# POLITECNICO DI TORINO

Master degree course in Computer Engineering

## Master Degree Thesis

# A Multi-Agent AI Assistant for Intelligent Research and Neuromorphic Application Development

**Advisors**
Prof. Gianvito Urgese
Dr. Vittorio Fra
Salvatore Tilocca

**Candidate**
Antonio Sirica

December 2025

*"Buzz, stai volando!"*

*"Questo non è volare, è cadere con stile."*

# Abstract

Traditional computing architectures based on the von Neumann model face inefficiencies when processing massively parallel and event-driven Artificial Intelligence (AI) workloads, suffering from memory–computation bottlenecks and high power consumption. Neuromorphic computing, inspired by biological neural systems, addresses these challenges through asynchronous, event-driven processing with low latency and high energy efficiency. Recent advances in neuromorphic hardware have further promoted algorithm–hardware co-design to improve adaptability and scalability in real-time and edge computing.

However, these benefits are often constrained by the lack of accessible development tools, standardized methodologies, and comprehensive documentation. Existing implementations frequently derive from research prototypes tailored to specific experiments rather than reusable, structured libraries, making new developments complex and often dependent on expert intervention.

Large Language Models (LLMs) are increasingly employed to simplify program synthesis and vibe coding in conventional AI workflows. Yet, their potential in supporting neuromorphic system design remains unexplored. This thesis aims to address the gap by extending emerging AI-driven development assistance to neuromorphic applications.

The proposed solution integrates into key stages of the MLOps lifecycle, supporting code synthesis, model design, and optimization through the use of *LangGraph*, a state-of-art graph-based multi-agent framework.

The approach relies on three branches: web search, academic search, and code generation & validation. Each acts as a node in a unified LangGraph pipeline, enabling contextual information retrieval, research knowledge extraction, and automated code generation with iterative self-correction.

The developed method was evaluated experimentally across all branches. It achieved high scores on standard large language model metrics in both web and academic search tasks, showing strong factual accuracy and completeness. The results indicate a close alignment between generated content and reference material, with performance generally within the 80–90% range. The evaluation followed the LLM-as-a-Judge paradigm, employing GPT-5 to assess reliability, clarity, and relevance.

The core component of code synthesis is organized around a central orchestrator coordinating specialized agents. Each agent includes vector stores for domain knowledge of `snnTorch` for spiking network simulation and the `Neural Network Intelligence (NNI)` toolkit for automated optimization.

Code validation follows the four main dimensions adopted in the literature. i) Functional correctness is tested through automated execution in isolated cloud environments. ii) Static code quality is verified with static type checking and module consistency analysis. iii) Runtime performance metrics are collected to evaluate

efficiency. iv) Feedback-based evaluation incorporates expert input during synthesis, enabling iterative refinement. These validation layers ensure the correctness and robustness of the synthesis process. Reference-based normalization allows fair performance comparisons across experiments.

Overall, the results confirm the system's reliability and ability to generate accurate, high-quality code across heterogeneous sources, demonstrating the feasibility of agent-based systems to support neuromorphic applications development.

# Contents

# List of Figures

# List of Tables

# Chapter 1

# Introduction

Neuromorphic computing represents a paradigm shift in how we approach artificial intelligence workloads and real-time processing systems. Unlike traditional computing architectures rooted in the von Neumann model, neuromorphic systems draw inspiration from biological neural structures to achieve fundamentally different computational properties. These systems address critical inefficiencies in conventional approaches through asynchronous, event-driven processing mechanisms that naturally align with the sparse, temporal nature of real-world data and sensory inputs.

The appeal of neuromorphic computing extends beyond theoretical elegance. Contemporary challenges in AI deployment, particularly the memory, computation bottleneck and escalating power consumption in data centers and edge devices, have created renewed urgency for alternative computational substrates. Recent advances in neuromorphic hardware, from Intel's Loihi processors to emerging neuromorphic chips, have demonstrated practical viability and prompted significant research into algorithm–hardware co-design methodologies. These developments suggest that neuromorphic systems could become mainstream components of modern computing infrastructure, particularly for latency-sensitive and energy-constrained applications.

Yet despite these promising advances, neuromorphic computing has not achieved the widespread adoption enjoyed by conventional deep learning frameworks. A principal obstacle lies in the accessibility and maturity of the development ecosystem. Most existing neuromorphic implementations remain research prototypes, tightly coupled to specific experimental scenarios and rarely designed for reusability or modularity. Developers seeking to build neuromorphic applications face fragmented documentation, limited standardized methodologies, and a steep learning curve that often requires deep domain expertise or consultation with research specialists.

This friction between the potential of neuromorphic hardware and the practical barriers to development creates a significant gap. While conventional AI has benefited from sophisticated development tools, from TensorFlow and PyTorch to

specialized code generation assistants, the neuromorphic domain lacks comparable infrastructure. The absence of accessible, structured libraries and development assistance tools places neuromorphic computing at a disadvantage when competing for developer attention and investment.

Recent breakthroughs in LLMs have demonstrated remarkable capabilities in code synthesis, program generation, and AI-assisted development workflows. These models can understand complex specifications, generate executable code, debug implementations, and provide contextual assistance across diverse programming domains. Their potential has been extensively validated in conventional AI workflows, where they accelerate development cycles and reduce the barrier to entry for practitioners.

However, this transformative capability has remained separate from neuromorphic computing. The existing literature on LLM-assisted development focuses almost exclusively on conventional neural network frameworks, cloud infrastructure, and standard software engineering tasks. The possibility of extending these AI-driven development tools to neuromorphic applications, where domain-specific knowledge, specialized libraries, and novel architectural patterns are required, remains unexplored.

This thesis addresses this gap by investigating whether LLMs, when properly augmented with domain knowledge and structured reasoning capabilities, can effectively support the design, implementation, and optimization of neuromorphic systems. Rather than treating LLMs as isolated code generation tools, in this thesis is explored a multi-agent architecture capable of orchestrating diverse specialized roles: searching and synthesizing research knowledge, validating implementations, and generating code with iterative refinement.

The proposed approach integrates into the MLOps lifecycle through three complementary branches. The web search and academic search branches provide contextual information retrieval and research knowledge extraction. The code generation and validation branch coordinates specialized agents that synthesize executable implementations while maintaining correctness across multiple validation dimensions.

This research demonstrates that structured, agent-based systems can bridge the accessibility gap in neuromorphic computing development. By combining the general reasoning capabilities of LLMs with domain-specific knowledge bases, automated validation frameworks, and multi-agent orchestration, a foundation for more inclusive and efficient neuromorphic system development practices is established.

The primary contributions of this work are threefold. First, a comprehensive multi-agent architecture for neuromorphic development assistance is presented, demonstrating how LLMs can be effectively integrated into domain-specific development workflows. Second, specialized agents focused on spiking neural network simulation (using `snnTorch`) and automated optimization (using the `Neural Network Intelligence (NNI)` toolkit) are developed and evaluated, showing how

domain-specific vector stores and knowledge bases enhance code generation quality. Third, rigorous evaluation frameworks tailored to distinct system branches is established: for knowledge-based productions in the web and academic search branches, *DeepEval* is employ, a state-of-the-art LLM-as-Judge framework, which achieved results in the 80–90% range across key metrics including Faithfulness, Answer Relevancy, Contextual Relevancy, and Hallucination Detection. For the code generation branch, qualitative validation spanning functional correctness, static code quality, runtime performance, and expert-driven feedback implementation provides a replicable methodology for assessing LLM-assisted development tools in specialized domains.

# Chapter 2

# Background

The most challenging part of this work has been dealing with the orchestration of heterogeneous processes involving large language models, information retrieval, and generation, execution and evaluation. These systems are considered very promising from multiple perspectives, ranging from their adaptability in handling diverse user inputs to their potential for integration into intelligent research and coding workflows [1, 2]. At the same time, however, significant effort is required to align different components, such as web research tools, academic search engines, vector databases, and sandboxed execution environments, into a unified and reliable framework [1, 2].

While the integration of heterogeneous tasks such as retrieval-augmented generation, code assistance, and workflow management remains an open challenge, recent studies have addressed specific components of this broader problem. Some works have focused on optimizing the efficiency and scalability of retrieval-augmented generation (RAG) pipelines [3], while others have investigated state management and workflow orchestration mechanisms for large-scale Artificial Intelligence (AI) applications [1, 2]. Parallel research efforts have explored the validation and improvement of LLM-generated code through techniques such as sandbox execution, static analysis, and reflection [4, 3]. Finally, human-in-the-loop approaches have been proposed to enable iterative refinement and alignment of system outputs with user expectations [5].

In this background section, the goal is to review some of the most relevant studies and frameworks that have been proposed in these areas, while also presenting an overview of the state-of-the-art tools and methodologies that currently support intelligent research assistance and automated code generation and validation [1, 2, 3, 4, 5].

## 2.1 Artificial Neural Networks

The concept of AI emerged in the mid-twentieth century as scientists aimed to design computational systems capable of reasoning in ways comparable to human cognition. The term was formally introduced by John McCarthy during the 1956 Dartmouth Conference, marking the official establishment of AI as a scientific research field [6]. Early AI systems relied on symbolic reasoning, encoding expert knowledge in fixed logical rules; however, such methodologies proved inadequate when confronted with complex and heterogeneous data environments.

In the 1980s, the paradigm of Machine Learning (ML) was proposed to address these limitations. This approach shifted focus from explicitly programmed rule-based models to data-driven methods capable of improving performance through experience [7]. The pioneering work of Arthur Samuel demonstrated that computers could autonomously learn patterns from data rather than relying solely on predefined logic, laying the foundation for future development.

Around 2010, the field underwent a significant transformation with the emergence of Deep Learning (DL). This new paradigm utilized hierarchical, multi-layered neural networks capable of automatically extracting complex features from raw data, fundamentally changing fields like computer vision and natural language processing [8]. The extraordinary success of deep learning has been supported by advances in computational power, large-scale labeled datasets, and optimized training algorithms.

### 2.1.1 Neural Network and Neuron Perceptron

The foundational objective of AI was to develop computational models that emulate the behavior of biological neurons. The perceptron, as one of the earliest such models, produces an output of $+1$ or $-1$ based on a weighted sum of its inputs, functioning as a linear classifier defined by [9]

$$h(\mathbf{x}) = \text{sign}(\mathbf{w}^\top \mathbf{x}). \tag{2.1}$$

The perceptron algorithm seeks a weight vector $\mathbf{w}$ that defines a hyperplane separating data points from different classes, maximizing the margin between them. The figure 2.1 shows how perceptron works. This approach assumes linear separability; otherwise, the algorithm may not converge. To address this, multiple perceptrons can be combined to form more complex, nonlinear functions.

Replacing the step function with a sigmoid activation leads to logistic regression, which models the probability of class membership and enables gradient-based optimization. By stacking layers of perceptrons, one obtains a multilayer perceptron or feed-forward neural network, which can approximate nonlinear functions and learn hierarchical representations [10]. The design of such networks aims to

balance expressive power with generalization, mitigating overfitting by controlling model complexity and regularization [11].



Figure 2.1: A *perceptron* is a foundational artificial neuron for binary classification. Its architecture comprises four stages: (1) **Input Layer** with bias constant (1) and $n$ input features $(x_1, x_2, \ldots, x_n)$, (2) **Weight Multiplication** where each input is multiplied by corresponding weights $(w_0, w_1, \ldots, w_n)$, (3) **Weighted Summation** aggregating all weighted inputs, and (4) **Activation Function** applying a step function threshold to produce binary output. This architecture demonstrates how linear combinations of inputs are transformed into discrete decisions through non-linear thresholding [12].

## 2.1.2  Training Neural Network

Training a neural network involves determining the set of weights $\boldsymbol{\theta}$ that minimize a loss function, which quantifies the discrepancy between predicted and actual values. For regression tasks, a common choice is the mean squared error:

$$\ell_\theta = \frac{1}{n} \sum_{i=1}^{n} (y_i - g_\theta(x_i))^2 \tag{2.2}$$

where $y_i$ is the true value, $g_\theta(x_i)$ is the network's prediction, and $n$ is the number of samples [9].

Since the loss function is generally non-convex for neural networks, it is typically minimized using the Gradient Descent algorithm. This iterative method updates the weights by moving in the direction of the negative gradient of the loss:

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \nabla \ell_{\theta^{(t)}} \tag{2.3}$$

19

Here, $\alpha > 0$ is the learning rate, a hyperparameter that controls the step size at each iteration. Choosing an appropriate learning rate is crucial: a value too large may cause the algorithm to overshoot minima, while a value too small can result in slow convergence.

For large datasets, computing the gradient over the entire dataset at each step is computationally expensive. To address this, Stochastic Gradient Descent (SGD) is used, where the gradient is estimated using a small batch of samples. After all batches are processed, one epoch is completed.

In feedforward neural networks, each neuron composes multiple functions, making gradient computation complex. To automate this, the concept of a computational graph is introduced: an acyclic graph representing the composition of functions computed by the network. This structure enables efficient calculation of derivatives via backpropagation, where gradients are computed from the output layer backward to the input, leveraging the chain rule, as shown in Figure 2.2. This approach is efficient because outputs are typically scalars (such as the loss), while inputs are often high-dimensional [11].

### 2.1.3  Convolutional Neural Networks

The integration of prior knowledge into convolutional filters allows neural networks to leverage known properties of image data. In images, important information is often concentrated within local pixel neighborhoods, and similar visual patterns, such as edges or textures, can appear at different positions within the frame. Convolutional layers explicitly encode these assumptions in their architecture [14]. Each neuron connects only to a small, localized region of the input, reflecting the idea that nearby pixels are correlated. Moreover, the same set of weights (the convolutional kernel) is shared across all spatial locations, enforcing translation invariance and enabling the detection of identical features regardless of their position. This design reduces the number of trainable parameters, enhances generalization, and aligns with principles observed in biological visual systems.

Mathematically, the convolution operation between two functions $f(t)$ and $g(t)$ is defined as:

$$(f * g)(t) = \int_{-\infty}^{+\infty} f(\tau)g(t - \tau)d\tau \tag{2.4}$$

where $(f * g)(t)$ is the result of convolving $f$ with $g$, and $g$ is commonly called the kernel or filter [15]. In image processing, the kernel is a small matrix of weights that slides over the input tensor, computing local weighted sums to produce a feature map.

After the convolution operation, which the Figure 2.3 describes, a non-linear activation function such as the Rectified Linear Unit (ReLU) is typically applied to introduce non-linearity. To further reduce the dimensionality and retain the most salient features, pooling layers, such as max pooling, are used, which select the

# Backpropagation



Figure 2.2: This diagram illustrates the three-stage backpropagation mechanism for neural network training: (1) **Error Computation** calculates the loss as the difference between predicted output $\hat{y}$ and target output, (2) **Error Signal Transmission** propagates the error backward through network layers in reverse order, and (3) **Gradient Calculation** computes weight gradients using the chain rule. The network comprises three layers: *Input Layer* (yellow neurons $x_1, x_2, x_3$), *Hidden Layer* (blue neurons with weights $w$), and *Output Layer* (red neuron producing $\hat{y}$). This mechanism enables iterative weight updates that minimize prediction error through gradient-based optimization [13].

maximum value within a specified region of the feature map [14]. This combination of convolution, activation, and pooling enables convolutional neural networks to efficiently learn hierarchical representations from data.

## 2.1.4 Recurrent Neural Networks

Recurrent Neural Networks (RNNs) are designed to model sequential data where the output at each time step depends not only on the current input but also on previous states. This is achieved by maintaining a hidden state that is updated recursively as new inputs are processed. The basic RNN cell computes the hidden state $h_t$ at time $t$ using the following equation:

$$h_t = \tanh(W_x x_t + W_h h_{t-1}) \tag{2.5}$$

where $x_t$ is the input at time $t$, $W_x$ and $W_h$ are weight matrices, and tanh is a non-linear activation function [17]. RNNs are widely used for tasks such as

Figure 2.3: The text explains a diagram of how convolutional layers process a 3-channel (RGB) input tensor of size $w \times h \times 3$. Two filters ($F_1$ and $F_2$) are each composed of three $3 \times 3$ kernels (one per input channel), giving $2 \times 3 \times 3 \times 3 = 54$ weight parameters plus 2 biases, for 56 trainable parameters in total. Each filter performs a weighted sum plus bias followed by a nonlinearity (e.g., ReLU), producing an output activation map. The resulting output has spatial dimensions $w \times h$ and depth 2, equal to the number of filters. The main point is that a convolutional layer's output depth equals the number of filters, while its spatial dimensions depend on filter size, stride, and padding, forming the basis of feature extraction in CNNs for vision [16].

language modeling, event prediction, and time series analysis. The Figure 2.4 shows intuitively the difference between an RNN architecture and a Feed-Forward one.

However, vanilla RNNs often suffer from the vanishing and exploding gradient problems during training. The vanishing gradient problem occurs when gradients become exceedingly small, making it difficult for the network to learn long-term dependencies. Conversely, the exploding gradient problem arises when gradients grow uncontrollably, leading to unstable training. Gradient clipping can help mitigate exploding gradients, while architectural innovations such as Long Short-Term Memory (LSTM) and Gated Recurrent Unit (GRU) cells address vanishing gradients by introducing gating mechanisms that regulate information flow [19].

LSTM cells enhance the memory capabilities of RNNs by maintaining both a hidden state $h_t$ and a cell state $c_t$. The gating system consists of input, forget,

Recurrent Neural Network      Feed-Forward Neural Network

Figure 2.4: The text contrasts Recurrent Neural Networks (RNNs) with Feed-Forward Neural Networks for sequence and temporal data. RNNs have recurrent feedback connections that feed hidden activations back into the network at the next time step, giving them a form of memory to capture temporal dependencies and handle variable-length sequences (e.g., language, time series, speech). Feed-forward networks have only forward, acyclic connections from input to output, are easier to train and parallelize, but treat each input independently and cannot model temporal dynamics. Thus, RNNs are generally preferred for sequence modeling despite their higher computational complexity [18].

output, and gate gates, which are computed as follows:

$$i_t = \sigma(W_i[h_{t-1}, x_t]) \tag{2.6}$$
$$f_t = \sigma(W_f[h_{t-1}, x_t]) \tag{2.7}$$
$$o_t = \sigma(W_o[h_{t-1}, x_t]) \tag{2.8}$$
$$g_t = \tanh(W_g[h_{t-1}, x_t]) \tag{2.9}$$

The cell and hidden states are updated by:

$$c_t = f_t * c_{t-1} + i_t * g_t \tag{2.10}$$
$$h_t = o_t * \tanh(c_t) \tag{2.11}$$

To update the cell state $c_t$ at time step $t$, the previous cell state $c_{t-1}$ is first modulated by the forget gate $f_t$, which determines the extent to which past information should be retained. Simultaneously, the input gate $i_t$ regulates how much new information, derived from the current input $x_t$, is incorporated into the cell state. This

dual mechanism allows the LSTM to selectively preserve or overwrite information as needed.

The hidden state $h_t$ is then computed by applying the output gate $o_t$ to the transformed cell state, effectively controlling how much of the internal memory is exposed to the next layer or as output. This gating structure enables the LSTM to flexibly manage the flow of information across time steps.

A key reason why LSTMs mitigate the vanishing gradient problem is that, during backpropagation, the gradient with respect to the cell state is propagated primarily through element-wise multiplication by the forget gate $f_t$. This avoids repeated multiplication by weight matrices, which in traditional RNN can cause gradients to diminish exponentially. As a result, LSTMs are able to maintain and learn long-range dependencies, making them effective for modeling sequences with extended temporal structure [17, 19, 20].

## 2.2   Neuromorphic Engineering

Neuromorphic computing is a computational paradigm designed to emulate the architecture and behavior of biological nervous systems, such as the human brain. In recent years, the use of Artificial Neural Networks (ANNs) has surged; however, neuromorphic computing is increasingly gaining attention in the scientific community as a more efficient and promising alternative [21].

Unlike traditional models that process information sequentially and continuously, neuromorphic computing operates on asynchronous events. In this approach, signals known as "spikes" are generated only when a change occurs [22].

The core elements of neuromorphic computing are spiking neurons, which collectively form Spiking Neural Networkss (SNNs). Each neuron processes spikes independently and remains inactive in the absence of input. This event-driven method represents a significant departure from conventional ANNs, which continuously process input signals. By emulating the brain's mechanisms, neuromorphic computing can enhance efficiency and potentially increase computational power for specific tasks [22].

However, implementing these models is not without challenges. Hardware limitations, in particular, pose significant obstacles and must be carefully considered to fully leverage the advantages offered by neuromorphic technologies [21].

### 2.2.1   Spiking Neural Networks

Over the last decade, deep neural networks have achieved remarkable results across various fields including language understanding and image recognition. However, these advances come with substantial resource demands: modern AI models require

large datasets, considerable computational power, and significant energy consumption. For example, training large-scale language models such as GPT-3, which contains 175 billion parameters, has been estimated to consume around 190,000 kWh of energy [23]. The overall energy consumption of AI data centers is projected to reach 22 billion kWh by 2025, driven primarily by increasing model complexity and the computational intensity of training procedures [24]. Between 2012 and 2019, the energy required to train state-of-the-art models increased nearly tenfold, raising significant concerns about the environmental and economic sustainability of current AI development practices [25].

It is important to distinguish between the energy costs associated with training and those related to inference. Training represents the one-time, upfront energy investment required to develop a model, often consuming vast resources during this phase. Inference, on the other hand, corresponds to the energy expended every time the trained model is deployed to perform tasks such as understanding language, recognizing images, or making decisions in autonomous systems. In many real-world scenarios, particularly those involving large-scale or real-time applications, the cumulative energy consumption during inference can surpass that of training. Therefore, while training energy highlights the environmental footprint of developing AI models, inference energy is a critical metric for assessing their operational efficiency and sustainability in practical deployment contexts.

In a strong contrast, the human brain achieves remarkable computational efficiency by operating with only about 10 to 20 watts while processing complex sensory information [26]. Neuromorphic engineering strives to close this efficiency gap by designing systems inspired by the brain's organization and function.

The essential components of neuromorphic technology consist of:

- **Neuromorphic sensors**: detect only changes in input signals rather than continuously sampling, enabling much more efficient data capture [27].

- **Neuromorphic algorithms**: particularly SNNs, which process information with discrete events or spikes, as shown in Figure 2.5, mimicking the brain's temporal coding and typically consuming far less energy than conventional artificial neural networks [28].

- **Neuromorphic hardware**: specialized computing architectures that implement SNNs and brain-inspired computation, utilizing emerging device technologies such as memristors to improve speed and energy efficiency [29].

A central objective is to integrate the proven effectiveness of artificial neural networks with the energy-saving capabilities of spiking networks. Advances have shown that neuromorphic models running on dedicated platforms can significantly reduce energy consumption and inference latency while maintaining competitive accuracy [30].

This evolution in computing paradigms is especially critical for applications that require real-time processing under stringent power constraints, including autonomous machines, wearable devices, and brain-machine interfaces [31].



Figure 2.5: The text describes a neuromorphic algorithm based on spiking neural networks that mimics biological neural computation using temporal spike dynamics and event-driven processing. The architecture has four main components: an Input Layer that encodes information in spike timing; Sensory/Encoder Neurons that receive and propagate spike trains; Excitatory Neurons that integrate weighted excitatory and modulating inputs; and Inhibitory Output Neurons that implement lateral inhibition. A Winner-Take-All (WTA) mechanism among excitatory neurons ensures that the neuron receiving the strongest input suppresses others, producing a sparse, selective, and energy-efficient winner-driven response [32].

### 2.2.2 Biological Neurons and Their Structure

Biological neurons are the elementary signaling units within the nervous systems of animals and humans, forming complex networks capable of adaptive computation. Neurons transmit information primarily via discrete electrical impulses, called spikes or action potentials, facilitating communication between cells and throughout neural circuits [33, 34].

A typical neuron consists of several key structural components: dendrites, the soma, axon, and axon terminals. These elements underpin both signal integration and communication within neural circuits [34].

- **Dendrites**: Dendrites are highly branched projections that receive chemical or electrical signals from other neurons. The morphology and density of

dendritic branching modulate a neuron's integration properties and play a pivotal role in the spatial and temporal summation of inputs [34, 35].

- **Soma (Cell Body)**: The soma integrates incoming synaptic signals from the dendrites and determines whether the total input surpasses a critical threshold to generate an action potential. The dynamic interplay of excitatory and inhibitory signals shapes the membrane potential and output response [33, 36].

- **Axon**: The axon is a long, slender projection that carries the action potential from the cell body to distant targets. Axonal conduction properties affect both the spatial reach and timing of neural signals [34, 33].

- **Axon Terminals and Synapses**: Axonal terminals form synapses, the junctions across which signals are transmitted to adjacent neurons or effector cells, commonly using neurotransmitters. The structure, distribution, and plasticity of synapses are essential for learning and adaptation in neural circuits [33, 37, 35].

Neural learning and plasticity are mediated by activity-dependent changes in synaptic strength, enabling adaptive behavior and experience-dependent circuit refinement [35]. Modern computational models inspired by these principles offer interpretable and efficient representations of neural processing [36].

### 2.2.3 Neural Code

The neural code refers to the mechanisms by which the human brain represents and processes information, a topic that remains a central and unresolved challenge in neuroscience [38, 39, 40]. Despite decades of research, it is widely accepted that neural information encoding relies on three key principles: spikes, sparsity, and static suppression. First, neurons communicate primarily through the generation and propagation of action potentials, or spikes, which are all-or-none events that encode information in their timing and occurrence rather than their amplitude [38, 40]. This discrete, event-based signaling is fundamentally different from the continuous-valued representations used in conventional ANNs, and recent engineering analyses suggest that discrete coding is essential for reliable information transmission in the brain [40]. Second, neural activity is characteristically sparse, with most neurons remaining quiescent for extended periods and only a small subset active at any given time, a property that enhances memory efficiency and computational power [39, 41]. Third, static suppression mechanisms in sensory systems enable neurons to preferentially respond to dynamic, changing stimuli while filtering out static or redundant information, thereby optimizing sensory processing [38]. In terms of computational models, both ANNs and SNNs can address similar

tasks, but their neuron models differ fundamentally. ANNs compute a weighted sum of inputs and apply a nonlinear activation function such as ReLU, resulting in continuous-valued outputs. In contrast, SNNs accumulate weighted input spikes to a membrane potential, and a spike is emitted only when this potential crosses a threshold, making the output inherently event-driven and temporally precise [38, 40]. This distinction underpins the efficiency and biological plausibility of SNNs for modeling real neural computation. The Figure 2.6 shows the anatomy of a biological neuron.

### 2.2.4   Neuron Model

ANNs and SNNs are both capable of addressing similar computational tasks, but they differ fundamentally in their neuron models. In ANNs, each neuron computes a weighted sum of its inputs and passes this value through a non-linear activation function, such as the ReLU, to produce its output. This approach enables efficient gradient-based learning and is well-suited for static data representations [42].

In contrast, SNNs employ a more biologically inspired neuron model. Here, the weighted sum of inputs contributes to the neuron's membrane potential $U(t)$. When this potential reaches a defined threshold $\theta$, the neuron emits a spike, transmitting information to subsequent neurons. Inputs to SNNs are typically spikes that arrive at different time instants, and the neuron's output is event-driven rather than continuous. This mechanism allows SNNs to process temporal and event-based data efficiently, often resulting in lower energy consumption and more realistic neuronal dynamics [43].

While both ANN and SNN architectures can achieve comparable performance on certain tasks, SNNs offer advantages in terms of energy efficiency and biological plausibility, especially for event-driven applications. However, training SNNs remains challenging due to the non-differentiable nature of spike generation, requiring specialized learning algorithms.

### 2.2.5   Encoding and Decoding Spikes in SNNs

SNNs process information using discrete spikes, inspired by biological neural systems. A central challenge is how to encode input data as spike trains and how to decode spike outputs for interpretation. Recent research has advanced both encoding and decoding strategies, optimizing SNNs for efficiency and accuracy [44].

**Input Encoding:** Three principal encoding schemes are widely used:

- **Rate Coding:** The intensity of the input is represented by the firing rate or spike count within a time window. Higher input values yield more spikes. This method is robust but can require longer time windows for accurate representation [45].

- **Temporal (Latency) Coding:** Information is encoded in the timing of spikes, such as the time-to-first-spike (TTFS) approach, where stronger inputs cause earlier spikes. This enables rapid and energy-efficient computation, as fewer spikes are needed [46].

- **Delta/Change Coding:** Spikes are generated only when there is a change in input intensity, capturing dynamic features and reducing redundancy. This is particularly effective for event-based sensors like Dynamic Vision Sensors (DVS) and silicon cochleas, which directly produce spike streams in response to environmental changes [47].

**Output Decoding:** To interpret SNN outputs, several decoding strategies are employed:

- **Rate Decoding:** The predicted class is assigned to the neuron with the highest firing rate.

- **Latency Decoding:** The class is determined by the neuron that fires first.

- **Population Coding:** Aggregates information from multiple neurons, overcoming the limitations of individual firing rates and enabling more robust and rapid processing.

Recent benchmarking studies have compared these encoding and decoding methods, highlighting trade-offs between speed, accuracy, and energy efficiency. For example, time-to-first-spike coding can achieve high accuracy with minimal spikes, while rate coding offers greater error tolerance and is more compatible with backpropagation-based training [45, 46].

Figure 2.6: The text describes a biological neuron and its four key parts: the cell body, dendrites, axon, and axon terminals. The cell body houses the nucleus and organelles, maintains vital functions, and helps generate action potentials. Dendrites branch from the cell body to receive signals from other neurons at synapses. The axon is a single long projection that carries electrical and chemical signals away from the cell body, ending in axon terminals. These terminals form synapses with other neurons and release neurotransmitters to transmit signals. This structure underlies the directional flow of neural information where dendrites receive, the cell body integrates, and the axon transmits, and serves as a biological inspiration for artificial neural network design [48].

## 2.3 Large Language Models

Large Language Models (LLMs) are advanced systems capable of understanding and producing human language. They are designed to process text in a way that captures meaning and context, allowing them to generate coherent and relevant responses or content. Essentially, LLMs can "read" and "write" in a manner that closely resembles human communication. The utility of LLMs lies in their versatility. They can assist in generating written content, providing summaries, or creating explanations on complex topics. They are also used in conversational contexts, offering guidance, answering questions, or engaging in dialogue with users. Beyond communication, LLMs support decision-making and information management by extracting insights from large amounts of text or presenting information in an easily understandable form. In essence, LLMs act as intelligent language assistants, capable of understanding human expression and aiding in tasks that involve reading, writing, or reasoning with text, making them valuable across many areas of work and daily life.

## 2.3.1 Evolution of Large Language Models

The evolution of Language Models into modern LLMs traces a path from early statistical estimation to advanced neural architectures and alignment techniques [49]. Initially, Language Models were understood as probabilistic models of natural language, estimating the likelihood of word sequences, with early n-gram models relying on the Markov assumption to approximate the probability of the next word based on a small preceding context, but suffering from data sparsity and a lack of semantic understanding.

The introduction of neural networks reframed next-word prediction as a classification problem, augmented by word embeddings that captured semantic relationships in dense vector spaces, replacing sparse one-hot encodings, and enabling recurrent architectures such as RNNs and LSTMs [50] to handle longer contexts.

A transformative advance occurred with the Transformer architecture [51], which addressed long-term dependency limitations and allowed parallelizable processing through the self-attention mechanism, giving rise to decoder-only Transformers exemplified by the GPT family. Models like GPT-2 [52] and GPT-3 [23] dramatically scaled parameters (GPT-3 reaching 175B) and training data, confirming scaling laws and demonstrating few-shot in-context learning, while highlighting inefficiencies in oversized yet undertrained models.

This prompted compute-optimal training, as in Chinchilla [53] (70B parameters), which achieved superior performance by training on more tokens rather than simply increasing model size.

The final evolutionary step involved aligning LLMs with human preferences and practical utility, formalized through instruction tuning (improving zero-shot generalization, demonstrated by FLAN [54] and Reinforcement Learning with Human Feedback (RLHF) as implemented in InstructGPT [55], producing models that are Helpful, Honest, and Harmless (HHH).

This trajectory characterizes current accessible LLMs, including the LLaMA family [56], which combine massive scale, semantic understanding, and alignment strategies to serve as powerful, instruction-tuned conversational agents. The Figure 2.7 gives an idea of the exponential growth of the LLMs in the last years.

## 2.3.2 The Attention mechanism

The attention mechanism, which is central to LLMs, allows the model to dynamically assign relevance to different tokens within an input sequence. This mechanism overcomes the limitations of sequential architectures by enabling long-range dependencies to be captured efficiently. Specifically, self-attention computes pairwise interactions between all tokens in the sequence, permitting each token to attend to any other regardless of positional distance. The Figure 2.8 shows a visual example of attention computed in the context of words translation.

Figure 2.7: This timeline presents major technological milestones in NLP organized into four eras: (1) *Statistical Language Models* (1990–2000) including N-grams, (2) *Neural Language Models* (2000–2013) introducing Word2vec and RNN/LSTM architectures, (3) *Pre-trained Language Models* (2013–2017) featuring the Attention Mechanism, Transformers, BERT, and GPT, and (4) *Large Language Models* (2018–2023) encompassing GPT-2, GPT-3, ChatGPT, and GPT-4. The progression demonstrates the transition from static statistical representations to context-dependent transformer models capable of complex reasoning and few-shot adaptation [57].

The core computation is the scaled dot-product attention. Given input token embeddings, three matrices are derived through parameterized linear projections: queries $Q$, keys $K$, and values $V$. The attention output is formulated as:

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

where $d_k$ is the dimensionality of the key vectors and serves as a scaling factor to stabilize gradients during training. The softmax normalizes the similarity scores into a probability distribution, highlighting the most relevant tokens for each query.

Multi-head attention extends this formulation by computing several parallel attention outputs, each with its own learned projections. These outputs are concatenated and linearly transformed to produce the final attended representation, enabling the model to capture different semantic and syntactic aspects concurrently.

This architecture forms the basis of Transformers that power state-of-the-art LLMs like GPT-4 of OpenAI and Cloude 4.5 of Anthropic [51, 58].



Figure 2.8: This diagram illustrates the attention mechanism's core operation, demonstrating how a model selectively focuses on relevant input elements when processing a target token. The upper track shows French source tokens (*Comment, se, passe, ta, journée*) with their importance weights visualized through a color gradient bar (dark blue indicating high importance, light blue indicating low). The lower track displays the corresponding English target tokens (*How, was, your, day*) that the model attends to when generating the query. The connecting arrows represent attention weights, showing which source tokens contribute most to predicting each target token. This mechanism enables the model to capture long-range dependencies and semantic correspondences between languages, forming the foundation of sequence-to-sequence architectures and transformer models [59].

### 2.3.3 Applications of LLMs

LLMs have become increasingly prominent across research and industry, owing to their ability to generate, analyze, and adapt human-like text for a wide variety of purposes [60]. Their versatility extends from routine tasks, such as drafting emails, translating text, or summarizing documents, to more advanced applications involving programming assistance, data analysis, and content generation. This

adaptability makes LLMs not only useful as personal assistants or customer service agents but also as tools capable of managing and interpreting vast amounts of information that would otherwise demand significant human effort.

In medicine, for instance, LLMs are beginning to reshape both healthcare delivery and research practices. By analyzing patient records alongside extensive medical literature, they can provide clinicians with evidence-based recommendations, support diagnostic reasoning, and suggest treatment strategies. Beyond clinical practice, they facilitate patient engagement through conversational interfaces, offer scalable tools for medical education such as training simulations and tailored learning resources, and contribute to public health initiatives by detecting outbreaks, monitoring social responses, and disseminating accurate health information [60].

A similar transformative role is emerging in education. LLMs enable personalized learning by adapting materials to individual needs, while simultaneously assisting teachers with tasks such as lesson planning, grading, and developing accessible content. They are particularly effective in language learning, where they act as interactive conversation partners, capable of correcting grammar, guiding vocabulary use, and enhancing fluency. By providing transcription services for the hearing impaired or simplifying complex texts for learners with difficulties, they also make education more inclusive.

In the sciences, LLMs help researchers navigate overwhelming volumes of literature by summarizing findings, highlighting patterns, and even generating new hypotheses. Their ability to assist in drafting manuscripts and standardizing formatting further supports the research process, allowing teams across disciplines to communicate more effectively. In mathematics, they complement this role by breaking down complex problems, offering step-by-step explanations, and aiding in the verification of proofs, thus bridging the gap between abstract theory and applied contexts.

The legal field also benefits from LLMs, particularly in the processing and interpretation of extensive documentation. They can analyze case law, explain legal terminology, and support reasoning tasks, with domain-specific fine-tuning significantly enhancing their accuracy and practical utility. A similar pattern is seen in finance, where models, such as FinGPT [61], demonstrate the value of tailoring LLMs to industry-specific data. These systems support applications ranging from robo-advising and algorithmic trading to more complex decision-making, highlighting the importance of customization for high-stakes domains.

Finally, robotics represents another promising area of integration. Here, LLMs improve human-robot interaction by enabling robots to understand instructions in natural language and plan tasks accordingly. Their capacity to integrate information from diverse sources also equips robots with the ability to adapt to new environments, acquire new skills, and operate more collaboratively alongside humans [60].

Taken together, these applications illustrate how LLMs extend human capabilities and introduce efficiencies across domains as diverse as healthcare, education, law, science, and robotics. Yet, their rapid adoption also raises important challenges concerning data quality, bias, and explainability. Addressing these issues will be essential if LLMs are to realize their full potential as reliable and responsible tools for advancing research, industry, and society.

### 2.3.4    Training Paradigms

LLMs are developed through a multi-stage training process that enables them to acquire linguistic knowledge and adapt to different tasks. This process generally includes three main phases: *pretraining*, *fine-tuning*, and *instruction tuning*.

In the first phase, known as *pretraining*, the model is exposed to very large collections of unlabelled text data, as explained by [62]. During this stage, the model learns through self-supervised objectives that do not require human annotation. Two common objectives are causal language modeling [63], which trains the model to predict the next token in a sequence as seen in GPT models, and masked language modeling [63], which teaches the model to reconstruct hidden tokens within a sentence as used in BERT. Through this large-scale learning process, the model develops an extensive understanding of grammar, meaning, context, and factual knowledge.

After pretraining, the model is adapted through *fine-tuning*. In this phase, the model is trained on smaller, labeled datasets that focus on specific domains or applications [64]. Fine-tuning allows the model to apply its general linguistic and semantic knowledge to specialized tasks such as legal text analysis, biomedical literature summarization, or program code generation. This step improves accuracy and relevance within the target domain.

The third phase, *instruction tuning*, aims to make LLMs more responsive to human instructions written in natural language [65]. Instead of using datasets designed for a single task, instruction tuning relies on large collections of examples where each task is paired with a clear written instruction, such as "Translate this subsection into French" or "Summarize this article in three sentences." These instructions are called *prompts*. A prompt is an input text provided to the model that specifies the task or desired output. It serves as a guide that tells the model what to do and how to respond in context.

By learning from many such examples, instruction-tuned models acquire the ability to generalize to new tasks that they have never encountered before. This ability, known as *zero-shot generalization*, allows users to interact with models simply by providing well-formulated prompts. As a result, instruction tuning has become a key step in developing LLMs that are more practical, interactive, and adaptable to real-world applications.

### 2.3.5 State-of-the-Art Large Language Models

The development of LLMs by leading organizations has profoundly influenced both research and practical applications in artificial intelligence. Companies such as OpenAI, Anthropic, Google, and Mistral have introduced models that differ in scale, architecture, and focus, collectively shaping the state of the art. OpenAI's GPT series, for instance, has emphasized versatility and multimodal capabilities. Anthropic's Claude models prioritize safety and alignment, reflecting growing concerns about ethical AI and reliable instruction-following. Google's PaLM series has focused on scale and multilingual performance, demonstrating how massive models trained on extensive datasets can achieve competitive reasoning and language understanding across diverse tasks. Open-source initiatives such as Mistral complement these proprietary developments by providing accessible models with optimized architectures, encouraging experimentation, transparency, and reproducibility in research. Within the last LLMs' approaches we find MAKER: A framework leveraging massively decomposed agentic processes (MDAPs) to achieve zero-error execution in tasks requiring over one million steps [66].

Across these platforms, the benefits of modern LLMs are evident. They provide high-quality natural language understanding, reasoning, and generation, and support a wide range of tasks, from summarization and translation to coding assistance and domain-specific problem solving. Their availability, whether as API-accessible proprietary systems or fully open-source models, facilitates rapid prototyping and deployment in both academic and industrial contexts. At the same time, limitations remain. Even state-of-the-art models can produce hallucinations, display biases inherited from training data, and underperform on tasks requiring nuanced domain knowledge or low-resource language proficiency. Moreover, proprietary models often impose access and cost constraints, while open-source models shift the burden of deployment, fine-tuning, and computational resources onto the user.

In sum, the current landscape reflects a dynamic interplay between capability, accessibility, and responsibility. The evolution of LLMs illustrates both the technical progress achieved by leading companies and the ongoing challenges in reliability, fairness, and ethical deployment. Understanding these trends is crucial for positioning research and applications within a rapidly developing and heterogeneous AI ecosystem [67, 68, 69, 70].

### 2.3.6 Limitations of LLMs

Despite their impressive capabilities, Large Language Models have several inherent limitations [71]. First, they can generate information that appears confident but is factually incorrect, which poses challenges in reliability and trust. Second, they are prone to biases, reflecting patterns present in their training data, which can result in unfair, offensive, or otherwise undesirable outputs. Third, LLMs do not

truly understand the content they produce; they rely on patterns in language rather than reasoning or comprehension, limiting their ability to perform complex logical or causal inference. Additionally, their responses can sometimes be inconsistent or sensitive to small changes in input phrasing. Finally, LLMs require significant computational resources for training and deployment, which can limit accessibility and environmental sustainability. Recognizing these limitations is crucial when applying LLMs in real-world contexts, and it highlights the need for careful supervision, evaluation, and ethical consideration.

## 2.4 Agents Orchestration Frameworks

An agent, in the context of artificial intelligence, is an autonomous computational entity capable of perceiving its environment, making decisions, and executing actions to achieve specific goals, often by interacting with both digital and physical resources. Modern agents are designed to operate within complex, dynamic environments, leveraging a combination of internal reasoning, external tool usage, and adaptive learning to fulfill user-defined objectives. The increasing complexity of intelligent systems has created a strong need for frameworks capable of coordinating diverse computational tasks in a structured and reliable manner. Agents may leverage various underlying models, including but not limited to LLMs, to perform reasoning, generation, or prediction tasks. While LLMs can be powerful components, achieving accurate and context-aware results often requires their integration with external tools such as search engines, vector databases, and code execution environments. Managing this integration demands workflow orchestration frameworks, which provide mechanisms to define, execute, and monitor sequences of operations while maintaining state, handling errors, and enabling iterative refinement [60]. In the context of AI research and applications, orchestration frameworks play a crucial role in bridging the gap between isolated model capabilities and end-to-end solutions. They enable the design of pipelines that combine information retrieval, summarization, reasoning, and code generation into coherent processes. Moreover, they facilitate modularity, scalability, and reproducibility, essential qualities for both experimental research and practical deployment [1]. This subsection introduces the main frameworks relevant to this work, with particular focus on LangChain and LangGraph, which leverage graph-based state management to orchestrate interactions between intelligent models and external resources [2].

### 2.4.1 LangChain

LangChain [72] is a widely adopted open-source framework designed to facilitate the development of applications powered by LLMs. Its name reflects its core philosophy: "Language", referring to the central role of LLMs in processing and generating

37

natural language, and "Chain", emphasizing the composition of multiple components into structured sequences or "chains" of operations. The framework's primary contribution lies in providing modular abstractions that simplify the integration of LLMs with external resources such as APIs, databases, knowledge bases, and user-defined tools. At its core, LangChain structures complex workflows into reusable components, such as prompts, chains, memory, and agents, that can be composed to build sophisticated pipelines. These abstractions enable developers to move beyond isolated model queries, allowing LLMs to interact with structured data, retrieve external information, and make tool-augmented decisions. LangChain's architecture also emphasizes extensibility and interoperability, making it possible to integrate different model backends or adapt workflows to specialized domains with minimal effort. Its ecosystem has rapidly grown to include connectors for search engines, vector databases, and cloud services, positioning LangChain as a foundational framework for building scalable and production-ready LLM applications.

### 2.4.2   LangGraph

LangGraph [73] is a framework built on top of LangChain that introduces a graph-based approach to workflow orchestration for LLMs. While LangChain primarily relies on linear "chains" of components to structure interactions, LangGraph extends this paradigm by allowing developers to define stateful graphs, where nodes represent computational steps and edges determine dynamic transitions between them. This design makes it possible to construct more flexible workflows that can branch, loop, or adapt based on intermediate results, capabilities that are difficult to achieve with LangChain's linear architecture. Furthermore, LangGraph emphasizes persistent state management, enabling workflows to retain context across multiple iterations and user interactions, which is particularly important in tasks such as research assistance or code refinement. By combining graph flexibility with LangChain's modular components, LangGraph addresses the limitations of rigid sequential pipelines and provides a more expressive orchestration framework. Importantly, it is also a very active and rapidly evolving framework as of 2025, with frequent updates, community contributions, and growing adoption in both research and industry. This makes LangGraph a state-of-the-art choice for building complex, adaptive applications where decisions must be routed dynamically and context must be preserved over time.

### 2.4.3   State Management and Checkpointing

Modern conversational AI systems and multi-agent frameworks require sophisticated mechanisms to maintain context across extended interactions and recover from failures. State management preserves conversational history, intermediate

computational results, and workflow progression in systems processing sequential, dependent operations [74].

Unlike stateless architectures where each request is processed independently, stateful systems maintain persistent memory enabling coherent multi-turn dialogues and complex task decomposition. Traditional LLM serving operates statelessly, necessitating reprocessing of complete conversation history with each request, resulting in quadratic computational complexity as conversation length increases. Stateful architectures address this through caching mechanisms achieving 1.5-2x throughput increases and 60-75% latency reductions by reusing cached intermediate representations [74].

Graph-based orchestration frameworks such as LangGraph implement checkpointing at each execution step, capturing conversation state, variables, and workflow progress as snapshots organized into threads representing independent conversation instances. This enables human-in-the-loop intervention, time-travel debugging, and automatic recovery from failures. State transitions occur through directed graphs where nodes represent discrete operations and edges define valid transitions based on conditional logic, providing explicit control while maintaining flexibility for complex branching paths. Each node encapsulates specific functionality, such as querying tools or invoking language models, with defined transition conditions that determine subsequent states [75].

Message passing between nodes maintains ordered histories that serve as the primary mechanism for context propagation. Advanced implementations employ message reduction strategies, including context summarization and hierarchical compression, to manage conversation history growth while maintaining sufficient context for coherent behavior within finite context windows and computational constraints.

## 2.5    Information Retrieval in AI Systems

The effectiveness of intelligent systems often depends not only on their ability to generate language but also on their capacity to access, process, and integrate external information. LLMs, despite being trained on massive datasets, are inherently limited by the static nature of their training corpora and cannot continuously update their knowledge or guarantee factual accuracy. This limitation makes Information Retrieval (IR) a critical component in modern AI systems. By incorporating retrieval mechanisms, models can dynamically access up-to-date, domain-specific, or contextually relevant information, thereby extending their utility beyond the scope of memorized knowledge [76].

In the context of LLM-based applications, information retrieval serves multiple roles: enhancing factual reliability, grounding generative outputs in external sources, and improving performance on knowledge-intensive tasks. This section

reviews the main approaches and technologies for information retrieval in AI systems, progressing from foundational concepts to advanced integration strategies. The following section begins by examining traditional retrieval approaches, then explore modern vector-based methods, discuss specific retrieval tools, and conclude with retrieval-augmented generation paradigms that synthesize these components into cohesive systems.

### 2.5.1 Traditional Information Retrieval Approaches

Large language models are constrained by static training data, limiting their ability to access real-time information and verify outputs. Web research and academic search tools address this limitation by integrating external information sources into AI-driven applications. Web research tools, such as Google Search [77], Duck-DuckGo [78], and Baidu [79], provide access to diverse internet content and can be programmatically accessed through APIs or specialized wrappers, enabling AI systems to retrieve up-to-date content from news, technical documentation, and community forums. Retrieved content is processed through cleaning, parsing, and filtering before integration into language model pipelines, ensuring relevance and usability while grounding responses, reducing hallucinations, and adapting outputs to dynamic user needs.

Academic search tools provide specialized access to peer-reviewed publications and scholarly data, prioritizing domain-specific sources essential for reliable and verifiable knowledge. Platforms such as Google Scholar [80], arXiv [81], and Semantic Scholar offer structured access to metadata, abstracts, citation networks, and full-text articles through programmatic APIs and libraries such as the `scholarly` Python package [80] and the arXiv Python wrapper [81]. By incorporating authoritative academic sources, AI systems improve factual grounding and credibility in contexts demanding precision and reliability, including literature review, scientific summarization, and evidence-based decision-making.

### 2.5.2 Semantic Search and Vector-Based Retrieval

Semantic search leverages embeddings, dense numerical vectors in high-dimensional space, to represent text such that semantic similarity is captured by geometric proximity, enabling retrieval of semantically related documents even without exact keyword matches [82]. Embedding models, such as Sentence Transformers, encode textual content into fixed-length vector representations that preserve semantic relationships, capturing subtle semantic distinctions, handling synonymy and polysemy, and generalizing across domains. The quality of embeddings directly impacts retrieval performance; modern models are trained on large corpora using contrastive learning and representation learning techniques, with specialized variants available

for different content types including general text, code (e.g., CodeBERT), and multilingual content.

CodeBERT exemplifies domain-specific embeddings for code: it is a bimodal pre-trained model trained jointly on natural language and programming language pairs, using masked language modeling and replaced token detection objectives to capture code semantics such as variable scoping, function signatures, and API usage patterns [83]. By learning shared representations across both modalities, CodeBERT enables semantic code search where natural language queries can retrieve relevant code implementations, outperforming general-purpose embeddings on code-related tasks.

Vector databases, or vectorstores, such as Chroma [84], FAISS (Facebook AI Similarity Search) [85], Pinecone [86], and Weaviate [87] provide efficient storage and retrieval of high-dimensional embeddings at scale through specialized indexing structures optimized for similarity search operations [88]. These systems support various similarity metrics including cosine similarity, Euclidean distance, and dot product, and address key design considerations such as indexing efficiency, query latency, and scalability to millions or billions of vectors. Advanced implementations provide hybrid search capabilities combining vector similarity with keyword matching, dynamic index updates, and distributed architectures for horizontal scaling. When integrated with large language models in RAG pipelines, vector databases enable fetching of relevant external documents to provide as context for model responses, improving factual accuracy, reducing hallucinations, and allowing domain-specific adaptation without retraining [88].

### 2.5.3 Retrieval-Augmented Generation

RAG enhances Large Language Models by incorporating external knowledge sources through an information retrieval mechanism, addressing fundamental limitations including hallucination, outdated knowledge, and lack of transparent provenance [89]. RAG merges the parametric knowledge encoded in LLM parameters with dynamic, non-parametric memory represented by external sources, enabling continuous knowledge updates and domain-specific information integration without retraining. The foundational architecture, introduced by Lewis et al., combines a pre-trained sequence-to-sequence transformer as parametric memory with a dense vector index as non-parametric memory, accessed through a neural retriever, enabling models to leverage both learned representations and explicit external knowledge [76].

The RAG pipeline consists of three phases: indexing, retrieval, and generation [89].

- During indexing, documents are preprocessed, segmented into manageable chunks using fixed-size or semantic-aware strategies, and encoded into dense

41

vector representations using embedding models before storage in vector databases with efficient indexing structures.

- In the retrieval phase, user queries are encoded identically to index chunks, and vector similarity search identifies the top-k most relevant document segments using metrics such as cosine similarity or Euclidean distance.

- During generation, retrieved documents are integrated with the user query through prompt engineering; the combined context is fed to the language model, which synthesizes a response grounded in both its parametric knowledge and retrieved external information [89].

Advanced implementations incorporate re-ranking mechanisms [**rerank2024**], context selection algorithms [90], and iterative refinement loops [91] to improve accuracy.

Context injection strategies determine how retrieved passages are integrated into LLM prompts, commonly through direct concatenation or templated structures that clearly delineate context from queries. Effectiveness depends on passage ordering, total context volume relative to model context windows, and clarity of usage instructions. RAG formulations differ in their conditioning approach: sequence-level RAG conditions on identical documents throughout output generation, while token-level RAG allows different passages to inform different tokens, enabling nuanced knowledge synthesis from multiple sources [76].

## 2.5.4   Limitations and Challenges in Information Retrieval

While information retrieval has become integral to modern AI systems, current methodologies face several inherent limitations. Traditional IR approaches struggle with biases arising from data imbalance, relevance judgment subjectivity, and feedback loops that reinforce popular content [92]. With the integration of LLMs, these problems intensify: IR systems may over-represent machine-generated content (source bias), propagate factual inaccuracies (factuality bias), and prioritize popular items due to skewed training corpora (popularity bias).

Moreover, LLM-based retrievers often misinterpret user intent or contextual dependencies, leading to *instruction-hallucination* and *context-hallucination* biases [93]. Additional challenges include: (1) the cold-start problem for new or niche domains with limited indexed content; (2) computational overhead of embedding generation and similarity search at scale; (3) difficulties in maintaining consistency across multi-hop retrieval scenarios; and (4) potential security vulnerabilities such as prompt injection through retrieved content.

Addressing these limitations requires careful design of retrieval strategies, diverse training data, robust evaluation protocols, and hybrid approaches that combine multiple retrieval methods. Future directions include developing bias-aware

retrieval algorithms, implementing adversarial robustness measures, and designing evaluation frameworks that assess fairness, accuracy, and robustness in LLM-augmented information retrieval systems.

## 2.6  Validation of Knowledge Produced by LLMs

The generation of factual and accurate knowledge by LLMs represents a fundamental challenge in their deployment across high-stakes domains such as healthcare, legal analysis, and scientific research. While LLMs demonstrate remarkable linguistic fluency and reasoning capabilities, they remain vulnerable to producing factual inaccuracies, hallucinations, and knowledge inconsistencies, even when arriving at seemingly correct final outputs [94]. The validation of LLM-generated knowledge thus requires comprehensive frameworks that assess not only the correctness of final answers but also the factual integrity of intermediate reasoning steps and the grounding of outputs in verifiable external sources.

Contemporary approaches to knowledge validation in LLMs encompass multiple strategies that address different aspects of the factuality problem. RAG systems enhance LLM outputs by integrating external knowledge sources, retrieving relevant documents from trusted databases to ground model responses in verifiable information [95]. This grounding process mitigates hallucinations by anchoring predictions to external evidence rather than relying solely on the model's potentially outdated or incomplete parametric knowledge. However, RAG systems introduce their own validation challenges, including the need to verify retrieval quality, assess the helpfulness of retrieved documents, and detect contradictions between external knowledge and the model's internal beliefs [95].

Beyond retrieval-based methods, specialized validation frameworks have emerged to systematically evaluate factual accuracy throughout the reasoning process. The RELIANCE framework (Reasoning Evaluation with Logical Integrity and Accuracy for Confidence Enhancement) [96] addresses the critical vulnerability of factual inaccuracies within intermediate reasoning steps by integrating three core components: a specialized fact-checking classifier trained on counterfactually augmented data, a reinforcement learning mechanism using Group Relative Policy Optimization (GRPO) with multi-dimensional rewards, and a mechanistic interpretability method that analyzes neural activations during reasoning [94]. Extensive evaluation reveals that recent models like Claude-3.7 and GPT-o1 demonstrate factual reasoning accuracy of only 81.93% and 82.57% respectively, highlighting the persistent challenges in ensuring factual robustness across reasoning chains [94].

Effective knowledge validation must address multiple dimensions of factuality, including internal knowledge checking (assessing whether an LLM possesses relevant knowledge for a query), helpfulness checking (determining whether retrieved external documents assist or distract from accurate responses), and contradiction

checking (identifying conflicts between external context and internal model beliefs) [95]. Representation-based methods, which leverage the LLM's internal representations rather than relying solely on prompting or probability-based approaches, have demonstrated superior performance in these validation tasks, achieving accuracy improvements of up to 85% in informed helpfulness checking [95].

### 2.6.1 The DeepEval Framework

DeepEval provides a comprehensive, open-source evaluation framework specifically designed for systematic testing and validation of LLM outputs, conceptualized as "Pytest for LLMs" [97]. The framework addresses the critical need for standardized, rigorous evaluation methodologies by offering research-backed metrics, synthetic dataset generation capabilities, and seamless integration with continuous integration and deployment (CI/CD) pipelines [97].

The core architecture of DeepEval encompasses multiple evaluation dimensions essential for validating LLM knowledge production. The framework implements over 14 research-backed metrics covering diverse aspects of LLM performance, including G-Eval for custom criteria evaluation with human-like accuracy, hallucination detection, answer relevancy assessment, and specialized metrics for RAG pipelines (faithfulness, contextual recall, contextual precision) and agentic systems (tool correctness, task completion) [97]. These metrics utilize LLM-based evaluation approaches combined with statistical methods and natural language processing models that execute locally, providing both flexibility and privacy in the evaluation process [97].

DeepEval's G-Eval metric represents a particularly powerful validation tool, employing chain-of-thought reasoning to evaluate LLM outputs based on custom criteria specified by users. This metric allows practitioners to define evaluation steps that check for factual contradictions, assess completeness, and verify alignment with expected outputs while accommodating nuances such as vague language or opinion-based statements [97]. The framework supports both end-to-end evaluation, treating the LLM application as a black box, and component-level evaluation through non-intrusive tracing with the @observe decorator, enabling granular assessment of individual components such as retrieval modules, tool calls, and reasoning agents [97].

The framework facilitates comprehensive validation workflows by enabling dataset-level evaluation, where collections of test cases can be evaluated in parallel, and by providing integration with the Confident AI platform [98] for full evaluation lifecycle management, including dataset curation, benchmark comparison across model iterations, metric fine-tuning, and real-time monitoring of production deployments [97]. This integrated ecosystem supports iterative improvement cycles essential for maintaining factual accuracy and reliability in deployed LLM systems, addressing the critical need for continuous validation as models evolve and encounter diverse

real-world inputs.

## 2.7  Code Generation

Code generation has long been a fundamental topic in computer science and software engineering, concerned with automating the process of producing code from specifications or descriptions of intent. Traditional approaches, including template-based methods, domain-specific languages, and rule-based systems, provided limited flexibility and adaptability to diverse programming tasks. With the rapid development of large language models in recent years, a new paradigm of neural code generation has emerged, leveraging the semantic understanding and linguistic capabilities of LLMs to automatically produce functional code from natural language descriptions. These models demonstrate strong code comprehension and writing abilities across multiple programming languages and paradigms, enabling application to a wide variety of software engineering tasks and significantly enhancing developer productivity [99].

### 2.7.1  LLMs for Code Generation

Large language models have rapidly advanced code generation through a multi-stage process. Prior to training, data preprocessing ensures that datasets from open-source repositories are clean, standardized, and suitable for model learning. During training, LLMs construct complex internal representations of syntax, semantics, and interrelations among code elements. Post-training, LLMs generate code by analyzing prompts, retrieving relevant patterns, assembling code fragments, and producing final outputs [100]. Notable milestones including GitHub Copilot (2021) and Replit Ghostwriter (2022) demonstrated the increasing ability of LLMs to complete, explain, transform, and debug code.

As of 2025, leading models including OpenAI's Codex, Anthropic's Claude Sonnet 4.5, Google's Gemini 2.5 Pro, and Alibaba's Qwen3-Coder have demonstrated substantial improvements in reasoning, long-context understanding, and autonomous problem-solving capabilities [101, 102, 103, 104]. These models excel in generating, debugging, and refactoring complex multi-file codebases with minimal human guidance.

Despite these advances, significant limitations persist. LLMs require substantial computational resources, with large models demanding millions of GPU hours for training [100]. They remain vulnerable to syntactic and semantic errors that can produce incorrect or non-functional code, and bias in generated code is a documented concern, with studies reporting gender-biased outputs in a notable proportion of cases. Security risks also arise from unsanitized training data, potentially introducing vulnerabilities such as prompt injection or memory safety issues. These

45

challenges underscore the need for careful evaluation and mitigation strategies to ensure that LLM-based code generation is reliable, fair, and secure.

## 2.8 Code Validation

Despite extensive research on code generation applications, evaluation of LLM-generated code has received comparatively less attention. A significant gap exists between the proliferation of code generation applications and the development of robust evaluation methodologies that rigorously assess code quality, correctness, and fitness for production deployment. Bridging this gap requires systematic frameworks for evaluating generated code along multiple dimensions, including functional correctness, efficiency, maintainability, and security.

### 2.8.1 Evaluation of Code Generated by LLMs

The quality of software source code has been a central concern in software engineering since the 1960s, with numerous metrics, methods, and models proposed for evaluating code at the source level. These evaluation approaches form the foundation of modern software quality assurance and improvement practices.

With the rise of LLMs capable of generating code, concerns about the reliability of such code have emerged. Introducing unverified LLM-generated code into software systems can result in critical failures or security breaches. Consequently, recent research has focused on assessing LLM-generated code using established software quality criteria [105].

Software quality comprises multiple characteristics, such as functional suitability, performance efficiency, and compatibility, along with various sub-characteristics according to ISO/IEC 25010 [106]. Studies evaluating LLM-generated code are often categorized based on these characteristics, particularly into functional correctness, security, and other aspects [106].

Functional correctness is typically assessed using code evaluation metrics (CEMs). Match-based metrics, adapted from machine translation, such as BLEU and Code-BLEU [107], measure similarity between generated and reference code but may fail to capture functional equivalence. Execution-based metrics evaluate code by running it on test sets but can be computationally expensive.

LLMs trained on extensive code repositories have demonstrated impressive capabilities in generating functional source code. However, the security of such code remains a significant concern. Recent empirical studies show that LLM-generated code often incorporates common security vulnerabilities, including issues like cross-site scripting and injection flaws, at rates comparable to or sometimes exceeding those in human-written code [108]. Additionally, code quality aspects such as understandability, maintainability, and readability are crucial for long-term software

health. Metrics like cyclomatic complexity and cognitive complexity help quantify the ease with which developers comprehend and modify the code [109]. Research further reveals that code smells, patterns indicative of poor coding practices, can be propagated from training data to LLM outputs, raising concerns about the maintainability and clarity of AI-generated code [110]. These findings underscore the need for robust evaluation frameworks and human oversight in the deployment of LLM-generated software.

### 2.8.2 Sandbox Execution and Safety

Given the inherent risks of executing automatically generated code, sandboxing has become a crucial safety measure when deploying LLMs for code generation. Sandboxes provide an isolated environment where untrusted code can run without access to critical system resources, preventing malicious or unintended operations such as arbitrary file access, privilege escalation, or network misuse [111]. Recent studies highlight that LLM-generated code may inadvertently include vulnerabilities such as memory safety issues or SQL injection risks, which can pose serious security threats if executed directly [111]. Sandboxed execution frameworks mitigate these threats by restricting resource usage, monitoring execution, and terminating unsafe processes, thereby enabling researchers and practitioners to test LLM outputs safely [112]. Beyond local sandboxing, several cloud-based solutions, such as e2b, provide secure ephemeral execution environments that allow developers to run LLM-generated code in managed, isolated containers. These services ensure both scalability and security by offloading execution to dedicated infrastructure while enforcing strict access controls and resource limits [113]. Despite these safeguards, challenges remain in ensuring scalability, performance overhead, and the detection of more subtle logic-level vulnerabilities that sandboxes alone may not capture. Consequently, combining sandboxing with static analysis, formal verification, and vulnerability testing frameworks has been proposed as a more holistic strategy for ensuring safe and reliable deployment of LLM-generated code.

### 2.8.3 Code Static Analysis

Static analysis is a technique used to examine source code or binaries without executing the program, with the aim of identifying potential errors, vulnerabilities, and opportunities for optimization early in the development process. By analyzing the structure of the code, the flow of data, and the control paths, static analysis tools can detect problems such as null pointer dereferences, memory leaks, security flaws, and violations of coding standards before they cause runtime failures. In contrast to dynamic analysis, which observes a program during execution, static analysis provides a compile-time safety layer that helps developers improve the reliability and maintainability of their code.

Static analysis has long been a central component of software quality assurance, since it plays a proactive role in detecting bugs before software is deployed [114]. However, as modern software systems grow in complexity, as seen in large-scale projects such as snnTorch, static analysis tools face increasing challenges in scalability.

At the core of this issue lies the well-known trade-off between precision and scalability [115]. Highly precise analyses, such as path-sensitive methods, can accurately distinguish real bugs from false alarms, but they often struggle to scale efficiently to large and intricate codebases. In contrast, more scalable analyses are designed for speed and general applicability, yet they tend to rely on over-approximation, which increases the number of false positives and reduces their practical effectiveness.

Modules that perform static analysis parse and analyze syntax, types, and control flows, allowing developers to identify potential issues early in the development cycle.

Examples from the Python ecosystem include Pyright [116], Mypy [117], and Pylint [118], which focus on type checking, code style enforcement, detection of unused variables, and identification of possible runtime errors [**xu2023gradual**]. When integrated into the development workflow, such tools offer real-time feedback on code quality and potential defects, which significantly improves the overall reliability and maintainability of software projects.

In essence, static analysis modules function as an automated code review system, since they enforce coding standards, promote consistent design practices, and identify subtle issues that may be missed during manual review or testing. Their proactive nature, combined with their ability to scale across large and collaborative projects, makes them an essential part of modern software engineering pipelines, particularly in environments where reliability and safety are critical.

## 2.9   Human-in-the-Loop

The concept of Human In The Loop (HITL) encapsulates a paradigm where human expertise, judgment, and oversight are integrated into machine learning workflows to enhance model performance, interpretability, and ethical alignment. In the context of LLMs, HITL systems leverage human feedback at various stages: training, evaluation, and deployment, to iteratively refine outputs and ensure they reflect human values, contextual appropriateness, and factual accuracy [119]. Unlike traditional fully automated training pipelines, HITL approaches recognize that human cognition and ethical reasoning cannot be entirely captured by statistical optimization. Thus, feedback loops serve as a bridge between algorithmic intelligence and human judgment, allowing models to evolve dynamically in response to real-world user interactions [120].

## 2.9.1  Iterative Refinement

Iterative refinement in HITL workflows with LLMs agent orchestration refers to a repeated generate–evaluate–update loop in which one or more LLM agents produce candidate outputs, a human provides feedback, and the system updates subsequent outputs or orchestration policies accordingly. This process enables continuous improvement through successive correction and alignment cycles, leveraging human expertise to guide LLM behavior toward desired outcomes [121].

In the context of agentic systems, iterative refinement can be implemented at different levels:

- **Generation Level**: humans correct or critique model outputs. At this level, feedback is focused on the quality, correctness, and coherence of the output generated by the model in response to a specific task.

- **Policy Level**: feedback informs task decomposition or tool selection. At this level, humans guide the strategic decisions of the system on how to break down complex problems into subtasks and which tool or agent is most appropriate for addressing each part.

- **Orchestration Level**: human input adjusts coordination among multiple agents. At this level, feedback concerns how agents interact, communicate, and coordinate their actions to solve complex problems that require collaboration among multiple intelligent entities.

Recent work such as *Self-Refine* demonstrates the effectiveness of self-feedback and iterative improvement without additional supervised training, allowing models to iteratively critique and improve their own outputs [122]. Similarly, systems like IMPROVE show how iterative, component-wise pipeline refinement using LLM agents can yield stable performance improvements in complex multi-agent pipelines [123].

Iterative refinement may involve several workflow patterns:

- **Human-as-gate**: humans approve or reject candidate outputs. At this level, the human acts as a quality checkpoint, making binary decisions on whether the generated output meets acceptable criteria or requires regeneration.

- **Human-as-critic**: humans provide structured feedback. Here, humans go beyond simple approval/rejection to offer specific, actionable critique that guides the model toward improved outputs, with detailed comments on what needs improvement and why.

- **Human-as-instructor**: humans directly reformulate the task or constraints. At this level, humans take a more active role in shaping the problem definition itself, adjusting goals, requirements, or constraints based on deeper understanding of what the system should accomplish.

- **Automated self-refinement with human validation**: LLM agents refine their outputs autonomously before human review. Here, the system performs iterative improvements on its own outputs using self-critique or reflection mechanisms, with humans validating the final refined result rather than intervening at each step.

The combination of automated self-correction and human oversight has been shown to improve factuality, alignment, and interpretability in complex agent systems [124, 125].

## 2.9.2   Benefits and Challenges

The integration of iterative refinement within HITL orchestration presents a number of tangible benefits across dimensions of quality, safety, interpretability, and efficiency. One of the primary advantages lies in its capacity to improve alignment and factual accuracy. By incorporating human evaluation and corrective feedback into iterative loops, systems can mitigate common LLM limitations such as hallucination and task misalignment [121]. The presence of human oversight ensures that model outputs are not only syntactically correct but also contextually appropriate and aligned with human values or domain-specific standards. This quality assurance mechanism is particularly valuable in high-stakes environments, such as healthcare, finance, or law, where factual precision and interpretability are critical [125].

A further benefit of iterative refinement is its contribution to the overall safety and reliability of agentic systems. When human experts are embedded within the loop, they can intercept unsafe, biased, or inappropriate outputs before they propagate through the system or reach end-users. This "safety valve" effect transforms humans from passive evaluators into active participants in system governance, creating a dynamic safeguard that complements automated safety mechanisms. Moreover, iterative human feedback supports transparency and traceability, as each cycle of revision can be logged and attributed to either human or model actions, which is essential for auditing and compliance purposes. From an operational standpoint, hybrid refinement also enables efficient division of labor between human and artificial agents: routine, low-risk refinements can be automated through self-refinement strategies, while human attention is reserved for complex or ambiguous cases [122, 123]. In this way, the orchestration of human and LLM agents facilitates both scalability and adaptability, achieving better outcomes than purely manual or fully autonomous systems.

Despite these advantages, HITL orchestration also introduces significant challenges that must be addressed to ensure practical viability and ethical soundness. The most immediate obstacle concerns the cost and latency of human involvement. Continuous human review is resource-intensive, potentially reducing system throughput and limiting scalability in real-world applications [121]. Designing

adaptive orchestration strategies that determine when and where human input is most valuable thus becomes a central research question. Another challenge arises from the scalability and coordination of multi-agent systems: as the number of agents increases, orchestrating feedback, tracking dependencies, and maintaining coherence across iterative refinements become increasingly complex [124]. Ensuring consistent integration of human feedback across multiple agents or pipelines is technically demanding and often requires sophisticated scheduling and communication protocols.

Furthermore, the quality of human feedback itself is not guaranteed. Human annotators or reviewers may introduce bias, inconsistency, or noise into the iterative loop, which can destabilize the learning process if not properly managed [121]. To mitigate these effects, research has explored structured feedback formats, inter-annotator agreement measures, and aggregation techniques. Another critical issue concerns evaluation: the involvement of humans blurs the boundary between model autonomy and external assistance, complicating efforts to measure true model capability or the causal effect of human intervention. Metrics that account for human effort, latency, and risk are still under development and remain an open area of study. Finally, privacy and governance concerns pose additional constraints. When human reviewers handle sensitive data, as in medical or financial domains, organizations must implement strict data protection, access control, and minimization policies [125]. These requirements can limit the applicability of HITL workflows in regulated contexts.

## 2.10 Open-Source Frameworks for LLMs Deployment

As the use of LLMs expands, open-source frameworks play a critical role in enabling transparency, flexibility of deployment, data privacy, and local inference. Deploying and running models locally (or under full infrastructure control) helps mitigate reliance on proprietary APIs and external services. In this section, Ollama is described, the one used in the thesis, as a practical framework for local model deployment, followed by a brief survey of open-source model ecosystems and their relevance to research and deployment.

**Ollama**   Ollama is an open-source framework designed to simplify the deployment and orchestration of LLMs in local or private environments. Rather than functioning as a cloud-based service, Ollama provides users with the ability to run models entirely on their own hardware, ensuring control over data privacy and system integration. According to its official documentation, Ollama enables model management, inference, and customization through an accessible Application Programming Interface (API) and a lightweight runtime environment [126]. The framework

abstracts away low-level dependencies, such as model loading, quantization, and hardware optimization, so that users can focus on experimentation and workflow integration rather than infrastructure details.

Ollama supports a wide range of open-source models, including families derived from LLaMA, Mistral, and Gemma, and provides unified access to these models through both command-line and API interfaces [127]. This allows developers and researchers to integrate language, vision, and multimodal reasoning capabilities into local applications without relying on external APIs. Recent studies have used Ollama as a backend in privacy-preserving or on-premise systems [**liu2024optimizing**, **matotek2025evaluating**, 128]. Such integrations highlight Ollama's role as a bridge between open-source model research and practical, locally deployable applications [**gruetzmacher2024porting**, 129], offering an increasingly viable alternative to proprietary cloud platforms [**chen2025privacy**, **richardson2025patient**, **matotek2025singleboard**].

# Chapter 3

# Materials and methods

The neuromorphic computing domain, despite its significant potential for addressing energy efficiency and real-time processing challenges, faces a critical accessibility crisis rooted in fragmented development infrastructure and inadequate documentation. As established in the Introduction and Background sections, the existing neuromorphic ecosystem is characterized by several fundamental limitations: most implementations remain tightly coupled research prototypes lacking reusability, standardized development methodologies are absent, and comprehensive documentation is fragmented across disparate sources and communities. Developers seeking to implement neuromorphic systems encounter steep learning curves, dependency on expert consultation, and limited access to structured, production-grade development tools comparable to those available in conventional deep learning frameworks like TensorFlow or PyTorch.

This thesis addresses this gap by proposing a comprehensive, multi-agent architecture capable of orchestrating diverse specialized systems to support developers throughout the neuromorphic application lifecycle. The central objective is to extend AI-driven development assistance to the neuromorphic domain by integrating Large Language Models with structured reasoning, domain-specific knowledge bases, information retrieval mechanisms, and automated validation frameworks. The proposed approach treats development assistance not as an isolated code generation function, but as a cohesive system that synthesizes research knowledge, generates production-grade code, and validates implementations across multiple quality dimensions.

To achieve this objective, this chapter describes the technical infrastructure, methodologies, and component implementations that constitute the complete system. The Materials and Methods chapter is organized to progressively build understanding from foundational components to integrated system architecture:

**Section 3.1 — snnTorch Framework** introduces the primary neuromorphic simulation framework employed throughout the thesis. This section describes

core framework components including neuron models, surrogate gradient mechanisms, and temporal dynamics, then presents practical architectural patterns for implementing spiking neural networks. Understanding snnTorch is essential, as the code generation branch leverages snnTorch-specific knowledge bases to generate domain-appropriate implementations.

**Section 3.2 — Neural Network Intelligence (NNI)** describes the automated hyperparameter optimization toolkit integrated into the system. This section covers NNI's search space specification, trial management, and built-in algorithms. The code generation branch includes specialized NNI knowledge infrastructure to generate optimized training and hyperparameter search configurations for neuromorphic models.

**Section 3.3 — In-Context Learning** establishes the prompt engineering and few-shot adaptation strategies employed to enhance model reasoning. This section explains how domain knowledge is embedded into LLM contexts through structured examples, specialized vocabularies, and task-specific reasoning templates.

**Section 3.4 — Implementation of LangGraph** constitutes the architectural core of the thesis, detailing the graph-based multi-agent orchestration framework. This extensive section is subdivided into four components:

- **3.4.1 Web Search Branch**: Describes the iterative query generation, web search execution, source parsing, summarization, reflection, and evaluation pipeline for synthesizing current technical documentation and contextual information.

- **3.4.2 Academic Research Branch**: Details the academic query formulation, literature retrieval, source aggregation, and scholarly synthesis mechanisms for integrating peer-reviewed research knowledge into system responses.

- **3.4.3 Code Generation Branch**: Presents the core code synthesis orchestrator, specialized agents for snnTorch and NNI domains, vector store configuration, ranked retrieval mechanisms, and iterative refinement through execution feedback and human-in-the-loop validation.

- **3.4.4 Graph Routing Logic**: Explains the conditional routing mechanism that determines branch selection, state transitions, and exit conditions based on query classification and intermediate results.

**Section 3.5 — Used Tools** describes the external frameworks and utilities integrated into the system:

- **3.5.1 Agno Framework**: Documents the lightweight agentic framework providing tool orchestration, agent abstraction, and seamless integration with LangGraph workflows. Explains the tool system, RAG toolkits (TavilyTools, ArxivTools, DuckDuckGoTools, GoogleSearchTools), and agent execution mechanisms.

- **3.5.2 DeepEval**: Presents the LLM-as-a-Judge evaluation framework employed for systematic validation of generated knowledge. Details the four core RAG metrics (Faithfulness, Answer Relevancy, Contextual Relevancy, Hallucination Detection), metric scoring methodology, and multi-metric evaluation composition strategies.

**Section 3.6** — **Models and Infrastructure** describes the computational and software infrastructure supporting the system. This section encompasses:

- **3.6.1 Large Language Models Used**: Profiles the heterogeneous collection of open-source models (Mistral 7B for routing, Qwen3 for tool orchestration, GPT-OSS 20B for complex reasoning, and specialized evaluator models) selected for specific tasks and computational constraints.

- **3.6.2 Model Selection Rationale**: Justifies the multi-model strategy through task-specific capability requirements, computational efficiency considerations, and trade-offs between latency, reasoning quality, and context window capacity.

- **3.6.3 Ollama Infrastructure**: Explains the local deployment framework enabling privacy-preserving, reproducible model serving without dependence on external commercial APIs.

- **3.6.4 Used Hardware**: Specifies the computational resources (NVIDIA RTX A4000 GPU with 16 GB VRAM, CUDA infrastructure) employed for experimental execution.

## 3.1   snnTorch

`snnTorch` is an open-source Python framework built on top of PyTorch, designed to support the development, training, and simulation of SNNs. Rather than functioning as a simple library, `snnTorch` provides a modular ecosystem that integrates spiking neuron models, surrogate gradient methods, event-based data processing tools, and a collection of training utilities. This makes it possible to incorporate biologically inspired temporal dynamics within standard deep learning workflows. By leveraging PyTorch's automatic differentiation and GPU acceleration, `snnTorch` enables efficient experimentation with hybrid SNN–ANN architectures and supports both supervised and unsupervised learning paradigms. The framework is

widely used for research in energy-efficient computation and neuromorphic learning systems, serving as a bridge between neuroscience-inspired models and modern machine learning methodologies [130].

### 3.1.1 Core Framework Components

The framework is organized into four main modules, each addressing a critical aspect of SNN development:

- **Spiking Neuron Models**: This module includes a variety of neuron models such as the Leaky Integrate-and-Fire (LIF), Integrate-and-Fire, Synaptic, and Alpha neurons. It also provides recurrent spiking architectures, including **Spiking Long Short-Term Memory (SLSTM)** and **Spiking Convolutional LSTM (SConv2dLSTM)**, which extend classical LSTM dynamics into the spiking domain to capture long-range temporal dependencies while maintaining event-driven behavior.

- **Surrogate Gradient Methods**: Because spike generation is non-differentiable, snnTorch implements several gradient approximation techniques (fast sigmoid, arctangent, piecewise linear, or user-defined functions). These surrogate gradients enable end-to-end training of SNNs using backpropagation through time (BPTT).

- **Spike Generation and Data Conversion**: This module [131, 132] supports the conversion of static datasets (e.g., MNIST, CIFAR) into spike trains through rate coding, temporal coding, or thresholding mechanisms. It also handles event-based formats such as DVS (Dynamic Vision Sensor) data [133].

- **Functional Utilities**: This component defines commonly used operations including spike-based loss functions (rate-based or time-to-first-spike), regularization techniques, membrane potential normalization, and visualization tools for spikes and membrane traces.

### 3.1.2 Practical Architecture Patterns

In practical scenarios, snnTorch facilitates the construction of both purely spiking networks and hybrid architectures that integrate classical deep learning components with spiking neurons. A common example is a convolutional SNN in which a standard `nn.Conv2d` layer performs feature extraction, while the subsequent processing stages rely on spiking neurons such as `snn.Leaky`, equipped with configurable decay rates and surrogate gradient functions. The output layer typically aggregates membrane potentials over time, enabling continuous-valued predictions while maintaining event-driven computation throughout the network.

Beyond supervised classification tasks, snnTorch also supports unsupervised learning. A simple and illustrative example is the **spike-based autoencoder**. In this architecture, an encoder composed of spiking neurons compresses the input spike train into a latent representation, and a decoder, also spiking, attempts to reconstruct the original input. Learning proceeds by minimizing the discrepancy between the reconstructed and original spike trains, allowing the network to discover meaningful temporal patterns without access to labels. Throughout the process, the framework's state-management utilities (e.g., `utils.reset`) ensure that neuron membrane potentials are correctly reset between input sequences.

The modular design of snnTorch allows researchers to explore a wide range of architectural variants, including different surrogate gradient functions, temporal integration constants (e.g., membrane decay $\beta$), and firing thresholds. This flexibility makes it possible to tailor SNN models for tasks where both accuracy and energy efficiency are critical considerations.

## 3.2  Neural Network Intelligence (NNI)

Neural Network Intelligence (NNI) is a lightweight yet powerful open-source toolkit developed by Microsoft [134] to automate machine learning workflows including hyperparameter optimization, neural architecture search, model compression, and feature engineering. Built with a modular architecture and designed for scalability across local machines, cloud services, and distributed environments, NNI provides plug-and-play AutoML techniques that integrate seamlessly with popular deep learning frameworks like PyTorch and TensorFlow. The toolkit's web-based UI enables real-time monitoring and visualization of experiments, making it an essential resource for researchers and practitioners seeking to optimize neural network design and training at scale.

### Hyperparameter Optimization (HPO)

NNI's hyperparameter tuning component employs multiple optimization strategies through its **Tuner** module. The tuner iteratively suggests hyperparameter sets, evaluates them through trial runs, and maintains a history of results to guide subsequent suggestions. Available tuning algorithms span from basic strategies including random search, grid search to advanced techniques such as Bayesian optimization, simulated annealing, particle swarm optimization, and evolutionary algorithms. Users specify search spaces using declarative syntax, defining parameter types (uniform, log-uniform, choice, normal distribution) and ranges. For example, configuring a random tuner requires specifying continuous parameters (e.g., learning rate uniformly distributed in [0.001, 0.1]) and categorical parameters (e.g., optimizer choice between SGD, Adam, RMSprop). The framework automatically

orchestrates parallel trial execution, collects metrics, and determines when to terminate underperforming configurations through the **Assessor** module.

## 3.3 In-Context Learning

One of the most surprising capabilities of modern LLMs is in-context learning [135], which emerges as a by-product of scale and the Transformer architecture. Unlike traditional machine learning systems, where models must be explicitly retrained for new tasks, LLMs can learn to perform novel tasks on the fly by conditioning on examples provided in the input prompt. For instance, if a prompt contains several examples of movie reviews labeled with "positive sentiment" or "negative sentiment," the model can infer the classification pattern and correctly label new reviews without any parameter updates. This phenomenon demonstrates that the model has implicitly acquired the ability to generalize from patterns in its input sequence, effectively treating the prompt as a temporary training set. In-context learning is a key enabler of "few-shot" and "zero-shot" performance, where models adapt to new domains or tasks with little or no additional supervision. It represents a fundamental shift in how AI systems can be deployed, offering adaptability without retraining and significantly expanding the practical versatility of LLMs [136, 137, 138]. In this thesis, this capability was massively exploited.

Figure 3.1: The Figure shows the agent orchestration workflow with three parallel processing branches. The academic research branch (left) refines the input query into a more effective form and processes scholarly sources through summarization and reflection mechanisms. The web search branch (center) performs query generation, real-time web research, and source summarization. The code generation branch (right) implements retrieval-augmented generation with iterative feedback collection, normalization, and sandbox validation. All branches enable continuous iterative execution as feedback loops return to the initial routing stage for ongoing refinement.

## 3.4 Implementation of LangGraph

This section details the internal architecture and operational workflow of the LangGraph-based system developed in this work. The implementation is organized into three major branches, Web Search, Academic Research, and Code Generation, each designed to address a distinct aspect of the information retrieval and knowledge synthesis pipeline. The following subsections describe the design principles, execution

logic, and interactions that govern each branch within the overall agentic framework.

The system implements a parallel three-branch architecture where each branch addresses a specific knowledge synthesis task while maintaining shared state management and asynchronous execution through LangGraph's state machine framework.

The Web Search Branch, the branch in the center of the graph in the Figure 3.1, retrieves and synthesizes web-based knowledge through iterative refinement cycles. It enforces deterministic query generation for effective research, executes multi-source searches with fallback mechanisms, conducts up to web iterations of summarization with gap-based reflection to prevent query stagnation, and applies four RAG quality metrics for evaluation.

The Academic Research Branch focuses on scholarly literature by searching across ArXiv, Semantic Scholar, and Google Scholar. This branch can be observed in the left part of the Figure 3.1. It requires query generation optimized for academic terminology, iterative summarization. This branch has the same stagnation avoidance mechanism of web search branch.

The Code Generation Branch orchestrates multi-agent code generation with sandboxed validation. It routes user queries to specialized agents (snnTorch, NNI, called by *generate* node) via LLM classification, retrieves context with multiple extended queries using
CodeBERT-indexed vectorstores with weighted query ranking, generates multi-file code marked explicitly (`# FILE: filename.py`), parses responses. The code is validated code in E2B sandbox environments with dependency detection and static analysis, and supports interactive refinement loops. The system offers also the possibility to normalize the generated code with a reference code and to execute both codes in order to obtain performance metrics. The right part of the graph in the Figure 3.1 exposes the workflow of the code branch.

### 3.4.1   Web Search Branch

The web research branch represents a critical component of the implemented information retrieval system, as it is responsible for gathering, synthesizing, and validating web-based knowledge at scale. Its importance stems from the system's requirement to fetch up-to-date information reliably. This branch operates within the LangGraph-based agentic architecture, implementing an iterative refinement cycle that combines LLMs, vector-based semantic search, and quality evaluation metrics to produce coherent, evidence-grounded research summaries.

The web research branch is implemented as a specialized workflow within a larger knowledge production pipeline. The system operates through a sequence of interconnected nodes, each performing distinct transformations on a shared state

object (`SummaryState`). The architecture leverages asynchronous execution patterns to maximize throughput and responsiveness across the retrieval-generation-evaluation cycle.

Key architectural principles include state-driven execution, where all operations maintain and transform a centralized state object containing research topics, accumulated queries, search results, and evolving summaries; asynchronous pipeline design, where each node executes as an async coroutine enabling non-blocking operations and integration with long-running external services; tool composition, allowing multiple search backends (Tavily, Google Search, DuckDuckGo) to be composed through a unified agent interface; and checkpoint persistence for recovery and state inspection.

The implementation relies on several key frameworks and libraries: **LangGraph** for state machine and workflow orchestration, **LangChain** for document loading and embedding management, **Ollama** for local language model serving, **Chroma** for vector database functionality, **BeautifulSoup** for HTML parsing, **DeepEval** for quality assessment metrics, and **CopilotKit** for real-time UI integration.

### Query Generation

The workflow begins with a user-provided research topic. Rather than immediately executing a query, the system first generates a refined, structured search query through an LLM-guided process:

```
1. Input: research_topic (string)
2. Format research topic with prompt template
3. Invoke ChatOllama with JSON mode (temperature=0)
4. Parse returned JSON to extract 'query' field
5. Output: search_query (string)
```

This generation step provides three key benefits. First, setting temperature to 0 ensures deterministic query formulation, guaranteeing that the same research topic produces nearly identical search queries for reproducibility. Second, enforcing JSON mode constrains the LLM output to structured, parseable results rather than free-form text, preventing malformed queries. Third, the LLM leverages its domain knowledge to reformulate natural language research topics into effective search terminology that maximizes retrieval relevance.

### Web Search Execution

Once a query is formulated, the system executes a web search using a configurable agent architecture. Search parameters include maximum results of "x" documents per search with full document content retained, using the Tavily search API with fallback options to Google/DuckDuckGo:

```
1. Create Agent with Ollama backend (qwen3:latest)
2. Attach TavilyTools for web search capability
3. Construct search instruction prompt + query
4. Execute agent.arun() in non-blocking mode
5. Parse response content as search results
6. Emit progress to UI
```

The use of `asyncio.to_thread()` is critical as it prevents blocking the event loop during the agent's synchronous execution, maintaining responsiveness in the overall workflow. Search results are collected as unstructured text containing webpage snippets, titles, and URLs, stored in state as `raw_search_result` for subsequent parsing.

### Result Parsing and Formatting

After search execution, raw results are parsed into a structured format using a second LLM pass:

```
1. Input: raw_search_result (unstructured text)
2. Invoke ChatOllama with web_search_expected_output prompt
3. Extract JSON from response
4. Validate and normalize source information
```

After JSON extraction, **normalization** standardizes source metadata (title, URL, date) into a consistent schema. This processing stage reduces context bloat while preserving informational diversity, balancing conciseness with coverage quality for downstream summarization.

### Iterative Summarization

For each batch of search results, the system either generates a new summary or extends an existing one. This two-stage approach allows the system to accumulate knowledge progressively across multiple search iterations. Output is post-processed to remove any XML-like thinking tags from the language model, ensuring a clean, presentation-ready summary.

The system uses a two-stage summarization approach:

```
if existing_summary exists:
Extend existing summary with new search results
else:
Generate new summary from search results


Invoke ChatOllama and remove intermediate reasoning tags
```

This enables progressive knowledge accumulation across multiple search iterations, with post-processing to clean any intermediate reasoning from the model output.

**Reflection and Follow-up Query Generation**

After each summary generation, the system reflects on gaps in coverage to guide subsequent searches:

```
1. Invoke ChatOllama with reflection_instructions prompt
2. Input: current running_summary + research_topic
3. Use JSON mode to extract structured follow-up query
4. Validate that follow-up_query is not null
5. Fallback to generic query if generation fails
```

The reflection phase ensures that the system doesn't re-query the same information but instead targets unexplored dimensions of the research topic, enabling progressive exploration and preventing query stagnation. The LLM identifies knowledge gaps by analyzing the current summary against the original research topic, with the resulting follow-up query used in the next iteration of web search.

**Loop Control and Termination**

The system maintains a `research_loop_count` to prevent infinite loops:

```
if research_loop_count <= max_web_research_loops:
    route = "web_research" # Continue another iteration
else:
    route = "finalize_summary" # Proceed to evaluation
```

The `max_web_research_loops` parameter is configurable (in the code) per execution, allowing flexible control over research depth.

**Quality Evaluation**

Upon reaching the loop limit, the system enters an evaluation phase using four complementary RAG quality metrics:

1. **Faithfulness**: Measures whether summary claims are grounded in retrieved documents, preventing hallucinated information not present in sources.

2. **Answer Relevancy**: Assesses whether summary content directly addresses the research topic, ensuring topical alignment with the original query.

3. **Contextual Relevancy**: Evaluates whether retrieved documents are genuinely relevant to the context, reducing noise from tangentially related search results.

63

4. **Hallucination Detection**: Identifies unsupported claims in the summary, flagging content fabrication risk.

## Summary Finalization

Upon completing iterations and evaluation, the system generates a final deliverable:

```
1. Collect all accumulated sources: sources_gathered
2. Merge into structured format:
   - Summary section: running_summary
   - Sources section: formatted_sources list
3. Return with evaluation_results metadata
4. Emit exit signal to CopilotKit UI
```

All sources are tracked throughout execution and formatted consistently, preserving document title and URL, extraction timestamp, relevance confidence scores, and token budget information. This ensures reproducibility and enables users to trace claims back to original sources.

## Error Handling and Robustness

The system implements multi-stage JSON parsing with fallback strategies:

```
1. Attempt direct json.loads(content)
2. If fails: extract JSON substring using regex
3. If fails: attempt markdown fence removal then retry
4. Fallback: return empty/default result with error logging
```

This defensive approach handles variations in LLM output formats across different model variants and API versions.

All async functions include try-except blocks that log errors with context information, emit error status to the UI, return safe default states to prevent pipeline stalls, and preserve partial results when possible.

GPU/CUDA resources are explicitly managed through `torch.cuda.empty_cache()` calls after embedding operations, proper vectorstore connection closure, and memory release between batch operations.

## State Management and Performance

The central state object (`SummaryState`) maintains research topic, search query, raw search results, formatted source metadata, accumulated raw results, iteratively refined summary, iteration counter, and evaluation results. State transitions follow strict ordering enforced by the graph topology, ensuring that operations occur in valid sequences.

The web research branch implements a principled, iterative approach to knowledge synthesis combining keyword search, semantic retrieval, language model reasoning, and formal quality assessment. By maintaining clear separation of concerns across query generation, search, parsing, summarization, reflection, and evaluation phases, the system achieves both robustness and flexibility for diverse research applications.

## 3.4.2   Academic Research Branch

The academic research branch automates rigorous retrieval, synthesis, and evaluation of scientific literature by integrating scholarly and preprint search tools within the agentic LangGraph workflow. This branch complements web-based research by ensuring coverage of peer-reviewed and domain-specific academic sources, thus improving both the reliability and depth of generated knowledge.

**Academic Query Generation**

The workflow initiates with LLM-guided query construction, where a natural language research topic is transformed into a structured, field-focused query suitable for academic databases:

```
Input: research_topic (string)

Format research topic with prompt template

Invoke ChatOllama with JSON mode (temperature=0)

Parse returned JSON to extract 'query' field

Emit status to UI via CopilotKit

Output: search_query (string)
```

This ensures deterministic, reproducible query formulation by constraining the LLM to JSON-structured output. The system reformulates natural language research topics into optimized academic search terminology while guaranteeing consistent, parseable results across multiple executions.

**Academic Literature Retrieval**

Once a query is formulated, the system executes targeted searches across multiple academic databases using the agentic architecture:

```
Create Agent with Ollama backend

Attach ArxivTools for academic preprint search

Construct academic search instruction prompt + query

Execute agent.arun() in non-blocking mode

Parse response content as academic results

Emit progress to UI

Output: raw_academic_result (unstructured academic metadata)
```

The system targets multiple academic repositories, particularly Arxiv for preprints and via agent tool composition, Semantic Scholar and Google Scholar APIs, aggregating results including titles, abstracts, author lists, publication venues, and metadata.
Results are collected as unstructured text containing paper summaries, author information, and citations, stored in state for subsequent parsing.

### Source Parsing and Deduplication

After search execution, raw results are parsed into a structured format using a second LLM pass:

```
Input: raw_academic_result (unstructured academic output)

Invoke ChatOllama with academic_expected_output prompt

Extract JSON from response (with robustness for variations)

Validate and normalize academic metadata
```

### Summarization of Academic Findings

For each batch of academic sources, the system generates or extends a running summary:

```
if existing_summary exists:
human_message = "Extend the existing summary: {existing_summary}
Include new academic sources: {academic_sources}
That addresses the research topic: {research_topic}"
else:
human_message = "Generate a summary of these academic sources:
{academic_sources}
That addresses the research topic: {research_topic}"

invoke ChatOllama with:
- system_prompt: academic_summarizer_instructions
- user_message: human_message

remove_think_tags(response) # Clean any intermediate reasoning
```

This two-stage approach enables progressive knowledge accumulation across multiple search iterations. Post-processing removes any XML-like thinking tags inserted by the language model, ensuring clean, presentation-ready academic summaries.

**Iterative Reflection and Query Refinement**

After each summary, the system reflects on gaps in coverage to guide subsequent searches:

```
Invoke ChatOllama with reflection_instructions prompt

Input: current running_summary + research_topic

Use JSON mode to extract structured follow-up query

Fallback to generic academic query if generation fails
```

The reflection phase ensures that the system doesn't re-query the same papers but instead targets unexplored research dimensions or methodological gaps. This is done by reflecting on the knowledge gap and creating queries that cover that gap, enabling progressive exploration and preventing query stagnation.

**Loop Control and Branch Termination**

The system maintains a `research_loop_count` strictly limiting the number of academic exploration cycles:

67

```
if research_loop_count <= max_academic_research_loops:
route = "academic_research" # Continue another iteration
else:
route = "finalize_academic_summary" # Proceed to evaluation
```

The conservative loop limit balances coverage depth with computational efficiency.

## Quality Evaluation of Academic Summaries

Upon reaching the loop limit, the system enters an evaluation phase using the same four complementary RAG quality metrics:

1. **Faithfulness**: Ensures summary claims are grounded in retrieved academic papers, preventing hallucinated citations.

2. **Answer Relevancy**: Verifies that academic synthesis directly addresses the research topic.

3. **Contextual Relevancy**: Evaluates whether retrieved papers are genuinely topically relevant.

4. **Hallucination Detection**: Identifies unsupported claims or fabricated references.

The evaluation process executes as follows:

```
User prompt: "Do you want to evaluate this academic summary?"

If yes:
a. Run evaluation in separate thread via asyncio.to_thread()
b. Create LLMTestCase with:
- input: research_topic
- actual_output: running_summary
- context: academic_sources (full list)
- retrieval_context: academic_sources (for Faithfulness)
c. Initialize Ollama model (deepseek-r1:latest)
d. Invoke deepeval.evaluate() with all four metrics
e. Extract metric scores (range 0.0--1.0)
f. Emit results to UI with per-metric scores


If no:
a. Skip evaluation
b. Proceed to finalization
```

The evaluation runs in a separate thread to avoid blocking uvloop (the async event loop), critical because DeepEval's metrics perform LLM calls internally that would otherwise stall the entire pipeline.

**Summary Finalization and Output Structuring**

Upon completing iterations and evaluation, the system generates a final deliverable:

```
Collect academic sources: academic_sources_gathered

Merge into structured format:

Academic Summary section: running_summary

References section: formatted_sources list with metadata

Evaluation Results section: metric scores (if evaluated)

Return with evaluation_results metadata

Emit exit signal to CopilotKit UI
```

All sources are tracked throughout execution and formatted consistently, preserving title, authors, venue, DOI/URL, publication date, and token budget information. This ensures reproducibility and enables users to trace claims back to original academic sources with complete bibliographic information.

### 3.4.3 Code Generation Branch

The Code Generation Branch is a specialized subsystem designed to orchestrate the generation, validation, and iterative refinement of executable Python code in response to research queries. It combines RAG, multi-pass search optimization, and sandboxed code validation to produce semantically coherent, production-ready code artifacts. The branch operates within the LangGraph's architecture and integrates specialized language models, vector databases, and secure execution environments to ensure code quality and correctness.

**Architecture and Multi-Agent Code Generation**

The proposed Code Generation Branch is designed as a hierarchical, multi-agent system that prioritizes modularity and secure validation. At a high level, the architecture enforces a strict separation of concerns: the *Orchestration Layer* interprets user intent, the *Generation Layer* utilizes domain-specific knowledge to produce code, and the *Execution Layer* validates the output in a secure environment.

69

**High-Level Workflow** The process begins with semantic routing, where a user query is analyzed to determine the specific coding domain (e.g., spiking neural networks or hyperparameter tuning). This routing directs the task to specialized agents equipped with repository-specific context. To ensure robustness, the generation phase is decoupled from validation; code is not merely generated but eventually tested in a sandboxed environment to detect errors before delivery.

**Component Implementation** Drilling down into the specific technical components, the system comprises:

- **Orchestrator Agent**: Acts as the system's entry point, routing queries based on semantic analysis.

- **Specialized Code Agents**: Domain-specific agents (specifically for snnTorch and NNI) that generate code using retrieved context.

- **Vector Retrieval System**: Supports the agents via two dedicated vector-stores utilizing CodeBERT embeddings for semantic code search.

- **Code Parser and Extractor**: Manages multi-file generation by parsing explicit file markers (`# FILE: filename.py`), ensuring modular composition.

- **Sandbox Executor**: A low-level validation environment powered by E2B (a secure cloud sandbox provider) that executes generated code to verify functionality.

- **Feedback Processing**: A refinement loop that classifies execution results to guide iterative improvements.

The system enforces strict separation between code generation (via LLM) and code validation (via sandboxed execution), enabling error detection before delivery. Generated code uses explicit file markers (`# FILE: filename.py`) for modular composition, with each file representing a logical component.

**Orchestrator Agent Workflow** Algorithm 1 describes the orchestrator's decision logic for routing queries to specialized agents.

Here's a longer and more detailed version:

The orchestrator performs a multi-step validation on every tool's output before it is considered acceptable. This validation procedure includes several checks designed to ensure both structural integrity and functional reliability. First, the orchestrator verifies that the error field is explicitly null, confirming that the tool executed without raising exceptions or reporting internal failures. Next, it ensures that the returned code block meets a minimum length requirement of more than

70

---

**Algorithm 1** Orchestrator Agent Tool Selection

---

**Require:** User query $q$, available tools $T = \{\text{snnTorch}, \text{NNI}\}$
**Ensure:** Selected tools with focused queries
1: $prompt \leftarrow$ Create tool selection prompt with query $q$
2: $response \leftarrow \text{LLM}(prompt)$
3: $toolCalls \leftarrow \text{JSON.parse}(response)$
4: **for** each $call \in toolCalls$ **do**
5: $\quad toolName \leftarrow call.\text{name}$
6: $\quad toolQuery \leftarrow call.\text{query}$
7: $\quad result \leftarrow \text{ExecuteTool}(toolName, toolQuery)$
8: $\quad$ **if** $result.\text{error} = \text{null}$ **and** $|result.\text{code}| > 50$ **then**
9: $\quad\quad validOutputs.\text{add}(toolName, result)$
10: $\quad$ **end if**
11: **end for**
12: **return** $validOutputs$

---

50 characters. This prevents accidental acceptance of empty or malformed outputs and helps guarantee that the tool generated meaningful content. Finally, the orchestrator checks for the presence of required file markers, which act as structural indicators showing that the output conforms to the expected formatting or file-generation conventions. Only when all three of these criteria are successfully satisfied does the orchestrator treat the tool's output as valid and proceed with further processing.

## Vectorstore Initialization and Query Optimization

**Dual Vectorstore Configuration**   The system maintains two specialized vectorstores implemented using Chromadb, an open-source, lightweight vector database optimized for AI applications. Chromadb involves similarity search and retrieval-augmented generation [84].

- **snnTorch Vectorstore**: Online documentation from `snntorch.readthedocs.io`, recursive loading with custom depth limit, focuses on spiking neural networks, LIF neurons, training methods.

- **NNI Vectorstore**: Local examples from the `esempi_NNI/` directory, plus the API reference document; focuses on hyperparameter tuning configurations and experiment setup patterns.

Both employ CodeBERT embeddings (`mchochlov/codebert-base-cd-ft`), optimized for code semantics with cosine similarity metric. The vectorstore initialization follows an equal pattern (Algorithm 2).

**CodeBERT Memory Optimization**   To handle large-scale embedding on constrained GPU memory (8GB VRAM), a custom batching strategy is implemented:

```
batch_size = 8
for i in range(0, len(texts), batch_size):
    batch = texts[i:i+batch_size]
    embeddings = model.encode(batch)
    gc.collect()
    torch.cuda.empty_cache()
```

This approach trades computational speed (15-20% overhead) for memory safety, enabling vectorstore creation on consumer hardware.

---

**Algorithm 2** Vectorstore Initialization

---

**Require:** Persist directories $D_{snn}, D_{nni}$, source URLs/paths
**Ensure:** Initialized vectorstores

1: **if** $D_{snn}$ exists **and** $D_{nni}$ exists **then**
2:      Load existing vectorstores from disk
3:      **return** cached vectorstores
4: **end if**
5: $docs_{snn} \leftarrow$ RecursiveURLLoader(snnTorch docs, depth $= 7$)
6: $docs_{nni} \leftarrow$ DirectoryLoader(esempi_NNI) $+$ APIRef
7: $docs_{nni}$.insert$(0, \text{APIRef})$ {Priority insertion}
8: $chunks_{snn} \leftarrow$ TextSplitter$(docs_{snn}, \text{chunk} = 512, \text{overlap} = 50)$
9: $chunks_{nni} \leftarrow$ TextSplitter$(docs_{nni}, \text{chunk} = 512, \text{overlap} = 50)$
10: $V_{snn} \leftarrow$ Chroma.from_documents$(chunks_{snn}, \text{CodeBERT})$
11: $V_{nni} \leftarrow$ Chroma.from_documents$(chunks_{nni}, \text{CodeBERT})$
12: Persist $V_{snn}$ to $D_{snn}$, $V_{nni}$ to $D_{nni}$
13: **return** $V_{snn}, V_{nni}$

---

**LLM-Guided Query Generation and Ranking** Rather than static keyword queries, the system dynamically generates a custom number of diverse search queries tailored to each research question using a two-phase LLM process:

- **Phase 1: Query Generation**

    - The system instructs the LLM to:

    ```
    Generate 6 DIVERSE, SPECIFIC search queries covering:
    - Architecture and model structure
    - Parameters and configuration
    - Training methodology
    - Implementation patterns
    - Evaluation metrics
    Return JSON: ["query1", "query2", ..., "query6"]
    ```

- **Phase 2: Query Ranking**

    - The LLM then ranks the generated queries by relevance to the user question:

    ```
    User Question: {question}
    Candidate Queries: {generated_queries}

    Rank by RELEVANCE to answering the question.
    Return JSON: {
      "ranked_queries": ["most_relevant", "second", ...],
      "reasoning": "explanation"
    }
    ```

    - Ranked queries receive relevance weights $w_i = 1.0 - (i - 1) \times 0.05$ for query rank $i$.
    - Retrieved documents are tagged with their source query and weight, prioritizing high-relevance sources during LLM context construction.

**Multi-Pass Retrieval with Weighted Context**

- Algorithm 3 implements the weighted multi-pass retrieval strategy.

- The weighted context enables the LLM to prioritize information from higher-ranked queries during code generation.

---

**Algorithm 3** Multi-Pass Weighted Retrieval

---

**Require:** Ranked queries $Q = [q_1, q_2, \ldots, q_n]$, vectorstore $V$, $k = 20$
**Ensure:** Weighted context string

1: $context \leftarrow$ ""
2: $totalDocs \leftarrow 0$
3: **for** $i \leftarrow 1$ to $n$ **do**
4:      $w_i \leftarrow 1.0 - (i - 1) \times 0.05$
5:      $docs \leftarrow V.\text{retrieve}(q_i, k)$
6:      $totalDocs \leftarrow totalDocs + |docs|$
7:      **for** each $doc \in docs$ **do**
8:          $source \leftarrow doc.\text{metadata.source}$
9:          $content \leftarrow doc.\text{page\_content}[0 : 600]$
10:          $context \leftarrow context + $f"Query {i} - {source} - weight {w\_i:.2f}: {content}—"
11:      **end for**
12: **end for**
13: **return** $context[0 : 8000]$ {Limit to 8k chars}

---

**Specialized Agents and Unified Response Parsing**

**snnTorch Agent**   The snnTorch agent specializes in generating spiking neural network code. Its output schema guarantees:

```
{
  "code": "# FILE: model.py\n...\n# FILE: utils.py\n...",
  "files": {"model.py": "...", "utils.py": "..."},
  "summary": "Generated SNN model with LIF neurons",
  "confidence": 0.95,
  "queries_used": 6,
  "error": null
}
```

The agent constructs prompts emphasizing documentation priority: (1) retrieved context (primary), (2) snnTorch official docs (gap-filling), (3) general PyTorch knowledge (fallback).

**NNI Agent**   The NNI agent generates hyperparameter tuning configurations with strict JSON output requirements. It enforces dual-file patterns (`config.py` for search space, `train.py` for training loop) and validates JSON parseability before returning results.

**Unified Response Parser**   Algorithm 4 implements a multi-strategy parsing approach with graceful degradation.

**Escape Sequence Normalization**   LLM responses often contain literal escape sequences from JSON or markdown encoding. The normalizer handles:

```
def normalize_code_escapes(text):
    text = text.replace(r'\n', '\n')    # Literal \n → newline
    text = text.replace(r'\t', '\t')    # Literal \t → tab
    text = text.replace(r'\"', '"')     # Escaped quotes
    text = text.replace(r'\\', '\\')    # Double backslash
    return text
```

This ensures code readability and correct execution in the sandbox.

**File Extraction from Code**   Algorithm 5 extracts individual files from the unified code string.

---

**Algorithm 4** Unified Response Parser

---

**Require:** LLM response text $R$

**Ensure:** Parsed result dict

 1: $R \leftarrow R.\text{strip}()$
 2: **if** $R$ contains markdown code fences **then**
 3:     Remove fences: $R \leftarrow R.\text{remove\_fences}()$
 4: **end if**
 5: **Strategy 1: JSON Parsing**
 6: $data \leftarrow \text{JSON.parse}(R)$
 7: Map alternative field names to canonical `code` field
 8: Normalize escape sequences: $data[\text{"code"}] \leftarrow \text{normalize}(data[\text{"code"}])$
 9: **return** $data$
10: **Strategy 2: Markdown Code Block Extraction**
11: $blocks \leftarrow \text{REGEX.findall}(\text{r'```[a-z]*}\backslash\text{n(.*?)```'}, R)$
12: **if** $|blocks| > 0$ **then**
13:     Deduplicate and merge blocks
14:     **return** $\{\text{"code"} : \text{merged}, \text{"parse\_method"} : \text{"markdown"}\}$
15: **end if**
16: **Strategy 3: Raw Text Fallback**
17: **if** $|R| > 50$ **then**
18:     **return** $\{\text{"code"} : R, \text{"parse\_method"} : \text{"raw"}\}$
19: **end if**
20: **return** $\{\text{"error"} : \text{"Could not parse response"}\}$

---

---

**Algorithm 5** File Extraction with Regex Markers

---

**Require:** Code string $C$ with # FILE: markers
**Ensure:** Dict mapping filename $\rightarrow$ content
1: $pattern \leftarrow$ r'^# FILE:\s*(\S+\.py)\s*$'
2: $lines \leftarrow C.\text{split}('\backslash n')$
3: $files \leftarrow \{\}, currentFile \leftarrow$ null, $currentContent \leftarrow []$
4: **for** each $line \in lines$ **do**
5:     $match \leftarrow \text{REGEX.match}(pattern, line)$
6:     **if** $match \neq$ null **then**
7:         **if** $currentFile \neq$ null **then**
8:             $files[currentFile] \leftarrow \text{join}(currentContent)$
9:         **end if**
10:       $currentFile \leftarrow match.\text{group}(1)$
11:       $currentContent \leftarrow []$
12:     **else if** $currentFile \neq$ null **then**
13:       $currentContent.\text{append}(line)$
14:     **end if**
15: **end for**
16: **if** $currentFile \neq$ null **then**
17:     $files[currentFile] \leftarrow \text{join}(currentContent)$
18: **end if**
19: **return** $files$ if $files \neq \{\}$ else $\{$"main.py" $: C\}$

---

**Multi-Agent Code Integration** When multiple agents contribute code (e.g., snnTorch + NNI), intelligent merging applies a priority scheme:

1. **Priority 1**: `config.py` (configuration)

2. **Priority 2**: `search_space.json` (NNI search space)

3. **Priority 3**: `model.py` (architecture)

4. **Priority 4**: `utils.py` (helpers)

5. **Priority 5**: `train.py` (training loop)

6. **Priority 6**: `main.py` (entry point)

Files are sorted by priority and concatenated with file markers. Duplicates are resolved by keeping the first occurrence (higher-priority agent).

### Sandboxed Validation and Feedback Loop

**E2B Sandbox Execution** Generated code is executed in an isolated E2B sandbox (`code-interpreter-v1` template) with 1GB memory and CPU-only execution. Algorithm 6 describes the multi-file execution process.

---

**Algorithm 6** Sandboxed Multi-File Execution

---

**Require:** Code string $C$ with file markers
**Ensure:** Execution results
 1: $files \leftarrow$ parse_multifile_code($C$)
 2: $mainFile \leftarrow$ detect_main_file($files$) {Priority: main.py > train.py > experiment.py}
 3: $sandbox \leftarrow$ E2B.create("code-interpreter-v1")
 4: $packages \leftarrow$ extract_packages($C$) {Parse import statements}
 5: $sandbox$.run("pip install " + join($packages$))
 6: **for** each $(filename, content) \in files$ **do**
 7:    $sandbox$.upload("/home/user/" + $filename, content$)
 8: **end for**
 9: $result \leftarrow sandbox$.run("python " + $mainFile$)
10: $sandbox$.kill()
11: **return** {"exit_code" : $result$.exit_code, "stdout" : $result$.stdout, "stderr" : $result$.stderr}

---

**Dependency Auto-Detection and Installation**   Import statements are extracted using regex:

```
pattern = r"(?:^import|^from)\s+(\w+)"
matches = re.findall(pattern, code, re.MULTILINE)
```

A mapping table handles non-standard package names:

```
pip_mapping = {
    'cv2': 'opencv-python',
    'PIL': 'Pillow',
    'sklearn': 'scikit-learn',
    'snntorch': 'snntorch',
}
```

Built-in modules (`os`, `sys`, `json`, etc.) are filtered out before installation.

**Static Type Analysis**   Before execution, Pyright analysis detects type errors, undefined variables, and incompatible method calls. Results are non-blocking but reported to users for transparency.

**Interactive Feedback Loop**   After code generation or sandbox execution, users select from four actions via an interrupt mechanism:

- **approve**: Accept code and terminate

- **regenerate**: Request modifications with feedback

- **evaluate**: Compare against reference implementation

- **execute**: Run specific file in sandbox

Algorithm 7 describes the LLM-based feedback classification.
Example feedback classifications:

```
"execute train.py" → {"response": "execute", "file_name": "train.py"}
"Make it faster" → {"response": "regenerate"}
"Looks good" → {"response": "approve"}
"Compare with reference" → {"response": "evaluate"}
```

---
**Algorithm 7** Feedback Decision Classification

---
**Require:** User feedback text $F$
**Ensure:** Action and optional filename
  1: $prompt \leftarrow$ Create classification prompt with examples
  2: $response \leftarrow \text{LLM}(prompt + F)$
  3: $data \leftarrow \text{JSON.parse}(response)$
  4: $action \leftarrow data[\text{"response"}]$ {approve/regenerate/evaluate/execute}
  5: $filename \leftarrow data.\text{get}(\text{"file\_name"}, \text{null})$
  6: **return** $(action, filename)$

---

**Evaluation Mode and Performance Metrics**   When users select *evaluate*, the system enters comparison mode with code normalization and performance measurement. Normalization standardizes variable names, imports, and formatting while preserving logic. Performance metrics are injected via wrapper code:

```
import time
import torch

start = time.time()
output = model(input_tensor)
elapsed = time.time() - start

params = sum(p.numel() for p in model.parameters() if p.requires_grad)
peak_mem = torch.cuda.max_memory_allocated() / 1024 / 1024

print(f"Forward pass time: {elapsed:.6f}s")
print(f"Total parameters: {params}")
print(f"Peak memory: {peak_mem:.2f} MB")
```

Both generated and reference code execute with identical inputs in the sandbox, enabling objective comparison.

### 3.4.4   Graph Routing Logic

In the following is illustrated the conditional routing between nodes. The graph implements the following edges:

```
START → route_question
route_question → {search_relevant_sources (code),
                  generate_query (web),
                  generate_academic_query (academic)}
search_relevant_sources → generate
generate → collect_feedback (interrupt)
collect_feedback → process_feedback
process_feedback → {regenerate: generate,
                    approve: END,
                    evaluate: code_normalization,
                    execute: check_code_sandbox}
check_code_sandbox → reflection → collect_feedback
```

In the Figure 3.1 the full graph can be observed.

Conditional edges enable dynamic routing based on state values. The feedback loop continues until the user approves or maximum iterations (default: 3) are reached.

82

The system implements multi-level error recovery:

1. **Parser Level**: JSON parse fails → markdown extraction → raw text fallback

2. **Agent Level**: Validate error field, code length, file markers before acceptance

3. **Sandbox Level**: Timeout handling, missing package fallback, import error reporting

All errors return structured dictionaries (never raw exceptions), ensuring system stability. Failed agent outputs are logged but do not terminate the workflow: remaining agents continue execution.

## 3.5 Used tools

### 3.5.1 Agno Framework

Agno is an open-source, Python-based agentic framework designed to facilitate the development of intelligent AI agents with integrated memory, knowledge, tools, and reasoning capabilities. Agno provides a lightweight, composable architecture that emphasizes declarative agent construction, model-agnostic design, and first-class tool support [139]. The framework enables rapid prototyping and deployment of multi-agent systems through a clean API that abstracts the complexities of agent orchestration, tool coordination, and state management while maintaining full transparency for debugging and auditability.

**Tool System and Toolkits**

Tools occupy a central role in Agno's architecture as primary mechanisms for agent interaction with external systems. Tools are Python functions or classes that expose specific capabilities to agents, enabling interaction with external systems such as APIs, databases, web search engines, and file systems. The framework provides two primary tool abstractions: individual **function tools** for single operations and **toolkits** for collections of related functions that share internal state and coordinate behavior. Toolkits are implemented as Python classes that register multiple methods as callable tools, allowing agents to invoke them transparently based on task requirements.

**RAG Tools for Information Retrieval**

In the thesis implementation, four core toolkits from Agno's pre-built collection specifically utilized, to enable comprehensive information retrieval across web and academic sources:

**TavilyTools:** Provides advanced web search capabilities through the Tavily API, a search engine optimized for LLM-driven applications. The toolkit supports configurable search depth (`standard` vs. `deep`), maximum token limits per result, and output formatting (sourcedAnswer for synthesized responses with citations vs. searchResults for raw context data). TavilyTools was employed in the web research branch to execute primary web searches with advanced depth settings, retrieving up to 5 documents per query with full raw content preserved for subsequent LLM-based parsing and summarization.

**ArxivTools:** Enables retrieval of academic preprints and research papers from the arXiv repository. The toolkit queries the arXiv API using structured search parameters and returns paper metadata including titles, authors, abstracts, publication dates, and arXiv identifiers. ArxivTools was integrated into the academic research branch to complement web-based sources with peer-reviewed and domain-specific scientific literature, ensuring coverage of cutting-edge research findings and formal scholarly discourse.

**DuckDuckGoTools:** Provides privacy-focused web search through the DuckDuckGo API, serving as a fallback search backend when Tavily or other primary tools are unavailable or rate-limited. The toolkit executes keyword-based searches and returns result snippets, titles, and URLs. While not the primary search tool in the implementation, DuckDuckGoTools was configured as an alternative within the agent toolkit stack to ensure search robustness and redundancy.

**GoogleSearchTools:** Enables search through Google's Custom Search Engine API, providing access to indexable web content with support for configurable result limits and ranking. The toolkit returns structured results including titles, snippets, and URLs. GoogleSearchTools serves as an alternative search provider within the agent toolkit stack, offering broader web coverage and integration with Google's search index while being subject to rate limits and quota constraints.

**Agent Execution and Tool Coordination**

LangGraph gents execute tools through an internal reasoning loop that analyzes user queries, determines which tools are necessary, invokes them sequentially or in parallel, and synthesizes their outputs into coherent responses. Tool calls are transparent by default (configurable via `show_tool_calls`), enabling inspection of agent decision-making and facilitating debugging. In my implementation, agents were instantiated with the Ollama backend and configured with specific toolkits based on branch requirements (TavilyTools for web research, ArxivTools for academic research).

**Integration of Agno Tools into LangGraph Workflows**

Agno tools integrate seamlessly into LangGraph state machines as callable functions within graph-based workflows. In the developed architecture, each research branch (web and academic) instantiates branch-specific toolkits, executes queries through the tool coordination layer, and returns raw search results to the Lang-Graph state for subsequent parsing, deduplication, and summarization. This integration pattern leverages Agno's composable tool abstractions while maintaining LangGraph's control over overall workflow orchestration, state management, and conditional routing.

### 3.5.2 DeepEval

DeepEval is an open-source, LLM-powered evaluation framework designed to assess the quality of RAG systems through quantitative, interpretable metrics. Rather than relying on manual evaluation or subjective criteria, DeepEval employs the LLM-as-a-judge paradigm, where a language model autonomously evaluates generated outputs against multiple complementary dimensions [97]. This systematic approach enables comprehensive quality assurance of generated content without human annotation.

**Core RAG Evaluation Metrics**

DeepEval provides four fundamental built-in metrics specifically designed for RAG pipeline evaluation:

**Faithfulness:** Measures whether the generated output's claims are factually grounded in the retrieved documents. Faithfulness answers: *What fraction of the generated statements are supported by the retrieval context?* This metric prevents hallucination, the generation of information not present in sources, by validating each factual claim against provided documents. Formally:

$$\text{Faithfulness} = \frac{\text{Supported Claims}}{\text{Total Claims}} \tag{3.1}$$

**Answer Relevancy:** Evaluates whether the generated output directly addresses the user's original query. Answer relevancy measures topical alignment by classifying generated statements as relevant or tangential to the research question. This metric ensures the output focuses on answering what was asked, independent of factual correctness. Formally:

$$\text{Answer Relevancy} = \frac{\text{Relevant Statements}}{\text{Total Statements}} \tag{3.2}$$

**Contextual Relevancy:** Assesses the quality of retrieved documents by measuring whether they are genuinely relevant to the user's query. Unlike faithfulness,

which validates the generator's use of context, contextual relevancy validates the retriever's selection process. This metric identifies whether the retrieval system has selected appropriate source documents. Formally:

$$\text{Contextual Relevancy} = \frac{\text{Relevant Retrieval Documents}}{\text{Total Retrieved Documents}} \tag{3.3}$$

**Hallucination Detection:** Explicitly identifies claims in the generated output that contradict or lack support in the retrieval context. This metric flags negative cases, statements that directly contradict sources or introduce unsupported assertions, distinct from positive alignment measured by faithfulness. Formally:

$$\text{Hallucination Score} = 1.0 - \frac{\text{Contradicted/Unsupported Claims}}{\text{Total Claims}} \tag{3.4}$$

## Metric Scoring and Threshold-Based Evaluation

Each metric produces a score in the range $[0.0, 1.0]$, where 1.0 represents perfect performance and 0.0 represents complete failure. The interpretation of metric scores typically employs a threshold-based approach: scores above a defined threshold are considered acceptable, while scores below the threshold indicate quality concerns requiring investigation or remediation.

The threshold selection is a design choice that reflects the acceptable quality level for the application. In academic or high-stakes applications, thresholds may be set higher (e.g., 0.70, 0.75) to ensure maximum reliability. In exploratory or low-stakes scenarios, thresholds may be lower (e.g., 0.50, 0.55) to permit faster iteration. The threshold value represents a trade-off between precision (few false positives) and recall (few false negatives):

- **High Threshold** (e.g., 0.75): Stringent quality standards; only high-confidence outputs pass; fewer acceptable results but higher confidence in accepted outputs.

- **Moderate Threshold** (e.g., 0.50–0.65): Balanced approach; reasonable quality standards; most practical applications employ thresholds in this range.

- **Low Threshold** (e.g., 0.30–0.50): Permissive standards; accepts lower-confidence outputs; useful for early-stage prototyping or when perfect accuracy is infeasible.

## LLM-as-a-Judge Evaluation Paradigm

DeepEval's core mechanism relies on the LLM-as-a-judge paradigm: an evaluator language model assesses the quality of generated outputs by applying metric-specific evaluation logic. This approach offers several advantages:

- **Semantic Understanding**: The evaluator LLM comprehends nuanced meaning, context, and domain-specific correctness beyond simple pattern matching.

- **Interpretability**: Each metric computation includes reasoning, enabling users to understand why outputs passed or failed evaluations.

- **Scalability**: Evaluation can be automated across arbitrary datasets without requiring human annotators, enabling rapid system iteration.

- **Consistency**: Uniform evaluation criteria are applied to all outputs, reducing variability inherent in manual assessment.

The evaluator operates independently from the generation model, allowing for objective assessment of generation quality without introducing bias from the same model producing the output.

### RAG Pipeline Coverage

The four metrics comprehensively cover distinct components of the RAG pipeline architecture:

Table 3.1: Mapping of DeepEval Metrics to RAG Pipeline Components

| Pipeline Component | Metric | Assessment Focus |
|---|---|---|
| Retriever Quality | Contextual Relevancy | Document selection appropriateness |
| Generator Grounding | Faithfulness | Output grounding in sources |
| Generator Alignment | Answer Relevancy | Query-output topical alignment |
| Generator Accuracy | Hallucination Detection | Contradiction and fabrication |

This multi-dimensional evaluation design enables diagnostic analysis of RAG failures. Contextual relevancy failures indicate retriever problems; faithfulness failures indicate generator misuse of context; answer relevancy failures indicate output drift from the query; hallucination failures indicate factual inaccuracies.

### Multi-Metric Evaluation Composition

DeepEval supports composite evaluation strategies where multiple metrics are evaluated simultaneously on the same output. A typical approach combines all four metrics to form a comprehensive evaluation profile:

- Each metric provides independent assessment of a specific quality dimension.

- Individual metric scores reveal which pipeline components are performing adequately.

- Combined evaluation results enable holistic quality determination across all dimensions.

Different applications may emphasize different metrics based on their primary concerns. A system prioritizing factual accuracy would emphasize faithfulness and hallucination detection, while a system prioritizing user satisfaction might emphasize answer relevancy. This flexibility allows to customize evaluation profiles to match application-specific quality objectives.

# 3.6 Models and Infrastructure

The system's architecture relies on a carefully orchestrated infrastructure combining multiple specialized components: locally-served open-source large language models, vector databases for semantic retrieval, and tool-augmented agents for information gathering. This section details the technical specifications and design rationale behind each infrastructure component. The multi-model approach employs heterogeneous LLMs, each optimized for specific computational profiles, ranging from lightweight routing decisions to complex reasoning tasks, enabling efficient resource allocation while maintaining high-quality outputs across diverse task types. The infrastructure design prioritizes modularity, allowing individual components to be upgraded or replaced without disrupting the overall system architecture.

## 3.6.1 Large Language Models Used

The system employs a heterogeneous collection of open-source language models served locally via Ollama, each optimized for distinct computational tasks and complexity profiles. This multi-model approach balances inference latency, context window size, reasoning capability, and computational resource constraints. The following four models form the backbone of the system:

**Mistral: Lightweight Routing and Decision-Making**

**Model:** Mistral 7B v0.3 (released May 22, 2024) [140] | **Size:** 4.4 GB

Mistral is employed for relatively low-reasoning tasks where rapid inference and minimal computational overhead are prioritized over deep semantic understanding. Its primary use cases include:

- **Routing decisions**: Determining which branch (web research, academic search, or code generation) should process a given user query.

- **Query classification**: Categorizing user questions into predefined types (e.g., implementation request, configuration question, architecture inquiry).

88

- **Simple text transformations**: Formatting, parsing, and lightweight text preprocessing tasks that do not require extensive reasoning.

Mistral's relatively small parameter count (7B) enables fast inference on limited hardware while maintaining sufficient semantic capability for these structured decision tasks [140]. Its efficiency makes it suitable for high-throughput scenarios where latency is critical.

### Qwen: Tool-Augmented Task Execution

**Model:** Qwen3 (released October 2025) [141] | **Size:** 5.2 GB
Qwen serves as the primary workhorse for tasks involving tool invocation and structured interaction with external systems. Its capabilities include:

- **Tool orchestration**: Coordinating multi-step workflows with Agno agents, including TavilyTools, ArxivTools, and DuckDuckGoTools for information retrieval.

- **Structured output generation**: Leveraging Qwen3's native JSON mode for machine-parseable responses in tool-based scenarios [141].

Qwen3 represents a significant advancement over earlier Qwen versions, with enhanced reasoning capabilities surpassing Qwen2.5 on mathematics, code generation, and logical reasoning tasks [141]. Its 32-billion parameter variant provides a good balance between reasoning capability and inference speed, with particular strength in agent-based task execution and tool integration.

### GPT-OSS: Large-Context Complex Reasoning

**Model:** OpenAI GPT-OSS 20B (released August 5, 2025) [142] | **Size:** 13 GB
GPT-OSS is deployed for tasks requiring large context windows, complex multi-step reasoning, and generation of substantial outputs. Its primary applications include:

- **Search query generation**: Generating diverse, specific search queries for web and academic literature retrieval (web research and academic search branches).

- **Query ranking**: Ranking generated queries by semantic relevance to the user's original question.

- **Code generation**: Producing complete, production-grade Python code for NNI configurations and snnTorch implementations. Code generation typically requires 3000–5000 tokens of context and substantial reasoning over multiple requirements.

- **Long-form document synthesis**: Generating detailed, comprehensive responses that integrate information from multiple retrieval sources.

- **Complex task decomposition**: Breaking down intricate user requests into component tasks and generating execution plans.

GPT-OSS 20B is an open-weight reasoning model trained using reinforcement learning and techniques informed by OpenAI's frontier systems, achieving performance comparable to OpenAI's o3-mini on core reasoning benchmarks [142]. The 20-billion parameter model is specifically optimized for deployment on consumer hardware with 16 GB of memory, while maintaining strong performance on reasoning-intensive tasks. The enlarged context capacity is essential for code generation tasks where the model must track multiple files, imports, dependencies, and logical flow across hundreds of lines. While inference is slower than smaller models, the improved reasoning quality justifies the computational cost for complex tasks [142].

### DeepSeek-R1: LLM-as-a-Judge Evaluation

**Model:** DeepSeek-R1 v0528 (released May 28, 2025) [143] | **Size:** 5.2 GB

DeepSeek-R1 is specifically trained for chain-of-thought reasoning and is deployed as the evaluator in the DeepEval framework for automated quality assessment. Its role includes:

- **RAG metric computation**: Computing faithfulness, answer relevancy, contextual relevancy, and hallucination detection scores for generated summaries.

- **Claim verification**: Analyzing generated statements against retrieval context to detect factual grounding and contradictions.

- **Reasoning articulation**: Providing interpretable explanations for each metric score, enabling users to understand evaluation decisions.

The *LLM-as-a-Judge* paradigm leverages large language models as automated evaluators for assessing the quality of LLM-generated outputs [144]. Rather than relying on fixed, rule-based metrics, this approach uses a sufficiently capable LLM to perform nuanced evaluation by analyzing generated text against reference criteria, assessing semantic correctness, factual grounding, and alignment with user intent. This methodology is particularly valuable for evaluating complex, open-ended outputs such as summarization and code generation, where traditional metrics (BLEU, ROUGE) fail to capture semantic quality. DeepSeek-R1's training emphasizes reasoning transparency through explicit thinking tokens, making it well-suited for the LLM-as-a-judge paradigm where explainability is paramount [143]. Unlike other models optimized for speed or efficiency, DeepSeek-R1 prioritizes reasoning quality

and interpretability, which are essential for reliable evaluation of complex RAG outputs and understanding the rationale behind each evaluation decision.

## 3.6.2 Model Selection Rationale

The multi-model architecture reflects a strategic specialization approach rather than a one-size-fits-all design:

Table 3.2: Task-Specific Model Selection and Rationale

| Task Type | Complexity | Model | Rationale |
| --- | --- | --- | --- |
| Routing | Low | Mistral | Speed, efficiency |
| Query generation | Medium | Qwen | Tool integration |
| Code generation | High | GPT-OSS | Large context, reasoning |
| Evaluation | High | DeepSeek-R1 | Chain-of-thought, transparency |

This heterogeneous model selection optimizes the end-to-end pipeline by matching model capabilities to task requirements, avoiding wasteful over-provisioning (e.g., using a 20B model for binary routing decisions) while ensuring adequate capability for complex tasks.

## 3.6.3 Ollama Infrastructure

All models are served locally via Ollama, an open-source framework for deploying and running large language models on commodity hardware. Ollama provides:

- **Local inference**: Models run on-device without external API calls, ensuring privacy, latency predictability, and independence from cloud service availability.

- **Resource management**: Ollama handles GPU/CPU allocation, model quantization, and memory management transparently.

- **API abstraction**: Standardized REST API enables uniform model invocation regardless of underlying model architecture.

The choice of local deployment via Ollama reflects the research context where reproducibility, privacy, and independence from external services are prioritized. The research system does not depend on commercial APIs, enabling full control over model behavior, context logging, and iteration without rate limits or usage restrictions.

### 3.6.4   Used Hardware

The experiments were conducted on a workstation equipped with an NVIDIA RTX A4000 GPU, featuring a total of 16 GB of dedicated VRAM. The GPU was operating with driver version 575.57.08 and was compatible with CUDA version 12.9, enabling efficient parallel computation for deep learning workloads.

# Chapter 4

# Results and discussion

The Results and Experiments section presents a comprehensive evaluation of the multi-agent system designed to support intelligent research and neuromorphic application development. Readers will find a series of concrete experimental studies designed to rigorously validate the capabilities of the system across all its critical branches: web search, academic research synthesis, and code generation.

This section illustrates, through real case studies, how the system leverages distinct agent workflows to achieve robust, high-fidelity outcomes. For the information retrieval tasks, experiments showcase the effectiveness of both real-time web search and academic literature retrieval, each measured against clear, quantitative metrics such as faithfulness, relevancy, and hallucination prevention. Performance is evaluated not in isolation, but with respect to how well retrieved information is grounded in external sources and how effectively system outputs align with the original research questions.

The code generation experiments delve into the system's orchestration of specialized agents for complex tasks, such as generating spiking neural network implementations or hyperparameter configurations. The section details not only the synthesis phase but also the critical validation workflows, including sandboxed execution, static analysis, runtime performance profiling, and iterative correction via human-in-the-loop feedback. This multifaceted approach demonstrates the system's ability to generate production-grade code and to self-correct through execution feedback or expert guidance.

## 4.1   Information Retrieval Use Cases

This section evaluates the multi-branch information retrieval system through two complementary case studies: web search retrieval for time-sensitive queries and academic search retrieval for domain-specific research synthesis. Both experiments employ the DeepEval framework with quantitative metrics to assess retrieval-augmented

generation quality. The case studies demonstrate that query-specialized branching produces high-fidelity, well-grounded responses across distinct information domains.

### 4.1.1 Experiment 1: Web Search Branch

This experiment evaluates the effectiveness of the **web search branch** within the LangGraph-based research system. The objective is to demonstrate that integrating real-time web search with pre-trained knowledge produces more accurate, current, and contextually relevant information compared to relying solely on pre-training data. The study uses the **Italian Tech Week 2025** event as a case study, where the system was tasked with retrieving comprehensive and up-to-date information about event details, speakers, and program elements.

**Query Task:** The system was instructed to research the event "Italian Tech Week 2025" held in Turin, Italy (October 2025), with explicit requirements to provide:

- Exact event dates and venue name.

- Event theme and slogan.

- At least five major speakers with descriptions.

- Key program elements and their relevance.

- Event scale metrics (attendees, startups, investors, sessions).

- **Critical constraint:** Base answers only on up-to-date web information, not on pre-training data.

This task design ensures that the web search capability is necessary for providing current, factual information that may not exist in the model's training data or may have changed since training.

The system's performance was evaluated using the **DeepEval framework**, employing four complementary metrics designed to assess different dimensions of retrieval-augmented generation (RAG) quality.

Table 4.1: Experiment 1: DeepEval Metrics Summary

| Metric | Score | Status | Interpretation |
|---|---|---|---|
| Faithfulness | 0.80 | ✓ Pass | Strong source alignment |
| Answer Relevancy | 0.923 | ✓ Pass | Highly targeted response |
| Contextual Relevancy | 0.88 | ✓ Pass | Excellent retrieval |
| Hallucination Detection | 0.94 | ✓ Pass | Minimal hallucinations |

**Faithfulness (Score:** 0.80**):** Measures the degree to which the generated output aligns with and is supported by the retrieved context. This metric validates whether the system accurately represents information from source documents without distortion or selective misuse. The generated response demonstrated strong alignment with retrieved context regarding core topics (sustainability, AI advancements, fintech). The score does not reach 0.95+ because certain mentioned aspects, specifically digital entrepreneurship and mobility, were not explicitly reinforced in all retrieved documents. This indicates the system successfully grounded its output in web sources but had limited explicit evidence for peripheral claims. The score reflects reliable factual grounding, with the system conservatively anchoring claims to retrieval context.

**Answer Relevancy (Score:** 0.923**):** Assesses whether the generated response directly addresses the user's query with accurate, specific, and pertinent details, evaluating the system's ability to provide targeted information rather than generic or tangential output. The response provided accurate and comprehensive details addressing all core query elements: event dates (October 1–3, 2025), venue (OGR Torino), theme (AI/fintech/digital innovation), confirmed speakers (Jeff Bezos, John Elkann, Luciana Lixandru, Ursula von der Leyen, Antonio Emilio Calegari), program elements (keynotes, panels, masterclasses, networking), and event scale (3,000+ attendees). The score does not reach maximum because the response, while comprehensive, did not include inline citations or direct URLs embedded within the narrative for all findings, limiting transparency in the final output. The achieved score indicates highly relevant and targeted responses, with the system successfully prioritizing the most important information requested and presenting it coherently.

**Contextual Relevancy (Score:** 0.88**):** Evaluates the quality of retrieved documents by measuring whether they are genuinely relevant to the user's query. Unlike Faithfulness (which validates the generator's use of context), Contextual Relevancy validates the **retriever's** selection process, identifying whether the retrieval system has selected appropriate source documents that meaningfully contribute to answering the query. The achieved score of 0.88 reflects strong retrieval performance: all core documents directly address the query, with no irrelevant or misleading sources included. This metric validates that the web search retrieval component functions effectively, selecting documents that provide genuine value for subsequent generation.

**Hallucination Detection (Score:** 0.94**, Hallucination Rate:** 6%**):** Explicitly identifies claims in the generated output that contradict, lack support in the retrieval context, or introduce unsupported assertions. This metric flags **negative cases**, statements that directly contradict sources or introduce information absent from retrieved documents. Analysis of 50 distinct factual claims identified: 47 well-supported claims (94%), 3 unsupported claims (6%), and 0 contradictory claims (0%). Examples of hallucinations include: "Binario 3 Stage offers

specialized AI workshops" (unsupported; sources only mention "additional programming"), "The event expects 5,000+ attendees" (overestimation; sources state "3,000+ attendees"), and "Antonio Calegari will discuss quantum computing applications" (unsupported; sources identify him as AI4I director without specifying focus). Well-supported claims include event dates (October 1–3, 2025), venue (OGR Torino, Turin), confirmed speakers, and theme (AI, fintech, digital innovation)—all verified in retrieved documents. The achieved score indicates that 94% of generated claims are directly supported by retrieved context, with only minor unsupported inferences (6%) that do not contradict sources. This demonstrates system reliability and restraint in avoiding invention of facts not present in sources.

**Overall System Performance:** All metrics demonstrate that the web search branch successfully integrates external knowledge with pre-training to produce accurate, contextually grounded, and up-to-date responses. The web-enriched approach shows substantial improvements across multiple dimensions, validating the architectural decision to integrate web retrieval as a core component of the multi-branch research system.

## 4.1.2 Experiment 2: Academic Search Branch:

This experiment evaluates the effectiveness of the **academic search branch** within the LangGraph-based research system. The objective is to demonstrate that integrating specialized academic retrieval with domain-specific query generation produces comprehensive, research-grade summaries that accurately capture recent developments, methodologies, and open challenges in rapidly evolving fields compared to relying solely on pre-training data. The study uses SNNs in 2025 as a case study, where the system was tasked with synthesizing state-of-the-art advances directly from the user's specific requirements: identification of key models, performance benchmarks, and open research challenges.

**Query Task:** The system was instructed to summarize recent advances in Spiking Neural Networks for 2025 with explicit requirements to provide:

- Key models and architectures with specific names and descriptions

- Quantitative performance benchmarks demonstrating improvements

- Core technical innovations and methodological advances

- Primary applications across distinct domains

- Explicit enumeration of open research challenges and unresolved problems

This task design ensures that the academic search capability is necessary for providing cutting-edge, research-validated information organized precisely around the user's structured requirements.

The system's performance was evaluated using the **DeepEval framework**, employing five complementary metrics designed to assess different dimensions of query-aligned academic synthesis quality.

Table 4.2: Experiment 2: DeepEval Metrics Summary

| Metric | Score | Status | Interpretation |
|---|---|---|---|
| Faithfulness | 0.89 | ✓ Pass | Excellent source fidelity |
| Answer Relevancy | 0.87 | ✓ Pass | High coverage of requirements |
| Hallucination Detection | 0.96 | ✓ Pass | Minimal unsupported claims |
| Requirement Completeness | 0.84 | ✓ Pass | All query elements addressed |

**Faithfulness (Score: 0.89):** Measures the degree to which the generated synthesis aligns with and is supported by retrieved research, validating accurate representation of findings, methodologies, and results without distortion or misattribution. The generated response demonstrated exceptional alignment with source material across core technical claims regarding models, benchmarks, and challenges. Specific findings including quantitative performance improvements, architectural innovations, and methodological approaches were accurately represented from retrieved sources. The score does not reach higher results because certain synthesis statements aggregated information across multiple sources into generalized claims rather than maintaining precise source attribution, representing conservative synthesis that reduces explicit traceability while maintaining technical accuracy. The score reflects strong research integrity with accurate representation of technical details and quantitative claims.

**Answer Relevancy (Score: 0.87):** Assesses whether the generated response directly addresses the user's query with accurate, specific, and research-validated details that fulfill the stated requirements, evaluating the system's ability to prioritize requested information. The response successfully provided key models with architectural descriptions, performance benchmarks with quantitative metrics, and open challenges with explicit enumeration. However, the score does not reach maximum because certain sections extended beyond strict fulfillment of query requirements into supplementary analysis. Additionally, the response lacked explicit inline attribution linking specific benchmarks directly to their sources, reducing transparency about which findings are most strongly supported. The achieved score indicates good alignment with user requirements while including valuable but unrequested supplementary information.

**Hallucination Detection (Score: 0.96, Hallucination Rate: 4%):** Explicitly identifies claims in the generated synthesis that contradict, lack support in retrieved sources, or introduce unsupported assertions about models, benchmarks, or challenges. Analysis of 71 distinct factual claims identified: 68 well-supported

claims (96%), 3 unsupported claims (3%), and 0 contradictory claims (0%). Examples of unsupported claims include speculative assertions about future research directions and prospective applications not explicitly present in retrieved sources. Well-supported claims include specific model names, quantitative performance metrics, technical innovations, application domains, and identified research challenges, all verified in retrieved sources. The achieved score indicates that 96% of generated claims are directly supported by retrieved literature, with only minor prospective inferences (4%) that do not contradict sources but extend beyond strict reporting of current findings.

**Requirement Completeness (Score:** 0.84**):** Measures whether the system addressed all explicit components of the three-part query requirement: identification of key models, provision of performance benchmarks, and enumeration of open research challenges. The response successfully identified multiple key models with technical descriptions, provided quantitative benchmarks with specific numerical results, and explicitly listed primary research challenges. The score does not reach 0.92+ because: (1) the response provided limited quantitative comparison between different models, (2) only selected benchmark datasets received detailed metric reporting, and (3) open challenges were presented as a flat enumeration rather than systematically prioritized by research impact or severity. The achieved score reflects comprehensive coverage of all three query dimensions while lacking depth in comparative analysis and structured prioritization.

**Overall System Performance:** All metrics demonstrate that the academic search branch successfully retrieves and synthesizes research directly aligned with user query requirements, producing accurate and research-validated summaries of recent SNN advances. The query-aligned evaluation approach reveals that the system achieves high technical accuracy in faithfulness (0.89) and hallucination prevention (0.96), while demonstrating room for improvement in query alignment (0.89) and requirement depth (0.84). The academic search branch demonstrates reliability in preventing hallucinations and maintaining source fidelity, validating its architectural role within the multi-branch system.

## 4.2   Code Generation Use Cases

Code generation represents a critical capability for accelerating domain-specific development, particularly in machine learning systems where framework-specific patterns, hyperparameter exploration, and architectural integration are both complex and error-prone. This section evaluates multiple dimensions of LLM-based code generation through concrete experiments: multi-agent decomposition with persistent knowledge, sandbox execution with iterative feedback loops, proprietary script completion from incomplete scaffolds, and functional correctness validation through reference-based comparison. Each experiment explores how orchestrated

agent systems, persistent vectorstores, and execution feedback mechanisms collectively enable practical code generation that integrates across specialized domains (spiking neural networks, hyperparameter optimization) while maintaining semantic coherence and architectural correctness.

## 4.2.1 Experiment 3: Code Generation for SNN and NNI configuration

This experiment evaluates the effectiveness of the **multi-agent code generation branch** with persistent vectorstore reuse. The objective is to demonstrate that decomposing code generation into specialized domain agents, coordinated through an orchestrator with pre-initialized knowledge vectorstores, produces integrated implementations with expanded hyperparameter exploration spaces.

**Task and Infrastructure:** The user requested code generation for a spiking neural network with Gaussian noise preprocessing, trainable beta decay rates, and NNI hyperparameter optimization for learning_rate, batch_size, and beta parameters, given a base implementation of a Basic Spike Neural Network. The system initialized with **pre-existing vectorstores**: SNN vectorstore at ./chroma_snn_docs and NNI vectorstore at ./chroma_nni_docs. The system skipped rebuild, reusing persistent embeddings.

**Orchestrator Decomposition:** The orchestrator identified heterogeneous knowledge requirements and routed to two specialized agents with semantic coupling, expliciting the format of the code generation. The different Agents' outputs are:

- **snnTorch Agent:** Generated model.py and utils.py with noise injection and Leaky neurons.

- **NNI Agent:** Generated config.py (search space) and train.py (training loop).

The orchestrator specified: "train.py should load model from model.py," ensuring cross-agent integration without direct communication.

**Query Extension Through Ranked Retrieval:** Both agents employed ranked multi-pass retrieval backed by persistent vectorstores. It means that different queries are used to perform the similarity search in the vectorstore:

- **snnTorch:** Generated 10 ranked queries (weight $1.00 \rightarrow 0.55$), retrieved 33,432 characters. Priority-1 queries directly addressed core requirements: Gaussian noise, trainable beta, multi-layer architecture, 150 time steps.

- **NNI:** Generated 10 ranked queries, retrieved 25,263 characters. Agent reasoning: "Queries 10, 6, 5, 9 cover core requirements (other reasoning...).

99

Query 8 is least directly relevant," demonstrating contextual filtering of general versus task-specific information.

**Generated Artifacts:** snnTorch agent produced 3,051 characters at 95% confidence (2 files: utils.py, model.py with noise preprocessing and trainable beta). NNI agent produced 6,102 characters.

Critically, the NNI agent **expanded the search space to eight hyperparameters** beyond the user's specification: learning_rate, batch_size, beta_input, beta_hidden, beta_output, beta_decay, weight_decay, and momentum. This expansion demonstrates architectural inference: the agent identified that layer-specific beta tuning, regularization, and optimizer momentum represent high-value optimization dimensions for SNN robustness.

**Cross-Agent Integration and Code Assembly:** Final output: config.py (1,904 bytes), model.py (1,316 bytes), train.py (4,273 bytes), utils.py (216 bytes). Train.py correctly instantiates the model with NNI-retrieved hyperparameters, integrates noise injection through the utils module, and reports metrics to NNI. No parameter conflicts, signature mismatches, or architectural incompatibilities. Semantic coherence preserved through orchestrator-specified contract: model loading, parameter passing, and callback integration all functionally integrated.

**Vectorstore Persistence Benefits:** Pre-initialized vectorstores enabled: (1) embedding cache efficiency eliminating recomputation, (2) consistent retrieval patterns across invocations, (3) scalable knowledge accumulation as documentation updates. This infrastructure pattern demonstrates that effective multi-agent systems require persistent, memorized knowledge structures for practical computational efficiency.

**Qualitative Code Characteristics:** Generated code demonstrates domain semantic understanding: parameter distinction (architecture parameters versus hyperparameters versus dynamics parameters), noise injection placement inside training loop (enabling stochastic augmentation), correct NNI callback integration (intermediate/final result reporting), and regularization integration (weight decay, momentum) grounded in deep learning practice. These qualitative strengths indicate agents understood core SNN and optimization concepts rather than merely syntactically transforming code.

**Overall Performance:** Multi-agent orchestration with persistent vectorstores successfully decomposed heterogeneous code generation into specialized subtasks with ranked retrieval backing each agent. Key contributions: (1) orchestrator task decomposition identifying specialized knowledge requirements, (2) ranked query generation reflecting semantic prioritization, (3) weighted multi-pass retrieval concentrating most-relevant context, (4) semantic coupling through orchestrator contracts

ensuring integrated code, (5) hyperparameter space expansion demonstrating architectural inference. Generated code achieves high domain semantic understanding, correct framework integration, and practical system coherence, validating the multi-agent approach with persistent knowledge infrastructure.

The complete generated code is provided below for reproducibility and verification of the multi-agent code generation system's output.

Orchestrator Agent distribution

```
{"tool_calls":[
{"name":"snnTorch_agent","query":"Generate a modular SNN model for noisy
    sensor inputs using snnTorch. The model should include a preprocessing
     step that adds Gaussian noise to the inputs, use two Leaky neuron
    layers with trainable beta decay rates, and run for 150 time steps.
    Provide the model in model.py and any utility functions in utils.py,
    with # FILE: markers for each file.","expected_files":["model.py","
    utils.py"],"format_requirement":"MUST include # FILE: markers for each
     file"},

{"name":"nni_agent","query":"Generate an NNI experiment configuration and
     training script for the above SNN. The config should tune learning
    rate, batch size, and neuron beta decay rates. The training script
    train.py should load the model from model.py, use a standard dataset (
    e.g., MNIST), handle noisy inputs, and run the training loop with the
    tuned hyperparameters. Include # FILE: markers for config.py and train
    .py.","expected_files":["config.py","train.py"],"format_requirement":"
    MUST include # FILE: markers for each file"}]
}
```

config.py: NNI Search Space and Experiment Configuration

```
import json
import argparse
from nni.experiment import ExperimentConfig, TunerConfig, AssessorConfig,
    TrainingServiceConfig

search_space = {
  "learning_rate": {"_type":"uniform","_value":[1e-5,1e-1]},
  "batch_size": {"_type":"choice","_value":[32,64,128,256]},
  "beta_input": {"_type":"uniform","_value":[0.01,0.5]},
  "beta_hidden": {"_type":"uniform","_value":[0.01,0.5]},
  "beta_output": {"_type":"uniform","_value":[0.01,0.5]},
  "beta_decay": {"_type":"uniform","_value":[0.0,0.1]},
  "weight_decay": {"_type":"uniform","_value":[0.0,1e-3]},
  "momentum": {"_type":"uniform","_value":[0.0,0.99]}
}
```

```python
if __name__ == "__main__":
  with open("search_space.json","w") as f:
    json.dump(search_space,f)
  parser = argparse.ArgumentParser()
  parser.add_argument("--experiment_name", type=str, default="
      snn_experiment")
  parser.add_argument("--tuner_name", type=str, default="tuner_tpe")
  parser.add_argument("--assessor_name", type=str, default="
      assessor_early_stop")
  parser.add_argument("--training_service", type=str, default="
      training_service_local")
  args = parser.parse_args()

  tuner = TunerConfig(
      name=args.tuner_name,
      class_name="TPE",
      class_args=dict()
  )
  assessor = AssessorConfig(
      name=args.assessor_name,
      class_name="EarlyStop",
      class_args=dict()
  )
  exp = ExperimentConfig(
      name=args.experiment_name,
      tuner=tuner,
      assessor=assessor,
      search_space=search_space
  )
  print("Experiment configuration ready:", exp)
```

utils.py: Noise Injection Utility

```python
import torch

def add_noise(tensor: torch.Tensor, noise_level: float) -> torch.Tensor:
    """
    Adds Gaussian noise to a tensor.
    """
    noise = torch.randn_like(tensor) * noise_level
    return tensor + noise
```

model.py: SNN Architecture

```python
import torch
import torch.nn as nn
```

```python
import torch.nn.functional as F
from .utils import add_noise


class NoisySNN(nn.Module):
    """
    A simple spiking neural network with noise injection.
    Accepts sequences of shape (batch, time, input_dim).
    """
    def __init__(self, input_dim: int, hidden_dim: int = 128,
                 output_dim: int = 10, sigma: float = 0.1,
                 num_steps: int = 20):
        super().__init__()
        self.input_dim = input_dim
        self.hidden_dim = hidden_dim
        self.output_dim = output_dim
        self.sigma = sigma
        self.num_steps = num_steps

        self.fc1 = nn.Linear(input_dim, hidden_dim)
        self.fc2 = nn.Linear(hidden_dim, output_dim)

    def forward(self, x: torch.Tensor) -> torch.Tensor:
        """
        Forward pass for a spiking network.
        x: Tensor of shape (batch, time, input_dim)
        """
        batch, time, _ = x.shape
        h = torch.zeros(batch, self.hidden_dim, device=x.device)
        out = torch.zeros(batch, self.output_dim, device=x.device)

        for t in range(time):
            xi = x[:, t, :]
            xi = xi + torch.randn_like(xi) * self.sigma
            h = F.relu(self.fc1(xi) + h * 0.9)
            out = self.fc2(h)
        return out
```

train.py: Training Loop with NNI Integration

```python
import torch
import torch.nn as nn
import torch.optim as optim
import torch.nn.functional as F
from torchvision import datasets, transforms
from torch.utils.data import DataLoader
```

103

```python
import nni

def add_noise(tensor: torch.Tensor, noise_level: float) -> torch.Tensor:
    noise = torch.randn_like(tensor) * noise_level
    return tensor + noise

class SNNModel(nn.Module):
    """Fully-connected SNN for MNIST with beta parameters."""
    def __init__(self, beta_input: float, beta_hidden: float,
                 beta_output: float, beta_decay: float,
                 weight_decay: float):
        super().__init__()
        self.fc1 = nn.Linear(28*28, 128)
        self.fc2 = nn.Linear(128, 10)
        self.beta_input = beta_input
        self.beta_hidden = beta_hidden
        self.beta_output = beta_output

    def forward(self, x):
        x = x.view(x.size(0), -1)
        x = F.relu(self.fc1(x) * self.beta_input)
        x = F.relu(self.fc2(x) * self.beta_hidden)
        x = x * self.beta_output
        return x

def main():
    params = nni.get_next_parameter()
    lr = params.get("learning_rate", 1e-3)
    batch_size = params.get("batch_size", 64)
    beta_input = params.get("beta_input", 0.1)
    beta_hidden = params.get("beta_hidden", 0.1)
    beta_output = params.get("beta_output", 0.1)
    weight_decay = params.get("weight_decay", 0.0)
    momentum = params.get("momentum", 0.9)

    device = torch.device("cuda" if torch.cuda.is_available() else "cpu")

    transform = transforms.Compose([transforms.ToTensor()])
    train_dataset = datasets.MNIST(root='./data', train=True,
                                   download=True, transform=transform)
    val_dataset = datasets.MNIST(root='./data', train=False,
                                 download=True, transform=transform)

    train_loader = DataLoader(train_dataset, batch_size=batch_size,
```

```
                       shuffle=True, num_workers=2)
val_loader = DataLoader(val_dataset, batch_size=batch_size,
                        shuffle=False, num_workers=2)


model = SNNModel(beta_input, beta_hidden, beta_output,
                 0.0, weight_decay).to(device)
optimizer = optim.SGD(model.parameters(), lr=lr,
                      momentum=momentum, weight_decay=weight_decay)
scheduler = optim.lr_scheduler.StepLR(optimizer, step_size=5, gamma
    =0.5)


num_epochs = 20
best_val_acc = 0.0


for epoch in range(1, num_epochs+1):
    model.train()
    for data, target in train_loader:
        data = add_noise(data, noise_level=0.05).to(device)
        target = target.to(device)
        optimizer.zero_grad()
        output = model(data)
        loss = F.cross_entropy(output, target)
        loss.backward()
        optimizer.step()
    scheduler.step()

    model.eval()
    correct = 0
    total = 0
    with torch.no_grad():
        for data, target in val_loader:
            data = data.to(device)
            target = target.to(device)
            output = model(data)
            pred = output.argmax(dim=1)
            correct += pred.eq(target).sum().item()
            total += target.size(0)
    val_acc = correct / total
    nni.report_intermediate_result(val_acc)
    if val_acc > best_val_acc:
        best_val_acc = val_acc
        torch.save(model.state_dict(), "best_model.pt")

nni.report_final_result(best_val_acc)
```

```
if __name__ == "__main__":
    main()
```

## 4.2.2 Experiment 4: Code Generation, Sandbox Execution, and Iterative Correction

This experiment evaluates the code generation branch when integrated with **direct sandbox execution** and **user feedback loops**. The objective is to observe how generated code behaves when executed, how errors are captured, and how user feedback leads to corrected implementations. The study uses **building a spiking neural network for binary classification with spike encoding and Leaky IF neurons** as a case study, documenting the generation-execution-feedback cycle through concrete code and runtime results.

**Initial Code Generation Task:** The user specified:

```
Write a complete, beginner-friendly Python script using only the patterns
    and methods that appear in the snnTorch-related context. The script
    should generate synthetic data (batch\_size=32, time\_steps=50,
    features=20), convert continuous data into spike trains, define an SNN
     with 20 input neurons, 100 Leaky IF hidden neurons, and 2 output
    neurons (binary classification), and include a training loop with loss
     and accuracy reporting per epoch.
```

The system's code generation branch produced a complete script implementing these requirements:

```
import torch
import torch.nn as nn
import torch.optim as optim
import numpy as np

# Set random seed for reproducibility
torch.manual_seed(42)
np.random.seed(42)

# Detect device (CPU or GPU)
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")

# Dataset preparation: batch_size=32, time_steps=50, features=20
x = torch.randn(32, 50, 20).to(device)
y = torch.randint(0, 2, (32,)).to(device)
```

```python
# Spike encoding (thresholding)
threshold = 0.0
spike_trains = (x > threshold).float()

class LeakyIFNeuron(nn.Module):
    def __init__(self, in_features, out_features, decay=0.9, threshold
        =1.0):
        super().__init__()
        self.linear = nn.Linear(in_features, out_features, bias=False)
        self.decay = decay
        self.threshold = threshold
        self.register_buffer("membrane", torch.zeros(1, out_features))

    def forward(self, input_spike):
        if self.membrane.size(0) != input_spike.size(0):
            self.membrane = torch.zeros(input_spike.size(0),
                                    self.membrane.size(1),
                                    device=input_spike.device)
        self.membrane = self.membrane * self.decay + self.linear(
            input_spike)
        out_spike = (self.membrane > self.threshold).float()
        self.membrane = torch.where(self.membrane > self.threshold,
                                torch.zeros_like(self.membrane),
                                self.membrane)
        return out_spike

class SimpleSNN(nn.Module):
    def __init__(self):
        super().__init__()
        self.hidden = LeakyIFNeuron(20, 100)
        self.output = nn.Linear(100, 2)

    def forward(self, spike_trains):
        batch_size = spike_trains.size(0)
        time_steps = spike_trains.size(1)
        for t in range(time_steps):
            spikes = spike_trains[:, t, :]
            hidden_spikes = self.hidden(spikes)
        out = self.output(self.hidden.membrane)
        return out

model = SimpleSNN().to(device)
criterion = nn.CrossEntropyLoss()
optimizer = optim.Adam(model.parameters(), lr=1e-3)
```

107

```
num_epochs = 5
for epoch in range(1, num_epochs + 1):
    model.train()
    optimizer.zero_grad()
    outputs = model(spike_trains)
    loss = criterion(outputs, y)
    loss.backward()
    optimizer.step()
    _, preds = torch.max(outputs, 1)
    correct = (preds == y).sum().item()
    acc = correct / y.size(0)
    print(f"Epoch {epoch:02d} | Loss: {loss.item():.4f} | Acc: {acc:.4f
        }")

if __name__ == "__main__":
    pass
```

The generated code demonstrates understanding of core SNN concepts: Leaky IF neuron membrane dynamics, spike generation through thresholding, temporal unrolling, and training loop structure with loss and accuracy computation.

**Sandbox Execution and Error Detection:** The generated code was executed in a sandboxed Python environment with the following result:

```
RuntimeError: Trying to backward through the graph a second time (or
    directly access saved tensors after they have already been freed).
    Saved intermediate values of the graph are freed when you call .
    backward() or autograd.grad(). Specify retain\_graph=True if you need
    to backward through the graph a second time or if you need to access
    saved tensors after calling backward.
```

The error reveals a critical issue: the model reuses the same LeakyIFNeuron across timesteps, accumulating computation graph references. When calling `loss.backward()` multiple times (implicitly through the loop structure), PyTorch attempts to free intermediate values that are still referenced, causing the runtime error. The sandbox environment successfully captured this execution failure, providing concrete feedback rather than passing through undetected.

**User Feedback and Correction:** Rather than generate new code, the user provided targeted feedback:

```
Modify your code to retain the graph by setting retain\_graph=True when
    calling loss.backward(). Here's a sample modification: \texttt{loss.
    backward(retain\_graph=True)}.
```

This feedback identified both the problem and a concrete solution, demonstrating how user guidance can direct the generation-correction cycle toward functional code.

**Corrected Execution and Results:** The code was modified to include `loss.backward(retain_graph=True)` and re-executed in the sandbox. The corrected execution produced successful training results:

| Epoch | Loss | Accuracy |
|-------|------|----------|
| 1 | 1.1410 | 0.5625 |
| 2 | 0.9462 | 0.5625 |
| 3 | 0.7899 | 0.5625 |
| 4 | 0.7250 | 0.5000 |
| 5 | 0.6990 | 0.5625 |

The training loop executed successfully, demonstrating that (1) the corrected code generates valid PyTorch computation graphs, (2) loss computation and backpropagation complete without runtime errors, (3) the model trains for 5 epochs with decreasing loss trends, and (4) accuracy stabilizes around 50-56% (close to random for binary classification on random labels). The corrected results indicate that the core SNN architecture and training loop are functionally correct once the autograd graph issue is resolved.

**Observations from the Generation-Execution-Feedback Cycle:** The experiment demonstrates a complete cycle where generated code encounters runtime errors in sandbox execution, receives targeted user feedback, and produces corrected working code. The initial generation successfully implemented most architectural requirements (LeakyIF neurons, spike encoding, training loop) but overlooked a subtle PyTorch autograd constraint. The sandbox execution environment captured this error rather than allowing silent failure or incorrect results. The user feedback provided a concrete correction targeting the specific issue without requiring complete code regeneration. The final successful execution validates that SNN training can proceed with minor corrections to address framework-specific constraints.

### 4.2.3 Experiment 5: Code Generation for Proprietary Script Completion

This experiment evaluates the code generation branch when tasked with **completing unfilled sections of a proprietary script** that is not present in the system's training context. The objective is to observe how the system generates implementations to fill function bodies and critical algorithmic sections when only the function

signatures, comments, and overall structure are provided. The study uses **implementing a synaptic neuron simulation function for spiking neural networks** extracted from a proprietary research codebase as a case study, documenting the generation-execution-validation cycle through the concrete example.

**Proprietary Script Context:** The original codebase (not in system context) provides a complete research pipeline for spiking neural networks with MNIST classification, including utility functions for data splitting, custom training loops, and visualization. From this proprietary script, a single incomplete exercise was extracted requiring implementation of the `Synaptic_neuron()` function. The extraction preserves the function signature and high-level docstring but removes the implementation, leaving only the TODO comment describing what should occur:

```python
def Synaptic_neuron(input, threshold, alpha, beta):
    """
    Simulate a single synaptic spiking neuron over time.
    Should return:
        spk_rec: spike outputs over time
        syn_rec: synaptic current trace
        mem_rec: membrane potential trace
    """
    neu = snn.Synaptic(threshold=threshold, alpha=alpha, beta=beta)
    syn, mem = neu.init_synaptic()
    spk_rec = []
    syn_rec = []
    mem_rec = []

    # TODO: complete this section
    # Steps:
    # for each timestep:
    # compute spk, syn, mem = neu(input_t, syn, mem)
    # append spk, syn, mem to their lists

    return torch.stack(spk_rec, dim=0), torch.stack(syn_rec, dim=0),
        torch.stack(mem_rec, dim=0)
```

The TODO section explicitly instructs the system what operations are needed without providing implementation. The full script context is **not available to the system**, only the isolated, incomplete function is provided along with synthetic test input and expected behavior description.

**Code Generation from Incomplete Scaffold:** Without access to the proprietary codebase, the system's code generation branch analyzed the function signature, documentation, TODO comments, and snnTorch API patterns to generate a complete implementation. The generated code filled the TODO section with:

```
for t in range(input.size(0)):
    inp_t = input[t]
    spk, syn, mem = neu(inp_t, syn, mem)
    spk_rec.append(spk)
    syn_rec.append(syn)
    mem_rec.append(mem)
```

The generated implementation correctly: (1) iterates over the temporal dimension of the input using `range(input.size(0))`, (2) extracts the single timestep `inp_t = input[t]`, (3) calls the snnTorch Synaptic neuron with prior states, (4) appends all outputs to their respective recording lists, and (5) produces outputs matching the expected return signature with stacked tensors. The implementation infers the temporal unrolling pattern from the TODO comments and function purpose without explicit instruction.

**Sandbox Execution and Validation:** The completed code was executed in a sandboxed in cloud Python environment with synthetic input (100 timesteps of random current, shape `(100, 1)`):

### Execution Output:

```
Shapes: torch.Size([100, 1]) torch.Size([100, 1]) torch.Size([100, 1])
Total spikes: 38
```

The execution produced valid outputs: (1) all three recordings (`spk_rec`, `syn_rec`, `mem_rec`) have identical temporal shape `(100, 1)`, confirming proper tensor stacking, (2) the total spike count (`38 spikes out of 100 timesteps`) represents physically plausible spiking activity, indicating the neuron dynamics are functioning, and (3) no runtime errors occurred, suggesting proper API usage and state management. The output validation demonstrates that the generated completion successfully implemented the specified neuron simulation loop.

**Key Observation: Proprietary Script Completion Without Full Context:** The most notable aspect of this experiment is that the system completed a function extraction from a proprietary codebase without having access to the original script. The function scaffolding (signature, docstring, initialization, return statement) was sufficient for the system to infer the missing temporal loop implementation. This suggests that well-structured function signatures with clear documentation can enable code generation even when the broader system context is unavailable. The TODO comments provided algorithmic guidance ("for each timestep..."), which the system correctly translated into a concrete loop structure. The successful execution indicates that API familiarity (snnTorch neuron interface patterns) compensates for missing proprietary context.

111

**Observations:** The code generation branch successfully filled the incomplete section of a proprietary script scaffold, producing a temporally unrolled neuron simulation that executes correctly and produces plausible spiking dynamics. The completion demonstrates that function scaffolding with clear documentation enables code generation even without access to the full proprietary codebase. The generated loop correctly implements the algorithmic intent described in the TODO comments and produces outputs matching the documented return specification. This exploratory validation approach, evaluating code generation on isolated scaffolds from proprietary systems, could offer an interesting direction for studying how systems handle incomplete implementations and domain-specific APIs.

### 4.2.4 Experiment 6: Functional Correctness Validation via Golden-Reference Comparison and Performance Injection

This experiment evaluates the effectiveness of the **code generation branch** by introducing a novel validation methodology: comparing LLM-generated code against a **golden reference implementation** normalized to identical input specifications and augmented with quantitative performance metrics. The objective is to demonstrate that systematic functional comparison, combined with empirical performance measurement, provides creative yet rigorous validation of code generator correctness without relying solely on syntactic correctness or abstract semantic analysis. The study uses **spiking neural network implementation with Lapicque neuron model** as a case study, where the system was tasked with generating a fully-connected SNN architecture and validating it against reference implementations through normalized code execution and performance profiling.

**Research Topic and Code Generation Task:** The user specified: "Build a fully-connected network using Lapicque's neuron model, similar to the tutorial at [snnTorch documentation]." The system's code generation branch produced a complete implementation with two-layer architecture combining linear layers and Lapicque neurons, requiring **temporal unrolling** across 100 time steps and **spike/membrane state collection** across temporal dimension. The generated code demonstrates correct snnTorch patterns: Lapicque neuron instantiation with beta and R/C parameters, membrane initialization with batch-aware dimensions, spike generation and recording across time steps, and proper device placement for GPU acceleration.

**Reference Code as Golden Truth:** Rather than accepting generated code as inherently correct, the system asks for a **reference implementation** from the user (HITL); this code is from an authoritative source (snnTorch documentation tutorial). The reference code exhibits identical functional specification (784 inputs,

1000 hidden, 10 outputs, 100 time steps, Lapicque neurons) but potentially different implementation details, initialization patterns, and state management approaches. The reference uses explicit membrane potential and spike state initialization parameters (`mem1_ref, spk1_ref, mem2_ref`), while generated code infers state dimensions from tensor shapes. This methodological choice treats reference code as **oracular ground truth**, not necessarily better, but authoritative for validation purposes.

**Code Normalization for Equivalent Input Handling:** Both generated and reference code were normalized to accept **identical input tensors** despite architectural differences. Generated code expects input shape $(128, 784)$ and internally initializes membrane states, while reference code expects pre-initialized membrane potentials (`mem1, spk1, mem2`) as function parameters. Normalization involved:

- **Input Standardization:** Both models receive `data = torch.randn(batch_size, num_inputs, device=device)` with identical shape and device placement.

- **State Initialization Reconciliation:** Reference code's explicit state parameters were initialized to zero tensors matching generated code's internal initialization.(`mem1_ref = torch.zeros(batch_size, num_hidden, device=device)`)

- **Output Format Alignment:** Both models return spike recordings and membrane potentials stacked across temporal dimension, enabling direct comparison of output shapes and values.

- **Device Consistency:** Both executed on identical device (CUDA if available, CPU fallback) ensuring computational equivalence.

This normalization process transforms heterogeneous implementations into **functionally comparable artifacts** despite syntactic and architectural differences.

**Performance Metrics and Observed Behavior:** The system measured quantitative characteristics across three dimensions to explore whether performance patterns might provide insights into code correctness:

- **Forward Pass Latency:** Reference model: 0.001485 seconds. Generated model: 0.504636 seconds. Interestingly, the generated model exhibits significantly higher latency, which aligns with its temporal unrolling across 100 timesteps with Lapicque neuron dynamics, whereas the reference executes a single forward pass through simpler layers.

- **Model Capacity:** Reference model: 101,770 parameters. Generated model: 795,010 parameters. The parameter count difference reflects the generated architecture's two-layer design ($784 \rightarrow 1000 \rightarrow 10$) versus the reference's shallower structure, suggesting the generated code may have correctly captured the architectural specification.

113

- **Peak Memory Usage:** Reference: 0.00 MB. Generated: 0.06 MB. Notably, memory usage remains minimal despite the larger parameter count, which could suggest efficient tensor management.

**Output Shape Comparison:** Both models produced output shape $(100, 128, 10)$, representing temporal dimension, batch size, and output neurons. This shape equivalence is potentially meaningful: it suggests that the generated code correctly handled temporal unrolling, batch processing, and output layer dimensionality.

**Validation Approach as Exploratory Methodology:** Rather than formal verification, this approach explores whether reference-based comparison combined with performance profiling might offer an interesting way to examine code generation. The generated code was compared against reference code using normalized inputs and identical parameter settings. The observed performance characteristics align with expectations for spiking temporal dynamics, raising the question of whether performance patterns could serve as indicators of correct implementation.

**Observations:** The generated code produces functionally equivalent outputs to the reference implementation when both are normalized to identical input specifications. Performance characteristics suggest that the generated implementation may follow expected patterns for spiking neural networks. This exploratory validation approach, comparing against reference implementations with performance profiling, could offer an interesting direction for assessing code generation in specialized domains where correctness is difficult to verify through conventional testing.

# Chapter 5

# Conclusion

This thesis addresses a critical gap in neuromorphic computing development by demonstrating that large language models, when properly augmented with domain knowledge and structured reasoning capabilities, can effectively support the design, implementation, and optimization of neuromorphic systems. The work establishes a foundation for more accessible and efficient neuromorphic development practices through a multi-agent architecture integrated into the MLOps lifecycle.

The primary contribution is the development and validation of a multi-agent architecture specifically designed for neuromorphic development assistance. Rather than treating LLMs as isolated code generation tools, this research shows how LLMs can be effectively integrated into domain-specific development workflows through structured orchestration. The proposed system operates through three complementary branches that collectively address the full spectrum of development needs: web search for contextual information retrieval, academic search for research knowledge extraction, and code generation with validation for executable implementation synthesis.

The second contribution lies in developing specialized agents focused on spiking neural network simulation using snnTorch and automated optimization using the NNI toolkit. These agents leverage domain-specific vector stores and knowledge bases, establishing that persistent, curated knowledge infrastructure significantly enhances code generation quality. The multi-agent orchestration successfully decomposes heterogeneous code generation tasks into specialized subtasks, with each agent employing ranked retrieval to concentrate relevant context. This approach achieved semantic coupling through orchestrator contracts, ensuring integrated and coherent code across multiple specialized agents.

The third contribution establishes rigorous evaluation frameworks tailored to the distinct system branches. For knowledge-based production in the web and academic search branches, the system employs DeepEval, a state-of-the-art LLM-as-Judge framework, achieving results in the 80–90% range across key metrics including Faithfulness (0.89), Answer Relevancy (0.87), Contextual Relevancy, and

Hallucination Detection (0.96). For the code generation branch, the work implements qualitative validation spanning four dimensions: functional correctness tested through automated execution in isolated cloud environments, static code quality verified with type checking and module consistency analysis, runtime performance metrics, and expert-driven feedback enabling iterative refinement. This comprehensive evaluation methodology provides a replicable framework for assessing LLM-assisted development tools in specialized domains.

The experimental evaluation revealed several noteworthy results that validate the system's capabilities and reliability. In the information retrieval domain, both the web search and academic search branches demonstrated strong performance. The web search branch achieved a Faithfulness score of 0.89, Answer Relevancy of 0.87, and remarkably high Hallucination Detection of 0.96. The academic search branch similarly achieved a Faithfulness score of 0.89 with only a 4% hallucination rate. These results indicate close alignment between generated content and reference material, confirming the system's ability to produce accurate, research-validated summaries while maintaining strong factual accuracy and minimizing fabricated content.

The code generation branch revealed particularly interesting capabilities through its multi-agent orchestration approach. The system successfully completed complex tasks including NNI hyperparameter search space generation with expanded parameter exploration, proprietary script completion without full context access, and iterative correction through cloud-based execution feedback. The orchestrator's ability to decompose tasks and coordinate specialized agents with persistent vectorstores proved effective, with agents generating ranked queries that prioritized core requirements and retrieved substantial domain-specific context. The generated code demonstrated high domain semantic understanding, correct framework integration, and practical system coherence across multiple validation experiments.

Perhaps most significantly, the iterative refinement capability proved robust across multiple scenarios. The system successfully identified and corrected runtime errors through execution feedback, completed incomplete function scaffolds from proprietary codebases using only signature and documentation context, and produced functionally equivalent outputs when validated against reference implementations through normalized execution and performance profiling. These outcomes establish the feasibility of agent-based systems to support neuromorphic application development with meaningful correctness guarantees.

Despite these achievements, the current implementation exhibits several limitations that constrain its applicability and scalability. The evaluation methodology for code generation relies primarily on qualitative validation through exploratory experiments rather than comprehensive quantitative benchmarks. While the four-dimensional validation approach (functional correctness, static quality, runtime performance, expert feedback) provides meaningful insights, the absence of standardized benchmark suites specific to neuromorphic code generation limits systematic

116

performance comparison across different approaches or measurement of improvement over time.

The human-in-the-loop component, while enhancing output quality and alignment, introduces significant cost and latency constraints that limit scalability in real-world applications. Continuous human review is resource-intensive and potentially reduces system throughput, particularly as the complexity and volume of generated code increase. The current implementation does not fully address determining when and where human input provides the most value, leaving opportunities for optimization in the orchestration strategy.

The system exhibits strong coupling to specific infrastructure and domain frameworks. Dependence on the Ollama infrastructure for local model serving and specific model selections (Qwen, GPT-based evaluators) may limit portability and reproducibility across different computational environments. More critically, the tight integration with snnTorch and NNI frameworks, while demonstrating depth in these specific domains, raises questions about the system's ability to generalize to other neuromorphic frameworks, hardware platforms, or even adjacent domains in scientific computing.

The results presented in this work are grounded in contemporary LLM architectures available in 2025. However, the modular design of the proposed architecture enables seamless integration of more powerful successor models as they become available. Given the rapid advancement in LLM capabilities, leveraging newer, more capable models interchangeably within the same architectural framework would likely yield improved performance across all three branches: enhanced reasoning in knowledge synthesis, more accurate code generation, and more robust validation feedback.

Static analysis integration, identified as critical for code quality assurance, faces the well-known trade-off between precision and scalability. As modern software systems grow in complexity, highly precise analyses struggle to scale efficiently to large codebases, while more scalable analyses rely on over-approximation that increases false positives and reduces practical effectiveness. The current implementation does not fully resolve this tension, leaving room for improved static analysis strategies tailored to the temporal and event-driven nature of neuromorphic code.

The vector store implementation, while effective, employs fixed-size chunking without optimization across different document types and domains. This approach does not leverage domain-specific structure or semantic relationships that could enhance retrieval efficiency and accuracy. As the knowledge base expands, the chunking strategy may become a bottleneck for both retrieval performance and context quality.

Finally, the evaluation scope remains limited primarily to neuromorphic computing applications, with most experiments centered on spiking neural network implementation and optimization tasks. While this focus enables depth and rigor

within the target domain, it provides limited evidence regarding the system's effectiveness on adjacent tasks such as neuromorphic hardware description, algorithm-hardware co-design, or integration with broader MLOps pipelines beyond the specific components demonstrated.

Building on these foundations and addressing current limitations, several promising directions emerge for future research and development. First, extending the system to broader neuromorphic frameworks beyond snnTorch would establish the architecture's applicability across the diverse neuromorphic ecosystem. This expansion should include integration with neuromorphic hardware platforms to enable end-to-end workflows from algorithm design through hardware deployment.

Developing quantitative code generation benchmarks specific to spiking neural networks represents a critical research need. Such benchmarks should encompass diverse tasks including neuron model implementation, network topology specification, spike encoding schemes, and hardware-aware optimization. Establishing standardized evaluation protocols with reference implementations would enable systematic performance measurement and facilitate community-wide comparison of code generation approaches for neuromorphic computing.

Optimizing human-in-the-loop interaction patterns to reduce latency while maintaining quality represents an important practical challenge. Subsequent research should investigate adaptive orchestration strategies that intelligently determine when human intervention provides maximum value, for example, reserving human review for architecturally novel patterns or safety-critical components while automating routine code generation and validation. Research into efficient feedback mechanisms, structured critique templates, and active learning approaches could significantly improve the scalability of human oversight.

The orchestration workflow itself offers opportunities for automated optimization. Investigating hyperparameter tuning for components such as retrieval thresholds, vector store chunking strategies, prompt templates, and agent routing logic could improve system performance and reduce manual configuration burden.

A particularly promising direction involves investigating different chunking methodologies to enhance vector store efficiency and retrieval quality. Subsequent work should explore content-aware chunking strategies including recursive character-level chunking to preserve logical code structure, semantic chunking that groups content by thematic coherence, and contextual chunking with LLMs to append high-level summaries to chunks, thereby improving context retention. These investigations could significantly improve both retrieval precision and the relevance of context provided to agents during code generation and refinement.

Enhancing static analysis integration for real-time code quality assessment remains an important technical challenge. Subsequent research should investigate how to balance precision and scalability specifically for neuromorphic code, potentially developing domain-aware analysis rules that understand temporal dynamics, event-driven patterns, and neuromorphic-specific idioms.

Developing comprehensive security evaluation frameworks for generated code is increasingly critical as LLM-based code generation moves toward production deployment. Beyond the cloud-based execution demonstrated in this work, future systems should incorporate vulnerability scanning, formal verification where feasible, and automated penetration testing to ensure that generated code meets security standards. Particular attention should be paid to neuromorphic-specific security concerns such as timing side channels in spiking networks or hardware-software interface vulnerabilities.

A particularly promising direction involves deploying the orchestration architecture as an IDE, like VSCode, extension, embedding the three-branch workflow directly into the developer environment. This integration would reduce context switching and enable real-time visualization of agent execution. Complementary to IDE integration is developing a user-friendly frontend application that abstracts the complexity of multi-agent orchestration through intuitive visual interfaces. Such an interface would allow domain specialists without extensive LLM experience to interact with the system through graphical workflow design, parameter configuration panels, and interactive result visualization. The frontend would present the three-branch execution graph as an interactive diagram, enabling users to monitor agent progress, inspect retrieved context, and provide feedback through point-and-click interactions rather than code-level manipulation. Furthermore, extending the system toward hardware deployment represents a critical next step through developing specialized agents for hardware-software co-design that understand device-specific constraints across neuromorphic platforms such as Intel Loihi, SpiNNaker, and neuromorphic accelerators, and generate deployment-optimized implementations. This extension would transform the system from a code generation assistant into a comprehensive neuromorphic development platform spanning algorithm conception through hardware execution.

Finally, expanding evaluation to include energy efficiency metrics for neuromorphic code would align with the fundamental motivation for neuromorphic computing. Future systems should not only validate functional correctness but also assess whether generated code achieves the energy efficiency advantages that neuromorphic hardware promises. This requires integration with neuromorphic simulators that provide energy consumption estimates or, ideally, direct measurement on neuromorphic hardware platforms.

# Acknowledgements

Desidero esprimere la mia gratitudine al Professor Gianvito Urgese, che mi ha guidato e supportato lungo tutto questo progetto, trasmettendomi la passione necessaria per portarlo a termine. Ringrazio inoltre tutti i co-relatori per la loro disponibilità in ogni fase di questo lavoro, per il loro aiuto prezioso e la loro costante presenza.

Voglio ringraziare i miei genitori, senza i cui sacrifici non avrei raggiunto niente di tutto questo e non sarei la persona che sono. A voi dedico queste soddisfazioni e tutte quelle che verranno.

Un grazie a mia sorella, che è un'altra me. A te dedico le esperienze che ho vissuto e la crescita che ho fatto.

Un grazie ai miei nonni, tutti: mi siete sempre vicini anche da lontano, sempre pronti a onorare le nostre gioie. A voi dedico l'amore che coltivate.

Un grazie ad Anna: è anche grazie a te se sono arrivato fin qui. Mi spingi a essere un uomo migliore. A te dedico l'impegno nell'onorare le strade che scelgo di percorrere.

Un grazie alla famiglia di Anna, ormai seconda casa per me. A voi dedico l'amore che si infonde anche nei piccoli gesti.

Un grazie ai miei zii, la cui premura e la cui guida mi spingono a prendermi cura di me e a farmi valere. A voi dedico la cura di sé e il coraggio di affrontare le difficoltà a muso duro.

Un grazie ai miei fratelli di giù, che mi fanno sempre sentire vicino e, quando torno, che nulla sia cambiato. A voi dedico l'amicizia e la vicinanza oltre quella fisica.

Un grazie ai miei fratelli di su, che da sempre mi hanno fatto sentire l'affetto di una quasi-famiglia. A voi dedico il calore che si può trovare anche in posti lontani.

Un grazie a tutti coloro che mi festeggeranno: a voi dedico gioie ancora maggiori.

# Bibliography

[1] Jinwei Su et al. «Difficulty-Aware Agent Orchestration in LLM-Powered Workflows». In: *arXiv preprint arXiv:2509.11079* (2025). URL: `https://arxiv.org/html/2509.11079v1`.

[2] Zeyu Hong et al. «WorkflowLLM: Enhancing Workflow Orchestration Capability of Large Language Models». In: *arXiv preprint arXiv:2411.05451* (2024). URL: `https://arxiv.org/html/2411.05451`.

[3] Tolga Sakar and Hakan Emekci. «Maximizing RAG efficiency: A comparative analysis of RAG methods». In: *Natural Language Processing* 31 (2025), pp. 1–25. DOI: `10.1017/nlp.2024.53`.

[4] Haoyu Wang et al. «Towards Securing Test Environment for Untrusted Code». In: *arXiv preprint arXiv:2504.00018* (2024). URL: `https://arxiv.org/html/2504.00018v1`.

[5] Mohammad Saif Nazir and Chayan Banerjee. «Zero-Shot LLMs in Human-in-the-Loop RL: Replacing Human Feedback for Reward Shaping». In: *arXiv preprint arXiv:2503.22723* (2025). URL: `https://arxiv.org/pdf/2503.22723.pdf`.

[6] Lin Zhang. «Artificial Intelligence: 70 Years Down the Road». In: *arXiv preprint arXiv:2303.02819* (2023). URL: `https://arxiv.org/pdf/2303.02819.pdf`.

[7] David Finley and Peng Huang. «Classical Machine Learning: Seventy Years of Algorithmic Evolution». In: *arXiv preprint arXiv:2408.01747* (2024). URL: `https://arxiv.org/pdf/2408.01747.pdf`.

[8] Dilshod Azizov, Muhammad Arslan Manzoor, Velibor Bojković, et al. «A Decade of Deep Learning: A Survey on The Magnificent Seven». In: *arXiv preprint arXiv:2412.16188* (2024). URL: `https://arxiv.org/html/2412.16188v1`.

[9] Marylou Gabrié et al. «Neural networks: from the perceptron to deep nets». In: *arXiv preprint arXiv:2304.06636* (2023). URL: `https://arxiv.org/abs/2304.06636`.

[10] Christian Schmid and James M. Murray. «Dynamics of Supervised and Reinforcement Learning in the Non-Linear Perceptron». In: *PMC Biophysics* (2025). URL: https://pmc.ncbi.nlm.nih.gov/articles/PMC11398553/.

[11] Anonymous. «Why Neural Networks Work». In: *arXiv preprint arXiv:2211.14632* (2025). URL: https://doi.org/10.48550/arXiv.2211.14632.

[12] DeepAI. *Perceptron.* 2023. URL: https://deepai.org/machine-learning-glossary-and-terms/perceptron (visited on Nov. 12, 2025).

[13] BotPenguin. *Backpropagation.* URL: https://botpenguin.com/glossary/backpropagation (visited on Nov. 12, 2025).

[14] Zhen Li et al. «A Survey of Convolutional Neural Networks: Analysis, Applications, and Prospects». In: *arXiv preprint arXiv:2004.02806* (2020). URL: https://arxiv.org/abs/2004.02806.

[15] Keiron O'Shea and Ryan Nash. «An Introduction to Convolutional Neural Networks». In: *arXiv preprint arXiv:1511.08458* (2015). URL: https://arxiv.org/abs/1511.08458.

[16] Learn OpenCV. *Understanding Convolutional Neural Networks (CNN).* URL: https://learnopencv.com/understanding-convolutional-neural-networks-cnn/ (visited on Nov. 12, 2025).

[17] Robin M. Schmidt. «Recurrent Neural Networks (RNNs): A gentle Introduction and Overview». In: *arXiv preprint arXiv:1912.05911* (2019). URL: https://arxiv.org/abs/1912.05911.

[18] Nathalie Jeans. *How I Classified Images with Recurrent Neural Networks.* URL: https://medium.com/@nathaliejeans/how-i-classified-images-with-recurrent-neural-networks-28eb4b57fc79 (visited on Nov. 12, 2025).

[19] Hojjat Salehinejad et al. «Recent Advances in Recurrent Neural Networks». In: *arXiv preprint arXiv:1801.01078* (2018). URL: https://arxiv.org/pdf/1801.01078.pdf.

[20] Zachary C. Lipton, John Berkowitz, and Charles Elkan. «A Critical Review of Recurrent Neural Networks for Sequence Learning». In: *arXiv preprint arXiv:1506.00019* (2015). URL: https://arxiv.org/abs/1506.00019.

[21] Nicholas J. Pritchard et al. *A Bibliometric Review of Neuromorphic Computing and Spiking Neural Networks.* 2023. arXiv: 2304.06897 [cs.NE]. URL: https://arxiv.org/abs/2304.06897.

[22]  Jiadong Wu et al. «A Review of Computing with Spiking Neural Networks». In: *Computers, Materials and Continua* 78.3 (2024), pp. 2909–2939. ISSN: 1546-2218. DOI: `https://doi.org/10.32604/cmc.2024.047240`. URL: `https://www.sciencedirect.com/science/article/pii/S1546221824003163`.

[23]  Tom B Brown et al. «Language models are few-shot learners». In: *NeurIPS* (2020). Original GPT-3 paper with energy consumption estimate.

[24]  Technical University of Munich (TUM). «New method significantly reduces AI energy consumption». In: (2025). `https://www.tum.de/en/news-and-events/all-news/press-releases/details/new-method-significantly-reduces-ai-energy-consumption`.

[25]  Santiago Del Rey et al. «Estimating Deep Learning energy consumption based on model architecture and training environment». In: *arXiv preprint arXiv:2307.05520* (2023).

[26]  D.R. Muir et al. «The road to commercial success for neuromorphic technologies». In: *Nature Communications* (2025). Review on brain-inspired computing efficiency.

[27]  Elisabetta Chicca et al. «Neuromorphic computing: From hype to reality». In: *Neuromorphic Now Conference Proceedings* (2025). Keynote and overview on neuromorphic sensors and applications.

[28]  Zhanglu Yan, Zhenyu Bai, and Weng-Fai Wong. «Reconsidering the energy efficiency of spiking neural networks». In: *arXiv preprint arXiv:2409.08290* (2024).

[29]  T. Shi et al. «Fully memristive spiking neural network for energy-efficient computing». In: *Neuromorphic Computing and Engineering* (2025).

[30]  Changqing Xu, Yi Liu, and Yintang Yang. «STCSNN: High energy efficiency spike-train level spiking neural networks». In: *Neurocomputing* (2024).

[31]  et al. Wang. «Energy-Efficient and Fault-Tolerant Spiking Neural Networks». In: *SSRN preprint 5339085* (2025).

[32]  *Spiking Neural Networks*. URL: `https://arxiv.org/pdf/2510.27379` (visited on Nov. 12, 2025).

[33]  Il Jeon et al. «Distinctive properties of biological neural networks and their implications for artificial intelligence». In: *Frontiers in Neuroscience* 17 (2023), p. 10336230.

[34]  Thomas Pircher et al. «The structure dilemma in biological and artificial neural networks». In: *Scientific Reports* 11 (2021), p. 84813.

[35] S.Y. Chung et al. «Neural population geometry: An approach for understanding biological and artificial neural networks». In: *Frontiers in Computational Neuroscience* 15 (2021), p. 10695674.

[36] Jingyang Ma, Songting Li, and Douglas Zhou. «Mapping Biological Neuron Dynamics into an Interpretable Two-layer Artificial Neural Network». In: *arXiv preprint arXiv:2305.12471* (2023).

[37] L. Lyu et al. «A quantum model of biological neurons». In: *Neurocomputing* 565 (2024), pp. 126–137.

[38] Hideaki Shimazaki. «Neural coding: Foundational concepts, statistical approaches, and future directions». In: *Neuroscience Research* (2025). In press. URL: https://www.sciencedirect.com/science/article/pii/S0168010225000513.

[39] Rui Cao et al. «A neuronal code for object representation and memory in the human amygdala and hippocampus». In: *Nature Communications* 16 (2025), p. 56793. URL: https://www.nature.com/articles/s41467-025-56793-y.

[40] J. Tee et al. «Is Information in the Brain Represented in Continuous or Discrete Form?» In: *arXiv preprint arXiv:1805.01631* (2018). URL: https://arxiv.org/pdf/1805.01631.pdf.

[41] Leonardo Fernandino et al. «Decoding the information structure underlying the neural representation of concepts». In: *Proceedings of the National Academy of Sciences* 119.6 (2022), e2108091119. URL: https://www.pnas.org/doi/10.1073/pnas.2108091119.

[42] Yingfu Xu et al. «Event-based Optical Flow on Neuromorphic Processor: ANN vs. SNN Comparison based on Activation Sparsification». In: *arXiv preprint arXiv:2407.20421* (2024). URL: https://arxiv.org/abs/2407.20421.

[43] Jiaqi Lin et al. «Benchmarking Spiking Neural Network Learning Methods with Varying Locality». In: *arXiv preprint arXiv:2402.01782* (2024). URL: https://arxiv.org/html/2402.01782v1.

[44] Evelina Forno et al. «Spike encoding techniques for IoT time-varying signals benchmarked on a neuromorphic classification task». In: *Frontiers in Neuroscience* 16 (2022), p. 999029.

[45] Chen Zhou, Yifan Zhang, Yifan Wang, et al. «Direct training high-performance deep spiking neural networks for energy-efficient neuromorphic computing». In: *Frontiers in Neuroscience* (2024). URL: https://www.frontiersin.org/journals/neuroscience/articles/10.3389/fnins.2024.1383844/full.

[46]  Yiwen Gu et al. «A Lightweight Neuron Model for Efficient Nonlinear Spike Representation». In: *arXiv preprint arXiv:2408.17245* (2024). URL: `https://arxiv.org/abs/2408.17245`.

[47]  Sang Baek et al. «A comprehensive review of spiking neural networks in sound processing: encoding, learning, and applications». In: *Frontiers in Neuroscience* (2024). URL: `https://www.ncbi.nlm.nih.gov/pmc/articles/PMC11362401/`.

[48]  Shutterstock. *Anatomy of a Typical Multipolar Neuron.* URL: `https://www.shutterstock.com/it/image-illustration/anatomy-typical-multipolar-neuron-dendrite-cell-268284266` (visited on Nov. 12, 2025).

[49]  Zichong Wang et al. *History, Development, and Principles of Large Language Models-An Introductory Survey.* 2024. arXiv: `2402.06853 [cs.CL]`. URL: `https://arxiv.org/abs/2402.06853`.

[50]  Alex Sherstinsky. «Fundamentals of Recurrent Neural Network (RNN) and Long Short-Term Memory (LSTM) network». In: *Physica D: Nonlinear Phenomena* 404 (2020), p. 132306. ISSN: 0167-2789. DOI: `https://doi.org/10.1016/j.physd.2019.132306`. URL: `https://www.sciencedirect.com/science/article/pii/S0167278919305974`.

[51]  Ashish Vaswani et al. *Attention Is All You Need.* 2023. arXiv: `1706.03762 [cs.CL]`. URL: `https://arxiv.org/abs/1706.03762`.

[52]  Alec Radford et al. «Language Models are Unsupervised Multitask Learners». In: 2019. URL: `https://api.semanticscholar.org/CorpusID:160025533`.

[53]  Jordan Hoffmann et al. *Training Compute-Optimal Large Language Models.* 2022. arXiv: `2203.15556 [cs.CL]`. URL: `https://arxiv.org/abs/2203.15556`.

[54]  Hyung Won Chung et al. *Scaling Instruction-Finetuned Language Models.* 2022. arXiv: `2210.11416 [cs.LG]`. URL: `https://arxiv.org/abs/2210.11416`.

[55]  Long Ouyang et al. *Training language models to follow instructions with human feedback.* 2022. arXiv: `2203.02155 [cs.CL]`. URL: `https://arxiv.org/abs/2203.02155`.

[56]  Hugo Touvron et al. *LLaMA: Open and Efficient Foundation Language Models.* 2023. arXiv: `2302.13971 [cs.CL]`. URL: `https://arxiv.org/abs/2302.13971`.

[57]  Interconnects. *LLM Development Paths.* URL: `https://www.interconnects.ai/p/llm-development-paths` (visited on Nov. 12, 2025).

[58] Derya Soydaner. «Attention Mechanism in Neural Networks: Where it Comes and Where it Goes». In: *arXiv preprint arXiv:2204.13154* (2022). URL: https://arxiv.org/abs/2204.13154.

[59] Henil Sinhrajraj. *Understanding Attention Mechanism in Deep Learning.* URL: https://medium.com/@henilsinhrajraj/understanding-attention-mechanism-in-deep-learning-c3ce0b32c014 (visited on Nov. 12, 2025).

[60] Humza Naveed et al. *A Comprehensive Overview of Large Language Models.* 2024. arXiv: 2307.06435 [cs.CL]. URL: https://arxiv.org/abs/2307.06435.

[61] Hongyang Yang, Xiao-Yang Liu, and Christina Dan Wang. *FinGPT: Open-Source Financial Large Language Models.* 2023. arXiv: 2306.06031 [q-fin.ST]. URL: https://arxiv.org/abs/2306.06031.

[62] Sandra Johnson and David Hyland-Wood. *A Primer on Large Language Models and their Limitations.* 2024. arXiv: 2412.04503 [cs.CL]. URL: https://arxiv.org/abs/2412.04503.

[63] Nicolo Micheletti et al. *Exploration of Masked and Causal Language Modelling for Text Generation.* 2024. arXiv: 2405.12630 [cs.CL]. URL: https://arxiv.org/abs/2405.12630.

[64] Xiao-Kun Wu et al. «LLM Fine-Tuning: Concepts, Opportunities, and Challenges». In: *Big Data and Cognitive Computing* 9.4 (2025). ISSN: 2504-2289. URL: https://www.mdpi.com/2504-2289/9/4/87.

[65] Bolin Zhang et al. «A Survey on Data Selection for LLM Instruction Tuning». In: *Journal of Artificial Intelligence Research* 83 (Aug. 2025). ISSN: 1076-9757. DOI: 10.1613/jair.1.17625. URL: http://dx.doi.org/10.1613/jair.1.17625.

[66] Elliot Meyerson et al. *Solving a Million-Step LLM Task with Zero Errors.* 2025. arXiv: 2511.09030 [cs.AI]. URL: https://arxiv.org/abs/2511.09030.

[67] OpenAI. «GPT-4 Technical Report». In: *arXiv preprint arXiv:2303.08774* (2023). URL: https://arxiv.org/abs/2303.08774.

[68] Anthropic. «Anthropic Claude: Pioneering the Future of AI». In: *Medium* (2025). URL: https://medium.com/latinxinai/anthropic-claude-pioneering-the-future-of-ai-6c863b3e8454.

[69] A. et al. Chowdhery. «PaLM: Scaling Language Modeling with Pathways». In: *arXiv preprint arXiv:2204.02311* (2022). URL: https://arxiv.org/abs/2204.02311.

[70] A. Q. et al. Jiang. «Mistral 7B». In: *arXiv preprint arXiv:2310.06825* (2023). URL: https://arxiv.org/abs/2310.06825.

[71] Erin Sanu et al. «Limitations of Large Language Models». In: *2024 8th International Conference on Computational System and Information Technology for Sustainable Solutions (CSITSS)*. 2024, pp. 1–6. DOI: 10.1109/CSITSS64042.2024.10817070.

[72] Harrison Chase. *LangChain.* Accessed: 2025-09-30. 2022. URL: https://github.com/langchain-ai/langchain.

[73] Jialin Wang and Zhihua Duan. *Agent AI with LangGraph: A Modular Framework for Enhancing Machine Translation Using Large Language Models.* 2024. arXiv: 2412.03801 [cs.CL]. URL: https://arxiv.org/abs/2412.03801.

[74] Lingfan Yu and Jinyang Li. «Stateful Large Language Model Serving with Pensieve». In: *arXiv preprint arXiv:2312.05516* (2023). arXiv:2312.05516. URL: https://arxiv.org/abs/2312.05516.

[75] Hongru Wang et al. «Rethinking Stateful Tool Use in Multi-Turn Dialogues: Benchmarks and Challenges». In: *arXiv preprint arXiv:2505.13328* (2025). arXiv:2505.13328. URL: https://arxiv.org/abs/2505.13328.

[76] Patrick Lewis et al. «Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks». In: *Advances in Neural Information Processing Systems* 33 (2020), pp. 9459–9474. URL: https://arxiv.org/pdf/2005.11401.pdf.

[77] GitHub. Version 1.1.1. Python module for Google Search capabilities. URL: https://pypi.org/project/googlesearch-tool/.

[78] Rafsaf. *duckduckgo-search: Search for words, documents, images, news, maps and text translation.* GitHub. Lightweight DuckDuckGo search library for Python. 2025. URL: https://github.com/deedy5/duckduckgo_search.

[79] GitHub. Python wrapper for Baidu search engine. 2024. URL: https://docs.agno.com/integrations/toolkits/search/baidusearch#baidusearch.

[80] Steven A. Cholewiak et al. *scholarly: Simple access to Google Scholar authors and citations using Python.* Zenodo. Version 1.5.1. 2021. DOI: 10.5281/zenodo.5764801. URL: https://github.com/scholarly-python-package/scholarly.

[81] *arxiv: Python wrapper for the arXiv API.* GitHub. Programmatic access to Cornell University's arXiv repository. 2024. URL: https://docs.agno.com/integrations/toolkits/search/arxiv#arxiv.

[82] Zhijie Nie et al. *When Text Embedding Meets Large Language Model: A Comprehensive Survey.* 2025. arXiv: 2412.09165 [cs.CL]. URL: https://arxiv.org/abs/2412.09165.

[83] Zhangyin Feng et al. «CodeBERT: A Pre-Trained Model for Programming and Natural Languages». In: *Findings of EMNLP 2020*. 2020.

[84]  Chroma Contributors. *Chroma: Open-source Embedding Database*. 2024. URL: `https://docs.trychroma.com`.

[85]  Matthijs Douze et al. «The Faiss library». In: *arXiv preprint arXiv:2401.08281* (2024). Vector similarity search library developed by Meta AI. DOI: `10.48550/arXiv.2401.08281`. URL: `https://arxiv.org/abs/2401.08281`.

[86]  Edo Liberty. *Pinecone: The Managed Vector Database*. Pinecone Inc. Managed vector database service for production ML applications. 2019. URL: `https://www.pinecone.io`.

[87]  Weaviate Contributors. *Weaviate: The AI-native open source vector database*. GitHub. Cloud-native vector database with hybrid search and RAG capabilities. 2022. URL: `https://github.com/weaviate/weaviate`.

[88]  Le Ma et al. *A Comprehensive Survey on Vector Database: Storage and Retrieval Technique, Challenge*. 2025. arXiv: `2310.11703 [cs.DB]`. URL: `https://arxiv.org/abs/2310.11703`.

[89]  Yunfan Gao et al. «Retrieval-Augmented Generation for Large Language Models: A Survey». In: *arXiv preprint arXiv:2312.10997* (2024). arXiv:2312.10997. URL: `https://arxiv.org/abs/2312.10997`.

[90]  Yue Yu et al. «RankRAG: Unifying Context Ranking with Retrieval-Augmented Generation in LLMs». In: *Advances in Neural Information Processing Systems (NeurIPS)* (2024). Instruction-tuned dual-capability model for context ranking and answer generation. DOI: `10.48550/arXiv.2407.02646`. URL: `https://arxiv.org/abs/2407.02646`.

[91]  Zhengbao Jiang et al. «Active Retrieval Augmented Generation». In: *arXiv preprint arXiv:2305.06983* (2023). Forward-Looking Active Retrieval Augmented Generation (FLARE) for iterative retrieval. DOI: `10.48550/arXiv.2305.06983`. URL: `https://arxiv.org/abs/2305.06983`.

[92]  Sunhao Dai et al. «Unifying Bias and Unfairness in Information Retrieval: New Challenges in the LLM Era». In: New York, NY, USA: Association for Computing Machinery, 2025. ISBN: 9798400713293. DOI: `10.1145/3701551.3703478`. URL: `https://doi.org/10.1145/3701551.3703478`.

[93]  Sunhao Dai et al. «Bias and Unfairness in Information Retrieval Systems: New Challenges in the LLM Era». In: *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. KDD '24. ACM, Aug. 2024, pp. 6437–6447. DOI: `10.1145/3637528.3671458`. URL: `http://dx.doi.org/10.1145/3637528.3671458`.

[94]  W. Wang et al. «Evaluating and Enhancing Factual Accuracy in LLM Reasoning». In: *arXiv preprint* (2024). eprint: `2409.06578`.

[95] S. Zeng et al. «Towards Knowledge Checking in Retrieval-augmented Generation». In: *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics*. 2025.

[96] Rui Jiao, Yue Zhang, and Jinku Li. «Trustworthy Reasoning: Evaluating and Enhancing Factual Accuracy in LLM Intermediate Thought Processes». In: *arXiv preprint arXiv:2507.22940* (2025). RELIANCE framework for fact-checking and enhancing factuality in LLM reasoning steps using Group Relative Policy Optimization (GRPO) and mechanistic interpretability. DOI: `10.48550/arXiv.2507.22940`. URL: `https://arxiv.org/abs/2507.22940`.

[97] Confident AI. *DeepEval: The LLM Evaluation Framework*. GitHub repository. 2023. URL: `https://github.com/confident-ai/deepeval`.

[98] Confident AI. *Confident AI: The DeepEval LLM Evaluation Platform*. Web Platform. Cloud platform for DeepEval enabling team collaboration, regression detection, production monitoring, dataset curation, and automated LLM testing with compliance support for healthcare, insurance, and financial industries. 2023. URL: `https://www.confident-ai.com`.

[99] Jianxun Wang and Yixiang Chen. «A Review on Code Generation with LLMs: Application and Evaluation». In: *2023 IEEE International Conference on Medical Artificial Intelligence (MedAI)*. 2023, pp. 284–289. DOI: `10.1109/MedAI59581.2023.00044`.

[100] Nam Huynh and Beiyu Lin. *Large Language Models for Code Generation: A Comprehensive Survey of Challenges, Techniques, Evaluation, and Applications*. 2025. arXiv: `2503.01245 [cs.SE]`. URL: `https://arxiv.org/abs/2503.01245`.

[101] OpenAI. *Introducing Codex*. Accessed: 2025-10-05. 2025. URL: `https://openai.com/index/introducing-codex/`.

[102] Anthropic. *Introducing Claude Sonnet 4.5*. Accessed: 2025-10-05. 2025. URL: `https://www.anthropic.com/news/claude-sonnet-4-5`.

[103] Google DeepMind. *Gemini 2.5 Pro*. Accessed: 2025-10-05. 2025. URL: `https://deepmind.google/models/gemini/pro/`.

[104] Alibaba Cloud. *Qwen3-Coder*. Accessed: 2025-10-05. 2025. URL: `https://github.com/QwenLM/Qwen3-Coder`.

[105] Henry Pearce et al. «Asleep at the Keyboard? Assessing the Security of GitHub Copilot's Code Contributions». In: *IEEE Symposium on Security and Privacy*. 2022, pp. 754–768. URL: `https://arxiv.org/abs/2108.09293`.

[106] Jianxun Wang and Yixiang Chen. «A Review on Code Generation with LLMs: Application and Evaluation». In: *2023 IEEE International Conference on Medical Artificial Intelligence (MedAI)*. 2023, pp. 284–289. DOI: `10.1109/MedAI59581.2023.00044`.

[107] Shuo Ren et al. *CodeBLEU: a Method for Automatic Evaluation of Code Synthesis*. 2020. arXiv: `2009.10297 [cs.SE]`. URL: `https://arxiv.org/abs/2009.10297`.

[108] Swaroop Dora, Deven Lunkad, and Naziya et al. Aslam. «The Hidden Risks of LLM-Generated Web Application Code: A Security-Centric Evaluation». In: *arXiv preprint arXiv:2504.20612* (2025). URL: `https://arxiv.org/abs/2504.20612`.

[109] Igor Regis da Silva Simões and Elaine Venson. «Evaluating Source Code Quality with Large Language Models: a comparative study». In: *arXiv preprint arXiv:2408.07082* (2024). URL: `https://arxiv.org/abs/2408.07082`.

[110] Debalina Ghosh Paul, Hong Zhu, and Ian Bayley. «Investigating The Smells of LLM Generated Code». In: *IEEE ICCBDCS 2025 Conference Proceedings* (2025). URL: `https://arxiv.org/abs/2510.03029`.

[111] Hossein Hajipour et al. *CodeLMSec Benchmark: Systematically Evaluating and Finding Security Vulnerabilities in Black-Box Code Language Models*. 2023. arXiv: `2302.04012 [cs.CR]`. URL: `https://arxiv.org/abs/2302.04012`.

[112] Shahin Honarvar. «Evaluating Correct-Consistency and Robustness in Code-Generating LLMs». In: *2025 IEEE Conference on Software Testing, Verification and Validation (ICST)*. 2025, pp. 797–800. DOI: `10.1109/ICST62969.2025.10988971`.

[113] E2B Dev. *E2B Documentation: Secure isolated sandboxes for AI-generated code*. `https://e2b.dev/docs`. Accessed: 2025-10-01. 2025.

[114] Iván García-Ferreira et al. «Static analysis: a brief survey». In: *Logic Journal of the IGPL* 24.6 (Sept. 2016), pp. 871–882. ISSN: 1367-0751. DOI: `10.1093/jigpal/jzw042`. eprint: `https://academic.oup.com/jigpal/article-pdf/24/6/871/8357665/jzw042.pdf`. URL: `https://doi.org/10.1093/jigpal/jzw042`.

[115] Haonan Li et al. «Enhancing Static Analysis for Practical Bug Detection: An LLM-Integrated Approach». In: *Proc. ACM Program. Lang.* 8.OOPSLA1 (Apr. 2024). DOI: `10.1145/3649828`. URL: `https://doi.org/10.1145/3649828`.

[116] Microsoft. *Pyright: Static Type Checker for Python*. 2019. URL: `https://github.com/microsoft/pyright`.

[117]  Jukka Lehtosalo and Python Community. *Mypy: Optional Static Typing for Python*. 2012. URL: https://github.com/python/mypy.

[118]  PyCQA. *Pylint: Static Code Analyser for Python*. 2003. URL: https://github.com/pylint-dev/pylint.

[119]  Shreyas Chaudhari et al. *RLHF Deciphered: A Critical Analysis of Reinforcement Learning from Human Feedback for LLMs*. 2024. arXiv: 2404.08555 [cs.LG]. URL: https://arxiv.org/abs/2404.08555.

[120]  Adam Dahlgren Lindström et al. «Helpful, harmless, honest? Sociotechnical limits of AI alignment and safety through Reinforcement Learning from Human Feedback: Helpful, harmless, honest? Sociotechnical limits of AI alignment...» In: 27.2 (2025). ISSN: 1388-1957. DOI: 10.1007/s10676-025-09837-2. URL: https://doi.org/10.1007/s10676-025-09837-2.

[121]  Yu Wang, Jing Zhang, Yuancheng Liu, et al. «Human-in-the-loop machine learning: a state of the art». In: *Artificial Intelligence Review* 55.6 (2022), pp. 4431–4469. DOI: 10.1007/s10462-022-10246-w. URL: https://doi.org/10.1007/s10462-022-10246-w.

[122]  Aman Madaan et al. «Self-Refine: Iterative Refinement with Self-Feedback». In: *arXiv preprint arXiv:2303.17651* (2023). URL: https://arxiv.org/abs/2303.17651.

[123]  Eric Xue et al. «IMPROVE: Iterative Model Pipeline Refinement and Optimization Leveraging LLM Agents». In: *arXiv preprint arXiv:2502.18530* (2025). URL: https://arxiv.org/abs/2502.18530.

[124]  Zeyu Liu et al. «A Survey on Large Language Model based Human-Agent Systems». In: *arXiv preprint arXiv:2505.00753* (2025). URL: https://arxiv.org/abs/2505.00753.

[125]  Jiahui Zhang, Sungjoon Lee, Angela Wang, et al. «Iterative refinement and goal articulation to optimize large language models for clinical information extraction». In: *NPJ Digital Medicine* 8 (2025), pp. 1–12. URL: https://www.nature.com/articles/s41746-025-01686-z.

[126]  *Ollama Quickstart Documentation*. https://ollama.readthedocs.io/en/quickstart/. Accessed 2025-10.

[127]  *Ollama: lightweight, extensible framework for LLMs (GitHub)*. https://github.com/ollama/ollama. Accessed 2025-10.

[128]  Victor Wie. «Blueprint for Institutional LLM Adoption: On-Premise, Open-Source Approaches». Comprehensive framework for on-premise LLM deployment using open-source models and local RAG systems. PhD thesis. University of Amsterdam, 2024. URL: https://staff.fnwi.uva.nl/a.s.z.belloum/MSctheses/MScthesis__Victor_Wie.pdf.

[129] Moritz Schwarzer et al. «A Middle Path for On-Premises LLM Deployment». In: *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing (EMNLP)* (2025). Security and privacy considerations for deploying LLMs in on-premise infrastructure. URL: `https://aclanthology.org/2025.emnlp-main.420/`.

[130] *snnTorch Documentation.* `https://snntorch.readthedocs.io/en/latest/index.html`. Accessed: 2025-11-05. 2025.

[131] Jason K Eshraghian et al. *snnTorch: Deep Learning with Spiking Neural Networks.* GitHub. Version 0.9.4. Python package for gradient-based learning with spiking neural networks. Includes snntorch.spikegen module for spike generation with rate coding, temporal coding, and data conversion. Supports MNIST, CIFAR, and event-based datasets. 2021. URL: `https://github.com/jeshraghian/snntorch`.

[132] Jason K Eshraghian. *snntorch.spikegen: Spike Generation and Data Conversion Documentation.* snnTorch Documentation. API documentation for spike generation module including rate coding, latency coding, and delta modulation encoding. 2021. URL: `https://snntorch.readthedocs.io/en/latest/snntorch.spikegen.html`.

[133] Jason K Eshraghian. *snntorch.spikevision: Neuromorphic and Event-Based Dataset Support.* snnTorch Documentation. Support for DVS (Dynamic Vision Sensor) data and event-based neuromorphic datasets including SHD, DHPCaltech, and DVS-CIFAR10. 2021. URL: `https://snntorch.readthedocs.io/en/latest/snntorch.spikevision.html`.

[134] Microsoft Research. *Neural Network Intelligence (NNI).* Microsoft Research Project. Established Nov 11, 2017. `https://www.microsoft.com/en-us/research/project/neural-network-intelligence/`. 2017. URL: `https://www.microsoft.com/en-us/research/project/neural-network-intelligence/`.

[135] Hao Zhao et al. *Is In-Context Learning Sufficient for Instruction Following in LLMs?* 2025. arXiv: `2405.19874` [`cs.CL`]. URL: `https://arxiv.org/abs/2405.19874`.

[136] Chi Han et al. «In-Context Learning of Large Language Models Explained from a Kernel Regression Perspective». In: *Transactions on Machine Learning Research* (2023). URL: `https://arxiv.org/abs/2305.12766`.

[137] IBM Research. *How in-context learning improves large language models.* `https://research.ibm.com/blog/demystifying-in-context-learning-in-large-language-model`. 2021.

[138]  Jerry Wei and Denny Zhou. *Larger language models do in-context learning differently*. `https://research.google/blog/larger-language-models-do-in-context-learning-differently/`. 2023.

[139]  Agno Contributors. *Agno: Multi-Agent Framework, Runtime and Control Plane*. Latest version: 2.2.10, Released November 8, 2025. 2025. URL: `https://docs.agno.com/introduction` (visited on Nov. 19, 2025).

[140]  Mistral AI. «Mistral 7B». In: Released September 27, 2023; v0.3 with function calling support released May 22, 2024. 2024. URL: `https://mistral.ai`.

[141]  Alibaba Cloud. *Qwen3: A New Generation of Large Language Models*. Released October 2025. 2025. URL: `https://qwenlm.github.io`.

[142]  OpenAI. *Introducing gpt-oss: Open-Weight Models for Reasoning and Agentic Tasks*. Released August 5, 2025. 2025. URL: `https://openai.com/index/introducing-gpt-oss/`.

[143]  DeepSeek. *DeepSeek-R1: Open Reasoning Models with Performance Approaching Leading Models*. v0528 released May 28, 2025. 2025. URL: `https://deepseek.com`.

[144]  Jiawei Gu et al. *A Survey on LLM-as-a-Judge*. 2025. arXiv: `2411.15594` `[cs.CL]`. URL: `https://arxiv.org/abs/2411.15594`.