
A Survey of Neurosymbolic Answer Set Programming

Journal Title
XX(X):1–50
©The Author(s) 2025
Reprints and permission:
sagepub.co.uk/journalsPermissions.nav
DOI: 10.1177/ToBeAssigned
www.sagepub.com/

SAGE

Alexander Philipp Rader¹ and Alessandra Russo¹

Abstract

Neurosymbolic artificial intelligence (AI) combines neural networks and symbolic methods to create robust and explainable frameworks. This survey provides an overview of the literature on neurosymbolic AI that uses answer set programming (ASP) as its symbolic language of choice. ASP is a logical formalism that can represent expressive rules and common-sense reasoning in a compact and human-readable form. Bridging the gap between neural representations and categorical symbols is a difficult task, especially when the learning of knowledge is involved. Many approaches have been proposed in the field to overcome these challenges and we categorise them based on which components are hard-coded or learned. We provide illustrations and explanations of the different types of frameworks and compare them with each other. We discuss the advantages of such hybrid models in terms of explainability and logical robustness. Lastly, we explore the limits of the field, including the simplicity of tasks, the extensive use of hard-coded knowledge, and the limited scalability of methods. We argue that better benchmarks, improvements in scalability and novel ways of propagating the learning signal through ASP components are needed to propel the field forward.

Keywords

Answer Set Programming, Neurosymbolic AI

Introduction

As AI systems are deployed more widely, issues of trust, safety and interpretability become ever more important. Modern neural AI models have made strides in many areas such as holding conversations (Liu et al. 2023b), passing difficult exams (OpenAI 2023) and generating images from text prompts (Podell et al. 2024). They are remarkably capable of processing real-world data and autonomously learning

¹Department of Computing, Imperial College London, United Kingdom

Corresponding author:

Alexander Philipp Rader, Imperial College London, Exhibition Rd, London SW7 2AZ, UK.
Email: apr20@ic.ac.uk

10 knowledge from examples. However, they lack explicit reasoning and logic capabilities, which can lead
11 to inconsistent outputs and hallucinations (Farquhar et al. 2024). Their black-box nature also means they
12 are not explainable and lack formal guarantees (John-Mathews 2021).

13 Marcus (2020) argues that symbolic representations are necessary for AI to achieve robust reasoning.
14 Unlike neural networks, symbolic AI acquires an internal model to represent abstract knowledge and can
15 reason with it logically. This model is human-readable, making its decisions transparent and explainable.
16 Depending on the formal representation of the model, guarantees can be made about its behaviour as well.
17 However, symbolic methods struggle to deal with real-world, noisy data and often scale exponentially
18 with the size of the problem domain.

19 The intersection of neural networks and symbolic methods is known as neurosymbolic AI and aims
20 to combine the best of both worlds (Garcez and Lamb 2023). There are countless ways of integrating
21 these two paradigms. In this survey, we focus on the use of answer set programming (ASP), a type of
22 symbolic AI that belongs to the field of logic programming. ASP is more expressive than languages
23 like Prolog, while still being relatively efficient to compute (Lifschitz 2019). It is therefore well suited
24 for representing complex tasks and a popular choice for neurosymbolic AI. We explore the capabilities,
25 advantages and drawbacks of frameworks combining ASP and neural networks throughout this survey.

26 We start by positioning this survey among other overview papers for similar fields in the related work
27 section. Additionally, we outline the use of other symbolic languages in neurosymbolic AI and compare
28 them with ASP, with a particular focus on formalisms within inductive logic programming. We establish
29 that this survey focuses on frameworks combining neural networks with ASP that were published within
30 the last six years. The remainder of the related work section includes pointers to work combining ASP
31 with reinforcement learning and application papers. In the background section, we formally define ASP,
32 explain how to learn ASP rules and provide a quick overview of neural networks.

33 The *Neurosymbolic ASP frameworks* section contains a discussion of papers in the field of
34 neurosymbolic ASP. It is split into four parts, which represent different categories of frameworks. We
35 categorise frameworks based on whether their neural component is pre-trained, their symbolic component
36 is hard-coded, or either of them is learned. We briefly describe each framework and discuss its strengths
37 and weaknesses. The descriptions are accompanied by illustrations allowing the reader to compare and
38 contrast them.

39 In the *Analysis of neurosymbolic ASP* section, we evaluate the performance of papers and identify the
40 current strengths, limits and open challenges in the field. The section is split into four parts: assessing the
41 difficulty of perception tasks, comparing accuracies on datasets, identifying the limits of ASP generation
42 and analysing scalability issues. Overall, the neurosymbolic frameworks perform well, often beating fully
43 neural baselines in terms of accuracy in addition to explainability and robustness. However, most datasets
44 contain simple inputs that do not reflect real-world scenarios. Another weakness is the reliance on hand-
45 written background knowledge and the need to heavily restrict the search space for finding ASP rules.
46 Lastly, the field suffers from timeouts and scalability issues that limit the ability of frameworks to scale
47 to real-world tasks.

48 We conclude that a great deal of research effort has been devoted to neurosymbolic ASP in recent
49 years. To continue pushing the field forward, more challenging datasets and novel methods for learning
50 are needed. A promising new direction is the combination of foundation models and ASP, where the
51 two paradigms complement each other well. ASP can help models move beyond summarising existing
52 knowledge towards making new discoveries while being robust and explainable.

53 Related work

54 In this section, we position our work among existing survey papers and related fields. We discuss the focus
55 of our survey on ASP and examine its strengths and weaknesses compared to other logic programming
56 paradigms. We also review related literature and provide pointers to adjacent work.

57 *Neurosymbolic AI*

58 This survey is situated within the wider field of neurosymbolic AI. [Belle and Marcus \(2026\)](#) provide a
59 historical perspective on the subject, highlight key developments, and conclude with an outlook into the
60 future. High-level overviews of approaches can be found in [Hitzler et al. \(2022\)](#) and [Sheth et al. \(2023\)](#),
61 which serve as a broad introduction.

62 For more comprehensive reviews, there are a variety of surveys that cover the field in detail: [Colelough
63 and Regli \(2025\)](#) conduct a systematic search of all papers published about neurosymbolic AI between
64 1970 and 2024. They uncover a large increase in publications on this topic beginning in 2020. [Gibaut
65 et al. \(2023\)](#) focus on six methods in neurosymbolic AI and survey the associated papers. [Yu et al. \(2023\)](#)
66 give a comprehensive overview of the domain, including its taxonomy, techniques and applications. [Wan
67 et al. \(2024a\)](#) group systems into five categories and profile them based on metrics such as underlying
68 operations. [Hitzler and Sarker \(2021\)](#) include a collection of 17 overview papers, encompassing different
69 branches of the field. The methods include graph reasoning, boolean circuits or logic tensor networks.
70 Last but not least, [Bhuyan et al. \(2024\)](#) create a conceptual map of the papers in the field, categorising
71 them based on domain, type and properties. They also delve deeper into individual frameworks, breaking
72 down the kinds of reasoning, representations and logics they use.

73 Many survey papers focus on specific aspects of the field. [Marra et al. \(2024\)](#) explore neurosymbolic
74 AI together with statistical relational AI. They identify shared dimensions between them and position
75 frameworks along these dimensions. [Acharya and Song \(2025\)](#) look at the field through the lens of
76 robustness, uncertainty quantification and intervenability. They classify how different techniques enable
77 improvements in these areas and outline current challenges. [Wan et al. \(2024b\)](#) experimentally evaluate
78 a suite of neurosymbolic algorithms in terms of metrics such as runtime and memory usage. They
79 uncover bottlenecks caused by factors such as symbolic operations, data dependencies, or complex flow
80 control. [Bouneffouf and Aggarwal \(2022\)](#) focus on applications of neurosymbolic AI in fields including
81 healthcare, finance, and information retrieval. [Manhaeve et al. \(2026\)](#) provide an overview of popular
82 benchmarks used to evaluate frameworks. They study the strengths and limitations of these tasks used
83 and propose desirable features that new benchmarks should have. [Lamb et al. \(2021\)](#) explore the use of
84 graph neural networks in neurosymbolic domains, including developments and applications. [Odense and
85 d'Avila Garcez \(2025\)](#) take a first step towards a formalisation of neurosymbolic AI. They establish
86 a formal definition of semantic encodings and show that many methods already fall under it. [Smet
87 and Raedt \(2025\)](#) contribute definitions for neurosymbolic models and inference. They build on the
88 observation that such systems combine logic with beliefs and show that many classes of models can
89 be cast within their definitions. These recent papers demonstrate growing efforts to establish a unifying
90 formalisation of the field and its many approaches.

91 We discussed surveys that provide a comprehensive overview of different methods in neurosymbolic
92 AI. Due to their broad scope, they tend to focus on major streams of research. There is no survey that
93 details the structures, strengths and limitations of frameworks that use ASP. The closest contender is the

94 work by [Borroto et al. \(2025\)](#), which reports on the intersection between ASP and neurosymbolic AI.
95 However, the authors limit their focus on a selection of frameworks and applications, leaving room for a
96 more comprehensive survey. Our work aims to fill this gap and provide a detailed account of the state of
97 the art in neurosymbolic ASP.

98 *Symbolic languages*

99 The *symbolic* part of neurosymbolic AI encompasses a wide range of different languages and paradigms.
100 In this section, we highlight some of the most common symbolic methods and discuss how they have been
101 combined with neural networks. For a more thorough overview, we direct you to the surveys discussed
102 in the previous section. Finally, we position ASP within the paradigm of logic programming and contrast
103 it with other languages in that category.

104 *Propositional and first-order logic.* A fundamental branch of symbolic logic is propositional logic,
105 which deals with two-valued variables combined using conjunctions, disjunctions and negation. There
106 are many frameworks encoding propositional formulas within layers of neural networks to make them
107 more explainable, such as Pix2Rule ([Cingillioglu and Russo 2021](#); [Baugh et al. 2023](#)), logic gate
108 networks ([Petersen et al. 2022](#)) or TELL ([Ragno et al. 2024](#)). First-order logic increases expressivity
109 by adding quantifiers, predicates and variables for representing complex knowledge. There are many
110 efforts to integrate first-order logic with neural networks, including logic explained networks ([Ciravegna
111 et al. 2023](#)), logical neural networks ([Riegel et al. 2020](#)) and neural logic machines ([Dong et al. 2019](#)).
112 [Badreddine et al. \(2022\)](#) propose logic tensor networks and define their own language called *real logic*
113 for it. It is a variant of first-order logic with fuzzy semantics, where domains are represented by tensors
114 and reasoning is implemented with real-valued tensor operations.

115 Symbolic sentences can take different forms that are useful for modelling and solving problems. In
116 decision trees, logical propositions are found in nodes and their truth values determine which edges to
117 take up to a leaf node that represents a classification. Frameworks that propose neurosymbolic decision
118 trees are [Möller et al. \(2025\)](#) and [Kairgeldin and Carreira-Perpiñán \(2025\)](#). Probabilistic circuits can
119 encode logic formulas into computational graphs that represent probability distributions and support
120 exact inference ([Choi et al. 2020](#)). They have been integrated into neural networks as a probabilistic layer
121 that ensures that predictions satisfy logical constraints ([Ahmed et al. 2022](#)). Other logical formalisms with
122 neurosymbolic frameworks include deterministic finite automata ([Zhang et al. 2024](#)), planning domain
123 definition language ([Liu et al. 2023a](#)) and action language ([Ishay and Lee 2025](#)).

124 *Logic programming.* A widely-used symbolic paradigm is logic programming, which models
125 computations as inference through a program. Such a program consists of facts and rules, in which the
126 head holds true whenever the body holds true. A query can be solved by checking whether it holds true
127 given the program. ASP belongs to the family of logic programming, along with languages like Prolog,
128 Datalog and abductive logic. There has been great interest in the research community to combine logic
129 programming with neural networks.

130 Prolog is the quintessential logic programming language. It supports representing knowledge as
131 definite clauses (with one head literal) and solves queries in a top-down, goal-directed fashion ([Clocksin
132 and Mellish 1981](#)). ProbLog is an extension to Prolog that allows facts to be annotated with probabilities
133 and therefore model uncertainty ([De Raedt et al. 2007](#)). By specifying these probabilities dynamically
134 using neural networks, [Manhaeve et al. \(2021\)](#) create the neurosymbolic framework DeepProbLog.

135 Given a program, DeepProbLog enables end-to-end training of the neural network by calculating
136 the gradients of probabilistic facts and backpropagating them to the network. The result is a unified
137 framework that integrates symbolic knowledge with the capabilities of neural networks to classify real-
138 world data. There are numerous variants and extensions of DeepProbLog that overcome shortcomings or
139 enhance the capabilities of the framework: DeepSeaProbLog adds continuous probability distributions
140 by incorporating weighted model integration (De Smet et al. 2023). DeepStochLog integrates neural
141 networks into stochastic definite clause grammars, which scales better than DeepProbLog as it does not
142 have to enumerate all possible worlds (Winters et al. 2022). VAEI combines ProbLog with variational
143 autoencoders, introducing symbolic and neural atoms to the latent space to aid generating the desired
144 output (Misino et al. 2022).

145 Datalog is a function-free subset of Prolog originally built for handling databases. Its restrictions
146 guarantee that queries will terminate and allow it to apply efficient bottom-up evaluation strategies for
147 better scalability. Scallop uses a probabilistic database as the interface between the neural and symbolic
148 components. The neural network predictions are stored as facts in the database and then queried over to
149 find the downstream label. By restricting logical evaluation to the top-k proofs, Scallop can calculate
150 gradients for the neural component efficiently (Huang et al. 2021). The apperception engine uses a
151 temporal variation of Datalog, which incorporates causal rules to model transitions. The framework builds
152 a theory that predicts future states given temporal data (Evans et al. 2021). Other variations include neural
153 Markov Prolog, which compiles logic rules into a Markov network (Thomson and Page 2023), and CR-
154 Prolog, a non-monotonic extension of Prolog that has been used for visual question answering (Riley and
155 Sridharan 2019).

156 Logical abduction extends standard logic programming by inferring facts and hypotheses that best
157 explain an observation. The Abductive Learning (ABL) framework implements abduction to bridge the
158 gap between neural predictions and a logical program. Given background knowledge and pseudo-labels
159 from the neural network, the abductive procedure evaluates whether the neural predictions are consistent
160 with the observation and provides corrected labels if not (Dai et al. 2019). MetaABD is an extension to
161 ABL that trains a neural network and induces the logical theory jointly. The background knowledge is not
162 given in full, instead the framework uses a combination of abduction and induction to learn recursive first-
163 order theories with predicate invention (Dai and Muggleton 2021). Han et al. (2023) employ abductive
164 logical reasoning to induce hypotheses for unpartitioned data. These hypothesis capture the relationship
165 between the target and subconcepts, which in turn are used to train neural networks. Similar to VAEI,
166 MetaAbd has been applied to increase control over visual generation with logical symbols (Peng et al.
167 2025). The method has been improved with pre-training to reduce the search space and achieve speed-
168 ups (Jin et al. 2025).

169 ASP introduces constructs into logic programming that help model commonsense reasoning with
170 incomplete knowledge. It supports default negation, where an atom is assumed false unless there is
171 evidence for it, and non-monotonic reasoning, where additional evidence can retract existing conclusions.
172 For these reasons, it is more expressive than languages like Prolog, where rules are limited to definite
173 logic, or MetaAbd, which cannot learn hard constraints. Law et al. (2018) prove that the learning of
174 answer set programs subsumes the other paradigms in inductive logic programming, making learning
175 in ASP the most general inductive framework. The high expressivity, combined with the availability
176 computationally efficient solvers, makes ASP a popular choice among practitioners.

177 However, expressivity is not the only relevant factor for neurosymbolic frameworks. More restricted
178 languages tend to run faster and can be easier to use. For example, Datalog can use efficient database
179 algorithms to scale better than Prolog. Abductive methods can conduct a guided search over the logic
180 program space to find latent labels, rather than enumerating all world possibilities. The query-driven
181 nature of Prolog can make framework design more intuitive than the stable model semantics of ASP.
182 Different frameworks also contain unique features, such as modelling probabilistic facts in ProbLog
183 or continuous probability distributions in DeepSeaProbLog. Many tasks do not make use of the higher
184 expressivity that ASP offers, often containing simple programs such as addition. In such cases, it is
185 desirable to choose frameworks with sufficient expressivity and superior scalability. There is no single
186 best language for neurosymbolic frameworks, as they all involve trade-offs.

187 *Survey scope*

188 This survey focuses on neurosymbolic frameworks that utilise ASP as their logical language. There are
189 two main reasons for this scope: First, it enables us to cover each framework in sufficient detail, complete
190 with illustrations and comparative analyses of their specific differences. Second, the properties of ASP
191 present unique opportunities and challenges to neurosymbolic AI that are worth studying in detail.

192 ASP's stable model semantics compute all valid solutions of a task at once. This paradigm aligns
193 well with the uncertain nature of neural networks, which assign probability mass to all possible concepts
194 given an input. Methods can consider all valid answer sets and assign them probabilities based on neural
195 predictions, rather than choosing one explanation that might be incorrect. However, these semantics
196 also increase the complexity of calculating solutions, limiting scalability. Furthermore, the existence
197 of multiple explanations poses a challenge for propagating gradients. There is no longer one single path
198 from the downstream solution to the input, but many paths through each answer set. For large problem
199 spaces with hundreds of thousands of answer sets, this can result in uninformative gradients. Bridging
200 this gap and providing a useful learning signal to the neural components is one of the main challenges in
201 neurosymbolic ASP.

202 ASP also lends itself to modelling logic problems, thanks to its syntactic constructs and non-monotonic
203 reasoning capabilities. A growing number of papers have combined large language models with ASP
204 to improve their logical reasoning performance. The authors justify their choice of ASP on the basis
205 of its suitability to model the commonsense knowledge found in natural language puzzles (see for
206 example [Ishay et al. 2023](#); [Rajasekharan et al. 2023b](#); [Alviano et al. 2024](#)). In particular, ASP is well-
207 suited to modelling problems using the generate-define-test pattern: Choice rules generate the search
208 space, auxiliary rules define the problem structure and constraints prune invalid solutions. Negation as
209 failure can handle default rules common in logic puzzles, i.e. general statements with specific exceptions.
210 Monotonic languages like Prolog are unable to model such problems.

211 Another advantage is the mature ecosystem of programs from which ASP-based systems can draw.
212 There are efficient solvers like Clingo ([Gebser et al. 2017](#)) or DLV ([Alviano et al. 2017](#)) that can
213 solve problems with large search spaces. Specialized variations provide utility for different use cases,
214 such as s(CASP), which solves programs top-down and can generate justifications ([Arias et al. 2018](#)).
215 Inductive learning of answer set programs from examples can be done with ILASP ([Law et al. 2020b](#)) or
216 FastLAS ([Law et al. 2020a](#)). All these systems are integrated into numerous neurosymbolic frameworks,
217 where they are essential for their functionality and efficiency.

218 To keep the content of this survey relevant, we only discuss papers released within the last six years,
219 i.e. from 2020 onwards. Moreover, we restrict our focus to “classical” learning of neural networks with
220 supervised or semi-supervised methods. For prominent work on combining reinforcement learning and
221 ASP, we provide a brief overview in the next section.

222 *Reinforcement learning*

223 ASP is also a popular specification language for rewards in reinforcement learning (RL). A variety of
224 papers define goals in ASP and some even learn rules from traces.

225 [Agostinelli et al. \(2024\)](#) specify the set of goal states with ASP, rather than listing each state one by
226 one. A deep reinforcement learning algorithm then estimates a heuristic function of the distance from the
227 current state to this set of goal. [Albilani and Bouzeghoub \(2023\)](#) use ASP rules in two ways: to generate
228 traces for learning low-level policies and as a backup policy in safety-critical scenarios. [Tudor and Gupta
\(2024\)](#) specify rules in the goal-directed ASP variant s(CASP) for RL and use a dependency graph to
229 prune the rules for faster execution. In this way, ASP helps train the neural component.

231 [Leonetti et al. \(2016\)](#) use ASP to represent the transition model of an environment and calculate plans
232 using an ASP solver. The plans represent partial policies, which restrict an agent to reasonable actions
233 during execution. Among these actions, the agent learns the expected cumulative reward using RL. In
234 principle, any planner can be used, but the authors chose ASP due to its ability to represent defaults. It
235 allows them to compactly represent optimistic assumptions, e.g. that all doors are open unless proven
236 otherwise.

237 [Furelos-Blanco et al. \(2021\)](#) introduce ISA, a framework for learning automata for subgoals in RL.
238 The automata are presented in ASP and have states for achieving or failing high-level goals. They use
239 ILASP to learn these automata from RL traces. [Parać et al. \(2024\)](#) extend this work and introduce the
240 ability to handle noisy data.

241 *Applications*

242 Neurosymbolic ASP has been applied in a variety of settings and we highlight some of them in this
243 section. For a more thorough review of applications in neurosymbolic AI more generally, refer to
244 [Bouneffouf and Aggarwal \(2022\)](#).

245 [Suchan et al. \(2021\)](#) combine object detection with spatial-temporal prediction and ASP for visual
246 sensemaking in autonomous driving. The environment is modelled with relational representations
247 pertaining to space and motion, while ASP computes commonsense interpretations about safety, space
248 and change. As it is applied to self-driving, the framework works online in a perceive-interpret-decide
249 cycle. [Rajasekharan et al. \(2023a\)](#) employ LLMs and ASP for argumentation analysis. They extract the
250 argumentation structure of a set of documents and represent it as an answer set program, which is used to
251 prove a claim. As an application, they demonstrate their framework on the topic of the MH17 Malaysian
252 Airlines flight downing. [Barbara et al. \(2023\)](#) propose an industrial application for neurosymbolic ASP
253 in the compliance checking of electrical control panels. They use deep learning to recognize electrical
254 components in images of panels and reconstructing its scheme. ASP is then employed to compare the
255 reconstructed scheme with the original schematic to discover problems. The authors prove the systems
256 by deploying it on a real test case in electrical control panel construction. Lastly, [Chu-Carroll et al. \(2024\)](#)
257 have used neuro-symbolic ASP in real-world applications at their company Elemental AI. They use it for

258 solving constraint satisfaction and optimisation problems. ASP serves as the logical reasoning engine and
 259 LLMs are used for knowledge acquisition and user interaction, in what they call an “LLM sandwich”.
 260 For knowledge acquisition, the LLM translates user inputs into an intermediate language, called Cogent,
 261 which is a constrained subset of English. For user interaction, the LLM takes the output of the ASP
 262 reasoner and presents it in a natural language interface.

263 Background

264 Answer set programming (ASP)

265 We give a brief overview of the syntax and semantics of ASP. For detailed definitions, please refer
 266 to [Gelfond and Kahl \(2014\)](#) and [Lifschitz \(2019\)](#).

267 *Syntax.* The language of ASP consists of constants, functions, predicates and variables. A **term** can be
 268 constructed as follows:

- 269 • A variable or constant is a term,
- 270 • If t_1, \dots, t_n are terms and f is a function of arity n , then $f(t_1, \dots, t_n)$ is a term.

271 An **atom** is an expression of the form $p(t_1, \dots, t_n)$ where p is a predicate of arity n . If $n = 0$, we omit
 272 the parentheses and just write p . Terms without variables are called **ground** and an atom is ground if
 273 every term in it is ground. A set of ground atoms is called an interpretation.

274 A general **rule** consists of atoms h_1, \dots, h_k and b_1, \dots, b_n and has the form:

$$h_1 \vee \dots \vee h_k :- b_1, \dots, b_m, \text{not } b_{m+1}, \dots, \text{not } b_n$$

275 The left part of the rule is the **head** and the right part is the **body**. The head is a disjunction of atoms and
 276 the symbol \vee is read as *or*. If the head only contains one atom, it is a **normal** rule. If the head is empty, the
 277 rule is a **constraint**. The body contains positive and negative atoms, the latter are indicated by the symbol
 278 **not** in front of them. Unlike classical negation, the symbol **not** denotes **negation as failure** and means
 279 that an atom is not believed to be true. If there are no negative atoms in the body, it is a **definite** rule.
 280 A **program** is a collection of rules. A definite program consists of only definite rules. A tight program
 281 contains no positive cycles and a stratified program contains no cycles through negation.

282 *Semantics.* A set S of ground atoms satisfies

- 283 1. atom p if $p \in S$,
- 284 2. **not** p if $p \notin S$,
- 285 3. $h_1 \vee \dots \vee h_k$ if for some $1 \leq i \leq k$, $h_i \in S$
- 286 4. b_1, \dots, b_n if S satisfies every atom in it,
- 287 5. a rule r if, whenever S satisfies the body of r , it satisfies the head of r .

288 The solutions of an answer set program are defined as sets of ground atoms, called **answer sets**. To
 289 determine whether a candidate set S is an answer set of a program Π , the reduct of the program, Π^S , needs
 290 to be constructed. This is done in multiple steps: First, ground the program by replacing all variables with
 291 all possible ground constants mentioned in the program. Then, remove all rules containing **not** p such

292 that $p \in S$. Last, remove all remaining body atoms containing `not`. S is an answer set of Π if it satisfies
 293 the rules of Π^S and is minimal (i.e. there is no proper subset of S satisfying the rules of Π^S).

294 The language contains more constructs to facilitate the modelling of complex problems. The main
 295 constructs include cardinality constraints, aggregates and optimisation statements.

296 A **cardinality constraint** is of the form

$$l\{h_1, \dots, h_k\}u :- b_1, \dots, b_m, \text{not } b_{m+1}, \dots, \text{not } b_n$$

297 where l and u are integer values such that $l \leq u$. Whenever the body holds, between l and u atoms in
 298 the head of the rule must be included in the answer set of the program. By leaving out the bounding
 299 numbers in a cardinality constraint, any combination of atoms in that set can be included in an answer
 300 set. Cardinality constraints are also known as **choice rules**.

301 **Aggregates** perform operations on sets. For example, the expression $\#\text{count}\{X : p(X)\}$ represents the
 302 number of elements in p . Other operations include sums, maxima and minima. Aggregates are often used
 303 in conjunction with comparison operators to model complex relationships.

304 **Optimisation statements** instruct the answer set solver to find solutions that either minimise or
 305 maximise certain properties. They are written using directives such as $\#\text{maximize}\{X : p(X)\}$. Rule
 306 bodies can also be annotated with weights, so-called weak constraints. The solver finds the optimal
 307 solution and ranks the answer sets based on these criteria.

308 $s(\text{ASP})$ (Marple et al. 2017) is an extension of the language, which provides solutions top-down in a
 309 goal-driven manner. Given a query, the system computes a partial answer set that contains it, bypassing
 310 the need to compute the entire answer set. Moreover, it does not need to ground the program, leading
 311 to performance improvements for some problems. A further extension is $s(\text{CASP})$ (Arias et al. 2018),
 312 which introduces constraints on variables, including over dense and unbounded domains.

313 ASP is a modelling language, meaning that problems are represented in a declarative way. To solve
 314 a problem, you typically model it in ASP and a solver calculates the possible solutions. One widely
 315 used solver is Clingo (Gebser et al. 2017). An answer set program can have multiple solutions and can
 316 represent problems up to the second level of the polynomial hierarchy (Law et al. 2018). ASP is also
 317 non-monotonic, meaning that adding new rules can invalidate previously held conclusions. This enables
 318 commonsense reasoning to be modelled through default rules, which usually hold unless disproven by
 319 new evidence.

320 Learning from answer sets

321 Rather than defining answer set programs by hand, the Learning from answer sets (LAS) task aims to
 322 automatically learn an answer set program from examples. This process, known as inductive learning,
 323 tries to find general rules that explain the given data. In this section, we define how to set up and solve a
 324 LAS task, using the definitions from Law et al. (2019).

325 Examples are represented in the form of weighted context-dependent partial interpretations (WCDPIs).
 326 A **WCDPI** is a tuple $e = \langle e_{id}, e_{pen}, e_{pi}, e_{ctx} \rangle$, where e_{id} is a unique identifier, e_{pen} is a penalty value,
 327 e_{pi} is a partial interpretation and e_{ctx} is the context. A partial interpretation e_{pi} consists of a pair of atom
 328 sets $\langle e^{inc}, e^{exc} \rangle$, called the inclusion and exclusion sets. The context is written in the form of an answer
 329 set program. A program P accepts a WCDPI e , iff there exists at least one answer set $I \in AS(P \cup e_{ctx})$,
 330 such that $e^{inc} \subseteq I$ and $e^{exc} \cap I = \emptyset$.

331 A **LAS task** is a tuple $T = \langle B, S_M, \langle E^+, E^- \rangle \rangle$, where B is the background knowledge, S_M is
332 the hypothesis space and E^+, E^- are sets of positive and negative WCDPIs. B is simply represented
333 as an answer set program and S_M is a set of ASP rules. The goal is to learn an optimal hypothesis
334 $H \subseteq S_M$, such that $B \cup H$ accepts as many positive WCDPIs and as few negative WCDPIs as possible.
335 The optimality of H is determined by summing up the penalties of negative/positive WCDPIs that are
336 accepted/not accepted respectively. In addition, a penalty for the length of H is applied to encourage
337 shorter hypotheses. Since it is infeasible to define S_M as a list of all possible rules that H could contain,
338 it is declared using a mode bias in practice. The **mode bias** consists of declarations and other constructs
339 that specify the hypothesis space.

340 Learning a hypothesis that defines concepts already observed in the examples is known as observational
341 predicate learning (OPL). This is computationally easier than non-OPL tasks, which require learning
342 concepts that are not directly observed (Law et al. 2021).

343 The two main LAS frameworks are ILASP (Law et al. 2020b) and FastLAS (Law et al. 2020a). ILASP
344 is capable of learning full answer set programs. FastLAS is more scalable but also more restricted. It
345 does not learn constructs such as choice rules or recursive rules and cannot invent new predicates. Both
346 systems are purely symbolic and cannot natively integrate neural networks.

347 *Neural networks*

348 Neural networks are composed of multiple layers of nodes, which are connected by weighted vertices.
349 They process an input by pushing it from one layer to the next, transforming it with linear and non-linear
350 functions. The transformed data that is output from the last layer represents the result. A dataset of input-
351 label pairs is used to learn a task. For each dataset example, the neural network output is compared with
352 the label and a loss is calculated. The weights of the neural connections are changed with regards to
353 the loss through the process of backpropagation (LeCun et al. 2015). Unlike symbolic methods, neural
354 networks represent knowledge sub-symbolically via their structure and the values of the weights applied
355 to connections.

356 *CNNs and GNNs.* A class of neural networks for image classification are convolutional neural networks
357 (CNNs). They contain specialist convolutional layers, which detect local features, and pooling layers,
358 which merge the features into higher-level concepts. These operations are particularly suited for
359 images, as the different convolutional functions act like filters and extract different properties from the
360 input (LeCun et al. 2015).

361 Graph neural networks (GNNs) use a specialized architecture for encoding graph structure directly in
362 neural space. They consist of nodes and edges to model entities and relationships. The main differentiator
363 compared to typical neural networks is message passing, whereby the internal states of nodes are updated
364 based on neighbours in the graph structure. GNNs are particularly well suited to data that can be
365 represented as graphs, such as molecules or social networks (Lamb et al. 2021).

366 *Foundation models.* Neural networks with billions of parameters that are trained on vast amounts of
367 broad data are known as foundation models (Bommasani et al. 2022). They are capable of solving a
368 wide range of problems through the use of **in-context learning**. The model learns how to solve a task
369 from examples provided in the prompt as a natural language description. No weights are changed in this
370 process, instead the model is only conditioned to utilise existing parameters. This type of adaptation is
371 known as few-shot learning. When no examples are provided in the prompt, the model performs zero-shot

learning. Performance can be improved by encouraging the model to break down its reasoning into steps, a process known as chain-of-thought-prompting. **Finetuning** is used to adapt a model to task-specific data by changing its weights.

The two main types of foundation models are large language models (LLMs) and vision language models (VLMs). LLMs are based on the transformer architecture and process text through the mechanism of attention. VLMs combine an LLM with a vision encoder to process multimodal input in the form of text and images (Bordes et al. 2024).

Neurosymbolic ASP frameworks

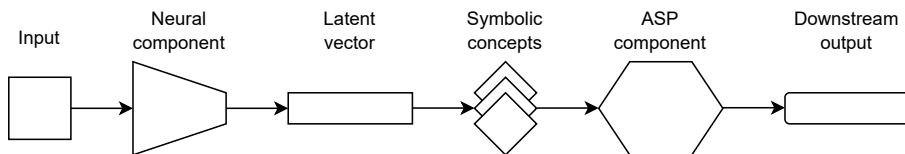


Figure 1. High-level depiction of inference through the components of neurosymbolic ASP frameworks.

There are a wide variety of frameworks which combine neural networks and ASP. In all approaches, the neural component processes raw inputs, while the ASP component performs logical reasoning to create an output, as illustrated in Figure 1. Combining neural and symbolic methods requires a translation between the latent vector representation of neural outputs and the symbolic concept representation of symbolic reasoners. We illustrate a typical inference procedure through such a neurosymbolic architecture with an example.

Example. In the MNIST Addition task, each input consists of two images of handwritten digits from the MNIST dataset (Deng 2012). Each downstream output is a single number, representing the sum of the two input numbers. In a neurosymbolic framework, the neural component can be a simple CNN which processes one image at a time and produces a latent vector of size 10. The i th entry in the vector represents the probability of the input being number i , for $i \in \{0, \dots, 9\}$. By choosing the argmax , i.e. the index of the entry with the highest probability, each latent vector can be translated into a symbolic concept. The ASP component can include a rule for adding up the two symbolic concepts: $\text{result}(Z) :- \text{digit}(1, X), \text{digit}(2, Y), Z = X + Y$. The downstream output is the number Z in $\text{result}(Z)$.

Compared to fully neural methods, a neurosymbolic approach has the advantages of robustness and explainability. The ASP component provides a decision based on human-readable rules, in contrast to an opaque neural network, which uses layers of nodes and weights. Any conclusion output by the ASP component is also logically robust given these set of rules, which is not guaranteed with a neural network. However, the increased transparency comes at the cost of complexity.

First, the translation between the continuous vector space of neural networks and the discrete symbols and rules of ASP is not always trivial. Symbols have to be extracted from raw data, such as natural language text or images. Depending on the problem, the translation may involve an unknown number of concepts or complex perception tasks.

404 Second, providing a learning signal for the neural and/or symbolic component is difficult. In a
 405 traditional neural network task, the model is trained end-to-end using the input and labels. For
 406 neurosymbolic ASP frameworks, the neural network outputs latent concepts, for which labels are often
 407 unavailable. Instead, the learning signal comes from the downstream labels, which have to be propagated
 408 through the non-differentiable ASP component. In many tasks, different combinations of latent symbols
 409 can result in the same downstream output, providing a noisy learning signal to the neural network.
 410 Learning the ASP component itself is challenging as well, as it receives noisy inputs from the neural
 411 network.

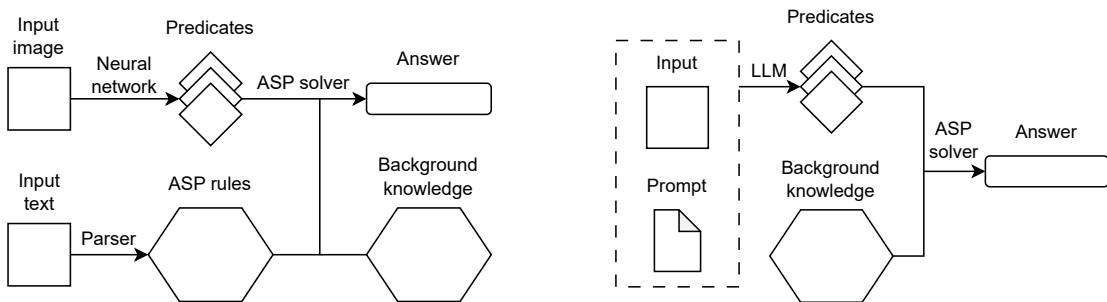
412 The proposed frameworks in the literature are all structured differently and deal with their own set of
 413 challenges. They can be split up into four broad categories:

- 414 1. Frameworks with a pre-trained neural and hard-coded ASP component. Their main challenge lies
 415 in the translation of neural outputs into symbols.
- 416 2. Frameworks with a hard-coded ASP component that train a neural network. Their main challenge
 417 lies in the propagation of the downstream learning signal through the ASP component.
- 418 3. Frameworks with a pre-trained neural component that learn an answer set program. Their main
 419 challenge lies in the learning of ASP rules with noisy neural predictions.
- 420 4. Frameworks that learn the neural and ASP component jointly. Their main challenge is a
 421 combination of all the problems above.

422 In this section, we discuss each category and illustrate the frameworks within it.

423 *Pre-trained neural and hard-coded symbolic component*

424 When both the neural and symbolic components are already given, the main focus lies on bridging the
 425 gap between them. Papers in this area tend to choose challenging tasks with natural language texts and
 426 complex images to demonstrate the usefulness of their framework. While earlier works use hand-crafted
 427 parsing pipelines, more recent papers experiment with LLMs.



(a) High-level depiction of the inference procedure in ASP-VQA (Eiter et al. 2022), AQuA (Basu et al. 2020) and NSGRAPH (Bauer et al. 2025).

(b) High-level depiction of the inference procedure in [LLM]+ASP (Yang et al. 2023), STAR (Rajasekharan et al. 2023b) and LLMASP (Alviano et al. 2024).

Figure 2. High-level depictions of frameworks with pre-trained neural and hard-coded symbolic components.

428 Figure 2a depicts the general structure of frameworks for the problem of visual question answering
429 (VQA). In VQA, the task involves answering questions about an image, such as “How many blue objects
430 are in the scene?”

431 *ASP-VQA and AQuA.* Both the ASP-VQA (Eiter et al. 2022) and AQuA (Basu et al. 2020) frameworks
432 use a YOLO network (Redmon and Farhadi 2018) to extract predicates from the image. YOLO is a neural
433 network architecture that outputs bounding boxes for objects in an image. Each row in its output vector
434 represents an object, and the columns correspond to class probabilities. Converting these neural vectors
435 into symbols is done simply by picking the class with the highest probability. On top of that, ASP-VQA
436 uses thresholding to select multiple likely classes per object and aggregates them into a choice rule. The
437 ASP component can choose any of these top predictions, allowing room for error. As both approaches
438 pre-train the YOLO network directly on given latent labels, no learning is occurring, just inference.

439 To extract the query and knowledge from text questions, both frameworks use a parsing pipeline. ASP-
440 VQA does not use the natural language text directly, but its functional representation, which the dataset
441 provides. The functional representation is a structured format made up of function symbols, predicates
442 and relations. The translation into ASP can therefore be done by a straightforward set of parsing rules.
443 AQuA, on the other hand, parses the natural language text and converts it into ASP by utilising a part-
444 of-speech tagger and dependency parser from CoreNLP (Manning et al. 2014). Both approaches also
445 include extensive background knowledge of the task hard-coded in ASP. An ASP solver then calculates
446 the answer by combining the extracted predicates and background knowledge.

447 By adding a symbolic layer on top of the YOLO network, the frameworks achieve the capability to
448 answer complex queries. AQuA even exceeds human baseline performance on one of their datasets. The
449 symbolic layer also adds robustness, as ASP-VQA reports good results even when the network is poorly
450 trained.

451 Eiter et al. (2023) extend ASP-VQA with the ability to provide contrastive explanations. The enhanced
452 framework is able to explain why the answer is P, rather than foil F, by showing how the input would
453 need to change to yield F. This abduction problem is encoded ASP by spanning the search space with
454 choice rules, adding constraints to ensure the answer contains the foil and adding weak constraints that
455 prefer minimal changes. For example, a question might ask whether there are three purple objects to the
456 left of the sphere in a given image of 3D objects. A contrastive explanation could include the answer *no*
457 and the explanation that shifting the sphere 20 pixels to the right would change the answer to *yes*. The
458 authors use the CLEVR dataset (Johnson et al. 2017) as a testbed, where the each input is a scene of 3D
459 objects and a corresponding question. They expand the dataset with 20 questions specific to contrastive
460 explanations and achieve an accuracy of 99%.

461 *NSGRAPH.* Similar techniques have been applied to graph problems in (Bauer et al. 2025). The input
462 images are pictures of graphs inspired by metro maps and the input text contains questions like “How
463 many stations are between A and B?”. To parse the input picture, NSGRAPH uses an optical graph
464 recognition network for the nodes and edges, and an optical character recognition network for the label
465 text. The questions are parsed using regular expressions, or LLMs for harder tasks. The dataset is based on
466 graph images generated from CLEGR (Mack and Jefferson 2018), which the authors call CLEGR^V. They
467 also present two further extensions: CLEGR⁺ and CLEGR-HUMAN. The former adds reformulated
468 questions by replacing words with synonyms and rephrasing sentences. The latter adds questions that

469 have been hand-crafted by humans using an online survey. NSGRAPH achieves accuracies of 85% and
470 94% respectively.

471 Recently, LLMs have taken over the process of parsing natural language input. Their remarkable ability
472 to produce structured language output from natural language makes them well suited for the task of
473 extracting ASP facts from text. They also require much less hand-crafting than pipelines with taggers
474 and parsers. Figure 2b illustrates the use of LLMs in bridging the gap between text and ASP predicates.

475 *[LLM]+ASP*. Yang et al. (2023) use LLMs for extracting ASP facts from natural language text puzzles
476 in their [LLM]+ASP framework. As LLMs are general-purpose models, in-context learning is employed
477 through examples of correct extractions in the prompt. The extracted facts are combined with hard-coded
478 background knowledge and an ASP solver arrives at the answer. The background knowledge comes in
479 the form of knowledge modules, which are written in a general way and therefore reusable for different
480 datasets. For example, the *location* module includes ASP rules for calculating an object's location using
481 offsets and is used for spatial reasoning, navigation and path-finding problems. The authors show that the
482 framework can solve robot planning tasks that the LLM alone fails at, thereby enhancing its capabilities.
483 Moreover, their method has uncovered errors in datasets such as StepGame (Shi et al. 2021), where over
484 10% of datapoints contain conflicting information. While neural-only methods learn to fit to the errors,
485 [LLM]+ASP flags them for being inconsistent.

486 Nguyen et al. (2025) have built a framework dedicated to detecting such errors. They integrate
487 LLMs with the explainable ASP solver XClingo (Cabalar et al. 2020) to find misleading information
488 in the CLUTRR dataset (Sinha et al. 2019). The LLM parses atoms from the input text and Clingo
489 finds the answer sets. If they contain conflicting information, e.g. *mother(theresa, darnell)* and
490 *sister(theresa, darnell)*, then XClingo generates an explanation tree that traces back to the
491 source of the inconsistency. Their experiments show that about 15% of examples in CLUTRR contain
492 inconsistent or ambiguous information.

493 *STAR*. The STAR framework (Rajasekharan et al. 2023b) extracts predicates from language inputs
494 using an LLM and reasons over them with ASP. The authors make use of both in-context learning and
495 finetuning to improve performance. Unlike [LLM]+ASP, they reason with the s(CASP) variation, which
496 is query-driven and more scalable. s(CASP) adds the ability to justify an answer in form of a proof tree,
497 which enhances explainability compared to using an LLM directly. Hard-coded background knowledge
498 is once again needed to represent the commonsense knowledge necessary for solving the problems.

499 In further work, the authors use the STAR framework to create a neurosymbolic chatbot. In this
500 scenario, the background knowledge is a conversational template that includes rules for staying on topic,
501 gathering relevant information from the user and answering their questions. To create a natural flow of
502 conversation, the LLM translates the ASP output into natural language again. Zeng et al. (2023) use this
503 system to create a conversational agent called AutoConcierge using the LLM GPT3 (Brown et al. 2020).
504 The aim is to provide restaurant recommendations based on nine relevant properties, e.g. location and
505 price preferences, that are extracted from the user with natural language dialogues. Zeng et al. (2024) use
506 the same system to create AutoCompanion, which is a social conversational bot for movies.

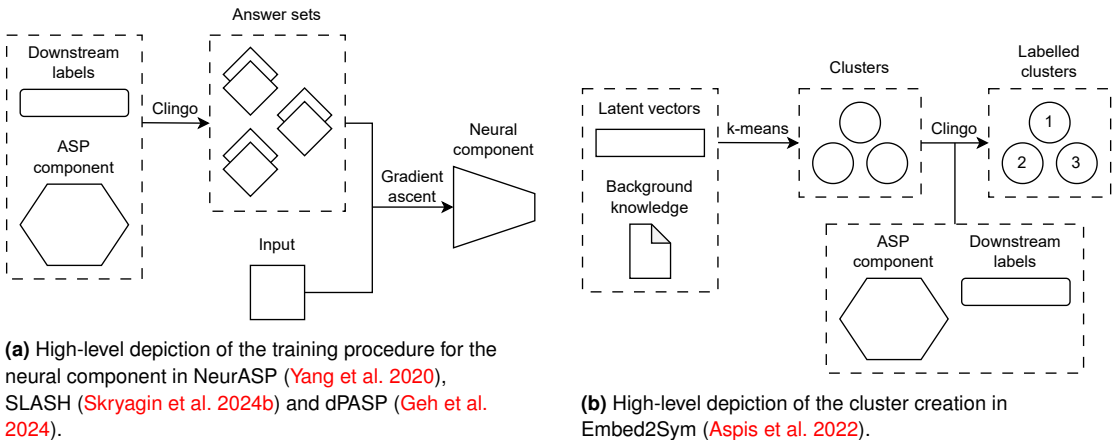
507 *LLMASP*. Alviano et al. (2024) use the YAML format to formalise the structure of prompts and
508 background knowledge. The LLMASP system makes use of two hand-written YAML files to generate the
509 prompt: The behaviour file includes generic instructions for the LLM to translate the input into ASP facts.

510 The application file includes domain-specific knowledge that describes the context of the task and what
 511 kind of ASP facts should be extracted. Both files are combined into one prompt and the LLM extracts the
 512 relevant facts from the user input. An ASP solver then finds the answer, making use of the extracted facts
 513 and a hard-coded ASP knowledge base. Finally, the LLM translates the ASP facts into natural language
 514 to answer the user’s questions. By providing the prompt structures in YAML, the system can be adapted
 515 to new domains in an efficient and rigorous manner.

516 Overall, the frameworks in this section use symbolic components to improve the capabilities of neural
 517 networks and LLMs. They successfully apply parsers or LLM predictions to bridge the gap between raw
 518 data and ASP facts. Answers from such hybrid networks are more robust than from the neural networks
 519 alone and can even find errors in the dataset or include contrastive explanations. However, the scope
 520 has been limited to the interplay between neural and ASP methods, where no neurosymbolic learning is
 521 involved. The absence of automatic learning procedures also necessitates labour-intensive hard-coding,
 522 thereby limiting adaptability to new problems.

523 *Neural training with hard-coded symbolic component*

524 The frameworks in this section contain hand-written answer set programs and tackle the issue of training
 525 neural networks indirectly. This type of task is also referred to as neurosymbolic *reasoning*. The neural
 526 component must learn latent symbols without access to latent labels, instead relying on downstream
 527 labels. Thus, the main challenge lies in propagating the learning signal through the non-differentiable
 528 symbolic component to the neural network.



(a) High-level depiction of the training procedure for the neural component in NeurASP (Yang et al. 2020), SLASH (Skryagin et al. 2024b) and dPASP (Geh et al. 2024).

(b) High-level depiction of the cluster creation in Embed2Sym (Aspis et al. 2022).

Figure 3. High-level depictions of frameworks with hard-coded symbolic components that train the neural network.

529 Many frameworks in this section follow the same high-level procedure for training, which Figure 3a
 530 illustrates. For each example, they calculate all answer sets and use them as noisy labels for training
 531 the neural network. Where they differ is the loss functions they use and the structure of their neural
 532 components and learning algorithms.

533 *NeurASP*. The first framework of this kind is *NeurASP* (Yang et al. 2020), which trains a neural network
534 given an answer set program and downstream labels. The neural component outputs a latent vector,
535 which acts as a probability distribution over ASP concepts. A hard-coded ASP ruleset then calculates the
536 downstream prediction using the most probable symbolic concepts from the neural network output. To
537 propagate the learning signal through the symbolic program to the neural network, *NeurASP* first finds
538 all answer sets that satisfy the downstream label. Each answer set contains a set of so-called neural atoms,
539 which are the concepts that the neural component predicts. *NeurASP* calculates the probability of each
540 answer set by multiplying the probabilities of the neural atoms in it. Finally, it calculates the gradient
541 of the loss with respect to each neural output and performs gradient ascent. Given an entry in the neural
542 latent vector, its gradient is increased for each answer set that contains it and decreased for each answer
543 set that does not contain it, weighted by the probability of the answer set. In effect, the answer sets act as
544 a weighted set of noisy latent labels.

545 **Example.** We revisit the MNIST Addition example from the beginning of this section to illustrate
546 the training procedure. In *NeurASP*, the neural network predicts the value of a single digit and the sum
547 operation is hard-coded in the ASP component. At training time, *NeurASP* calculates all answer sets for
548 a given downstream label. For example, if the downstream label is 11, then the answer sets contain the
549 following combinations of neural atoms: $\{(2, 9), (3, 8), (4, 7), (5, 6), (6, 5), (7, 4), (8, 3), (9, 2)\}$. Each
550 answer set is assigned a probability based on neural network confidences. If the neural network is highly
551 confident that the first number is 6 and the second number is 5, then the answer set (6, 5) will have
552 a higher probability than, say, (8, 3). The gradient ascent operation will then increase the weights for
553 predicting numbers 2 to 9, as they appear in the answer sets, and decrease the weights for 0 and 1, which
554 do not appear in any answer set. The increases and decreases are weighted by the answer set probabilities.

555 Splitting up a task in this way alleviates pressure on the neural network. Instead of solving the entire
556 task, it only has to learn latent concepts. Existing knowledge can then be utilised in the form of ASP rules
557 to find the downstream solution.

558 *SLASH*. Skryagin et al. (2022) introduce *SLASH*, an extension of *NeurASP* that integrates more
559 sophisticated probability estimations. In *NeurASP*, the perception component is a neural network and
560 is only capable of estimating conditional probabilities for each symbol C given data X : $P(C|X)$.
561 This is typically done using a Softmax function on its last layer. *SLASH* extends this notion with
562 neural-probabilistic predicates (NPPs). NPPs can learn the probability distribution of the latent concepts,
563 allowing *SLASH* to estimate $P(X|C)$ and $P(X, C)$ as well. Through density estimation, *SLASH* can
564 also handle missing data points and regenerate them. In the paper, NPPs are realised using probabilistic
565 circuits, but they can be replaced by any other component that estimates probabilities, including neural
566 networks.

567 The scalability of *SLASH* is improved in Skryagin et al. (2024b), where the authors introduce a
568 method called SAME to prune insignificant answer sets and speed up learning. As the neural predictions
569 improve during training, SAME gradually eliminates symbolic latent concepts with low probabilities
570 when generating answer sets. Over time, each epoch gradually speeds up, as the gradients are calculated
571 using fewer and fewer answer sets. Further speed improvements come in the form of answer set networks
572 (ASNs), which calculate answer sets by leveraging the GPU (Skryagin et al. 2024a). ASNs are answer
573 set programs encoded into the form of a GNN, where nodes and edges represent the atoms and relations.
574 This representation can compute answer sets in parallel on GPU nodes, unlike solvers like Clingo,
575 which are CPU-bound. The authors use ASNs as a replacement for Clingo in *SLASH* and report speed

576 improvements of two orders of magnitude for complex problems, allowing them to finetune LLMs with
577 SLASH. However, ASNs require answer set programs to be tight, which require an exponential number
578 of additional formulas in the worst case [Lin and Zhao \(2004\)](#).

579 *dPASP*. [Geh et al. \(2024\)](#) introduce a more powerful specification language with dPASP. It extends
580 the capabilities of NeurASP and SLASH by introducing interval-valued facts and disjunctions that are
581 annotated with probabilities. They implement two semantics for their framework: maxent and credal.
582 The former assigns probabilities based on maximising entropy, while the latter is more conservative and
583 assigns tighter bounds. The syntax of dPASP allows for the seamless integration of Python code and
584 an interface between raw data and program constants. The learning function is based on a Lagrange
585 multiplier derivation for gradient ascent and ends up being very similar to NeurASP’s learning rule. It
586 calculates the same terms as NeurASP, but multiplies them with weight factors $\frac{1}{m}$ and $1 - \frac{1}{m}$, where m
587 is the number of possible atoms that the neural network output represents.

588 *Embed2Sym*. Unlike the previous frameworks, Embed2Sym ([Aspis et al. 2022](#)) fully pre-trains the
589 neural component on the downstream labels. The neural network is structured to contain a separate
590 perception and reasoning component, both of which are neural. The perception component processes
591 the input and projects it into a latent dimension. The reasoning component takes the concatenated latent
592 vectors and predicts the downstream output. This structure allows the entire neural network to be trained
593 end-to-end with downstream labels, while producing a representation of latent concepts. To bridge the
594 gap between latent vectors and symbolic concepts, the framework uses k-means clustering, as shown in
595 [Figure 3b](#). The number of clusters is equal to the number of values a latent concept can take and is hard-
596 coded in the background knowledge. In the case of MNIST Addition, there are 10 clusters, one for each
597 single-digit number. Matching each cluster with the correct concept label is done with Clingo using a
598 hand-crafted answer set program. This program includes rules for reaching the downstream answer from
599 latent concepts and chooses a cluster/concept matching that maximises the number of correct downstream
600 predictions. In the MNIST Addition example, the answer set program would include rules for summing
601 up the two latent concepts and assigning each digit the cluster that leads to the maximum number of
602 correct sum predictions. At inference time, the framework uses the neural perception component to
603 create latent vectors. It assigns each vector the symbolic label corresponding to the nearest cluster and
604 then solves the task using the hard-coded ASP rules. The ASP component can be modified to solve new
605 problems, such as subtraction, without retraining the neural network. This makes the model transferable
606 to new domains, unlike a purely neural solution.

607 The latent embedding space of the neural component might not produce perfect clusters, leading
608 to misclassified concepts. [Rader and Russo \(2023\)](#) alleviate this issue by extending the framework
609 with active learning. They use the clusters to create a dataset of latent labels and finetune the
610 perception component to predict latent concepts directly. For each example where the downstream
611 prediction is correct, the cluster assignment is used as the the latent label. For examples with incorrect
612 downstream predictions, an oracle provides *active* latent labels for the dataset. The extension improves
613 the performance of the framework, enabling it to classify concepts more accurately than the clusters. As
614 only a small number of datapoints need to be labelled and the network is not retrained from scratch, this
615 extension is both data- and time-efficient.

616 Overall, there are two approaches to propagate the learning signal through a hard-coded ASP component.
617 NeurASP, SLASH and dPASP use answer set calculations to enumerate all possible latent concepts

618 for a downstream label. Embed2Sym trains a neural network end-to-end and uses clustering in the
619 embedding space to generate latent labels. Either approach creates a more explainable framework and
620 enables generalization to new tasks without the need of retraining, as only the ASP component has to be
621 changed.

622 *Symbolic learning with pre-trained neural component*

623 Pre-training a neural component on latent labels alleviates the challenge of propagating the downstream
624 learning signal back to the neural network. Papers in this section instead focus on translating neural
625 outputs into symbols and learning an answer set program to solve the task. They use two main approaches
626 for learning ASP rules: Either extracting symbolic concepts and then using a LAS solver, or generating
627 rules directly with an LLM. Setting up a LAS task is not trivial, as the search space of possible rules
628 is large and neural outputs are noisy. Generating rules with LLMs is challenging as well, because the
629 free-form output has to be syntactically and semantically correct. In this section, we discuss the different
630 strategies that have been proposed to overcome these issues.

631 Figure 4a illustrates frameworks that deploy an off-the-shelf LAS solver like ILASP or FastLAS to
632 find ASP rules. They differ in what neural methods they use to create the symbolic concepts for the LAS
633 task.

634 *FFNSL*. The framework FFNSL (Cunnington et al. 2023a) uses a standard neural network that is
635 pre-trained using latent labels. Their so-called data-to-knowledge generator translates the latent vector
636 outputs of the neural component into ASP atoms by selecting the index of the maximum value in the
637 vector. This translation is hard-coded, so it is known which vector entries correspond to which symbolic
638 concepts. The symbolic component then learns an answer set program that solves the downstream task.
639 It uses ILASP or FastLAS to find ASP rules based on the predicted symbolic concepts, the downstream
640 labels and hard-coded background knowledge, which includes the search space for the possible rules.
641 The paper investigates the effect of distributional shifts in the dataset, which cause the accuracy of the
642 neural networks to plummet. However, the symbolic learning remains robust to noisy predictions and the
643 generated rules outperform fully neural baselines.

644 *NeSyGPT*. Instead of training a traditional neural network, NeSyGPT (Cunnington et al. 2024) extracts
645 symbols from the data using a VLM. The framework feeds the input image into BLIP (Li et al. 2022)
646 alongside a question designed to extract the latent concept. For example, the authors set the question
647 “What number is this?” in the MNIST Addition task. For more complex perception tasks, such as
648 extracting the suit and rank of a playing card, they additionally finetune BLIP with latent labels. Since
649 there are no guarantees on the BLIP output, they use a text distance metric to map the VLM output to
650 a predefined set of symbolic concepts. This interface between neural and symbolic components, i.e. the
651 set of symbolic concepts and what questions to ask the VLM, is not necessarily manually engineered.
652 The authors present a way to programmatically generate it using LLMs. After all examples have been
653 converted into symbolic concepts, NeSyGPT learns ASP rules with ILASP.

654 *Embed2Rule*. To reduce the number of calls to a VLM, Embed2Rule (Aspis et al. 2024) uses BLIP
655 to label clusters rather than each individual example. It follows the same procedure as Embed2Sym
656 to generate the clusters, which we illustrated in Figure 3b. Images from each cluster are then sampled
657 and weakly labelled using BLIP. As BLIP might assign different labels to images in the same cluster,

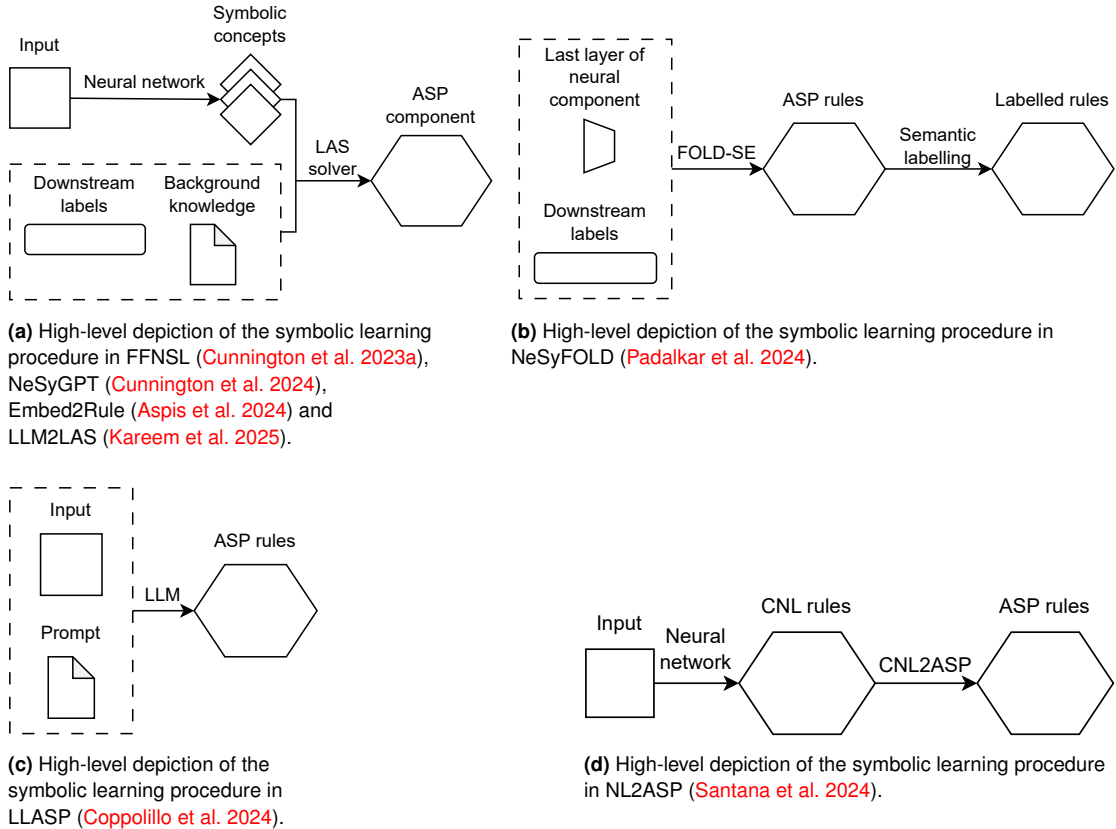


Figure 4. High-level depictions of symbolic learning procedures in frameworks with pre-trained neural components.

658 an optimisation algorithm finds the cluster-label assignment that maximises agreement with the BLP
 659 labels. Lastly, the learned symbolic concepts are used to find rules with ILASP. Only requiring the VLM
 660 to label a few datapoints per cluster enhances the data and compute efficiency of this method.

661 *LLM2LAS*. The domain of LLM2LAS (Kareem et al. 2025) consists of story-based questions written
 662 in natural language. The framework utilises an LLM to generate predicates from the story using few-
 663 shot prompting. A parsing pipeline then converts these predicates into an event calculus representation
 664 depicting actions, effects and timepoints. Moreover, the pipeline creates mode biases and CDPIs for
 665 the LAS task automatically. Only a minimal amount of background knowledge, two lines of ASP, is
 666 hard-coded. The pipeline for creating predicates and mode declarations is also quite general, as long as
 667 the task can be represented in event calculus. This is the case for 13 out of 20 task types in the bAbI
 668 dataset (Weston et al. 2015) that the authors test LLM2LAS with.

669 In an updated version, [Borroto Santana et al. \(2025\)](#) replace the parsing pipeline with an LLM that
670 generates mode biases for ILASP. This enables the framework to solve 15 rather than 13 tasks in bAbI.
671 An additional two tasks can be solved with the help of extended background knowledge. For the rest, the
672 hypothesis space remains too large for ILASP to find a solution.

673 *NeSyFOLD*. The aim of NeSyFOLD ([Padalkar et al. 2024](#)) is not to solve a task, but to explain decision-
674 making in CNNs. The CNN is pre-trained on downstream labels and NeSyFOLD turns its final layer into
675 ASP rules, as Figure 4b illustrates. The rationale is that filters in the final layer tend to represent high-
676 level concepts that can be formalised in logic. NeSyFOLD first binarises the last layer by thresholding
677 the activation of each filter given an input, creating a tabular dataset. Then, the framework makes use of
678 the FOLD-SE algorithm, which turns tabular data into default ASP rules ([Wang and Gupta 2023](#)). The
679 resulting program approximates the decision-making of the last layer of the CNN by treating each filter
680 as an atom that is used in the rule set. These atoms do not have human-readable names, which is why a
681 semantic labelling step is necessary. An oracle, such as a human or foundation model, gives each atom a
682 name by looking at the parts of the images that are activated by the filter that the atom represents. This
683 results in a neuro-symbolic model that mostly maintains the predictive power of the CNN while being
684 explainable and human-readable.

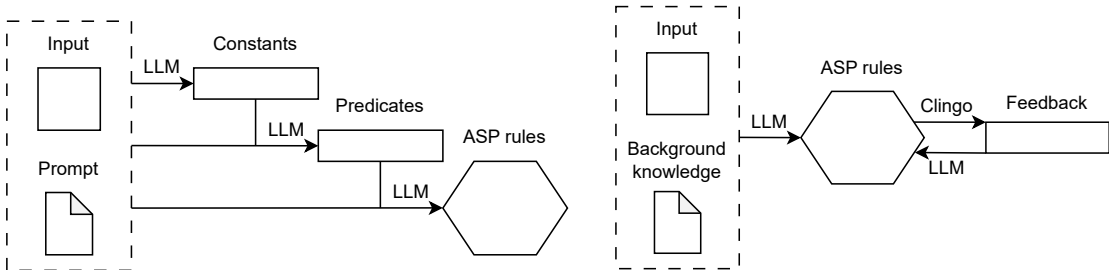
685 The authors have expanded the framework multiple times. NeSyFOLD-G ([Padalkar et al. 2023](#)) is
686 a variant which groups similar kernels together before binarising them. This reduces the number of
687 generated rules and therefore increases interpretability. NeSyBiCor ([Padalkar et al. 2025](#)) introduces the
688 ability to remove biases in a CNN based on the rules extracted by NeSyFOLD. The user can tag undesired
689 concepts in those rules that should not be used to make decisions. For example, the CNN might use the
690 colour of the sky for predicting the type of road in an image, which is irrelevant. The framework then
691 finetunes the CNN using a semantic similarity loss to push it away from making predictions with such
692 undesired concepts. This process largely maintains the accuracy of the rule set and often reduces the
693 number of rules.

694 The remaining frameworks in this section utilise LLMs to create ASP rules. Their goal is to augment
695 the reasoning capabilities of LLMs by encoding tasks in ASP instead of solving them directly. The
696 remarkable ability of LLMs to produce structured languages such as Python code has been widely
697 demonstrated in literature. However, as their training sets include much less ASP than Python, LLMs
698 struggle to generate correct answer set programs from scratch. Different approaches therefore propose
699 various methods to improve LLM capabilities.

700 *LLASP*. The conceptually simplest approach involves finetuning the LLM to output an ASP program
701 directly, as is done in LLASP ([Coppolillo et al. 2024](#)). The authors first systematically evaluate the ASP
702 capabilities of LLMs and find them to be inadequate in terms of syntactic and semantic correctness. To
703 alleviate this problem, they perform supervised finetuning on lightweight LLMs using an ad-hoc dataset.
704 This dataset consists of fundamental ASP programming patterns, such as constraints, joins, preferences
705 or filtering. Despite being smaller in model size, LLASP is able to generate ASP rules in one shot, as
706 illustrated in Figure 4c. On basic tasks, it achieves 89% semantic and syntactic accuracy on the ad-
707 hoc dataset, outperforming much larger LLMs. When including combined problems, however, there are
708 mixed results, indicating room for improvement in the field.

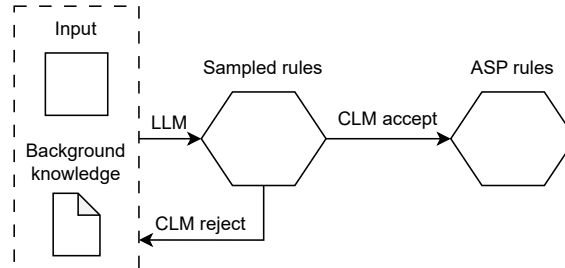
709 *NL2ASP*. Rather than outputting ASP directly, [Santana et al. \(2024\)](#) use an intermediate representation
710 called controlled natural language (CNL). Figure 4d illustrates the two-step process. First, a neural

711 network such as BART (Lewis et al. 2020) transforms natural language input into CNL. Second, the tool
 712 CNL2ASP (Caruso et al. 2023) translates the CNL sentences into ASP. The authors create a dataset to test
 713 out NL2ASP, consisting of ASP encodings taken from competitions and online resources and manually
 714 translated into CNL. To improve performance, they finetune the pre-trained neural network on CNL
 715 labels. NL2ASP consistently achieves BLEU scores over 0.9 and is able to produce 99% syntactically
 716 correct CNL statements with an F1 score of 0.93. The intermediate translation into CNL is an easier task
 717 than directly outputting ASP, as CNL is a higher-level language. Nevertheless, it still supports the main
 718 ASP constructs, such as facts, strong and weak constraints or choice rules.



(a) High-level depiction of the symbolic learning procedure in GPT-ASP (Ishay et al. 2023).

(b) High-level depiction of the symbolic learning procedure in DSPy-ASP (Wang et al. 2024) and LLM-ARC (Kalyanpur et al. 2024).



(c) High-level depiction of the symbolic learning procedure in CLM-ASP (Kaur et al. 2025).

Figure 5. More high-level depictions of symbolic learning procedures in frameworks with pre-trained neural components.

719 **GPT-ASP.** Ishay et al. (2023) devise a four-step method for converting natural language logic puzzles
 720 into answer set programs. Their framework GPT-ASP first generates constants, then predicates and then
 721 rules, as Figure 5a illustrates. At each step, the LLM has access to both the input and the ASP generated
 722 in the previous steps. The rule generation step is split up into two parts. First, the LLM generates choice
 723 rules, which increase the number of answer sets. Then, it creates constraints, which limit the number of
 724 answer sets again. This process is similar to how humans model problems in ASP. The pipeline allows

725 you to spot and correct errors easily by inspecting the constructed constants, predicates and rules. This is
726 not possible with LLM-only models that simply output the answer.

727 The capabilities of LLMs to generate ASP code can be further improved by introducing feedback loops.
728 Two frameworks make use of this technique, which is shown in Figure 5b.

729 *DSPy-ASP*. In the DSPy-ASP framework (Wang et al. 2024), the LLM can revise the ASP rules for
730 three iterations. First, it generates ASP predicates and queries, which are input to the Clingo solver
731 together with predefined knowledge modules. Second, the LLM then revises the answer set program
732 based on feedback from the solver, such as error messages. This process is repeated three times. While
733 the LLM only generates predicates at first, it does have the ability to revise and generate new rules during
734 the feedback loops. The authors use the DSPy Python framework (Khattab et al. 2024) to automate
735 the prompt engineering process and show that adding feedback loops further improves task success.
736 Compared to direct-prompting, the addition of ASP results in significant accuracy increases of up to 50%
737 in spatial reasoning tasks.

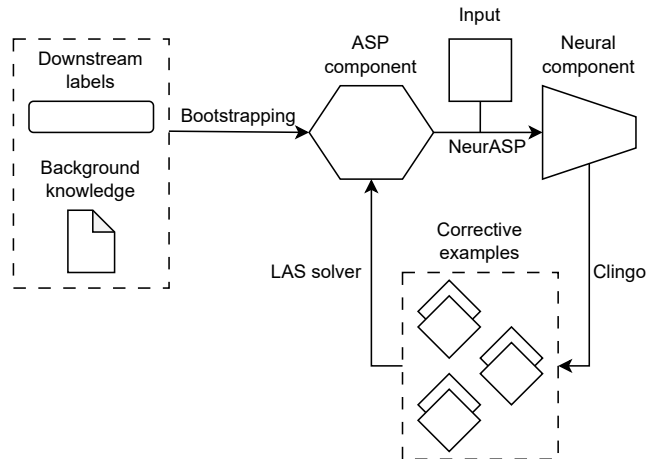
738 *LLM-ARC*. Kalyanpur et al. (2024) go further and use LLMs to generate tests in addition to ASP rules
739 in their LLM-ARC framework. These tests are meant to verify the semantic correctness of the code. Just
740 like the generated ASP rules, they are run through Clingo, which provides feedback through its error
741 messages. The authors provide a simple schema for specifying tests, including mechanisms for checking
742 that a proposition is true in any, all or no answer sets. The LLM can then correct the code and tests in
743 an iterative manner, until everything compiles and all tests pass. The prompt includes few-shot examples
744 of how to solve questions from the benchmark, including how to write tests and correct errors. Even
745 though there are no guarantees that the generated tests are semantically sound, the authors show that they
746 improve the correctness of the generated ASP rules in practice.

747 *CLM-ASP*. Kaur et al. (2025) apply conformal language modelling (CLM) to improve the capability of
748 LLMs to produce answer set programs. CLM is a rejection sampling method with statistical guarantees
749 that assesses the output of LLMs based on certain criteria. If the criteria are not met, CLM rejects the
750 output and samples from the LLM again, as Figure 5c illustrates. The CLM check consists of two steps:
751 First, all unacceptable samples that do not pass the admission function are filtered out. In the case of
752 CLM-ASP, the admission function checks that the output is syntactically correct using Clingo. Second,
753 the output is evaluated through the metrics of quality, diversity and confidence. The authors experiment
754 with transition scores and $ROUGE_L$ metrics, as well as leveraging another LLM as the judge. The latter
755 method involves finetuning an LLM to assess the quality of ASP and leads to better results than the
756 other metrics. Finally, in-context learning is used to guide the creation of ASP. The prompt includes
757 hand-written ASP rules that are needed to solve the given tasks.

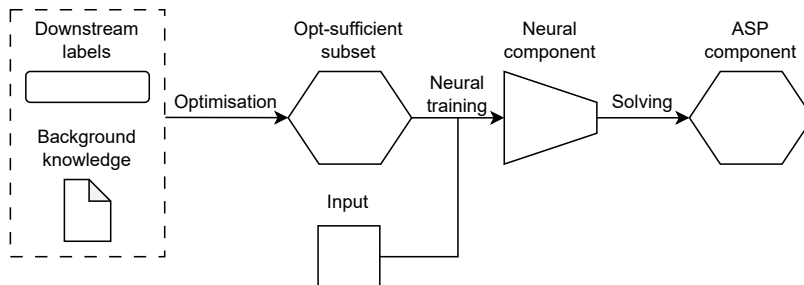
758 To sum up, there are two main approaches to learning ASP rules from raw data: Converting the data into
759 symbolic examples to use with off-the-shelf solvers or generating rules directly. Both strands have been
760 influenced by the rise of foundation models. In the former, VLMs act as the perception component,
761 extracting predicates from images. In the latter, LLMs additionally act as the reasoning component,
762 writing and improving ASP rules to solve a task using in-context learning. Foundation models have
763 improved the scope and accuracy of neurosymbolic ASP methods. In turn, ASP has improved the logical
764 capabilities of LLMs, which by themselves struggle with reasoning.

765 *Joint learning of neural and symbolic components*

766 Training the neural component and learning ASP rules at the same time is a very challenging task, because
 767 it resembles a chicken-and-egg problem. The neural network does not have latent labels to train with, as
 768 there is no answer set program that can generate them. And the ASP component does not have latent
 769 concepts to learn from, as the neural network is not trained yet. There are two papers in the literature that
 770 have tried to overcome these challenges without using any pre-trained components and we will discuss
 771 them in this section.



(a) High-level depiction of the learning procedure in NSIL (Cunnington et al. 2023b).



(b) High-level depiction of the learning procedure in NeuralFastLAS (Charalambous et al. 2023).

Figure 6. High-level depictions of frameworks that jointly learn the neural and ASP component.

772 *NSIL*. [Cunnington et al. \(2023b\)](#) address the chicken-and-egg problem in their NSIL framework by
 773 bootstrapping a hypothesis, as Figure 6a illustrates. The bootstrapping task takes only the downstream
 774 labels and background knowledge into account. It is set up using WCDPIs, each of which contains a

775 downstream label in the inclusion set and choice rules for the neural atoms in the context. A LAS solver
776 then finds ASP rules that cover as many WCDPIs as possible.

777 For tasks like MNIST Addition, bootstrapping works well. In the paper, the mode bias includes
778 functions for addition, subtraction and multiplication. To maximise coverage of WCDPIs, the LAS solver
779 would choose the addition function, because it is the only one that can cover all 20 downstream labels.
780 A subtraction function of two positive digits cannot cover the labels 10 to 19, while a multiplication
781 function cannot arrive at prime numbers like 13. For more complex tasks, the initial hypothesis might
782 be wrong or incomplete, which is where the iterative nature of NSIL comes into play. The bootstrapped
783 hypothesis is used to train the neural component with NeurASP. The freshly trained neural predictions
784 are turned into corrective examples to learn a better hypothesis using FastLAS or ILASP. At this point,
785 the loop starts again by training the neural network using the new hypothesis.

786 Splitting up the process into a neural and symbolic component allows NSIL to solve NP-complete
787 tasks like the hitting set problem. It plays to the strengths of both paradigms, at the expense of a difficult
788 learning regime.

789 *NeuralFastLAS*. Charalambous et al. (2023) introduce NeuralFastLAS, which learns ASP rules and
790 trains a neural network jointly in one iteration, as shown in Figure 6b. First, the framework constructs
791 a set of ASP rules that can prove the downstream labels for each example, given all possible choices
792 of neural atoms. Using the background knowledge and constraints such as symmetry, this set of rules is
793 pruned in the optimisation step to obtain the opt-sufficient subset. Crucially, the opt-sufficient subset is
794 proven to contain the optimal symbolic solution. The answer sets stemming from the opt-sufficient subset
795 are then used as noisy latent labels to train the neural network. Since there are many different rules in the
796 opt-sufficient subset, the neural component has a second head that computes a posterior probability for
797 each rule. After the network has been trained, an optimal hypothesis is found given the neural network
798 predictions and rule posteriors. The paper proves the theoretic correctness of the method and provides
799 conditions for the guaranteed convergence of the neural network. As the name suggests, the framework is
800 modelled after FastLAS and thereby inherits its expressive power, which is limited to stratified programs.

801 The papers in this section tackle true neurosymbolic learning without any pre-training and very limited
802 background knowledge. Both frameworks break down complex problems into a neural and symbolic
803 part and try to solve both simultaneously - a very difficult task. They start the process by bootstrapping
804 rules, either computing a single hypothesis or a space of possible hypotheses. While this approach works
805 for simple examples, it tends to get stuck in local minima and is limited in its scalability. Much more
806 research is needed to find algorithms that efficiently traverse the search space of possible rules while
807 training neural components at the same time.

808 **Analysis of neurosymbolic ASP**

809 The field of neurosymbolic ASP contains a diverse set of frameworks, datasets and metrics. In this
810 section, we aim to synthesize results across the landscape to identify the capabilities, strengths and
811 limits of current methods. We start by evaluating the perception tasks used in benchmarks and conduct a
812 comparative analysis of the performance of different frameworks. We find that neurosymbolic methods
813 are often able to outperform fully neural methods not only in explainability, but also accuracy. However,
814 we identify some barriers to progress, including simple perception inputs, a gap in challenging but

815 solvable datasets, a limited capacity to generate ASP and scalability issues. Throughout this section,
816 we propose ways forward to address these limitations.

817 *Perception tasks*

Synthetic images

MNIST (Deng 2012)	Embed2Sym, NeurASP, SLASH, dPASP, FFNSL, NeSyGPT, Embed2Rule, NSIL, NeuralFastLAS
ShapeWorld (Kuhle and Copestake 2017)	SLASH
CLEVR (Johnson et al. 2017)	ASP-VQA, AQuA, SLASH, NeSyGPT
CLEGR ^V (Bauer et al. 2025)	NSGRAPH

Real-world images

CIFAR-10 (Krizhevsky and Hinton 2009)	Embed2Sym
VQAR (Huang et al. 2021)	SLASH
Playing cards (Cunnington et al. 2023a)	FFNSL, NeSyGPT, Embed2Rule
PlantVillage (Hughes and Salathe 2016)	FFNSL
Indoor scenes (Quattoni and Torralba 2009)	FFNSL
PlantDoc (Singh et al. 2020)	NeSyGPT
Places (Zhou et al. 2018)	NeSyFOLD
German traffic signs (Stallkamp et al. 2012)	NeSyFOLD

Natural language text

bAbI (Weston et al. 2015)	[LLM]+ASP, LLM2LAS
CLUTRR (Sinha et al. 2019)	[LLM]+ASP
gSCAN (Ruis et al. 2020)	[LLM]+ASP
StepGame (Shi et al. 2021)	[LLM]+ASP, DSPy-ASP, CLM-ASP
SpartQA (Mirzaee et al. 2021)	DSPy-ASP
Logic grid puzzles (Mitra and Baral 2015)	GPT-ASP
FOLIO (Han et al. 2024)	LLM-ARC
QuaRel (Tafjord et al. 2019)	STAR
ASP patterns (Coppolillo et al. 2024)	LLASP
NL2CNL (Santana et al. 2024)	NL2ASP

Table 1. Datasets used to train the perception components and the corresponding frameworks.

818 Table 1 provides a summary of all types of input that papers have used to test their frameworks. The
819 datasets can be split up into three categories: synthetic images, real-world pictures, and natural language
820 text. Within these categories, there are discrepancies about the difficulty of the perception task. In general,
821 frameworks that incorporate more learning tend to use simpler perception tasks, while frameworks with
822 pre-trained components can handle more realistic inputs.

823 *Synthetic images.* The MNIST dataset consists of 28x28 pixel greyscale images of handwritten
824 digits (Deng 2012). It was created in 1994 and formed the basis for testing one of the first convolutional
825 neural networks, LeNet-5, which already achieved 99% accuracy (Lecun et al. 1995). As such, it is
826 considered one of the easiest perception datasets and even very small neural networks can learn to
827 predict it perfectly. Nine of the 13 frameworks that take images as inputs are tested on MNIST images,
828 despite its simplicity. The main reason is that the image classification forms only the first part of a more
829 complex neurosymbolic task, such as addition or set membership. To convincingly demonstrate the real-
830 life viability of these neurosymbolic frameworks, however, more realistic perception tasks are needed.

831 The ShapeWorld dataset is a step above MNIST, consisting of two-dimensional shapes with various
832 orientations and colours against a white background (Kuhnle and Copestake 2017). The SLASH
833 framework uses a variant with up to four shapes and trains a CNN to recognise properties of the objects
834 (colour, shape, shade and size). The downstream label is the combination of all properties, for example
835 `has_attributes(object1, red, circle, bright, small)`. Unlike in tasks like MNIST Addition, the
836 downstream label therefore includes all four latent labels without obfuscating them. Therefore, the neural
837 component is trained directly on the latent labels and the symbolic component only collects all latent
838 attributes into one predicate. While the perception task with ShapeWorld is more difficult, the actual
839 neurosymbolic task is easier than for MNIST tasks.

840 In the CLEVR dataset, the shapes are three-dimensional and can partially occlude each other (Johnson
841 et al. 2017). Again, SLASH trains the neural component directly on the latent attributes of the objects.
842 ASP-VQA and AQUA use a pre-trained YOLO network and do not perform any learning at all. NeSyGPT
843 uses a VLM instead, which is pre-trained on a large corpus of general images. In addition, the authors
844 finetune it with a small number of latent labels. Similarly, the CLEGR^V dataset generates images of
845 graphs deterministically from their specifications. This results in limited variety and a straightforward
846 perception task, which NSGRAPH solves using pre-trained graph recognition models.

847 In the category of synthetic images, only MNIST is truly used for neurosymbolic training of the
848 perception component. Most MNIST tasks include just the downstream labels and therefore only provide
849 noisy signals for the latent classification task. ShapeWorld, CLEVR and CLEGR^V, while embodying a
850 more complex perception task, provide direct latent labels in all cases.

851 *Real-world images.* CIFAR-10 is a dataset of real-world images compressed to 32x32 pixels that depict
852 one of ten object categories (Krizhevsky and Hinton 2009). Embed2Sym uses it for the CIFAR-10
853 Addition task, where each image category is arbitrarily assigned a number and the goal is to find the
854 sum of two images. Just like in MNIST Addition, no latent labels are given, but the perception task is
855 more difficult.

856 The Visual Question Answering and Reasoning (VQAR) dataset contains diverse real-world images
857 with a much higher resolution than CIFAR-10 (Huang et al. 2021). But the SLASH framework uses a
858 pre-trained network to find the bounding boxes of objects, sidestepping the difficult task of training the
859 neural component.

860 The Playing cards dataset consists of photos of real playing cards with a resolution of 523x831
861 pixels (Cunnington et al. 2023a). It is used to learn the rules for determining the winner of card games.
862 None of the three papers that use the dataset actually train the neural component from downstream labels.
863 FFNSL pre-trains the neural component on the latent playing card labels directly, while Embed2Rule and
864 NeSyGPT use a VLM with latent finetuning.

865 The same is the case for the PlantVillage (Hughes and Salathe 2016) and Indoor scenes (Quattoni
866 and Torralba 2009) datasets. While they contain real-world images of diseased crops and varied indoor
867 rooms, FFNSL uses a pre-trained neural network. For the PlantDoc dataset, which contains images of
868 diseased plants (Singh et al. 2020), NeSyGPT uses a VLM with latent finetuning.

869 The Places dataset contains images of indoor and outdoor scenes (Zhou et al. 2018) and is used by
870 the NeSyFOLD and NeSyBiCor frameworks. The downstream labels represent scenes, while the latent
871 concepts are objects in the image. NeSyFOLD first trains a CNN to predict the scene category and then
872 extracts rules from the CNN’s last layer. Each atom in those rules represents a (set of) filter activations,
873 which roughly correspond to objects in the image. The names of the objects are provided through manual
874 annotation of segmentation masks. Therefore, the framework is not able to classify latent concepts
875 autonomously. These papers use the same techniques for the German traffic sign dataset (Stallkamp
876 et al. 2012).

877 While many of the datasets in this section display realistic scenes and objects, they are not used in
878 a neurosymbolic training regime. Instead, almost all frameworks either pre-train or finetune the neural
879 components on image labels directly, which amounts to a basic classification task. The notable exception
880 is CIFAR-10, where latent symbols are extracted automatically without labels in the Embed2Sym
881 framework.

882 *Natural language text.* The last category of datasets comprises collections of natural language texts,
883 which are used by the LLM-based neurosymbolic frameworks. They all take the form of logical tasks or
884 puzzles, making them well-suited for translating into ASP.

885 bAbI is a collection of questions that involve skills such as counting, path-finding or negation to solve.
886 The questions are generated from a simulation of entities and actions, which are turned into natural
887 language using a simple automated grammar (Weston et al. 2015). Therefore, the dataset is not too
888 complex. CLUTRR poses the task of inferring family relations from short stories. It is more realistic
889 than bAbI, as the stories were written by crowd-workers, who generated narratives from the generated
890 kinship facts (Sinha et al. 2019). gSCAN represents a grid world in JSON format and asks natural
891 language questions about how to achieve a goal. The questions are very direct instructions without much
892 language variability and the answer comes in the form of a sequence of actions (Ruis et al. 2020). In all
893 these datasets, [LLM]+ASP uses additional hand-written knowledge modules to solve the tasks. Thus,
894 even though many of the tasks are complex, the framework requires substantial manual engineering.
895 LLM2LAS manages to solve bAbI questions with more minimal background knowledge, but the pipeline
896 only works for 15 out of the 20 tasks.

897 StepGame contains questions that require multi-hop spatial reasoning. It consists of descriptions of
898 entities and their spatial relationships in a grid-based world and asks queries about their relative positions.
899 These descriptions are generated automatically, but utilise different ways to describe spatial relations
900 from a set of crowdsourced synonyms (Shi et al. 2021). Both [LLM]+ASP and DSPy-ASP make use of
901 hand-written knowledge modules to solve the task, limiting their applicability in real-world scenarios.
902 While CLM-ASP makes the LLM generate the entire answer set program, all the rules necessary are
903 included in the prompt. This leads to the curious scenario that most errors are caused by the LLM
904 incorrectly copying the rules from the prompt. A more complex benchmark is SpartQA, which consists
905 of quantifier-based reasoning around blocks and objects, generated automatically (Mirzaee et al. 2021).
906 DSPy-ASP beats LLMs in terms of accuracy, but again at the expense of requiring manually specified
907 ASP knowledge modules.

908 The logic grid puzzles dataset provides a set of categories, each containing an equal number of
909 elements. The aim is to match elements based on clues given in the question. Since the dataset was
910 compiled from a puzzle website, the questions are presumably human-made (Mitra and Baral 2015). GPT-
911 ASP manages to achieve a high accuracy, unlike LLM-only methods, without relying on any hand-crafted
912 ASP. Instead, all the necessary knowledge is encoded within the prompt, which contains instructions and
913 a few solved examples from the dataset.

914 FOLIO consists of a set of premises and a conclusion. The task is to determine whether the conclusion
915 is true, false or uncertain. It was created by experts in 2024 to challenge the state-of-the-art language
916 models of the time and includes logically complex tasks written in natural language (Han et al. 2024).
917 LLM-ARC manages to outperform LLM-only models without needing any hand-written ASP rules.

918 QuaRel is a set of commonsense physics questions based on properties like friction, speed or time.
919 The questions were crowdsourced by asking people to come up with imaginative scenarios for the given
920 relations (Taffjord et al. 2019). The STAR framework uses LLMs to extract predicates, while modelling
921 the commonsense knowledge by hand in ASP. Both the ASP patterns (Coppolillo et al. 2024) and
922 the NL2CNL (Santana et al. 2024) tasks are ad-hoc datasets created for their respective framework.
923 The former represents classic ASP patterns, which are encoded as templates and can be combined for
924 more complex questions. Patterns include guessing, constraints, joins, transitive closure, preferences and
925 filtering. For the latter, the authors collected ASP problems from competitions, lecture notes and online
926 resources. They hand-crafted their CNL and natural language representation to create the questions and
927 labels. In both cases, the tasks represent standard toy problems, rather than real-world applications.

928 In summary, most natural language datasets were generated synthetically, with some using crowd-
929 workers to enrich the questions. Since they are based on simulations or structured graphs, transforming
930 them into ASP is more straightforward than with true natural language inputs. The notable exception is
931 FOLIO, which was written by experts with the goal of providing a challenging and varied benchmark.
932 It is therefore impressive that LLM-ARC achieves state-of-the-art results on it without using any hand-
933 written ASP knowledge. Another potential issue is that all datasets other than FOLIO were created before
934 the advent of LLMs and have been released publicly. As LLMs are trained on large amounts of publicly
935 available data, it is possible that they have seen these datasets in their training procedure. For a fair
936 analysis of LLM-based methods, authors should make sure to use datasets that the model could not have
937 encountered before.

938 All in all, methods with more complex neurosymbolic requirements are limited to more rudimentary
939 perception datasets. Only the two simplest visual datasets, MNIST and CIFAR-10, are used for training
940 the neural component without latent labels. For any more advanced perception tasks, neurosymbolic
941 frameworks either train their neural component directly or finetune a VLM with latent labels. The
942 current limit for true neurosymbolic learning with ASP therefore lies with 32x32 pixel images with 10
943 categories. For textual inputs, most datasets are synthetically generated and the majority of frameworks
944 need additional hard-coded ASP knowledge. Only LLM-ARC generates all ASP autonomously and has
945 been tested on a challenging, real-world dataset with FOLIO. However, as LLMs are pre-trained on vast
946 amounts of data, they might have encountered the same or similar problems during training.

947 *Performance analysis*

948 The experimental sections of the literature offer insights into the viability of proposed methods on
 949 different datasets. In this section, we consolidate the results from numerous papers and tasks and provide
 950 a comparative analysis. We have compiled these results from the original papers running their own
 951 framework, as well as papers running other frameworks for comparison. We find that neurosymbolic
 952 frameworks vary in terms of capabilities, but most are capable of achieving accuracies of 80–100% on
 953 benchmark tasks. Combining ASP with LLMs improves performance on logical tasks across the board,
 954 but the literature lacks comparisons between neurosymbolic methods.

Task	Neur ASP	SLASH	dPASP	Embed 2Sym	NSIL	Neural FastLAS	NeSy GPT	FFNSL
Add 2	98	99	94	98	98	–	–	–
Add 3	98	99	–	–	8	93	–	–
Add 4	98	99	–	94	–	–	–	–
Add 6	T/O	–	–	94	–	–	–	–
Add 8	T/O	–	–	92	–	–	–	–
Add 30	T/O	–	–	66	–	–	–	–
$a \times b + c$	–	–	–	–	87	97	–	–
Even9Plus	91	81	–	89	90	98	95	90
1k labels	90	85	–	72	88	–	95	88

Table 2. Accuracy (%) comparison of different frameworks on MNIST arithmetic tasks. T/O = Timeout.

955 *MNIST arithmetic.* Table 2 presents accuracies on MNIST arithmetic tasks, which take MNIST digits
 956 as inputs and perform arithmetic operations in the symbolic component. Among them, addition is the
 957 most basic variation and involves calculating the sum of 2 up to 30 digits. Another variation includes
 958 multiplication of the first two digits added to the third. In Even9Plus, the result is the second value if the
 959 first value is even, or 9 plus the second value otherwise. The 1k labels task involves two-digit addition
 960 with training data limited to 1000 examples.

961 Eight frameworks report results on MNIST arithmetic tasks, demonstrating its role as an established
 962 benchmark for neurosymbolic ASP. As the first row in Table 2 shows, two-digit addition is a solved
 963 problem and all frameworks manage to achieve over 90% accuracy. Even9Plus and $a \times b + c$ are more
 964 challenging, with some frameworks dropping below 90%. Both NeurASP and SLASH do very well up
 965 until 4 digits in MNIST addition. From 6 digits onwards, NeurASP times out, whereas Embed2Sym
 966 manages to scale all the way up to 30. The reason lies in its structure: Embed2Sym trains a neural
 967 network end-to-end and then clusters the latent space. While the number of inputs increases, the number
 968 of clusters stays at 10, which remains manageable. However, the reliance on neural networks makes
 969 Embed2Sym less data-efficient, as it achieves the lowest accuracy among all frameworks tested on 1,000
 970 labels.

971 Both NSIL and NeuralFastLAS learn the rules of the task while training their neural component
 972 simultaneously. This is a much harder challenge, but they manage to keep up with other frameworks
 973 on most tasks. The exception is 3 digit addition, where NSIL only achieves 8% accuracy. In general,
 974 NeuralFastLAS is superior, but it is more limited than NSIL in the rules it can learn. NeSyGPT and
 975 FFNSL are pre-trained on latent labels and therefore unsurprisingly perform well.

Task	NeurASP	Embed2Sym
Member 3	97	97
Member 4	T/O	97
Member 5	T/O	98
Member 20	T/O	93

Table 3. Accuracy (%) comparison of different frameworks on member tasks with specified number of entries. T/O = Timeout.

976 *Member.* Table 3 presents accuracies for the member task. Given a list of MNIST images and a symbolic
 977 digit, it involves stating whether the digit is included in the list of images. The downstream label is
 978 boolean and therefore only provides a sparse signal: If *false*, each image could represent any value except
 979 the given digit. If *true*, each image could represent any value, but at least one of them must be the given
 980 digit. Every added image in the list increases the number of possible latent labels by a factor of 10.
 981 NeurASP times out at a length of 4, where there are about $10^4 = 10,000$ answer sets per example.
 982 Embed2Sym scales much better again because it does not need to enumerate all latent possibilities. This
 983 illustrates the inherent scaling advantage of fully neural training well.

Task	SLASH	NSIL	Embed2Rule	NeurASP	Embed 2Sym	NeSyGPT	FFNSL
MNIST 5/4	50	100	99	–	–	–	–
MNIST 10/6	T/O	T/O	93	–	–	–	–
FashionMNIST	–	88	–	–	–	–	–
PlantDoc	–	–	–	64	77	99	98

Table 4. Accuracy (%) comparison of different frameworks on hitting set tasks. T/O = Timeout.

984 *Hitting sets.* Table 4 presents accuracies for different variations of the hitting set problem. In this task,
 985 you are given a universe of elements U and a set of sets S , where each $T \in S$ only contains elements
 986 from U . A hitting set $H \subseteq U$ contains elements such that for every $T \in S, T \cap H \neq \emptyset$. Given an integer
 987 $k \geq 1$, the hitting set task requires determining whether there exists a hitting set H with $|H| \leq k$ (Aspis
 988 et al. 2024). This task showcases the benefit of ASP’s added expressivity, as it requires default negation
 989 to solve. Table 4 includes four variants: The first two use MNIST images as inputs with 5 elements in U
 990 and 4 input images, or 10 elements in U and 6 input images respectively, where $k = 2$. FashionMNIST
 991 contains greyscale images of clothing items and its hitting set task contains 5 elements in U with up to 4
 992 sets per example and $k = 2$. Lastly, the PlantDoc hitting set task contain 38 elements in U with up to 5
 993 sets per example and $k = 2$.

994 For this task, NeSyGPT, FFNSL and Embed2Rule fare much better than NeurASP, SLASH,
 995 Embed2Sym or NSIL. There is a clear divide between the former frameworks, which train their neural
 996 component, and the latter, which use pre-trained neural networks but learn the rules. This suggests that
 997 the ruleset itself is not too complicated, but does not provide an effective learning signal for the neural
 998 component. Therefore, hitting sets can serve as a good benchmark for further research.

Task	NeurASP	SLASH	Embed 2Sym	NSIL	NeSyGPT	FFNSL	Embed 2Rule	ASP- VQA	AQuA
CIFAR10 add	T/O	–	85	–	–	–	–	–	–
ShapeWorld	–	85	–	–	–	–	–	–	–
Crop yield	–	–	–	–	–	100	–	–	–
Indoor scenes	–	–	–	–	–	100	–	–	–
Sudoku 4x4	–	98	–	50	–	100	98	–	–
Sudoku 9x9	–	T/O	–	50	–	100	78	–	–
Follow suit 4	25	25	–	25	100	100	86	–	–
Follow suit 10	10	10	17	10	100	100	–	–	–
CLEVR	99	90	–	–	99	–	–	97	94

Table 5. Accuracy (%) comparison of different frameworks on miscellaneous benchmarks. T/O = Timeout.

999 *Miscellaneous tasks.* Other prominent tasks from the literature are summarised in Table 5. CIFAR10
1000 addition is reminiscent of two-digit MNIST addition, but the input images come from the CIFAR10
1001 dataset. Each category is arbitrarily assigned a digit. The symbolic task in ShapeWorld amounts to
1002 predicting the attributes of each object, which is equivalent to a supervised classification problem. In crop
1003 yield prediction, the symbolic task defines the quality of the yield as poor, moderate or strong depending
1004 on the the crop’s location, species and health. While the location is given, the species and health must
1005 be predicted from the input image. Indoor scene classification maps room classes (e.g. bathroom) into
1006 higher-level superclasses (e.g. home). The two Sudoku variants encode the grid validity task on boards
1007 of size 4x4 and 9x9. Given a sequence of MNIST images representing the board entries, the downstream
1008 label indicates whether the board is valid. In the follow suit game, the framework must determine the
1009 winner given 4 or 10 images of playing cards. The winner is the card that matches the suit of card 1 and
1010 has the highest rank among all matches. CLEVR involves determining the attributes of objects in a scene,
1011 such as like colour or material, and answering reasoning questions about them.

1012 Despite the similarity to MNIST addition, NeurASP fails to learn CIFAR10 addition and times out.
1013 SLASH is able to achieve a high accuracy for ShapeWorld and FFNSL gets 100% on crop yield
1014 and indoor scenes. Note that FFNSL reports the accuracy of the learned hypothesis, rather than the
1015 downstream accuracy, which would take neural prediction errors into account. While NeurASP can learn
1016 4x4 Sudoku grid validity, it times out for the larger 9x9 version. NSIL struggles in both cases and is
1017 unable to learn a valid hypothesis or train the neural network. Only FFNSL and Embed2Rule, which
1018 make use of pre-trained neural components, are able to achieve high accuracies on Sudoku.

1019 A similar picture emerges for follow suit, where the predictions of NeurASP, SLASH and NSIL all
1020 amount to random guessing. This happens because the learning signal from the downstream label is
1021 extremely limited. It only provides the index of the winner, and there can be millions of card combinations
1022 leading to the same winner in the 4-player version alone. In the 10-player version, there are up to $\binom{52}{10}$
1023 possible latent label combinations per downstream label. Embed2Sym accomplishes a slightly better
1024 accuracy of 17% for Follow suit 10, but this is still much lower than for Member 20 or MNIST addition
1025 30. That might be surprising, given that Follow Suit 10 has only has 10 inputs compared to 20 and 30
1026 for the other tasks. The main difference is that playing cards have 52 unique values, one for each rank-
1027 suit combination. This requires Embed2Sym to create 52 clusters, as opposed to 10 for MNIST digits.

1028 Embed2Sym scales well with regard to the number of inputs, but struggles when the number of clusters
 1029 increases. Finally, all frameworks achieve at least 90% accuracy on CLEVR across the board.

1030 Neurosymbolic methods offer advantages beyond increases in accuracy. FFNSL maintains a better
 1031 performance under distributional shifts than a neural network alone. NeSyFOLD obtains accuracies
 1032 of 92%, 67% and 44% for increasingly difficult variations of Places and 78% for the German traffic
 1033 signs. This is actually a slight drop in performance to the original CNN. However, NeSyFOLD enables
 1034 explanations of decisions and can correct biases.

Task	[LLM]+ASP	LLM2LAS	DSPy-ASP	CLM-ASP
bAbI	100	100	–	–
StepGame 1	93	–	94	80
StepGame 5	93	–	88	65
StepGame 10	88	–	80	–
StepGame 15	–	–	–	1

Table 6. Accuracy (%) comparison on bAbI and StepGame reasoning benchmarks with varying steps.

1035 *Natural language tasks.* Table 6 compares results for LLM-based frameworks on natural language
 1036 tasks. [LLM]+ASP and LLM2LAS both achieve 100% accuracy on bAbI, indicating that this dataset
 1037 now easily solvable by modern architectures. In StepGame, the task becomes more difficult with an
 1038 increasing number of hops. Interestingly, [LLM]+ASP, which employs a straightforward fact extraction
 1039 using a single prompt, outperforms the more involved DSPy-ASP. CLM-ASP struggles even more, but
 1040 unlike the other two it does not make use of hard-coded ASP rules.

1041 bAbI and StepGame are the only two datasets that have been used by multiple papers. In all other
 1042 cases, each paper has selected a different task on which to test their method, which illustrates a lack
 1043 of agreed-upon baselines. Moreover, no LLM-based framework has compared its method with those of
 1044 others in the literature. Instead, they report improvements over prompting the LLMs directly without use
 1045 of ASP.

1046 Yang et al. (2023) show that GPT3 alone can solve 80% of bAbI tasks with few-shot prompting, and
 1047 86% with chain-of-thought prompting. While GPT-3 achieves 32% fewer accuracy points for 1 reasoning
 1048 hop in StepGame, the difference grows to 57% at 10 hops. Wang et al. (2024) compare their DSPy-based
 1049 approach with Deepseek, Llama3 and GPT-4.0 mini, reporting similar gaps for StepGame. Kaur et al.
 1050 (2025) report 45%, 29% and 0% accuracy on StepGames 1, 5 and 15 respectively for a simple prompt-
 1051 based approach with Llama3, much lower than CLM-ASP. The reason CLM-ASP fails on 15 hops is that
 1052 the authors only calibrated the LLM with examples of up to 5 hops, indicating a lack of generalization
 1053 capabilities in LLM-based ASP generation.

1054 Other baselines exhibit similar differences. [LLM]+ASP clearly outperforms non-LLM-based methods
 1055 by 20–30%, achieving 91% accuracy on CLUTRR and 100% on gScan. For SpartQA, the results are
 1056 closer together: Deepseek achieves 60%, Llama3 55% and GPT4.0 mini 56%. Rajasekharan et al.
 1057 (2023b) run tests with two versions of GPT3 with 7 billion and 175 billion parameters. The larger
 1058 GPT model actually achieves almost the same accuracy as STAR on QuaRel (91 vs 94%). For the 7B
 1059 version, however, STAR improves GPT3’s accuracy from 53 to 86%. This suggests that neurosymbolic
 1060 methods can help smaller models to match the performance of bigger models. The grid puzzles dataset
 1061 poses the greatest challenge for vanilla LLMs, with accuracies of only 21% for GPT4, compared to

1062 the 92% achieved by GPT-ASP (Ishay et al. 2023). For FOLIO, GPT4 obtains 74% using chain-of-
1063 thought prompting, considerably lower than the 88% achieved by LLM-ARC (Kalyanpur et al. 2024).
1064 Rajasekharan et al. (2023b) have conducted qualitative testing using real user input on their chatbots
1065 based on STAR. They conclude that STAR performs better than vanilla LLMs in metrics such as staying
1066 on topic or providing relevant responses. STAR can also generate justifications for each answer in form
1067 of a proof tree, a feature lacking in fully neural methods. Lastly, NSGRAPH achieves 81% accuracy on
1068 small graphs (median of 10 nodes and 8 edges), decreasing to 67% for large graphs (median of 24 nodes
1069 and 26 edges).

1070 Our comparative analysis has revealed the relative strengths and weaknesses of different approaches.
1071 Traditional neural networks scale better, but are less explainable and data-efficient than their
1072 neurosymbolic counterparts. For LLM-based approaches, the integration with ASP brings added value to
1073 reasoning tasks and outperforms chain-of-thought prompting. Common baselines like MNIST arithmetic
1074 pose little challenge to most frameworks, whereas harder tasks like follow suit lead to time-outs. Apart
1075 from hitting sets, there is a lack of challenging benchmarks where neurosymbolic frameworks achieve
1076 solid, but not near-perfect results.

1077 *ASP generation*

1078 ASP is able to express all NP-search problems and includes a variety of useful constructs such as
1079 disjunctions, choices and negation as failure to efficiently model statements Brewka et al. (2011).
1080 However, many neurosymbolic framework can only learn a subset of ASP. Furthermore, they often
1081 require extensive background knowledge and mode biases to limit the search space. Even LLMs, which
1082 can in theory generate any answer set program, are often limited to just producing predicates and have
1083 only been shown to generate deductive proofs rather than general knowledge. In this section, we will
1084 discuss these limitations, focussing on frameworks that learn at least part of an answer set program.

1085 *Expressivity limits.* Predicates are the most basic form of ASP that a framework can produce and can
1086 be extracted directly from the input using a translation procedure. A few frameworks only generate
1087 predicates: ASP-VQA, AQuA and NSGRAPH utilise neural networks, while [LLM]+ASP, STAR and
1088 LLMASP take advantage of LLMs. The rest of the answer set program is either extracted using a fixed
1089 parsing pipeline or hard-coded in the form of knowledge modules. Some of these knowledge modules
1090 are general enough to be reusable for different tasks. Initially, the DSPy-ASP framework also generates
1091 just facts and relies on hard-coded ASP rules to form a program. However, it can change those rules and
1092 generate new ones in the iterative refinement stage based on Clingo feedback.

1093 Learning rules is a step up from predicates, but not all algorithms can generate rules that exploit the
1094 full expressivity of ASP. For example, NeSyFOLD uses the FOLD-SE algorithm, which can only learn
1095 default theories. These are equivalent to stratified answer set programs and cannot contain any cycles
1096 through negation. NeuralFastLAS is modelled after the first version of FastLAS, which is limited to
1097 observational predicate learning (OPL) and cannot learn recursive rules (Law et al. 2020a). Three further
1098 frameworks use FastLAS directly and are compatible with newer versions, which support non-OPL
1099 learning (Law et al. 2021): FFNSL, NSIL and NeSyGPT. Together with LLM2LAS and Embed2Rule,
1100 they also support ILASP for tasks that require higher expressive power. NSIL uses ILASP's capability
1101 to learn choice rules in the hitting set task. FFNSL, NeSyGPT and Embed2Rule require its predicate
1102 invention capability for the Follow suit task. Figure 7 provides a summary of the levels of expressivity and

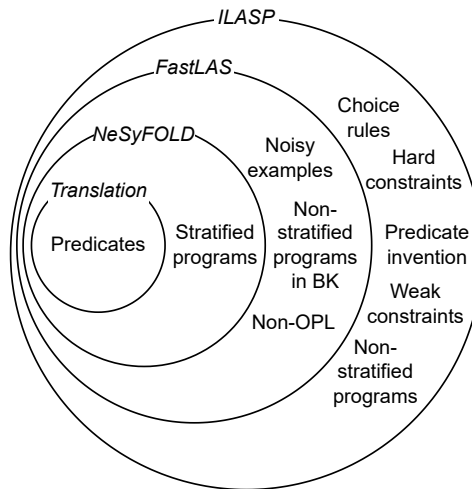


Figure 7. Levels of expressivity of different LAS frameworks.

1103 features that different frameworks can reach. Notably, FastLAS allows the use of non-stratified programs
 1104 in the background knowledge, while ILASP supports it outright. ILASP can also learn higher-level ASP
 1105 constructs such as hard and weak constraints, as well as choice rules.

1106 Lastly, there are five frameworks that use LLMs to produce entire ASP rulesets: LLASP and CLM-ASP
 1107 generate rules directly, with the latter using rejection sampling to pick the best program. GPT-ASP has a
 1108 four step process and LLM-ARC generates both ASP rules and tests, refining them iteratively. NL2ASP
 1109 generates sentences in CNL, which are then deterministically translated into ASP. As LLMs can output
 1110 any combination of letters, they are in theory capable of producing any unrestricted answer set program,
 1111 using all the syntactic constructs available in the language. This remains the case for NL2ASP, as CNL
 1112 supports all standard constructs of ASP.

1113 In summary, the expressive capabilities of generated answer set programs vary between frameworks.
 1114 Seven frameworks only generate predicates, the lowest level of expressivity. NeuralFastLAS can learn
 1115 rules, but is restricted to OPL tasks. Five frameworks use ILASP for generation, which gives them the
 1116 ability to generate complex programs, including constructs such as negation as failure. Another five LLM-
 1117 based frameworks generate rules directly, rather than just extracting predicates, allowing them to create
 1118 rules of any level of expressivity in principle.

1119 *Background knowledge.* Out of the 23 frameworks discussed in this survey, 11 hard-code the ASP
 1120 component. Out of the remaining frameworks, six use FastLAS or ILASP and therefore make extensive
 1121 use of background knowledge. FFNSL, NeSyGPT, Embed2Rule, NSIL and NeuralFastLAS use rule
 1122 templates to restrict the search space and make the task tractable for FastLAS or ILASP.

1123 We illustrate the extent of the background knowledge with the MNIST addition task that has been a
 1124 running example throughout this survey. The following is an example of a typical rule template for this
 1125 task, which we have adapted from NSIL:

```

1126 num(0..18).
1127 digit_type(0..9).
1128 result(Y) :- digit(1,X0), digit(2,X1), solution(X0,X1,Y).
1129 :- digit(1,X0), digit(2,X1), result(Y1), result(Y2), Y1 != Y2.
1130 #modeh(solution(var(digit_type),var(digit_type),var(num))).
1131 #modeb(var(num) = var(digit_type)).
1132 #modeb(var(num) = var(digit_type) + var(digit_type)).
1133 #maxv(3).
1134 #bias("penalty(1, head(X)) :- in_head(X).").
1135 #bias("penalty(1, body(X)) :- in_body(X).").

```

1136 The first two lines restrict the domain of the labels (`num`) and digits (`digit_type`). The next line
1137 specifies that the result is the solution of the two input digits. The constraint ensures that there is only
1138 one result. The rest of the background knowledge consists of mode biases. `#modeh` specifies that the
1139 head of the learned rule must include a solution predicate with two digits and a number as its arguments.
1140 The `#modeb` lines tell the LAS solver that only a digit or the sum of two digits can be in the body of the
1141 learned rule. `#maxv(3)` restricts the LAS solver to a maximum of three variables in each rule. The last
1142 two lines specify that a penalty is given for each example that is not covered. The penalties instruct the
1143 LAS solver to find an optimal solution that covers as much of the data as possible.

1144 As this example demonstrates, much of the ASP structure is still hand-crafted, even for a simple task
1145 like MNIST addition. NSIL also runs experiments with superfluous functions, such as subtraction or
1146 multiplication, to increase the hypothesis space. However, this still presents a vastly restricted search
1147 space compared to the universe of all possible functions. For more complex tasks, such as playing card
1148 games, even more information is needed to make the learning task tractable. Scaling up to real-world
1149 tasks would therefore require a substantial amount of manual engineering. This limits the usefulness of
1150 LAS solvers in the real world, because background knowledge is expensive to codify and sometimes
1151 unattainable. LLM2LAS is the only approach that automates some of the mode bias generation. An LLM
1152 generates the fluents in a sentence, which are then transformed into mode body and head atoms using a
1153 hard-coded pipeline. The use of LLMs to generate the ILASP search space instead of rules holds promise.
1154 It saves a lot of manual work and still results in a structured search for rules that provably fit the data.

1155 Six frameworks generate 100% of their ASP rules themselves: NeSyFOLD, LLASP, NL2ASP, GPT-
1156 ASP, LLM-ARC and CLM-ASP. NeSyFOLD restricts the search space by only considering stratified
1157 ASP rules. As NL2ASP employs a traditional language model, it trains on the dataset using k-fold cross
1158 validation. All other papers use LLMs and guide the search through in-context learning. In these cases,
1159 some background knowledge is implicitly included in the prompt through solved examples.

1160 *Induction vs deduction.* Learning answer set programs has traditionally been conducted through the
1161 lens of inductive logic programming (ILP). The goal of ILP is to learn general facts and rules from
1162 examples (Muggleton 1991). All neurosymbolic ASP frameworks in this survey that use traditional neural
1163 networks perform inductive learning. They generate one answer set program using multiple examples that
1164 models the entire dataset. The frameworks using LLMs, however, work differently. They generate a new
1165 answer set program for every example, with the aim of solving the puzzle in the question. The program
1166 models the natural language question in ASP and the solver then performs deductive inference to arrive

1167 at the solution. This is more akin to a translation and deduction task, rather than discovering new, general
 1168 knowledge through induction.

1169 LLMs in general struggle much more with inductive reasoning than deductive reasoning (Hua
 1170 et al. 2025). Therefore, it remains to be seen if frameworks using LLMs are able to generate answer
 1171 set programs inductively. More research and experiments are needed to develop this capability in
 1172 neurosymbolic ASP.

1173 Scalability

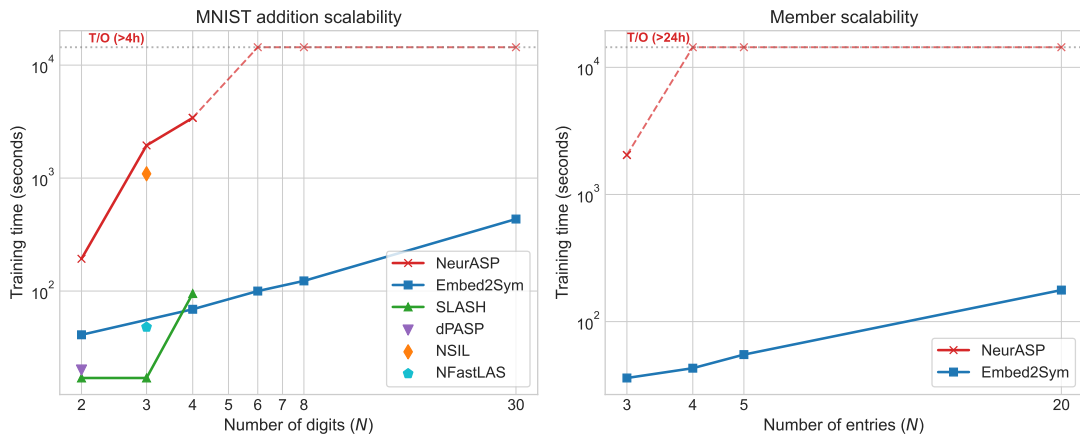


Figure 8. Scalability comparison of different frameworks on MNIST addition and member tasks. NFastLAS = NeuralFastLAS.

1174 Figure 8 compares the training times of different approaches on the MNIST addition and member
 1175 tasks with increasing complexity. As the input sizes increase, NeurASP eventually times out on both
 1176 tasks. Embed2Sym scales much better because it utilizes a fully neural training procedure. The gap to
 1177 the symbolic space is crossed using clustering, which is faster than computing all answer sets since there
 1178 are only 10 clusters for MNIST inputs. A similar picture is painted by SLASH. It scales much better
 1179 than NeurASP because it employs the SAME method, which prunes unlikely answer sets as the network
 1180 grows more confident during training. But the paper only provides numbers for up to 4 digits. dPASP
 1181 runs similarly quickly for 2 digits, but no more data points are available. NSIL calls NeurASP as part of
 1182 its training loop, albeit with an improved implementation. This allows it to be quicker on MNIST add 3,
 1183 even though it has to learn the rules as well. NeuralFastLAS vastly improves on NSIL by eliminating the
 1184 need to use NeurASP and calculating a solution in a single iteration. In addition to its faster runtime in
 1185 MNIST add 3, it solves the Even9Plus task in 18 seconds compared to NSIL's 14 minutes.

1186 When moving to more complex tasks, frameworks like SLASH and NSIL hit their limits, as depicted
 1187 in Figure 9. Apart from implementation improvements, the training procedures of NeurASP, SLASH,
 1188 dPASP and NSIL all boil down to the same computations: For each example, they enumerate all valid
 1189 answer sets, calculate their probabilities and compute the gradients. This is feasible for tasks like MNIST



Figure 9. Scalability comparison of different frameworks on hitting set, sudoku and follow suit tasks. T/O = timeout after 24 hours.

1190 addition or 4x4 sudoku with a few hundred possible answer sets. However, on more complex tasks
 1191 like hitting set 10/6 and follow suit, they time out. Follow suit 4 in particular yields anywhere from
 1192 800,000 to 5 million answer sets. Embed2Rule makes use of the same fully neural training procedure
 1193 as Embed2Sym, allowing it to learn tasks where SLASH and NSIL fail. However, training times remain
 1194 high for tasks like Sudoku 9x9, since it learns rules using ILASP.

1195 Searching for solutions in the space of answer set programs is a difficult task to scale. For example,
 1196 the complexity for ILASP to decide whether a hypothesis is an optimal inductive solution is Σ_2^P -
 1197 complete (Law et al. 2018). This inherently limits the scalability of frameworks like Embed2Rule or
 1198 NeSyGPT. FastLAS was created to be a more scalable LAS solver, but at a cost of features such as
 1199 learning recursive rules or predicate invention (Law et al. 2021).

1200 Such scalability issues can be bypassed with the use of foundation models, which do not calculate
 1201 rules exactly, but rather predict the right words and symbols. ASP-VQA answers questions much faster
 1202 than NeurASP on CLEVR through the use of an LLM, reporting 89 seconds of runtime compared to
 1203 NeurASP’s 1 hour. However, LLMs require a lot of resources to be trained in the first place, far surpassing
 1204 the cost of training frameworks like NeurASP. Moreover, LLMs such as OpenAI’s GPT models are
 1205 proprietary and require dedicated server architecture to run, usually incurring a per-token cost. Unlike
 1206 traditional methods, they cannot be run on local machines. Even smaller models like BLIP, which is used
 1207 in NeSyGPT and Embed2Rule, require a lot of resources to finetune. Training them from scratch would
 1208 be prohibitively expensive for researchers, which is why pre-trained models are used instead.

1209 There is still a lot of work to be done to speed up ASP methods. While neural networks benefit from
 1210 GPU parallelisation and very efficient Python frameworks, most ASP methods like Clingo, ILASP or
 1211 NeurASP run on the CPU. There is progress on parallelisation, as with answer set networks, which can
 1212 compute answer sets on the GPU (Skryagin et al. 2024a). More such methods are needed to cover all steps

1213 of neurosymbolic pipelines. Only with faster implementations, in addition to theoretical discoveries, can
1214 neurosymbolic ASP methods become more viable in real-world tasks.

1215 Conclusion

1216 In this survey, we have discussed the current literature on neurosymbolic ASP, highlighting the
1217 achievements and limitations of the field. We categorised the wide array of different frameworks
1218 according to which components are learned or hard-coded. Frameworks with fully pre-trained neural
1219 networks and hard-coded ASP components focus on the translation between the components. The
1220 addition ASP transforms decisions into transparent and verifiable processes and enables further features
1221 like contrastive explanations. When the neural components need to be trained, the main challenge lies in
1222 propagating the learning signal through the non-differentiable ASP component. Papers in this category
1223 either enumerate answer sets to serve as noisy labels or use clustering. There is room for further research
1224 into different methods of providing learning signals that scale better to the number of answer sets or
1225 clusters. The predictions of pre-trained neural networks can serve as noisy examples for LAS solvers,
1226 enabling rule learning from noisy data. Recent work has started utilizing LLMs to predict rules directly.
1227 Their results show that integrating ASP improves the reasoning capabilities of LLMs compared to
1228 traditional prompting. Lastly, jointly learning the neural and symbolic component serves as the most
1229 difficult challenge, as there are no reliable learning signals for either component at the start. Only two
1230 frameworks have been proposed so far and have made the first steps with simple problems.

1231 The biggest challenge for neurosymbolic methods is scalability. The current limit for joint learning in
1232 perception tasks is CIFAR-10 and approaches using traditional neural networks time out when the task
1233 size increases. One reason lies in suboptimal implementations. Commonly-used frameworks like Clingo
1234 or ILASP do not have GPU support and cannot be parallelized easily. Code built on top of them, such as
1235 NeurASP, is often a research prototype and therefore not optimised. More efficient implementations are
1236 necessary for scaling up to larger tasks, but that is not enough. Novel methods are needed for problems
1237 like efficiently traversing the search space of ASP rules and propagating learning signals through ASP
1238 components. A tighter integration between neural and symbolic representations could overcome some
1239 inefficiencies.

1240 The need for extensive hard-coding also holds the field back. Hand-writing rules can be a very
1241 expensive processes and leads to poor generalizability. Even frameworks that learn rules needs to limit
1242 the search space with manually engineered rule templates. Recent work has explored the use of LLMs
1243 to automate rule generation, which represents a promising direction. LLMs can generate ASP rules
1244 directly and reduce the need for hard-coding, while ASP brings added value to their predictions in
1245 terms of reasoning ability and explainability. Further research is necessary to learn inductive knowledge
1246 with LLMs, for example by integrating them with LAS solvers. Another untapped research area is the
1247 combination ASP and foundation models for multi-modal inference, such as voice and text.

1248 An opportunity for the field is development of more realistic datasets. The most popular benchmarks
1249 for traditional neurosymbolic methods are tasks based on MNIST images and frameworks routinely
1250 report accuracies of 80 or 90%. Most datasets for LLM-based methods use synthetically generated
1251 sentences and tasks with clear-cut symbolic representations. More realistic tasks are needed to
1252 demonstrate the applicability of neurosymbolic ASP in real-world settings. There is a need to create
1253 challenging but solvable tasks that can be widely adopted to measure different approaches against.

1254 As neural models continue to struggle with complex logical reasoning, integrating efficient and
1255 search-based symbolic methods can achieve better performance, robustness and explainability. With
1256 its combination of expressiveness and readability, ASP is well-placed to fulfill this role. We have
1257 the opportunity to create trustworthy neurosymbolic frameworks that go beyond summarising existing
1258 information and towards discovering new knowledge.

1259 Acknowledgements

1260 This work was supported by the UKRI Centre for Doctoral Training in Safe and Trusted Artificial Intelligence
1261 [EP/S0233356/1].

1262 References

- 1263 Acharya K and Song H (2025) A comprehensive review of neuro-symbolic ai for robustness, uncertainty
1264 quantification, and intervenability. *Arabian Journal for Science and Engineering* DOI:10.1007/
1265 s13369-025-10887-3.
- 1266 Agostinelli F, Panta R and Khandelwal V (2024) Specifying goals to deep neural networks with answer set
1267 programming. *Proceedings of the International Conference on Automated Planning and Scheduling* 34(1): 2–
1268 10. DOI:10.1609/icaps.v34i1.31454. URL [https://ojs.aaai.org/index.php/ICAPS/article/
1269 view/31454](https://ojs.aaai.org/index.php/ICAPS/article/view/31454).
- 1270 Ahmed K, Teso S, Chang KW, den Broeck GV and Vergari A (2022) Semantic probabilistic layers for neuro-
1271 symbolic learning. In: *Neural Information Processing Systems*. pp. 29944–29959. URL [https://api.
1272 semanticscholar.org/CorpusID:249240063](https://api.semanticscholar.org/CorpusID:249240063).
- 1273 Albilani M and Bouzeghoub A (2023) Guided hierarchical reinforcement learning for safe urban driving. In:
1274 *2023 IEEE 35th International Conference on Tools with Artificial Intelligence (ICTAI)*. pp. 746–753. DOI:
1275 10.1109/ICTAI59109.2023.00115.
- 1276 Alviano M, Calimeri F, Dodaro C, Fuscà D, Leone N, Perri S, Ricca F, Veltri P and Zangari J (2017) The asp
1277 system dlv2. In: Balduccini M and Janhunen T (eds.) *Logic Programming and Nonmonotonic Reasoning*. Cham:
1278 Springer International Publishing, pp. 215–221.
- 1279 Alviano M, Scudo FL, Grillo L and Reiners LAR (2024) Answer set programming and large language models
1280 interaction with yam! : Second report. In: *KoDis+CAKR+SYNERGY@KR*. pp. 4330–4338. URL [https:
1281 //api.semanticscholar.org/CorpusID:274981259](https://api.semanticscholar.org/CorpusID:274981259).
- 1282 Arias J, Carro M, Salazar E, Marple K and Gupta G (2018) Constraint answer set programming without grounding.
1283 *Theory and Practice of Logic Programming* 18(3–4): 337–354. DOI:10.1017/S1471068418000285.
- 1284 Aspis Y, Albinhassan M, Lobo J and Russo A (2024) Embed2rule scalable neuro-symbolic learning via latent space
1285 weak-labelling. In: *Neural-Symbolic Learning and Reasoning: 18th International Conference, NeSy 2024,
1286 Barcelona, Spain, September 9–12, 2024, Proceedings, Part I*. Berlin, Heidelberg: Springer-Verlag. ISBN 978-
1287 3-031-71166-4, p. 195–218. DOI:10.1007/978-3-031-71167-1_11. URL [https://doi.org/10.1007/
1288 978-3-031-71167-1_11](https://doi.org/10.1007/978-3-031-71167-1_11).
- 1289 Aspis Y, Broda K, Lobo J and Russo A (2022) Embed2Sym - Scalable Neuro-Symbolic Reasoning via Clustered
1290 Embeddings. In: *Proceedings of the 19th International Conference on Principles of Knowledge Representation
1291 and Reasoning*. pp. 421–431. DOI:10.24963/kr.2022/44. URL [https://doi.org/10.24963/kr.
1292 2022/44](https://doi.org/10.24963/kr.2022/44).

- 1293 Badreddine S, d'Avila Garcez A, Serafini L and Spranger M (2022) Logic tensor networks. *Artificial Intelligence*
1294 303: 103649. DOI:<https://doi.org/10.1016/j.artint.2021.103649>. URL <https://www.sciencedirect.com/science/article/pii/S0004370221002009>.
- 1295
- 1296 Barbara V, Guarascio M, Leone N, Manco G, Quarta A, Ricca F and Ritacco E (2023) Neuro-symbolic ai for
1297 compliance checking of electrical control panels. URL <https://arxiv.org/abs/2305.10113>.
- 1298 Basu K, Shakerin F and Gupta G (2020) Aqua: Asp-based visual question answering. In: *Practical Aspects*
1299 *of Declarative Languages: 22nd International Symposium, PADL 2020, New Orleans, LA, USA, January*
1300 *20–21, 2020, Proceedings*. Berlin, Heidelberg: Springer-Verlag. ISBN 978-3-030-39196-6, p. 57–72. DOI:
1301 10.1007/978-3-030-39197-3_4. URL https://doi.org/10.1007/978-3-030-39197-3_4.
- 1302 Bauer JJ, Eiter T, Higuera Ruiz N and Oetsch J (2025) Visual graph question answering with asp and llms for
1303 language parsing. *Electronic Proceedings in Theoretical Computer Science* 416: 15–28. DOI:10.4204/eptcs.
1304 416.2. URL <http://dx.doi.org/10.4204/EPTCS.416.2>.
- 1305 Baugh KG, Cingillioglu N and Russo A (2023) Neuro-symbolic rule learning in real-world classification tasks. URL
1306 <https://arxiv.org/abs/2303.16674>.
- 1307 Belle V and Marcus G (2026) The future is neuro-symbolic: Where has it been, and where is it going? In: *Proceedings*
1308 *of the 40th AAAI Conference on Artificial Intelligence*, Proceedings of the AAAI Conference on Artificial
1309 Intelligence. AAAI Press, pp. 1–8. URL <https://aaai.org/conference/aaai/aaai-26/>. The
1310 40th Annual AAAI Conference on Artificial Intelligence, AAAI-26 ; Conference date: 20-01-2026 Through
1311 27-01-2026.
- 1312 Bhuyan BP, Ramdane-Cherif A, Tomar R and Singh TP (2024) Neuro-symbolic artificial intelligence: a survey.
1313 *Neural Comput. Appl.* 36(21): 12809–12844. DOI:10.1007/s00521-024-09960-z. URL [https://doi.org/](https://doi.org/10.1007/s00521-024-09960-z)
1314 [10.1007/s00521-024-09960-z](https://doi.org/10.1007/s00521-024-09960-z).
- 1315 Bommasani R, Hudson DA, Adeli E, Altman R, Arora S, von Arx S, Bernstein MS, Bohg J, Bosselut A, Brunskill E,
1316 Brynjolfsson E, Buch S, Card D, Castellon R, Chatterji N, Chen A, Creel K, Davis JQ, Demszky D, Donahue C,
1317 Doumbouya M, Durmus E, Ermon S, Etchemendy J, Ethayarajh K, Fei-Fei L, Finn C, Gale T, Gillespie L, Goel
1318 K, Goodman N, Grossman S, Guha N, Hashimoto T, Henderson P, Hewitt J, Ho DE, Hong J, Hsu K, Huang J,
1319 Icard T, Jain S, Jurafsky D, Kalluri P, Karamcheti S, Keeling G, Khani F, Khattab O, Koh PW, Krass M, Krishna
1320 R, Kuditipudi R, Kumar A, Ladhak F, Lee M, Lee T, Leskovec J, Levent I, Li XL, Li X, Ma T, Malik A, Manning
1321 CD, Mirchandani S, Mitchell E, Munyikwa Z, Nair S, Narayan A, Narayanan D, Newman B, Nie A, Niebles JC,
1322 Nilforoshan H, Nyarko J, Ogut G, Orr L, Papadimitriou I, Park JS, Piech C, Portelance E, Potts C, Raghunathan
1323 A, Reich R, Ren H, Rong F, Roohani Y, Ruiz C, Ryan J, Ré C, Sadigh D, Sagawa S, Santhanam K, Shih A,
1324 Srinivasan K, Tamkin A, Taori R, Thomas AW, Tramèr F, Wang RE, Wang W, Wu B, Wu J, Wu Y, Xie SM,
1325 Yasunaga M, You J, Zaharia M, Zhang M, Zhang T, Zhang X, Zhang Y, Zheng L, Zhou K and Liang P (2022)
1326 On the opportunities and risks of foundation models. URL <https://arxiv.org/abs/2108.07258>.
- 1327 Bordes F, Pang RY, Ajay A, Li AC, Bardes A, Petryk S, Mañas O, Lin Z, Mahmoud A, Jayaraman B, Ibrahim M, Hall
1328 M, Xiong Y, Lebensold J, Ross C, Jayakumar S, Guo C, Bouchacourt D, Al-Tahan H, Padthe K, Sharma V, Xu
1329 H, Tan XE, Richards M, Lavoie S, Astolfi P, Hemmat RA, Chen J, Tirumala K, Assouel R, Moayeri M, Talattof
1330 A, Chaudhuri K, Liu Z, Chen X, Garrido Q, Ullrich K, Agrawal A, Saenko K, Celikyilmaz A and Chandra V
1331 (2024) An introduction to vision-language modeling. URL <https://arxiv.org/abs/2405.17247>.
- 1332 Borroto M, Ielo A, Mazzotta G and Ricca F (2025) Answer set programming and neurosymbolic ai: Applications
1333 and future perspectives (invited talk). In: Bruno P, Calimeri F, Cauteruccio F and Terracina G (eds.) *Hybrid*

- 1334 *Models for Coupling Deductive and Inductive Reasoning*. Cham: Springer Nature Switzerland. ISBN 978-3-
1335 031-89366-7, pp. 1–21.
- 1336 Borroto Santana MA, GALLAGHER K, IELO A, KAREEM I, RICCA F and RUSSO A (2025) Question
1337 answering with llms and learning from answer sets. *Theory and Practice of Logic Programming* : 1–25DOI:
1338 10.1017/S1471068425100343.
- 1339 Bouneffouf D and Aggarwal CC (2022) Survey on applications of neurosymbolic artificial intelligence. URL
1340 <https://arxiv.org/abs/2209.12618>.
- 1341 Brewka G, Eiter T and Truszczynski M (2011) Answer set programming at a glance. *Commun. ACM* 54: 92–103.
1342 DOI:10.1145/2043174.2043195.
- 1343 Brown T, Mann B, Ryder N, Subbiah M, Kaplan JD, Dhariwal P, Neelakantan A, Shyam P, Sastry G, Askell A,
1344 Agarwal S, Herbert-Voss A, Krueger G, Henighan T, Child R, Ramesh A, Ziegler D, Wu J, Winter C, Hesse
1345 C, Chen M, Sigler E, Litwin M, Gray S, Chess B, Clark J, Berner C, McCandlish S, Radford A, Sutskever I
1346 and Amodei D (2020) Language models are few-shot learners. In: Larochelle H, Ranzato M, Hadsell R, Balcan
1347 M and Lin H (eds.) *Advances in Neural Information Processing Systems*, volume 33. Curran Associates, Inc.,
1348 pp. 1877–1901. URL [https://proceedings.neurips.cc/paper_files/paper/2020/file/
1349 1457c0d6bfbcb4967418bfb8ac142f64a-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfbcb4967418bfb8ac142f64a-Paper.pdf).
- 1350 Cabalar P, Fandinno J and Muñiz B (2020) A system for explainable answer set programming. *Electronic*
1351 *Proceedings in Theoretical Computer Science* 325: 124–136. DOI:10.4204/EPTCS.325.19.
- 1352 Caruso S, Dodaro C, Maratea M, Mochi M and RICCIO F (2023) Cnl2asp: Converting controlled natural language
1353 sentences into asp. *Theory and Practice of Logic Programming* 24. DOI:10.1017/S1471068423000388.
- 1354 Charalambous T, Aspis Y and Russo A (2023) Neurafastlas: Fast logic-based learning from raw data. URL
1355 <https://arxiv.org/abs/2310.05145>.
- 1356 Choi Y, Vergari A and den Broeck GV (2020) Probabilistic circuits: A unifying framework for tractable probabilistic
1357 models. URL <https://api.semanticscholar.org/CorpusID:221997880>.
- 1358 Chu-Carroll J, Beck A, Burnham G, Melville DO, Nachman D, Özcan AE and Ferrucci D (2024) Beyond llms:
1359 Advancing the landscape of complex reasoning. URL <https://arxiv.org/abs/2402.08064>.
- 1360 Cingillioglu N and Russo A (2021) pix2rule: End-to-end neuro-symbolic rule learning. *International Workshop on*
1361 *Neural-Symbolic Learning and Reasoning* URL [https://api.semanticscholar.org/CorpusID:
1362 235422026](https://api.semanticscholar.org/CorpusID:235422026).
- 1363 Ciravegna G, Barbiero P, Giannini F, Gori M, Liò P, Maggini M and Melacci S (2023) Logic explained networks.
1364 *Artificial Intelligence* 314: 103822. DOI:<https://doi.org/10.1016/j.artint.2022.103822>. URL [https://www.
1365 sciencedirect.com/science/article/pii/S000437022200162X](https://www.sciencedirect.com/science/article/pii/S000437022200162X).
- 1366 Clocksin WF and Mellish C (1981) *Programming in Prolog*. Springer. URL [https://api.
1367 semanticscholar.org/CorpusID:10092527](https://api.semanticscholar.org/CorpusID:10092527).
- 1368 Colelough BC and Regli W (2025) Neuro-symbolic ai in 2024: A systematic review. *ArXiv abs/2501.05435*. URL
1369 <https://api.semanticscholar.org/CorpusID:274180938>.
- 1370 Coppolillo E, Calimeri F, Manco G, Perri S and Ricca F (2024) Llaspl: fine-tuning large language models for
1371 answer set programming. In: *Proceedings of the 21st International Conference on Principles of Knowledge*
1372 *Representation and Reasoning*, KR '24. ISBN 978-1-956792-05-8, pp. 834–844. DOI:10.24963/kr.2024/78.
1373 URL <https://doi.org/10.24963/kr.2024/78>.
- 1374 Cunnington D, Law M, Lobo J and Russo A (2023a) Ffnsl: Feed-forward neural-symbolic learner. *Mach.*
1375 *Learn.* 112(2): 515–569. DOI:10.1007/s10994-022-06278-6. URL <https://doi.org/10.1007/>

- 1376 [s10994-022-06278-6](#).
- 1377 Cunnington D, Law M, Lobo J and Russo A (2023b) Neuro-symbolic learning of answer set programs from raw
1378 data.
- 1379 Cunnington D, Law M, Lobo J and Russo A (2024) The role of foundation models in neuro-symbolic learning
1380 and reasoning. In: Besold TR, d'Avila Garcez A, Jimenez-Ruiz E, Confalonieri R, Madhyastha P and Wagner B
1381 (eds.) *Neural-Symbolic Learning and Reasoning*. Cham: Springer Nature Switzerland. ISBN 978-3-031-71167-
1382 1, pp. 84–100.
- 1383 Dai WZ and Muggleton S (2021) Abductive knowledge induction from raw data. In: Zhou ZH (ed.) *Proceedings of*
1384 *the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*. International Joint Conferences
1385 on Artificial Intelligence Organization, pp. 1845–1851. DOI:10.24963/ijcai.2021/254. URL <https://doi.org/10.24963/ijcai.2021/254>. Main Track.
- 1386
- 1387 Dai WZ, Xu Q, Yu Y and Zhou ZH (2019) Bridging machine learning and logical reasoning by abductive learning.
1388 In: *Proceedings of the 33rd International Conference on Neural Information Processing Systems*. Red Hook,
1389 NY, USA: Curran Associates Inc., pp. 2815–2826.
- 1390 De Raedt L, Kimmig A and Toivonen H (2007) Problog: a probabilistic prolog and its application in link discovery.
1391 In: *Proceedings of the 20th International Joint Conference on Artificial Intelligence, IJCAI'07*. San Francisco,
1392 CA, USA: Morgan Kaufmann Publishers Inc., p. 2468–2473.
- 1393 De Smet L, Zuidberg Dos Martires P, Manhaeve R, Marra G, Kimmig A and De Raedt L (2023) Neural probabilistic
1394 logic programming in discrete-continuous domains. In: Evans RJ and Shpitser I (eds.) *Proceedings of the Thirty-*
1395 *Ninth Conference on Uncertainty in Artificial Intelligence, Proceedings of Machine Learning Research*, volume
1396 216. PMLR, pp. 529–538. URL <https://proceedings.mlr.press/v216/de-smet23a.html>.
- 1397 Deng L (2012) The mnist database of handwritten digit images for machine learning research. *IEEE Signal*
1398 *Processing Magazine* 29(6): 141–142.
- 1399 Dong H, Mao J, Lin T, Wang C, Li L and Zhou D (2019) Neural logic machines. In: *International*
1400 *Conference on Learning Representations*. pp. 2574–2595. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=BlxY-hRctX)
1401 [BlxY-hRctX](https://openreview.net/forum?id=BlxY-hRctX).
- 1402 Eiter T, Geibinger T, Higuera N and Oetsch J (2023) A logic-based approach to contrastive explainability for
1403 neurosymbolic visual question answering. In: Elkind E (ed.) *Proceedings of the Thirty-Second International*
1404 *Joint Conference on Artificial Intelligence, IJCAI-23*. International Joint Conferences on Artificial Intelligence
1405 Organization, pp. 3668–3676. DOI:10.24963/ijcai.2023/408. URL [https://doi.org/10.24963/](https://doi.org/10.24963/ijcai.2023/408)
1406 [ijcai.2023/408](https://doi.org/10.24963/ijcai.2023/408). Main Track.
- 1407 Eiter T, Higuera N, Oetsch J and Pritz M (2022) A neuro-symbolic asp pipeline for visual question answering. URL
1408 <https://arxiv.org/abs/2205.07548>.
- 1409 Evans R, Bošnjak M, Buesing L, Ellis K, Pfau D, Kohli P and Sergot M (2021) Making sense of raw input.
1410 *Artificial Intelligence* 299: 103521. DOI:<https://doi.org/10.1016/j.artint.2021.103521>. URL [https://www.](https://www.sciencedirect.com/science/article/pii/S0004370221000722)
1411 [sciencedirect.com/science/article/pii/S0004370221000722](https://www.sciencedirect.com/science/article/pii/S0004370221000722).
- 1412 Farquhar S, Kossen J, Kuhn L and Gal Y (2024) Detecting hallucinations in large language models using
1413 semantic entropy. *Nature* 630: 625–630. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:270615909)
1414 [270615909](https://api.semanticscholar.org/CorpusID:270615909).
- 1415 Furelos-Blanco D, Law M, Jonsson A, Broda K and Russo A (2021) Induction and exploitation of subgoal
1416 automata for reinforcement learning. *J. Artif. Int. Res.* 70: 1031–1116. DOI:10.1613/jair.1.12372. URL
1417 <https://doi.org/10.1613/jair.1.12372>.

- 1418 Garcez A and Lamb L (2023) Neurosymbolic ai: the 3rd wave. *Artificial Intelligence Review* : 1–20DOI:
1419 10.1007/s10462-023-10448-w.
- 1420 Gebser M, Kaminski R, Kaufmann B and Schaub T (2017) Multi-shot ASP solving with clingo. *CoRR*
1421 abs/1705.09811.
- 1422 Geh RL, Gonçalves J, Silveira IC, Mauá DD and Cozman FG (2024) dPASP: A Probabilistic Logic Programming
1423 Environment For Neurosymbolic Learning and Reasoning. In: *Proceedings of the 21st International Conference*
1424 *on Principles of Knowledge Representation and Reasoning*. pp. 731–742. DOI:10.24963/kr.2024/69. URL
1425 <https://doi.org/10.24963/kr.2024/69>.
- 1426 Gelfond M and Kahl Y (2014) *Knowledge Representation, Reasoning, and the Design of Intelligent Agents: The*
1427 *Answer-Set Programming Approach*. Cambridge University Press.
- 1428 Gibaut W, Pereira L, Grassioto F, Osorio A, Gadioli E, Munoz A, Gomes S and Santos Cd (2023) Neurosymbolic ai
1429 and its taxonomy: a survey. *ArXiv* DOI:10.48550/ARXIV.2305.08876. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2305.08876)
1430 [2305.08876](https://arxiv.org/abs/2305.08876).
- 1431 Han S, Schoelkopf H, Zhao Y, Qi Z, Riddell M, Zhou W, Coady J, Peng D, Qiao Y, Benson L, Sun L, Wardle-Solano
1432 A, Szabó H, Zubova E, Burtell M, Fan J, Liu Y, Wong B, Sailor M, Ni A, Nan L, Kasai J, Yu T, Zhang R, Fabbri
1433 A, Kryscinski WM, Yavuz S, Liu Y, Lin XV, Joty S, Zhou Y, Xiong C, Ying R, Cohan A and Radev D (2024)
1434 FOLIO: Natural language reasoning with first-order logic. In: AI-Onaizan Y, Bansal M and Chen YN (eds.)
1435 *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. Miami, Florida,
1436 USA: Association for Computational Linguistics, pp. 22017–22031. DOI:10.18653/v1/2024.emnlp-main.1229.
1437 URL <https://aclanthology.org/2024.emnlp-main.1229/>.
- 1438 Han Z, Cai LW, Huang YX, Wei B, Wang W and Yin Y (2023) Abductive subconcept learning. *Science China*
1439 *Information Sciences* 66. DOI:10.1007/s11432-020-3569-0.
- 1440 Hitzler P, Eberhart A, Ebrahimi M, Sarker MK and Zhou L (2022) Neuro-symbolic approaches in artificial
1441 intelligence. *National Science Review* 9(6). DOI:10.1093/nsr/nwac035. URL [https://doi.org/10.](https://doi.org/10.1093/nsr/nwac035)
1442 [1093/nsr/nwac035](https://doi.org/10.1093/nsr/nwac035).
- 1443 Hitzler P and Sarker MK (2021) *Neuro-Symbolic Artificial Intelligence: The State of the Art*. IOS Press. URL
1444 <https://api.semanticscholar.org/CorpusID:274983612>.
- 1445 Hua W, Sun F, Pan L, Jardine A and Wang WY (2025) Inductionbench: LLMs fail in the simplest complexity class.
1446 *Workshop on Reasoning and Planning for Large Language Models* URL [https://openreview.net/](https://openreview.net/forum?id=brw11PScQM)
1447 [forum?id=brw11PScQM](https://openreview.net/forum?id=brw11PScQM).
- 1448 Huang J, Li Z, Chen B, Samel K, Naik M, Song L and Si X (2021) Scallop: From probabilistic deductive
1449 databases to scalable differentiable reasoning. In: Ranzato M, Beygelzimer A, Dauphin Y, Liang P and
1450 Vaughan JW (eds.) *Advances in Neural Information Processing Systems*, volume 34. Curran Associates, Inc., pp.
1451 25134–25145. URL [https://proceedings.neurips.cc/paper_files/paper/2021/file/](https://proceedings.neurips.cc/paper_files/paper/2021/file/d367eef13f90793bd8121e2f675f0dc2-Paper.pdf)
1452 [d367eef13f90793bd8121e2f675f0dc2-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/d367eef13f90793bd8121e2f675f0dc2-Paper.pdf).
- 1453 Hughes DP and Salathe M (2016) An open access repository of images on plant health to enable the development of
1454 mobile disease diagnostics. URL <https://arxiv.org/abs/1511.08060>.
- 1455 Ishay A and Lee J (2025) Llm+al: bridging large language models and action languages for complex reasoning about
1456 actions. In: *Proceedings of the Thirty-Ninth AAAI Conference on Artificial Intelligence and Thirty-Seventh*
1457 *Conference on Innovative Applications of Artificial Intelligence and Fifteenth Symposium on Educational*
1458 *Advances in Artificial Intelligence*, AAAI'25/IAAI'25/EAAI'25. AAAI Press. ISBN 978-1-57735-897-8, pp.
1459 24212–24220. DOI:10.1609/aaai.v39i23.34597. URL [https://doi.org/10.1609/aaai.v39i23.](https://doi.org/10.1609/aaai.v39i23.34597)

- 1460 34597.
- 1461 Ishay A, Yang Z and Lee J (2023) Leveraging large language models to generate answer set programs. In:
1462 *Proceedings of the 20th International Conference on Principles of Knowledge Representation and Reasoning*,
1463 KR '23. ISBN 978-1-956792-02-7, pp. 374–383. DOI:10.24963/kr.2023/37. URL [https://doi.org/10.](https://doi.org/10.24963/kr.2023/37)
1464 [24963/kr.2023/37](https://doi.org/10.24963/kr.2023/37).
- 1465 Jin Y, Liu J, Luo Z, Peng Y, Qin Z, Dai WZ and Ding YX (2025) Pre-training meta-rule selection policy for visual
1466 generative abductive learning. *International Joint Conference on Learning and Reasoning 2024* : 163–180DOI:
1467 10.1007/978-3-032-09087-4_11.
- 1468 John-Mathews JM (2021) Critical empirical study on black-box explanations in ai. URL [https://arxiv.org/](https://arxiv.org/abs/2109.15067)
1469 [abs/2109.15067](https://arxiv.org/abs/2109.15067).
- 1470 Johnson J, Hariharan B, van der Maaten L, Fei-Fei L, Zitnick CL and Girshick R (2017) Clevr: A diagnostic dataset
1471 for compositional language and elementary visual reasoning. In: *2017 IEEE Conference on Computer Vision*
1472 *and Pattern Recognition (CVPR)*. pp. 1988–1997. DOI:10.1109/CVPR.2017.215.
- 1473 Kairgeldin R and Carreira-Perpiñán MA (2025) Neurosymbolic models based on hybrids of convolutional neural
1474 networks and decision trees. In: H Gilpin L, Giunchiglia E, Hitzler P and van Krieken E (eds.) *Proceedings*
1475 *of The 19th International Conference on Neurosymbolic Learning and Reasoning, Proceedings of Machine*
1476 *Learning Research*, volume 284. PMLR, pp. 796–813. URL [https://proceedings.mlr.press/](https://proceedings.mlr.press/v284/kairgeldin25a.html)
1477 [v284/kairgeldin25a.html](https://proceedings.mlr.press/v284/kairgeldin25a.html).
- 1478 Kalyanpur A, Saravanakumar KK, Barres V, Chu-Carroll J, Melville D and Ferrucci D (2024) Llm-arc: Enhancing
1479 llms with an automated reasoning critic. URL <https://arxiv.org/abs/2406.17663>.
- 1480 Kareem I, Gallagher K, Borroto M, Ricca F and Russo A (2025) Using learning from answer sets for robust question
1481 answering with llm. In: Dodaro C, Gupta G and Martinez MV (eds.) *Logic Programming and Nonmonotonic*
1482 *Reasoning*. Cham: Springer Nature Switzerland. ISBN 978-3-031-74209-5, pp. 112–125.
- 1483 Kaur N, McPheat L, Russo A, Cohn AG and Madhyastha P (2025) An empirical study of conformal prediction in
1484 llm with asp scaffolds for robust reasoning. URL <https://arxiv.org/abs/2503.05439>.
- 1485 Khattab O, Singhvi A, Maheshwari P, Zhang Z, Santhanam K, A SV, Haq S, Sharma A, Joshi TT, Moazam H,
1486 Miller H, Zaharia M and Potts C (2024) DSPy: Compiling declarative language model calls into state-of-the-art
1487 pipelines. *The Twelfth International Conference on Learning Representations* URL [https://openreview.](https://openreview.net/forum?id=sY5N0zyY5Od)
1488 [net/forum?id=sY5N0zyY5Od](https://openreview.net/forum?id=sY5N0zyY5Od).
- 1489 Krizhevsky A and Hinton G (2009) Learning multiple layers of features from tiny images. Technical
1490 Report 0, University of Toronto, Toronto, Ontario. URL [https://www.cs.toronto.edu/~kriz/](https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf)
1491 [learning-features-2009-TR.pdf](https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf).
- 1492 Kuhnle A and Copestake A (2017) Shapeworld - a new test methodology for multimodal language understanding.
1493 URL <https://arxiv.org/abs/1704.04517>.
- 1494 Lamb LC, d'Avila Garcez A, Gori M, Prates MO, Avelar PH and Vardi MY (2021) Graph neural networks meet
1495 neural-symbolic computing: a survey and perspective. In: *Proceedings of the Twenty-Ninth International Joint*
1496 *Conference on Artificial Intelligence, IJCAI'20*. ISBN 9780999241165, pp. 4877–4884.
- 1497 Law M, Russo A, Bertino E, Broda K and Lobo J (2020a) Fastlas: Scalable inductive logic programming
1498 incorporating domain-specific optimisation criteria. *Proceedings of the AAAI Conference on Artificial*
1499 *Intelligence* 34(03): 2877–2885. DOI:10.1609/aaai.v34i03.5678. URL [https://ojs.aaai.org/index.](https://ojs.aaai.org/index.php/AAAI/article/view/5678)
1500 [php/AAAI/article/view/5678](https://ojs.aaai.org/index.php/AAAI/article/view/5678).

- 1501 Law M, Russo A and Broda K (2018) The complexity and generality of learning answer set programs.
1502 *Artificial Intelligence* 259: 110–146. DOI:<https://doi.org/10.1016/j.artint.2018.03.005>. URL <https://www.sciencedirect.com/science/article/pii/S000437021830105X>.
- 1503
1504 Law M, Russo A and Broda K (2019) Logic-based learning of answer set programs. *Reasoning Web. Explainable*
1505 *Artificial Intelligence: 15th International Summer School 2019, Bolzano, Italy, September 20–24, 2019,*
1506 *Tutorial Lectures* : 196–231 DOI:10.1007/978-3-030-31423-1_6. URL [https://doi.org/10.1007/](https://doi.org/10.1007/978-3-030-31423-1_6)
1507 [978-3-030-31423-1_6](https://doi.org/10.1007/978-3-030-31423-1_6).
- 1508 Law M, Russo A and Broda K (2020b) The ilasp system for inductive learning of answer set programs.
- 1509 Law M, Russo A, Broda K and Bertino E (2021) Scalable non-observational predicate learning in asp. In: Zhou
1510 ZH (ed.) *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21.*
1511 *International Joint Conferences on Artificial Intelligence Organization*, pp. 1936–1943. DOI:10.24963/ijcai.
1512 2021/267. URL <https://doi.org/10.24963/ijcai.2021/267>. Main Track.
- 1513 LeCun Y, Bengio Y and Hinton G (2015) Deep learning. *Nature* 521: 436–44. DOI:10.1038/nature14539.
- 1514 Lecun Y, Jackel L, Bottou L, Cortes C, Denker J, Drucker H, Guyon I, Muller U, Sackinger E, Simard P and Vapnik V
1515 (1995) Learning algorithms for classification: A comparison on handwritten digit recognition. *Neural networks*
1516 : 261–276.
- 1517 Leonetti M, Iocchi L and Stone P (2016) A synthesis of automated planning and reinforcement learning for efficient,
1518 robust decision-making. *Artificial Intelligence* 241: 103–130. DOI:<https://doi.org/10.1016/j.artint.2016.07.004>.
1519 URL <https://www.sciencedirect.com/science/article/pii/S0004370216300819>.
- 1520 Lewis M, Liu Y, Goyal N, Ghazvininejad M, Mohamed A, Levy O, Stoyanov V and Zettlemoyer L (2020) BART:
1521 Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension.
1522 In: Jurafsky D, Chai J, Schluter N and Tetreault J (eds.) *Proceedings of the 58th Annual Meeting of the*
1523 *Association for Computational Linguistics*. Online: Association for Computational Linguistics, pp. 7871–7880.
1524 DOI:10.18653/v1/2020.acl-main.703. URL <https://aclanthology.org/2020.acl-main.703/>.
- 1525 Li J, Li D, Xiong C and Hoi S (2022) BLIP: Bootstrapping language-image pre-training for unified vision-language
1526 understanding and generation. In: Chaudhuri K, Jegelka S, Song L, Szepesvari C, Niu G and Sabato S (eds.)
1527 *Proceedings of the 39th International Conference on Machine Learning, Proceedings of Machine Learning*
1528 *Research*, volume 162. PMLR, pp. 12888–12900. URL [https://proceedings.mlr.press/v162/](https://proceedings.mlr.press/v162/li22n.html)
1529 [li22n.html](https://proceedings.mlr.press/v162/li22n.html).
- 1530 Lifschitz V (2019) *Answer Set Programming*. 1st edition. Springer Publishing Company, Incorporated. ISBN
1531 3030246574.
- 1532 Lin F and Zhao Y (2004) Assat: computing answer sets of a logic program by sat solvers. *Artificial*
1533 *Intelligence* 157(1): 115–137. DOI:<https://doi.org/10.1016/j.artint.2004.04.004>. URL <https://www.sciencedirect.com/science/article/pii/S0004370204000578>. Nonmonotonic Reasoning.
- 1534
1535 Liu B, Jiang Y, Zhang X, Liu Q, Zhang S, Biswas J and Stone P (2023a) Llm+p: Empowering large language models
1536 with optimal planning proficiency. URL <https://arxiv.org/abs/2304.11477>.
- 1537 Liu Y, Han T, Ma S, Zhang J, Yang Y, Tian J, He H, Li A, He M, Liu Z, Wu Z, Zhu D, Li X, Qiang N, Shen D, Liu
1538 T and Ge B (2023b) Summary of chatgpt/gpt-4 research and perspective towards the future of large language
1539 models.
- 1540 Mack D and Jefferson A (2018) Clevr graph: A dataset for graph question answering. URL [github.com/](https://github.com/Octavian-ai/clevr-graph)
1541 [Octavian-ai/clevr-graph](https://github.com/Octavian-ai/clevr-graph).

- 1542 Manhaeve R, Dumančić S, Kimmig A, Demeester T and De Raedt L (2021) Neural probabilistic logic programming
1543 in deepprolog. *Artificial Intelligence* 298: 103504. DOI:<https://doi.org/10.1016/j.artint.2021.103504>. URL
1544 <https://www.sciencedirect.com/science/article/pii/S0004370221000552>.
- 1545 Manhaeve R, Giannini F, Ali M, Azzolini D, Bizzarri A, Borghesi A, Bortolotti S, De Raedt L, Dhami D, Diligenti M,
1546 Dumančić S, Faltings B, Gentili E, Gerevini A, Gori M, Guns T, Homola M, Kersting K, Lehmann J, Lombardi
1547 M, Lorello L, Marconato E, Melacci S, Passerini A, Paul D, Riguzzi F, Teso S, Yorke-Smith N and Lippi M
1548 (2026) Benchmarking in neuro-symbolic ai. In: Dai WZ (ed.) *Learning and Reasoning*. Cham: Springer Nature
1549 Switzerland. ISBN 978-3-032-09087-4, pp. 238–249.
- 1550 Manning C, Surdeanu M, Bauer J, Finkel J, Bethard S and McClosky D (2014) The Stanford CoreNLP natural
1551 language processing toolkit. In: Bontcheva K and Zhu J (eds.) *Proceedings of 52nd Annual Meeting of
1552 the Association for Computational Linguistics: System Demonstrations*. Baltimore, Maryland: Association for
1553 Computational Linguistics, pp. 55–60. DOI:10.3115/v1/P14-5010. URL [https://aclanthology.org/
1554 P14-5010/](https://aclanthology.org/P14-5010/).
- 1555 Marcus G (2020) The next decade in ai: Four steps towards robust artificial intelligence.
- 1556 Marple K, Salazar E, Chen Z and Gupta G (2017) The s(asp) predicate answer set programming system. URL
1557 <https://api.semanticscholar.org/CorpusID:195834851>.
- 1558 Marra G, Dumančić S, Manhaeve R and De Raedt L (2024) From statistical relational to neurosymbolic artificial
1559 intelligence: A survey. *Artificial Intelligence* 328: 104062. DOI:<https://doi.org/10.1016/j.artint.2023.104062>.
1560 URL <https://www.sciencedirect.com/science/article/pii/S0004370223002084>.
- 1561 Mirzaee R, Rajaby Faghihi H, Ning Q and Kordjamshidi P (2021) SPARTQA: A textual question answering
1562 benchmark for spatial reasoning. In: *Proceedings of the 2021 Conference of the North American Chapter
1563 of the Association for Computational Linguistics: Human Language Technologies*. Online: Association for
1564 Computational Linguistics, pp. 4582–4598. DOI:10.18653/v1/2021.naacl-main.364. URL [https://
1565 aclanthology.org/2021.naacl-main.364/](https://aclanthology.org/2021.naacl-main.364/).
- 1566 Misino E, Marra G and Sansone E (2022) Vael: Bridging variational autoencoders and probabilistic logic
1567 programming. *ArXiv* abs/2202.04178. URL [https://api.semanticscholar.org/CorpusID:
1568 246679982](https://api.semanticscholar.org/CorpusID:246679982).
- 1569 Mitra A and Baral C (2015) Learning to automatically solve logic grid puzzles. In: Márquez L, Callison-Burch C
1570 and Su J (eds.) *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*.
1571 Lisbon, Portugal: Association for Computational Linguistics, pp. 1023–1033. DOI:10.18653/v1/D15-1118.
1572 URL <https://aclanthology.org/D15-1118/>.
- 1573 Muggleton S (1991) Inductive logic programming. *New Gen. Comput.* 8(4): 295–318. DOI:10.1007/BF03037089.
1574 URL <https://doi.org/10.1007/BF03037089>.
- 1575 Möller M, Norlander A, Martires PZD and Raedt LD (2025) Neurosymbolic decision trees. URL [https://
1576 arxiv.org/abs/2503.08762](https://arxiv.org/abs/2503.08762).
- 1577 Nguyen QA, Pham TT, Vuong THY, Trinh VG and Thanh NH (2025) Detecting Misleading Information with
1578 LLMs and Explainable ASP. In: *ICAART 2025 - 17th International Conference on Agents and Artificial
1579 Intelligence*. Porto, Portugal: SCITEPRESS - Science and Technology Publications, pp. 1327–1334. DOI:
1580 10.5220/0013357400003890. URL <https://inria.hal.science/hal-04920600>.
- 1581 Odense S and d'Avila Garcez A (2025) A semantic framework for neurosymbolic computation. *Artificial Intelligence*
1582 340: 104273. DOI:<https://doi.org/10.1016/j.artint.2024.104273>. URL [https://www.sciencedirect.
1583 com/science/article/pii/S0004370224002091](https://www.sciencedirect.com/science/article/pii/S0004370224002091).

- 1584 OpenAI (2023) Gpt-4 technical report.
- 1585 Padalkar P, Wang H and Gupta G (2023) Using logic programming and kernel-grouping for improving interpretability
1586 of convolutional neural networks. In: Gebser M and Sergey I (eds.) *Practical Aspects of Declarative Languages*.
1587 Cham: Springer Nature Switzerland. ISBN 978-3-031-52038-9, pp. 134–150.
- 1588 Padalkar P, Wang H and Gupta G (2024) Nesyfold: a framework for interpretable image classification. In:
1589 *Proceedings of the Thirty-Eighth AAAI Conference on Artificial Intelligence and Thirty-Sixth Conference*
1590 *on Innovative Applications of Artificial Intelligence and Fourteenth Symposium on Educational Advances in*
1591 *Artificial Intelligence, AAAI'24/IAAI'24/EAAI'24*. AAAI Press. ISBN 978-1-57735-887-9, pp. 4378–4387.
1592 DOI:10.1609/aaai.v38i5.28235. URL <https://doi.org/10.1609/aaai.v38i5.28235>.
- 1593 Padalkar P, Ślusarz N, Komendantskaya E and Gupta G (2025) A neurosymbolic framework for bias correction
1594 in convolutional neural networks. *Theory and Practice of Logic Programming* 24: 644–662. DOI:10.1017/
1595 S1471068424000322.
- 1596 Parać R, Nodari L, Ardon L, Furelos-Blanco D, Cerutti F and Russo A (2024) Learning robust reward machines from
1597 noisy labels. In: *Proceedings of the 21st International Conference on Principles of Knowledge Representation*
1598 *and Reasoning, KR '24*. ISBN 978-1-956792-05-8, pp. 909–919. DOI:10.24963/kr.2024/85. URL <https://doi.org/10.24963/kr.2024/85>.
- 1600 Peng Y, Zha Z, Jin Y, Luo Z, Dai WZ, Ren Z, Ding YX and Zhou K (2025) Generating by understanding: Neural
1601 visual generation with logical symbol groundings. In: *Proceedings of the 31st ACM SIGKDD Conference*
1602 *on Knowledge Discovery and Data Mining V.2, KDD '25*. New York, NY, USA: Association for Computing
1603 Machinery. ISBN 9798400714542, p. 2291–2302. DOI:10.1145/3711896.3736978. URL <https://doi.org/10.1145/3711896.3736978>.
- 1605 Petersen F, Borgelt C, Kuehne H and Deussen O (2022) Deep differentiable logic gate networks. In: Oh AH,
1606 Agarwal A, Belgrave D and Cho K (eds.) *Advances in Neural Information Processing Systems*. pp. 2006–2018.
1607 URL <https://openreview.net/forum?id=vF3WefcoePW>.
- 1608 Podell D, English Z, Lacey K, Blattmann A, Dockhorn T, Müller J, Penna J and Rombach R (2024) SDXL: Improving
1609 latent diffusion models for high-resolution image synthesis. *The Twelfth International Conference on Learning*
1610 *Representations* URL <https://openreview.net/forum?id=di52zR8xgf>.
- 1611 Quattoni A and Torralba A (2009) Recognizing indoor scenes. In: *2009 IEEE Conference on Computer Vision and*
1612 *Pattern Recognition*. pp. 413–420. DOI:10.1109/CVPR.2009.5206537.
- 1613 Rader AP and Russo A (2023) Active learning in neurosymbolic ai with embed2sym. *COGAI@IJCLR* URL
1614 <https://api.semanticscholar.org/CorpusID:269089085>.
- 1615 Ragno A, Plantevit M, Robardet C and Capobianco R (2024) Transparent explainable logic layers. In: *European*
1616 *Conference on Artificial Intelligence*. ISBN 9781643685489, pp. 914–921. DOI:10.3233/FAIA240579.
- 1617 Rajasekharan A, Zeng Y and Gupta G (2023a) Argument analysis using answer set programming and semantics-
1618 guided large language models. *ICLP Workshops 2023* URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:259503415)
1619 [CorpusID:259503415](https://api.semanticscholar.org/CorpusID:259503415).
- 1620 Rajasekharan A, Zeng Y, Padalkar P and Gupta G (2023b) Reliable natural language understanding with large
1621 language models and answer set programming. In: Pontelli E, Costantini S, Dodaro C, Gaggl S, Calegari
1622 R, D'Avila Garcez A, Fabiano F, Mileo A, Russo A and Toni F (eds.) *Proceedings 39th International*
1623 *Conference on Logic Programming, Imperial College London, UK, 9th July 2023 - 15th July 2023, Electronic*
1624 *Proceedings in Theoretical Computer Science*, volume 385. Open Publishing Association, pp. 274–287. DOI:
1625 10.4204/EPTCS.385.27.

- 1626 Redmon J and Farhadi A (2018) Yolov3: An incremental improvement. URL [https://arxiv.org/abs/](https://arxiv.org/abs/1804.02767)
1627 [1804.02767](https://arxiv.org/abs/1804.02767).
- 1628 Riegel R, Gray A, Luus F, Khan N, Makondo N, Akhalwaya IY, Qian H, Fagin R, Barahona F, Sharma U, Ikbal
1629 S, Karanam H, Neelam S, Likhanyi A and Srivastava S (2020) Logical neural networks. URL [https:](https://arxiv.org/abs/2006.13155)
1630 [//arxiv.org/abs/2006.13155](https://arxiv.org/abs/2006.13155).
- 1631 Riley H and Sridharan M (2019) Integrating non-monotonic logical reasoning and inductive learning with deep
1632 learning for explainable visual question answering. *Frontiers in Robotics and AI* Volume 6 - 2019. DOI:10.
1633 3389/frobt.2019.00125. URL [https://www.frontiersin.org/journals/robotics-and-ai/](https://www.frontiersin.org/journals/robotics-and-ai/articles/10.3389/frobt.2019.00125)
1634 [articles/10.3389/frobt.2019.00125](https://www.frontiersin.org/journals/robotics-and-ai/articles/10.3389/frobt.2019.00125).
- 1635 Ruis L, Andreas J, Baroni M, Bouchacourt D and Lake BM (2020) A benchmark for systematic generalization in
1636 grounded language understanding. In: *Proceedings of the 34th International Conference on Neural Information*
1637 *Processing Systems, NIPS '20*. Red Hook, NY, USA: Curran Associates Inc. ISBN 9781713829546, pp. 19861–
1638 19872.
- 1639 Santana MB, Kareem I and Ricca F (2024) Towards automatic composition of asp programs from natural language
1640 specifications. In: *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence,*
1641 *IJCAI '24*. ISBN 978-1-956792-04-1, pp. 6198–6206. DOI:10.24963/ijcai.2024/685. URL [https://doi.](https://doi.org/10.24963/ijcai.2024/685)
1642 [org/10.24963/ijcai.2024/685](https://doi.org/10.24963/ijcai.2024/685).
- 1643 Sheth A, Roy K and Gaur M (2023) Neurosymbolic Artificial Intelligence (Why, What, and How) .
1644 *IEEE Intelligent Systems* 38(03): 56–62. DOI:10.1109/MIS.2023.3268724. URL [https://doi.](https://doi.ieeeecomputersociety.org/10.1109/MIS.2023.3268724)
1645 [ieeecomputersociety.org/10.1109/MIS.2023.3268724](https://doi.ieeeecomputersociety.org/10.1109/MIS.2023.3268724).
- 1646 Shi Z, Zhang Q and Lipani A (2021) Stepgame: A new benchmark for robust multi-hop spatial reasoning in
1647 texts. In: *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36. pp. 11321–11329. DOI:
1648 10.1609/aaai.v36i10.21383.
- 1649 Singh D, Jain N, Jain P, Kayal P, Kumawat S and Batra N (2020) Plantdoc: A dataset for visual plant disease detection.
1650 In: *Proceedings of the 7th ACM IKDD CoDS and 25th COMAD, CoDS COMAD 2020*. New York, NY, USA:
1651 Association for Computing Machinery. ISBN 9781450377386, p. 249–253. DOI:10.1145/3371158.3371196.
1652 URL <https://doi.org/10.1145/3371158.3371196>.
- 1653 Sinha K, Sodhani S, Dong J, Pineau J and Hamilton WL (2019) CLUTRR: A diagnostic benchmark for inductive
1654 reasoning from text. In: Inui K, Jiang J, Ng V and Wan X (eds.) *Proceedings of the 2019 Conference on*
1655 *Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural*
1656 *Language Processing (EMNLP-IJCNLP)*. Hong Kong, China: Association for Computational Linguistics, pp.
1657 4506–4515. DOI:10.18653/v1/D19-1458. URL <https://aclanthology.org/D19-1458/>.
- 1658 Skryagin A, Ochs D, Deibert P, Kohaut S, Dhimi DS and Kersting K (2024a) Answer set networks: Casting answer
1659 set programming into deep learning. URL <https://arxiv.org/abs/2412.14814>.
- 1660 Skryagin A, Ochs D, Dhimi DS and Kersting K (2024b) Scalable neural-probabilistic answer set programming. *J.*
1661 *Artif. Int. Res.* 78. DOI:10.1613/jair.1.15027. URL <https://doi.org/10.1613/jair.1.15027>.
- 1662 Skryagin A, Stammer W, Ochs D, Dhimi DS and Kersting K (2022) Neural-Probabilistic Answer Set Programming.
1663 *Proceedings of the 19th International Conference on Principles of Knowledge Representation and Reasoning :*
1664 463–473 DOI:10.24963/kr.2022/48. URL <https://doi.org/10.24963/kr.2022/48>.
- 1665 Smet LD and Raedt LD (2025) Defining neurosymbolic ai. URL <https://arxiv.org/abs/2507.11127>.
- 1666 Stallkamp J, Schlipsing M, Salmen J and Igel C (2012) Man vs. computer: Benchmarking machine
1667 learning algorithms for traffic sign recognition. *Neural Networks* 32: 323–332. DOI:[https://doi.org/10.](https://doi.org/10.1016/j.neunet.2011.12.014)

- 1668 1016/j.neunet.2012.02.016. URL <https://www.sciencedirect.com/science/article/pii/S0893608012000457>. Selected Papers from IJCNN 2011.
- 1669
- 1670 Suchan J, Bhatt M and Varadarajan S (2021) Commonsense visual sensemaking for autonomous driving – on
1671 generalised neurosymbolic online abduction integrating vision and semantics. *Artificial Intelligence* 299:
1672 103522. DOI:<https://doi.org/10.1016/j.artint.2021.103522>. URL <https://www.sciencedirect.com/science/article/pii/S0004370221000734>.
- 1673
- 1674 Tafford O, Clark P, Gardner M, Yih Wt and Sabharwal A (2019) Quarel: a dataset and models for answering questions
1675 about qualitative relationships. In: *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence*
1676 *and Thirty-First Innovative Applications of Artificial Intelligence Conference and Ninth AAAI Symposium on*
1677 *Educational Advances in Artificial Intelligence*, AAAI'19/IAAI'19/EAAI'19. AAAI Press. ISBN 978-1-57735-
1678 809-1, pp. 7063–7071. DOI:10.1609/aaai.v33i01.33017063. URL <https://doi.org/10.1609/aaai.v33i01.33017063>.
- 1679
- 1680 Thomson A and Page D (2023) Neural markov prolog. URL <https://arxiv.org/abs/2312.01521>.
- 1681 Tudor A and Gupta G (2024) Autonomous task completion based on goal-directed answer set programming. *ICLP*
1682 *Workshops* URL <https://api.semanticscholar.org/CorpusID:276307594>.
- 1683 Wan Z, Liu CK, Yang H, Li C, You H, Fu Y, Wan C, Krishna T, Lin Y and Raychowdhury A (2024a) Towards
1684 cognitive ai systems: a survey and prospective on neuro-symbolic ai. URL <https://arxiv.org/abs/2401.01040>.
- 1685
- 1686 Wan Z, Liu CK, Yang H, Raj R, Li C, You H, Fu Y, Wan C, Samajdar A, Lin YC, Krishna T and Raychowdhury A
1687 (2024b) Towards cognitive ai systems: Workload and characterization of neuro-symbolic ai. In: *2024 IEEE*
1688 *International Symposium on Performance Analysis of Systems and Software (ISPASS)*. pp. 268–279. DOI:
1689 10.1109/ISPASS61541.2024.00033.
- 1690 Wang H and Gupta G (2023) Fold-se: An efficient rule-based machine learning algorithm with scalable explainability.
1691 In: Gebser M and Sergey I (eds.) *Practical Aspects of Declarative Languages*. Cham: Springer Nature
1692 Switzerland, pp. 37–53.
- 1693 Wang R, Sun K and Kuhn J (2024) Dspy-based neural-symbolic pipeline to enhance spatial reasoning in llms. URL
1694 <https://arxiv.org/abs/2411.18564>.
- 1695 Weston J, Bordes A, Chopra S, Rush AM, van Merriënboer B, Joulin A and Mikolov T (2015) Towards ai-complete
1696 question answering: A set of prerequisite toy tasks. URL <https://arxiv.org/abs/1502.05698>.
- 1697 Winters T, Marra G, Manhaeve R and Raedt LD (2022) Deepstochlog: Neural stochastic logic programming.
1698 *Proceedings of the AAAI Conference on Artificial Intelligence* 36(9): 10090–10100. DOI:10.1609/aaai.v36i9.
1699 21248. URL <https://ojs.aaai.org/index.php/AAAI/article/view/21248>.
- 1700 Yang Z, Ishay A and Lee J (2020) Neurasp: Embracing neural networks into answer set programming. In:
1701 Bessiere C (ed.) *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence,*
1702 *IJCAI-20*. International Joint Conferences on Artificial Intelligence Organization, pp. 1755–1762. DOI:
1703 10.24963/ijcai.2020/243. URL <https://doi.org/10.24963/ijcai.2020/243>. Main track.
- 1704 Yang Z, Ishay A and Lee J (2023) Coupling large language models with logic programming for robust and general
1705 reasoning from text. In: *Findings of the Association for Computational Linguistics, ACL 2023*, Proceedings
1706 of the Annual Meeting of the Association for Computational Linguistics. Association for Computational
1707 Linguistics (ACL), pp. 5186–5219. DOI:10.18653/v1/2023.findings-acl.321. Publisher Copyright: © 2023
1708 Association for Computational Linguistics.; 61st Annual Meeting of the Association for Computational
1709 Linguistics, ACL 2023 ; Conference date: 09-07-2023 Through 14-07-2023.

- 1710 Yu D, Yang B, Liu D, Wang H and Pan S (2023) A survey on neural-symbolic learning systems. *Neural Networks*
1711 166: 105–126. DOI:<https://doi.org/10.1016/j.neunet.2023.06.028>. URL <https://www.sciencedirect.com/science/article/pii/S0893608023003398>.
- 1713 Zeng Y, Rajasekharan A, Basu K, Wang H, Arias J and Gupta G (2024) A reliable common-sense reasoning socialbot
1714 built using llms and goal-directed asp. *Theory and Practice of Logic Programming* 24(4): 606–627. DOI:
1715 10.1017/s147106842400022x. URL <http://dx.doi.org/10.1017/S147106842400022X>.
- 1716 Zeng Y, Rajasekharan A, Padalkar P, Basu K, Arias J and Gupta G (2023) Automated interactive domain-specific
1717 conversational agents that understand human dialogs. In: Gebser M and Sergey I (eds.) *Practical Aspects of*
1718 *Declarative Languages*. Cham: Springer Nature Switzerland. ISBN 978-3-031-52038-9, pp. 204–222.
- 1719 Zhang H, Kung PN, Yoshida M, den Broeck GV and Peng N (2024) Adaptable logical control for large language
1720 models. URL <https://arxiv.org/abs/2406.13892>.
- 1721 Zhou B, Lapedriza A, Khosla A, Oliva A and Torralba A (2018) Places: A 10 million image database for scene
1722 recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 40(6): 1452–1464. DOI:
1723 10.1109/TPAMI.2017.2723009.