



**UNIMORE**

UNIVERSITÀ DEGLI STUDI DI  
MODENA E REGGIO EMILIA

Dipartimento di Scienze Fisiche, Informatiche e Matematiche

---

**Corso di Laurea Magistrale in Informatica**

# **RaTA-Tool: Retrieval-based Tool Selection with Multimodal Large Language Models**

**Candidato:**

Gabriele Mattioli

**Relatori:**

Prof. Lorenzo Baraldi

Prof.ssa Marcella Cornia

**Correlatori:**

Dott. Lorenzo Baraldi

Dott.ssa Sara Sarto

Anno Accademico 2024/2025



## Riassunto in lingua italiana

I recenti progressi dei Large Language Models (LLMs) hanno stravolto numerosi settori, evolvendosi da modelli specializzati nella generazione di testo a strumenti general-purpose capaci di ragionare. La loro influenza pervasiva nella vita odierna rimarca un significativo cambio di traiettoria dell'innovazione tecnologica, definendo nuovi standard nella collaborazione tra uomo e macchina. Nonostante tali capacità, i Large Language Models sono ancora principalmente utilizzati come strumenti di conversazione, che richiedono la compilazione di un prompt con una richiesta ben precisa e producono in output una risposta generata token per token. Allo stesso tempo, la ricerca scientifica si sta rapidamente muovendo verso un nuovo paradigma: l'Agentic AI; l'obiettivo è quello di equipaggiare i Large Language Models e altri modelli di AI con abilità di pianificazione, utilizzo di tool esterni e collaborazione con altri modelli, allo scopo di assolvere a compiti sfaccettati e composti da più fasi, senza la necessità di un intervento umano costante. In tale contesto, questa tesi esamina le capacità dei Multimodal Large Language Models (MLLMs) di usare tool esterni al fine di affrontare task multimodali complessi. Le soluzioni esistenti sono generalmente limitate a input testuali e faticano a comprendere istruzioni multimodali. Esse tendono anche a funzionare solamente in impostazioni sperimentali *closed-world*, ovvero non sono in grado di generalizzare a tool non osservati durante la fase di training. Questi problemi li rendono significativamente meno utili per applicazioni del mondo reale. Il framework sviluppato, RaTA-Tool, converte input e tool multimodali in descrizioni strutturate, e trova il tool migliore per un task comparando quanto la richiesta dell'utente corrisponda alla descrizione di ciascun tool. Inoltre, si propone un dataset per l'utilizzo di tool multimodali in un setting *open-world*, che contiene descrizioni strutturate di tool generate a partire dalle model card di HuggingFace.



# Abstract

Recent advancements in Large Language Models (LLMs) have revolutionized numerous domains, evolving from specialized text generators into general-purpose reasoning engines. Their pervasive influence on contemporary life highlights a significant shift in the trajectory of technological innovation, establishing a new standard for human-machine collaboration. Despite these abilities, LLMs are still mainly used as conversational instruments, prompting requests and waiting for a token by token generated response. Meanwhile, research is quickly moving towards a new paradigm: Agentic AI. The goal is to give LLMs and other AI models the ability to plan, use external tools and collaborate with other models, in order to accomplish multi-step multifaceted tasks without the need of constant human intervention. Within this framework, this thesis examines how Multimodal Large Language Models (MLLMs) use external tools to tackle complex multimodal tasks. Current solutions are usually limited to text-only inputs and struggle to understand multimodal instructions. They also tend to work only in closed-world settings, meaning they are not able to generalize to new tools not seen during training. These issues make them much less useful for real-world applications. The framework developed, RaTA-Tool, turns multimodal input and tools into structured descriptions, and finds the best tool for a task by comparing how well the user's request matches each tool's description. Furthermore, a custom dataset for open-world multimodal tool use is proposed, which contains structured tool descriptions derived from Hugging Face model cards.



# Ringraziamenti



# Contents

<b>Introduction</b>	<b>1</b>
<b>1 Background</b>	<b>3</b>
1.1 Neural networks . . . . .	3
1.2 Multilayer perceptron . . . . .	4
1.3 Transformers . . . . .	5
1.3.1 Attention . . . . .	6
Masked attention . . . . .	7
Multi-head attention . . . . .	8
1.3.2 Fully connected feed-forward network . . . . .	8
1.3.3 Transformer architecture . . . . .	9
1.4 Large Language Models (LLMs) . . . . .	10
1.4.1 Tokenization . . . . .	11
1.4.2 Embedding layer . . . . .	11
1.4.3 Training . . . . .	12
Pre-training . . . . .	12
Supervised Fine-Tuning (SFT) . . . . .	12
Reinforcement Learning from Human Feedback (RLHF) . . . . .	12
1.5 Multimodal Large Language Models (MLLMs) . . . . .	13
1.5.1 Vision Transformer (ViT) . . . . .	14
Classification . . . . .	14

1.5.2	Contrastive Language-Image Pre-training (CLIP)	15
1.6	Parameter Efficient Fine-Tuning (PEFT)	16
1.6.1	LoRA	17
1.7	Direct Preference Optimization (DPO)	18
1.8	Agentic AI	19
1.8.1	Tool learning	19
	Tool categorization	20
	Tool learning framework	21
1.8.2	Multimodal tool learning	22
<b>2</b>	<b>Proposed method</b>	<b>25</b>
2.1	Dataset creation	26
2.1.1	ToolMMBench	26
2.1.2	Our dataset	27
	Tools characterization	27
	Splitting, merging and cleaning	28
2.2	Pipeline	29
2.2.1	Query structured description	30
2.2.2	Retrieval	30
2.2.3	Supervised Fine-Tuning (SFT)	31
2.3	Direct Preference Optimization (DPO)	32
2.3.1	DPO dataset	32
<b>3</b>	<b>Experiments</b>	<b>35</b>
3.1	Experimental setup	35
3.1.1	SFT dataset	35
3.1.2	DPO dataset	36
3.2	Training strategies	36
3.2.1	Supervised Fine-Tuning (SFT)	36
3.2.2	Direct Preference Optimization (DPO)	38
3.3	Inference paradigms	38

---

3.3.1	Query task description generation . . . . .	38
	Zero-shot . . . . .	38
	Few-shot . . . . .	38
	SFT and DPO . . . . .	38
3.3.2	Retrieval . . . . .	39
3.4	Experimental results . . . . .	39
3.4.1	MLLM backbone and Embedding backbone . . . . .	40
	Raw models: few-shot inference . . . . .	40
	SFT models . . . . .	41
	DPO models . . . . .	41
	Overall considerations . . . . .	41
3.4.2	Description formats . . . . .	41
3.4.3	In-context learning and supervised fine-tuning . . . . .	43
3.5	Qualitative results . . . . .	44
<b>4</b>	<b>Conclusions</b>	<b>45</b>
4.1	Key findings and implications . . . . .	45
4.2	Limitations and future work . . . . .	46
4.3	Final remarks . . . . .	46
<b>Appendix</b>		
<b>A</b>	<b>Qualitative examples</b>	<b>47</b>
<b>B</b>	<b>Prompts</b>	<b>51</b>
<b>C</b>	<b>Publications</b>	<b>57</b>
	<b>Bibliography</b>	<b>76</b>



# List of Figures

1.1	Architecture of an MLP . . . . .	5
1.2	Scaled masked attention. Figure from [1]. . . . .	6
1.3	Multi-head attention. Figure from [1]. . . . .	9
1.4	Vision Transformer (ViT) architecture. Figure from [3]. . . . .	15
1.5	Direct Preference Optimization (DPO) vs. Reinforcement Learning from Human Feedback (RLHF). Figure from [19]. . . . .	18
1.6	Tool categorization. Figure from [20]. . . . .	20
1.7	Tool learning framework. Figure from [20] . . . . .	22
1.8	MLLM-Tool architecture. Figure from [21]. . . . .	22
2.1	Dataset creation pipeline. . . . .	28
2.2	Qualitative dataset examples. . . . .	29
2.3	RaTA-Tool pipeline. . . . .	30
3.1	Qualitative results of RaTA-Tool under different input modalities. . . . .	44
A.1	Qualitative examples of RaTA-Tool for text-only queries. . . . .	48
A.2	Qualitative examples of RaTA-Tool for multimodal queries with image inputs. . . . .	48
A.3	Qualitative comparison of zero-shot task description and RaTA-Tool one, along with the ground-truth description. . . . .	49
B.1	Example of prompt for tool description generation in JSON. . . . .	52

B.2	Example of prompt for tool description generation in natural language. . . . .	53
B.3	Example of prompt for task-description generation at inference time in JSON. . . . .	54
B.4	Example of prompt for task-description generation at inference time in natural language. . . . .	55

# List of Tables

1	Distribution of training and testing sets. Table from [21]. . . . .	27
2	Dataset statistics. . . . .	35
3	Decoding strategies used to sample responses for the DPO preference dataset. . . . .	37
4	Accuracy results on aggregate metrics: average per query ( $Avg_q$ ) and average per modality ( $Avg_m$ ). Comparative results between Qwen2.5-Omni-3B [31] and Qwen2.5-Omni-7B [31] as the MLLM backbone. Comparative results between Qwen3-Embedding [32] and Contriever [34] as the embedding backbone. . . . .	40
5	Experiments on the effect of query and tool descriptions format and prompting or fine-tuning strategy on tool-selection performance, reported for both embedding backbones (Contriever and Qwen3-Embedding). Results are shown for text, image, and audio queries, together with per query ( $Avg_q$ ) and per modality ( $Avg_m$ ) aggregate metrics. . . . .	42



# Introduction

The influential paper “Attention Is All You Need” by Vaswani *et al.* [1] forever changed the deep learning landscape. The presented transformer architecture quickly became the de facto standard for a wide range of models, spanning from machine translation to image understanding and generation. Subsequently, the research community focused heavily on this new architecture, resulting in a vast number of derivative models and solutions.

Transformers rely the attention mechanism [1], which allows the model to process all words in a sequence simultaneously, deriving meaning not solely from their independent significance, but by establishing meaningful connections with the rest of the sentence. Consequently, the model understands the broader context and meaning of the sequence, rather than analyzing each word in isolation. Furthermore, the architecture employs a multi-head attention mechanism, enabling the model to jointly attend to different representation subspaces simultaneously, yielding a more comprehensive understanding of the text.

Shortly after, GPT-1 [2] introduced the concept of Generative Pre-trained Transformer. This decoder-only transformer was trained in a self-supervised fashion on a massive corpus of data, yielding a model capable of generating coherent, extended text. More precisely, it was trained to predict the next word (or token) in a sequence. This type of model is known as autoregressive, meaning

it uses its previous outputs as inputs for subsequent generation steps. GPT-1 demonstrated that decoder-only transformers could achieve state-of-the-art results in natural language generation.

Contemporary LLMs, descended from GPT-1 and similar architectures, master not only text generation, but also excel in coding, reasoning, and information extraction. Modern LLMs are widely deployed as virtual assistants across a broad spectrum of tasks. Moreover, the Vision Transformer [3] enabled these architectures to jointly process text and image sequences, giving rise to Multimodal Large Language Models (MLLMs). Subsequent works expanded this capability to other input modalities, such as audio and video [4]–[6].

Given the success of these models, the research community sought to equip them with new capabilities to solve increasingly difficult problems and achieve more human-like autonomy. This concept evolved into what is known today as agentic AI.

The goal of agentic AI is to make models capable of pursuing complex, multi-task, and multifaceted objectives autonomously. While traditional LLMs are prompt-oriented, requiring detailed user instructions for each sequential step, agentic systems are goal-oriented, capable of autonomously determining the necessary steps to accomplish a user-defined goal.

The literature identifies four foundational patterns for building agents: reflection, external tool use, planning and multi-agent collaboration [7].

The goal of this thesis is to explore the current state-of-the-art results achieved in external tool use and to develop a new framework to allow further advancements in this field.

# Chapter 1

# Background

## 1.1 Neural networks

Neural networks, also known as Artificial Neural Networks (ANNs), are computational models that resemble the structure and functions of biological neural networks. ANNs emerged as an attempt to leverage the architecture of the human brain to perform tasks that conventional algorithms struggled to solve. The foundational element of such a network is the neuron, a component that, in the biological brain, transmits information using electrical and chemical signals. [8]

An artificial neural network consists of interconnected simulated neurons. Each neuron is connected to other nodes via links analogous to biological axon-synapse-dendrite connections. [8] These nodes receive input data and process it to perform specific operations. Crucially, each link has an associated weight that determines the strength of one node's influence on another, thereby modulating the signal transmitted between neurons.

These neurons are grouped into layers: an input layer that receives the initial features, an output layer that computes the final results, and one or more hidden layers between them. The more hidden layers a network has, the deeper it is

considered to be. Historically, ANNs with at least two hidden layers were classified as deep neural networks.

Each artificial neuron aggregates signals from the units in the preceding layer to generate a single scalar output. This output is then propagated to multiple neurons in the subsequent layer of the network.

Mathematically, the output of each artificial neuron is expressed as follows:

$$z = \sigma\left(\sum_{i=1}^n w_i x_i + b\right) \quad (1.1)$$

where  $\sigma$  indicates a non-linear activation function,  $w_i$  denotes the weight associated with the  $i$ -th input  $x_i$ , and  $b$  represents the scalar bias.

## 1.2 Multilayer perceptron

The multilayer perceptron (MLP) is a foundational class of feedforward artificial neural network, wherein information propagates in a single direction: from the input nodes, through the hidden nodes, to the output nodes. An MLP comprises at least three layers, and each node connects to every node in the subsequent layer, forming a fully connected network.

MLPs are distinguished by their capacity to learn complex, non-linear mappings within data, thanks to non-linear activation functions. They are typically trained in a supervised learning paradigm using the backpropagation algorithm.

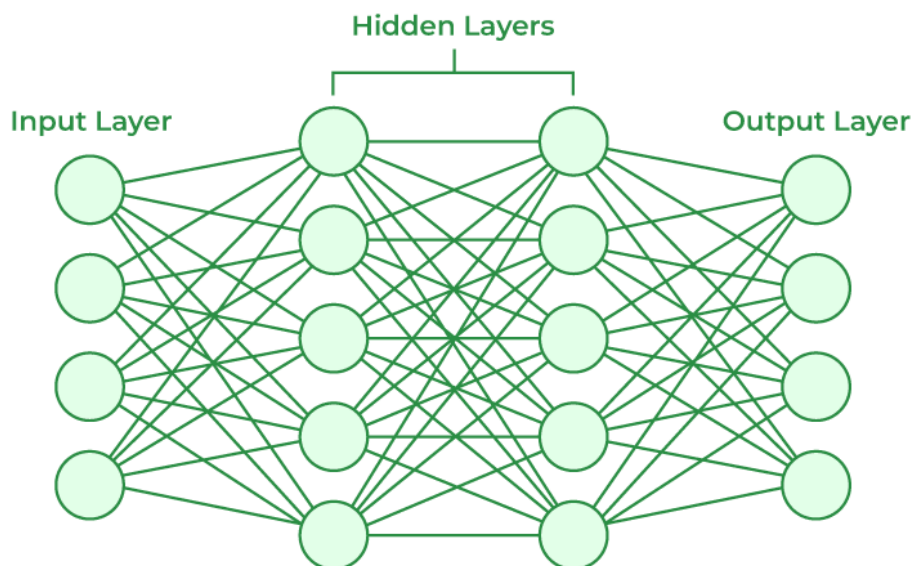


Figure 1.1: Architecture of an MLP

## 1.3 Transformers

Contemporary LLMs predominantly leverage the Transformer architecture. [1] While the original paradigm used an encoder-decoder structure, autoregressive generative models often use decoder-only variants. For this reason, the following description will focus on the decoder-only Transformer and its operational principles.

The decoder, given a sequence of continuous representations  $\mathbf{Z} = (z_1, \dots, z_n)$  with  $z_i \in \mathbb{R}^{d_{model}}$ , generates an output sequence  $\mathbf{y} = (y_1, \dots, y_m)$ . The decoder is composed of a stack of  $N$  identical layers, each consisting of a masked multi-head attention mechanism and a position-wise fully connected feed-forward network. After the last layer, a linear projection and a softmax function produce the final output probabilities.

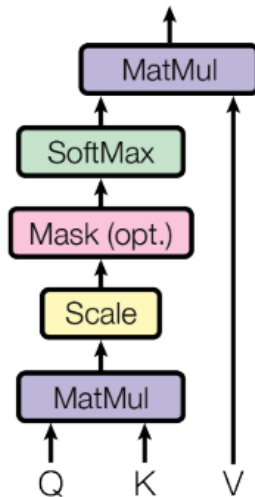


Figure 1.2: Scaled masked attention. Figure from [1].

### 1.3.1 Attention

The attention mechanism is central to the Transformer and enables the modeling of the representation of a sequence by relating elements at different positions.

The input sequence  $\mathbf{Z}$  is first projected into three distinct latent spaces: queries ( $Q$ ), keys ( $K$ ) and values ( $V$ ). Each projection is the result of a learned linear transformation, parametrized by the weight matrices  $W_q, W_k \in \mathbb{R}^{d_{model} \times d_k}$  and  $W_v \in \mathbb{R}^{d_{model} \times d_v}$  as follows:

$$Q = \mathbf{Z}W_q \quad (1.2)$$

$$K = \mathbf{Z}W_k \quad (1.3)$$

$$V = \mathbf{Z}W_v \quad (1.4)$$

As the projections are all linear transformations of the same input  $\mathbf{Z}$ , this mechanism is formally referred to as self-attention. This differs from cross-attention, where queries are generated from one sequence, whereas keys and values derive from another input stream.

Given  $Q$ ,  $K$  and  $V$ , the particular type of attention applied is called scaled dot-product attention (Figure 1.2) and is formally defined as:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (1.5)$$

**Masked attention** The result of  $QK^T$  relates each input token with itself and all others in the sequence. In generative tasks, however, this representation must be causally masked to prevent a token from attending to subsequent tokens. This is achieved by applying a look-ahead mask before the softmax operation. Specifically, the attention energy  $e_{ij}$  is set to  $-\infty$  for all  $j > i$ , nullifying the influence of subsequent tokens on the current hidden state. Practically, this is obtained via an element-wise addition of a mask matrix  $\mathbf{M}$  to the attention score matrix  $QK^T$ . The mask matrix is obtained from a binary lower triangular matrix  $L$ :

$$L = \begin{bmatrix} 1 & 0 & 0 & \cdots & 0 & 0 \\ 1 & 1 & 0 & \cdots & 0 & 0 \\ 1 & 1 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & 1 & 1 & \cdots & 1 & 0 \\ 1 & 1 & 1 & \cdots & 1 & 1 \end{bmatrix} \quad (1.6)$$

Subsequently, an element-wise mapping  $f(L_{ij})$  is applied to this binary matrix to construct the functional mask:

$$\mathbf{M} = f(L_{ij}) = \begin{cases} 0, & \text{if } L_{ij} = 1 \\ -\infty, & \text{if } L_{ij} = 0 \end{cases} \quad (1.7)$$

Upon adding this mask, the subsequent softmax operation suppresses the  $-\infty$  values to zero ( $e^{-\infty} \rightarrow 0$ ). This effectively eliminates the contribution of future tokens.

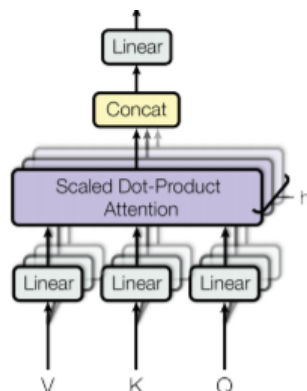
The result of  $QK^T$  is also scaled by  $\frac{1}{\sqrt{d_k}}$ . The magnitude of the dot-product attention can increase significantly for large values of  $d_k$ , pushing the softmax function into regions with vanishing gradients. This scaling prevents this effect.

The final stage involves the multiplication of the attention scores by the value matrix  $\mathbf{V}$ . This operation computes a weighted sum of the values, where the weights represent the relative importance of each token within the masked context.

**Multi-head attention** Rather than computing a single attention function, the architecture employs Multi-Head Attention (MHA), where multiple attention heads operate in parallel. As shown in Figure 1.3, each head has distinct learned linear projections, enabling the extraction of features at different levels of abstraction. The outputs from these individual heads are then concatenated and linearly projected to form the final aggregated representation for each token.

### 1.3.2 Fully connected feed-forward network

Following the multi-head attention sub-layer, the representations are passed through a position-wise feed-forward network. This component introduces non-linearities into the model, expanding its representational capacity beyond simple linear aggregations. The function is defined as follows:



**Figure 1.3:** Multi-head attention. Figure from [1].

$$FFN(x) = \sigma(xW_1 + b_1)W_2 + b_2 \quad (1.8)$$

It consists of two linear transformations with a non-linear activation function applied in between. The inner dimension is typically significantly larger than  $d_{model}$ , enabling the model to learn much more complex relationships than it could within the original, narrower dimension  $d_{model}$ .

Unlike the attention mechanism, the feed-forward layers process token's information individually. There is no communication between tokens; each token's features are processed in isolation to yield a deeper representation of the token itself.

### 1.3.3 Transformer architecture

Following the multi-head attention and the feed-forward network of each Transformer layer, residual connections and layer normalization are employed. The output of each sub-layer is then:

$$LayerNorm(x + Sublayer(x)) \quad (1.9)$$

As stated in the architectural overview, the Transformer decoder is made up of a stack of  $N$  identical layers. Following the final layer, a linear layer  $\mathbf{W}_{out} \in \mathbb{R}^{d_{model} \times |\mathcal{V}|}$  projects the dense vectors into the vocabulary space. The resulting *logits*, which are raw unnormalized scores, undergo a softmax function, yielding a probability distribution over the entire lexicon. Each element of this output vector represents the predicted probability that the  $i$ -th token in the vocabulary is the next token in the sequence.

## 1.4 Large Language Models (LLMs)

Large Language Models (LLMs) are a class of models designed to process and generate human-like text. They are typically categorized into three main architectural branches:

**Decoder-only (Generative):** Examples include GPT [9], Llama [10] and Qwen[11]; they predict the next token based on the input context and previously generated tokens.

**Encoder-only (Discriminative):** Examples include BERT [12]; they build contextual representations by processing the entire input sequence bidirectionally and can predict a specific label. This category of models is often employed for sentiment analysis or general classification tasks.

**Encoder-Decoder (Hybrid):** Examples include T5 [13]; they translate or transform the input sequence. These architectures are typically employed for sequence-to-sequence tasks such as machine translation or summarization.

Generative LLMs have demonstrated extraordinary capabilities in human-like text generation and are frequently leveraged as few-shot learners across a diverse set of downstream tasks. We will focus on this class of language models.

### 1.4.1 Tokenization

The vocabulary, denoted as  $\mathcal{V}$ , is the finite set of all tokens that a specific model is capable of recognizing and processing. During the tokenization process, the input text is decomposed into these units, each of which is represented by a unique index. Consequently, the raw textual input is transformed into a vector of discrete indices  $\mathbf{x} = (x_1, x_2, \dots, x_n)$  where  $x_i$  is the  $i$ -th token index.

While these tokens may represent characters, subwords, or full words depending on the specific tokenizer used, the terms “token” and “word” are used interchangeably throughout this work for convenience.

### 1.4.2 Embedding layer

Splitting input text into a sequence of token indices is the essential first step to move from human language to machine-understandable content. However, token indices fail to encode semantic relationships between tokens; two similar words may be assigned widely distant index values. For this reason, they are unsuitable for direct processing by the neural network and must be projected into a continuous vector space before being used as input.

This is achieved through an embedding layer, where each input token index is mapped to a  $d_{model}$ -dimensional vector. This mapping is performed via an embedding matrix  $W_e \in \mathbb{R}^{|\mathcal{V}| \times d_{model}}$ . Here,  $|\mathcal{V}|$  denotes the vocabulary cardinality, representing the total number of unique tokens in the model’s lexicon. Conceptually, the embedding matrix functions as a high-dimensional lookup table that maps discrete token indices to their corresponding dense vector representations. The vector representations within the embedding matrix are learnable parameters.

The resulting sequence  $\mathbf{X} \in \mathbb{R}^{n \times d_{model}}$  captures the static semantic properties of each individual token. Given that the Transformer architecture is inherently

permutation-invariant, it treats the input as a “bag of words” and does not account for the sequential order that characterizes natural language unless explicitly provided with positional information. Consequently, positional embeddings are summed element-wise with  $\mathbf{X}$  to integrate sequential information into the input before it enters the attention layers.

### 1.4.3 Training

Modern LLMs undergo a series of training stages.

**Pre-training** The first step is known as “pre-training” and adopts a self-supervised approach, using unlabeled data and a next token prediction objective. During this phase, the model is exposed to a vast dataset scraped from the internet and learns the underlying statistical structures and semantic nuances of language. This phase serves as the basis for all subsequent training, providing the model with the fundamental ability to generate syntactically and semantically correct text.

**Supervised Fine-Tuning (SFT)** The second step is Supervised Fine-Tuning (SFT), where the model is further conditioned on datasets of curated, high-quality prompt-response pairs. At this stage, the model can be fine-tuned for diverse downstream tasks, such as summarization and question answering. Rather than optimizing for a single task, modern generative models are fine-tuned on a multitude of instruction sets concurrently. This results in a highly flexible system capable of performing diverse tasks, as often seen in modern virtual assistant models.

**Reinforcement Learning from Human Feedback (RLHF)** The final stage is Reinforcement Learning from Human Feedback (RLHF) [14]. Although

SFT provides a structural foundation, models frequently retain residual behaviors that diverge from established safety and utility constraints. These issues can manifest as the generation of toxic output, harmful instructions, or hallucinated content. To mitigate these risks, RLHF is employed to more closely align the model's responses with human values and safety guidelines.

## 1.5 Multimodal Large Language Models (MLLMs)

While traditional LLMs were restricted to textual processing, recent advancements have enabled multimodal dialogue within these models. Multimodal Large Language Models (MLLMs) are capable of processing heterogeneous data streams, including images, video, and audio, and generating text that aligns with both text instructions and multimodal information.

In this theoretical analysis, we will focus on MLLMs supporting image inputs. The underlying architecture and principles are applicable, with some adjustments, to other modalities. The essential components of a Multimodal Large Language Model are:

**Visual Encoder** : This component extracts feature vectors from images. A standard choice is CLIP [15] with a Vision Transformer (ViT) [3] image encoder.

**Adapter Module** : This module performs dimensional alignment between the visual features extracted by the visual encoder and the LLM input space. Usually, a 2-layer MLP or a linear projector is employed.

**Language Model Backbone** : A standard Transformer-based LLM responsible for language processing and generation.

### 1.5.1 Vision Transformer (ViT)

The Vision Transformer (ViT) [3] first pioneered the application of the traditional Transformer architecture to image processing. This approach moved beyond the traditional convolutional methods that were dominant in this field, such as the ResNet-like architectures [16] that previously defined the state of the art in large-scale image recognition.

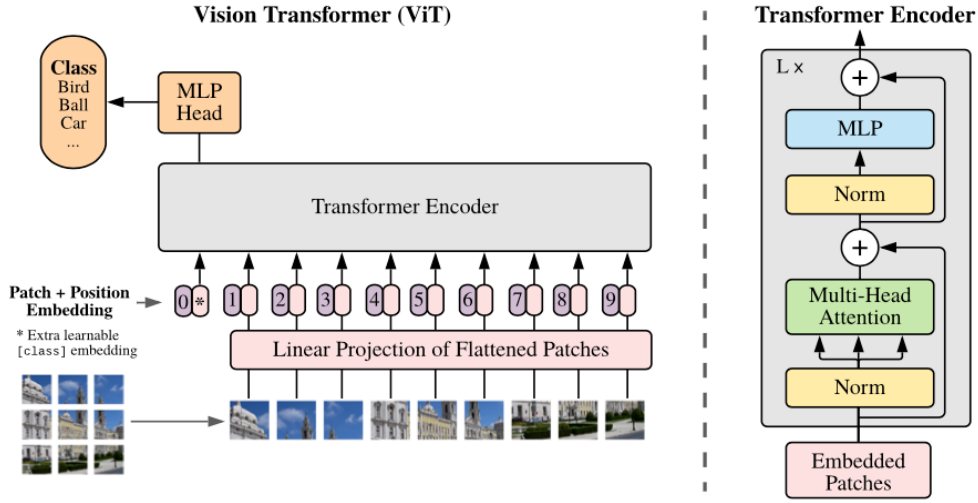
The main idea behind ViT is to split an image into fixed-size, non-overlapping patches and feed the sequence of these patches as input to a Transformer, analogous to tokens in a sentence. The architecture is illustrated in Figure 1.4.

To accomplish this, the original image  $\mathbf{x} \in \mathbb{R}^{H \times W \times C}$ , where  $H$  and  $W$  indicate the original image resolution (height and width) and  $C$  the number of channels, is reshaped into a sequence of flattened 2D patches  $\mathbf{x}_p \in \mathbb{R}^{N \times (P^2 \cdot C)}$ .  $P \times P$  indicates the resolution of each image patch, and  $N = \frac{HW}{P^2}$  is the total number of resulting patches. The 2D flattened patches are treated as the token embeddings, where  $P^2 \cdot C$  is the embedding vector length and  $N$  the sequence length.

Since Transformers do not inherently account for spatial information, positional embeddings are added to the patches, similar to the approach used in traditional LLMs.

**Classification** To perform classification tasks, a specialized learnable embedding vector, denoted as the class token ( $\mathbf{x}_{class}$ ), is prepended to the sequence of patches, similar to the  $[CLS]$  token in BERT [12]. This vector serves as a global information aggregator: due to the attention mechanism, it attends to every other patch in the image and integrates information from all of them. At the final layer, it has encoded sufficient relevant features to be used as a global image representation and serves as the input for the classification head.

Formally, the input to the Vision Transformer is represented as:



**Figure 1.4:** Vision Transformer (ViT) architecture. Figure from [3].

$$\mathbf{z}_0 = [\mathbf{x}_{class}; \mathbf{x}_p^1 \mathbf{E}; \mathbf{x}_p^2 \mathbf{E}; \dots; \mathbf{x}_p^N \mathbf{E}] + \mathbf{E}_{pos} \quad (1.10)$$

where  $\mathbf{x}_{class}$  is the special learnable embedding vector,  $\mathbf{x}_p^i$  denotes the  $i$ -th image patch and  $\mathbf{E} \in \mathbb{R}^{(P^2 \cdot C) \times d_{model}}$  is a linear projection mapping the patches to the Transformer hidden dimension.

Subsequently, a series of traditional Transformer encoder layers is stacked, and after the final layer, the class token  $\mathbf{x}_{class}$  can be used for classification.

## 1.5.2 Contrastive Language-Image Pre-training (CLIP)

Contrastive Language-Image Pre-training (CLIP) [15] bridges the semantic gap between vision and language by mapping both modalities into a unified embedding space, enabling the model to interpret images through the lens of natural language.

CLIP uses a contrastive learning objective, where the model is tasked with predicting the correct associations between images and captions.

The architecture is composed of an image encoder, usually a Vision Transformer, and a text encoder, a standard Transformer. The Vision Transformer here is not used as a classifier; instead, it is used to extract high-dimensional latent representations from raw image data. Instead of prepending a special class token to the sequence of patches and extracting it after the final Transformer layer, the entire matrix  $\mathbf{z}_L \in \mathbb{R}^{N \times d_{model}}$  is used as a dense representation of the image. The text encoder brackets the sequence with  $[SOS]$  and  $[EOS]$  tokens; the activations at the  $[EOS]$  token index from the final Transformer layer serve as the global feature representation of the text input.

Image and text latent representations are linearly projected into the same multimodal embedding space. The model takes a batch of  $N$  image embeddings and  $N$  text embeddings and computes the cosine similarity between each possible pair. The model’s goal is to maximize the similarity values for correct pairs while minimizing those of wrong pairs.

After training, CLIP is able to encode visual inputs into a latent space aligned with linguistic semantics. By anchoring images to linguistic concepts, the framework provides the required continuous vector representations to feed visual information into the Transformer backbone of an MLLM.

## 1.6 Parameter Efficient Fine-Tuning (PEFT)

While Supervised Fine-Tuning (SFT) is an essential methodology for steering foundation models towards specialized downstream tasks without the need to pre-train from scratch, it is not without drawbacks. Full-parameter fine-tuning is often prohibitively expensive and may lead to catastrophic forgetting [17], a phenomenon where the model’s pre-existing knowledge is degraded during adaptation.

Parameter-Efficient Fine-Tuning (PEFT) aims to mitigate these issues by fine-

tuning only a tiny subset of the model’s parameters while achieving performance comparable to that of full-parameter fine-tuning.

### 1.6.1 LoRA

Low-Rank Adaptation (LoRA) [18] proposes a framework for Parameter-Efficient Fine-Tuning (PEFT). Specifically, instead of fine-tuning all the parameters of a pre-trained weight matrix  $\mathbf{W}_0 \in \mathbb{R}^{d \times k}$ , LoRA approximates its weight update  $\delta\mathbf{W}$  using two low-rank matrices  $\mathbf{A} \in \mathbb{R}^{r \times k}$  and  $\mathbf{B} \in \mathbb{R}^{d \times r}$ , with  $r \ll \min(d, k)$ .  $\mathbf{W}_0$  is frozen during training, whereas the A and B parameters are updated. The fine-tuned weight matrix  $\mathbf{W} \in \mathbb{R}^{d \times k}$  is obtained by summing the original pre-trained weight matrix  $\mathbf{W}_0$  and the product of the low-rank matrices  $\mathbf{A}$  and  $\mathbf{B}$ :

$$\mathbf{W} = \mathbf{W}_0 + \mathbf{B}\mathbf{A} \tag{1.11}$$

The parameter  $\alpha$  controls the magnitude of the LoRA update. The adapter values are scaled by  $\frac{\alpha}{r}$  before being added to the original weights. The update formula with the scaling factor becomes:

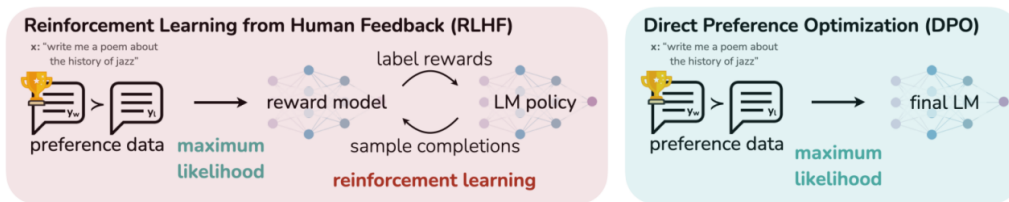
$$\mathbf{W} = \mathbf{W}_0 + \frac{\alpha}{r}\mathbf{B}\mathbf{A} \tag{1.12}$$

After training, the LoRA weights can be merged and stored with the original matrices, thus eliminating any additional inference latency.

Such an approach significantly reduce the memory and storage used for training, allowing for fine-tuning with fewer resources. Another advantage is that a fine-tuned model can be easily adapted to a different downstream task

by simply swapping the LoRA weights.

## 1.7 Direct Preference Optimization (DPO)



**Figure 1.5:** Direct Preference Optimization (DPO) vs. Reinforcement Learning from Human Feedback (RLHF). Figure from [19].

While teaching LLMs broad language knowledge is a conceptually straightforward task driven by self-supervised learning, aligning a model to specific behavioral patterns is a non-trivial challenge.

Reinforcement Learning from Human Feedback (RLHF) [14] remains the de facto standard for this objective: a separate reward model is trained to predict human preferences based on collected comparative data. The policy model is then optimized to maximize the reward model’s score. However, this procedure is notoriously complex and unstable, driving the development of more streamlined frameworks.

Direct Preference Optimization (DPO) [19] addresses these challenges by reformulating the alignment objective. By deriving a closed-form solution for the optimal policy, DPO avoids the instabilities of reinforcement learning and instead optimizes the model using a straightforward classification loss on preference data.

The DPO objective is defined as:

$$\mathcal{L}_{\text{DPO}}(\pi_{\theta}; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \log \frac{\pi_{\theta}(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \beta \log \frac{\pi_{\theta}(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right) \right] \quad (1.13)$$

## 1.8 Agentic AI

Agentic AI refers to autonomous systems capable of reasoning, planning, and using external tools to achieve complex, multi-step goals with minimal human intervention. According to research, four design patterns must be followed to develop Agentic AI systems [7]:

**Reflection** : The model is able to iteratively review its own work, adjusting and improving it without human intervention.

**Tool use** : The model can effectively decide when and which external tools use depending on the task. Tools can include web search, code execution, or other specialized models.

**Planning** : Given a multi-step task, the model is able to plan the steps needed to solve it and carry out the action plan autonomously.

**Multi-agent collaboration** : Rather than a single AI agent solving all tasks independently, a team of agents collaborates to refine ideas and plans, divide the work, and review final results.

### 1.8.1 Tool learning

In human evolutionary history, the development of specialized tools has served as a primary mechanism for augmenting innate physiological and cognitive capacities, allowing the species to transcend its biological constraints. [20]

A similar trajectory is evident in the evolution of AI. Original foundation

models were only able to perform next-token prediction, operating exclusively within the limits of their internal parameters.

Subsequently, the paradigm shifted as LLMs were augmented with the capability to interface with external resources, such as web search engines. This integration of external tools transformed LLMs from isolated predictors into interactive assistants capable of real-world interactions.

The capability to use external tools is crucial when creating agentic systems. LLMs are indeed extremely good at next-token prediction but lack human-level ability in many areas, and do not have access to real-time information that is not encoded in their weights. Moreover, many specialized models are outperforming generalist LLMs in many domains, making the use of the latter less effective in those contexts.

Tool learning with foundation models refers to this phenomenon of teaching foundation models to use external tools, mirroring the trajectory of human tool use. Despite the rapid advancements in this field, it remains characterized by a lack of unified systematic frameworks and significant unresolved research challenges.

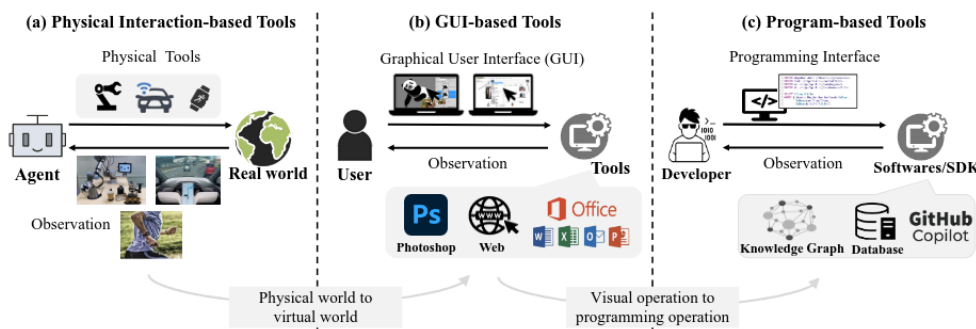


Figure 1.6: Tool categorization. Figure from [20].

**Tool categorization** Tools can be categorized into a three tier taxonomy (Figure 1.6) based on the nature of their interaction. [20]. The three categories are:

- **Physical Interaction-based Tools:** This includes all tools that require physical interaction, such as robots and wearables.
- **GUI-based Tools:** This type of tool involves a virtual interaction via a programmatic interface. The user usually interacts with the system through digital components like buttons, text fields, and so on.
- **Program-based Tools:** Tools at this level can be used via a programming interface, offering a high degree of flexibility. For instance, programming libraries, SDKs, or neural network based tools.

The goal of tool learning is to bridge the architectural divide between disparate interaction modalities, empowering agents to navigate a continuum of tool environments, from programmatic to physical, within a single framework.

**Tool learning framework** A comprehensive tool learning framework is essential to establish a standardized conceptual baseline, defining both the main components and the interdependencies between them. The framework illustrated in Figure 1.7 comprises four main components:

- **Tool set:** The tool set  $\mathcal{T} = \{\mathcal{T}_1, \mathcal{T}_2, \dots\}$  is a heterogeneous set of tools with diverse functional capabilities. The foundation model has the capability, when answering a particular query, to exploit some of these to better frame its reasoning and provide a more accurate response.
- **Environment:** The environment serves as the operational substrate, whether virtual or physical, that provides the necessary infrastructure for tool execution and returns feedback to the model.
- **Perceiver:** The perceiver sends an elaborated summary to the controller based on the environmental and human feedback.
- **Controller:** The controller, usually a foundation model, learns to map user requests to specific tools when interaction with them is considered helpful.

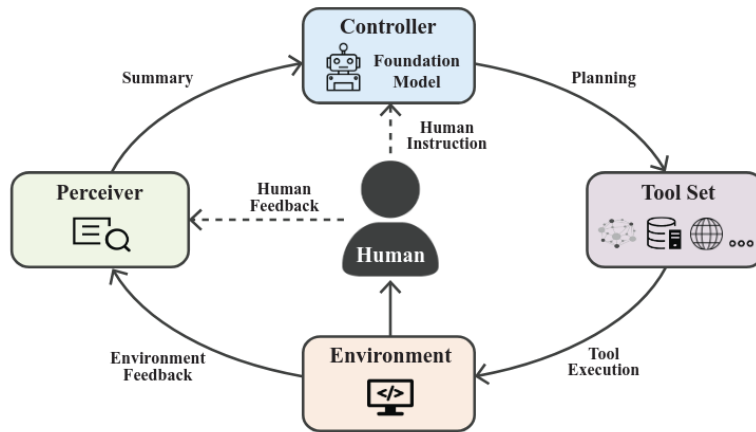


Figure 1.7: Tool learning framework. Figure from [20]

## 1.8.2 Multimodal tool learning

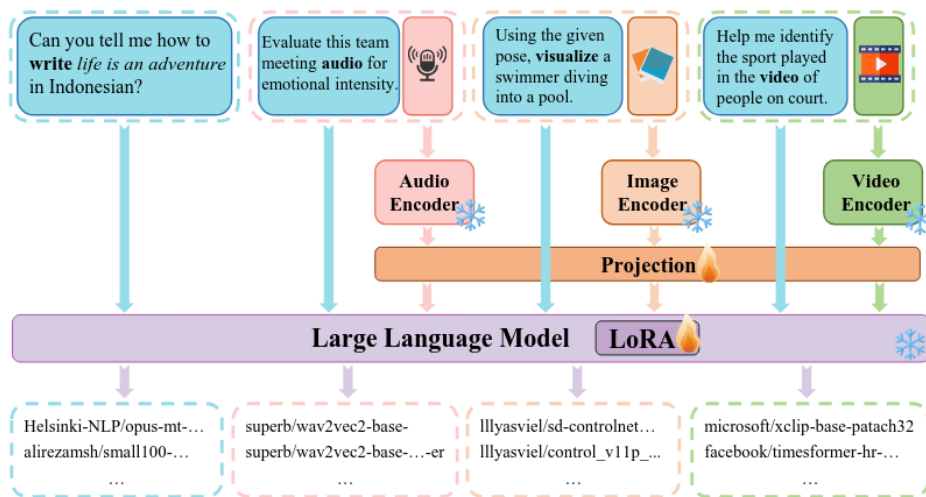


Figure 1.8: MLLM-Tool architecture. Figure from [21].

Recent studies have demonstrated that text-only queries can frequently lead to ambiguity when used as the exclusive information source for the controller’s tool-selection logic.

Consequently, solutions that extend tool-use capabilities to Multimodal Large Language Models (MLLMs) are emerging. These solutions enable the controller to process cross-modal sensory information, thereby providing the necessary contextual grounding to accurately interpret user requests.

MLLM-Tool [21] represent an initial step towards this direction. Its architectural approach, depicted in Figure 1.8, combines a multimodal encoder with a traditional LLM backbone to jointly map visual/auditory and textual features into the same space.

The study uses ToolMMBench, a novel dataset designed to associate multimodal queries with their corresponding set of relevant tools. The tool set was sourced from the Hugging Face repository.

ImageBind [22] was used as the multimodal encoder, with its weights frozen. It extracts features from the multimodal input (image, audio, or video) and linearly projects these features into the LLM’s feature space.

The LLM component was trained for the tool selection task using Low-Rank Adaption (LoRA) [18]. By jointly processing the prompt and the aligned multimodal embeddings, the model is capable of accurately predicting the tool that best fits the task. This methodology was applied to a series of open-source foundation models: Vicuna [23], Llama [24], Llama2 [24], and Llama2-Chat [24], with experiments conducted on both 7B and 13B variants.

The LLM fine-tuning leveraged ToolMMBench, which contains instruction-answer pairs. In the training set, one-to-many pairs were split into multiple one-to-one pairs to allow for the proper calculation of the objective function.

The fine-tuning process enables the LLM to act as the controller within the architecture illustrated in Figure 1.7. However, as the training objective was specifically designed for tool prediction, the controller is limited to operating within a closed set of tools encountered during the fine-tuning process.



## Chapter 2

# Proposed method

This thesis aims to analyze Multimodal Large Language Models (MLLMs) serving as central controllers within an agentic system framework.

Although pioneering works like Toolformer [25], Gorilla [26] and ToolLLM [27] have successfully demonstrated tool-use capabilities, these models remain functionally constrained to text-based inputs, often leading to interpretative ambiguity and a lack of multimodal perception.

HuggingGPT [28], Visual ChatGPT [29] and GPT-4Tools [30] prioritize task orchestration by decomposing multifaceted tasks into subtasks. These subproblems are subsequently addressed by a team of specialized models, allowing the central LLM to serve as a high-level manager.

While these solutions lay the foundation for LLM-driven agentic systems, they are functionally constrained by a text-only input modality. This imposed limitation frequently hinders the model’s ability to discern subtle nuances within contextually similar requests, thereby undermining its capacity to accurately identify specialized tools.

To mitigate the interpretative constraints of unimodal inputs, recent research has leveraged MLLMs to achieve more robust contextual grounding. MLLM-Tool [21] is an initial step in this direction, enabling the retrieval of external

APIs based on multimodal queries.

However, it operates under a closed-world assumption; its predictive accuracy is bound to the specific API corpus used during the training phase. In dynamic environments where tools are frequently updated, such an approach severely limits its practical applicability, as the need for constant re-training can hardly be cost-effective.

This work introduces a novel methodology for open-world multimodal tool selection. Rather than training the controller to predict discrete tool labels, our approach uses an embedding-based retrieval system, which project both multimodal queries and tool specifications into a high-dimensional semantic space. Tool matching is performed using vector similarity measures, bypassing the need for continual fine-tuning and enabling the model to naturally generalize to new tools.

## 2.1 Dataset creation

### 2.1.1 ToolMMBench

We started from the publicly available dataset from MLLM-Tool [21], ToolMM-Bench.

The dataset comprises four different input modalities: *text-only*, *text-image*, *text-audio* and *text-video*. It is composed of one-to-many instruction-answer pairs, where the instruction is the user query along with possible multimodal input, whereas the answer is a collection of the ground truth external tools for that task. This one-to-many mapping emulate real tasks, where many tools can potentially be employed for the same task.

External tools are represented by Hugging Face scraped open-source models. Overall, the dataset contains 932 different models: 654 with text inputs, 191 with mixed text and image inputs, 79 with text and audio, and 9 with text and

video.

Following a standard 80/20 split, the dataset was partitioned into a training set comprising 44,533 instructions and a test set totaling 2,397 instructions. The distribution of each set across the input modalities is represented in Table 1.

Each one-to-many query in the training set is split into multiple one-to-one instances: this is done to ensure that each tool is individually represented during the training phase, allowing the model to learn the different possible answers for each given query.

**Table 1:** Distribution of training and testing sets. Table from [21].

<b>Settings</b>	<b>Text</b>	<b>Audio</b>	<b>Image</b>	<b>Video</b>
Training Set (w/ split)	31,280	2,944	10,085	224
Training Set (w/o split)	5,563	562	3,938	82
Testing Set	1,452	164	760	21

Note: The first line represents the number of instruction-answer pairs after splitting one-to-many queries into multiple one-to-one instructions. The second line represents distinct query numbers for the four modalities.

### 2.1.2 Our dataset

While ToolMMBench served as a solid foundation for our work, its supervised formulation limited its applicability. In fact, it only allows to train a model to predict tools that are seen during the training phase, limiting its ability to handle user queries requiring unseen tools. For this reason, instead of using the raw instruction-tool pairs to teach the model some statistical distribution, we employed structured tool descriptions to represent each tool. Moreover, further pre-processing operations allowed to adjust the original dataset structure to our training objective.

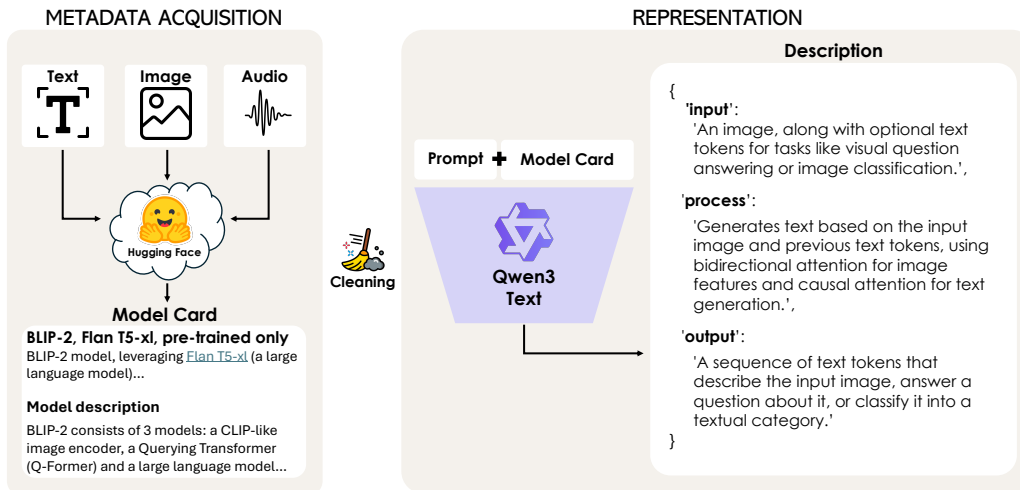


Figure 2.1: Dataset creation pipeline.

**Tools characterization** Each tool in the dataset is a Hugging Face model. To build high-dimensional tool representations, we programmatically scraped the Hugging Face model card for each tool, containing a natural language description of functionalities, inputs and outputs, fragments of code, and training details.

The model card contains heterogeneous content: not all the information is useful at identifying the problem it addresses, and description length may largely vary between models. Thereby, the scraped model cards for each tool were used as input to generate a structured JSON description. For this phase, an LLM (*i.e.*, Qwen3 [11]) was involved: a customized prompt along with the tool’s model card was used to prompt the LLM to generate the description. An explanatory scheme is represented in Figure 2.1.

The JSON description is composed of three fields:

- **Input:** The expected input format and modalities expected by the model.
- **Process:** A functional summary of the task performed.
- **Output:** The expected output format.

Some qualitative examples of the dataset are shown in Figure 2.2.

**Splitting, merging and cleaning** We first split all one-to-many mappings into one-to-one instances and merged the original training and test splits to

KHA-WHITE/MANGA-OCR-BASE	
Model Card	Description
<p><b>Manga OCR</b> Optical character recognition for Japanese text, with the main focus being Japanese manga.</p> <p>It uses <a href="#">Vision Encoder Decoder</a> framework.</p> <p>Manga OCR can be used as a general purpose printed Japanese OCR, but its main goal was to provide a high quality text recognition, robust against various scenarios specific to manga:</p> <ul style="list-style-type: none"> <li>• both vertical and horizontal text</li> <li>• text with furigana</li> <li>• text overlaid on images</li> <li>• wide variety of fonts and font styles</li> <li>• low quality images</li> </ul> <p>Code is available <a href="#">here</a>.</p>	<pre>{   'input':     'An image, along with optional text tokens for tasks like visual question answering or image classification.',   'process':     'Generates text based on the input image and previous text tokens, using bidirectional attention for image features and causal attention for text generation.',   'output':     'A sequence of text tokens that describe the input image, answer a question about it, or classify it into a textual category.' }</pre>
VALHALLA/T5-BASE-QG-HL	
Model Card	Description
<p><b>T5 for question-generation</b> This is <a href="#">t5-base</a> model trained for answer aware question generation task. The answer spans are highlighted within the text with special highlight tokens.</p> <p>You can play with the model using the inference API, just highlight the answer spans with <code>&lt;h1&gt;</code> tokens and end the text with <code>&lt;/s&gt;</code>.</p> <p>For example <code>&lt;h1&gt; 42 &lt;/h1&gt; is the answer to life, the universe and everything. &lt;/s&gt;</code></p> <p>For more details see <a href="#">this</a> repo. Model in action 🚀 You'll need to clone the <a href="#">repo</a>.</p>	<pre>{   'input':     'Text containing answer spans highlighted with &lt;h1&gt; tokens, ending with &lt;/s&gt; token..',   'process':     'Generates a question based on the provided text and the highlighted answer span.',   'output':     'A dictionary containing the generated question and the corresponding answer.' }</pre>

Figure 2.2: Qualitative dataset examples.

form a unified dataset of one-to-one instruction-tool pairs. During this step, we performed dataset-wide cleaning to remove duplicated pairs<sup>1</sup> and inconsistent samples<sup>2</sup>.

## 2.2 Pipeline

This section outlines the RaTA-Tool pipeline. Given a user query  $q \in \mathcal{Q}$ , the pipeline’s objective is to retrieve the most appropriate tool from a set  $\mathcal{T} = \{T_1, T_2, \dots, T_n\}$ . The main components are an MLLM (*i.e.*, Qwen2.5-Omni [31]) that takes the user query as input to generate a structured description, and a retrieval module to select the tool. An overview of the pipeline is shown in Figure 2.3.

<sup>1</sup>By duplicated pairs, we refer to two or more occurrences of the same  $(text\_input, tool)$  pair for text-only queries or  $(text\_input, multimodal\_input, tool)$  set for multimodal queries.

<sup>2</sup>Models no longer supported on Hugging Face.

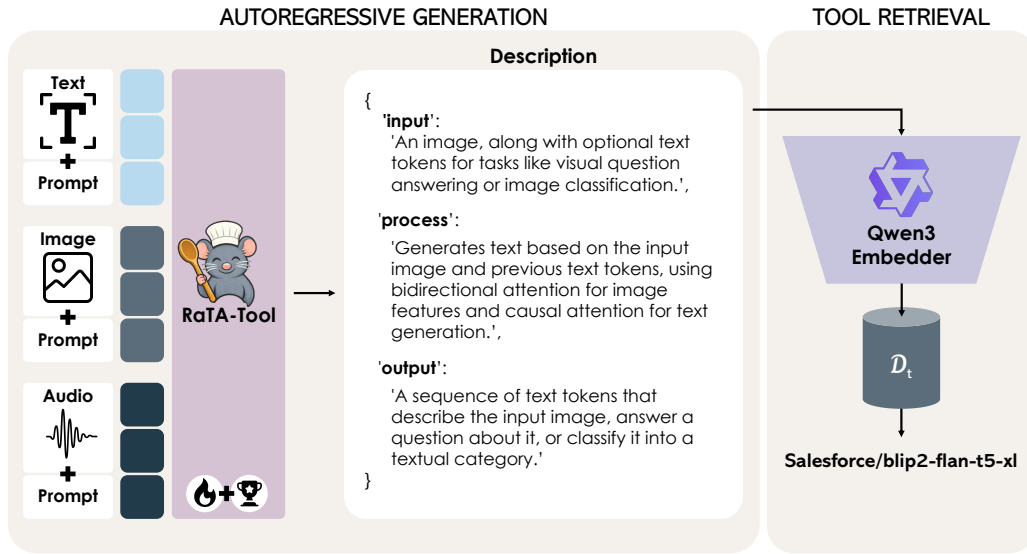


Figure 2.3: RaTA-Tool pipeline.

### 2.2.1 Query structured description

The user query  $q$  is first transformed into a structured JSON description to make it comparable to the tool descriptions. Given the input query and a custom prompt, the MLLM is prompted to generate a structured description following the standardized JSON format of the dataset:

- **Input:** The input modality and a brief description.
- **Process:** A summary of the user requested task.
- **Output:** The expected task output.

### 2.2.2 Retrieval

To perform the retrieval, both the query and the tools JSON descriptions are projected into a high-dimensional semantic space through a pre-trained embedding model (*i.e.*, Qwen3-Embedding [32]).

Let denote with  $E(\cdot)$  the embedding function, the similarity between a query description  $d_q$  and a tool description  $d_t$  is computed using cosine similarity:

$$\text{sim}(d_q, d_t) = \frac{E(d_q)E(d_t)}{\|E(d_q)\| \cdot \|E(d_t)\|} \quad (2.1)$$

Given the set of external tools  $\mathcal{T} = \{t_1, t_2, \dots, t_n\}$ , the selected tool  $t^*$  is the one maximizing the similarity score:

$$t^* = \arg \max_{t \in \mathcal{T}} \text{sim}(d_q, d_t) \quad (2.2)$$

### 2.2.3 Supervised Fine-Tuning (SFT)

The MLLM responsible for query task description generation was fine-tuned using our custom dataset. Specifically, for each query-tool pair, the multimodal query was included in the input prompt, while the corresponding JSON tool description served as the target response that the model was trained to generate. As a result, the MLLM learns to write task descriptions that closely match the actual descriptions of the tools. This training helps the model translate a user’s complex request into a clear summary designed to find and activate the appropriate external tool.

To ensure the model focuses exclusively on generating high-quality task descriptions, we employ label masking. Specifically, the loss is only calculated on the tokens corresponding to the task description, while the tokens belonging to the multimodal query and the system prompt are masked using a standard ignore index (*e.g.*, -100).

To perform the training stage, a standard autoregressive cross-entropy loss was employed:

$$\mathcal{L}_{\text{SFT}} = - \sum_{t=1}^T \log p_G(z_t | q, z_{<t}) \quad (2.3)$$

where  $p_G(z_t | q, z_{<t})$  represents the conditional probability of the next token  $z_t$  given the user query  $q$  and the sequence of previously generated tokens  $z_{<t}$ .

## 2.3 Direct Preference Optimization (DPO)

Direct Preference Optimization (DPO) optimizes the policy model by increasing the relative likelihood of a preferred response (*chosen*) over a dispreferred one (*rejected*), anchored by a static reference model. Within the RaTA-Tool framework, the task-description generator serves as the policy model, while the multimodal user query  $q$  is the input.

### 2.3.1 DPO dataset

DPO requires a specific preference dataset containing a prompt and two corresponding responses: *chosen* and *rejected*.

To construct this preference dataset, for each query  $q$ , the SFT model is prompted to generate multiple candidate task descriptions using different decoding strategies, including greedy decoding, beam-search and temperature-based sampling. For each generated task description  $s_i \in \mathcal{S} = \{s_1, \dots, s_k\}$  a *loss score* is computed. The similarity function shown in Equation 2.1 is computed between the embeddings of each task description  $s_i$  and the embeddings of the tools in  $\mathcal{T}$ . Then, for each query, the tools are ranked in descending order of similarity, where a lower rank indicates higher similarity.

Let  $\mathcal{T}_i^{\text{ranked}}$  denote the ordered set of ranked tools for task description  $s_i$ . Considering  $\mathcal{T}_q^* = \{t_{q1}^*, \dots, t_{qn}^*\}$  as the set of ground truth tools for query  $q$ , the

*loss score* of task description  $s_i$  can be defined as:

$$\mathcal{L}_i = \arg \min_{t_q \in \mathcal{T}_q^*} \mathcal{T}_i^{ranked} \quad (2.4)$$

The sampled response with the lowest *loss score* is selected as the *chosen* response. From the other sampled responses that have a different *loss score*, a random one is selected as the *rejected* option. If no difference in scores between the sampled responses exist, the query is discarded from the preference dataset.



## Chapter 3

# Experiments

## 3.1 Experimental setup

### 3.1.1 SFT dataset

To create a custom dataset to fine-tune the MLLM, we started from ToolMM-Bench and applied the pre-processing steps described in Sec. 2.1.2.

To generate the structured tool descriptions, we employed Qwen3-8B [11]. Both JSON and natural language (NL) descriptions were generated to test both formats. A one-shot example was injected into the prompt to align the model with the desired description structure. Refer to Appendix B, Figure B.1 and Figure B.2 for the prompts used to instruct Qwen3 to generate the descriptions.

The resulting data was divided into training, validation, and test splits.

**Table 2:** Dataset statistics.

	Unique Queries				Unique Tools				Total Items
	Text	Image	Audio	All	Text	Image	Audio	All	All
Training split	5,742	3,239	571	9,552	585	170	71	826	42,105
Validation split	639	336	54	1,029	499	140	51	690	4,436
Test split	3,231	922	190	4,343	66	20	8	94	5,396
Overall	9,612	4,497	815	14,924	651	190	79	920	51,937

Dataset statistics are reported in Table 2.

The first step was splitting the dataset into training and test sets with a 90/10 ratio, at the tool level. This matches the objective of open-world learning. In fact, by dividing the training and test sets based on tools, the test set is composed entirely of tools not seen during training, thus allowing us to measure the accuracy of predicting unseen tools. During this operation, particular caution was taken to maintain input modality balance: to avoid imbalances between the training and test sets, the 90/10 ratio split was applied to each input modality separately.

Subsequently, the training set was further divided into two splits: training and validation. Following the standard 90/10 ratio, this step involved splitting the training set at the query level. The resulting validation set was used during the supervised fine-tuning (SFT) and direct preference optimization (DPO) phases to compute validation metrics. Tools in the validation set were already seen during SFT, whereas the queries were not.

### **3.1.2 DPO dataset**

To perform direct preference optimization (DPO), a custom preference dataset was created as outlined in Section 2.3.1. The dataset is composed of 2,981 training items and 409 validation items. Table 3 shows the five different decoding strategy configurations.

## **3.2 Training strategies**

### **3.2.1 Supervised Fine-Tuning (SFT)**

Supervised Fine-Tuning (SFT) allowed to align MLLM generated task descriptions with tool descriptions, empirically increasing the semantic similarity

**Table 3:** Decoding strategies used to sample responses for the DPO preference dataset.

	Temperature	Beams
Greedy decoding	0	1
Beam search	0	5
Medium-temperature sampling	0.7	1
High-temperature sampling	1	1
Sampling + Beam search	1	3

between the query and the appropriate tool for solving that task.

The 3B and 7B parameter variants of Qwen2.5-Omni [31] were adopted as the foundational MLLMs, fine-tuned using Low-Rank Adaptation (LoRA) [18] with a rank  $r = 8$ , scaling factor  $\alpha = 16$  and dropout = 0.1. LoRA adapters were applied across all linear layers, including the attention projections  $(q, k, v, o)$ , as well as MLP components (gate, up, down).

LoRA reduced the number of trainable parameters from 8,953,244,160 to 21,430,272. Consequently, only 0.24% of all parameters were trained. This significantly improved memory efficiency, allowing to work with bigger models.

The fine-tuning was performed for 4 epochs on 2 NVIDIA L40S/A40 GPUs, with a learning rate of  $2 \times 10^{-4}$  and a weight decay of 0.01. Fused AdamW, a high-performance implementation of AdamW, was employed as the optimizer.

Brain Floating Point (BF16) [33] precision was used during training. BF16 uses 16 bits to represent each parameter as FP16, but unlike the latter, it employs 8 bits for the exponent, as in FP32. This is really useful in deep learning training, since it allows to represent the same range of numbers as FP32, while using half the memory.

### 3.2.2 Direct Preference Optimization (DPO)

Direct Preference Optimization (DPO) further aligned the model’s generation with tool matching. It was applied using LoRA with the same parameters used in SFT.

Both the 3B and 7B variants of Qwen2.5-Omni [31] were trained using DPO for 8 epochs on a single NVIDIA L40S/A40 GPU, with a learning rate of  $5 \times 10^{-5}$  and a weight decay of 0.01. BF16 [33] numerical precision was applied, as in SFT.

## 3.3 Inference paradigms

### 3.3.1 Query task description generation

**Zero-shot** Qwen2.5-Omni [31] without any adaptation was used in the first pipeline. This established a baseline from which to proceed with subsequent training and adaptations.

In the zero-shot pipeline, the model was instructed through a custom prompt to generate a description of the user query following either a JSON or NL structure.

**Few-shot** To align the MLLM’s query task description generation with tool descriptions, few-shot examples were employed. In this second version of the pipeline, two examples were used inside the prompt: one text-only and one image multimodal example. Each example is composed of a user query and its corresponding task description, derived from the ground-truth tool description.

**SFT and DPO** While the few-shot pipeline improved task description generation and subsequent prediction accuracy, training was a fundamental step to

adapt the model to our downstream task. Training strategies were consistently applied using few-shot prompts.

### 3.3.2 Retrieval

Given a task description generated from a user query, the retrieval phase predicts the best tool based on the semantic similarity between the task description and tool descriptions within the candidate tool set.

To measure the similarity, the descriptions are projected into high-dimensional semantic spaces, as outlined in Section 2.2.2.

Two pre-trained models were employed to create these embedding vectors: Contriever [34] and Qwen3-Embedding-8B [32].

The resulting embedding vectors have a dimensionality of 768 for Contriever and 4096 for Qwen3-Embedding.

## 3.4 Experimental results

A series of experiments was conducted to evaluate different RaTA-Tool configurations. In particular, the components compared were:

- **MLLM backbone:** Qwen2.5-Omni-3B [31] and Qwen2.5-Omni-7B [31].
- **Embedding backbone:** Contriever [34] and Qwen3-Embedding [32].
- **Stage of MLLM training:** Zero-shot, few-shot, supervised fine-tuning (SFT) and Direct Preference Optimization (DPO).
- **Description format:** Natural language (NL) and JSON.

Comparative results are detailed in the subsequent sections.

### 3.4.1 MLLM backbone and Embedding backbone

Qwen2.5-Omni-3B and Qwen2.5-Omni-7B [31] are employed as the MLLM backbones for structured task description generation. Comparative accuracy results are shown in Table 4, providing single-modality metrics (text, image, audio), average per query ( $\text{Avg}_q$ ), and average per modality ( $\text{Avg}_m$ ) across input modalities. The average per query is weighted by the number of items of each input modality in the test set (Table 2), whereas the average per modality is obtained by dividing the total by the number of modalities (*i.e.*, 3).

Moreover, two embedding models are compared: Qwen3-Embedding [32] which is used in the main configuration, and Contriever [34], to assess the influence of the embedding model on tool selection accuracy.

**Table 4:** Accuracy results on aggregate metrics: average per query ( $\text{Avg}_q$ ) and average per modality ( $\text{Avg}_m$ ). Comparative results between Qwen2.5-Omni-3B [31] and Qwen2.5-Omni-7B [31] as the MLLM backbone. Comparative results between Qwen3-Embedding [32] and Contriever [34] as the embedding backbone.

	Retriever	Training		Modalities			$\text{Avg}_q$	$\text{Avg}_m$
		SFT	DPO	Text	Image	Audio		
Qwen2.5-Omni-3B [31]	Contriever	-	-	17.7	8.0	36.3	16.5	20.7
Qwen2.5-Omni-7B [31]	Contriever	-	-	18.3	6.5	27.9	16.2	17.6
Qwen2.5-Omni-3B [31]	Qwen3-Embedding	-	-	24.9	31.3	69.0	28.2	41.7
Qwen2.5-Omni-7B [31]	Qwen3-Embedding	-	-	25.1	32.1	67.4	28.4	41.5
	Contriever	✓	-	57.1	37.0	80.5	53.8	58.2
	Qwen3-Embedding	✓	-	71.0	55.7	82.5	68.3	69.7
<b>RaTA-Tool-3B</b>	Qwen3-Embedding	✓	✓	<b>71.6</b>	<b>56.1</b>	<b>83.7</b>	<b>68.8</b>	<b>70.4</b>
	Contriever	✓	-	51.2	31.8	79.0	48.3	54.0
	Qwen3-Embedding	✓	-	69.3	<b>58.8</b>	<b>83.7</b>	67.7	70.6
<b>RaTA-Tool-7B</b>	Qwen3-Embedding	✓	✓	<b>71.6</b>	57.1	<b>83.7</b>	<b>69.1</b>	<b>70.8</b>

**Raw models: few-shot inference** First, the raw models are compared, with the Qwen2.5-Omni-7B variant achieving comparable results with respect to its

3B counterpart, using either Contriever or Qwen3-Embedding as the embedding backbone. Qwen3-Embedding demonstrates instead a substantial advantage over Contriever, with a +11.7%  $\text{Avg}_q$  for the 3B parameter model and a +12.2% increase for the 7B variant.

**SFT models** Subsequently, RaTA-Tool-3B and RaTA-Tool-7B, representing the framework with adapted MLLMs via SFT and DPO, are compared. Here, for the SFT-only version, Qwen3-Embedding clearly overperform Contriever again, with an  $\text{Avg}_q$  of 68.3 against 53.8 for Contriever (+14.5%) for the 3B parameter model, and 67.7 against 48.3 (+19.4%) for the 7B variant.

**DPO models** The models were further aligned with DPO. Specifically, given the clear superiority of Qwen3-Embedding over Contriever, the preference dataset was constructed based on loss scores computed against tool embedding vectors generated by Qwen3-Embedding. Through DPO, a slight increase is noticeable in both backbone variants, with a +0.5%  $\text{Avg}_q$  gain for Qwen2.5-Omni-3B and +1.4% for Qwen2.5-Omni-7B. RaTA-Tool-7B shows a slightly better overall performance compared to the 3B counterpart.

**Overall considerations** Experimental results outline important factors influencing performance: a strong embedding model is essential for retrieval-based tool selection, embedding dimension represents a foundational value, and preference-based alignment provides additional performance gains over supervised fine-tuning alone.

### 3.4.2 Description formats

The main RaTA-Tool configuration uses a structured JSON description with three fields (*input*, *process*, *output*) to represent both tools and queries.

## Experiments

Experiments were conducted on natural language (NL) descriptions to examine how description format can affect retrieval performance.

**Table 5:** Experiments on the effect of query and tool descriptions format and prompting or fine-tuning strategy on tool-selection performance, reported for both embedding backbones (Contriever and Qwen3-Embedding). Results are shown for text, image, and audio queries, together with per query ( $Avg_q$ ) and per modality ( $Avg_m$ ) aggregate metrics.

	Retriever	Descriptions		Modalities			$Avg_q$	$Avg_m$
		Type	Mode	Text	Image	Audio		
	Contriever	NL	Zero-shot	16.8	17.2	43.7	18.1	25.9
	Contriever	NL	Few-shot	13.3	13.7	49.5	15.0	25.5
	Contriever	JSON	Zero-shot	18.3	6.5	27.9	16.2	17.6
	Contriever	JSON	Few-shot	21.0	15.1	35.8	20.4	24.0
	Contriever	NL	SFT	34.3	<b>50.3</b>	79.0	39.6	<b>54.5</b>
<b>RaTA-Tool-7B</b>	Contriever	JSON	SFT	<b>51.2</b>	31.8	<b>79.0</b>	<b>48.3</b>	54.0
	Qwen3-Emb	NL	Zero-shot	27.6	32.0	72.1	30.5	43.9
	Qwen3-Emb	NL	Few-shot	31.5	22.1	72.1	31.3	41.9
	Qwen3-Emb	JSON	Zero-shot	25.1	32.1	67.4	28.4	41.5
	Qwen3-Emb	JSON	Few-shot	25.9	38.7	70.0	30.6	44.9
	Qwen3-Emb	NL	SFT	41.6	44.0	<b>95.8</b>	44.5	60.5
<b>RaTA-Tool-7B</b>	Qwen3-Emb	JSON	SFT	<b>69.3</b>	<b>58.8</b>	83.7	<b>67.7</b>	<b>70.6</b>

Results reported in Table 5 show the effect of different description formats under zero-shot and few-shot prompting, as well as with the fine-tuned models. The performance of both embedding models is reported.

In the zero-shot and few-shot prompting scenarios, JSON and NL descriptions perform similarly. Indeed, with the Contriever embedding, NL ( $Avg_q = 18.1$ ) outperforms JSON ( $Avg_q = 16.2$ ) by +1.9% in the zero-shot setting, whereas the opposite happens in the few-shot setting, where JSON ( $Avg_q = 20.4$ ) significantly outperforms NL ( $Avg_q = 15.0$ ) by +5.4%. With Qwen3-Embedding, NL outperforms JSON in both zero-shot and few-shot settings, with a +2.1% and +0.7% overall average per query, respectively.

Nevertheless, this tendency reverses with the SFT models of both sizes. Here, the Contriever JSON configuration has an average accuracy per query of 48.3 versus 39.6 for NL (+ 8.7%). Similarly, with Qwen3-Embedding, the JSON Avg<sub>q</sub> stands at 67.7 while NL stands at 44.5 (+ 23.2%). This result suggests that a stronger embedding model (Qwen3-Embedding) can create deeper semantic representations of already well structured descriptions with key information that helps in the tool retrieval goal. However, when the description lacks information or contains unnecessary details, the model’s dimension and embedding power cannot compensate.

### 3.4.3 In-context learning and supervised fine-tuning

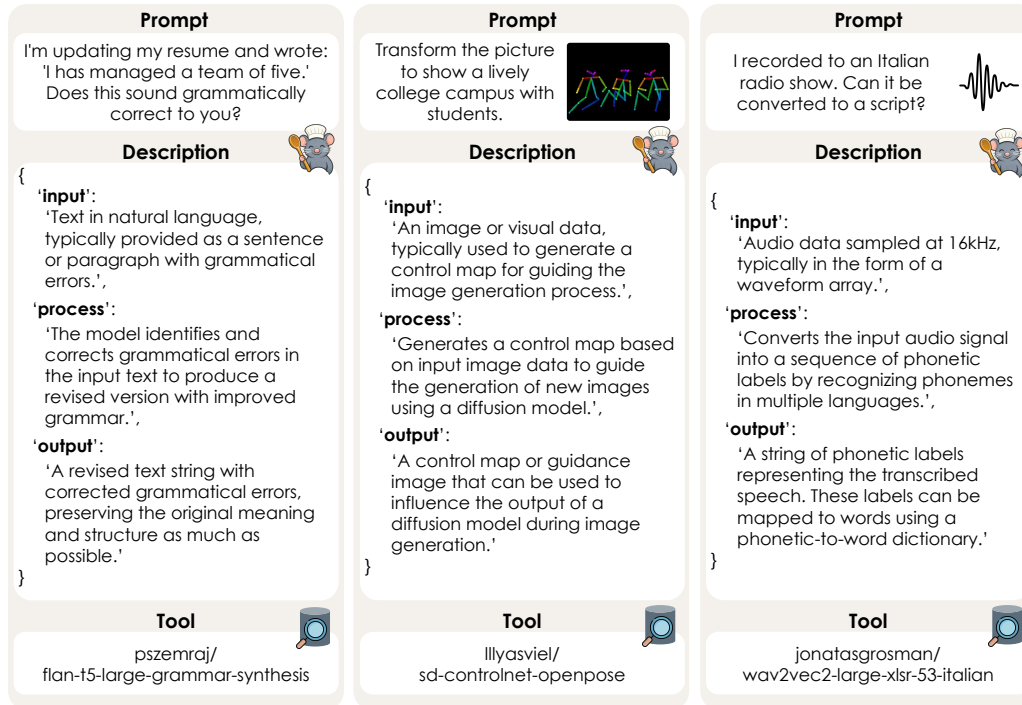
The results presented in Table 5 outline the performance across zero-shot, few-shot, and SFT generated query task descriptions.

Few-shot prompting almost always improves retrieval performance: for instance, JSON descriptions embedded with Contriever achieve an average accuracy per query of 16.2 in the zero-shot setting and 20.4 with in the few-shot setting, improving by +4.2%. Similarly, Qwen3-Embedding with JSON descriptions yields an accuracy of 28.4 in the zero-shot setting versus 30.6 in the few-shot setting, reflecting an increase of +2.2%.

SFT substantially improves model performance with both embedding models. In particular, with JSON descriptions SFT yields a +27.9% gain with Contriever compared to the previous best average accuracy per query, and a +37.1% gain using Qwen3-Embedding.

These results highlight the importance of in-context examples, as well as a specific training strategy, in aligning the model with downstream task needs, leading to the generation of more discriminative task descriptions for retrieval. Overall, the best results are achieved by combining SFT with structured JSON descriptions, consistently across both retrieval backbones.

### 3.5 Qualitative results



**Figure 3.1:** Qualitative results of RaTA-Tool under different input modalities.

Qualitative results for RaTA-Tool are presented in Figure 3.1. Examples include both text-only queries and multimodal queries with image and audio inputs. The model generates a structured JSON description of the user query, including *input*, *process*, and *output* fields. The embedding vector of this description is compared with all tool vectors, to retrieve the highest similarity tool reported in the figure. Across all examples, the retrieved tool corresponds to the ground truth tool. Additional qualitative examples can be found in Appendix A.

## Chapter 4

# Conclusions

### 4.1 Key findings and implications

The objective of this thesis was to lay a solid foundation for the use of external tools by foundation models that supported both text and multimodal inputs. In contrast to existing solutions, this work focused on retrieval-based tool selection, allowing controllers in an agentic framework to retrieve unseen tools. To achieve this, we proposed a novel framework, RaTA-Tool, composed of an MLLM aligned with SFT and DPO, and a retrieval backbone. To fine-tune the model, a custom dataset was developed.

The experimental results outline the effectiveness of this approach. RaTA-Tool demonstrated strong capabilities in tool selection for all input modalities involved. Moreover, qualitative results showed the effectiveness of the model training, with structured descriptions that capture more semantic details when generated by the aligned model compared to the zero-shot baseline.

## 4.2 Limitations and future work

While RaTA-Tool exhibits great capabilities in tool selection, it is not without limitations. Its use is currently limited to single-turn conversations, lacking the ability to adjust its selection based on subsequent user requests or additional constraints.

Furthermore, RaTA-Tool focuses on tools that are Hugging Face models, leaving out traditional tools or APIs.

Another limitation is its restriction to single-tool solutions. RaTA-Tool cannot plan a series of actions involving multiple tools to solve a single complex task. This is a fundamental requirement for a controller within an agentic framework, in order to tackle practical problems that often require a series of tools to be executed in sequence.

Future work may concentrate on expanding RaTA-Tool’s capabilities beyond tool selection and towards a fully autonomous agent.

## 4.3 Final remarks

In conclusion, this thesis presented a novel framework, RaTA-Tool, for effective open-world retrieval-based multimodal tool selection that reasons over structured descriptions. To support evaluation in this setting, we introduced the first benchmark for open-world multimodal tool use with standardized tool descriptions. Our results demonstrate that semantic, description-driven tool reasoning is a promising direction for scalable and robust tool-learning systems.

## Appendix A

# Qualitative examples

This section presents additional qualitative examples, including both text-only and multimodal queries. Figure A.1 and Figure A.2 show examples of text-only and image multimodal queries respectively, along with their corresponding generated JSON description and the retrieved tool. In all cases, RaTA-Tool retrieved the correct tool.

Figure A.3 compares the zero-shot, RaTA-Tool, and ground truth task descriptions (*i.e.*, the ground truth tool description) for a given user query. This example highlights how RaTA-Tool not only correctly structures the description but also captures deeper semantic details, thereby improving retrieval accuracy.

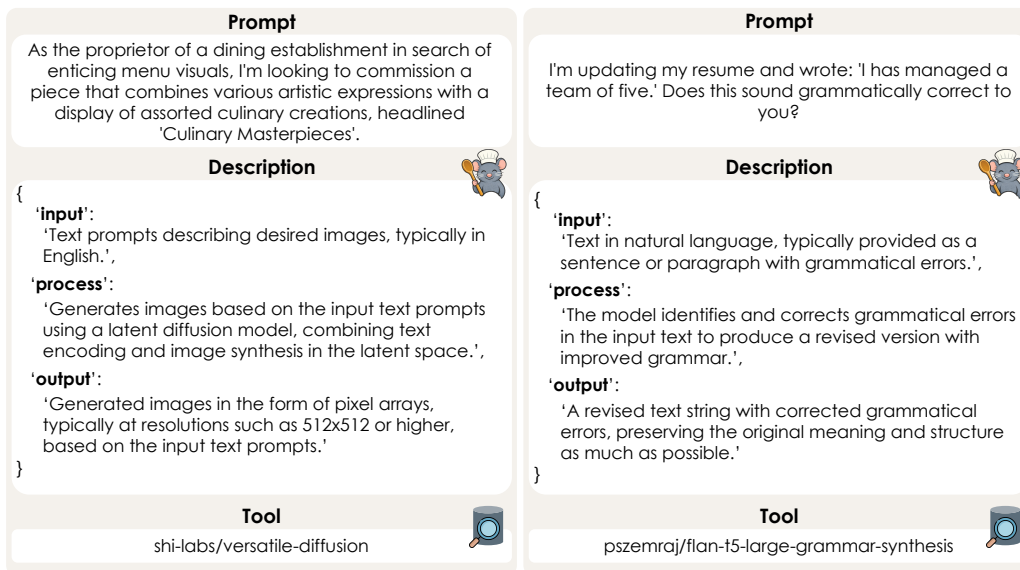


Figure A.1: Qualitative examples of RaTA-Tool for text-only queries.

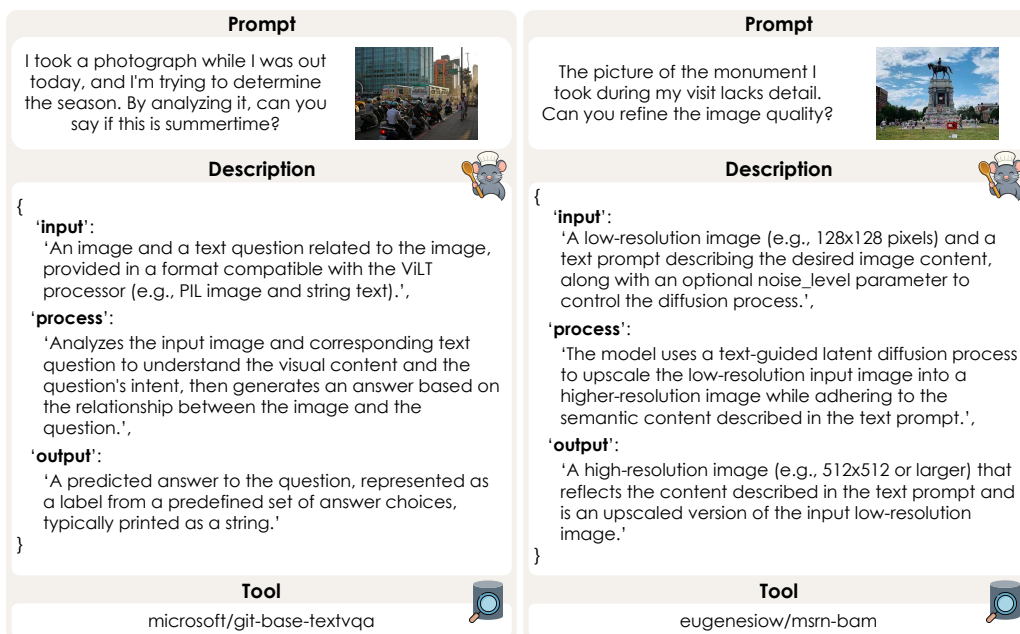
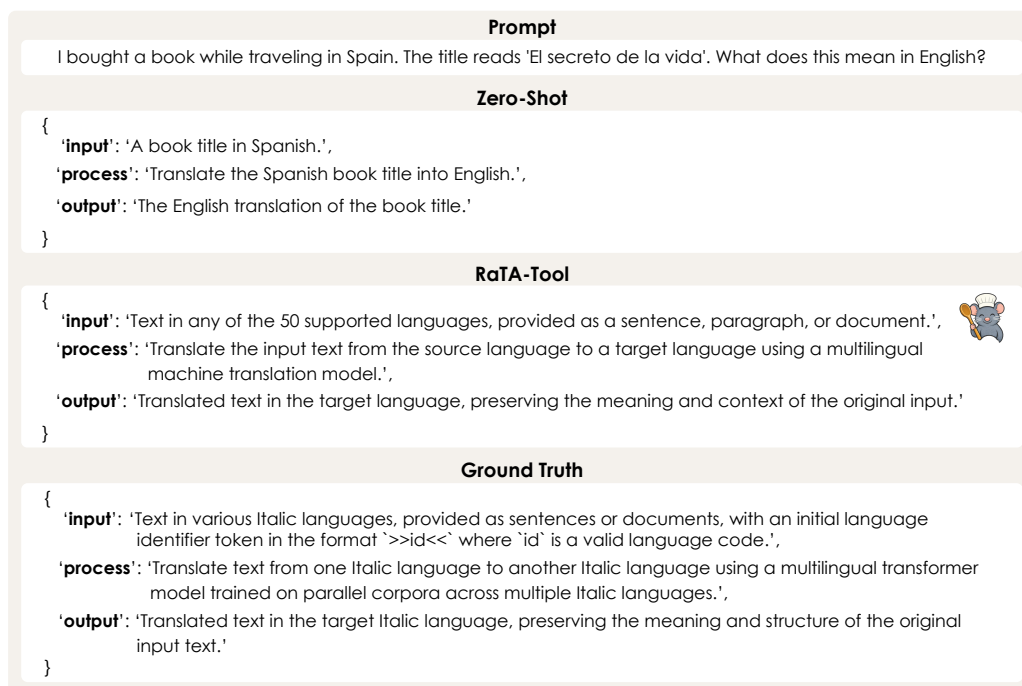


Figure A.2: Qualitative examples of RaTA-Tool for multimodal queries with image inputs.



**Figure A.3:** Qualitative comparison of zero-shot task description and RaTA-Tool one, along with the ground-truth description.



## Appendix B

# Prompts

This section provides examples of the different prompts used for generation. Both JSON and natural language (NL) versions are shown for each prompt.

Figure B.1 and Figure B.2 present examples of prompts for tool description generation. They are used to instruct the model to generate tool descriptions from the information contained in the tool’s Hugging Face model card.

Figure B.3 and Figure B.4 show examples of prompts used for query task description generation.

## JSON Tool Description Prompt

```
You are an AI assistant tasked with extracting structured information from Hugging Face model cards.
Given the model card below, analyze it thoroughly and generate a JSON object that describes the model in terms of:
1. Input: What type of data the model takes in (e.g., text, image, audio, structured data, etc.), including format or shape if mentioned.
2. Process: What the model does with the input (e.g., sentiment classification, translation, object detection, etc.), described concisely.
3. Output: What the model produces as output, including the format, labels or value types where applicable.

Your output must be:
- Model-agnostic (do not mention brand names or model IDs)
- Written strictly in English
- A valid JSON object in the following format:
{
  "input": "<description of the expected input>",
  "process": "<brief description of what the model does>",
  "output": "<description of the output produced>"
}

EXAMPLE
Model card: DeepPavlov/rubert-base-cased HuggingFace model card

Model description:
{
  "input": "Text written in Russian, provided as a sentence, paragraph, or document.",
  "process": "Analyze the linguistic and semantic content of the Russian text and transform it into a multidimensional vector representation that captures meaning, syntax, and contextual nuances.",
  "output": "A numerical embedding vector that encodes the semantic information of the input text, suitable for downstream tasks such as similarity comparison, clustering, or classification."
}

Now do the same for the model card below:
Model card: {model_card}
Model description:
```

**Figure B.1:** Example of prompt for tool description generation in JSON.

---

## NL Tool Description Prompt

You are an AI assistant tasked with summarizing Hugging Face model cards. Given the Model Card below, analyze it thoroughly and write a concise, natural language description that explains what kind of input the model expects, what it does with that input, and what it produces as output.

Write the description in clear and complete sentences, using natural language. Avoid technical formatting such as JSON or bullet points. Focus only on information that is explicitly stated or clearly implied in the model card. The description must be written strictly in English.

### EXAMPLE

Model card:

DeepPavlov/rubert-base-cased HuggingFace model card

Model description:

This model processes Russian language text by converting it into a dense semantic vector representation, allowing the meaning of the text to be captured and expressed as an embedding.

Now do the same for the model card below:

Model card: {model\_card}

Model description:

**Figure B.2:** Example of prompt for tool description generation in natural language.

## JSON Inference Prompt

Given the user query below, write a complete and precise description of the task requested by the user in JSON format. The JSON object must have three fields:

- input: The nature and format of the input the model will receive.
- process: The process or transformation the model is expected to perform.
- output: The nature and format of the output the model should produce.

Each field must contain short, clear sentences written in English. Do not include specific details about the current query; keep the description applicable to similar tasks of this type.

Your output must be a valid JSON object in the following format:

```
{
  "input": "<general description of the input>"
  "process": "<general description of the reasoning or steps>"
  "output": "<general description of the expected output>"
}
```

### EXAMPLES

Query 1: I'm watching a WWII movie with a German phrase 'Hände hoch. Ich kann dir nicht wehtun.' Can you translate this into Spanish for me?

Task Description 1:

```
{
  "input": "Text in any of the West Germanic languages, including Afrikaans, Dutch, German, and others, provided as sentences or documents. The input must include a language identifier token in the format '>id<' where 'id' corresponds to a valid target language.",
  "process": "Translate text from one West Germanic language to another West Germanic language using a multilingual transformer model trained on parallel corpora.",
  "output": "Translated text in the target West Germanic language, corresponding to the input text. The output format is text, and the translation quality is measured using metrics like BLEU and chr-F."
}
```

Query 2: Show this picture as a group of friends gathered around a campfire, telling stories.

Query 2 Image: {query\_2\_image}

Task Description 2:

```
{
  "input": "An image in the form of a normal map, which represents surface orientation information using a 3D vector for each pixel, typically in a grayscale or RGB format.",
  "process": "The model processes the normal map input to guide the diffusion model in generating an image that aligns with the surface orientation information provided by the normal map.",
  "output": "A generated image that adheres to the surface orientation details specified in the normal map input, typically in RGB format."
}
```

Now do the same for the query below:

Query: {query}

**Figure B.3:** Example of prompt for task-description generation at inference time in JSON.

---

## NL Inference Prompt

Given the user query below, write a complete and precise description of the task requested by the user. Your description must clearly specify:

- The nature and format of the input the model will receive
- The process or transformation the model is expected to perform
- The nature and format of the output the model should produce

Write the description in clear and direct natural language using full sentences. Do not use code, lists, or technical notation.

Write strictly in English.

### EXAMPLES

Query 1: I'm watching a WWII movie with a German phrase 'Hände hoch. Ich kann dir nicht wehtun.' Can you translate this into Spanish for me?

Task description 1: This model processes text in various West Germanic languages by translating input text from one West Germanic language into another West Germanic language.

The input text must include a language identifier at the beginning of each sentence to specify the source language. The model transforms the input by generating a corresponding translation in the target language. The output is the translated text in the specified target West Germanic language.

Query 2: Show this picture as a group of friends gathered around a campfire, telling stories.

Query 2 image: {query\_2\_image}

Task description 2: This model processes visual input in the form of normal maps, which are images that encode surface orientation information. It is designed to work in conjunction with a text-to-image diffusion model to guide the generation of images based on the provided normal map. The model interprets the normal map to influence the texture and structure of the generated image, producing a realistic visual output that aligns with the input conditions. The output is a fully rendered image that reflects both the textual description and the spatial orientation details provided by the normal map.

Now do the same for the query below:

Query: {query}

**Figure B.4:** Example of prompt for task-description generation at inference time in natural language.



## Appendix C

# Publications

The research conducted during the development of this thesis resulted in the paper: Gabriele Mattioli, Evelyn Turri, Sara Sarto, Lorenzo Baraldi, Marcella Cornia, Lorenzo Baraldi and Rita Cucchiara. “RaTA-Tool: Retrieval-based Tool Selection with Multimodal Large Language Models”. Submitted to the International Conference on Pattern Recognition (ICPR) 2026.

# RaTA-Tool: Retrieval-based Tool Selection with Multimodal Large Language Models

Gabriele Mattioli<sup>1</sup>, Evelyn Turri<sup>1</sup>[0009-0005-5668-0839],  
Sara Sarto<sup>1</sup>[0000-0003-1057-3374], Lorenzo Baraldi<sup>2</sup>[0009-0000-4658-8928],  
Marcella Cornia<sup>1</sup>[0000-0001-9640-9385], Lorenzo Baraldi<sup>1</sup>[0000-0001-5125-4957],  
and Rita Cucchiara<sup>1</sup>[0000-0002-2239-283X]

<sup>1</sup> University of Modena and Reggio Emilia, Italy

{name.surname}@unimore.it

<sup>2</sup> University of Pisa, Italy

{name.surname}@phd.unipi.it

**Abstract.** Tool learning with foundation models aims to endow AI systems with the ability to invoke external resources – such as APIs, computational utilities, and specialized models – to solve complex tasks beyond the reach of standalone language generation. While recent advances in Large Language Models (LLMs) and Multimodal Large Language Models (MLLMs) have expanded their reasoning and perception capabilities, existing tool-use methods are predominantly limited to text-only inputs and closed-world settings. Consequently, they struggle to interpret multimodal user instructions and cannot generalize to tools unseen during training. In this work, we introduce RaTA-Tool, a novel framework for *open-world multimodal tool selection*. Rather than learning direct mappings from user queries to fixed tool identifiers, our approach enables an MLLM to convert a multimodal query into a structured task description and subsequently retrieve the most appropriate tool by matching this representation against semantically rich, machine-readable tool descriptions. This retrieval-based formulation naturally supports extensibility to new tools without retraining. To further improve alignment between task descriptions and tool selection, we incorporate a preference-based optimization stage using Direct Preference Optimization (DPO). To support research in this setting, we also introduce the first dataset for open-world multimodal tool use, featuring standardized tool descriptions derived from Hugging Face model cards. Extensive experiments demonstrate that our approach significantly improves tool-selection performance, particularly in open-world, multimodal scenarios.

**Keywords:** Tool Agent Learning · Retrieval-Augmented Selection · Multimodal Large Language Models.

## 1 Introduction

Humans possess a remarkable ability to invent, adapt, and employ tools, allowing them to solve problems far beyond their innate physical capabilities. As foundation models continue to evolve [3,24], AI systems are starting to acquire

comparable capabilities. This emerging direction, commonly referred to as *tool learning with foundation models* [28], seeks to merge the general reasoning abilities of large models with the precision and specialization of external tools. By granting models access to dynamic knowledge bases [36,34,26] and computational utilities [36], tool-augmented systems can perform tasks far beyond the scope of standalone language generation.

Exploiting the already extended capabilities of Large Language Models (LLMs) [1,4,13,41], some works integrate in these models tool-use mechanisms to support real-world tool use at scale. This extension has opened the door to a broad range of applications. Within this broader landscape, Multimodal Large Language Models (MLLMs) [5,8,22] play a pivotal role. By jointly processing text, images, audio, and other modalities [7,12,14,19], MLLMs unify diverse tasks under a single architecture and have achieved impressive performance across conversational, analytical, and reasoning applications, particularly when trained with instruction-following paradigms.

Despite these advances, most existing tool-use frameworks operate under two restrictive assumptions. First, they largely rely on textual inputs [29], limiting their ability to interpret multimodal instructions where visual or audio information is essential for disambiguation. Second, they adopt a closed-world formulation: tool selection is treated as a supervised classification problem over a fixed set of tools observed during training. As a result, these systems struggle to generalize to newly introduced tools or to tasks whose required capabilities were not explicitly represented in the training distribution.

Recent works have started to explore multimodal tool selection by training MLLMs to map multimodal inputs directly to tool identifiers [38]. While effective in controlled evaluation settings, such approaches fundamentally inherit the limitations of closed-set learning. In contrast, real-world tool ecosystems are dynamic, continually evolving as new tools are introduced and existing ones are updated. This gap highlights the need for an *open-world* formulation of multimodal tool selection, where models can reason over previously unseen tools without retraining.

To address this challenge, we introduce RaTA-Tool, a novel retrieval-based framework for open-world multimodal tool selection. Rather than learning direct mappings from user queries to tool identifiers, our approach enables an MLLM to reinterpret a multimodal user query into a structured task description that captures the underlying intent and required capabilities. This task representation is then matched against a collection of semantically rich, machine-readable tool descriptions using embedding-based retrieval, allowing the system to naturally generalize to unseen tools. To further enhance performance and refine generated task descriptions, we incorporate a preference-based optimization stage using Direct Preference Optimization (DPO) [30].

Finally, to support systematic evaluation in this setting, we introduce the first dataset for open-world multimodal tool selection. It features standardized, structured tool descriptions derived from Hugging Face model cards, enabling reproducible and extensible evaluation beyond standard classification pipelines.

**Contributions.** In summary, our contributions are as follows:

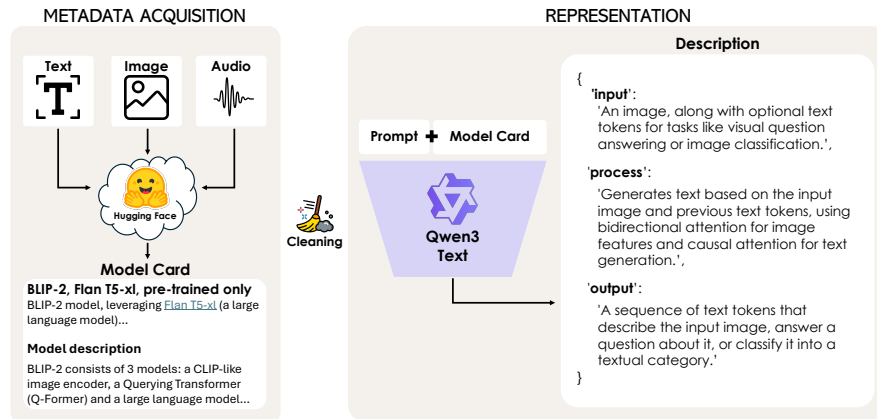
- We introduce a novel framework for open-world multimodal tool selection that reasons over structured tool descriptions instead of relying on closed-set, supervised query-tool mappings.
- We propose a retrieval pipeline in which an MLLM converts user queries into task descriptions that are matched to tools via embedding-based similarity.
- We enhance tool-selection performance through a preference-based alignment stage using DPO.
- We construct the first benchmark tailored for open-world multimodal tool selection, featuring standardized, machine-readable tool descriptions.

## 2 Related Work

**Multimodal Large Language Models.** Recent years have witnessed rapid progress in extending LLMs beyond text to support multimodal inputs such as images, audio, and video. Early approaches align pretrained language backbones with modality-specific encoders through cross-modal projection and instruction tuning. Representative models such as LLaVA [23,22], BLIP [21], and Flamingo [2] demonstrate strong performance in multimodal perception, reasoning, and generation while largely preserving frozen or lightly adapted encoders.

Subsequent works have extended MLLMs beyond images to additional modalities such as audio and video by aligning pretrained language backbones with modality-specific encoders through adapter or connector modules [6,45]. Complementary efforts focus on video-centric modeling, addressing temporal dynamics and audio-visual fusion via specialized encoders and spatiotemporal alignment mechanisms [7,43]. More recent approaches move toward unified multimodal perception and generation within a single architecture, supporting tightly coupled reasoning across text, vision, audio, and speech [12,19]. In parallel, alternative designs explore modality-agnostic representations by removing modality-specific components altogether and introducing universal encoders capable of handling diverse input types [14]. While these models significantly enhance multimodal understanding, they primarily focus on improving perception and reasoning within the model itself. They do not explicitly address how multimodal inputs can be leveraged to select and invoke external tools, nor do they consider open-world settings where tools are unseen during training.

**Tool Use with Language Models.** Parallel to advances in LLM capabilities, a growing body of work explores the use of language models as agents that invoke external tools to solve complex tasks. Toolformer [32] introduces the idea of teaching LLMs to decide when and how to call APIs through self-supervised annotation. Subsequent systems extend this paradigm to broader application domains such as code generation, web search, healthcare, and task automation. Several works treat tool usage as either retrieve an API or planning actions from a predefined set of tools. Gorilla [27] focuses on accurate API selection in machine learning workflows, while ToolLLM [29] emphasizes multi-step task execution using predefined tool inventories. On a different line, systems such



**Fig. 1.** Pipeline of dataset creation. We first perform a metadata acquisition step by scraping model cards from Hugging Face. The collected data then undergoes a cleaning stage. Finally, given a user query and its associated prompt, we feed this information into an LLM to generate a structured tool description in a standardized JSON format.

as HuggingGPT [33], Visual GPT [39] and GPT-4Tools [42] decompose user requests into subtasks and route them to specialized models. Despite their effectiveness, these methods primarily operate on textual inputs and therefore cannot directly leverage multimodal cues present in many real-world user instructions.

**Multimodal Tool Agents.** More recently, multimodal tool agents extend tool-use paradigms to MLLMs, enabling direct perception of visual or audio inputs when selecting tools. MLLM-Tool [38] represents a first step in this direction by training MLLMs to select tools based on multimodal queries. However, existing approaches largely rely on direct mappings from queries to predefined tool identifiers, inheriting the limitations of closed-set learning and restricting generalization to unseen tools. In contrast, our work addresses *open-world* multimodal tool selection, where relevant tools may not be observed during training. Rather than predicting tool names or IDs, we transform multimodal user queries into structured task descriptions and retrieve suitable tools via semantic matching against rich, machine-readable tool representations. This retrieval-based formulation decouples tool selection from fixed inventories and enables scalable generalization to newly introduced tools.

### 3 Proposed Method

#### 3.1 Dataset Creation

**Preliminaries.** As an initial step, we build upon ToolMMBench [38], a publicly available dataset that links multimodal user queries to Hugging Face models. The dataset is designed to evaluate LLM ability to select external tools when faced with text-ambiguous queries.

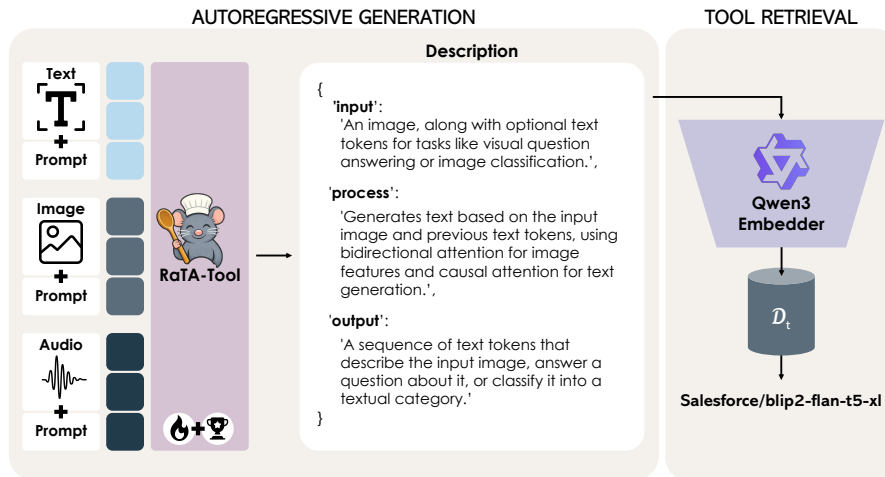
KHA-WHITE/MANGA-OCR-BASE	
Model Card	Description
<p><b>Manga OCR</b> Optical character recognition for Japanese text, with the main focus being Japanese manga.</p> <p>It uses <a href="#">Vision Encoder Decoder</a> framework.</p> <p>Manga OCR can be used as a general purpose printed Japanese OCR, but its main goal was to provide a high quality text recognition, robust against various scenarios specific to manga:</p> <ul style="list-style-type: none"> <li>• both vertical and horizontal text</li> <li>• text with furigana</li> <li>• text overlaid on images</li> <li>• wide variety of fonts and font styles</li> <li>• low quality images</li> </ul> <p>Code is available <a href="#">here</a>.</p>	<pre>{   'input':     'An image, along with optional text tokens for tasks like visual question answering or image classification.',   'process':     'Generates text based on the input image and previous text tokens, using bidirectional attention for image features and causal attention for text generation.',   'output':     'A sequence of text tokens that describe the input image, answer a question about it, or classify it into a textual category.' }</pre>
VALHALLA/T5-BASE-QG-HL	
Model Card	Description
<p><b>T5 for question-generation</b> This is <a href="#">t5-base</a> model trained for answer aware question generation task. The answer spans are highlighted within the text with special highlight tokens. You can play with the model using the inference API, just highlight the answer spans with <code>&lt;h1&gt;</code> tokens and end the text with <code>&lt;/s&gt;</code>.</p> <p>For example <code>&lt;h1&gt; 42 &lt;h1&gt; is the answer to life, the universe and everything. &lt;/s&gt;</code></p> <p>For more details see <a href="#">this repo</a>. Model in action 🚀 You'll need to clone the <a href="#">repo</a>.</p>	<pre>{   'input':     'Text containing answer spans highlighted with &lt;h1&gt; tokens, ending with &lt;/s&gt; token.',   'process':     'Generates a question based on the provided text and the highlighted answer span.',   'output':     'A dictionary containing the generated question and the corresponding answer.' }</pre>

**Fig. 2.** Qualitative dataset examples illustrating Hugging Face model cards and the corresponding tool descriptions represented in the standardized JSON format.

Let  $\mathcal{Q} = \{q_1, q_2, \dots, q_M\}$  denote the set of user queries and  $\mathcal{T} = \{t_1, t_2, \dots, t_N\}$  the set of candidate tools. ToolMMBench provides an initial collection of query-tool associations organized into modality-specific splits, including *text-only*, *text-image*, and *text-audio* inputs. In each split, a user textual query  $q \in \mathcal{Q}$ , optionally accompanied by multimodal input, is associated with a single external tool  $t \in \mathcal{T}$ .

Although this supervised formulation enables controlled evaluation of tool selection, it inherently restricts generalization to tools observed during training. In particular, models trained under this setting can only select from a fixed, closed set of tools, limiting their ability to handle user queries that require previously unseen tools or newly introduced capabilities. This closed-world assumption motivates the need for methods that generalize beyond the training tool set by reasoning over tool descriptions rather than relying on query-tool associations.

Therefore, we use this resource solely as an initial foundation, upon which we perform extensive pre-processing and adaptation to construct a dataset aligned with our experimental setting. To align the original dataset with our experimental setting, we first standardize the structure of all query-tool pairs to match our required format. We then perform dataset-wide cleaning procedures to remove duplicated or inconsistent samples, ensuring that each query is associated with a single, well-defined tool description.



**Fig. 3.** Overview of our pipeline. RaTA-Tool supports multimodal user queries by jointly processing the input prompt and associated modalities. A fine-tuned LLM encodes the combined inputs to generate a structured task description of the user request, further aligned via Direct Preference Optimization (DPO). This description is embedded and used to retrieve the most relevant tool from an external tool collection.

**Tool Metadata Acquisition and Representation.** Each tool in our dataset corresponds to a Hugging Face model. To obtain descriptive metadata, we programmatically scraped the model card associated with each tool. Model cards typically contain semi-structured natural language descriptions of the model functionality, intended use cases, input modalities, and expected outputs. However, this information is not directly suitable for automated reasoning or structured tool selection.

To address this limitation, we employ an LLM (*i.e.*, Qwen3 [41]) to transform each model card into a standardized JSON representation. Specifically, given the user query and the corresponding model card as input, the LLM is prompted to extract and generate a structured description comprising three fields: ① **input**, the expected input format and modalities supported by the tool; ② **output**, the type and structure of the generated output; and ③ **process**, a concise functional summary of the tool capabilities.

This structured representation enables consistent and machine-readable tool descriptions, facilitating downstream tool selection and reasoning tasks. An overview of the dataset creation pipeline is shown in Fig. 1, while some qualitative examples from our dataset are shown in Fig. 2.

### 3.2 The RaTA-Tool Framework

Given an input query  $q \in \mathcal{Q}$ , our objective is to retrieve the most appropriate tool for a user-requested task, from a collection of candidate descriptions. Overall, RaTA-Tool comprises two components: an MLLM that takes as input the user

query and a retrieval module which selects the most relevant tool. An overview of the proposed framework is shown in Fig. 3.

**Multimodal Integration and Autoregressive Generation.** A MLLM typically receives a query composed of one or more input modalities (for example, text, image, or audio) and generates a textual output in an autoregressive manner. The textual prompt usually includes a pre-defined system-level prompt and a question related to the input modality, given by the user. Specifically for our use case, given the input query and a specific prompt, an MLLM  $\mathcal{G}$  is trained via a supervised fine-tuning (SFT) stage to generate a structured task description, following the standardized JSON format of our dataset. Formally,  $\tilde{d} = \mathcal{G}(q)$ .

Let  $\tilde{d} = \{z_1, z_2, \dots, z_T\}$  denote the sequence of tokens composing the task description. Training is performed using a standard autoregressive cross-entropy loss, defined as

$$\mathcal{L}_{\text{SFT}} = - \sum_{t=1}^T \log p_{\mathcal{G}}(z_t | q, z_{<t}), \quad (1)$$

where  $p_{\mathcal{G}}(z_t | q, z_{<t})$  denotes the conditional probability of generating the next token given the input query and the previously generated tokens.

The generated descriptions capture the expected input, output, and functional intent of the user’s task, forming the basis for subsequent tool retrieval.

**Tool Collection Retrieval.** The external tool collection comprises a set of tool description  $\mathcal{D}_i = \{d_{t_1}, d_{t_2}, \dots, d_{t_N}\}$ , each represented in the standardized JSON format. The goal of the retrieval module is to identify the tool whose description best aligns with the model-generated task description  $\tilde{d}$ .

To compare descriptions, we embed both  $\tilde{d}$  and each  $d_t \in \mathcal{D}_i$  using a pre-trained embedding model (*i.e.*, Qwen3-Embedding [44]). Let  $E(\cdot)$  denote the embedding function. The similarity between the generated query description and a candidate tool description is modeled as the inner product between their respective embeddings:

$$\text{sim}(\tilde{d}, d_t) = E(\tilde{d}) \cdot E(d_t). \quad (2)$$

The selected tool is the one maximizing this similarity score:

$$t^* = \arg \max_{t \in \mathcal{T}} \text{sim}(\tilde{d}, d_t). \quad (3)$$

### 3.3 Tool Selection via Direct Preference Optimization

**Preliminaries.** Direct Preference Optimization (DPO) [30] has emerged as an effective alternative to Reinforcement Learning from Human Feedback (RLHF) [20] for aligning language models with preference signals. Unlike RLHF, which relies on explicit reward modeling and policy optimization, DPO directly optimizes the model by contrasting preferred and dispreferred responses. Although initially proposed for text-only LLMs, DPO and its variants [15,35,40] have been successfully extended to multimodal models [9,18] and a range of

generative tasks, including diffusion models [37], image captioning [25,31], and visual question answering [10].

In DPO, the policy model is trained to assign higher likelihood to a preferred response  $y_w$  than to a less-preferred response  $y_l$ , relative to a fixed reference model. In our setting, the policy corresponds to the task-description generator  $\mathcal{G}$ , and the input corresponds to the multimodal user query  $q$ . The DPO objective is defined as:

$$\mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(q, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \log \frac{\pi_\theta(y_w | q)}{\pi_{\text{ref}}(y_w | q)} - \beta \log \frac{\pi_\theta(y_l | q)}{\pi_{\text{ref}}(y_l | q)} \right) \right], \quad (4)$$

where  $\beta$  controls the strength of the regularization that ties the policy  $\pi_\theta$  to the frozen reference model  $\pi_{\text{ref}}$ . In practice,  $\pi_\theta$  is initialized from  $\pi_{\text{ref}}$ .

**Preference Dataset Construction.** To apply DPO in our framework, we construct a preference dataset derived from the SFT-trained model. Given a multimodal query  $q$ , we generate multiple candidate task descriptions using the SFT model under diverse decoding strategies, including greedy decoding, beam search, and temperature-based sampling. Let  $\mathcal{S} = \{s_1, \dots, s_K\}$  denote the resulting set of candidate descriptions.

Let  $\mathcal{T} = \{t_1, t_2, \dots, t_N\}$  be the tool collection, and let  $t^* \in \mathcal{T}$  denote the ground-truth tool associated with query  $q$ . Each candidate description  $s_k$  is embedded and compared against all tool descriptions using the similarity function defined in Eq. 2. Tools are ranked in descending order of similarity, where a lower rank indicates better alignment with the ground-truth tool.

For each query, the candidate description yielding the best (lowest) rank for  $t^*$  is selected as the preferred response  $y_w$ . A non-preferred response  $y_l$  is sampled uniformly from the remaining candidates. Queries for which all candidates achieve identical ranks are discarded to avoid ambiguous preference signals. The resulting preference pairs are then used to fine-tune the model via the DPO objective defined in Eq. 4.

## 4 Experiments

### 4.1 Experimental Setting

**Dataset Details.** For the SFT stage, we build upon ToolMMBench [38] and apply the pre-processing and standardization steps described in Sec. 3.1. To generate the ground-truth descriptions in JSON format starting from the Hugging Face model cards, we use Qwen3-8B<sup>1</sup> [41]. The resulting dataset is partitioned into training and test splits. Dataset statistics are reported in Table 1.

The split is performed at the *tool level* to reflect the open-world evaluation setting. Specifically, 90% of the tools are assigned to the training split and the

<sup>1</sup> Qwen/Qwen3-8B

**Table 1.** Dataset statistics.

	Unique Queries				Unique Tools			
	Text	Image	Audio	All	Text	Image	Audio	All
Training split	6,381	3,575	625	10,581	585	170	71	826
Test split	3,231	922	190	4,343	66	20	8	94
Overall	9,612	4,497	815	14,924	651	190	79	920

remaining 10% to the test split, yielding 826 training tools and 94 test tools, respectively. Queries are assigned to a split based on the associated tool, ensuring that tools appearing in the test set are never observed during training. Particular care is taken to preserve the distribution of input modalities (text-only, text-image, and text-audio) across splits.

For DPO training, we construct a preference dataset comprising 2,981 training items and 409 evaluation items, following the procedure described in Sec. 3.3. For each query, we generate five candidate responses using diverse decoding strategies to encourage output variability. Specifically, we employ greedy decoding, beam search with a beam width of 5, and three sampling-based configurations: medium-temperature sampling (using temperature equal to 0.7), high-temperature sampling (with temperature equal to 1.0), and sampling combined with beam search using three beams. This mix of deterministic and stochastic decoding strategies yields a diverse set of candidate responses from which preference pairs are constructed.

**Training and Evaluation Protocols.** Each dataset item is defined as a pair consisting of a user query and a tool. A query may be purely textual or accompanied by an additional modality, such as an image or audio input. As a consequence, multiple items may share the same textual prompt while differing in their associated multimodal inputs, and are therefore treated as distinct queries. Moreover, the dataset naturally exhibits a many-to-many relationship: a single tool can be associated with multiple queries, and conversely, a single query may correspond to different tools depending on the intended task.

During training, this multiplicity is preserved, allowing the model to observe multiple valid tool associations for similar or identical queries. At test time, each query-tool pair is evaluated independently and exactly once. The model predicts a single tool for each query, and performance is measured by whether the retrieved tool matches the ground-truth tool associated with that query instance. Detailed statistics on the number of queries and tools in each split are provided in Table 1. Overall, the dataset consists of 14,924 unique queries.

**Implementation and Training Details.** For structured task description generation, we adopt Qwen2.5-Omni models [19] with two different scales, namely 3B and 7B parameters<sup>2</sup>, as the underlying MLLM. Both variants are fine-tuned using Low-Rank Adaptation (LoRA) [16] with rank  $r = 8$  and scaling factor  $\alpha = 16$ . Supervised fine-tuning is performed for 4 epochs on 2 NVIDIA L40S/A40

<sup>2</sup> Qwen/Qwen2.5-Omni-3B and Qwen/Qwen2.5-Omni-7B

**Table 2.** Accuracy results on aggregate metrics: unweighted average across modalities ( $\text{Avg}_q$ ) and weighted average ( $\text{Avg}_m$ ).

	Retriever	Training		Modalities			$\text{Avg}_q$	$\text{Avg}_m$
		SFT	DPO	Text	Image	Audio		
Qwen2.5-Omni-3B [19]	Contriever	-	-	17.7	8.0	36.3	16.5	20.7
Qwen2.5-Omni-7B [19]	Contriever	-	-	18.3	6.5	27.9	16.2	17.6
Qwen2.5-Omni-3B [19]	Qwen3-Embedding	-	-	24.9	31.3	69.0	28.2	41.7
Qwen2.5-Omni-7B [19]	Qwen3-Embedding	-	-	25.1	32.1	67.4	28.4	41.5
	Contriever	✓	-	57.1	37.0	80.5	53.8	58.2
	Qwen3-Embedding	✓	-	71.0	55.7	82.5	68.3	69.7
<b>RaTA-Tool-3B</b>	Qwen3-Embedding	✓	✓	<b>71.6</b>	<b>56.1</b>	<b>83.7</b>	<b>68.8</b>	<b>70.4</b>
	Contriever	✓	-	51.2	31.8	79.0	48.3	54.0
	Qwen3-Embedding	✓	-	69.3	<b>58.8</b>	<b>83.7</b>	67.7	70.6
<b>RaTA-Tool-7B</b>	Qwen3-Embedding	✓	✓	<b>71.6</b>	57.1	<b>83.7</b>	<b>69.1</b>	<b>70.8</b>

GPUs, using a learning rate of  $2 \times 10^{-4}$  and a weight decay of 0.01. We further align the models via DPO, again employing LoRA with the same configuration as in SFT. DPO training is conducted for 8 epochs on a single NVIDIA L40S/A40 GPU, with a learning rate of  $5 \times 10^{-5}$  and a weight decay of 0.01. For the retrieval stage, we use Qwen3-Embedding-8B<sup>3</sup> [44] as the text embedding model to encode both generated task descriptions and tool descriptions.

## 4.2 Experimental Results

Table 2 reports accuracy results on the test set of the introduced dataset using the standardized JSON tool-description format. Results are shown separately for text, image, and audio queries, together with two aggregate metrics: an unweighted average across modalities ( $\text{Avg}_q$ ) and a weighted average ( $\text{Avg}_m$ ), where weights reflect the number of samples per modality. This evaluation protocol allows us to assess modality-specific behavior while also capturing overall system performance. We evaluate two embedding backbones for the retrieval stage. Our main configuration relies on Qwen3-Embedding [44], while we additionally consider the Contriever model [17], based on the BERT base version [11], to assess the impact of the embedding model on tool selection.

In the zero-shot setting, performance is generally limited, particularly when using Contriever. With this retrieval model, both Qwen2.5-Omni backbones achieve low aggregate accuracy, with  $\text{Avg}_m$  below 21 for the 3B model and even lower for the 7B variant (*i.e.*, 17.6). The modality breakdown reveals uneven behavior: audio queries reach moderate accuracy, whereas text and image queries remain substantially lower. Replacing Contriever with Qwen3-Embedding yields a large improvement, increasing  $\text{Avg}_m$  to 41.7 (3B) and 41.5 (7B). The gains are consistent across modalities and especially pronounced for image and audio, indicating that retrieval quality and modality coverage of the embedding space are critical even without task-description fine-tuning.

<sup>3</sup> Qwen/Qwen3-Embedding-8B

**Table 3.** Ablation study on the effect of tool-description format and prompting or fine-tuning strategy on tool-selection performance, reported for both retrieval backbones (Contriever and Qwen3-Embedding). Results are shown for text, image, and audio queries, together with unweighted ( $\text{Avg}_q$ ) and weighted ( $\text{Avg}_m$ ) aggregate metrics.

Retriever	Descriptions		Modalities			$\text{Avg}_q$	$\text{Avg}_m$
	Type	Mode	Text	Image	Audio		
Contriever	NL	Zero-shot	16.8	17.2	43.7	18.1	25.9
Contriever	NL	Few-shot	13.3	13.7	49.5	15.0	25.5
Contriever	JSON	Zero-shot	18.3	6.5	27.9	16.2	17.6
Contriever	JSON	Few-shot	21.0	15.1	35.8	20.4	24.0
Contriever	NL	SFT	34.3	<b>50.3</b>	79.0	39.6	<b>54.5</b>
<b>RaTA-Tool-7B</b>	Contriever	JSON	<b>51.2</b>	31.8	<b>79.0</b>	<b>48.3</b>	54.0
Qwen3-Emb	NL	Zero-shot	27.6	32.0	72.1	30.5	43.9
Qwen3-Emb	NL	Few-shot	31.5	22.1	72.1	31.3	41.9
Qwen3-Emb	JSON	Zero-shot	25.1	32.1	67.4	28.4	41.5
Qwen3-Emb	JSON	Few-shot	25.9	38.7	70.0	30.6	44.9
Qwen3-Emb	NL	SFT	41.6	44.0	<b>95.8</b>	44.5	60.5
<b>RaTA-Tool-7B</b>	Qwen3-Emb	JSON	<b>69.3</b>	<b>58.8</b>	83.7	<b>67.7</b>	<b>70.6</b>

SFT substantially improves performance for both model sizes, confirming the importance of aligning the generator to produce structured task descriptions in the standardized JSON format. Again, the retriever choice remains critical: for the 3B setting, Qwen3-Embedding improves  $\text{Avg}_m$  from 58.2 (Contriever) to 69.7; for 7B, it increases  $\text{Avg}_m$  from 54.0 to 70.6. These differences are reflected across modalities, with particularly gains for text and image queries.

Finally, applying DPO yields further improvements. For RaTA-Tool-3B, DPO increases performance to  $\text{Avg}_q = 68.8$  and  $\text{Avg}_m = 70.4$ , reaching 71.6 (text), 56.1 (image), and 83.7 (audio). For RaTA-Tool-7B, DPO produces the best overall results with  $\text{Avg}_q = 69.1$  and  $\text{Avg}_m = 70.8$ , and consistently strong accuracy across modalities. Overall, these results show that (i) a strong embedding model is essential for retrieval-based tool selection, and (ii) preference-based alignment provides complementary gains beyond supervised fine-tuning.

### 4.3 Ablation and Analysis

Table 3 presents an ablation study using Qwen2.5-Omni-7B as the task-description generator, focusing on the impact of the tool-description format and the prompting or fine-tuning strategy. Results are reported for both retrieval backbones previously considered, namely Contriever and Qwen3-Embedding.

We first analyze the effect of the description format under zero-shot and few-shot prompting. Across both retrievers, structured JSON descriptions outperform natural-language (NL) ones. For example, with Qwen3-Embedding in the zero-shot setting, JSON achieves an  $\text{Avg}_m$  of 41.5, compared to 43.9 for NL, while under few-shot prompting JSON improves to 44.9, narrowing and reversing

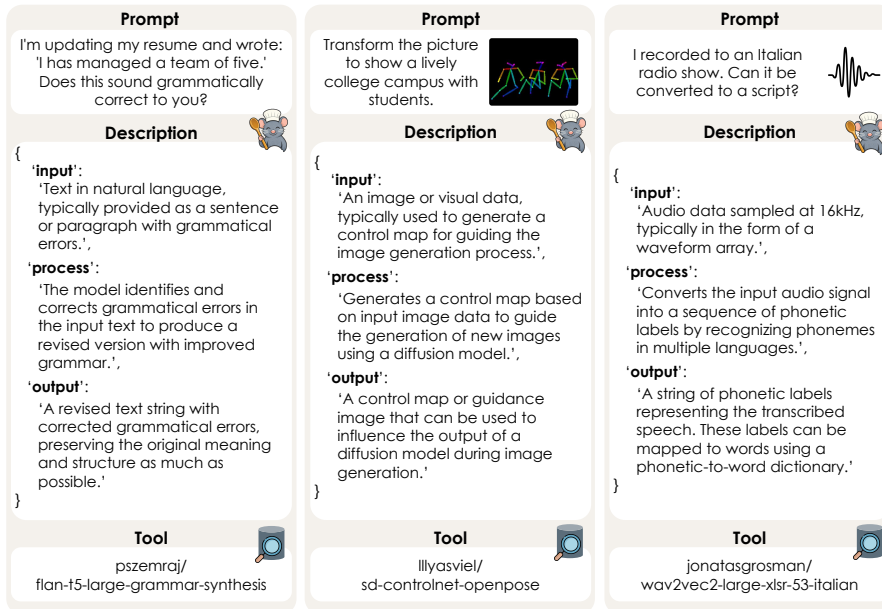


Fig. 4. Qualitative results of RaTA-Tool under different input modalities.

the gap relative to NL (41.9). A similar trend is observed with Contriever, where JSON under few-shot prompting ( $Avg_m = 24.0$ ) outperforms NL ( $Avg_m = 25.5$ ).

Comparing prompting strategies, few-shot prompting consistently improves over zero-shot inference for both NL and JSON formats, with especially notable gains for image and audio queries. For instance, when using Qwen3-Embedding with JSON descriptions, image accuracy increases from 32.1 to 38.7 when moving from zero-shot to few-shot prompting. These improvements suggest that in-context examples help the model better capture task-relevant constraints and generate more discriminative task descriptions for retrieval. SFT leads to the largest performance improvements across both retrieval backbones. With Contriever, SFT with JSON descriptions improves  $Avg_m$  from 24.0 (few-shot) to 54.0, while with Qwen3-Embedding it increases from 44.9 to 70.6. In both cases, JSON descriptions outperform NL under SFT, confirming the advantage of standardized, machine-readable formats.

Overall, these results indicate that while retrieval quality remains important, the choice of description format and training strategy plays a central role in effective multimodal tool selection. The best results are achieved by combining SFT with structured JSON descriptions, consistently across both retrieval backbones.

#### 4.4 Qualitative Results

We present qualitative results of RaTA-Tool in Fig. 4. The examples include purely textual prompts (first column) as well as multimodal prompts (second

and third columns). In all cases, the model generates a structured JSON description of the user request, which is then used to retrieve the most relevant tool from the external collection. The generated descriptions accurately capture the expected inputs, functional intent, and outputs of the task, demonstrating effective understanding of both textual and multimodal cues. Across all shown examples, RaTA-Tool retrieves the correct tool. Additional qualitative examples, together with the exact prompts used in our experiments, are provided in the supplementary material.

## 5 Conclusion

We presented RaTA-Tool, a retrieval-based framework for open-world multimodal tool selection that reasons over structured, machine-readable tool descriptions instead of relying on closed-set supervision. By converting multimodal user queries into structured task descriptions and retrieving tools via semantic similarity, our approach naturally generalizes to unseen tools. We further showed that preference-based alignment with DPO complements standard SFT, improving robustness and consistency across modalities. To support evaluation in this setting, we introduced the first benchmark for open-world multimodal tool use with standardized tool descriptions. Our results demonstrate that semantic, description-driven tool reasoning is a promising direction for scalable and robust tool-learning systems.

## References

1. Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F.L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al.: GPT-4 Technical Report. arXiv preprint arXiv:2303.08774 (2023)
2. Alayrac, J.B., Donahue, J., Luc, P., Miech, A., Barr, I., Hasson, Y., Lenc, K., Mensch, A., Millican, K., Reynolds, M., et al.: Flamingo: a Visual Language Model for Few-Shot Learning. In: NeurIPS (2022)
3. Awais, M., Naseer, M., Khan, S., Anwer, R.M., Cholakkal, H., Shah, M., Yang, M.H., Khan, F.S.: Foundation Models Defining a New Era in Vision: A Survey and Outlook. IEEE Trans. PAMI (2025)
4. Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al.: Language Models are Few-Shot Learners. In: NeurIPS (2020)
5. Caffagni, D., Cocchi, F., Barsellotti, L., Moratelli, N., Sarto, S., Baraldi, L., Cornia, M., Cucchiara, R.: The Revolution of Multimodal Large Language Models: A Survey. In: ACL Findings (2024)
6. Chen, F., Han, M., Zhao, H., Zhang, Q., Shi, J., Xu, S., , Xu, B.: X-LLM: Bootstrapping Advanced Large Language Models by Treating Multi-Modalities as Foreign Languages. arXiv preprint arXiv:2305.04160 (2023)
7. Cheng, Z., Leng, S., Zhang, H., Xin, Y., Li, X., Chen, G., Zhu, Y., Zhang, W., Luo, Z., Zhao, D., et al.: VideoLLaMA 2: Advancing Spatial-Temporal Modeling and Audio Understanding in Video-LLMs. arXiv preprint arXiv:2406.07476 (2024)

8. Cocchi, F., Moratelli, N., Caffagni, D., Sarto, S., Baraldi, L., Cornia, M., Cucchiara, R.: LLaVA-MORE: A Comparative Study of LLMs and Visual Backbones for Enhanced Visual Instruction Tuning. In: ICCV Workshops (2025)
9. Compagnoni, A., Caffagni, D., Moratelli, N., Baraldi, L., Cornia, M., Cucchiara, R.: Mitigating Hallucinations in Multimodal LLMs via Object-aware Preference Optimization. In: BMVC (2025)
10. Compagnoni, A., Morini, M., Sarto, S., Cocchi, F., Caffagni, D., Cornia, M., Baraldi, L., Cucchiara, R.: ReAG: Reasoning-Augmented Generation for Knowledge-based Visual Question Answering. arXiv preprint arXiv:2511.22715 (2025)
11. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In: NAACL (2019)
12. Fu, C., Lin, H., Long, Z., Shen, Y., Dai, Y., Zhao, M., Zhang, Y.F., Dong, S., Li, Y., Wang, X., et al.: VITA: Towards Open-Source Interactive Omni Multimodal LLM. arXiv preprint arXiv:2408.05211 (2024)
13. Grattafiori, A., Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Vaughan, A., et al.: The Llama 3 Herd of Models. arXiv preprint arXiv:2407.21783 (2024)
14. Han, J., Gong, K., Zhang, Y., Wang, J., Zhang, K., Lin, D., Qiao, Y., Gao, P., Yue, X.: OneLLM: One Framework to Align All Modalities with Language. In: CVPR (2024)
15. Hong, J., Lee, N., Thorne, J.: ORPO: Monolithic Preference Optimization without Reference Model. In: EMNLP (2024)
16. Hu, E.J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W., et al.: LoRA: Low-Rank Adaptation of Large Language Models. In: ICLR (2022)
17. Izacard, G., Caron, M., Hosseini, L., Riedel, S., Bojanowski, P., Joulin, A., Grave, E.: Unsupervised Dense Information Retrieval with Contrastive Learning. TMLR (2021)
18. Jia, H., Jiang, C., Xu, H., Ye, W., Dong, M., Yan, M., Zhang, J., Huang, F., Zhang, S.: SymDPO: Boosting In-Context Learning of Large Multimodal Models with Symbol Demonstration Direct Preference Optimization. In: CVPR (2025)
19. Jin Xu, Zhifang Guo, Jinzheng He, Hangrui Hu, Ting He, Shuai Bai, Keqin Chen, Jialin Wang, Yang Fan, Kai Dang, Bin Zhang, Xiong Wang, Yunfei Chu, Junyang Lin: Qwen2.5-Omni Technical Report. arXiv preprint arXiv:2503.20215 (2025)
20. Lambert, N.: Reinforcement Learning from Human Feedback. arXiv preprint arXiv:2504.12501 (2025)
21. Li, J., Li, D., Savarese, S., Hoi, S.: BLIP-2: Bootstrapping Language-Image Pre-training with Frozen Image Encoders and Large Language Models. In: ICML (2023)
22. Liu, H., Li, C., Li, Y., Lee, Y.J.: Improved Baselines with Visual Instruction Tuning. In: CVPR (2024)
23. Liu, H., Li, C., Wu, Q., Lee, Y.J.: Visual Instruction Tuning. In: NeurIPS (2023)
24. Minaee, S., Mikolov, T., Nikzad, N., Chenaghlu, M., Socher, R., Amatriain, X., Gao, J.: Large Language Models: A Survey. arXiv preprint arXiv:2402.06196 (2024)
25. Moratelli, N., Caffagni, D., Cornia, M., Baraldi, L., Cucchiara, R.: Revisiting Image Captioning Training Paradigm via Direct CLIP-based Optimization. In: BMVC (2024)
26. Nakano, R., Hilton, J., Balaji, S., Wu, J., Ouyang, L., Kim, C., Hesse, C., Jain, S., Kosaraju, V., Saunders, W., et al.: WebGPT: Browser-assisted question-answering with human feedback. arXiv preprint arXiv:2112.09332 (2021)
27. Patil, S.G., Zhang, T., Wang, X., Gonzalez, J.E.: Gorilla: Large Language Model Connected with Massive APIs. In: NeurIPS (2024)

28. Qin, Y., Hu, S., Lin, Y., Chen, W., Ding, N., Cui, G., Zeng, Z., Zhou, X., Huang, Y., Xiao, C., et al.: Tool Learning with Foundation Models. *ACM Computing Surveys* **57**(4), 1–40 (2024)
29. Qin, Y., Liang, S., Ye, Y., Zhu, K., Yan, L., Lu, Y., Lin, Y., Cong, X., Tang, X., Qian, B., et al.: ToolLLM: Facilitating Large Language Models to Master 16000+ Real-world APIs. In: *ICLR* (2024)
30. Rafailov, R., Sharma, A., Mitchell, E., Manning, C.D., Ermon, S., Finn, C.: Direct Preference Optimization: Your Language Model is Secretly a Reward Model. In: *NeurIPS* (2023)
31. Sarto, S., Moratelli, N., Cornia, M., Baraldi, L., Cucchiara, R.: Positive-Augmented Contrastive Learning for Vision-and-Language Evaluation and Trainin. *IJCV* **133**(11), 7647–7671 (2025)
32. Schick, T., Dwivedi-Yu, J., Dessì, R., Raileanu, R., Lomeli, M., Hambro, E., Zettlemoyer, L., Cancedda, N., Scialom, T.: Toolformer: Language Models Can Teach Themselves to Use Tools. In: *NeurIPS* (2023)
33. Shen, Y., Song, K., Tan, X., Li, D., Lu, W., Zhuang, Y.: HuggingGPT: Solving AI Tasks with ChatGPT and its Friends in Hugging Face. In: *NeurIPS* (2023)
34. Shuster, K., Xu, J., Komeili, M., Ju, D., Smith, E.M., Roller, S., Ung, M., Chen, M., Arora, K., Lane, J., et al.: BlenderBot 3: a deployed conversational agent that continually learns to responsibly engage. *arXiv preprint arXiv:2208.03188* (2022)
35. Song, F., Yu, B., Li, M., Yu, H., Huang, F., Li, Y., Wang, H.: Preference Ranking Optimization for Human Alignment. In: *AAAI* (2024)
36. Thoppilan, R., De Freitas, D., Hall, J., Shazeer, N., Kulshreshtha, A., Cheng, H.T., Jin, A., Bos, T., Baker, L., Du, Y., et al.: LaMDA: Language Models for Dialog Applications. *arXiv preprint arXiv:2201.08239* (2022)
37. Wallace, B., Dang, M., Rafailov, R., Zhou, L., Lou, A., Purushwalkam, S., Ermon, S., Xiong, C., Joty, S., Naik, N.: Diffusion Model Alignment Using Direct Preference Optimization. In: *CVPR* (2024)
38. Wang, C., Luo, W., Dong, S., Xuan, X., Li, Z., Ma, L., Gao, S.: MLLM-Tool: A Multimodal Large Language Model For Tool Agent Learning. In: *WACV* (2025)
39. Wu, C., Yin, S., Qi, W., Wang, X., Tang, Z., Duan, N.: Visual ChatGPT: Talking, Drawing and Editing with Visual Foundation Models. *arXiv preprint arXiv:2303.04671* (2023)
40. Wu, J., Xie, Y., Yang, Z., Wu, J., Gao, J., Ding, B., Wang, X., He, X.:  $\beta$ -DPO: Direct Preference Optimization with Dynamic  $\beta$ . In: *NeurIPS* (2024)
41. Yang, A., Li, A., Yang, B., Zhang, B., Hui, B., Zheng, B., Yu, B., Gao, C., Huang, C., Lv, C., et al.: Qwen3 Technical Report. *arXiv preprint arXiv:2505.09388* (2025)
42. Yang, R., Song, L., Li, Y., Zhao, S., Ge, Y., Li, X., Shan, Y.: GPT4Tools: Teaching Large Language Model to Use Tools via Self-instruction. In: *NeurIPS* (2023)
43. Zhang, H., Li, X., Bing, L.: Video-LLaMA: An Instruction-tuned Audio-Visual Language Model for Video Understanding. In: *EMNLP Demos* (2023)
44. Zhang, Y., Li, M., Long, D., Zhang, X., Lin, H., Yang, B., Xie, P., Yang, A., Liu, D., Lin, J., Huang, F., Zhou, J.: Qwen3 Embedding: Advancing Text Embedding and Reranking Through Foundation Models. *arXiv preprint arXiv:2506.05176* (2025)
45. Zhao, Z., Guo, L., Yue, T., Chen, S., Shao, S., Zhu, X., Yuan, Z., Liu, J.: Chat-Bridge: Bridging Modalities with Large Language Model as a Language Catalyst. *arXiv preprint arXiv:2305.16103* (2023)

# Bibliography

- [1] A. Vaswani, N. Shazeer, N. Parmar, *et al.*, «Attention is all you need», *Advances in neural information processing systems*, vol. 30, 2017.
- [2] A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever, «Improving language understanding by generative pre-training», *OpenAI Technical Report*, 2018.
- [3] A. Dosovitskiy, L. Beyer, A. Kolesnikov, *et al.*, «An image is worth 16x16 words: Transformers for image recognition at scale», *arXiv preprint arXiv:2010.11929*, 2020.
- [4] G. Bertasius, H. Wang, and L. Torresani, «Is space-time attention all you need for video understanding?», in *Icml*, vol. 2, 2021, p. 4.
- [5] A. Arnab, M. Dehghani, G. Heigold, C. Sun, M. Lučić, and C. Schmid, «Vivit: A video vision transformer», in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 6836–6846.
- [6] Y. Gong, Y.-A. Chung, and J. Glass, «Ast: Audio spectrogram transformer», *arXiv preprint arXiv:2104.01778*, 2021.
- [7] A. Ng, *How agents can improve LLM performance*, <https://www.deeplearning.ai/the-batch/how-agents-can-improve-llm-performance/>, The Batch (DeepLearning.AI). Accesso: 22 Maggio 2024, 2024.
- [8] Wikipedia, *Neural network (machine learning)* — *Wikipedia, the free encyclopedia*, [Online; accessed 26-February-2026], 2026. [Online]. Available: [https://en.wikipedia.org/w/index.php?title=Neural\\_network\\_\(machine\\_learning\)&oldid=1339866777](https://en.wikipedia.org/w/index.php?title=Neural_network_(machine_learning)&oldid=1339866777).

- [9] A. Singh, A. Fry, A. Perelman, *et al.*, «Openai gpt-5 system card», *arXiv preprint arXiv:2601.03267*, 2025.
- [10] A. Grattafiori, A. Dubey, A. Jauhri, *et al.*, «The llama 3 herd of models», *arXiv preprint arXiv:2407.21783*, 2024.
- [11] A. Yang, A. Li, B. Yang, *et al.*, «Qwen3 technical report», *arXiv preprint arXiv:2505.09388*, 2025.
- [12] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, «Bert: Pre-training of deep bidirectional transformers for language understanding», in *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, 2019, pp. 4171–4186.
- [13] C. Raffel, N. Shazeer, A. Roberts, *et al.*, «Exploring the limits of transfer learning with a unified text-to-text transformer», *Journal of machine learning research*, vol. 21, no. 140, pp. 1–67, 2020.
- [14] L. Ouyang, J. Wu, X. Jiang, *et al.*, «Training language models to follow instructions with human feedback», *Advances in neural information processing systems*, vol. 35, pp. 27 730–27 744, 2022.
- [15] A. Radford, J. W. Kim, C. Hallacy, *et al.*, «Learning transferable visual models from natural language supervision», in *International conference on machine learning*, PmLR, 2021, pp. 8748–8763.
- [16] K. He, X. Zhang, S. Ren, and J. Sun, «Deep residual learning for image recognition», in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [17] Y. Luo, Z. Yang, F. Meng, Y. Li, J. Zhou, and Y. Zhang, «An empirical study of catastrophic forgetting in large language models during continual fine-tuning», *IEEE Transactions on Audio, Speech and Language Processing*, 2025.
- [18] E. J. Hu, Y. Shen, P. Wallis, *et al.*, «Lora: Low-rank adaptation of large language models.», *Iclr*, vol. 1, no. 2, p. 3, 2022.

- 
- [19] R. Rafailov, A. Sharma, E. Mitchell, C. D. Manning, S. Ermon, and C. Finn, «Direct preference optimization: Your language model is secretly a reward model», *Advances in neural information processing systems*, vol. 36, pp. 53 728–53 741, 2023.
- [20] Y. Qin, S. Hu, Y. Lin, *et al.*, «Tool learning with foundation models», *ACM Computing Surveys*, vol. 57, no. 4, pp. 1–40, 2024.
- [21] C. Wang, W. Luo, S. Dong, *et al.*, «Mllm-tool: A multimodal large language model for tool agent learning», in *2025 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, IEEE, 2025, pp. 6678–6687.
- [22] R. Girdhar, A. El-Nouby, Z. Liu, *et al.*, «Imagebind: One embedding space to bind them all», in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2023, pp. 15 180–15 190.
- [23] W.-L. Chiang, Z. Li, Z. Lin, *et al.*, *Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality*, Mar. 2023. [Online]. Available: <https://lmsys.org/blog/2023-03-30-vicuna/>.
- [24] H. Touvron, T. Lavril, G. Izacard, *et al.*, «Llama: Open and efficient foundation language models», *arXiv preprint arXiv:2302.13971*, 2023.
- [25] T. Schick, J. Dwivedi-Yu, R. Dessi, *et al.*, «Toolformer: Language models can teach themselves to use tools», *Advances in neural information processing systems*, vol. 36, pp. 68 539–68 551, 2023.
- [26] S. G. Patil, T. Zhang, X. Wang, and J. E. Gonzalez, «Gorilla: Large language model connected with massive apis», *Advances in Neural Information Processing Systems*, vol. 37, pp. 126 544–126 565, 2024.
- [27] Y. Qin, S. Liang, Y. Ye, *et al.*, «Toollm: Facilitating large language models to master 16000+ real-world apis», *arXiv preprint arXiv:2307.16789*, 2023.
- [28] Y. Shen, K. Song, X. Tan, D. Li, W. Lu, and Y. Zhuang, «Hugginggpt: Solving ai tasks with chatgpt and its friends in hugging face», *Advances in Neural Information Processing Systems*, vol. 36, pp. 38 154–38 180, 2023.

- [29] C. Wu, S. Yin, W. Qi, X. Wang, Z. Tang, and N. Duan, «Visual chatgpt: Talking, drawing and editing with visual foundation models», *arXiv preprint arXiv:2303.04671*, 2023.
- [30] R. Yang, L. Song, Y. Li, *et al.*, «Gpt4tools: Teaching large language model to use tools via self-instruction», *Advances in Neural Information Processing Systems*, vol. 36, pp. 71 995–72 007, 2023.
- [31] Jin Xu, Zhifang Guo, Jinzheng He, Hangrui Hu, Ting He, Shuai Bai, Keqin Chen, Jialin Wang, Yang Fan, Kai Dang, Bin Zhang, Xiong Wang, Yunfei Chu, Junyang Lin, «Qwen2.5-Omni Technical Report», *arXiv preprint arXiv:2503.20215*, 2025.
- [32] Y. Zhang, M. Li, D. Long, *et al.*, «Qwen3 embedding: Advancing text embedding and reranking through foundation models», *arXiv preprint arXiv:2506.05176*, 2025.
- [33] D. Kalamkar, D. Mudigere, N. Mellempudi, *et al.*, «A study of bfloat16 for deep learning training», *arXiv preprint arXiv:1905.12322*, 2019.
- [34] G. Izacard, M. Caron, L. Hosseini, *et al.*, «Unsupervised dense information retrieval with contrastive learning», *arXiv preprint arXiv:2112.09118*, 2021.