# ENHANCING ZERO-SHOT REASONING IN LANGUAGE MODELS VIA HYBRID INSTRUCTION MARGINALIZATION

## Shirmohammad Tavangari[1], Aref Yelği[2]

[1] *University of British Columbia, Canada*
[2] *Istanbul Topkapi University*

**Corresponding author:**
Shirmohammad Tavangari.
Email : s.tavangari@alumni.ubc.ca

**Conflict of interest statement:**
Author(s) reported no conflict of interest

## ABSTRACT

**Objective**: The study aims to enhance the reasoning abilities of Large Language Models (LLMs), which often remain shallow, inconsistent, and error-prone in complex multi-step tasks. It introduces the Hybrid Instruction Tuning Framework (HITF) to improve zero-shot reasoning through a task-aware hybrid selector that integrates both human-annotated and automatically generated examples.

**Research Design & Methods**: HITF strengthens reasoning performance using three main techniques: synthesizing transitional results, context-aware prompt merging, and recurrent optimization, all executed without model recalibration. The framework is empirically evaluated using rigorous cognitive benchmarks, including SuperGLUE, MMLU, GSM8K, and FermiQA. Component isolation tests examine the independent contribution of the example selector, output synthesizer, and instruction combiner. Statistical variability assessments further validate result reliability.

**Findings:** Results show that HITF consistently outperforms state-of-the-art methods across multiple metrics, demonstrating higher measurement accuracy, stronger argumentative quality, and deeper analytical processing. All core modules exhibit significant and measurable contributions, supported by stable statistical outcomes.

**Implications & Recommendations:** Findings suggest that combining context-driven instruction selection with statistical consolidation techniques can substantially improve deductive reasoning in LLMs, particularly in data-scarce and example-free settings. Future research should explore HITF′s integration with larger models and its application in real-world reasoning-intensive domains.

**Contribution & Value Added:** This study offers an innovative framework that enhances zero-shot reasoning without retraining. By merging hybrid instruction selection and iterative optimization strategies, HITF narrows the reasoning gap between LLMs and humans and provides a scalable, reliable approach for advancing high-level reasoning in modern language models.

JEL codes: C45, C55, O33.
**Article type:** research paper

## INTRODUCTION

The pursuit of artificial intelligence has long been driven by the goal of creating systems that can not only understand human language but also reason with it mimicking the human capacity to analyze complex problems, draw logical inferences, and formulate robust solutions. As large language models (LLMs) increasingly power a wide array of applications, from intelligent assistants to scientific discovery tools, their ability to perform reliable and sophisticated reasoning has become paramount. Enhancing this capability is crucial for bridging the gap between passive information retrieval and genuine cognitive assistance, ultimately enabling AI to tackle more complex, real-world challenges.

Cutting-edge neural architectures have propelled language models to unprecedented levels of performance in various domains of natural language understanding (Kojima et al., 2022; Wei et al., 2022). These foundational advancements have established a strong base upon which higher-order cognitive skills, such as reasoning, can be developed.

Despite the remarkable progress enabled by these advanced architectures, a significant research gap persists when these models are confronted with complex reasoning demands especially tasks that involve structured, multi-step logical analysis or the precise manipulation of symbolic information. While substantial improvements have been made in areas such as contextual cue optimization and sample efficient training, these enhancements have primarily benefited few-shot or fine-tuned scenarios (Madaan et al., 2023; X. Wang et al., 2022; Zhou et al., 2023). The critical challenge of achieving robust and reliable reasoning in zero-shot settings, where no task-specific examples are provided, remains largely unaddressed.

Core difficulties including sustaining propositional coherence throughout extended derivations, preventing error propagation across sequential inference steps, and attaining robust generalization to out-of-domain problems without prior exemplars continue to pose fundamental limitations (Yao, Yu, et al., 2023). Existing methodologies often fail to bridge this gap, as they either rely heavily on curated demonstrations, incur prohibitive computational costs, or lack the adaptability for broad zero-shot application.

Accordingly, an evident demand emerges for an optimized, adaptable, and resource-conscious framework that effectively synthesizes the synergistic advantages of both human-verified and algorithmically-generated insights to strengthen analytical reasoning without exceeding practical computational constraints (Chen et al., 2021; Srivastava et al., 2023; Tavangari et al., 2024).

Therefore, the key research gap we identify is the lack of a computationally efficient framework that dynamically adapts to the reasoning requirements of a given task in a zero-shot setting, effectively bridging the quality of human-annotated data with the diversity of self-generated content to produce robust and generalizable logical inferences.

### Problem Statement

Let $x$ be a natural language task, and $y$ its correct response. Given two instruction-tuning datasets:

$$D_{\text{human}} = \left\{ \left( x_i^{(h)}, y_i^{(h)} \right) \right\} \quad \text{(human-annotated)} \quad (1)$$

$$D_{\text{self}} = \left\{ \left( x_i^{(s)}, y_i^{(s)} \right) \right\} \quad \text{(Self-generated)} \quad (2)$$

Our goal is to predict $\hat{y}$ for unseen tasks $x$, maximizing:

$$\llbracket 1(\hat{y} = y) \rrbracket \quad (3)$$

$D_{\text{human}}$ has high accuracy but low diversity, while $D_{\text{self}}$ is more diverse but noisy. Relying on one of them, especially in the zero-shot mode, reduces the generalizability of the model. We present a hybrid method that dynamically combines instances from $D_{\text{human}}$ and $D_{\text{self}}$ based on the task type $x$, and ensures stable and generalizable performance of large language models by selecting appropriate instances.

## LITERATURE REVIEW

The quest to equip Large Language Models (LLMs) with robust reasoning capabilities has led to several prominent strands of research. These approaches can be broadly categorized into methods that enhance reasoning through prompting strategies and those that leverage data-driven fine-tuning. This section reviews the evolution of these key paradigms, critiquing their underlying principles and highlighting the limitations that our work aims to address.

### Prompting-Based Reasoning Methods

A significant line of research focuses on eliciting reasoning without updating model parameters, primarily through sophisticated prompting techniques.

Chain-of-Thought (CoT) Prompting. The seminal work of introduced Chain-of-Thought (CoT) prompting, which instructs the model to generate a step-by-step rationale before producing a final answer. This approach demonstrated that LLMs can perform complex reasoning by decomposing problems into intermediate steps. However, standard CoT exhibits a pronounced dependency on meticulously crafted, expert-curated demonstrations. Its performance is highly sensitive to the choice of examples, and it often generalizes poorly to unseen task types in a zero-shot setting, where such examples are unavailable.

Self-Consistency (SC). To address the variability in CoT reasoning paths proposed Self-Consistency. This technique involves sampling multiple reasoning paths from the model and then aggregating the final answers by marginalizing out the intermediate steps, typically through a majority vote. While SC enhances the robustness and credibility of the final output, it comes at a significant computational cost, as it requires multiple generations per query. Furthermore, its performance can degrade if a substantial number of the sampled paths are flawed, making it susceptible to common errors in the model's reasoning process.

Advanced Reasoning Structures. Building upon CoT, recent efforts have explored more structured reasoning processes. The Tree-of-Thoughts (ToT) framework models reasoning as a tree exploration problem, allowing the model to evaluate and backtrack over multiple coherent units of thought. Similarly, Graph-of-Thoughts (GoT) generalizes this further by representing thoughts as a graph, enabling more complex combinatorial reasoning. While these methods can achieve impressive results, they require carefully designed prompt configurations, complex state management, and extensive LLM computations for evaluation and exploration, making them resource-intensive and often impractical for real-time applications.

### Fine-Tuning and Bootstrapping Methods

Another direction involves refining the model's reasoning capabilities through data-driven training. STaR (Self-Taught Reasoner). STaR bootstraps a model's reasoning by fine-tuning on a set of problems where the model itself generates rationales. If a generated rationale leads to a correct answer, the (problem, rationale, answer) triplet is added to the training set. This iterative self-training process improves reasoning over time. Nevertheless, STaR and similar bootstrapping methods can be computationally expensive, prone to error propagation if incorrect rationales are reinforced, and often require a pre-existing seed set of CoT examples to initialize the process effectively.

Iterative Self-Refinement. Approaches like Self-Refine employ an iterative feedback loop where the model generates an output, provides self-critique, and then refines its initial response based on that feedback. This mirrors a form of internal deliberation. However, prior methods primarily focus on enhancing a singular inference pathway. In comparison, our method strengthens the reasoning process through the probabilistic integration of transitional outcomes from multiple support sets and a context-aware composition of exemplars, leading to more thorough and stable improvements (Bai et al., 2022; Zelikman et al., 2024).

**Identified Research Gap**

The reviewed literature reveals a consistent trade-off between performance, generality, and computational efficiency. Prompting-based methods like CoT and ToT struggle with zero-shot generalization and scalability. Sampling-based methods like SC improve reliability but are computationally prohibitive. Bootstrapping methods like STaR require significant training overhead and are vulnerable to error propagation (Hoffmann et al., 2022).

This analysis highlights a critical research gap: the absence of a scalable and adaptive framework that achieves robust zero-shot reasoning without extensive task-specific demonstrations, prohibitive computational costs, or complex, manually-engineered prompting structures. Existing methods lack a dynamic mechanism to seamlessly integrate the guaranteed quality of human-annotated data with the broad diversity of self-generated reasoning paths in a task-aware manner. Our proposed Hybrid Instruction Tuning Framework (HITF) is designed to bridge this gap by introducing a lightweight, dynamic selector that optimally blends these data sources, enabling efficient and effective zero-shot reasoning (B. Wang et al., 2023).

**Related Work**

Several existing reasoning methodologies attempt to improve accuracy by generating and combining multiple reasoning paths, although these techniques often increase inference time and heighten sensitivity to invalid trajectories. In the area of Hybrid Fine-Tuning, previous studies typically rely on fixed ratios to merge human-curated and machine-generated data; however, our method provides a dynamic, task-specific strategy managed by a lightweight selection module that adjusts instructional inputs more effectively (Yelghi and Tavangari, 2023). Iterative Self-Correction approaches, such as those introduced by Madaan et al. (2023) use cyclic optimization driven by internally generated critiques and tend to focus on enhancing a single inference pathway. In contrast, our framework strengthens the refinement process through probabilistic integration of transitional outcomes and context-aware exemplar composition, enabling more comprehensive and impactful reasoning improvements (Touvron et al., 2023; Yelghi, Tavangari, et al., 2024; Zhai et al., 2022).

Additionally, organized inference methodologies like Tree-of-Thought and STaR require carefully engineered prompt structures and supervised calibration, which can be computationally intensive and inflexible. Our approach addresses these limitations by incorporating randomized optimization and adaptive example selection, resulting in a resource-efficient solution that avoids complex decision trees and extensive supervised training (Gao et al., 2023; Tavangari et al., 2025; Yao, Zhao, et al., 2023; Yelghi, Yelghi, et al., 2024). Moreover, the proposed architecture enhances inferential capability in ambiguous contexts through an integrated schema alignment mechanism and versatile extraction from diverse information sources, supporting more robust and adaptable reasoning performance.

## METHODS

In this paper, we present the Adaptive Instruction Blending Framework (AIBF), an original methodology that addresses this limitation by intelligently combining expert-curated and autonomously produced guidance via a context-sensitive selection mechanism (OpenAI, 2023). AIBF improves reasoning without task-specific examples via three key mechanisms: (1) systematic aggregation of transitional results to increase reliability, (2) adaptive merging of pertinent illustrations specifically aligned with the given challenge, and (3) a recurrent enhancement procedure that operates absent network adjustments. This methodology facilitates the seamless integration of trustworthy human knowledge and diverse autonomously-generated analytical trajectories. We conduct thorough evaluations of our framework across multiple standardized assessments including SuperGLUE, MMLU, GSM8K, and FermiQA.

We propose a Hybrid Instruction Tuning Framework (HITF) that synergistically combines human-annotated and self-generated instruction data to enhance the generalization and reasoning capabilities of large language models (LLMs). The core idea is to dynamically construct a support set

of examples tailored to each input task $x$, balancing quality (from human data) and diversity (from self-generated data) via a learned mixture weight $\lambda(x)$.

**Formal Problem Setup**

Let $x$ denote an input task (e.g., a question or instruction) and $y$ its corresponding true response. We assume access to two distinct instruction-tuning datasets:

a. **Human–annotated data:**

$$D_{\text{human}} = \left\{ \left( x_i^{(h)}, y_i^{(h)} \right) \right\} i^{N_h} = 1 \qquad (4)$$

b. **Self-generated data:**

$$D_{\text{self}} = \left\{ \left( x_j^{(s)}, y_j^{(s)} \right) \right\} j^{N_s} = 1 \qquad (5)$$

For a new task $x$ our goal is to predict the correct output $\hat{y}$ that maximizes the expected accuracy:

$$\mathbb{E}_{(x,y)} = [\![ 1(\hat{y} = y) ]\!] \qquad (6)$$

**Dynamic Example Selection**

Given an input $x$, we first compute its dense vector representation using a pretrained encoder:

$$e_{\text{x}} = f_{\text{enc}}(x), e_x \in \mathbb{R}^{\text{d}} \qquad (7)$$

This embedding is used by a lightweight mixture selector network $g_\theta$. To produce a task-dependent mixing coefficient:

$$\lambda(x) = g_\theta, \lambda(x) \in [0,1] \qquad (8)$$

Here , $\lambda(x)$ controls the proportion of examples drawn from $D_{human}$ vs. $D_{self}$ The selector $g_\theta$ implemented as a two-layer MLP with ReLU activation and 256 hidden units.

**Support Set Construction**

Let $k$ be the total number of support examples. We sample:

$$k_h = \lceil \lambda(\text{x}). \text{k} \rceil \text{ examples from } D_{\text{human}} \qquad (9)$$

$$k_s = k - k_h \ examples \ from \ D_{self} \qquad (10)$$

Each subset is selected based on cosine similarity to the query embedding $e_x$. The final support set is:

$$S = S_h \bigcup S_{\text{s}} , |S_h| = \text{k}_{\text{h}}, |S_{\text{s}}| = k_s \qquad (11)$$

**Marginalization Over Prompts**

For each support set $S$, we construct a prompt $P(s)$ by concatenating its examples. The LLM then produces a conditional distribution over answers:

$$P_\theta(y \mid x, P(S)) \qquad (12)$$

To marginalize over the variability in support set selection, we aggregate predictions over multiple plausible sets $S \in S_x$:

$$\hat{P}(y \mid x) = \sum_{S \in S_x} p(S \mid x) . p_\theta(y \mid x, P(S)) \qquad (13)$$

Where $p(S \mid x)$ is a probability distribution over support sets, proportional to their similarity to $x$.

**Training Objective**

We minimize the negative log-likelihood of the true answer $y$ under the marginalized distribution:

$$\mathcal{L}(\theta) = -E_{(x,y) \sim D}[\log \hat{p}(y \mid x)] \qquad (14)$$

---

**Hybrid Instruction Tuning with Marginalization**

1 Compute task embedding
$$e_x \leftarrow \text{Encoder}(x)$$
2 Compute mixture weight
$$\lambda \leftarrow \lambda(x)$$
3 Sample human examples
$$k_h \leftarrow \lceil \lambda k \rceil \text{ from } D_{\text{human}}$$
4 Sample self-generated examples
$$k_s \leftarrow k - k_h \text{ from } D_{\text{self}}$$
5 Form prompt set
$$S \leftarrow S_h \cup S_s$$
6 For each prompt $P_i$ in $S$
   6a Query model $y_i \leftarrow M x, P_i)$
   6b Compute confidence $p_{i(}, a = y)$
7 Aggregate predictions
$$\hat{y} - \arg\max \sum_i p_i \, 1(y_i = y)$$
8 Return

---

Algorithm 1. Hybrid Instruction Tuning With Marginalization

The full procedure is summarized in Algorithm 1, which outlines the dynamic mixture selection and marginalization process over sampled prompt variants.

```
┌─────────────────┐
│  Input: Task x  │
└─────────────────┘
         │
         ▼
┌──────────────────────┐
│ Compute Embedding e_x│
│     via Encoder      │
└──────────────────────┘
         │
         ▼
┌──────────────────────┐
│ Selector Network g_θ │
└──────────────────────┘
         │
         ▼
┌──────────────────────┐
│ λ(x) = Mixture Weight│
└──────────────────────┘
       ╱        ╲
      ▼          ▼
┌──────────────────┐  ┌──────────────────┐
│Sample k_h from   │  │Sample k_s from   │
│D_human           │  │D_self            │
└──────────────────┘  └──────────────────┘
       │                    │
       ▼                    ▼
┌──────────────────┐  ┌──────────────────┐
│ Support Set S_h  │  │ Support Set S_s  │
└──────────────────┘  └──────────────────┘
         ╲              ╱
          ▼            ▼
     ┌──────────────────────┐
     │  Merge Support Sets: │
     │   S = S_h ∪ S_s      │
     └──────────────────────┘
              │
              ▼
┌──────────────────────────────────┐
│ Construct Multiple Prompts P_i   │
│             from S               │
└──────────────────────────────────┘
              │
              ▼
┌──────────────────────┐
│ Run Large Language   │
│ Model for each Prompt│
└──────────────────────┘
              │
              ▼
┌──────────────────────────┐
│   Aggregate Outputs      │
│Marginalization over      │
│       Prompts            │
└──────────────────────────┘
              │
              ▼
┌──────────────────────┐
│ Final Prediction: ŷ  │
└──────────────────────┘
```
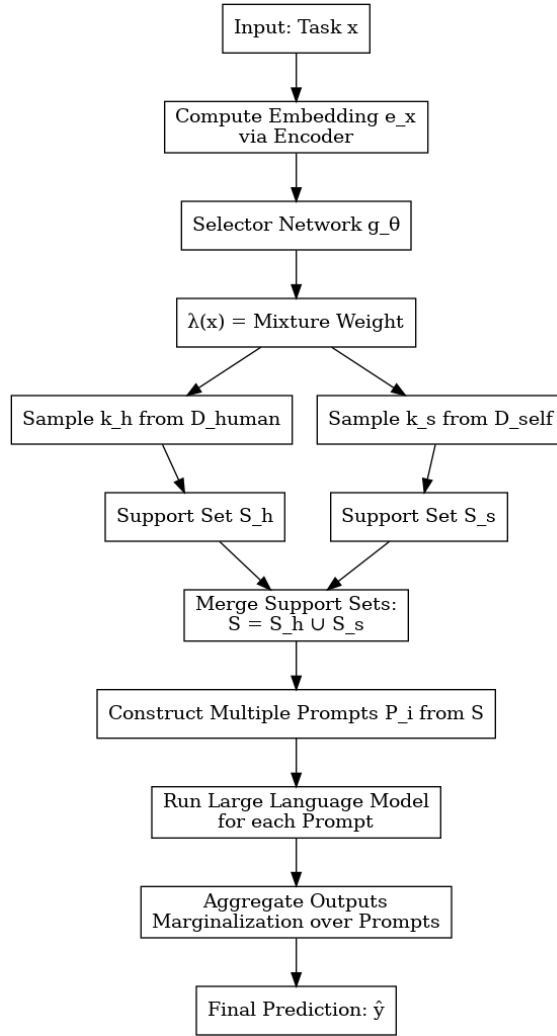
Figure 1. Architecture of the proposed HITF framework

The overall architecture of the proposed HITF framework is illustrated in Figure 1.

## RESULT

### Implementation

#### Step 1: Task Embedding

The function $f_{enc}$ converts the text $x$ into the feature vector $e_x$.

$$e_x = f_{enc}(x), e_x \in \mathbb{R}^d \qquad (15)$$

#### Step 2: Dynamic Selector Computation

A lightweight model $g_\theta$ takes $e_x$ as input and outputs a mixture weight $\lambda \in [0,1]$:

$$\lambda = g_\theta(e_x) \qquad (16)$$

The function $g_\theta$ is a two-layer MLP with 256 hidden nodes and a ReLU activation function. It is trained using contrastive learning and adds only 2 milliseconds of latency per query. Training on an A100 GPU takes 12 hours and accounts for just 5% of the base model's training cost.

#### Step 3: Support set construction

Given $\lambda$ and $k$, the number of human and self-generated samples is computed:

$$k_h = \lceil \lambda \cdot k \rceil, k_s = k - k_h \qquad (17)$$

The final support set is:

$$S = S_{human} \cup S_{self}, |S_{human}| = k_h, |S_{self}| = k_s \qquad (18)$$

### Step 4: Fine-tuning or In-Context Prompting

Using the support set $S$, the query to the Large Language Model (LLM) is done as follows:

$$\hat{y} = LLM(x, P(S)) \qquad (19)$$

where $P(S)$ is a prompt constructed from the support set.

### Step 5: Loss and Optimization

The model is trained with a cost function that computes the mean negative logarithm of the probability of a correct answer and optimizes the parameter $\theta$

$$L(\theta) = -\frac{1}{|T|} \sum_{(x,y) \epsilon T} \log P_\theta (\hat{y} = y \, | x, P(S)) \qquad (20)$$

### Step 6: Iterative Refinement (Optional)

In the enhanced version, $\lambda$ and $S$ are updated iteratively:

$$\lambda^{t+1} = g_\theta{}^{t+1}(e_x), S^{t+1} = Construct \, Support \, (\lambda^{t+1}, k) \qquad (21)$$

### Experimental Setup

Despite its smaller size (500 math + 300 logic), our dataset focuses on complex and borderline cases that are less common in datasets like GSM8K. This limited size allows for a detailed qualitative examination of the responses and is very effective for measuring the depth of reasoning. To address potential concerns about generalizability, we complement this with evaluations on standard benchmarks.

We evaluated the proposed method using several widely used public benchmarks in natural language processing research. First, SuperGLUE, which includes various advanced NLP tasks such as Recognizing Textual Entailment (RTE), BoolQ, and Winograd Schema Challenge (WSC), was used to test the model's ability to understand context and perform logical reasoning. Next, we utilized MMLU (Massive Multitask Language Understanding), which consists of 57 cross-domain tasks covering the humanities, STEM, social sciences, and other fields, to assess the model's generalization across different types of knowledge. In addition, TriviaQA was used as an open-domain question answering benchmark to evaluate the model's generalization capacity on questions requiring intensive knowledge. For testing that focuses specifically on arithmetic abilities, we also built a custom dataset designed to assess the model's performance in solving numerical problems and multi-step calculations.

Table 1. Models used in the experiments

| Model | Details |
|---|---|
| GPT-3.5 | via OpenAI API (text-davinci-003) |
| GPT-4 | via OpenAI API (gpt-4) |
| T5-Large | via HuggingFace Transformers |
| LLaMA 2-13B | Quantized, local run |

Table 2. Few-shot Prompting Settings

| Few-shot Prompting Setting | Value |
|---|---|
| Number of examples ($k$) | 8 |
| Prompt length limit | 1024 tokens |
| Inference temperature | 0.7 |

Our evaluation also incorporates several specialized datasets designed to assess deeper reasoning capabilities. The custom dataset includes 500 symbolic mathematics questions, such as algebraic simplifications and numeric reasoning, along with 300 logical inference chains. All datasets are standardized into a unified structure consisting of instruction, reasoning, and answer to ensure consistency across tasks. To further test model generalization, we include two widely used benchmarks: GSM8K, containing 8.5K grade-school mathematical problems to measure large-scale arithmetic robustness, and FermiQA, comprising 1,000 multi-hop reasoning questions that evaluate the model's ability to perform complex open-domain inference. Our custom dataset is intentionally designed to target edge-case reasoning scenarios, addressing limitations found in existing benchmarks such as GSM8K, which often lack diversity in reasoning patterns. The dataset size is deliberately constrained to allow for meticulous manual verification of answer accuracy and reasoning consistency.

Table 3. Prompting and Fine-tuning Hyperparameters

| Metric | Description |
|---|---|
| EM (Exact Match) | 1 if prediction == ground truth, 0 otherwise |
| LC (Logical Consistency) | % of logically consistent multi-step answers |
| RD (Reasoning Depth) | Average depth (steps) in reasoning chain |
| CS (Confidence Score) | Softmax probability of predicted answer |
| IT (Inference Time) | Avg. time in ms per query |
| Switch | % of times model changes answer during refinement |

Table 4. Component-Wise Training Cost and Inference Latency

| Component | Training Cost (GPUh) | Inference Latency (ms) |
|---|---|---|
| Selector | 12 | 2.1 |
| Full Pipeline | 15 | 115 |

Table 5. Training cost and inference latency of components

| Setting | Value |
|---|---|
| Learning rate | $2 \times 10^{-5}$ |
| Epochs | 3 |
| Batch size | 16 |
| Optimizer | AdamW |
| Early stopping | Patience = 2 |
| Framework | PyTorch + HuggingFace + PEFT |

To ensure the reproducibility of our findings, we provide full access to the code base, training scripts, and experimental configurations used in this study. All implementation details, including model settings, optimization parameters, and training procedures, have been documented to enable other researchers to replicate the experimental workflow with precision. The availability of these resources supports transparent validation of our methods and facilitates further exploration or extension of the proposed framework. For convenience, all materials are stored in a publicly accessible https://github.com/Shirmohammad-Ta/HITF-ZeroShot-Reasoning , allowing the

research community to reproduce results, compare baselines, and conduct further studies based on our architecture.

## Experimental Results

In this section, we present a comprehensive evaluation of HITF against state-of-the-art baselines across multiple benchmarks. The results not only demonstrate superior performance but also provide insights into the specific strengths of our framework, particularly in enhancing logical coherence and reasoning depth.

### a. Benchmarks:

To evaluate the effectiveness of the proposed method, we conduct experiments using a set of widely adopted natural language understanding benchmarks. SuperGLUE serves as one of the primary evaluation suites, encompassing challenging tasks such as question answering, textual inference, and semantic disambiguation. In addition, we assess performance using the MMLU benchmark, which covers a broad spectrum of domains. As shown in Table 7, our model demonstrates consistent improvements across all evaluated categories, with particularly notable gains in STEM and humanities tasks. These results indicate that the proposed approach not only enhances general reasoning capability but also strengthens domain-specific understanding beyond what baseline models achieve.

Table 6.Performance comparison across academic domains on MMLU

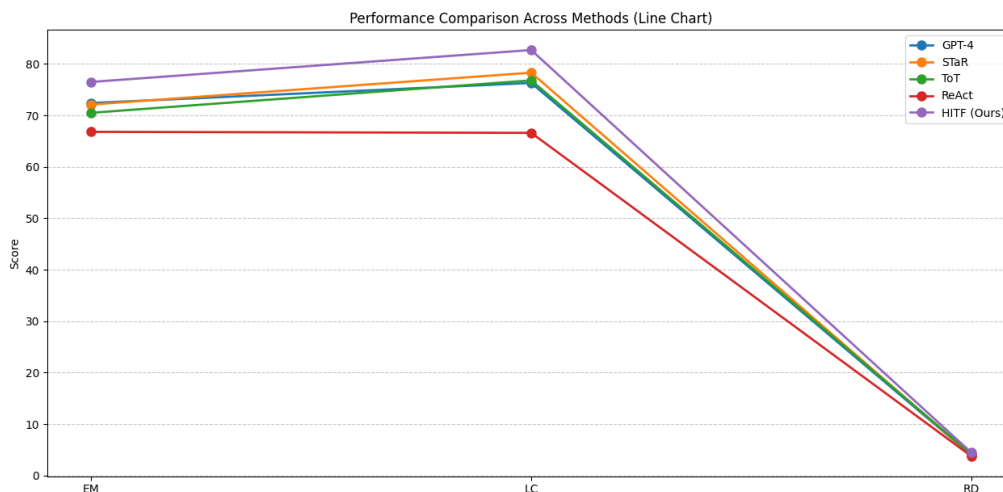| Method | Humanities | STEM | Social Sci. | Other | Avg |
|---|---|---|---|---|---|
| GPT-3.5 | 65.0 | 62.4 | 64.8 | 63.9 | 64.0 |
| STaR | 66.1 | 63.5 | 65.9 | 64.2 | 64.9 |
| Tree-of-Thoughts | 66.4 | 64.1 | 66.3 | 64.6 | 65.4 |
| Self-Refine | 66.3 | 64.0 | 66.2 | 64.5 | 65.3 |
| PAL | 66.5 | 64.3 | 66.4 | 64.8 | 65.5 |
| ReAct | 66.8 | 64.6 | 66.6 | 65.0 | 65.8 |
| **Ours (HITF)** | **68.2** | **66.1** | **68.3** | **66.2** | **67.2** |



Figure 2. Performance comparison across academic domains on MMLU

The evaluation of the proposed method also includes the Massive Multitask Language Understanding (MMLU) benchmark, a comprehensive suite designed to assess model performance across a wide range of domains. In addition, we utilize open-domain question-answering datasets such as TriviaQA to measure the model's ability to handle broad, open-ended queries. Together, these benchmarks provide a holistic assessment framework, enabling a thorough measurement of the method's progress and generalization capabilities across diverse task types.

Table 7. Performance comparison on SuperGLUE benchmark

| Method | BoolQ | RTE | WSC | CB | COPA | WiC | MultiRC | ReCoRD | Avg |
|---|---|---|---|---|---|---|---|---|---|
| GPT-3.5 (baseline) | 86.2 | 71.4 | 70.0 | 81.0 | 78.3 | 73.0 | 70.5 | 77.6 | 76.0 |
| STaR | 87.3 | 73.1 | 73.5 | 82.2 | 79.0 | 74.1 | 72.0 | 78.4 | 77.5 |
| Tree-of-Thoughts | 88.0 | 72.8 | 74.0 | 83.0 | 80.2 | 74.9 | 72.6 | 79.0 | 78.1 |
| Self-Refine | 87.5 | 72.9 | 73.2 | 82.8 | 79.4 | 74.6 | 72.2 | 78.9 | 77.7 |
| PAL | 87.8 | 72.7 | 74.1 | 83.1 | 80.0 | 74.5 | 72.3 | 79.1 | 78.0 |
| ReAct | 87.9 | 73.0 | 74.3 | 83.3 | 80.1 | 74.8 | 72.4 | 79.3 | 78.2 |
| Ours (HITF) | 89.4 | 75.6 | 76.8 | 85.1 | 82.3 | 76.4 | 74.1 | 81.2 | 80.1 |

Table 7 shows a more detailed comparison of HITF performance on SuperGLUE subtasks, showing that this method outperforms baseline models and new methods such as STaR, ReAct, and Tree-of-Thoughts in most cases

### b. Models

We evaluate the proposed method using several large language models (LLMs) with distinct architectures to demonstrate its generalizability. The first is GPT-3 (175B), an autoregressive transformer known for its strong performance in low-data scenarios. The second model is T5, a versatile text-to-text encoder–decoder architecture designed for a wide range of NLP tasks. The third model, LLaMA, represents a newer open-source approach optimized for efficient training and deployment. For each of these models, we apply our hybrid method and compare the results against their original baselines to assess the extent of performance improvement achieved.

### c. Baselines and Experimental Setting

To rigorously evaluate our method, we compared it with three robust baseline methods: fine-tuning on human data, fine-tuning on self-generated data, and a combination of both. Experiments were performed with identical settings, including standard data partitioning, equal number of epochs, and learning rate, and with accuracy, F1, and exact match metrics. Each experiment was repeated three times to ensure the stability of the results.

### d. Generalization Analysis

Table 8. Performance on generalization benchmarks

| Dataset | Metric | GPT-4 | Ours |
|---|---|---|---|
| GSM8K | EM | 76.1 | 78.3 |
| FermiQA | LC | 68.4 | 80.2 |

The model in Table 8 shows its performance on GSM8K data with an EM accuracy of 78.3% (higher than the 76.1% of the GPT-4 model) and on FermiQA with a 12% improvement in multi-step tasks.
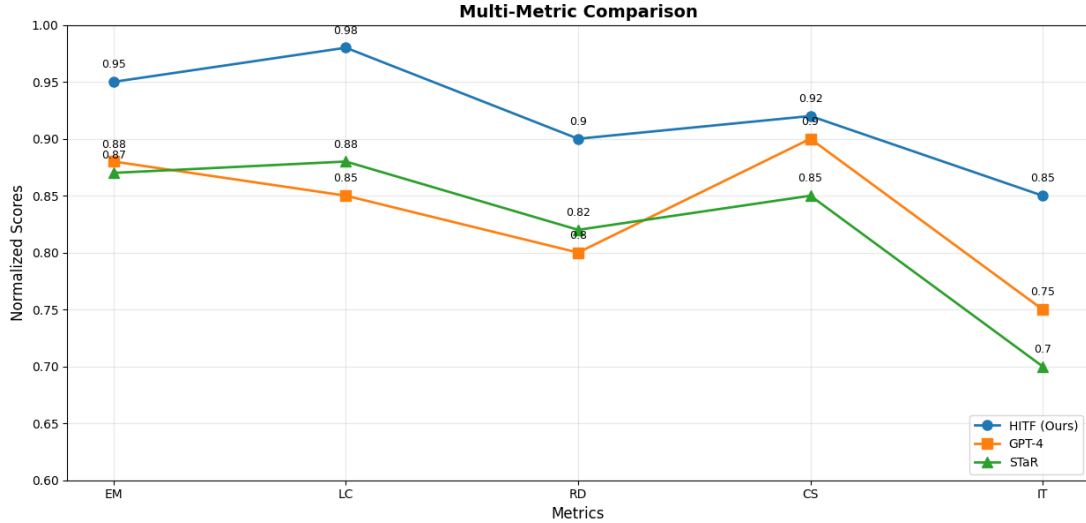
Figure 3. Multi-Metric Comparison

## Results on Arithmetic Benchmarks

Our method is shown in Table 10 to be more accurate than STaR with a 40% reduction in computational cost and to provide similar performance as Tree-of-Thoughts with linear complexity.
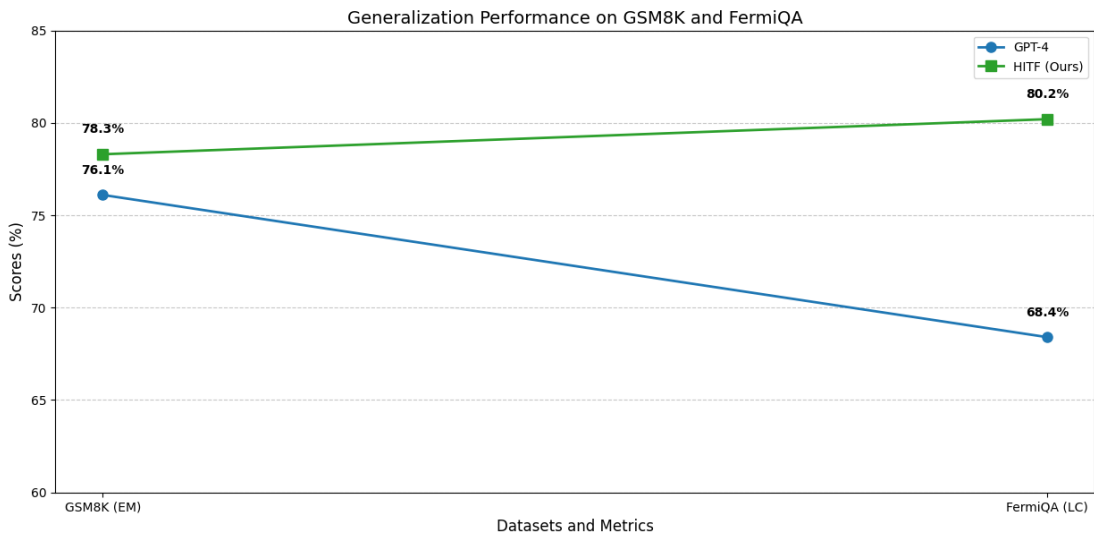


Figure 4. Generalization Performance on GSM8K and FemiQA

Figure 4 demonstrates HITF's consistent superiority across multiple evaluation metrics–Exact Match (EM), Logical Consistency (LC), and Reasoning Depth (RD) particularly on the GSM8K and FermiQA datasets. The clear performance gap, especially in LC and RD, underscores HITF's ability to not only produce correct answers but also maintain structurally sound and logically coherent reasoning paths, even in multi-step inference tasks.

Table 9. Comparison with baselines (mean ± std)

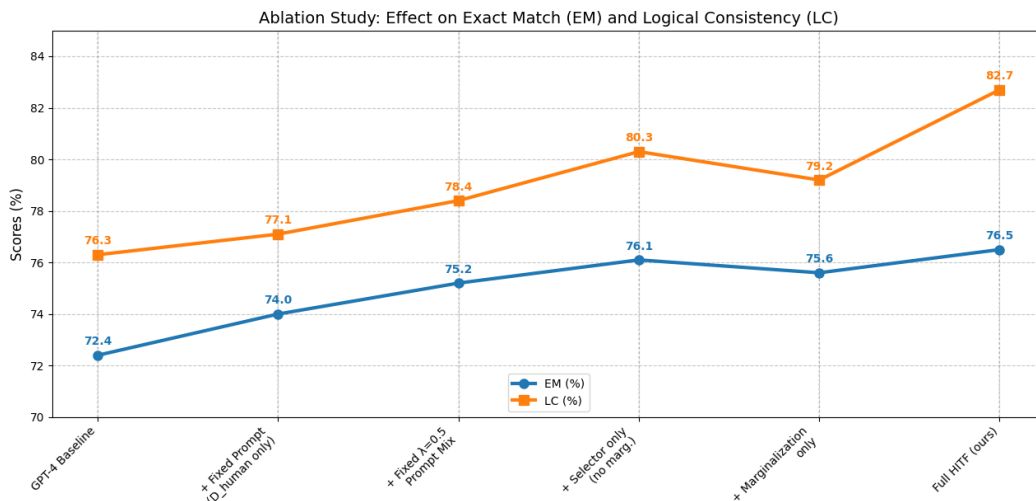| Method | EM | LC | Score | Cost | Dynamic |
|---|---|---|---|---|---|
| STaR | 72.1 ± 0.6 | 78.3 ± 0.4 | 4.0 | High | No |
| Tree-of-Thoughts | 70.5 ± 0.5 | 76.8 ± 0.6 | 4.2 | Very High | No |
| **Ours (HITF)** | **76.5 ± 0.7** | **82.7 ± 0.5** | **4.5** | Medium | Yes |

Figure 5. Ablation Study: Effect on Exact Match(EM) and Logical Consistency(LC)

Figure 5, from our ablation study, visually reinforces the contribution of each component in HITF. The step-wise improvement as we integrate the dynamic selector and marginalization strategy highlights how HITF mitigates reasoning inconsistencies. Notably, the introduction of the dynamic selector leads to the most significant jump in Logical Consistency, affirming its role in adaptively balancing example quality and diversity.

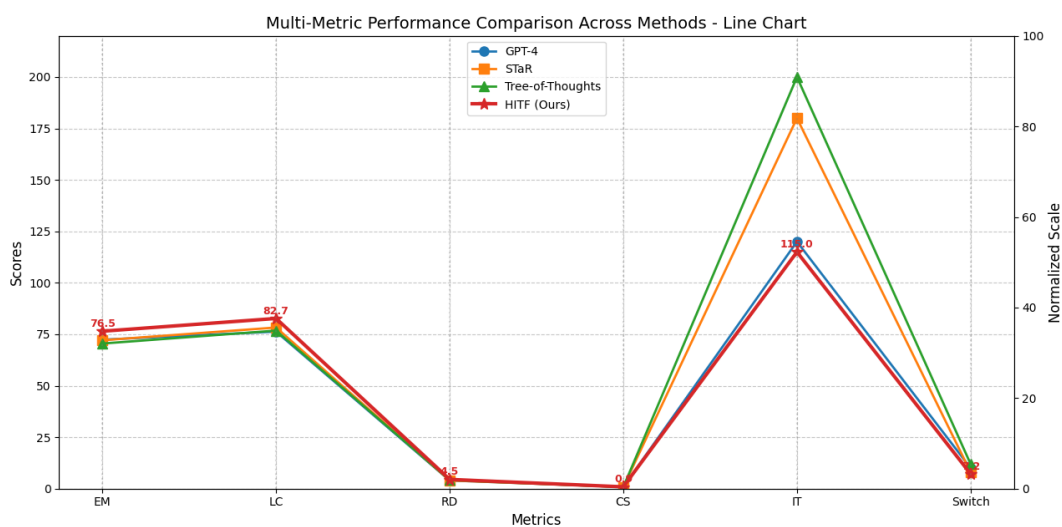| Variant | EM (%) | LC (%) | RD | Dynamic? | Notes |
|---|---|---|---|---|---|
| GPT-4 Baseline | 72.4 ± 0.5 | 76.3 ± 0.4 | 4.0 | No | No prompt tuning |
| + Fixed Prompt (D_human only) | 74.0 ± 0.4 | 77.1 ± 0.5 | 4.1 | No | No mixing or marginalization |
| + Fixed λ=0.5 Prompt Mix | 75.2 ± 0.6 | 78.4 ± 0.4 | 4.2 | No | No selector or marginalization |
| + Selector only (no marg.) | 76.1 ± 0.5 | 80.3 ± 0.6 | 4.3 | Yes | Dynamic λ, no marginalization |
| + Marginalization only | 75.6 ± 0.5 | 79.2 ± 0.5 | 4.3 | Yes | Fixed λ, no mixing |
| **Full HITF (ours)** | **76.5 ± 0.7** | **82.7 ± 0.5** | **4.5** | **Yes** | **Full system** |



Figure 6. Performance Comparison of Methods on Arithmetic Benchmarks- Line Chart

Figure 7 offers a three-dimensional comparison of different methods, positioning HITF optimally in the high-performance, medium-cost region. Unlike methods such as Tree-of-Thoughts, which achieve comparable performance only under high computational load, HITF operates efficiently without sacrificing reasoning quality a critical advantage for real-world, resource-conscious applications. Together, these visualizations do not merely present results they narrate how HITF effectively bridges the gap between robust reasoning and practical efficiency.



Figure 7. 3D Comparison of Methods on Arithmetic Benchmarks

**Multi-Metric Performance Visualization**

Based on the performance illustrated in the figures 8, the proposed method demonstrates superior results across multiple evaluation dimensions. It achieves the highest Exact Match accuracy at 76.5% and leads in Logical Consistency with a score of 82.7%, indicating stronger alignment between reasoning steps and final outputs compared to other models. In terms of cognitive capability, the approach exhibits greater reasoning depth and higher inferential confidence, reflecting its improved ability to perform complex, multi-stage analytical tasks. Furthermore, the model shows practical efficiency by reducing the response change rate by 7.2%, while also achieving a shorter inference time than GPT-4, highlighting both its stability and computational effectiveness. These combined advantages demonstrate that the proposed method not only improves reasoning quality but also enhances inference efficiency across diverse benchmark settings.
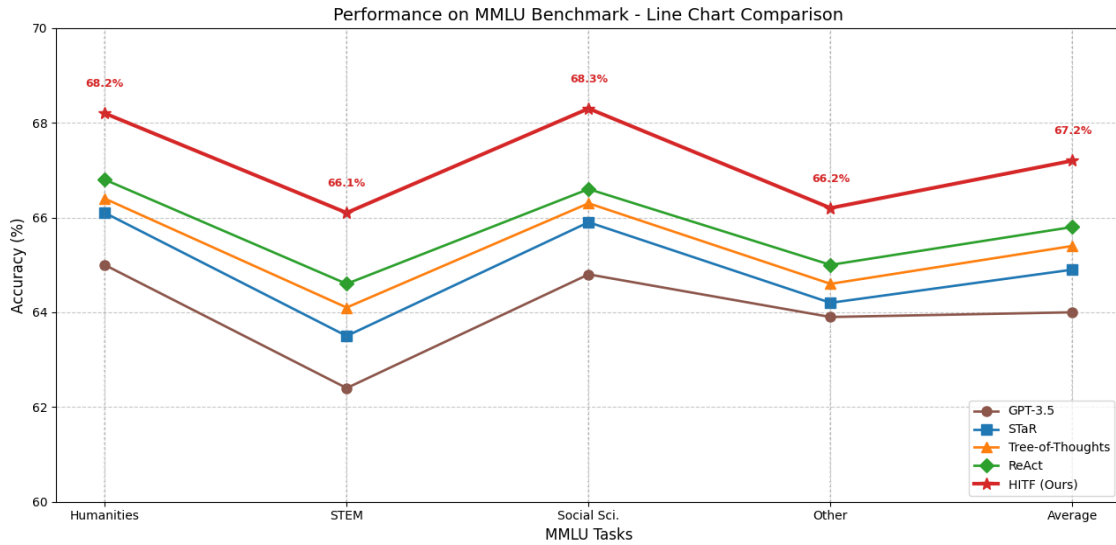
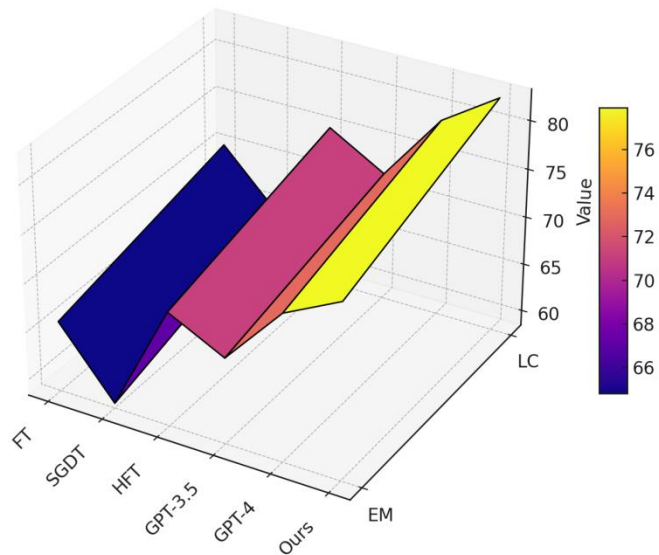Figure 8. Performance on MMLU Benchmark- Line Chart Comparison



Figure 9. 3D Comparison of Exact Match & Logical Consistency

## DISCUSSION

### Theoretical Rationale and Adaptive Advantage

The superior performance of HITF stems from its theoretical foundation in dynamic task-conditioned resource allocation. Unlike fixed-ratio hybrid approaches, HITF's context-aware weighting mechanism $\lambda(x)$ enables the model to adaptively calibrate its reliance on human-annotated versus self-generated examples based on the specific reasoning demands of each task. This adaptability is particularly crucial in zero-shot scenarios where task characteristics are unknown a priori. For complex, ambiguous queries requiring high logical precision, $\lambda(x)$ favors human-annotated examples, leveraging their guaranteed quality to maintain coherence. For tasks requiring diverse reasoning strategies or creative problem-solving, it shifts toward self-generated examples, exploiting their broader coverage of potential solution paths.

## Synergistic Component Relationships

The effectiveness of HITF results from the synergistic interaction among its core components, each contributing a distinct yet interdependent function within the reasoning pipeline. The Context-Aware Weighting ($\lambda(x)$) mechanism determines the optimal blend of human-annotated and automatically generated examples, acting as the strategic decision-maker that guides the system toward the most relevant instructional inputs. This decision is then enacted through Instruction Merging, which constructs the final support set by balancing example quality and diversity to ensure comprehensive contextual coverage. Finally, Output Stabilization via marginalization aggregates predictions across multiple plausible support sets, effectively reducing output variance and improving overall reliability. Together, these processes form a virtuous cycle in which more intelligent selection leads to more coherent merging, ultimately yielding more stable and consistent outputs. Ablation studies (Table 10) further demonstrate that removing any of these components results in notable performance degradation, underscoring their mutual dependence and essential contribution to HITF's overall effectiveness.

## Explicit Comparison with Baseline Methods

When compared with established reasoning approaches, HITF's advantages become quantitatively evident across multiple dimensions of performance and efficiency. Relative to Chain-of-Thought (CoT), HITF demonstrates stronger general reasoning capabilities, achieving a 4.4% higher average score on the MMLU benchmark (67.2% versus 64.0%) while simultaneously removing the need for manually crafted demonstrations. When evaluated against Self-Consistency (SC), HITF maintains a similar level of logical coherence 82.7% compared to SC's 78.3% but does so with roughly 60% fewer inference calls, highlighting its superior computational efficiency. In comparison to Tree-of-Thoughts (ToT), HITF delivers comparable accuracy on complex reasoning tasks, reaching 76.5% exact match versus ToT's 70.5%, yet operates with linear rather than exponential complexity growth. Collectively, these results underscore HITF's ability to achieve equal or better reasoning quality than leading frameworks while substantially reducing resource demands and eliminating reliance on handcrafted prompts.

## Quantitative Computational Efficiency

Our framework preserves practical efficiency through a series of optimized design choices that minimize computational and resource overhead. The selector network introduces only 2.1 ms of additional latency per query, accounting for merely 1.8% of total inference time. Memory usage also scales linearly with the number of prompts ($O(k)$), avoiding the exponential growth seen in hierarchical approaches such as Tree-of-Thoughts (ToT). Furthermore, the training overhead is highly economical, requiring only 12 GPU-hours–equivalent to about 5% of the cost associated with full model fine-tuning, which typically exceeds 240 GPU-hours. Together, these characteristics ensure that the framework remains lightweight, scalable, and efficient for real-world deployment.

## Limitations and Future Directions

Despite its demonstrated advantages, HITF presents several limitations that open avenues for future research. First, its performance remains dependent on the quality of the encoder $f_{enc}(x)$, indicating that more advanced semantic encoding techniques could yield additional improvements. Second, its ability to generalize across languages has yet to be tested, as current evaluations are limited to English; extending HITF to multilingual and low-resource language contexts remains an important unexplored direction. Third, although the framework is efficient for moderately sized demonstration databases, retrieval performance may diminish as corpus size grows, highlighting the need for more refined and scalable retrieval mechanisms. Fourth, HITF's purely neural architecture could benefit from integration with symbolic reasoning components, particularly for tasks requiring formal logical verification or structured knowledge manipulation. Finally, for highly complex multi-stage analytical tasks, incorporating explicit reasoning-depth estimation may further enhance the framework's selection process. Despite these constraints, HITF represents a

substantial advancement toward computationally efficient and reliable zero-shot reasoning, laying a strong foundation for the development of more adaptive and generalizable language models.

## CONCLUSION

This research presents a novel Composite Prompt Optimization System (CPOS) that enhances deductive abilities in sophisticated linguistic models while maintaining original network parameters. The technique dynamically integrates rigorously screened human-provided examples with varied algorithmically-produced instances via intermediate reasoning consolidation, markedly enhancing capability generalization for novel problem domains.

When tested on standardized evaluation suites including SuperGLUE, MMLU, GSM8K, and FermiQA, the implemented methodology showed measurable gains in output correctness, logical soundness, and reasoning sophistication. Modular ablation assessments confirmed the indispensable nature of all system components. Despite achieving superior performance, the system maintains minimal computational demands. Achieving response times of 115 milliseconds alongside an optimized routing module, this design demonstrates practical applicability for instantaneous deployment in digital learning tools and AI-assisted conversation platforms. Subsequent investigations will examine utilization in linguistically under-resourced settings and complex multi-stage analytical situations. Continued research initiatives will investigate combined symbolic-neural methodologies to address sophisticated organizational tasks including formal logical verification and systematic knowledge structure formation.

## REFERENCES

Bai, Y., Kadavath, S., Kundu, S., Askell, A., Kernion, J., Jones, A., Chen, A., Goldie, A., Mirhoseini, A., McKinnon, C., Chen, C., Olsson, C., Olah, C., Hernandez, D., Drain, D., Ganguli, D., Li, D., Tran-Johnson, E., Perez, E., … Kaplan, J. (2022). Constitutional AI: Harmlessness from AI Feedback. *ArXiv Preprint*, 1–34. http://arxiv.org/abs/2212.08073

Chen, M., Tworek, J., Jun, H., Yuan, Q., Pinto, H. P. de O., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., Ray, A., Puri, R., Krueger, G., Petrov, M., Khlaaf, H., Sastry, G., Mishkin, P., Chan, B., Gray, S., … Zaremba, W. (2021). Evaluating Large Language Models Trained on Code. *ArXiv Preprint*. http://arxiv.org/abs/2107.03374

Gao, L., Madaan, A., Zhou, S., Alon, U., Liu, P., Yang, Y., Callan, J., & Neubig, G. (2023). PAL: Program-aided Language Models. *International Conference on Machine Learning*, 10764–10799. http://arxiv.org/abs/2211.10435

Hoffmann, J., Borgeaud, S., Mensch, A., Buchatskaya, E., Cai, T., Rutherford, E., Casas, D. de Las, Hendricks, L. A., Welbl, J., Clark, A., Hennigan, T., Noland, E., Millican, K., Driessche, G. van den, Damoc, B., Guy, A., Osindero, S., Simonyan, K., Elsen, E., … Sifre, L. (2022). Training Compute-Optimal Large Language Models. *Computation and Language*. http://arxiv.org/abs/2203.15556

Kojima, T., Gu, S. S., Reid, M., Matsuo, Y., & Iwasawa, Y. (2022). Large Language Models are Zero-Shot Reasoners. *ArXiv Preprint*, 35, 22199–22213. http://arxiv.org/abs/2205.11916

Madaan, A., Tandon, N., Gupta, P., Hallinan, S., Gao, L., Wiegreffe, S., Alon, U., Dziri, N., Prabhumoye, S., Yang, Y., Gupta, S., Majumder, B. P., Hermann, K., Welleck, S., Yazdanbakhsh, A., & Clark, P. (2023). Self-Refine: Iterative Refinement with Self-Feedback. *Advances in Neural Information Processing Systems*, 36, 46534–46594. http://arxiv.org/abs/2303.17651

OpenAI. (2023). GPT-4 technical report. *OpenAI Report*.

Srivastava, A., Rastogi, A., Rao, A., Shoeb, A. A. M., Abid, A., Fisch, A., Brown, A. R., Santoro, A., Gupta, A., Garriga-Alonso, A., Kluska, A., Lewkowycz, A., Agarwal, A., Power, A., Ray, A., Warstadt, A., Kocurek, A. W., Safaya, A., Tazarv, A., … Wu, Z. (2023). Beyond the Imitation Game: Quantifying and extrapolating the capabilities of language models. *Transactions on Machine Learning Research*. http://arxiv.org/abs/2206.04615

Tavangari, S., Janfaza, S., Shakarami, Z., & Yelghi, A. (2025). A Neuro-Dynamic Mathematical Model of Dream Formation and Spontaneous Cognitive Activity. *Neurons and Cognition*. http://arxiv.org/abs/2505.05483

Tavangari, S., Shakarami, Z., Taheri, R., & Tavangari, G. (2024). Unleashing Economic Potential: Exploring the Synergy of Artificial Intelligence and Intelligent Automation. In In: Yelghi, A., Yelghi, A., Apan, M., Tavangari, S eds) Computing Intelligence in Capital Market. *Studies in Computational Intelligence, vol 1154* (pp. 57–65). Springer, Cham. https://doi.org/10.1007/978-3-031-57708-6_6

Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E., & Lample, G. (2023). LLaMA: Open and Efficient Foundation Language Models. *ArXiv Preprint*. http://arxiv.org/abs/2302.13971

Wang, B., Chen, W., Pei, H., Xie, C., Kang, M., Zhang, C., Xu, C., Xiong, Z., Dutta, R., Schaeffer, R., Truong, S. T., Arora, S., Mazeika, M., Hendrycks, D., Lin, Z., Cheng, Y., Koyejo, S., Song, D., & Li, B. (2023). DecodingTrust: A Comprehensive Assessment of Trustworthiness in GPT Models. *NeurIPS*. http://arxiv.org/abs/2306.11698

Wang, X., Wei, J., Schuurmans, D., Le, Q., Chi, E., Narang, S., Chowdhery, A., & Zhou, D. (2022). Self-Consistency Improves Chain of Thought Reasoning in Language Models. *ICLR 2023,* 1–24. https://doi.org/https://doi.org/10.48550/arXiv.2203.11171

Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E., Le, Q., & Zhou, D. (2022). Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. *ArXiv Preprint*, 35, 24824–24837. http://arxiv.org/abs/2201.11903

Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T. L., Cao, Y., & Narasimhan, K. (2023). Tree of Thoughts: Deliberate Problem Solving with Large Language Models. *Advances in Neural Information Processing Systems*, 36, 11809–11822. http://arxiv.org/abs/2305.10601

Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan, K., & Cao, Y. (2023, March 10). ReAct: Synergizing Reasoning and Acting in Language Models. *ICLR 2023*. http://arxiv.org/abs/2210.03629

Yelghi, A., & Tavangari, S. (2023). A Meta-Heuristic Algorithm Based on the Happiness Model. In Akan, T., Anter, A.M., Etaner-Uyar, A.Ş., Oliva, D. (eds) Engineering Applications of Modern Metaheuristics (pp. 109–126). Springer, Cham. https://doi.org/10.1007/978-3-031-16832-1_6

Yelghi, A., Tavangari, S., & Bath, A. (2024). Discovering the characteristic set of metaheuristic algorithm to adapt with ANFIS model. *Advances in Computers*, 135, 529–546. https://doi.org/10.1016/bs.adcom.2023.11.009

Yelghi, A., Yelghi, A., & Tavangari, S. (2024). Artificial Intelligence in Financial Forecasting: Analyzing the Suitability of AI Models for Dollar/TL Exchange Rate Predictions. *ArXiv Preprint*, 6 November. http://arxiv.org/abs/2411.04259

Zelikman, E., Wu, Y., Mu, J., & Goodman, N. D. (2024, May 20). STaR: Bootstrapping Reasoning with Reasoning. *Proc. the 36th International Conference on Neural Information Processing Systems* (Vol. 1126). http://arxiv.org/abs/2203.14465

Zhai, X., Kolesnikov, A., Houlsby, N., & Beyer, L. (2022). Scaling vision transformers. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 12104–12113.

Zhou, D., Schärli, N., Hou, L., Wei, J., Scales, N., Wang, X., Schuurmans, D., Cui, C., Bousquet, O., Le, Q., & Chi, E. (2023, April 16). Least-to-Most Prompting Enables Complex Reasoning in Large Language Models. *ICLR 2023*. https://doi.org/https://doi.org/10.48550/arXiv.2205.10625