
Meaning without reference in large language models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The widespread success of large language models (LLMs) has been met with
2 skepticism that they possess anything like human concepts or meanings. Contrary to
3 claims that LLMs possess no meaning whatsoever, we argue that they likely capture
4 important aspects of meaning, and moreover work in a way that approximates a
5 compelling account of human cognition in which meaning arises from *conceptual*
6 *role*. Because conceptual role is defined by the relationships between internal
7 representational states, meaning cannot be determined from a model’s architecture,
8 training data, or objective function, but only by examination of how its internal
9 states relate to each other. This approach may clarify why and how LLMs are so
10 successful and suggest how they can be made more human-like.

11 **L**ARGE language models (LLMs) have begun to display an array of competencies that were long
12 thought to be out of the reach of neural networks. At the same time, critics have been vocal
13 that existing methods, model structures, and training paradigms will be insufficient for anything
14 like human language use. In particular, it is argued that LLMs will never achieve “meaning” or
15 “understanding” due to either the objective function they optimize, the format of their internal
16 representations, or the type of training data that they receive. This short comment aims to contest the
17 claims that LLMs are necessarily incapable of acquiring meaning, and suggest that LLM meanings
18 may already be similar to human-like meaning in many (but not all) ways. We argue that LLMs
19 have likely already achieved some key aspects of meaning which, while imperfect, mirror the way
20 meanings work in cognitive theories, as well as approaches to the philosophy of language.

21 LLMs are trained to predict words from massive datasets of text from the internet. They typically
22 contain billions of parameters that are jointly optimized—at great computational, energy, and financial
23 expense [Bender et al., 2021]—to make predictions about word occurrence from surrounding context.
24 This setup differs from human language acquisition in data scale and format. On the other hand, the
25 core objective of word prediction is a central piece of human language processing [e.g. Altmann and
26 Kamide, 1999, Hale, 2001, Levy, 2008] and has long been shown capable of providing a learning
27 signal from which linguistic structures and semantic categories can emerge [Elman, 1990].

28 A key debate in recent literature has been about *what* models trained in this manner come to know. A
29 prominent view is that models trained only on text cannot acquire realistic meanings because they
30 lack reference, or connection to objects in the real world. Bender and Koller [2020] illustrate this
31 with an “octopus test.” They imagine an octopus that learns to use words correctly by eavesdropping
32 on a conversation between two people on land. If the octopus has no access to the referents of the
33 words then there are gaps in its meaning for the words. For example, if the octopus must suddenly
34 determine which *object* is a coconut, then its expertise in using the word “coconut” won’t help. Its
35 knowledge of co-occurrence statistics between “coconut” and other words won’t help either since it
36 also knows nothing about the referents of other words. The octopus simply does not have the required
37 meanings to find the coconut. Humans learners don’t have this problem because their input—and
38 consequently representations—are tied to real-world referents. Bender & Koller’s position is that no

39 amount of predictive linguistic savvy can give models that are trained on text alone the knowledge of
40 reference they need to acquire meanings.

41 **Meaning and reference** The octopus test assumes that reference determines meaning, but in fact
42 cognitive scientists and philosophers have found a variety of problems with this view. One is that
43 there are many terms that are meaningful to us but have no discernible referent at all, such as abstract
44 words like “justice” and “wit.” People can think of new concepts like “aphid-sized accordion” that
45 don’t exist (thus no referent), or even terms that have no *possible* referent like “perpetual motion
46 machine” or “imaginary cup of tea.” We can think of concepts like “king of San Francisco” that pick
47 out nobody, but are at least meaningful enough to reason about (for example: “If there was a King of
48 San Francisco, he’d live in The Presidio.”). Other examples show reference can be quite decoupled
49 from meaning. Terms like “treaty” and “contract” are often thought to have a concrete referent, but
50 what is important is actually an abstract entity: a treaty is still valid if the piece of paper is destroyed.
51 Frege’s example is of the “morning star” and the “evening star.” Both are terms for the planet Venus,
52 but were once conceived of as different entities without knowing that they are the same object.

53 This problem is not solely an issue with abstract concepts. Many concrete terms seem not to get
54 their most important semantic properties from reference either. Consider an example like a “postage
55 stamp.” Everyone can conjure up an image of a typical, physical, postage stamp. But with further
56 scrutiny, none of the concrete features of typical stamps seem necessary. We can imagine a country
57 whose postage stamps were made of clear glass, for example, or stamps that were microscopic so that
58 they could not be seen, or stamps that were paid for and tracked entirely online, or others that were
59 larger than a house (for mailing very large packages). We could think of postage stamps that were
60 *RFID* tags that went inside the envelope. Maybe intelligent ants would use postage stamps that were
61 pheromones rather than paper; in the future we might have postage stamps that sprout wings and fly
62 your letter to its intended location. Already we have stamps that only have barcodes and stamps that
63 you can draw yourself. You know the term “postage stamp” but there is no way that you could have
64 considered all of the possible referents yet, so reference cannot be what determines the concept.

65 Similar arguments were made by Wittgenstein [1953] when considering the concept of “water”,
66 leading him to conclude that reference plays little or no role in determining meaning in the general
67 case. One way to explain these observations is to assume that our meaning for terms like “postage
68 stamp” (or “water”) may be primarily determined by the role these concepts play in some greater
69 mental theory. Very roughly, people call something a “postage stamp” if you pay for it and then
70 attach it to a letter in order to have the letter to be delivered. In this view, the meaning of the word is
71 intrinsically intertwined with other concepts like “payment”, “letter” and “delivery.” The interrelation
72 is key [see, e.g. Deacon, 1998, Santoro et al., 2021] because when the associated terms shift in
73 meaning even slightly (e.g. credit cards were developed as wholly new a form of “payment”) the
74 meaning of “postage stamp” comes along for the ride (we know right away that they can be paid for
75 by a credit card). Such relationships between concepts are *the* essential, defining, aspects of meaning
76 and, in fact, possessing the appropriate relationships allows you to determine the reference. This view
77 makes sense of the puzzling examples above.

78 **Conceptual role theory** In philosophy of mind, versions of this approach go by the name “con-
79 ceptual role theory.” Following Block [1998], consider statements in physics like $f = m \cdot a$. This
80 equation is not exactly a definition of f (*force*), nor is it a definition of *mass* (m) or *acceleration* (a),
81 but, is a statement about the interrelationship of f , m , and a . Most of us can’t give much more detail
82 about the ultimate physical reference of these terms—forces have something to do with interacting
83 elementary particles (whatever those are), masses have something to do with the Higgs field (whatever
84 that is)—but our thoughts about $f = m \cdot a$ certainly do not seem meaningless. Many believe that
85 conceptual roles are one of the most promising ways to characterize human concepts *in general* [for
86 an overview of this and competing theories, see, e.g., Margolis and Laurence, 1999]. Murphy and
87 Medin [1985] for example argue that our organization of categories is based on entire theories of
88 structured conceptual domains (rather than simple features or similarities), an idea they trace to Quine
89 [1977]. Murphy and Medin note how we might reflexively consider a composite category such as
90 “prime numbers or apples” to be an unnatural or incoherent set of entities. However, if we know
91 someone called Wilma who is a number theorist grew up on an apple farm, then the category “topics
92 of conversation with Wilma”, involving the same constituent entities, seems perfectly reasonable.
93 What is a natural category or concept depends on our mental conception of how the underlying pieces
94 relate, and concepts can even be assembled fluidly, in an ad hoc manner or context-dependent manner

95 [Barsalou, 1983, Casasanto and Lupyan, 2015]. Theories in cognition have also been explored in
96 learning models [e.g. Goodman et al., 2011, Ullman et al., 2012] and experimentally shown to shape
97 how children explore the world [e.g. Gopnik et al., 1999, Gopnik and Schulz, 2004, Bonawitz et al.,
98 2012].

99 **Conceptual role in LLMs** If anything like this view is correct, then the search for meaning in
100 learning models—or brains—should focus on understanding the way that the systems’ internal
101 representational states *relate to each other*. Once a learning model finds it probable that “postage
102 stamps” are “affixed” to “letters” so they can be “delivered,” then it has acquired some important
103 pieces of conceptual role for these terms. It would not be possible to conclude anything about what
104 meanings a system does and does not possess from its training data or architecture because these may
105 not be informative about how the internal states relate to each other.

106 Relations between internal states have been long emphasized in cognitive theories [Shepard and
107 Chipman, 1970, Deacon, 1998, Fodor and Pylyshyn, 1988], for example early attempts to discover
108 the geometry of psychological space [Shepard, 1980] and more recent analyses of brain data based
109 on representational similarity [Kriegeskorte et al., 2008]. Elman [2004] argues for a closely related
110 view of the mental lexicon in which a word’s meaning is the effect it has on other mental states.
111 In deep learning models, the relational geometry of vector representations have been examined
112 for instance in analogy problems [Mikolov et al., 2013], match to human similarities [Hill et al.,
113 2017], and encoding of humanlike gradient distinctions [Vulić et al., 2017]. Grand et al. [2022]
114 show that semantic embeddings from these models capture gradient scales of multiple features, like
115 from “small” to “big” or “safe” to “dangerous.” It is even possible to align the word representations
116 acquired by text-based models across languages to translate between them effectively with no prior
117 knowledge of which words or phrases should have the same meaning [Lample et al., 2017]. Similarly,
118 Abdou et al. [2021] show that a model trained on text can recover key geometry of color space, even
119 without grounding; with a few examples of grounding, LLMs are able to align their structure with
120 the real grounded one, suggesting that they already possess the right relations Patel and Pavlick
121 [2021]. Importantly, as the performance of LLMs has improved in recent years the extent to which
122 their relational geometry reflects human data has also consistently increased [Peters et al., 2018,
123 Devlin et al., 2018, Brown et al., 2020]. Larger models also better reflect the human tendency for
124 semantic or mental models to influence formal reasoning behaviour [Wason and Johnson-Laird,
125 1972], on challenging logic problems that are not observed in their training data [Dasgupta et al.,
126 2022]. Moreover, this increasing correspondence between LLMs and human data is not observed
127 only behaviourally. Recent fMRI studies show that the semantic models that best account for the
128 representational geometry and processing activity of human brains are precisely the neural network
129 LLMs which are trained on the largest amount of data [Schrimpf et al., 2021, Goldstein et al., 2022,
130 Kumar et al., 2022].

131 Many of the tasks that LLMs succeed on are ones that require maintaining the right relationships
132 between concepts. Impressively, the largest models can now devise coherent narratives [Brown
133 et al., 2020], extend stories [Xu et al., 2020, Li et al., 2021], answer factual questions [Jiang et al.,
134 2021] solve Winograd Schema [Kocijan et al., 2020] and resolve complex quantitative reasoning
135 problems [Lewkowycz et al., 2022]. Increasingly, such models are even aware of the likelihood that
136 they can answer a given question correctly; i.e. they have an explicit sense of the extent of their
137 own knowledge [Kadavath et al., 2022]. Each of these capacities requires, in some way or another,
138 sensitivity to conceptual roles because the required words and concepts must be used jointly together
139 in a coherent way that mimics how humans would.

140 Despite these empirical successes, there are many places where these models can still be improved (for
141 detailed analyses, see, Lake and Murphy [2021], McClelland et al. [2020], Pavlick [2022]). Lake and
142 Murphy [2021] emphasize the need for reference, inference, better and more robust compositionality,
143 more structure and more consistent abstract reasoning. Models trained on multimodal datasets show
144 better match to human judgements than those trained on text alone [Hill et al., 2016, De Deyne et al.,
145 2021]; at the same time, even multimodal LLMs are missing many aspects of a complete theory of
146 semantics, including the ability to simulate situations in which their physical or linguistic behaviour
147 affects their environment [McClelland et al., 2020] as well as knowledge of the goals and desires that
148 drive how people use words [Bisk et al., 2020, Lake and Murphy, 2021].

149 Our claim, then, is not that LLMs perfectly capture human concepts or perfectly reflect human
150 meaning. Unlike the more radical perspectives entertained by Wittgenstein or Quine, we also do not

151 consider that reference should play *no* role in a principled treatment of meaning. Instead, we find it
152 productive to consider reference as just one (optional) aspect of a word’s full conceptual role. It is
153 relevant for some concepts [see Putnam, 1974] and not others—just like color, valence, or teleology is
154 relevant for some concepts and not others. Experience of both agency and a perceptual environment
155 similar to our own may lead to the richest, most human-like understanding of language in machines
156 [Bisk et al., 2020, McClelland et al., 2020].

157 As these improvements are made, the models will come into closer alignment with humans, and each
158 such improvement will enrich the model’s sense of meaning. This process of progressive enrichment
159 is also found in human concepts. When people discovered that H_2O was the chemical composition of
160 water, they grew their conceptual network and even revised their reference for the term. There was no
161 hard transition from a meaningless concept of “water” to a meaningful one. Some meaning was there
162 all along because “water” had a conceptual role even before its chemical composition was known.
163 What changed was the richness and interconnection of this concept—the way in which it was related
164 to other concepts like “hydrogen” and “oxygen.” In much the same way, we see no reason to assume
165 that the world of a system that receives input from a single sensory modality is meaningless, even if
166 the addition of further sensors provides clear enrichment. When thinking about improving LLMs it
167 we should therefore consider ways to enrich the internal conceptual roles of these systems, including
168 to better reflect the structure, inference and algorithmic sophistication of humans [Tenenbaum et al.,
169 2011, Lake and Murphy, 2021, Rule et al., 2020].

170 **Discovering conceptual role** Conceptual role theory also provides a compelling way to understand
171 learning, including the way in which neural networks may come to represent symbolic processes. A
172 symbol like *AND* only means logical conjunction if it interacts (composes) with others symbols like
173 *TRUE* and *FALSE* in the appropriate way—i.e. when it has the right conceptual role in the broader
174 system of symbols. The technique of *church-encoding* in mathematical logic [see Pierce, 2002]
175 provides a way to understand how such roles may be learned within neural networks or dynamical
176 systems Piantadosi [2021]. In church-encoding, a representation is constructed in one system (e.g.
177 lambda calculus or a neural network) in order to mimic the behavior of another system (e.g. boolean
178 logic) in the sense that the representations in the first system interact with each other in a way that
179 yields the desired conceptual roles of the second. Piantadosi [2021] shows how a church-encoding
180 learner could acquire structures like logic, lists, trees, hierarchies, numbers, quantifiers, and recursion
181 without possessing them to start, and how this metaphor may provide an “assembly language” that
182 translates from symbolic computational or cognitive theories into underlying implementations. In
183 this view, neural networks would train their parameters so that their internal, intrinsic dynamics
184 church-encode the conceptual roles of a targeted domain.

185 The key question for LLMs is whether training to predict text could actually support discovery of
186 conceptual roles. To us the question is empirical, and we believe has been answered in a promising,
187 partial affirmative by studies showing success on tasks that require knowledge of relationships
188 between concepts. Text provides such clues to conceptual role because human conceptual roles
189 generated the text. Analogously, it is possible to build a theory of gravity from measurements of the
190 moon’s movement because gravity *generated* these movements; the goal of essentially all inductive
191 learning techniques is to invert from observations to likely generating processes or parameters.
192 Moreover, the entailment relationships between sentences are often intrinsically related to analogous
193 patterns in thought [e.g. Fodor and Pylyshyn, 1988], meaning that a model which captures how
194 sentences relate to each other might indeed capture how thoughts relate to each other. One helpful
195 analogy is that of *embedding theorems* in dynamical systems [e.g. Packard et al., 1980, Takens, 1981]
196 which allow some properties of systems to be recovered from seemingly impoverished representations
197 of their state. In the paper “Geometry from a Time Series”, Packard et al. [1980] show, remarkably,
198 that one can sometimes reconstruct the geometry of a multi-dimensional dynamical system from a
199 *one-dimensional* projection of its state. Thus, information about high-dimensional state (sometimes
200 essentially all of it) can be decoded from the trajectory of low-dimensional projection. People use
201 concepts in thinking and reasoning based on their meaning, and text is a low-dimensional projection
202 of some of these patterns of use, so it is plausible that some properties of the real meaning could be
203 inferred from text. At the very least, embedding theorems illustrate that there may not be a simple
204 way to intuit what a learner can or cannot deduce about the underlying mental states from text alone.

205 The protein folding neural network AlphaFold [Jumper et al., 2021] provides further evidence that
206 transformer-based networks can infer and generalise complex latent multi-entity structures. AlphaFold

207 is trained to predict the configuration of single proteins only, but acquires actionable knowledge of
208 concepts not explicitly present in its training data. Despite never seeing a zinc ion, AlphaFold often
209 perfectly infers its location and places all the protein side chains correctly right around it. Some
210 proteins only fold with multiple copies of themselves (homomers). Again, AlphaFold has never seen
211 more than one copy of a protein, but often infers both the number of required copies and their relative
212 placement correctly. This suggests that the process of predicting the structure of single proteins
213 enabled the AlphaFold network to infer non-trivial facts about chemistry and biology.

214 **Communicative intentions** Separable from meaning and reference, many have also rejected the
215 idea that LLMs produce language with *intention* [Bender et al., 2021, Bender and Koller, 2020].
216 Because LLMs are trained only on sequence prediction, they are argued to be, “stochastic parrots”
217 Bender et al. [2021] or just sophisticated “auto-complete” algorithms¹ that cannot access the intentions
218 of those who produced their training data, and themselves produce language without intending
219 anything in particular. But in our view, a key difference between autocompleting parrots and LLMs is
220 that the latter have rich, causal, and structured internal states.

221 One view of intent is semantic, corresponding to whether the language they produce arises from an
222 internal representation of (intended) meaning. The conversion of internal states into language and
223 back *is* the essential function of LLMs, embedded in their architecture and training. We have argued
224 that LLM’s internal state has some notions of conceptual role, so LLM’s utterances have the semantic
225 intent corresponding to these roles.

226 Another view of intent is pragmatic, asking what might be achieved by producing a sentence.
227 This corresponds to asking whether they engage in any goal-directed *planning* [Russell, 2010] when
228 producing language. We consider it probable that multi-layer LLMs do, in an emergent, implicit sense,
229 execute a form of planning as part of the process of repeated (self-attention-based) analyses of current
230 and past inputs. These computations likely involve representation of the current situation and at least
231 implicit evaluation of consequences of utterances. Recent work has argued that LLMs possess model-
232 like belief structures Hase et al. [2021] and update representations of dynamic semantics, objects
233 and situations, throughout a discourse Li et al. [2021]. These emergent semantics causally determine
234 LLM output. Of course, differences between humans and LLMs in their training experience and
235 objectives mean that the planning process in LLMs is less explicit and sophisticated than in humans
236 (making errors more likely, for instance, in cases of hypothetical or counterfactual reasoning [Ortega
237 et al., 2021]).

238 **Conclusion** Bender & Koller argue that text-based LLMs will never have meaning because these
239 models lack reference. However, they do not demonstrate that reference is the key to meaning—
240 instead they assume it. As we have argued, this assumption is hard to reconcile with theories of
241 cognition and the phenomena that motivate them. People are happy to think about concepts without
242 referents and otherwise often don’t know many details of reference. Meaning instead seems to come
243 from the way concepts relate to *each other*. It is these interrelations that LLMs know something
244 about since their internal geometries and trajectories approximate those of humans. Like people who
245 don’t know that water is H_2O and so could not pick it out based on chemical composition, Bender
246 & Koller’s octopus lacks some aspects of conceptual role like physical appearance. But, both the
247 octopus and people know other parts of conceptual role that are sophisticated in their own right. If
248 theories about conceptual role are the correct account, then LLMs likely already share the foundation
249 of how our own concepts get their meaning.

250 References

- 251 M. Abdou, A. Kulmizev, D. Hershcovich, S. Frank, E. Pavlick, and A. Søgaard. Can language
252 models encode perceptual structure without grounding? A case study in color. *arXiv preprint*
253 *arXiv:2109.06129*, 2021.
- 254 G. T. Altmann and Y. Kamide. Incremental interpretation at verbs: Restricting the domain of
255 subsequent reference. *Cognition*, 73(3):247–264, 1999.
- 256 L. W. Barsalou. Ad hoc categories. *Memory & cognition*, 11(3):211–227, 1983.

¹<https://garymarcus.substack.com/p/nonsense-on-stilts>

- 257 E. M. Bender and A. Koller. Climbing towards NLU: On meaning, form, and understanding in the
258 age of data. In *Proceedings of the 58th Annual Meeting of the Association for Computational*
259 *Linguistics*, pages 5185–5198, Online, July 2020. Association for Computational Linguistics. doi:
260 10.18653/v1/2020.acl-main.463. URL <https://aclanthology.org/2020.acl-main.463>.
- 261 E. M. Bender, T. Gebru, A. McMillan-Major, and S. Shmitchell. On the dangers of stochastic parrots:
262 Can language models be too big? In *Proceedings of the 2021 ACM Conference on Fairness,*
263 *Accountability, and Transparency*, pages 610–623, 2021.
- 264 Y. Bisk, A. Holtzman, J. Thomason, J. Andreas, Y. Bengio, J. Chai, M. Lapata, A. Lazaridou, J. May,
265 A. Nisnevich, et al. Experience grounds language. *arXiv preprint arXiv:2004.10151*, 2020.
- 266 N. Block. Semantics, conceptual role. *Routledge encyclopedia of philosophy*, 8:652–657, 1998.
- 267 E. B. Bonawitz, T. J. van Schijndel, D. Friel, and L. Schulz. Children balance theories and evidence
268 in exploration, explanation, and learning. *Cognitive psychology*, 64(4):215–234, 2012.
- 269 T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam,
270 G. Sastry, A. Askell, et al. Language models are few-shot learners. *Advances in neural information*
271 *processing systems*, 33:1877–1901, 2020.
- 272 D. Casasanto and G. Lupyan. All Concepts Are Ad Hoc Concepts. In E. Margolis and S. Laurence,
273 editors, *The Conceptual Mind: New directions in the study of concepts*, pages 543–566. Cambridge:
274 MIT Press, 2015.
- 275 I. Dasgupta, A. K. Lampinen, S. C. Chan, A. Creswell, D. Kumaran, J. L. McClelland, and F. Hill.
276 Language models show human-like content effects on reasoning. *arXiv preprint arXiv:2207.07051*,
277 2022.
- 278 S. De Deyne, D. J. Navarro, G. Collell, and A. Perfors. Visual and affective multimodal models of
279 word meaning in language and mind. *Cognitive Science*, 45(1):e12922, 2021.
- 280 T. W. Deacon. *The symbolic species: The co-evolution of language and the brain*. Number 202. WW
281 Norton & Company, 1998.
- 282 J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional
283 transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- 284 J. L. Elman. Finding structure in time. *Cognitive science*, 14(2):179–211, 1990.
- 285 J. L. Elman. An alternative view of the mental lexicon. *Trends in cognitive sciences*, 8(7):301–306,
286 2004.
- 287 J. A. Fodor and Z. W. Pylyshyn. Connectionism and cognitive architecture: A critical analysis.
288 *Cognition*, 28(1-2):3–71, 1988.
- 289 A. Goldstein, Z. Zada, E. Buchnik, M. Schain, A. Price, B. Aubrey, S. A. Nastase, A. Feder,
290 D. Emanuel, A. Cohen, et al. Shared computational principles for language processing in humans
291 and deep language models. *Nature neuroscience*, 25(3):369–380, 2022.
- 292 N. D. Goodman, T. D. Ullman, and J. B. Tenenbaum. Learning a theory of causality. *Psychological*
293 *review*, 118(1):110, 2011.
- 294 A. Gopnik and L. Schulz. Mechanisms of theory formation in young children. *Trends in cognitive*
295 *sciences*, 8(8):371–377, 2004.
- 296 A. Gopnik, A. N. Meltzoff, and P. K. Kuhl. *The scientist in the crib: Minds, brains, and how children*
297 *learn*. William Morrow & Co, 1999.
- 298 G. Grand, I. A. Blank, F. Pereira, and E. Fedorenko. Semantic projection recovers rich human
299 knowledge of multiple object features from word embeddings. *Nature Human Behaviour*, pages
300 1–13, 2022.
- 301 J. Hale. A probabilistic earley parser as a psycholinguistic model. In *Second meeting of the north*
302 *american chapter of the association for computational linguistics*, 2001.

- 303 P. Hase, M. Diab, A. Celikyilmaz, X. Li, Z. Kozareva, V. Stoyanov, M. Bansal, and S. Iyer. Do
304 language models have beliefs? methods for detecting, updating, and visualizing model beliefs.
305 *arXiv preprint arXiv:2111.13654*, 2021.
- 306 F. Hill, K. Cho, and A. Korhonen. Learning distributed representations of sentences from unlabelled
307 data. *arXiv preprint arXiv:1602.03483*, 2016.
- 308 F. Hill, K. Cho, S. Jean, and Y. Bengio. The representational geometry of word meanings acquired by
309 neural machine translation models. *Machine Translation*, 31(1):3–18, 2017.
- 310 Z. Jiang, J. Araki, H. Ding, and G. Neubig. How can we know when language models know? on
311 the calibration of language models for question answering. *Transactions of the Association for*
312 *Computational Linguistics*, 9:962–977, 2021.
- 313 J. Jumper, R. Evans, A. Pritzel, T. Green, M. Figurnov, O. Ronneberger, K. Tunyasuvunakool,
314 R. Bates, A. Žídek, A. Potapenko, et al. Highly accurate protein structure prediction with alphafold.
315 *Nature*, 596(7873):583–589, 2021.
- 316 S. Kadavath, T. Conerly, A. Askell, T. Henighan, D. Drain, E. Perez, N. Schiefer, Z. H. Dodds,
317 N. DasSarma, E. Tran-Johnson, S. Johnston, S. El-Showk, A. Jones, N. Elhage, T. Hume, A. Chen,
318 Y. Bai, S. Bowman, S. Fort, D. Ganguli, D. Hernandez, J. Jacobson, J. Kernion, S. Kravec,
319 L. Lovitt, K. Ndousse, C. Olsson, S. Ringer, D. Amodei, T. Brown, J. Clark, N. Joseph, B. Mann,
320 S. McCandlish, C. Olah, and J. Kaplan. Language models (mostly) know what they know, 2022.
321 URL <https://arxiv.org/abs/2207.05221>.
- 322 V. Kocijan, T. Lukasiewicz, E. Davis, G. Marcus, and L. Morgenstern. A review of winograd schema
323 challenge datasets and approaches. *arXiv preprint arXiv:2004.13831*, 2020.
- 324 N. Kriegeskorte, M. Mur, and P. A. Bandettini. Representational similarity analysis-connecting the
325 branches of systems neuroscience. *Frontiers in systems neuroscience*, page 4, 2008.
- 326 S. Kumar, T. R. Sumers, T. Yamakoshi, A. Goldstein, U. Hasson, K. A. Norman, T. L. Griffiths, R. D.
327 Hawkins, and S. A. Nastase. Reconstructing the cascade of language processing in the brain using
328 the internal computations of a transformer-based language model. *bioRxiv*, 2022.
- 329 B. M. Lake and G. L. Murphy. Word meaning in minds and machines. *Psychological review*, 2021.
- 330 G. Lample, A. Conneau, L. Denoyer, and M. Ranzato. Unsupervised machine translation using
331 monolingual corpora only. *arXiv preprint arXiv:1711.00043*, 2017.
- 332 R. Levy. Expectation-based syntactic comprehension. *Cognition*, 106(3):1126–1177, 2008.
- 333 A. Lewkowycz, A. Andreassen, D. Dohan, E. Dyer, H. Michalewski, V. Ramasesh, A. Slone, C. Anil,
334 I. Schlag, T. Gutman-Solo, et al. Solving quantitative reasoning problems with language models.
335 *arXiv preprint arXiv:2206.14858*, 2022.
- 336 B. Z. Li, M. Nye, and J. Andreas. Implicit representations of meaning in neural language models.
337 *arXiv preprint arXiv:2106.00737*, 2021.
- 338 E. E. Margolis and S. E. Laurence. *Concepts: Core Readings*. The MIT Press, 1999.
- 339 J. L. McClelland, F. Hill, M. Rudolph, J. Baldridge, and H. Schütze. Placing language in an integrated
340 understanding system: Next steps toward human-level performance in neural language models.
341 *Proceedings of the National Academy of Sciences*, 117(42):25966–25974, 2020.
- 342 T. Mikolov, W.-t. Yih, and G. Zweig. Linguistic regularities in continuous space word representations.
343 In *Proceedings of the 2013 conference of the north american chapter of the association for*
344 *computational linguistics: Human language technologies*, pages 746–751, 2013.
- 345 G. L. Murphy and D. L. Medin. The role of theories in conceptual coherence. *Psychological review*,
346 92(3):289, 1985.
- 347 P. A. Ortega, M. Kunesch, G. Delétang, T. Genewein, J. Grau-Moya, J. Veness, J. Buchli, J. Degraeve,
348 B. Piot, J. Perolat, et al. Shaking the foundations: delusions in sequence models for interaction and
349 control. *arXiv preprint arXiv:2110.10819*, 2021.

- 350 N. H. Packard, J. P. Crutchfield, J. D. Farmer, and R. S. Shaw. Geometry from a time series. *Physical*
351 *review letters*, 45(9):712, 1980.
- 352 R. Patel and E. Pavlick. Mapping language models to grounded conceptual spaces. In *International*
353 *Conference on Learning Representations*, 2021.
- 354 E. Pavlick. Semantic structure in deep learning. *Annual Review of Linguistics*, 8:447–471, 2022.
- 355 M. Peters, M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee, and L. Zettlemoyer. Deep
356 contextualized word representations. arxiv 2018. *arXiv preprint arXiv:1802.05365*, 12, 2018.
- 357 S. T. Piantadosi. The computational origin of representation. *Minds and machines*, 31(1):1–58, 2021.
- 358 B. C. Pierce. *Types and programming languages*. MIT press, 2002.
- 359 H. Putnam. Meaning and reference. *The journal of philosophy*, 70(19):699–711, 1974.
- 360 W. Quine. Natural kinds. In S. Schwartz, editor, *Naming, necessity, and natural kinds*, pages 155–175.
361 Ithaca, NY: Cornell University Press, 1977.
- 362 J. S. Rule, J. B. Tenenbaum, and S. T. Piantadosi. The child as hacker. *Trends in cognitive sciences*,
363 24(11):900–915, 2020.
- 364 S. J. Russell. *Artificial intelligence a modern approach*. Pearson Education, Inc., 2010.
- 365 A. Santoro, A. Lampinen, K. Mathewson, T. Lillicrap, and D. Raposo. Symbolic behaviour in
366 artificial intelligence. *arXiv preprint arXiv:2102.03406*, 2021.
- 367 M. Schrimpf, I. A. Blank, G. Tuckute, C. Kauf, E. A. Hosseini, N. Kanwisher, J. B. Tenenbaum, and
368 E. Fedorenko. The neural architecture of language: Integrative modeling converges on predictive
369 processing. *Proceedings of the National Academy of Sciences*, 118(45):e2105646118, 2021.
- 370 R. N. Shepard. Multidimensional scaling, tree-fitting, and clustering. *Science*, 210(4468):390–398,
371 1980.
- 372 R. N. Shepard and S. Chipman. Second-order isomorphism of internal representations: Shapes of
373 states. *Cognitive psychology*, 1(1):1–17, 1970.
- 374 F. Takens. Detecting strange attractors in turbulence. In *Dynamical systems and turbulence, Warwick*
375 *1980*, pages 366–381. Springer, 1981.
- 376 J. B. Tenenbaum, C. Kemp, T. L. Griffiths, and N. D. Goodman. How to grow a mind: Statistics,
377 structure, and abstraction. *science*, 331(6022):1279–1285, 2011.
- 378 T. D. Ullman, N. D. Goodman, and J. B. Tenenbaum. Theory learning as stochastic search in the
379 language of thought. *Cognitive Development*, 27(4):455–480, 2012.
- 380 I. Vulić, D. Gerz, D. Kiela, F. Hill, and A. Korhonen. Hyperlex: A large-scale evaluation of graded
381 lexical entailment. *Computational Linguistics*, 43(4):781–835, 2017.
- 382 P. C. Wason and P. N. Johnson-Laird. *Psychology of reasoning: Structure and content*, volume 86.
383 Harvard University Press, 1972.
- 384 L. Wittgenstein. *Philosophical investigations*. John Wiley & Sons, 1953.
- 385 P. Xu, M. Patwary, M. Shoeybi, R. Puri, P. Fung, A. Anandkumar, and B. Catanzaro. Megatron-cntrl:
386 Controllable story generation with external knowledge using large-scale language models. *arXiv*
387 *preprint arXiv:2010.00840*, 2020.