

Automatic Construction of Error Taxonomies with Large Language Models

Anonymous ACL submission

Abstract

Error taxonomies organize error instances into structured categories and provide an important foundation for error analysis. However, most existing taxonomies are manually designed by experts, making them difficult to adapt to new domains and data distributions. In this paper, we propose an LLM-based framework for automatically constructing error taxonomies directly from raw error instances. The framework generates semantic error type labels with a large language model, builds hybrid representations that combine label semantics with transformation signals between erroneous and corrected expressions, and induces taxonomy structures through hierarchical clustering, followed by LLM-based semantic naming and structure refinement. Experiments on both grammatical errors from learner corpora and mathematical reasoning errors generated by large language models show the automatically constructed taxonomies achieve strong performance across multiple evaluation metrics, including exclusivity, coverage, consistency, and usability, and achieve competitive performance compared with manually designed taxonomies, outperforming them in several settings.

1 Introduction

Error analysis plays an important role in many natural language processing and educational applications, including grammatical error correction, automated essay scoring and intelligent tutoring systems (Dulay and Burt, 1972; Dodigovic, 2007; Heift and Schulze, 2007; Ye et al., 2022, 2023a,b; Huang et al., 2023; Zhang et al., 2023a; Li et al., 2023, 2025b,c). A common approach for understanding error patterns is through an *error taxonomy*, which organizes error instances into structured categories according to their underlying linguistic or reasoning mechanisms (Ma et al., 2022; Ye et al., 2025; Fei et al., 2023). Figure 1 illustrates a simple example of classifying an error instance

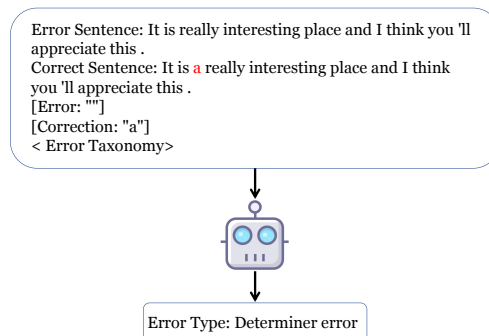


Figure 1: Example of mapping an error instance to an abstract error type using an error taxonomy.

into an abstract error type using an error taxonomy. Such taxonomies provide an essential bridge between individual error instances and higher-level explanations, enabling systematic error analysis, model evaluation, and targeted feedback generation (Corder, 2015; Bialystok et al., 1982).

Error taxonomies have been widely used across different domains. For example, grammatical error taxonomies are commonly adopted in learner corpus research to analyze language acquisition patterns (Ma et al., 2022; Ye et al., 2025; Fei et al., 2023), while reasoning error taxonomies have recently been used to study the failure modes of large language models in mathematical word problem solving and reasoning tasks (Sun et al., 2025; McNichols et al., 2023; Li et al., 2024). However, most existing error taxonomies are manually designed by experts. Although expert-designed taxonomies often provide valuable insights, they suffer from several practical limitations. First, manual construction is labor-intensive and difficult to scale when large datasets or multiple domains are involved. Second, category definitions often rely heavily on expert intuition, which may introduce subjective inconsistencies across datasets and research communities (Politzer and Ramirez, 1973; Tucker et al., 1974; Bryant et al., 2017). Third, as the distribution

| | | | |
|-----|--|--|-----|
| 070 | of errors evolves with new models and tasks, manu- | Our contributions are summarized as follows ¹ : | 120 |
| 071 | ally maintaining and extending such taxonomies | | |
| 072 | becomes increasingly challenging. | | |
| 073 | Recent advances in large language models | | |
| 074 | (LLMs) have demonstrated strong capabilities in | | |
| 075 | semantic abstraction and structure induction (Li | | |
| 076 | et al., 2025a; Wang et al., 2023; Zeng et al., 2024), | • We formulate the task of automatic error tax- | 121 |
| 077 | enabling new possibilities for automatically orga- | onomy construction and propose an LLM- | 122 |
| 078 | nizing textual information. However, automatically | based framework that derives hierarchical er- | 123 |
| 079 | constructing error taxonomies remains challenging. | ror taxonomies directly from error instances. | 124 |
| 080 | Unlike traditional concept taxonomy construction, | • We introduce an error-oriented representa- | 125 |
| 081 | where hierarchical relations are defined over exist- | tion that captures error semantics and error- | 126 |
| 082 | ing concepts, error taxonomy construction requires | correction transformation signals, enabling | 127 |
| 083 | abstracting category definitions directly from di- | more reliable clustering of error patterns. | 128 |
| 084 | verse error instances. Furthermore, error data of- | | |
| 085 | ten contains paired information describing the rela- | • We conduct experiments on English gram- | 129 |
| 086 | tionship between erroneous outputs and their cor- | matical and mathematical reasoning errors , | 130 |
| 087 | rected counterparts, which provides important sig- | demonstrating that the proposed framework | 131 |
| 088 | nals about the underlying causes of errors but is | constructs high-quality taxonomies competi- | 132 |
| 089 | rarely utilized in existing approaches. | tive with expert-designed ones. | 133 |
| 090 | In this work, we propose an LLM-based frame- | | |
| 091 | work for the automatic construction of error tax- | 2 Related Work | 134 |
| 092 | onomies . Our framework is primarily designed for | | |
| 093 | constructing taxonomies of English grammatical | 2.1 Taxonomy Induction | 135 |
| 094 | errors, where large-scale learner data is available | | |
| 095 | and fine-grained error categorization is essential for | Automatic taxonomy induction has long been stud- | 136 |
| 096 | language learning and evaluation. Our approach | ied for knowledge organization and semantic struc- | 137 |
| 097 | first leverages a large language model to abstract | ture learning. Early work mainly relied on lexi- | 138 |
| 098 | concise error type labels from error instances. We | cal-syntactic patterns and graph-based methods to | 139 |
| 099 | then construct semantic representations that jointly | extract hypernym-hyponym relations from large | 140 |
| 100 | encode the semantics of the generated error types | text corpora (Kozareva and Hovy, 2010). Later | 141 |
| 101 | and the transformation signals between erroneous | studies modeled hierarchical relations in continu- | 142 |
| 102 | and correct expressions. Based on these represen- | ous embedding spaces, for example by learning | 143 |
| 103 | tations, an adaptive hierarchical clustering algo- | transformations between hyponym and hypernym | 144 |
| 104 | rithm is used to construct the taxonomy structure, | embeddings (Fu et al., 2014), or by formulating | 145 |
| 105 | while an LLM-based naming module generates | taxonomy construction as a global structure opti- | 146 |
| 106 | interpretable category labels for each cluster, fol- | mization problem such as reinforcement learning | 147 |
| 107 | lowed by a structure optimization stage that refines | based taxonomy generation (Mao et al., 2018). | 148 |
| 108 | the hierarchy and improves semantic consistency. | Recent work has begun to explore the use of | 149 |
| 109 | To evaluate the generalization ability of the | large language models for taxonomy construction. | 150 |
| 110 | framework, we further apply it to a different error | Some approaches use language models to predict | 151 |
| 111 | analysis scenario, namely mathematical reasoning | candidate hierarchical relations before building | 152 |
| 112 | errors. Experimental results show that the automa- | taxonomy structures under structural constraints | 153 |
| 113 | tically constructed taxonomies achieve strong per- | (Chen et al., 2023). Others directly generate taxon- | 154 |
| 114 | formance across multiple evaluation dimensions. | omy structures with LLMs. For example, Taxon- | 155 |
| 115 | These results demonstrate that the proposed ap- | omyGPT formulates taxonomy construction as a | 156 |
| 116 | proach can effectively construct high-quality er- | conditional text generation task (Li et al., 2025a), | 157 |
| 117 | ror taxonomies across different domains, despite | while Chain-of-Layer incrementally expands tax- | 158 |
| 118 | variations in the optimal taxonomy structure, with | onomy structures with layer-wise generation and | 159 |
| 119 | minimal human intervention. | validation mechanisms (Zeng et al., 2024). | 160 |
| | | However, most existing studies focus on organiz- | 161 |
| | | ing entities or topical concepts. Taxonomy induc- | 162 |
| | | tion for <i>error classification</i> , where categories must | 163 |
| | | be abstracted directly from error instances and may | 164 |
| | | involve error-correction signals, remains relatively | 165 |
| | | underexplored in the literature. | 166 |

¹All codes and data will be released after the review.

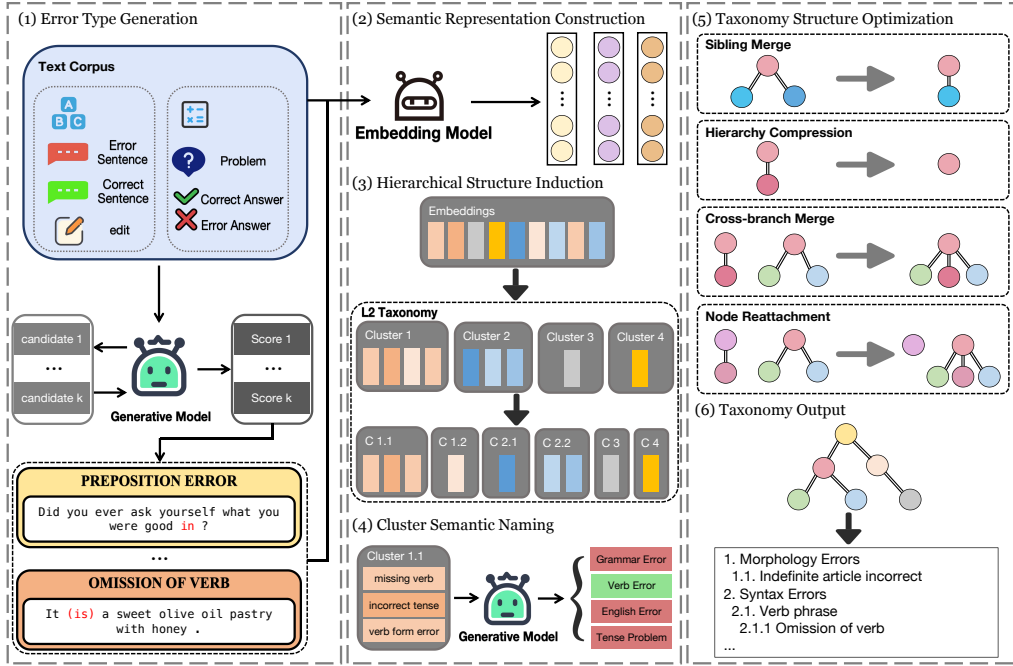


Figure 2: Overview of the proposed LLM-based error taxonomy construction framework

2.2 LLM-based Text Clustering

Text clustering is a widely used approach for discovering latent semantic structures from unlabeled data and is often adopted as a basis for taxonomy construction. Recent work incorporates LLMs into clustering in several complementary ways. One line of research improves text representations using high-quality LLM-generated embeddings before applying traditional clustering algorithms (Petukhova et al., 2025). CLUSTERLLM further leverages LLMs to generate semantic triplet constraints that guide contrastive representation learning for clustering (Zhang et al., 2023b). Another line of work performs clustering directly with LLMs. For example, LLM-MemCluster treats clustering as an iterative labeling process with dynamically maintained cluster labels (Zhu et al., 2025), while other approaches adopt a two-stage framework where LLMs first generate cluster labels and then assign texts to them accordingly (Huang and He, 2025). LLMs have also been used to refine clustering structures, such as boundary-aware reassignment with LLM reasoning (Feng et al., 2024), integrating cluster-level semantic descriptions (Patnaik et al., 2024), or guiding clustering with natural language objectives (Wang et al., 2023).

Despite these advances, most work focuses on topic discovery or semantic similarity grouping, while clustering *error instances* for structured error taxonomy construction remains underexplored.

3 Method

We propose an LLM-based framework for automatically constructing error taxonomies from raw error instances. As illustrated in Figure 2, the framework first generates concise error type labels using an LLM, and then constructs semantic representations that integrate label semantics with transformation signals between erroneous and corrected expressions. Hierarchical clustering is subsequently applied to induce the taxonomy structure, followed by LLM-based cluster naming and a structure optimization stage that removes redundant nodes and improves overall hierarchical consistency.

3.1 Error Type Generation

To construct the taxonomy, we first generate an error type label for each error instance. Given an error instance represented as a pair consisting of an erroneous expression and its corresponding correction, the goal is to summarize the underlying error pattern into a concise semantic label.

Leveraging the semantic reasoning ability of large language models, we generate candidate error type labels for each instance. To improve robustness, we adopt a multi-candidate generation strategy with LLM-based self-evaluation, enabling more reliable label selection. For each error instance x , multiple candidate labels are generated:

$$T(x) = \{t_1, t_2, \dots, t_k\} \quad (1)$$

Each candidate label is evaluated according to four criteria: **Alignment**, **Specificity**, **Reusability**, and **Consistency**. Alignment measures whether the label accurately captures the semantic difference between the erroneous expression and its correction, reflecting the underlying cause of the error. Specificity evaluates whether the label has an appropriate level of abstraction, avoiding overly generic or overly specific descriptions. Reusability measures whether the label can serve as a stable category name that can be reused across multiple error instances. Consistency evaluates whether the label appears consistently across multiple generations, indicating higher semantic stability. The overall quality score of a candidate label is defined as

$$q(t) = s_{\text{align}} + s_{\text{spec}} + s_{\text{reuse}} + s_{\text{cons}} \quad (2)$$

The final error type label is selected as

$$t^* = \arg \max_{t \in T(x)} q(t) \quad (3)$$

3.2 Semantic Representation Construction

After generating error type labels, we construct semantic representations for error instances to support similarity-based taxonomy induction. The representation encodes two complementary signals: the semantic meaning of error type labels and the transformation signals between erroneous and corrected expressions. First, the generated error type label is embedded into a semantic vector

$$\mathbf{e}_{\text{type}} = \text{Embedding}(t^*) \quad (4)$$

Next, we encode both the erroneous expression and the corrected expression and compute their semantic difference

$$\mathbf{e}_{\text{diff}} = \text{Embedding}(x_{\text{correct}}) - \text{Embedding}(x_{\text{error}}) \quad (5)$$

which captures the semantic transformation during error correction. The final representation is obtained by concatenating these two vectors

$$\mathbf{e} = \text{concat}(\mathbf{e}_{\text{type}}, \mathbf{e}_{\text{diff}}) \quad (6)$$

This representation jointly captures abstract error semantics and correction transformations, enabling more accurate and robust measurement of similarity between error instances.

3.3 Hierarchical Structure Induction

Given the semantic representations of error instances, we construct the taxonomy using an adaptive hierarchical clustering procedure. Starting from a root cluster containing all instances, clusters are recursively partitioned to form a hierarchical taxonomy with progressively refined structure. For a cluster C containing n samples, the number of subclusters is determined adaptively as

$$k = \min(K, \max(2, \lfloor n/3 \rfloor)) \quad (7)$$

where K denotes the maximum number of clusters allowed at each split. Cluster compactness is measured using the sum of squared errors (SSE)

$$SSE(C) = \sum_{x \in C} \|x - \mu_C\|^2 \quad (8)$$

where μ_C is the centroid of cluster C . After splitting C into k subclusters C_1, \dots, C_k , the total within-cluster SSE becomes

$$SSE_{\text{split}} = \sum_{i=1}^k SSE(C_i) \quad (9)$$

The improvement ratio is defined as

$$\Delta = \frac{SSE(C) - SSE_{\text{split}}}{SSE(C)} \quad (10)$$

If $\Delta < \tau$, the cluster is considered sufficiently coherent and the splitting process stops; otherwise the resulting subclusters are added as child nodes in the taxonomy tree. The procedure continues recursively until the maximum depth is reached or clusters become stable. The detailed clustering algorithm is provided in Appendix A.

3.4 Cluster Semantic Naming

After clustering, each cluster contains a set of semantically similar error types but lacks an interpretable category name. To transform the clustering structure into a readable taxonomy, we use an LLM to generate candidate semantic labels for each cluster. For a cluster C , the set of error type labels within the cluster is provided as input to the LLM, which produces a set of candidate category names

$$L(C) = \{l_1, l_2, \dots, l_m\} \quad (11)$$

To select the most appropriate label, we evaluate each candidate using two criteria: *coverage* and

304 *confusion*. Coverage measures how well a label
305 represents the semantics of the current cluster

$$306 \quad Cov(l, C) = \frac{1}{|C|} \sum_{t \in C} sim(l, t) \quad (12)$$

307 where $sim(\cdot)$ denotes semantic similarity. Con-
308 fusion measures how well the label also describes
309 other clusters at the same level

$$310 \quad Conf(l, C) = \frac{1}{|S(C)|} \sum_{C' \in S(C)} Cov(l, C') \quad (13)$$

311 where $S(C)$ denotes the set of sibling clusters
312 of C . The final score is defined as

$$313 \quad Score(l, C) = Cov(l, C) - Conf(l, C) \quad (14)$$

314 and the final cluster label is selected as

$$315 \quad l^* = \arg \max_{l \in L(C)} Score(l, C) \quad (15)$$

316 This procedure assigns interpretable and repre-
317 sentative category names to clusters while preserv-
318 ing clear semantic boundaries and minimizing over-
319 lap between sibling clusters.

3.5 Taxonomy Structure Optimization

320 The taxonomy obtained from hierarchical cluster-
321 ing may still contain redundant nodes or struc-
322 turally inconsistent relations, since the hierarchy
323 is constructed primarily from local semantic simi-
324 larity. To further improve the structural quality of
325 the taxonomy, we introduce a structure optimiza-
326 tion stage that iteratively refines the hierarchy. We
327 consider a set of candidate structural operations
328 that modify the taxonomy while preserving its hi-
329 erarchical organization, including sibling merging,
330 hierarchy compression, cross-branch merging, and
331 node reattachment. Details of these operations are
332 provided in Appendix B.

333 To evaluate the quality of a taxonomy T , we
334 define the following objective function
335

$$336 \quad J(T) = \frac{1}{|E(T)|} \sum_{(p,c) \in E(T)} d(p, c) \\ 337 \quad + \frac{1}{|Sib(T)|} \sum_{(u,v) \in Sib(T)} \max(0, \tau_s - d(u, v)) \\ 338 \quad + \frac{|T|}{|T_0|} \quad (16)$$

339 The three terms respectively encourage parent-
child semantic coherence, sibling separability, and
structural compactness, jointly defining a balanced

340 objective for hierarchy refinement. The distance
341 function is defined as

$$342 \quad d(u, v) = \frac{1 - \cos(u, v)}{2} \quad (17)$$

343 where $E(T)$ denotes the set of parent-child
344 edges and $Sib(T)$ denotes the set of sibling node
345 pairs. During optimization, a candidate operation
346 producing taxonomy T' is accepted if it reduces
347 the objective value

$$348 \quad \Delta J = J(T') - J(T) < 0 \quad (18)$$

349 This greedy procedure is applied iteratively un-
350 til no further improvement is achieved, yielding a
351 more compact and semantically coherent taxonomy
352 across the hierarchy.

3.6 Taxonomy Output

353 Finally, the optimized taxonomy is converted into
354 a standardized hierarchical representation. Each
355 node is assigned a hierarchical index through depth-
356 first traversal (e.g., 1, 1.1, 1.1.1), which preserves
357 the parent-child relations within the taxonomy.
358 This representation facilitates downstream appli-
359 cations such as automatic error classification, error
360 analysis, and educational feedback generation.

4 Experiments

4.1 Experimental Setup

362 To evaluate the effectiveness of the proposed auto-
363 matic error taxonomy construction method, experi-
364 ments are conducted on two types of error datasets:
365 grammatical errors from learner corpora and math-
366 ematical reasoning errors generated by large lan-
367 guage models. These datasets provide diverse error
368 instances from different domains, enabling us to
369 construct error taxonomies and evaluate their struc-
370 tural quality using the proposed evaluation metrics.

371 **Datasets.** We conduct experiments on two
372 datasets: a grammatical error dataset from learner
373 corpora and a mathematical reasoning error dataset.
374

375 For grammatical errors, we use the dataset intro-
376 duced in (Zou et al., 2025), derived from the
377 W&I+LOCNESS corpus (Bryant et al., 2019).
378 It contains 487 single-error instances annotated
379 under four representative grammatical error tax-
380 onomies: POL73 (Politzer and Ramirez, 1973),
381 TUC74 (Tucker et al., 1974), BRY17 (Bryant et al.,
382 2017), and FEI23 (Fei et al., 2023), which serve as
383 reference taxonomies in our evaluation.
384

| Taxonomy | Exclusivity \uparrow | Coverage \uparrow | Consistency (ITC / ITS) \uparrow | Usability (Macro / Micro) \uparrow |
|---|------------------------|---------------------|------------------------------------|--------------------------------------|
| (A) Existing Error Taxonomies | | | | |
| POL73 | 0.752 | 0.698 | 0.215 / 0.894 | 0.562 / 0.534 |
| TUC74 | 0.250 | 0.160 | 0.158 / 0.950 | 0.315 / 0.119 |
| BRY17 | 0.928 | 0.980 | 0.288 / 0.931 | 0.626 / 0.821 |
| FEI23 | 0.853 | 0.924 | 0.221 / 0.824 | 0.645 / 0.772 |
| (B1) Proposed Method | | | | |
| Flat KMeans (L1K50) | 0.784 | 0.971 | 0.270 / 0.895 | 0.580 / 0.737 |
| Hierarchical KMeans (L2K15) | 0.820 | 0.985 | 0.310 / 0.937 | 0.565 / 0.719 |
| Hierarchical KMeans (L3K6) | 0.875 | 0.995 | 0.318 / <u>0.945</u> | 0.710 / <u>0.820</u> |
| Hierarchical KMeans (L4K4) | 0.645 | 0.688 | 0.221 / 0.895 | 0.547 / 0.511 |
| (B2) Module Ablation | | | | |
| L3K6 (Single Generation) | 0.842 | 0.972 | 0.311 / 0.938 | 0.689 / 0.804 |
| L3K6 (Only Type) | 0.802 | 0.951 | 0.305 / 0.930 | 0.656 / 0.760 |
| L3K6 (Only Diff) | 0.550 | 0.480 | 0.222 / 0.887 | 0.595 / 0.721 |
| L3K6 (Direct Naming) | 0.868 | 0.993 | <u>0.316</u> / 0.942 | 0.701 / 0.815 |
| L3K6 (w/o Structural Refinement) | 0.740 | 0.998 | 0.300 / 0.935 | 0.633 / 0.721 |
| (B3) Backbone Variants (Embedding / LLM) | | | | |
| L3K6 + KALM Embedding | 0.847 | 0.992 | 0.289 / 0.930 | 0.686 / 0.797 |
| L3K6 + LLaMA Embedding | 0.855 | 0.978 | 0.307 / 0.936 | 0.691 / 0.780 |
| L3K6 + Gemini Generation | <u>0.880</u> | 0.983 | 0.314 / 0.950 | <u>0.703</u> / 0.803 |
| L3K6 + Claude Generation | 0.877 | <u>0.997</u> | <u>0.316</u> / 0.941 | 0.695 / 0.799 |

Table 1: Performance comparison of English error taxonomies. All results are obtained using GPT-5-mini as the evaluator. LxKy denotes a hierarchical clustering configuration with x levels and at most y sub-clusters per parent node. **Bold** indicates the best result and underline indicates the second-best.

For mathematical reasoning errors, we construct an experimental dataset based on MWPEs-300k (Sun et al., 2025), which consists of incorrect reasoning outputs generated by large language models on several math word problem benchmarks, including SVAMP (Patel et al., 2021), GSM8K (Cobbe et al., 2021), MATH (Hendrycks et al., 2021), and AQuA (Ling et al., 2017). Each instance contains a math problem, a model-generated incorrect solution, and the corresponding ground-truth answer.

To keep the experiments computationally manageable, we randomly sample 500 erroneous instances from this dataset and annotate them following the procedure described in (Zou et al., 2025).

Models. Two types of models are used in the taxonomy construction framework: embedding models for semantic representation and generative models for error type and cluster name generation.

For semantic representation, we primarily use Qwen3-Embedding-8B (Zhang et al., 2025). To analyze the impact of embedding models on taxonomy construction, two additional embedding models are included in ablation experiments: KaLM-Embedding (Hu et al., 2025) and LLaMA-Embed-Nemotron-8B (Babakhin et al., 2025).

For error type generation and cluster semantic naming, we primarily use GPT-5-mini (OpenAI,

2024). To evaluate the robustness of the framework across different generative models, we additionally test Claude-Haiku-4.5 (Anthropic, 2025) and Gemini-3-Flash-preview (Google DeepMind, 2025). All generation experiments use a unified prompting template and identical decoding parameters to ensure comparability.

Evaluation Metrics. To evaluate the quality of the constructed error taxonomies, we adopt the taxonomy evaluation framework proposed in (Zou et al., 2025), which includes four dimensions: exclusivity, coverage, consistency, and usability.

Exclusivity evaluates whether category boundaries are clearly defined, coverage measures how well the taxonomy captures observed error instances, consistency assesses structural stability in semantic space through intra-type consistency (ITC) and inter-type separability (ITS), and usability reflects the practical applicability of the taxonomy for automatic classification and human annotation. Detailed definitions of these metrics are provided in Appendix C.

4.2 Main Results and Analysis

The main experimental results are presented and summarized in Table 1.

Comparison with Existing Taxonomies. Compared with manually designed taxonomies, the au-

| Taxonomy | Exclusivity \uparrow | Coverage \uparrow | Consistency (ITC / ITS) \uparrow | Usability (Macro / Micro) \uparrow |
|------------------------------|------------------------|---------------------|------------------------------------|--------------------------------------|
| (A) Original Taxonomy | | | | |
| Original | <u>0.572</u> | 0.618 | <u>0.765</u> / 0.250 | <u>0.206</u> / 0.364 |
| (B) Ours | | | | |
| Flat KMeans (L1K50) | 0.647 | 0.992 | 0.793 / 0.192 | 0.540 / 0.572 |
| Hierarchical KMeans (L2K15) | 0.541 | <u>0.880</u> | 0.724 / <u>0.320</u> | 0.144 / 0.320 |
| Hierarchical KMeans (L3K6) | 0.473 | 0.764 | 0.730 / 0.318 | 0.128 / 0.336 |
| Hierarchical KMeans (L4K4) | 0.193 | 0.526 | 0.733 / 0.323 | 0.110 / 0.180 |

Table 2: Performance comparison of Math error taxonomies. All results are obtained using GPT-5-mini as the evaluator. **Bold** indicates the best result and underline indicates the second-best.

| Setting | Exclusivity | Coverage | Consistency (ITC / ITS) | Usability (Macro / Micro) |
|-----------------------------|-------------------|-------------------|---------------------------------------|---------------------------------------|
| Flat KMeans (L1K50) | 0.783 \pm 0.006 | 0.970 \pm 0.004 | 0.269 \pm 0.005 / 0.894 \pm 0.006 | 0.579 \pm 0.009 / 0.736 \pm 0.010 |
| Hierarchical KMeans (L2K15) | 0.819 \pm 0.008 | 0.984 \pm 0.006 | 0.309 \pm 0.007 / 0.936 \pm 0.005 | 0.563 \pm 0.010 / 0.717 \pm 0.011 |
| Hierarchical KMeans (L3K6) | 0.874 \pm 0.006 | 0.994 \pm 0.003 | 0.317 \pm 0.005 / 0.944 \pm 0.004 | 0.708 \pm 0.008 / 0.819 \pm 0.009 |
| Hierarchical KMeans (L4K4) | 0.642 \pm 0.015 | 0.686 \pm 0.017 | 0.220 \pm 0.012 / 0.893 \pm 0.010 | 0.545 \pm 0.018 / 0.509 \pm 0.020 |

Table 3: Robustness analysis under different random seeds (mean \pm standard deviation).

439 tomatically constructed taxonomy achieves more
440 balanced performance across different metrics. In
441 particular, the L3K6 configuration achieves the
442 highest intra-type consistency (ITC = 0.318) and
443 the best Macro F1 score (0.710), while maintain-
444 ing strong performance on other metrics. It also
445 achieves higher coverage (0.995) than the best man-
446 ual taxonomy BRY17 (0.980), indicating that data-
447 driven semantic grouping can capture a broader
448 range of error patterns. Although BRY17 achieves
449 the highest exclusivity and Micro F1, its perfor-
450 mance is less balanced across different metrics.

451 **Effect of Hierarchical Structure.** Further
452 analysis shows that hierarchical structure signif-
453 icantly influences taxonomy quality. Flat cluster-
454 ing (L1K50) achieves high coverage but performs
455 worse in exclusivity and usability, suggesting that
456 purely flat structures fail to capture meaningful se-
457 mantic organization. Introducing hierarchical struc-
458 tures improves consistency and usability. Among
459 all configurations, the three-level structure (L3K6)
460 provides the best balance between semantic granu-
461 larity and structural stability. However, increasing
462 the hierarchy depth further (L4K4) leads to per-
463 formance degradation across most metrics, indi-
464 cating that excessive hierarchy depth may cause
465 over-fragmentation of categories.

466 **Module Ablation.** Ablation experiments further
467 reveal the contributions of different components.
468 Removing multi-candidate generation (Single Gen-
469 eration) slightly reduces performance, demonstrat-

470 ing the importance of generation diversity and self-
471 evaluation filtering. Using only semantic error type
472 labels (Only Type) leads to moderate degradation,
473 while relying only on difference embeddings (Only
474 Diff) causes a substantial performance drop across
475 metrics. This indicates that semantic error labels
476 provide the primary information, while transforma-
477 tion signals serve as complementary cues. Remov-
478 ing structural refinement slightly increases cover-
479 age but decreases exclusivity and consistency, sug-
480 gesting that structure optimization helps eliminate
481 redundant nodes and clarify category boundaries.

482 **Backbone Robustness.** Replacing the embed-
483 ding model or generation model only causes mi-
484 nor performance changes. Across all backbone
485 variants, the L3K6 structure consistently achieves
486 strong results, indicating that the effectiveness of
487 the proposed framework mainly comes from the
488 overall taxonomy construction pipeline rather than
489 relying on a specific model.

490 4.3 Mathematical Error Taxonomy 491 Evaluation

492 To further evaluate the generalization ability of the
493 proposed method, experiments are conducted on
494 the mathematical error dataset. The results are
495 shown in Table 2 for detailed comparison.

496 Compared with grammatical error analysis, the
497 mathematical task exhibits a different performance
498 pattern. Flat KMeans achieves the best overall
499 results across most metrics, including coverage
500 (0.992), exclusivity (0.647), intra-type consistency

(ITC = 0.793), and usability (0.540 / 0.572). These results indicate that flat clustering is more suitable for organizing mathematical reasoning errors. This difference arises from the nature of reasoning errors in mathematical problem solving. Such errors are typically local computation or reasoning mistakes rather than hierarchical conceptual relations. As a result, introducing deeper hierarchical structures tends to fragment semantically related errors and leads to lower exclusivity and usability.

Although Flat KMeans does not utilize the full hierarchical structure of our framework, the automatically constructed taxonomies still significantly outperform the manually designed taxonomy across most metrics. In particular, the coverage improves from 0.618 to 0.992 and usability increases substantially. This demonstrates that the main improvements come from the semantic representation learning and error type generation stages in the proposed pipeline, which enable the taxonomy to better capture the distribution of reasoning error patterns than manually designed error taxonomy.

4.4 Robustness to Random Initialization

Since the hierarchical clustering stage uses KMeans, different random seeds may affect the results. To evaluate robustness, we repeat the taxonomy construction process five times with different seeds and report the mean and standard deviation of the evaluation metrics in Table 3.

Overall, the variance across different seeds is small, indicating that the proposed method is robust to initialization randomness. In particular, the best-performing configuration L3K6 shows very stable results with minimal standard deviations across all metrics. Although deeper hierarchies such as L4K4 exhibit slightly larger fluctuations, the variations remain within a reasonable range.

These results demonstrate that the performance of the proposed taxonomy construction framework mainly arises from the overall pipeline design rather than specific initialization conditions.

4.5 Human Annotation Agreement

To further evaluate the usability of automatically constructed taxonomies, we measure inter-annotator agreement using Cohen’s Kappa (κ) (Cohen, 1960), which quantifies agreement beyond chance:

$$\kappa = \frac{p_o - p_e}{1 - p_e}, \quad (19)$$

| Taxonomy | A1 & A2 | A2 & A3 | A1 & A3 | Average |
|--------------------------------------|--------------|--------------|--------------|--------------|
| (A) Existing Error Taxonomies | | | | |
| POL73 | 0.663 | 0.619 | 0.632 | 0.638 |
| TUC74 | 0.641 | 0.531 | 0.557 | 0.576 |
| BRY17 | 0.817 | 0.767 | 0.803 | 0.796 |
| FEI23 | 0.703 | 0.770 | 0.711 | 0.728 |
| (B) Proposed Method | | | | |
| Flat KMeans (L1K50) | 0.777 | 0.749 | 0.733 | 0.753 |
| Hierarchical KMeans (L2K15) | 0.783 | 0.737 | 0.748 | 0.756 |
| Hierarchical KMeans (L3K6) | 0.861 | 0.795 | 0.820 | 0.825 |
| Hierarchical KMeans (L4K4) | 0.703 | 0.695 | 0.711 | 0.703 |

Table 4: Inter-annotator agreement (Cohen’s κ) for English error taxonomies.

where p_o denotes the observed agreement between annotators and p_e denotes the expected agreement estimated from marginal label distributions across all categories.

As shown in Table 4, the automatically constructed taxonomy L3K6 achieves the highest agreement score (0.825), outperforming both other automatic structures and manually designed taxonomies. This indicates that the proposed taxonomy provides clearer semantic boundaries and more interpretable category definitions. Flat KMeans and L2K15 show slightly lower agreement, while deeper hierarchies such as L4K4 yield lower agreement due to increased structural complexity. These results are consistent with the usability metrics observed in the main experiments.

5 Conclusion

We present an LLM-based framework for automatically constructing error taxonomies from raw error instances. Our method generates semantic error type labels using a large language model, builds hybrid representations that integrate label semantics with error-correction transformation signals, and induces taxonomy structures through hierarchical clustering followed by semantic naming and structure refinement. Experiments on grammatical errors and mathematical reasoning errors show that the automatically constructed taxonomies achieve strong performance across multiple evaluation metrics and remain competitive with manually designed taxonomies. These results demonstrate that large language models can effectively support the automatic construction of structured error taxonomies across different domains with minimal human intervention.

584 Limitations

585 Although the proposed framework shows strong
586 performance in automatically constructing error
587 taxonomies, several limitations remain.

588 First, the hierarchical clustering stage still re-
589 quires predefined structural parameters, such as the
590 maximum hierarchy depth and the upper bound on
591 the number of clusters per split. While the adaptive
592 splitting strategy helps mitigate this issue, the over-
593 all taxonomy structure still depends on these preset
594 parameters. In future work, it would be interesting
595 to explore fully adaptive approaches where both
596 the hierarchy depth and the number of clusters can
597 be determined automatically from the data.

598 Second, the current framework is primarily de-
599 signed and evaluated on English grammatical er-
600 rors. When applied to other domains, such as math-
601 ematical reasoning errors, certain modules in the
602 pipeline become less effective due to differences in
603 the structural properties of error patterns. Although
604 the overall pipeline can still produce reasonable
605 taxonomies in such scenarios, developing a more
606 domain-general taxonomy construction framework
607 remains an important direction for future research.

608 Ethical Considerations

609 We conduct our experiments using publicly avail-
610 able datasets, including learner corpus data de-
611 rived from W&I-LOCNESS and mathematical rea-
612 soning error data constructed from MWPEs-300k.
613 These datasets do not contain personally identi-
614 fiable or sensitive information, and we use them
615 solely for research purposes in accordance with
616 their original licenses and intended usage.

617 We use publicly available language models and
618 embedding models in our experiments and properly
619 acknowledge their sources. All models are applied
620 following their recommended usage guidelines.

621 For the annotation-related parts of this study, we
622 recruited three postgraduate students with back-
623 grounds in linguistics and related fields as annota-
624 tors. We compensated them at a standard hourly
625 rate and ensured that the workload remained rea-
626 sonable. The annotation process involved routine
627 research tasks and did not expose annotators to
628 sensitive or harmful content.

629 Despite these considerations, potential risks re-
630 main. LLM-based label generation may introduce
631 biases or inconsistencies, and the resulting tax-
632 onomies may be over-interpreted in downstream
633 applications. Our framework mitigates these risks

through repeated label generation, hybrid represen- 634
tations, and structure refinement, which improve 635
consistency and robustness. Nevertheless, the con- 636
structed taxonomies should be treated as supportive 637
analytical tools rather than definitive standards. 638

References 639

- Anthropic. 2025. Introducing claude haiku 4.5. <https://www.anthropic.com/news/claude-haiku-4-5>.
640
641
642
Published: Oct. 15, 2025.
- Yauhen Babakhin, Radek Osmulski, Ronay Ak, 643
Gabriel Moreira, Mengyao Xu, Benedikt Schifferer, 644
Bo Liu, and Even Oldridge. 2025. [Llama-embed- 645
nemotron-8b: A universal text embedding model 646
for multilingual and cross-lingual tasks](#). *Preprint*, 647
arXiv:2511.07025. 648
- Ellen Bialystok, Heidi Dulay, Marina Burt, and Stephen 649
Krashen. 1982. [Language two](#). *The Modern Lan- 650
guage Journal*, 67:273. 651
- Christopher Bryant, Mariano Felice, Øistein E Ander- 652
sen, and Ted Briscoe. 2019. The bea-2019 shared 653
task on grammatical error correction. In *Proceedings 654
of the fourteenth workshop on innovative use of NLP 655
for building educational applications*, pages 52–75. 656
- Christopher Bryant, Mariano Felice, and Ted Briscoe. 657
2017. Automatic annotation and evaluation of error 658
types for grammatical error correction. In *Proceed- 659
ings of the 55th annual meeting of the association for 660
computational linguistics (Volume 1: Long Papers)*, 661
pages 793–805. 662
- Boqi Chen, Fandi Yi, and Dániel Varró. 2023. [Prompt- 663
ing or fine-tuning? a comparative study of large 664
language models for taxonomy construction](#). In 665
*2023 ACM/IEEE International Conference on Model 666
Driven Engineering Languages and Systems Com- 667
panion (MODELS-C)*, page 588–596. IEEE Press. 668
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, 669
Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias 670
Plappert, Jerry Tworek, Jacob Hilton, Reiichiro 671
Nakano, Christopher Hesse, and John Schulman. 672
2021. [Training verifiers to solve math word prob- 673
lems](#). *Preprint*, arXiv:2110.14168. 674
- Jacob Cohen. 1960. A coefficient of agreement for 675
nominal scales. *Educational and psychological mea- 676
surement*, 20(1):37–46. 677
- Stephen Pit Corder. 2015. The significance of learners’ 678
errors. In *Error analysis*, pages 19–27. Routledge. 679
- Marina Dodigovic. 2007. [Artificial intelligence and 680
second language learning: An efficient approach to 681
error remediation](#). *Language Awareness*, 16(2):99– 682
113. 683
- Heidi C Dulay and Marina K Burt. 1972. Goofing: 684
An indicator of children’s second language learning 685
strategies 1. *Language learning*, 22(2):235–252. 686

| | | | |
|-----|--|---|-----|
| 687 | Yuejiao Fei, Leyang Cui, Sen Yang, and 1 others. 2023. | induction using llms: An enhanced framework | 743 |
| 688 | Enhancing grammatical error correction systems with | by integrating doubly-checked mechanism and self- | 744 |
| 689 | explanations. In <i>Proceedings of the 61st Annual</i> | evaluation strategy. In <i>China Conference on Knowl-</i> | 745 |
| 690 | <i>Meeting of the Association for Computational Lin-</i> | <i>edge Graph and Semantic Computing and Interna-</i> | 746 |
| 691 | <i>guistics</i> , pages 7489–7501. | <i>tional Joint Conference on Knowledge Graphs</i> , pages | 747 |
| 692 | Zijin Feng, Luyang Lin, Lingzhi Wang, Hong Cheng, | 134–146, Singapore. Springer Nature Singapore. | 748 |
| 693 | and Kam-Fai Wong. 2024. LLMEdgeRefine: En- | Xiaoyuan Li, Wenjie Wang, Moxin Li, Junrong Guo, | 749 |
| 694 | hancing text clustering with LLM-based boundary | Yang Zhang, and Fuli Feng. 2024. Evaluating mathe- | 750 |
| 695 | point refinement . In <i>Proceedings of the 2024 Confer-</i> | matical reasoning of large language models: A focus | 751 |
| 696 | <i>ence on Empirical Methods in Natural Language Pro-</i> | on error identification and correction . In <i>Findings of</i> | 752 |
| 697 | <i>cessing</i> , pages 18455–18462, Miami, Florida, USA. | <i>the Association for Computational Linguistics: ACL</i> | 753 |
| 698 | Association for Computational Linguistics. | 2024, pages 11316–11360, Bangkok, Thailand. As- | 754 |
| 699 | Ruiji Fu, Jiang Guo, Bing Qin, Wanxiang Che, Haifeng | sociation for Computational Linguistics. | 755 |
| 700 | Wang, and Ting Liu. 2014. Learning semantic hier- | Yinghui Li, Haojing Huang, Shirong Ma, Yong Jiang, | 756 |
| 701 | archies via word embeddings . In <i>Proceedings of the</i> | Yangning Li, Feng Zhou, Hai-Tao Zheng, and Qingyu | 757 |
| 702 | <i>52nd Annual Meeting of the Association for Compu-</i> | Zhou. 2023. On the (in) effectiveness of large lan- | 758 |
| 703 | <i>tational Linguistics (Volume 1: Long Papers)</i> , pages | guage models for chinese text correction. <i>arXiv</i> | 759 |
| 704 | 1199–1209, Baltimore, Maryland. Association for | <i>preprint arXiv:2307.09007</i> . | 760 |
| 705 | Computational Linguistics. | Yinghui Li, Shirong Ma, Shaoshen Chen, Haojing | 761 |
| 706 | Google DeepMind. 2025. Best for frontier intelligence | Huang, Shulin Huang, Yangning Li, Hai-Tao Zheng, | 762 |
| 707 | at speed. https://deepmind.google/models/ | and Ying Shen. 2025b. Correct like humans: Progres- | 763 |
| 708 | gemini/flash/ . | sive learning framework for chinese text error correc- | 764 |
| 709 | Trude Heift and Mathias Schulze. 2007. <i>Errors and In-</i> | tion. <i>Expert Systems with Applications</i> , 265:126039. | 765 |
| 710 | <i>telligence in Computer-Assisted Language Learning:</i> | Yinghui Li, Shang Qin, Jingheng Ye, Haojing Huang, | 766 |
| 711 | <i>Parsers and Pedagogues</i> . Routledge. | Yangning Li, Shu-Yu Guo, Libo Qin, Xuming Hu, | 767 |
| 712 | Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul | Wenhao Jiang, Hai-Tao Zheng, and 1 others. 2025c. | 768 |
| 713 | Arora, Steven Basart, Eric Tang, Dawn Song, and | Rethinking the roles of large language models in chi- | 769 |
| 714 | Jacob Steinhardt. 2021. Measuring mathematical | nese grammatical error correction. In <i>Proceedings of</i> | 770 |
| 715 | problem solving with the math dataset . | <i>the 63rd Annual Meeting of the Association for Com-</i> | 771 |
| 716 | Xinshuo Hu, Zifei Shan, Xinping Zhao, Zetian Sun, | <i>putational Linguistics (Volume 6: Industry Track)</i> , | 772 |
| 717 | Zhenyu Liu, Dongfang Li, Shaolin Ye, Xinyuan | pages 553–567. | 773 |
| 718 | Wei, Qian Chen, Baotian Hu, Haofen Wang, Jun | Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blun- | 774 |
| 719 | Yu, and Min Zhang. 2025. Kalm-embedding: Supe- | som. 2017. Program induction by rationale genera- | 775 |
| 720 | rior training data brings a stronger embedding model. | tion: Learning to solve and explain algebraic word | 776 |
| 721 | <i>Preprint</i> , arXiv:2501.01028. | problems . In <i>Proceedings of the 55th Annual Meet-</i> | 777 |
| 722 | Chen Huang and Guoxiu He. 2025. Text clustering | <i>ing of the Association for Computational Linguistics</i> | 778 |
| 723 | as classification with llms . In <i>Proceedings of the</i> | <i>(Volume 1: Long Papers)</i> , pages 158–167, Vancouver, | 779 |
| 724 | <i>2025 Annual International ACM SIGIR Conference</i> | Canada. Association for Computational Linguistics. | 780 |
| 725 | <i>on Research and Development in Information Re-</i> | Shirong Ma, Yinghui Li, Rongyi Sun, Qingyu Zhou, | 781 |
| 726 | <i>trieval in the Asia Pacific Region, SIGIR-AP 2025,</i> | Shulin Huang, Ding Zhang, Li Yangning, Ruiyang | 782 |
| 727 | page 374–384, New York, NY, USA. Association for | Liu, Zhongli Li, Yunbo Cao, and 1 others. 2022. Lin- | 783 |
| 728 | Computing Machinery. | guistic rules-based corpus generation for native chi- | 784 |
| 729 | Haojing Huang, Jingheng Ye, Qingyu Zhou, Yinghui | nese grammatical error correction. In <i>Findings of the</i> | 785 |
| 730 | Li, Yangning Li, Feng Zhou, and Hai-Tao Zheng. | <i>Association for Computational Linguistics: EMNLP</i> | 786 |
| 731 | 2023. A frustratingly easy plug-and-play detection- | 2022, pages 576–589. | 787 |
| 732 | and-reasoning module for chinese spelling check. In | Yuning Mao, Xiang Ren, Jiaming Shen, Xiaotao Gu, | 788 |
| 733 | <i>Findings of the Association for Computational Lin-</i> | and Jiawei Han. 2018. End-to-end reinforcement | 789 |
| 734 | <i>guistics: EMNLP 2023</i> , pages 11514–11525. | learning for automatic taxonomy induction . In <i>Pro-</i> | 790 |
| 735 | Zornitsa Kozareva and Eduard Hovy. 2010. A semi- | <i>ceedings of the 56th Annual Meeting of the Associa-</i> | 791 |
| 736 | supervised method to learn and construct taxonomies | <i>tion for Computational Linguistics (Volume 1: Long</i> | 792 |
| 737 | using the web . In <i>Proceedings of the 2010 Confer-</i> | <i>Papers)</i> , pages 2462–2472, Melbourne, Australia. As- | 793 |
| 738 | <i>ence on Empirical Methods in Natural Language</i> | sociation for Computational Linguistics. | 794 |
| 739 | <i>Processing</i> , pages 1110–1118, Cambridge, MA. As- | Hunter McNichols, Mengxue Zhang, and Andrew Lan. | 795 |
| 740 | sociation for Computational Linguistics. | 2023. Algebra error classification with large lan- | 796 |
| 741 | Jiaye Li, Yuan Meng, Lijun Wang, Tianhao Qian, | guage models . In <i>Artificial Intelligence in Education:</i> | 797 |
| 742 | Songlin Zhai, and Guilin Qi. 2025a. Taxonomy | <i>24th International Conference, AIED 2023, Tokyo,</i> | 798 |
| | | <i>Japan, July 3–7, 2023, Proceedings</i> , page 365–376, | 799 |
| | | Berlin, Heidelberg. Springer-Verlag. | 800 |

Algorithm 1 Adaptive Hierarchical Clustering for Error Taxonomy Construction

Require: Embeddings \mathbf{X} , maximum depth L , maximum clusters per split K , minimum cluster size n_{\min} , threshold τ

Ensure: Taxonomy tree \mathcal{T}

- 1: Initialize root cluster $C^{(0)} \leftarrow \mathbf{X}$ and taxonomy tree \mathcal{T}
- 2: **for** $\ell = 1$ to L **do**
- 3: **for** each cluster C at level $\ell - 1$ **do**
- 4: **if** $|C| < n_{\min}$ **then**
- 5: mark C as leaf
- 6: **continue**
- 7: **end if**
- 8: $k \leftarrow \min(K, \max(2, \lfloor |C|/3 \rfloor))$
- 9: $\{C_1, \dots, C_k\} \leftarrow \text{KMEANSPLIT}(C, k)$
- 10: compute $\Delta = \frac{SSE(C) - SSE_{split}}{SSE(C)}$
- 11: **if** $\Delta < \tau$ **then**
- 12: mark C as leaf
- 13: **continue**
- 14: **end if**
- 15: **for** $i = 1$ to k **do**
- 16: add C_i as child of C in \mathcal{T}
- 17: **end for**
- 18: **end for**
- 19: **end for**
- 20: **return** \mathcal{T}

B Structural Optimization Operations

To refine the taxonomy structure, we apply several structural operations that modify the hierarchy while preserving the overall tree organization. These operations aim to remove redundant nodes and improve semantic consistency within the taxonomy.

- **Sibling Merge.** If two sibling nodes under the same parent exhibit high semantic similarity, they are merged into a single node. The merged node inherits the union of their subtrees, reducing redundant categories.
- **Hierarchy Compression.** If a parent node contains only one child and the two nodes are semantically similar, the intermediate hierarchy level is removed by absorbing the child into the parent node.
- **Cross-branch Merge.** If nodes located in different branches but at the same hierarchy level are semantically similar, they may be merged into a single node, and their subtrees are unified.
- **Node Reattachment.** If a node is semantically more similar to another potential parent than its current parent, it can be reattached to the more appropriate parent while preserving its subtree structure.

These operations are evaluated using the objective function described in Section 3.5, and are applied iteratively until no further improvement can be achieved.

C Detailed definitions of Evaluation Metrics

This section briefly describes the evaluation metrics used in our experiments. All metrics follow the taxonomy evaluation framework proposed in (Zou et al., 2025). Here we provide a high-level description of how each metric is computed.

Exclusivity Exclusivity measures whether the error types in a taxonomy are mutually exclusive. The metric analyzes the confidence scores produced by large language models when assigning error categories to instances. If multiple error types receive high confidence for the same instance, this indicates overlapping category boundaries and reduces the exclusivity score. The final score is obtained by aggregating overlap statistics across the dataset.

Coverage Coverage evaluates how well the taxonomy accounts for the observed error instances. It is defined as the proportion of instances that can be assigned to predefined error categories (excluding “Other” categories). Higher coverage indicates that the taxonomy captures a larger portion of the observed error space.

Consistency Consistency measures the structural coherence of the taxonomy in semantic space. Error instances are mapped into an embedding space, and two complementary properties are evaluated: intra-type consistency (ITC) and inter-type separability (ITS). ITC measures the compactness of instances belonging to the same error type, while ITS evaluates the separability between different error types.

Usability Usability reflects the practical applicability of a taxonomy in real-world error classification scenarios. It is evaluated from two perspectives: model-level usability and human-level usability. Model-level usability measures how effectively models can apply the taxonomy for automatic classification, while human-level usability evaluates the consistency of human annotations when applying the taxonomy.