

# Modular Machine Learning: An Indispensable Path towards New-Generation Large Language Models

Xin Wang, *Member, IEEE*, Haoyang Li, Zeyang Zhang, Haibo Chen and Wenwu Zhu, *Fellow, IEEE*,

**Abstract**—Large language models (LLMs) have dramatically advanced machine learning research including natural language processing, computer vision, data mining, etc., yet they still exhibit critical limitations in reasoning, factual consistency, and interpretability. In this paper, we introduce a novel learning paradigm—Modular Machine Learning (MML)—as an essential approach toward new-generation LLMs. MML decomposes the complex structure of LLMs into three interdependent components: modular representation, modular model, and modular reasoning, aiming to enhance LLMs' capability of counterfactual reasoning, mitigating hallucinations, as well as promoting fairness, safety, and transparency. Specifically, the proposed MML paradigm can: i) clarify the internal working mechanism of LLMs through the disentanglement of semantic components; ii) allow for flexible and task-adaptive model design; iii) enable interpretable and logic-driven decision-making process. We present a feasible implementation of MML-based LLMs via leveraging advanced techniques such as disentangled representation learning, neural architecture search and neuro-symbolic learning. We critically identify key challenges, such as the integration of continuous neural and discrete symbolic processes, joint optimization, and computational scalability, present promising future research directions that deserve further exploration. Ultimately, the integration of the MML paradigm with LLMs has the potential to bridge the gap between statistical (deep) learning and formal (logical) reasoning, thereby paving the way for robust, adaptable, and trustworthy AI systems across a wide range of real-world applications.

**Index Terms**—Large Language Model, Neuro-Symbolic Learning, Disentangled Representation Learning, Neural Architecture Search.



## 1 INTRODUCTION

THE advent of Large Language Models (LLMs) [1], epitomized by the likes of ChatGPT, has undeniably been a watershed moment in the evolution of the AI landscape. These models have astounded the research community and the general public alike with their seemingly superhuman language-processing capabilities. In a plethora of tasks, they have managed to mimic human-level language understanding and generation with remarkable fidelity. From drafting eloquent essays and engaging in fluent conversations to providing detailed summaries of complex texts, LLMs have proven their capabilities across various scenarios. This remarkable performance has led some researchers to advocate the “bigger is better” principle. The underlying rationale is that as the model size and volume of training data grow exponentially (i.e., the scaling law), so does the breadth and depth of knowledge that the LLM can encapsulate. This enables the model to handle a wide range of language-related challenges with an accuracy and fluency that were previously unimaginable.

However, as with any technological marvel, LLMs have their own *Achilles' heels*. A major limitation of LLMs surfaces when it comes to quantitative reasoning tasks. For in-

stance, LLMs often struggle with simple two-digit arithmetic problems. Without human interaction, they have trouble applying basic mathematical principles to obtain the right answers. This deficiency highlights the fact that, despite their remarkable linguistic capabilities, their cognitive and logical reasoning abilities, particularly in domains requiring formal and rule-based understanding, remain in need of substantial enhancement. Consequently, the research community has witnessed a growing trend for augmenting LLMs with additional tools, where the additional tools are designed to fill the gaps in reasoning capabilities. For example, coupling an LLM with a symbolic solver can enhance LLM's ability to handle mathematical and logical problems. By integrating external knowledge bases, LLMs can access domain-specific information that might not have been part of their original training data as well, thereby enhancing their adaptable problem-solving capabilities. Such hybrid systems could potentially release the full potential of LLMs, allowing them to overcome current limitations and become more versatile and reliable in a wide range of AI applications.

LLMs are expected to be explainable, reliable, adaptable, and extendable in real-world applications. For example, in many critical domains such as healthcare, finance, and legal systems, the decisions made by LLMs can have far-reaching consequences. When a model provides a diagnosis or financial advice, patients and stakeholders need to understand the underlying rationale. Without explainability, it becomes difficult to validate the accuracy and reliability of the output. In a medical context, if an LLM suggests a particular treatment plan, both doctors and patients must clearly understand the reasoning behind that recommendation. This not only

• All authors were with the Department of Computer Science and Technology, Tsinghua University, China, 100084. E-mail: {xin\_wang, wwzhu}@tsinghua.edu.cn, lihy218@gmail.com {zy-zhang20, chb24}@mails.tsinghua.edu.cn

This work was supported in part by National Natural Science Foundation of China No. 62222209, Beijing National Research Center for Information Science and Technology (BNRist) under Grant No. BNR2023RC01003, Beijing Key Lab of Networked Multimedia, China.

fosters trust but also enables error detection and correction. Furthermore, as the demands of applications evolve and new tasks emerge, LLMs need to be adaptable to unseen scenarios via progressively incorporating new knowledge, domains, and functions. In the fast-paced world of technology, a static LLM would quickly become obsolete and outdated. Additionally, as new scientific discoveries arise, an extensible LLM should be capable of incorporating these advancements to enhance its capabilities.

In this paper, we propose a novel learning paradigm, i.e., Modular Machine Learning (MML), which includes i) *modular representation*, ii) *modular model* and iii) *modular reasoning* for deep learning architectures such as LLMs, presenting a viable path to attaining the above crucial characteristics required in a wide range of real-world applications. For ease of understanding, we utilize a visual question answering (VQA) task as an example to illustrate the process of MML for LLMs in Fig. 1. To further validate the methodology, in Fig. 2 we provide an instantiation of MML under VQA task with a feasible implementation. 1) Disentangled Representation Learning (DRL) for modular representation allows for a more organized and interpretable structure of information within the LLMs. By disentangling complex representations into independent semantic dimensions, it becomes easier to understand and control the model’s internal working process, thus enhancing explainability. 2) Neural Architecture Search (NAS) for modular model enables the creation of modularized neural architectures with various modules in charge of different functions. Different modules can be optimized with respect to specific tasks, making the overall LLM capable of reliably completing complex tasks with an adaptive combination of different prerequisite tasks. The modularized architecture also facilitates parameter maintenance and updating within different modules, contributing to easy extension with new modules regarding new functions. 3) Neuro-Symbolic Learning (NSL) for modular reasoning bridges the gap between the black-box neural network inference process and human-understandable symbolic reasoning. The modularized design formalizes the inference steps, improving explainability and reliability as the decision-making process becomes more transparent and auditable, as well as promoting fairness via enabling the detection and correction of potential biases in the reasoning chain. As such, *modular machine learning holds the potential to transform LLMs into more powerful, trustworthy, and versatile tools for real-world applications*. More discussions regarding the implementation will be presented in Section 4.

We will first present the methodology of MML in Section 2, then discuss the significance of MML for LLMs in Section 3, followed by a detailed description of a feasible implementation of MML in Section 4. Additionally, we will also explore the future directions in Section 5. In sum, we aim to provide a comprehensive understanding of the potential and limitations of current machine learning techniques, with the hope that insights from this paper may contribute to the advancement of Artificial General Intelligence (AGI).

## 2 METHODOLOGY OF MODULAR MACHINE LEARNING (MML)

Modular Machine Learning (MML) is a learning paradigm that decomposes the complex structure of large language models (LLMs) into three interdependent components: modular representation, modular model, and modular reasoning. This decomposition allows for improved reasoning, interpretability, and adaptability of LLMs in various applications.

### 2.1 Mathematical Definition

Let  $\mathcal{X}$  be the input space and  $\mathcal{Y}$  the output space. A Modular Machine Learning model  $M$  can be defined as a composition of three modular functions:

$$\mathcal{M}(x) = \mathcal{M}_{\theta_R} \left( \mathcal{M}_{\theta_M} \left( \mathcal{M}_{\theta_D}(x) \right) \right), \quad (1)$$

where:

- $\mathcal{M}_{\theta_D}$ : Modular representation function that extracts disentangled and semantically meaningful features from the input  $x \in \mathcal{X}$ .
- $\mathcal{M}_{\theta_M}$ : Modular model function that dynamically adapts the model architecture to the input features.
- $\mathcal{M}_{\theta_R}$ : Modular reasoning function that applies symbolic reasoning on the structured representation obtained from the modular model.

### 2.2 Optimization Objective

The training objective for MML can be formulated as:

$$\min_{\theta_D, \theta_M, \theta_R} \mathcal{L}(M(x), y), \quad (2)$$

where  $\theta_D$ ,  $\theta_M$ , and  $\theta_R$  represent the parameters of the modular representation, modular model, and modular reasoning functions respectively, and  $\mathcal{L}$  denotes the loss function for the task. This formal definition highlights how MML integrates modular representation, modular model, and modular reasoning into a unified framework, enabling LLMs to perform complex reasoning with enhanced flexibility and interpretability.

## 3 THE KEY SIGNIFICANCE OF MML FOR LLMs

**Enable Counterfactual Reasoning.** One of the most profound limitations of current LLMs is their inability to effectively reason about counterfactuals [2]. Counterfactual reasoning involves exploring hypothetical scenarios that differ from the observed reality to answer “what if” questions. For example, in medicine, counterfactual reasoning might explore the outcome which assumes a patient had received a different treatment. While LLMs excel at predicting texts based on patterns in large datasets, their statistical foundations make them ill-suited for tasks requiring logical consistency and causal reasoning. MML bridges this gap by integrating modular reasoning systems capable of counterfactual reasoning. These systems rely on explicitly defined rules or causal graphs to model relationships between variables, enabling the generation of logically coherent hypotheses about unobserved scenarios. By non-trivially integrating such modules into LLMs, it is possible to enhance the ability of LLMs to handle tasks requiring robust logical reasoning [3].

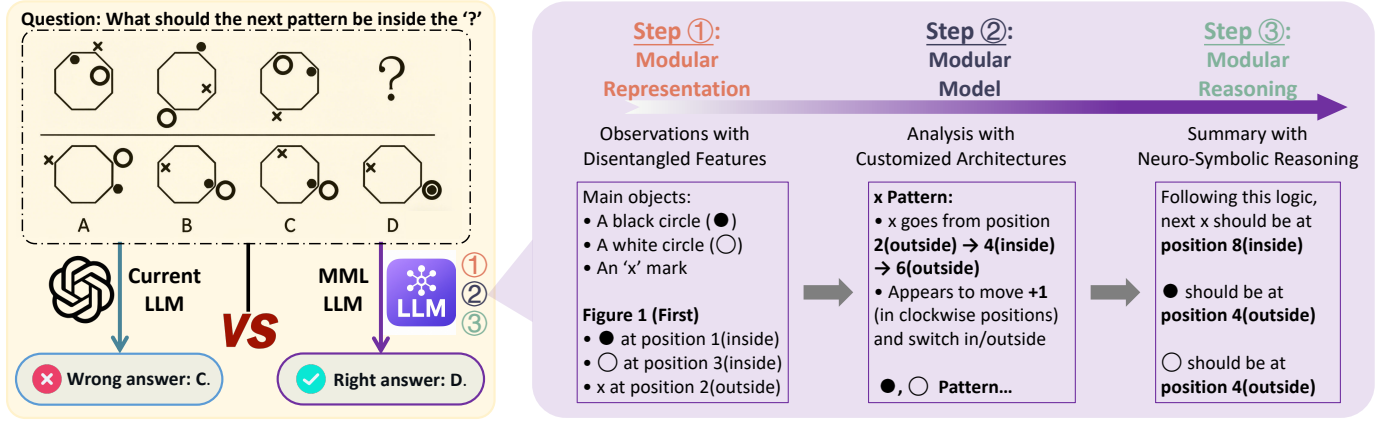


Fig. 1: We employ a visual question answering (VQA) task to compare traditional LLMs with our proposed Modular Machine Learning (MML) framework for LLMs. Current LLMs (e.g., ChatGPT-4o) often fail on tasks requiring complex reasoning. MML is able to overcome the weakness by adopting a modular approach as follows, **Step ①**: Modular Representation first disentangles the visual content to separate and extract key features from image based on the textual question (e.g., identifying the target objects **X**, **Solid Circle** and **Open Circle**); **Step ②**: Modular Model then utilizes customized neural architectures tailored for different functionalities to analyze sequential position patterns of objects (e.g., tracking the sequential positions of the relevant objects **X**, **Solid Circle** and **Open Circle**); **Step ③**: Modular Reasoning finally performs logical reasoning to infer the next pattern and summarize the answer (e.g., predict the positions of **X**, **Solid Circle** and **Open Circle** inside the question mark).

The integration can take various forms. One approach is to use neural modules as perception engines to parse and preprocess complex, high-dimensional data into structured representations compatible with modular systems. These modular systems then perform counterfactual reasoning by manipulating these structured representations according to well-defined logical or causal principles. For example, in policy decision-making, an MML-enhanced LLM could analyze historical data to suggest alternative outcomes under different policy scenarios, offering interpretable insights for potential causal mechanisms. Beyond applications, MML-based frameworks can provide architectural advantages by separating perception from reasoning. This modular design ensures that the reasoning layer operates independently of the statistical biases inherent in neural networks. For instance, while a neural network might overfit correlations from training data, an MML-based module is able to enforce constraints that align with known causal relationships, mitigating the risk of generating spurious counterfactuals. Thus, MML could serve as a foundation for counterfactual reasoning within LLMs [4].

**Conduct Rule-Based Reasoning.** One of the main challenges of LLMs is their reliance on probabilistic text generation, which can lead to factually incorrect outputs [5]. MML enhances LLMs by introducing structured rule-based reasoning mechanisms that ensure factual consistency. Modules within MML frameworks can extract rules from LLM-generated content and validate them against external knowledge bases, ontologies, or predefined constraints. For example, in the medical domain, an LLM could synthesize clinical guidelines derived from scientific literature, while an MML-integrated modular reasoning module can cross-check these recommendations against structured knowledge, identifying inconsistencies or inaccuracies. This approach enhances the reliability of LLM outputs, making them suitable for high-stakes applications where precision is critical.

Additionally, by continuously updating its rule base, MML ensures that LLMs remain aligned with up-to-date domain knowledge, thereby reducing the risk of outdated or biased recommendations.

**Unify Perception and Reasoning.** The integration of MML and LLM into a unified framework bridges the gap between perception and reasoning, offering a scalable and adaptable AI paradigm. By leveraging MML's modular nature, this unified system can efficiently handle complex multimodal data, combining LLM-driven perception with modular reasoning to enable comprehensive decision-making. One of the key benefits of this approach lies in its flexibility. While LLMs can be fine-tuned for domain-specific applications, MML's modular design ensures that different reasoning modules can be updated independently. For example, in an autonomous system, the perceptual module (LLM) can process sensor data and natural language commands, while the reasoning module (MML) ensures compliance with pre-defined safety constraints. This separation enhances system robustness and reduces the need for extensive retraining when adapting to new tasks.

**Eliminate Hallucinations.** A critical challenge for LLMs is the phenomenon of hallucinations—the generation of plausible yet factually incorrect or inconsistent information [6]. Hallucinations stem from the probabilistic nature of LLMs, which are trained to optimize likelihood functions rather than ensure factual correctness. This limitation poses significant risks in high-stakes domains such as healthcare, legal systems, scientific research, etc. MML offers a transformative solution by introducing modular learning and reasoning modules that enforce factuality. These modules operate as logical evaluators, cross-referencing LLM outputs with structured knowledge bases, ontologies, or predefined rules. For instance, in a medical diagnosis application, modular systems could validate treatment recommendations generated by an LLM against established clinical guidelines,

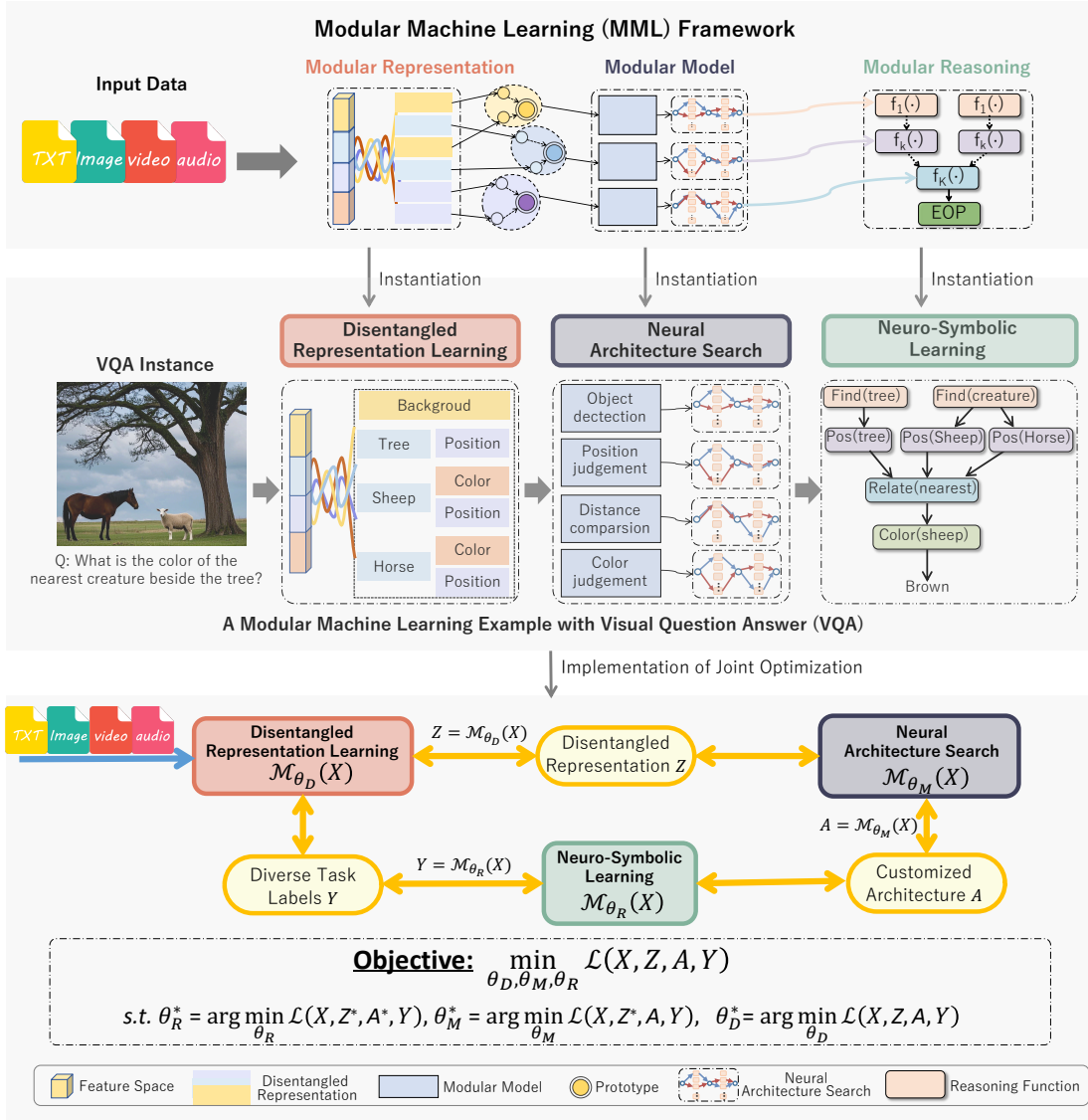


Fig. 2: The general framework of our proposed MML (Modular Representation, Modular Model, and Modular Reasoning) with a feasible implementation (Disentangled Representation Learning, Neural Architecture Search, and Neuro-Symbolic Learning) to solve a practical task. To maintain consistency with Fig. 1, we again take a visual question answering (VQA) task as an example to illustrate the instantiation of MML. We utilize the VQA example to demonstrate how MML disentangles the input into modular components, processing through a sequential combination of disentangled representation learning, neural architecture search, and neuro-symbolic learning. Specifically, the disentangled representation learning component is responsible for extracting a set of disentangled features from the input data, which are subsequently fed into the neural architecture search component. The neural architecture search component then identifies the optimal architecture that best captures the disentangled representations. Following this, the neuro-symbolic learning component further refines the output through structured, logic-driven reasoning. In the VQA example, the input image that contains a tree, a horse, and a sheep is processed by MML to extract disentangled representations indicating background, position, and color etc. Based on these disentangled representations, MML then discovers the optimal neural architecture to perform various tasks, including object detection, position judgment, distance comparison, and color identification. Finally, neuro-symbolic learning integrates the outputs in a logical reasoning manner to answer the question: “What is the color of the nearest creature beside the tree?” The answer, “Brown” correctly identifies the nearest creature (the sheep) and its color, using symbolic logic to establish the relationship between the entities. Additionally, we formally present the joint optimization process at the bottom of the figure, which involves the simultaneous optimization of the three modules: disentangled representation learning, neural architecture search, and neuro-symbolic learning. The objective function aims to minimize the loss  $\mathcal{L}(X, Z, A, Y)$  through a coordinated end-to-end learning process. Specifically, the optimal parameters for each module ( $\theta_D$ ,  $\theta_M$ ,  $\theta_R$ ) are determined by minimizing the respective loss functions in a joint manner: first optimizing the neuro-symbolic learning component, followed by the neural architecture search component, and finally the disentangled representation learning component. This joint optimization framework enables the seamless integration of different modular components, achieving a more efficient and effective learning process that can be adapted to diverse tasks.

Notations:  $\mathcal{L}$ : loss function (e.g., cross-entropy loss, MSE loss, etc.);  $X$ : input data;  $Z$ : disentangled representation;  $A$ : architecture of the modules;  $Y$ : target output, i.e., label;  $\theta_D$ ,  $\theta_M$ ,  $\theta_R$ : parameters of the representation, architecture and reasoning modules, respectively.

flagging any inconsistencies. Furthermore, MML can enable iterative refinement of LLM outputs through modular post-processing. Consider a scenario where an LLM generates a scientific summary. An MML-based module could validate the summary by applying logical rules derived from domain-specific knowledge, such as consistency checks for units, measurements, or causal relationships. This iterative process ensures that the final output is both factually accurate and logically coherent. Another promising direction is the use of MML for real-time monitoring and correction of LLM outputs. MML-based modules could act as dynamic filters, evaluating intermediate outputs during the generation process. For example, when generating legal documents, modular systems could enforce compliance with statutory requirements by rejecting outputs that violate predefined legal constraints. This dynamic interaction significantly enhances the reliability and trustworthiness of LLMs, making them suitable for applications where accuracy is paramount.

**Encourage Fairness and De-Biasing.** The issue of fairness and bias in LLMs arises from the biases hidden in their training data. These biases can manifest as discriminatory behaviors in applications ranging from recruitment to lending [7]. While traditional methods focus on mitigating bias at the data preprocessing or model training stage, MML introduces a novel paradigm by embedding fairness constraints directly into the reasoning process. MML-based modules within the framework can explicitly define fairness criteria, ensuring that decisions align with ethical guidelines and societal norms. For instance, in a job recommendation system, modular rules could enforce demographic parity by evaluating whether candidates from different groups receive equitable recommendations. Similarly, in lending scenarios, modular reasoning could validate creditworthiness decisions to prevent discrimination based on race or gender. In addition to enforcing fairness, MML can further facilitate bias detection and correction. Neural modules can preprocess large datasets to identify patterns indicative of bias, while modular reasoning can analyze these patterns to infer the underlying causations. This dual capability enables a comprehensive approach to addressing bias, combining the scalability of neural networks with the interpretability of modular systems. Moreover, the modularity of MML allows for the continuous evolution of fairness criteria. As societal norms change, human-understandable rules can be updated without retraining the entire system. This adaptability ensures that MML-enhanced LLMs remain aligned with contemporary ethical standards, setting a new benchmark for responsible AI development.

**Enforce Robustness.** Safety is a fundamental requirement for deploying AI systems in real-world scenarios, particularly in critical domains such as healthcare, finance, and autonomous systems. Traditional LLMs, while powerful, lack intrinsic mechanisms to ensure safety, making them susceptible to adversarial attacks, out-of-distribution scenarios, and catastrophic errors [8], [9]. MML addresses these vulnerabilities by embedding safety constraints within MML-based modules. These modules act as “guardrails”, enforcing rules that prevent unsafe behaviors. For instance, in autonomous driving applications, modular systems could validate decisions made by LLMs against predefined safety protocols, such as maintaining a safe distance from other

vehicles or adhering to traffic laws. The robustness of MML-enhanced LLMs extends beyond safety constraints. MML-based modules can perform real-time anomaly detection by analyzing outputs to check inconsistencies or deviations from expected patterns. In the domain of cybersecurity, for example, modular reasoning could identify unusual network activity indicative of a potential breach, enabling proactive intervention.

**Enhance Interpretability.** The interpretability of AI systems is a pressing concern, particularly as they are increasingly deployed in domains where accountability is paramount [10]. Traditional LLMs, often described as “black box”, struggle to provide explanations for their outputs, limiting their adoption in critical applications. MML offers a solution by integrating modular reasoning modules that generate interpretable outputs. These modules translate complex neural representations into human-readable logical explanations, bridging the gap between machine learning and human understanding. For instance, in a legal AI system, MML-based modules could explain how specific statutes and precedents influenced a recommendation, providing a clear audit trail for decisions. The modular architecture of MML also facilitates domain-specific customization of explanations. By tailoring modular rules to the requirements of a particular domain, MML can generate contextually relevant explanations. In healthcare, for example, modular systems could elucidate the clinical pathways leading to a diagnosis, enabling practitioners to validate and refine their decisions. Moreover, MML-enhanced LLMs support interactive explainability, allowing users to query and explore the reasoning behind outputs. For instance, in a financial application, users could ask why a loan application was rejected and receive a detailed explanation of the factors considered. This interactive capability not only enhances transparency but also fosters user trust and engagement. Last but not least, interpretability can also facilitate trust in AI systems by providing a clear rationale for decisions. For instance, in medical diagnostics, modular systems could generate detailed explanations for treatment recommendations, enabling clinicians to verify their validity. This transparency is essential for trust through fostering confidence in AI systems, particularly in high-stakes applications.

Therefore, LLMs have demonstrated remarkable capabilities in processing unstructured data, understanding natural language, and capturing semantic nuances across various modalities. However, their limitations in reasoning, factual consistency, fairness, and interpretability necessitate an enhanced framework that integrates MML principles. MML provides a structured backbone that augments LLMs, ensuring logical coherence, adaptability, and robustness in real-world applications. One of the critical advantages of MML is its modular architecture, which enables the decoupling of perception and reasoning. LLMs, with their superior ability to parse complex, high-dimensional data, can serve as the perceptual layer, converting raw inputs into structured representations. These structured representations can then be processed by dedicated reasoning modules within the MML framework, ensuring consistency, interpretability, and logical soundness. Of course, developing a unified framework requires addressing the technical challenge of integrating continuous representations produced by LLMs

with discrete modular reasoning. Differentiable reasoning architectures, such as neuro-modular integration pipelines, facilitate seamless interaction between these components. These architectures approximate modular logic with continuous operations, enabling gradient-based optimization while preserving logical consistency.

To summarize, the synergy between MML and LLMs leads to emergent properties that neither system can achieve independently. Inspired by cognitive dual-process theories [11], this hybrid system enables intuitive, fast decision-making through LLMs while ensuring deliberate, logical reasoning via MML-based modules. This dual-layered approach significantly improves decision accuracy and adaptability across diverse applications.

In addition, the interplay between LLM-driven content generation and MML-based validation is able to foster a robust iterative process for hypothesis generation and verification.

By structuring LLMs within the MML framework, AI systems can achieve enhanced robustness, interpretability, and adaptability, setting the foundation for the next generation of unified intelligent systems [12].

## 4 A FEASIBLE IMPLEMENTATION OF MML

We present the methodology of MML, and show the general framework with VQA as an example for instantiation in Fig. 2. In this section, we continue to elaborate on the feasible implementation via introducing three complementary strategies: disentangled representation learning (DRL), neural architecture search (NAS), and neuro-symbolic learning (NSL). First, DRL isolates independent semantic factors into modular representations, promoting transparency, generalizability, and controllability. The disentanglement enables the separation of complex data attributes, making the learned representations more interpretable and robust. Second, NAS automates the design of modular networks by efficiently exploring large search spaces, and identifying optimal operator configurations and interconnections that strike a balance between performance and computational efficiency. This automation significantly reduces human effort in model design while enhancing flexibility. Third, NSL integrates symbolic logic into the inference process within LLMs in a modular manner, enabling structured, rule-based validation and iterative refinement of neural outputs. This symbolic design mitigates issues such as hallucinations and reinforces logical consistency, becoming crucial for high-stakes applications. Collectively, these interconnected strategies exemplify the practical implementation of MML. They not only address the inherent limitations of conventional LLMs but also pave the way for the next generation of robust, interpretable, and logic-capable models in various real-world applications. Fig. 3 illustrates the overall optimization details of this MML implementation.

### 4.1 Disentangled Representation Learning for Modular Representation

Disentangled Representation Learning (DRL) for modular representation aims at separating and representing the underlying factors of variation in data independently and

meaningfully. DRL seeks to break down complex data into distinct components or modules, each focusing on a specific aspect, such as color, shape, or size in images. Disentanglement is crucial in unsupervised and reinforcement learning as it facilitates effective decision-making through clear and meaningful representations. Unlike traditional end-to-end deep learning models that directly learn data representations capturing entangled features, DRL strives to extract latent variables that correspond to individual factors, fostering human-like generalization [13]. The modular representation enables models to learn structured, explainable features, enhancing their generalization ability and adaptability to new situations [14].

**4.1.0.1 Definition: Disentangled Representation Learning (DRL):** Let  $\mathcal{M}_{\theta_D} = \{D_1, D_2, \dots, D_n\}$  represent a set of disentangled representation modules, where each module  $D_i$  maps the input data  $x \in \mathcal{X}$  to a latent variable  $z_i$  through a function  $d_{\theta_i}$ , i.e.,

$$z_i = d_{\theta_i}(x), \quad \forall i \in \{1, 2, \dots, n\}. \quad (3)$$

A *disentangled embedding* (a.k.a. *representation*)  $z = (z_1, z_2, \dots, z_n) \in \mathcal{Z}$  denotes the latent variables  $z_i$  correspond to statistically independent factors of variation  $v_i$  in the input data, such that:

$$p(z | x) = \prod_{i=1}^n p(z_i | v_i). \quad (4)$$

This indicates that each latent variable  $z_i$  is responsible for capturing only one underlying factor  $v_i$  and remains invariant to changes in other factors  $v_j$  ( $j \neq i$ ). The modular representation function  $\mathcal{M}_{\theta_D}$  can then be expressed as:

$$\mathcal{M}_{\theta_D}(x) = (d_{\theta_1}(x), d_{\theta_2}(x), \dots, d_{\theta_n}(x)). \quad (5)$$

This function outputs a set of disentangled and semantically meaningful features.

The objective of DRL is to learn disentangled representations by minimizing the following loss function:

$$\theta_D^* = \arg \min_{\theta_D} \mathcal{L}(\mathbf{X}, \mathbf{Z}, \mathbf{A}, \mathbf{Y}), \quad (6)$$

$$\mathbf{Z}^* = \mathcal{M}_{\theta_D^*}, \quad (7)$$

where  $\mathbf{Z}$  is the disentangled representation and  $\mathbf{Z}^*$  is the optimal learned representation,  $\mathcal{L}$  denotes the loss function that measures the discrepancy between the predicted and the true outputs,  $\mathbf{X}$  is the input data,  $\mathbf{A}$  is the architecture of the modules, and  $\mathbf{Y}$  is the target output. The optimization process aims to learn the parameters  $\theta_D$  of the disentangled representation modules such that the learned representations are semantically meaningful and independent.

This modular design highlights how DRL decomposes complex representations into independent, meaningful components, facilitating robust and interpretable learning. Specifically, DRL enables machine learning models to identify and disentangle hidden factors within observed data, thereby aligning learned representations with real-world semantics. These representations are invariant to external semantic changes [15], [16], aligned with real semantics, and robust against confounding or biased information [17], making them suitable for diverse downstream tasks.



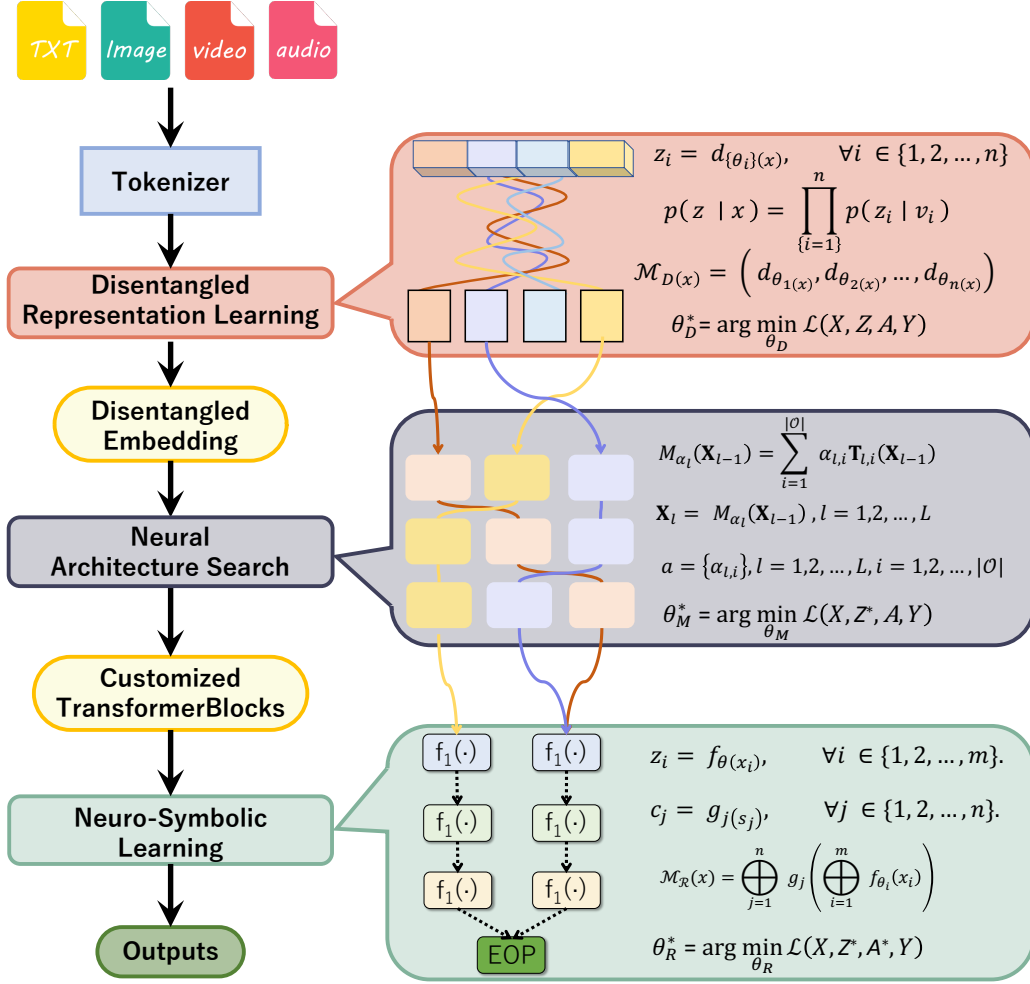


Fig. 3: The overall optimization of our proposed feasible MML implementation: a joint optimization of: 1) Disentangled Representation Learning (DRL) module, which is responsible for learning a set of disentangled representations from the input data; 2) Neural Architecture Search (NAS) module, which is responsible for searching customized transformerblocks for different functions; and 3) Neuro-Symbolic Learning (NSL) module, which is responsible for refining the outputs of the NAS module through neuro-symbolic reasoning.

Furthermore, DRL contributes to enhanced interpretability and efficiency in foundation models, such as ChatGPT and Stable Diffusion, by disentangling task-relevant knowledge from redundant components [18]. This capability aids in making foundation models more transparent and adaptable, addressing challenges related to task specificity and interpretability.

#### 4.2 Neural Architecture Search for Modular Model

Neural Architecture Search (NAS) automates the design of neural network architectures, significantly reducing human effort while achieving state-of-the-art performance [19]. It is particularly valuable for modular model, where complex networks are broken down into specialized substructures or blocks that can be independently designed and optimized. NAS aims to identify the most effective combination of these modules and their interconnections, thus improving overall performance while maintaining computational efficiency. The modular model allows NAS to adapt to various tasks, such as image classification, object detection, semantic segmentation,

and language processing, by tailoring network components to specific requirements [20], [21], [22].

**Definition 2.** (Neural Architecture Search (NAS)) Let  $\mathbf{M} = \{M_1, M_2, \dots, M_K\}$  represent a set of modular components, where each module  $M_i$  corresponds to a specific function. These modules are selected through routing from a supernet. We define each layer in this supernet as follows:

$$M_{\alpha_l}(\mathbf{X}_{l-1}) = \sum_{i=1}^{|\mathcal{O}|} \alpha_{l,i} \mathbf{T}_{l,i}(\mathbf{X}_{l-1}), \quad (8)$$

where  $\mathbf{T}_{l,i}$  is the  $i$ -th block of the  $l$ -th layer,  $\alpha_{l,i}$  is the weight of the  $i$ -th block of the  $l$ -th layer, and  $\mathbf{X}_{l-1}$  is the input to the  $l$ -th layer,  $|\mathcal{O}|$  is the number of blocks in the each layer. We can adopt a stack of layers to compose a module,

$$\mathbf{X}_l = M_{\alpha_l}(\mathbf{X}_{l-1}), \quad l = 1, 2, \dots, L, \quad (9)$$

where  $L$  is the total number of layers.

Let the architecture search space  $\mathcal{A}$  consist of all possible configurations formed by combining these blocks. Let  $a_j \in \mathcal{A}$

denote a specific architecture to form a module  $M_j$  to achieve a specific function.

$$a_j = \{\alpha_{l,i}\}, \quad l = 1, 2, \dots, L, i = 1, 2, \dots, |\mathcal{O}| \quad (10)$$

where  $\alpha_{l,i}$  determines how the module  $M_j$  weighs the blocks in each layer, and the group of them, i.e.  $a_j$ , is the configuration of the module  $M_j$  that determines how it selects different calculation paths within the supernet to provide the flexibility of the model design.

The objective of NAS  $\mathcal{M}_{\theta_M}$  is to find the different optimal architecture to form modular *w.r.t.* the task. Let  $\mathbf{A}$  denote the architecture set  $\{a_1, a_2, \dots, a_K\}$ , where  $a_j$  architecture for modular  $M_j$ . The optimization objective can be mathematically formulated as:

$$\theta_M^* = \arg \min_{\theta_M} \mathcal{L}(\mathbf{X}, \mathbf{Z}^*, \mathbf{A}, \mathbf{Y}), \quad (11)$$

$$\mathbf{A}^* = \mathcal{M}_{\theta_M^*}, \quad (12)$$

where  $\mathbf{A}$  is the chosen architecture and  $\mathbf{A}^*$  is the optimal learned architecture,  $\mathcal{L}$  denotes the loss function that measures the discrepancy between the predicted and the true outputs. The search process optimizes both the modules' design and their interconnections using search algorithms such as reinforcement learning, evolutionary algorithms, or gradient-based optimization.

This modular design allows for the flexible combination of different neural components while efficiently identifying the optimal architecture, balancing between performance and computational cost. The NAS process typically comprises three key components: search space, search strategy, and evaluation strategy [19].

The search space in modular NAS is divided into macro and micro spaces. The macro space addresses the overall architecture, while the micro space focuses on individual modules or blocks. The modular design enables the integration of diverse operators such as convolution, attention, or aggregation, making it suitable for different data types like images, sequences, and graphs [23]. The search strategies include reinforcement learning (RL), evolutionary algorithms (EA), and gradient-based methods. RL-based NAS treats the search process as a decision-making task, selecting modular components based on expected performance gains [24]. EA-based NAS evolves modular architectures by simulating biological evolution, iteratively selecting and mutating modules to enhance performance [25]. Gradient-based NAS, such as DARTS [26], relaxes the search space into a continuous domain, allowing for efficient gradient optimization of modular connections and configurations. Evaluation strategies aim to assess model performance efficiently, given the extensive search space. Techniques like weight sharing [27], predictor-based estimation [28], and zero-shot methods [29] reduce the cost of evaluating different modular combinations. One-shot NAS, for instance, trains a super-network where all modular sub-networks share weights, enabling rapid evaluation of individual architectures without retraining [30].

### 4.3 Neuro-Symbolic Learning for Modular Reasoning

Neuro-Symbolic Learning (NSL) integrates the strengths of neural networks and symbolic reasoning to develop modular AI systems that are both flexible and interpretable [31]. In

modular reasoning, complex tasks are divided into smaller, specialized components, each responsible for a specific aspect of the problem. By leveraging neural networks' ability to learn from data and symbolic reasoning's structured logic, NSL enhances the strengths of both dynamic adaptability and logical reasoning, making it suitable for complex tasks in robotics, natural language processing, and autonomous systems [32], [33].

**Definition 3.** (Neuro-Symbolic Learning (NSL))

Let  $\mathcal{M}_{\theta_R} = \{N_1, N_2, \dots, N_m\}$  represent a set of neural network modules, where each module  $N_i$  maps input data  $x_i$  to a latent representation  $z_i$  through a function  $f_{\theta_i}$ , i.e.,

$$z_i = f_{\theta_i}(x_i), \quad \forall i \in \{1, 2, \dots, m\}. \quad (13)$$

Let  $\mathcal{S} = \{S_1, S_2, \dots, S_n\}$  represent a set of symbolic reasoning modules, where each symbolic module  $S_j$  processes a symbolic representation  $s_j$  to produce a logical conclusion  $c_j$  through a reasoning function  $g_j$ , i.e.,

$$c_j = g_j(s_j), \quad \forall j \in \{1, 2, \dots, n\}. \quad (14)$$

The objective of NSL is to learn a joint modular model  $\mathcal{M}$  that combines neural modules  $\mathcal{N}$  and symbolic modules  $\mathcal{S}$  as follows:

$$\mathcal{M}_{\mathcal{R}}(x) = \bigoplus_{j=1}^n g_j \left( \bigoplus_{i=1}^m f_{\theta_i}(x_i) \right). \quad (15)$$

Here, the symbol  $\bigoplus$  denotes the modular aggregation operation, which combines the outputs from multiple modules into a unified representation. Specifically, the inner aggregation  $\bigoplus_{i=1}^m$  collects the outputs from the neural modules, while the outer aggregation  $\bigoplus_{j=1}^n$  combines the logical conclusions from the symbolic reasoning modules.

The model's optimization objective is to minimize the loss function:

$$\theta_R^* = \arg \min_{\theta_R} \mathcal{L}(\mathbf{X}, \mathbf{Z}^*, \mathbf{A}, \mathbf{Y}), \quad (16)$$

$$\mathbf{Y}^* = \mathcal{M}_{\theta_R^*}, \quad (17)$$

where  $\mathbf{Y}$  is the output of the model,  $\mathbf{X}$  is the input data,  $\mathbf{A}$  is the architecture of the modules, and  $\mathbf{Y}^*$  is the optimal learned output. The optimization process aims to learn the parameters  $\theta_R$  of the neural and symbolic modules such that the combined model produces accurate predictions while maintaining interpretability and logical consistency.

This modular design fosters a balanced interaction between neural and symbolic systems. For instance, the "System 1 and System 2" framework by Yoshua Bengio<sup>1</sup> outlines a collaborative model, where neural modules perform fast pattern recognition while symbolic modules conduct logical reasoning [31]. The modular reasoning is able to optimize both neural networks and symbolic systems, enhancing the efficiency and interpretability of complex reasoning tasks [34]. NSL also shows potential in improving LLMs by integrating logical consistency, guiding inference strategies, and enhancing prompt engineering, making LLMs more robust and interpretable [35], [36].

1. <https://yoshuabengio.org/2023/03/21/scaling-in-the-service-of-reasoning-model-based-ml/>



## 5 CHALLENGES AND FUTURE DIRECTIONS

Although hybrid systems that integrate MML into LLMs promise enhanced reasoning, safety, and explainability, their development confronts several challenges illuminating pathways for future research.

**Integration and Configurable Interface.** One of the foremost challenges is the seamless interface between neural networks—which operate in a continuous, gradient-driven space—and MML-based modules that execute discrete, logic-based operations. Current methods often depend on heuristic interfaces or differentiable relaxations to enable communication between these components. Moving forward, research should explore adaptive and dynamically configurable interfaces that can automatically modulate the degree of modular intervention based on input complexity or task demands. Advances in differentiable programming and co-training strategies may pave the way for truly unified architectures.

**Differentiability and Joint Optimization.** The inherent mismatch between differentiable neural computations and non-differentiable modular reasoning creates significant obstacles for end-to-end training. Current solutions rely on surrogate loss functions or approximations, which can introduce biases and limit performance. Developing novel training methodologies such as meta-learning, reinforcement learning-guided structure search, or even fully differentiable modular reasoning paradigms can help to bridge this gap. Such approaches would enable more accurate backpropagation through modular components, ensuring that the combined system is able to learn coherently.

**Scalability and Computational Efficiency.** LLMs are inherently computationally intensive, and the addition of modular reasoning layers may further escalate resource demands. This integration poses challenges for both training and real-time inference, particularly in resource-constrained environments or when processing high-volume data streams. Future research should prioritize the development of scalable architectures and efficient computational techniques for MML-based LLMs. Advances such as parameter-efficient fine-tuning, model quantization, and parallel processing frameworks can mitigate computational burdens while retaining the model performance.

**Evaluation and Benchmarking.** Traditional evaluation metrics for LLMs—centered on language fluency and task-specific accuracy—fall short of capturing the logical consistency and interpretability gains afforded by MML. Without standardized benchmarks that assess both performance and reasoning quality, it is difficult to measure progress across different systems. Establishing comprehensive evaluation protocols that include metrics for logical coherence, safety, fairness, and explanation quality will be essential. Such benchmarks will not only facilitate rigorous comparisons but also drive the development of models that are robust in safety-critical and high-stakes applications.

**Balancing Interpretability and Performance.** While integrating MML-based modules significantly enhances interpretability by providing clear audit trails and logical explanations, it often comes with trade-offs in raw predictive performance on data-intensive tasks. Striking a balance between maintaining high performance and achieving human-understandable reasoning remains an open research ques-

tion. Future work might investigate adaptive systems that dynamically adjust the level of modular processing based on the confidence or complexity of the task, such as inference strategies that modulate reasoning depth in response to uncertainty measures.

**Modular World Models for Adaptive Reasoning.** Another key future direction is the integration of Modular World Models (MWMs) to enhance reasoning, adaptability, and generalization with LLMs. Unlike traditional end-to-end deep learning approaches, MWMs decompose world knowledge into structured, reusable components that can dynamically interact with LLMs for causal inference, counterfactual reasoning, and real-world simulation. This modular approach enables efficient adaptation to new tasks, reducing the need for extensive retraining while improving interpretability and robustness. Future research should focus on a modular architecture that allows LLMs to query and update modular world models in the real-time world, leveraging advances in graph-based reasoning, differentiable programming, and structured knowledge distillation. By embedding MWMs, next-generation AI systems can move from static pattern recognition toward adaptive, compositional, and transparent decision-making in complex environments.

In summary, addressing these challenges requires a holistic research agenda that improves both the integration mechanisms and the training methodologies of MML-based LLMs. By developing adaptive interfaces, novel joint optimization techniques, scalable architectures, and standardized evaluation benchmarks, as well as finding effective ways to balance interpretability with performance, the community can unlock the full potential of MML systems. Such advancements will be critical in paving the way toward more robust, transparent, and versatile foundation models that can better meet the demands of real-world applications.

## 6 CONCLUSIONS

In this paper, we propose a new concept of Modular Machine Learning (MML), elaborating the pivotal role of MML in advancing LLMs. By dissecting the challenges inherent to current LLMs which range from limitations in logical reasoning and factual consistency to issues of interpretability and scalability, we illustrate how a modular learning paradigm can bridge the gap between raw (continuous) neural perception and structured (discrete) logic reasoning. We show that the integration of MML and LLMs not only mitigates issues such as hallucinations and bias, but also paves the way for more explainable and accountable AI applications in critical domains.

Looking ahead, addressing the challenges of configurable interface, joint optimization, and computational efficiency is crucial. Future research should prioritize the development of adaptive interfaces and advanced hybrid training methodologies, as well as the establishment of standardized benchmarks for rigorously evaluating both linguistic and logical competencies. By continuing to refine these approaches, the AI community can harness the full potential of MML, steering large foundation models toward greater reliability and real-world utility. As the final conclusion, we believe that MML ultimately stands as an indispensable strategy for the evolution of new-generation LLMs.

## REFERENCES

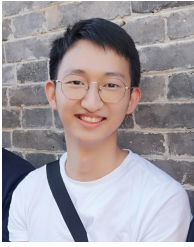
- [1] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat *et al.*, “Gpt-4 technical report,” *arXiv preprint arXiv:2303.08774*, 2023.
- [2] S. J. Hoch, “Counterfactual reasoning and accuracy in predicting personal events,” *Journal of Experimental Psychology: Learning, Memory, and Cognition*, vol. 11, no. 4, p. 719, 1985.
- [3] H. Dong, J. Mao, T. Lin, C. Wang, L. Li, and D. Zhou, “Neural logic machines,” *arXiv preprint arXiv:1904.11694*, 2019.
- [4] E. Cambria, L. Malandri, F. Mercurio, N. Nobani, and A. Seveso, “Xai meets llms: A survey of the relation between explainable ai and large language models,” *arXiv preprint arXiv:2407.15248*, 2024.
- [5] Z. Zeng, Q. Cheng, and Y. Si, “Logical rule-based knowledge graph reasoning: A comprehensive survey,” *Mathematics*, vol. 11, no. 21, p. 4486, 2023.
- [6] F. Larøi, T. M. Luhrmann, V. Bell, W. A. Christian Jr, S. Deshpande, C. Fernyhough, J. Jenkins, and A. Woods, “Culture and hallucinations: overview and future directions,” *Schizophrenia bulletin*, vol. 40, no. Suppl\_4, pp. S213–S220, 2014.
- [7] Y. Li, M. Du, R. Song, X. Wang, and Y. Wang, “A survey on fairness in large language models,” *arXiv preprint arXiv:2308.10149*, 2023.
- [8] A. Kumar, C. Agarwal, S. Srinivas, A. J. Li, S. Feizi, and H. Lakkaraju, “Certifying llm safety against adversarial prompting,” *arXiv preprint arXiv:2309.02705*, 2023.
- [9] A. Wei, N. Haghtalab, and J. Steinhardt, “Jailbroken: How does llm safety training fail?” *Advances in Neural Information Processing Systems*, vol. 36, pp. 80 079–80 110, 2023.
- [10] S. Huang, S. Mamidanna, S. Jangam, Y. Zhou, and L. H. Gilpin, “Can large language models explain themselves? a study of llm-generated self-explanations,” *arXiv preprint arXiv:2310.11207*, 2023.
- [11] C. G. Beevers, “Cognitive vulnerability to depression: A dual process model,” *Clinical psychology review*, vol. 25, no. 7, pp. 975–1002, 2005.
- [12] S. Wu, H. Fei, L. Qu, W. Ji, and T.-S. Chua, “Next-gpt: Any-to-any multimodal llm,” *Forty-first International Conference on Machine Learning*, 2024.
- [13] R. Geirhos, J.-H. Jacobsen, C. Michaelis, R. Zemel, W. Brendel, M. Bethge, and F. A. Wichmann, “Shortcut learning in deep neural networks,” *Nature Machine Intelligence*, vol. 2, no. 11, pp. 665–673, 2020.
- [14] Y. Bengio, A. Courville, and P. Vincent, “Representation learning: A review and new perspectives,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 35, no. 8, pp. 1798–1828, 2013.
- [15] H. Kim and A. Mnih, “Disentangling by factorising,” pp. 2649–2658, 2018.
- [16] J. Lee, E. Kim, J. Lee, J. Lee, and J. Choo, “Learning debiased representation via disentangled feature augmentation,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 25 123–25 133, 2021.
- [17] R. Suter, D. Miladinovic, B. Schölkopf, and S. Bauer, “Robustly disentangled causal mechanisms: Validating deep representations for interventional robustness,” pp. 6056–6065, 2019.
- [18] J. Zeng, Y. Jiang, S. Wu, Y. Yin, and M. Li, “Task-guided disentangled tuning for pretrained language models,” pp. 3126–3137, 2022.
- [19] T. Elsken, J. H. Metzen, and F. Hutter, “Neural architecture search: A survey,” *The Journal of Machine Learning Research*, vol. 20, no. 1, pp. 1997–2017, 2019.
- [20] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le, “Learning transferable architectures for scalable image recognition,” pp. 8697–8710, 2018.
- [21] Y. Chen, T. Yang, X. Zhang, G. Meng, X. Xiao, and J. Sun, “Detnas: Backbone search for object detection,” *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [22] Y. Wang, Y. Yang, Y. Chen, J. Bai, C. Zhang, G. Su, X. Kou, Y. Tong, M. Yang, and L. Zhou, “Textnas: A neural architecture search space tailored for text representation,” vol. 34, no. 05, pp. 9242–9249, 2020.
- [23] B. Zoph and Q. V. Le, “Neural architecture search with reinforcement learning,” *arXiv preprint arXiv:1611.01578*, 2016.
- [24] Y. Jaafra, J. L. Laurent, A. Deruyver, and M. S. Naceur, “Reinforcement learning for neural architecture search: A review,” *Image and Vision Computing*, vol. 89, pp. 57–66, 2019.
- [25] K. De Jong, “Evolutionary computation: a unified approach,” pp. 185–199, 2016.
- [26] H. Liu, K. Simonyan, and Y. Yang, “Darts: Differentiable architecture search,” *arXiv preprint arXiv:1806.09055*, 2018.
- [27] L. Xie, X. Chen, K. Bi, L. Wei, Y. Xu, L. Wang, Z. Chen, A. Xiao, J. Chang, X. Zhang *et al.*, “Weight-sharing neural architecture search: A battle to shrink the optimization gap,” *ACM Computing Surveys (CSUR)*, vol. 54, no. 9, pp. 1–37, 2021.
- [28] C. White, A. Zela, B. Ru, Y. Liu, and F. Hutter, “How powerful are performance predictors in neural architecture search?” *arXiv preprint arXiv:2104.01177*, 2021.
- [29] H. Chen, M. Lin, X. Sun, and H. Li, “Nas-bench-zero: A large scale dataset for understanding zero-shot neural architecture search,” 2021.
- [30] H. Pham, M. Guan, B. Zoph, Q. Le, and J. Dean, “Efficient neural architecture search via parameters sharing,” *ICML*, pp. 4095–4104, 2018.
- [31] Y. Bengio, “From system 1 deep learning to system 2 deep learning,” 2019.
- [32] K. Yi, J. Wu, C. Gan, A. Torralba, P. Kohli, and J. B. Tenenbaum, “Neural-symbolic vqa: Disentangling reasoning from vision and language understanding,” *arXiv preprint arXiv:1810.02338*, 2018.
- [33] W.-Z. Dai, Q. Xu, Y. Yu, and Z.-H. Zhou, “Bridging machine learning and logical reasoning by abductive learning,” pp. 2811–2822, 2019.
- [34] J. Pfeiffer, S. Ruder, I. Vulić, and E. M. Ponti, “Modular deep learning,” *arXiv preprint arXiv:2302.11529*, 2023.
- [35] Y. Zhang, Y. Li, L. Cui, D. Cai, L. Liu, T. Fu, X. Huang, E. Zhao, Y. Zhang, Y. Chen *et al.*, “Siren’s song in the ai ocean: A survey on hallucination in large language models,” *arXiv preprint arXiv:2309.01219*, 2023.
- [36] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. L. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray *et al.*, “Training language models to follow instructions with human feedback,” *arXiv preprint arXiv:2203.02155*, 2022.



**Xin Wang** is currently an Associate Professor at the Department of Computer Science and Technology, Tsinghua University. He got both of his Ph.D. and B.E degrees in Computer Science and Technology from Zhejiang University, China. He also holds a Ph.D. degree in Computing Science from Simon Fraser University, Canada. His research interests include multimedia intelligence, machine learning and its applications. He has published over 200 high-quality research papers in ICML, NeurIPS, IEEE TPAMI, IEEE TKDE, ACM KDD, WWW, ACM SIGIR, ACM Multimedia etc., winning three best paper awards including ACM Multimedia Asia. He is the recipient of ACM China Rising Star Award, IEEE TCMC Rising Star Award and DAMO Academy Young Fellow.



**Haoyang Li** is currently a postdoc researcher in the Department of Population Health Sciences at Weill Cornell Medicine of Cornell University. He received his Ph.D. from the Department of Computer Science and Technology of Tsinghua University in 2023. He received his B.E. from the Department of Computer Science and Technology of Tsinghua University in 2018. His research interests are mainly in machine learning on graphs and out-of-distribution generalization. He has published high-quality papers in prestigious journals and conferences, e.g., IEEE TKDE, ACM TOIS, NeurIPS, ACM KDD, ACM Web Conference, AAAI, IJCAI, ICLR, ACM Multimedia, IEEE ICDE, IEEE ICDM, etc., winning one best paper award.



**Zeyang Zhang** received his B.E. from the Department of Computer Science and Technology, Tsinghua University in 2020. He is a Ph.D. candidate in the Department of Computer Science and Technology of Tsinghua University. His main research interests focus on graph representation learning, automated machine learning and out-of-distribution generalization. He has published several papers in prestigious conferences, e.g., NeurIPS, AAAI, etc.



**Haibo Chen** is currently a Ph.D. student at the Department of Computer Science and Technology, Tsinghua University. He received his B.E. degree from the School of Computer Science and Engineering, Central South University. His main research interests include graph machine learning, out-of-distribution learning, and multi-modal graphs.



**Wenwu Zhu** is currently a Professor in the Department of Computer Science and Technology at Tsinghua University. He also serves as the Vice Dean of National Research Center for Information Science and Technology, and the Vice Director of Tsinghua Center for Big Data. Prior to his current post, he was a Senior Researcher and Research Manager at Microsoft Research Asia. He was the Chief Scientist and Director at Intel Research China from 2004 to 2008. He worked at Bell Labs, New Jersey as Member of Technical Staff during

1996-1999. He received his Ph.D. degree from New York University in 1996. His research interests are in the area of data-driven multimedia networking and Cross-media big data computing. He has published over 400 referred papers and is the inventor or co-inventor of over 100 patents. He received eight Best Paper Awards, including ACM Multimedia 2012 and IEEE Transactions on Circuits and Systems for Video Technology in 2001 and 2019.

He served as EiC for IEEE Transactions on Multimedia from 2017-2019. He served in the steering committee for IEEE Transactions on Multimedia (2015-2016) and IEEE Transactions on Mobile Computing (2007-2010), respectively. He serves as General Co-Chair for ACM Multimedia 2018 and ACM CIKM 2019, respectively. He is an AAAS Fellow, IEEE Fellow, SPIE Fellow, and a member of The Academy of Europe (Academia Europaea).