

Aligning Generalisation Between Humans and Machines

Filip Ilievski^{1*}, Barbara Hammer², Frank van Harmelen¹,
Benjamin Paassen², Sascha Saralajew³, Ute Schmid⁴,
Michael Biehl⁵, Marianna Bolognesi⁶, Xin Luna Dong⁷,
Kiril Gashteovski^{3,8}, Pascal Hitzler⁹, Giuseppe Marra¹⁰,
Pasquale Minervini^{11,12}, Martin Mundt¹³,
Axel-Cyrille Ngonga Ngomo¹⁴, Alessandro Oltramari¹⁵,
Gabriella Pasi¹⁶, Zeynep G. Saribatur¹⁷, Luciano Serafini¹⁸,
John Shawe-Taylor¹⁹, Vered Shwartz^{20,21}, Gabriella Skitalinskaya²²,
Clemens Stachl²³, Gido M. van de Ven¹⁰, Thomas Villmann^{24,25}

¹Vrije Universiteit Amsterdam, Amsterdam, The Netherlands.

²University of Bielefeld, Bielefeld, Germany.

³NEC Laboratories Europe, Heidelberg, Germany.

⁴University of Bamberg, Bamberg, Germany.

⁵University of Groningen, Groningen, The Netherlands.

⁶Università di Bologna, Bologna, Italy.

⁷Meta Reality Labs, Redmond, WA, USA.

⁸CAIR, Ss. Cyril and Methodius University, Skopje, North Macedonia.

⁹Kansas State University, Manhattan, KS, USA.

¹⁰KU Leuven, Leuven, Belgium.

¹¹University of Edinburgh, Edinburgh, UK.

¹²Miniml.AI, Edinburgh, UK.

¹³University of Bremen, Bremen, Germany.

¹⁴Paderborn University, Paderborn, Germany.

¹⁵Carnegie Bosch Institute, Pittsburgh, PA, USA.

¹⁶Università degli Studi di Milano Bicocca, Milan, Italy.

¹⁷TU Wien, Vienna, Austria.

¹⁸Fondazione Bruno Kessler, Trento, Italy.

¹⁹University College London, London, UK.

²⁰University of British Columbia, Vancouver, Canada.

²¹Vector Institute, Toronto, Canada.

²²Duolingo, Pittsburgh, PA, USA.

²³University of St. Gallen, Institute of Behavioral Science and Technology, St. Gallen, Switzerland.

²⁴University of Applied Sciences Mittweida, Mittweida, Germany.

²⁵Technical University Freiberg, Freiberg, Germany.

*Corresponding author(s). E-mail(s): f.ilievski@vu.nl;

Abstract

Recent advances in AI—including generative approaches—have resulted in technology that can support humans in scientific discovery and forming decisions, but may also disrupt democracies and target individuals. The responsible use of AI and its participation in human-AI teams increasingly shows the need for *AI alignment*, that is, to make AI systems act according to our preferences. A crucial yet often overlooked aspect of these interactions is the different ways in which humans and machines *generalise*. In cognitive science, human generalisation commonly involves abstraction and concept learning. In contrast, AI generalisation encompasses out-of-domain generalisation in machine learning, rule-based reasoning in symbolic AI, and abstraction in neurosymbolic AI. In this perspective paper, we combine insights from AI and cognitive science to identify key commonalities and differences across three dimensions: notions of, methods for, and evaluation of generalisation. We map the different conceptualisations of generalisation in AI and cognitive science along these three dimensions and consider their role for alignment in human-AI teaming. This results in interdisciplinary challenges across AI and cognitive science that must be tackled to support effective and cognitively supported alignment in human-AI teaming scenarios.

Keywords: generalisation, human-AI teaming, alignment

Contents

1	Introduction	4
2	Parallels in Generalisation by Humans and Machines	5
3	Notions of Generalisation	8
3.1	Generalisation as a process	8
3.2	Generalisation as a product	8
3.3	Generalisation as an operator	8
3.4	Alignment of human and machine notions of generalisation	9
4	Methods	9
4.1	Statistical generalisation methods in AI	10
4.2	Knowledge-informed generalisation methods in AI	11
4.3	Instance-based translation in AI	12
4.4	Aligning machine generalisation methods and human expectations	12
5	Evaluation	13
5.1	Measuring distributional shifts	14
5.2	Determining under- and overgeneralisation	14
5.3	Distinguishing memorisation and generalisation	15
5.4	Alignment of machine evaluation of generalisation to humans	15
6	Emerging Directions	16

1 Introduction

Recent advances in artificial intelligence (AI) enable meaningful support for humans in complex tasks, such as scientific discovery and decision-making [71]. Conversely, AI can also potentially disrupt democracies and target individuals [39]. The responsible use of AI increasingly motivates *AI alignment*, which aims to “make AI systems act according to our preferences” [21]. AI alignment is essential for effective human-AI teaming in complex scenarios where neither humans nor AI perform well on their own [99]. For example, AI can help invent novel biomedical application hypotheses following the objectives and guidance of a scientist [51]. Alternatively, humans can iteratively edit samples provided by an AI model to improve its accuracy and trustworthiness, for example, in classifying skin cancer [35]. Besides human-AI teaming, AI alignment is necessary for its safe use [8] and demonstrable adherence to accountability, privacy, and transparency requirements in legal frameworks such as the EU’s AI Act [6].

A crucial, yet often overlooked aspect of this alignment in interactive scenarios is the complementary ways in which humans and machines *generalise* (Figure 1). Generalisation is typically defined as “the process of transferring knowledge or skills from specific instances or exemplars to new contexts” [123]. Following cognitive science, human generalisation commonly involves concept learning and the abstraction of general characteristics to a collection of entities [56]. Humans excel at generalising from a few examples, compositionality, and robust generalisation to noise, shifts, and Out-Of-Distribution (OOD) data [60]. Humans can learn from little data and seemingly generalise beyond the observed distribution largely because, through evolution, experience, or both, they have access to strong causality-driven *common sense* priors at multiple hierarchical levels that characterise physical principles in nature and human behaviour in interactions.

In sharp contrast, data-driven (*statistical*) AI systems struggle to generalise beyond their training distribution and to abstract effectively. Although some neural architectures might display better alignment with physical laws [88], the generalisability of statistical machines, averaged over all distributions and in the absence of prior knowledge, is constant (no-free-lunch theorem [138]).

The goal of *human-machine teaming* [128] is that each side addresses the limitations of the other while aligning on similar goals. For example, some generalisation capabilities of large language models (LLMs), like the quick production of rhetorically polished texts on any topic, are beyond those of most humans. However, their overgeneralisation errors (“hallucinations” [68]), like replacing specific facts with nonfactual information, can be easily caught by a human expert. Effective teaming requires that humans can assess AI responses and access its underlying rationales (“explanations”) (Figure 1).

The complementarity of humans and AI, and the requirements for effective human-AI teaming, shed new light on traditional knowledge-informed (*analytical*) and *instance-based* AI paradigms. Analytical methods provide compositionality and accessible semantics, albeit in limited scenarios [77], whereas instance-based models are robust to distributional shifts when adequate representation is available [91]. Combining the strengths of various machine methods has inspired emerging research directions under the umbrella of neurosymbolic AI [59].

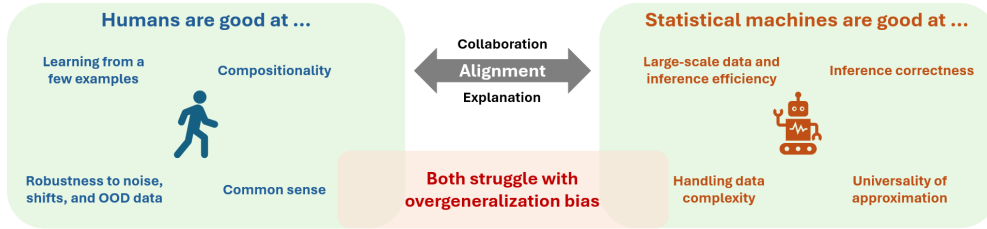


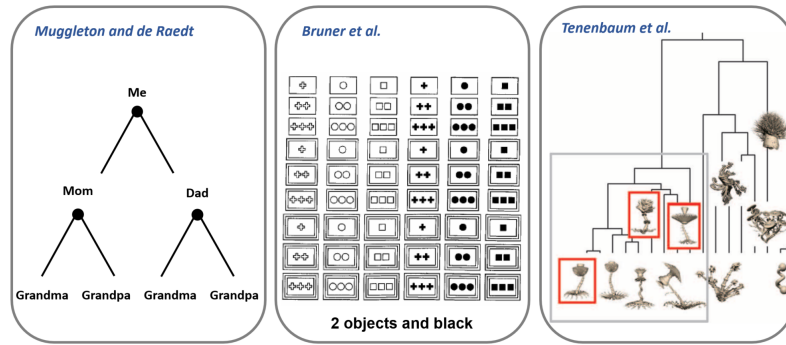
Fig. 1: Comparison of the strengths of humans and statistical machines, illustrating their complementary *generalisation* in human-AI teaming scenarios. Humans excel at compositionality, common sense, abstraction from a few examples, and robustness. Statistical ML excels at large-scale data and inference efficiency, inference correctness, handling data complexity, and the universality of approximation. Overgeneralisation biases remain challenging for both humans and machines. Collaborative and explainable mechanisms are key to achieving alignment in human-AI teaming. See Table 3 for a complete overview of the properties of machine methods, including instance-based and analytical machines.

46 This perspective paper draws on insights about the generalisation of humans and
 47 machines from AI and cognitive science. We analyse three dimensions from the per-
 48 spective of AI alignment: *notions* of, *methods* for, and *evaluation* of generalisation. We
 49 focus on the following questions: What are the known notions of generalisation? What
 50 are the strengths and weaknesses of the generalisation of AI methods? How is generali-
 51 sation evaluated today? What is the impact of current trends in AI, such as foundation
 52 models, on generalisation theories, methodological frameworks, and evaluation prac-
 53 tices? Addressing these questions highlights the need for interdisciplinary approaches
 54 for effective and cognitively supported alignment of human and AI generalisation.

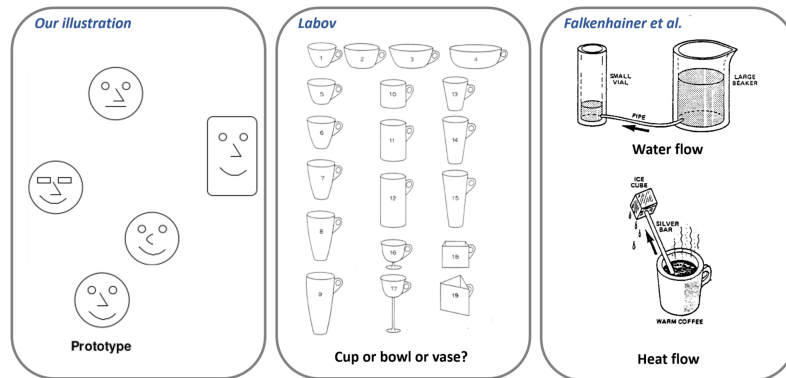
55 2 Parallels in Generalisation by Humans and 56 Machines

57 Approaches to generalisation have been proposed in AI as well as in cognitive psychol-
 58 ogy, and they often mutually inspired each other. This holds for all types of approaches,
 59 whether rule-based, symbolic and knowledge-informed, case- and analogy-based, as
 60 well as neural and statistical. In the following, the mutual influence between AI meth-
 61 ods and cognitive psychology will be illustrated by selected historical milestones,
 62 summarised in Figure 2.

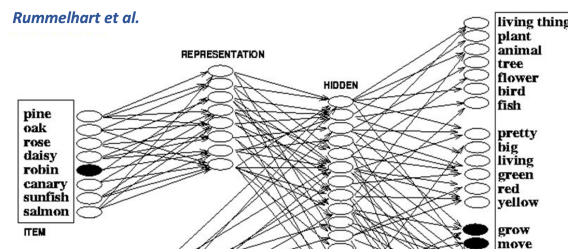
63 In the early days of cognitive psychology, Bruner et al. [17] empirically investigated
 64 human concept learning (Figure 2a), which inspired the first decision tree learn-
 65 ing algorithms [63]. Observations on human learning of relational concepts inspired
 66 early machine learning (ML) approaches to learning from structural representations
 67 [102, 137] and recursive concepts [116]. This class of approaches, often referred to as
 68 inductive programming, allows learning from few examples and taking into account
 69 background knowledge for model induction [54]. Rule learning approaches have been
 70 extended to statistical relational learning [32] to overcome the brittleness of symbolic



(a) Learning the relational rule ‘grandparent’ using a background theory ‘parent’ in Muggleton and De Raedt [102], conjunctive rules in Bruner et al. [17], and names of alien objects modelled as Bayesian inference over a tree-structured domain representation in Tenenbaum et al. [126].



(b) Example-based prototypical representations (cf. Rosch and Mervis [113], and Medin et al. [98]), context-effects (Labov [79]), and analogy (Falkenhainer et al. [38]).



(c) Statistical generalisation: neural network model of semantic memory (Rumelhart and Todd [114]; crop of the illustration).

Fig. 2: Illustrative examples of human generalisation and its inspiration of rule-based (top), example-based (middle), and statistical ML approaches (bottom).

71 learning. Bayesian approaches to rule learning have been introduced as a plausible
 72 framework to model human learning in complex domains such as language acquisition
 73 [126].

74 Rule learning as an explicit approach (system 2 [72]) is apparent for domains of
75 high-level cognition where relevant features can be verbalised [80]. Other cognitive
76 theories have been proposed for domains where knowledge is not (entirely) available
77 in explicit form. For example, prototype theory was proposed as a similarity-based
78 approach where entities are grouped into categories for which similarity within cat-
79 egory borders is maximised and between categories minimised [113], see Figure 2b.
80 This approach is reflected in similarity-based methods to ML, especially k-nearest
81 neighbours [27]. Exemplar theories [108] have been proposed to address the flexibil-
82 ity of human categorisation. For instance, context-dependence of the classification of
83 visual objects has been shown by Labov [79] where a cup might be classified as a
84 bowl or a vase when different contexts, e.g., soup or flowers, are introduced. Another
85 similarity-based approach is analogical reasoning, addressing knowledge transfer from
86 one situation to another, often from a different domain [38, 42, 46]. In contrast to
87 other methods, the analogy is not based on feature but on structural similarity.

88 Neural network approaches were proposed by cognitive scientists [115] as a method
89 of generalisation learning that overcomes the brittleness of symbolic approaches; see
90 Figure 2c. Despite strong arguments from researchers in symbolic AI [41], neural
91 networks and other statistical approaches became the dominant branch of ML due to
92 their superior performance in increasingly large datasets. However, some core concerns
93 by Fodor and Pylyshyn [41] about the relationship between statistical ML and human
94 cognition remain. First, data are separated from a semantic model. While humans
95 who have learnt a concept robustly recognise OOD inputs, ML has only addressed this
96 problem more recently [139]. Second, even if knowledge is primarily implicit in many
97 domains, humans can verbally describe at least part of what constitutes a concept.
98 This observation has recently been reflected in research on explainable AI, where novel
99 approaches to explain black-box models have been proposed focusing on explanations
100 based on concepts and relations [1]. Third, human explanations are typically based
101 on the causal history of an event and a causal explanation for the generalisation itself
102 [100]. Substantial empirical evidence has demonstrated that humans do not focus on
103 the superficial level of event covariations, but reason and learn based on deeper causal
104 representations [132]. In ML, discovering high-level causal variables from low-level
105 observations remains a significant challenge [117].

106 The combination of implicit, neural learning, and explicit, symbolic approaches
107 is addressed in neurosymbolic AI research [59]. Neurosymbolic approaches promise
108 to combine the complementary strengths of neural and analytical techniques, pre-
109 serving robustness while enabling data-model separation and integrating available
110 background knowledge [97]. Combining neural and symbolic approaches to general-
111 isation and explainability within human-centric AI reflects many aspects of human
112 learning flexibility [64]. For effective joint decision-making and problem-solving, the
113 human-AI interface must align with human information processing [67]. Alignment
114 must be established for different aspects, including knowledge state, current informa-
115 tion needs, or values [19]. In cognitive modelling, researchers aim to align algorithmic
116 approaches to learning and human learning [82]. Recent results show that, to date,
117 the performance of human-AI teams lags behind that of the best AI model or the best
118 human alone in many domains [128].

119 3 Notions of Generalisation

120 In the broader context of cognitive science and AI research, there are three different
121 notions of generalisation, which we will cover in turn.

122 3.1 Generalisation as a process

123 Generalisation as a *process* refers to constructing concepts or rules from example
124 data. In cognitive science, the process is typically called *abstraction*, either through
125 associative learning, generalisation by similarity, or the transformation of schemas
126 from lower to higher stages of cognitive development [25]. More broadly, French [43]
127 distinguishes (1) generalisation of concrete instances into an abstract schema, which we
128 call *abstraction*, (2) generalisation through the application or extension of the schema
129 to various situations, which we call *extension*, and (3) generalisation involving the
130 transformation/adaptation of the schema to fit a new context, which we call *analogy*.
131 In AI methods, abstraction (1) corresponds to concept or rule mining in knowledge-
132 informed AI or the learning of models from data in statistical AI [11]; extension (2)
133 relates to online, multi-task, few-shot, or continual learning schemes in statistical and
134 instance-based AI [92, 131, 145]; and analogy (3) relates to transfer learning [146] in
135 statistical AI and reasoning by analogy in analytical AI.

136 Importantly, a generalisation process does not have to start from example data
137 but may abstract, extend, or transfer a pre-existing model beyond its original scope.

138 3.2 Generalisation as a product

139 Generalisation may also refer to the *products* of a generalisation process, such as cat-
140 egories, concepts, rules, and models, in their various representations. Generalisations
141 of *categories and concepts* may be represented using a symbolic definition, as a list of
142 attributes and their bounds (cognitive science: [66]; or decision trees in AI [10]), as a
143 prototype (cognitive science: [112]; AI: [12]), or as a set of exemplars of the category
144 (cognitive science: [108]; or k-nearest neighbor in AI: [110]). Categories or concepts
145 may also be represented as a probability distribution from which examples of this cat-
146 egory can be drawn, which is the notion implicitly used by generative AI models [7].
147 Beyond categories or concepts, products may also be *rules or relations*, represented,
148 for instance, via functions or graphs in parametric or non-parametric form.

149 3.3 Generalisation as an operator

150 Finally, we refer to generalisation as an *operator*, namely the application of a product
151 to new data. The successful application to new data is at the core of generalisation
152 in statistical AI [119]. Under the assumption that training data and new data are
153 independently sampled from the same distribution (IID), one can mathematically
154 prove generalisation via one of three theories: (1) The *Probably Approximately Correct*
155 (*PAC*) framework analyses whether a model (i. e., a product) derived via a machine
156 learning algorithm (i. e., a process) from a random sample of data can be expected to
157 achieve a low prediction error on new data from the same distribution in most cases
158 [119]. (2) *Statistical physics of learning* aims to understand the typical properties of

159 learning algorithms (i. e., processes) with many adaptive parameters [33]. (3) *Vapnik–*
160 *Chervonenkis (VC) dimension theory* focuses on the storage capacity of model classes
161 and their subsequent ability to make accurate predictions on new data [130].

162 Generalisation theory in ML is limited in several ways. It typically predicts that
163 generalisation only occurs if the available data is large enough to not just be memorised
164 [28]. In contrast, humans can generalise from a few samples for a specific task, as
165 generalisation in humans is not a singular event but based on lifelong experience
166 of regularities observed in nature. Few-shot learning addresses this to some extent
167 [16, 103].

168 Ultimately, all three theories have been primarily applied to abstraction processes.
169 For model extension or analogy, the mathematical theory is less well-established. Of
170 particular interest is generalisation across compositions—for instance, in language
171—which has been addressed in analytical AI [48], but is limited by the undecidability
172 or high complexity of inference [144]. Despite such fundamental restrictions, humans
173 operate with compositions similar to those found in language and vision.

174 3.4 Alignment of human and machine notions of generalisation

175 We observe that human and machine notions of generalisation are misaligned. While
176 humans tend toward sparse abstractions and conceptual representations that can be
177 composed or transferred to new domains via analogical reasoning, generalisations in
178 statistical AI tend to be statistical patterns and probability distributions, which can
179 sometimes be extended but still fail to generalise to tasks and domains that are too
180 far removed from the training data. In other words, because humans and machines use
181 different *processes* (e. g., abstraction vs. data-driven learning), they arrive at different
182 *products* (e. g., categories and rules vs. probability distributions) that generalise dif-
183 ferently; if we wish to achieve human-like generalisation ability (as an *operator*), we
184 need new methods for machine generalisation.

185 4 Machine Methods for Generalisation

186 Humans excel in systematic generalisation across representations, contexts, and tasks
187 based on a few observations. Although specific machine methods have shown remark-
188 able results in solving compositional tasks [81], the underlying mechanisms of human
189 generalisation are not sufficiently understood to mimic them in artificial systems [126].

190 AI methods are usually structured according to algorithmic aspects rather than
191 their generalisability (see Table 1). Although these impact generalisation behaviour
192—for instance, symbolic methods often implement compositionality—there is no sim-
193 ple mapping to the form of generalisation that an AI model can achieve. Therefore,
194 we focus on another categorisation, the interplay of observational data (i. e., single
195 instances) and models (i. e., principles that apply to a whole population), as this corre-
196 lates to the generalisability of the model and its suitable evaluation (Section 5). Three
197 categories can be distinguished: (1) The transfer of individual observations to a popula-
198 tion is the basis for *statistical generalisation methods*. (2) The search for observational
199 evidence of an explicit theory is done in *knowledge-informed methods*. (3) *Instance-*
200 *based methods* focus on individuals concerning the source and target of generalisation.

Table 1: Common categories to structure AI methods algorithmically centred. These categories are not uniquely related to their type of generalisation.

Category	Attributes
Training signal	supervised, unsupervised, reinforcement, semi-supervised, self-supervised
Data type	tabular data, data structures (e. g., text, graph), prior knowledge
Model representation	parametric/non-parametric, symbolic/sub-symbolic, black-/white-/grey-box
Training objective	Bayesian inference, maximum likelihood principle, rule learning, mean squared error minimization

201 These choices have different characteristics in terms of their generalisability and
 202 alignment with human generalisation.

203 4.1 Statistical generalisation methods in AI

204 Many modern methods, including deep learning, aim at statistical generalisation:
 205 observational data (i. e., training data) serve as input to an inference mechanism that
 206 extracts a model for the entire population (i. e., the underlying distribution). Gen-
 207 eralisation refers to the ability of the inferred patterns to be successfully applied to
 208 new data (Section 3.3). Typically, algorithmic methods are expressed as optimisation
 209 methods for a model’s loss function, such as the prediction error. As it cannot be
 210 evaluated in the entire population, it is approximated in a given training set, known
 211 as *empirical risk minimisation* [129]. Although evaluation in an independent test set
 212 constitutes an unbiased estimator of generalisation ability (Section 5), the empirical
 213 loss in the training set systematically underestimates the model loss —*generalisation*
 214 *needs to be accounted for explicitly*. Popular strategies regularize models towards a
 215 better generalisation behaviour, such as maximum margin or stability [14]. Even heav-
 216 ily over-parametrised deep models display surprising generalisation capabilities due to
 217 intrinsic regularisation [52].

218 Statistical methods excel in inference correctness and efficiency. Yet, they typically
 219 require large training data sets. This challenge can be partially overcome by technolo-
 220 gies that build their inference on already learnt representations and instance-based
 221 translations, such as few- or zero-shot mechanisms [103]. Still, empirical risk minimi-
 222 sation has a fundamental limitation compared to human generalisation: *generalisation*
 223 *can only be expected in areas covered by observations*, but not for out-of-sample events,
 224 novel contexts, or distributional changes [142]. Indeed, machine behaviour for OOD
 225 settings can significantly deviate from human expectation, with adversarial attacks as
 226 prominent examples of this phenomenon [109].

227 In recent years, many vital settings, including LLMs, have been targeted that
 228 do not allow for a simple analytic expression of human’s intention. Thus, surrogate
 229 losses, such as next token probability, are used as a proxy. With massive training
 230 data, instruction tuning, or human feedback, impressive generalisability arises [16].
 231 However, *the emerging generalisation abilities are only partially understood and do*
 232 *not necessarily align with human expectation*, necessitating a downstream evaluation
 233 (Section 5) if possible at all —human intentions are not necessarily well-formed or

234 static, and the type of information an AI provides could influence human objectives
235 in interactions [21].

236 Statistical generalisation methods are often based on model families with univer-
237 sal approximation capability to account for the lack of domain-specific knowledge.
238 Deep models, for example, can deal with high degrees of nonlinearity and multimodal
239 signals [70]. Yet, the product is typically a black box, which does not reveal insight
240 into its generalisation behaviour; hence, partially unintended generalisation behaviour
241 can easily occur. Recently developed post hoc *explanation methods* allow for a closer
242 inspection of the underlying rationale and its impact on the generalisation behaviour
243 of the model [29].

244 4.2 Knowledge-informed generalisation methods in AI

245 Knowledge-informed generalisation methods aim to find empirical evidence of a theo-
246 ry, resulting in a meaningful representation confirmed by the data. Popular methods
247 include mechanistic models [5], causal models [141], or functional programs [76]. As
248 their semantics is directly accessible, *humans can inspect how these models gener-*
249 *alise to previously not encountered scenarios*. Generalisation is often well aligned with
250 human expectations. Yet, model parameters require semantic grounding, which is chal-
251 lenging to realise with low-level sensor data. Neurosymbolic integration can partially
252 overcome such limitations [96].

253 *Learning the optimal model structure is demanding*, and fundamental limitations
254 such as non-identifiability of structural components might exist [77]. Hence, many
255 methods are restricted to simple schemes, such as description logic, rather than uni-
256 versal approximators. Learning methods such as semantic clustering, probabilistic rule
257 mining, subsumption, or analogies mirror specific representations. Since limited noise
258 robustness and inference efficiency pose significant challenges, hybrid approaches have
259 emerged, such as a transfer of symbolic models to a real-valued embedding space where
260 efficient numeric inference is possible [20]. Knowledge-based approaches enable the
261 explicit inspection and manipulation of information, allowing generalisation based on
262 a few examples, as accompanying rules can ensure valid generalisation. Yet, they are
263 restricted to domains where a theory can be formalised with reasonable effort. As this
264 is often limited, current models do not reach the impressive capabilities of statistical
265 approaches trained on massive datasets.

266 *Systematic compositionality* refers to the ability to generalise and produce novel
267 combinations from known components. It has been fundamental in the design of tradi-
268 tional, logic-based systems; yet statistical methods have struggled with compositional
269 generalisation [41]. Compositionality seems to be a universal principle in nature, since
270 it has been observed in many species [9]. In recent years, significant progress has
271 been made in improving compositional generalisation in deep learning, typically by
272 adding components that mirror the compositional structure of the domain, such as
273 structure-processing neural networks [55] or metalearning for compositional generali-
274 sation [81]. Although these efforts provide a pathway for neural networks to generalise
275 systematically, most of the results are only empirical [136]. There remains a *significant*
276 *gap between the systematic generalisation capabilities of knowledge-informed models*

277 *and the representation learning techniques of deep models*, with neurosymbolic AI
278 promising a viable bridge [59].

279 4.3 Instance-based translation in AI

280 Instance-based non-parametric techniques, such as nearest neighbour methods or case-
281 based reasoning [2], rely on local inference, which is computed when needed based on
282 similar cases encountered previously. They are among the most popular ML meth-
283 ods, showing high flexibility when combined with complex representations [75]. Since
284 they adjust their complexity as needed, they offer universal approximation capabil-
285 ity. Suitable data structures enable efficient training and inference in large data sets.
286 Furthermore, local inference usually allows human inspection of individual decisions
287 —although not of the entire model. Instance-based methods closely resemble concepts
288 in cognitive linguistics, such as a graded degree of belonging to a category, which can
289 be represented by a prototype [112].

290 Instance-based models have shown great promise for incremental learning of distri-
291 butional shifts [91]: They can identify out-of-sample instances based on their similarity
292 to previously encountered data, and they can naturally deal with the challenge of
293 catastrophic forgetting in continual or *lifelong learning* as they can memorize possi-
294 bly relevant data points. This principle also suggests possible solutions to catastrophic
295 forgetting in continual learning using deep statistical models [22]. Conversely, it is pos-
296 sible to implement forget mechanisms if older instances become invalid. The reliance
297 of instance-based methods on similarity means that *a suitable representation is key*
298 *to support generalisation* [78], as it directly influences the model’s ability to evolve
299 patterns across diverse datasets and tasks. Recent work investigates how to achieve
300 representations to support generalisations across tasks or domains [57].

301 *Context has a unique role* as generalisation requires adapting knowledge learnt in
302 one setting to fit a novel, unseen one. Humans can cope with the challenge of acquir-
303 ing context knowledge and assessing the similarity of two contextual representations
304 [120]. ML techniques such as transfer learning, prompting, or retrieval augmented gen-
305 eralisation mimic parts of this process [45]. In this realm, LLMs have demonstrated
306 remarkable capabilities for few-shot or in-context learning [16], often still relying on
307 implicit contextual information. An explicit representation of contextual knowledge
308 and its compositionality, for example, through neurosymbolic AI, is the subject of
309 ongoing research.

310 4.4 Aligning machine generalisation methods and human 311 expectations

312 While *statistical methods* offer powerful universal approximation, their generalisation
313 behaviour does not match human generalisation well, lacking generalisation to OOD
314 samples and compositionality. Another challenge is their black-box nature, where
315 post hoc explanations provide solutions for specific cases. In contrast, *knowledge-*
316 *based methods* enable human insight and compositionality by design, but often at the
317 expense of universality and algorithmic efficiency in the face of structure learning.
318 Emerging neurosymbolic approaches aim to combine these methodological principles.

Table 2: Characterisation of AI generalisation methods.

Pros	Cons
<i>Statistical: generalisation from observations to a population</i>	
universal approximation, surprising generalisation of deep models	black boxes, generalisation failures outside of training distribution
<i>Knowledge-informed: confirm/adapt hypothesis based on observations</i>	
meaningful models, identifiable parameters, generalisation in the limit, compositionality	restriction to simple scenarios, optimization/structure identification computationally demanding
<i>Instance-based: translation from previous observations to a new observation</i>	
flexible to change/distributional shift	rely on suitable representation

319 *Instance-based methods* try a different approach, focusing on generalisation from and
 320 to individual instances. This principle is well-aligned with human generalisation and
 321 enables learning from a few data points and lifelong learning scenarios; yet, results
 322 depend strongly on the choice of representation and context. Recent statistical meth-
 323 ods for representation learning offer promising directions that enable machines to
 324 generalise from a few examples, much like humans. Ultimately, making claims about
 325 the generalisation properties of various machine approaches requires a meaningful
 326 evaluation, which we discuss next.

327 5 Evaluation of Generalisation

328 The theoretical generalisation strengths and weaknesses of the machine method
 329 families are summarised in Table 2. Statistical methods enable *universality of approx-*
 330 *imation* and *inference correctness*, excel at handling *data complexity* and *large-scale*
 331 *data*, and ensure *inference efficiency*. Analytical methods support *compositionality*,
 332 *explainable predictions*, and perform *explicit knowledge manipulation*. Instance-based
 333 methods support *robustness to noise, shifts, and OOD data*, *memorise training*
 334 *samples* reliably, and can *learn from a few samples*.

335 Deriving provable robustness and generalisation guarantees is necessary to define
 336 the theoretical limits of models. Meanwhile, empirically evaluating the machine’s
 337 generalisation is also desirable. From a statistical learning perspective, evaluating
 338 the generalisation of supervised approaches estimates their applicability to new data
 339 (Section 3.3). This formalisation of measuring *inference correctness* on IID data is the-
 340 oretically grounded and remains relevant when assessing systems. It allows measuring
 341 the *universality of approximation* and the ability to handle *large-scale data and per-*
 342 *form efficient inference* on more complex datasets by increasing the task complexity.
 343 However, with increasing task complexity and system opaqueness, it becomes challeng-
 344 ing to guarantee the generalisability assumptions: IID and task-representative data.
 345 For example, Li and Flanigan [86] found that ChatGPT performed well on benchmarks
 346 released before its launch but poorly on those published afterwards. Here, test set leak-
 347 age into ChatGPT violates the IID assumption, invalidating generalisation estimates

348 on pre-release benchmarks. The lack of transparency about the data used to train
349 LLMs makes it challenging to create novel test sets, leading to possible overestima-
350 tions of the model’s generalisation [64, 85]. Namely, while the emergence of foundation
351 models have enabled evaluation of *learning from a few examples* via zero- and few-
352 shot tests [103], the risk of test set memorisation means that test data may appear
353 partially in their training set, invalidating the findings (data contamination [34]).

354 The following discusses key areas related to the evaluation of generalisation and
355 its role in AI applications.

356 5.1 Measuring distributional shifts

357 Distribution shifts can be estimated using statistical distance measures such as the
358 Kullback-Leibler divergence between the feature distributions of the training and test
359 sets [90]. Generative models produce an explicit likelihood estimate $p(x)$ that indicates
360 how typical a sample is to the training distribution. Since discriminative models do
361 not offer this possibility, proxy techniques include calculating cosine similarity between
362 embedding vectors and using nearest-neighbour distances in a transformed feature
363 space. For LLMs, a standard proxy measure for familiarity is perplexity. When the
364 model’s internal representations cannot be directly accessed, the layers of non-linear
365 abstractions in modern (deep) machine learning models allow for gauging relations
366 through intermittent embeddings. Learning evaluation in the context of *drift* can be
367 done using tailored benchmarks [26]. The model’s *robustness to noise, distributional*
368 *shifts, and OOD data* can be studied using adversarial and counterfactual techniques.
369 Adversarial techniques alter data features, such as syntax, semantics, or context, while
370 preserving the underlying task and the original label [109]. In contrast, counterfac-
371 tual techniques create data samples that alter target prediction with minimal input
372 changes [85].

373 5.2 Determining under- and overgeneralisation

374 AI models are created to provide value to humans and are thus assessed against
375 human generalisation, often using the human-centric concepts of under- and overgen-
376 eralisation. *Undergeneralisation* occurs when a change in the input, perceptible or
377 imperceptible, causes a considerable modification within a model. Examples of under-
378 generalisation include model performance degradation under various natural changes,
379 such as environmental perturbations in computer vision [58]. For foundation models,
380 prompt engineering substantially affects performance [49]. In contrast, models *over-*
381 *generalise*, which means that they over-confidently make false predictions for (known
382 or novel) concepts because critical differences are ignored in prediction [13]. A well-
383 known overgeneralisation phenomenon is hallucination, which refers to models that
384 deviate from the source of the information [68]. Other overgeneralisations are biased
385 predictions, for example, when a model predicts a property of an individual from the
386 statistical properties of a demographic group to which they belong [61], and logical
387 fallacies [124].

388 Characterising the model’s under- and overgeneralisation requires an appropri-
389 ate metric, defining its use, and establishing a mechanism to interpret the metric’s

390 score in terms of generalisability beyond the particular test examples. This procedure
391 is susceptible to three caveats. First, the choice between discriminative and genera-
392 tive models determines which representational basis is used to infer similarity [104].
393 For instance, a discriminative model will only learn representations useful to opti-
394 mize classification accuracy, whereas a generative approach would additionally learn
395 representations necessary to describe the data distribution. Second, deep models are
396 prone to learning decision shortcuts and ignoring meaningful features [83], sometimes
397 also called simplicity bias [118]. Third, modern models are often proprietary and fre-
398 quently updated, partly based on user interactions through reinforcement learning,
399 which increases their exposure to datasets and hinders reproducibility. To protect
400 against these caveats, besides using open-source models, it is critical to evaluate across
401 different levels of abstraction, from surface forms to semantic similarity and higher-
402 level structural mappings, and explicitly consider the application context and limits.
403 Consequently, machines are increasingly tested for their ability to *handle complex data*,
404 such as multimodal datasets [143], and to exhibit *compositionality* in tasks such as
405 commonsense reasoning [30], analogies [125], and concept induction [106].

406 5.3 Distinguishing memorisation and generalisation

407 In AI, *memorisation* refers to learning details from the training data, including facts
408 and noise. Memorisation may be beneficial in some cases (e.g., Paris is the capital of
409 France), but detrimental in others (e.g., Biden is the president of the United States).
410 This observation raises a question: *When should models generalise, when should they*
411 *memorise, and how can this distinction be evaluated?* Whether the models should
412 generalise or memorise is set a priori. When learning from experience, generalisation
413 is crucial, for instance, in recognising a new manifestation of a vase [79]. Conse-
414 quently, generalisation setups include cross-domain validation and robustness testing.
415 In contrast, factual knowledge is often memorised: Paris is the capital of France, and
416 mosquitoes fly. Tasks such as answering factual questions [133] and reasoning about
417 legal precedents [53] require *memorisation*. While evaluating memorisation and gen-
418 eralisation separately is informative, many tasks, including causal reasoning [107],
419 argumentation [4], and theorem proving [140], require holistic integration of generalisa-
420 tion and memorisation, centred around *explicit knowledge manipulation*. The explicit
421 knowledge use allows testing a model’s *explainable predictions* through human studies
422 and faithfulness evaluation [105].

423 5.4 Alignment of machine evaluation of generalisation to 424 humans

425 Effective alignment of AI with humans requires a principled evaluation of its strengths
426 and weaknesses in generalisation. Evaluating AI generalisability in the context of
427 its alignment comprises: (1) deriving provable guarantees about AI’s properties, like
428 compositionality and learning from a few samples, and (2) performing empirical eval-
429 uation by leveraging task-specific benchmarks and metrics. The evaluation of AI
430 generalisability measures its performance concerning distribution shifts, determining
431 its undergeneralisation (adaptability to task variations such as camera perturbations)

432 and overgeneralisation (e. g., hallucinations), and its ability to memorise and gener-
 433 alise adequately when necessary. These three aspects are essential for alignment, since
 434 distributional shifts, task variations, and both facts and noise are natural in human-AI
 435 teaming scenarios. We summarise typical approaches for evaluating key generalisation
 436 properties in Table 3. In the next section, we discuss open challenges for measuring
 437 generalisation in the context of human-AI alignment.

Table 3: Desired properties of generalisation that emerge Sections 3 and 4. The properties are listed in order of their appearance in Section 5, split by a subsection using a horizontal line. Each of the properties is supported by statistical (S), analytical (A), and instance-based (I) AI methods, as discussed in Section 4. The evaluation methods for each property are discussed and substantiated by the references in Section 5. Achievement is indicated by ”+” and non-achievement by ”-”; we use a strict evaluation to avoid partial scoring.

Property	Method			Evaluation
	S	A	I	
inference correctness	+	-	-	train-test splits
universality of approximation	+	-	+	increasing task complexity
large-scale data and inference efficiency	+	-	+	large/complex datasets
learning from a few samples	-	+	+	zero/few-shot tests
robustness to noise, shifts, OOD data	-	-	+	adversarial, shifted, counterfactual tasks
compositionality	-	+	-	analogy, abstraction, concept induction
handling data complexity	+	-	-	multimodal datasets
memorisation of training samples	-	-	+	factuality datasets, precedents
explicit knowledge manipulation	-	+	-	causality, argumentation, theorem proving
explainable predictions	-	+	-	human studies, faithfulness

438 6 Emerging Directions

439 The previous sections addressed the challenges of aligning human and machine
 440 intelligence, emphasising AI’s potential to enhance human generalisation. Table 3
 441 summarises the properties, methods and evaluation practices, and highlights the com-
 442 plementarity of the three methods in achieving important properties for aligning
 443 human and AI generalisation. Statistical approaches enable universality of approx-
 444 imation and inference correctness, instance-based methods enable robustness and
 445 memorisation, while analytical techniques are designed for compositionality and
 446 explainable predictions. The evaluation column displays the range of approaches used
 447 to assess various methods. Next, we discuss future research directions for novel gen-
 448 eralisation theories, hybrid methods, evaluation practices, and alignment in future
 449 human-AI teams.

450 ***Generalisation theory in the era of foundation models***

451 Recent zero-shot and in-context learning approaches in LLMs and LRMs (large
452 reasoning models, [73]) implicitly generalise to tasks unrelated to their training with-
453 out explicit similarity [18]. In other words, model builders assume that LLMs have
454 implicitly generalised (process, Section 3.1) to generalisations (product, Section 3.2)
455 that allow generalisation (operator, Section 3.3) to entirely new tasks and domains.
456 However, this assumption remains unsubstantiated, leading to an overestimation of
457 the generalisability of foundation models [73]. This highlights the need for further
458 research. First, new generalisations *processes and products* are needed to provide guar-
459 antees (or at least reasons to believe) that zero-shot application to new tasks is viable,
460 potentially through encoding invariances or equivariances [24], as used in complex
461 architectures such as AlphaFold; or through cognitively inspired representations, such
462 as prototypes, which have proven efficient for domain generalisation [87]. Second,
463 a new *theory* is required to define when few- or zero-shot applications are feasible.
464 It has recently been shown that —unlike in classical learning theory—high dimen-
465 sionality of the signals might be key to the generalisability of few-shot learners and
466 over-parametrised deep networks [127].

467 ***Generalisable neurosymbolic methods***

468 Neurosymbolic AI combines robust, data-driven latent models with the precision of
469 explicit compositional models [59]. However, several challenges remain: Defining *prov-*
470 *able generalisation properties*, including worst-case bounds, is crucial. Recent work
471 exploits the compositionality of neurosymbolic systems to derive upper and lower
472 bounds for the robustness of generalisations [95]. Verifying correctness instead of just
473 robustness remains an unaddressed problem.

474 Current symbolic representations in neurosymbolic systems are typically of low
475 *expressivity* (knowledge graphs, propositional logic), allowing for only limited forms of
476 generalisation. Recent works explore the use of richer formalisms. The use of descrip-
477 tion logics in neurosymbolic systems is particularly relevant since description logics
478 are designed to capture forms of generalisation [122].

479 While a theory about the *compositionality* of neurosymbolic *systems* is emerg-
480 ing [31], a theory on how to compose the generalisations themselves is lacking. For
481 symbolic representations, a theory on abstractions [47] exists, but the question of
482 composing latent representations, such as embeddings, remains open.

483 Finally, handling the *context dependency* of generalisations remains a challenge.
484 How do we measure the distance between contexts and apply generalisation across
485 contexts? How do we know when a context is too novel and we overgeneralise? A
486 possible direction is formal modelling of contextual dimensions such as time and space,
487 following prior research on axiomatising commonsense knowledge [50]. CYC’s notion of
488 hierarchical microtheories, each containing a collection of axioms [84], can be revisited
489 from a neurosymbolic perspective, for instance, by expressing axioms of microtheories
490 in flexible natural language representations [134].

491 ***Generalisation in continual learning***

492 The enormous effort required to train foundational models, and the increasing avail-
493 ability of pre-trained models, has led to a transition from models which are trained
494 from scratch to systems based on foundation models that undergo continuous adapta-
495 tion to new data or tasks. However, naive approaches carry a high risk of *catastrophic*
496 *forgetting*, which, surprisingly, seems to be higher for larger LLMs [93].

497 The relation between continual learning and generalisation is complex: continual
498 learning methods are designed to combat catastrophic forgetting when generalising to
499 novel tasks, while within-task generalisation facilitates faster learning and improved
500 performance in subsequent continual learning tasks [121]. Thus, *concept drift* requires
501 learning of the underlying features that can be readily applied to novel tasks. This
502 can be achieved by an extension of statistical methods to hybrid approaches, which
503 attend to the preservation of learnt signals: for example, formalising domain rules as
504 ontologies or symbolic constraints enables a system to detect drift whenever incoming
505 data or model outputs violate these constraints, serving as an early warning signal for
506 distributional change that may disrupt the generalisation capabilities of the system.
507 Recent work on graph streams [94] employs neurosymbolic prototypes, where repre-
508 sentative subgraphs are embedded in vector spaces. Another remedy can be based on
509 data-driven approaches, such as (possibly self-supervised) rehearsal technologies, for a
510 robust memorisation of important information, albeit at increased computational costs
511 [62]. More efficient alternatives aim for architectural solutions such as the incorpora-
512 tion of instance-based representations into statistical models [44]. However, theoretical
513 insight on the effect of overparametrisation or task similarity on the generalisability
514 and forgetting of a model currently exists for very simple models only [89].

515 ***Evaluation of generalisation in foundation models***

516 Several directions have emerged to address data contamination, spurious correlations,
517 and overfitting in state-of-the-art models. Abstraction benchmarks for visual reason-
518 ing [23], analogy [125], and lateral thinking [69] are gaining popularity. Crowdsourcing
519 can be used to create and scale benchmarks, but it can also introduce cognitive and
520 cultural biases by annotators [36], which remains poorly understood. On the other
521 hand, evaluation servers and public leaderboards with private test datasets prevent
522 overfitting but lack standardisation and are costly to maintain. A final direction is sim-
523 ulation environments and synthetic data generators [37], though they often suffer from
524 a sim-to-real gap. To address reproducibility, researchers proposed the model [101] and
525 data cards [111] to report the details of the experiment, and reproducibility checklists
526 based on a broad consensus [74], albeit with limited coverage of generalisability.

527 ***Aligning generalisation in future human-AI teams***

528 The goal of effective human-AI teaming and the appearance of legal frameworks such
529 as the EU AI Act [6, 8] require transparent *collaboration* workflows, with *explanations*
530 bridging the gaps between human and AI reasoning [3, 64]. Section 2 discussed that
531 this alignment must occur at the *output level*. However, when misalignments occur,
532 (e. g., AI predicts tumour type 1 and the doctor diagnoses tumour type 3), mechanisms
533 for realignment and error correction become critical. Such mechanisms pose stricter

534 requirements for collaboration on the *process level* through concepts and relations
535 [82]. Examples of realignment techniques include language games, where realignment
536 emerges from interaction, and physics-informed models that refine predictions on
537 object permanence.

538 A critical challenge of human-AI teaming is *reconciling the fundamentally differ-*
539 *ent reasoning paradigms* of humans and AI, like human causal models and AI’s deep
540 learning associations. Efforts such as concept-based explanations [135] and those con-
541 sidering relationships [40] suggest the potential for intertranslatability into a common
542 explanatory language.

543 Furthermore, robust evaluation frameworks should consider both objective task-
544 related outcomes and subjective process-related experiences, as well as the long-term
545 ramifications of the collaboration, taking into account each party’s contributions
546 and responsibilities. Despite the emergence of evaluation frameworks and metrics for
547 humans that augment AI (e.g., in manual data labelling), AI that helps humans
548 (e.g., conversational question answering), and balanced collaborations where both
549 contribute equally (e.g., medical decision-making) [15], there is little research on
550 evaluating the generalisation capabilities of such teams.

551 **Acknowledgements.** The manuscript resulted from the May 2024 Dagstuhl semi-
552 nar: Generalization by People and Machines (24192). Ken Forbus, Piek Vossen, Dafna
553 Shahaf, Wael Abd-Almageed, and Michael Waldmann provided valuable insights dur-
554 ing the seminar. FI is funded by the NWO AiNed project “Human-Centric AI Agents
555 with Common Sense”. BH, BP, A-CNN gratefully acknowledge funding by the Min-
556 istry of Culture and Science of North Rhine-Westphalia (MKW NRW) through the
557 project SAIL (grant no. NW21-059A-D).

558 **Competing interests.** The authors declare no competing interests.

559 References

- 560 [1] Reduan Achtibat, Maximilian Dreyer, Ilona Eisenbraun, Sebastian Bosse,
561 Thomas Wiegand, Wojciech Samek, and Sebastian Lapuschkin. From Attribution
562 Maps to Human-Understandable Explanations through Concept Relevance
563 Propagation. *Nature Machine Intelligence*, 5:1006–1019, 2023.
- 564 [2] D.W. Aha. *Lazy Learning*. Springer Netherlands, 2013.
- 565 [3] Zeynep Akata, Dan Balliet, Maarten de Rijke, Frank Dignum, Virginia Dignum,
566 Gusztai Eiben, Antske Fokkens, Davide Grossi, Koen V. Hindriks, Holger H.
567 Hoos, and et al. A Research Agenda for Hybrid Intelligence: Augmenting Human
568 Intellect With Collaborative, Adaptive, Responsible, and Explainable Artificial
569 Intelligence. *IEEE Computer*, 53(8):18–28, 2020.
- 570 [4] Katie Atkinson, Pietro Baroni, Massimiliano Giacomin, Anthony Hunter, Henry
571 Prakken, Chris Reed, Guillermo Simari, Matthias Thimm, and Serena Villata.
572 Towards Artificial Argumentation. *AI Magazine*, 38(3):25–36, 2017.

- 573 [5] Ruth E Baker, Jose-Maria Peña, Jayaratnam Jayamohan, and Antoine
574 Jérusalem. Mechanistic models versus machine learning, a fight worth fighting
575 for the biological community? *Biology Letters*, 14(5), 2018.
- 576 [6] Alejandro Bellogín, Oliver Grau, Stefan Larsson, Gerhard Schimpf, Biswa Sen-
577 gupta, and Gürkan Solmaz. The EU AI Act and the Wager on Trustworthy AI.
578 *Communications of the ACM*, 67(12):58–65, 2024.
- 579 [7] Staphord Bengesi, Hoda El-Sayed, Md Kamruzzaman Sarker, Yao Houkpati,
580 John Irungu, and Timothy Oladunni. Advancements in Generative AI: A
581 Comprehensive Review of GANs, GPT, Autoencoders, Diffusion Model, and
582 Transformers. *IEEE Access*, 2024.
- 583 [8] Yoshua Bengio, Geoffrey Hinton, Andrew Yao, Dawn Song, Pieter Abbeel,
584 Trevor Darrell, Yuval Noah Harari, Ya-Qin Zhang, Lan Xue, Shai Shalev-
585 Shwartz, et al. Managing extreme AI risks amid rapid progress. *Science*, 384
586 (6698):842–845, 2024.
- 587 [9] M. Berthet, M. Surbeck, and S. W. Townsend. Extensive compositionality in
588 the vocal system of bonobos. *Science*, 388(6742):104–108, 2025.
- 589 [10] Dimitris Bertsimas and Jack Dunn. Optimal classification trees. *Machine*
590 *Learning*, 106:1039–1082, 2017.
- 591 [11] Michael Biehl. *The Shallow and the Deep: A biased introduction to neural*
592 *networks and old school machine learning*. University of Groningen Press, 2023.
- 593 [12] Jacob Bien and Robert Tibshirani. Prototype selection for interpretable
594 classification. *The Annals of Applied Statistics*, 5(4):2403–2424, 2011.
- 595 [13] Terrance E. Boulton, Steve Cruz, Akshay R. Dhamija, Manuel Gunther, James
596 Henrydoss, and Walter J. Scheirer. Learning and the Unknown : Surveying Steps
597 Toward Open World Recognition. *AAAI Conference on Artificial Intelligence*,
598 2019.
- 599 [14] Olivier Bousquet and André Elisseeff. Algorithmic Stability and Generalization
600 Performance. In *Advances in Neural Information Processing Systems*, volume 13.
601 MIT Press, 2000.
- 602 [15] Marvin Braun, Maike Greve, and Ulrich Gnewuch. The New Dream Team?: A
603 Review of Human-AI Collaboration Research From a Human Teamwork Per-
604 spective. In *International Conference on Information Systems*, 2023. Association
605 for Information Systems.
- 606 [16] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan,
607 Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda
608 Askell, et al. Language Models are Few-Shot Learners. In *Advances in*

- 609 *Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran
610 Associates, Inc., 2020.
- 611 [17] Jerome Bruner, Jacqueline J. Goodnow, and George A. Austin. *A Study of*
612 *Thinking*. Wiley, 1956.
- 613 [18] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric
614 Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al.
615 Sparks of Artificial General Intelligence: Early experiments with GPT-4, 2023.
- 616 [19] Patrick Butlin. AI Alignment and Human Reward. In *2021 AAAI/ACM*
617 *Conference on AI, Ethics, and Society*, pages 437–445, 2021.
- 618 [20] Jiahang Cao, Jinyuan Fang, Zaiqiao Meng, and Shangsong Liang. Knowledge
619 Graph Embedding: A Survey from the Perspective of Representation Spaces.
620 *ACM Computing Surveys*, 56(6), 2024.
- 621 [21] Micah Carroll, Davis Foote, Anand Siththaranjan, Stuart Russell, and Anca
622 Dragan. AI alignment with changing and influenceable reward functions. In *41st*
623 *International Conference on Machine Learning*. JMLR, 2024.
- 624 [22] Hung-Jen Chen, An-Chieh Cheng, Da-Cheng Juan, Wei Wei, and Min Sun.
625 Mitigating Forgetting in Online Continual Learning via Instance-Aware Param-
626 eterization. In *Advances in Neural Information Processing Systems*, volume 33,
627 pages 17466–17477. Curran Associates, Inc., 2020.
- 628 [23] François Chollet. On the measure of intelligence. *arXiv preprint*
629 *arXiv:1911.01547*, 2019.
- 630 [24] Taco Cohen and Max Welling. Group Equivariant Convolutional Networks. In
631 *33rd International Conference on Machine Learning*, volume 48 of *Proceedings*
632 *of Machine Learning Research*, pages 2990–2999, 2016. PMLR.
- 633 [25] Eliana Colunga and Linda B Smith. The emergence of abstract ideas: Evidence
634 from networks and babies. *Philosophical Transactions of the Royal Society of*
635 *London. Series B: Biological Sciences*, 358(1435):1205–1214, 2003.
- 636 [26] Andrea Cossu, Davide Bacciu, Alessio Bernardo, Emanuele Della Valle, Alexan-
637 der Gepperth, Federico Giannini, Barbara Hammer, and Giacomo Ziffer. Don’t
638 Drift Away: Advances and Applications of Streaming and Continual Learning.
639 *European Symposium on Artificial Neural Networks*, 2025.
- 640 [27] Thomas Cover and Peter Hart. Nearest neighbor pattern classification. *IEEE*
641 *Transactions on Information Theory*, 13(1):21–27, 1967.
- 642 [28] Thomas M Cover. Geometrical and Statistical Properties of Systems of Linear
643 Inequalities with Applications in Pattern Recognition. *IEEE Transactions on*

- 644 *Electronic Computers*, (3):326–334, 1965.
- 645 [29] Abhilekha Dalal, Rushrukh Rayan, Adrita Barua, Eugene Y. Vasserman,
646 Md. Kamruzzaman Sarker, and Pascal Hitzler. On the Value of Labeled Data
647 and Symbolic Methods for Hidden Neuron Activation Analysis. In *Neural-
648 Symbolic Learning and Reasoning Proceedings, Part II*, volume 14980 of *Lecture
649 Notes in Computer Science*, pages 109–131, 2024. Springer.
- 650 [30] Ernest Davis. Benchmarks for Automated Commonsense Reasoning: A Survey.
651 *ACM Computing Surveys*, 56(4):1–41, 2023.
- 652 [31] Maaïke de Boer, Quirine Smit, Michael van Bekkum, André Meyer-Vitali, and
653 Thomas Schmid. Design Patterns for LLM-based Neuro-Symbolic Systems.
654 *Neurosymbolic Artificial Intelligence (to appear)*, 2025.
- 655 [32] Luc De Raedt and Kristian Kersting. Probabilistic inductive logic programming.
656 In *Probabilistic inductive logic programming: theory and applications*, pages 1–27.
657 Springer, 2008.
- 658 [33] Aurélien Decelle. An introduction to machine learning: a perspective from sta-
659 tistical physics. *Physica A: Statistical Mechanics and its Applications*, page
660 128154, 2022.
- 661 [34] Jesse Dodge, Maarten Sap, Ana Marasović, William Agnew, Gabriel Ilharco,
662 Dirk Groeneveld, Margaret Mitchell, and Matt Gardner. Documenting Large
663 Webtext Corpora: A Case Study on the Colossal Clean Crawled Corpus. In
664 *2021 Conference on Empirical Methods in Natural Language Processing*, pages
665 1286–1305, 2021. Association for Computational Linguistics.
- 666 [35] Jon Donnelly, Zhicheng Guo, Alina Jade Barnett, Hayden McTavish, Chaofan
667 Chen, and Cynthia Rudin. Rashomon Sets for Prototypical-Part Networks:
668 Editing Interpretable Models in Real-Time. In *Computer Vision and Pattern
669 Recognition Conference (CVPR)*, pages 4528–4538, 2025.
- 670 [36] Tim Draws, Alisa Rieger, Oana Inel, Ujwal Gadiraju, and Nava Tintarev. A
671 Checklist to Combat Cognitive Biases in Crowdsourcing. In *AAAI Conference
672 on Human Computation and Crowdsourcing*, volume 9, pages 48–59, 2021.
- 673 [37] Jiafei Duan, Samson Yu, Hui Li Tan, Hongyuan Zhu, and Cheston Tan. A Survey
674 of Embodied AI: From Simulators to Research Tasks. *IEEE Transactions on
675 Emerging Topics in Computational Intelligence*, 6(2):230–244, 2022.
- 676 [38] Brian Falkenhainer, Kenneth D Forbus, and Dedre Gentner. The structure-
677 mapping engine: Algorithm and examples. *Artificial intelligence*, 41(1):1–63,
678 1989.

- 679 [39] Emilio Ferrara. GenAI against humanity: Nefarious applications of generative
680 artificial intelligence and large language models. *Journal of Computational Social*
681 *Science*, pages 1–21, 2024.
- 682 [40] Bettina Finzel, Patrick Hilme, Johannes Rabold, and Ute Schmid. When a Rela-
683 tion Tells More Than a Concept: Exploring and Evaluating Classifier Decisions
684 with CoReX. *arXiv preprint arXiv:2405.01661*, 2024.
- 685 [41] Jerry A Fodor and Zenon W Pylyshyn. Connectionism and cognitive architec-
686 ture: A critical analysis. *Cognition*, 28(1-2):3–71, 1988.
- 687 [42] Kenneth D Forbus, Ronald W Ferguson, Andrew Lovett, and Dedre Gentner.
688 Extending SME to handle large-scale cognitive modeling. *Cognitive Science*, 41
689 (5):1152–1201, 2017.
- 690 [43] Robert M. French. *The Subtlety of Sameness: A Theory and Computer Model*
691 *of Analogy-Making*. MIT Press, 1995.
- 692 [44] Neil De La Fuente, Maria Pilligua, Daniel Vidal, Albin Soutiff, Cecilia Curreli,
693 Daniel Cremers, and Andrey Barsky. Prototype Augmented Hypernetworks for
694 Continual Learning. *arXiv preprint, arXiv:2505.07450*, 2025.
- 695 [45] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai,
696 Jiawei Sun, Meng Wang, and Haofen Wang. Retrieval-Augmented Generation
697 for Large Language Models: A Survey. *arXiv preprint, arXiv:2312.10997*, 2024.
- 698 [46] Dedre Gentner. Structure-Mapping: A Theoretical Framework for Analogy.
699 *Cognitive Science*, 7(2):155–170, 1983.
- 700 [47] Fausto Giunchiglia, Adolfo Villafiorita, and Toby Walsh. Theories of abstraction.
701 *AI Communications*, 10(3-4):167–176, 1997.
- 702 [48] E Mark Gold. Language identification in the limit. *Information and Control*, 10
703 (5):447–474, 1967.
- 704 [49] Hila Gonen, Srini Iyer, Terra Blevins, Noah Smith, and Luke Zettlemoyer.
705 Demystifying Prompts in Language Models via Perplexity Estimation. In *Find-*
706 *ings of the Conference on Empirical Methods in Natural Language Processing*,
707 pages 10136–10148, 2023. Association for Computational Linguistics.
- 708 [50] Andrew S Gordon and Jerry R Hobbs. *A Formal Theory of Commonsense*
709 *Psychology: How People Think People Think*. Cambridge University Press, 2017.
- 710 [51] JuraJ Gottweis, Wei-Hung Weng, Alexander N. Daryin, Tao Tu, Anil Palepu,
711 Petar Sirkovic, Artiom Myaskovsky, Felix Weissenberger, Keran Rong, Ryutaro
712 Tanno, et al. Towards an AI co-scientist. *CoRR*, abs/2502.18864, 2025.

- 713 [52] P. Grohs and G. Kutyniok, editors. *Mathematical Aspects of Deep Learning*.
714 Cambridge University Press, 2022.
- 715 [53] Neel Guha, Julian Nyarko, Daniel Ho, Christopher Ré, Adam Chilton, Alex
716 Chohlas-Wood, Austin Peters, Brandon Waldon, Daniel Rockmore, Diego Zam-
717 brano, et al. LegalBench: A Collaboratively Built Benchmark for Measuring
718 Legal Reasoning in Large Language Models. *Advances in Neural Information
719 Processing Systems*, 36:44123–44279, 2023.
- 720 [54] Sumit Gulwani, José Hernández-Orallo, Emanuel Kitzelmann, Stephen H Mug-
721 gleton, Ute Schmid, and Benjamin Zorn. Inductive programming meets the real
722 world. *Communications of the ACM*, 58(11):90–99, 2015.
- 723 [55] Barbara Hammer. *Learning with Recurrent Neural Networks*. Springer-Verlag,
724 2000.
- 725 [56] S. Harnad. To cognize is to categorize: Cognition is categorization. In *Hand-
726 book of Categorization in Cognitive Science (2nd edition)*, pages 21–54. Elsevier
727 Academic Press, 2017.
- 728 [57] Jerry Zhi-Yang He, Zackory Erickson, Daniel S. Brown, Aditi Raghunathan, and
729 Anca Dragan. Learning Representations that Enable Generalization in Assistive
730 Tasks. In *6th Annual Conference on Robot Learning*, 2022.
- 731 [58] Dan Hendrycks and Thomas Dietterich. Benchmarking Neural Network Robust-
732 ness to Common Corruptions and Perturbations. *International Conference on
733 Learning Representations (ICLR)*, 2019.
- 734 [59] Pascal Hitzler, Md. Kamruzzaman Sarker, and Aaron Eberhart, editors. *Com-
735 pendium of Neurosymbolic Artificial Intelligence*, volume 369 of *Frontiers in
736 Artificial Intelligence and Applications*. IOS Press, 2023.
- 737 [60] Andreas Holzinger, Anna Saranti, Alessa Angerschmid, Bettina Finzel, Ute
738 Schmid, and Heimo Mueller. Toward human-level concept learning: Pattern
739 benchmarking for AI algorithms. *Patterns*, 4:100788, 2023.
- 740 [61] Dirk Hovy and Shannon L. Spruit. The Social Impact of Natural Language
741 Processing. In *54th Annual Meeting of the Association for Computational
742 Linguistics (Volume 2: Short Papers)*, pages 591–598, 2016. Association for
743 Computational Linguistics.
- 744 [62] Jianheng Huang, Leyang Cui, Ante Wang, Chengyi Yang, Xinting Liao, Linfeng
745 Song, Junfeng Yao, and Jinsong Su. Mitigating Catastrophic Forgetting in Large
746 Language Models with Self-Synthesized Rehearsal. In *62nd Annual Meeting of
747 the Association for Computational Linguistics (Volume 1: Long Papers)*, pages
748 1416–1428, 2024. Association for Computational Linguistics.

- 749 [63] Earl B Hunt, Janet Marin, and Philip J Stone. *Experiments in Induction*.
750 Academic Press, 1966.
- 751 [64] Filip Ilievski. *Human-Centric AI with Common Sense*. Springer, 2025.
- 752 [65] Filip Ilievski, Barbara Hammer, Frank van Harmelen, Benjamin Paassen, Sascha
753 Saralajew, Ute Schmid, Michael Biehl, Marianna Bolognesi, Xin Luna Dong,
754 Kiril Gashteovski, et al. Aligning generalization between humans and machines.
755 *Nature Machine Intelligence*, 7(9):1378–1389, 2025.
- 756 [66] Ray S Jackendoff. *Semantics and cognition*, volume 8. MIT press, 1985.
- 757 [67] Jiaming Ji, Tianyi Qiu, Boyuan Chen, Borong Zhang, Hantao Lou, Kaile Wang,
758 Yawen Duan, Zhonghao He, Jiayi Zhou, Zhaowei Zhang, et al. AI Alignment: A
759 Comprehensive Survey. *arXiv preprint arXiv:2310.19852*, 2023.
- 760 [68] Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii,
761 Ye Jin Bang, Andrea Madotto, and Pascale Fung. Survey of Hallucination in
762 Natural Language Generation. *ACM Computing Surveys*, 55(12), 2023.
- 763 [69] Yifan Jiang, Filip Ilievski, Kaixin Ma, and Zhivar Sourati. BRAINTEASER:
764 Lateral Thinking Puzzles for Large Language Models. In *Conference on*
765 *Empirical Methods in Natural Language Processing*, pages 14317–14332, 2023.
766 Association for Computational Linguistics.
- 767 [70] Tianzhe Jiao, Chaopeng Guo, Xiaoyue Feng, Yuming Chen, and Jie Song. A
768 Comprehensive Survey on Deep Learning Multi-Modal Fusion: Methods, Tech-
769 nologies and Applications. *Computers, Materials and Continua*, 80(1):1–35,
770 2024.
- 771 [71] John M. Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Fig-
772 urnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin
773 Žídek, Anna Potapenko, et al. Highly accurate protein structure prediction with
774 AlphaFold. *Nature*, 596(7873):583–589, 2021.
- 775 [72] Daniel Kahneman. *Thinking, Fast and Slow*. Macmillan, 2011.
- 776 [73] Subbarao Kambhampati, Kaya Stechly, and Karthik Valmeekam. (How) Do
777 reasoning models reason? *Annals of the New York Academy of Sciences*, 2025.
- 778 [74] Sayash Kapoor, Emily M. Cantrell, Kenny Peng, Thanh Hien Pham, Christo-
779 pher A. Bail, Odd Erik Gundersen, Jake M. Hofman, Jessica Hullman,
780 Michael A. Lones, Momin M. Malik, et al. REFORMS: Consensus-based Rec-
781 ommendations for Machine-learning-based Science. *Science Advances*, 10(18):
782 eadk3452, 2024.

- 783 [75] Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike
784 Lewis. Generalization through Memorization: Nearest Neighbor Language
785 Models. In *International Conference on Learning Representations*, 2020.
- 786 [76] Emanuel Kitzelmann and Ute Schmid. Inductive Synthesis of Functional Pro-
787 grams: An Explanation Based Generalization Approach. *Journal of Machine*
788 *Learning Research*, 7(15):429–454, 2006.
- 789 [77] D. Koller and N. Friedman. *Probabilistic Graphical Models: Principles and*
790 *Techniques*. Adaptive computation and machine learning. MIT Press, 2009.
- 791 [78] Brian Kulis. Metric Learning: A Survey. *Foundations and Trends® in Machine*
792 *Learning*, 5(4):287–364, 2013.
- 793 [79] William Labov. The boundaries of words and their meanings. In *New Ways*
794 *of Analyzing Variation in English*, pages 67–90. Georgetown University Press,
795 1973.
- 796 [80] Daniel Lafond, Yves Lacouture, and Andrew L Cohen. Decision-tree models
797 of categorization response times, choice proportions, and typicality judgments.
798 *Psychological Review*, 116(4):833, 2009.
- 799 [81] Brenden M. Lake and Marco Baroni. Human-like systematic generalization
800 through a meta-learning neural network. *Nature*, 623(7985):115–121, 2023.
- 801 [82] Pat Langley and Herbert A Simon. The Central Role of Learning in Cognition. In
802 *Cognitive Skills and Their Acquisition*, pages 361–380. Psychology Press, 2013.
- 803 [83] Sebastian Lapuschkin, Stephan Wäldchen, Alexander Binder, Grégoire Mon-
804 tavon, Wojciech Samek, and Klaus-Robert Müller. Unmasking Clever Hans Pre-
805 dictors and Assessing What Machines Really Learn. *Nature Communications*,
806 10, 2019.
- 807 [84] Douglas B Lenat. CYC: A large-scale investment in knowledge infrastructure.
808 *Communications of the ACM*, 38(11):33–38, 1995.
- 809 [85] Martha Lewis and Melanie Mitchell. Using counterfactual tasks to evaluate
810 the generality of analogical reasoning in large language models. *arXiv preprint*
811 *arXiv:2402.08955*, 2024.
- 812 [86] Changmao Li and Jeffrey Flanigan. Task contamination: Language models
813 may not be few-shot anymore. In *AAAI Conference on Artificial Intelligence*,
814 volume 38, pages 18471–18480, 2024.
- 815 [87] Muxin Liao, Shishun Tian, Yuhang Zhang, Guoguang Hua, Wenbin Zou, and
816 Xia Li. Calibration-Based Multi-Prototype Contrastive Learning for Domain
817 Generalization Semantic Segmentation in Traffic Scenes. *IEEE Transactions on*

- 818 *Intelligent Transportation Systems*, 25(12):20985–21001, 2024.
- 819 [88] Henry W. Lin, Max Tegmark, and David Rolnick. Why Does Deep and Cheap
820 Learning Work So Well? *Journal of Statistical Physics*, 168(6):1223–1247, 2017.
- 821 [89] Sen Lin, Peizhong Ju, Yingbin Liang, and Ness Shroff. Theory on forgetting and
822 generalization of continual learning. In *International Conference on Machine*
823 *Learning*. JMLR, 2023.
- 824 [90] Jiashuo Liu, Zheyang Shen, Yue He, Xingxuan Zhang, Renzhe Xu, Han Yu, and
825 Peng Cui. Towards out-of-distribution generalization: A survey. *arXiv preprint*
826 *arXiv:2108.13624*, 2021.
- 827 [91] Viktor Losing, Barbara Hammer, and Heiko Wersing. KNN Classifier with Self
828 Adjusting Memory for Heterogeneous Concept Drift. In *IEEE International*
829 *Conference on Data Mining (ICDM)*, pages 291–300, 2016.
- 830 [92] Jie Lu, Anjin Liu, Fan Dong, Feng Gu, João Gama, and Guangquan Zhang.
831 Learning under Concept Drift: A Review. *IEEE Transactions on Knowledge*
832 *and Data Engineering*, 31:2346–2363, 2019.
- 833 [93] Yun Luo, Zhen Yang, Fandong Meng, Yafu Li, Jie Zhou, and Yue Zhang. An
834 Empirical Study of Catastrophic Forgetting in Large Language Models During
835 Continual Fine-tuning. *arXiv preprint arXiv:2308.08747*, 2025.
- 836 [94] Kleantes Malialis, Jin Li, Christos G Panayiotou, and Marios M Polycar-
837 pou. Incremental Learning with Concept Drift Detection and Prototype-based
838 Embeddings for Graph Stream Classification. In *International Joint Conference*
839 *on Neural Networks (IJCNN)*, pages 1–7. IEEE, 2024.
- 840 [95] Vasileios Manginas, Nikolaos Manginas, Edward Stevinson, Sherwin Vargh-
841 ese, Nikos Katzouris, Georgios Paliouras, and Alessio Lomuscio. A Scal-
842 able Approach to Probabilistic Neuro-Symbolic Verification. *arXiv preprint*
843 *arXiv:2502.03274*, 2025.
- 844 [96] Robin Manhaeve, Sebastijan Dumančić, Angelika Kimmig, Thomas Demeester,
845 and Luc De Raedt. Neural probabilistic logic programming in DeepProbLog.
846 *Artificial Intelligence*, 298:103504, 2021.
- 847 [97] Gary F Marcus. *The Algebraic Mind: Integrating Connectionism and Cognitive*
848 *Science*. MIT press, 2003.
- 849 [98] Douglas L Medin, William D Wattenmaker, and Sarah E Hampson. Fam-
850 ily resemblance, conceptual cohesiveness, and category construction. *Cognitive*
851 *Psychology*, 19(2):242–279, 1987.

- 852 [99] Jason S Metcalfe, Brandon S Perelman, David L Boothe, and Kaleb Mcdowell.
853 Systemic oversimplification limits the potential for human-AI partnership. *IEEE*
854 *Access*, 9:70242–70260, 2021.
- 855 [100] Tim Miller. Explanation in artificial intelligence: Insights from the social
856 sciences. *Artificial Intelligence*, 267:1–38, 2019.
- 857 [101] Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasser-
858 man, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru.
859 Model Cards for Model Reporting. In *Conference on Fairness, Accountability,*
860 *and Transparency*, pages 220–229, 2019.
- 861 [102] Stephen Muggleton and Luc De Raedt. Inductive logic programming: Theory
862 and methods. *The Journal of Logic Programming*, 19:629–679, 1994.
- 863 [103] Thomas Müller, Guillermo Pérez-Torró, and Marc Franco-Salvador. Few-Shot
864 Learning with Siamese Networks and Label Tuning. In *60th Annual Meeting of*
865 *the Association for Computational Linguistics (Volume 1: Long Papers)*, pages
866 8532–8545, 2022. Association for Computational Linguistics.
- 867 [104] Martin Mundt, Yongwon Hong, Iuliia Pliushch, and Visvanathan Ramesh. A
868 wholistic view of continual learning with deep neural networks: Forgotten lessons
869 and the bridge to active and open world learning. *Neural Networks*, 160:306–336,
870 2023.
- 871 [105] Meike Nauta, Jan Trienes, Shreyasi Pathak, Elisa Nguyen, Michelle Peters, Yas-
872 min Schmitt, Jörg Schlötterer, Maurice van Keulen, and Christin Seifert. From
873 Anecdotal Evidence to Quantitative Evaluation Methods: A Systematic Review
874 on Evaluating Explainable AI. *ACM Computing Surveys*, 55(13s), 2023.
- 875 [106] Weili Nie, Zhiding Yu, Lei Mao, Ankit B Patel, Yuke Zhu, and Anima Anand-
876 kumar. Bongard-logo: A new benchmark for human-level concept learning and
877 reasoning. *Advances in Neural Information Processing Systems*, 33:16468–16480,
878 2020.
- 879 [107] Ana Rita Nogueira, Andrea Pugnana, Salvatore Ruggieri, Dino Pedreschi, and
880 João Gama. Methods and tools for causal discovery and causal inference. *Wiley*
881 *interdisciplinary reviews: data mining and knowledge discovery*, 12(2):e1449,
882 2022.
- 883 [108] Robert M Nosofsky. Exemplar-based accounts of relations between classifica-
884 tion, recognition, and typicality. *Journal of Experimental Psychology: Learning,*
885 *Memory, and Cognition*, 14(4):700, 1988.
- 886 [109] Nicolas Papernot, Patrick McDaniel, Somesh Jha, Matt Fredrikson, Z. Berkay
887 Celik, and Ananthram Swami. The Limitations of Deep Learning in Adversarial
888 Settings. In *IEEE European Symposium on Security and Privacy (EuroS&P)*,

- 889 pages 372–387, 2016.
- 890 [110] Leif E Peterson. K-nearest neighbor. *Scholarpedia*, 4(2):1883, 2009.
- 891 [111] Mahima Pushkarna, Andrew Zaldivar, and Oddur Kjartansson. Data Cards:
892 Purposeful and Transparent Dataset Documentation for Responsible AI. In
893 *2022 ACM Conference on Fairness, Accountability, and Transparency*, pages
894 1776–1826, 2022.
- 895 [112] Eleanor Rosch. Natural categories. *Cognitive Psychology*, 4(3):328–350, 1973.
- 896 [113] Eleanor Rosch and Carolyn B Mervis. Family resemblances: Studies in the
897 internal structure of categories. *Cognitive Psychology*, 7(4):573–605, 1975.
- 898 [114] David E Rumelhart and Peter M Todd. Learning and connectionist representa-
899 tions. In *Attention and performance XIV: Synergies in experimental psychology,*
900 *artificial intelligence, and cognitive neuroscience*, pages 3–30. MIT Press, 1993.
- 901 [115] David E Rumelhart, James L McClelland, and PDP Research Group. *Parallel*
902 *distributed processing, volume 1: Explorations in the microstructure of cognition:*
903 *Foundations*. MIT press, 1986.
- 904 [116] Ute Schmid and Emanuel Kitzelmann. Inductive rule learning on the knowledge
905 level. *Cognitive Systems Research*, 12(3-4):237–248, 2011.
- 906 [117] Bernhard Schölkopf, Francesco Locatello, Stefan Bauer, Nan Rosemary Ke, Nal
907 Kalchbrenner, Anirudh Goyal, and Yoshua Bengio. Toward causal representation
908 learning. *Proceedings of the IEEE*, 109(5):612–634, 2021.
- 909 [118] Harshay Shah, Kaustav Tamuly, Aditi Raghunathan, Prateek Jain, and Pra-
910 neeth Netrapalli. The Pitfalls of Simplicity Bias in Neural Networks. *Neural*
911 *Information Processing Systems*, 2020.
- 912 [119] Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning:*
913 *From theory to algorithms*. Cambridge University Press, 2014.
- 914 [120] Roger N. Shepard. Toward a Universal Law of Generalization for Psychological
915 Science. *Science*, 237(4820):1317–1323, 1987.
- 916 [121] Zenglin Shi, Jie Jing, Ying Sun, Joo-Hwee Lim, and Mengmi Zhang. Unveil-
917 ing the Tapestry: The Interplay of Generalization and Forgetting in Continual
918 Learning. *IEEE Transactions on Neural Networks and Learning Systems*, 2025.
- 919 [122] Gunjan Singh, Riccardo Tommasini, Sumit Bhatia, and Raghava Mutharaju.
920 Benchmarking Neuro-Symbolic Description Logic Reasoners: Existing Chal-
921 lenges and A Way Forward. *Neurosymbolic Artificial Intelligence (to appear)*,
922 2025.

- 923 [123] Ji Y Son, Linda B Smith, and Robert L Goldstone. Simplicity and generalization:
924 Short-cutting abstraction in children’s object categorizations. *Cognition*, 108(3):
925 626–638, 2008.
- 926 [124] Zhivar Sourati, Vishnu Priya Prasanna Venkatesh, Darshan Deshpande, Himan-
927 shu Rawlani, Filip Ilievski, Hông-Ân Sandlin, and Alain Mermoud. Robust
928 and explainable identification of logical fallacies in natural language arguments.
929 *Knowledge-Based Systems*, 266:110418, 2023.
- 930 [125] Zhivar Sourati, Filip Ilievski, Pia Sommerauer, and Yifan Jiang. ARN: Analogi-
931 cal Reasoning on Narratives. *Transactions of the Association for Computational*
932 *Linguistics*, 12:1063–1086, 2024.
- 933 [126] Joshua B Tenenbaum, Charles Kemp, Thomas L Griffiths, and Noah D Good-
934 man. How to grow a mind: Statistics, structure, and abstraction. *Science*, 331
935 (6022):1279–1285, 2011.
- 936 [127] Ivan Y. Tyukin, Alexander N. Gorban, Muhammad H. Alkhudaydi, and Qinghua
937 Zhou. Demystification of Few-shot and One-shot Learning. In *International*
938 *Joint Conference on Neural Networks (IJCNN)*, pages 1–7, 2021.
- 939 [128] Michelle Vaccaro, Abdullah Almaatouq, and Thomas Malone. When combi-
940 nations of humans and AI are useful: A systematic review and meta-analysis.
941 *Nature Human Behaviour*, pages 1–11, 2024.
- 942 [129] Vladimir N. Vapnik. *The nature of statistical learning theory*. Springer-Verlag
943 New York, Inc., 1995.
- 944 [130] Vladimir N Vapnik and A Ya Chervonenkis. On the uniform convergence of
945 relative frequencies of events to their probabilities. In *Measures of complexity:*
946 *festschrift for alexey chervonenkis*, pages 11–30. Springer, 2015.
- 947 [131] Eli Verwimp, Rahaf Aljundi, Shai Ben-David, Matthias Bethge, Andrea Cossu,
948 Alexander Gepperth, Tyler L Hayes, Eyke Hüllermeier, Christopher Kanan,
949 Dhireesha Kudithipudi, et al. Continual learning: Applications and the road
950 forward. *Transactions on Machine Learning Research*, 2024.
- 951 [132] Michael R Waldmann, York Hagmayer, and Aaron P Blaisdell. Beyond the
952 information given: Causal models in learning and reasoning. *Current Directions*
953 *in Psychological Science*, 15(6):307–311, 2006.
- 954 [133] Cunxiang Wang, Xiaoze Liu, Yuanhao Yue, Xiangru Tang, Tianhang Zhang,
955 Cheng Jiayang, Yunzhi Yao, Wenyang Gao, Xuming Hu, Zehan Qi, et al.
956 Survey on factuality in large language models: Knowledge, retrieval and domain-
957 specificity. *arXiv preprint arXiv:2310.07521*, 2023.

- 958 [134] Nathaniel Weir, Bhavana Dalvi Mishra, Orion Weller, Oyvind Tafjord, Sam
959 Hornstein, Alexander Sabol, Peter Jansen, Benjamin Van Durme, and Peter
960 Clark. From Models to Microtheories: Distilling a Model’s Topical Knowledge
961 for Grounded Question Answering. *arXiv preprint arXiv:2412.17701*, 2024.
- 962 [135] Cara Leigh Widmer, Md Kamruzzaman Sarker, Srikanth Nadella, Joshua
963 Fiechter, Ion Juvina, Brandon Minnery, Pascal Hitzler, Joshua Schwartz, and
964 Michael Raymer. Towards Human-Compatible XAI: Explaining Data Differ-
965 entials with Concept Induction over Background Knowledge. *Journal of Web
966 Semantics*, 79:100807, 2023.
- 967 [136] Thaddäus Wiedemer, Jack Brady, Alexander Panfilov, Attila Juhos, Matthias
968 Bethge, and Wieland Brendel. Provable Compositional Generalization for
969 Object-Centric Learning. In *The International Conference on Learning Repre-
970 sentations*, 2024.
- 971 [137] Patrick H Winston. Learning structural descriptions from examples. Technical
972 report, MIT, AI Technical Reports, 1970.
- 973 [138] David H Wolpert and William G Macready. No free lunch theorems for
974 optimization. *IEEE Transactions on Evolutionary Computation*, 1(1):67–82,
975 1997.
- 976 [139] Jingyang Yang, Kaiyang Zhou, Yixuan Li, and Ziwei Liu. Generalized out-of-
977 distribution detection: A survey. *International Journal of Computer Vision*,
978 pages 1–28, 2024.
- 979 [140] Kaiyu Yang, Aidan Swope, Alex Gu, Rahul Chalamala, Peiyang Song, Shixing
980 Yu, Saad Godil, Ryan J Prenger, and Animashree Anandkumar. LeanDojo:
981 Theorem Proving with Retrieval-Augmented Language Models. *Advances in
982 Neural Information Processing Systems*, 36:21573–21612, 2023.
- 983 [141] Liuyi Yao, Zhixuan Chu, Sheng Li, Yaliang Li, Jing Gao, and Aidong Zhang. A
984 Survey on Causal Inference. *ACM Transactions on Knowledge Discovery from
985 Data (TKDD)*, 15(5), 2021.
- 986 [142] Haotian Ye, Chuanlong Xie, Tianle Cai, Ruichen Li, Zhenguo Li, and Liwei
987 Wang. Towards a Theoretical Framework of Out-of-Distribution Generalization.
988 In *Advances in Neural Information Processing Systems*, 2021.
- 989 [143] Yuan Yuan, Zhaojian Li, and Bin Zhao. A Survey of Multimodal Learning:
990 Methods, Applications, and Future. *ACM Computing Surveys*, 57(7), 2025.
- 991 [144] Thomas Zeugmann. Can Learning in the Limit Be Done Efficiently? In
992 *Algorithmic Learning Theory*, pages 17–38, 2003. Springer Berlin Heidelberg.

- ⁹⁹³ [145] Yu Zhang and Qiang Yang. A survey on multi-task learning. *IEEE Transactions*
⁹⁹⁴ *on Knowledge and Data Engineering*, 34(12):5586–5609, 2021.
- ⁹⁹⁵ [146] Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu
⁹⁹⁶ Zhu, Hui Xiong, and Qing He. A comprehensive survey on transfer learning.
⁹⁹⁷ *Proceedings of the IEEE*, 109(1):43–76, 2020.