054

000

# Position: Fodor and Pylyshyn's Legacy — Still No Human-like Systematic Compositionality in Neural Networks

Anonymous Authors<sup>1</sup>

# Abstract

The strength of human language and thought lies in their ability for systematic compositionality: the meaning of a unit (semantics) can be inferred from its structure (syntax). While Fodor and Pylyshyn famously posited that neural networks inherently lack this capacity and in turn are no viable model of the human mind, Lake and Baroni recently presented meta-learning as a pathway to compositionality. In this position paper, we critically evaluate this claim, highlighting limitations in the proposed framework of metalearning for compositionality (MLC). Specifically, we identify a class of test cases compatible with Lake and Baroni's setup that consistently provoke transduction errors despite falling well within the scope of human-like abilities. We further identify overlooked yet essential elements required for substantive claims of systematic generalization. Therefore, despite the success of neural models in mimicking human behavior, it seems premature to claim that modern architectures have overcome the limitations raised by Fodor and Pylyshyn. This issue is pivotal to the AGI debate, as systematic generalization is crucial for human-like reasoning and adaptability.

# 1. Introduction

Meta-learning, or *learning to learn* from different situations, is an interesting challenge closely related to human intelligence. It is a core element of our educational system that we learn *how to learn* without explicit prior knowledge about each situation in life, as their variations are manifold. Similarly, the use of language embodies this adaptability, requiring the integration of learned rules with contextual nuances



Figure 1. "A conceptual illustration visualizing rule-learning, without any text" by DALL·E; some semantics of *text* seem to be misunderstood.

to navigate both familiar and novel scenarios. Language exemplifies how humans apply systematic generalization, seamlessly combining learned grammatical structures and vocabulary to create and interpret new expressions. This dynamic interplay between rules and context bridges the abstract principles of meta-learning with the practical mechanisms that underlie communication and cognitive reasoning.

The principle of compositionality is a key challenge for artificial neural networks, as it requires the ability to develop systematic representations and behavior. Unlike humans, neural models often struggle to generalize such rules (Nezhurina et al. 2024, Wüst et al. 2024a, Bayat et al. 2025) across contexts, reflecting fundamental gaps in their representational and operational frameworks. As artificial neural networks are constrained by their reliance on finite representational spaces and distributed encoding schemes, these limitations manifest in their difficulty in consistently applying composition rules across different scenarios. While humans can effortlessly recombine learned concepts to interpret novel sentences or solve unique problems, neural networks lack the inherent transparency, flexibility, and reflexivity to perform similar feats. Their opacity, driven by distributed representations, hinders their ability to systematically manipulate components and infer relationships.

Lake & Baroni (2023) introduced a meta-learning frame-

<sup>&</sup>lt;sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

work attempting to mitigate these challenges by introducing episodic training tasks that require rule inference. The 057 framework involves presenting neural networks with sup-058 port examples governed by hidden grammars and testing 059 their ability to generalize these rules. This episodic ap-060 proach aims to train networks for systematic generalization, using meta-learning principles to approximate human-like 061 062 reasoning. They claimed to overcome some fundamental 063 limitations of neural networks, prominently stated by Fodor 064 & Pylyshyn (1988). However, there is also plenty of evi-065 dence of the limitations of modern deep learning models 066 with human-like capabilities in language understanding that 067 rely on systematic compositional reasoning (Deletang et al. 068 2023, Zhang et al. 2023, Dziri et al. 2024, Mészáros et al. 069 2024, Bayat et al. 2025), and we provide further insights that 070 even Lake and Baroni's model fails to prove its systematic 071 behavior in several instances.

Despite its potential, the framework's reliance on learned distributions and predefined rule forms limits its scope. Gen-074 eralization remains limited to permutations of known rules 075 rather than discovering entirely new principles. The diffi-076 culty of scaling to complex tasks with deeper nesting further 077 underscores the persistent gaps in achieving true human 078 compositional reasoning. Lake and Baroni's framework 079 provides valuable insights, but also highlights the need for innovation in training and evaluating neural networks to 081 overcome these limitations, since behavioral similarities 082 may mask fundamental differences in underlying mecha-083 nisms. 084

Thus, we argue in this paper that: Neural networks havenot yet achieved learning systematic compositional skills.

We derive this position as follows:

087

088

089

090

091

092

093

094 095

096

097

098

099

100

104

105

106

109

- We locate the nature of Lake and Baroni's approach in refuting Fodor and Pylyshyn's claims that neural networks cannot reliably develop *compositional representations* and *structure-sensitive operations*.
- We show that within their setup, the model exhibits various non-systematic behaviors that can not be considered human-like.
- We argue for several necessary aspects of training and evaluation of meta-learning systems to achieve and assess their systematicity.
- We adapt Fodor and Pylyshyn's arguments in light of the modern development of deep-learning systems to argue for a future of models capable of learning symbolic representations for artificial cognition and representation learning.

# 2. The Challange of Compositionality

# 2.1. Fodor and Pylyshyn's Legacy

In their influential 1988 paper, *Connectionism and cognitive architecture: A critical analysis*<sup>1</sup>, Fodor and Pylyshyn claim that artificial neural models are unsuitable for modeling the human mind on a cognitive level. They review several arguments for the combinatorial structure of mental representations, highlighting the systematicity of these representations that follow the compositional nature of cognitive capabilities; the ability to understand some given thoughts implies the ability to understand various thoughts not only with semantically related content but also of more combinatorial complex structure. Nevertheless, they also consider the possibility that artificial neural networks may play a role in *implementing* cognition.

Differentiating neural networks and symbolic systems. Their work starts with discussing the disagreement about the nature of mental processes and mental representations between the so-called Connectionist approach, which focuses on artificial neural networks, and the Classical approach, favoring symbolic systems like Turing Machines for modeling cognitive capacities. They stress that it is neither about the explicitness of rules, nor the reality of representational states, nor about nonrepresentational architecture, as a "Connectionist neural network can perfectly well implement a Classical architecture at the cognitive level."<sup>2</sup> While both "assign semantic content to *something*"<sup>3</sup>, it is identified as the central difference that they disagree about what primitive relations hold among these content-bearing entities. The sole importance of causal connectedness in neural networks is contrasted with a range of semantic and structural relations in symbolic systems. Only the sensitivity for both semantic and structural relations is expected to allow for commitment to the compositionality of mental representations with combinatorial syntax and semantics. Furthermore, the operations the models perform in transforming one representation into another are sensitive to the structure of these representations and not only their semantics.

**Productivity, compositionality and systematicity of cognitive ability.** The need for these two properties of symbol systems, *compositional representations* and *structuresensitive operations*, is justified by "three closely related features of cognition: its productivity, its compositionality, and its inferential coherence."<sup>4</sup> Only structure-sensitive operations in combination with a combinatorial structure and semantics of representations can explain the (under appropriate idealization) *unbounded capacities* of a representational

<sup>&</sup>lt;sup>1</sup>(Fodor & Pylyshyn, 1988).

<sup>&</sup>lt;sup>2</sup>Ibid., p.11.

<sup>&</sup>lt;sup>3</sup>Ibid., p.12, emphasis in original.

<sup>&</sup>lt;sup>4</sup>Ibid., p.33.

system. Similarly, cognitive capacities are systematic in that 111 the capability for producing or processing some represen-112 tations is syntactically connected to the capability for pro-113 ducing or processing certain other representations without 114 relying on processing every specific semantics, e.g. under-115 standing the form of the expression  $A \wedge B \Rightarrow A$  implies the 116 capability to understand the expression for any substituents 117 of A or B. In fact, systematicity makes a stronger by using 118 a weaker assumption as, "[p]roductivity arguments infer 119 the internal structure of mental representations from the 120 presumed fact that nobody has a *finite* intellectual compe-121 tence [and by] contrast, systematicity arguments infer the 122 internal structure of mental representations from the patent 123 fact that nobody has a *punctuate* competence." <sup>5</sup> Closely 124 related to systematicity is the compositionality of mental 125 representations since capabilities for producing or process-126 ing representations can be linked not only syntactically but 127 also semantically. Here, it is important to note that not ev-128 ery mental representation is expected to be compositional, 129 e.g., the understanding of some expressions in natural lan-130 guage, as the "similarity of constituent structure accounts 131 for the semantic relatedness between systematically related 132 sentences only to the extent that the semantical properties of the shared constituents are context-independent."6 A last 133 134 cognitive feature is the systematicity of inference. Recalling 135 the example of conjunction  $A \wedge B$  entailing its constituent 136 A, it is not only the mental representation of understanding 137 this rule that is systematic but also its application for coher-138 ent inference between thoughts, demanding again for the 139 structure-sensitivity of operations in symbolic systems.

140 Neural networks for implementing symbol systems. Fi-141 nally, Fodor and Pylyshyn comment on treating Connection-142 ism as an implementation theory for cognitive architecture. 143 They "have no principled objection to this view"7. Still, 144 they stress that when neural networks are only a method 145 for implementing cognitive architecture, their internal states 146 are useless for understanding the nature of mental repre-147 sentations and, therefore, "irrelevant for the psychological 148 theory"8; neural networks would be just of neurological 149 means, and the need for and relevance of symbol systems 150 for modeling cognition would stay untouched. 151

#### 2.2. Lake and Baroni's Objection

Compositional seq2sec tasks. Lake and Baroni present
their work as evidence against Fodor and Pylyshyn's claims.
They introduce a meta-learning framework that they claim
achieves or exceeds human-level systematic generalization
across their evaluations. Their experimental setup is based

164

152

153

on sequence-to-sequence (sec2sec) transduction tasks, They consider sequences generated over 8 pseudolanguage tokens  $u \in U$  for the input domain  $X = U^*$ , while the output domain  $Y = C^*$  comprises sequences generated over 6 different color tokens  $c \in C$ . Both domains are linked by a transduction grammar, i.e., a set of production rules that define how a sequence of input tokens is translated into a color sequence. Each rule is of two sorts; it can state a primitive transduction rule  $u \rightarrow c$ , simply mapping an input token to an output token; otherwise, it states a unary operation  $v_1 u \to f_u(v_1)$  or binary operation  $v_1 u v_2 \to g_u(v_1, v_2)$ where any f is some n-fold (n < 8) repetition, any q is some combination of repetition, permutation and concatenation. Each  $v_i$  is either a single token  $u_i$  or the whole proceeding or succeeding token string  $x_i$ . By the iterative composition of these rules, such a grammar generates a set of translatable input sequences  $\bar{X} \subseteq X$ .

Seq2seq meta-learning framework for evaluation of human systematic generalization. With these transduction tasks, Lake and Baroni set up a meta-learning framework with EPISODES associated with different transduction grammars. Each episode combines a SUPPORT set of inputoutput transduction pairs and a OUERY set of input-output pairs, while any pair is consistent with the associated grammar. The query outputs are hidden, and it is the task to replicate them with the support examples as the only information given; the underlying transduction grammar also remains hidden. In this way, explicitly inferring the grammar rule is unnecessary. Still, the capability to implicitly extract or hypothesize the actual grammar rules is expected to be essential for reliably deriving the correct query outputs. A standard seq2seq transformer network is now trained on query examples of various episodes. The transformer encoder processes a query input combined with the support pairs of its episode as context, and the transformer decoder generates an output sequence.

# 3. Systematicity through Meta-Learning

Locating Lake and Baroni's approach. In order to evaluate the proposed framework for systematic generalization by meta-learning neural networks with respect to Fodor and Pylyshyn's claims, we will first clarify which of Fodor and Pylyshyn's arguments Lake and Baroni are referring to, since they primarily present an implementation of what they claim is a human-like systematic capability, but directly address a challenge. They themselves situate their work as a contribution to the line of argument that Fodor and Pylyshyn's statements no longer apply to current model architectures; they are not criticizing the properties of human cognition, but the alleged inability of neural networks to reliably develop *compositional representations* and *structuresensitive operations*. By focusing on behavioral tests rather

<sup>&</sup>lt;sup>5</sup>Ibid., p.40, emphasis in original.

<sup>&</sup>lt;sup>3</sup>Ibid., p.40, <sup>6</sup>Ibid., p.42.

<sup>&</sup>lt;sup>161</sup><sup>7</sup>Ibid., p.42.

<sup>162 &</sup>lt;sup>1</sup>01d., p.67. 163 <sup>8</sup>Ibid., p.65.

165 than ablation studies that directly examine the structure of learned representations, Lake and Baroni emphasize 167 the structure sensitivity and systematicity of their model, 168 which is crucial for demonstrating compositional abilities 169 and coherent behavior. Furthermore, they present their meta-170 learning framework for compositionality to systematically 171 train neural networks with these abilities. While a single 172 neural network with compositional abilities would not con-173 tradict Fodor and Pylyshyn, who did not claim any limits on 174 implementability of cognitive abilities, a framework that reli-175 ably archives compositional abilities by stochastic learning 176 methods would actually contradict their main point of crit-177 icism. Unfortunately, we will see in the following section 178 that the model trained on meta-learning still fails to reliably 179 demonstrate compositional ability in several examples. 180

# 181 **3.1. Examining the Lack of Compositionality**

182 Lake and Baroni mention that generalization beyond train-183 ing only occurs with respect to new combinations of three 184 grammar rules out of the same set of grammar rules used 185 during training. However, when we account for invariance 186 under the atomic assignments of colors to language tokens 187 and the mere labeling of operations, we find that 179/200 188 validation episodes have a combination of non-primitive 189 grammar operations that were already within the 100000 190 training episodes. So, even if the model would achieve highly systematic results on the test episodes, it could be just due to memorization of the experienced operation pat-193 terns and learning to extract the correct labels out of the episode's support examples. However, we can even show 195 that there is non-systematic behavior within their repos-196 itory of testing episodes; we reevaluate their pre-trained 197 'net - BIML - top' model on the same set of 'algebraic' test episodes with the mere difference that we did 10 evalu-199 ations of all query examples for every testing episode, for 200 statistical purposes, similar to the one episode they further 201 evaluated against human performance. We find that the 202 model performs worst on episodes #133, #32, and #122, 203 with accuracies of only 41%, 52%, and 54% on the query 204 examples, respectively. (See next paragraph and Appendix 205 7 for details.) 206

Failure in rule extraction. Further investigation of Episode 208 #133 (see Table 1) reveals that the model struggles to cor-209 rectly process the semantics of the language token  $\langle fep \rangle$ 210 with the hidden grammar rule  $x_1$  fep  $\rightarrow x_1 x_1$  and will 211 therefore refer to as (twice) and mixes it up with the to-212 ken  $\langle gazzer \rangle$  (with  $x_1 gazzer \rightarrow x_1 x_1 x_1$ ) we will call 213 (thrice). It seems to have a problem with the sole example 214 featuring  $\langle twice \rangle$ ,  $\langle \blacksquare$  thrice twice  $\rightarrow \bullet \bullet \bullet \bullet \bullet \bullet \rangle$ , which 215 also happens to contain  $\langle \text{thrice} \rangle$ . But since  $\langle \text{thrice} \rangle$ 216 has several iconic examples in the support, it is expected 217 that a reasoner with compositional skills will be able to 218 systematically use a single example and remain consistent 219



Table 1. Episode #133 with 10 evaluations for each query example; SUPPORT and QUERY are decoded for better readability. Expected outputs backed with green. The model shows *incoherent processing* and *systematically mistakes* twice for thrice. Further results can be found in Appendix A.1. (Best viewed in color.)

273

274

with the rest of the support information. By considering the examples  $\langle \blacksquare \rightarrow \bullet \rangle$ ,  $\langle \blacksquare \rightarrow \bullet \rangle$ ,  $\langle \blacksquare \text{ thrice} \rightarrow \bullet \bullet \bullet \rangle$ , human systematicity would at least suspect some semantics of  $\langle \text{twice} \rangle$  that are different to those of  $\langle \text{thrice} \rangle$ .

Non-systematic parsing. Interestingly, the hidden grammar allows for an ambiguous interpretation of nested transduction queries, which would actually be a challenge for a systematic reasoner. For instance, the query ( before twice) could be parsed as either ( before ( twice)) (marked as target by Lake and Baroni ) or  $\langle (\square \text{ before } \square) \text{ twice} \rangle$ , and similarly for a query with (thrice). But the support example  $\langle \blacksquare$  before  $\blacksquare$  thrice  $\rightarrow \bullet \bullet \bullet \bullet \rangle$  should at least induce a bias toward the intended processing. But the responses to this challenge also lack systematicity; while the frequent mistakes ( $\blacksquare$  before  $\blacksquare$  before  $\blacksquare$  twice  $\rightarrow \bullet \bullet \bullet \bullet \bullet$ ) and ( before before twice  $\rightarrow$ processing ••••• could be explained by  $\langle u_1 \text{ before } (u_2 \text{ before } (u_3 \text{ thrice})) \rangle$  while, in contrast, a similar explanation to the the error  $\langle \blacksquare$  before  $\blacksquare$  twice  $\rightarrow \bullet \bullet \bullet \bullet \bullet \bullet \rangle$  would be the parsing  $\langle (u_1 \text{ before } u_2) \text{ thrice} \rangle$ . We will further discuss the importance of systematicity for meta-learning systems in Section 3.2.

Violating structure-sensitivity. Besides both previous failure modes that are related to incompetence in extracting information from the support examples, we also found query examples for episode #1 that reveal additional non-systematicity (see Table 2 in Appendix 7 for extended version). For queries with patterns  $\langle u_1 \text{ thrice around } u_2 u_3 \rangle$  and  $\langle u_1 \text{ around } u_2 u_3 \text{ twice} \rangle$ we first see that the model never parses (around) as intended. Instead of  $\langle ((u_1 \text{ thrice}) \text{ around } u_2) u_3 \rangle$ and  $((\langle u_1 \text{ around } u_2) \ u_3)$  twice $\rangle$  the stable outputs can be explain with parsing  $\langle \texttt{around} \rangle$  as intended. Instead of  $\langle (u_1 \text{ thrice}) \text{ around } (u_2 u_3) \rangle$  and  $(\langle u_1 \text{ around } (u_2 \ u_3)) \text{ twice} \rangle$  — except for the cases,  $\langle \blacksquare$  thrice around  $\blacksquare \blacksquare \rangle$ ,  $\langle \blacksquare$  thrice around  $\blacksquare \blacksquare \rangle$ ,  $\langle \square$  around  $\square$   $\square$  twice $\rangle$ , where it would not make any difference! Only the (also ambiguous) case around twice is correctly processed in 6/10 cases - however, with even worse performance than on the unambiguous examples. Despite the structural similarities to the other query examples up to the color combination, we see a non-systematic deviation in responding that leaves compositional skills in doubt.

Limits in productivity. Finally, we want to point out that Lake and Baroni's setup only enables the model to process input sequences of up to 10 tokens and generate output sequences of up to 8 color tokens (which further restricts the admissible input sequences). This limits the possibilities for testing more complex input sequences and thus assessing



*Table 2.* Episode #1 with our own query examples and with 10 evaluations for each input; SUPPORT and QUERY are decoded for better readability. Expected outputs backed with green. Further results can be found in Appendix A.4 (Best viewed in color.)

the productivity for the model's skill.

#### 3.2. Our Position on Meta-Learning Systems

We now discuss whether meta-learning, beyond Baroni's framework, could be a promising approach towards humanlike compositional skills, despite the demonstrated limitations in the specific setup. Meta-learning systems aim to emulate human-like learning by incorporating systematic275 ity and flexibility into their architectures. These systems 276 aim to (1) generalize beyond training examples by infer-277 ring composition rules from limited examples, (2) adapt to 278 novel contexts with flexibility as a key expectation, allowing 279 systems to quickly transferring skills to new domains with 280 minimal retraining, and (3) mirror human-like cognition 281 by ensuring that error patterns and reasoning paths are still 282 systematic, explainable, or even self-correcting.

283 Weakness of non-reflective training. One of the primary 284 shortcomings in Lake and Baroni's work is the employment 285 of a one-shot prediction approach. Models are trained to 286 perform a direct transduction without intermediate reflection 287 or validation steps on the presented support examples. To 288 guarantee systematic production of outputs, we claim that 289 it is of primary importance for meta-learning models to 290 iteratively extract, validate, and correct their current belief 291 in the extracted rules. In the previous section we showed 292 that the models of Lake and Baroni fail to validate extracted 293 rules against the support and, therefore, systematically fail 294 to correctly extract (and in consequence apply), e.g., the 295 twice rule. 296

297 Focus on systematicity rather than productivity. Consid-298 ering the role that underlying grammars play with regard to 299 meta-learning or specifically non-meta-learning problems, 300 any of today's modern transformer systems can be broken 301 by providing them with more and more complex problems, 302 up to a point where the models are no longer expressive 303 enough to comprehend the problem as a whole. This could 304 be, for example, due to the depth of rule nesting or simply 305 due to the length of the input. While the general ability to 306 learn to transcribe rules certainly is a prior requirement of 307 meta-learning systems in the particular discussed setting, 308 one would not necessarily deny such systems the ability 309 to perform meta-learning reasoning even when failing to accomplish such tasks due to the aforementioned reasons. 311 When talking about meta-learning tasks, one is not so much 312 interested in the ability to derive rules of arbitrary com-313 plexity —which rather constitutes a problem of classical 314 machine learning- but in the ability of these models to sys-315 tematically discover, verify, apply, and combine said rules 316 or to systematically learn from its mistakes. When com-317 paring to human reasoning (Nezhurina et al., 2024; Wüst 318 et al., 2024b), meta-reasoning abilities are not judged by the 319 ability to produce transductions in a one-shot fashion, but 320 rather with a focus on the result being correct in the final output. We, therefore, formulate the following claim: 322

*Claim* 1. A characteristic of *successful* meta-learning systems is the ability to consistently abstain from *non-systematic* errors.

323

324

325

327

328

329

Our claim primarily regards the consistency of model behavior and we, therefore, differentiate between systematic and non-systematic errors. Systematic errors might arise due to wrong inherent assumptions of the model that, then however, get systematically applied in consequence. Such errors might stem from wrong assumptions on the general task setup. In our setting, this might involve assumptions about the unique interpretation of rules -see, e.g., our discussion on possibly ambiguous rule interpretations in Lake and Baroni-, and in general might be due to exogenous factors and implicit assumptions that where not be captured during training phase. While such errors might not produce the desired output they follow a systematicity that give rise to the assumption that the model might have been able to learn the right rules, given the correct underlying assumptions. The absence of systematicity, however, poses a much larger fault. Here, models might expose erratic 'glitches', resulting in a non-human-like behavior, that is absent of any systematicity. As the underlying reasons for such behavior might not be understood in general, it is unclear how to treat and correct such errors.

Last, we derive two positions regarding essential aspects of evaluation and training of successful meta-learning systems:

**Position on Evaluation.** Assessing and postulating systematic or compositional skills in neural networks requires either the direct evaluation of the model's internal representations, which would require an inspectable or explainable network architecture, or the use of comprehensive ablation studies that systematically testing a model's behavior in out-of-distribution situations.

Position on Implementation and Training. In order to obtain compositionality and systematicity within the discussed meta-learning tasks, the presence of symbolic representations within neural networks is vital to ensure consistent application and composition of rules. We want to emphasize that while Fodor and Pylyshyn remain unrefuted in the general analysis, today's discussion of modern neural network architectures continuously evolve to develop symbolic representations e.g. in the form of circuits (Olah et al., 2020; Wang et al., 2022; Conmy et al., 2023; Hanna et al., 2024). These explicit representations are important building blocks that promote consistent behavior and allow for explicit reflection and iterative correction of possible inconsistencies of the extracted rule sets. Last, it is important to mention that reflective behavior is likely to not evolve from training on one-shot transduction tasks, but requires models to have the possibility iterate, validate and correct over the extracted rule sets. Most recently important breakthroughs in this direction have been achieved in RL training of language reasoning models (Stiennon et al., 2020; Ouyang et al., 2022; Bai et al., 2022; Lee et al., 2023; DeepSeek-AI et al., 2025).

# 330 4. Related Work

Human-like compositionality. Regarding the importance of compositionality for cognitive skills Fodor & Lepore 333 (2001) and Fodor (2001) extend the disscusion of Fodor & 334 Pylyshyn (1988) on the compositional nature of language 335 and thought. While (natural) language incorporates some 336 non-compositional structures due to context sensitivity, com-337 positionality is argued to be essential for (a language of) 338 thought. This falls in line with more recent work of Fe-339 dorenko et al. (2024) trying to find evidence for how lan-340 guage is primarily a tool for communication rather than 341 thought. 342

343 Compositionality in neural networks. Besides Lake & 344 Baroni (2023), there is older as well as recent work trying to 345 demonstrate compositional or meta-learning skills achieved 346 with neural network architecture (Botvinick & Plaut, 2004; 347 Santoro et al., 2016; Park et al., 2024; DeepSeek-AI et al., 2025). Other work is proposing frameworks for learning and 349 assessing compositional skills (Petrache & Trivedi, 2024; 350 Sinha et al., 2024) or other intelligent behavior (Chollet, 351 2019) and Bayat et al. (2025) is introducing memorization-352 aware training to tackle overfitting to spurious correlation 353 encountered in training. 354

Limitations in systematicity. Several works evaluate and 355 demonstrate the limitations of modern AI models in com-356 positional or systematic generalization tasks (Bender et al., 357 2021; Deletang et al., 2023; Dziri et al., 2024; Mészáros 358 et al., 2024; Nezhurina et al., 2024; Zhang et al., 2024; 359 Wüst et al., 2024b) and there is also another targeted re-360 sponse to Lake and Baroni's work, presenting problems of 361 non-systematic behovior (Goodale & Mascarenhas, 2023). 362

363 Importance of symbolics. There is also more recent work 364 that stresses the importance of symbolics. Ellis et al. (2020) 365 introduces a machine learning system that utilisez neural 366 guided program synthesis to learns to solve problems. Wüst 367 et al. (2024a) furhter demonstrates the advantages of using 368 program synthesis for unsupervised learning of complex, 369 relational concepts from images, focusing on the benefits 370 in terms of generalization, interpretability, and revisability. 371 Stammer et al. (2024b), on the other hand investigated the 372 benefits of symbolic representations for improved gener-373 alization and interpretability of low-level visual concepts. 374 The position of the importance of symbols for AI expla-375 nations is further discussed by Kambhampati et al. (2022). 376 The approach of Dinu et al. (2024) combines generative 377 models and solvers by use of large language models as se-378 mantic parsers. Shindo et al. (2025) model human ability to 379 combine symbolic reasoning with intuitive reactions by a 380 neuro-symbolic reinforcement learning framework.

- 381
- 382
- 383 384

# **5.** Alternative Views

Historically, Fodor & Pylyshyn (Fodor & Pylyshyn, 1988) argued for the emergence or implementation of symbolically reasoning structures within neural networks as a necessary aspect for achieving human-like meta-learning. However, the considerations for meta-learning discussed in their and our paper strongly focus on the learning of logical and arithmetic rules where concepts can be reduced onto symbolic expressions. These representations, therefore, naturally align well with the abilities of symbolic reasoners, but leave out other possible forms of meta-learning systems. Considering different modalities, for example for composing visual patterns or motion sequences, might pose a strong hurdle for classical symbolic systems. Such domains that do not operate over discrete 'crystallized' symbols, but rather operate on abstract 'fluid' concepts, still lack a well defined notion of what constitutes meta-learning within them. As a consequence it is unclear how to measure and systematically asses the abilities of models with regard to meta-learning in possible benchmarks.

Untargeted Emergence of Systematic Reasoning. Even without targeted training towards meta-learning models, LLM have shown to exhibit emergent abilities for various tasks (Brown et al., 2020; Wei et al., 2022a; Schaeffer et al., 2024). While 'true' understanding of the world might only be achieved via (embodied) interaction (Lipson & Pollack, 2000; Gupta et al., 2021; Zečević et al., 2023), some works have argued that such abilities might even be learned through mere passive observation (Lampinen et al., 2024), while other approaches argue for the value of selfexplanation guided learning (Stammer et al., 2024a). Considering the underlying aspect of systematic learning and reasoning, several works already where able to distill symbolically acting *circuits* that emerged during training from LLM (Olah et al., 2020; Wang et al., 2022; Conmy et al., 2023; Hanna et al., 2024). In light of these results, it stands yet to to be seen whether meta-learning abilities of language reasoning models might also emerge as a consequence of pure scaling laws (Sutton, 2019; Kaplan et al., 2020; Bubeck et al., 2023).

# 6. Discussion and Conclusion

For this final section, we will revisit the key points that constitute our position (c.f. Sec. 1) and that we believe to form important aspects towards the goal of achieving meta-learning models capable of performing human-like systematic compositionality:

(I) Criteria for Systematic Compositionality. The main criteria for models with productive, systematic and compositional skills remain compositional representation and structure-sentitive operations.

(II) Non-systematic Behavior. As Fodor and Pylyshyn's
trained model exhibits various non-systematic behaviors, it
failed to demonstrate human-like compositional learning
capacities and, furthermore, refutes the presented claims
that their meta-learning framework is achieving human-like
systematic generalization.

(III) Assessing Compositionality. Systematic testing of
 several types of out-of-distribution episodes is necessary to
 assess compositional skills.

(IV) Emergence and Learning of Symbolic Representations. Meta-learning systems have to encourage the emergence of symbolic representations during training. For that,
we expect training tasks and model architecture that makes
iteration, self-validation, and self-correction over the extracted rule sets possible as well as necessary.

Before concluding we now summarize the key considerations required for achieving truly meta-learning systems.

404 Aspects of meta-learning. The limitations of current neural 405 models emphasize the importance of hybrid architectures 406 that integrate the strengths of symbolic and connectionist 407 paradigms. Neuro-symbolic models offer a compelling solu-408 tion, combining explicit rule representation, error correction 409 mechanisms, and dynamic scalability. Key advancements in 410 this direction include, (1) Systematicity: Embedding mech-411 anisms for representing and manipulating compositional 412 rules within neural architectures; (2) Reflective Reasoning: 413 Incorporating iterative self-correction processes to emulate 414 human-like adaptability; (3) Scalability: Enabling models 415 to dynamically expand rule sets and adapt to novel tasks, 416 mirroring human flexibility; and we will discuss those in 417 the following:

418
419
420
420
421
421
421
422
421
421
423
424
425
426
427
427
427
428
429
429
429
420
420
420
421
420
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421
421

422 Reflection and iterative refinement. A distinct ability of 423 human reasoning is the ability to reflect and develop a set of 424 currently hypothesized rules. The extraction and validation 425 of rules from provided support examples might pose a com-426 plex task which can scale with exponential complexity with 427 number of provided support examples. Models exhibiting 428 one-shot behavior might be able to perform such tasks up to 429 a particular problem size, but are ultimately limited by their 430 own model capacity. In this paper we, therefore, provide 431 arguments towards the use of reflective learners, as a par-432 ticular class of models, capable of iterative refinement and 433 self-correction of their current beliefs, a practice already 434 adopted with great success for general language reasoning 435 models (Wei et al., 2022b; Stammer et al., 2024a; Yao et al., 436 2024; DeepSeek-AI et al., 2025). This iterative behavior 437 allows for repeated validation of the conjectures rules and, 438 therefore, fundamentally stands in contrast to models trained

439

to provide answers in a one-shot fashion.

**Scalability, memory and context.** An often overlooked part of learning to reflecting upon ones beliefs is the requirement of learning to store and operate on suitable representations of the models' beliefs. Particularly, this includes the presence of some sort of memory that can be read and updated. Upon the application of rules to a given query a model might additionally want to track its current context (e.g. the nesting depth of current rules), which, again, might require some sort of memory to generalize to arbitrary problem sizes and overcome the limitations of a static number model parameters.

**Conclusion.** Models with the discussed properties have the potential to address foundational critiques of connectionism while advancing the capabilities of artificial cognition. By bridging the gap between symbolic and connectionist principles, hybrid architectures could achieve systematicity and productivity, paving the way for truly human-like reasoning.

The enduring relevance of Fodor and Pylyshyn's critique underscores the challenges in developing systems capable of systematic generalization and compositional reasoning. While meta-learning frameworks represent significant progress, they fall short of resolving foundational limitations. Future advancements must embrace integrative approaches that merge the strengths of symbolic and connectionist paradigms, paving the way for a more robust understanding of artificial cognition. By addressing these challenges, we can move closer to realizing the vision of human-like artificial intelligence.

# 7. Impact Statement.

Strong meta-learning abilities show as an important skill to navigate the complex and changing tasks of today's world. When presenting models that aim to robustly adapt to novel environment, it is important to refrain from making unsolidified claims about the achievement of meta-learning systems which ultimately do not hold true upon closer inspection. Hiding behind details of what technically constitutes as meta-learning systems does not help the general discussion, but might enhance the public trust in such models, possibly leading into a false reliance in them. Our analysis showed that modern neural meta-learning systems can only archive such tasks, if at all, only under a very narrow and restricted definition of a meta-learning setup. In this paper we promote the systematic evaluation of meta-learning systems beyond their training distribution in order to truthfully assess their ability of performing compositional reasoning. We furthermore. As a result, we claim that 'Fodor and Pylyshyn's Legacy' persists and we conclude that there is still no human-like systematic compositionality learned in neural networks as of today.

# 440 **References**

458

- Bai, Y., Kadavath, S., Kundu, S., Askell, A., Kernion, J., Jones, A., Chen, A., Goldie, A., Mirhoseini, A., McKinnon, C., et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.
- Bayat, R., Pezeshki, M., Dohmatob, E., Lopez-Paz, D.,
  and Vincent, P. The pitfalls of memorization: When
  memorization hurts generalization. In *The Thirteenth International Conference on Learning Representations*,
  2025. URL https://openreview.net/forum?
  id=vVhZh9ZpIM.
- Bender, E. M., Gebru, T., McMillan-Major, A., and
  Shmitchell, S. On the dangers of stochastic parrots: Can
  language models be too big? In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, pp. 610–623, 2021.
- Botvinick, M. and Plaut, D. C. Doing without schema hierarchies: a recurrent connectionist approach to normal and impaired routine sequential action. *Psychological review*, 111(2):395–429, 2004. doi: 10.1037/0033-295X.111.
  2.395. URL https://pubmed.ncbi.nlm.nih.gov/15065915/.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D.,
  Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G.,
  Askell, A., et al. Language models are few-shot learners.
  Advances in neural information processing systems, 33:
  1877–1901, 2020.
- Bubeck, S., Chandrasekaran, V., Eldan, R., Gehrke, J.,
  Horvitz, E., Kamar, E., Lee, P., Lee, Y. T., Li, Y.,
  Lundberg, S., et al. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- 477 Chollet, F. On the measure of intelligence, 2019. URL
  478 https://arxiv.org/abs/1911.01547.
- 479
  480
  481
  481
  482
  482
  483
  483
  484
  484
  484
  485
  486
  487
  487
  488
  488
  489
  489
  480
  480
  480
  481
  481
  481
  482
  483
  484
  484
  484
  485
  484
  485
  484
  485
  484
  485
  484
  485
  486
  487
  487
  487
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
  488
- 485 DeepSeek-AI, Guo, D., Yang, D., Zhang, H., Song, J., 486 Zhang, R., Xu, R., Zhu, Q., Ma, S., Wang, P., Bi, X., 487 Zhang, X., Yu, X., Wu, Y., Wu, Z. F., Gou, Z., Shao, 488 Z., Li, Z., Gao, Z., Liu, A., Xue, B., Wang, B., Wu, B., 489 Feng, B., Lu, C., Zhao, C., Deng, C., Zhang, C., Ruan, 490 C., Dai, D., Chen, D., Ji, D., Li, E., Lin, F., Dai, F., Luo, 491 F., Hao, G., Chen, G., Li, G., Zhang, H., Bao, H., Xu, 492 H., Wang, H., Ding, H., Xin, H., Gao, H., Qu, H., Li, 493 H., Guo, J., Li, J., Wang, J., Chen, J., Yuan, J., Qiu, J., 494

Li, J., Cai, J. L., Ni, J., Liang, J., Chen, J., Dong, K., Hu, K., Gao, K., Guan, K., Huang, K., Yu, K., Wang, L., Zhang, L., Zhao, L., Wang, L., Zhang, L., Xu, L., Xia, L., Zhang, M., Zhang, M., Tang, M., Li, M., Wang, M., Li, M., Tian, N., Huang, P., Zhang, P., Wang, Q., Chen, Q., Du, Q., Ge, R., Zhang, R., Pan, R., Wang, R., Chen, R. J., Jin, R. L., Chen, R., Lu, S., Zhou, S., Chen, S., Ye, S., Wang, S., Yu, S., Zhou, S., Pan, S., Li, S. S., Zhou, S., Wu, S., Ye, S., Yun, T., Pei, T., Sun, T., Wang, T., Zeng, W., Zhao, W., Liu, W., Liang, W., Gao, W., Yu, W., Zhang, W., Xiao, W. L., An, W., Liu, X., Wang, X., Chen, X., Nie, X., Cheng, X., Liu, X., Xie, X., Liu, X., Yang, X., Li, X., Su, X., Lin, X., Li, X. Q., Jin, X., Shen, X., Chen, X., Sun, X., Wang, X., Song, X., Zhou, X., Wang, X., Shan, X., Li, Y. K., Wang, Y. Q., Wei, Y. X., Zhang, Y., Xu, Y., Li, Y., Zhao, Y., Sun, Y., Wang, Y., Yu, Y., Zhang, Y., Shi, Y., Xiong, Y., He, Y., Piao, Y., Wang, Y., Tan, Y., Ma, Y., Liu, Y., Guo, Y., Ou, Y., Wang, Y., Gong, Y., Zou, Y., He, Y., Xiong, Y., Luo, Y., You, Y., Liu, Y., Zhou, Y., Zhu, Y. X., Xu, Y., Huang, Y., Li, Y., Zheng, Y., Zhu, Y., Ma, Y., Tang, Y., Zha, Y., Yan, Y., Ren, Z. Z., Ren, Z., Sha, Z., Fu, Z., Xu, Z., Xie, Z., Zhang, Z., Hao, Z., Ma, Z., Yan, Z., Wu, Z., Gu, Z., Zhu, Z., Liu, Z., Li, Z., Xie, Z., Song, Z., Pan, Z., Huang, Z., Xu, Z., Zhang, Z., and Zhang, Z. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. URL https://arxiv.org/abs/2501.12948.

- Deletang, G., Ruoss, A., Grau-Moya, J., Genewein, T., Wenliang, L. K., Catt, E., Cundy, C., Hutter, M., Legg, S., Veness, J., and Ortega, P. A. Neural networks and the chomsky hierarchy. In *The Eleventh International Conference on Learning Representations*, 2023. URL https: //openreview.net/forum?id=WbxHAzkeQcn.
- Dinu, M.-C., Leoveanu-Condrei, C., Holzleitner, M., Zellinger, W., and Hochreiter, S. Symbolicai: A framework for logic-based approaches combining generative models and solvers, 2024. URL https://arxiv. org/abs/2402.00854.
- Dziri, N., Lu, X., Sclar, M., Li, X. L., Jiang, L., Lin, B. Y., Welleck, S., West, P., Bhagavatula, C., Le Bras, R., et al. Faith and fate: Limits of transformers on compositionality. *Advances in Neural Information Processing Systems*, 36, 2024.
- 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, 2020. URL https://arxiv.org/abs/ 2006.08381.
- Fedorenko, E., Piantadosi, S. T., and Gibson, E. A. F. Language is primarily a tool for commu-

495 nication rather than thought. Nature, 630:575–
496 586, 2024. doi: 10.1038/s41586-024-07522-w.
497 URL https://www.nature.com/articles/
498 s41586-024-07522-w#citeas.

- Fodor, J. and Lepore, E. Brandom's burdens: Compositionality and inferentialism. *Philosophy and Phenomenological Research, Vol. 63, No. 2, 465-481, 2001.* URL https://www.jstor.org/stable/3071079.
- Fodor, J. A. Language, thought and compositionality. *Mind & Language*, 16(1):1–15, 2001. doi: https://doi.org/10.1111/1468-0017.00153. URL
   https://onlinelibrary.wiley.com/doi/ abs/10.1111/1468-0017.00153.
- Fodor, J. A. and Pylyshyn, Z. W. Connectionism and cognitive architecture: a critical analysis. *Cognition 28, 3–71* (1988), 1988.
- Goodale, M. and Mascarenhas, S. Fodor and pylyshyn's systematicity challenge still stands: A reply to lake and baroni (2023), 2023. URL https://lingbuzz.net/lingbuzz/007759.

518

519

520

521

526

- Gupta, A., Savarese, S., Ganguli, S., and Fei-Fei, L. Embodied intelligence via learning and evolution. *Nature communications*, 12(1):5721, 2021.
- Hanna, M., Liu, O., and Variengien, A. How does gpt-2
  compute greater-than?: Interpreting mathematical abilities in a pre-trained language model. *Advances in Neural Information Processing Systems*, 36, 2024.
- Kambhampati, S., Sreedharan, S., Verma, M., Zha, Y., and
  Guan, L. Symbols as a lingua franca for bridging humanai chasm for explainable and advisable AI systems. In *Conference on Artificial Intelligence, (AAAI)*, pp. 12262– 12267. AAAI Press, 2022.
- Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Child, R., Gray, S., Radford, A., Wu, J., and Amodei, D. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- Lake, B. M. and Baroni, M. Human-like systematic generalization through a meta-learning neural network. *Nature* 623, 115–121, 2023.
- Lampinen, A., Chan, S., Dasgupta, I., Nam, A., and Wang,
  J. Passive learning of active causal strategies in agents
  and language models. *Advances in Neural Information Processing Systems*, 36, 2024.
- Lee, H., Phatale, S., Mansoor, H., Lu, K. R., Mesnard, T.,
  Ferret, J., Bishop, C., Hall, E., Carbune, V., and Rastogi,
  A. Rlaif: Scaling reinforcement learning from human feedback with ai feedback. 2023.

- Lipson, H. and Pollack, J. B. Automatic design and manufacture of robotic lifeforms. *Nature*, 406(6799):974–978, 2000.
- Mészáros, A., Ujváry, S., Brendel, W., Reizinger, P., and Huszár, F. Rule extrapolation in language modeling: A study of compositional generalization on OOD prompts. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https: //openreview.net/forum?id=Li2rpRZWjy.
- Nezhurina, M., Cipolina-Kun, L., Cherti, M., and Jitsev, J. Alice in wonderland: Simple tasks showing complete reasoning breakdown in state-of-the-art large language models, 2024. URL https://arxiv.org/abs/2406. 02061.
- Olah, C., Cammarata, N., Schubert, L., Goh, G., Petrov, M., and Carter, S. Zoom in: An introduction to circuits. *Distill*, 5(3):e00024–001, 2020.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al. Training language models to follow instructions with human feedback. *Advances in neural information* processing systems, 35:27730–27744, 2022.
- Park, C. F., Okawa, M., Lee, A., Lubana, E. S., and Tanaka, H. Emergence of hidden capabilities: Exploring learning dynamics in concept space. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum? id=owuEcT6BT1.
- Petrache, M. and Trivedi, S. Position paper: Generalized grammar rules and structure-based generalization beyond classical equivariance for lexical tasks and transduction, 2024. URL https://arxiv.org/abs/2402. 01629.
- Santoro, A., Bartunov, S., Botvinick, M., Wierstra, D., and Lillicrap, T. Meta-learning with memory-augmented neural networks. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48*, ICML'16, pp. 1842–1850. JMLR.org, 2016.
- Schaeffer, R., Miranda, B., and Koyejo, S. Are emergent abilities of large language models a mirage? *Advances in Neural Information Processing Systems*, 36, 2024.
- Shindo, H., Delfosse, Q., Dhami, D. S., and Kersting, K. Blendrl: A framework for merging symbolic and neural policy learning. In *Proceedings of the International Conference on Representation Learning (ICLR)*, 2025.

- Sinha, S., Premsri, T., and Kordjamshidi, P. A survey on compositional learning of AI models: Theoretical and experimental practices. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL https:// openreview.net/forum?id=BXDxwItNqQ. Survey Certification.
- Stammer, W., Friedrich, F., Steinmann, D., Brack, M.,
  Shindo, H., and Kersting, K. Learning by self-explaining. *Transactions on Machine Learning Research*, 2024a.
  ISSN 2835-8856. URL https://openreview.
  net/forum?id=bpjU7rLjJ7.
- Stammer, W., Wüst, A., Steinmann, D., and Kersting, K.
  Neural concept binder. *Advances in Neural Information Processing Systems*, 2024b.
- Stiennon, N., Ouyang, L., Wu, J., Ziegler, D., Lowe, R.,
  Voss, C., Radford, A., Amodei, D., and Christiano,
  P. F. Learning to summarize with human feedback. *Advances in Neural Information Processing Systems*, 33: 3008–3021, 2020.
- 572 Sutton, R. The bitter lesson. *Incomplete Ideas (blog)*, 13(1):573 38, 2019.
- Wang, K., Variengien, A., Conmy, A., Shlegeris, B., and
  Steinhardt, J. Interpretability in the wild: a circuit for
  indirect object identification in gpt-2 small. *arXiv preprint arXiv:2211.00593*, 2022.
- Wei, J., Tay, Y., Bommasani, R., Raffel, C., Zoph, B.,
  Borgeaud, S., Yogatama, D., Bosma, M., Zhou, D., Metzler, D., et al. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*, 2022a.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi,
  E., Le, Q. V., Zhou, D., et al. Chain-of-thought prompting
  elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837,
  2022b.
- Wüst, A., Stammer, W., Delfosse, Q., Dhami, D. S., and Kersting, K. Pix2code: Learning to compose neural visual concepts as programs. In *The 40th Conference on Uncertainty in Artificial Intelligence*, 2024a. URL https://openreview.net/forum?id=EE4ikEQnOT.
- Wüst, A., Tobiasch, T., Helff, L., Dhami, D. S., Rothkopf,
  C. A., and Kersting, K. Bongard in wonderland: Visual puzzles that still make ai go mad? *arXiv preprint arXiv:2410.19546*, 2024b.
- Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T., Cao, Y., and Narasimhan, K. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36, 2024.

- Zečević, M., Willig, M., Dhami, D. S., and Kersting, K. Causal parrots: Large language models may talk causality but are not causal. *Transactions on Machine Learning Research*, 2023.
- Zhang, D., Tigges, C., Zhang, Z., Biderman, S., Raginsky, M., and Ringer, T. Transformer-based models are not yet perfect at learning to emulate structural recursion. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL https:// openreview.net/forum?id=Ry5CXXmlsf.
- Zhang, S. D., Tigges, C., Biderman, S., Raginsky, M., and Ringer, T. Can transformers learn to solve problems recursively?, 2023. URL https://arxiv.org/abs/ 2305.14699.

# A. APPENDIX: Position: Fodor and Pylyshyn's Legacy — Still No Human-like Systematic Compositionality in Neural Networks

This appendix contains the full set of grammar rules, support examples and query examples of Lake and Baroni's meta-learning. We reevaluated their pre-trained 'net - BIML - top' model on the same set of 'algebraic' testing-episodes. Here, we reported the outputs for #133, #132, #122, and modified #1.



*Table 3.* GRAMMAR and SUPPORT for Episode #133; decodedfor better readability. (Best viewed in color.)



Table 4. Episode #133 with 10 evaluations for each query example; decoded for better readability. Expected outputs backed with green. The model shows *incoherent processing* and *systematically mistakes* twice for thrice. (Best viewed in color.)

OUT	COU 8 1 1 8
•	8 1 1 8
•	$\frac{1}{1}$
•	1 8
••	8
••	0
	2
	5
••••	3
	1
••••	1
••••	8
•	1
	1
	4
	4
••••	1
•••••	1
••	2
•••••	1
•••••	1
••••	1
	1
	1
••••	1
	1
•••	1
	_
••	3
••••	2
	1
	1
	1
••••	1
	1
	_
	9
	1
	4
	+ 9
	1
	1
	1
•••••	1
••••••	
•••••	10

3. Complete responses for Lake and Baroni's QUERY (Lake		nd Baroni)	
meta-learning testing-episode #122.	IN	OUT	C
	four times	••••	
GRAMMAR #122 (Lake and Baroni)		••	
blicket $\rightarrow \bullet$ , kiki $\rightarrow \bullet$ , zup $\rightarrow \bullet$ , lug $\rightarrow \bullet$ ,		•••	
$x_1 \operatorname{dax} \to x_1 x_1 x_1 x_1 x_1,$		••••••	
$x_1 \text{ fep } x_1 \to x_1 \ u_1 \ u_1 \ x_1,$		••••••	
$x_1 \text{ gazzer} \to x_1 x_1$		••••	
DECODING (this paper; for readability)		•••••	
blicket : $\blacksquare$ , kiki : $\blacksquare$ , zup : $\blacksquare$ , lug : $\blacksquare$ ,		•••••	
$\operatorname{dax}: \texttt{four times},$	📕 📕 twice within 📒 🗖	•••••	
$\mathrm{fep}: \texttt{twice within},$		•••	
gazzer:twice		•••••	
SUPPORT (Lake and Baroni)		••••••	
$\blacksquare \to \bullet, \blacksquare \to \bullet, \blacksquare \to \bullet,$		••••••	
$\blacksquare \ - \bullet \bullet \bullet,$		•••••	
$\blacksquare \to \bullet \bullet,$		••	
$\blacksquare$ twice $ o$ $ulletullet$ ,		•	
$\blacksquare$ twice $ o$ $ulletullet$ ,		••••	
four times $\rightarrow \bullet \bullet \bullet \bullet$ ,		•••••	
four times $\rightarrow \bullet \bullet \bullet \bullet$ ,		•••••	
<b>t</b> wice $\rightarrow \bullet \bullet \bullet \bullet$ ,		•••••	
<b>t</b> wice $\rightarrow \bullet \bullet \bullet \bullet$ ,	<b>.</b>	•••••	
<b>E E I</b> twice $\rightarrow \bullet \bullet \bullet \bullet \bullet \bullet \bullet \bullet \bullet$ ,	twice within twice	•••••	
<b>t</b> wice within <b>t</b> wice $\rightarrow \bullet \bullet \bullet \bullet \bullet \bullet$ ,		••••	
		•••••	
		•••••	
<i>Fable 7.</i> GRAMMAR and SUPPORT for Episode #122; decoded			
or better readability. (Best viewed in color.)		•••	
		••••	
		•••	
		••	
	twice within	•••	
		••••	
		•••	
		••••	
		••••	
		••	
		•	
	-	•	
		••••	
	four times	•••	
		••	
	<pre>twice</pre>		
	four times	••••	
		•••	
		•••	

770	A.4. Complete responses for our modified version of				
772	Lake and Darom's testing-episode #1.				
773	GRAMMAR #1 (Lake and Baroni)				
774	tufa $\rightarrow \bullet$ , wif $\rightarrow \bullet$ , lug $\rightarrow \bullet$ , fep $\rightarrow \bullet$ .				
775	$u_1 \text{ gazzer} \rightarrow x_1 x_1 x_1.$				
776	$x_1$ kiki $u_1 \rightarrow x_1 u_1 x_1$ ,				
777	$x_1 \operatorname{zup} \to x_1 x_1$				
778	DECODING (this paper; for readability)				
779	tufa : $\blacksquare$ , wif : $\blacksquare$ , lug : $\blacksquare$ , fep : $\blacksquare$ ,				
780	gazzer: thrice,				
781	kiki : around,				
782	zup:twice				
783	SUPPORT (Lake and Baroni)				
784	$\blacksquare \rightarrow \bullet, \blacksquare \rightarrow \bullet,$				
785	$\blacksquare$ $\rightarrow$ ••,				
786	$\blacksquare$ twice $ o$ $ullet$ ,				
787	<b>thrice</b> $\rightarrow \bullet \bullet \bullet$ ,				
788	<b>E L</b> twice $\rightarrow \bullet \bullet \bullet \bullet$ ,				
789	<b>thrice</b> $\rightarrow \bullet \bullet \bullet \bullet \bullet \bullet$ ,				
790	<b>The second sec</b>				
791	$\blacksquare \texttt{thrice} \to \bullet \bullet \bullet \bullet \bullet \bullet,$				
792	$\blacksquare$ around $\blacksquare \rightarrow \bullet \bullet \bullet$ ,				
793	$\blacksquare$ around $\blacksquare \rightarrow \bullet \bullet \bullet \bullet$ ,				
794	<b>a</b> round <b>b</b> around <b>b</b> $\rightarrow$ ••••••••,				
795	$\blacksquare$ around $\blacksquare$ twice $\rightarrow \bullet \bullet \bullet \bullet \bullet \bullet$ ,				
796	$\blacksquare$ thrice around $\blacksquare  o \bullet \bullet \bullet \bullet \bullet \bullet \bullet$				
797					
798	Table 9. GRAMMAR and SUPPORT for Episode #1; decoded for				
799	bener readability. (Dest viewed in color.)				
800					

QUERY (new; this paper)							
IN	OUT	COUNT					
thrice around	•••••	10					
	•••••						
	••••	8					
+hrico pround	••••	1					
	••	1					
	•••••	_					
	•••••	8					
thrice around	•••••	1					
	••••	1					
	•••••	_					
	••••	9					
🗖 thrice around 📒 📕	••••	1					
	••••	_					
	•••••	8					
thrico pround	••••	1					
	••••	1					
	•••••	_					
	•••••	9					
thrice around	••	1					
	•••••	_					
	•••••	8					
around twice		1					
	••••	1					
	•••••						
	•••••	8					
around twice	•••••	1					
	••••	1					
	•••••						
	•••••	6					
	••••	1					
🗖 around 📕 🗖 twice	•••••	1					
	••••	1					
	•••••	1					
	•••••	7					
	••••	1					
🗖 around 🗖 📕 twice	•••••	1					
	••••	1					
	•••••						
	•••••	5					
	•••••	2					
around	•••••	1					
	••••	1					
	••••	1					

Table 10. Our own query examples for Episode #1 with 10 evaluations each; decoded for better readability. Expected outputs backed with green. (Best viewed in color.)