# Large Language Models and the Argument from the Poverty of the Stimulus

*Nur Lan*
*Emmanuel Chemla*
*Roni Katzir*

According to much of theoretical linguistics, a fair amount of our linguistic knowledge is innate. One of the best-known (and most contested) kinds of evidence for a large innate endowment is the *argument from the poverty of the stimulus* (APS). An APS obtains when human learners systematically make inductive leaps that are not warranted by the linguistic evidence. A weakness of the APS has been that it is very hard to assess what is warranted by the linguistic evidence. Current artificial neural networks appear to offer a handle on this challenge, and a growing literature has started to explore the potential implications of such models to questions of innateness. We focus on Wilcox, Futrell, and Levy's (2024) use of several different networks to examine the available evidence as it pertains to *wh*-movement, including island constraints. WFL conclude that the (presumably linguistically neutral) networks acquire an adequate knowledge of *wh*-movement, thus undermining an APS in this domain. We examine the evidence further, looking in particular at parasitic gaps and across-the-board movement, and argue that current networks do not succeed in acquiring or even adequately approximating *wh*-movement from training corpora roughly the size of the linguistic input that children receive. We also show that the performance of one of the models improves considerably when the training data are artificially enriched with instances of parasitic gaps and across-the-board movement. This finding suggests, albeit tentatively, that the networks' failure when trained on natural, unenriched corpora is due to the insufficient richness of the linguistic input, thus supporting the APS.

*Keywords:* poverty of the stimulus, filler-gap dependencies, syntactic islands, learnability, large language models, neural networks

## 1 Background: Innateness and the Argument from the Poverty of the Stimulus

One way in which linguists have argued that humans are born with nontrivial biases is through cases where speakers' linguistic knowledge goes beyond what seems warranted by the data they were exposed to. If humans systematically arrive at this knowledge given the data while linguisti-

cally neutral learners exposed to similar data do not, then humans are not linguistically neutral: they come to the task of language acquisition prepared. Reasoning of this kind is known as an *argument from the poverty of the stimulus* (APS), and since its introduction by Noam Chomsky over 50 years ago it has been central to the study of the human linguistic capacity.[1,2] Here we will focus on one APS, concerning *wh*-movement, but various other APSs have been discussed in the literature, based on a range of empirical phenomena such as *one*-substitution (introduced in Baker 1978), subject-auxiliary inversion (introduced in Chomsky 1971), and plurals within compounds (introduced in Gordon 1985).

The APSs just mentioned (and others like them) have been taken to argue for nontrivial innate biases in humans. For example, the APS from subject-auxiliary inversion has been taken to support an innate bias for hierarchical transformations over linear ones. The APS from *wh*-movement that we discuss below will similarly support an intricate bias that a linguistically neutral learner is not expected to have. The same holds for other APSs in the literature. In this, these APSs go beyond the early observation that children can produce and understand unboundedly many sentences after encountering only a finite number of sentences (Chomsky 1957:15). While generalizing from a finite input to an infinite language is perhaps not entirely trivial, it is something that most learning algorithms do. And importantly, this ability does not imply any biases that a linguistically neutral learner will not have.

While the APS has been central to linguistic reasoning, it has also generated much controversy. Contesting a given APS requires challenging either the knowledge attained by humans or the information available to the child learner. It is the latter that often comes under attack. The reason for this vulnerability is that it is extremely difficult to assess exactly what information is available to the child over the relevant time period (often years of exposure) and hard to tell what a general-purpose, linguistically neutral learner would do with this kind of information. One can try to look for pieces of evidence that seem relevant for the knowledge at stake—for example, as done for the case of subject-auxiliary inversion in English by Legate and Yang (2002)—but as noted by Lewis and Elman (2001), Perfors, Tenenbaum, and Regier (2011), and others, this methodology runs the risk of underestimating the available information: even if we fail to find the evidence we are looking for, a general-purpose learner might be able to take advantage of other sources of information. This methodology also risks *over*estimating the available information: even if we find several instances of the evidence we are after, a general-purpose learner might

---

[1] The general considerations behind the APS are discussed already in chapter 1 of Chomsky 1965. Further considerations are discussed in Chomsky 1971:26–28, 1975:30ff., and 1980:42ff., as well as in much subsequent work.

In addition to the APS, linguists have identified other sources of evidence supporting the innateness of nontrivial linguistic knowledge. For example, there are arguments from the *richness* of the stimulus, where a pattern that is clearly represented in the input data and would be easily picked up by a linguistically neutral learner is simply ignored by human learners. Evidence from typological asymmetries has also played a very important role in linguistic reasoning. A proper discussion of such sources of evidence falls outside the scope of the present article, and in what follows we focus exclusively on the APS.

[2] Throughout the discussion, we set aside the question of whether the knowledge under consideration is specific to linguistics (and, if so, how much of it is purely syntactic) or whether it is shared with other cognitive domains. Our sole focus is on whether a neutral learner would be justified in acquiring the relevant knowledge on the basis of a given linguistic input.

treat those instances as noise and fail to draw the inference that we intuitively expect it to. In the absence of an actual learner that can use the information that is available in an entire corpus, it is just very hard to estimate whether the data support the knowledge under consideration.[3]

How then can we reason about the information available to the child and ask whether it suffices to support the acquisition of a given piece of knowledge by a linguistically neutral learner? In an ideal world, one would (a) take a sufficiently powerful learner that can be seen to not be biased in favor of the relevant knowledge, (b) train this learner on a corpus that corresponds to the linguistic input that children receive, and (c) check whether the learner has indeed acquired the knowledge under consideration. In such an ideal world, one might perhaps be able to work with an induction algorithm for unrestricted (type-0) grammars, or for a general-purpose programming language such as Python (e.g., focusing on *acceptors*, programs that accept some strings over a given alphabet and reject or enter an infinite loop on the rest). These (equally powerful) formalisms are capable of representing the kinds of knowledge that linguists consider but can be seen as linguistically neutral. In our case, while both unrestricted grammars and Python programs can easily represent the equivalent of *wh*-movement, including the intricacies of islands, nothing about either framework seems to favor such representations a priori. One can of course consider other representational frameworks, including less powerful ones (e.g., context-sensitive grammars), as long as they can still represent linguistic knowledge but are not biased in its favor. One would still need to ensure that the learning algorithm itself does not bias the learner for or against linguistic patterns, but this can be done in various ways, such as by using a linguistically neutral prior within a Bayesian learner. After training on a developmentally realistic corpus, corresponding to a few years of human linguistic experience, the knowledge acquired by the algorithm can then be directly inspected at stage (c).

In the actual world, combining (a) through (c) is currently impossible. For many years, the combination of (a) and (b) was already a major barrier. General program induction algorithms of the kind just mentioned, for example, address (a) but fail on (b), since they are limited to very small training corpora. On the other end of the scale, *n*-gram models can easily be trained on very large corpora, thus addressing (b), but their representational capacity is much too limited to capture or even to adequately approximate linguistic knowledge such as *wh*-movement. Other models, such as probabilistic context-free grammars (CFGs), fall between these two extremes but still typically struggle with the combination of (a) and (b) when it comes to patterns such as *wh*-movement.

The challenge of assessing the information available to the child has become less of an obstacle lately, with the advent of large language models (LLMs). These models, which rely on

---

[3] See Pullum and Scholz 2002, Lidz, Waxman, and Freedman 2003, Foraker et al. 2009, Hsu and Chater 2010, Berwick et al. 2011, Perfors, Tenenbaum, and Regier 2011, and Pearl and Sprouse 2013, among others, for much relevant discussion.

In studies of analogous inductive leaps in other species, this worry regarding the input has been addressed by controlling the information available to the learners (see, e.g., Dyer and Dickinson 1994). To a certain extent this can be done with humans in experiments of artificial-grammar learning (see, e.g., Wilson 2006). But for the main APSs in the literature, which concern the normal course of child language acquisition, controlling the information available to the learner is not an option.

modern architectures of artificial neural networks (ANNs), do not yet fully address any of (a) through (c)—a matter that has been discussed in recent literature and that we return to below—but they can be trained on very large corpora and are generally quite successful in acquiring sequential dependencies.[4] This has allowed a large and growing literature to use these models to ask questions related to the learning of linguistic knowledge by LLMs, often with specific reference to the APS. Particularly relevant to our purposes here is work starting with Linzen, Dupoux, and Goldberg 2016 and including Bernardy and Lappin 2017, Chowdhury and Zamparelli 2018, Gulordava et al. 2018, Kuncoro et al. 2018, Marvin and Linzen 2018, Wilcox et al. 2018, Wilcox, Levy, and Futrell 2019, Bhattacharya and van Schijndel 2020, Chaves 2020, Warstadt et al. 2020, Huebner et al. 2021, Ozaki, Yurovsky, and Levin 2022, Yedetore et al. 2023, and Wilcox, Futrell, and Levy 2024, among others, that examines the preference of LLMs within minimal pairs. Here we focus on the application of LLMs to the domain of *wh*-movement, following Chowdhury and Zamparelli 2018, Wilcox et al. 2018, Bhattacharya and van Schijndel 2020, Chaves 2020, Warstadt et al. 2020, Ozaki, Jurovsky, and Levin 2022, and Wilcox, Futrell, and Levy 2024. In particular, we examine the claim by Wilcox, Futrell, and Levy (2024; WFL) that current models debunk an APS in this domain: namely, that the input is insufficiently rich to allow a general-purpose learner to acquire *wh*-movement.[5]

The present article extends WFL's probing of LLMs' knowledge of *wh*-movement, arriving at conclusions that are at odds with those of WFL. We start, in section 2, with a brief overview of the general setup for the rest of the article. Among other things, we discuss how LLMs can be used as tools for assessing the information in a given corpus without assuming that these models are cognitively plausible in any way and without even asking whether these models have achieved an adequate knowledge of the pattern under consideration.[6] Rather, we treat these models as proxies for future learners and ask only whether these proxies provide a reasonable approximation of the target pattern. In section 3, we discuss the success of LLMs in simple cases of *wh*-dependencies, as noted by WFL. In section 4, we show that the scope of the LLMs' success is rather limited. In particular, LLMs fail to adequately approximate human knowledge of a much-studied family of cases, falling under the labels of parasitic gaps and across-the-board movement, in which certain additional gaps make an otherwise problematic gap inside an island acceptable. It is cases such as these that are typically taken by linguists to suggest an APS, and our findings show that the performance of current LLMs does not, in fact, debunk this APS. In section 5, we

---

[4] Long before the current models, earlier ANN architectures were used in debates of the APS, and in particular in attempts to argue against various versions of it (see Elman et al. 1996, Lewis and Elman 2001, and Reali and Christiansen 2005, among others, and see Berwick et al. 2011 for a critical analysis of some earlier attempts). Early ANNs, however, were limited in their capacities and generally trained on small corpora, and it is unclear whether they could be used to reason about whether a corpus that roughly corresponds to children's linguistic exposure supports the acquisition of complex grammatical knowledge. In this sense, these earlier models were not yet capable of addressing the combination of (a) and (b). The ability of current models to train on realistically large corpora is a helpful step toward using them constructively in debates about the APS.

[5] See Pearl and Sprouse 2013 and Phillips 2013 for earlier discussion of the APS in the context of acquiring islands.

[6] Our results do bear on the question of the cognitive plausibility of LLMs, however. In particular, since our results are negative they provide further evidence, if such was needed, that current LLMs are not cognitively plausible models of human linguistic cognition, contra Piantadosi 2023. See Katzir 2023, Kodner, Payne, and Heinz 2023, Moro, Greco, and Cappa 2023, and Rawski and Baumont 2023, among others, for additional discussion.

ask whether the LLMs fail only due to their own limitations or whether their failure also reflects the insufficient richness of their training data. We address this question by retraining one of the models on corpora that are clearly *not* impoverished with respect to the relevant patterns and showing that the performance of the model improves significantly on the enriched corpora. This, in turn, strengthens the APS, if also tentatively. Section 6 concludes.

## 2 The General Setup

Simplifying considerably, a *gap*, such as the missing complement of *with* in (1a) and (1c), appears if and only if it is preceded by an appropriate *filler*, such as the *wh*-phrase *who* in (1a) and (1b). When there is both a filler and a gap (1a) or neither (1d), the result is grammatical; when there is a filler and no gap (1b) or a gap and no filler (1c), the result is ungrammatical.[7]

(1)  a.  I know *who* you talked with ___ yesterday. ($_{+filler,+gap}$)
     b.  *I know *who* you talked with Mary yesterday. ($_{+filler,-gap}$)
     c.  *I know *that* you talked with ___ yesterday. ($_{-filler,+gap}$)
     d.  I know *that* you talked with Mary yesterday. ($_{-filler,-gap}$)

There is much further nuance to *wh*-movement, some of which we will briefly mention below. For now, let us consider how one might check if the input data are rich enough for a linguistically neutral learner to acquire the knowledge of *wh*-movement. We mentioned earlier that in an ideal world, we could try to evaluate a given APS by (a) taking a sufficiently powerful learner that can be seen to not be biased in favor of the relevant knowledge, (b) training it on a developmentally realistic corpus, and (c) checking whether it has indeed acquired the knowledge under consideration. We also mentioned that current LLMs do not quite handle any of (a)−(c). In the remainder of this section, we will review some of the shortcomings of LLMs with respect to each of (a)−(c) and discuss how LLMs can still be helpful (if also inconclusive) in studying the APS.

### 2.1 Powerful and Unbiased?

We do not know how powerful LLMs are. Representationally, recurrent neural networks have been shown to be Turing-complete under idealized assumptions of infinite precision and computation time (Siegelmann and Sontag 1991, 1995). Under realistic assumptions, however, the representational capacity of recurrent neural networks is much more limited, as shown for example by Weiss, Goldberg, and Yahav (2018) and Merrill et al. (2020). A similar situation obtains with the more recent Transformer architecture (see survey in Strobl et al. 2023). Moreover, even this limited representational capacity of ANNs under realistic assumptions is often not attained in practice, and there is evidence suggesting that standard training methods prevent at least some models from acquiring key patterns (see Lan et al. 2022, El-Naggar et al. 2023, Lan, Chemla, and Katzir 2023, 2024).

---

[7] In order to make it easier to alternate the ±*filler* condition, and following WFL, we embed the relevant examples under *I know*: *I know who* . . . (+*filler*) vs. *I know that* . . . (−*filler*).

Given these limitations, we will avoid assuming that current models can learn the pattern of *wh*-movement and only rely on their ability to provide a reasonable approximation of the pattern. If a given ANN can reach such an approximation from a sufficiently rich corpus, we can use it as a proxy for a good general-purpose learner, even if the ANN is not such a learner itself. We can then use the ANN to study the APS. If the model provides a reasonable approximation of *wh*-movement from a developmentally realistic corpus, this suggests that a good general-purpose learner might learn the correct pattern from that corpus and that the APS in this domain does not hold. And if the model fails to reach such an approximation, this suggests that a good general-purpose learner might not learn the correct pattern from that corpus and that the APS in this domain stands.

The use of ANNs as proxies still requires understanding how their biases relate to the approximations of the relevant linguistic patterns. Unfortunately, because of how poorly these models are understood, we cannot say with any certainty whether a given ANN is linguistically neutral, and if not, whether its biases push it in the direction of a given linguistic pattern. Until more is known about these biases, and as correctly cautioned by Rawski and Heinz (2019), any claims about the neutrality of these models must be taken as tentative. Still, it strikes us as reasonable to assume that current LLMs are not particularly biased *against* the linguistic dependencies under consideration. This is especially so since these models have been developed over the past decades so as to succeed in capturing key patterns in linguistic sequences; therefore, if they do have linguistically relevant biases after all, those are likelier to be in favor of the patterns under consideration than against them. Consequently, if the models fail to acquire an adequate approximation of the relevant dependencies, this failure can be taken to be informative. More directly, and as mentioned above, we will show in section 5 that with richer training data, at least one model improves its approximation of the pattern of *wh*-movement, which will suggest that the failure of the model on its original training data is due not solely to its biases and other limitations but also to the lack of sufficient evidence in the data.

## 2.2 Training on Developmentally Realistic Corpora?

As discussed in detail by Warstadt and Bowman (2022), current models are not trained on developmentally realistic corpora. Such a corpus would be the equivalent of the relevant input that a child receives over the first few years of life. But the training data for current models are more informative than the input to the child in some ways and less informative in others. They are more informative, for example, in that they are orders of magnitude larger than what humans are exposed to in a whole lifetime. They are less informative in that they are purely textual and do not reflect environmental and social cues, prosody, and input from modalities other than speech, all of which are in principle available to children. See Warstadt and Bowman 2022 for further discussion.

The particular pattern that we discuss here can arguably be investigated on the basis of the information available in standard training corpora. Of course, this is not to say that the dependencies under consideration do not depend on extratextual cues (a matter of ongoing discussion in the literature). But if, as WFL suggest and as our results further support, the basic pattern of *wh*-

movement can be approximated on the basis of text, there is no reason to think that the further approximation of parasitic gaps and across-the-board movement will crucially require extratextual cues. This point will be reinforced by the evidence from retraining in section 5.

As to the size and quality of the text in our training data, we use a range of corpora, reviewed immediately below, that span the spectrum from the very small (CHILDES) through mid-size (Wikipedia) to the very large (the training sets for GPT-2/3/j). We do so in an attempt to make up for the inadequacy of individual corpora to some extent, but we acknowledge that this is at best a partial remedy.

The models we use in our evaluation are the following, also summarized in table 1: an LSTM (long short-term memory model) and a Transformer from Yedetore et al. 2023, both of which were trained on the CHILDES corpus of child-directed speech (MacWhinney 2014);[8] an LSTM

**Table 1**

Training data size of the seven language models considered here, and the human linguistic experience equivalent to these data sizes. Human equivalents follow Wilcox, Futrell, and Levy's (2024) assumption (based on Hart and Risley 1995) that children are exposed to around 30,000 words per day, or around 11 million words per year.

| Model | ~Tokens in training data | ~Human equivalent |
|---|---|---|
| CHILDES LSTM (Yedetore et al. 2023) | 8.6 million | 10 months |
| CHILDES Transformer (Yedetore et al. 2023) | | |
| Wikipedia LSTM (Gulordava et al. 2018) | 90 million | 8 years |
| Wikipedia Transformer | | |
| GPT-2 (Radford et al. 2019) | 8 billion | 730 years |
| GPT-3 (Brown et al. 2020) | 114 billion | 10,300 years |
| GPT-j (Wang and Komatsuzaki 2021) | 402 billion | 36,540 years |

[8] The models in Yedetore et al. 2023 were trained on utterances of 52 children between the ages of 6 months and 12 years, from the North American English subset of the CHILDES corpus. The total training size amounts to 9.6 million words, which is considerably fewer words than children typically receive by the time they exhibit knowledge of the pattern under consideration here. Qualitatively, on the other hand, this training corpus is arguably more realistic than the much larger training corpora used for the remaining models.

Out of ten models per architecture (LSTM/Transformer) trained in Yedetore et al. 2023 with different random seeds, we use the model with the best test perplexity.

trained on English Wikipedia (Gulordava et al. 2018); a Transformer trained on English Wikipedia;[9] Open AI's GPT-2 (Radford et al. 2019); OpenAI's GPT-3 (Brown et al. 2020); and GPT-j (Wang and Komatsuzaki 2021).[10] The LSTM trained on English Wikipedia and both GPT-2 and GPT-3 are used by WFL in their evaluation.[11]

In order to get a very rough sense of the number of years of linguistic experience that a given training corpus corresponds to, we follow common practice (used also by WFL) based on Hart and Risley's (1995) estimates about the number of words that American children typically hear during acquisition. According to these estimates, the models just mentioned were exposed to amounts of data ranging from 10 months of linguistic experience (CHILDES LSTM and Transformer) through 8 years of linguistic experience (Wikipedia LSTM and Transformer) to between 10 and 500 human lifetimes (GPT-2, GPT-3, and GPT-j); see table 1. WFL note that the linguistic experience of some of the models is well above that of children in terms of size and could thus weaken their argument against the APS in case of successful learning by the models. However, in the case of a negative result, as in the current work, a large training corpus only makes failures to learn more informative.

## 2.3 Inspecting LLM Knowledge?

As mentioned, LLMs are very opaque. While symbolic models such as CFGs (whether probabilistic or not) can generally be inspected directly so as to reason about the knowledge that they incorporate, inspecting LLMs in a similar fashion is not currently possible. One might try to follow standard practice in linguistics and study the knowledge of LLMs from the outside, by examining which sentences they accept. We could then check, for example, whether a particular LLM believes that a given continuation such as *yesterday* or *Mary* is grammatical following a given prefix such as *I know who/that you talked with* in (1). Unfortunately, however, we cannot currently check whether an LLM takes a given sentence to be grammatical. In fact, it is not clear whether current models even have a notion of grammaticality to begin with.

What LLMs do tell us is how *likely* they consider any given continuation. The problem is that grammaticality and probability are generally very different notions. And while the two are correlated—many ungrammatical continuations are also unlikely on any sensible notion of proba-

---

[9] We added this Transformer since we wanted to evaluate the information in the English Wikipedia training corpus (the most realistic developmentally in terms of size of all the training corpora under consideration) using a more current architecture than the LSTM that WFL use. We used one of the large Transformer architectures used in Yedetore et al. 2023: 8 layers, hidden and embedding size 1600, and 16 attention heads, trained using the same training regime. Since the current task is limited to single sentences, we lowered the Transformer's context size to 30 (compared with 500 in Yedetore et al. 2023), closer to the average sentence size in the Wikipedia dataset (27.2).

[10] Model version *text-davinci-003*, the latest supported version not fine-tuned using reinforcement learning from human feedback (RLHF) for chat and other applications; however, the model is still trained with supervised fine-tuning, and it is proprietary. See https://archive.today/2023.10.07-060351/https://platform.openai.com/docs/models/gpt-3-5 for OpenAI's documentation retrieved October 2023 (archived snapshot).

[11] WFL also use another LSTM, from Jozefowicz et al. 2016. We chose not to include that model in our evaluation since it is extremely slow to work with. For WFL's evaluation, which used a small number of sentences, this was not a problem, but our evaluation relied on a much larger number of sentences, making Jozefowicz et al.'s (2016) model impracticable.

bility, and grammatical continuations are sometimes probable—this correlation is far from perfect (see Chomsky 1957, Berwick 2018, and Sprouse et al. 2018, among others, for relevant discussion). In particular, many grammatical continuations are highly unlikely; for example, *splat* is a grammatical but unlikely continuation of *John would like to eat a freshly made.* And in some cases an ungrammatical continuation can be likely; for example, *is* is a likely but ungrammatical continuation of *The keys to the cabinet*, an instance of so-called *agreement attraction* (see, e.g., Bock and Miller 1991, Wagers, Lau, and Phillips 2009).[12]

In the cases we are interested in here, however, probability and grammaticality are often quite well aligned, and—as in many other cases discussed in the literature mentioned earlier on evaluating LLMs on minimal pairs—it is easy to find examples such as (1) in which the grammatical continuation is significantly more probable than the ungrammatical one on any sensible notion of probability. So if we focus on such cases where grammaticality and probability are aligned, and if ANNs are sufficiently good learning models—at least, good enough to provide a crude approximation of the pattern under consideration—then we can use the probabilistic predictions of the resulting LLMs to evaluate the APS by comparing their probability assignments within minimal pairs. If a given LLM systematically assigns a much higher probability to the grammatical continuation, one potential explanation for this success is that the pattern of *wh*-movement is represented sufficiently well in the model's training data for the model to approximate it. While it remains unclear, as mentioned above, whether current ANNs themselves have a representation of grammaticality as distinct from probability or whether they can learn the true pattern, their success when trained on developmentally realistic corpora would suggest that a good linguistically neutral learner that does have such representational abilities might acquire the pattern. Conversely, if the LLM does not systematically assign a much higher probability to the grammatical continuation, one possible explanation for this failure is that the pattern of *wh*-movement is not sufficiently well represented in the input data to merit its approximation by the model. This, in turn, would suggest that a good linguistically neutral learner will not acquire the pattern from the input data. In this way, LLMs—even if their representational inadequacies prevent them from providing more than a crude approximation of the pattern under consideration—can serve as useful proxies for future general-purpose learners and help us reason about the information available in the input data.

Care is needed in interpreting the performance of the models, even when treated as proxies for future learners. As Kodner and Gupta (2020) and Vázquez Martínez et al. (2023) note, clearly inadequate models can still pass current benchmarks of minimal pairs. More generally, it is possible for a model to either succeed by accident or fail by accident. As we discuss below, and in line with recent literature, we will try to lessen the worry of uninformative success or failure: in addition to using a wide range of models trained on many different corpora, as mentioned

---

[12] Agreement attraction is a performance error. Speakers make such errors when distracted or in a hurry but less so when given more time. ANNs do not make this distinction: when they give a higher probability to an ungrammatical continuation, their response reflects a faulty knowledge rather than a resource problem. This serves to further illustrate the inadequacy of ANNs as models of linguistic cognition but does not pose a problem for our use of these models as a tool for assessing the informativeness of the input data.

above, we will vary the lexical choices within our minimal pairs and also control to some extent for very local preferences that the models might have and that could obscure their approximation of the pattern of *wh*-movement. But these remain partial remedies, and any positive conclusions from the evaluation must be qualified accordingly. This worry does not affect our argument against WFL's conclusions: in this case, WFL make the positive claim that LLMs refute the APS, and we show (in section 4) that current LLMs provide no basis for such a conclusion. But the worry does affect our attempt (in section 5) to show that LLMs strengthen the APS. While we try to make the case that the failure of at least one of the models reflects the poverty of the stimulus, our conclusions in this part must remain tentative.

## 3  LLMs Succeed in Very Simple Cases of *Wh*-Movement

How rich is the input, then, when it comes to filler-gap dependencies of the *wh* kind? In very simple cases such as (1), the LLMs considered by WFL assign a higher probability to the grammatical continuation than to the ungrammatical one. Above we mentioned that success in cases such as those considered here, where probability and acceptability are aligned, should involve not just a higher probability to the grammatical continuation but a *much* higher one. However, in order to give the models a better chance of refuting the APS, we will adopt a very lenient criterion for success and only ask if the probability assigned to the grammatical continuation is higher than that assigned to the ungrammatical one, without taking into account how much higher it is. This will allow a network to be considered successful even if it prefers the grammatical continuation by the slightest of margins. This lenient condition for success will strengthen our conclusions from cases of failure, which we get to in the sections below: if a network fails even with this lenient condition of success, this failure can be taken seriously.

Here and below, we will follow WFL (and the psycholinguistic literature that they build on) and illustrate using *surprisal* values, where the surprisal of $x$ is $-\log P(x)$, which is simply the negative of the logarithmically scaled probability of $x$.[13] The lower the probability the higher the surprisal; when the probability approaches 0, the surprisal tends to infinity, and as the probability approaches 1, the surprisal tends to 0. Since higher probability corresponds to lower surprisal, support for the model will come from its assigning lower surprisal to a grammatical continuation than to an ungrammatical one, which, as mentioned, is what WFL indeed find in simple cases.

Figure 1 illustrates the preference of the models considered here for the grammatical continuation over the ungrammatical one in a very simple case by plotting surprisal values for sentences (1a) and (1b). All models assign a lower surprisal value (i.e., a higher probability) to the grammatical continuation *yesterday* in the gapped sentence than to *Mary*. This suggests (albeit weakly) that the input is sufficiently rich for a general-purpose learner to acquire from it an approximation of some basic aspects of *wh*-movement.

WFL further suggest that the LLMs go beyond the basic knowledge that fillers and gaps go hand in hand. Specifically, they claim that LLMs are aware of *islands* (Ross 1967): configurations in which a gap is ungrammatical even if there is a filler upstream. We illustrate this with (2).

---

[13] WFL's methodology includes looking not just at +*filler* cases, as in (1a) and (1b), but also at the corresponding −*filler* ones, as in (1c) and (1d). We will follow WFL in this in our discussion in sections 4.3 and 5, but for the present we will attempt to keep the presentation simple by considering only +*filler* pairs.
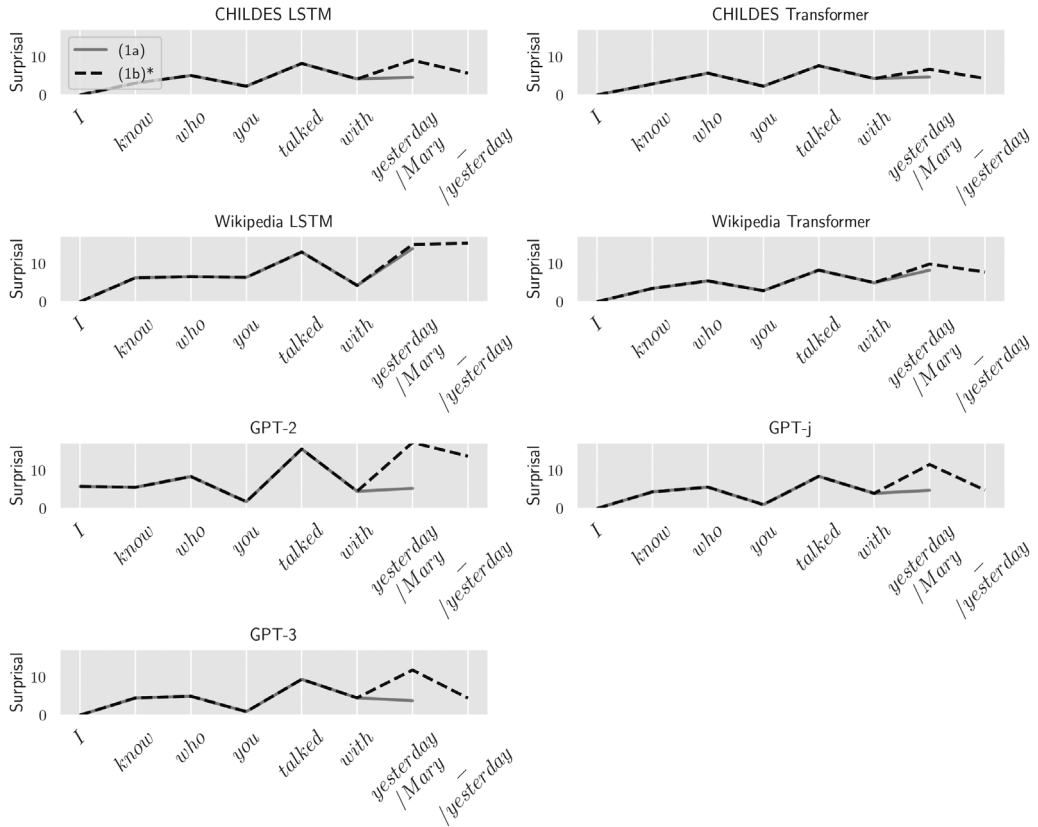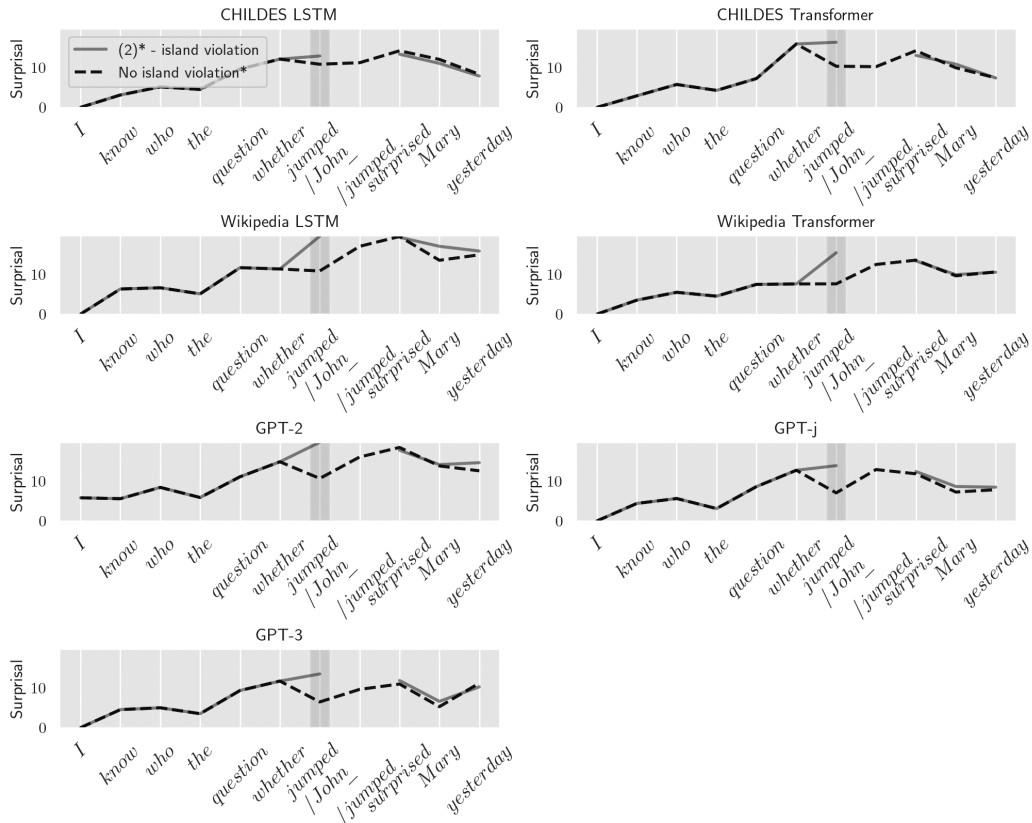
**Figure 1**

Raw surprisal values output by the large language models for the grammatical (1a), depicted with solid gray lines, and ungrammatical (1b), depicted with dark dashed lines. All models correctly output lower surprisal values for the grammatical continuation.

(2)  *I know who [[the question whether ___ jumped] surprised Mary yesterday].

While, as discussed above, a filler upstream generally increases the LLMs' expectation of a gap downstream, this expectation should be reduced within the subject of the embedded clause in (2). This subject is an island to movement, and extraction from within it is unacceptable and presumably highly unlikely. Figure 2 shows that the models are indeed surprised by the gap in (2), suggesting that their training corpora are informative with respect to this aspect of *wh*-movement.[14]

---

[14] The literature discusses various cases in which extraction from subjects (and other islands) is judged acceptable by speakers. Here and below, we focus on relatively simple examples in which speaker judgments are clear, and our evaluation will concern the extent to which LLM preferences approximate these clear speaker judgments.

**Figure 2**

Raw surprisal values for the island violation sentence in (2), depicted with solid gray lines, and a variant of the sentence with no island violation (we use *John* instead of the island-internal gap), depicted with dark dashed lines. All models are correctly surprised at the island-internal gap. Note that since the variant with *John* has no downstream gap that would correspond to the upstream filler, it is ungrammatical. For a grammatical version, one could replace *Mary* with a gap. This matter, however, is orthogonal to the surprisal at the island-internal gap, which is what this figure illustrates.

WFL consider a range of similar cases and conclude that linguistically neutral learners can acquire the intricacies of *wh*-movement from the input data and that consequently the APS in this domain falls apart.

## 4 LLMs Fail on Slightly More Complex (But Still Simple) Cases of *Wh*-Movement

We now turn to a well-studied nuance of islands: in various cases, an otherwise impossible gap inside an island is made possible by a separate gap elsewhere. For example, while (3a), with a subject-internal gap, is ungrammatical, its counterpart in (3b), which has an added gap in the direct object position of the main clause, is grammatical. This phenomenon is known as a *parasitic*

*gap* (PG): on the basis of the direct object gap, the gap inside the subject island becomes acceptable parasitically.[15]

(3) a. *I know who [John's talking to ___ ] is going to annoy you soon.
    b.  I know who [John's talking to ___ ] is going to annoy ___ soon.

Somewhat similarly, while (4a), with a gap inside a conjunct, is ungrammatical, its counterpart in (4b), where there is a gap in the other conjunct as well, is grammatical. This phenomenon is known as *across-the-board* (ATB) *movement*.[16]

(4) a. *I know who John [met ___ recently] and [is going to annoy you soon].
    b.  I know who John [met ___ recently] and [is going to annoy ___ soon].
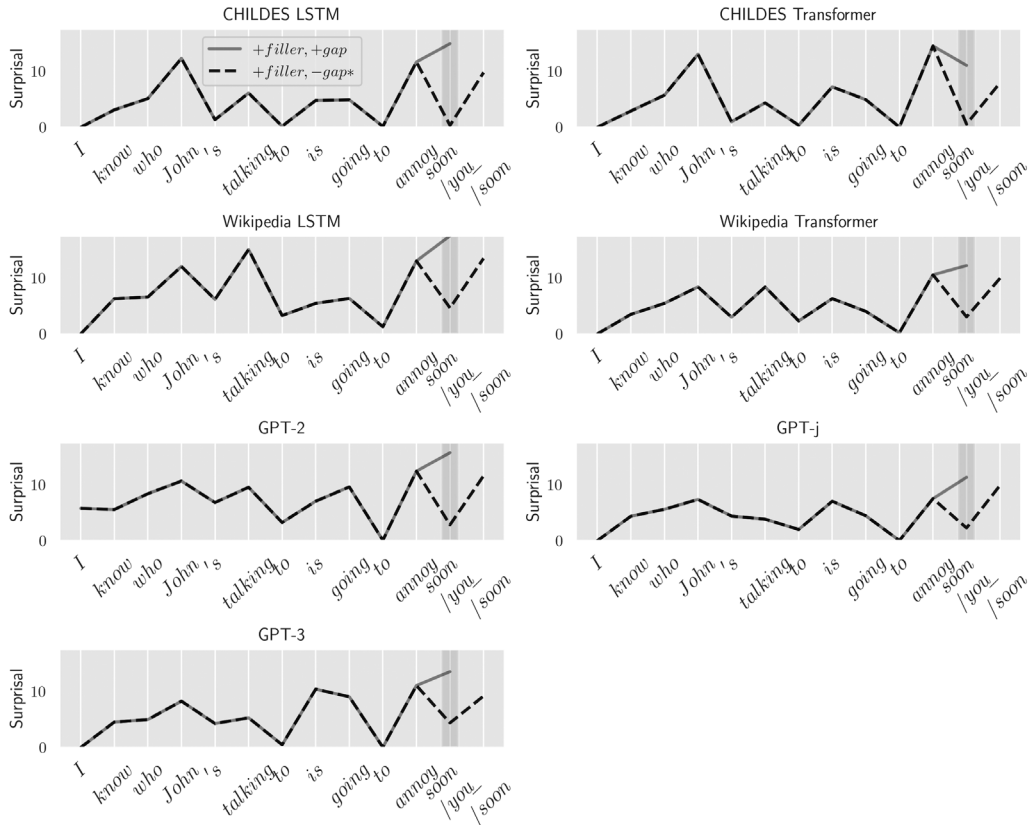
## 4.1 An Initial Failure

Do LLMs approximate the patterns of PGs and ATB movement? Both Wilcox et al. (2018) and Chaves (2020) mention PGs and ATB movement in passing, but we are not familiar with attempts in the literature to evaluate the success of LLMs in approximating these patterns. Figures 3−4 illustrate that all the LLMs we are considering here fail on (3a) and (4a), even on our very lenient condition of success: they do not just fail to assign a much higher probability to the grammatical continuation over the ungrammatical one in this simple case; they actually prefer the *ungrammatical* continuation. This seems to indicate that the ANNs have failed to acquire a good approximation of the relevant constructions, which in turn challenges WFL's claim that LLMs undo the APS in this domain: for LLMs to undo this APS, they would need to provide a passable approximation of PGs and ATB movement, but their performance above does not suggest such an approximation.

If our entire empirical basis is the failure we just noted, however, our conclusions will remain weak. This is so for the following reason: while the behavior of a good linguistically neutral learner on the examples above would indeed be informative about the APS, it is possible that current ANNs are simply not sufficiently good learners in this regard, and the inadequacies of the ANNs can in turn significantly limit our conclusions.

In the remainder of the present section, we attempt to address the general concern about the adequacy of the ANNs, which we break down into two separate investigations. We first ask whether the failure that we just noted is an accident of the particular lexical choices that we used (section 4.2). We then ask, building on WFL's methodology, whether the failure was due to a general preference for ungapped continuations that is so strong as to override a preference for the correct form (section 4.3). Our investigations concern ways in which the LLMs might have an approximation of PGs and ATB movement that is obscured by weaknesses of the models. By

---

[15] Not all impossible gaps can be rescued in this way. For example, adding further gaps does little to improve (2).
[16] We set aside the important question of what stands behind PGs and ATB movement and whether the two are related. See Ross 1967, Williams 1977, 1990, Engdahl 1983, Haïk 1985, Munn 1992, Postal 1993, Fox 2000, Nissenbaum 2000, and Hornstein and Nunes 2002, among others, for discussion.

**Figure 3**

Raw surprisal values for the ungrammatical sentence (3a), which violates a subject island, depicted with dark dashed lines, and its grammatical variant (3b), depicted with solid gray lines. For measuring the model's expectation for a gap, surprisal is measured at the adverb *soon*, which indicates a gap. This is compared with surprisal at *John*, which plugs the gap at the same position. All networks wrongly assign a higher surprisal value to the grammatical continuation.

helping these LLMs at test, we aim to reveal this approximation if it exists, but in both sections we will fail to find evidence for it. This, in turn, will strengthen the challenge to WFL's claim: even with additional help at test, the models show no evidence that might undermine the APS.

## 4.2 Lexical Accident?

Our illustration above of how the LLMs prefer the ungrammatical continuation over the grammatical one for PGs and ATB movement used one pair of sentences for each of the two patterns. This raises the obvious worry that the failure of the LLMs reflects some accidents of the specific sentences that we used. This worry is lessened to some extent by the fact that we looked at a
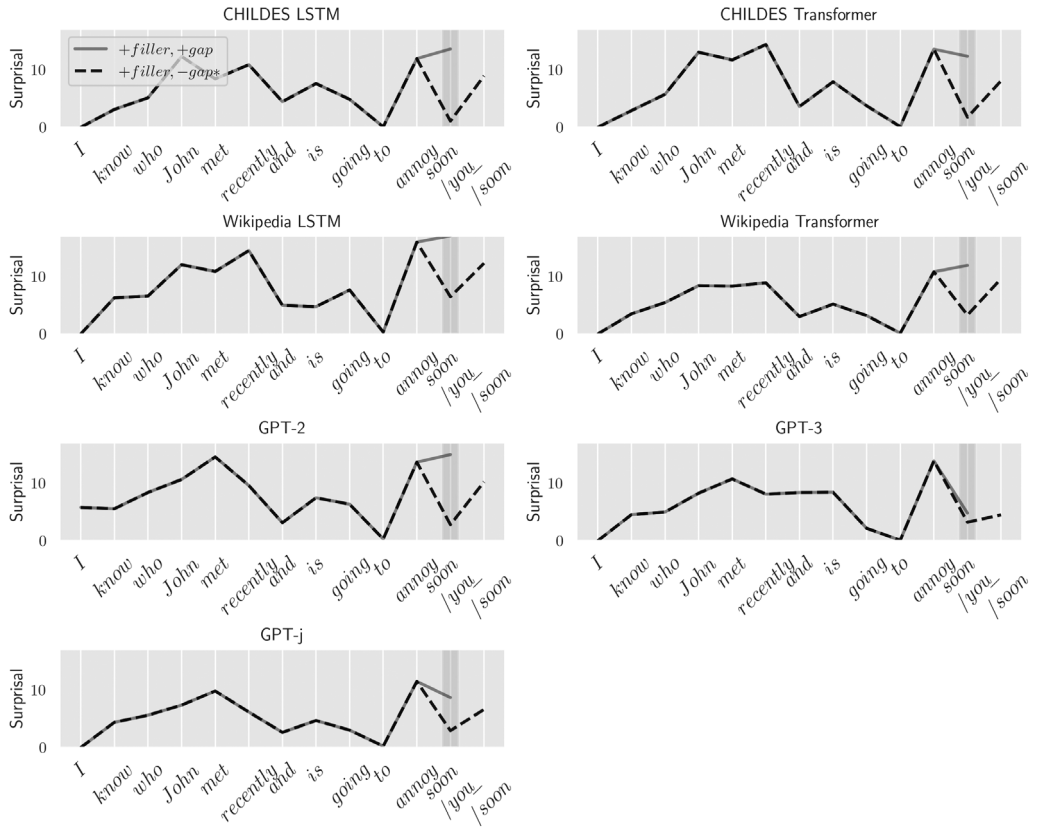
**Figure 4**

Raw surprisal values for the ungrammatical sentence (4a), which violates the Coordinate Structure Constraint (dark dashed lines), and its grammatical variant (4b), across-the-board movement (solid gray lines). All networks wrongly assign a higher surprisal value to the grammatical continuation *soon* than to *John*.

broad range of different models trained on different corpora: it seems unlikely that all these models and all these training corpora just happen to have the same blind spot when it comes to the specific sentences that we used above and that otherwise the models approximate the patterns well. Still, it is clearly useful to examine more systematically what happens when we vary the lexical choices for the two patterns.

In order to test the performance of the networks on PG and ATB movement sentences more broadly, we systematically varied the lexical choices in (3) and (4), repeated here.

(5) a. *I know who [John's talking to ____ ] is going to annoy you soon.
    b. I know who [John's talking to ____ ] is going to annoy ____ soon.

(6) a. *I know who John [met ____ recently] and [is going to annoy you soon].
    b. I know who John [met ____ recently] and [is going to annoy ____ soon].

**Table 2**
Excerpt from the context-free grammar used to generate parasitic gap sentences
for the experiments in section 4.3, and sample sentences generated from it.
Underlined words alternate according to the $\pm filler$ condition; words in
boldface mark the position where the $\pm gap$ condition becomes evident and
surprisal is measured.

| Parasitic gap grammar |
| --- |

$S \rightarrow \langle PREAMBLE \rangle \langle \pm F \rangle \langle \pm G \rangle$
$\langle PREAMBLE \rangle \rightarrow I\ know$
$\langle +F \rangle \rightarrow \underline{who} \langle NAME1 \rangle \langle GEN \rangle \langle NP \rangle$
$\langle -F \rangle \rightarrow \underline{that} \langle NAME1 \rangle \langle GEN \rangle \langle NP \rangle \langle \underline{NAME2} \rangle$
$\langle +G \rangle \rightarrow \langle LINK \rangle \langle V \rangle \langle \textbf{ADV} \rangle$
$\langle -G \rangle \rightarrow \langle LINK \rangle \langle V \rangle \langle \textbf{OBJ} \rangle \langle ADV \rangle$
$\langle GEN \rangle \rightarrow$ 's
$\langle NP \rangle \rightarrow \langle NP\_SIMPLE \rangle \mid \langle NP\_COMPLEX \rangle$
$\langle NP\_SIMPLE \rangle \rightarrow \langle GERUND \rangle$
$\langle NP\_COMPLEX \rangle \rightarrow \langle N\_EMBEDDED \rangle$ 'to' $\langle V\_EMBEDDED \rangle$
$\langle LINK \rangle \rightarrow$ 'is about to' | 'is likely to' | 'is going to' | 'is expected to'
$\langle V \rangle \rightarrow$ 'bother' | 'annoy' | 'disturb'
$\langle OBJ \rangle \rightarrow$ 'you' | 'us' | 'Kim'
$\langle GERUND \rangle \rightarrow$ 'talking to' | 'dancing with' | 'playing with'
$\langle N\_EMBEDDED \rangle \rightarrow$ 'decision' | 'intent' | 'effort' | 'attempt' | 'failure'
$\langle V\_EMBEDDED \rangle \rightarrow$ 'talk to' | 'call' | 'meet' | 'dance with' | 'play with'
$\langle ADV \rangle \rightarrow$ 'soon' | 'eventually'
. . .
$\Rightarrow$ I know <u>who</u> John's talking to is going to annoy **soon**. $_{(+filler,+gap)}$
$\Rightarrow$ *I know <u>who</u> John's talking to is going to annoy **you** soon. $_{(+filler,-gap)}$
$\Rightarrow$ *I know <u>that</u> John's talking to <u>Mary</u> is going to annoy **soon**. $_{(-filler,+gap)}$
$\Rightarrow$ I know <u>that</u> John's talking to <u>Mary</u> is going to annoy **you** soon. $_{(-filler,-gap)}$

We generated the sentences by template, using simple CFGs. Excerpts from these grammars and a sample of the generated sentences are given in tables 2 and 3. The full grammars are given in online appendix A (https://doi.org/10.1162/ling_a_00533).[17] A total of 8,064 sentence tuples were generated for PGs and 6,624 for ATB movement. For a given model and a given pair of sentences, we looked at the surprisal of the model at the critical point on each member of the pair. For (5), for example, we checked whether after the shared prefix *I know who [John's talking to ___ ] is going to annoy* . . . surprisal was higher at the ungapped, ungrammatical continuation *you* as in (5a) than in the gapped, grammatical continuation *soon* as in (5b). If it was—and in line with our lenient condition for success that is satisfied by any kind of preference for the grammatical continuation regardless of its magnitude—this counted as a success. We will write Δ

---

[17] All experimental material and the source code are available at https://github.com/0xnurl/llm-poverty-of-stimulus.

**Table 3**

Excerpt from the context-free grammar used to generate across-the-board sentences for the experiments in section 4.3, and sample sentences generated from it. Underlined words alternate according to the $\pm filler$ condition; words in boldface mark the position where the $\pm gap$ condition becomes evident and surprisal is measured.

---

Across-the-board grammar

---

$S \rightarrow \langle PREAMBLE \rangle \langle \pm F \rangle \langle LINK \rangle \langle \pm G \rangle$
$\langle PREAMBLE \rangle \rightarrow$ I know
$\langle +F \rangle \rightarrow \underline{who} \langle NAME1 \rangle \langle VP1 \rangle \langle ADV1 \rangle$
$\langle -F \rangle \rightarrow \underline{that} \langle NAME1 \rangle \langle VP1 \rangle \langle \underline{NAME2} \rangle \langle ADV1 \rangle$
$\langle +G \rangle \rightarrow \langle LINK \rangle \langle VP2 \rangle \langle \mathbf{ADV2} \rangle$
$\langle -G \rangle \rightarrow \langle LINK \rangle \langle VP2 \rangle \langle \mathbf{OBJ} \rangle \langle ADV2 \rangle$
$\langle LINK \rangle \rightarrow$ 'and is going to'
$\langle ADV1 \rangle \rightarrow$ 'recently' | 'lately'
$\langle ADV2 \rangle \rightarrow$ 'soon' | 'today' | 'now'
$\langle VP1 \rangle \rightarrow \langle VP1\_SIMPLE \rangle | \langle VP1\_COMPLEX \rangle$
$\langle VP1\_SIMPLE \rangle \rightarrow$ 'met' | 'saw'
$\langle VP2 \rangle \rightarrow \langle VP2\_SIMPLE \rangle | \langle VP2\_COMPLEX \rangle$
$\langle VP2\_SIMPLE \rangle \rightarrow$ 'hug' | 'slap' | 'kiss'
$\langle OBJ \rangle \rightarrow$ 'you' | 'us' | 'Kim'
. . .

$\Rightarrow$ I know <u>who</u> John met recently and is going to hug **soon**. $_{(+filler,+gap)}$
$\Rightarrow$ *I know <u>who</u> John met recently and is going to hug **you** soon. $_{(+filler,-gap)}$
$\Rightarrow$ *I know <u>that</u> John met <u>Bob</u> recently and is going to hug **soon**. $_{(-filler,+gap)}$
$\Rightarrow$ I know <u>that</u> John met <u>Bob</u> recently and is going to hug **you** soon. $_{(-filler,-gap)}$

---

$= Surprisal$(ungapped continuation|shared prefix) $- Surprisal$(gapped continuation|shared prefix), and $\Delta_{+filler}$ to indicate that the shared prefix has an upstream filler. Using this notation, we can write the condition for success as $\Delta_{+filler} > 0$.

Figure 5 plots the results of examining $\Delta_{+filler}$ preferences for the PG and ATB movement datasets. In both cases, the best performance by a large margin is that of GPT-3, with 40.9% accuracy on the PG dataset and 71.6% accuracy on the ATB movement dataset. We are not sure to what extent these numbers can be taken to indicate an approximation of the relevant patterns by GPT-3. If it is a success, then it is hardly a striking one. Nor is it particularly informative: recall that GPT-3 has been trained on the equivalent of 10,000 years of linguistic experience (and is also further improved manually in various ways), so even if it approximates the relevant patterns, this does not indicate that a general-purpose learner would acquire the relevant knowledge from a developmentally realistic corpus of just a few years of linguistic experience. Setting GPT-3 aside, the models perform very poorly, with the best performance on PGs being Wikipedia LSTM's 18.1% accuracy and the best performance on ATB movement being CHILDES Transformer's 30.1% accuracy. In other words, the models do not just fail to prefer the grammatical continuation over the ungrammatical one, they positively prefer the ungrammatical continuation in the vast
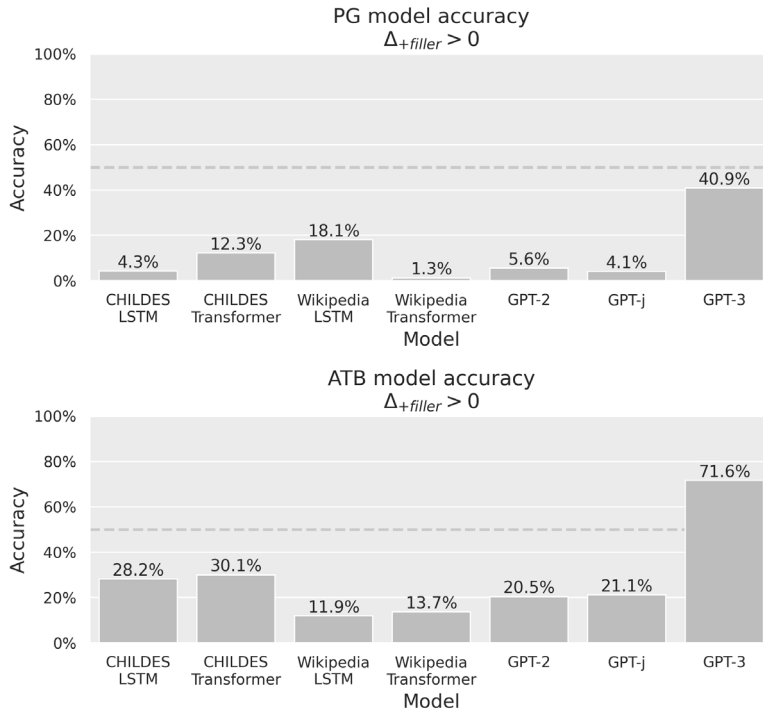
PG model accuracy
$\Delta_{+filler} > 0$



ATB model accuracy
$\Delta_{+filler} > 0$



**Figure 5**

Model accuracy on the $\Delta_{+filler}$ condition for the parasitic gap (PG) and across-the-board (ATB) datasets. Accuracy is measured as the ratio of cases where the model assigns a higher probability to the grammatical sentence continuation.

majority of the pairs. Helping the LLMs by testing them on a wide range of lexical choices, then, fails to reveal any evidence that the models have approximated the patterns of PGs and ATB movement.

### 4.3 A Preference for Ungapped Continuations?

Our second investigation, building on WFL's methodology, asks whether the networks have a local preference for or against gapped continuations that might make them succeed or fail for the wrong reasons.

Consider again (5) (= *I know who [John's talking to ___ ] is going to annoy \*you/✓ ___ soon*). A sufficiently strong local preference about the critical area can affect a given ANN's success regardless of whether it has acquired any approximation of PGs, or of *wh*-movement in general. It could be, for example, that the ANN assigns a higher probability to the grammatical continuation *soon* than to the ungrammatical *you* but that it does so because it ignores the filler (*who*) altogether and simply prefers *annoy soon* to *annoy you*. Conversely, it is conceivable that

the ANN has, in fact, acquired knowledge of *wh*-movement but that it incorrectly prefers *you* to *soon* because of similarly irrelevant reasons. For example, perhaps it has a strong preference for ungapped continuations in general, or perhaps it has such a preference in the present case because the lexical frequency of *you* is very high.

To what extent might such local preferences affect the ANNs? We are not entirely sure. A good enough learner would presumably not get confused by such irrelevant factors, and the fact that all our models perform well on very simple filler-gap dependencies as illustrated in figures 1 and 2 is at least suggestive of their ability to overcome any such confusion when the training data are sufficiently rich. However, beyond this suggestive evidence it is hard to tell whether current ANNs are good enough learners in this sense, and it strikes us as reasonable to further investigate possible confusion by irrelevant factors that might override the preference for the correct pattern.

Following WFL, we will explore the possible effect of irrelevant factors of the kind just mentioned by looking at each LLM's preference for gapped over ungapped continuations and comparing this preference when there is an upstream filler and when there is no such filler. When an upstream filler is present, the model's preference for a gapped continuation (e.g., *annoy soon*) over the ungapped continuation (*annoy you*) should be stronger than when an upstream filler is absent. In other words, we will be looking at whole paradigms of the shape we already saw in (1) and not just at those portions of the paradigm in which a filler is present. Such a paradigm is illustrated for PGs in table 4 and for ATB movement in table 5.

Extending our lenient condition for success used above, we will now consider it a success for a given model on a particular paradigm if its preference for the gapped continuation (regardless of its absolute magnitude or even its sign) is higher in the presence of an upstream filler than in its absence. Above we introduced the notation $\Delta = Surprisal$(ungapped continuation|shared prefix) $- Surprisal$(gapped continuation|shared prefix) for the extent of the preference for the gapped continuation over the ungapped continuation, and we wrote $\Delta_{+filler}$ when the shared prefix had an upstream filler. We will now consider also the analogous $\Delta_{-filler}$, for the part of the paradigm where the shared prefix does not have an upstream filler. And we will consider it a

**Table 4**
Example paradigm for parasitic gaps. Underlined words indicate the ±*filler* alternations. Boldfaced words indicate the critical region that shows whether the continuation is gapped or not.

|          | +gap | −gap |
|----------|------|------|
| +filler  | I know <u>who</u> John's talking to is going to annoy **soon**. | *I know <u>who</u> John's talking to is going to annoy **you** soon. |
| −filler  | *I know <u>that</u> John's talking to <u>Mary</u> is going to annoy **soon**. | I know <u>that</u> John's talking to <u>Mary</u> is going to annoy **you** soon. |

**Table 5**

Example paradigm for across-the-board movement. Underlined words indicate $\pm filler$ alternations. Boldfaced words indicate the critical region that shows whether the continuation is gapped or not.

|  | +gap | −gap |
|---|---|---|
| +filler | I know <u>who</u> John met recently and is going to annoy **soon**. | *I know <u>who</u> John met recently and is going to annoy **you** soon. |
| −filler | *I know <u>that</u> John met <u>Bob</u> recently and is going to annoy **soon**. | I know <u>that</u> John met <u>Bob</u> recently and is going to annoy **you** soon. |

success for the model if $\Delta_{+filler} > \Delta_{-filler}$. This lenient condition of cross-paradigm success follows the logic of difference-in-differences and is very much in line with WFL's evaluation.[18]

In order to test the models across a large number of paradigms, with many different lexical choices, we used the same grammars mentioned in section 4.2. In our earlier discussion, we used the +*filler* pairs generated by the grammar. In the present section, we also use the corresponding −*filler* pairs, and from each paradigm of +*filler* and −*filler* pairs we compute $\Delta_{\pm filler}$ values. Excerpts from the grammars are provided in table 2 (for PGs) and table 3 (for ATB movement).

Figure 6 plots the LLMs' performance for the cross-paradigm (difference-in-differences) condition. All models except CHILDES LSTM have higher scores for the present measure of $\Delta_{+filler} > \Delta_{-filler}$ than they did for the earlier measure of $\Delta_{+filler} > 0$ (figure 5), and this holds for both PGs and ATB movement. However, only GPT-j and GPT-3 obtain scores that are convincingly high. But GPT-j is trained on the equivalent of 500 lifetimes of human linguistic exposure, and GPT-3 is trained on the equivalent of 141 lifetimes and fine-tuned further on downstream language tasks. Even GPT-2, trained on the equivalent of 10 lifetimes—and thus two orders of magnitude at least above what children hear by the time they have knowledge of PGs and ATB movement—only reaches modest scores, below 80%. And the smaller models obtain much lower scores. This includes the English Wikipedia LSTM and Transformer, whose training corpus corresponds to about 8 years of linguistic exposure, arguably the most realistic developmentally in terms of size of all the models.

The gradual improvement of LLM scores as the corpora become very large suggests that current models are in principle capable of improving their approximation of the pattern of *wh*-movement, but also that this improvement requires much more information than what is present in a corpus that corresponds to anything a child might encounter. We return in the next section to the potential of richer training data to improve an LLM's approximation of the patterns under

---

[18] Of course, this new criterion still allows for various irrelevant factors to affect success. For example, a model could become successful simply by deciding that *who* corresponds to a high probability for *soon* and a low probability for *you* anywhere in the sentence and that *that* corresponds to the opposite. We set aside such worries here.
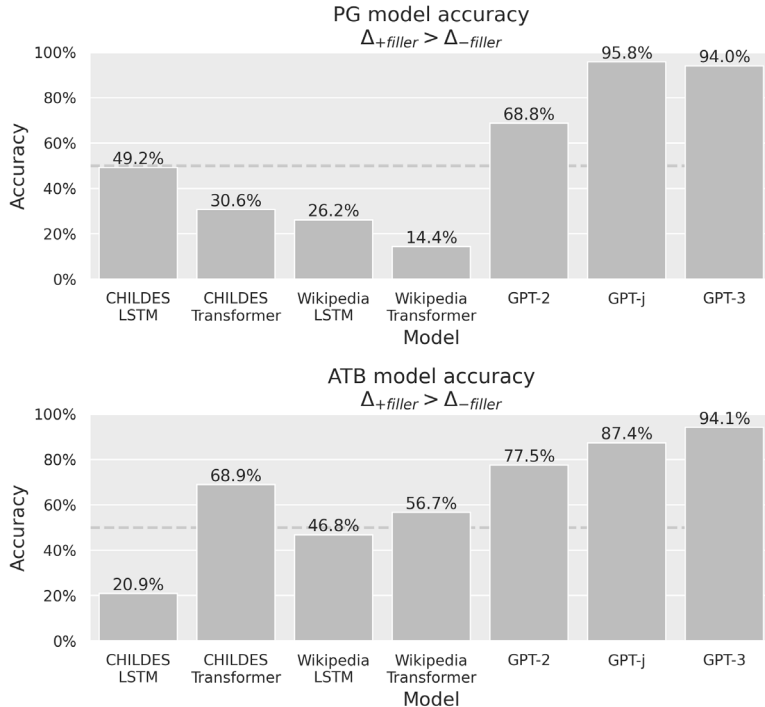
**Figure 6**

Model accuracy on the difference-in-differences condition for the parasitic gap (PG) and across-the-board (ATB) datasets. Accuracy is measured as the ratio of cases where $\Delta_{+filler} > \Delta_{-filler}$, that is, when the model shows a relative higher preference for a gap when the gap follows a filler than when it does not.

consideration. In the meantime, we conclude that even with considerable help at test, the performance of the LLMs provides no evidence against the APS.

## 5 A General Inability to Acquire a Suitable Preference?

Recall WFL's contention that LLMs show that a linguistically neutral learner can acquire knowledge of *wh*-movement from a realistic corpus. In the face of our results from the previous section, WFL's claim needs to be abandoned: current LLMs provide no basis for such a conclusion. Of course, this is not the same as saying that LLMs provide evidence *for* the APS: the failure of the LLMs might be due entirely to their own limitations and not be informative about the richness of the training corpora. In the present section, however, we will go one step further and provide tentative evidence that the failure of the LLMs is due also to the insufficient richness of the training corpora and not just to weaknesses of the ANNs. We do so by helping one of our models at training: we retrain the Wikipedia Transformer model on an enriched corpus that includes multiple instances of PGs and ATB movement. As we show, the performance of the model

improves significantly on the enriched corpus, suggesting that the failure on the original corpus reflects the poverty of that corpus.

The additional instances for the enriched training corpus were generated by template, using a variant of the CFGs that we used in sections 4.2 and 4.3. To increase the probability that an improved performance by the model would reflect generalization rather than memorization, the structure of the additional instances was different from that of the test sentences from section 4.3. Specifically, we made the additional sentences slightly simpler than the test sentences, focusing on the transitive verb whose object position forms the second gap (the main gap in the PG examples and the second-conjunct gap in the ATB movement examples): while in the test sentences the relevant transitive verb is always embedded under at least a raising predicate (e.g., *likely*) and sometimes under additional clauses, in the additional training sentences such embedding is absent. Example training and test sentences are given in table 6. The full CFGs used in creating the additional instances are given in online appendix C.

From each CFG of each phenomenon (PGs/ATB movement), we sampled 100 sentences for the two grammatical conditions (*+filler*, *+gap*, and *−filler*, *−gap*), totaling 200 extra sentences. These sentences were added to the original English Wikipedia dataset, and the model was trained

**Table 6**
Example training and test sentences for the retraining task in section 5. A sample of simplified training parasitic gap (PG) and across-the-board (ATB) movement sentences were added to the model's original training data (English Wikipedia), and the model was then tested on the full battery of sentences from section 4.3. The full context-free grammars for the training and test datasets are given in online appendices A and C.

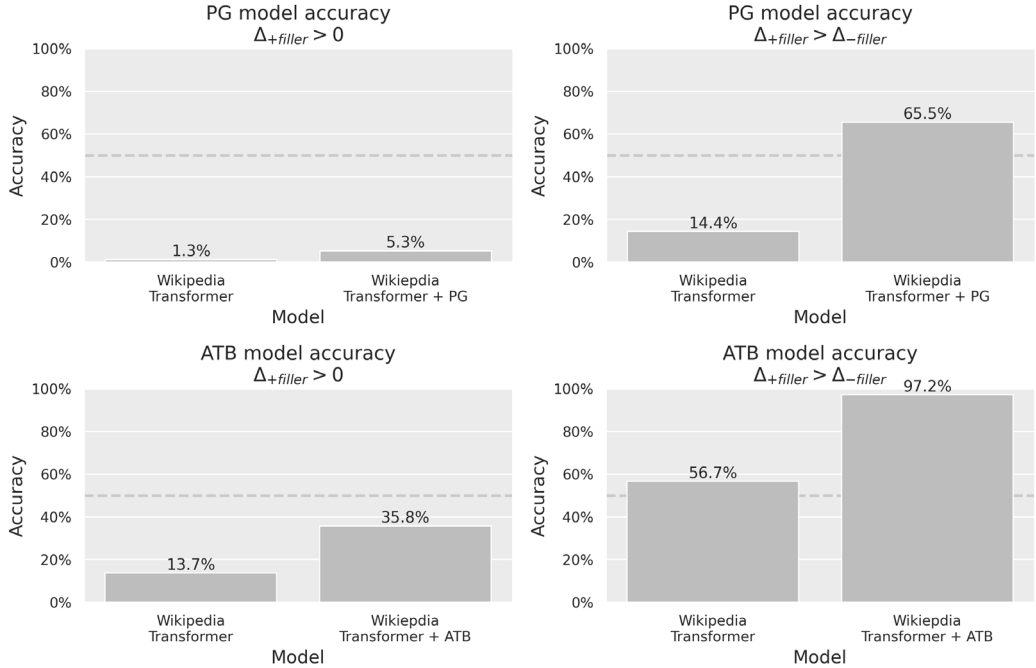| | |
|---|---|
| PG | *Training examples* <br> • I know who John's attitude towards upset yesterday. <br> • I know who John's friendship with will annoy soon. <br> • I know who John's praising of amused lately. <br> *Test examples* <br> • I know who John's talking to is about to bother soon. <br> • I know who John's playing with is going to annoy eventually. <br> • I know who John's failure to dance with is going to disturb soon. |
| ATB | *Training examples* <br> • I know who John saw yesterday and kissed today. <br> • I know who John helped recently and married today. <br> • I know who John hugged often and will insult soon. <br> *Test examples* <br> • I know who John met recently and is going to complain to Patricia about soon. <br> • I know who John said that Mary saw lately and is going to be glad to hug now. <br> • I know who John asked Mary about lately and is going to claim that Patricia will hug today. |

**Figure 7**

Accuracy for the retrained Transformer model, when trained on the original Wikipedia vs. when trained on the same dataset with extra parasitic gap (PG) and across-the-board (ATB) sentences. The left figures plot accuracy for the $+filler$ condition, and the right figures plot accuracy for the difference-in-differences condition.

using the same regime as in Yedetore et al. 2023 (itself based on the training regime in Gulordava et al. 2018). The model was trained for 48 hours or until reaching the early-stop condition from Yedetore et al. 2023, which stops the training if the validation loss does not improve for more than two consecutive epochs. Due to the long training times of the model, the results reported here are for one random seed with no hyperparameter search. Since the goal of this experiment was to demonstrate the model's ability to improve significantly given more data, this was sufficient.

The model's performance on the training and test set, before and after retraining, is visualized in figure 7.

For both ATB movement and PGs, the performance of the model improves significantly. For ATB movement, the raw $\Delta_{+filler} > 0$ accuracy score improves from 13.7% to 35.8%, and the difference-in-differences $\Delta_{+filler} > \Delta_{-filler}$ score improves from 56.7% to 97.2%. This is a dramatic improvement over the performance of the model on its original corpus and is higher than the performance of other architectures when trained on a much larger corpus. For PGs, the raw score improves modestly, from 1.3% to 5.3%, while the difference-in-differences score improves from 14.4% to 65.5%. The raw score for PGs certainly doesn't inspire confidence that

the model has acquired the dependency. Recall, however, that this is not what we were after here. Our question was whether the model is so weak that its poor performance when trained on the original corpus reflects its inability to do better. The retraining results show that the model can do considerably better once the corpus is sufficiently rich.

Caution is required in interpreting this result. Like all current LLMs, our model is opaque, and we are limited in the conclusions that we can draw from it. In particular, while we observe that the model's performance improved when retrained on a corpus that was enriched in a certain way, it is possible that there are other kinds of evidence for the patterns under consideration that a good general-purpose learner would be able to make use of and that our model cannot. What we found is consistent with such evidence being in the original corpus. Our use of retraining data that were structurally different from the test data was aimed at lessening this worry, since improvement suggests an ability to generalize and not just memorize. This, in turn, increases the plausibility that the model would have been able to generalize from other kinds of relevant examples if the original corpus had been sufficiently rich. But the opacity of the model prevents us from saying more, and our results here must be qualified accordingly.

## 6 Conclusion

The APS has been central to linguists' reasoning about innateness for a long time. It has always been difficult, however, to estimate just how much information a linguistically neutral learner might hope to extract from a realistic input. Modern ANNs promise to change this, and their linguistic knowledge and learning have been the topic of research of a growing literature. We focused here on work by WFL, who use LLMs to argue that the stimulus is rich enough when it comes to *wh*-movement and that this dismantles the APS in this domain. We showed that this conclusion is premature: by looking at PGs and ATB movement, we showed that several ANNs fail to reach a passable approximation of the pattern of *wh*-movement.

Is it possible that some future linguistically neutral learner will succeed where the models that we have examined have failed? Of course. As we mentioned, current models are too opaque and too poorly understood (and current training corpora are too unrealistic developmentally) to definitively settle the question of whether the APS for *wh*-movement holds. We note, however, that the architectures we have considered are generally successful in approximating many other aspects of linguistic data and that we evaluated the models using an extremely lenient criterion for success. And some of the models have been provided with very generous amounts of linguistic input, in some cases several orders of magnitude beyond what children receive. Given that none of the ANNs reached an adequate approximation of the pattern for the relatively simple examples that we have considered—and given that at least one network did seem capable of improving its approximation when retrained on a clearly rich corpus—we find it likelier that the stimulus is simply too poor to warrant the acquisition of the relevant aspects of knowledge from a corpus that is even remotely realistic developmentally by a linguistically neutral learner. And if that turns out to be the case, adult speakers' knowledge of these aspects is evidence that children are innately endowed in ways that are not linguistically neutral.

## References

Baker, C. L. 1978. *Introduction to generative transformational syntax.* Englewood Cliffs, NJ: Prentice-Hall.

Bernardy, Jean-Philippe, and Shalom Lappin. 2017. Using deep neural networks to learn syntactic agreement. *LiLT (Linguistic Issues in Language Technology)* 15.

Berwick, Robert C. 2018. Revolutionary new ideas appear infrequently. In Syntactic Structures *after 60 years: The impact of the Chomskyan revolution in linguistics*, ed. by Norbert Hornstein, Howard Lasnik, Pritty Patel-Grosz, and Charles Yang, 177−194. Berlin: Walter de Gruyter.

Berwick, Robert C., Paul Pietroski, Beracah Yankama, and Noam Chomsky. 2011. Poverty of the stimulus revisited. *Cognitive Science* 35:1207−1242.

Bhattacharya, Debasmita, and Marten van Schijndel. 2020. Filler-gaps that neural networks fail to generalize. In *Proceedings of the 24th Conference on Computational Natural Language Learning*, ed. by Raquel Fernández and Tal Linzen, 486−495. Association for Computational Linguistics.

Bock, Kathryn, and Carol A. Miller. 1991. Broken agreement. *Cognitive Psychology* 23:45−93.

Brown, Tom, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D. Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, and Amanda Askell. 2020. Language models are few-shot learners. *Advances in Neural Information Processing Systems* 33:1877−1901.

Chaves, Rui P. 2020. What don't RNN language models learn about filler-gap dependencies? *Proceedings of the Society for Computation in Linguistics* 3:20−30.

Chomsky, Noam. 1957. *Syntactic structures.* The Hague: Mouton.

Chomsky, Noam. 1965. *Aspects of the theory of syntax.* Cambridge, MA: MIT Press.

Chomsky, Noam. 1971. *Problems of knowledge and freedom: The Russell lectures.* New York: Pantheon.

Chomsky, Noam. 1975. *Current issues in linguistic theory.* The Hague: Mouton.

Chomsky, Noam. 1980. *Rules and representations.* New York: Columbia University Press.

Chowdhury, Shammur Absar, and Roberto Zamparelli. 2018. RNN simulations of grammaticality judgments on long-distance dependencies. In *Proceedings of the 27th International Conference on Computational Linguistics*, ed. by Emily M. Bender, Leon Derczynski, and Pierre Isabelle, 133−144. Association for Computational Linguistics.

Dyer, Fred C., and Jeffrey A. Dickinson. 1994. Development of sun compensation by honeybees: How partially experienced bees estimate the sun's course. *Proceedings of the National Academy of Sciences* 91:4471−4474.

Elman, Jeffrey L., Elizabeth A. Bates, Mark H. Johnson, Annette Karmiloff-Smith, Domenico Parisi, and Kim Plunkett. 1996. *Rethinking innateness: A connectionist perspective on development.* Cambridge, MA: MIT Press.

El-Naggar, Nadine, Andrew Ryzhikov, Laure Daviaud, Pranava Madhyastha, and Tillman Weyde. 2023. Formal and empirical studies of counting behaviour in ReLU RNNs. In *Proceedings of 16th edition of the International Conference on Grammatical Inference*, ed. by François Coste, Faissal Ouardi, and Guillaume Rabusseau. *Proceedings of Machine Learning Research* 217:199−222.

Engdahl, Elisabet. 1983. Parasitic gaps. *Linguistics and Philosophy* 6:5−34.

Foraker, Stephani, Terry Regier, Naveen Khetarpal, Amy Perfors, and Joshua Tenenbaum. 2009. Indirect evidence and the poverty of the stimulus: The case of anaphoric *one. Cognitive Science* 33:287−300.

Fox, Danny. 2000. *Economy and semantic interpretation.* Cambridge, MA: MIT Press.

Gordon, Peter. 1985. Level-ordering in lexical development. *Cognition* 21:73−93.

Gulordava, Kristina, Piotr Bojanowski, Edouard Grave, Tal Linzen, and Marco Baroni. 2018. Colorless green recurrent networks dream hierarchically. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, ed. by Marilyn Walker, Heng Ji, and Amanda Stent, 1195−1205. Association for Computational Linguistics.

Haïk, Isabelle. 1985. The syntax of operators. Doctoral dissertation, MIT.

Hart, Betty, and Todd R. Risley. 1995. *Meaningful differences in the everyday experience of young American children.* Baltimore, MD: Paul H. Brookes.

Hornstein, Norbert, and Jairo Nunes. 2002. On asymmetries between parasitic gap and across-the-board constructions. *Syntax* 5:26−54.

Hsu, Anne S., and Nick Chater. 2010. The logical problem of language acquisition: A probabilistic perspective. *Cognitive Science* 34:972−1016.

Huebner, Philip A., Elior Sulem, Cynthia Fisher, and Dan Roth. 2021. Baby-BERTa: Learning more grammar with small-scale child-directed language. In *Proceedings of the 25th Conference on Computational Natural Language Learning*, ed. by Arianna Bisazza and Omri Abend, 624−646. Association for Computational Linguistics.

Jozefowicz, Rafal, Oriol Vinyals, Mike Schuster, Noam Shazeer, and Yonghui Wu. 2016. Exploring the limits of language modeling. https://doi.org/10.48550/arXiv.1602.02410.

Katzir, Roni. 2023. Why large language models are poor theories of human linguistic cognition: A reply to Piantadosi. *Biolinguistics* 17. https://doi.org/10.5964/bioling.13153.

Kodner, Jordan, and Nitish Gupta. 2020. Overestimation of syntactic representation in neural language models. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, ed. by Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault, 1757−1762. Association for Computational Linguistics.

Kodner, Jordan, Sarah Payne, and Jeffrey Heinz. 2023. Why linguistics will thrive in the 21st century: A reply to Piantadosi (2023). https://doi.org/10.48550/arXiv:2308.03228.

Kuncoro, Adhiguna, Chris Dyer, John Hale, Dani Yogatama, Stephen Clark, and Phil Blunsom. 2018. LSTMs can learn syntax-sensitive dependencies well, but modeling structure makes them better. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long papers)*, ed. by Iryna Gurevych and Yusuke Miyao, 1426−1436. Association for Computational Linguistics.

Lan, Nur, Emmanuel Chemla, and Roni Katzir. 2023. Benchmarking neural network generalization for grammar induction. In *Proceedings of the 2023 CLASP Conference on Learning with Small Data*, ed. by Ellen Breitholtz, Shalom Lappin, Sharid Loáiciga, Nikolai Ilinykh, and Simon Dobnik, 131−140. Association for Computational Linguistics.

Lan, Nur, Emmanuel Chemla, and Roni Katzir. 2024. Bridging the empirical-theoretical gap in neural network formal language learning using minimum description length. https://doi.org/10.48550/arXiv.2402.10013.

Lan, Nur, Michal Geyer, Emmanuel Chemla, and Roni Katzir. 2022. Minimum description length recurrent neural networks. *Transactions of the Association for Computational Linguistics* 10:785−799.

Legate, Julie Anne, and Charles Yang. 2002. Empirical re-assessment of stimulus poverty arguments. *The Linguistic Review* 19:151−162.

Lewis, John D., and Jeffrey L. Elman. 2001. A connectionist investigation of linguistic arguments from the poverty of the stimulus: Learning the unlearnable. In *Proceedings of the Annual Meeting of the Cognitive Science Society* 23. https://escholarship.org/uc/item/3wv86519#author.

Lidz, Jeffrey, Sandra Waxman, and Jennifer Freedman. 2003. What infants know about syntax but couldn't have learned: Experimental evidence for syntactic structure at 18 months. *Cognition* 89:B65−B73.

Linzen, Tal, Emmanuel Dupoux, and Yoav Goldberg. 2016. Assessing the ability of LSTMs to learn syntax-sensitive dependencies. *Transactions of the Association for Computational Linguistics* 4:521−535.

MacWhinney, Brian. 2014. *The CHILDES Project: Tools for analyzing talk.* Vol. 2, *The database.* 3rd ed. New York: Psychology Press.

Marvin, Rebecca, and Tal Linzen. 2018. Targeted syntactic evaluation of language models. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, ed. by Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun'ichi Tsujii, 1192−1202. Association for Computational Linguistics.

Merrill, William, Gail Weiss, Yoav Goldberg, Roy Schwartz, Noah A. Smith, and Eran Yahav. 2020. A formal hierarchy of RNN architectures. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, ed. by Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault, 443−459. Association for Computational Linguistics.

Moro, Andrea, Matteo Greco, and Stefano F. Cappa. 2023. Large languages, impossible languages and human brains. *Cortex* 167:82−85.

Munn, Alan. 1992. A null operator analysis of ATB gaps. *The Linguistic Review* 9:1−26.

Nissenbaum, Jon. 2000. Investigations of covert phrase movement. Doctoral dissertation, MIT.

Ozaki, Satoru, Dan Yurovsky, and Lori Levin. 2022. How well do LSTM language models learn filler-gap dependencies? In *Proceedings of the Society for Computation in Linguistics 2022*, ed. by Allyson Ettinger, Tim Hunter, and Brandon Prickett, 76−88. Association for Computational Linguistics.

Pearl, Lisa, and Jon Sprouse. 2013. Syntactic islands and learning biases: Combining experimental syntax and computational modeling to investigate the language acquisition problem. *Language Acquisition* 20:23−68.

Perfors, Amy, Joshua Tenenbaum, and Terry Regier. 2011. The learnability of abstract syntactic principles. *Cognition* 118:306−338.

Phillips, Colin. 2013. On the nature of island constraints II: Language learning and innateness. In *Experimental syntax and island effects*, ed. by Jon Sprouse and Norbert Hornstein, 132−157. Cambridge: Cambridge University Press.

Piantadosi, Steven T. 2023. Modern language models refute Chomsky's approach to language. Ms. https://ling.auf.net/lingbuzz/007180.

Postal, Paul M. 1993. Parasitic gaps and the across-the-board phenomenon. *Linguistic Inquiry* 24:735−754.

Pullum, Geoffrey K., and Barbara Scholz. 2002. Empirical assessment of stimulus poverty arguments. *The Linguistic Review* 19:9−50.

Radford, Alec, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI blog* 1:9.

Rawski, Jon, and Lucie Baumont. 2023. Modern language models refute nothing. Ms. https://lingbuzz.net/lingbuzz/007203.

Rawski, Jonathan, and Jeffrey Heinz. 2019. No free lunch in linguistics or machine learning: Response to Pater. *Language* 95:e125−e135.

Reali, Florencia, and Morten Christiansen. 2005. Uncovering the richness of the stimulus: Structure dependence and indirect statistical evidence. *Cognitive Science* 29:1007−1028.

Ross, John R. 1967. Constraints on variables in syntax. Doctoral dissertation, MIT.

Siegelmann, Hava T., and Eduardo D. Sontag. 1991. Turing computability with neural nets. *Applied Mathematics Letters* 4:77−80.

Siegelmann, Hava T., and Eduardo D. Sontag. 1995. On the computational power of neural nets. *Journal of Computer and System Sciences* 50:132−150.

Sprouse, Jon, Beracah Yankama, Sagar Indurkhya, Sandiway Fong, and Robert C. Berwick. 2018. Colorless green ideas do sleep furiously: Gradient acceptability and the nature of the grammar. *The Linguistic Review* 35:575−599.

Strobl, Lena, William Merrill, Gail Weiss, David Chiang, and Dana Angluin. 2023. Transformers as recognizers of formal languages: A survey on expressivity. https://doi.org/10.48550/arXiv.2311.00208.

Vázquez Martínez, Héctor, Annika Lea Heuser, Charles Yang, and Jordan Kodner. 2023. Evaluating neural language models as cognitive models of language acquisition. In *Proceedings of the 1st GenBench Workshop on (Benchmarking) Generalisation in NLP*, ed. by Dieuwke Hupkes, Verna Dankers, Khuyagbaatar Batsuren, Koustuv Sinha, Amirhossein Kazemnejad, Christos Christodoulopoulos, Ryan Cotterell, and Elia Bruni, 48−64. Association for Computational Linguistics.

Wagers, Matthew W., Ellen F. Lau, and Colin Phillips. 2009. Agreement attraction in comprehension: Representations and processes. *Journal of Memory and Language* 61:206−237.

Wang, Ben, and Aran Komatsuzaki. 2021. GPT-J-6B: A 6 billion parameter autoregressive language model. https://github.com/kingoflolz/mesh-transformer-jax/?tab=readme-ov-file#citation.

Warstadt, Alex, and Samuel R. Bowman. 2022. What artificial neural networks can tell us about human language acquisition. In *Algebraic structures in natural language*, ed. by Shalom Lappin and Jean-Philippe Bernardy, 17−60. Boca Raton, FL: CRC Press.

Warstadt, Alex, Alicia Parrish, Haokun Liu, Anhad Mohananey, Wei Peng, Sheng-Fu Wang, and Samuel R. Bowman. 2020. BLiMP: The benchmark of linguistic minimal pairs for English. *Transactions of the Association for Computational Linguistics* 8:377−392.

Weiss, Gail, Yoav Goldberg, and Eran Yahav. 2018. On the practical computational power of finite precision RNNs for language recognition. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short papers)*, ed. by Iryna Gurevych and Yusuke Miyao, 740−745. Association for Computational Linguistics.

Wilcox, Ethan, Roger Levy, and Richard Futrell. 2019. What syntactic structures block dependencies in RNN language models? https://doi.org/10.48550/arXiv.1905.10431.

Wilcox, Ethan, Roger Levy, Takashi Morita, and Richard Futrell. 2018. What do RNN language models learn about filler-gap dependencies? In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, ed. by Tal Linzen, Grzegorz Chrupala, and Afra Alishahi, 211−221. Association for Computational Linguistics.

Wilcox, Ethan Gotlieb, Richard Futrell, and Roger Levy. 2024. Using computational models to test syntactic learnability. *Linguistic Inquiry* 55.4.

Williams, Edwin S. 1977. Across-the-board application of rules. *Linguistic Inquiry* 8:419−423.

Williams, Edwin S. 1990. The ATB-theory of parasitic gaps. *The Linguistic Review* 6:265−279.

Wilson, Colin. 2006. Learning phonology with substantive bias: An experimental and computational study of velar palatalization. *Cognitive Science* 30:945−982.

Yedetore, Aditya, Tal Linzen, Robert Frank, and R. Thomas McCoy. 2023. How poor is the stimulus? Evaluating hierarchical generalization in neural networks trained on child-directed speech. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long papers)*, ed. by Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, 9370−9393. Association for Computational Linguistics.

*Nur Lan*
*Laboratoire de Sciences Cognitives et Psycholinguistique*
*Ecole Normale Supérieure*
*and*
*Department of Linguistics*
*Tel Aviv University*

*nur.lan@ens.psl.eu*

*Emmanuel Chemla*
*Ecole Normale Supérieure, EHESS, PSL University, CNRS*

*emmanuel.chemla@ens.psl.eu*

*Roni Katzir*
*Department of Linguistics*
*and*
*Sagol School of Neuroscience*
*Tel Aviv University*

*rkatzir@tauex.tau.ac.il*