



**Politecnico
di Torino**

Politecnico di Torino

Master's Degree in Data Science and Engineering

A.a. 2024/2025

Graduation Session July, 2025

Neuro-Symbolic Artificial for Visual Reasoning via Dynamic Logic Tensor Networks

Supervisor:

Lia Morra

Candidate:

Homayoun Afshari

Summary

Recent advancements in Large Language Models (LLMs) have demonstrated impressive abilities in contextual understanding and chain-of-thought reasoning. However, while these neural networks excel at inductive pattern recognition, they still struggle with deductive rule chaining and abductive hypothesis generation. As a promising solution, Neuro-Symbolic (NeSy) AI integrates this neural perception with symbolic logic to equip AI with more reasoning capabilities. This thesis focuses on visual reasoning through NeSy AI, a domain where perception and logic must jointly operate to interpret partially observable visually encoded structures. Despite progress in this field, current NeSy visual reasoning approaches often suffer from three key limitations: task-specific rigidity (lack of flexibility), the generation of non-interpretable rules (lack of explainability), and the production of informal rules that are difficult to verify automatically (lack of formality). To address these challenges, this work proposes a novel NeSy visual reasoning framework designed to simultaneously achieve flexibility, explainability, and formality. The architecture takes both textual and visual descriptions of a reasoning task and processes them through a context generator. The resulting context is then used by a rule generator to produce symbolic rules, while a visual processor independently converts the visual input into symbolic visual atoms. Consequently, these rules and symbols are evaluated by a rule verifier, which checks their coherence against visual facts. The system uses this feedback to iteratively refine both context and rule generation, enabling continuous improvement in a manner that is similar to reinforcement learning. To evaluate this system, we use the ViSudo-PC benchmark, which is a symbolic visual reasoning task that distinguishes valid from corrupted Sudoku boards using digits drawn from various visual domains. Results show that the system effectively induces symbolic rules from multi-modal inputs, verifies them automatically, and improves its reasoning performance through iterative feedback. Furthermore, when the rule generator is bypassed to retain only the visual encoder and rule verifier and ensure compatibility with existing benchmark protocols, the system achieves significant improvements over state-of-the-art methods. Accordingly, the proposed approach not only satisfies all its predefined objectives, but also surpasses existing methods in performance.

Acknowledgements

I would like to express my deepest gratitude to my advisor, Professor Lia Morra, for her continuous guidance, support, and encouragement throughout this research. I am also thankful to the members of her team, particularly Ph.D candidate Alessandro Russo, for his valuable feedback and insights. Special thanks to my classmates and friends for the stimulating discussions and moral support during the course of this thesis. Finally, I am deeply grateful to my family and my beloved wife, Fatemeh, for their unwavering encouragement and patience.

Table of Contents

List of Figures	VIII
1 Introduction	1
1.1 What is Reasoning?	1
1.2 How can AI Reason?	2
1.3 A Solution: NeSy AI	2
1.4 NeSy AI Empowered by Fuzzy Logic	3
1.5 Visual Reasoning through NeSy AI	4
1.6 Issues of NeSy Visual Reasoning	4
1.7 Novelties of the Thesis	5
1.8 Structure of the Thesis	7
2 Literature Review	8
2.1 Foundational Concepts	8
2.2 Different Tasks and Benchmarks	10
2.3 Basic Visual Reasoning Frameworks	12
2.4 Visual Reasoning with VSAs	13
2.5 Visual Reasoning with LTNs	14
2.6 Visual Reasoning with LLMs	15
3 Methodology	17
3.1 Proposed Method	17
3.2 Rule Generation	18
3.2.1 Prompt Engineering	18
3.2.2 Symbolic Rule Generation	21
3.3 Rule Verification	21
3.3.1 Syntax Tree Generation	22
3.3.2 LTN Implementation	26
3.4 Visual Processing	30
3.4.1 CNN Design	31
3.4.2 Visual Symbol Definition	32

3.5	Feedback Handling	32
3.5.1	LTN Development	33
3.5.2	Loop Creation	33
4	Experimental Environment	36
4.1	Setup	36
4.1.1	Benchmark and Dataset	36
4.1.2	Programming Tools	38
4.2	Specifications	39
4.2.1	VLM	39
4.2.2	D-LTN	41
4.2.3	CNN Block	43
4.3	Development Strategy	43
4.3.1	System Loops	43
4.3.2	Data Usage	44
4.4	Evaluation Metrics	44
4.4.1	Loss Function	44
4.4.2	Performance Metrics	45
5	Results	46
5.1	Experiments	46
5.1.1	Hyperparameter Tuning	46
5.1.2	Rule Generation	47
5.1.3	Experimental Tests	49
5.2	Discussions	50
5.2.1	Analyses	50
5.2.2	Comparison with Prior Literature	51
6	Conclusion	52
6.1	Findings	52
6.2	Future Directions	53
	Bibliography	55

List of Figures

1.1	Conceptual diagram of the proposed NeSy system.	6
3.1	Functional diagram of the proposed NeSy system.	17
3.2	Reasoning with the help of VLM.	18
3.3	Automatic generation of LTN formulas.	22
3.4	Perception with the help of CNN.	31
3.5	Detailed functional diagram of the proposed NeSy system.	32
4.1	First pair of the training subset of the 11-split of the benchmark. .	37

Chapter 1

Introduction

1.1 What is Reasoning?

Recently, there has been a critical shift in Artificial Intelligence (AI) research: the pursuit of models that not only excel in **Contextual Understanding**, as done by Large Language Models (LLMs), but also demonstrate powerful reasoning capabilities, as if a human brain is pulling the strings behind the scene [1]. Not from a long ago, we have observed how OpenAI’s ChatGPT or Anthropic’s Claude have been equipped with the ability to generate **Chains of Thought** in order to resolve analytical problems, along with their remarkable abilities in pattern recognition, language generation, and zero-shot generalization [2]. These behaviors have sparked numerous discussions in the AI community about how these **Trainable Machines** are able to perform such tasks [3]. Are they *truly* capable of reasoning, or are they merely *simulating* reasoning behaviors? This question, however, is a serious dilemma that may never be answered, but it poses a fundamental question; **What Is Reasoning?**

As defined by cognitive science, reasoning is the process of *drawing conclusions* from available information using three primary modes: **Induction**, **Deduction**, and **Abduction** [4]. Induction involves *generalizing* rules from specific observations [5]. For instance, noticing that the sun rises every morning, one can infer that it will do so again tomorrow. Deduction, on the other hand, uses established *rules* to derive specific *conclusions* [6]. For example, if all humans are mortal and Socrates is human, then one can deduce that Socrates must also be mortal. Abduction, however, is more subtle and complex. It involves inferring the *most likely explanation* from incomplete evidence, as if a detective is hypothesizing the cause of a mysterious event [7]. For example, observing a wet floor under a cloudy sky, one can believe that it must have been raining lately. Understanding these distinct forms of reasoning is essential if we want to assess how far AI models have

come and how far they still have to go in achieving or mimicking true reasoning.

1.2 How can AI Reason?

At the heart of AI, **Neural Networks**, including LLMs, have demonstrated strong capabilities in induction, primarily due to their *training on large datasets* and their proficiency in pattern extraction [8]. In other words, considering that inductive reasoning involves generalizing from specific examples to broader patterns, neural networks excel in this regard because they are optimized to minimize *prediction error* over millions of instances, which stems from their ability to implicitly encode associations within high-dimensional vector spaces and capturing complex correlations *without* the need for explicit knowledge representations [5]. This property allows these models excel at generalizing from specific instances, capturing statistical regularities, and producing coherent outputs across many tasks such as language modeling, image captioning, and even basic question-answering [9].

Nevertheless, when it comes to deduction, the reasoning performance of neural networks is often *contingent* upon the presence of implicitly encoded logic within the data distribution [10]. In other words, since these AI models rely heavily on statistical associations rather than formal rule-following, they tend to lose strength when reasoning requires precise logical consistency or the chaining of abstract rules. More critically, abduction also remains a persistent reasoning challenge. Effective abductive reasoning in AI requires *integrating* world knowledge, contextual awareness, and causal inference to construct plausible but unobserved hypotheses [7]. Neural networks, *by design*, lack explicit mechanisms to support this kind of inferential leap, which makes them unreliable for tasks that require *systematical exploration* of reasoning paths beyond what the data directly supports [5].

1.3 A Solution: NeSy AI

Early paradigms in AI, known as Good Old-Fashioned AI (**GOF AI**), focused on *symbolic reasoning*, where knowledge was represented explicitly in logic-based systems and manipulated through rule-based inference [11]. Unlike neural networks, GOF AI systems were designed to perform all three major forms of reasoning by leveraging structured representations, ontologies, and well-defined search procedures [10]. However, while these symbolic systems offered explainability and formal guarantees, but often *lacked* scalability and robustness in dealing with noisy or unstructured data [5]. With the rise of neural networks, however, modern AI systems began to prioritize statistical learning over symbolic manipulation, achieving impressive results in perception and pattern recognition but *sacrificing*

core reasoning capabilities in the process. This shift has created a discontinuity: we now expect neural networks to handle reasoning tasks that were once *explicitly* addressed by symbolic AI, but *without* equipping them with the tools to do so.

As a consequence, the so-called issues have *prompted* AI researchers to explore hybrid approaches that combine the inductive strength of neural networks with deductive and abductive potentials of GOF AI, aiming to develop systems that can reason more robustly across all three reasoning modes [5]. One such promising direction is Neuro-Symbolic (**NeSy**) AI, which seeks to integrate the sub-symbolic power of neural networks with the explicit and structured reasoning capabilities of symbolic systems [8]. This hybridization aims to leverage the scalability and flexibility of neural models while incorporating the transparency, explainability, and logical rigor of symbolic representations.

This architectural duality originates from two-system theory of reasoning in cognitive science: **System 1** and **System 2** [3]. The first refers to fast, automatic, and cognitive reasoning processes that are typically unconscious, effortless, and operate in parallel, which is similar to the pattern-based memory-driven nature of neural networks optimized for inductive reasoning [1]. On the other hand, system 2 represents slow, deliberate, and analytical thought through explicitly manipulating structured representations like rules, logic, and symbols, which resembles deduction and abduction in symbolic reasoning [3]. As in neural networks, system 1 enables rapid recognition and generalization, but it often *lacks* precision [12]. System 2, by contrast, enables methodical and consistent reasoning, *albeit* at a higher computational cost [5]. A successful NeSy system, therefore, aims to simulate the complementary functions of both systems, enabling NeSy models to handle a broad range of reasoning tasks with greater versatility and depth.

1.4 NeSy AI Empowered by Fuzzy Logic

In the context of NeSy AI, **Symbols** refer to abstract representations within a specific domain that carry explicit meanings, such as constants, variables, numbers, or fixed semantic labels [5]. These symbols can be combined and manipulated using formal languages like First-Order Logic (FOL), allowing for the representation of structured knowledge and reasoning over it [5]. For instance, a rule such as $\forall x : \text{isBrid}(x) \rightarrow \text{canFly}(x)$ encodes a general statement about birds' ability to fly, where each component, including quantifiers, predicates, and logical operators, is a symbolic element. This structured framework enables symbolic components of NeSy systems to draw logical conclusions and generate plausible hypotheses from incomplete observations, i.e., deduction and abduction through symbolic reasoning, thereby overcoming some of the fundamental reasoning limitations faced by neural networks [3].

Considering this importance of logical formalisms in NeSy AI, recent research has explored moving toward **Real Logic**, which is a differentiable form of fuzzy logic [5]. Unlike binary logic systems that treat statements as either entirely true or false, fuzzy logic introduces degrees of truth, allowing the system to represent and reason with uncertainty and vagueness that often characterize real-world data. Real Logic enables soft constraints to be imposed over continuous values and supports gradient-based optimization, which makes it compatible with neural networks during training [5]. As a result, symbolic rules are no longer fixed or hand-crafted, but can be learned or fine-tuned using stochastic gradient descent based on observed data [13]. This ability to merge differentiability with symbolic reasoning significantly enhances the adaptability of NeSy models in ambiguous or noisy environments, where classical logic systems typically struggle to operate effectively [8]. For instance, the rule $\text{isBird}(x) \rightarrow \text{canFly}(x)$ can be evaluated to a degree, rather than as strictly true or false, allowing for exceptions (e.g., penguins) to be incorporated naturally.

1.5 Visual Reasoning through NeSy AI

The soft interpretation of logic allows NeSy systems to achieve more robust generalization and contextual reasoning across a variety of domains, which sets the stage for applications in dynamic and perceptually complex settings like **Visual Reasoning** [14]. In such settings, NeSy systems must integrate perception from partial information with reasoning about the relationships between different visual objects, which often involves partial observations, occlusions, or ambiguous scenes [12]; the conditions under which, pure neural networks may fail to generalize and purely symbolic systems may lack flexibility [3]. Accordingly, in this thesis, we will also focus on the class of NeSy systems that are built on top of fuzzy logic for performing visual reasoning tasks.

1.6 Issues of NeSy Visual Reasoning

Regardless of all the opportunities NeSy AI provides for promising visual reasoning, several limitations still persist across different methodological paradigms. These limitations highlight ongoing trade-offs between flexibility, explainability, and formality in current approaches [12, 15]. We will discuss these approaches in detail in Section ??; here, we only focus on high-level descriptions of their shortcomings. In this respect, we can categorize the issues into three groups: Flexible But Non-Explainable (**FBN**), Explainable But Rigid (**ENR**), and Flexible But Verbose (**FBV**) approaches.

FBN approaches emphasize learning symbolic reasoning patterns directly from

visual data without hard-coded rules, often in flexible and open-ended environments. These methods embed reasoning mechanisms within neural network weights or architectures [3]. While this enables scalability and end-to-end training, it often results in latent rules that are non-explainable [15]. In other words, the visual reasoning processes become entangled with the model’s internal representations, making it difficult to extract, verify, or validate the learned logic. As a result, despite solving the task effectively, these models sacrifice transparency and explainability.

In response to this, ENR methods address the explainability gap by explicitly embedding symbolic rules into the system prior to training. Therefore, by integrating known symbolic rules as inductive biases or constraints, these approaches enforce explainability and logical consistency during both training and inference [9]. However, such methods often suffer from limited generalizability. Their performance and applicability hinge on domain-specific rule sets that do not easily transfer to new tasks, unseen visual domains, or more open-ended reasoning problems.

FBV methods, on the other hand, aim to overcome both limitations by generating rules dynamically at runtime, often in the form of free-form text or structured outputs [16]. For example, LLMs with proper prompting can produce reasoning traces that explain their decisions [16]. While this introduces flexibility and an element of explainability, the resulting explanations are typically verbose, informal, or non-machine-readable [15]. This restricts their utility in downstream symbolic processing and makes verification challenging. Additionally, such free-form reasoning is highly dependent on prompt engineering and lacks the formal structure needed for robust generalization.

Taken together, these limitations suggest that current NeSy visual reasoning systems have not yet reached the peak of their promised potential. Therefore, fully addressing these problems remains an open challenge, and future work will need to investigate architectures that can learn modular, interpretable, and grounded rules that generalize across tasks without requiring manual encoding or opaque abstraction. That is what we also seek to explore in this thesis.

1.7 Novelties of the Thesis

Considering the challenges highlighted in previous sections, this thesis proposes a novel NeSy visual reasoning framework grounded in the flexibility, explainability, and formality principles. As shown in Figure 1.1, the system processes two modalities: textual and visual. The **Textual Input** can contain optional contextual information, while the **Visual Input** comprises images representing the reasoning task. These inputs, along with an *intermediate feedback*, are integrated via an **Context Generator** to produce a semantically enriched **Context**. This input is then passed to a neural **Rule Generator**, which proposes symbolic rules in a

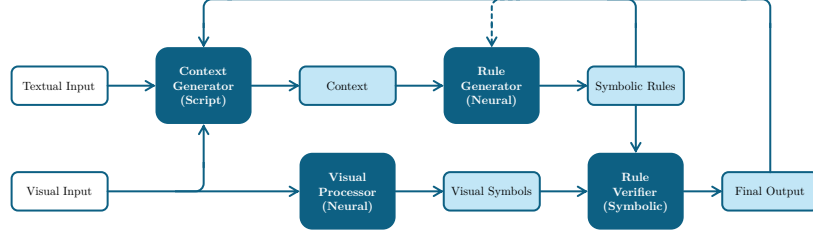


Figure 1.1: Conceptual diagram of the proposed NeSy system.

formal language such as FOL. Concurrently, the visual input is also transformed into **Visual Symbols** through a **Visual Processor**. After that, both the symbolic rules and visual symbols are processed by a symbolic **Rule Verifier**, which assesses their alignment and determines the **Final Output**.

On top of containing the results of the described experiment, both the generated symbolic rules and the final output serve as the source of the so-called intermediate feedback. This feedback is *primarily* designed to enhance the overall performance of the system by informing the context generator about possible failures or inconsistencies in previous experiments. By incorporating this feedback, the system can dynamically adjust the semantic context to better support rule generation. Additionally, the feedback can also be *directly* passed to the rule generator, which can enable it to immediately fine-tune its symbolic hypotheses based on empirical evidence. This optional feedback also increases the robustness and *expressive power* of the generated rules and promotes more efficient convergence accurate reasoning.

Consequently, the proposed system operates in an **Iterative Loop** that resembles Reinforcement Learning (**RL**) dynamics, as it iteratively uses the intermediate feedback to gradually refine its internal blocks. This process continues until a predefined **Termination Criterion** is met, such as when the generated symbolic rules adequately explain the majority of the images included in the visual input. This iterative design not only enables **Continuous Refinement of Symbolic Hypotheses** but also facilitates **Feedback-Based Learning and Convergence**. Through this mechanism, the system addresses common NeSy limitations by:

- **Providing Flexibility** in visual reasoning tasks, imposing minimal constraints on the types of problems it can handle.
- **Ensuring Explainable** both regarding its architectural design and the rules it generates at the end.
- **Upholding Formality** in rule generation, enabling the output to be understood and verified by both humans and machines.

1.8 Structure of the Thesis

This thesis is organized into six main chapters, including the current one. Chapter 2 presents a comprehensive literature review, covering foundational concepts, relevant tasks and benchmarks, and various approaches to NeSy visual reasoning. Chapter 3 outlines the proposed method and explains each component in details, from symbolic rule generation and verification to visual processing and feedback integration. Chapter 4 describes the experimental environment, including component specifications, the development strategy, and the evaluation metrics used. Chapter 5 reports the experimental results and provides a critical analysis of the outcomes. Finally, Chapter 6 summarizes the key contributions of the thesis and suggests potential directions for future research.

Chapter 2

Literature Review

In this chapter, we review foundational concepts, benchmarks, and architectural paradigms within NeSy AI with a focus on visual reasoning tasks, where we will cover different approaches for integrating neural and symbolic methods.

2.1 Foundational Concepts

As discussed earlier, NeSy AI aims to *synthesize* the complementary capabilities of neural networks and symbolic reasoning systems, moving beyond the traditional dichotomy between *connectionist* and *symbolic* paradigms [3]. A growing body of work now explores how these paradigms can be integrated to form hybrid models that are both data-efficient and cognitively plausible [1]. Importantly, NeSy systems *are not* monolithic, instead, they vary in architecture, areas of integration, and the degree to which symbolic rules influence learning or inference [17]. This diversity is not merely technical but reflects deeper assumptions about cognition and representation, *increasingly* shaped by interdisciplinary insights from cognitive science, linguistics, and formal logic [3].

To bring structure to this growing body of work, **Professor Henry Kautz** proposed a detailed taxonomy that *classifies* NeSy systems based on the nature and depth of interaction between neural and symbolic components. The taxonomy includes six types [3, 9]:

1. **symbolic Neuro symbolic**: This category resembles the standard pipeline of deep learning, where symbolic inputs (e.g., words) are *transformed* into vector embeddings, processed through neural architectures, and then mapped back into symbolic outputs. Most Natural Language Processing (**NLP**) systems fall into this category.
2. **Symbolic[Neuro]**: In this configuration, a symbolic system governs the overall

logic and control, while neural networks are used internally as *subroutines*. A well-known example is **AlphaGo**, which is an AI playing the game of GO through symbolic Monte Carlo tree search enhanced by neural networks to estimate game states.

3. **Neuro;Symbolic**: This architecture features a pipelined integration in which neural and symbolic systems operate on *distinct components* of the task. Therefore, these systems communicate in a loosely coupled manner, such as passing information back and forth to enhance mutual performance.
4. **Neuro:Symbolic** \rightarrow **Neuro**: Here, symbolic knowledge is *embedded* directly into the structure of the neural network itself. In other words, rather than merely using symbolic inputs, these systems compile symbolic rules into architectural priors or weight initialization, influencing the neural learning process in a deeper model-intrinsic way.
5. **NeuroSymbolic**: This category involves *tensorizing* symbolic structures (e.g., FOL representations) so that neural networks can perform reasoning tasks over them. These approaches maintain an interplay between formal logic and neural computation in a fully differentiable manner [5].
6. **Neuro[Symbolic]**: This is the *transpose* of the second category, where a neural architecture performs symbolic reasoning by learning structural relations or by attending to symbolic elements when required. As a prototypical example, we can consider Graph Neural Networks (**GNNs**), where the neural model effectively learns and reasons over symbolic graph structures.

The taxonomy provides a conceptual foundation for analyzing and comparing the many flavors of NeSy architectures and reflects both *practical implementation concerns* such as modularity, explainability, or scalability and *theoretical insights* into how symbolic and sub-symbolic representations can co-evolve [5]. Building on this framework, the survey **Towards Cognitive AI Systems** extends Kautz’s taxonomy by introducing dimensions such as probabilistic integration, types of intermediate representations, and degrees of explainability, which offer *more comprehensive mapping* of recent developments in NeSy AI [9]. In this survey, rather than framing integration as a binary property, it is tried to treat NeSy AI as a *spectrum* that ranges from loosely coupled ensembles to deeply fused models with joint optimization. This perspective underscores how design decisions at the neural or symbolic interface affect the cognitive fidelity and explanatory power of the resulting systems [3]. Importantly, the survey highlights that these trade-offs are *particularly* impactful in domains like visual reasoning, where perception and logic must interact seamlessly to achieve meaningful interpretations.

As mentioned earlier, again, visual reasoning, in particular, serves as a compelling testbed for NeSy architectures, which requires *integrating* bottom-up perception with top-down logical inference to interpret scenes, detect causal relations, and apply structured reasoning over visual inputs [18]. However, unlike conventional visual recognition tasks, visual reasoning demands interpreting spatial and temporal patterns, analogies, and counterfactuals [7]. We will cover the examples later, but here, as an overview, we can consider tasks such as how VAR challenges models to infer plausible causes for incomplete visual sequences or how RPM tests the ability to detect and apply abstract visual rules [7, 18]. NeSy models address these challenges by *jointly* learning from visual data and structured symbolic domains, such as spatial logic, temporal constraints, or analogical schemes [8]. These capabilities not only benchmark NeSy systems but also drive innovation in architectures that balance learned representations with interpretable symbolic reasoning [1].

2.2 Different Tasks and Benchmarks

To evaluate NeSy visual reasoning systems, a range of tasks and benchmarks have been proposed, each emphasizing different aspects of reasoning capabilities [19]. One notable task is Visual Abductive Reasoning (**VAR**), which, as mentioned before, focuses on *uncovering* hidden causal relationships in visual narratives, even when the connections between events appear non-obvious or disjoint [7]. This task challenges AI systems to reason holistically over visual scenes and synthesize contextual information to infer the most plausible underlying causes. Additionally, as a task closely related to VAR, Dense Video Captioning (**DVC**) seeks to *produce* rich, multi-sentence descriptions of untrimmed videos, which demands temporally grounded, comprehensive narrative understanding [7]. Some approaches address this task by first parsing events and then generating text, while others reverse the process or unify both stages.

Building upon the VAR and DVC tasks, the **VideoABC** benchmark introduces a challenging procedural visual reasoning task [20]. Designed to assess a model’s ability to *interpret* and explain physical processes in instructional videos, VideoABC emphasizes long-term dependencies and common-sense reasoning [20]. The primary task involves *selecting* the most plausible action or step to complete a visual procedure given an initial and final state. The secondary task, on the other hand, is a more demanding task that asks the model to *justify* why alternative choices are less plausible, thus discouraging superficial pattern matching. Notably, VideoABC avoids reliance on natural language inputs to provide a *purely* visual benchmark that requires genuine high-level reasoning over visual sequences [20].

Another foundational benchmark is Raven’s Progressive Matrices (**RPM**), along

with its variants Relational and Analogical Visual Reasoning (**RAVEN**) and Impartial RAVEN (**I-RAVEN**), which are modeled after classic IQ tests [18]. As mentioned before, these benchmarks assess abstract and analogical reasoning by requiring the model to *infer* the missing piece in a visual grid based on underlying rules of symmetry, transformation, or pattern continuation [21]. NeSy methods have shown strong performance on RPM tasks by combining visual perception with symbolic rule-based reasoning [18]. For instance, some approaches model the task as a form of probabilistic abductive reasoning, where solutions are inferred within a *symbolically* structured space constrained by prior background knowledge [1]. While this approach offers explainability and generalizability, it often incurs a high computational cost due to the complexity of the symbolic search space [21].

Beyond structured grid-based reasoning, Visual Question Answering (**VQA**) represents a more open-ended and linguistically grounded task, which challenges systems to *answer* natural language questions based on images through integrating of visual scene understanding, semantic reasoning, and multi-modal inference [19]. Successful models in this task must handle relational reasoning across spatial and semantic dimensions, combining features from both the visual input and the textual query [6]. However, early VQA benchmarks suffered from dataset biases and superficial *shortcuts* that allowed models to perform well without truly understanding the visual input [20]. To overcome these limitations, the Compositional Language and Elementary Visual Reasoning (**CLEVR**) dataset was introduced, featuring images of 3D objects and explicitly relational questions, where previous powerful models struggled with the relational aspects [6]. CLEVR has become a key benchmark for evaluating whether systems genuinely perform reasoning rather than exploit statistical patterns in language.

Finally, a benchmark tailored specifically for NeSy visual reasoning paradigms is Visual Sudoku Puzzle Classification (**ViSudo-PC**), which blends visual perception with symbolic relational constraints, requiring systems to *determine* whether a visually rendered Sudoku grid is correctly solved [22]. Unlike traditional numeric input, ViSudo-PC puzzles are built from images drawn from datasets such as Modified National Institute of Standards and Technology (**MNIST**), Extended MNIST (**EMNIST**), Kuzushiji MNIST (**KMNIST**), and Fashion MNIST **FMNIST**, which introduces visual variability and noise [10]. The benchmark supports both 4×4 and 9×9 grid sizes to test the model’s ability to integrate low-level visual recognition with high-level reasoning about Sudoku rules [22]. As such, ViSudo-PC represents a holistic NeSy challenge that demands *joint* learning of visual perception and reasoning.

2.3 Basic Visual Reasoning Frameworks

Building on the foundational concepts of NeSy AI and the tasks and benchmarks introduced previously, this section discusses several core frameworks for NeSy visual reasoning. Among these methods, one of the most influential neural components tailored for reasoning tasks is the Relation Network (**RN**) [6]. These networks are *modular* neural architectures designed explicitly for relational reasoning, which is an essential skill for many visual reasoning tasks [6]. When applied to datasets such as CLEVR, RNs significantly enhances performance by enabling models to reason about object relations, an issue that convolutional architectures alone often *struggled* to capture [6]. By augmenting perception pipelines with RNs, models can gain the ability to *implicitly* learn logical structures such as comparison, counting, and spatial relationships.

Moving from modular perception-enhancing components to integrated reasoning architectures, the Deep Symbolic Learning (**DSL**) framework represents a crucial step toward unifying perception and symbolic inference [13]. DSL enables *end-to-end* learning of symbolic representations directly from raw perceptual inputs, and simultaneously discovers the underlying symbolic rules [13]. By embedding discrete symbolic choices within differentiable neural layers, DSL offers a compositional approach to reasoning that supports generalization across tasks and domains [13]. This makes it particularly relevant in contexts where the structure of symbolic reasoning must emerge from perception, such as visual classification with latent symbolic structure [12]. However, while DSL integrates reasoning and perception tightly, another key challenge for hybrid systems is efficient training, especially when symbolic supervision is weak or indirect [12]. A practical and broadly applicable strategy is the use of Transfer Learning (**TL**) for the neural perception modules [12]. In this approach, the perception component of a NeSy system is *pre-trained* on the downstream task using standard supervision, before being integrated with the symbolic reasoning module [12]. This technique mitigates issues such as slow convergence and local minima by ensuring that the perception model already maps inputs to semantically informative representations [12]. It has shown consistent improvements in both performance and training efficiency across various NeSy setups and complex visual reasoning tasks [12].

Among more formalized integrations of symbolic logic and neural perception, Neural Probabilistic Soft Logic (**NeuPSL**) exemplifies a *general-purpose* NeSy framework that employs neural capabilities to extend Probabilistic Soft Logic (**PSL**), which is a statistical relational learning framework that, similarly to fuzzy FOL, represents logical rules as soft constraints using continuous truth values [10]. NeuPSL inherits PSL’s flexible probabilistic reasoning framework but *introduces* differentiable components that interface directly with raw data via deep neural networks [10]. By supporting joint learning and inference over symbolic rules and

perceptual inputs, NeuPSL demonstrates the utility of energy-based modeling for tasks like MNIST-Addition, where visual and symbolic aspects interact [10]. As part of the broader family of Neuro-Symbolic Energy-Based Models (**NeSy-EBMs**), NeuPSL contributes to a growing line of work seeking probabilistic formalism within hybrid reasoning [10].

Lastly, the challenge of symbol grounding in generative tasks is addressed by frameworks like Abductive Visual Generation (**AbdGen**), which *combine* neural visual generation models with logical rule learning via abductive reasoning [23]. This generative visual reasoning requires the system to *assign* semantic symbols to latent neural factors and to infer rules that guide the generation process [23]. AbdGen achieves this purpose through a combination of quantized abduction and contrastive meta-abduction, effectively grounding symbols in visual data while maintaining a logic-based generative structure [23]. This approach enables it to learn from *limited* data and supports precise and explainable generation based on symbolic conditions. [23]

2.4 Visual Reasoning with VSAs

Vector-Symbolic Architectures (**VSAs**) are computational models that *represent* both atomic and composite concepts using high-dimensional distributed vectors [18]. These representations are manipulated through a defined set of algebraic operations, such as binding (e.g., Hadamard product), unbinding, bundling (additive superposition), permutation, and cleanup via associative memory, which together allow for the combination of symbolic structure with distributed and connectionist processing [1]. VSAs support key cognitive properties like compositionality, transparency, and robustness, which makes them well-suited for tasks that require reasoning over structured representations [24].

Early applications of VSAs to visual reasoning focused on *analogical reasoning*, often assuming symbolic inputs *without* integrating visual perception [1]. These limitations have led to the development of Neuro-Vector Symbolic Architectures (**NVSAs**), which incorporate deep neural networks for visual input processing and use VSA-based mechanisms for reasoning [1]. The NVSA framework *transforms* raw images into fixed-width VSA vectors that preserve perceptual uncertainty, which enables symbolic-like reasoning while retaining neural flexibility [1]. Then, the vector-based backend enables efficient probabilistic reasoning, such as probabilistic abduction, without resorting to *exhaustive* symbolic search [21]. In doing so, NVSA models allow perceptual processes to be shaped by the demands of downstream reasoning [1].

Building on the above ideas, the Relational Reasoning with Symbolic and Object-Level Features Using VSA processing (**RESOLVE**) architecture enhances the

NVSA approach by introducing an *attention mechanism* in high-dimensional bipolar vector spaces [24]. RESOLVE encodes both object-level features and inter-object relations through VSA operations such as binding and bundling, while its attention mechanism improves relational extraction efficiency, which is a notable challenge for transformer-based models [24]. RESOLVE demonstrates *strong* generalization in both fully relational and mixed-relational visual reasoning benchmarks [24].

Another notable approach is Probabilistic Abduction for Visual Abstract Reasoning via Learning Rules in VSAs (**Learn-VRF**), which uses VSAs to learn rule-based formulations for solving abstract reasoning problems, specifically RPM [18]. Unlike prior models that rely on *hand-crafted rule* representations, Learn-VRF learns rule structures directly from data while maintaining symbolic explainability [21]. It performs well on Out-of-Distribution (**OOD**) samples, demonstrating generalization to unseen combinations of attributes and rules [21]. More recently, the Abductive Rule Learner with Context-awareness (**ARLC**) has also extended the Learn-VRF by introducing a new training objective *specifically* tailored for abductive reasoning [18]. ARLC not only learns the rules underlying abstract reasoning tasks but also allows for the incorporation of domain knowledge through programmatic constraints [18]. It addresses prior limitations of Learn-VRF, such as its limited rule expressiveness and sub-optimal selection mechanisms, offering a more flexible and context-aware framework for rule induction [18].

Expanding further, the Abductive Visual Generation (**AbdGen**) framework also integrates logic programming with neural generative models under an abductive learning paradigm [23]. While primarily designed for visual generation tasks, AbdGen employs symbolic inference techniques, such as quantized abduction, which uses *nearest-neighbor lookups* in semantic code books to ground symbolic hypotheses in perceptual data [23]. This demonstrates the potential of VSAs in facilitating symbol grounding and abductive reasoning in visually rich contexts [23].

2.5 Visual Reasoning with LTNs

Logic Tensor Networks (**LTNs**) are presented as a significant NeSy framework designed to effectively *integrate* deep learning capabilities with symbolic reasoning, particularly relevant for visual reasoning tasks that demand processing both rich perceptual data and abstract knowledge [5]. The foundational element of LTNs is Real Logic, which, as mentioned, facilitates learning, querying, and reasoning by grounding the elements of a FOL signature onto data through neural computational graphs [5]. This framework employs fuzzy FOL semantics to transform the typically binary constraints of classical logic into continuous, differentiable operations, ultimately allowing logical knowledge to function as a differentiable regularizer

in the loss function [12]. Consequently, LTNs offer a *unified* language capable of representing and computing a variety of AI tasks crucial for visual understanding and reasoning, including multi-label classification, relational learning, embedding learning, and query answering [5]. Notably, LTNs can perform reasoning on combinations of axioms not explicitly trained upon and provide high explainability, which are valuable characteristics when tackling complex visual scenarios [3].

A compelling demonstration of LTNs applied to visual reasoning tasks is explored in their application to the ViSudo-PC benchmark [8]. An LTN-based approach addresses this by integrating Convolutional Neural Networks (**CNNs**), which act as the perceptual module for tasks like digit recognition or whole-board classification, with an LTN responsible for *enforcing* the logical constraints of Sudoku [8]. Various methods for this integration and for formulating the logical constraints are investigated, which include **Indirect Solutions**, where a non-trainable predicate encodes the Sudoku rules to verify digits detected by a trainable predicate, and **Direct Solutions**, where a predicate is involved that directly calculates the validity of the entire board and, along with an auxiliary digit detection and rule enforcement predicate, is grounded by CNNs (either separate or with a shared backbone) [8]. This dual-level integration of perceptual recognition and logical reasoning exemplifies how LTNs can effectively *bridge* sub-symbolic and symbolic components in visual reasoning tasks [5].

2.6 Visual Reasoning with LLMs

As introduced in Section 1, LLMs have shown notable promise in performing reasoning tasks across multiple domains [14]. Extending this ability into the realm of visual reasoning, we can mention Vision-Language Models (VLMs), which combine natural language understanding with visual perception [14]. These models can be adapted for visual reasoning by *encoding* visual content into structured symbolic or textual representations [21]. In many current approaches, the raw visual data is preprocessed externally and transformed into a format that the LLM can interpret, such as object attributes, spatial relationships, or transformation rules [14]. Once this symbolic encoding is complete, the VLM is *prompted* to reason over the visual structure. This setup has enabled zero-shot performance on abstract reasoning benchmarks like RPMs, where VLMs can identify visual patterns and analogies without task-specific fine-tuning [21]. While the model itself does not directly interpret pixels, its general language-based reasoning abilities can be leveraged through careful design of symbolic inputs and prompts [3].

Recent research has further explored how post-training methods influence the reasoning capabilities of LLMs and VLMs, particularly comparing Supervised Fine-Tuning (**SFT**) and RL: while SFT encourages memorization of training examples,

RL fosters the emergence of generalized reasoning skills [14]. For instance, in visual tasks like RPMs, SFT-tuned models often struggle with OOD examples, whereas RL-trained models demonstrate improved adaptability and abstraction [14]. This insight is central to the development of systems like DeepSeek-R1, which applies RL to incentivize structured reasoning [2]. The DeepSeek-R1-Zero model, notably trained without prior SFT, reveals that reasoning capabilities can emerge solely from reinforcement signals [2]. Meanwhile, the full DeepSeek-R1 pipeline, which combines SFT and RL, balances stability and performance, guiding LLMs toward generating more accurate and interpretable Chain-of-Thought (**CoT**) reasoning [2].

Beyond DeepSeek, other works such as **MathPrompter** and **Auto-CoT** represent earlier or parallel efforts to enhance LLM reasoning through prompting strategies [25]. MathPrompter augments reasoning by retrieving relevant equations and concepts to scaffold mathematical problem solving, while Auto-CoT automates the generation of reasoning chains to improve model performance across logical tasks [16]. Though these methods primarily target textual reasoning, they share conceptual ties with visual reasoning when symbolic representations of images are treated analogously to structured text [21]. Together, these methods underscore a broader trend: effective visual reasoning with LLMs often hinges on external symbolic encoding and internal alignment via prompt engineering, fine-tuning, or reinforcement [21, 19, 14]. The ongoing research reflects a shift from passive language modeling to active reasoning architectures, where models are trained not just to generate plausible text, but to follow reasoning trajectories across modalities [25, 21].

Chapter 3

Methodology

In this chapter, we discuss our methodology through four main modules: rule generation via prompt engineering and symbolic formulation, rule verification using syntactic parsing and LTN implementation, visual processing using CNNs, and feedback handling with the help of looped training mechanisms.

3.1 Proposed Method

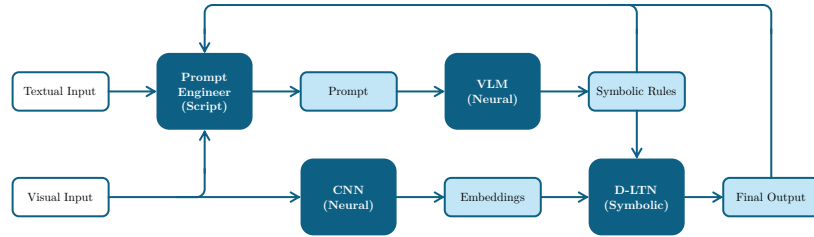


Figure 3.1: Functional diagram of the proposed NeSy system.

As discussed in Section 1.7, our objective is to overcome the limitations of the methods presented in the previous chapter by emphasizing flexibility, explainability, and formality. To this end, we transition from the high-level conceptual design shown in Figure 1.1 to a more concrete functional architecture illustrated in Figure 3.1. This refined design introduces the following key components:

- A **Prompt Engineer** serves as the context generator, combining the textual and visual inputs with intermediate feedback to construct a coherent **Prompt**. This prompt is then fed to the next processing block as the required context.
- A **VLM** acts as the rule generator. The idea behind this choice is to leverage

the advanced reasoning capabilities of these modern language models for generating symbolic rules. Additionally, to obtain formal rules, we prompt the VLM to compose them in FOL and Python, which are necessary for constructing the symbolic knowledge base and performing logic-based reasoning.

- A **CNN** functions as the visual processor by converting the visual input into **Embeddings**, which are directly interpreted as visual symbols. By using a CNN, we map the images from a two-dimensional complex space into a simpler space of embedding vectors.
- An Dynamic LTN (**D-LTN**) performs rule verification. With the help of this block, the generated symbolic rules are automatically converted into an LTN, which is then tasked with applying the extracted rules to the embeddings.

In the following, we will explain the complexities and the detailed operations performed by each block of the functional diagram.

3.2 Rule Generation

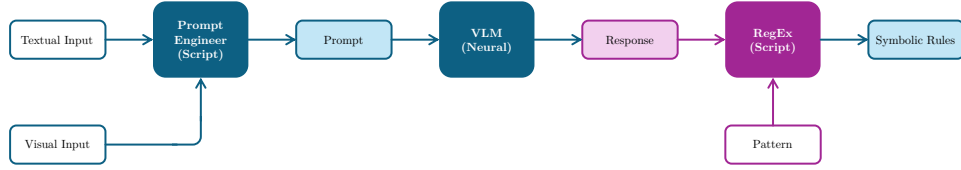


Figure 3.2: Reasoning with the help of VLM.

This section outlines the complete flow from the system inputs to the prompt engineering stage, through the VLM, and ultimately to the generation of symbolic rules. As illustrated in Figure 3.2, the process involves three additional components, the raw **Response** produced by the VLM, a Regular Expression (**RegEx**) used to extract symbolic rules from the response, and a specific **Pattern** that guides the extraction process.

3.2.1 Prompt Engineering

The following script presents our initial prompt engineering approach. For clarity, the template follows the style recommended by OpenAI in their documentation for using VLMs via Python, where the `client` object is used to access their Application Programming Interface (**API**) [26].

```

system_role = '''
You are a helpful assistant that can extract the First-OrderLogic
(FOL) from images. The grammar of the FOL is as follows:
  1. Constants: <PLACE HOLDER>.
  2. Variable: <PLACE HOLDER>.
  3. Functions: <PLACE HOLDER>.
  4. Predicates: <PLACE HOLDER>.
  6. The symbols used for AND, OR, and NOT: `<PLACE HOLDER>`,
    `<PLACE HOLDER>`, and `<PLACE HOLDER>`, respectively.
  7. The symbols used for implication and equivalence:
    `<PLACE HOLDER>`, and `<PLACE HOLDER>`, respectively.
  8. The symbols for universal and existential quantifiers:
    `<PLACE HOLDER>` and `<PLACE HOLDER>`, respectively.
  9. Use parentheses for preserving operation precedence.
Act based on the following:
  1. Before FOL rule generation, deeply analyze the images.
  2. Consider that all the images must follow the same FOL rule.
  3. The FOL rule applies to the visual objects inside each image.
  4. At the end of your chain of thought, use the following
    template to present the extracted rules:
    ```JSON
 {
 "rule_1": "first possible rule",
 "rule_2": "second possible rule",
 ...
 }
    ```

  5. Then, provide the groundings of constants, functions, and
    predicates in the following template:
    ```Python
 the groundings
    ```
'''

prompt = [{
  'type': 'text',
  'text': 'These are images you can use as reference:'
}]
for base64_image in base64_image_list:
  prompt.append({
    'type': 'image_url',
    'image_url': {
      'url': f'data:image/png;base64,{base64_image}'
    }
  })

```

```

chat_completion = client.chat.completions.create(
    messages=[
        {'role': 'system', 'content': system_role},
        {'role': 'user', 'content': prompt}
    ]
)

response = chat_completion.choices[0].message.content

```

In this template, the placeholder '`<PLACE_HOLDER>`' indicates where specific information must be inserted into the `system_role` variable, which will be specified in the following sections once the grammar for the FOL language has been defined. Also, as shown in the template, the textual input of the system is provided via the `system_role` variable, while the visual input is passed to the VLM as part of the prompt, using images stored in the list `base64_image_list`. These images serve as reference examples from the training data, guiding the VLM in generating the corresponding symbolic rules. As previously mentioned, these rules are produced in two forms:

- **FOL.** According to the prompt, the VLM is first instructed to generate a JavaScript Object Notation (**JSON**) dictionary containing the possible FOL rules that describe the underlying relationships among the visual objects. The JSON format provides the flexibility for the VLM to extract and organize as many relevant rules as it identifies.
- **Python.** The VLM is also asked to produce a Python script containing the **Groundings**, which are the contextual meanings of the symbols appearing in the FOL rules. A grounding can be value or tensor assigned to a constant or a variable, a Python function or neural network represented by an FOL function or predicate, or a tensor operation performing the abstract idea of a universal quantifier. Even though every symbol of an FOL rule requires grounding in practice, the VLM is prompted to specify the groundings only for constants, functions, and predicates. We explain the details in subsequent sections.

The above approach assumes the VLM is already familiar with the concept of FOL and is capable of performing analytical reasoning. In fact, in our proposed architecture, the VLM effectively serves as the reasoning brain of the system, which allows it to operate with minimal hard-coded logic or rule-based intervention. Therefore, the reasoning process is delegated entirely to the VLM and the sophistication of its reasoning capabilities directly impacts the quality of the results. In other words, the more advanced the model, the more reliable, nuanced, and expressive the extracted rules will be.

3.2.2 Symbolic Rule Generation

The raw response generated by the VLM is returned in plain text, which most likely includes the CoT reasoning steps it followed before producing the symbolic rules. This inclusion is typical behavior of modern language models, especially when prompted to reason step-by-step or operate under system instructions that prioritize transparency and explainability. Also, the CoT reasoning is inevitable, as it provides insight into the intermediate steps the model uses to reach its conclusions. However, from a system integration perspective, this verbose output becomes noise. What we actually require for downstream processing are only the two formal parts of the response, the JSON dictionary that holds the extracted FOL rules, and the Python script that lists the groundings. To separate these relevant segments from the surrounding explanatory text, we employ RegEx, which is a robust pattern-matching tool widely used in text processing. RegEx allows us to define flexible yet precise search patterns that can locate and extract the relevant blocks within unstructured textual output. Therefore, to extract the symbolic rule components effectively, we define two RegEx patterns, one for the FOL rules and another for the groundings:

```
fol_pattern = r'```JSON\s*(\{[^\}]*\})\s*```'
python_pattern = r'```Python\s*([^\`]*)\s*```'
```

The first pattern, `fol_pattern`, searches for a block of text that begins with the literal prefix ````JSON`, followed by any amount of whitespace, and then captures a JSON object enclosed in curly braces `{...}`. The pattern ensures that it terminates at the matching closing brace before the final triple backticks `````. Similarly, the second pattern, `python_pattern`, captures the Python code block. It searches for the prefix ````Python`, then greedily captures all content until it encounters the closing backticks `````. These carefully crafted expressions ensure the correct segments of the response are isolated from any explanatory text or CoT reasoning.

3.3 Rule Verification

With the VLM response decoded into an JSON dictionary with possible FOL rules and a Python script defining the groundings, we are now able to convert those symbolic rules into proper formats that finally lead to the implementation of a D-LTN. In summary, as illustrated in Figure 3.3, we can use a **Rule Parser** to convert the symbolic rules, stated in a pre-defined **Grammar**, into a **Syntax Tree**, and use it to generate an **LTN Formula** with the help of an **LTN Builder** by converting the raw nodes of the syntax tree into the customized **Node Classes** we already defined. In this section, we discuss the details of this process.

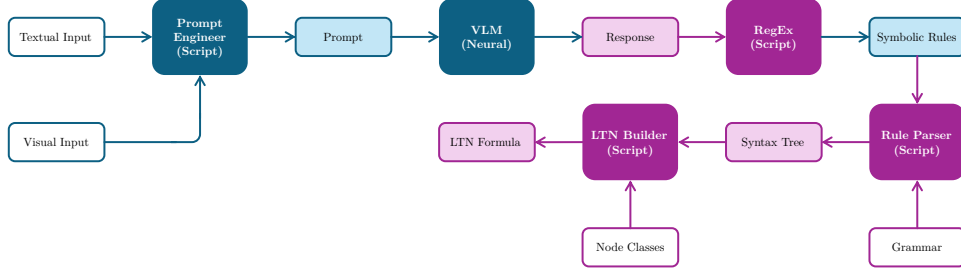


Figure 3.3: Automatic generation of LTN formulas.

3.3.1 Syntax Tree Generation

A syntax tree is a hierarchy of nodes that represents the abstract syntactic structure of a set of instructions written in a formal language, in our case, FOL or Python. Each node in the tree denotes a construct occurring in the logical expression or the source code. Accordingly, in our methodology, a syntax tree serves two distinct purposes:

- The syntax tree for the FOL rules represents their logical structure, with nodes corresponding to predicates, constants, variables, functions, logical connectives, and quantifiers. This tree provides a canonical and unambiguous representation of the rule, abstracting away superficial details of the specific textual syntax.
- Similarly, a syntax tree for the Python script represents its programmatic structure, with nodes corresponding to expressions, statements, function definitions, and variable assignments. This tree is crucial for programmatically accessing and interpreting the groundings that the VLM provides.

With the use of these syntax trees, we can create well-defined formulas for our D-LTN block, ensuring that the symbolic knowledge is precisely interpreted and integrated. Accordingly, to convert the extracted FOL rules into a machine-readable format suitable for our D-LTN implementation, we leverage **Lark**, which is a powerful and flexible parsing toolkit for Python that allows for building custom parsers [27]. For the Python grounding script, on the other hand, we simply make use of Python’s built-in **ast** module, which provides tools for working with syntax trees of Python source code directly. This simplifies the extraction of groundings by allowing direct traversal and inspection of the script’s structure. For the FOL rules, however, the parsing process is more intricate.

As mentioned, Lark is a powerful parser generator. It uses context-free grammars to define the structure of a language and employs various parsing algorithms (such as **LALR**, **Earley**, or **CYK**) to convert plain text into syntax trees [27].

This process is similar to how development environments like Visual Studio Code (**VSCode**) utilize internal syntax parsers or textual grammars (such as **TextMate** grammars) to transform a plain Python script into a colored, syntactically highlighted visualization, regardless of its semantic meaning. In fact, a parser's sole responsibility is to construct a syntax tree based on the defined grammar. It does not concern itself with the semantic correctness or logical validity of the content. Therefore, just as VSCode does not inherently validate the runtime behavior or logical soundness of a Python script, our FOL parser similarly focuses exclusively on the grammatical correctness of the symbolic rules. The semantic validation, which involves assessing whether a syntactically correct FOL rule accurately describes the relationships within a specific domain, is a separate and crucial step that will be investigated later.

To define a context-free grammar for the conversion of plain FOL expressions into syntax trees using Lark, we primarily focus on balancing the **generalizability** of the grammar and its inherent **formality**. In other words, we want to define a grammar robust enough to encompass a broad spectrum of common FOL rules while simultaneously enforcing a strict formal structure, ensuring that all generated rules consistently adhere to the same syntactic conventions. Also, this allows the parser to render the rules unambiguously machine-readable for seamless integration into our D-LTN. According to Lark documentations, we need to define this grammar in Extended Backus-Naur Form (**EBNF**) [27], which in our case, is as follows:

```
//// Explanations ////

// constant identifiers always start with "C"
// variable identifiers always start with "x"
// function identifiers always start with "f"
// predicate identifiers always start with "P"
// wrapper symbol is "(" and ")"
// negation symbol is "!"
// conjunction symbol is "&"
// disjunction symbol is "|"
// implication symbol is "implies"
// equivalence symbol is "iff"
// universal quantifier symbol is "forall"
// existential quantifier symbol is "exists"

//// Initialization ////

// imports
%import common.WS
%ignore WS
```

```

// entry point
?start: expression

//// Term-Level Terminal Definitions ////

// Tree Structure:
// term
// atom
// constant, variable
// mapper
// function

// Abstract Terminal (no precedence)
?term: constant | variable | function

// Concrete Terminals (no precedence)
constant: /C[a-z0-9_]*/
variable: /x[a-z0-9_]*/
function: /f[a-z0-9_]*/ "(" term ("," term)* ")"

//// Expression-Level Terminal Definitions ////

// Tree Structure:
// expression
// evaluator
// predicate
// unary_connective
// wrapper, logical_not
// binary_connective
// logical_and, logical_or, iff, implies
// quantifier
// exists, forall

// Abstract Terminal (ascending precedence)
?expression: level_0
?level_0: level_1 | exists | forall
?level_1: level_2 | iff | implies
?level_2: level_3 | logical_or
?level_3: level_4 | logical_and
?level_4: level_5 | logical_not | wrapper
?level_5: predicate

// Concrete Terminals (ascending precedence)
exists: "exists" variable ("," variable)* expression
forall: "forall" variable ("," variable)* expression

```

```

iff: level_1 "iff" level_2
implies: level_1 "implies" level_2
logical_or: level_2 "|" level_3
logical_and: level_4 "&" level_3
logical_not: "!" level_4
wrapper: "(" expression ")"
predicate: /P[a-z0-9_]* / "(" term ("," term)* ")"

```

The provided grammar defines the formal syntax for the FOL language used in the proposed system, along with a clear hierarchy of operator precedence. At the base level, constants, variables, and functions are identified by the prefixes **C**, **x**, and **f**, respectively. On top of them, predicates, marked by the prefix **P**, act as the core evaluable units. Next, logical operators are applied in a strict order of precedence: negation (!) binds most tightly, followed by conjunction (&), disjunction (|), then implication (implies) and equivalence (iff), and finally universal (forall) and existential (exists) quantifiers operating at the highest level and introducing scoped variable bindings. To override this default precedence and enforce specific groupings, the grammar also includes a wrapper using parentheses, which enables unambiguous parsing of nested logical forms. According to this grammar, we can now update the prompt used by the VLM:

```

system_role = '''
You are a helpful assistant that can extract the First-OrderLogic
(FOL) from images. The grammar of the FOL is as follows:
1. Constants: always starting with "C", e.g., "C", "C1" etc.
2. Variable: always starting with "x", e.g., "x", "x_2", etc.
3. Functions: always starting with "f", e.g., "f", "f_get", etc.
4. Predicates: always starting with "P", e.g., "P", "P_equal", etc.
6. The symbols used for AND, OR, and NOT: `&`, `|`, and `!`,
   respectively.
7. The symbols used for implication and equivalence: `implies`
   and `iff` respectively.
8. The symbols for universal and existential quantifiers: `forall`
   and `exists`, respectively.
9. Use parentheses for preserving operation precedence.
Act based on the following:
1. Before FOL rule generation, deeply analyze the images.
2. Consider that all the images must follow the same FOL rule.
3. The FOL rule applies to the visual objects inside each image.
4. At the end of your chain of thought, use the following
   template to present the extracted rules:
   ```JSON
 {
 "rule_1": "first possible rule",

```

---

```

 "rule_2": "second possible rule",
 ...
 }
 ...

5. Then, provide the groundings of constants, functions, and
 predicates in the following template:
   ```Python
   the groundings
   ```
'''

prompt = [{
 'type': 'text',
 'text': 'These are images you can use as reference:'
}]
for base64_image in base64_image_list:
 prompt.append({
 'type': 'image_url',
 'image_url': {
 'url': f'data:image/png;base64,{base64_image}'
 }
 })

chat_completion = client.chat.completions.create(
 messages=[
 {'role': 'system', 'content': system_role},
 {'role': 'user', 'content': prompt}
]
)

response = chat_completion.choices[0].message.content

```

### 3.3.2 LTN Implementation

Following the successful parsing of the symbolic rules, encompassing both the FOL interaction descriptions and their corresponding Python groundings, we can now proceed with their conversion into LTNs. To understand this process, it is essential to delve deeper into the fundamental nature of an LTN. At its core, an LTN represents a *fuzzified grounding* of an FOL rule, leveraging functions and neural networks to enable differentiability and trainability [5]. Therefore, the implementation of an LTN essentially involves translating this grounding process into a computational framework. For instance, consider an FOL rule such as  $\forall x_1 \exists x_2 : \text{Pred}(\text{func}(x_1), x_2)$ . The conversion starts by defining how the function `func` and the predicate `Pred` behave within a fuzzy logical space, which is actually

why we previously prompted the VLM to generate Python scripts containing these concrete groundings for constants, functions, and predicates. Subsequently, we must also define the fuzzy groundings for the logical quantifiers,  $\exists$  (existential) and  $\forall$  (universal). With these fuzzy groundings established for all components of the rule, an LTN is effectively constructed.

Consequently, given the groundings of  $x_1$  and  $x_2$ , which are the input embeddings of the D-LTN block, we can then *trigger* this instantiated LTN, similar to evaluating the truth value of the original FOL rule, but in a continuous and differentiable manner. However, given our system’s automatic nature, this conversion process simply cannot be manual. So, how do we transform the generated symbolic rules into an executable LTN formula? The answer lies within the syntax trees we have already discussed. In summary, we use Lark again to convert the FOL-based syntax trees into LTN-based syntax trees. Next, these LTN-based trees are grounded with the Python objects that we have prepared earlier.

The core idea behind the above approach is that traversing these newly formed LTN-based trees is functionally equivalent to feeding inputs directly into the LTN. In other words, the structure of the LTN-based syntax tree *is* the computational graph of the LTN itself. Considering our example rule,  $\forall x_1 \exists x_2 : \text{Pred}(\text{func}(x_1), x_2)$ , once converted, this rule will be represented as a tree where the root node is  $\forall x_1$ , its child is  $\exists x_2$ , whose child is  $\text{Pred}$ , and  $\text{Pred}$  then has  $\text{func}$  and  $x_2$  as its children, with  $\text{func}$  in turn having  $x_1$  as its child. When we *traverse* this tree, we start from the innermost node, which is  $x_1$ . Then, the node representing  $\text{func}$ , which is now a fuzzy, differentiable function defined in our Python groundings, receives the embedding for  $x_1$  as input and computes its output. After that, the node representing  $\text{Pred}$ , which is grounded by a fuzzy predicate, takes the computed output from  $\text{Pred}(\text{func}(x_1), x_2)$  and the embedding for  $x_2$  as its inputs, computing its fuzzy truth value. Next, the  $\exists x_2$  node, grounded as a fuzzy existential quantifier, processes the truth value from  $\text{Pred}$  by considering all possible instantiations of  $x_2$  in the domain to find the existentially-aggregated truth value. Finally, the  $\forall x_1$  node, grounded as a fuzzy universal quantifier, takes this result, considers all possible instantiations of  $x_1$ , and finds the universally-aggregated truth value.

The above sequential evaluation of an LTN-based syntax tree *is* the exact process of feeding the embeddings directly into the LTN to compute its final truth value for the rule. However, the fundamental question is how to we create an LTN-based syntax tree. Along with the ability to create custom parsers, Lark also provides transformers to convert raw syntax trees into customized ones, in our case, from an FOL-based syntax tree into an LTN-based syntax tree. For this purpose, we need to define custom node classes for the tree so that we achieve the functionality described in the previous example. To begin, we define the `Base` class as the superclass of all other nodes. Then, considering the grammar we defined previously, we proceed to define the remaining node classes according the hierarchy provided

in the grammar.

Below the **Base** class level, we also have other high-level superclasses, which describe the fundamental roles of the nodes in the LTN-based syntax tree by grouping them into abstract categories based on their functionality. These classes are as follows:

- **Term**: extends **Base**. This abstract class serves as a superclass for all nodes representing logical terms, such as variables, constants, and function applications. The instances of this class are tensors that are fed into predicates or other expressions.
- **Expression**: extends **Base**. This abstract class serves as a superclass for all nodes representing logical expressions, including predicates, connectives, and quantified formulas. The instances of this class are responsible for computing the truth values of the expressions they represents.

Also, below the above classes, we have middle-level superclasses, which provide more specialized abstractions that further refine the structure and semantics of the LTN-based syntax tree. Each class defines a logical category of behavior that will be inherited by concrete node types. They are as follows:

- **Atom**: extends **Term**. This abstract class represents the basic terms, i.e., constants and variables.
- **Mapper**: extends **Term**. This abstract class represents functions that map one or more terms to another term, including both built-in operations (e.g., vector summation) and user-defined functions.
- **Evaluator**: extends **Expression**. This abstract class represents logical predicates that evaluate terms to a truth value, including both built-in relations (e.g., equality, inequality) and user-defined predicates.
- **UnaryConnective**: extends **Expression**. This abstract class represents logical connectives that operate on a single expression, i.e., negations and groupings with parentheses.
- **BinaryConnective**: extends **Expression**. This abstract class represents logical connectives that operate on two expressions, i.e., conjunctions, disjunctions, equivalences, and implications.
- **Quantifier**: extends **Expression**. This abstract class represents quantified logical expressions that operate over domains of variables, i.e., universal or existential quantifiers.

Finally, the low-level concrete classes are the actual building blocks that instantiate the nodes of the LTN-based syntax tree. These classes implement specific logic or structure corresponding to the elements of FOL. They are as follows:

- **Constant**: extends **Atom**. It represents a constant term.
- **Variable**: extends **Atom**. It represents a variable term.
- **Function**: extends **Mapper**. It represents a user-defined function.
- **Predicate**: extends **Evaluator**. It represents a user-defined predicate.
- **Wrapper**: extends **UnaryConnective**. It represents a syntactic wrapper for grouping expressions using parentheses.
- **LogicalNot**: extends **UnaryConnective**. It represents negation of a single expression.
- **LogicalAnd**: extends **BinaryConnective**. It represents conjunction between two expressions.
- **LogicalOr**: extends **BinaryConnective**. It represents disjunction between two expressions.
- **Implies**: extends **BinaryConnective**. It represents implication between two expressions.
- **Iff**: extends **BinaryConnective**. It represents equivalence between two expressions.
- **ForAll**: extends **Quantifier**. It represents a universal quantifier.
- **Exists**: extends **Quantifier**. It represents an existential quantifier.

We designed the above classes as general as possible to enable Lark to convert any FOL-based syntax tree into an LTN-based syntax tree. However, in the following, we assume that every LTN formula is always derived from an FOL expression and never a term. For example, a term like  $f(x)$  is not an expression and it cannot serve as a valid LTN formula on its own. Therefore, the root of any LTN-based syntax tree must always be an expression, which ensures that the output of the LTN evaluation is always a truth value, and not a multi-dimensional tensor. Moreover, we intentionally did not implement concrete classes for built-in functions or predicates. If such functionality is required, it must be added explicitly by extending either **Mapper** for functions or **Evaluator** for predicates. Alternatively, possible users of these classes can define their own. For example, although the framework does not

include a built-in predicate for comparing variables (e.g., greater-than), a user could extend `Evaluator` to define a custom predicate class such as `P_gt` to implement this logic.

To finalize the LTN formula, we need to equip the `Base` class with **Groundability** and **Callability** capabilities. Groundability refers to the ability to assign specific groundings to the nodes, while callability allows input values to be passed to the formula for evaluation. By adding them to the base class, all subclasses automatically inherit these features. To implement them, we define a `ground` method that accepts grounding values as named arguments, and a `call` method that also accepts inputs via named arguments. Since the resulting LTN formula is structured as a tree, these methods are only invoked at the root node. In other words, for any LTN formula, we only need to call the `ground` method once to propagate groundings to all nodes, and similarly, we call the formula once to provide input values to the entire structure.

It is worth noting that both methods described above serve the purpose of grounding the LTN formula, but in different ways. The `ground` method is used to ground all nodes that are static in nature, such as constants, functions, predicates, and any other node that is not an instance of the `Variable` class. In contrast, the `call` method is specifically used to ground variables. This separation is deliberate, as it highlights the dynamic behavior of variables in contrast to the fixed nature of other nodes. In other words, non-variable nodes are grounded once and remain unchanged throughout the evaluation, while variables may be grounded multiple times, reflecting their role as placeholders for changing input values. This dynamic nature is also the reason we have referred to variable groundings as *inputs*. Additionally, when defining the node classes, we can optionally assign default groundings to connectives and quantifiers. These elements are formal and domain-independent, so their semantics can be predefined. In contrast, constants, functions, and predicates are user-defined and domain-specific. Therefore, they must be explicitly grounded by the user or provided externally. This is the main reason we prompt the VLM to supply them along with the FOL rules.

## 3.4 Visual Processing

The next major component in our functional diagram is where the visual input is processed. As shown in Figure 3.4, we do not explicitly expand the diagram with additional visual processing details, but this section will describe the internal workings of the CNN block. Specifically, we explain how it processes raw visual input into embeddings and how these embeddings are subsequently used as visual symbols within our system.

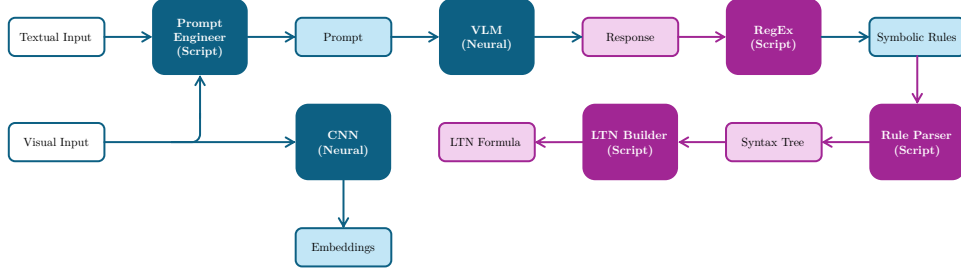


Figure 3.4: Perception with the help of CNN.

### 3.4.1 CNN Design

The visual input consists of a batch of images, each containing multiple visual objects. Assuming a batch size of  $B$ , a maximum of  $O$  objects per image, object dimensions of  $W \times H$ , and  $C$  color channels, the input tensor has the shape  $B \times O \times C \times W \times H$ . The CNN processes this input through a series of *convolutional* layers, followed by a *linear projection* that maps each object to an embedding of dimension  $E$ , resulting in an output tensor of shape  $B \times O \times E$ . As previously mentioned, one key benefit of this CNN block is that it transforms complex, high-dimensional visual data into a compact and semantically meaningful embedding space. Additionally, this visual encoder plays a foundational role in the system: it is responsible for perceiving and abstracting visual information, while the rest of the architecture focuses on symbolic reasoning over these representations. For this reason, it is commonly referred to in the NeSy literature as **Perceptor**.

The perceptor can be integrated into the system in two principal ways, depending on the desired trade-off between modularity and adaptability:

- **Modular Setting (Frozen Perceptor):** In this setup, the CNN is pre-trained independently on a task related to the visual domain, such as object classification or detection. Once trained, the perceptor is frozen and used as a fixed feature extractor during the reasoning phase. This approach can offer reduced training complexity and modularity, which allows researchers to analyze the symbolic reasoning component in isolation.
- **End-to-End Setting (Trainable Perceptor):** Here, the perceptor is not pre-trained but instead trained jointly with the reasoning module using supervision from final reasoning outcomes (e.g., logical inference labels or task-specific decisions). This setup allows the system to adapt visual representations to better align with symbolic tasks, potentially improving performance.

Both settings are valid and useful. The modular setting is often favored when pre-trained vision models are available. The end-to-end setting, on the other

hand, is more suitable when domain-specific visual features are not easily captured by existing models, or when tight coupling between perception and reasoning is necessary for optimal performance.

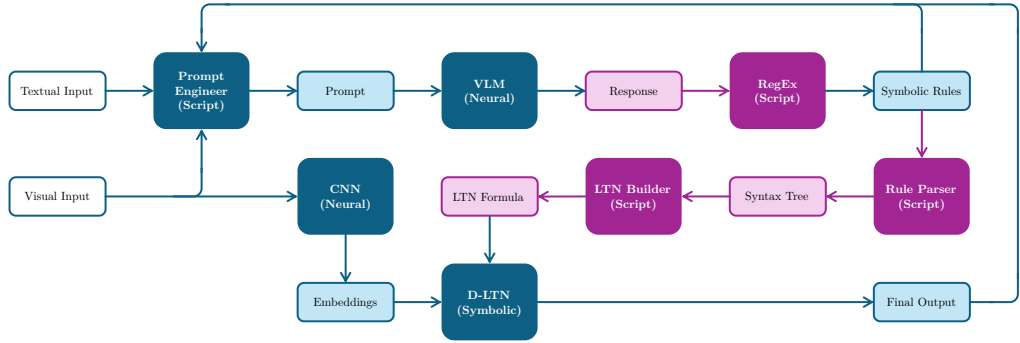
### 3.4.2 Visual Symbol Definition

As discussed earlier, the embeddings produced by the CNN are used directly as visual symbols. These embeddings encode contextual properties of each object, such as shape, color, and texture, into a tensor of shape  $B \times O \times E$ . However, in addition to these semantic features, each visual object is also associated with positional information, such as bounding boxes or spatial indices, which play a crucial role in visual reasoning. To incorporate this information into the embeddings, we can consider two approaches:

- **Projection-Based Fusion:** Apply a separate learnable projection layer to the positional features to map them into the same embedding space, and then add the result to the semantic embeddings.
- **Concatenation-Based Fusion:** Directly concatenate the positional features to the semantic embeddings, resulting in an augmented embedding vector.

Both strategies aim to enrich the visual symbols with spatial context, enabling the reasoning component to exploit both semantic and positional cues.

## 3.5 Feedback Handling



**Figure 3.5:** Detailed functional diagram of the proposed NeSy system.

Figure 3.5 contains the full functional diagram of the proposed system with its final feedback loop. As discussed, this loop is to provide the prompt engineer with the intermediate feedback both from generated symbolic rules and the final output of the D-LTN. In this section, we will discuss how this feedback is prepared.

### 3.5.1 LTN Development

Given our indirect approach to implementing LTNs (traversing the LTN-based syntax trees via the `ground` and `call` methods), we define a D-LTN as trainable if it contains at least one node that is grounded by a trainable object such as a neural network. For example, consider the formula  $\forall x_1 \exists x_2 : \text{Pred}(\text{func}(x_1), x_2)$ . If either `Pred` or `func` is grounded by a neural model within the corresponding LTN-based syntax tree, then the entire D-LTN becomes trainable. This capability allows us to optimize the D-LTN and shape its behavior through learning.

Therefore, to train such a system, we require a dataset consisting of appropriate textual inputs and corresponding labeled visual inputs. In addition, in order to evaluate the model's performance, we must also define a suitable metric that reflects its symbolic reasoning capabilities grounded in visual perception. Since our system is designed to extract and verify symbolic rules from images, we focus the evaluation on binary decision-making, which is determining whether the given label for an image is correct. This approach aligns with the semantics of fuzzified FOL rules, which produce a truth value  $t \in [0,1]$ . By setting a decision threshold  $\tau \in [0,1]$ , we can interpret the output as follows as if  $t \geq \tau$ , the system considers the rule to be satisfied, otherwise, it does not.

### 3.5.2 Loop Creation

Once the final output of the system is produced, the learning loop can be completed. At this stage, the system generates intermediate feedback consisting of the most recently derived symbolic rules along with the corresponding performance evaluation. This feedback is then incorporated into the next prompt, enabling the system to iteratively refine its symbolic understanding and improve performance over time. This prompt is what we have provided below, where '`<PREVIOUS_FOL_RULE>`', '`<PREVIOUS_GROUNDINGS>`', '`<METRIC_NAME>`', and '`<METRIC_VALUE>`' denote the most recent FOL rule with the highest performance, the most recent groundings, the evaluation metric used to evaluate the system, and the numerical performance score associated with it.

```
system_role = '''
You are a helpful assistant that can extract the First-OrderLogic
(FOL) from images. The grammar of the FOL is as follows:
1. Constants: always starting with "C", e.g., "C", "C1" etc.
2. Variable: always starting with "x", e.g., "x", "x_2", etc.
3. Functions: always starting with "f", e.g., "f", "f_get", etc.
4. Predicates: always starting with "P", e.g., "P", "P_equal", etc.
6. The symbols used for AND, OR, and NOT: `&`, `|`, and `!`,
 respectively.
7. The symbols used for implication and equivalence: `implies`
```

```

 and `iff` respectively.
 8. The symbols for universal and existential quantifiers: `forall`
 and `exists`, respectively.
 9. Use parentheses for preserving operation precedence.
Act based on the following:
 1. Before FOL rule generation, deeply analyze the images.
 2. Consider that all the images must follow the same FOL rule.
 3. The FOL rule applies to the visual objects inside each image.
 4. You performed the same operation previously, where you extracted
 <PREVIOUS_FOL_RULE>
 and these groundings in Python
 <PREVIOUS_GROUNDINGS>
 where you achieved <METRIC_NAME> at <METRIC_VALUE>
 5. At the end of your chain of thought, use the following
 template to present the extracted rules:
        ```JSON
        {
            "rule_1": "first possible rule",
            "rule_2": "second possible rule",
            ...
        }
        ```

 6. Then, provide the groundings of constants, functions, and
 predicates in the following template:
        ```Python
        the groundings
        ```

'''

prompt = [{
 'type': 'text',
 'text': 'These are images you can use as reference:'
}]
for base64_image in base64_image_list:
 prompt.append({
 'type': 'image_url',
 'image_url': {
 'url': f'data:image/png;base64,{base64_image}'
 }
 })

chat_completion = client.chat.completions.create(
 messages=[
 {'role': 'system', 'content': system_role},
 {'role': 'user', 'content': prompt}
]
)

```

```
]
)

response = chat_completion.choices[0].message.content
```

## Chapter 4

# Experimental Environment

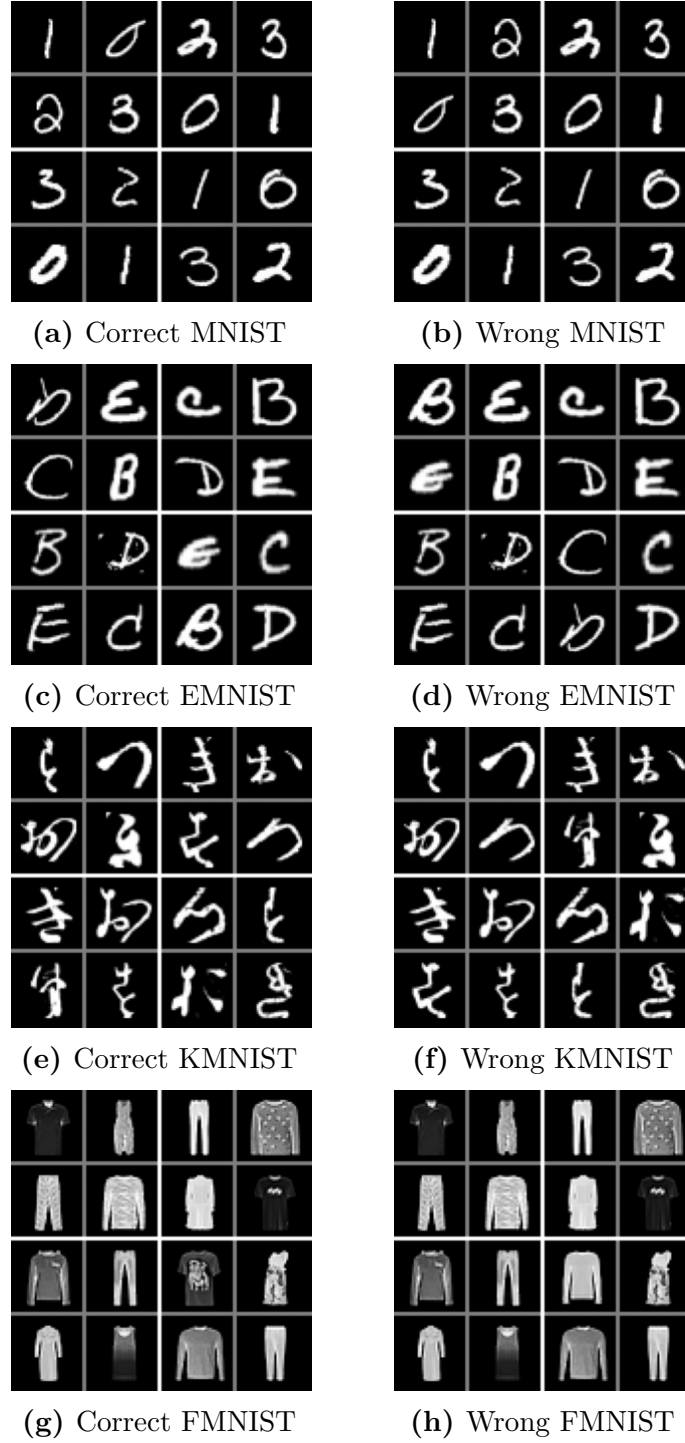
In this chapter, we describe the experimental setup, including the evaluation benchmark, datasets, programming environment, system specifications, development strategy, and evaluation metrics used to validate our proposed approach.

### 4.1 Setup

In this section, we detail the benchmark task and datasets used to evaluate our model, along with the programming environment and software tools that supported the development and experimentation process.

#### 4.1.1 Benchmark and Dataset

To assess the performance of our system, we utilized the **ViSudo-PC** benchmark. As described in Chapter 2, this benchmark evaluates a system’s ability to verify the correctness of visually encoded Sudoku boards. The dataset comprises 11 splits: the first 10 are used for scoring, while the 11-th is reserved for experimentation. Each split can be configured to contain either  $4 \times 4$  or  $9 \times 9$  boards. Additionally, each split includes its own training (100 pairs), validation (100 pairs), and test (100 pairs) subsets. As illustrated in Figure 4.1, each pair consists of two visually identical boards, except that the second board is deliberately corrupted to violate the Sudoku rules. Consequently, each subset provides 200 instances: 100 labeled as correct and 100 as incorrect [22]. Moreover, each cell in every board is represented by a  $28 \times 28$  grayscale image containing a handwritten digit (from MNIST), a handwritten letter representing a digit (from EMNIST), a Japanese character used as a digit (from KMNIST), or a fashion item standing in for a digit (from FMNIST). As discussed, this diverse visual representation requires the system to simultaneously address both perception and reasoning.



**Figure 4.1:** First pair of the training subset of the 11-split of the benchmark.

It is worth noting that our method to access the dataset was guided by the official repository of the benchmark [28]. Specifically, we used the following command:

```
python <ROOT>/generate-split.py
 --dataset <DATASET>
 --split <SPLIT>
 --out-dir <OUT_DIR>
```

where the capitalized placeholders, respectively, denote the directory containing the script, the data source to be accessed (`mnist`, `emnist`, `kmnist`, and `fmnist`), the split index (from 1 to 11), and the destination path for the generated data. The arguments mentioned in this commands, however, are the ones we explicitly set to access the dataset. The other arguments were left at their default settings. All other available arguments were left at their default values, which are listed below:

- **corrupt-chance**: The probability of applying an additional corruption after one has already been made, with 0.5.
- **dimension**: The size of the square puzzle grid, with 4 as the default value.
- **force**: If set to 1, allows overwriting of existing output directories, with 0 as the default value.
- **num-train**: The number of correct training puzzles to generate, with an equal number of incorrect puzzles also created, with 100 as the default value.
- **num-valid**: Similar to the previous argument but for validation.
- **num-test**: Similar to the previous argument but for test.
- **overlap-percent**: The fraction of additional puzzles generated by overlapping with the base dataset, with 0.0 as the default value.
- **seed**: A random seed to ensure reproducibility, with `None` as the default value.
- **strategy**: The strategy of puzzle generation, with `simple` as the default value.

### 4.1.2 Programming Tools

The experiments were conducted using the computational resources provided by **Google Colaboratory**, which is a cloud-based Jupyter notebook environment offering access to Graphical Processing Units (**GPUs**) such as NVIDIA T4 and storage space for efficient model training and inference [29]. The implementation was developed in **Python 3.11.13** with the use of **PyTorch 2.6.0** for tensor operations and neural network computations due to its flexibility and support for GPU

acceleration [30]. For the VLM, we utilized `llama-4-maverick-17b-128e-instruct`, which is a multi-modal language model that is composed of 17 billion activated parameters and supports a 128K token context window. This VLM has demonstrated robust performance in text and image reasoning tasks and aligns well with the system’s requirements for scalable and responsive rule generation [31]. We accessed it via an API key provided by Groq, which is a high-throughput inference engine running on cloud [32].

## 4.2 Specifications

Here, we describe the technical specifications and architectures of our key system components, including the VLM, the D-LTN, and the CNN block responsible for visual feature extraction.

### 4.2.1 VLM

During preliminary experiments, we identified the need to impose certain constraints on the prompt design. Initially, we allowed the VLM to generate Python-based groundings directly, as described in Section 3.2. However, we observed that without guidance, the VLM struggled to infer the correct FOL rules, even after multiple iterations. To address this, we restructured the prompt to present a predefined set of grounding alternatives, from which the VLM could choose. This modification improved rule accuracy and convergence. Consequently, the symbolic rules generated by the VLM were restricted to expressing only the relations between visual objects in FOL, with all groundings (except for variables) predefined in advance. We also instructed the VLM to generate only a single FOL rule per iteration, thereby reducing the complexity of the rule verification process and enhancing system stability. Accordingly, we used the following prompt:

```
system_role = '''
 You are a helpful assistant that can extract the First-Order
 Logic (FOL) rule from images.
 THE GRAMMAR OF FOL:
 - Constants: Not allowed in the rule.
 - Variables: Your options are `x1`, `x2`, ..., which
 represent visual objects.
 - Functions: Not allowed in the rule.
 - Predicates: Your options are `P_same_row`, `P_same_col`,
 `P_same_block`, `P_same_loc`, and `P_same_value`.
 - To compare variables, only use predicates.
 - The symbols used for logical AND, OR, and NOT are
 respectively `&`, `|`, and `!`.
```

- The symbols used for implication and equivalence are respectively ``implies`` and ``iff``.
  - The symbols used for universal and existential quantifiers are respectively ``forall`` and ``exists``.
  - Use parentheses for preserving operation precedence.
- WHAT YOU MUST CONSIDER:
- Use your own knowledge to analyze and deeply think about the images provided as your reference.
  - All the images must follow the same rule that you extract.
  - The rule applies to the visual objects within each image.
  - The visual objects may represent numbers rather than what they really are.
  - At the end of your chain of thought, put the extracted rule in the following template:  
`EXTRACTED_RULE: "the rule you extracted"`

```

'''
if len(history_list) > 0:
 n_extracted_rules = 0
 system_role += 'HISTORY OF PREVIOUS TRIALS:'
 for trial, incident in enumerate(history_list):
 error_message, extracted_fol_rule, ratio = incident
 system_role += (
 f'- Trial {trial+1} -> '
)
 if error_message != '':
 system_role += (
 f'error: "{error_message}"'
)
 else:
 n_extracted_rules += 1
 system_role += (
 f'extracted rule: "{extracted_fol_rule}", '
 f'conforming images: {100 * ratio:.2f}%'
)
 if ratio < termination_threshold:
 system_role += 'IMPORTANT LESSON FROM HISTORY:'
 if n_extracted_rules == 0:
 system_role += (
 '- Pay attention to the the instructions!'
)
 else:
 system_role += (
 '- The next FOL rule must be an improved version'
 ' of the above!'
)

```

```

prompt = [{
 'type': 'text',
 'text': 'These are the reference images:'
}]
for base64_image in base64_image_list:
 prompt.append({
 'type': 'image_url',
 'image_url': {
 'url': f'data:image/png;base64,{base64_image}'
 }
 })

chat_completion = client.chat.completions.create(
 messages=[
 {'role': 'system', 'content': system_role},
 {'role': 'user', 'content': prompt}
]
)

response = chat_completion.choices[0].message.content

```

In the above prompt, while the VLM is asked to abstractly define the predicates in FOL, they are grounded in executable Python functions tailored to our benchmark. Each predicate captures a specific type of relationship between visual objects and is implemented using either exact or approximate similarity. Specifically, `P_same_row`, `P_same_col`, and `P_same_block` evaluate whether two objects belong to the same row, column, or a block respectively, `P_same_loc` verifies whether two objects occupy the same grid position, and `P_same_value` checks whether two visual objects represent the same semantic concept (e.g., digit or letter). This design ensures that semantic relationships are modeled flexibly, while the VLM decides how to insert them into the FOL rule that it extracts.

### 4.2.2 D-LTN

Since the prompt used for the VLM eliminated the need for a Python script to define groundings, we instead needed to define them manually. For this reason and to ensure fair comparison against state-of-the-art methods, we adopted the same set of groundings as in [8]. These groundings are detailed as follows:

- `P_same_row`: grounded via binary similarity, this predicate checks whether two visual objects share the same row attribute:

$$P_{\text{same\_row}}(x_1, x_2) := 1_{x_1^{\text{row}}=x_2^{\text{row}}} \quad (4.1)$$

- **P\_same\_col**: grounded via binary similarity, this predicate evaluates whether two visual objects share the same column attribute:

$$P_{\text{same\_col}}(x_1, x_2) := 1_{x_1^{\text{col}}=x_2^{\text{col}}} \quad (4.2)$$

- **P\_same\_block**: grounded via binary similarity, this predicate assesses whether two visual objects belong to the same block:

$$P_{\text{same\_block}}(x_1, x_2) := 1_{x_1^{\text{block}}=x_2^{\text{block}}} \quad (4.3)$$

- **P\_same\_loc**: grounded via binary similarity, this predicate checks whether two visual objects share the same location:

$$P_{\text{same\_loc}}(x_1, x_2) := 1_{x_1^{\text{loc}}=x_2^{\text{loc}}} \quad (4.4)$$

- **P\_same\_value**: grounded via exponential similarity, this predicate measures the similarity of the contents of two visual objects:

$$P_{\text{same\_value}}(x_1, x_2) := \exp\left(-\text{relu}\left(\|x_1^{\text{value}} - x_2^{\text{value}}\|_p\right)\right) \quad (4.5)$$

- **Wrapper**: grounded by the identity function, i.e.,

$$g_{\text{wrapper}}(b) := b \quad (4.6)$$

- **LogicalNot**: grounded by the complement function, i.e.,

$$g_{\text{not}}(b) := 1 - b \quad (4.7)$$

- **LogicalAnd**: grounded using Goguen's product t-norm:

$$g_{\text{and}}(b_1, b_2) := b_1 b_2 \quad (4.8)$$

- **LogicalOr**: grounded using Goguen's t-conorm:

$$g_{\text{or}}(b_1, b_2) := 1 - (1 - b_1)(1 - b_2) = b_1 + b_2 - b_1 b_2 \quad (4.9)$$

- **Implies**: grounded using Reichenbach's implication:

$$g_{\text{implies}}(b_1, b_2) := 1 - b_1 + b_1 b_2 \quad (4.10)$$

- **Iff**: grounded using linear similarity:

$$g_{\text{iff}}(b_1, b_2) := 1 - |b_1 - b_2| \quad (4.11)$$

- **ForAll**: grounded via the power mean of the complements:

$$g_{\text{iff}}(b_1, b_2, \dots, b_k) := 1 - \sqrt[p]{\frac{1}{k} \sum_{i=1}^k (1 - b_i)^p} \quad (4.12)$$

- **Exists**: grounded using the power mean:

$$g_{\text{exists}}(b_1, b_2, \dots, b_k) := \sqrt[p]{\frac{1}{k} \sum_{i=1}^k b_i^p} \quad (4.13)$$

### 4.2.3 CNN Block

Regarding our CNN block, we employed a neural network composed of a stack of **Convolutional** layers followed by a stack of **Fully Connected** layers. Each stage of the convolutional stack begins with a 2D convolutional layer, where the number of output channels and the kernel size are respectively determined by the `cnn_dims` and `kernel_dims` hyperparameters. This convolutional layer is followed by a **ReLU** activation to introduce non-linearity, a **Group Normalization** layer with 2 groups to promote stable training, and a 2D **max pooling** operation with a kernel size of 2 for spatial downsampling. To improve generalization, a 2D **dropout** layer is also applied at the end of the stack, with dropout probability specified by the `drop_prob` parameter.

The resulting feature maps are then flattened and passed to the fully connected stack. Each stage of this second stack consists of a linear projection layer configured according to the `embed_dims` hyperparameter. With the exception of the final layer, each projection is followed by a ReLU activation and a dropout layer with the same dropout probability as that of the previous stack. In addition, depending on the `use_softmax` flag, the final stage may optionally include a **softmax** activation. This pipeline transforms each visual symbol into a compact embedding that captures both local visual features and high-level semantic abstractions.

To preserve spatial information, we also generate positional indices for each board cell using the **Cartesian product** of row and column coordinates. While these positional coordinates remain separate from the visual embeddings within the CNN block, we later integrate them using the concatenation-based fusion strategy described in Section 3.4. Specifically, the visual embedding for each object is concatenated with its corresponding positional coordinates to form the final representation.

## 4.3 Development Strategy

This section explains our development approach, focusing on the design of system loops for iterative training and feedback, as well as the strategies employed for effective data utilization during the learning process.

### 4.3.1 System Loops

The primary process driving our system is the **Generation Loop**, which is a term we introduce here to refer to the overall iterative cycle described in the previous chapters and illustrated in Figure 3.5. As discussed before, in each iteration of this loop, the prompt engineer combines the textual and visual inputs with the intermediate feedback and feeds the result into the VLM, which generates candidate

FOL rules. These rules are then used to configure the D-LTN, which attempts to verify them against the visual embeddings produced by the CNN block. For each iteration of this loop, we randomly sampled a small portion of the training subset from a given split of the ViSudo-PC dataset to serve as a reference for the VLM. The full training and validation subsets of the split were then used to develop the CNN block and D-LTN, and the test subset was reserved for the final performance evaluation and early stopping the generation loop in case no further improvement was observed.

As we can see, nested within the generation loop is the **Development Loop**, which handles the training and validation of the CNN block and D-LTN. Similar to the generation loop, we also implemented early stopping for this loop, but with the help the validation subset, as done typically in ML workflows. Accordingly, the development loop can be executed independently of the VLM to evaluate the CNN block and D-LTN independently. In such cases, FOL rules must be supplied to the D-LTN, either created manually or taken from previous VLM outputs, using the grammar we defined in the previous chapter. This setup allows for repeated and controlled testing of the system components without re-engaging the VLM.

### 4.3.2 Data Usage

To ensure compliance with the ViSudo-PC benchmark and maintain the integrity of our evaluation, we restricted the generation loop to operate exclusively on the 11-th data split. However, as this loop is centered on generating FOL rules, it was necessary to first tune the hyperparameters of the CNN block and D-LTN prior to rule generation, for which we also used the 11-th split. All other dataset splits were reserved for our experimental tests to enable a fair and standardized comparison of our CNN block and D-LTN against state-of-the-art methods evaluated within the ViSudo-PC benchmark. It is worth noting that this benchmark is designed not to assess rule generation but rather to evaluate the ability of NeSy systems to verify the correctness of synthesized Sudoku boards.

## 4.4 Evaluation Metrics

In this section, we present the loss function employed to guide training and the performance metrics used to quantitatively assess the reasoning capabilities and overall effectiveness of our model.

### 4.4.1 Loss Function

Since the CNN block is a neural network, the development loop requires a loss function to guide training. However, given that the primary focus of this thesis is

not visual perception and that the ViSudo-PC benchmark only permits supervision based on the final truth value produced by the D-LTN, we adopted an end-to-end evaluation strategy, as outlined in Section 3.4. In this strategy, rather than developing and evaluating the CNN block independently, we measure its effectiveness indirectly by observing how it contributes to the overall system performance. Accordingly, instead of designing complex loss functions, we opted for the simplest viable alternative, which was using the D-LTN’s truth value as the loss we want to reduce. Specifically, if we denote a single Sudoku board as  $x$ , its label as  $y \in \{0,1\}$ , and the composed function of the CNN block followed by the D-LTN as  $f$ , the loss function  $l$  can be defined as:

$$l(x, y) = y + (1 - 2y)f(x) = \begin{cases} f(x) & \text{if } y = 0, \\ 1 - f(x) & \text{if } y = 1 \end{cases} \quad (4.14)$$

#### 4.4.2 Performance Metrics

The loss function described in (4.14) penalizes the system proportionally to how far its output is from the correct truth value, and is minimized when the system assigns high confidence to correct evaluations. However, while this function is effective for supervising the system from the perspective of D-LTN rule verification, it does not fully capture the individual contributions of the VLM or CNN block. Therefore, to account for varying levels of perceptual fidelity and reasoning capability, we introduce the following set of complementary evaluation metrics:

- **VLM Load:** The number of training instances provided to the VLM during the generation loop. A lower value indicates better data efficiency and generalization capability of the VLM.
- **Total Iterations:** The total number of iterations required by the generation loop before convergence. Fewer iterations suggest quicker convergence and more effective rule generation.
- **Area Under the Curve (AUC):** The area under the **Receiver Operating Characteristic** (ROC) curve computed over D-LTN’s truth values. A higher AUC indicates better discriminative performance across different decision thresholds.
- **Accuracy:** The proportion of Sudoku boards for which the D-LTN correctly predicts the label, using a decision threshold of 0.5 on the produced truth value. This metric captures the overall effectiveness of the system in verifying board correctness based on the learned rules.

# Chapter 5

## Results

In this chapter, we present experimental outcomes from hyperparameter tuning, rule generation, and performance tests, followed by analyses and comparisons with state-of-the-art methods.

### 5.1 Experiments

This section covers the results of our experiments. First, we provide the details of our hyperparameter tuning process. Next, we present the rules obtained by the system. In the final part, we outline the results obtained in testing the generated rules on the CNN block and D-LTN assuming the VLM is bypassed.

#### 5.1.1 Hyperparameter Tuning

As discussed earlier, the rule generation loop depends on having both the CNN block and the D-LTN tuned in advance. To this end, we constructed the hyperparameter grid shown in Table 5.1, which outlines 192 different configurations tested for the CNN block’s hyperparameters introduced in the previous chapter. Regarding the D-LTN, while it is not a neural module, we needed to consider its interaction with varying datasets. For this reason, we evaluated each combination of hyperparameters through four separate runs of the development loop using the 11-th split of the Visudo-PC dataset with digits sourced from MNIST, EMNIST, KMNIST, and FMNIST, increasing the total number of separate runs to 768. The average test performance across these four sources was then used to identify the optimal configuration. Additionally, as previously discussed, the D-LTN in this experiment must be initialized with a manually crafted FOL rule that encodes the core semantic constraint of a valid Sudoku board. To this end, we focused on the fundamental principle that no two identical digits may appear in the same row,

column, or block, unless they refer to the same physical cell. To formally express this constraint, we designed a rule stating that if any two cells share the same value, then they must either correspond to the same location or not belong to the same row, column, or block. In other words, identical values are only permitted when referring to the same position; otherwise, their coexistence would violate Sudoku consistency. We believe this rule effectively enforces the uniqueness of values across structural dimensions and captures the logical essence of Sudoku validity. Considering the FOL grammar defined in Chapter 3, the rule we chose to initialize the D-LTN is as follows:

```
forall x1, x2
 P_same_value(x1, x2) implies
 P_same_loc(x1, x2) | (!P_same_row(x1, x2) &
 !P_same_col(x1, x2) &
 !P_same_block(x1, x2))
```

| Hyperparameter | Alternatives                   |
|----------------|--------------------------------|
| cnn_dims       | (8, 16), (16, 32), (32, 64)    |
| kernel_dims    | (4, 4), (2, 2)                 |
| embed_dims     | (4,), (16,), (64, 4), (64, 16) |
| drop_prob      | 0.1, 0.2, 0.3                  |
| use_softmax    | <b>False</b> , <b>True</b>     |

**Table 5.1:** Hyperparameters and possible alternatives for each one.

The hyperparameter search strategy defined for the CNN block and D-LTN provides a solid and exhaustive exploration of the relevant configuration space, as it spans diverse datasets, architectural setups, embedding sizes, regularization levels, and output activation strategies. This comprehensive setup enhances the chances of identifying a robust and generalizable model. Based on this strategy, Table 5.2 presents the results of our grid search, showing the top five configurations. While both the average test AUC and accuracy are reported by table, the rows are sorted by the first metric, as it is generally favored in the literature [8]. Consequently, we selected the configuration with the highest average test AUC, where the values for `cnn_dims`, `kernel_dims`, `embed_dims`, `drop_prob`, and `use_softmax` are (32, 64), (4, 4), (64,), 0.2, and **True**, respectively.

### 5.1.2 Rule Generation

With the CNN block and D-LTN tuned, we proceeded to the rule generation experiment. The results are presented in Table 5.3. For this experiment, as in the previous hyperparameter tuning experiment, we used the 11-th split of the

| ID | Hyperparameters |             |            |           |             | Average Test |          |
|----|-----------------|-------------|------------|-----------|-------------|--------------|----------|
|    | cnn_dims        | kernel_dims | embed_dims | drop_prob | use_softmax | AUC          | Accuracy |
| 1  | (32, 64)        | (4, 4)      | (64,)      | 0.2       | True        | 0.9560       | 85.50%   |
| 2  | (32, 64)        | (4, 4)      | (64, 4)    | 0.2       | True        | 0.9528       | 86.25%   |
| 3  | (32, 64)        | (4, 4)      | (64,)      | 0.1       | True        | 0.9523       | 86.12%   |
| 4  | (16, 32)        | (4, 4)      | (64,)      | 0.2       | True        | 0.9477       | 84.87%   |
| 5  | (16, 32)        | (4, 4)      | (4,)       | 0.2       | True        | 0.9458       | 82.50%   |

**Table 5.2:** Top five combinations of the hyperparameters.

Visudo-PC dataset and ran the generation loop separately on MNIST, EMNIST, KMNIST, and FMNIST. As shown, the loop successfully concluded with a rule for the first three sources, while it terminated with an error for FMNIST. We will analyze this issue in the next section, but it is worth noting here that the failure occurred due to the computational cost associated with the VLM, as we limited the loop to a maximum of 20 iterations. While it is possible that the loop could have succeeded on FMNIST with more iterations, we consider this limit acceptable, as the results already demonstrate the system’s ability to generate valid and semantically meaningful rules, which are detailed in the following.

1. **Rule 1:**

```
forall x1, x2
 !P_same_loc(x1, x2) & (P_same_row(x1, x2) |
 P_same_col(x1, x2) |
 P_same_block(x1, x2))
 implies !P_same_value(x1, x2)
```

2. **Rule 2:**

```
forall x1 forall x2
 (P_same_row(x1, x2) |
 P_same_col(x1, x2) |
 P_same_block(x1, x2)) & !P_same_loc(x1, x2)
 implies !P_same_value(x1, x2)
```

3. **Rule 3:**

```
forall x1, x2
 !P_same_loc(x1, x2) & P_same_value(x1, x2) implies (
 !P_same_row(x1, x2) &
 !P_same_col(x1, x2) &
 !P_same_block(x1, x2))
```

| Data Source                | VLM Load | Generated Rule | Total Iterations | Test   |          |
|----------------------------|----------|----------------|------------------|--------|----------|
|                            |          |                |                  | AUC    | Accuracy |
| MNIST-11 ( $4 \times 4$ )  | 3        | Rule 1         | 19               | 0.9974 | 00.99%   |
| EMNIST-11 ( $4 \times 4$ ) | 3        | Rule 2         | 13               | 0.9012 | 00.99%   |
| KMNIST-11 ( $4 \times 4$ ) | 3        | Rule 3         | 9                | 0.9517 | 00.99%   |
| FMNIST-11 ( $4 \times 4$ ) | 3        | Error          | 20               | -      | -        |

**Table 5.3:** Outcomes of running the rule generation loop.

### 5.1.3 Experimental Tests

Considering the hyperparameters and FOL rules obtained from our previous experiments, we initialized the D-LTN and ran the development loop across additional splits of the Visudo-PC dataset to evaluate our approach against other benchmark methods. Specifically, we conducted a comprehensive experiment using all three extracted rules, all ten dataset splits, and all four data sources (MNIST, EMNIST, KMNIST, and FMNIST), resulting in a total of 120 separate runs of the development loop. In the following section, we will aggregate the outcomes in accordance with standard practices in the literature and use them to compare the performance of our proposed method against state-of-the-art alternatives. For reference in this section, however, the seven best and seven worst non-aggregated outcomes based on the test AUC are also reported in Table 5.4.

| Data Source                | Utilized Rule | Training |          | Validation |          | Test   |          |
|----------------------------|---------------|----------|----------|------------|----------|--------|----------|
|                            |               | AUC      | Accuracy | AUC        | Accuracy | AUC    | Accuracy |
| MNIST-3 ( $4 \times 4$ )   | Rule 1        | 0.9988   | 99.50%   | 0.9699     | 90.50%   | 0.9949 | 91.00%   |
| MNIST-7 ( $4 \times 4$ )   | Rule 2        | 0.9988   | 99.00%   | 0.9665     | 90.00%   | 0.9917 | 93.50%   |
| MNIST-6 ( $4 \times 4$ )   | Rule 2        | 1.0000   | 100.00%  | 0.9856     | 94.00%   | 0.9892 | 97.00%   |
| MNIST-3 ( $4 \times 4$ )   | Rule 2        | 0.9988   | 98.50%   | 0.9668     | 92.00%   | 0.9862 | 92.50%   |
| MNIST-4 ( $4 \times 4$ )   | Rule 3        | 0.9975   | 99.50%   | 0.9800     | 94.00%   | 0.9849 | 94.50%   |
| MNIST-3 ( $4 \times 4$ )   | Rule 3        | 0.9975   | 99.50%   | 0.9752     | 92.50%   | 0.9846 | 92.50%   |
| MNIST-7 ( $4 \times 4$ )   | Rule 1        | 0.9988   | 99.50%   | 0.9791     | 91.00%   | 0.9845 | 94.00%   |
| FMNIST-9 ( $4 \times 4$ )  | Rule 2        | 0.9713   | 93.50%   | 0.8625     | 70.50%   | 0.8637 | 66.50%   |
| FMNIST-8 ( $4 \times 4$ )  | Rule 2        | 1.0000   | 99.50%   | 0.8725     | 77.00%   | 0.8571 | 74.00%   |
| EMNIST-4 ( $4 \times 4$ )  | Rule 3        | 0.5000   | 50.00%   | 0.5000     | 50.00%   | 0.5000 | 50.00%   |
| EMNIST-6 ( $4 \times 4$ )  | Rule 2        | 0.5000   | 50.00%   | 0.5000     | 50.00%   | 0.5000 | 50.00%   |
| EMNIST-10 ( $4 \times 4$ ) | Rule 2        | 0.5000   | 50.00%   | 0.5000     | 50.00%   | 0.5000 | 50.00%   |
| EMNIST-10 ( $4 \times 4$ ) | Rule 1        | 0.5000   | 50.00%   | 0.5000     | 50.00%   | 0.5000 | 50.00%   |
| FMNIST-3 ( $4 \times 4$ )  | Rule 3        | 0.5000   | 50.00%   | 0.5000     | 50.00%   | 0.5000 | 50.00%   |

**Table 5.4:** The seven best and worst outcomes of running the development loop.

## 5.2 Discussions

In our final section in this chapter, we will analyze the results of our experiments and compare our system against the state-of-the-art.

### 5.2.1 Analyses

Regarding hyperparameter tuning, Table 5.2 indicate that the convolutional stack performs best with larger layers, which is likely because larger kernels and more filters enable the model to capture more complex spatial patterns. Additionally, the optimal dropout probability in this layer is minimal, suggesting that the model does not require strong regularization at this point, possibly due to the structured nature of the visual data and the relatively low risk of overfitting early on. In contrast, the linear stack performs best with a single high-dimensional embedding layer. This suggests that deeper transformations beyond this point do not significantly improve performance and may even degrade it, likely because the essential visual features have already been extracted. It also implies that additional ReLU activations and dropout layers offer limited benefit at this stage, as introducing further non-linearity or regularization may disrupt rather than enhance the learned representation. Lastly, using a softmax activation at the output of the CNN block proved effective, likely because it stabilizes the feature distribution and emphasizes the most informative components, thereby facilitating downstream symbolic reasoning.

Turning to the rule generation experiment, Table 5.3 validates the effectiveness of our NeSy system in learning formal rules from textual and visual inputs. As established in Chapter 1, our goals were flexibility, explainability, and formality. The generated rules, which are expressed in FOL, directly address the latter two, providing interpretable outcome for the visual reasoning process in a formal structure. Regarding flexibility, although this chapter focused on the Sudoku consistency task, the functional architecture described in Figure 3.1 imposes no restrictions on other visual reasoning problems, provided they can be framed within the same formal framework. That said, the system is not without limitations. As shown in Table 5.3, while rule generation succeeded for MNIST, EMNIST, and KMNIST, it failed for FMNIST. This discrepancy likely stems from the nature of the visual data. The first three data sources contain digits and letters, which provide semantically meaningful patterns that the VLM can leverage to infer logical rules. In contrast, FMNIST consists of clothing items, whose visual features lack inherent symbolic meaning in the context of Sudoku, making it significantly more challenging for the system to extract consistent logical patterns.

Regarding the experimental tests on the development loop, the results in Table 5.4 generally demonstrate strong and consistent performance. Out of the 120

runs, which cover all combinations of data sources, benchmark splits, and extracted rules, only 5 outcomes (4.17%) exhibit noticeably low AUC or accuracy values in the test subset. While these under-performing cases indicate that the system is not flawless, they are not concentrated in any particular configuration. Instead, the failures are scattered across different data sources, splits, and rules, suggesting that the system’s performance is not overly sensitive to specific experimental conditions, which highlights the system’s robustness and generalizability. In other words, despite occasional deviations, the system consistently converges toward high-quality outcomes, demonstrating its ability to extract meaningful patterns and apply them effectively across diverse settings.

### 5.2.2 Comparison with Prior Literature

As previously discussed, benchmarking our proposed method requires aggregating the results from our experimental tests. Following the literature, we begin this process by grouping the results according to the data sources and computing the average system performance over the first ten splits of the Visudo-PC dataset [8]. However, since we evaluated three distinct rules for each data source and each split, a customized aggregation strategy was necessary. Specifically, for each data source, we computed the average performance across splits after selecting the average performance across the rules. This approach remains consistent with benchmarking conventions and allows us to showcase the best-performing configuration of our system when compared to state-of-the-art methods. The aggregated results are summarized in Table 5.5, which compares our method against the NeuPSL [10] and LTN-IND [8], with the latter tested under three variations A, B, and C. Each cell in the table shows the aggregated test AUC along with its standard deviation. As the results clearly indicate, the proposed method outperforms all baselines by a significant margin across all four data sources (MNIST, EMNIST, KMNIST, and FMNIST) with almost the lowest possible standard deviation. This consistent performance highlights both the effectiveness and robustness of our approach in handling the Visudo-PC task under varying input distributions.

| Data Source             | NeuPSL [10]     | LTN-IND [8]     |                 |                 | Proposed Method                   |
|-------------------------|-----------------|-----------------|-----------------|-----------------|-----------------------------------|
|                         |                 | A               | B               | C               |                                   |
| MNIST ( $4 \times 4$ )  | $0.88 \pm 0.02$ | $0.83 \pm 0.18$ | $0.84 \pm 0.14$ | $0.94 \pm 0.10$ | <b><math>0.97 \pm 0.01</math></b> |
| EMNIST ( $4 \times 4$ ) | $0.79 \pm 0.09$ | $0.58 \pm 0.04$ | $0.58 \pm 0.06$ | $0.65 \pm 0.14$ | <b><math>0.91 \pm 0.10</math></b> |
| KMNIST ( $4 \times 4$ ) | $0.65 \pm 0.12$ | $0.83 \pm 0.09$ | $0.85 \pm 0.11$ | $0.87 \pm 0.09$ | <b><math>0.93 \pm 0.01</math></b> |
| FMNIST ( $4 \times 4$ ) | $0.74 \pm 0.04$ | $0.67 \pm 0.11$ | $0.76 \pm 0.15$ | $0.83 \pm 0.11$ | <b><math>0.89 \pm 0.04</math></b> |

**Table 5.5:** Comparison of the proposed method to state-of-the-art.

# Chapter 6

## Conclusion

In this chapter, we summarize the key findings and contributions of the thesis and outline future research directions to improve and extend the proposed method.

### 6.1 Findings

During the course of this thesis, we proposed and developed a novel NeSy visual reasoning system. We began by investigating the limitations of existing methods through categorizing them into FBN, ENR, and FBV approaches in Chapter 1. Using this categorization, we demonstrated that none of the existing approaches simultaneously satisfy the three key goals of flexibility in addressing diverse visual reasoning tasks, explainability in articulating the reasoning process, and formality in expressing rules through a well-defined logical language. To address this challenge, we defined these principles as the core goals of our system in Chapter 3 and designed an iterative RL-inspired architecture by integrating a rule-generating VLM, a rule-verifying D-LTN, and a visual encoding CNN block. Next, focusing on the Visudo-PC benchmark, we conducted various experiments on our proposed method in Chapter 5 and proved its success in achieving all three goals by flexibly generating explainable and formal FOL rules.

Beyond this, we showed that, by bypassing the VLM, we could independently evaluate the combination of the CNN block and the D-LTN against state-of-the-art methods in the benchmark, achieving superior performance. To explain this improvement, we need to refer to the method proposed in [8], which forms the core inspiration behind our system. As discussed in Chapter 2, this method is also an LTN-based NeSy visual reasoning approach specifically designed for the Visudo-PC benchmark. However, while it also couples a CNN and an LTN similar to our system’s coupling of the CNN block and D-LTN, the roles of these components differ significantly. In this method, the CNN is used primarily for visual classification,

and the LTN is implemented manually using the LTNTorch framework [33]. These two differences help explain why our proposed method offers better performance:

- **Classification vs. Embedding:** While the original CNN performs classification, our CNN block produces embeddings. For a  $4 \times 4$  Sudoku, this means that the original method outputs a single scalar per visual object, whereas our embeddings are high-dimensional vectors. Consequently, our system captures more nuanced visual features, which are richer and more informative for downstream reasoning.
- **LTNTorch vs. FOL Grammar:** The LTNTorch framework supports the conversion of a wide range of arbitrary FOL rules into LTN implementations. While this generality is advantageous in principle, it can reduce practical effectiveness due to increased complexity and weaker optimization structure. In contrast, our D-LTN is automatically built from FOL rules that conform to a predefined grammar. Aside from automating the implementation of LTN formulas, this simplification enables more efficient learning and a tighter alignment between the logical rules and the network architecture. As a result, our method achieves more reliable rule verification and better integration with the learned visual features.

## 6.2 Future Directions

Despite the promising results and strengths of our proposed method, several opportunities remain for further improvement and expansion. Referring back to the conceptual design of our system in Figure 1.1, we outlined a blueprint for a flexible and modular platform, which has been shown to be effective based on the criteria introduced in Chapter 1. In this framework, we instantiated the rule generator as a VLM, the rule verifier as a D-LTN, and the visual encoder as a CNN block. However, each of these components can be independently studied and potentially replaced with alternative subsystems, as long as the overarching RL-inspired structure is preserved. Therefore, our final implementation, illustrated in Figure 3.1, represents just one of many possible realizations. However, even within this practical design, several promising directions merit further investigation:

- **VLM Improvement:** The VLM used in our system was not necessarily state-of-the-art; rather, it was the best option available to us for experimental purposes. By integrating a more powerful and context-aware VLM, we can increase the level of abstraction in prompts and broaden the range of visual reasoning tasks supported by our framework. Practically, this enhancement may allow us to progressively lift some of the constraints we applied to the

prompt engineer in Chapter 4, bringing us closer to realizing the complete FOL grammar and rule generation capabilities described in Chapter 3.

- **Fine-Tuned VLM:** In our current setup, the FOL grammar is enforced through in-context learning, where constraints are embedded within the prompt itself [14]. A natural progression would be to fine-tune the VLM directly on this grammar. With the right dataset, the VLM can be optimized to generate outputs strictly within the specified FOL syntax. This fine-tuning process means that even a relatively small VLM could be sufficient, since it would no longer need to support free-form text generation. Instead, its capacity could be focused entirely on formal logical reasoning. Techniques such as SFT and RL offer practical approaches for this direction [14].
- **True RL Loop:** Our system currently mimics RL through its iterative structure, but a future version could incorporate a genuine RL loop by introducing an explicit reward function. Instead of routing intermediate feedback solely to the prompt engineer, it could also be directed back to the VLM, enabling it to learn how to refine its rule generation autonomously. One practical implementation would be to define a reward based on the accuracy of rule verification over a validation subset and apply policy optimization techniques such as **Proximal Policy Optimization** (PPO), which would enable a full RL-based adaptation using machine feedback.

# Bibliography

- [1] Michael Hersche, Mustafa Zeqiri, Luca Benini, Abu Sebastian, and Abbas Rahimi. «A neuro-vector-symbolic architecture for solving Raven’s progressive matrices». In: *Nature Machine Intelligence* 5.4 (2023), pp. 363–375 (cit. on pp. 1, 3, 8, 10, 11, 13).
- [2] Daya Guo et al. «Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning». In: *arXiv preprint arXiv:2501.12948* (2025) (cit. on pp. 1, 16).
- [3] Kyle Hamilton, Aparna Nayak, Bojan Božić, and Luca Longo. «Is neuro-symbolic AI meeting its promises in natural language processing? A structured review». In: *Semantic Web* 15.4 (2024), pp. 1265–1306 (cit. on pp. 1, 3–5, 8, 9, 15).
- [4] ‘Deduction’ vs. ‘Induction’ vs. ‘Abduction’ — *merriam-webster.com*. <https://www.merriam-webster.com/grammar/deduction-vs-induction-vs-abduction>. (Visited on 12/12/2024) (cit. on p. 1).
- [5] Samy Badreddine, Artur d’Avila Garcez, Luciano Serafini, and Michael Spranger. «Logic tensor networks». In: *Artificial Intelligence* 303 (2022), p. 103649 (cit. on pp. 1–4, 9, 14, 15, 26).
- [6] Adam Santoro, David Raposo, David G Barrett, Mateusz Malinowski, Razvan Pascanu, Peter Battaglia, and Timothy Lillicrap. «A simple neural network module for relational reasoning». In: *Advances in neural information processing systems* 30 (2017) (cit. on pp. 1, 11, 12).
- [7] Chen Liang, Wenguan Wang, Tianfei Zhou, and Yi Yang. «Visual abductive reasoning». In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2022, pp. 15565–15575 (cit. on pp. 1, 2, 10).
- [8] Lia Morra et al. «Designing Logic Tensor Networks for Visual Sudoku Puzzle Classification.» In: *NeSy*. 2023, pp. 223–232 (cit. on pp. 2–4, 10, 15, 41, 47, 51, 52).

- [9] Zishen Wan et al. «Towards cognitive ai systems: a survey and prospective on neuro-symbolic ai». In: *arXiv preprint arXiv:2401.01040* (2024) (cit. on pp. 2, 5, 8, 9).
- [10] Connor Pryor, Charles Dickens, Eriq Augustine, Alon Albalak, William Wang, and Lise Getoor. «Neupsl: Neural probabilistic soft logic». In: *arXiv preprint arXiv:2205.14268* (2022) (cit. on pp. 2, 11–13, 51).
- [11] Laura Von Rueden et al. «Informed machine learning—a taxonomy and survey of integrating prior knowledge into learning systems». In: *IEEE Transactions on Knowledge and Data Engineering* 35.1 (2021), pp. 614–633 (cit. on p. 2).
- [12] Alessandro Daniele, Tommaso Campari, Sagar Malhotra, and Luciano Serafini. «Simple and Effective Transfer Learning for Neuro-Symbolic Integration». In: *International Conference on Neural-Symbolic Learning and Reasoning*. Springer. 2024, pp. 166–179 (cit. on pp. 3, 4, 12, 15).
- [13] Alessandro Daniele, Tommaso Campari, Sagar Malhotra, and Luciano Serafini. «Deep symbolic learning: Discovering symbols and rules from perceptions». In: *arXiv preprint arXiv:2208.11561* (2022) (cit. on pp. 4, 12).
- [14] Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Shengbang Tong, Saining Xie, Dale Schuurmans, Quoc V Le, Sergey Levine, and Yi Ma. «Sft memorizes, rl generalizes: A comparative study of foundation model post-training». In: *arXiv preprint arXiv:2501.17161* (2025) (cit. on pp. 4, 15, 16, 54).
- [15] Xin Zhang and Victor S Sheng. «Neuro-Symbolic AI: Explainability, Challenges, and Future Trends». In: *arXiv preprint arXiv:2411.04383* (2024) (cit. on pp. 4, 5).
- [16] Shima Imani, Liang Du, and Harsh Shrivastava. «Mathprompter: Mathematical reasoning using large language models». In: *arXiv preprint arXiv:2303.05398* (2023) (cit. on pp. 5, 16).
- [17] Nicholas Watters, Daniel Zoran, Theophane Weber, Peter Battaglia, Razvan Pascanu, and Andrea Tacchetti. «Visual interaction networks: Learning a physics simulator from video». In: *Advances in neural information processing systems* 30 (2017) (cit. on p. 8).
- [18] Giacomo Camposampiero, Michael Hersche, Aleksandar Terzić, Roger Wattenhofer, Abu Sebastian, and Abbas Rahimi. «Towards learning abductive reasoning using vsa distributed representations». In: *International Conference on Neural-Symbolic Learning and Reasoning*. Springer. 2024, pp. 370–385 (cit. on pp. 10, 11, 13, 14).

- [19] Yeongbin Kim, Gautam Singh, Junyeong Park, Caglar Gulcehre, and Sungjin Ahn. «Imagine the unseen world: a benchmark for systematic generalization in visual world models». In: *Advances in Neural Information Processing Systems* 36 (2023), pp. 27880–27896 (cit. on pp. 10, 11, 16).
- [20] Wenliang Zhao, Yongming Rao, Yansong Tang, Jie Zhou, and Jiwen Lu. «Videoabc: A real-world video dataset for abductive visual reasoning». In: *IEEE Transactions on Image Processing* 31 (2022), pp. 6048–6061 (cit. on pp. 10, 11).
- [21] Michael Hersche, Francesco Di Stefano, Thomas Hofmann, Abu Sebastian, and Abbas Rahimi. «Probabilistic abduction for visual abstract reasoning via learning rules in vector-symbolic architectures». In: *arXiv preprint arXiv:2401.16024* (2024) (cit. on pp. 11, 13–16).
- [22] Eriq Augustine, Connor Pryor, Charles Dickens, Jay Pujara, William Wang, and Lise Getoor. «Visual sudoku puzzle classification: A suite of collective neuro-symbolic tasks». In: *International Workshop on Neural-Symbolic Learning and Reasoning (NeSy)*. 2022 (cit. on pp. 11, 36).
- [23] Yifei Peng, Yu Jin, Zhexu Luo, Yao-Xiang Ding, Wang-Zhou Dai, Zhong Ren, and Kun Zhou. «Generating by Understanding: Neural Visual Generation with Logical Symbol Groundings». In: *arXiv preprint arXiv:2310.17451* (2023) (cit. on pp. 13, 14).
- [24] Mohamed Mejri, Chandramouli Amarnath, and Abhijit Chatterjee. «RE-SOLVE: Relational Reasoning with Symbolic and Object-Level Features Using Vector Symbolic Processing». In: *arXiv preprint arXiv:2411.08290* (2024) (cit. on pp. 13, 14).
- [25] Paul Tarau. «On LLM-generated Logic Programs and their Inference Execution Methods». In: *arXiv preprint arXiv:2502.09209* (2025) (cit. on p. 16).
- [26] OpenAI. *Images and vision: Learn how to understand or generate images*. publisher: OpenAI. URL: <https://platform.openai.com/docs/guides/images-vision?api-mode=responses> (visited on 06/09/2025) (cit. on p. 18).
- [27] Erez Shinan. *Lark - a parsing toolkit for Python*. URL: <https://github.com/lark-parser/lark> (cit. on pp. 22, 23).
- [28] Eriq Augustine. *Visual Sudoku Puzzle Classification*. Version 1.0.0. May 2022. DOI: <https://github.com/linqs/visual-sudoku-puzzle-classification>. URL: <https://github.com/linqs/visual-sudoku-puzzle-classification> (cit. on p. 38).

- [29] Google Research. *Google Colaboratory: A Hosted Jupyter Notebook Service for Machine Learning and Data Science*. <https://colab.research.google.com>. Accessed: 2025-06-17. 2025 (cit. on p. 38).
- [30] PyTorch Contributors. *PyTorch: An Open Source Machine Learning Framework*. <https://pytorch.org>. Version 2.5.1, Accessed: 2025-06-17. 2025 (cit. on p. 39).
- [31] Meta AI. *Llama-4-Maverick-17B-128E-Instruct: A Multimodal Mixture-of-Experts Model*. <https://huggingface.co/meta-llama/Llama-4-Maverick-17B-128E-Instruct>. Accessed: 2025-06-17. 2025 (cit. on p. 39).
- [32] Groq, Inc. *Groq: Fast Inference for Large Language Models*. <https://groq.com>. Accessed: 2025-06-17. 2025 (cit. on p. 39).
- [33] Tommaso Carraro. *LTNtorch: PyTorch implementation of Logic Tensor Networks*. Version 1.0.0. Mar. 2022. DOI: 10.5281/zenodo.6394282. URL: <https://doi.org/10.5281/zenodo.6394282> (cit. on p. 53).