581. Cambridge University Press Cambridge.
- Terence Odlin. 1989. *Language Transfer*. Cambridge Applied Linguistics. Cambridge University Press.
- Timothy Osborne, Michael Putnam, and Thomas Groß. 2012. [Catenae: Introducing a novel unit of syntactic analysis](#). *Syntax*, 15(4):354–396.
- Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning. 2020. [Stanza: A python natural language processing toolkit for many human languages](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 101–108, Online. Association for Computational Linguistics.
- Jack C. Richards. 1971. [A non-contrastive approach to error analysis](#)¹. *ELT Journal*, XXV(3):204–219.
- Jack C. Richards and Richard W. Schmidt. 2011. *Longman Dictionary of Language Teaching and Applied Linguistics*, 4th edition. Routledge, London.
- Yixiao Song, Kalpesh Krishna, Rajesh Bhatt, Kevin Gimpel, and Mohit Iyyer. 2024. [GEE! grammar error explanation with large language models](#). In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 754–781, Mexico City, Mexico. Association for Computational Linguistics.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Ahmet Yavuz Uluslu and Gerold Schneider. 2022. [Scaling native language identification with transformer adapters](#). In *Proceedings of the 5th International Conference on Natural Language and Speech Processing (ICNLSP 2022)*, pages 298–302, Trento, Italy. Association for Computational Linguistics.
- Helen Yannakoudakis, Ted Briscoe, and Ben Medlock. 2011. [A new dataset and method for automatically grading ESOL texts](#). In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 180–189, Portland, Oregon, USA. Association for Computational Linguistics.
- Wei Zhang and Alexandre Salle. 2023. Native language identification with large language models. *arXiv preprint arXiv:2312.07819*.
- Gustavo Zomer and Ana Frankenberg-Garcia. 2021. [Beyond grammatical error correction: Improving L1-influenced research writing in English using pre-trained encoder-decoder models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2534–2540, Punta Cana, Dominican Republic. Association for Computational Linguistics.

A Prompt

LLM Annotation Prompt

Task: You are an expert at identifying and classifying spelling and language errors made by English learners. Your highest priority is to identify errors that may be due to L1 (native language) interference and provide a brief but **specific** explanation of how the L1 could cause such an error. Your explanation should include:

- A **concrete linguistic example** from the L1 (e.g., a word or phrase in the learner's native language) or a well-known phonological, orthographic, or syntactic feature of the L1 that contributes to the error.
- A short discussion of how that L1 feature leads the learner to produce the erroneous English form.

If there is **no** L1 interference, classify the error into one of the following categories: orthographic (including typos), lexical, morphological, or grammatical.

Steps to follow for each erroneous word:

1. **Determine if L1 interference is involved.**
 - If **yes**, select the appropriate L1 interference subcategory and provide a "l1_interference_reason" that:
 - Identifies the specific L1 feature (e.g., a Spanish prefix rule, an Arabic root pattern, a Japanese phonological constraint).
 - Explains how that feature maps to the incorrect English form.
 - If **no**, classify under other subcategories: orthographic (including typos), lexical, morphological, or grammatical.
2. **Return the errors in the order they appear in the text.**

Error Categories and Descriptions

1. Orthography Subcategories

- **Phonetic Errors**
 - Definition: Words spelled purely by sound, ignoring English orthographic norms.
 - Examples:
 - * *fone* → *phone*
 - * *nife* → *knife*
- **Vowel Substitution and Omission**
 - Definition: Substituting or omitting vowels incorrectly.
 - Examples:
 - * *hop* → *hope*
 - * *beter* → *better*
- **Silent Letters and Irregular Spelling**
 - Definition: Ignoring or mishandling silent letters or irregular spelling patterns.
 - Examples:
 - * *clim* → *climb*
 - * *writting* → *writing*
- **Consonant Substitution Errors**
 - Definition: Replacing one consonant with another.
 - Examples:
 - * *shose* → *chose*
 - * *joke* → *yoke*
- **Hyphenation, Compound Words, and Spacing Errors**
 - Definition: Errors in spacing or hyphenation of compound words.
 - Examples:
 - * *infact* → *in fact*
 - * *some where* → *somewhere*

2. Lexical Subcategories

- **Homophone Confusion**
 - Definition: Mixing up words that sound alike but differ in spelling and meaning.
 - Examples:
 - * *their* → *there*
 - * *peace* → *piece*
- **Lexical Errors**

- Definition: Errors involving incorrect word choice due to misunderstanding of meaning.
- Examples:
 - * *among* → *below*
 - * *borrow* → *lend*
- **Phonological Confusion**
 - Definition: Errors where words are confused due to phonological similarities, often involving metathesis, substitution of similar phonemes, or confusion between near-homophones.
 - Examples:
 - * *aboard* → *abroad* (Metathesis: reversed phonemes)
 - * *form* → *from* (Transposition of adjacent sounds)
 - * *claps* → *class* (Substitution of "p" for "s")

3. Morphological Subcategories

- **Morphemic Errors with Affixes**
 - Definition: Incorrect handling of prefixes or suffixes.
 - Examples:
 - * *beautifull* → *beautiful*
 - * *hoping* → *hopping*
- **Overgeneralization of Spelling Rules**
 - Definition: Applying English morphological or spelling rules too broadly.
 - Examples:
 - * *buyed* → *bought*
 - * *goed* → *went*

4. L1 Interference Subcategories

- **Orthographic Interference**
 - Definition: Applying L1 spelling conventions to English.
 - Examples:
 - * *esplendid* → *splendid* (Spanish: adding "e" before "s" clusters)
 - * *colur* → *colour* (British vs. American orthography confusion)
- **Lexical Interference**
 - Definition: Using L1-based lexical forms or cognates in English.
 - Examples:
 - * *telefon* → *telephone* (Spanish or German influence)
 - * *faciliter* → *facilitate* (French influence)
- **Grammatical Interference**
 - Definition: Applying L1 grammatical patterns to English.
 - Examples:
 - * *She has 24 years* → *She is 24 years old* (Spanish: "Ella tiene 24 años")
 - * *He doesn't know nothing* → *He doesn't know anything* (Negative concord in some L1s)
- **Syntactic Interference**
 - Definition: Applying L1 syntactic structures to English.
 - Examples:
 - * *He to the store goes* → *He goes to the store* (German word order influence)
 - * *Beautiful is she* → *She is beautiful* (Japanese syntax influence)

5. Grammatical Subcategories

- **Grammatical Errors**
 - Definition: Errors in grammar, syntax, word order, or agreement.
 - Examples:
 - * *She go yesterday* → *She went yesterday*
 - * *He like apples* → *He likes apples*

Categories and Subcategories:

We define a hierarchical categorization system using Python enums for clarity and consistency:

```
from enum import Enum

class OrthographySubcategory(Enum):
    PHONETIC = "Phonetic Errors"
    VOWEL_SUBSTITUTION_OMISSION = "Vowel Substitution and Omission"
```

```

SILENT_LETTERS_IRREGULAR = "Silent Letters and Irregular Spelling"
CONSONANT_SUBSTITUTION = "Consonant Substitution Errors"
HYPHENATION_SPACING = "Hyphenation, Compound Words, and Spacing Errors"
CONSONANT_DOUBLING = "Consonant Doubling and Dropping"
CAPITALIZATION_PUNCTUATION = "Capitalization and Punctuation Errors"
TYPO = "Typo"

```

```

class LexicalSubcategory(Enum):
    HOMOPHONE_CONFUSION = "Homophone Confusion"
    LEXICAL = "Lexical Errors"
    PHONOLOGICAL_CONFUSION = "Phonological Confusion"

```

```

class MorphologicalSubcategory(Enum):
    MORPHEMIC_AFFIX = "Morphemic Errors with Affixes"
    OVERGENERALIZATION = "Overgeneralization of Spelling Rules"
    CONSONANT_DOUBLING = "Morphological Consonant Doubling and Dropping"

```

```

class L1InterferenceSubcategory(Enum):
    ORTHOGRAPHIC_INTERFERENCE = "Orthographic Interference"
    LEXICAL_INTERFERENCE = "Lexical Interference"
    GRAMMATICAL_INTERFERENCE = "Grammatical Interference"
    SYNTACTIC_INTERFERENCE = "Syntactic Interference"

```

```

class GrammaticalSubcategory(Enum):
    GRAMMATICAL = "Grammatical Errors"

```

Probabilities:

- For each error, provide a "type" field as an object where keys are the enum names (e.g., "OrthographySubcategory.PHONETIC") and values are probabilities (floats).
- Probabilities must sum to 1.0 for that error.

If L1 Interference is detected:

- Include "l1_interference_reason" explaining how the L1 caused the error.

Output Format:

Return a JSON array of objects. Each object should contain:

- "incorrect": the misspelled or erroneous word.
- "correct": the correct form.
- "type": a dictionary of {error_type: probability} where probabilities sum to 1.0.
- "l1_interference_reason": a string if L1 Interference applies.

Format strictly as JSON, with no additional commentary.

Few-Shot Examples:

Example Input:

L1: Spanish

Text: After the long *fly* and waiting two hours, I saw a *plant* arrive, which I thought was the right one because it looked so *esplendid* even though I felt *beter* knowing I had finally gotten there. The *clim* was tough, but I *buyed* a ticket, carrying my *childs* with rain, my friend said he'd *shose* a seat for me, but *infact* issues we had. *im* sad.

Example Output:

```

[
  {
    "incorrect": "plant",
    "correct": "plane",
    "type": {
      "OrthographySubcategory.PHONETIC": 0.8,
      "OrthographySubcategory.CONSONANT_SUBSTITUTION": 0.2
    }
  }
]

```

```
    }  
  },  
  {  
    "incorrect": "esplendid",  
    "correct": "splendid",  
    "type": {  
      "L1InterferenceSubcategory.ORTHOGRAPHIC_INTERFERENCE": 0.7,  
      "OrthographySubcategory.PHONETIC": 0.3  
    },  
    "l1_interference_reason": "Spanish speakers often add an 'e' before 's' clusters due to L1  
      orthographic habits."  
  },  
  ...  
]
```

Note: This is a truncated example. The full prompt can be found in the GitHub repository.

B Error Annotation Analysis

Error Type	Capitalization	Consonant Doubling	Consonant Substitution	Grammatical	Hyphenation	L1 Interference	Lexical	Morphological	Punctuation	Typo
Capitalization	-	(-5.050, 0.000)	(-9.758, 0.000)	(-10.483, 0.000)	(-9.363, 0.000)	(-10.490, 0.000)	(-14.618, 0.000)	(-2.804, 0.006)	(2.660, 0.008)	(-2.488, 0.017)
Consonant Doubling	(5.050, 0.000)	-	(-1.345, 0.185)	(-3.970, 0.000)	(-2.436, 0.017)	(-6.166, 0.000)	(-5.177, 0.000)	(2.781, 0.008)	(6.250, 0.000)	(0.669, 0.506)
Consonant Substitution	(9.758, 0.000)	(1.345, 0.185)	-	(-3.426, 0.001)	(-1.531, 0.128)	(-5.862, 0.000)	(-5.039, 0.000)	(5.478, 0.000)	(12.613, 0.000)	(1.657, 0.104)
Grammatical	(10.483, 0.000)	(3.970, 0.000)	(3.426, 0.001)	-	(1.836, 0.067)	(-3.183, 0.002)	(-0.683, 0.495)	(7.564, 0.000)	(12.031, 0.000)	(3.633, 0.001)
Hyphenation	(9.363, 0.000)	(2.436, 0.017)	(1.531, 0.128)	(-1.836, 0.067)	-	(-4.647, 0.000)	(-2.885, 0.004)	(6.128, 0.000)	(11.186, 0.000)	(2.458, 0.017)
L1 Interference	(10.490, 0.000)	(6.166, 0.000)	(5.862, 0.000)	(3.183, 0.002)	(4.647, 0.000)	-	(2.926, 0.004)	(8.731, 0.000)	(11.323, 0.000)	(5.676, 0.000)
Lexical	(14.618, 0.000)	(5.177, 0.000)	(5.039, 0.000)	(0.683, 0.495)	(2.885, 0.004)	(-2.926, 0.004)	-	(10.021, 0.000)	(17.420, 0.000)	(4.356, 0.000)
Morphological	(2.804, 0.006)	(-2.781, 0.008)	(-5.478, 0.000)	(-7.564, 0.000)	(-6.128, 0.000)	(-8.731, 0.000)	(-10.021, 0.000)	-	(4.605, 0.000)	(-1.182, 0.243)
Punctuation	(-2.660, 0.008)	(-6.250, 0.000)	(-12.613, 0.000)	(-12.031, 0.000)	(-11.186, 0.000)	(-11.323, 0.000)	(-17.420, 0.000)	(-4.605, 0.000)	-	(-3.113, 0.004)
Typo	(2.488, 0.017)	(-0.669, 0.506)	(-1.657, 0.104)	(-3.633, 0.001)	(-2.458, 0.017)	(-5.676, 0.000)	(-4.356, 0.000)	(1.182, 0.243)	(3.113, 0.004)	-

Table 8: Keyboard distance analysis: Pairwise T-tests between error types. Each cell shows the T-statistic and P-value for the corresponding pair.

Error Category	Precision	Recall	F1 Score
Orthography: Vowel Substitution/Omission	0.333	1.000	0.500
Grammatical: Grammatical	0.483	1.000	0.651
L1 Interference: Grammatical Interference	0.882	1.000	0.938
Orthography: Consonant Substitution	0.600	1.000	0.750
Orthography: Phonetic	0.333	1.000	0.500
Orthography: Typo	1.000	1.000	1.000
Orthography: Capitalization/Punctuation	0.778	1.000	0.875
Orthography: Hyphenation/Spacing	0.750	1.000	0.857
L1 Interference: Orthographic Interference	0.000	0.000	0.000
Lexical: Lexical	0.833	1.000	0.909
L1 Interference: Lexical Interference	1.000	1.000	1.000
Orthography: Silent Letters/Irregular	0.400	1.000	0.571
L1 Interference: Syntactic Interference	0.800	1.000	0.889
Lexical: Phonological Confusion	1.000	1.000	1.000
Morphological: Overgeneralization	0.000	0.000	0.000
Morphological: Morphemic/Affix	0.000	0.000	0.000

Table 9: Type-wise performance metrics for LLM annotations compared to human annotations. Rows with all zero values indicate the model didn't produce the given error at all.

C Error Trends by L1 and Year

In this section, we present the aggregated error trends for each L1 group across different years. Each plot shows the distribution of top-level error categories normalized by text length.

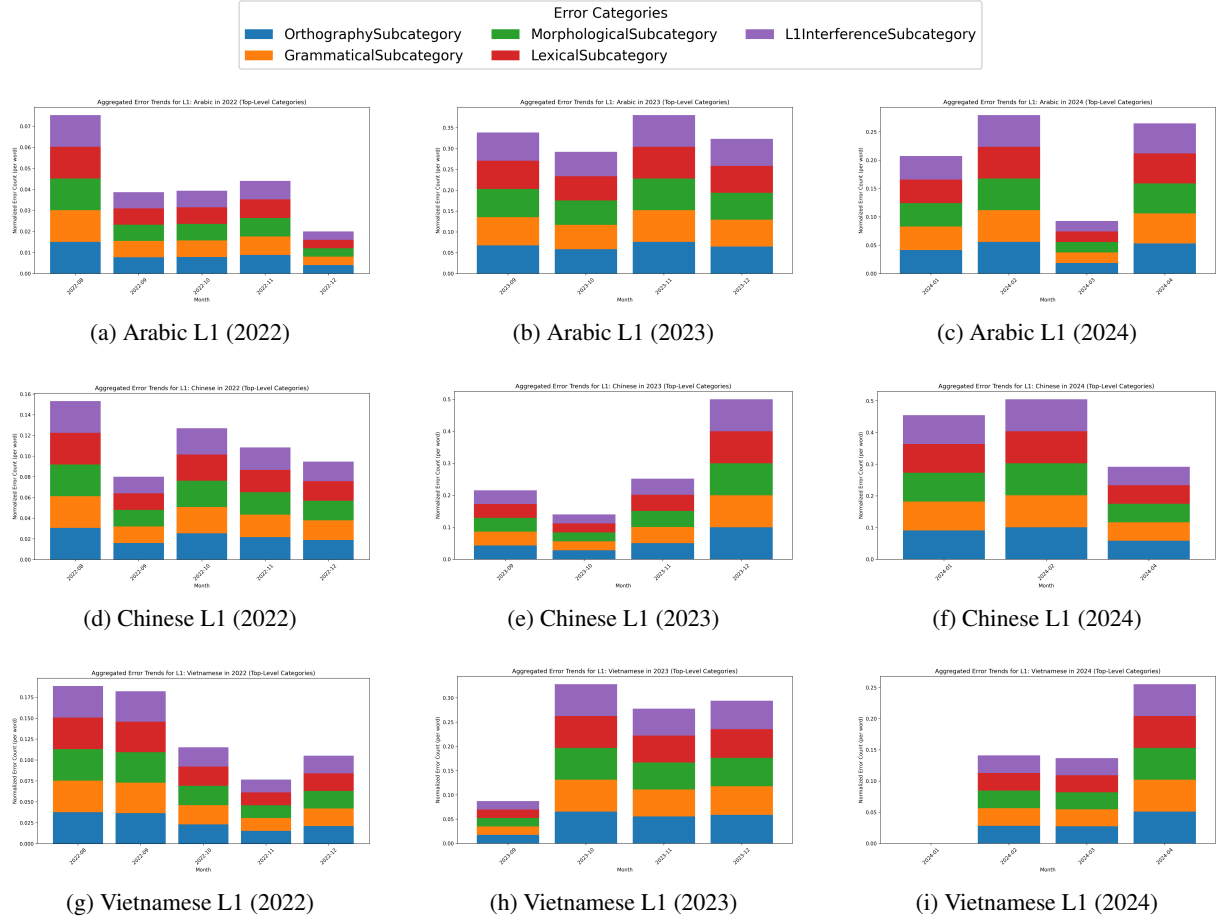


Figure 4: Aggregated error trends by L1 and year. Each subfigure represents a different L1-year combination.

D Corpus Composition

L1	Learners	Docs	Tokens	Med. tok/doc	Entries per learner	Span(wks)	Count		Proportion	
							Long	Short	Long	Short
Arabic	35	345	63090	80.0	9.83	10.00	121	223	0.35	0.65
Azerbaijani	2	12	1934	113.0	6.00	5.00	5	7	0.42	0.58
Bengali	1	3	1990	415.0	3.00	4.00	3	0	1.00	0.00
Chinese	18	132	28678	86.5	7.33	4.00	57	5	0.43	0.57
Dari	2	26	12118	332.0	13.00	9.86	26	0	1.00	0.00
French	1	17	9139	393.0	17.00	13.86	17	0	1.00	0.00
Indonesian	1	8	820	89.0	8.00	11.00	1	7	0.13	0.87
Korean	1	13	1933	140.0	13.00	9.00	8	5	0.62	0.38
Kyrgyz	1	3	339	83.0	3.00	2.00	1	2	0.33	0.67
Portuguese	1	3	1301	281.0	3.00	4.00	3	0	1.00	0.00
Russian	1	19	12024	490.0	19.00	13.86	19	0	1.00	0.00
Sindhi	1	16	9611	567.0	13.00	13.86	14	2	0.87	0.13
Telugu	2	36	19493	416.5	18.00	13.86	35	1	0.97	0.03
Urdu	1	2	384	192.0	2.00	0.29	2	0	1.00	0.00
Vietnamese	4	47	12471	199	11.75	11.9	35	12	0.70	0.30
Total	72	682	175325	-	-	-	-	-	-	-

Table 10: Corpus composition and per-L1 breakdown, including the total number of documents, tokens, learners, and document types analyzed in this paper.