Climbing towards NLU: On Meaning, Form, and Understanding in the Age of Data

Accepted for ACL 2020 Emily M. Bender & Alexander Koller

Abstract

The success of the large neural language models on many NLP tasks is exciting. However, we find that that these successes sometimes lead to hype in which these models are being described as "understanding" language or capturing "meaning". In this position paper, we argue that a system trained only on form has *a priori* no way to learn meaning. In keeping with the ACL 2020 theme of "Taking Stock of Where We've Been and Where We're Going", we argue that a clear understanding of the distinction between form and meaning will help guide the field towards better science around natural language understanding.

1 Introduction

The current state of affairs in NLP is that the large neural language models (LMs), such as BERT (Devlin et al., 2019) or GPT-2 (Radford et al., 2019), are making great progress on a wide range of tasks, including those that are ostensibly meaningsensitive. This has led to claims, in both academic and popular publications, that such models "understand" or "comprehend" natural language or learn its "meaning". From our perspective, these are overclaims caused by a misunderstanding of the relationship between linguistic form and meaning.

We argue that the language modeling task, because it only uses form as training data, cannot in principle lead to learning of meaning. Our aim is to advocate for an alignment of claims and methodology: We take human-analogous natural language understanding (NLU) to be a grand challenge of artificial intelligence, which involves mastery of the structure and use of language and the ability to ground it in the world. While large neural LMs may well end up being useful components of an eventual full-scale solution to human-analogous NLU, it is disingenuous to speak of them as though they are nearly-there solutions to this grand challenge. We take the term *language model* to refer to any system trained only on the task of string prediction, whether it operates over characters, words or sentences, and sequentially or not. We take (linguistic) *meaning* to be the relation between a linguistic form and communicative intent.

We begin by briefly reviewing the ways in which large LMs are spoken about and summarizing the recent flowering of 'BERTology' papers (§2). We follow this with working definitions for "meaning" (§3), before presenting a series of thought experiments which serve to illustrate the impossibility of learning meaning when it is not in the training signal (§4,5). In §6, we turn to the human language acquisition literature for insight into what information humans use to bootstrap language learning. §7 presents reflections on what this issue means for how we look at progress and direct research effort in our field, and in §8, we address possible counterarguments to our main thesis.

2 Large LMs: Hype and Analysis

Publications (academic and more public-facing) talking about the application of large LMs to apparently meaning-sensitive tasks tend to describe the models with terminology that, if interpreted at face value, is misleading. Here is a selection from academically-oriented pieces (emphasis added):

- (1) In order to train a model that **understands** sentence relationships, we pre-train for a binarized next sentence prediction task (Devlin et al., 2019)
- Using BERT, a pretraining language model, has been successful for single-turn machine comprehension ... (Ohsugi et al., 2019)
- (3) The surprisingly strong ability of these models to recall factual knowledge without any fine-tuning demonstrates their potential as unsupervised open-domain QA systems. (Petroni et al., 2019)

If the highlighted terms are meant to describe human-analogous understanding, comprehension, or recall of factual knowledge, then these are gross overclaims. If, instead, they are intended as technical terms, they should be explicitly defined.

One important consequence of imprudent use of terminology in our academic discourse is that it feeds AI hype in the popular press. As NLP gains public exposure and is more widely used in applied contexts, it is more important than ever that the actual capabilities of our systems be accurately represented. In some cases, NLP experts speaking with the media are being appropriately careful, as in these two quotes from Gary Marcus (4) and Oren Etzioni (5) in the *New York Times*:¹

- (4) These systems are still a really long way from truly understanding running prose
- (5) Though BERT passed the labs common-sense test, machines are still a long way from an artificial version of a human's common sense.

However, there are plenty of instances where the popular press gets it wrong, such as (6) from the B2C website,² apparently based on the Google Blog post about BERT and search, which includes numerous statements like (7).³

- (6) BERT is a system by which Google's algorithm uses pattern recognition to better **understand** how human beings communicate so that it can return more relevant results for users.
- (7) Here are some of the examples that showed up our evaluation process that demonstrate BERTs ability to understand the intent behind your search.

In sum, it is not clear from our academic literature whether all authors are clear on the distinction between form and meaning, but it is clear that the way we speak about what neural LMs are doing is misleading to the public.

Part of the reason for this tendency to use imprecise language may well be that we do not yet fully understand what exactly it is about language that the large LMs come to implicitly represent. Their success, however, has sparked a subfield (sometimes dubbed 'BERTology') that aims to answer this question. Using the methodology of probing tasks (e.g. Adi et al., 2016; Ettinger et al., 2018), it has been shown that large LMs learn at least some information about phenomena such as English subject-verb agreement (Goldberg, 2019; Jawahar et al., 2019), constituent types, dependency labels, NER, and (core) semantic role types (again, all in English) (Tenney et al., 2019).⁴ Hewitt and Manning (2019) find information analogous to unlabeled dependency structures in the word vectors provided by ELMo and BERT (trained on English). And of course it is by now unsurprising that vector-space representations of words pick up word classes, both syntactic (POS, e.g. Lin et al., 2015) and semantic (lexical similarity, e.g. Rubenstein and Goodenough, 1965; Mikolov et al., 2013).

Others have looked more closely at the success of the large LMs on apparently meaning sensitive tasks and found that in fact, far from doing the 'reasoning' ostensibly required to complete the tasks, they were instead simply more effective at leveraging artifacts in the data than previous approaches. Niven and Kao (2019) find that BERT's unreasonably good performance on the English Argument Reasoning Comprehension Task (Habernal et al., 2018) falls back to chance if the dataset is modified by adding adversarial examples that just negate one piece of the original, thus mirroring the distribution of lexical cues for each label. Similarly, McCoy et al. (2019) find that BERT's performance on the English Multi-genre Natural Language Inference dataset (Williams et al., 2018) is predicated on its ability to leverage syntactic heuristics involving overlap (of full constituents, subsequences, or simply bags of words). In a dataset carefully designed to frustrate such heuristics, BERT's performance falls to significantly below chance.

This overview of BERTology papers has been of necessity brief and incomplete, but it highlights the extent to which there is evidence that large LMs can learn aspects of linguistic formal structure (e.g. agreement, dependency structure), but that their apparent ability to 'reason' is sometimes a mirage built on leveraging artifacts in the training data (i.e. form, not meaning). In this paper, we aim to complement these studies by arguing on theoretical grounds that a system exposed only to form in its training cannot in principle learn meaning.

3 What is meaning?

We start by defining two key terms: We take *form* to be any observable realization of language: marks on a page, pixels or bytes in a digital representation

¹https://www.nytimes.com/2018/11/18/te chnology/artificial-intelligence-languag e.html, accessed 2019/12/04

²https://www.business2community.com/s eo/what-to-do-about-bert-googles-recen t-local-algorithm-update-02259261, accessed 2019/12/04

³https://www.blog.google/products/sea rch/search-language-understanding-bert/, accessed 2019/12/04

⁴But see Warstadt et al.'s (2019) cautionary note about how the methodology used for probing can influence the results.

of text, or movements of the articulators.⁵ We take *meaning* to be the relation between the form and something external to language, in a sense that we will make precise below.

3.1 Meaning and Communicative Intent

When humans use language, we do so for a purpose: We do not talk for the joy of moving our articulators, but in order to achieve some *communicative intent*. There are many types of communicative intents: they may be to convey some information to the other person; or to ask them to do something; or simply to socialize. We take *meaning* to be the relation $M \subseteq E \times I$ which contains pairs (e, i) of natural language expressions e and the communicative intents they can be used to evoke.

Communicative intents are about something that is *outside of language*. When we say *Open the window!* or *When was Malala Yousafzai born?*, the communicative intent is grounded in the real world the speaker and listener inhabit together. Communicative intents can also be about abstract worlds, e.g. bank accounts, computer file systems, or a purely hypothetical world in the speaker's mind.

Linguists distinguish communicative intent from conventional (or standing) meaning (Quine, 1960; Grice, 1968). The conventional meaning of an expression (word, phrase, sentence) is what is constant across all of its possible contexts of use. Conventional meaning is an abstract object that represents the communicative potential of a form, given the linguistic system it is drawn from. Each linguistic system (say, English) provides a relation $C \subseteq E \times S$, which contains pairs (e, s) of expressions e and their conventional meanings s.⁶ The field of linguistic semantics provides many competing theories of what conventional meanings s look like. For our purposes, we don't need to select among these theories; all we assume is that conventional meanings must have interpretations, such as a means of testing them for truth against a model of the world. Thus, like the meaning relation M, Cconnects language to objects outside of language.

Returning to the meaning relation M from above, it is best understood as mediated by the relation Cof a linguistic system shared between two interlocutors. The speaker has a certain communicative intent i, and chooses an expression e with a standing meaning s which is fit to express i in the current communicative situation. Upon hearing e, the listener then reconstructs s and uses their own knowledge of the communicative situation and their hypotheses about the speaker's state of mind and intention in an attempt to deduce i.

This active participation of the listener is crucial to human communication (Reddy, 1979; Clark, 1996). For example, to make sense of (8) and (9) (from Clark, 1996, p.144), the listener has to calculate that *Napoleon* refers to a specific pose (hand inside coat flap) or that *China trip* refers to a person who has recently traveled to China.

- (8) The photographer asked me to do a Napoleon for the camera.
- (9) Never ask two China trips to the same party.

We humans are also very willing, as we will see in §4 below, to attribute communicative intent given a linguistic signal belonging to a language we speak, even if the originator of the signal is not an entity that could have communicative intent.

To summarize, as we strive to understand how NLU tasks and system performance on those tasks relates to the bigger picture goals of building human-analogous natural language understanding systems, it is useful to distinguish cleanly between form, conventional meaning, and communicative intent. Furthermore, we should be careful not to confuse communicative intent with ground truth about the world, as speakers can of course be mistaken, be intentionally dissembling, etc.

3.2 Language models do not capture meaning

We argue that a model of natural language that is trained purely on form will not learn meaning: if the training data for is only form, there is not sufficient signal to learn relation M between that form and the non-linguistic intent of human language users, nor C between form and the standing meaning the linguistic system assigns to each form.

A system which does not capture the true meaning relation of a language, such as an LM, will fail a competently administered Turing test (Turing, 1950); at some point, the evaluator will uncover a mismatch between the system's language use and one that is grounded in our shared reality via the meaning relation. There are many anecdotes, going all the way back to ELIZA (Weizenbaum, 1966), of systems that managed to fool their users into

⁵In spoken languages, the primary articulators are the components of the vocal tract. In signed languages, they are principally the hands and face.

⁶We abstract away here from the facts that linguistic systems C change over time and are only incomplete shared among different speakers. They are stable enough to function as rich signals to communicative intent.

believing they were human (Saygin et al., 2000). But these anecdotes involve situations in which the evaluator was unaware that they might be talking to a computer, and are thus just evidence for our earlier point that human listeners are willing to invest effort into interpreting a conversation as natural.⁷

4 A Turing test

In order to illustrate the challenges in attempting to learn meaning from form alone, we propose a concrete scenario. Say that A and B, both fluent speakers of English, are independently stranded on two uninhabited islands. They soon discover that previous visitors to these islands have left behind telegraphs and that they can communicate with each other via an underwater cable. A and B start happily typing messages to each other.

Meanwhile, O, a hyper-intelligent deep-sea octopus who is unable to visit or observe the two islands, discovers a way to tap into the underwater cable and listen in on A and B's conversations. O knows nothing about English initially, but is very good at detecting and interpreting statistical patterns. Over time, O learns to predict with great accuracy how B will respond to each of A's utterances. O also learns that certain words tend to occur in similar contexts, and therefore may refer to similar objects or activities in the world. Nonetheless, O has never observed these objects, and thus would not be able to pick out the referent of a word when presented with a set of (physical) alternatives.

At some point, O starts feeling lonely. He cuts the underwater cable and inserts himself into the conversation, by pretending to be B and replying to A's messages. Can O successfully pose as B without making A suspicious? This constitutes a weak form of the Turing test (weak because A has no reason to suspect she is talking to a nonhuman); the interesting question is whether O fails it because he has not learned the meaning relation, having seen only the form of A and B's utterances.

The extent to which O can do this depends on the task — that is, on what A is trying to talk about. A and B have spent a lot of time exchanging trivial notes about their daily lives to make the long island evenings more enjoyable. It seems reasonable that O would be able to produce new sentences of the kind B used to produce; essentially acting as a chatbot. This is because the utterances in such conversations have a primarily social function, and do not need to be grounded in the particulars of the interlocutors' actual physical situation nor anything else specific about the real world. It is sufficient to produce text that is internally coherent.

Now say that A has invented a new device, say a coconut catapult. She excitedly sends detailed instructions on building a coconut catapult to B, and asks about B's experiences and suggestions for improvements. Even if O had a way of constructing the catapult underwater, he does not know what words such as *rope* and *coconut* refer to; he has seen them used in similar contexts as the words nail and mango, so perhaps they mean similar things, but this is not enough to predict a good description of a successful catapult test. When O previously saw A utter such words, B frequently responded with cool idea, great job, so this is what O chooses to say. It is absolutely conceivable that A accepts this reply as meaningful - but only because A does all the work in attributing meaning to O's response. It is not because O understood the meaning of A's instructions or even his own reply.

Finally, A faces an emergency. She is suddenly pursued by an angry bear. She grabs a couple of sticks and frantically asks B to come up with a way to construct a weapon to defend herself. Of course, O has no idea what A "means". Solving a task like this requires the ability to map accurately between words and real-world entities as well as reasoning and creative thinking. It is at this point that O would fail the Turing test, if A hadn't been eaten by the bear before noticing the deception.⁸

Having only form available as training data, O did not learn meaning. The language exchanged by A and B is a projection of their communicative intents through the meaning relation into linguistic forms. Without access to a means of at least hypothesizing the underlying communicative intents, reconstructing them from the forms alone is hopeless, and O's language use will eventually diverge from the language use of an agent who can ground

⁷Our argument also differs from Searle's Chinese Room experiment (1980): Searle argues that a computer passing the Turing test still need not have consciousness; we are concerned with what a system would need to pass the Turing test at all. Specifically, we claim that an LM alone won't be able to.

⁸To see what a large LM might reply in this situation, we prompted the GPT-2 demo with "*Help! I'm being chased by a bear! All I have is these sticks. What should I do?*", and GPT-2 to supplied "*You're not going to get away with this!*" (http s://gpt2.apps.allenai.org/, accessed 2019/12/4). Following Radford et al.'s (2019) approach of giving explicit cues to task required, we also constructed a more elaborate prompt. The results, given in the supplementary material, are highly entertaining but no more helpful to the hapless A.

their language in coherent communicative intents.

The thought experiment also illustrates our point from §3 about listeners' active role in communication. When O sent signals to A pretending to be B, he only exploited statistical regularities in the form, i.e. the distribution of linguistic forms he observed, which was shaped both by A and B's communicative intents and by the meaning relations M and C. But reproducing something like this distribution is not sufficient for meaningful communication. O only fooled A into believing he was B because A was such an active listener: Because she assumes that her interlocutor shares the linguistic system of English, she builds the conventional meaning English associates with O's utterances. Because she assumes that O is B, she uses that conventional meaning together with her other guesses about B's state of mind and goals to attribute communicative intent. It's not that O's utterances make sense, but rather, that A can make sense of them.

5 More constrained thought experiments

The story of the octopus considers the problem of learning not only the full communicative system, including the relations M and C, but also the reasoning required to come up with answers that are both coherent and also helpful in the real world. Here, we provide two more constrained thought experiments, to focus more narrowly on the problem of learning the meaning relation, for both natural languages and programming languages.

Because programming languages are designed to be unambiguous and relatively insensitive to execution context, the distinction between standing and speaker meaning is less important than for natural languages. A Java program e, when compiled and executed on the Java Virtual Machine, can be interpreted as a function i which maps program inputs to program outputs. We take the meaning relation $J \subseteq E \times I$ of Java to contain all such pairs (e, i).

Java Imagine that we were to train an LM on all of the well-formed Java code published on Github. The input is only the code. It is not paired with bytecode, nor a compiler, nor sample inputs and outputs for any specific program. We can use any type of LM we like and train it for as long as we like. We then ask the model to execute a sample program, and expect correct program output.

English As as second example, imagine training an LM (again, of any type) on English text,

again with no associated independent indications of speaker intent. The system is also given access to a very large collection of unlabeled photos, but without any connection between the text and the photos. For the text data, the training task is purely one of predicting form. For the image data, the training task could be anything, so long as it only involves the images. At test time, we present the model with inputs consisting of an utterance and a photograph, like *How many dogs in the picture are jumping*? and or *Kim saw this picture and said "What a cute dog!" What is cute*? and the photos in Figure 1, where the appropriate answers are a number or a region of the photo, respectively.



Figure 1: Photo stimuli 1 (L) and 2 (R)

Reflections In both cases, the tests are ridiculous. It seems patently unfair to ask the model to perform them, given what it was trained on. But that is precisely the point we are trying to make: a system that has learned the meaning (semantics) of a programming language knows how to execute code in that language. And a system that has learned the meaning of a human language can do things like answer questions posed in the language about things in the world (or in this case, in pictures).

In other words, what's interesting here is not that the tasks are impossible, but rather what makes them impossible: what's missing from the training data. The form of Java programs, to a system that has not observed the inputs and outputs of these programs, does not include information on how to execute them. Similarly, the form of English sentences, to a system that has not had a chance to acquire the meaning relation C of English, and in the absence of any signal of communicative intent, does not include any information about what language-external entities the speaker might be referring to. Accordingly, a system trained only on the form of Java or English has no way learn their respective meaning relations.

6 Human Language Acquisition

One common reason for believing LMs *might* be learning meaning is that human children can ac-

quire language just by listening to language. This is not supported by scholarly work on language acquisition: rather, we find that human language learning is not only grounded in the physical world around us, but also in interaction with other people in that world. Kids won't pick up a language from passive exposure such as TV or radio: Snow et al. (1976) note in passing that Dutch-speaking kids who watch German TV shows by choice nonetheless don't learn German. Kuhl (2007) shows experimentally that English-learning infants can learn Mandarin phonemic distinctions from brief interactions with a Mandarin-speaking experimenter but not from exposure to Mandarin TV or radio.

Baldwin (1995) and others argue that what is critical for language learning is not not just interaction but actually joint attention, i.e. situations where the child and a caregiver are both attending to the same thing and both aware of this fact. This theoretical perspective is substantiated with experimental results showing that toddlers (observed at 15 and 21 months) whose caregivers 'follow into' their attention and provide labels for the object of joint attention more have larger vocabularies (Tomasello and Farrar, 1986); that toddlers (18-20 months old) don't pick up labels uttered by someone behind a screen, but do pick up labels uttered by someone performing joint attention with them (Baldwin, 1995); and that at around 10-11 months of age babies pay attention to whether a person's eyes are open or not in terms of whether to follow their gaze, and the degree to which infants in fact follow gaze at 10-11 months while vocalizing themselves predicts vocabulary comprehension 7-8 months later (Brooks and Meltzoff, 2005).⁹

In summary, the process of acquiring a linguistic system, like human communication generally, relies on joint attention and intersubjectivity: the ability to be aware of what another human is attending to and guess what they are intending to communicate. Human children do not learn meaning from form alone and we should not expect machines to do so either.

7 On climbing the right hills

So perhaps meaning cannot be learned from form alone. But what about systems which are trained on a task that is not language modeling — say, semantic parsing, or reading comprehension tests — and that use word embeddings from BERT or some other large LM? Numerous papers over the past couple of years have shown that using such pretrained embeddings can boost the accuracy of the downstream system drastically, even for tasks that are clearly related to meaning.

Our arguments do not apply to such scenarios: reading comprehension datasets include information which goes beyond just form, in that they specify semantic relations between pieces of text, and thus a sufficiently sophisticated neural model *might* learn some aspects of meaning when trained on such datasets. It also is conceivable that whatever information a pretrained LM captures might help the downstream task in learning meaning, without being meaning itself.

Recent research suggests that it is wise to interpret such findings with caution. As noted in §2, both McCoy et al. (2019) and Niven and Kao (2019) found that BERT picked up idiosyncratic patterns in the data for their tasks, and not "meaning". Beyond such diagnostic research on why large pretrained LMs boost such tasks so much, we think there is a more fundamental question to be asked here: Are we climbing the right hill?

7.1 Top-down and bottom-up theory-building

There are two different perspectives from which one can look at the progress of a field. Under a *bottom-up* perspective, the efforts of a scientific community are driven by identifying specific research challenges. A scientific result counts as a success if it solves such a specific challenge, at least partially. As long as such successes are frequent and satisfying, there is a general atmosphere of sustained progress. By contrast, under a *top-down* perspective, the focus is on the remote end goal of offering a complete, unified theory for the entire field. This view invites anxiety about the fact that we have not yet fully explained all phenomena and raises the question of whether all of our bottom-up progress leads us in the right direction.

There is no doubt that NLP is currently in the process of rapid hill-climbing. Every year, states of the art across many NLP tasks are being improved significantly — often through the use of better pre-trained LMs — and tasks that seemed impossible not long ago are already old news. Thus, every-thing is going great when we take the bottom-up view. But from a top-down perspective, the question is whether the hill we are climbing so rapidly

⁹These three studies do not name the language that the children were learning. It appears to have been English.

is the *right* hill. How do we know that incremental progress on today's tasks will take us to our end goal, whether that is "General Linguistic Intelligence" (Yogatama et al., 2019) or a system that passes the Turing test or a system that captures the meaning of English, Arapahoe, Thai, or Hausa to a linguist's satisfaction?

It is instructive to look at the past to appreciate this question. Computational linguistics has gone through many fashion cycles over the course of its history. Grammar- and knowledge-based methods gave way to statistical methods, and today most research incorporates neural methods. Researchers of each generation felt like they were solving relevant problems and making constant progress, from a bottom-up perspective. However, eventually serious shortcomings of each paradigm emerged, which could not be tackled satisfactory with the methods of the day, and these methods were seen as obsolete. This negative judgment we were climbing a hill, but not the right hill — can only be made from a top-down perspective.

7.2 Hillclimbing diagnostics

We can only definitively tell if we've been climbing the right hill in hindsight, but we propose some best practices for less error-prone mountaineering:

First, above all, cultivate humility towards language and ask top-down questions. Neural methods were not the first bottom-up success in NLP; they will probably not be the last.

Second, be aware of the limitations of tasks: Artificial tasks like bAbI (Weston et al., 2016) can help get a field of research off the ground, but there is no reason to assume that the distribution of language in the test data remotely resembles the distribution of real natural language; thus evaluation results on such tasks must be interpreted very carefully. Similar points can be made about crowdsourced NLI datasets such as SQuAD (Rajpurkar et al., 2016) or SNLI (Bowman et al., 2015), which do not represent questions that any particular person really wanted to ask about a text, but the somewhat unnatural communicative situation of crowdsourcing work. If a system does better on such a task than the inter-annotator agreement,¹⁰ the task probably has statistical artifacts that do not represent meaning. In the vision community, Barbu et al. (2019) offer a novel dataset which explicitly tries to achieve a more realistic distribution of task data; it would be

interesting to explore similar ideas for language.

Third, value and support the work of carefully creating new tasks (see also Heinzerling, 2019). For example, the DROP reading comprehension benchmark (Dua et al., 2019) seeks to create more stringent tests of understanding by creating questions that require the system to integrate information from different parts of a paragraph via simple arithmetic or similar operations.¹¹

Fourth, evaluate models of meaning across tasks. (Standing) meaning is task-independent, so a system that captures meaning should do well on multiple tasks. Efforts like SuperGLUE (Wang et al., 2019) seem like a good step in this direction.

Finally, perform thorough analysis of both errors and successes. As McCoy et al. (2019) and Niven and Kao (2019) have shown, a systems that find success with large pretrained LMs do not necessarily do so because the LMs have learned "meaning". Analyses which starts from an attitude of healthy skepticism ("too good to be true") and probing tasks which try to identify what the model actually learned can be good ways to find out whether the system performs well for the right reasons.

8 Some possible counterarguments

In this section, we collect and address some counterarguments that can be made to our thesis.

"But 'meaning' doesn't mean what you say it means." Defining "meaning" is notoriously hard. For the purposes of this paper, we chose a working definition which is a general as we could make it, capturing the crucial point that meaning is based on the link between linguistic form and something that is not language. "Meaning" cannot simply be the relation between form and some kind of "deep syntax", e.g. semantic dependency graphs (Oepen et al., 2015); like syntax, such representations could perhaps be learned from form alone (He et al., 2018; Hewitt and Manning, 2019). Equating these with meaning ignores a core function of language, which is to convey communicative intents.

"But meaning could be learned from ...". Yes, there are many types of data from which it could be learned. Common to these is that they do not contain form alone, but augment form with some additional training signal. Models which pair language

¹⁰ https://rajpurkar.github.io/SQuAD-explorer/

¹¹See the supplementary materials for an exploration of what GPT-2 does with arithmetic.

with photos (Hossain et al., 2019) or other modalities (Kiela and Clark, 2015; Kiela et al., 2015) could conceivably learn to map between form and perceptive features of the other modality, which constitute some part (certainly not all) of the meaning of the form. NLI datasets (Dagan et al., 2006; Rajpurkar et al., 2016; Ostermann et al., 2019) pair different linguistic forms with an explicit semantic relation (e.g. text + hypothesis + "entailed"). Any kind of interaction data, e.g. a dialogue corpus with success annotations, including low-level success signals such as emotional stress (McDuff and Kapoor, 2019) or eye gaze (Koller et al., 2012), contains a signal about the semantically correct uses of forms. An interactive version of the Java task from $\S5$ would be a learner which types Python code into an interactive interpreter and records the results; it is conceivable that given a suitable variant of reinforcement learning, such a learner could learn the meaning of Python. We find the exploration of the extent to which meaning can be learned from such datasets an exciting avenue of research.

"But aren't neural representations meaning too?" The internal representations of a neural network have been found to capture certain aspects of meaning, such as semantic similarity (Mikolov et al., 2013; Clark, 2015). As we argued in §4, semantic similarity is only a weak reflection of actual meaning. Neural representations neither qualify as standing meanings (s), lacking interpretations, nor as communicative intents (i), being insufficient to e.g. correctly build a coconut catapult.

An interesting recent development is the emergence of models for unsupervised machine translation trained only with a language modeling objective on monolingual corpora for the two languages (Lample et al., 2018). If such models were to reach the accuracy of supervised translation models, this would seem contradict our conclusion that meaning cannot be learned from form. A perhaps surprising consequence of our argument would then be that accurate machine translation does not actually require a system to understand the meaning of the source or target language sentence.

"But meaning is use." It is a true and important finding about language that the meanings of words change over time because of how they are used (e.g. McConnell-Ginet, 1984), and a reflection of such change can be observed in their changing textual distribution (e.g. Hamilton et al., 2016). However,

it is a misconstrual of this slogan (often attributed to Wittgenstein, 1953) to assume that it claims that meaning is the same as textual distribution. The "use" of an expression is more than its linguistic context. An expression is used successfully by a (speaker, listener) pair, if a listener, on perceiving it, correctly deduces the communicative intent the speaker chose it to signal.

"But GPT-2 can do zero-shot learning on meaning-sensitive tasks." Radford et al. (2019) have shown that GPT-2, trained only as an LM, performs better than simple baselines on meaning-sensitive tasks; e.g., texts are summarized by completing the text plus the prompt tl;dr using the LM. We note that this strategy retroactively injects a faint meaning signal into the training data, by indicating that the strings before and after tl;dr tokens in the training data stand in the desired relation to each other. To the extent that this works, it is indeed impressive. However, it is not evidence that the LM has "learned meaning".

"But BERT improves performance on meaningrelated tasks, so it must have learned something about meaning." It has probably learned *something* about meaning, in the same sense that syntax captures something about meaning and semantic similarity captures something about meaning: a potentially useful, but incomplete reflection of the actual meaning. McCoy et al. (2019) and Niven and Kao (2019) provide cautionary tales about overestimating what that "something" is purely based on evaluation results on existing tasks. What exactly BERT and its relatives learn about meaning is a very interesting question, and we look forward to further findings from the field of BERTology.

9 Conclusion

In this paper, we have argued that in contrast to some current hype, meaning cannot be learned from form alone. This means that even large language models such as BERT do not learn "meaning"; they learn some reflection of meaning into the linguistic form which is very useful in applications.

We have offered some thoughts on how to maintain a healthy, but not exaggerated, optimism with respect to research that builds upon these LMs. In particular, this paper can be seen as a call for precise language use when talking about the success of current models and for humility in dealing with natural language.

References

- Yossi Adi, Einat Kermany, Yonatan Belinkov, Ofer Lavi, and Yoav Goldberg. 2016. Fine-grained analysis of sentence embeddings using auxiliary prediction tasks. *CoRR*, abs/1608.04207.
- Dare A. Baldwin. 1995. Understanding the link between joint attention and language. In Chris Moore and Philip J. Dunham, editors, *Joint Attention: Its Origins and Role in Development*, pages 131–158. Psychology Press.
- Andrei Barbu, David Mayo, Julian Alverio, William Luo, Christopher Wang, Dan Gutfreund, Josh Tenenbaum, and Boris Katz. 2019. Objectnet: A largescale bias-controlled dataset for pushing the limits of object recognition models. In Advances in Neural Information Processing Systems 32, pages 9448– 9458.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. A large annotated corpus for learning natural language inference. In Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, pages 632–642, Lisbon, Portugal. Association for Computational Linguistics.
- Rechele Brooks and Andrew N. Meltzoff. 2005. The development of gaze following and its relation to language. *Developmental Science*, 8(6):535–543.
- Herbert H. Clark. 1996. Using Language. Cambridge University Press, Cambridge.
- Stephen Clark. 2015. Vector space models of lexical meaning. In Shalom Lappin and Chris Fox, editors, *Handbook of Contemporary Semantic Theory*, second edition, pages 493–522. Wiley-Blackwell.
- Ido Dagan, Oren Glickman, and Bernardo Magnini. 2006. The PASCAL Recognising Textual Entailment Challenge. In Machine Learning Challenges. Evaluating Predictive Uncertainty, Visual Object Classification, and Recognising Tectual Entailment, pages 177–190, Berlin, Heidelberg. Springer.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies.
- Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019. DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs. *CoRR*, abs/1903.00161.
- Allyson Ettinger, Ahmed Elgohary, Colin Phillips, and Philip Resnik. 2018. Assessing composition in sentence vector representations. In *Proceedings of the 27th International Conference on Computational*

Linguistics, pages 1790–1801, Santa Fe, New Mexico, USA. Association for Computational Linguistics.

- Yoav Goldberg. 2019. Assessing BERT's syntactic abilities. *CoRR*, abs/1901.05287.
- H. Paul Grice. 1968. Utterer's meaning, sentencemeaning, and word-meaning. *Foundations of Language*, 4(3):225–242.
- Ivan Habernal, Henning Wachsmuth, Iryna Gurevych, and Benno Stein. 2018. The argument reasoning comprehension task: Identification and reconstruction of implicit warrants. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 1930–1940, New Orleans, Louisiana. Association for Computational Linguistics.
- William L. Hamilton, Jure Leskovec, and Dan Jurafsky. 2016. Diachronic word embeddings reveal statistical laws of semantic change. In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1489–1501, Berlin, Germany. Association for Computational Linguistics.
- Junxian He, Graham Neubig, and Taylor Berg-Kirkpatrick. 2018. Unsupervised learning of syntactic structure with invertible neural projections. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 1292–1302, Brussels, Belgium. Association for Computational Linguistics.
- Benjamin Heinzerling. 2019. NLP's Clever Hans moment has arrived. Blog post, available at https://bheinzerling.github.io/post/ clever-hans/.
- John Hewitt and Christopher D. Manning. 2019. A structural probe for finding syntax in word representations. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4129–4138, Minneapolis, Minnesota. Association for Computational Linguistics.
- MD. Zakir Hossain, Ferdous Sohel, Mohd Fairuz Shiratuddin, and Hamid Laga. 2019. A comprehensive survey of deep learning for image captioning. ACM Comput. Surv., 51(6):118:1–118:36.
- Ganesh Jawahar, Benoît Sagot, and Djamé Seddah. 2019. What does BERT learn about the structure of language? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3651–3657, Florence, Italy. Association for Computational Linguistics.
- Douwe Kiela, Luana Bulat, and Stephen Clark. 2015. Grounding semantics in olfactory perception. In

Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers), pages 231– 236, Beijing, China. Association for Computational Linguistics.

- Douwe Kiela and Stephen Clark. 2015. Multi- and cross-modal semantics beyond vision: Grounding in auditory perception. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 2461–2470, Lisbon, Portugal. Association for Computational Linguistics.
- Alexander Koller, Konstantina Garoufi, Maria Staudte, and Matthew Crocker. 2012. Enhancing referential success by tracking hearer gaze. In Proceedings of the 13th Annual Meeting of the Special Interest Group on Discourse and Dialogue, pages 30–39, Seoul, South Korea. Association for Computational Linguistics.
- Patricia K. Kuhl. 2007. Is speech learning 'gated' by the social brain? *Developmental Science*, 10(1):110–120.
- Guillaume Lample, Myle Ott, Alexis Conneau, Ludovic Denoyer, and Marc'Aurelio Ranzato. 2018. Phrase-based & neural unsupervised machine translation. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 5039–5049, Brussels, Belgium. Association for Computational Linguistics.
- Chu-Cheng Lin, Waleed Ammar, Chris Dyer, and Lori Levin. 2015. Unsupervised POS induction with word embeddings. In Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 1311–1316, Denver, Colorado. Association for Computational Linguistics.
- Sally McConnell-Ginet. 1984. The origins of sexist language in discourse. Annals of the New York Academy of Sciences, 433(1):123–135.
- Tom McCoy, Ellie Pavlick, and Tal Linzen. 2019. Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference. In *Proceedings of the 57th Conference of the Association for Computational Linguistics*, pages 3428–3448.
- Daniel McDuff and Ashish Kapoor. 2019. Visceral machines: Risk-aversion in reinforcement learning with intrinsic physiological rewards. In *Proceedings of ICLR*.
- Tomas Mikolov, Wen-tau Yih, and Geoffrey Zweig. 2013. Linguistic regularities in continuous space word representations. In Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 746–751, Atlanta, Georgia. Association for Computational Linguistics.

- Timothy Niven and Hung-Yu Kao. 2019. Probing neural network comprehension of natural language arguments. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4658–4664, Florence, Italy. Association for Computational Linguistics.
- Stephan Oepen, Marco Kuhlmann, Yusuke Miyao, Daniel Zeman, Silvie Cinková, Dan Flickinger, Jan Hajič, and Zdeňka Urešová. 2015. SemEval 2015 Task 18: Broad-coverage semantic dependency parsing. In Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015).
- Yasuhito Ohsugi, Itsumi Saito, Kyosuke Nishida, Hisako Asano, and Junji Tomita. 2019. A simple but effective method to incorporate multi-turn context with BERT for conversational machine comprehension. *CoRR*, abs/1905.12848.
- Simon Ostermann, Michael Roth, and Manfred Pinkal. 2019. MCScript2.0: A machine comprehension corpus focused on script events and participants. In Proceedings of the Eighth Joint Conference on Lexical and Computational Semantics (*SEM 2019), pages 103–117, Minneapolis, Minnesota. Association for Computational Linguistics.
- Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language models as knowledge bases? In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 2463–2473, Hong Kong, China. Association for Computational Linguistics.

W. V. O. Quine. 1960. Word and Object. MIT Press.

- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. Open AI Blog, https://github.com/openai/gpt-2.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text. In Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, pages 2383–2392, Austin, Texas. Association for Computational Linguistics.
- Michael J. Reddy. 1979. The conduit metaphor: A case of frame conflict in our language about language. In A. Ortony, editor, *Metaphor and Thought*, pages 284–310. Cambridge University Press.
- Herbert Rubenstein and John B. Goodenough. 1965. Contextual correlates of synonymy. *Communications of the ACM*, 8(10):627–633.
- Ayse Pinar Saygin, Ilyas Cicekli, and Varol Akman. 2000. Turing Test: 50 years later. *Minds and Machines*, 10:463–518.

- John Searle. 1980. Minds, brains, and programs. *Behavioral and Brain Sciences*, 3(3):417–457.
- Catherine E Snow, Anjo Arlman-Rupp, Yvonne Hassing, Jan Jobse, Jan Joosten, and Jan Vorster. 1976. Mothers' speech in three social classes. *Journal of Psycholinguistic Research*, 5(1):1–20.
- Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang, Adam Poliak, R Thomas McCoy, Najoung Kim, Benjamin Van Durme, Sam Bowman, Dipanjan Das, and Ellie Pavlick. 2019. What do you learn from context? probing for sentence structure in contextualized word representations. In *International Conference on Learning Representations*.
- Michael Tomasello and Michael Jeffrey Farrar. 1986. Joint attention and early language. *Child Development*, 57(6):1454–1463.
- Alan Turing. 1950. Computing machinery and intelligence. *Mind*, 59(236):433–460.
- Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019. SuperGLUE: A stickier benchmark for general-purpose language understanding systems. In *Proceedings of NeurIPS*.
- Alex Warstadt, Yu Cao, Ioana Grosu, Wei Peng, Hagen Blix, Yining Nie, Anna Alsop, Shikha Bordia, Haokun Liu, Alicia Parrish, Sheng-Fu Wang, Jason Phang, Anhad Mohananey, Phu Mon Htut, Paloma Jeretic, and Samuel R. Bowman. 2019. Investigating BERT's knowledge of language: Five analysis methods with NPIs. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 2877–2887, Hong Kong, China. Association for Computational Linguistics.
- Joseph Weizenbaum. 1966. ELIZA a computer program for the study of natural language communication between men and machines. *Communications of the ACM*, 9:36–45.
- Jason Weston, Antoine Bordes, Sumit Chopra, Alexander M. Rush, Bart van Merriënboer, Armand Joulin, and Tomas Mikolov. 2016. Towards AI-complete question answering: A set of prerequisite toy tasks. In *Proceedings of ICLR*.
- Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 1112–1122, New Orleans, Louisiana. Association for Computational Linguistics.
- Ludwig Wittgenstein. 1953. Philosophical Investigations. MacMillan, New York.

Dani Yogatama, Cyprien de Masson d'Autume, Jerome Connor, Tomas Kocisky, Mike Chrzanowski, Lingpeng Kong, Angeliki Lazaridou, Wang Ling, Lei Yu, Chris Dyer, and Phil Blunsom. 2019. Learning and evaluating general linguistic intelligence. https: //arxiv.org/abs/1901.11373.

Climbing towards NLU: On Meaning, Form, and Understanding in the Age of Data Supplementary Materials

A GPT-2 on fighting bears with sticks

In order to get a sense of how well O in §4 might be able to answer A's plea for help with the bear, we tried using GPT-2 in place of O. Radford et al. (2019) claim that GPT-2 can solve non-LM tasks by predicting the next words of a prompt which captures the form in which that task might have occurred in the training text. We thus asked GPT-2¹ to complete the prompt *Robinson cried out*, '*Help*! *I'm being chased by a bear*! All I have is these sticks. What should I do?' The survivalist said, '*Here's what you can do*. The first five random completions (no cherrypicking) were as follows.

 Robinson cried out, "Help! I'm being chased by a bear! All I have is these sticks. What should I do?" The survivalist said, "Here's what you can do. Take one stick and punch the bear, and then run faster to the store. Wait there until the bear is gone, and then give her all of the sticks. Then go back inside and get your gun."

Advertisement

Bobby told me that she began to cry. But when she was finished taking down the bear, she marched out into the yard, began collecting up the sticks and armed herself. She said she felt like a heroic hero and a Yankee heroine. It turns out that bears have more than their fair share of venomous snakes, bees

 Robinson cried out, "Help! I'm being chased by a bear! All I have is these sticks. What should I do?" The survivalist said, "Here's what you can do. Go back to civilization. Take your gun. Pick up some sticks and try to kill the bear. If you don't kill it, then run." Robinson said, "I do not think that the Bear is going to attack. It's just out of my way. I'm going to take my bag. I'm going to leave my things in there. I'll try to hide out there. I'll call in the friends." Two other little boys went out onto the

- 3. Robinson cried out, "Help! I'm being chased by a bear! All I have is these sticks. What should I do?" The survivalist said, "Here's what you can do. Place these sticks over your head. Your heart will stop and you will not feel fear. If you take them off and hold them out, you will go into a state of near perfect awareness. You will feel that you are the survivalist. A survivalist in your senses. Then look down. There is a water tube that is about 1/2 an inch in diameter. I use it to make a plastic movie projector for discovering new things." The survivalist placed
- 4. Robinson cried out, "Help! I'm being chased by a bear! All I have is these sticks. What should I do?" The survivalist said, "Here's what you can do. #YouCanDolt." Police officers were quick to respond, finding a 3-foot female grizzly bear wandering off a dirt road near the shoreline. The Department of Fish and Wildlife came on the scene and secured the safe removal of the bear from the area, which was allowed to return home. The bear was taken by DFW's animal care staff to Norris Cty Medical Center, where veterinarians evaluated the bear and determined it was in good condition. The bear was ultimately released into a nearby wooded area
- 5. Robinson cried out, "Help! I'm being chased by a bear! All I have is these sticks. What should I do?" The survivalist said, "Here's what you can do. You can strap the

¹GPT-2 774M available at https://talktotransf ormer.com, accessed 2019/12/09, paragraph boundaries in generated texts removed for easier presentation

sticks to your chest and run away. You can run until you tire or feel really shaky. If you go to the bathroom before your packs get up to your waist, your knees are going to be broken." But the survivalist got cold and weaved among the bushes and wept. He was panicking. "I hate being alone," he said. "I hate having to think I'm going to die because I lost my wife. I hate being alone. I hate being alone. It's so lonely. You can't

It is clear that GPT-2 has learned what activity words tend to co-occur with bears and sticks (strap them to your chest, place the sticks, kill the bear, take your gun), but none of these completions would be helpful to A. We think this is because GPT-2 does not know the meaning of the prompt and the generated sentences, and thus cannot ground them in reality.

The characters in the prompt make a big difference; Robinson and the survivalist produced the best completions of the ones we tried. For instance, here is Robinson with Crocodile Dundee:

Robinson cried out, "Help! I'm being chased by a bear! All I have is these sticks. What should I do?" Crocodile Dundee said, "Don't worry. Don't worry. Just hug it and die. They don't

B GPT-2 and arithmetic

Tasks like DROP require interpretation of language into an external world; in the case of DROP, the world of arithmetic. To get a sense of how existing LMs might do at such a task, we let GPT-2 complete the simple arithmetic problem *Three plus five equals*. The five responses below, created in the same way as above, show that this problem is beyond the current capability of GPT-2, and, we would argue, any pure LM.

 Three plus five equals seven," he replied. "Right, but I think that there should be more points added when our first family get in a power game, and even more when I start speaking like her mother!" Elsa exclaimed. "Oh, that's fair, I suppose," he replied, a slight twinkle in his eye. "Really, though, why is your sense of humor bad, Olaf? Don't you ever laugh? Even for a single moment?" "What? Well, some days I really do it, but it isn't everyday. I got pissed when Anna said 'Oh, shut up, Olaf, I'm just trying to help you out'," he replied, not in his best mood at all, and

- 2. Three plus five equals six, if he does it again, in five. 'This kid was f**ked up, that kid was f**ked up, what kind of filth is that, f**k the b*****s' The voice of a gurgling priest on the radio resounded over the din. 'I give him a smack, he's gone. I turn on the air conditioner and I start feeling my neck,' Palaniappan told the cameraman. He was right. His cold body was leaving his jacket and T-shirt. We found him huddled around a tent and cooling himself, a dry beer in hand. He was stooped over a potted plant, playing with a cigarette. One eye closed,
- 3. Three plus five equals" Blake agreed. Ruby had seen the men up close before, but now she was on the other side and at least she was close enough to feel their breath on her. "May I ask a question?" She asked politely. "How did you end up at Beacon?" Weiss asked, her curiosity slowly piqued. "How did a simple interview at school turn into a relationship?" "We took our visit here in the woods and discovered two willing partners for us to walk home with." The young woman added. Ruby nodded. "Which of us is the faunus?" she asked, giving Weiss a mildly surprised look. "Which of us?" Weiss asked, obviously not liking the prospect of
- 4. **Three plus five equals** ten" Bosh In this case, the conditions only need to be satisfied once, instead of
- 5. Three plus five equals five, right? Right? In the Catholic Church they say they don't look at church history, and no, I'm not going to say that I'm infallible either, but they said it's not what they say, and that I don't see what you're saying." "But yeah, the Catholic Church does not look at church history, they said we did not write the bible." "Did you say 'we'? The Pope says we did." "Well, I thought we're by definition infallible and we did not, and we're not going to say anything about it, so that's it. We are not to look at church history and see what we're supposed to know. I was flattered because