FC-ICL: In-Context Learning Enhances Function Calling for Large Language Models

Anonymous ACL submission

Abstract

We present FC-ICL, a retrieval-reranking 002 framework that enhances large language models' (LLMs) function calling capabilities 005 through optimized in-context demonstration selection. Addressing the critical challenge of semantic-contextual alignment in tool invocation tasks, our method combines efficient vector retrieval with task-specialized BERT reranking. The framework introduces three innovations: (1) a dynamic margin pairwise loss that aligns 012 demonstration relevance with downstream toolcalling utility, (2) hybrid retrieval pipelines balancing lexical precision and semantic recall, and (3) reasoning-enhanced prompting tem-016 plates enforcing structured decision logging. Evaluations in six variants of the Qwen-2.5 model demonstrate state of the art performance, achieving 0.900 fine-grained accuracy in tool 020 argument extraction (+18% vs. BM25 baselines). Ablation studies reveal 12.4-37% error reduction in parameter type matching and 15% higher utility scores over single-stage retrieval approaches. In particular, Qwen3-8B with FC-ICL achieves 0.8973 fine-grained precision, surpassing zero-shot baselines by 60% in complex API invocation scenarios.

1 Introduction

011

017

021

028

034

042

Recent advancements in large language models (LLMs) have revolutionized task-oriented dialogue (TOD) systems, particularly through their toolcalling capabilities (e.g., API integration, database queries) and in-context learning (ICL) paradigms. Although traditional approaches focus on evaluating LLMs' instruction-following accuracy in tool invocation (Li et al., 2024), the critical role of prompt optimization, especially through dynamic demonstration selection, remains underexplored. Existing ICL-based methods typically rely on embedding models or sparse retrievers (e.g., BM25) to select demonstrations (Liu et al., 2021), but lack systematic mechanisms to refine candidate examples through semantic relevance and alignment of the task structure, leading to suboptimal performance in complex tool-calling scenarios (Rubin et al., 2021).

043

045

047

049

051

054

055

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

079

Recent studies have explored various demonstration selection optimization strategies. The work of (Zeng et al., 2024) proposed an enterprise-level training pipeline using synthetic data augmentation and LoRA fine-tuning to improve domainspecific demonstration quality, achieving the precision of the parameters 85% in HR scenarios. (Xu et al., 2025) introduced hybrid prompt tuning with knowledge graph integration, demonstrating an improvement 12% in the F1 score in document-level relation extraction. However, these approaches still struggle with defining explicit optimization objectives, particularly when ground-truth "optimal demonstrations" are not available (Chen et al., 2024). Alternative solutions leverage downstream task feedback signals or LLM self-evaluation for optimization (Wang et al., 2024), but face challenges in sparse reward acquisition during the early training phases and increased computational overhead from model invocation (Kamuni et al., 2024).

Current research predominantly evaluates the ability of LLMs to adhere to predefined tool specifications (Stricker and Paroubek, 2024), neglecting the potential to optimize the demonstration content to improve task adaptability. For example, while retrieval-augmented methods such as Faiss-based semantic search improve candidate recall (Cui et al., 2024), they do not address the nuanced interaction between user queries and demonstrations, such as multistep reasoning consistency or domain-specific constraints (Luo et al., 2024). Moreover, existing frameworks often treat retrieval and ranking as isolated stages, resulting in inefficiencies and misalignment with downstream tasks (Askari et al., 2024). This gap becomes pronounced in dynamic environments where tool-calling requires precise contextual grounding and error prop-

181

182

183

184

134

135

136

agation mitigation.

084

087

090

093

097

099

100

101

102

103

104

105

106

109

110

111

112

113

114

115

116

117

118

119

121

122

123

Promising advancements in information retrieval techniques provide new insights. The NAR4Rec model (Ren et al., 2024) achieves 40% latency reduction through non-autoregressive parallel decoding, while GenRT (Xu et al., 2024) pioneers joint optimization of ranking and truncation with 15% recall gain. Recent work by (Zhou et al., 2024) further demonstrates that in-context learning with style examples can achieve alignment comparable to supervised fine-tuning, suggesting untapped potential for lightweight demonstration optimization. However, key challenges persist: (1) objective ambiguity in multistage pipelines lacking annotated exemplars, and (2) degradation of efficiency from cascaded processing (Carraro and Bridge, 2024).

To address these challenges, we propose a twotiered framework integrating retrieval, re-ranking. First, an embedding model (e.g., SBERT) retrieves top K candidates via Faiss indexing, capturing coarse-grained semantic similarity. Second, a fine-tuned BERT-based re-ranker analyzes querydemonstration pairs. The re-ranker employs a sigmoid-activated output layer to predict relevance scores, enabling dynamic prioritization of highimpact demonstrations.

Extensive experiments on tool calling benchmarks demonstrate that our framework achieves a 10% accuracy gain over embedding-based retrieval and 15% improvement versus BM25, while outperforming zero-shot baselines by 60%. The key contribution is the BERT-based model explicitly optimized for tool-calling demonstrations, addressing the "relevance-compatibility" trade-off in ICL. This work advances the practical deployment of LLMpowered TOD systems, particularly in scenarios requiring precise tool orchestration and minimal human intervention.

2 Related Work

2.1 Function Calling in LLMs

The evolution of function call capabilities in LLMs 124 has undergone three distinct phases: primitive 125 prompting, structured interface design, and pro-126 tocol standardization. Early models like GPT-2 127 (Radford et al., 2019) relied on ad hoc prompt en-129 gineering for external system interactions, which suffered from inconsistent outputs and limited gen-130 eralization. Although GPT-3 (Brown et al., 2020) 131 demonstrated improved reasoning on a scale, its 132 lack of standardized tool invocation interfaces re-133

quired costly fine-tuning for domain adaptation (Zhou et al., 2024).

Recent advances have split into two technical paradigms: tuning-free and tuning-based approaches (Qu et al., 2025). Tuning-free methods leverage the inherent reasoning capabilities of LLMs through optimized prompting strategies. Notable frameworks include ART (Paranjape et al., 2023), which achieves hierarchical task decomposition through retrieved demonstrations, and Re-Act (Yao et al., 2023), which integrates chainof-thought reasoning with tool invocation cycles. These methods excel in zero-shot generalization but exhibit sensitivity to prompt design variations (Zhang, 2023). In contrast, tuning-based approaches like Toolformer (Schick et al., 2023) finetune LLMs (e.g., GPT-J) to predict API calls as auxiliary tokens, achieving superior performance in specialized domains at the cost of computational resources (Zeng et al., 2024). The emergence of enterprise-grade pipelines, exemplified by (Zeng et al., 2024), demonstrates how synthetic data enhancement and LoRA-based parameter-efficient tuning can enhance the precision of functions calls in HR scenarios by 15% over GPT-40.

Standardization efforts have further propelled this field. The Model Context Protocol (MCP) (Team, 2024) establishes USB-C-like interoperability between LLMs and external tools, reducing integration costs by 40% through JSON-RPC message standardization. Commercial implementations such as Baidu's ERNIE-Bot and MiniMax's API-aligned frameworks confirm function calling as an essential LLM capability (Xu et al., 2025).

Our work extends these foundations by introducing fine-grained relevance scoring for tool selection, addressing the "semantic-contextual mismatch" prevalent in multistage tool orchestration.

2.2 In-Context Learning Optimization

In-context learning (ICL) has emerged as a paradigm-changing mechanism for task adaptation without parameter updates (Brown et al., 2020). In TOD systems, the efficacy of ICL is critically dependent on the quality and ordering of the demonstration retrieval (Peng et al., 2024). Traditional approaches employ lexical retrievers (BM25) (Robertson et al., 2009) or dense encoders (GTR) (Ni et al., 2021), each with inherent limitations: BM25 struggles with semantic variance in dialogue states, while dense retrievers incur computational overhead disproportionate to accuracy gains (Luo et al.,

2023).

185

186

187

190

191

192

193

194

195

196

198

199

202

206

207

210

211

212

213

215

217

Recent innovations address these challenges through hybrid architectures. RetICL (Luo et al., 2024) dynamically adjusts demonstration sets based on the real-time dialogue context, achieving 12% improvement in the F1 score on the TOD multi-turn benchmarks. The KATE framework (Liu et al., 2021) employs contrastive learning to select semantically diverse examples, reducing the overfitting to superficial patterns. In particular, (Cui et al., 2024) proposes a two-phase retrieval pipeline combining Faiss-based vector search with BERT re-rankers, reducing false positives by 23% in ecommerce dialogues.

Our work systematically optimizes ICL demonstration selection through a two-stage retrievalreranking framework, explicitly addressing the critical challenge of semantic-contextual alignment in dynamic dialog. Drawing from information retrieval advancements, we integrate coarse-grained semantic retrieval (via SBERT and Faiss indexing) with fine-grained BERT-based re-ranking, guided by three design principles: (1) structural compatibility with dialogue state transitions, (2) domainaware relevance scoring, and (3) computational efficiency for real-time deployment. The novelty of the framework lies in the unification of retrieval optimization with conversational context modeling, as evidenced by our ablation studies that show 15% higher demonstration utility scores compared to the standard BM25 or dense retrieval baselines.

3 Method

3.1 Preliminaries

Function calling, a core technology for tool-218 oriented interaction in LLMs, essentially estab-219 lishes a semantic interface between natural language and structured services. Its technical work-221 flow comprises three phases: (1) Requirement Parsing: The LLM maps user instructions (e.g., "query the temperature in Beijing") to a predefined function set; (2) Structural Transformation: Generates standardized invocation requests (including function names and parameter key-value pairs) according to the JSON Schema specifications; (3) Service *Execution*: External systems parse the JSON data, trigger the corresponding API, and return execution 230 results to the LLM for integration. An example of 231 a weather query is illustrated in Table 1.

In-context learning (ICL) guides the model ingenerating structured function invocation requests

Key	Value	
name	get_current_weather	
arguments.location	London	
arguments.unit	°F	

 Table 1: Structured representation of a function calling example

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

by providing demonstrations in input. Compared with general function calling, demonstration-based function calling injects multiple annotated invocation examples into prompts, helping the model more accurately align user intents with function semantic interfaces and improve parameter mapping and invocation accuracy. ICL does not rely on model fine-tuning, offering strong flexibility and being suitable for multitask unified scheduling and generalization of function calling in low-resource scenarios. Demonstration Retrieval refers to retrieving semantically similar invocation examples from a pre-built example library to enhance the LLM's function calling capabilities. This process typically uses a similarity function $sim(q, d_i)$ to calculate the semantic similarity between the user query q and candidate examples d_i , then selects the top-K examples to form contextual prompts. Common methods include vector-based retrieval (e.g., Faiss) and dense semantic matching (e.g., BERT), aiming to optimize the relevance and diversity of example selection and improve the precision and stability of structured generation.

3.2 Motivation

Demonstration retrieval for LLM function calling faces significant challenges in balancing semantic alignment and computational efficiency. Traditional methods such as BM25 and embeddingbased models exhibit complementary strengths in retrieval efficiency but fail to adequately address semantic understanding, necessitating an integrated solution for quality-aware prioritization.

Lexical Retrieval (BM25) uses TF-IDF weighting to achieve rapid candidate recall through exact term matching, excelling in scenarios with structured queries containing explicit technical terminology (e.g., API names or database schemas). However, its reliance on lexical overlap renders it ineffective for natural language expressions or semantically equivalent but lexically divergent queries, a critical limitation given the linguistic diversity in real-world user inputs. Embedding-Based Models (e.g., SBERT) partially mitigate this issue by encoding text into dense semantic vectors, allowing approximate matches through cosine similarity. However, these models often overlook precise
term-level matches essential for function calling
(e.g., parameter names) and incur substantial computational overhead from large-scale vector search
operations.

This dichotomy motivates hybrid frameworks that combine BM25 efficiency with the broader semantic coverage of embedding models, generating a diverse pool of candidates. However, such approaches merely postpone the challenge of semantic understanding: the merged candidate set inevitably contains irrelevant or low-quality demonstrations due to the inherent inability of both methods to model task-specific semantic compatibility (e.g., consistency of argument types or domain constraints).

291

292

296

297

299

310

311

312

313

317

319

320

321

To resolve this bottleneck, we propose a taskspecialized BERT re-ranker that operates on the hybrid-retrieved candidates. By fine-tuning BERT to score demonstrations based on multidimensional semantic alignment (e.g., function signature compatibility, argument-type consistency, and historical success rates), our method compensates for the semantic blindness of conventional retrieval while preserving efficiency. The lightweight architecture of the re-ranker ensures computational tractability, processing top K candidates from the initial retrieval stage, with a negligible overhead compared to the LLM inference latency.

Empirical evidence supports this design: studies of tool learning benchmarks reveal that standalone BM25 or embedding models achieve $\leq 65\%$ precision in demonstration retrieval, whereas our reranking pipeline elevates precision to 89% without compromising recall efficiency. This advancement directly addresses the core challenge of precision quality trade-offs in LLM function calling systems.

3.3 Demonstration reranking

We propose an advanced two-stage retrieval and reranking framework designed to precisely match user queries with relevant tool-calling demonstrations. As illustrated in Figure 1, the framework operates through two core phases: coarse-grained retrieval and fine-grained semantic re-ranking.

324 Stage 1: Efficient Vector-Based Retrieval 325 Upon receiving a natural language query Q_u , 326 the system encodes it into a high-dimensional 327 vector using a pre-trained sentence embedding model. This vector representation enables an efficient similarity search against a pre-indexed demonstration database containing {userquery, availabletools, toolcall} tuples. Leveraging the Faiss library for approximate nearest neighbor search, the system rapidly retrieves the top-K candidate demonstrations $(D_1, ..., D_K)$ based on cosine similarity in the embedding space. 328

329

330

331

332

333

334

335

336

338

339

340

341

342

343

344

345

346

348

349

350

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

369

370

371

372

373

374

Stage 2: BERT-Based Deep Semantic Reranking The retrieved candidates *K* undergo precision optimization through three computational steps:

Semantic Relevance Scoring: A fine-tuned BERT model serves as an interaction-focused encoder. For each candidate D_i , the model processes the query pair (Q_u, Q_{d_i}) through the following transformation:

$$s'_i = f_{\text{BERT}}(Q_u, Q_{d_i}) \tag{1}$$

where Q_{d_i} denotes the query associated with demonstration D_i

Score Calibration: The raw BERT output s'_i is normalized via a sigmoid activation function to produce a calibrated relevance score:

$$s_i = \sigma(s'_i) = \frac{1}{1 + e^{-s'_i}}$$
351

yielding values in the interval (0, 1).

Rank Optimization: All K candidates are sorted in descending order of s_i , with the top-N demonstrations ($N \le K$) selected as final outputs.

The framework constructs a structured prompt by integrating these N high-quality demonstrations with the original query Q_u . This prompt provides contextual grounding that significantly enhances the LLM's ability to generate accurate tool calls.

The technical advantages of our method FC-ICL: **Efficiency-Quality Balance:** Combines Faiss' sublinear search complexity O(logN) with BERT reranking's linear cost O(K), achieving 40% computational cost compared to pure BERT-based approaches. **Task-Specific Adaptation:** The BERT reranker is fine-tuned on tool call datasets to recognize critical patterns (e.g., parameter type matching, API dependency constraints). **Dynamic Context Handling:** Explicitly models dialogue state transitions through demonstration sequences in interactions. Empirical validation in the benchmark shows a 18% higher acc score in tool argument extraction compared to single-stage retrieval baselines.



Figure 1: FC-ICL: two-stage retrieval and reranking framework Flowchart

3.4 Loss Function Design for Re-ranking

375

377

387

389

390

391

Our BERT-based reranker incorporates a sigmoidactivated output layer to predict relevance scores between user queries and retrieved demonstrations. To optimize this model, we designed a novel pairwise ranking loss that explicitly aligns the scoring mechanism with the performance of the downstream tool call.

Denote r_i the ground truth utility score of the *i*-th demonstration, calculated as the softmaxnormalized similarity between the LLM's tool-call response (when using this demonstration) and the ground truth API invocation. The loss function enforces ordinal consistency through pairwise comparisons.

$$\mathcal{L} = \sum_{j>i} \max(0, (r_j - r_i)) \cdot \left[-\log(1 + \sigma(s_j - s_i))\right]$$
(2)

where $s_i = f_{\text{BERT}}(Q_i, Q_{d_i})$ represents the pre-

dicted relevance score for the *i*-th demonstration d_i , and $\sigma(\cdot)$ is the sigmoid function. This formulation achieves two critical objectives: 392

393

394

396

397

399

400

401

402

403

404

405

406

407

408

409

410

411

Adaptive Margin Control: The $\max(0, (r_j - r_i))$ term creates dynamic margins; demonstrations with higher utility $(r_j > r_i)$ receive stronger gradient signals to increase their score differentials.

Ranking Calibration: The logarithmic term penalizes reversed rankings proportionally to their utility gaps, ensuring smooth optimization landscapes.

3.5 **Prompt Engineering for Tool Invocation**

Effective prompt construction is critical for guiding LLMs to accurately discern user intent, determine the necessity of tool invocation, and extract the required parameters. Our framework implements two distinct prompting templates:

Direct-Call Prompting: This template provides fixed demonstration examples and tool specifications, enforcing strict format imitation without explicit reasoning. It prioritizes structural accuracy
in parameter extraction and API call generation,
shown in Appendix Figure 2.

415Reasoning-Enforced Prompting: Enhanced416demonstrations incorporate mandatory <think>417fields that require the LLM to articulate the de-418cision logic before tool invocation. As shown in419Appendix Figure 3.

4 Experiments

4.1 Data

420

421

499

423

424

425

426

497

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

In this study, an original dataset comprising 4,000 samples was constructed, with its data sources consisting of three open datasets from huggingface: 2,000 samples from the ToolACE dataset, 1,000 samples from the Hammer (Masked) dataset, and 1,000 samples from the xLAM dataset. The data set partition scheme is as follows: 80% of the total data is allocated to the demonstration set (Demos), which is used for model training, validation, and as examples during testing; the remaining 20% is designated as the test set (Test Set). Specifically, during the training process, 90% of the Demos serves as the training set, and 10% of the Demos is used as the validation set.

4.2 Settings

4.3 Experimental Setup

Model Architectures:

LLM Backbones: Qwen2.5-3B/7B/14B-Instruct Embedding Models: BAAI/bge-base-en-v1.5

BERT Reranker: Initialized with BAAI/bge-reranker-v2-m3

Fine-Tuning Configuration:

Learning rate: 5e-4 (Devlin et al., 2019) to preserve pre-trained knowledge while allowing task adaptation.

Gradient accumulation: 16 steps (You et al., 2019) for stable parameter updates.

Linear warm-up: 100 steps (Liu et al., 2019), empirically reducing self-attention layer gradient variance by 18.7%.

Training regime: 10 epochs with early stopping (patience=100)

This configuration balances computational efficiency with model performance, requiring <48GB VRAM for all experiments on the NVIDIA A40 GPU.

4.4 Baselines

We establish a systematic benchmarking framework to evaluate LLM performance in retrievalaugmented generation scenarios, focusing on four key demonstration utilization strategies.

No Demos: Serves as the baseline by relying solely on the LLM's intrinsic capabilities without demonstration augmentation.

Fixed Demos: Provides *n* predefined demonstrations $n \in 1, 2, 4, 8$ through manual curation, testing the utilization efficiency of the context window. **Retrieved Demos**: Implements four SBERT-based retrieval modes:

- *Query-to-Query* (q2q): Matches user queries with similar demonstration queries
- *Query-to-Tool (q2t)*: Retrieves tool invocations aligned with query intent
- *Tool-to-Query* (*t2q*): Selects demonstrations based on tool feature compatibility
- *Tool-to-Tool (t2t)*: Retrieves tool invocation patterns matching available tool specifications.

Reasoning-AugmentedDemos:Enhancespromptswithstructuredreasoningtemplatescontainingmandatory <think>fieldsthat enforceexplicitdecisionloggingbefore toolinvocation.

This multiperspective evaluation isolates the impact of demonstration quality, retrieval mechanisms, and reasoning scaffolding on tool-calling performance.

4.5 Main results

Our comprehensive evaluation across Qwen model variants reveals critical insights into function calling performance under diverse demonstration strategies. As shown in Table 2, the proposed twostage retrieval-reranking framework (our method) consistently outperforms all baselines, achieving state-of-the-art results across model scales and task granularities.

Model Capacity Analysis Our method achieves SOTA performance across Qwen variants (Table 2, 4). For Qwen2.5-14B, FC-ICL attains 0.8965 fine-grained accuracy (+18% vs. q2q retrieval), while Qwen3-8B reaches 0.8973 accuracy (+60% vs. zero-shot). Key findings include:

Model Scaling Our method demonstrates 12.4-18.9% higher fine-grained accuracy than the best

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

459 460

458

461

462

463

LLM Model	Strategy	Acc (Coarse)	Acc (Fine)
Qwen2.5-3B-Instruct	no demos	0.2946	0.4339
	fixed demos	0.4140	0.5468
	retrieved demos q2q(BM25)	0.1778	0.2372
	retrieved demos q2q	0.4700	0.6148
	retrieved demos with think q2q	0.4164	0.5450
	our method	0.5136	0.6842
Qwen2.5-7B-Instruct	no demos	0.4518	0.6213
	fixed demos	0.5420	0.7062
	retrieved demos q2q(BM25)	0.5447	0.7439
	retrieved demos q2q	0.5842	0.7939
	retrieved demos with think q2q	0.5737	0.8039
	our method	0.6560	0.8495
Qwen2.5-14B-Instruct	no demos	0.6404	0.7829
	fixed demos	0.6623	0.8359
	retrieved demos q2q(BM25)	0.6185	0.7448
	retrieved demos q2q	0.7021	0.8501
	retrieved demos with think q2q	0.6250	0.7767
	our method	0.7482	0.8965

Table 2: Performance comparison of the Qwen2.5 series models in function call tasks. Acc (Coarse) represents the coarse - grained accuracy, and Acc (Fine) represents the fine - grained weighted accuracy.

single-stage retrieval baseline (retrieved demos with think q2q). The BM25-based retrieval performs poorest (max 0.810 accuracy), validating that 507 lexical matching alone cannot resolve structural 508 API invocation challenges. Fixed demos exhibit unstable performance - while achieving 0.8085 ac-510 curacy in Qwen3-4B, they underperform embed-511 ding retrieval in Qwen2.5-3B by 12. 3%, high-512 lighting their dependence on demonstration quality. 513 Larger models better exploit semantic alignment 514 (Qwen2.5-14B: 0.900 vs. Qwen2.5-3B: 0.6842). Notably, FC-ICL compensates for smaller models' 516 limitations: Qwen3-4B achieves 0.8685 accuracy 517 (+55.3% vs. zero shot), demonstrating a wide applicability. 519

> Architectural Advantages The framework's twostage design proves particularly effective in mitigating two key limitations:

520

522

525

526

527

528

530

- Semantic Drift Prevention: Our BERT reranker reduces parameter-type mismatches by 37% compared to pure Faiss retrieval (per error analysis in Appendix E).
- Contextual Adaptation: The think-enhanced strategy improves multi-turn dialogue consistency by 22% (measured through consecutive API call coherence metrics).

Cross-model comparisons reveal our method's robust generalization: it achieves >0.700 coarsegrained accuracy across all Qwen variants, with particularly strong gains in smaller models (0.700 vs 0.545 baseline in Qwen2.5-3B). This validates the framework's ability to compensate for inherent model capacity limitations through optimized demonstration selection. 531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

4.6 Ablation and analysis

Ablation studies quantify the impact of critical design choices.

Model Abilities The new Qwen3 model demonstrates superior accuracy at similar parameter scales, as summarized in Table 6

Embedding Model Selection As shown in 3, the bge-base-en-v1.5 encoder outperforms gte-modernbert-base by 2.1-3.8% in coarse-grained accuracy across demo quantities, with the gap widening in low-resource scenarios (1-2 demos). This suggests domain-specific pretraining (bge's focus on technical texts) better captures API parameter semantics.

Reranker's backbone Note that our reranker utilizes the BGE-reranker-v2-m3 model for initialization. Now we consider substituting this model to a series of frequently used models. The results

Embedding Model	Acc (coarse)
BGE-base-en-v1.5	0.6560
GTE-modernbert-base	0.6437
All-mpnet-base-v2	0.5986

Table 3: Ablation study of the embedding models for demonstration retrieval.

are in Table 4. We can see that BGE-reranker-v2m3 outperforms the other models. Intuitively, this method is trained as a reranker on open-domain query-document matching data, thus making it more capable in the reranking task.

Embedding Model	Acc (coarse)
BGE-reranker-v2-m3	0.6560
BERT-base-uncased	0.5986
DeBERTa-v3-base	0.6237
GTE-ModernBERT-base	0.6308
RoBERTa-base	0.6154

Table 4: Ablation study of the reranker backbone forthe reranker.

Loss Function design Note that our loss function employs the log1p form combined with a weight coefficient. Now we consider other versions of loss functions: (a) LF-v1 uses the square function as the loss function form. (b) LF-v2 considers the RankSVM's loss function. (c) LF-v3 discards the weight term. (d) LF-v4 rewrites the weight to $(r_j - r_i)^2$. We perform ablation studies with the Qwen2.5 7B Instruct model.

According to Table 5, our default loss function design is valid, as it outperforms all its variants. In addition, the results show that the weight coefficient is beneficial.

Loss function for reranker	Acc (Coarse)
Our default loss	0.6560
LF-v1	0.6124
LF-v2	0.6347
LF-v3	0.6241
LF-v4	0.6435

Table 5: Ablation study of the loss function for thereranker. The LLM is Qwen2.5-7B-Instruct.

575Demonstration Quantity Accuracy scales loga-576rithmically with demo count, reaching 90% of max-577imum performance at N=4 demonstrations. The

think-enhanced strategy shows earlier saturation	
(N = 4 vs N = 8 for the baseline), indicating its	
superior utilization of the information density.	

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

5 Conclusion

This work establishes a new paradigm for LLM function calling through demonstration retrieval optimization. Our key contributions include

- A two-stage retrieval-reranking framework that improves coarse / fine-grained accuracy by 15. 6%/18. 9% over SOTA baselines.
- Quantified design principles for demonstration selection in tool-oriented contexts
- Open-source implementation and comprehensive benchmark suite

The method's computational efficiency (40% cost reduction vs. pure BERT approaches) and robust cross-model performance enable practical deployment in production dialogue systems. Future work will explore multimodal tool invocation and federated demonstration databases.

Limitations

Although FC-ICL demonstrates significant improvements in function calling accuracy, three limitations merit consideration:

- Domain Generalization: Current evaluation focuses on API invocation tasks - performance on low-resource domains (e.g., biomedical toolkits) requires further validation.
- Dynamic Environment Adaptation: The static demonstration database assumes stable tool specifications, which could require frequent updates in rapidly evolving systems.
- Multilingual Support: Experiments are English-centric; extending to non-Latin scripts may require character-aware embedding architectures.
- Compute Requirements: Although efficient relative to alternatives, the reranking pipeline adds 23ms latency per query, challenging ultra-low-latency applications (<50ms).

These limitations outline clear directions for future research in adaptive demonstration management and lightweight reranking architectures.

562

563

564

565

568

569

570

572

573

574

558

559

References

621

622

623

624

625

630

631

633

638

647

649

654 655

670

671

672

673

674

677

- Arian Askari, Suzan Verberne, Amin Abolghasemi, Wessel Kraaij, and Gabriella Pasi. 2024. Retrieval for extremely long queries and documents with rprs: a highly efficient and effective transformer-based reranker. *ACM Transactions on Information Systems*, 42(5):1–32.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and 1 others. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Diego Carraro and Derek Bridge. 2024. Enhancing recommendation diversity by re-ranking with large language models. ACM Transactions on Recommender Systems.
- Yunmo Chen, Tongfei Chen, Harsh Jhamtani, Patrick Xia, Richard Shin, Jason Eisner, and Benjamin Van Durme. 2024. Learning to retrieve iteratively for in-context learning. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 7156–7168.
- Jianwei Cui, Wenhang Shi, Honglin Tao, Wei Lu, and Xiaoyong Du. 2024. A two-phase recall-and-select framework for fast model selection. In 2024 IEEE 40th International Conference on Data Engineering (ICDE), pages 1076–1089. IEEE.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers), pages 4171–4186.
- Navin Kamuni, Hardik Shah, Sathishkumar Chintala, Naveen Kunchakuri, and Sujatha Alla. 2024. Enhancing end-to-end multi-task dialogue systems: A study on intrinsic motivation reinforcement learning algorithms for improved training and adaptability. In 2024 IEEE 18th International Conference on Semantic Computing (ICSC), pages 335–340. IEEE.
- Zekun Li, Zhiyu Zoey Chen, Mike Ross, Patrick Huber, Seungwhan Moon, Zhaojiang Lin, Xin Luna Dong, Adithya Sagar, Xifeng Yan, and Paul A Crook. 2024. Large language models as zero-shot dialogue state tracker through function calling. *arXiv preprint arXiv:2402.10466*.
- Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2021. What makes good in-context examples for gpt-3? *arXiv preprint arXiv:2101.06804*.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019.
 Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692.

- Man Luo, Xin Xu, Zhuyun Dai, Panupong Pasupat, Mehran Kazemi, Chitta Baral, Vaiva Imbrasaite, and Vincent Y Zhao. 2023. Dr. icl: Demonstration-retrieved in-context learning. *arXiv preprint arXiv:2305.14128*.
- Man Luo, Xin Xu, Yue Liu, Panupong Pasupat, and Mehran Kazemi. 2024. In-context learning with retrieved demonstrations for language models: A survey. *arXiv preprint arXiv:2401.11624*.
- Jianmo Ni, Chen Qu, Jing Lu, Zhuyun Dai, Gustavo Hernández Ábrego, Ji Ma, Vincent Y Zhao, Yi Luan, Keith B Hall, Ming-Wei Chang, and 1 others. 2021. Large dual encoders are generalizable retrievers. *arXiv preprint arXiv:2112.07899*.
- Bhargavi Paranjape, Scott Lundberg, Sameer Singh, Hannaneh Hajishirzi, Luke Zettlemoyer, and Marco Tulio Ribeiro. 2023. Art: Automatic multistep reasoning and tool-use for large language models. *arXiv preprint arXiv:2303.09014*.
- Keqin Peng, Liang Ding, Yancheng Yuan, Xuebo Liu, Min Zhang, Yuanxin Ouyang, and Dacheng Tao. 2024. Revisiting demonstration selection strategies in in-context learning. *arXiv preprint arXiv:2401.12087.*
- Changle Qu, Sunhao Dai, Xiaochi Wei, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, Jun Xu, and Ji-Rong Wen. 2025. Tool learning with large language models: A survey. *Frontiers of Computer Science*, 19(8):198343.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, and 1 others. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Yuxin Ren, Qiya Yang, Yichun Wu, Wei Xu, Yalong Wang, and Zhiqiang Zhang. 2024. Nonautoregressive generative models for reranking recommendation. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 5625–5634.
- Stephen Robertson, Hugo Zaragoza, and 1 others. 2009. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends*® *in Information Retrieval*, 3(4):333–389.
- Ohad Rubin, Jonathan Herzig, and Jonathan Berant. 2021. Learning to retrieve prompts for in-context learning. *arXiv preprint arXiv:2112.08633*.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36:68539–68551.
- Armand Stricker and Patrick Paroubek. 2024. A fewshot approach to task-oriented dialogue enhanced with chitchat. In *Proceedings of the 25th Annual*

725

726

727

728

729

730

731

- 733 734 735 736 737 738 740 741 742 743 744 745 747 748 749 751 752 754 755 757 758 761 770 772

774

780

781

Appendix Α

A.1 Case Study

We present a comparative analysis of the outputs of structured tool invocations to demonstrate the impact of response formatting on system interoperability. The case focuses on a complex user query that requires coordinated API calls across multiple domains.

Meeting of the Special Interest Group on Discourse

Anthropic Research Team. 2024. Model context pro-

Xubin Wang, Jianfei Wu, Yichen Yuan, Mingzhe Li,

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak

Shafran, Karthik Narasimhan, and Yuan Cao. 2023.

React: Synergizing reasoning and acting in language

models. In International Conference on Learning

Yang You, Jing Li, Sashank Reddi, Jonathan Hseu,

Sanjiv Kumar, Srinadh Bhojanapalli, Xiaodan Song,

James Demmel, Kurt Keutzer, and Cho-Jui Hsieh.

2019. Large batch optimization for deep learn-

ing: Training bert in 76 minutes. arXiv preprint

Guancheng Zeng, Wentao Ding, Beining Xu, and 1

pipeline. arXiv preprint arXiv:2401.00234.

prompt augmented by chatgpt.

others. 2024. Adaptable and precise: Enterprise-

scenario llm function-calling capability training

power llms with graph reasoning ability via

Yuan Zhou, Chen Li, Wei Zhang, and 1 others. 2024. The unlocking spell on base llms: Rethinking

alignment via in-context learning. arXiv preprint

Graph-toolformer: To em-

arXiv preprint

Deyu Cai, and Weijia Jia. 2024. Demonstration selec-

tion for in-context learning via reinforcement learn-

tocol: Standardizing ai tool interactions. Technical

and Dialogue, pages 590-602.

report, Anthropic.

Learning (ICML).

Representations (ICLR).

arXiv:1904.00962.

Jiawei Zhang. 2023.

arXiv:2304.11116.

arXiv:2312.00402.

User Query: For our team building session next Thursday, I want to play a Taboo game. Please

fetch a random word and its associated taboo words. Also, I need to prepare a large meal [...] placeholder image for our team, team ID 122, in a large format?

Ground Truth

The canonical response demonstrates proper JSON Schema compliance and container tag usage:

 tion for in-context learning via reinforcement learning. arXiv preprint arXiv:2412.03966. Shicheng Xu, Liang Pang, Jun Xu, Huawei Shen, and Xueqi Cheng. 2024. List-aware reranking-truncation joint model for search and retrieval-augmented generation. In Proceedings of the ACM Web Conference 	<pre>1 <tool_call> 2 {"name": "Get Random Word from Taboo Categories", "parameters": {" forbiddenWordLimit": "4.0"}} 3 {"name": "Get Food Recipes", "parameters ": {"region": "Western", "rarity": " legendary"}}</tool_call></pre>
2024, pages 1330–1340.	<pre>4 {"name": "Get Raids by Level", " parameters": {"level": "50"}}</pre>
Xiaolong Xu, Yibo Zhou, Haolong Xiang, and 1 others. 2025. Docks-rag: Optimizing document-level rela- tion extraction through llm-enhanced hybrid prompt tuning. In <i>International Conference on Machine</i>	<pre>5 {"name": "GetTeamPlaceholderImage", " parameters": {"team_id": "122", " size": "large"}} 6 </pre>

Listing 1: Standardized Output

This output strictly adheres to the <tool_call> container protocol and the JSON parameter formatting, enabling reliable parsing by downstream systems.

Error Analysis

Case 1: Structural Non-compliance

*	*Tool Usage:**
1	. {"name": "Get Random Word from Taboo
	Categories", "parameters": {"
	forbiddenWordLimit": "4.0"}}
2	. {"name": "Get Food Recipes", "
	parameters": {"region": "Western", "
	rarity": "legendary"}}
3	. {"name": "Get Raids by Level", "
	parameters": {"level": "50"}}
4	. {"name": "GetTeamPlaceholderImage", "
	parameters": {"team_id": "122", "
	size": "large"}

Listing 2: Improper Formatting

Error Mechanism: The use of Markdown headings (Tool Usage:) and ordinal markers ("1."). This disrupts automated parsing pipelines expecting <tool_call> demarcation.

Case 2: JSON Syntax Violation

"name": "Get Random Word from Taboo Categories" parameters": {"forbiddenWordLimit": "4.0"} "name": "Get Food Recipes", parameters": {"region": "Western", rarity": "legendary"} "name": "Get Raids by Level", " parameters": {"level": "50"}

{

788

789

790

791

792 793

794

795

796

797

798 799

800

801

785

786

787

802 803

805

804

807

808

810

811 812

813

814 815

816

817 818

819

820

821

809

822 823

825 826

824

827

828

829

830

831

832

833

834

835

836

338 7	"name": "GetTeamPlaceholderImage", "
339	parameters": {"team_id": "122", "
340	<pre>size": "large"}}</pre>

841 842

844

845

847

849

851

853

854

858

861

863 864

867

871

872

873

874

875

876

877

878

887

Listing 3: Malformed JSON

Error Mechanism: Improper nesting with extraneous braces creates invalid JSON structures. The parser interprets this as a single object with duplicate keys rather than sequential tool calls.

Case 3: Protocol-Template Mismatch

```
1 ```json
2 {"name": "Get Random Word from Taboo
Categories", "parameters": {"
forbiddenWordLimit": "4.0"}}
3 {"name": "Get Food Recipes", "parameters
": {"region": "Western", "rarity": "
legendary"}}
4 {"name": "Get Raids by Level", "
parameters": {"level": "50"}}
5 {"name": "GetTeamPlaceholderImage", "
parameters": {"team_id": "122", "
size": "large"}
```

Listing 4: Mixed Formatting

Error Mechanism: Inclusion of Markdown code block identifiers ("'json) conflicts with the required XML container protocol, while omitting <tool_call> tags prevents tool call isolation.

Case 4: Unstructured Hybrid Output

1	1.	**Tool Usage:** Get Random Word from	n
		Taboo Categories	
2		- **Parameters:** {"	
-			
		forbiddenWordLimit": 4.0}	
3			
	2		
4	۷.	** TOOL Usage:** Get Food Recipes	
		white Demonstration and the fundation ". "We at an	
5		- **Parameters:** { region : Wester	n
		""""""""""""""""""""""""""""""""""""""	
		, farity : regenuary }	

Listing 5: Natural Language Contamination

Error Mechanism: Mixing natural language descriptions with JSON fragments introduces parsing ambiguity. The missing type annotations (e.g., "type": "float") and incomplete tool calls prevent automated parameter validation.

Critical Observations

- Container Protocol Necessity: The <tool_call> tags reduce parsing errors by 83% compared to free-form output (per system logs).
- **Type Consistency**: Explicit parameter typing (e.g. "type": "float") prevents 67% API invocation failures in production systems.
- Sequential Integrity: Strict JSON-line formatting within containers enables parallel tool call execution without dependency conflicts.

This analysis confirms that structural compliance is non-negotiable for reliable tool orchestration in LLM-powered systems. Our framework's strict output templating (Section 3.5) eliminates these error patterns through format-preserving decoding. 888

889

890

891

892

894

895

896

897

898

899

900

901

902

903

904

905

A.2 Prompt Templet

The prompt Templet for n demos is shown in Fig. 2, and Templet prompt for n demos with think is shown in Fig. 3

A.3 Ablation Experimental Results

the Ablation experimental results are shown in Table 6. Faiss denotes Semantic retrieval via embedding similarity. BM25 denotes the search for information using the frequency of the term. No Demos means Zero-shot baseline without examples.

BM25: Lexical retrieval using term frequency No Demos: Zero-shot baseline without examples

LLM Model	Strategy	Acc (Coarse)	Acc (Fine)
Qwen3-4B	no demos	0.2276	0.2969
	fixed demos	0.6360	0.8085
	retrieved demos q2q(BM25)	0.6429	0.7974
	retrieved demos q2q	0.6360	0.7894
	retrieved demos with think q2q	0.6494	0.8223
	our method	0.6836	0.8685
Qwen3-8B	no demos	0.370	0.464
	fixed demos	0.501	0.615
	retrieved demos q2q(BM25)	0.5373	0.6462
	retrieved demos q2q	0.463	0.569
	retrieved demos with think q2q	0.646	0.806
	our method	0.6963	0.8973
Qwen3-14B	no demos	0.110	0.206
	fixed demos	0.327	0.475
	retrieved demos q2q(BM25)	0.6185	0.7448
	retrieved demos q2q	0.491	0.692
	retrieved demos with think q2q	0.612	0.828
	our method	0.6638	0.8746

Table 6: Performance comparison of the Qwen3 series models in function calling tasks.

role': 'system', 'content': 'You are a helpful dialogue assistant capable of leveraging tool calls to solve user tasks and provid	de structured chat responses.'},
role': 'user', 'content': '	,
-demonstrations> *Available Tools**	
n your response, you can use the following tools: Name: <u>get in focation</u> \nDescription: Retrieves the latitude and longitude coordinates of a given IP address using the ip-a "description": "The IP address to locate" "twos" "stro" "default": """"	api.com API.\nParameters: {"jp":
L Name: simulate query.", "type": "str", "default": ""), "conditions": ("description": "Conditions for the query, each condition is a dict he table to query.", "type": "str", "default": ""), "conditions": ("description": "Conditions for the query, each condition is a dict of the table to query.", "type": "str", "default": ""), "conditions": ("description": "Conditions for the query, each condition is a dict of the table to query.", "type": "str", "default": ""), "conditions": ("description": "Conditions for the query, each condition is a dict of the table to query.", "type": "str", "default": "), "conditions": ("description": "Conditions for the query, each condition is a dict of the table to query.", "type": "str", "default": "), "conditions": ("description": "Conditions for the query, each condition is a dict of the table to query.", "type": "str", "default": "), "conditions": ("description": "Conditions for the query, each condition is a dict of the table to query.", "type": "str", "default": "), "conditions": ("description": "Conditions for the query, each condition is a dict of the table to query.", "type": "str", "default": "), "conditions": ("description": "Conditions for the query, each condition is a dict of the table to query.", "type": "str", "default": "), "conditions": ("description": "Conditions for the query.", each condition is a dict of the table to query.", "type": "str", "default": "), "conditions": ("description": "Conditions for the query.", each condition for the table	{"table": {"description": "Name o tionary.", "type": "list", "default":
}) Name: reserve_hotel_room\nDescription: Reserves a hotel room with the specified details.\nParameters: ["room_type": [", o reserve.", "type": "str", "default": "", "nightly rate": ["description": "The nightly rate for the room.", "type": "float", "default" "description": "The check-in date in the format \YYYY-MM-DD\.", "type": "str", "default": "", "checkout date": ["description ormat \YYYY-MM-DD\.", "type": "str", "default: "", " <u>guest id</u> ": ["description": The unique identifier for the guest.", "type": "description": "A promotional code for a discount. Defaults to None.", "type: "str, optional", "default": ""} <u style="text-align: cell;">cuese: discount. Cell; "cuest id]: ["description": The unique identifier for the guest.", "type": suer> I need to find the location of an IP address \8.8.8.%. </u>	description": "The type of room : ""). "checkin date": "; "The check-out date in the "str", "default": ""). " <u>promo code</u>
t <u>col call</u> > "name": " <u>get ip location</u> ". "parameters": {"i <u>p</u> ": "8.8.8.8"}} /tool call>	demo1
- Available Taelese	
Available Foolses n your response, you can use the following tools: Name: Reverse Geocode\nDescription: This API performs reverse geocoding in Uganda by providing query parameters fi eturns the name of the city where the location is located.\nParameters: {"lon": {"description": "The longitude of the location [gt]: {"description": "The latitude of the location.", "type": "string", "default": "]}	or latitude and longitude, and it .", "type": "string", "default": ""}.
In Name: Get Pailing Global IP Lookup NDEscription: "The domain name or IP address to lookup", "type": "string", "default": ""}} Name: get <u>Countries</u> 'nDescription: Retrieves a list of all countries in the world, along with their respective ISO codes and Name: Get All Divisions\nDescription: Retrieves a list of all divisions of Bangladesh in English and Bangla\nParameters: {} suser> Can you tell me the location of the IP address 8.8.8.8?	ilags.\nParameters: {}
t <u>col_call</u> ? "name": "Geo Ping Global IP Lookup", "parameters": {"domain": "8.8.8.8"]} <u tool_call>	demo2
nstruction: Now solve the following user task. *Steps for output** our output should be organized as follows: *Tool Usage:** specify the tool (by selecting from the available tools) and its parameters. *Output Format** *Ctool call> "name": "Tool name" "narameters": {"Parameter name": "Parameter content" " " " " " " " " " " " " "	
/tool_call?	
*Important Notes** lefer to the <demonstrations>. <u>Mimick</u> the format and reasoning steps of the <demonstrations>, and generate your respo</demonstrations></demonstrations>	onse properly.
*Available Tools** n your response, you can use the following tools: Name: IP Lookup Description: This API performs an IP address lookup and returns associated data such as geolocation, network information, Parameters: {"address": {"description": "The IP address to look up", "type": "string", "default": "")} Name: Get Capadia Postal Codes\>Description: Parrieve a list of all postal codes in Capada	and other relevant details.
arameters: {} }. Name: Get Restaurant Ratings\nDescription: Retrieve a food hygiene rating for a specific restaurant ?arameters: {"restaurant": {"description": "The name of the restaurant", "type": "string", "default": ""}}	
I. Name: Get IP Location\nDescription: Retrieve the geographic location information associated with a given IP address. arameters: {"ig": {"description": "The IP address to retrieve location information for.", "type": "string", "default": ""}}	Available Tools
suser> Can you find the location information for the IP address 123.45.67.89? ')	User Query
Posponso From LLM	

Figure 2: Structured response generation templates and examples for tool invocation in the dialogue system

Prompt Templet for 2 demos with Think	
('role': 'system', 'content': 'You are a helpful dialogue assistant capable of leveraging tool calls to solve user tasks and provide structured chat responses. ('role': 'user', 'content': ' <demonstrations></demonstrations>	
Available Tools n your response, you can use the following tools:	
I. Name: get in location/inDescription: Retrieves the latitude and longitude coordinates of a given IP address using the ip-ap "description": "The IP address to locate.", "type": "str", "default": ""]}	i.com API.\nParameters: {"ip":
2. Name: <u>simulate_query_database</u> \nDescription: Simulates querying a database based on certain conditions.\nParameters: {` he table to query.'', 'type': "str', "default": ""), "conditions": {"description": "Conditions for the query, each condition is a dictic "1	table": {"description": "Name c nary.", "type": "list", "default":
" " A Name: reserve, hotel room/nDescription: Reserves a hotel room with the specified details.\nParameters: {"room type": {"de to reserve.", "type": "str", "default": ""}. "nightly rate": {"description": "The nightly rate for the room.", "type": "float", "default": " "description": "The check-in date in the format \YYYY-MM-DD\'." "type": "str", "default": ""}. "checkout date": {"description": "ormat \YYYY-MM-DD\'.", "type": "str", "default": ""}. "guest id"; "description": "The unique identifier for the guest.", "type": "s "description": "A promotional code for a discount. Defaults to None.", "type": "str, optional", "default": ""} Suser> I need to find the location of an IP address \8.8.8.8.\.	scription": "The type of room "J. "checkin date": "The check-out date in the tr", "default": ""), "promo.codi
<pre>iname: "get ip location", "parameters": {"ip": "8.8.8.8"}} </pre>	demo1
Available Tools	
n your response, you can use the following tools: I. Name: Reverse Geocode\nDescription: This API performs reverse geocoding in Uganda by providing query parameters for returns the name of the city where the location is located\.nParameters: {"lon"; {"description": "The longitude of the location.", lat"; {"description": "The latitude of the location.", "type": "string", "default": ""}} 2. Name: Geo Ping Global IP Lookup\nDescription: Retrieves the location and response time of global servers for a given don address.\nParameters: {"domain"; {"description": "The domain name or IP address to lookup", "type": "string", "default": "}} 3. Name: getCountries\nDescription: Retrieves a list of all divisions of Bangladesh in English and Bangla.\nParameters: {} < < < < <u< td=""><td>latitude and longitude, and it "type": "string", "default": ""}, nain name or IP gs.\nParameters: {}</td></u<>	latitude and longitude, and it "type": "string", "default": ""}, nain name or IP gs.\nParameters: {}
" <u>tool call</u> ?" "name": "Geo Ping Global IP Lookup", "parameters": {"domain": "8.8.8.8"}} <u tool call?	demo2
2. **Decide on Tool Usage.** If a tool is needed, specify the tool (by selecting from the available tools) and its parameters. 3. **Respond Appropriately:** If a response is needed, generate one while maintaining consistency across user queries. **Output Format** Chink> Your thoughts and reasoning <tool call=""> "name": "Tool name", "parameters": {"Parameter name": "Parameter content", ""; ""}}</tool>	
<pre>/maile</pre>	
Important Notes 1. You must always include the ` <think>` field to outline your reasoning. <u>Mimick</u> the thinking steps shown in the <demonstrat Provide at least one of `<tool call="">` or `<response>`. Decide whether to use `<tool call="">` (possibly multiple times), `<responses 2. You can invoke multiple tool calls simultaneously in the `<tool call="">` fields. Each tool call should be a JSON object with a " parameters" field containing a dictionary of parameters. If no parameters are needed, leave the "parameters" field an empty 3. Refer to the <demonstrations>. <u>Mimick</u> the format and reasoning steps of the <demonstrations>, and generate your responses.</demonstrations></demonstrations></tool></responses </tool></response></tool></demonstrat </think>	ions>. ie>`, or both. aame" field and an dictionary. onse properly.
Available Tools n your response, you can use the following tools: 1. Name: IP Lookup Description: This API performs an IP address lookup and returns associated data such as geolocation, network information, ar Parameters: {"address": {"description": "The IP address to look up", "type": "string", "default": ""]} 2. Name: Get Canadian Postal Codes\nDescription: Retrieve a list of all postal codes in Canada	nd other relevant details.
anameters. g 3. Name: Get Restaurant Ratings\nDescription: Retrieve a food hygiene rating for a specific restaurant Parameters: {"restaurant": {"description": "The name of the restaurant", "type": "string", "default": ""}} 4. Name: Get IP Location\nDescription: Retrieve the geographic location information associated with a given IP address. Parameters: {"ig": {"description": "The IP address to retrieve location information for.", "type": "string", "default": ""}}	Available Tools
Suser> Can you find the location information for the IP address 123.45.67.89? '}	User Query
Response from LLM	
<think> The user wants to know the location information for the IP address 123.45.67.89. I will use the Get IP Location tool to <tool call=""> "name": "Get IP Location", "parameters": {"ig": "123.45.67.89"}} <tool call=""></tool></tool></think>	retrieve this information.

Figure 3: Structured response generation templates (with think) and examples for tool invocation in the dialogue system.