

Learning from Near-Misses: Error-Aware Contrastive Few-Shot Learning for NL2Formula

Anonymous ACL submission

Abstract

Natural Language to Excel Formula (NL2Formula) translates user intent into executable spreadsheet formulas. However, current models often produce *near-miss* outputs—formulas that parse correctly yet fail at execution due to an incorrect function, operator, or reference. Through a systematic error analysis, we find that these errors repeatedly arise from a small set of structural decision points, motivating the need for typed error supervision rather than general error signals. To this end, we introduce an abstract syntax tree (AST)-based error taxonomy that organizes common error modes by the kind of decision that goes wrong in the parse tree. Building on this taxonomy, we propose Error-Aware Contrastive Few-Shot Learning (ECFL), an error-aware framework that unifies training and inference around typed error supervision. During offline training, ECFL mines near-miss errors, assigns error types under the taxonomy, and builds error-aware contrastive demonstrations for fine-tuning. During online inference, a lightweight predictor estimates likely error types and triggers targeted retrieval of contrastive demonstrations to guide single-pass decoding. Experiments show ECFL improves Exact Match (EM) by 6.4 points over supervised fine-tuning (SFT) and matches self-consistency (SC@5) accuracy at substantially lower inference cost.

1 Introduction

Spreadsheet formulas are a backbone of real-world data analysis and lightweight automation (Chen et al., 2024; Aivaloglou et al., 2017), but writing correct formulas remains a major barrier for non-expert users (Chen et al., 2024; Zhao et al., 2024; Powell et al., 2008). While large language models (LLMs) have lowered the barrier by translating natural language into spreadsheet formulas (Srinivasa Ragavan et al., 2022; Zhao et al., 2024), they often struggle with the strict structural constraints

of spreadsheet formulas (Aivaloglou et al., 2017). As a result, models often produce *near-miss* outputs that parse correctly but fail at execution due to a wrong function, operator, or reference. Such near-misses are particularly frustrating because they are usually hard to catch by inspection and can propagate to downstream analyses and decisions (Powell et al., 2008). The long-tail distribution of functions and complex nested compositions further amplify the problem (Chen et al., 2024; Aivaloglou et al., 2017; Zhao et al., 2024), making near-misses both common and costly (Powell et al., 2008).

Through a systematic error analysis (Section 3), we find that these errors are rarely arbitrary but repeatedly arise from a small set of structural decision points in the formula’s abstract syntax tree (AST). Existing methods either discard these errors entirely (standard SFT) or collapse them into a single undifferentiated negative signal (e.g., DPO; Rafailov et al., 2023), losing the structural distinctions needed for targeted correction. Self-correction methods iteratively refine outputs by generating a candidate, executing it (or checking syntax), diagnosing errors, and regenerating, but this loop incurs substantial inference cost (Chen et al., 2023; Zhang et al., 2023). Inference-time scaling methods such as self-consistency (Wang et al., 2023) improve accuracy by generating K candidates and selecting via majority vote, but pay a $K \times$ decoding cost. Overall, these methods do not effectively exploit near-miss errors as structured supervision. We therefore hypothesize that typed near-miss supervision can teach type-specific correction patterns that generic preference signals and brute-force sampling do not capture.

To operationalize this idea, we introduce an AST-based *error taxonomy* that types near-misses by the kind of structural choice that goes wrong in the parse tree. To operationalize this idea, we introduce an AST-based *error taxonomy* that types near-misses by the kind of structural choice that goes

wrong in the parse tree. Typing follows a skeleton-first rule that prioritizes structural mismatches in the AST skeleton over leaf-level substitutions (Section 3.3). Only when the skeletons match do we type fill-level errors such as reference or value mismatches. Building on this taxonomy, we propose ECFL (Error-Aware Contrastive Few-Shot Learning), an error-aware framework that exploits typed near-miss supervision for both training and inference. During offline training, ECFL mines near-miss predictions by Monte Carlo sampling, assigns each a deterministic error type via AST comparison, and stores typed wrong–correct pairs in an error bank. For Low-Rank Adaptation (LoRA) fine-tuning, ECFL constructs contrastive prompts by selecting type-aligned demonstrations from the error bank. During online inference, a lightweight predictor estimates likely error types and retrieves one demonstration per estimated type via query-embedding similarity. The model then decodes once, guided by these type-aligned demonstrations with minimal overhead.

To evaluate reliably, we clean a standard benchmark, NL2Formula (Zhao et al., 2024), and introduce TFBench, an expert-curated benchmark covering 215 functions and 12 operators with nested compositions. Experiments show consistent gains in both structural and execution correctness across standard and expert benchmarks, with larger improvements on long-tail functions, while matching or exceeding SC@5 accuracy using single-pass decoding (one generation plus a forward-only pass for type prediction), instead of $5\times$ full generations required by SC@5.

Our contributions are threefold:

- **A 9-type near-miss error taxonomy** that provides automatic and interpretable error types for NL2Formula results.
- **ECFL**, an error-aware framework that unifies training and inference around typed near-miss supervision.
- **TFBench**, an expert-curated benchmark and cleaned NL2Formula splits for reliable evaluation. Experiments show consistent improvements over strong baselines on both structural and execution correctness.

2 Related Work

Our work connects several research threads; we summarize key distinctions here and provide an extended discussion in Appendix C.

NL2Formula and Code Generation. Prior systems such as SpreadsheetCoder (Chen et al., 2021b) and FLAME (Joshi et al., 2023) frame formula generation as structured prediction over spreadsheet context, but remain brittle on long-tail functions and complex nesting. These methods rely on end-to-end memorization and lack explicit mechanisms to handle the systematic near-miss errors that arise from strict structural constraints.

Learning from Errors. Self-debugging approaches (Chen et al., 2023; Zhang et al., 2023) use execution feedback to iteratively refine outputs, but require multiple generation-execution cycles. Alternatively, preference learning methods like DPO (Rafailov et al., 2023) leverage collected errors as negative signals to sharpen decision boundaries. In contrast, we *mine* errors offline and type them via AST analysis, converting failure patterns into *error-aware* contrastive supervision for both training and inference, rather than treating all errors as uniform preference signals.

Inference-Time Scaling. Methods like self-consistency (Wang et al., 2023) and Tree of Thoughts (Yao et al., 2023) improve robustness via multi-sample aggregation or deliberate search, but incur prohibitive $K\times$ or more inference cost. Our approach avoids multi-sample generation by predicting likely error types and retrieving targeted contrastive demonstrations with a single extra forward-only pass, yielding comparable robustness at substantially lower cost.

3 Near-Miss Error Taxonomy

3.1 Error Analysis

During the analysis of existing NL2Formula benchmark (Zhao et al., 2024) and the construction of our expert-curated TFBench, annotators wrote gold formulas with model-generated candidates as references. Systematic comparison of these candidates against the final gold formulas revealed two consistent patterns.

First, most errors are *near-misses* rather than complete failures. The model typically understands the task but makes localized errors, such as adding or omitting a formula component, confusing similar functions (e.g., SUMIF vs. COUNTIF), or selecting the wrong operator (e.g., > vs. >=). While the overall formula structure remains reasonable, the difficulty lies in precise element selection and correct nesting depth.

Type	Condition	Example	
Skeleton	FuncMismatch	$ S^w = S^* $, diff at func	SUM→AVG
	FuncMissing	$ S^w < S^* $, missing func	IF(SUM)→SUM
	FuncRedundant	$ S^w > S^* $, extra func	A1→SUM(A1)
	OpMismatch	$ S^w = S^* $, diff at op	>→>=
	OpMissing	$ S^w < S^* $, missing op	A11+B11+C11→A11+B11
	OpRedundant	$ S^w > S^* $, extra op	A11*B11→A11*B12*2
Fill	RefMismatch	$S^w=S^*$, diff at ref	A1:A10→A1:A5
	ValMismatch	$S^w=S^*$, diff at val	"Yes"→"Y"
Miscellaneous	Otherwise	COUNT(A:A,"*")→COUNT(A:A)	

Table 1: The 9-type error taxonomy with formal conditions. Skeleton errors (top) arise from structural mismatches; fill errors (middle) arise from leaf-level substitutions under identical structure; miscellaneous errors (bottom) capture residual cases that do not fit the structural or fill definitions.

Second, references and values errors are overwhelmingly *mismatches*, rather than *missing* or *redundant*. Among approximately 10,000 error samples, only 19 cases (<0.2%) involve missing or redundant references and values. This asymmetry suggests that once the model commits to a function or operator, it generally follows the correct argument structure.

To sum up, the main challenge is choosing *which* function/operator to use and *what* reference/value to fill, rather than deciding *how many* arguments to provide. Moreover, errors are typically *structurally localized*: a single term in the formula is wrong while surrounding context remains correct. This locality makes an AST-based method a natural fit for both locating errors and assigning an interpretable error type by deterministic tree comparison. It therefore enables automatic error typing without manual annotation, enabling targeted retrieval of wrong–correct contrastive demonstrations.

3.2 Error Taxonomy

Based on the above analysis, we define a 9-type error taxonomy organized into three levels (Table 1). Figure 1 also illustrates a concrete example. *Skeleton-level* errors (6 types) reflect structural decision failures: function and operator nodes each admit Mismatch, Missing, and Redundant modes. *Fill-level* errors (2 types) occur when the skeleton is identical but leaf content differs, yielding RefMismatch and ValMismatch. Since fill-level Missing/Redundant errors are negligible (<0.2%), we group them with unparseable cases into *Miscellaneous* (1 type) (1 type), which is excluded from contrastive demonstration redemonstration retrieval.

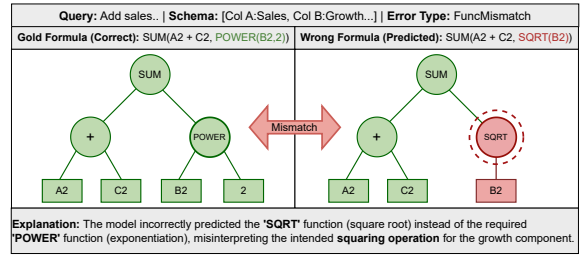


Figure 1: AST-based error typing for formula prediction. The predicted formula incorrectly uses SQRT instead of POWER, resulting in a FuncMismatch error type.

3.3 Error Typing

We formalize error typing as a deterministic function $\tau : (f^w, f^*) \mapsto e$. Given a formula f , its **Abstract Syntax Tree** $\mathcal{T}(f)$ has internal nodes for functions and operators, and leaf nodes for references and values. The **skeleton** $S(f)$ is the sequence of internal node labels obtained by preorder traversal; the **leaf sequence** $L(f)$ is defined analogously for leaves.

Our typing rule follows a **skeleton-first** principle: structural errors take precedence over fill-level errors, as they reflect more fundamental decision failures. Given $S^w = S(f^w)$ and $S^* = S(f^*)$, we first compare skeleton lengths to determine Missing ($|S^w| < |S^*|$), Redundant ($|S^w| > |S^*|$), or Mismatch (equal length but different node type), then identify whether the first differing node is a function or operator. Only when $S^w = S^*$ do we proceed to leaf comparison, classifying errors as RefMismatch or ValMismatch based on the first differing leaf type. Cases that do not fit any pattern are assigned Miscellaneous.

Since both skeleton and leaf sequence are generated by preorder traversal, left-to-right comparison effectively simulates tree traversal: the first mismatched position corresponds to the first erroneous decision point in the AST. This design aligns error localization with the top-down, left-to-right generation order of autoregressive language models, where the first incorrect token often causes subsequent errors. It is easy to show that this taxonomy is exhaustive (every parseable error receives exactly one type), mutually exclusive (skeleton-first prevents ambiguity), and deterministic (pure AST comparison, no execution or manual annotation).

4 Method

Building on the error taxonomy defined in Section 3, we now present ECFL, an error-aware frame-

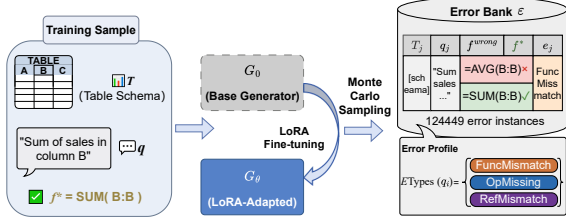


Figure 2: Data structures and notation. A training sample contains table schema T , query q , and gold formula f^* . The error bank \mathcal{E} stores typed error tuples, each pairing a wrong prediction with the correct answer and an AST-derived error label e_j . The error profile $ETypes(q_i)$ tracks which error types a query triggered during sampling.

work that transforms near-miss errors into supervisory signals. We organize ECFL into three components: error bank construction, offline training, and online inference (Figure 3).

4.1 Error Bank Construction

We first construct an error bank as a reusable repository of typed near-miss mistakes mined from a base generator \mathcal{G}_0 (e.g., an untuned LLM). This error bank supports type-conditioned selection of wrong–correct demonstrations for both offline training and online inference. Each training instance of NL2Formula task consists of a table schema T_i , a natural language query q_i , and a gold formula f_i^* . To mine near-miss errors, we perform Monte Carlo sampling on \mathcal{G}_0 by drawing K samples for each (T_i, q_i) under temperature τ and nucleus sampling:

$$f_i^{(k)} \sim p_{\mathcal{G}_0}(f | T_i, q_i; \tau, \text{top-}p), \quad k = 1, \dots, K. \quad (1)$$

We retain an instance only if $f_i^{(k)} \neq f_i^*$ and its error type is not MISCELLANEOUS. Each retained near-miss is stored in the error bank \mathcal{E} as a typed tuple $(T_i, q_i, f_i^{\text{wrong}}, f_i^*, e_i)$. We also update the per-query error profile $ETypes(q_i)$, which records the set of error types observed for q_i during sampling. Figure 2 illustrates these data structures.

4.2 Offline Training

Contrastive LoRA fine-tuning. After constructing the error bank, we build contrastive demonstrations to help the model identify and avoid near-miss errors. For each training query q_i , we select M demonstrations (e.g., $M=3$) from the error bank \mathcal{E} . Each demonstration is formatted as a wrong–correct pair under its table schema, ex-

plícitly highlighting the decision boundary (Figure 4). If $ETypes(q_i)$ is non-empty, we select one demonstration per observed error type by maximizing cosine similarity $\text{Sim}(q_i, q_j)$ between query embeddings (from a sentence encoder like MiniLM (Wang et al., 2020)):

$$d_e = \arg \max_{(T_j, q_j, \dots, e_j) \in \mathcal{E}, e_j = e} \text{Sim}(q_i, q_j). \quad (2)$$

If fewer than N demonstrations are retrieved, we pad the remaining slots with globally most similar demonstrations. We fine-tune the base generator \mathcal{G}_0 with LoRA adapters to obtain \mathcal{G}_θ on these contrastive demonstrations, minimizing standard autoregressive NLL:

$$\mathcal{L}_{\text{SFT}} = - \sum_{t=1}^{|f_i^*|} \log p_{\mathcal{G}_\theta}(f_{i,t}^* | f_{i,<t}^*, d_1, \dots, d_M). \quad (3)$$

Error type predictor. To enable error-aware retrieval at inference time (Section 4.3), we need to predict which error types a query is likely to trigger without access to the gold formula. We train a lightweight multilayer perceptron (MLP) classifier \mathcal{C} to estimate error types from the LLM’s hidden states. Specifically, for each tuple $(T_i, q_i, f_i^{\text{wrong}}, f_i^*, e_i)$ in the error bank \mathcal{E} , we run a forward-only pass through \mathcal{G}_θ with (T_i, q_i) and extract the last-token hidden state \mathbf{h}_i . Each (\mathbf{h}_i, e_i) pair becomes one training instance. The classifier predicts error types via:

$$p(e | q) = \text{softmax}(\mathcal{C}(\mathbf{h}_{\text{last}})), \quad (4)$$

trained with standard cross-entropy loss. This method achieves 78% macro-F1, verifying the effectiveness of our error-type predictor. Further training and evaluation details are provided in Appendix A.

4.3 Online Inference

Given a test instance (T, q) , we first run a forward-only pass through \mathcal{G}_θ to obtain the hidden state \mathbf{h} . The classifier \mathcal{C} outputs a distribution over error types, from which we select the top- M types $\hat{E} = \{\hat{e}_1, \dots, \hat{e}_M\}$. For each $\hat{e}_m \in \hat{E}$, we retrieve the most similar demonstration:

$$d_{\hat{e}_m} = \arg \max_{(T_j, q_j, \dots, e_j) \in \mathcal{E}, e_j = \hat{e}_m} \text{Sim}(q, q_j). \quad (5)$$

The retrieved demonstrations form a contrastive prompt, and \mathcal{G}_θ generates the final formula \hat{f} in a single pass.

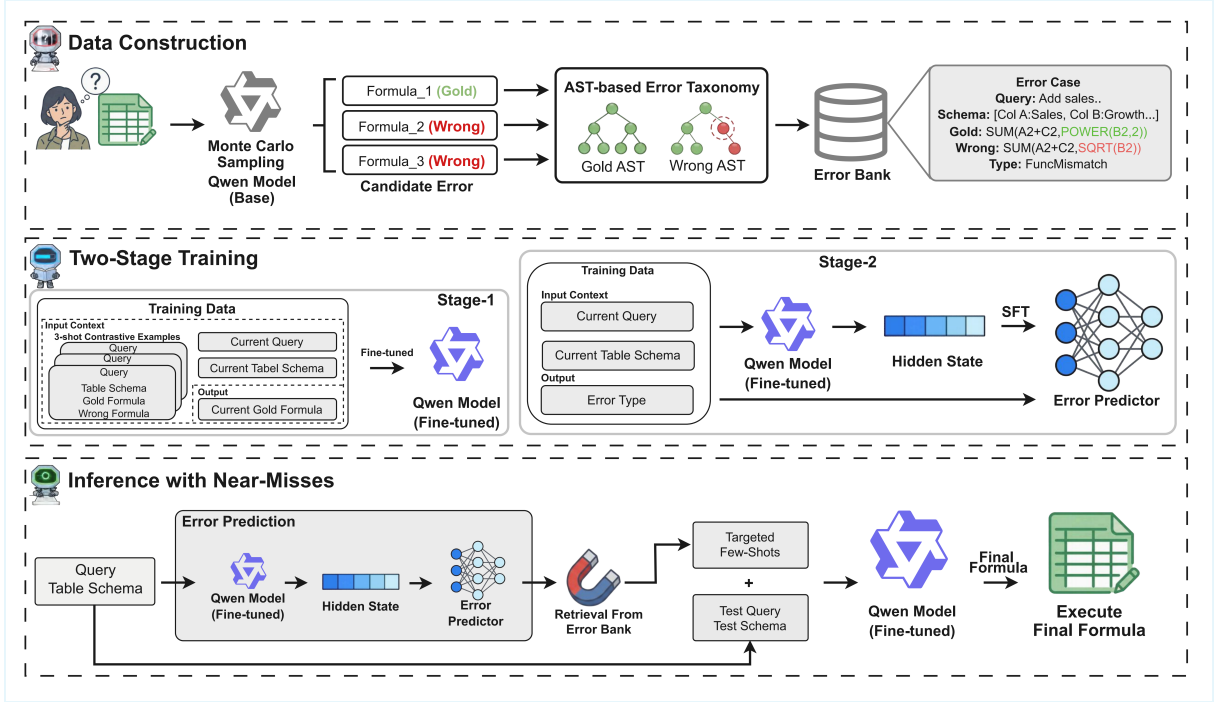


Figure 3: Framework overview. We first build an error bank by Monte Carlo sampling and assign AST-based types into 9 categories (Section 3). The error bank is used to construct error-aware contrastive demonstrations for LoRA fine-tuning, and train an error-type predictor. During inference, we predict error types, retrieve type-aligned demonstrations, and generate the final formula.

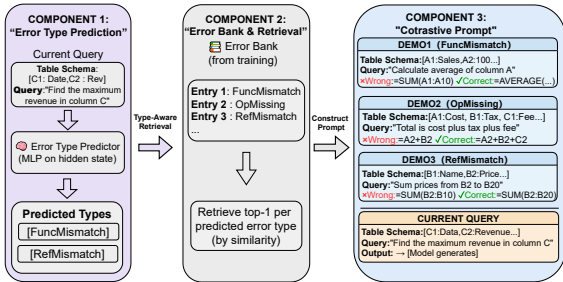


Figure 4: Error-aware contrastive prompt construction. Each one shows a wrong–correct pair under its table schema, followed by the current query.

This *predict–retrieve–generate* procedure adds minimal overhead (one forward pass for hidden states, plus fast retrieval) while providing targeted contrastive evidence aligned with the query’s likely failure modes. The classifier is a small MLP over one hidden state and the retrieval step is a single similarity lookup, so their cost is negligible compared to full autoregressive decoding of a formula.

Algorithm summary. Algorithm 1 summarizes the complete ECFL framework.

5 Experiments

We evaluate ECFL on two benchmarks and compare against different training-based methods, inference-time scaling methods, and API baselines.

5.1 Experimental Setup

5.1.1 Datasets

NL2Formula (public). We use NL2Formula (Zhao et al., 2024) as the in-domain benchmark. Public NL2Formula corpora contain substantial noise, so we systematically audit the data and remove low-quality samples, filtering 7,732 problematic instances in total: (i) *domain contamination*—formulas mixing Excel with Power BI/DAX-specific syntax (e.g., SUMMARIZE, SUMX); (ii) *redundant wrappers*—formulas that wrap dynamic-array expressions with SUM() to force a single scalar value; (iii) *annotation mismatches*—samples where the natural language question does not align with the formula. We split the cleaned data into training (5,412 samples) / validation (1,546 samples) / test (773 samples) sets with a ratio of 7:2:1 in a cross-table manner, ensuring no table appears in more than one split. All baselines are trained and evaluated on this cleaned data, ensuring fair comparison.

Algorithm 1 ECFL

Require: Data \mathcal{D} , base generator \mathcal{G}_0
Ensure: Generator \mathcal{G}_θ , classifier \mathcal{C} , error bank \mathcal{E}

- 1: /* Phase 1: Offline Training */
- 2: $\mathcal{E} \leftarrow \emptyset$
- 3: **for** $(T_i, q_i, f_i^*) \in \mathcal{D}$ **do**
- 4: **for** $k = 1$ to K **do**
- 5: $f_i^{(k)} \sim p_{\mathcal{G}_0}(f | T_i, q_i)$
- 6: **if** $f_i^{(k)} \neq f_i^*$ **and** parseable **then**
- 7: $e \leftarrow \text{ASTTYPE}(f_i^{(k)}, f_i^*)$
- 8: $\mathcal{E} \leftarrow \mathcal{E} \cup \{(T_i, q_i, f_i^{(k)}, f_i^*, e)\}$
- 9: **end if**
- 10: **end for**
- 11: **end for**
- 12: **for** $(T_i, q_i, f_i^*) \in \mathcal{D}$ **do**
- 13: $\mathcal{D}_i \leftarrow \text{ERRORAWARERETRIEVE}(q_i, \mathcal{E}, \text{ETypes}(q_i))$
- 14: Fine-tune \mathcal{G}_θ with contrastive prompt from \mathcal{D}_i
- 15: **end for**
- 16: Train \mathcal{C} on $\{(h_i, e)\}$ pairs from \mathcal{E}
- 17:
- 18: /* Phase 2: Online Inference */
- 19: **Input:** Test instance (T, q)
- 20: $h \leftarrow \mathcal{G}_\theta.\text{FORWARD}(T, q)$
- 21: $\hat{E} \leftarrow \text{TOPK}(\mathcal{C}(h), M)$
- 22: $\mathcal{D}_{\text{demo}} \leftarrow \text{ERRORAWARERETRIEVE}(q, \mathcal{E}, \hat{E})$
- 23: $\hat{f} \leftarrow \mathcal{G}_\theta.\text{GENERATE}(T, q, \mathcal{D}_{\text{demo}})$
- 24: **return** \hat{f}

TFBench (expert-curated). Existing benchmarks over-represent common functions like SUM and IF, which can mask brittleness on long-tail functions and complex compositions. To stress-test compositional generalization on long-tail functions, we construct TFBench from publicly sourced tables and expert-annotated gold formulas. TFBench includes 13,722 expert-annotated instances, covering 215 distinct Excel functions and 12 operators, including rarely used configurations (e.g., advanced XLOOKUP parameters) and newer constructs (LET, LAMBDA). These instances are stratified by nesting depth (1–11) and function rarity, enabling fine-grained analysis of where models fail. Figure 5 illustrates function co-occurrence patterns, showing how functions commonly nest together in complex formulas. Crucially, TFBench is strictly held out from all training, error mining, and retrieval bank construction, preventing any data leakage. Additional statistics and annotation details are provided in Appendix D and Appendix E.

5.1.2 Baselines

We compare ECFL against baselines spanning four categories (details in Appendix B):

Training-based methods. We consider both standard SFT and its variant with preference optimization. SFT fine-tunes the base LLM on cleaned

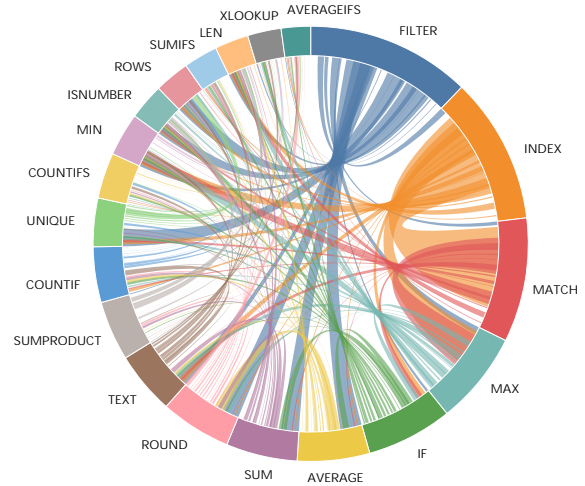


Figure 5: Chord diagram of function co-occurrence in TFBench. The 20 most frequent functions are arranged around the circle, with arc size proportional to frequency. Ribbons indicate co-occurrence within the same formula.

NL2Formula with standard instruction format. SFT + DPO (Rafailov et al., 2023) further applies preference optimization using wrong–correct pairs as preference data.

Inference-time scaling methods. These methods build on the SFT model and improve accuracy by inference-time scaling. Self-Consistency (Wang et al., 2023) (SC@ k) samples k outputs and selects via majority vote. Self-Debugging (Chen et al., 2023) iteratively refines outputs based on execution feedback.

In-Context Learning (ICL). ICL (positive-only) retrieves the most similar query–gold exemplars from the training set via embedding similarity, uses them as demonstrations for both training and inference, but without error typing or wrong–correct contrastive pairs. This tests whether retrieval alone suffices, or whether contrastive pairs are necessary.

API LLMs. We evaluate DeepSeek-R1, Gemini-2.5-Pro, Claude-3.7-Sonnet, and GPT-4o under zero-shot settings.

5.1.3 Metrics

We report two complementary metrics. **Exact Match (EM@1)** normalizes formulas by removing redundant whitespace, canonicalizing function-name case, and standardizing separators, then computes an exact string match against the gold formula. This metric is fully reproducible and aligns with our training-time verification. **Execution Ac-**

Table 2: Key hyperparameters of our method.

Component	Setting
MC samples	$K=3$, $\tau=0.7$, $\text{top-}p=0.9$
Demos	$N=3$
Predicted types	$M=3$ (top-3 from MLP)
Encoder	MiniLM (cosine sim.)
LoRA rank/ α /dropout	32/64/0.05
LR/Batch/Epoch	$2e-5/4/3$
Max length	2500
MLP hidden/layer	512/4

423 **curacy (Exec@1)**, when a spreadsheet execution
 424 engine and input tables are available, counts a pre-
 425 diction as correct if it produces the same output as
 426 the gold formula on the provided table.

427 5.1.4 Implementation

428 We use Qwen3-8B as the base LLM (\mathcal{G}_0) and adapt
 429 it with LoRA to obtain \mathcal{G}_θ (Section 4). For Monte
 430 Carlo error mining, we sample $K=3$ candidates per
 431 query with $\tau=0.7$ and $\text{top-}p=0.95$. For retrieval,
 432 we embed queries using MiniLM and compute co-
 433 sine similarity. Error-aware retrieval fetches the
 434 most similar demonstration per selected error type.
 435 Key hyperparameters are summarized in Table 2.

436 For API LLMs, we use identical input format
 437 (table schema + query) and enforce a strict out-
 438 put format (“Return *only* the Excel formula”). We
 439 report results under (i) zero-shot prompting and
 440 (ii) retrieval-augmented ICL using demonstrations
 441 from the NL2Formula training set. Unless other-
 442 wise stated, temperature is set to 0 (or the closest
 443 deterministic setting) to reduce variance, with out-
 444 put length capped to avoid verbose explanations.

445 5.2 Results

446 Table 3 reports results on NL2Formula and TF-
 447 Bench. On NL2Formula, ECFL performs the best
 448 and achieves the highest EM (75.32%) and Exec
 449 (86.24%), improving over SFT by +6.4 EM and
 450 +6.1 Exec. SFT+DPO improves EM over SFT
 451 but reduces Exec (80.18% \rightarrow 79.01%), suggesting
 452 overfitting to surface formatting. In comparison,
 453 our error-aware contrastive supervision provides a
 454 more reliable learning signal. On TFBench, ECFL
 455 improves over SFT from 38.46% to 42.05% EM
 456 (+3.6) and from 58.91% to 61.25% Exec (+2.3).
 457 The gain is larger on rare functions (+16.28 EM
 458). This suggests that ECFL is more effective on
 459 long-tail functions. We provide a more fine-grained
 460 breakdown by error type in Section 5.3.

461 Compared to inference-time scaling, ECFL out-

Table 3: Main results on NL2Formula and TFBench. EM = Exact Match, Exec = Execution Accuracy.

Method	NL2Formula		TFBench	
	EM	Exec	EM	Exec
<i>Training Methods</i>				
SFT	0.6894	0.8018	0.3846	0.5891
SFT + DPO	0.7006	0.7901	0.3876	0.5626
ECFL (Ours)	0.7532	0.8624	0.4205	0.6125
<i>Inference-Time Enhancement</i>				
SFT + SC@3	0.7036	0.8162	0.3897	0.5917
SFT + SC@5	0.7169	0.8319	0.4073	0.6095
SFT + Self-Debugging	0.7017	0.8007	0.4020	0.5798
<i>ICL Baselines</i>				
ICL (positive-only)	0.7254	0.8385	0.3948	0.5611
<i>API LLMs</i>				
DeepSeek-R1	0.3318	0.6103	0.2292	0.5456
Gemini-2.5-Pro	0.3497	0.6819	0.2841	0.6025
Claude-3.7-Sonnet	0.3439	0.6122	0.2440	0.5816
GPT-4o	0.3402	0.6612	0.2771	0.5924

performs SC@5 on both NL2Formula (75.32% vs. 462
 71.69% EM) and TFBench (42.05% vs. 40.73% 463
 EM), while avoiding multi-sample decoding. Self- 464
 Debugging underperforms SC@5 and ECFL on 465
 EM in both benchmarks, because formula errors 466
 are often silent, limiting the utility of iterative re- 467
 finement. 468

ICL (positive-only) is a strong baseline (72.54% 469
 EM and 83.85% Exec on NL2Formula), showing 470
 that retrieval alone provides substantial gains. How- 471
 ever, ECFL remains higher by +2.8 EM and +2.4 472
 Exec on NL2Formula, indicating that typed wrong- 473
 correct demonstrations add value beyond query- 474
 gold demonstrations alone. 475

API LLMs exhibit a large EM–Exec gap (e.g., 476
 GPT-4o: 34.02% EM vs. 66.12% Exec), suggesting 477
 frequent non-canonical yet executable outputs. Yet, 478
 ECFL still outperforms API LLMs, and the smaller 479
 gap between EM and Exec suggests that contrastive 480
 fine-tuning teaches canonical solutions matching 481
 both functionally and formally. 482

483 5.3 Ablation

Training vs. Inference Contribution. Table 4 484
 ablates whether contrastive demonstrations are 485
 used during training (**Train**) and inference (**Infer**). 486
 Row 1 uses no contrastive demonstrations in ei- 487
 ther stage (SFT). Row 2 uses contrastive demon- 488
 strations only at inference time. Row 3 uses con- 489
 trastive demonstrations only during training. Row 4 490
 uses contrastive demonstrations in both stages (full 491
 ECFL). Using both yields the best results (75.32% 492
 EM on NL2Formula), while using only training 493

Table 4: Ablation of training- and inference-time contributions. The Train/Infer columns indicate whether contrastive demonstrations are used, respectively.

Train	Infer	NL2Formula		TFBench	
		EM	Exec	EM	Exec
✗	✗	0.6894	0.8018	0.3846	0.5891
✗	✓	0.6271	0.7829	0.3615	0.4916
✓	✗	0.6852	0.7998	0.3940	0.5467
✓	✓	0.7532	0.8624	0.4205	0.6125

Table 5: Ablation of the usage of error types for contrastive demonstration retrieval.

Retrieval	NL2Formula		TFBench	
	EM	Exec	EM	Exec
Sim-only (no types)	0.6862	0.7982	0.3914	0.5470
Fixed-types (top-3 types)	0.7092	0.8126	0.4028	0.5567
Predicted-types (ECFL)	0.7532	0.8624	0.4205	0.6125

(68.52%) or only inference (62.71%) is substantially worse. Notably, inference-only contrastive demonstrations underperform the SFT baseline, indicating that the model must learn how to use contrastive evidence during training.

Value of Error-Aware Retrieval. Table 5 compares three contrastive demonstrations retrieval strategies. **Sim-only** ignores error types and retrieves demonstrations purely by query similarity. **Fixed-types** always retrieves demonstrations from the three most frequent error types. **Predicted-types** (ECFL) predicts likely error types for each query and retrieves one contrastive demonstration per predicted type. The results show that both fixed-types and predicted-types outperform sim-only, indicating that sampling demonstrations by error type is beneficial. Predicted-types further achieves the best performance, suggesting that selecting error types conditioned on the query is more effective.

6 Limitations

We acknowledge several limitations and outline directions for future work.

Error bank specificity. Our error bank is constructed offline from Qwen3-8B under specific sampling conditions. The optimal bank may differ across model families and deployment settings. In the future, we will explore adaptive bank construction that updates using target-model feedback or online error accumulation.

Predictor accuracy. Our ablation shows that retrieving demonstrations based on query-specific predicted error types improves performance. While our error-type predictor achieves 78% macro-F1, misclassifications still occur. These errors can route retrieval toward less relevant demonstrations. It is beneficial to incorporate uncertainty-aware prediction that abstains or retrieves broader demonstrations when confidence is low, or more generally improve the predictor to better support retrieval.

Taxonomy assumption. Our taxonomy is defined based on the error patterns observed in our error analysis. Due to resource constraints, we did not conduct a dedicated ablation to systematically compare alternative taxonomies or identify an optimal one. As a result, the current taxonomy may miss finer-grained failure modes. Moreover, we currently consider only the first error, which may underrepresent multi-step failures within a single formula. We will explore taxonomy design choices (e.g., granularity, hierarchy, and multi-label typing) and evaluate their impact on demonstration retrieval and formula generation.

7 Conclusion

We propose ECFL (Error-Aware Contrastive Few-Shot Learning), an error-aware training-inference framework that explicitly learns from model mistakes. By mining errors via Monte Carlo sampling, typing them with AST rules, and retrieving error-aware contrastive demonstrations, ECFL improves structural correctness on both the cleaned NL2Formula benchmark and our expert-curated TFBench, exceeding SC@5 accuracy at substantially lower inference cost. These results suggest that *targeted error correction* through typed retrieval offers a compelling alternative to blind memorization or costly multi-sample decoding. We believe this error-aware paradigm generalizes beyond spreadsheets: any structure-sensitive code generation task—where near-misses are systematic and AST-comparable—could benefit from similar typed supervision. We hope ECFL and TFBench serve as useful resources for the community.

References

Efthimia Aivaloglou, David Hoepelman, and Felienne Hermans. 2017. [Parsing excel formulas: A grammar and its application on 4 large datasets](#). *Journal of Software: Evolution and Process*, 29(12):e1895.

571	Mark Chen, Jerry Tworek, Heewoo Jun, and 1 others.	Sruti Srinivasa Ragavan, Zhitao Hou, Yun Wang, An-	623
572	2021a. Evaluating large language models trained on	drew D. Gordon, Haidong Zhang, and Dongmei	624
573	code. <i>arXiv preprint arXiv:2107.03374</i> .	Zhang. 2022. GridBook: Natural language formulas	625
574	Sibe Chen, Yeye He, Weiwei Cui, Ju Fan, Song Ge,	for the spreadsheet grid . In <i>Proceedings of the 27th</i>	626
575	Haidong Zhang, Dongmei Zhang, and Surajit Chaud-	<i>International Conference on Intelligent User Inter-</i>	627
576	huri. 2024. Auto-Formula: Recommend formulas	<i>faces (IUI '22)</i> , pages 345–368, New York, NY, USA.	628
577	in spreadsheets using contrastive learning for table	Association for Computing Machinery.	629
578	representations . <i>Proceedings of the ACM on Man-</i>	Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan	630
579	<i>agement of Data</i> , 2(3).	Yang, and Ming Zhou. 2020. Minilm: Deep self-	631
580	Xinyun Chen, Maxwell Lin, Nathanael Schärli, and	attention distillation for task-agnostic compression	632
581	Denny Zhou. 2023. Teaching large language models	of pre-trained transformers. <i>Advances in neural in-</i>	633
582	to self-debug. In <i>arXiv preprint arXiv:2304.05128</i> .	<i>formation processing systems</i> , 33:5776–5788.	634
583	Xinyun Chen, Petros Maniatis, Rishabh Singh, Charles	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le,	635
584	Sutton, Hanjun Dai, Max Lin, and Denny Zhou.	Ed H Chi, Sharan Narang, Aakanksha Chowdhery,	636
585	2021b. SpreadsheetCoder: Formula prediction from	and Denny Zhou. 2023. Self-consistency improves	637
586	semi-structured context. In <i>International Conference</i>	chain of thought reasoning in language models. <i>arXiv</i>	638
587	<i>on Machine Learning</i> , pages 1661–1672. PMLR.	<i>preprint arXiv:2203.11171</i> .	639
588	Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021.	Yue Wang, Weishi Wang, Shafiq Joty, and Steven CH	640
589	SimCSE: Simple contrastive learning of sentence em-	Hoi. 2021. CodeT5: Identifier-aware unified pre-	641
590	beddings. In <i>Proceedings of the 2021 Conference on</i>	trained encoder-decoder models for code understand-	642
591	<i>Empirical Methods in Natural Language Processing</i> ,	ing and generation. In <i>Proceedings of the 2021 Con-</i>	643
592	pages 6894–6910.	<i>ference on Empirical Methods in Natural Language</i>	644
593	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan	<i>Processing</i> , pages 8696–8708.	645
594	Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran,	646
595	Weizhu Chen. 2022. LoRA: Low-rank adaptation of	Tom Griffiths, Yuan Cao, and Karthik Narasimhan.	647
596	large language models. In <i>International Conference</i>	2023. Tree of thoughts: Deliberate problem solving	648
597	<i>on Learning Representations</i> .	with large language models. <i>Advances in neural</i>	649
598	Harshit Joshi, Chunning Zhao, Bing Li, Yu Liu,	<i>information processing systems</i> , 36:11809–11822.	650
599	and 1 others. 2023. FLAME: A small language	Kechi Zhang, Zhuo Li, Jia Li, and 1 others. 2023. Self-	651
600	model for spreadsheet formulas. <i>arXiv preprint</i>	Edit: Fault-aware code editor for code generation. In	652
601	<i>arXiv:2301.13779</i> .	<i>Proceedings of the 61st Annual Meeting of the Associ-</i>	653
602	Ansong Ni, Srinu Iyer, Dragomir Radev, and 1 others.	<i>ation for Computational Linguistics</i> , pages 769–787.	654
603	2023. LEVER: Learning to verify language-to-code	Wei Zhao, Zhitao Hou, Siyuan Wu, Yan Gao, Haoyu	655
604	generation with execution. In <i>International Confer-</i>	Dong, Yao Wan, Hongyu Zhang, Yulei Sui, and	656
605	<i>ence on Machine Learning</i> , pages 26106–26128.	Haidong Zhang. 2024. NL2Formula: Generating	657
606	Stephen G. Powell, Kenneth R. Baker, and Barry	spreadsheet formulas from natural language queries .	658
607	Lawson. 2008. A critical review of the literature	In <i>Findings of the Association for Computational Lin-</i>	659
608	on spreadsheet errors . <i>Decision Support Systems</i> ,	<i>guistics: EACL 2024</i> , pages 2377–2388, St. Julian’s,	660
609	46(1):128–138.	Malta. Association for Computational Linguistics.	661
610	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christo-	A Implementation Details	662
611	pher D. Manning, Stefano Ermon, and Chelsea Finn.	Monte Carlo Sampling. We use temperature	663
612	2023. Direct preference optimization: Your language	$\tau = 0.7$, top- $p = 0.9$, and maximum length 512	664
613	model is secretly a reward model . In <i>Advances in</i>	tokens, with each query sampled $K = 3$ times.	665
614	<i>Neural Information Processing Systems</i> .	Error Bank Statistics. The final error banks con-	666
615	Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, and 1	tain 124,449 and 28,645 unique error instances for	667
616	others. 2023. Code Llama: Open foundation models	NL2Formula and TFBench, respectively.	668
617	for code. In <i>arXiv preprint arXiv:2308.12950</i> .	Error-Type Predictor. The MLP classifier has	669
618	Freda Shi, Daniel Fried, Marjan Ghazvininejad, Luke	4096-dim input \rightarrow 1024-dim hidden (with residual	670
619	Zettlemoyer, and Sida I Wang. 2022. Natural lan-	connection) \rightarrow 512-dim hidden \rightarrow 8-dim output,	671
620	guage to code translation with execution. In <i>Proceed-</i>	trained with the AdamW optimizer (lr=1e-3, batch	672
621	<i>ings of the 2022 Conference on Empirical Methods</i>	size 12, 100 epochs).	673
622	<i>in Natural Language Processing</i> , pages 3533–3546.		

674	Software Versions. Python 3.10, PyTorch 2.9.1,	• SFT + DPO: After SFT, we apply Direct Prefer-	722
675	Transformers 4.57.1, PEFT 0.18.0, Sentence-	ence Optimization (Rafailov et al., 2023) using	723
676	Transformers 5.2.0.	wrong–correct pairs as preference data. This tests	724
677	Reproducibility. All results are from single runs	whether preference learning can implicitly lever-	725
678	with a fixed random seed (42). We verified stability	age error signals without explicit error typing.	726
679	on the SFT baseline with 3 independent runs, ob-		
680	serving standard deviation <0.3% on both EM and	Inference-Time Scaling Methods.	727
681	Exec.	• SFT + SC@ k : Self-Consistency (Wang et al.,	728
682	Computational Resources. We use Qwen3-8B	2023) samples k outputs from the SFT model	729
683	(8B parameters) as the base model. All experi-	and selects the most frequent one via majority	730
684	ments were conducted on 4 NVIDIA A800 80GB	vote. We report $k \in \{3, 5\}$.	731
685	GPUs. LoRA fine-tuning takes approximately 3	• SFT + Self-Debugging: Following Chen et al.	732
686	hours; Monte Carlo error mining takes approxi-	(2023), the model iteratively refines its output	733
687	mately 2 hours; total training time is under 6 GPU-	based on execution feedback or syntax errors.	734
688	hours per run.	We allow up to 3 refinement rounds.	735
689	Licenses. We use the following open-source		
690	resources: Qwen3-8B (Apache 2.0), MiniLM	ICL Baselines.	736
691	(MIT), NL2Formula dataset (CC-BY 4.0). TF-	• ICL (positive-only): Retrieves top- N most sim-	737
692	Bench is constructed from public spreadsheet	ilar query-gold exemplars from the training set	738
693	data sourced from Kaggle (https://www.kaggle.com/datasets/) and Alibaba Tianchi (https://	based on query embedding similarity. This rep-	739
694	tianchi.aliyun.com/dataset), which are re-	resents the standard ICL approach without error	740
695	leased under ODbL v1.0, DbCL v1.0, and CC0 1.0	awareness.	741
696	licenses. TFBench will be released under CC-BY	API LLMs. We evaluate four strong general-	742
697	4.0. Our code will be released under MIT license	purpose models: DeepSeek-R1, Gemini-2.5-Pro,	743
698	upon acceptance.	Claude-3.7-Sonnet, and GPT-4o. Each is tested	744
699		under (i) zero-shot (schema + query only). We	745
700	Use of AI Assistants. We used AI writing assis-	use identical input serialization and enforce a strict	746
701	tants (e.g., ChatGPT, Claude) for grammar check-	output format.	747
702	ing and formatting corrections during manuscript		
703	preparation. All scientific claims, experimental de-	C Extended Related Work	748
704	sign, methodology, and core contributions are the	This appendix provides an extended discussion of	749
705	authors’ original work.	related work, elaborating on the connections and	750
706	Ethical Considerations. Our work focuses on	distinctions summarized in Section 2.	751
707	translating natural language to spreadsheet formu-		
708	las, a low-risk application domain. The datasets	C.1 NL2Formula and Code Generation	752
709	used (NL2Formula and TFBench) contain only for-	The task of generating code from natural lan-	753
710	mula expressions and table schemas, with no per-	guage has seen remarkable progress with the ad-	754
711	sonally identifiable information or offensive con-	vent of large language models. Codex (Chen et al.,	755
712	tent. We do not foresee direct risks from misuse, as	2021a) demonstrated that LLMs pretrained on code	756
713	the generated formulas operate within spreadsheet	can achieve strong performance on programming	757
714	environments with limited external impact.	benchmarks. Subsequent work has explored spe-	758
715	B Baseline Details	cialized architectures such as CodeT5 (Wang et al.,	759
716	We describe each baseline category in detail.	2021) for code understanding and Code Llama	760
717	Training-based Methods.	(Roziere et al., 2023) for code-specific tasks.	761
718	• SFT (Supervised Fine-Tuning): Fine-tunes \mathcal{G}_0 on	NL2Formula represents a specialized instance	762
719	cleaned NL2Formula data with standard instruc-	of code generation where the target language	763
720	tion format. No contrastive demonstrations are	is Excel’s formula syntax. SpreadsheetCoder	764
721	used during training or inference.	(Chen et al., 2021b) introduced neural encoder-	765
		decoder models that incorporate spreadsheet con-	766
		text (row/column headers, cell values) as struc-	767
		tured input. FLAME (Joshi et al., 2023) explored	768

769 smaller, task-specialized models for formula gen-
770 eration. Despite progress, these systems struggle
771 with the strict structural constraints of Excel for-
772 mulas, where minor errors in function names, ar-
773 gument order, or cell references can completely
774 invalidate the output. Crucially, existing methods
775 rely on end-to-end memorization and lack explicit
776 mechanisms to handle the systematic *near-miss* er-
777 rors that arise from long-tail function distributions
778 and complex nesting patterns.

779 C.2 Learning from Model Errors

780 A growing body of work explores how models can
781 learn from their own mistakes, which is central to
782 our approach.

783 **Self-Debugging and Refinement.** Self-
784 debugging approaches use execution feedback to
785 iteratively refine outputs. [Chen et al. \(2023\)](#) teach
786 models to self-debug by generating and executing
787 test cases. Self-Edit ([Zhang et al., 2023](#)) trains fault-
788 aware editors to localize and fix errors. LEVER ([Ni
789 et al., 2023](#)) learns verifiers that predict execution
790 correctness to filter candidates. These methods typ-
791 ically operate at inference time and require multiple
792 generation-execution cycles. Our approach differs
793 by mining errors *offline* during training, building a
794 reusable error bank that enables single-pass gener-
795 ation at inference time.

796 **Preference Optimization and Error Signals.**
797 Direct Preference Optimization (DPO) ([Rafailov
798 et al., 2023](#)) and its variants have emerged as power-
799 ful alternatives to RLHF for aligning language mod-
800 els with human or model-generated preferences. In
801 the context of code generation, collected errors
802 can serve as the “rejected” samples to guide the
803 model away from failure modes. However, stan-
804 dard DPO treats each error as a monolithic fail-
805 ure, potentially missing the fine-grained structural
806 reasons behind a near-miss. Our framework com-
807 plements preference-based methods by explicitly
808 typing errors via AST analysis, enabling more gran-
809 ular contrastive supervision during both fine-tuning
810 and retrieval.

811 **Error Typing and Categorization.** Prior work
812 on error analysis in code generation often treats
813 errors as monolithic failures. We introduce a fine-
814 grained AST-based error taxonomy that categorizes
815 near-misses into 9 types (function/operator mis-
816 match, missing, redundant, fill-level errors, and
817 miscellaneous). This typed supervision enables
818 targeted retrieval of contrastive demonstrations
819 aligned with predicted error patterns.

C.3 Inference-Time Robustness 820

821 For code generation, inference-time strategies can
822 substantially improve robustness but often incur
823 high computational cost.

824 **Self-Consistency.** Self-consistency ([Wang et al.,
825 2023](#)) samples multiple outputs and selects by ma-
826 jority vote or execution agreement. While effective,
827 sampling K candidates requires K forward passes,
828 leading to multiplicative inference cost. Our ap-
829 proach achieves comparable robustness without
830 multi-sample generation by predicting likely er-
831 ror types and retrieving targeted demonstrations
832 with a single extra forward-only pass.

833 **Execution-Guided Methods.** Execution-guided
834 methods ([Shi et al., 2022](#)) use program execution
835 to filter or rank candidates. These require access to
836 an execution environment at inference time, which
837 may not always be available. Our error-type pre-
838 dictor anticipates likely failure modes *before* gener-
839 ation, enabling proactive rather than reactive error
840 handling.

C.4 Contrastive Learning for Generation 841

842 Contrastive learning has proven effective across
843 NLP tasks by learning representations that distin-
844 guish positive pairs from negative pairs ([Gao et al.,
845 2021](#)). In the context of generation, contrastive
846 signals can sharpen decision boundaries between
847 correct and incorrect outputs.

848 We adopt a simple but effective form of con-
849 trast: presenting wrong–correct pairs as in-context
850 demonstrations. Unlike representation-level con-
851 trastive learning, our approach operates at the out-
852 put level, directly showing the model examples of
853 common mistakes and their corrections. This in-
854 terpretable contrast is naturally aligned with the
855 structure of formula generation errors, where near-
856 misses differ from gold outputs by specific, identi-
857 fiable deviations.

C.5 Parameter-Efficient Fine-Tuning 858

859 Adapting large language models to specialized
860 tasks traditionally required full fine-tuning, which
861 is computationally expensive and risks catastrophic
862 forgetting. LoRA ([Hu et al., 2022](#)) introduces low-
863 rank adapter matrices that can be efficiently trained
864 and merged with pretrained weights.

865 In our framework, LoRA serves as the adap-
866 tation mechanism, but the key innovation lies
867 in *what* we fine-tune on: contrastive few-shot
868 prompts constructed from the error bank. By ex-

posing the model to wrong–correct pairs during fine-tuning, we inject error awareness directly into the adapted parameters, complementing inference-time retrieval.

C.6 Data Quality and Benchmark Design

The quality of training data significantly impacts code generation performance. For NL2Formula specifically, public datasets often contain non-Excel functions (from Power BI/DAX), redundant wrapper patterns, and inconsistent annotations. Our data governance pipeline addresses these issues through systematic cleaning: removing domain-mismatched samples, normalizing redundant structures via AST analysis, and canonicalizing references.

Existing NL2Formula benchmarks tend to over-emphasize common functions like SUM and IF, making it difficult to diagnose failures on rare functions and complex compositions. TFBench addresses this gap by deliberately covering 215 distinct functions and 12 operators, including long-tail functions and advanced constructs (XLOOKUP, LET, LAMBDA), with controlled difficulty stratification by nesting depth. This enables more nuanced evaluation of compositional generalization.

D TFBench Statistics

We construct TFBench by curating public spreadsheet tables sourced from Kaggle and Alibaba Tianchi. The experts are asked to write the corresponding Excel formula using standard Excel syntax. TFBench contains 13722 expert-annotated instances and is stratified by nesting depth and function rarity for fine-grained analysis. The annotation protocol and quality control procedures are described in Appendix E.

Statistic	Value
Total examples	13722
Unique functions	215
Unique operators	12
Avg. nesting depth	2.35
Max nesting depth	11
Avg. formula length	23.82 tokens

Table 6: TFBench statistics.

E TFBench Annotation Protocol

Annotator Instructions. Annotators were instructed to: (1) read the natural language query describing a spreadsheet operation; (2) write the

corresponding Excel formula using standard Excel syntax; (3) verify the formula by mental execution against the provided table schema. We did not provide step-by-step templates to avoid biasing annotators toward specific function choices.

Quality Control. Each formula was verified by a second annotator. Disagreements were resolved by a third senior annotator.

Annotator Recruitment and Compensation. Annotators were graduate students in mathematics and computer science with proficiency in spreadsheet applications, recruited from the research team. They were compensated at standard research assistant rates. All participation was voluntary, and annotators consented to their annotations being used for research purposes.

Annotator Demographics. Annotation was performed by 10 graduate students (6 male, 4 female) with 2+ years of Excel experience, located in China. No personally identifiable information about annotators is disclosed beyond aggregate demographics.

F Fine-Grained Analysis

F.1 Long-Tail Function Performance

To understand where ECFL gains come from, we stratify TFBench results by function frequency in the training set. Table 7 reports performance on queries containing low-frequency functions.

Thresh.	#Func	#Samp.	SFT	ECFL	Δ
Bot. 0.1%	39	86	20.93	37.21	+16.28
Bot. 0.05%	29	49	12.24	26.53	+14.29

Table 7: Performance on low-frequency functions (%). ECFL shows larger gains on rarer functions.

The results confirm that ECFL is particularly effective on long-tail functions, where error-aware contrastive demonstrations provide targeted guidance that standard SFT lacks.

F.2 Case Study

Table 8 shows examples where SFT produces incorrect formulas but ECFL succeeds.

Case 1: Calculate the Kurtosis of the trading Volume distribution for the year 2015.

Gold: =KURT(FILTER(G2:G15097, YEAR(A2:A15097)=2015))

SFT: =STDEV(...) X

ECFL: =KURT(...) ✓

Error: FuncMismatch

Case 2: Forecast the user count when the temperature is 0.5.

Gold: =FORECAST(0.5, P2:P732, J2:J732)

SFT: =FORECAST.ETS(...) X

ECFL: =FORECAST(...) ✓

Error: FuncMismatch

Table 8: Case study: SFT errors corrected by ECFL.