Position: Efficient General Intelligence requires Neuro-Symbolic Integration: Pillars, Benchmarks, and Beyond

Anonymous Authors¹

Abstract

Recent breakthroughs in Large Language Model (LLM) development have rekindled hopes for broadly capable artificial intelligence. Yet, these models still exhibit notable limitations - particularly in deductive reasoning and efficient skill acquisition. In contrast, neuro-symbolic approaches, which integrate sub-symbolic pattern extraction with explicit logical structures, offer more robust generalization across diverse tasks. We argue that additional factors - such as modular transparency, flexible representations, and targeted prior knowledge - are crucial to further enhance this generalization. Our analysis of both historical and contemporary AI methods suggests that a multi-component neuro-symbolic implementation strategy is necessary for efficient general intelligence. This position is reinforced by the latest performance gains on the ARC-AGI benchmark and by concrete case studies demonstrating how neuro-symbolic designs address gaps left by purely neural or purely symbolic systems.

1. Introduction

We posit that truly general artificial intelligence demands the unification of data-driven neural approaches with the interpretability and compositional expressiveness of symbolic methods. Across decades of progress, research on "artificial intelligence" has often centered on narrow tasks and small leaps in computational automation, without necessarily pursuing robust, human-like intelligence. This changed with the rise of large, monolithic neural networks - models that excel in pattern extraction and display intriguing emergent capacities (Bubeck et al., 2023). Yet, while these black-box approaches are remarkable in many respects, they also suffer from opaque decision-making processes and often exhibit only local forms of generalization. They thus provide limited insights into the core mechanisms 051 underlying *flexible*, *human-level* intelligence.

Motivated by these gaps, an increasing number of researchers suggest incorporating symbolic reasoning or "transparent" inductive biases into deep learning pipelines, giving rise to *neuro-symbolic* approaches (d'Avila Garcez & Lamb, 2023; Keber et al., 2024). By preserving the neural model's strengths in statistical pattern recognition and combining them with symbolic structures that allow for compositional logic, explainable decisions, and interpretability, neuro-symbolic methods promise broader skill-acquisition efficiency, deeper semantic understanding, and safer realworld deployment (Hernández-Orallo, 2020; Hassija et al., 2024).

However, merely layering symbolic modules on top of neural back-ends does not automatically confer general intelligence. To foster meaningful progress, we must (i) clarify the very notion of intelligence, especially in terms of skill acquisition and generalization difficulty, and (ii) pinpoint how best to evaluate a model's capacity to abstract knowledge from sparse data and adapt to novel tasks. Below, we outline why Chollet (2019)'s emphasis on "skill-acquisition efficiency" is particularly fruitful, how debates on behaviorism versus internalism push us to seek more transparent model mechanisms, and why benchmarks specifically designed for general intelligence play such a central role. We then highlight the Abstraction and Reasoning Corpus (ARC) and its latest ARC-AGI-1 challenge (Chollet et al., 2025) as a tangible testbed for demonstrating how neuro-symbolic methods could open doors to true breadth of reasoning ability.

1.1. Defining Intelligence

Despite centuries of study, intelligence remains notoriously difficult to define comprehensively (Legg & Hutter, 2007). We adopt the formulation by Chollet (2019) that views the intelligence of a system as "a measure of its skill-acquisition efficiency over a scope of tasks, with respect to priors, experience, and generalization difficulty." This perspective shifts attention from raw *performance* on a single task to the *ability to learn new* tasks under constraints – such as limited data, novel transformations, or minimal prior knowledge. Indeed, skill-acquisition efficiency is at the heart of what sets "general" intelligence apart from specialized or over-engineered solutions (Bober-Irizar & Banerjee, 2024).

055 1.2. Behaviorism vs. Internalism

A longstanding philosophical debate pertains to whether 057 only external behavior matters (behaviorism) or whether 058 the internal mechanisms of thought carry essential explana-059 tory value (internalism). In contemporary machine learning, 060 this tension appears as "functionality vs. interpretability" 061 or "black-box vs. transparent systems." High-performing 062 but opaque models - like many Large Language Models -063 demonstrate that achieving sophisticated outputs does not 064 necessarily illuminate the process by which the model rea-065 sons (Hernández-Orallo, 2020; Schlangen, 2021). 066

067 As these systems are deployed in sensitive or high-stakes en-068 vironments, interpretability and control become paramount 069 (Hassija et al., 2024). Post-hoc explanations often provide 070 only a partial window into massive parameter spaces, leav-071 ing significant uncertainties about why a particular decision was reached (Kenny et al., 2021; Slack et al., 2021; Lee-073 mann et al., 2023; Rong et al., 2023). By contrast, inherent 074 model transparency - via symbolic modules, meaningful 075 structured interfaces, or modular architectures - can yield 076 more reliable comprehension of internal processes, facili-077 tate debugging, and bolster trustworthiness. Consequently, 078 we argue that internalist considerations should shape the 079 development of any model that aspires to broader, more 080 systematic intelligence. 081

0820831.3. Generalization Efficiency083

Even when a model attains notable performance on a suite of tasks, it is crucial to distinguish between intrinsic gener-085 alization and engineered solutions. Many recent successes hinge on massive data curation, architectural tuning, or man-087 ual injection of priors - leading to impressive system-centric results, but not necessarily reflecting a model's capacity to 089 autonomously learn how to solve unseen tasks. As Chol-090 let (2019) notes, a "developer-aware" perspective on skill 091 acquisition controls for these extra-human interventions. 092 Without such a perspective, higher benchmark scores risk 093 being misread as general intelligence. 094

Hence, if the field's ambition is true *general* intelligence –
rather than a proliferation of specialized or heavily handcrafted solutions – then adopting metrics and methods highlighting *skill-acquisition efficiency* becomes indispensable.
This, in turn, requires reliable ways to evaluate how well a
model performs under low-data, unseen, or compositional
scenarios – where brute-force training or naive memorization is infeasible.

1041051.4. Benchmarking for Generality

The search for a benchmark that isolates genuine abstraction and reasoning from mere pattern fitting has led to the
Abstraction and Reasoning Corpus (ARC) (Chollet, 2019),

later extended into ARC-AGI-1 (Chollet et al., 2025). Unlike tasks saturated by large, curated datasets, ARC consists of small, diverse puzzles that test "core knowledge" concepts like spatial manipulation, color/object transformations, or compositional logic (Moskvichev et al., 2023).

Despite being straightforward for humans, ARC tasks have proven unexpectedly difficult for computational models, with only about half the tasks consistently solved (Bober-Irizar & Banerjee, 2024; ARC Prize, 2024). This difficulty emerges precisely because ARC demands *abstract generalization* over a minimal set of examples, thwarting superficial shortcuts. While ARC alone is not a perfect proxy for all human-level reasoning (Chollet, 2019), it remains a valuable gauge of small-data adaptability, creative knowledge transfer, and flexible problem solving. Table 1 summarizes the key dimensions for designing models with broad generalization capabilities.

In what follows, we leverage ARC-AGI-1 to motivate why **hybrid neuro-symbolic architectures** – infused with explicit mechanisms for transparency, compositional reasoning, and high-level knowledge abstraction – are integral to bridging the gap from narrow task competence to more truly *general* intelligence.

2. Alternative Views

In this position paper, we argue that *modular neuro-symbolic integration* is central for achieving efficient generalization. Nevertheless, it is important to recognize other significant perspectives in the field, especially since some researchers propose that *either* purely neural or purely symbolic methods might suffice if combined with enough data, computational resources, or engineering effort. Below, we discuss these alternatives – large language models (LLMs) as an archetype of purely neural approaches, and domain-specific languages (DSLs) or program-synthesis approaches as a representative of purely symbolic strategies – and evaluate why they each have notable strengths yet ultimately fall short when it comes to *broadly* efficient generalization.

2.1. Purely Neural Approaches: Large Language Models

Transformers and large language models (LLMs) have undeniably exhibited broad emergent capabilities, including surprising generalization and few-shot reasoning, across multiple domains (Bubeck et al., 2023; Webb et al., 2023). Remarkably, they can perform competitively even on the Abstraction and Reasoning Corpus (ARC) when equipped with skillful prompting, chain-of-thought techniques, and *selfimprovement* querying (Greenblatt, 2024; Berman, 2024; ARC Prize, 2025; Chollet et al., 2025). Indeed, GPT-4, Sonnet 3.5, and *o3* consistently achieve the highest ARC-AGI-1

Dimension	Importance for General Intelligence	Representative Works	
Skill-Acquisition Efficiency	Emphasizes how well a system adapts to new tasks with- out extensive retraining; penalizes overreliance on devel- oper engineering or huge datasets.	(Chollet, 2019), (Bober-Irizar & Banerjee, 2024)	
Transparency & Interpretability	Strengthens trust and debugging; post-hoc explanations(Hernánare often insufficient for large black-box models. Inher- ent transparency is crucial for real-world reliability.2020), (12024)		
Symbolic Reasoning	Allows compositional, logically coherent transforma- tions. Fosters human-level abstraction and provides ro- bust handling of discrete structures.(d'Avila G Lamb, 2023 et al., 2024)		
Neural Representations	Harnesses powerful pattern-extraction capabilities from raw data (images, text), enabling feature discovery and capturing nuanced correlations.	from (Bubeck et al., 2023) 7 and	
Small-Data Adaptation	Avoids brute-forcing solutions by demanding strong gen- eralization from very few examples (as in ARC tasks), exposing true abstraction capabilities.	(Moskvichev et al., 2023), (Chollet et al., 2025)	

Table 1. Key Dimensions for Designing Models with Broad Generalization

public scores when allowed large-scale test-time optimization. This level of success leads many researchers to view LLMs as the foundation for future general-purpose AI.

124 125

134

135

136

137

138

139

140

141

142

143

144

145

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161 162

163

164

Strengths of LLMs. Modern LLMs have several strengths ranging from a wide knowledge coverage to reasoning capabilities and flexibility when applied to downstream tasks.

- **Pre-training on Massive Corpora** allows for extensive self-supervised learning on diverse text sources. In this way, LLMs acquire a wealth of representations, effectively consolidating and covering wide-ranging knowledge (Bubeck et al., 2023).
- Flexible Transfer of Knowledge can be applied to handle various downstream tasks (including non-linguistic tasks expressed in language) with minimal fine-tuning, thanks to in-context learning capabilities and powerful embedding spaces (Dong et al., 2023; Berman, 2024).
- Emergent Reasoning Behaviors can be elicited through prompting strategies such as chain-of-thought or retrieval augmented generation. Such reasoning-like procedures within LLMs often improve the performance on complex tasks (Webb et al., 2023).

Challenges and Limitations. Despite impressive benchmarks, purely neural methods still exhibit significant hurdles regarding *efficient* generalization:

1. **Opaque and Brittle Emergence:** The extent to which LLMs can perform genuine abstract reasoning (versus

pattern matching) remains an open debate (Valmeekam et al., 2023; Kaddour et al., 2023; Dziri et al., 2023; Lewis & Mitchell, 2024; Wang et al., 2024; Lotfi et al., 2024; Schuurmans et al., 2024). Their "emergent" abilities can be unreliable, hard to interpret, and domainspecific (Bober-Irizar & Banerjee, 2024).

- Data-Hungry and Costly: Training large-scale transformers demands massive, human-generated corpora

 and some fear we are reaching the upper limit of high-quality data for further scaling this approach (Sutskever, 2024). In addition, fine-tuning or test-time brute forcing can be expensive and inefficient (Sachdeva et al., 2024).
- 3. **Developer vs. Model Intelligence:** Many LLM-based successes rely heavily on *engineered prompting* and human-coded heuristics. Thus, high-level performance may reflect *developer-centric* skill more than an intrinsic model capacity for generalization (Chollet, 2019; Dong et al., 2023; Yu et al., 2023; Bober-Irizar & Banerjee, 2024).
- 4. Lack of Transparency: Unlike modular designs, LLMs encode reasoning steps in vast weight matrices, limiting interpretability. This black-box nature impedes deeper analysis of the reasoning process and complicates improvements targeted at genuine compositional intelligence (Garcez & Lamb, 2023).

Moreover, recent ARC results reveal that while LLM-based approaches can outperform other methods on the public

165 benchmark, they do so through *massive prompt engineering* or resource-intensive test-time synthesis (Greenblatt, 2024; 167 Berman, 2024). Mahowald et al. (2024) draw parallels to the 168 human brain's specialized "language areas," cautioning that 169 forcing a language-dominant model to cover abstract non-170 linguistic tasks may be fundamentally inefficient. Hence, 171 even though LLMs are powerful in practice, they are less 172 suitable as an academic research framework for understand-173 ing the mechanisms behind generalization.

174 175

184

185

186

187

188

189

190

191

193

195

196

197

198

199

200

201

202

204

205

206

208

209

210

211

212

213

214

215

216

217

218

219

176 Conclusion for LLMs. While purely neural approaches have reshaped modern AI, purely monolithic LLMs appear 177 suboptimal as a basis for broad efficiency. Large data com-178 179 bined with sufficient computing resources can brute force solutions, but they do not illuminate the core processes un-180 derlying abstract reasoning. For those interested in deeper 181 interpretability, explainability, or developer-aware skill ac-182 183 quisition, neuro-symbolic integration seems indispensable.

2.2. Purely Symbolic Approaches: Domain-Specific Languages and Program Synthesis

Although overshadowed by neural methods in recent years, purely symbolic or logic-based AI once dominated the field and retains a devoted following (Kastner & Hong, 1984). Within the ARC domain, the most visible symbolic attempts revolve around exhaustive search in a *Domain-Specific Language* (DSL) or program-synthesis methods such as Dream-Coder (Ellis et al., 2020).

DSL-Based Methods. Early top-ranked solutions in the original ARC challenge relied on large, hand-crafted DSLs (icecuber, 2020; de Miquel, 2020; Larchenko, 2020). By systematically searching over a predefined set of transformations and heuristics, these approaches found valid transformations for specific puzzles. However, these DSL-based methods achieved only modest coverage due to the combinatorial explosion of possible transformations and the diversity of ARC tasks. They also demanded extensive human engineering to hard-code each concept, undermining *developer-aware* generalization measures (Bober-Irizar & Banerjee, 2024).

Program Synthesis Approaches. Program-synthesis frameworks like DreamCoder (Ellis et al., 2020) extend the DSL idea with higher-level constructs (e.g., control-flow operators, recursion). While this unlocks greater expressiveness, it can also inflate the search space. Adapting a fully general programming language for ARC tasks becomes cumbersome because ARC-AGI-1 is already quite challenging without further increasing the solution space (Bober-Irizar & Banerjee, 2024).

Symbolic Drawbacks. While symbolic approaches can offer strong interpretability (one can often track each logical step explicitly), they typically struggle to infer abstract "core concepts" from limited data without some learned inductive biases. Their purely top-down logic has trouble coping with the noisy, high-dimensional input distributions where data-driven feature extraction is crucial. Additionally, naive symbolic search tends to be fragile in the face of tasks requiring approximate or probabilistic reasoning.

Conclusion for Symbolic Methods. Historically, purely symbolic solutions have rarely scaled well across diverse tasks and have difficulty encoding robust priors for low-data settings (Kastner & Hong, 1984; Ellis et al., 2020). The *ARC experience* confirms that exhaustive or highly engineered symbolic DSLs rapidly reach diminishing returns. Hence, purely symbolic approaches, while valuable for interpretability and logic, alone are insufficient for broad or efficient generalization.

2.3. Synthesis of Both Views

In summary, purely neural approaches (e.g., LLMs) can demonstrate remarkable capabilities but often rely on extensive engineering, computational resources, and data, with limited inherent interpretability. Purely symbolic approaches retain logical clarity but cannot cope effectively with the complexity and ambiguity that general tasks demand. We, therefore, see *neuro-symbolic integration* – the systematic coupling of learnable neural components with explicit symbolic representations and reasoning – as the most promising route to truly efficient, transparent generalization, far beyond either paradigm alone.

3. Neuro-Symbolic Approaches

Neuro-symbolic methods stand at the intersection of statistical learning and explicit symbolic reasoning, offering a promising path toward *efficient* generalization. As discussed in Section 2, purely neural or purely symbolic methods each have strengths, but neither alone excels at developer-aware skill acquisition. Nevertheless, these two paradigms clearly complement each other (Bober-Irizar & Banerjee, 2024), already hinting at the power of *neuro-symbolic integration* to tackle a broader range of tasks (Bober-Irizar & Banerjee, 2024; Chollet et al., 2025).

Multiple works have surveyed the general advantages and disadvantages of neuro-symbolic approaches in depth (Hamilton et al., 2022; Hitzler et al., 2022; Garcez & Lamb, 2023; Keber et al., 2024; Bhuyan et al., 2024). Rather than revisiting all of these aspects, we focus here on the key *generalization* benefits, underscored by the most recent ARC-AGI-1 findings (Chollet et al., 2025). 220 **Terminology.** The term "neuro-symbolic" (sometimes ab-221 breviated "NeSy") can encompass a wide variety of hybrid 222 architectures and learning strategies. While the specific 223 mechanisms vary, the core idea is to marry symbolic struc-224 tures (e.g., logic programs, DSLs, knowledge graphs) with 225 neural components (e.g., deep networks or learned embed-226 dings) (Hitzler et al., 2022; Garcez & Lamb, 2023; Keber 227 et al., 2024). 228

Table 2 summarizes the main aspects of representative stateof-the-art NeSy approaches for generalization in ARC-like tasks.

229

230

231

232

247

Coming from the Neural Side. In Large Language Mod-233 els (LLMs), one could argue that chain-of-thought prompt-234 ing or structured "reasoning graphs" already hint at neuro-235 symbolic principles (Hitzler et al., 2022; Keber et al., 2024). 236 These techniques often wrap a neural transformer in a scaf-237 fold of symbolic instructions or constraints (Yu et al., 2023), 238 thus improving performance across multiple tasks. For ex-239 ample, Xu et al. (2024) demonstrate how extensive logical 240 orchestration around LLM calls boosts reliability on diverse 241 tasks. Notably, the top LLM-based ARC-AGI-1 approaches 242 also incorporate symbolic heuristics to stabilize generaliza-243 tion - further evidence that purely sub-symbolic solutions 244 remain insufficient (Franzen et al., 2024; Barbadillo, 2024; 245 Chollet et al., 2025). 246

248 Coming from the Symbolic Side. Conversely, the golden 249 era of symbolic AI faded in the late 1980s, giving way to 250 sub-symbolic (neural) approaches. However, the limitations 251 that once suffocated symbolic AI - such as brittle rule sys-252 tems or exponential search complexity – can be mitigated by 253 modern neural advances and computing power (Mira, 2008). 254 Rather than reviving purely symbolic methods, researchers 255 increasingly aim to harness the strengths of both paradigms: 256 explicit logic for interpretability and systematic abstraction, 257 and neural modules for data-driven feature extraction and 258 robustness (Hitzler et al., 2022; Garcez & Lamb, 2023). A 259 clear illustration is Bober-Irizar & Banerjee (2024), who build upon a DSL-based ARC solver by adding learnable 261 "concept formation" components, significantly boosting ef-262 ficiency and success rates. 263

264 Synergy in Practice. One prominent driver behind neuro-265 symbolic integration is generalization efficiency (Bhuyan 266 et al., 2024). Hybrid models can learn abstract concepts 267 more compactly, leveraging both (i) a neural module to han-268 dle noisy or high-dimensional inputs and (ii) a symbolic 269 module to enforce logical coherence and compositional rea-270 soning. This synergy is particularly relevant in low-data 271 tasks like ARC, where purely neural systems often overfit, 272 and purely symbolic systems lack robust inductive priors. 273 While recent work has demonstrated promising gains on 274

ARC (Moskvichev et al., 2023; Chollet et al., 2025; Bober-Irizar & Banerjee, 2024), open challenges remain – most notably:

- Exploding Search Spaces. Combining symbolic search with neural heuristics can mitigate the worst-case combinatorial complexity explosion, but designing these heuristics remains nontrivial (Bober-Irizar & Banerjee, 2024).
- Data Efficiency vs. Model Complexity. ARC-AGI-1 tasks demand strong reasoning from minimal examples, stressing the importance of balanced architectures that do not over-parameterize (Moskvichev et al., 2023).
- Formation of New Concepts. Handling ever-evolving domains requires neuro-symbolic methods that can *learn new concepts dynamically* rather than rely solely on a hard-coded DSL (Bober-Irizar & Banerjee, 2024).

Though these obstacles are significant, the ability of neurosymbolic methods to unify inductive and deductive reasoning is an especially potent strength – analogous to "System 1" vs. "System 2" thinking in human cognition (Kahneman, 2011; Garcez & Lamb, 2023). As computational and data constraints grow more urgent, *this marriage of neural and symbolic approaches* will likely become not just beneficial but *indispensable*.

Closing the Gap. Ultimately, the goal of neuro-symbolic research is to exploit each paradigm's complementary strengths: neural networks excel at fast, intuitive processing of raw data, whereas symbolic formalisms enable explicit logic and compositional abstraction. By weaving these together, a system can move beyond behavioristic success into genuine *skill-acquisition efficiency*, operating effectively with minimal data or developer engineering while staying transparent, interpretable, and controllable. In the following sections, we delve deeper into the specific design components and synergy effects that make neuro-symbolic architectures uniquely suited to achieving broader generalization.

4. Pillars of Efficient Neuro-Symbolic Generalization

Achieving robust generalization through neuro-symbolic methods requires more than simply pairing a neural module with a symbolic one. As Odense & Garcez (2022)(p. 38) argue, the key is to exploit the "complementary strengths and weaknesses" of both connectionist and symbolic paradigms – rather than letting one approximate or overshadow the other. Consequently, the central question in the years ahead is *how* to fuse these components so that they collectively yield effective skill acquisition across diverse tasks. We propose

Position: Efficient General Intelligence requires Neuro-Symbolic Integration

Approach	Neural Component	Symbolic Compo- nent	Key Mechanism & Insights
Bober-Irizar & Banerjee (2024) (Bober- Irizar & Banerjee, 2024)	Learned concept- formation module (e.g., CNN-like em- beddings to identify object features)	DSL-based program search for transforma- tions	Uses neural heuristics to guide symbolic search, significantly reducing the DSL's combinatorial explosion. Demonstrates no table gains on ARC tasks versus purely sym bolic baselines.
SearChain (Xu et al., 2024)	Large Language Model (transformer) for reasoning over prompts	Search framework with symbolic con- straints (e.g., BFS or rule-based expan- sions)	Combines "search in the chain" logic with LLM prompting; symbolic scaffolding con strains the neural model's proposed trans formations, improving reliability on diverse puzzle-solving tasks.
DreamCoder (Ellis et al., 2020)	Neural "wake-sleep" cycle that learns com- mon subroutines or concepts	Inductive program synthesis in a high- level language (with control-flow, recur- sion)	Iteratively refines a library of reusable functions—symbolic <i>abstractions</i> —guided by neural scoring. <i>DreamCoder</i> is not specifically designed for ARC but illustrates how learned domain knowledge can be symbolically encoded.
Neuro-Symbolic DSL Enhancements (various) (Hamilton et al., 2022; Hitzler et al., 2022; Garcez & Lamb, 2023; Bhuyan et al., 2024)	Neural embeddings for object detection, classification, or spa- tial feature extraction	Logic-based DSL or ontology enforcing compositional rules	General family of hybrid methods: neural modules handle perceptual tasks or fuzzy matches, while symbolic DSL enforces interpretability and constraint satisfaction Shown to improve data-efficiency and inter pretability on small "grid-world" or ARC like puzzles.

Table 2. Representative Neuro-Symbolic Approaches for Generalization in ARC-like Tasks

that the following fundamental pillars are indispensable for attaining efficient (developer-aware) generalization in neuro-symbolic systems.

4.1. Multi-Component Synergy Effects

275

305 306

307

308

309

311 Although merging neural and symbolic layers is crucial, 312 other "side problems" - such as representation strategies, 313 uncertainty handling, and knowledge encoding - are equally 314 significant for broad-scope generalization (Bhuyan et al., 315 2024). Many advanced methods overlook at least one di-316 mension (e.g. using trivial transformations or underpowered representations), losing potential flexibility (Franzen et al., 318 2024; Berman, 2024). In contrast, a systematic approach 319 that addresses each sub-component fosters powerful synergy 320 effects between them(Garcez & Lamb, 2023).

322 While modular integration entails substantial engineering 323 (Garcez & Lamb, 2023) and necessitates careful data-format 324 alignment, Bober-Irizar & Banerjee (2024) demonstrate its 325 worth: their neuro-symbolic concept-formation technique, inspired by Ellis et al. (2020), overshadow naive DSL search 327 by leveraging richer learned representations and heuristics 328 to prune the search space. 329

Takeaway: By holistically optimizing each component in the system, one transcends individual contributions and achieves system-wide synergy, enabling more capable and efficient generalization.

4.2. Model Specificity

Global ambitions need not result in "solve-all" monstrosities. Instead, concentrating on a well-defined domain – such as ARC's core priors – prevents runaway complexity (Chollet et al., 2025). Despite ARC's seemingly simple puzzles (e.g. shape manipulation), thorough mastery proves nontrivial without the "core knowledge" priors they were constructed from (Chollet, 2019; Ellis et al., 2020). Meanwhile, large foundation models often require careful prompting to isolate relevant priors (Greenblatt, 2024; Berman, 2024). Their strengths in specialized world knowledge or linguistic reasoning are not relevant for ARC(Chollet et al., 2025)(p.1).

Takeaway: A targeted model scope, with sufficient coverage of relevant key primitives yet focused capabilities, yields a broad solution space while still being feasible and tractable.

Position:	Efficient	General	Intelligence	requires	Neuro-Symbolic	Integration
-----------	-----------	---------	--------------	----------	----------------	-------------

Pillar	Why It Matters	Key Challenges	Representative Works
Multi-Component Synergy (4.1)	Achieves more powerful system-wide effects tran- scending capacities of individual components	Engineering complexity; aligning data formats and abstraction levels of modules (e.g. integrating neural (pattern extraction) and symbolic (logical composi- tion) capacities for broader coverage)	(Ellis et al., 2020; Garcez & Lamb, 2023; Bober-Irizar & Banerjee, 2024)
Model Specificity (4.2)	Prevents scope bloat by fo- cusing on well-defined core concepts; avoids extraneous features	Balancing breadth vs. tractability; ensuring funda- mental priors are covered	(Chollet, 2019; Chollet et al., 2025)
Knowledge Encod- ing (4.3)	Embeds abstract human in- sights, reducing the burden of brute-force or data-heavy engineering	Over-encoding task-specific solutions; selecting which concepts to "hard-code" vs. learn	(icecuber, 2020; de Miquel, 2020; Bober-Irizar & Banerjee, 2024)
Knowledge Acqui- sition & Transfer (4.4)	Captures new concepts from training data and reuses them adaptively at test time	Designing effective but flexi- ble training paradigms (e.g. curriculum learning, test- time fine-tuning)	(Ellis et al., 2020; Akyürek et al., 2024; Bober-Irizar & Banerjee, 2024; Chollet et al., 2025)
Representation (4.5)	Governs how data is parsed and manipulated (object- based, graph-based, etc.), greatly influencing efficiency and model scope	Deciding the optimal abstrac- tion level (e.g., pixels vs. ob- jects) and bridging neural embeddings with symbolic structures	(Xu et al., 2023a; Skean et al., 2024; Barbadillo, 2024)
Abstractions & Hi- erarchies (4.6)	Filters superfluous detail, en- abling compositional reason- ing and meaningful represen- tations	Choosing the number and granularity of layers; ensur- ing each abstraction captures meaningful transformations	(Krizhevsky et al., 2017; Xu et al., 2023b)

Table 3. Key Pillars of Efficient Neuro-Symbolic Generalization

4.3. Knowledge Encoding

Symbolic frameworks excel at instilling human knowledge, yet enumerating solution scripts for each task kills adaptability. Instead, defining process-level abstractions (e.g. "move(object, vector)") fosters reusability across countless tasks (Xu et al., 2023a). Ellis et al. (2020, p. 18) emphasize that "rich systems of built-in knowledge" radically accelerate learning – a stance aligning with the principle that broad competence arises from fundamental, composable operators. Takeaway: Injecting abstract human expertise (conceptlevel rather than solution-level) boosts data efficiency and encourages flexible reuse.

4.4. Knowledge Acquisition, Transfer, and Combination

No matter how thorough the initial knowledge encoding, new tasks inevitably appear. Thus, a neuro-symbolic system 384

must *learn* fresh concepts during training and *recombine* them spontaneously at inference (Chollet, 2019). ARC-AGI-1 showcases how test-time fine-tuning (TTFT) can be essential for unseen tasks (Akyürek et al., 2024; Chollet et al., 2025). Likewise, DreamCoder's "sleep-wake" cycle continuously refines a library of existing abstractions (Ellis et al., 2020), which Bober-Irizar & Banerjee (2024) adapt to handle ARC's diverse puzzle types.

Takeaway: Flexible generalization arises from continual concept formation plus dynamic adaptation at test time.

4.5. Representation

Representational design profoundly shapes a system's ability to generalize. While neural embeddings capture latent structure, they can be overly broad for specialized tasks like ARC (Garcez & Lamb, 2023; Skean et al., 2024). On the

385 other hand, graph- or object-centric representations simplify 386 transformations (Xu et al., 2023a), thus reducing search 387 complexity and clarifying model behavior. Replacing pixel-388 level manipulations with object-level reasoning, for instance, 389 can diminish the needed symbolic operator set by an order 390 of magnitude as respective ARC puzzles are situated on this abstraction level(Xu et al., 2023a). In contrast to abstraction 392 capabilities, which are more processing-focused, representation spaces reflect the model's perspectives on the world (i.e., world model)(Huh et al., 2024; Barbadillo, 2024).

Takeaway: Accurately aligning representations with the *natural granularity* of the domain maintains interpretability and computational efficiency while setting a meaningful scope for the model.

400 401 **4.6. Abstractions and Hierarchies**

402 Layered abstractions are foundational to both human cogni-403 tion and deep-network architectures (LeCun et al., 1989; 404 Riesenhuber & Poggio, 1999; Grill-Spector & Malach, 405 2004; Krizhevsky et al., 2017). In the ARC-AGI-1 context, 406 moving from pixel-level to object- or pattern-level oper-407 ations delivers major efficiency improvements (Xu et al., 408 2023a;b). Each abstracted layer or module discards noisy 409 details, accentuating shared structures across tasks while 410 bolstering interpretability. 411

Takeaway: *Hierarchical design* combines low-level perception and high-level logic, enabling compositional reasoning and meaningful explanations/representations.

4.7. Concluding Remarks on the Pillars

Collectively, the six pillars in Table 3 - multi-component synergy, model specificity, knowledge encoding, knowledge acquisition/transfer, representation, and hierarchical abstractions - constitute the blueprint for efficient neuro-symbolic generalization. Of course, they are not exclusive to neuro-symbolic approaches and are potentially useful in other domains, too. When each is addressed deliberately and woven together cohesively, the connectionist-symbolic merger achieves far more than either paradigm alone. From DSL-based solutions fortified by learned heuristics (Bober-Irizar & Banerjee, 2024) to LLM-driven systems guided by symbolic constraints (ARC Prize, 2025), such interplay has already advanced complex reasoning tasks in ARC. We posit, therefore, that fully engaging these pillars is indispensable for the next leap in developer-aware, data-efficient, and transparently interpretable general AI.

5. Conclusion

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

The pillars outlined in Section 4 – from multi-component synergy and model specificity to hierarchical abstraction – reinforce and depend upon one another. Their interplay is precisely what enables systems to generalize, adapt, and recombine knowledge in new settings with minimal data or developer engineering. We contend that *this synergy* lies at the heart of intelligence itself: carefully crafted representations, thoughtfully curated curricula, and dynamic transfer strategies collectively drive skill-acquisition efficiency.

Interestingly, such modular-yet-intertwined specialization echoes the structure of the human brain. While AI need not mimic biology directly, the brain's functional organization offers strong evidence that partitioning cognition into cooperating modules – symbolic and neural – is both more efficient and more transparent than a monolithic design (Kahneman, 2011). The Abstraction and Reasoning Corpus (ARC) challenges highlight the difficulties and the *promise* of this approach (Chollet et al., 2025).

Indeed, *proper* neuro-symbolic integration has been recognized as crucial yet technically daunting (Garcez & Lamb, 2023; Chollet et al., 2025). Initial development costs may exceed those of monolithic solutions, but the reward – developer-aware, data-efficient systems that can flexibly adapt – is immense. Although our pillars are not fundamentally new in isolation, acknowledging how they interconnect offers a fresh lens on the essence of generalization.

To move the field forward, we propose **four priority direc-tions**:

- **Benchmarks for Synergy**: More comprehensive test suites (i.e., by extending ARC-AGI-style tasks) to systematically measure how effective the particular strengths of symbolic and neural components are utilized.
- **Open-Source DSL-Neural Frameworks**: Facilitating modular experimentation to ensure that promising ideas can be tested swiftly and reproducibly.
- **Interpretable Module Integration**: Standardized protocols to visualize how learned representations and symbolic rules interact.
- **Safety and Trust**: Deploying neuro-symbolic designs in safety-critical domains (medical, autonomous systems) to enhance transparency and reliability.

We encourage the research community to intensify its focus on these *synergy-driven* neuro-symbolic methods. While achieving truly general intelligence is undeniably challenging, *starting now* and grappling with multi-component integration – rather than evading it – stands to unlock the next great leap in flexible, safe, and transparently interpretable general AI.

440 Impact Statement

441 By advocating a synergy of data-driven and symbolic ap-442 proaches, this work paves the way for AI systems that learn 443 efficiently, reason transparently, and adapt efficiently to new 444 challenges. While our primary goal is to advance machine 445 learning methodologies, the ripple effects of more robust 446 and interpretable AI could reshape various sectors - from 447 healthcare and finance to education - by bolstering reliabil-448 ity and trustworthiness. 449

450 Compared to black-box models, neuro-symbolic solutions 451 are inherently more amenable to auditing and control, miti-452 gating risks of hidden biases or unintended behaviors. Such 453 transparency is essential for aligning AI with societal val-454 ues. We thus believe that emphasizing developer-aware 455 intelligence, modular design, and clear human oversight of-456 fers a safer, more equitable foundation for the technology's 457 continued evolution. 458

References

459

460

- 461 Akyürek, E., Damani, M., Qiu, L., Guo, H., Kim,
 462 Y., and Andreas, J. The Surprising Effectiveness of
 463 Test-Time Training for Abstract Reasoning, Novem464 ber 2024. URL http://arxiv.org/abs/2411.
 465 07279. arXiv:2411.07279 [cs].
- 466
 467
 468
 469
 469
 469
 469
 460
 460
 461
 462
 463
 464
 465
 465
 466
 467
 468
 469
 469
 469
 469
 469
 469
 460
 460
 461
 462
 463
 464
 464
 465
 465
 466
 467
 468
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
 469
- 470 ARC Prize, I. OpenAI o3 Breakthrough High Score on ARC471 AGI-Pub, 2025. URL https://arcprize.org/
 472 blog/oai-o3-pub-breakthrough.
- 473
 474
 475
 475
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
 476
- Berman, J. How I came in first on ARC-AGI-Pub using Sonnet 3.5 with Evolutionary Test-time Compute, 2024. URL https://jeremyberman.substack.com/ p/how-i-got-a-record-536-on-arc-agi.
- Bhuyan, B. P., Ramdane-Cherif, A., Tomar, R., and Singh,
 T. P. Neuro-symbolic artificial intelligence: a survey. *Neural Computing and Applications*, 36(21):12809–
 12844, July 2024. ISSN 1433-3058. doi: 10.1007/
 s00521-024-09960-z. URL https://doi.org/10.
 1007/s00521-024-09960-z.
- Bober-Irizar, M. and Banerjee, S. Neural networks for abstraction and reasoning. *Scientific Reports*, 14(1): 27823, November 2024. ISSN 2045-2322. doi: 10.1038/ s41598-024-73582-7. URL https://www.nature. com/articles/s41598-024-73582-7. Publisher: Nature Publishing Group.

- Bubeck, S., Chandrasekaran, V., Eldan, R., Gehrke, J., Horvitz, E., Kamar, E., Lee, P., Lee, Y. T., Li, Y., Lundberg, S., Nori, H., Palangi, H., Ribeiro, M. T., and Zhang, Y. Sparks of Artificial General Intelligence: Early experiments with GPT-4, April 2023. URL http://arxiv. org/abs/2303.12712. arXiv:2303.12712 [cs].
- Chollet, F. On the Measure of Intelligence, November 2019. URL http://arxiv.org/abs/1911.01547. arXiv:1911.01547 [cs].
- Chollet, F., Knoop, M., Kamradt, G., and Landers, B. ARC Prize 2024: Technical Report, January 2025. URL http://arxiv.org/abs/2412. 04604. arXiv:2412.04604 [cs].
- d'Avila Garcez, A. S. and Lamb, L. C. Neurosymbolic AI: the 3rd wave. *Artificial Intelligence Review*, 56(11): 12387–12406, November 2023. ISSN 1573-7462. doi: 10.1007/s10462-023-10448-w.
- de Miquel, A. 2020 kaggle arc-agi-1 challange: 2nd place solution, 2020. URL https://www.kaggle.com/competitions/ abstraction-and-reasoning-challenge/ discussion/154391.
- Dong, Q., Li, L., Dai, D., Zheng, C., Wu, Z., Chang, B., Sun, X., Xu, J., Li, L., and Sui, Z. A Survey on Incontext Learning, June 2023. URL http://arxiv. org/abs/2301.00234. arXiv:2301.00234 [cs].
- Dziri, N., Lu, X., Sclar, M., Li, X. L., Jiang, L., Lin, B. Y., West, P., Bhagavatula, C., Bras, R. L., Hwang, J. D., Sanyal, S., Welleck, S., Ren, X., Ettinger, A., Harchaoui, Z., and Choi, Y. Faith and Fate: Limits of Transformers on Compositionality, June 2023. URL http://arxiv. org/abs/2305.18654. arXiv:2305.18654 [cs].
- Ellis, K., Wong, C., Nye, M., Sable-Meyer, M., Cary, L., Morales, L., Hewitt, L., Solar-Lezama, A., and Tenenbaum, J. B. DreamCoder: Growing generalizable, interpretable knowledge with wake-sleep Bayesian program learning, June 2020. URL http://arxiv.org/ abs/2006.08381. arXiv:2006.08381 [cs].
- Franzen, D., Disselhoff, J., and Hartmann, D. The LLM ARChitect: Solving ARC-AGI Is A Matter of Perspective, 2024. URL https: //github.com/da-fr/arc-prize-2024/ blob/main/the_architects.pdf.
- Garcez, A. d. and Lamb, L. C. Neurosymbolic AI: the 3rd wave. *Artificial Intelligence Review*, 56(11):12387– 12406, November 2023. ISSN 1573-7462. doi: 10.1007/ s10462-023-10448-w. URL https://doi.org/10. 1007/s10462-023-10448-w.

- 495 Greenblatt, R. Getting 50% (SoTA) on ARC496 AGI with GPT-40, 2024. URL https:
 497 //redwoodresearch.substack.com/p/
 498 getting-50-sota-on-arc-agi-with-gpt.
- 499 Grill-Spector, K. and Malach, R. THE HUMAN VISUAL 500 CORTEX. Annual Review of Neuroscience, 27(Volume 501 27, 2004):649-677, July 2004. ISSN 0147-006X, 1545-502 4126. doi: 10.1146/annurev.neuro.27.070203.144220. 503 URL https://www.annualreviews.org/ 504 content/journals/10.1146/annurev. 505 neuro.27.070203.144220. Publisher: An-506 nual Reviews. 507
- 508 Hamilton, K., Nayak, A., Božić, B., and Longo, L. Is 509 neuro-symbolic AI meeting its promises in natural 510 language processing? A structured review. 511 Semantic Web. Preprint(Preprint):1–42, January 512 2022. ISSN 1570-0844. doi: 10.3233/SW-223228. 513 URL https://content.iospress.com/ 514 articles/semantic-web/sw223228. Publisher: 515 IOS Press.
- Hassija, V., Chamola, V., Mahapatra, A., Singal, A., Goel,
 D., Huang, K., Scardapane, S., Spinelli, I., Mahmud,
 M., and Hussain, A. Interpreting Black-Box Models: A
 Review on Explainable Artificial Intelligence. *Cognitive Computation*, 16(1):45–74, January 2024. ISSN 18669964. doi: 10.1007/s12559-023-10179-8. URL https:
 //doi.org/10.1007/s12559-023-10179-8.
- Hernández-Orallo, J. Twenty Years Beyond the Turing
 Test: Moving Beyond the Human Judges Too. *Minds and Machines*, 30(4):533–562, December 2020. ISSN 1572 8641. doi: 10.1007/s11023-020-09549-0. URL https:
 //doi.org/10.1007/s11023-020-09549-0.
- Hitzler, P., Eberhart, A., Ebrahimi, M., Sarker, M. K., and
 Zhou, L. Neuro-symbolic approaches in artificial intelligence. *National Science Review*, 9(6):nwac035, June
 2022. ISSN 2095-5138. doi: 10.1093/nsr/nwac035. URL
 https://doi.org/10.1093/nsr/nwac035.
- Huh, M., Cheung, B., Wang, T., and Isola, P. The Platonic
 Representation Hypothesis, July 2024. URL http://
 arxiv.org/abs/2405.07987. arXiv:2405.07987
 [cs].
- icecuber. 2020 kaggle arc-agi-1 challange: 1st place
 solution + code and official documentation, 2020. URL
 https://www.kaggle.com/competitions/
 abstraction-and-reasoning-challenge/
 discussion/154597.
- Kaddour, J., Harris, J., Mozes, M., Bradley, H., Raileanu, R., and McHardy, R. Challenges and Applications of Large Language Models, July 2023. URL http://arxiv. org/abs/2307.10169. arXiv:2307.10169 [cs].

- Kahneman, D. Thinking, fast and slow. *Farrar, Straus and Giroux*, 2011.
- Kastner, J. K. and Hong, S. J. A review of expert systems. European Journal of Operational Research, 18(3):285–292, December 1984. ISSN 0377-2217. doi: 10.1016/0377-2217(84)90150-4. URL https://www.sciencedirect.com/ science/article/pii/0377221784901504.
- Keber, M., Grubišić, I., Barešić, A., and Jović, A. A Review on Neuro-symbolic AI Improvements to Natural Language Processing. In 2024 47th MIPRO ICT and Electronics Convention (MIPRO), pp. 66– 72, May 2024. doi: 10.1109/MIPRO60963.2024. 10569741. URL https://ieeexplore.ieee. org/abstract/document/10569741. ISSN: 2623-8764.
- Kenny, E. M., Ford, C., Quinn, M., and Keane, M. T. Explaining black-box classifiers using post-hoc explanations-by-example: The effect of explanations and error-rates in xai user studies. *Artificial Intelligence*, 294: 103459, 2021.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. ImageNet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, May 2017. ISSN 0001-0782, 1557-7317. doi: 10.1145/3065386. URL https://dl.acm.org/doi/10.1145/3065386.
- Larchenko, I. 2020 Kaggle ARC-AGI-1 Challange: My part of the 3rd place solution, 2020. URL https://www.kaggle.com/competitions/ abstraction-and-reasoning-challenge/ discussion/154409.
- LeCun, Y., Boser, B., Denker, J. S., Henderson, D., Howard, R. E., Hubbard, W., and Jackel, L. D. Backpropagation Applied to Handwritten Zip Code Recognition. *Neural Computation*, 1(4):541–551, December 1989. ISSN 0899-7667. doi: 10.1162/neco. 1989.1.4.541. URL https://ieeexplore.ieee. org/abstract/document/6795724. Conference Name: Neural Computation.
- Leemann, T., Kirchhof, M., Rong, Y., Kasneci, E., and Kasneci, G. When are post-hoc conceptual explanations identifiable? In *Uncertainty in Artificial Intelligence*, pp. 1207–1218. PMLR, 2023.
- Legg, S. and Hutter, M. Universal Intelligence: A Definition of Machine Intelligence, December 2007. URL http: //arxiv.org/abs/0712.3329. arXiv:0712.3329 [cs].

- Lewis, M. and Mitchell, M. Using Counterfactual Tasks to Evaluate the Generality of Analogical Reasoning in Large Language Models, February 2024. URL http:// arxiv.org/abs/2402.08955. arXiv:2402.08955.
- Lotfi, S., Finzi, M., Kuang, Y., Rudner, T. G. J.,
 Goldblum, M., and Wilson, A. G. Non-Vacuous
 Generalization Bounds for Large Language Models,
 July 2024. URL http://arxiv.org/abs/2312.
 17173. arXiv:2312.17173.
- 560 Mahowald, K., Ivanova, A. A., Blank, I. A., Kan-561 wisher, N., Tenenbaum, J. B., and Fedorenko, E. 562 Dissociating language and thought in large language 563 Trends in Cognitive Sciences, 28(6):517models. 564 540, June 2024. ISSN 1364-6613, 1879-307X. doi: 565 10.1016/j.tics.2024.01.011. URL https://www. 566 cell.com/trends/cognitive-sciences/ 567 abstract/S1364-6613(24)00027-5. Publisher: 568 Elsevier. 569
- 570 Mira, J. M. Symbols versus connections: 50 571 years of artificial intelligence. Neurocomput-572 ing, 71(4):671-680, January 2008. ISSN 0925-573 2312. 10.1016/j.neucom.2007.06.009. doi: 574 URL https://www.sciencedirect.com/ 575 science/article/pii/S0925231207003451. 576
- Moskvichev, A., Odouard, V. V., and Mitchell, M.
 The ConceptARC Benchmark: Evaluating Understanding and Generalization in the ARC Domain, May 2023. URL http://arxiv.org/abs/2305.
 07141. arXiv:2305.07141 [cs].
- Odense, S. and Garcez, A. d. A Semantic Framework for Neural-Symbolic Computing, December 2022. URL http://arxiv.org/abs/2212.
 12050. arXiv:2212.12050 [cs].
- 587 Riesenhuber, M. and Poggio, T. Hierarchical models of
 588 object recognition in cortex. *Nature Neuroscience*, 2
 589 (11):1019–1025, November 1999. ISSN 1546-1726. doi:
 590 10.1038/14819. URL https://www.nature.com/
 591 articles/nn1199_1019. Publisher: Nature Publishing Group.
- Rong, Y., Leemann, T., Nguyen, T.-T., Fiedler, L., Qian,
 P., Unhelkar, V., Seidel, T., Kasneci, G., and Kasneci, E.
 Towards human-centered explainable ai: A survey of user
 studies for model explanations. *IEEE transactions on pattern analysis and machine intelligence*, 2023.

593

Sachdeva, N., Coleman, B., Kang, W.-C., Ni, J., Hong, L., Chi, E. H., Caverlee, J., McAuley, J., and Cheng, D. Z. How to Train Data-Efficient LLMs, February 2024. URL http://arxiv.org/abs/2402. 09668. arXiv:2402.09668 [cs].

- Schlangen, D. Targeting the Benchmark: On Methodology in Current Natural Language Processing Research. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers), pp. 670–674, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-short.85. URL https: //aclanthology.org/2021.acl-short.85.
- Schuurmans, D., Dai, H., and Zanini, F. Autoregressive Large Language Models are Computationally Universal, October 2024. URL http://arxiv.org/abs/ 2410.03170. arXiv:2410.03170.
- Skean, O., Arefin, M. R., LeCun, Y., and Shwartz-Ziv, R. Does Representation Matter? Exploring Intermediate Layers in Large Language Models, December 2024. URL http://arxiv.org/abs/2412. 09563. arXiv:2412.09563 [cs].
- Slack, D., Hilgard, A., Singh, S., and Lakkaraju, H. Reliable post hoc explanations: Modeling uncertainty in explainability. *Advances in neural information processing systems*, 34:9391–9404, 2021.
- Sutskever, I. Sequence to sequence learning with neural networks: what a decade. NeurIPS 2024, 2024. URL https://www.youtube.com/watch?v= WQQdd6qGxNs.
- Valmeekam, K., Marquez, M., and Kambhampati, S. Can Large Language Models Really Improve by Self-critiquing Their Own Plans?, October 2023. URL http://arxiv.org/abs/2310. 08118. arXiv:2310.08118 [cs].
- Wang, B., Yue, X., Su, Y., and Sun, H. Grokked Transformers are Implicit Reasoners: A Mechanistic Journey to the Edge of Generalization, May 2024. URL http:// arxiv.org/abs/2405.15071. arXiv:2405.15071 [cs].
- Webb, T., Holyoak, K. J., and Lu, H. Emergent Analogical Reasoning in Large Language Models, August 2023. URL http://arxiv.org/abs/2212. 09196. arXiv:2212.09196.
- Xu, S., Pang, L., Shen, H., Cheng, X., and Chua, T.-S. Search-in-the-Chain: Interactively Enhancing Large Language Models with Search for Knowledge-intensive Tasks, February 2024. URL http://arxiv.org/ abs/2304.14732. arXiv:2304.14732 [cs].
- Xu, Y., Khalil, E. B., and Sanner, S. Graphs, Constraints, and Search for the Abstraction and Reasoning Corpus. *Proceedings of the AAAI Conference on*

605	Artificial Intelligence, 37(4):4115–4122, June 2023a.
606	ISSN 2374-3468, 2159-5399. doi: 10.1609/aaai.v37i4.
607	25527. URL https://ojs.aaai.org/index.
608	php/AAAI/article/view/25527.
609	
610	Xu, Y., Li, W., Vaezipoor, P., Sanner, S., and Khalil, E. B.
611	LLMs and the Abstraction and Reasoning Corpus: Suc-
612	cesses, Failures, and the Importance of Object-based Rep-
613	resentations, May 2023b. URL http://arxiv.org/
614	abs/2305.18354. arXiv:2305.18354 [cs].
615	Vu 7 He I Wu 7 Dai X and Chen I Towards Better
616	Chain of Thought Prompting Strategies: A Survey Oc
617	tober 2023 LIBL http://arviv.org/abs/2310
618	0.0001 2023. OKL http://alkiv.org/abs/2310.
619	04939. arXiv.2310.04939 [C3].
620	
621	
622	
623	
624	
625	
626	
627	
628	
629	
630	
631	
632	
633	
634	
635	
636	
637	
638	
639	
640	
641	
642	
643	
044	
04J	
647	
649	
640	
650	
651	
652	
653	
654	
655	
656	
657	
658	
659	
007	