# FARM: Field-Aware Resolution Model for Intelligent Trigger-Action Automation

KHUSRAV BADALOV, Neouly Co., Ltd., Republic of Korea

YOUNG YOON, Hongik University, Republic of Korea

Trigger-Action Programming (TAP) platforms such as IFTTT and Zapier enable Web of Things (WoT) automation by composing event-driven rules across heterogeneous services. A TAP applet links a *trigger* to an *action* and must bind trigger outputs (*ingredients*) to action inputs (*fields*) to be executable. Prior work largely treats TAP as service-level prediction from natural language, which often yields non-executable applets that still require manual configuration. We study the *function-level configuration* problem: generating complete applets with correct ingredient-to-field bindings. We propose FARM (Field-Aware Resolution Model), a two-stage architecture for automated applet generation with full configuration. Stage 1 trains contrastive dual encoders with selective layer freezing over schema-enriched representations, retrieving candidates from 1,724 trigger functions and 1,287 action functions (2.2M possible trigger-action pairs). Stage 2 performs selection and configuration using an LLM-based multi-agent pipeline. It includes intent analysis, trigger selection, action selection via cross-schema scoring, and configuration verification. Agents coordinate through shared state and agreement-based selection. FARM achieves 81% joint accuracy on Gold (62% Noisy, 70% One-shot) at the function level, where both trigger and action *functions* must match the ground truth. For comparison with service-level baselines, we map functions to their parent services and evaluate at the service level. FARM reaches 79% joint accuracy and improves over TARGE by 21 percentage points. FARM also generates ingredient-to-field bindings, producing executable automation configurations.

CCS Concepts: • **Computing methodologies** → *Rule learning*; **Knowledge representation and reasoning**; **Supervised learning**; **Multi-agent planning**; • **Information systems** → **Similarity measures**.

Additional Key Words and Phrases: Trigger-Action Programming, Natural Language Processing, Internet of Things (IoT), Web of Things (WoT), Multi-Agentic AI, Retrieval Augmentation Generation (RAG)

## 1 Introduction

Trigger-Action Programming (TAP) is a declarative programming paradigm that allows users to define event-driven automation rules. It is the essential foundation behind systems such as IFTTT (If This Then That), Zapier, and similar platforms including Microsoft Power Automate and Apple Shortcuts. Ur et al. [45] demonstrated the practicality of TAP in smart home contexts, while Rahmati et al. [38] provided a comparative analysis of IFTTT and Zapier. TAP has become a mainstream interface for end-user automation across Web of Things (WoT) ecosystems, connecting both IoT devices and web applications; Mi et al. [34] found that over 400 services spanning both categories are integrated on IFTTT alone. Ur et al. [44] showed that real-world TAP rule corpora exhibit substantial diversity in services, rule structures, and user-written natural language descriptions, making TAP a challenging target for robust intent understanding and program synthesis.

Historically, TAP interaction was associated with composing simple IF-THEN rules by manually selecting a trigger and an action and combining them into a mashup [15, 38]. Recently, LLM-driven assistants and agentic systems have

---

Authors' Contact Information: Khusrav Badalov, Neouly Co., Ltd., Seoul, Republic of Korea , khusravvvb99@gmail.com; Young Yoon, Hongik University, Computer Engineering, Seoul, Republic of Korea, young.yoon@hongik.ac.kr.

shifted this workflow toward natural language specifications, where a model is expected to (i) infer user intent, (ii) select services, and (iii) produce executable configurations [22, 42]. Figure 1 illustrates this paradigm shift, contrasting traditional manual configuration with agentic AI-driven mashup generation. However, in practice, AI-based TAP builders frequently fail on underspecified or ambiguous requests: users often phrase the same intent in multiple ways, provide incomplete context (e.g., location, device identity, thresholds), or implicitly encode preferences (e.g., comfort vs. energy saving) [21, 27]. In such cases, current systems may generate multiple inconsistent trigger-action variants, select incorrect parameters, or produce brittle applets that are syntactically valid but semantically misaligned with user goals.
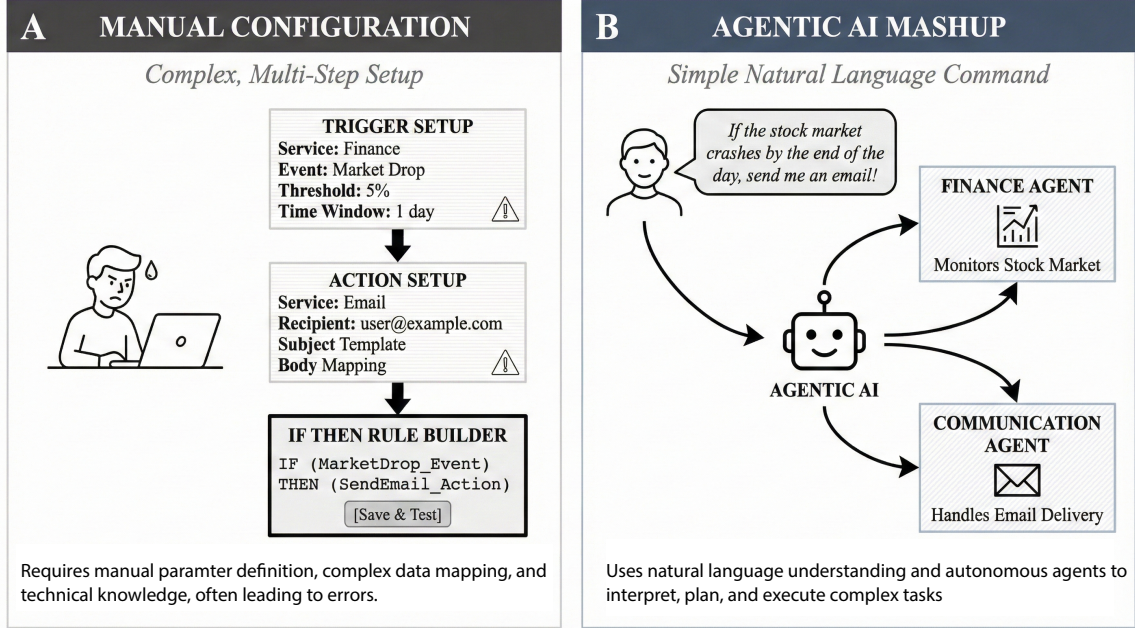


Fig. 1. An example of an TAP applet interface for manual configuration and our proposal for an automated applet generation.

## 1.1 Motivation

Consider a common scenario: Alice asks an AI assistant, "Notify me when it's going to rain tomorrow." Although seemingly straightforward, existing TAP platforms frequently mishandle such requests: alerts may trigger for the wrong location, apply an unintended threshold, or fail due to missing implicit parameters. This example highlights a central challenge in contemporary TAP systems: the widening semantic gap between informal user language and the structured, parameter-rich function-level invocations required for reliable automation. We use the term *function-level* to distinguish from service-level identification: while a service (e.g., "Google Sheets") may expose multiple triggers and actions, a function refers to a specific trigger (e.g., "New row added to spreadsheet") or action (e.g., "Update cell in spreadsheet"), each with its own input parameters, output data fields, and schema constraints. Prior work also shows that users struggle to debug such failures, because the root cause may lie in hidden parameters, platform semantics, or multi-rule interactions Zhang et al. [63].

Recent advances in TAP generation have made significant strides. Yusuf et al.[61] demonstrated that transformer seq2seq models can generate TAPs with improved service and field accuracy from natural language descriptions. Liu et al. [28] learning-based approaches, including attention and latent alignment formulations, framed TAP inference as service/action prediction under weak supervision . However, these state-of-the-art works largely focus on predicting the correct services and fields, and do not explicitly ensure that generated configurations are aligned with the user's higher-level intent or preferences.

Conversely, Wu et al. [55] pioneered implicit intention prediction through multi-view representation learning in their MvTAP framework, showing that users often create rules driven by latent goals (e.g., "energy saving" when turning off devices). Yet MvTAP addresses intention classification rather than the end-to-end generation of executable applet configurations with function-level parameters. In parallel, Hsu et al. [13] showed that even seemingly correct TAP rules can yield undesirable outcomes due to cross-rule interactions and hidden automation chains, reinforcing that correctness must be defined at the level of executable behavior, not only service accuracy.

This dichotomy between generation accuracy and intention understanding represents a critical gap. A system can generate syntactically valid TAPs while still missing the user's actual intent. On the other hand, a system can correctly infer, for instance, "energy saving", without being able to translate that goal into concrete trigger conditions and action parameters. In addition, the move toward agentic LLM tool use introduces a further constraint: systems must reliably decide *when* to call tools, *which* tools to call, and *how* to parameterize calls, under uncertainty in the user request Schick et al.[40].

To clearly illustrate these challenges, we performed a systematic study on IFTTT and identified recurrent failure patterns that neither intent-only nor generation-only approaches fully address:

- **Understanding the Core of Services:** Services vary widely in functional richness. Popular Web services such as Google Drive, X (formerly Twitter), and Spotify expose many triggers and actions, while WoT device services such as smart thermostats or light sensors are often sparse due to hardware constraints. This induces severe imbalance in training data and increases the risk of overfitting to richly represented services [44].
- **Intention-Structure Mismatch:** Classification-based approaches can identify latent intent, and generation-based approaches can produce valid structures. However, aligning intent with structural choices (trigger conditions, action parameters, and field values) remains unreliable [55, 61].
- **Context-Dependent Field Selection:** Field names alone are insufficient. The same services may require different threshold values or options depending on intent (e.g., comfort-focused vs. energy-saving automations), and current generators are weak at selecting values that reflect user goals [63].
- **Cross-Functional Reasoning:** Similar intents can couple disparate services (e.g., weather plus messaging, occupancy plus HVAC). Existing approaches often treat service selection as independent, missing intent-driven cross-service dependencies [55].

These observations motivate our core insight: effective TAP configuration requires simultaneously understanding user intentions and generating accurate technical implementations. Traditional approaches, whether intent-centric or generation-centric, address only one side of this dual requirement.

Through experimentation and analysis of 16k+ real IFTTT applets with enriched metadata, we find that neither intention understanding nor structure generation alone suffices. The remaining challenge is architectural: bridging intent comprehension with executable configuration generation. This motivates our multi-agent approach that decomposes the problem and coordinates complementary capabilities.

## 1.2   Research Questions

Our investigation was guided by the following research questions:

(1) **How can we handle severe service-level imbalance across TAP domains, where some platforms expose dozens of triggers and actions while others provide only a few?**
    The TAP ecosystem exhibits extreme distributional skew: popular services like Google Drive, X (formerly Twitter), or Spotify have rich function-level coverage, while niche WoT devices such as Philips Hue or Nest Thermostat may expose only one or two triggers. This imbalance leads to underrepresentation of sparse function-level metadata in learned representations, causing models to favor well-represented services regardless of user intent [44]. We investigate whether training separate encoders for triggers and actions—combined with a layer freezing strategy that preserves pretrained semantic knowledge—can achieve reliable retrieval across both common and rare services.

(2) **How can a multi-agent architecture decompose the TAP generation task to produce complete, executable applet configurations rather than just service name predictions?**
    Prior approaches predict only which services to use, leaving users to manually configure fields and parameters. We investigate whether separating the problem into specialized agents—one for understanding trigger events, another for reasoning about action requirements, and a third for verifying compatibility—can generate complete configurations including field bindings that map trigger outputs to action inputs.

## 1.3   Contributions

To summarize, this work makes the following contributions:

- We introduce **FARM**[1] (Field-Aware Resolution Model), a two-stage framework that combines contrastive-trained dual encoders for high-recall candidate retrieval with multi-agent LLM-based selection for precise configuration generation.
- We propose a **layer freezing strategy** for domain-specific encoder fine-tuning that preserves pretrained semantic knowledge while adapting to trigger-action retrieval, achieving over 90% recall at rank 5 for both triggers and actions.
- We design a **four-agent selection pipeline** where specialized agents (Intent Analyzer, Trigger Selector, Action Selector, Verifier) coordinate to generate complete function-level configurations with field bindings—not just service names. These configurations include trigger input parameters, action required fields, and ingredient-to-field mappings that specify how trigger outputs populate action inputs.
- We demonstrate that FARM achieves **79% joint accuracy** on clear queries, outperforming prior approaches by +21 points over TARGE [5], while additionally generating executable applet configurations that prior methods do not produce.
- We evaluate FARM across three test conditions—clear queries (Gold data), ambiguous queries (Noisy data), and rare function-level queries (One-shot data) —demonstrating robust performance even on challenging inputs where existing methods struggle.

---

[1]Code: https://github.com/DinaHongik/FARM. Dataset available upon request: young.yoon@hongik.ac.kr

## 2 Background and Related Work

Many existence methods have been conducted in the area of trigger-action programming (TAP), and we divided most of the relevant research into two groups: 1) Rule synthesis and generation 2) Rule searching and classification.

Early efforts on automatic trigger-action program (TAP) composition treated the problem as a multi-class classification task mapping natural language (NL) descriptions to pre-defined trigger-action functions. Quirk et al. [36] first framed TAP generation as classification by training a binary logistic regression classifier for each possible "if-then" rule component. Their system extracted linguistic features (unigrams, bigrams, character trigrams) from the NL description and predicted which trigger or action functions should appear in the resulting recipe. While effective, this approach required one classifier per candidate function, which did not scale as new IoT services and functions emerged. Subsequent classification-based approaches improved efficiency and accuracy. Yoon et al. [59] has been proposed a CRF-based learning method that identifies relevant trigger services and predicts trigger-action pairs for user requests. This approach combined information retrieval with parallel learning engines, outperforming earlier single-model classifiers in both accuracy and training time. Deep learning further enhanced classification methods: the Latent Attention Model (LAM) by Liu et al. [28] trained separate neural classifiers for triggers and actions however introduced a "latent" attention mechanism (LAM) to weight words in the description by importance. LAM achieved state-of-the-art accuracy at the time, yet still treated trigger and action prediction independently, failing to capture their inter dependencies. In general, pure classification approaches are fast and straightforward but tend to struggle when user descriptions are implicit or vague, since they ignore the joint semantics of triggers and actions. To address the limitations of disjoint classification, researchers turned to generative or semantic parsing techniques that model TAP creation as a structured prediction or sequence-to-sequence generation problem. Beltagy and Quirk [1] recast TAP synthesis as constructing a parse tree (sequence of rule components) rather than independent classifications. Their system predicted one component of the recipe at a time, conditioned on the NL input and previously generated components, using an ensemble of logistic regression and a multilayer perceptron. This structured approach outperformed the earlier purely binary classification by Quirk et al. [36] and mitigated the scalability issue, although it still relied on bag-of-words text features and thus ignored deeper semantic context. More advanced generative models leverage modern neural sequence learning. Yao et al. [58] developed an interactive semantic parsing framework where a conversational agent learns to synthesize if-then rules through dialogue, employing hierarchical reinforcement learning to decide on sub-goals and rule components. This interactive approach allowed the agent to ask clarifying questions and handle complex recipes, but training it required costly simulation of user interactions. In contrast, one-shot generation models aim to produce the rule in a single pass. Yusuf et al. [61] introduced RecipeGen, a Transformer-based sequence-to-sequence model that directly "translates" a NL description into a trigger-action script. By framing TAP creation as language generation rather than classification, their approach could learn the implicit relationships between trigger and action phrases (e.g., understanding that "photo tagged with me on Facebook" implies a Facebook file URL trigger and a Dropbox upload action). To boost accuracy, RecipeGen was warm-started with a pre-trained encoder, adapting it to the IFTTT domain vocabulary. This generative model significantly outperformed the LAM classifier on multiple real-world TAP datasets, achieving higher recall and BLEU scores for correct trigger-action predictions. The success of such sequence-learning methods demonstrates that modeling the joint structure of TAPs can handle unclear or complex user requests better than independent classifiers. However, generative models must ensure the validity of produced rules (respecting available services/fields), often addressed through constrained decoding or post-validation.

A third line of work uses hybrid approaches that combine data-driven learning with domain knowledge or that recast TAP creation as a recommendation problem. Instead of parsing free-form descriptions, these methods often assume partial user input (such as chosen trigger or a high-level goal) and leverage existing rule repositories or ontologies to suggest the rest of the rule. Corno et al. [7] developed RecRules, a recommendation system that represents IoT channels and functions in a semantic graph and uses collaborative filtering on past user-created rules. Given a user's current context or usage history, RecRules performs semantic reasoning over the graph (using manually crafted ontological rules) and then suggests likely trigger-action combinations based on similarity to what other users have done. This knowledge-driven system improved the relevance of suggestions but required maintaining an expert-defined semantic model of the domain.

To incorporate explicit user intentions in automation, Corno et al. [6] proposed TAPrec, which introduced an ontology (EUPont) of goal-oriented activities to guide rule creation. In TAPrec, when a user selects a trigger, the system infers the user's high-level goal (e.g. "personalize room lighting") by mapping the trigger to an OWL class in the ontology, and then recommends an action that fulfills that goal. This approach effectively casts action selection as a goal-conditioned classification task, constrained by expert knowledge of typical IoT goals. They later extended this idea with a conversational assistant called HeyTAP [8] which asks the user questions in natural language, extracts an abstract intention from the dialogue, and matches it to predefined goal classes to recommend an appropriate if-then rule. These ontology-backed systems emphasize semantic constraints (ensuring the recommended trigger and action logically serve the same goal), but they come with the overhead of manually constructing and updating the ontology and labeling rules with intention metadata. Recent research strives to reduce this manual effort by using learning on graph-structured representations. Huang et al. [16] introduced TAP-AHGNN, an attention-based heterogeneous graph neural network that learns embeddings for TAP components (devices, triggers, actions) and uses them to recommend services. TAP-AHGNN integrates a knowledge graph of the IoT domain with GNNs, automatically capturing relationships (e.g. which devices or services are often used together) and predicting the next action given a trigger context. This hybrid model leverages structured knowledge and data-driven patterns, outperforming purely collaborative or content-based recommenders by learning complex cross-service associations (while still requiring an initial knowledge graph schema). Another emerging trend is applying large language models: King et al. [22] present Sasha, a smart-home assistant that feeds ambiguous voice commands and home context into an LLM to generate possible automation rules. Sasha showed the viability of general-purpose LLMs for TAP recommendation, but still needed each training rule to be labeled with an intention (goal) to guide the model.

## 2.1 Comparison and Limitations of Previous Works

*2.1.1 Comparison.* The evolution from classification-based to generative and hybrid approaches demonstrates a steady increase in the capability of TAP automation systems. Early classification methods (e.g., Quirk et al. [36], Yoon et al. [59], Kuang et al. [24]) were effective for known trigger–action pairs but faced challenges in scalability and capturing semantic nuances. Generative approaches (e.g., Yusuf et al. [61], Yao et al. [58], Liu et al. [29], Liu et al. [30]), ranging from structured predictors to neural sequence-to-sequence models, introduced greater flexibility by jointly modeling triggers and actions while using context from natural language descriptions. These methods achieved improved accuracy on complex user requests and unseen combinations.

*2.1.2 Limitations of Previous Work.* While the methods mentioned above improve the efficiency and accuracy of user-created TAP rules, to the best of our knowledge, none of the existing TAP research explains how a user's natural

language query can directly specify which function-level parameters and fields are required, nor how these can be used to automatically mash up two independent function—trigger and action—into a complete applet. As illustrated in Fig. 2, IFTTT applets operate across *credential boundaries* (marked with ×) that isolate each service's authentication context—the trigger service's OAuth credentials never reach the action service, and vice versa. This security model means that only validated data (ingredients) flows between stages, requiring any automated system to understand both the schema validation requirements and the field-level mappings needed to bridge these isolated services.
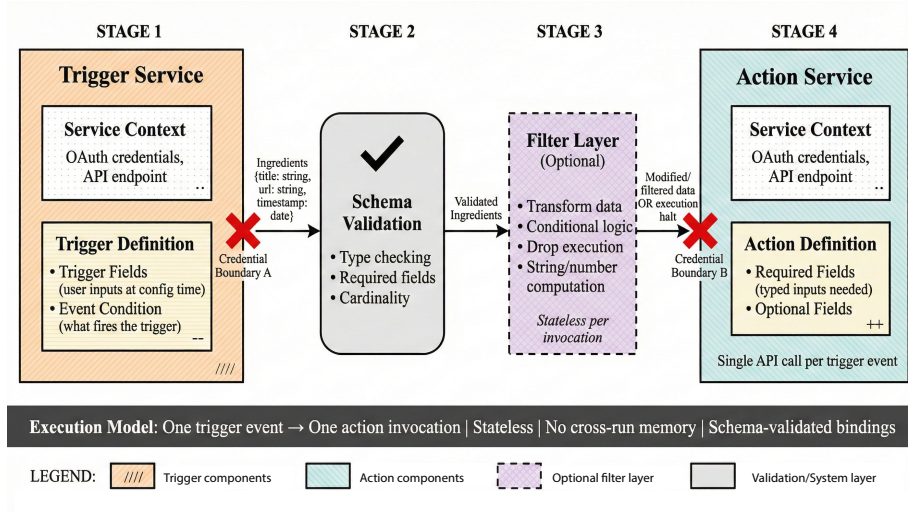


Fig. 2. IFTTT applet execution pipeline with credential boundaries.

Beyond this configuration gap, existing approaches suffer from limited expressivity in handling complex automation scenarios. The simple one-to-one trigger-action mapping requires an "impedance match" between trigger outputs and action inputs at a shared abstraction level. When deeper customization is needed—such as combining data from multiple sources or applying conditional transformations—users must write external "glue code" or resort to traditional programming methods. This limitation prevents platforms from supporting real-world scenarios involving multiple conditions, sequential triggers, or stateful computations, ultimately restricting the practical utility of automated TAP generation systems. Furthermore, current generation methods produce rules in isolation without considering inter-rule interactions or security implications. When multiple TAP rules are simultaneously enabled, complex system behaviors emerge that are difficult to diagnose. Wang et al. [49] demonstrated that 66% of synthetic IFTTT deployments exhibit potential inter-rule vulnerabilities, including action conflicts and infinite actuation loops. Existing TAP synthesis approaches focus solely on functional correctness for individual rules, lacking mechanisms to verify compatibility with a user's existing rule set or to ensure that generated configurations do not introduce unintended security risks. Table 1 provides a comprehensive comparison of existing approaches.

Table 1. Comparison of TAP Automation Approaches: Functional and Non-Functional Aspects

| Authors | Method | Input | Output | Functional Capabilities | | | Non-Functional Aspects | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | Compos. Reasoning | Iterative Refine | Service Extend. | Reported Accuracy | Dataset Scale | Data Require. |
| Quirk et al. [36] | Binary Classification | NL query | Executable AST | × (independent) | × | × (retrain) | Channel: 50% Function: 37% | 114K recipes 160 channels | High (labeled pairs) |
| Yoon et al. [59] | CRF-based Learning | NL query | Service names | × (independent) | × | × (retrain) | MA: 51.6% NE: 32.7% | 270K pages 1K-9K train | Medium-High (seq. labels) |
| Liu et al. [28] (LAM) | Neural Attention | NL query | Service names | × (independent) | × | × (retrain) | 87.5% (joint pred.) | 68K train 584 test | High (labeled pairs) |
| Beltagy & Quirk [1] | Structured Prediction | NL query | Executable JSON | Partial (sequential) | × (1-shot) | × (retrain) | Tree: 42.6% Channel: 54% | 77.5K train 5.2K dev, 4.3K test | High (labeled deriv.) |
| Yusuf et al. [61] (RecipeGen) | Transformer Seq2Seq | NL query | Executable Output | Implicit (decoder) | × (1-shot) | × (retrain) | BLEU: 59.1% SR@1: 50.6% MRR@3: 0.575 | 45K-120K merged sets | High (labeled pairs) |
| Yao et al. [58] | Hierarchical RL | Dialogue | Executable JSON | ✓ | ✓ (dialogue) | × (retrain) | Sim: 89.4% Human: 63.4% | 291K recipes + sim. dialogue | High (recipes + sim.) |
| Corno et al. [6] (TAPrec) | Ontology Reasoning | Partial selection | Service names | ✓ (goal-based) | × | Manual update | N/A (recommends) | Ur et al. scale + EUPont ont. | Medium (rules + ont.) |
| Huang et al. [16] (TAP-AHGNN) | GNN Embeddings | Context | Service names | Partial | × | Partial (graph) | HR@10: 0.937 NDCG@10: 0.952 MRR@10: 0.967 | 9,884 recipes 1,249 triggers 748 actions | Medium (graph struct.) |
| Wu et al. [55] (MvTAP) | Multi-view Learning | Multi-view (U/D/K) | User intentions | × (independent) | × | × (retrain) | Micro-F1: 0.783 Macro-F1: 0.751 | N/A N/A | Medium (multi-label) |
| Cimino et al. [5] (TARGE) | Cross-view Contrastive Learning + LLM + Perplexity Ranking | NL query (user intent) | Executable rules (channel + functionality) | ✓ (CRG: trigger-action conditioned) | × (1-shot, top-K available) | Partial (clustering, no full retrain) | EM: 61% (gold) EM: 41% (noisy) EM: 47% (1-shot) MRR@3: 0.65/0.45/0.50 | 34.8K rules 468 T-channels 446 A-channels 1.3K T-func 948 A-func | Medium-High (rules for contrastive learning + LLM LoRA) |
| **FARM (Ours)** | **Contrastive Learning + RAG + Multi-Agentic AI** | **NL query** | **Executable JSON** | **✓ (explicit)** | **✓ (fallback)** | **✓ (index)** | **R@1: T72%/A79%, R@5: 92% MRR@5: 0.81, MRR@3: 0.84 JM: 81%/62%/70% Faith: 0.44-0.48, Topic: 0.80-0.82** | **16.5K applets 1.7K trigger functions 1.3K action functions** | **Medium (12.6K pairs for encoder)** |

**Legend:** ✓ = Has capability; × = Does not have capability; N/A = Method does not produce executable outputs by design (recommendation systems).

**Data Requirements:** FARM requires 16.5K applets (12.6K training pairs) for contrastive encoder fine-tuning and 3,011 Function-level schema descriptions for indexing (Applets). TARGE requires 34.8K rules for cross-view contrastive learning and LLM fine-tuning via LoRA. Other supervised methods require 45K-291K labeled (query → trigger + action) pairs or complex annotations.

**FARM Metrics:** R@1 = Stage 1 retrieval (T=Trigger encoder, A=Action encoder reported separately); R@5 = Stage 1 retrieval (both encoders); MRR@5 = Stage 1 retrieval Mean Reciprocal Rank (macro-avg over trigger/action); MRR@3 = End-to-end Selection MRR; JM = Joint Accuracy on Gold/Noisy/One-shot sets; Faith = Faithfulness; Topic = Topic Adherence (Stage 2 quality).

**Compositional Reasoning:** Explicit verification that trigger outputs match action inputs (FARM, TAPrec); Implicit through decoder architecture (RecipeGen); Partial through goal-based reasoning (TAPrec, TAP-AHGNN).

**Service Extensibility:** Whether new services can be added without retraining the entire model. FARM supports this through simple index updates.

## 3 Approach

In this section, we will introduce our purposed method, which comprises the following 2 stages which is illustrated in Figure 3:

### 3.1 System Architecture Overview

Stage 1 uses contrastive-trained dual encoders for high-recall candidate retrieval. Stage 2 applies multi-agent selection with LLM-based verification to produce executable bindings. This leverages complementary strengths: neural speed and recall, plus LLM precision and schema-awareness.

Formally, given a user query $q$, Stage 1 retrieves candidate sets:

$$\mathcal{T}_k = \underset{t_i \in \mathcal{T}}{\text{top-}k} \left(\text{sim}(\mathcal{E}_T(q), \mathcal{E}_T(t_i))\right) \tag{1}$$

$$\mathcal{A}_k = \underset{a_j \in \mathcal{A}}{\text{top-}k} \left(\text{sim}(\mathcal{E}_A(q), \mathcal{E}_A(a_j))\right) \tag{2}$$

where $\mathcal{E}_T$ and $\mathcal{E}_A$ are the trigger and action encoders respectively, $\mathcal{T}$ and $\mathcal{A}$ are the full function catalogs, and $\text{sim}(\cdot, \cdot)$ denotes cosine similarity. Stage 2 then selects the optimal pair $(t^*, a^*) \in \mathcal{T}_k \times \mathcal{A}_k$ through multi-agent selection and generates the binding function $\beta : \mathcal{F}_a \to \mathcal{I}_t \cup \mathcal{S}$, mapping action input fields to trigger output ingredients or static values.

### 3.2 Stage 1: Contrastive Dual-Encoder Retrieval

The retrieval stage addresses the scalability challenge by reducing the search space from $O(|\mathcal{T}| \times |\mathcal{A}|)$ to $O(k^2)$ where $k \ll |\mathcal{T}|, |\mathcal{A}|$. We train separate encoders for triggers and actions because the same query requires identifying different semantic aspects: the *event* to detect (trigger) versus the *action* to perform.

*3.2.1 Schema-Enriched Text Representation.* Unlike prior work that embeds only service names, we encode complete function schemas to enable field-level retrieval. Each function $x \in \mathcal{T} \cup \mathcal{A}$ is represented as:

$$\text{text}(x) = [\text{channel}] \, [\text{category}] \, \text{name. desc} \, || \, \text{schema}(x) \tag{3}$$

where $||$ denotes concatenation and $\text{schema}(x)$ encodes the function's data interface:

**Trigger Schema** (data provider):

$$\text{schema}(t) = \text{"Provides: "} \bigoplus_{i \in \mathcal{I}_t} (\text{name}_i, \text{type}_i) \tag{4}$$

**Action Schema** (data consumer):

$$\text{schema}(a) = \text{"Requires: "} \bigoplus_{f \in \mathcal{F}_a} (\text{name}_f, \text{required}_f) \tag{5}$$

where $\mathcal{I}_t$ denotes the ingredient set of trigger $t$, $\mathcal{F}_a$ denotes the field set of action $a$, and $\bigoplus$ represents formatted concatenation. This representation enables the model to learn that, for example, "log sensor data" should retrieve actions requiring data fields, not just services with "log" in the name.

*Contrastive Learning with InfoNCE.* We train encoders using the InfoNCE objective [46], which learns to discriminate positive query-function pairs from in-batch negatives:

$$\mathcal{L}_{\text{InfoNCE}} = -\mathbb{E}_{(q,d^+) \sim \mathcal{D}} \left[ \log \frac{\exp(s(q, d^+)/\tau)}{\sum_{i=1}^{B} \exp(s(q, d_i)/\tau)} \right] \tag{6}$$

where $s(q, d) = \cos(\mathcal{E}(q), \mathcal{E}(d))$ is the cosine similarity between query and document embeddings, $\tau$ is the temperature hyperparameter, $B$ is the batch size, and $d^+$ is the positive (ground-truth) function for query $q$. The denominator sums over all $B$ documents in the batch, treating co-occurring samples as hard negatives.

The temperature $\tau = 0.05$ creates a sharp softmax distribution that strongly penalizes confusion between similar candidates. With batch size $B = 16$, each training step provides 15 hard negatives sampled from the same domain, encouraging fine-grained discrimination.

*3.2.2 Semantic Preservation via Layer Freezing.* A critical challenge in domain-specific training is catastrophic forgetting [9, 23, 26, 62]: the model loses pretrained semantic knowledge while adapting to the target task. We observed that full training caused the model to memorize exact service names while forgetting semantic equivalences (e.g., "log" ≈ "record" ≈ "add row").

We address this through selective layer freezing. Let $\theta = \{\theta^{(0)}, \theta^{(1)}, ..., \theta^{(L)}\}$ denote the parameters of an $L$-layer transformer. We partition layers into frozen and trainable sets:
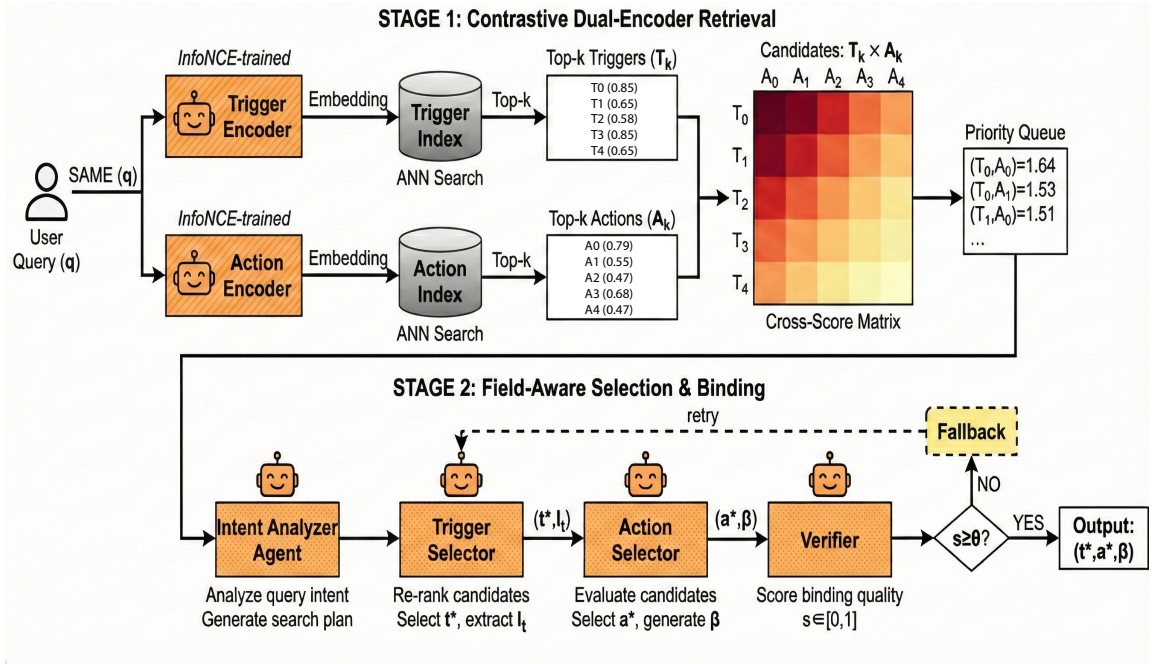


Fig. 3. FARM two-stage architecture. **Stage 1**: Dual contrastive encoders retrieve top-$k$ trigger and action candidates; cross-score matrix ranks $k \times k$ pairs by compatibility scores (70% ingredient-field coverage, 30% retrieval quality). **Stage 2**: Field-aware selection pipeline performs deep ingredient-to-field analysis—Intent Analyzer decomposes query intent, Trigger Selector identifies optimal trigger $t^*$ and extracts its ingredient schema $I_t$ (available output fields), Action Selector evaluates field compatibility and generates bindings $\beta$ mapping ingredients to action input fields, and Verifier scores binding completeness and semantic coherence. Fallback retries with next candidate pair when verification score $s < \theta$.
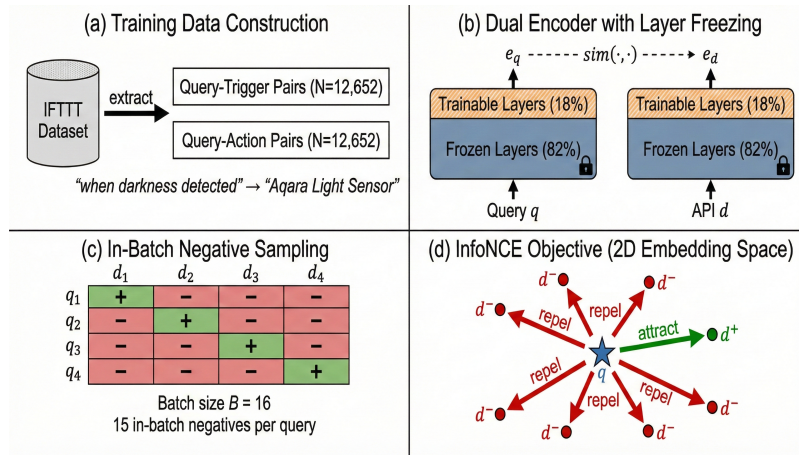
Fig. 4. Contrastive training pipeline for dual encoders. **(a)** Training data construction: query-function pairs extracted from IFTTT dataset. **(b)** Dual encoder architecture with layer freezing—trainable layers (orange, 18%) adapt to domain-specific patterns while frozen layers (blue-gray, 82%) preserve pretrained semantics. **(c)** In-batch negative sampling: diagonal entries are positive pairs, off-diagonal entries serve as hard negatives. **(d)** InfoNCE objective in embedding space: queries are pulled toward positive documents and pushed away from negatives.

$$\theta_{\text{frozen}} = \{\theta^{(0)}, ..., \theta^{(\ell)}\}, \quad \theta_{\text{train}} = \{\theta^{(\ell+1)}, ..., \theta^{(L)}\} \tag{7}$$

During training, gradients are computed only for $\theta_{\text{train}}$:

$$\theta_{\text{train}}^{(t+1)} = \theta_{\text{train}}^{(t)} - \eta \nabla_{\theta_{\text{train}}} \mathcal{L}_{\text{InfoNCE}} \tag{8}$$

For our final transformer encoder, we have chosen EmbeddingGemma [39] with 24 layers. We froze layers 0-11 (including the embedding layer), resulting in 82% of the parameters (252M of 307M) frozen. This preserves the pretrained model's semantic generalization in lower layers while allowing upper layers to learn domain-specific retrieval patterns.

**Theoretical Motivation.** Research on transformer representations [19, 31, 43] shows that lower layers encode lexical and syntactic features while upper layers encode task-specific semantics. By freezing lower layers, we preserve the model's ability to recognize semantic equivalences ("darkness detected" $\approx$ "low light sensed") while training upper layers for trigger-action discrimination.

*3.2.3  Training Configuration.* We train separate encoders for triggers and actions with identical architectures but independent weights. Table 2 summarizes the key hyperparameters. We use a low temperature ($\tau = 0.05$) to create sharp softmax distributions that strongly penalize confusion between similar candidates. With batch size 16, each training step provides 15 hard negatives sampled from the same domain, encouraging fine-grained discrimination between semantically similar functions.

*3.2.4  Trigger-Action Programming Formulation.* Trigger Action Programming TAP enables end users to create automation rules in the form IF trigger THEN action, following the paradigm of Ur et al. [44]. Formally, given a user query $q$ expressed in natural language, the system must:

(1) **Select** a trigger $t^* \in \mathcal{T}$ from the trigger catalog

| Parameter | Value |
|---|---|
| Base Model | Pretrained Text Encoder |
| Frozen Layers | 0–11 (+ embeddings) |
| Loss Function | InfoNCE |
| Temperature $\tau$ | 0.05 |
| Batch Size | 16 |
| Learning Rate | $2 \times 10^{-5}$ |
| Epochs | 3 |
| Training Pairs | 12,652 |

Table 2. Encoder training configuration.

(2) **Select** an action $a^* \in \mathcal{A}$ from the action catalog

(3) **Generate** bindings $\beta : \mathcal{F}_a \to \mathcal{I}_t \cup \mathcal{S}$ mapping action input fields to trigger output ingredients or static values

Each trigger $t$ exposes a set of *ingredients* $\mathcal{I}_t = \{i_1, i_2, ...\}$—typed output fields produced when the trigger fires (e.g., Subject, Body for an SMS trigger). Each action $a$ requires a set of *fields* $\mathcal{F}_a = \{f_1, f_2, ...\}$—typed input parameters needed for execution (e.g., To, Subject for an email action). The binding function $\beta$ creates the data flow that makes applets executable.

*3.2.5 Why Retrieval Alone is Insufficient.* While retrieval-augmented generation (RAG) has achieved success in many NLP tasks [18, 25], several factors make pure retrieval insufficient for TAP:

**Selection Ambiguity.** High recall ensures correct functions appear in the candidate set, but does not identify *which* candidate is correct. With high retrieval recall (Section 4.4.1 shows R@5 > 92% for both encoders), selecting the correct pair from $k \times k$ combinations—25 candidate pairs when $k = 5$ (Fig. 14)—requires reasoning beyond embedding similarity. Naive selection that simply takes the rank-1 result from each encoder evaluates only 1 pair and achieves 58% joint accuracy, whereas multi-agent search over all 25 pairs achieves 81% (Section 4.6).

**Schema Compatibility.** Retrieval operates on text similarity, not schema compatibility. A trigger providing [temperature, humidity] may rank highly for "weather logging" but be incompatible with an action requiring [image_url, caption]. Understanding data flow between functions requires schema-level reasoning.

**Binding Generation.** Even with correct trigger-action pairs, generating executable applets requires mapping trigger ingredients to action fields—determining that StockName should fill the row_content field. This semantic alignment task is beyond retrieval's capability.

**Ambiguity Resolution.** User queries often admit multiple valid interpretations. "Turn on lights when I leave" could use location triggers (GPS-based) or calendar triggers (schedule-based); an LLM can reason about context and user intent to select appropriately.

These limitations motivate our two-stage architecture: retrieval for efficient candidate generation, followed by multi-agent selection for reasoning and binding.

*3.2.6 Design Consideration: Why Not LoRA?.* Low-Rank Adaptation (LoRA) [14] has become the dominant parameter-efficient fine-tuning method, achieving impressive results across many tasks by learning low-rank update matrices $\Delta W = BA$ where $B \in \mathbb{R}^{d \times r}$ and $A \in \mathbb{R}^{r \times d}$ with rank $r \ll d$. Extensions like QLoRA [10] further reduce memory requirements through quantization.

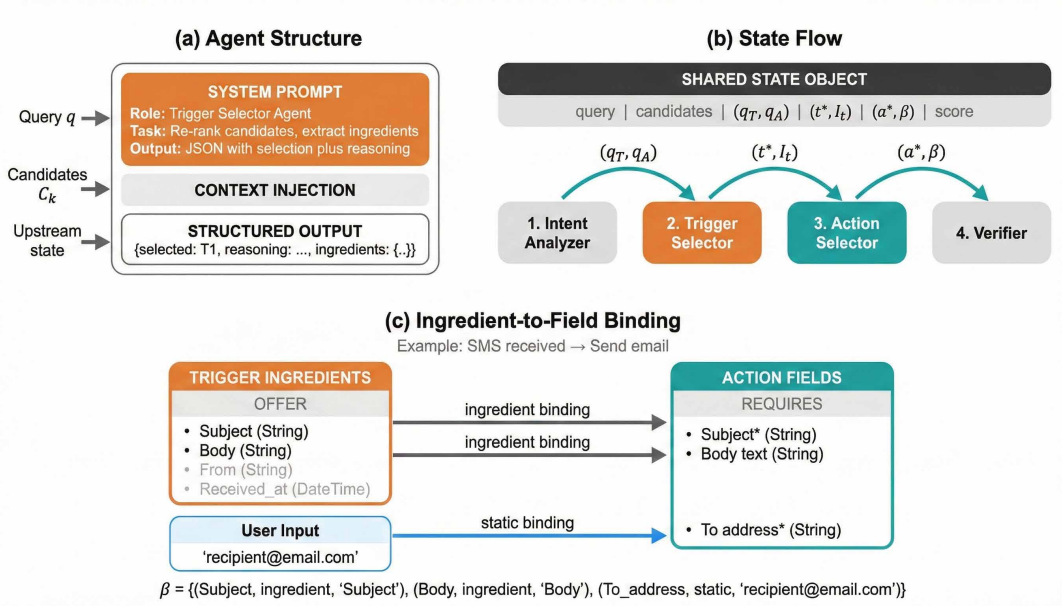However, for contrastive retrieval training, we found layer freezing superior to LoRA for three reasons:

Fig. 5. Stage 2 multi-agent selection components. (a) Agent structure showing input query, LLM processing with retrieved candidates $C_k$, and structured JSON output. (b) State-based pipeline where Intent Analyzer, Trigger Selector, Action Selector, and Verifier communicate through shared state. (c) Ingredient-to-field binding illustrating how trigger outputs (Subject, Body) dynamically populate action inputs, with static user-provided values for required fields.

**Semantic Preservation.** LoRA modifies all layers simultaneously through low-rank updates, which can subtly alter the pretrained semantic space. Layer freezing ensures that frozen layers retain pretrained representations, which we validate empirically in Section 4.6.

**Representation Stability.** Contrastive learning requires stable representations for effective in-batch negative sampling. LoRA's distributed updates across all layers can destabilize the embedding space during training, whereas layer freezing maintains a stable foundation in lower layers while adapting upper layers.

**Computational Simplicity.** Layer freezing requires no additional adapter modules, reducing implementation complexity and inference overhead. The frozen/trainable partition is straightforward: gradients simply are not computed for frozen parameters.

Recent work by Biderman et al. [2] has shown that LoRA can underperform full fine tuning on certain tasks, and our domain, contrastive retrieval with small datasets, appears to be one where the inductive bias of layer freezing that preserves lower layer semantics is particularly beneficial.

*3.2.7 State-Based Agent Communication.* Unlike message-passing multi-agent systems where agents exchange discrete messages, our agents communicate through a *shared state object* that accumulates information as execution proceeds. This design pattern, common in dataflow architectures, provides several advantages: (1) any downstream agent can access any upstream decision, (2) state can be serialized for debugging and replay, and (3) conditional routing decisions can inspect any state field.

Table 3 defines the key state fields and their producers/consumers.

| State Field | Producer | Consumer |
|---|---|---|
| query $q$ | Input | All agents |
| trigger_candidates $\mathcal{T}_k$ | Stage 1 | Trigger Selector |
| action_candidates $\mathcal{A}_k$ | Stage 1 | Action Selector |
| search_intents $(q_T, q_A)$ | Intent Analyzer | Selectors |
| selected_trigger $(t^*, I_t)$ | Trigger Selector | Action Selector |
| bindings $(a^*, \beta)$ | Action Selector | Verifier |
| verifier_score $s$ | Verifier | Routing logic |
| llm_overrode_rag | Selectors | Fallback logic |

Table 3. State schema for inter-agent communication. Each agent reads upstream fields and writes to designated output fields.

*3.2.8 Intent Analyzer Agent.* The Analyzer Agent initiates Stage 2 by analyzing the user query using chain-of-thought (CoT) reasoning [54]. Given query $q$, the agent performs:

**Intent Decomposition.** The analyzer identifies two distinct intents embedded in the query: (1) the trigger event to detect, and (2) the action to execute. For example, given "Change the light to green if stock price rises," the analyzer extracts:

- Trigger intent: "stock price rises" (financial market event)
- Action intent: "change light to green" (smart home control)

**Search Query Generation.** The analyzer generates optimized search queries $(q_T, q_A)$ to clarify trigger and action intents. These queries may expand abbreviations ("temp" $\rightarrow$ "temperature"), resolve ambiguities, or reformulate colloquial expressions into function-compatible terminology.

**Reasoning Trace.** Following the CoT paradigm [50, 53], the agent produces an explicit THINKING section followed by a Intent Analyze:

```
THINKING: 1. User wants to detect stock price changes (TRIGGER).
2. User wants to control smart light color (ACTION).
3. Need financial function for trigger, IoT function for action.
PLAN: Search for stock/price triggers and light/color actions.
```

*3.2.9 Trigger Selector Agent.* The Trigger Selector applies a two-stage selection process combining neural retrieval with LLM re-ranking, following the retrieve-then-rerank paradigm established in modern information retrieval [35].

**Candidate Pair Ranking.** Rather than selecting trigger and action independently, we pre-compute compatibility scores for all $k \times k$ candidate pairs:

$$\text{score}(t_i, a_j) = 0.7 \cdot \text{coverage}(t_i, a_j) + 0.3 \cdot \frac{\text{sim}(q, t_i) + \text{sim}(q, a_j)}{2} \tag{9}$$

where $\text{coverage}(t_i, a_j)$ measures ingredient-to-field compatibility (defined below in Equation 11) and the retrieval similarity term provides ranking based on Stage 1 encoder quality. This weighted combination prioritizes schema compatibility (70%) while using retrieval confidence (30%) as a tie-breaker. Pairs are sorted by this score to create a priority queue enabling best-first search through the candidate space.

**Agreement-Based Selection.** The agent presents top-$k$ trigger candidates to the LLM with full function descriptions. The LLM selects its preferred match $t_{LLM}$ with reasoning. To prevent hallucination from degrading retrieval quality, we

employ an *agreement-based* mechanism:

$$t^* = \begin{cases} t_{LLM} & \text{if } \dfrac{s_{RAG}(t_{LLM})}{s_{RAG}(t_0)} \geq \tau_T \\ t_0 & \text{otherwise} \end{cases} \tag{10}$$

where $t_0$ is RAG's top candidate, $s_{RAG}(\cdot)$ denotes retrieval similarity, and $\tau_T = 0.95$ is the trigger agreement threshold. This ensures LLM can only override RAG when its choice has comparable retrieval confidence, balancing reasoning capability with retrieval reliability.

**Ingredient Extraction.** After selection, the agent extracts available ingredients $I_t = \{(name, type, example)\}$ from the trigger's function schema. Ingredients represent data fields emitted when the trigger fires (e.g., `StockName`, `Price`, `PercentageChange`).

*3.2.10 Action Selector Agent.* The Action Selector receives the selected trigger $(t^*, I_t)$ and evaluates action candidates for schema compatibility.

**Cross-Scoring.** The agent computes a coverage score measuring how well trigger ingredients can satisfy action requirements:

$$\text{coverage}(t, a) = \frac{|\{f \in \mathcal{F}_a^{req} : \exists i \in I_t, \, \text{match}(i, f)\}|}{|\mathcal{F}_a^{req}|} \tag{11}$$

where $\mathcal{F}_a^{req}$ denotes required action fields. The $\text{match}(i, f)$ function checks compatibility using:

- **Direct matching**: `temperature` $\rightarrow$ `temperature`
- **Substring matching**: `stock_name` $\rightarrow$ `name`
- **Semantic mapping**: `message` $\approx$ `body` $\approx$ `content`

**LLM Re-ranking.** Similar to trigger selection, the LLM re-ranks action candidates with full context about the selected trigger and its ingredients. The agreement threshold $\tau_A = 0.80$ is lower than triggers, as action selection benefits more from LLM reasoning about functional intent.

**Binding Generation.** Upon selecting action $a^*$, the agent generates bindings $\beta$ mapping each action field to a data source:

$$\beta = \{(f, \text{src}, \text{val}) : f \in \mathcal{F}_a\} \tag{12}$$

where $\text{src} \in \{\text{ingredient}, \text{static}\}$. Ingredient bindings use trigger data (e.g., `StockName` $\rightarrow$ `row_content`); static bindings use placeholder values for fields without matching ingredients.

*3.2.11 Verifier Agent.* The Verifier implements the LLM-as-Judge pattern [64] to evaluate the complete applet configuration $(t^*, a^*, \beta)$:

- **Binding Quality**: Are bindings semantically appropriate?
- **Completeness**: Are all required action fields bound?
- **Executability**: Can this configuration execute without runtime errors?

The Verifier produces a quality score $s \in [0, 1]$ and natural language critique. A rule-based fallback handles LLM parsing failures by checking trigger/action presence, binding count, and required field coverage.

*3.2.12 Quality-Gated Fallback.* When the verifier score falls below threshold $\theta_v = 0.5$, the system triggers fallback logic:

(1) **LLM Fallback**: If the LLM overrode RAG's selection (`llm_overrode_rag` = true), retry with RAG's original top choice. This catches cases where LLM reasoning was incorrect.

(2) **Pair Iteration**: If still failing, advance to the next candidate pair from the priority queue and restart from Trigger Selector.

The system attempts up to $k^2$ pairs before declaring failure. In our evaluation (Section 4.4.2), most queries succeed within two attempts due to Stage 1's high recall.

*3.2.13 Execution Example.* We illustrate the complete pipeline with query: "Change light to green if stock price rises."

**Step 1: Analyzer.** Analyzes query and decomposes: $q_T$ = "stock price rises," $q_A$ = "change light color." Produces reasoning trace explaining the decomposition.

**Step 2: Trigger Selector.** Receives candidates $\mathcal{T}_5$ from RAG. Top results: $T_0$ = "Price rises above" (0.753), $T_1$ = "Today's price rises by percentage" (0.744). LLM selects $T_1$ with reasoning about percentage-based detection being more appropriate. Agreement ratio $0.744/0.753 = 0.99 \geq 0.95$, so LLM override accepted. Extracts $I_t$ = [`StockName`, `Price`, `PercentageChange`, `CheckTime`].

**Step 3: Action Selector.** Evaluates $\mathcal{A}_5$ against $I_t$. Computes coverage scores. Selects $A_0$ = "Turn on / change light mode" (0.627). Generates bindings: $\beta$ = [(`color`, static, "green"), (`light`, static, "Living room")].

**Step 4: Verifier.** Evaluates $(t^*, a^*, \beta)$. Confirms all required fields bound. Assigns score $s = 0.85 \geq 0.5$. Output accepted.

*3.2.14 Chain-of-Thought Prompting.* All agents employ chain-of-thought (CoT) prompting [4] to produce interpretable reasoning traces. Each prompt instructs the LLM to:

(1) Analyze user intent explicitly before making decisions
(2) Evaluate each candidate's fitness with stated criteria
(3) Check schema compatibility between data sources and targets
(4) Justify the final selection with concrete reasoning

This structured deliberation improves decision accuracy through explicit reasoning and provides interpretability for system debugging and error analysis.

## 3.3  Embedding Model Selection

The choice of base encoder significantly impacts retrieval quality. We identify three key selection criteria:

**Semantic Understanding.** The encoder should be trained on diverse corpora to provide robust semantic representations that generalize to WoT service descriptions. Models trained primarily on narrow domains may struggle with the mixed technical vocabulary common in function documentation.

**Efficiency.** The model should balance performance with computational tractability for domain-specific training. Larger models often offer diminishing returns, as pretrained general knowledge matters less than task-specific adaptation.

**Embedding Dimensionality.** The embedding dimension should balance expressiveness with retrieval efficiency, enabling fast approximate nearest neighbor search over large function catalogs.

For our implementation, we select Gemma Embedding [39] ( 307M parameters, 768-dimensional embeddings), which satisfies these criteria while remaining efficient for contrastive training.

## 4 Experiments

We evaluate FARM through comprehensive experiments designed to assess both stages of our architecture. Our evaluation addresses three key questions: (1) Does contrastive training with layer freezing improve retrieval quality while preserving semantic generalization? (2) Does multi-agent selection improve end-to-end accuracy over retrieval alone? (3) How does FARM compare to existing trigger-action programming approaches?

### 4.1 Function-Level TAP Dataset

Our dataset comprises 1,724 trigger functions and 1,287 action functions from the IFTTT platform, spanning 19 categories including Smart Home, Finance, Health, Developer Tools, Social Media, and others. This creates a search space of $1,724 \times 1,287 = 2,218,788$ possible trigger-action pairs. Critically, our dataset operates at the *function level* with complete schema specifications, which differs fundamentally from traditional service-level identification tasks.

*4.1.1 Service-Level vs. Function-Level Granularity.* To understand this distinction, consider the difference between identifying a service and selecting a specific trigger or action:

- **Service-Level (Service):** Identify the service providing functionality (e.g., "The New York Times", "Google Sheets", "Evernote")
- **function-Level (Trigger/Action):** Select the specific trigger or action within that service (e.g., "New article from search", "Add row to spreadsheet", "Append to note")

A single service may expose dozens of distinct trigger and action functions. For instance, the New York Times service provides multiple trigger functions ("New article from search", "New popular article", "New article in section", etc.), each with different functionality, parameters, and data outputs. Service-level identification only determines *which service* to use. Functions-level selection determines *which specific trigger or action* to configure and execute.

*4.1.2 Schema Structure and Data Interface Specifications.* Each trigger or action function in our dataset includes complete schema specifications that define its data interface. Table 4 shows a concrete example from our dataset, illustrating the full schema structure for a trigger-action pair.:

**Trigger Schema:**

- `service_name`: Function identifier (e.g., "New article from search")
- `category`: Domain category (e.g., "News & information")
- `description`: Natural language explanation of trigger behavior
- `Trigger fields`: Input parameters required to configure the trigger (e.g., "Search for" with type String)
- `Ingredients`: Output data fields produced when the trigger fires, with:
  - `Slug`: Programmatic identifier (e.g., `Title`, `ArticleUrl`)
  - `Type`: Data type (String, Date, Number, Boolean, etc.)
  - `Filter code`: Access path used by the IFTTT filter code environment
  - `Example`: Representative value showing expected format

**Action Schema:**

- `service_name`: Function identifier (e.g., "Append to note")
- `category`: Domain category (e.g., "Popular services")
- `description`: Natural language explanation of action behavior

| Trigger Function: "New article from search" | Action Function: "Append to note" |
|---|---|
| **Category:** News & information | **Category:** Popular services |
| **Description:** This Trigger fires every time a new article that is published by The New York Times matches a search query you specify. | **Description:** This Action will append to a note as determined by its title and notebook. |
| **Trigger Fields:**<br>• Search for (String, required) | **Action Fields:**<br>• Title (String, required)<br>• Body (String, required)<br>• Notebook (String, optional)<br>• Tags (String, optional) |
| **Ingredients (Outputs):**<br>• Title (String)<br>• Author (String)<br>• Blurb (String)<br>• ArticleUrl (String)<br>• ImageUrl (String)<br>• Source (String)<br>• Section (String)<br>• Keywords (String)<br>• PublishedDate (Date with time) | **Data Flow Mapping:**<br>Ingredients → Fields<br><br>Example binding:<br>• Ingredient `Title` → Field `Title`<br>• Ingredient `Blurb` → Field `Body`<br>• Ingredient `Keywords` → Field `Tags` |

Table 4. Example function-Level Schema from Dataset. The trigger provides 9 typed ingredients (Title, Author, Blurb, ArticleUrl, ImageUrl, Source, Section, Keywords, PublishedDate); the action requires 4 fields (Title, Body, Notebook, Tags), with Title and Body marked as required. This schema information enables reasoning about data flow compatibility and automatic generation of ingredient-to-field bindings.

- `Action fields`: Input parameters required to execute the action, with:
  - `Label`: Human-readable field name (e.g., "Title", "Body")
  - `Slug`: Programmatic identifier (e.g., `title`, `body`)
  - `Required`: Boolean indicating if field must be provided
  - `Helper text`: Usage instructions or constraints
  - `Filter code method`: Method signature used by the IFTTT filter code environment to set the field value

*4.1.3 Why Function-Level Schemas Enable Executable Applets.* The distinction between service-level and function-level is not merely semantic. It fundamentally changes the task complexity and system requirements.

**Service-Level Limitation:** Prior work treats TAP as a classification problem over service names. A system that outputs "The New York Times" and "Evernote" provides no actionable information. Which specific trigger should be used. Which specific action should be selected. How to configure the parameters. Which ingredient values should fill which action fields. The output requires extensive manual configuration before execution.

**Function-Level Advantage:** Our dataset's function-level schemas enable fully automated applet generation. By including:

- **Function specifications**: The exact trigger and action functions to select ("New article from search", "Append to note")
- **Type information**: Data type constraints for validation (String, Date, Boolean)
- **Required field markers**: Which parameters must be provided versus optional

- **Ingredient definitions**: Available data outputs from triggers
- **Field requirements**: Expected inputs for actions

The system can reason about *schema compatibility* (Does the trigger produce data the action needs?), generate *ingredient-to-field bindings* (Which trigger outputs map to which action inputs?), and produce *executable configurations* ready for deployment without human intervention.

For example, given the query "Send New York Times articles about recipes to Evernote", a service-level system outputs service names, while our function-level system produces:

- Trigger: "New article from search" with field `Search for` = "recipe"
- Action: "Append to note" with bindings:
  - `Title` ← Ingredient `Title`
  - `Body` ← Ingredient `Blurb`
  - `Tags` ← Ingredient `Keywords`

This complete specification is executable within the IFTTT execution model without further configuration. The function-level granularity with schema information is what makes our task and dataset fundamentally different from prior service identification approaches.

We evaluate FARM on the IFTTT platform, a widely-used trigger-action programming service for WoT (Web of Things) automation. The service catalog comprises:

- **Triggers**: 1,724 distinct functions across 19 categories (Smart Home, Finance, Health, Developer Tools, Social Media, etc.)
- **Actions**: 1,287 distinct functions across 19 categories
- **Search Space**: $1,724 \times 1,287 = 2,218,788$ possible trigger-action function pairs

Each function entry contains structured schema information: function name, category, natural language description, and data interface specifications (ingredients for triggers, required fields for actions). Following IFTTT's internal terminology, the dataset labels specific functions as `service_name` (e.g., "New article from search", "Add row to spreadsheet"), which we evaluate at the function-level rather than the platform-level (e.g., "The New York Times", "Google Sheets").

*Terminology.* IFTTT refers to individual trigger/action functions as `service_name` (e.g., "Add row to spreadsheet"). We refer to these as *functions* throughout. We reserve *platform service* for the parent application (e.g., "Google Sheets"). Unless stated otherwise (e.g., Table 5), all Stage 1 and Stage 2 metrics are computed at the function level.

## 4.2 Experimental Setup

*4.2.1 Evaluation Sets.* Following the evaluation protocol established by Cimino et al. [5], we construct three evaluation sets to test different aspects of system robustness:

- **Gold Set**: Clear, well-formed automation requests with unambiguous intent (e.g., "When darkness detected, log to spreadsheet")
- **Noisy Set**: Vague or ambiguous descriptions requiring inference (e.g., "track my fitness stuff")
- **One-Shot Set**: Queries involving rare or specialized functions with limited training exposure

*4.2.2 Model Configuration.* **Stage 1 - Contrastive Encoders:**

- **Base Model**: EmbeddingGemma [39] ( 307M parameters, 768-dimensional embeddings)
- **Architecture**: 24 transformer layers, 12 attention heads, 2,048 token context
- **Layer Freezing**: Layers 0–11 frozen, layers 12–23 trainable (18% of total parameters trainable)
- **Training**: InfoNCE loss [46], $\tau = 0.05$, batch size 16, 3 epochs, learning rate $2 \times 10^{-5}$

**Stage 2 - Multi-Agent Selection:**

- **LLM**: IBM Granite 4.0 Small [17] via local inference
- **Temperature**: 0.0 (deterministic outputs)
- **Note**: The multi-agent system operates through *prompting only*—no additional training is performed on the LLM. As Brown et al. [3] demonstrated, large-scale language models can achieve strong task performance without fine-tuning, and Wei et al. [52] showed that emergent abilities arise at sufficient scale through prompting alone. Agents use carefully designed system prompts with chain-of-thought reasoning [53] to elicit step-by-step analysis before producing structured outputs.

### 4.3 Evaluation Metrics

We employ metrics from established benchmarks to enable fair comparison with prior work.

*4.3.1 Stage 1: Retrieval Metrics.* We evaluate retrieval effectiveness using Recall at cutoff $K$ (Recall@K) and Mean Reciprocal Rank at cutoff $K$ (MRR@K), two standard ranking metrics in information retrieval [32, 47].

**Recall@K** measures retrieval coverage: what fraction of queries have the correct answer in the top-$K$ results? Given a query set $Q$:

$$\text{Recall@K} = \frac{1}{|Q|} \sum_{q \in Q} \mathbb{I}\left( \exists d \in \mathcal{D}_q^K \ \text{s.t.} \ d \in \mathcal{R}_q \right), \tag{13}$$

where $\mathcal{D}_q^K$ denotes the top $K$ retrieved documents for query $q$, $\mathcal{R}_q$ is the set of relevant documents for $q$, and $\mathbb{I}(\cdot)$ is the indicator function. Practically, R@1 asks: "Is the correct answer ranked first?" while R@5 asks: "Is the correct answer somewhere in the top 5?" We report R@1 and R@5 to evaluate both precision (top result quality) and recall (candidate set coverage for Stage 2).

**Joint Recall@K** extends this to dual-encoder retrieval, measuring when *both* components succeed simultaneously:

$$\text{Joint Recall@K} = \frac{1}{|Q|} \sum_{q \in Q} \mathbb{I}\left( t_{\text{true}} \in \mathcal{D}_q^K(T) \wedge a_{\text{true}} \in \mathcal{D}_q^K(A) \right), \tag{14}$$

where $\mathcal{D}_q^K(T)$ and $\mathcal{D}_q^K(A)$ denote the top $K$ triggers and actions retrieved for query $q$. For example, Joint R@5 = 85% means that for 85% of queries, the correct trigger appears in the trigger encoder's top-5 *and* the correct action appears in the action encoder's top-5. Critically, Joint R@5 establishes the *theoretical ceiling* for Stage 2 performance: Stage 2 selection cannot choose a correct pair that Stage 1 did not retrieve (Fig. 14).

**Mean Reciprocal Rank (MRR@K)** measures ranking quality by considering *where* the correct answer appears within the top-$K$ results:

$$\text{MRR@K} = \frac{1}{|Q|} \sum_{q \in Q} \begin{cases} \frac{1}{\text{rank}_q} & \text{if } \text{rank}_q \leq K, \\ 0 & \text{if } \text{rank}_q > K, \end{cases} \tag{15}$$

where $\text{rank}_q$ is the rank position of the first relevant document for query $q$. MRR@K gives partial credit based on position: rank 1 receives score 1.0, rank 2 receives 0.5, rank 5 receives 0.2, while anything beyond rank $K$ receives 0. The cutoff $K$ defines the evaluation scope and we use two values: (1) **MRR@5** for Stage 1 retrieval evaluation, measuring

how well individual encoders rank trigger and action candidates (macro-averaged over both encoders, reported in Fig. 7); (2) **MRR@3** for function-level prediction, applied *twice*—first to evaluate FARM's own trigger-action pair selection performance, then to compare FARM against prior work (LAM [28], RecipeGen++ [60], TARGE [5]) using the same metric for fair comparison (Table 5). MRR@3 follows established evaluation protocols in the TAP literature for function-level prediction tasks.

*4.3.2 Stage 2: Selection and Generation Metrics.* We define end-to-end evaluation metrics that measure system performance on complete applet generation:

**Goal Accuracy** measures task completion with partial credit, evaluating whether the system achieves the user's automation intent:

$$\text{Goal Accuracy}(q) = \begin{cases} 1.0 & \text{if } t_{\text{pred}} = t_{\text{true}} \wedge a_{\text{pred}} = a_{\text{true}} \\ 0.5 & \text{if } (t_{\text{pred}} = t_{\text{true}}) \oplus (a_{\text{pred}} = a_{\text{true}}) \\ 0.0 & \text{otherwise} \end{cases} \tag{16}$$

where $\oplus$ denotes exclusive OR (exactly one match). Practically, Goal Accuracy = 1.0 means the applet is fully correct, 0.5 means partial success (user achieves either trigger monitoring or action execution, but not the complete automation), and 0.0 means complete failure. We consider a query successful when Goal Accuracy $\geq$ 0.5, used to compute the Success Rate metric.

**Joint Accuracy** measures service-pair correctness, requiring both trigger and action to be correctly selected:

$$\text{Joint Accuracy}(q) = \mathbb{I}(t_{\text{pred}} = t_{\text{true}} \wedge a_{\text{pred}} = a_{\text{true}}),$$

where $t_{\text{pred}}$ and $a_{\text{pred}}$ are the predicted trigger and action service names, $t_{\text{true}}$ and $a_{\text{true}}$ are the ground truth services, and $\mathbb{I}(\cdot)$ is the indicator function. Joint Accuracy returns 1 only when both services are correct, and 0 otherwise. It is computed from the *full system* (Stage 1 + Stage 2 multi-agent selection), allowing Stage 2 LLM reasoning to refine Stage 1 candidates. Unlike Goal Accuracy which provides partial credit (0.5 for one service correct), Joint Accuracy is binary (0/1), making it stricter and comparable to metrics used in prior work (LAM, RecipeGen++, TARGE).

**Trigger/Action Accuracy** measures individual component correctness:

$$\text{Trigger Accuracy} = \frac{1}{|Q|} \sum_{q \in Q} \mathbb{I}(t_{\text{pred}} = t_{\text{true}}), \tag{17}$$

and similarly for Action Accuracy with $a_{\text{pred}}$ and $a_{\text{true}}$. These metrics isolate which component (trigger selection vs. action selection) causes failures, informing system debugging and improvement.

**Success Rate** measures task-level success with partial credit, defined as the fraction of queries where the system achieves at least partial automation:

$$\text{Success Rate} = \frac{1}{|Q|} \sum_{q \in Q} \mathbb{I}(\text{Goal Accuracy}(q) \geq 0.5). \tag{18}$$

Since Goal Accuracy$(q) \in \{0.0, 0.5, 1.0\}$ and Joint Accuracy$(q) \in \{0, 1\}$, the following ordering holds pointwise and therefore also in expectation:

$$\text{Success Rate} \geq \mathbb{E}[\text{Goal Accuracy}] \geq \mathbb{E}[\text{Joint Accuracy}]. \tag{19}$$

Practically, Success Rate counts queries where at least one service (trigger or action) is correctly identified. Success Rate < 100% indicates failures in service identification or system crashes.

We adopt generation quality metrics from the RAGAS evaluation framework [11], following the functional definitions provided in the official documentation [37].

**Faithfulness** measures whether generated applet configurations are grounded in retrieved function schemas, preventing hallucination:

$$\text{Faithfulness} = \frac{|\text{Claims verifiable in retrieved schemas}|}{|\text{Total claims in generated applet}|}. \tag{20}$$

A "claim" is an assertion about service capabilities, field names, or data types in the generated applet. For example, if the system generates a binding `trigger.ingredient_name` → `action.field_name`, Faithfulness checks whether both `ingredient_name` exists in the trigger schema and `field_name` exists in the action schema. Faithfulness = 1.0 means all generated content is verifiable from retrieved documentation; lower scores indicate hallucinated field names or capabilities not present in the actual function specifications.

**Topic Adherence** evaluates whether selected services match the user's intended automation domain:

$$\text{Topic Adherence} = \frac{|\text{Applets matching reference automation domain}|}{|\text{Total applets}|}. \tag{21}$$

Practically, if a user's query mentions "fitness tracking," Topic Adherence checks whether the selected trigger and action belong to health/fitness services (e.g., Fitbit, Google Fit, MyFitnessPal) rather than unrelated domains (e.g., social media, smart home). This metric ensures the system respects user intent boundaries and does not inappropriately cross automation contexts.

## 4.4 Experiment results

*4.4.1 Stage 1: Contrastive Learning Results.* We first evaluate the retrieval stage in isolation to understand the impact of contrastive training and layer freezing. Figure 6 illustrates the training progression over three epochs. Figure 7 presents retrieval performance comparing pretrained baseline against contrastive-trained encoders on the Gold evaluation set.

Layer freezing with contrastive training achieves substantial improvements across all retrieval metrics. Trigger encoder: R@1 improves from 30% to 72% (+140%), R@5 from 59% to 92% (+56%), and MRR@5 from 42% to 79% (+88%). Action encoder: R@1 improves from 39% to 79% (+103%), R@5 from 57% to 92% (+61%), and MRR@5 from 48% to 83% (+73%). Both encoders achieve >92% R@5, ensuring the target trigger/action almost always appears in the top-5 candidates for Stage 2 selection.

Critically, these individual R@5 scores correspond to a measured Joint R@5 of 85%, meaning both the correct trigger *and* correct action appear in their respective top-5 results for 85% of queries. This Joint R@5 establishes the theoretical upper bound for Stage 2 performance—our multi-agent selection can only choose from pairs that Stage 1 successfully retrieved.
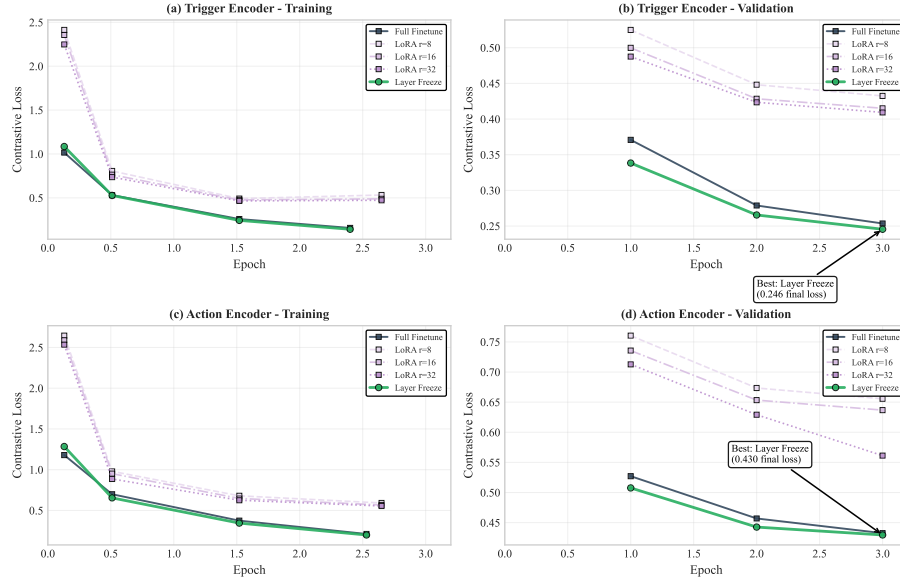
Fig. 6. Contrastive training curves showing loss convergence and R@1 improvement. Both encoders converge within 3 epochs, with the majority of improvement in the first epoch.
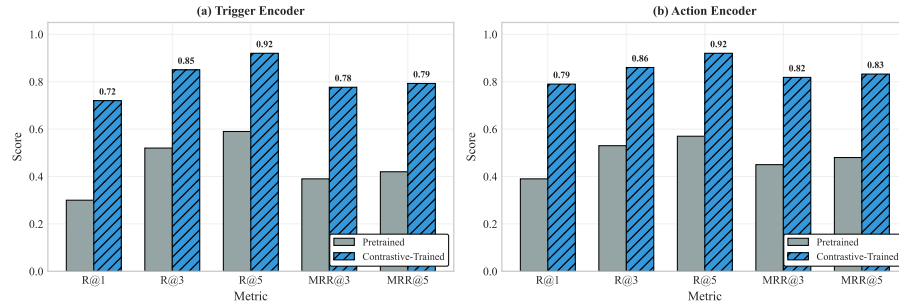


Fig. 7. Stage 1 retrieval quality. Contrastive training with layer freezing significantly improves trigger encoder (R@1: 30%→72%, R@5: 59%→92%, MRR@5: 42%→79%) and action encoder (R@1: 39%→79%, R@5: 57%→92%, MRR@5: 48%→83%).

*4.4.2 Stage 2: Multi-Agent Selection Results.* We evaluate the complete two-stage system, measuring end-to-end selection accuracy and generation quality across all three evaluation sets (Gold, Noisy, and One-Shot). Figure 8 presents selection accuracy across all three evaluation sets.

**Gold Set:** Clear queries enable high-confidence selection. The gap between individual accuracy (89-90%) and joint accuracy (81%) reflects the compounding difficulty of selecting both services correctly.

**Noisy Set:** Ambiguous queries challenge the system, particularly for action selection (72%). Vague descriptions require significant inference.

**One-Shot Set:** The system maintains reasonable performance on rare functions, demonstrating robust generalization from contrastive learning.
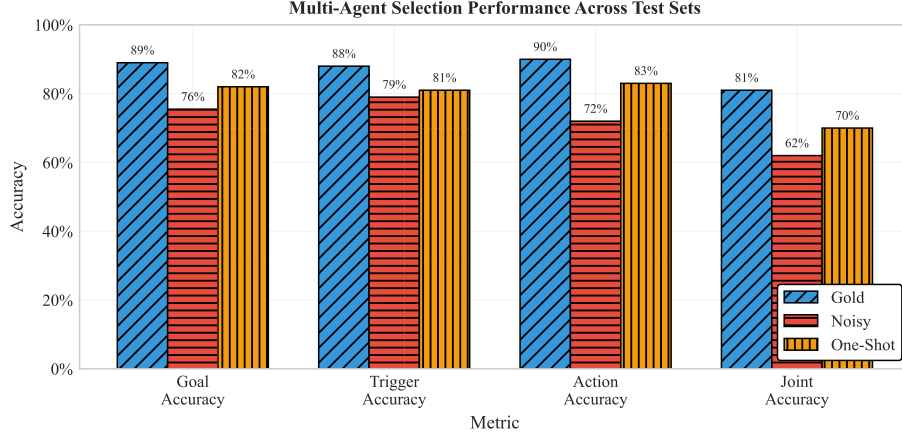
Fig. 8. Stage 2 multi-agent selection performance across evaluation sets. The system achieves 81% joint accuracy on Gold, with graceful degradation on challenging inputs.

*4.4.3 Generation Quality.* Figure 9 presents RAGAS [11] quality metrics for generated applet configurations.
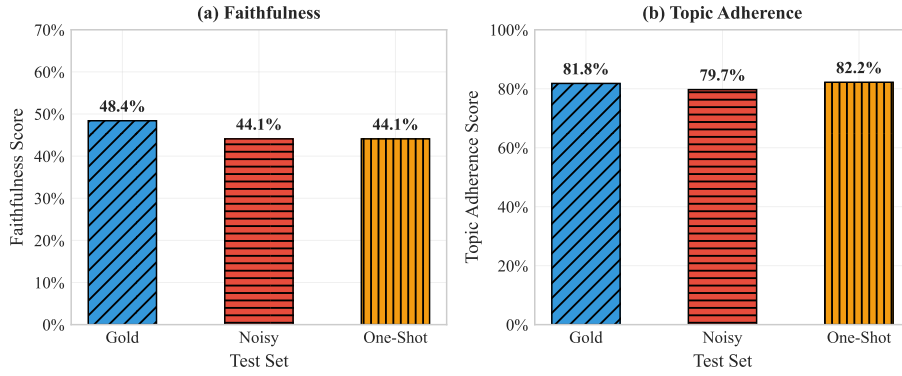


Fig. 9. Generation quality metrics across evaluation sets. Faithfulness (0.44–0.48) measures grounding in retrieved schemas; Topic Adherence (0.80–0.82) measures alignment with user's automation domain.

High topic adherence across all sets indicates that selected services consistently align with the user's stated automation domain. Moderate faithfulness scores reflect the challenging nature of generating precise field bindings that require inference beyond explicit context.

*4.4.4 Comprehensive Multi-Dataset Evaluation.* Figure 10 presents a comprehensive breakdown of system performance across all three evaluation sets, revealing both overall effectiveness and robustness characteristics.

**Gold Set Results:** Clear, well-formed queries enable high performance across all metrics. The system achieves 97% success rate (at least partial automation achieved), 89% goal accuracy (average function-level performance with partial credit), and 81% joint accuracy (both trigger and action functions correctly selected). Individual component accuracy reaches 88% for triggers and 90% for actions, with the gap to joint accuracy (81%) reflecting the compounding difficulty of perfect service-pair selection.
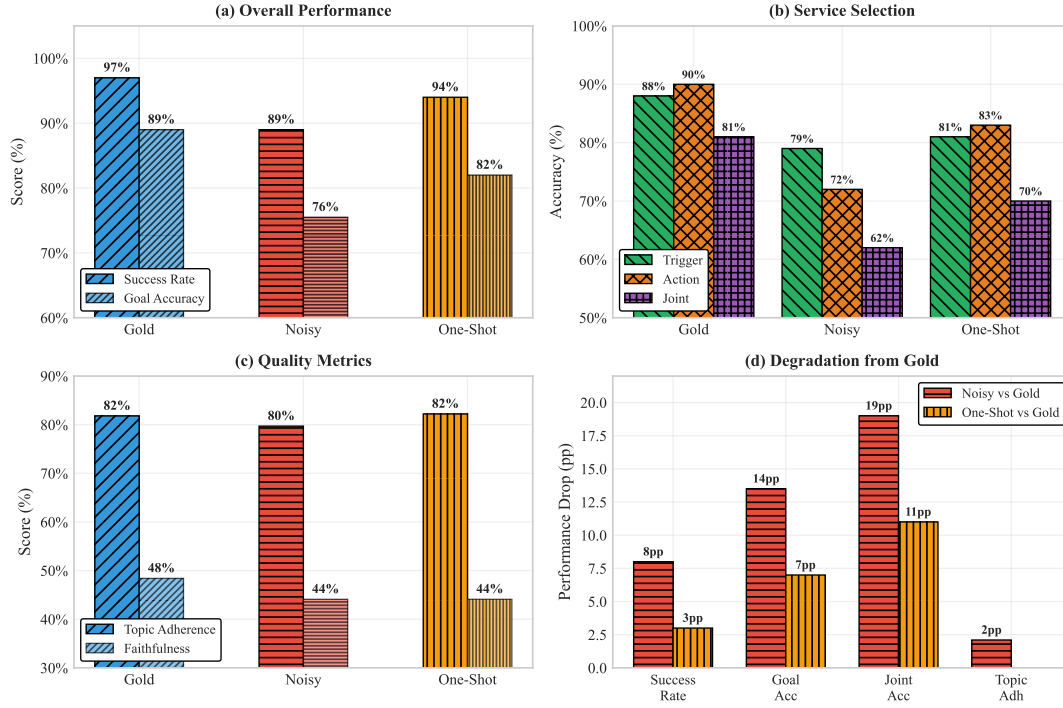
Fig. 10. Multi-agentic evaluation across test sets. **(a)** Overall performance: Gold achieves 97% success rate (at least one service correct) and 89% goal accuracy (average with partial credit); Noisy shows 8pp degradation in success rate; One-Shot maintains 94% success. The ordering Success Rate > Goal Accuracy reflects metric definitions: Success Rate = binary threshold (goal_accuracy ≥ 0.5), Goal Accuracy = arithmetic mean. **(b)** Service selection: Trigger accuracy (88%/79%/81%), Action accuracy (90%/72%/83%), Joint accuracy (81%/62%/70%) for Gold/Noisy/One-Shot respectively. **(c)** Quality metrics: Topic adherence remains stable (82%) across sets; Faithfulness moderate (44–48%). **(d)** Performance degradation: Noisy queries show largest drops in joint accuracy (19pp) while One-Shot maintains resilience (11pp drop).

**Noisy Set Results:** Vague or ambiguous descriptions present significant challenges. Success rate drops to 89% (8pp degradation from Gold), goal accuracy to 75.5% (13.5pp drop), and joint accuracy to 62% (19pp drop). The ordering Success Rate (89%) > Goal Accuracy (75.5%) > Joint Accuracy (62%) reflects the metric hierarchy: Success Rate requires only partial correctness (one service), Goal Accuracy averages partial credit, while Joint Accuracy requires both services correct. Action selection proves particularly challenging (72% accuracy vs 90% on Gold), as ambiguous queries require substantial inference to identify appropriate actions. Topic adherence remains stable at 79.7%, indicating robust domain alignment despite query ambiguity.

**One-Shot Set Results:** Queries involving rare functions demonstrate robust generalization from contrastive learning. The system maintains 94% success rate (3pp drop from Gold) and 82% goal accuracy (7pp drop), with 70% joint accuracy (11pp drop). Both trigger (81%) and action (83%) selection remain competitive, showing that layer freezing preserves sufficient pretrained semantic knowledge to handle low-resource scenarios.

**Quality Metrics Analysis:** Faithfulness scores range from 44% to 48% across all sets, reflecting the challenging nature of schema-grounded binding generation that requires inference beyond explicit context. Topic adherence remains

consistently high (80–82%), indicating that the multi-agent system reliably selects services within the user's intended automation domain across all test conditions.

## 4.5  Comparison with Related Work

We compare FARM against three established baselines for trigger-action programming: Liu et al. [28], Yusuf et al. [60], and Cimino et al. [5]. Table 5 presents quantitative results on a service-level prediction task across all three evaluation sets. Although FARM operates at function-level granularity with full schema specifications, for this table we map each predicted function to its parent platform service to enable fair comparison with baselines. To ensure a controlled comparison, we train and evaluate each baseline on the same service-level splits derived from our dataset.

| | Test Gold | | Test Noisy | | Test One-shot | |
|---|---|---|---|---|---|---|
| **Method** | Joint Acc | MRR@3 | Joint Acc | MRR@3 | Joint Acc | MRR@3 |
| LAM [28] | 0.24 | 0.27 | 0.23 | 0.26 | 0.12 | 0.19 |
| RecipeGen++ [60] | 0.31 | 0.35 | 0.29 | 0.34 | 0.17 | 0.24 |
| TARGE [5] | 0.58 | 0.63 | 0.39 | 0.44 | 0.45 | 0.49 |
| **FARM (Ours)** | **0.79** | **0.84** | **0.62** | **0.69** | **0.70** | **0.77** |

Table 5. Comparison with prior trigger-action programming approaches on a **service-level** benchmark. We derive service-level labels by mapping each trigger/action *function* in our dataset to its parent platform service. We re-ran LAM, RecipeGen++, and TARGE using the authors' released code on this service-level dataset, using identical train/validation/test splits across all methods. FARM operates at function granularity, but for this table we map FARM's predicted functions to their parent services to enable direct comparison. **Joint Accuracy** is computed on the top-1 predicted trigger-service and action-service pair. **MRR@3** is computed as the mean reciprocal rank of the ground-truth trigger-action *service pair* within the top-3 ranked pairs produced by each method.

## 4.6  Ablation Studies

We conduct ablation studies to understand the contribution of each system component.

*4.6.1  Layer Freezing Ablation.* Figure 11 compares full training against layer freezing on retrieval accuracy.

Layer freezing preserves the pretrained encoder's understanding in frozen layers (0–11) while training upper layers (12–23) for domain-specific patterns.

*4.6.2  LoRA vs. Layer Freezing.* We compare our layer freezing approach against Low-Rank Adaptation (LoRA) [14], the dominant parameter-efficient fine-tuning method. Figure 12 presents results across retrieval accuracy and parameter efficiency.

Layer freezing achieves 85% Joint R@5 with 18% trainable parameters, outperforming full fine-tuning (83% Joint R@5, 100% params) and all LoRA variants. LoRA's low-rank constraint limits its ability to adapt to the TAP domain, with increasing rank ($r=8{\rightarrow}32$) providing only marginal improvements. Layer freezing provides the best trade-off: highest retrieval performance with moderate parameter cost.

*4.6.3  Component Ablation.* **Component Contributions:**

(1) **Contrastive Training (+49 pts Joint R@1):** Improves retrieval top-1 dual success from near-random (7%) to viable (56%).
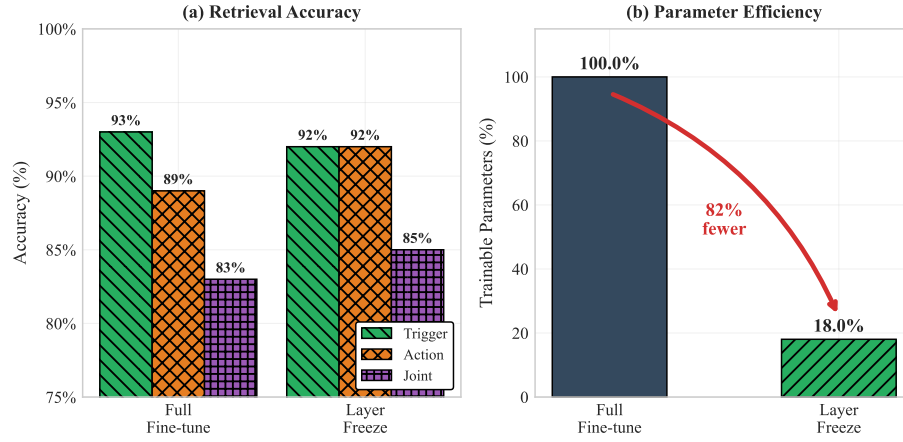
Fig. 11. Layer freezing ablation using R@5 metric (consistent with Figure 12). Layer freezing achieves slightly better joint R@5 accuracy (85%) with 92%/92% for trigger/action compared to full fine-tuning (83% joint) with 93%/89% for trigger/action, while using only 18% of trainable parameters. Note: R@5 is used here instead of R@1 to match the LoRA ablation evaluation protocol, which uses $k$=5 for retrieval evaluation.
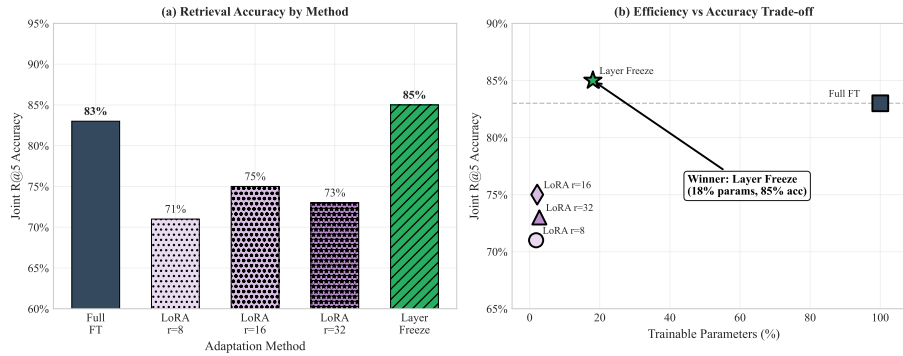


Fig. 12. LoRA vs. Layer Freezing ablation. (a) Joint R@5 accuracy by adaptation method. Layer freezing achieves the highest accuracy (85%) while LoRA variants underperform (69–73%). (b) Efficiency vs. accuracy trade-off showing layer freezing provides the best balance of parameter efficiency (18%) and retrieval performance.

(2) **Layer Freezing (+2 pts Joint R@1):** Improves retrieval top-1 dual success from 56% to 58% while using only 18% trainable parameters.

(3) **Multi-Agent Selection (+23 pts Joint Accuracy):** This substantial improvement (58% → 81%) stems from expanding the search space: naive selection evaluates only 1 pair (rank-1 from each encoder), while multi-agent selection evaluates up to $k^2 = 25$ pairs through cross-compatibility scoring (Fig. 13, 14).
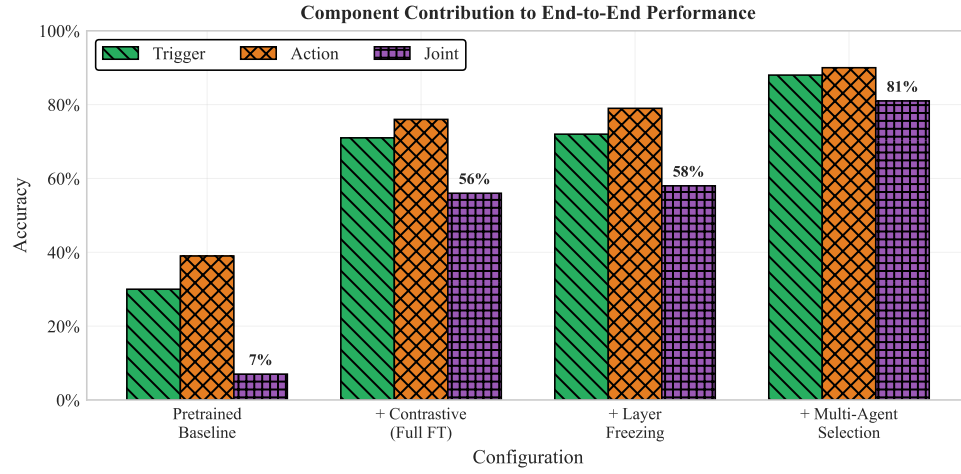
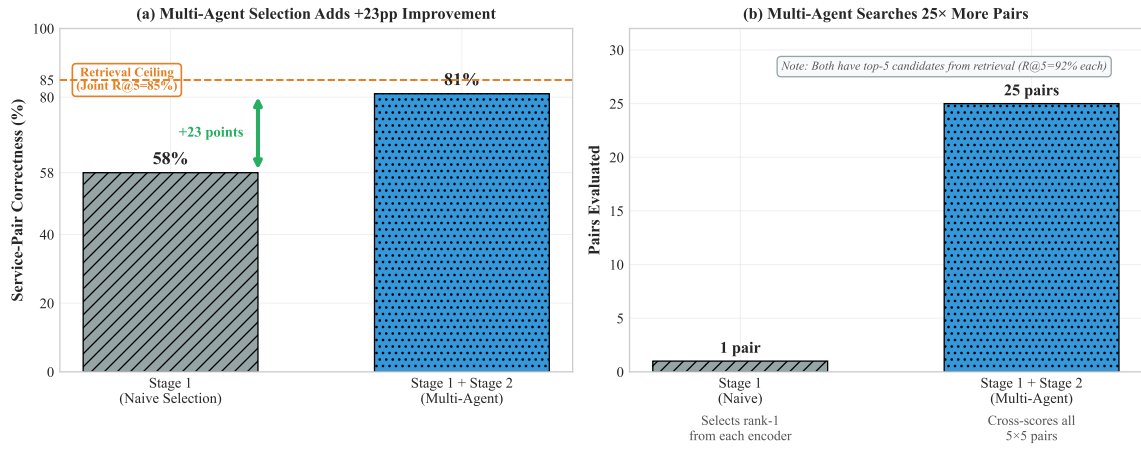Fig. 13. Component ablation showing incremental contributions.



Fig. 14. Stage 2 improvement mechanism. (a) Multi-agent selection achieves 81% joint accuracy, a +23 point improvement over naive rank-1 selection (58%), approaching the 85% retrieval ceiling (Joint R@5). (b) The improvement stems from evaluating 25 trigger-action pairs (5×5) with LLM-based cross-compatibility scoring, compared to naive selection which evaluates only 1 pair (rank-1 from each encoder).

*4.6.4  LLM Selection Ablation.* We evaluate the impact of different large language models on the multi-agent selection pipeline. Table 6 presents generation quality metrics for various LLMs suitable for tool calling and agentic workflows.

IBM Granite 4.0 achieves the highest scores on both metrics, with +9–17 points improvement in Faithfulness and +6–14 points improvement in Topic Adherence compared to alternatives. This performance advantage stems from Granite's optimization for enterprise agentic workflows with native tool calling support, enabling more reliable structured JSON output generation and better instruction following for our agent prompting strategy.

| LLM | Faithfulness | Topic Adherence |
|---|---|---|
| LLaMA 3 70B [33] | 0.38 | 0.74 |
| Qwen2 72B [57] | 0.41 | 0.76 |
| Mistral Large [20] | 0.36 | 0.71 |
| DeepSeek-Coder [12] | 0.33 | 0.68 |
| **IBM Granite 4.0 Small [17]** | **0.50** | **0.82** |

Table 6. LLM Ablation for Multi-Agent Selection. Faithfulness measures factual grounding in retrieved schemas; Topic Adherence measures alignment with automation domains. Higher is better. Note: These metrics represent best-case LLM performance on the Gold set; end-to-end system faithfulness across all datasets (Gold/Noisy/One-Shot) ranges from 0.44–0.48 as shown in Fig. 9.

*4.6.5 Base Encoder Selection.* To evaluate the impact of base encoder choice on our contrastive learning framework, we compare performance across four state-of-the-art embedding models with diverse architectures and training objectives. We apply our layer freezing strategy (freezing layers 0–11, training 12–23) to EmbeddingGemma and compare against alternative models to assess both absolute performance and the generalizability of our approach.

EmbeddingGemma with layer freezing achieves the best performance across all datasets and configurations (58% Gold, 50% Noisy, 54% One-Shot Joint R@1; Trigger R@5 = 92%, Action R@5 = 92% on Gold), validating our ablation findings and demonstrating consistent robustness to dataset quality variations. Different base encoders respond differently to training strategies: ModernBERT models benefit from full training (e.g., ModernBERT-large on Gold improves from 50% to 52%), whereas BGE and E5 models achieve higher performance under layer freezing, similar to EmbeddingGemma. This indicates that layer freezing is particularly effective for encoders with strong multilingual pretraining (EmbeddingGemma, BGE, E5), while architectures optimized for long-context processing (ModernBERT) benefit from full parameter updates. Performance degradation from Gold to Noisy datasets follows similar trends across models (7–10 percentage point drops), indicating that dataset quality affects all architectures comparably. Notably, even under their preferred training strategies, no alternative model surpasses EmbeddingGemma with layer freezing on any dataset split, confirming both the effectiveness of our approach and the importance of base encoder selection for trigger-action retrieval (Table 7).

| Base Model | Params | Layer Freezing (18%) | | | | Full Training (100%) | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | *Joint R@1* | | | R@5 | *Joint R@1* | | | R@5 |
| | | Gold | Noisy | One-Shot | | Gold | Noisy | One-Shot | |
| **EmbeddingGemma [39]** | 307M | **0.58** | **0.50** | **0.54** | **0.92** | **0.56** | 0.48 | **0.52** | 0.91 |
| ModernBERT-base [51] | 149M | 0.45 | 0.38 | 0.42 | 0.85 | 0.47 | 0.40 | 0.44 | 0.87 |
| ModernBERT-large [51] | 395M | 0.50 | 0.42 | 0.48 | 0.89 | 0.52 | 0.44 | 0.50 | **0.91** |
| BGE-base-en-v1.5 [56] | 109M | 0.48 | 0.40 | 0.45 | 0.88 | 0.46 | 0.38 | 0.43 | 0.86 |
| E5-large-v2 [48] | 335M | 0.52 | 0.44 | 0.49 | 0.90 | 0.50 | **0.49** | 0.47 | 0.89 |

Table 7. Base encoder comparison.

## 4.7 Efficiency Analysis

We analyze FARM's computational efficiency across retrieval, selection, and training phases. Table 8 summarizes the latency and resource requirements for each system component. All experiments conducted on NVIDIA Tesla V100-SXM2-32GB GPUs (32GB memory).

Table 8. System Efficiency Metrics

| Metric | Value |
|---|---|
| *Stage 1 (Retrieval)* | |
| Encoding Latency | ~50ms per query |
| Vector Search (Qdrant) | ~5ms for top-5 |
| *Stage 2 (Selection)* | |
| LLM Calls per Query | 2–4 |
| Average Total Latency | ~3s (local inference) |
| LLM Configuration | 3× V100 GPUs |
| Fallback Rate | <5% |
| *Training* | |
| Encoder Training Time | ~20 minutes |
| Training Hardware | Single V100 GPU |
| Training Data Size | 12,652 pairs |

The two-stage architecture provides computational efficiency: Stage 1 reduces the search space from 2.2M candidate pairs to $k^2 = 25$ combinations (with $k = 5$), enabling Stage 2's LLM reasoning to focus on high-quality candidates. This design is necessary because naive selection (taking rank-1 from each encoder, evaluating only 1 pair) achieves just 58% joint accuracy, whereas evaluating the full $k \times k$ search space with multi-agent reasoning achieves 81%. The system uses $k = 5$ throughout, matching the R@5 evaluation metrics reported in ablation studies.

## 4.8 Discussion

*4.8.1 Why Two Stages?* Our two-stage architecture addresses a fundamental trade-off between efficiency and accuracy. Stage 1 provides fast (~55ms) retrieval that reduces the search space from 2.2 million possible pairs to a manageable candidate set. However, *retrieval ranking alone is insufficient for accurate selection.*

To understand why, consider the selection problem: given top-$k$ candidates from each encoder, we must identify the correct trigger-action pair. A naive approach would take the rank-1 result from each encoder (evaluating only 1 pair), but this achieves just 58% joint accuracy because it requires *both* encoders to rank their correct items first simultaneously. Individual R@1 scores (Trigger: 72%, Action: 79%) compound to approximately 58% joint success, and when either encoder ranks the correct item at position 2–5, the naive approach fails.

Stage 2 addresses this limitation by evaluating the cross-product of candidates: with $k = 5$, the system considers up to $k^2 = 25$ trigger-action pairs (Fig. 14b). Multi-agent selection uses LLM reasoning to score pairs based on schema compatibility, ingredient-to-field matching, and intent alignment—factors invisible to embedding-based retrieval. This expanded search achieves 81% joint accuracy (Fig. 14a), approaching the retrieval ceiling of 85% (Joint R@5). The 4-point

gap between achieved (81%) and ceiling (85%) reflects LLM reasoning errors on queries where correct pairs exist in the candidate set but are incorrectly scored.

Put quantitatively: Stage 2 succeeds on 95% of retrievable queries (81% / 85%), demonstrating that the multi-agent search effectively exploits Stage 1's high-recall candidate sets. Neither stage alone could achieve this: Stage 1 lacks reasoning capabilities, while Stage 2 requires the computational efficiency of Stage 1 to focus LLM evaluation on high-quality candidates.

*4.8.2 Prompt-Based Multi-Agent Design.* Stage 2 operates through *prompting only*—no task-specific fine-tuning. This provides flexibility (agents modified via prompts), interpretability (explicit reasoning traces), and generalization (pre-trained LLM knowledge handles edge cases). We mitigate potential output inconsistency through structured schemas and agreement-based override mechanisms.

*4.8.3 Limitations.*

(1) **Retrieval Ceiling:** At $k$=5, individual recall reaches 92% and Joint R@5 reaches 85%. Stage 2 cannot recover if the correct service is not retrieved in Stage 1's top-5, establishing an 85% theoretical upper bound on joint accuracy.
(2) **Complex Queries:** Multi-intent queries (e.g., "turn off lights AND lock door") require decomposition not yet implemented.
(3) **Faithfulness:** Moderate scores (0.44–0.48) indicate room for improvement in explicit schema grounding.
(4) **LLM Dependency:** Stage 2 quality depends on the underlying LLM capability; smaller models may underperform.

Several promising directions emerge from this work:

(1) **Execution Validation:** Developing automated evaluation of ingredient-to-field binding correctness through integration with live IFTTT API endpoints. This would enable measuring end-to-end applet execution success rates and identifying failure modes in binding generation.
(2) **Query Decomposition:** Implementing multi-intent query parsing to handle complex automation requests requiring multiple trigger-action pairs.
(3) **Reinforcement Learning from Human Feedback:** Training a cross-encoder reranker using Group Relative Policy Optimization (GRPO) [41], which eliminates the need for a critic model by estimating baselines from group scores. This would enable learning nuanced scoring functions from user acceptance/rejection signals while remaining memory-efficient. GRPO's compatibility with existing LLM infrastructure makes it particularly suitable for improving Stage 2 selection quality through iterative refinement.
(4) **Cross-Platform Generalization:** Extending evaluation to other TAP platforms (Zapier, Make, Home Assistant) to validate approach generalizability across different service ecosystems and schema formats.
(5) **User Studies:** Conducting human evaluation to assess practical utility and user satisfaction with generated applets in real-world scenarios, including longitudinal studies of applet usage patterns.

**Acknowledgment**

## References

[1] I. Beltagy, K. Lo, and E. H. Hovy. 2016. Improved Semantic Parsers for If-Then Statements. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 1–10. doi:10.18653/v1/P16-1069

[2] Dan Biderman, Jacob Portes, Jose Javier Gonzalez Ortiz, Mansheej Paul, Philip Greengard, Connor Jennings, Daniel King, Sam Havens, Vitaliy Chiley, Jonathan Frankle, Cody Blakeney, and John P. Cunningham. 2024. LoRA Learns Less and Forgets Less. *arXiv preprint arXiv:2405.09673* (2024). https://arxiv.org/abs/2405.09673

[3] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, et al. 2020. Language Models are Few-Shot Learners. In *Advances in Neural Information Processing Systems (NeurIPS)*, Vol. 33. 1877–1901.

[4] Wei-Lin Chiang, Liang Zheng, Ying Sheng, Soheil Zhuang, Ziqi Wu, Yu Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, and Eric P. Xing. 2023. Chatbot Arena: An Open Platform for Evaluating LLMs by Human and AI Judges. *arXiv preprint arXiv:2305.14387* (2023).

[5] Gaetano Cimino, Vincenzo Deufemia, and Mattia Limone. 2025. IoT Rule Generation With Cross-View Contrastive Learning and Perplexity-Based Ranking. *IEEE Internet of Things Journal* 12, 17 (2025), 35828–35847.

[6] Fulvio Corno, L. De Russis, and A. Monge Roffarello. 2020. TapRec: Supporting the Composition of Trigger-Action Rules Through Dynamic Recommendations. In *Proceedings of the 25th International Conference on Intelligent User Interfaces*. 579–588.

[7] Fulvio Corno, L. De Russis, and A. Monge Roffarello. 2019. RecRules: Recommending If-Then Rules for End-User Development. *ACM Transactions on Intelligent Systems and Technology* 10, 5 (2019), 1–27.

[8] F. Corno, L. De Russis, and A. Monge Roffarello. 2021. From Users' Intentions to If-Then Rules in the Internet of Things. *ACM Transactions on Information Systems* 39, 4 (2021), 1–33.

[9] Matias Delange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Ales Leonardis, Gregory Slabaugh, and Tinne Tuytelaars. 2021. A Continual Learning Survey: Defying Forgetting in Classification Tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2021).

[10] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. QLoRA: Efficient Finetuning of Quantized LLMs. *arXiv preprint arXiv:2305.14314* (2023). https://arxiv.org/abs/2305.14314

[11] Shahul Es, Jithin James, Luis Espinosa-Anke, and Steven Schockaert. 2024. RAGAS: Automated Evaluation of Retrieval Augmented Generation. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations*. Association for Computational Linguistics, St. Julians, Malta, 150–158. doi:10.18653/v1/2024.eacl-demo.16

[12] Daya Guo et al. 2024. DeepSeek-Coder: When the Large Language Model Meets Programming. *arXiv* (2024). https://arxiv.org/abs/2401.14196

[13] Kai-Hsiang Hsu, Yu-Hsi Chiang, and Hsu-Chun Hsiao. 2019. SafeChain: Securing Trigger-Action Programming from Attack Chains. *IEEE Transactions on Information Forensics and Security* 14, 10 (2019), 2607–2622. doi:10.1109/TIFS.2019.2907062

[14] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu Chen. 2022. LoRA: Low-Rank Adaptation of Large Language Models. *arXiv preprint arXiv:2106.09685* (2022). https://arxiv.org/abs/2106.09685

[15] Justin Huang and Maya Cakmak. 2015. Supporting Mental Model Accuracy in Trigger-Action Programming. In *Proceedings of the 2015 ACM International Joint Conference on Pervasive and Ubiquitous Computing (UbiComp)*. ACM, 215–225. doi:10.1145/2750858.2805830

[16] Zijun Huang and Jiangfeng Li. 2023. TAP-AHGNN: An Attention-Based Heterogeneous Graph Neural Network for Service Recommendation on Trigger-Action Programming Platform. In *Proceedings of the 2023 International Conference on Data Science and Information Technology*. Springer, 157–170. doi:10.1007/978-981-99-4752-2_12

[17] IBM. 2025. IBM Granite 4.0: Hyper-Efficient, High-Performance Hybrid Models for Enterprise. https://www.ibm.com//new/announcements/ibm-granite-4-0-hyper-efficient-high-performance-hybrid-models. Official IBM announcement page.

[18] Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. 2022. Atlas: Few-shot Learning with Retrieval Augmented Language Models. *arXiv preprint arXiv:2208.03299* (2022). https://arxiv.org/abs/2208.03299

[19] Ganesh Jawahar, Benoît Sagot, and Djamé Seddah. 2019. What Does BERT Learn About the Structure of Language?. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. 3651–3657.

[20] Albert Jiang et al. 2023. Mistral 7B. *arXiv* (2023). https://arxiv.org/abs/2310.06825

[21] Seonghyeon Kim et al. 2025. Leveraging LLMs for Efficient and Personalized Smart Home Automation. *arXiv preprint arXiv:2601.04680* (2025). Demonstrates LLM failures on underspecified IoT instructions.

[22] E. King. 2023. Sasha: Creative Goal-Oriented Reasoning in Smart Homes with Large Language Models. arXiv:2305.09802 [cs.HC] https://arxiv.org/abs/2305.09802

[23] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. 2017. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences* 114, 13 (2017), 3521–3526.

[24] Zhejun Kuang, Yusheng Zhu, Dawen Sun, Jian Zhao, Yongheng Xing, Feng Wang, and Lei Sun. 2025. SAFE-TAP: Semantic-aware and fused embedding for TAP rule security detection. *Neurocomputing* 656 (2025), 131529. doi:10.1016/j.neucom.2025.131529

[25] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. In *Advances in Neural Information Processing Systems*, Vol. 33. https://arxiv.org/abs/2005.11401

[26] Zhizhong Li and Derek Hoiem. 2016. Learning without Forgetting. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 614–629.

[27] Alisa Liu, Zhaofeng Wu, Julian Michael, Alane Suhr, Peter West, Alexander Koller, Swabha Swayamdipta, Noah Smith, and Yejin Choi. 2023. We're Afraid Language Models Aren't Modeling Ambiguity. (2023), 790–807.

[28] Chang Liu, Xinyun Chen, Eui Chul Shin, Minghui Chen, and Dawn Song. 2016. Latent Attention for If-Then Program Synthesis. In *Advances in Neural Information Processing Systems*, Vol. 29.

[29] Liwei Liu, Wei Chen, Lu Liu, Kangkang Zhang, Jun Wei, and Yan Yang. 2021. TAGen: Generating Trigger-Action Rules for Smart Homes by Mining Event Traces. In *Proceedings of the 19th International Conference on Service-Oriented Computing (ICSOC 2021)*. Springer, 652–662. doi:10.1007/978-3-030-91431-8_41

[30] Liwei Liu, Wei Chen, Tao Wang, Wei Wang, Guoquan Wu, and Jun Wei. 2023. Generating Scenario-Centric TAP Rules for Smart Homes by Mining Historical Event Logs. In *2023 IEEE International Conference on Web Services (ICWS)*. IEEE, Conference Location: Chicago, IL, USA. doi:10.1109/ICWS60048.2023.00015

[31] Nelson F. Liu, Matt Gardner, Yonatan Belinkov, Matthew Peters, and Noah A. Smith. 2019. Linguistic Knowledge and Transferability of Contextual Representations. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 1073–1094.

[32] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. 2008. *Introduction to Information Retrieval*. Cambridge University Press.

[33] Meta AI. 2024. The LLaMA 3 Model Family. *arXiv* (2024). https://arxiv.org/abs/2404.09437

[34] Xianghang Mi, Feng Qian, Ying Zhang, and XiaoFeng Wang. 2017. An Empirical Characterization of IFTTT: Ecosystem, Usage, and Performance. In *Proceedings of the 2017 Internet Measurement Conference (IMC)*. ACM, 398–404. doi:10.1145/3131365.3131369

[35] Rodrigo Nogueira and Kyunghyun Cho. 2019. Passage Re-ranking with BERT. *arXiv preprint arXiv:1901.04085* (2019). https://arxiv.org/abs/1901.04085

[36] Chris Quirk, Raymond Mooney, and Michel Galley. 2015. Language to Code: Learning Semantic Parsers for If-Then Recipes. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. 878–888.

[37] RAGAS Contributors. 2024. RAGAS Documentation: Metrics for Retrieval-Augmented Generation. https://docs.ragas.io/en/stable/. Accessed: 2025-12-19.

[38] Amir Rahmati, Earlence Fernandes, Jaeyeon Jung, and Atul Prakash. 2017. IFTTT vs. Zapier: A Comparative Study of Trigger-Action Programming Frameworks. *arXiv preprint arXiv:1709.02788* (2017).

[39] Henrique Schechter Vera, Sahil Dua, Biao Zhang, Daniel Salz, Ryan Mullins, Sindhu Raghuram Panyam, Sara Smoot, Iftekhar Naim, Joe Zou, Feiyang Chen, Daniel Cer, Alice Lisak, Min Choi, Lucas Gonzalez, Omar Sanseviero, Glenn Cameron, Ian Ballantyne, Kat Black, Kaifeng Chen, Weiyi Wang, Zhe Li, Gus Martins, Jinhyuk Lee, Mark Sherwood, Juyeong Ji, Renjie Wu, Jingxiao Zheng, Jyotinder Singh, Abheesht Sharma, Divyashree Sreepathihalli, Aashi Jain, Adham Elarabawy, AJ Co, Andreas Doumanoglou, Babak Samari, Ben Hora, Brian Potetz, Dahun Kim, Enrique Alfonseca, Fedor Moiseev, Feng Han, Frank Palma Gomez, Gustavo Hernández Ábrego, Hesen Zhang, Hui Hui, Jay Han, Karan Gill, Ke Chen, Koert Chen, Madhuri Shanbhogue, Michael Boratko, Paul Suganthan, Sai Meher Karthik Duddu, Sandeep Mariserla, Setareh Ariafar, Shanfeng Zhang, Shijie Zhang, Simon Baumgartner, Sonam Goenka, Steve Qiu, Tanmaya Dabral, Trevor Walker, Vikram Rao, Waleed Khawaja, Wenlei Zhou, Xiaoqi Ren, Ye Xia, Yichang Chen, Yi-Ting Chen, Zhe Dong, Zhongli Ding, Francesco Visin, Gaël Liu, Jiageng Zhang, Kathleen Kenealy, Michelle Casbon, Ravin Kumar, Thomas Mesnard, Zach Gleicher, Cormac Brick, Olivier Lacombe, Adam Roberts, Qin Yin, Yunhsuan Sung, Raphael Hoffmann, Tris Warkentin, Armand Joulin, Tom Duerig, and Mojtaba Seyedhosseini. 2025. EmbeddingGemma: Powerful and Lightweight Text Representations. *arXiv preprint arXiv:2509.20354* (2025). https://arxiv.org/abs/2509.20354

[40] Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language Models Can Teach Themselves to Use Tools. In *Advances in Neural Information Processing Systems (NeurIPS)*. Proceedings publication.

[41] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. 2024. DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models. *arXiv preprint arXiv:2402.03300* (2024).

[42] Yingtian Shi et al. 2024. Bridging the Gap Between Natural User Expression with Complex Automation Programming in Smart Homes. *arXiv preprint arXiv:2408.12687* (2024).

[43] Ian Tenney, Dipanjan Das, and Ellie Pavlick. 2019. BERT Rediscovers the Classical NLP Pipeline. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. 4593–4601.

[44] Blase Ur, Melwyn Pak Yong Ho, Stephen Brawner, Jiyun Lee, Sarah Mennicken, Noah Picard, Diane Schulze, and Michael L. Littman. 2016. Trigger-Action Programming in the Wild: An Analysis of 200,000 IFTTT Recipes. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*. Association for Computing Machinery, New York, NY, USA. doi:10.1145/2858036.2858556

[45] Blase Ur, Elyse McManus, Melwyn Pak Yong Ho, and Michael L. Littman. 2014. Practical Trigger-Action Programming in the Smart Home. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. Association for Computing Machinery, New York, NY, USA.

doi:10.1145/2556288.2557420

[46] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation Learning with Contrastive Predictive Coding. In *arXiv preprint arXiv:1807.03748*.

[47] Ellen M. Voorhees. 1999. The TREC-8 Question Answering Track Report. In *Proceedings of the Eighth Text Retrieval Conference (TREC-8)*.

[48] Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. 2022. Text Embeddings by Weakly-Supervised Contrastive Pre-training. *arXiv preprint arXiv:2212.03533* (2022).

[49] Qi Wang, Pubali Datta, Wei Yang, Si Liu, Adam Bates, and Carl A. Gunter. 2019. Charting the Attack Surface of Trigger-Action IoT Platforms. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security (CCS)*. ACM, 1439–1453.

[50] Xuezhi Wang and Denny Zhou. 2024. Chain-of-Thought Reasoning Without Prompting. *arXiv preprint arXiv:2402.10200* (2024). arXiv:2402.10200 [cs.CL] doi:10.48550/arXiv.2402.10200

[51] Benjamin Warner, Antoine Chaffin, Benjamin Clavié, Orion Weller, Oskar Hallström, Said Taghadouini, Alexis Gallagher, Raja Biswas, Faisal Ladhak, Tom Aarsen, Nathan Cooper, Griffin Adams, Jeremy Howard, and Iacopo Poli. 2024. Smarter, Better, Faster, Longer: A Modern Bidirectional Encoder for Fast, Memory Efficient, and Long Context Finetuning and Inference. *arXiv preprint arXiv:2412.13663* (2024).

[52] Jason Wei, Yi Tay, Rishi Bommasani, et al. 2022. Emergent Abilities of Large Language Models. *Transactions on Machine Learning Research (TMLR)* (2022).

[53] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc Le, and Denny Zhou. 2022. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. *arXiv preprint arXiv:2201.11903* (2022). arXiv:2201.11903 [cs.CL] doi:10.48550/arXiv.2201.11903

[54] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. In *Advances in Neural Information Processing Systems*, Vol. 35. Curran Associates, Inc., 24824–24837.

[55] G. Wu et al. 2025. User Intention Prediction for Trigger-Action Programming via Multi-view Representation Learning. *Expert Systems with Applications* (2025). Article information available via ScienceDirect.

[56] Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas Muennighoff. 2023. C-Pack: Packaged Resources To Advance General Chinese Embedding. *arXiv preprint arXiv:2309.07597* (2023).

[57] An Yang et al. 2024. Qwen2 Technical Report. *arXiv* (2024). https://arxiv.org/abs/2407.10671

[58] Zijun Yao, Jiangfeng Li, Huijuan Zhang, Chenxi Zhang, and Gang Yu. 2019. Interactive User Interface Design for Trigger-Action Programming. In *Proceedings of the International Conference on Smart Internet of Things (SmartIoT)*.

[59] Y. Yoon. 2016. Performance analysis of CRF-based learning for processing WoT application requests expressed in natural language. *SpringerPlus* 5 (2016), 1324. doi:10.1186/s40064-016-3012-9

[60] I. N. B. Yusuf, D. B. A. Jamal, Lingxiao Jiang, and David Lo. 2022. RecipeGen++: An Automated Trigger-Action Programs Generator. In *Proceedings of the ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 1672–1676.

[61] I. N. B. Yusuf, L. Jiang, and D. Lo. 2022. Accurate Generation of Trigger-Action Programs with Domain-Adapted Sequence-to-Sequence Learning. In *Proceedings of the 30th International Conference on Program Comprehension (ICPC)*. 99–110.

[62] Friedemann Zenke, Ben Poole, and Surya Ganguli. 2017. Continual Learning Through Synaptic Intelligence. In *Proceedings of the 34th International Conference on Machine Learning (ICML)*. 3987–3995.

[63] Lefan Zhang, Cyrus Zhou, Hannah Yu, Blase Ur, and Michael L. Littman. 2023. Helping Users Debug Trigger-Action Programs. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* (2023). doi:10.1145/3569506

[64] Liang Zheng, Wei-Lin Chiang, Ying Sheng, Soheil Zhuang, Ziqi Wu, Yu Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, et al. 2023. Judging LLMs by LLMs: A Benchmarking Framework for Evaluating Large Language Models. *arXiv preprint arXiv:2306.05685* (2023).