Searching for a Sign: *A Semiological Analysis of LLMs*.

Ella Harrigan Swarthmore College December 2024 ABSTRACT

# 2. BACKGROUND ON LARGE LANGUAGE MODELS

- 3. SIGNIFIER AND SIGNIFIED
  - 3.1 Sound-image and Concept
  - 3.2 Convention and the Sign
  - 3.3 Linear and Non-Linear Dimensions
  - 3.4 Form and Substance

# 4. MEANING and COMMUNICATIVE INTENT

- 4.1 Applying Form and Substance
- <u>4.2 Communicative intent</u>
- 4.3 Understanding
- 5. MY EXPERIMENT
  - 5.1 Vocabulary and phrasing
  - 5.2 Methods
  - 5.3 Rationale behind my experimentation
  - 5.4 Limitations
  - 5.5 Results
    - A. Certainty-Level and Bias
    - B. Certainty Level and Gaslight Level
    - C. Incoherency Levels and Gaslight Level
    - D. Rationales Frequency and Answer Type
    - E. Incoherency Levels across Various Axes
- 6. DISCUSSION
  - 6.1 Experiment Discussion
  - 6.2 Conclusion
- <u>REFERENCES</u>

#### ABSTRACT

Large Language Models, also known as LLMs, are language models which work with an incomplete sign: the signifier without the signified. According to a Saussurian method of analysis, they operate on the strata of expression, but not on that of concept. We can call this the expression-concept gap. This is because the LLM lacks hierarchical structure and an ability to access extralinguistic information, both of which are necessary for the signified and concept of the sign. This loss of the signified causes LLM outputs to behave in sometimes strange ways, outputting responses that appear to be nonsensical, though grammatical, utterances. In my experimentation, I analyze some of these responses. Specifically, using a four-tiered "gaslighting" system, I look at the "justifications" generated by the LLM after giving me its initial response. These responses differ in interesting ways from human speech. Specifically, they show a higher rate of gendered bias, increased incoherency per gaslight-level, and an inability to demonstrate "learning" from past errors. These are all effects of the expression-concept gap.

#### **INTRODUCTION**

The transformation of objects in space, or objects in time, To objects outside either, but still tactile, still precise... It's always the same problem— Nothing's more abstract, more unreal, than what we actually see. Our job is to make it otherwise.

-Charles Wright, 1998

I started thinking about what became this thesis when I realized that I kept noticing connections between ideas in different places and fields — in late nineteenth century Switzerland with Saussure, in mid-twentieth century Germany with Wittgenstein, then France with Barthes, and in twenty-first century America with Emily Bender. These ideas cropped up in Continental and Analytic philosophy, in semiology, in NLP, and got digested into my introductory linguistics curriculum. They all circled around what we would call in linguistics and semiotics *the sign*, though the vocabulary changes depending on who you ask. And they were all about, roughly, what Saussure would call the distinction between signifier and signified (the word-image and concept). I wanted to compare the different ways Saussure's original ideas about signification have been woven into various parts of meaning and language analysis over the years.

While I was having this revelation about Saussure and Barthes and Wittgenstein, I was on summer vacation. I had a grant to write poetry. I spent my days writing sonnets and reading

semiotics, and then at night I went to a lot of strange little work-parties and gatherings with a very different crowd than I was used to. I was in San Francisco that summer, living with my girlfriend who, along with everyone we knew, was all in on the AI gold rush.

At those parties there was a lot of tipsy, half-feverish talk about important-sounding words like AGI and p-doom and so on, which were quite a lot of fun to say but began to drive me a little crazy. Everyone was excited, but we weren't sure quite what it was that we were so excited about, though we knew we liked the term *black box*, because it was mysterious, alliterative, and scientific. That box seemed capable of holding anything we wanted it to — AGI could burst out of it any day now, we thought, like Athena had done once, out of Zeus' skull.

A problem with LLMs is that, though their outputs look really good when you first read them, they often fall apart under scrutiny. Sentences that appear well-formed are revealed as nonsensical, citations lead to papers that don't exist, etc. The LLM is excellent at imitating the form of written human speech, at least grammatically — things like discourse markers are a bit more difficult. We can say that they are very good at expressing what appear to be well-formed utterances. But the actual substance of what those utterances say is often totally bizarre.

This is what we can call the expression-content gap.<sup>1</sup> This is the fact that large language models such as ChatGPT appear to have perfect (technically, near-perfect) expression. This means they produce outputs which take the expression of grammatical sentences, with words in the sentence which make sense together and in relative context. They take the form of our own discourse. However, when you look at the actual meanings of the words in the LLM outputs, you'll start to find what appear to be mistakes. Everyone who has been even slightly paying attention has heard grumblings about LLM "hallucinations," or else seen examples of ChatGPT failing to answer seemingly simple questions — while the outputs still have all the discourse markers of confidence.

To try and understand this phenomenon, I decided to apply a sign-based analysis of LLMs, drawing on work from semiology as well as related fields such as NLP (natural language processing), philosophy of language, and formal semantics. Then, I ran an experiment where I asked Chat-GPT two variants on the same question — a question which involved pronominal ambiguity — and analyzed its responses, looking at how that data reflected the expression-concept gap.

In section 2 of this paper, I provide background on language models, so that we have the technical framework for the rest of the paper. In section 3, I provide that same technical framework for the theories of signification I draw on throughout. Then, in section 4, I compare

<sup>&</sup>lt;sup>1</sup> This is similar but not identical to the idea of the *form-meaning* gap as proposed by Ted Gibson (Fridman, 2024).

different theories across time which all have to do with questions of expression and concept, and apply those generally to large language models. In section 5, I present my experiment and show results. And in section 6, I discuss my experimental results, situating them in the general body of research of modern NLP and AI work being done on LLM 'errors', and connect them to my previous semiological analysis.

# 2. BACKGROUND ON LARGE LANGUAGE MODELS

This paper is about large language models, also known as LLMs. LLMs are transformer-based language models which function via layers of neural networks. To explain what all of those words mean, we can start at the beginning.

Sometime in the twentieth century, we started working on natural language processing, also known as NLP. According to Wolf (2018), in natural language processing, a language model is a model which, given context, can generate text and predict the next word in a sentence. The language models I will be discussing in this paper are those that are trained on and output text. Raw text in and of itself codes no meaning; it is a collection of text units without their signifieds. And so, generally, the question for language models is this: given a dataset of raw text — what we can think of semiologically as signifiers without their signifieds — can meaning be learned?

The history of language models in NLP can be roughly summarized as a move from symbolic NLP to statistical NLP, and then, in the world of statistical processing, moving from classic statistical processing to "neural network" statistical processing, and then, as a subset of neural network language models, we arrive at transformer-based neural network models, also known as generative pre-trained transformers, also known as GPT. These GPT models are so-called large language models (LLMs) because they're trained on huge amounts of data; they are the type of model I work with in this paper. Specifically, my experiment is on Open AI's language model, known as ChatGPT, version 3.5.

To understand large language models properly — or at least, to understand them insofar as is relevant for a semiological analysis — I found it was important to have a bit of a historical background. LLMs are often shrouded in mystery — everyone is always talking about the *black box*, how no one knows how they work, et cetera, and while this is true to an extent, looking at the development of language models can help clear away the feeling of "magic." This is important because many of the debates and frameworks laid out in the mid and late twentieth century around language models are still relevant today; although much of the layout and formal mechanisms have changed, the fundamental nature of the inputs — text without a referent — have stayed the same.

So, a brief overview: the advent of language models can be debated depending on what you count as NLP and what you count as LMs; if you think about a language model as anything modeling language, then the early linguists and semioticians such as Saussure and Barthes count. So do all the generative syntax models, rules-based and OT based phonology models, even diachronic rules-based orderings in historical linguistics. But for practical purposes, people usually use "language model" to mean a computer model, which may or may not incorporate the structures and formal grammars from non-computational theoretical linguistics.

With that in mind, we can place the beginning of language models in the mid twentieth century. These language models were symbolic, which means they were rules-based. As an example, we can look at some of the rules used in the early chatbot ELIZA, which was made by Joseph Weizenbaum in the 1960s, and which simulated a Rogerian psychiatrist. (Weizenbaum, 1966) ELIZA is famous because it was one of the first LMs to somewhat pass the Turing Test. (Jurafsky and Martin, 2024; Turing, 1950) I pulled the following example from the book Chatbots and Dialogue Systems by Daniel Jurafsky & James H. Martin, but the code itself is obviously made by Joseph Weizenbaum.

(1)

Rule: (.\*) YOU (.\*) ME -> WHAT MAKES YOU THINK I \2 YOU User sentence: You hate me

System response. What makes you think I hate you

This isn't a general syntactic rule like the ones we have in Chomskyan syntax, for example. It's specific to the specific user sentence type. These rules are also encoded by hand one-by-one; Weizenbaum had to type in this specific rule to get this specific system response with this specific type of user sentence. There's no generative grammar. But symbolic NLP also includes language models where more general grammar and syntax rules were attempted to be encoded. When people made symbolic NLPs, they tried to encode syntactic and semantic information directly into the functioning mechanisms of the model.

The problem with symbolic NLPs is simply that they did not work very well. As computer processing power grew, symbolic NLP strategies started to be phased out and replaced with statistical ones. In addition, with the advent of the world-wide-web, the amount of data available for datasets grew tremendously. A lot of this work was going on at IBM in the late 20th century. In statistical NLP, specific rules were phased out in favor of "machine learning", where you give the machine data and then the machine generalizes from that data. The algorithm is learning pattern recognition and then extrapolating from those patterns. The word learn here means "A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P if its performance at tasks in T, as measured by P, improves with experience E." (Mitchell, 1997) This is an operational definition; the measure of learning is based on how the machine operates as opposed to any specific "cognitive" function.

Word n-gram models are the simplest type of statistical LM, and although LMs have gotten much (much) more complex since then, their core idea — predicting the next word based on the previous one(s) — has followed the field ever since. n-gram models decide the likelihood of a given word coming next in the sequence based on the context of previous words. A bigram model decides the likelihood of a given word coming next based on the context of just the following word; a trigram model decides the likelihood of a given word s, et cetera. So an n-gram model looks n-1 words into the past. (Jurafsky and Martin, 2024)

Every language model that's not symbolic is based on statistics and probability, but the ways in which they do it vary. Neural-network based language models still work with pattern recognition and data-crunching, but they use a different toolkit from the simpler word n-gram models. The toolkit is artificial neurons. These are functions which receive a group of inputs, then weight and sum the input(s) along with a bias term, then output them as one output. (Jurasky and Martin, 2024) There are various layers of these neural-networks — there's always going to be an initial input layer and then an output layer, but there are also "hidden layers" in between (any layer that's not the input or output layers). Raw data has to pass through layer and layer of neural networks before it can reach the final output; this process is called "deep learning."

Figure 1. Figure taken from Jurafsky & Martin Chapter 7.



Here we see a simple two-layer neural network (we don't count the input layer in the number of layers.) 'h' stands for 'hidden.'

Moving forward, we arrive at transformer-based language models, which are a type of neural-network-based statistical language model. This is what ChatGPT is.



Figure 2.

Transformers are a specific type of neural network which incorporate attention as a focusing framework. They were introduced in the paper *Attention is All You Need* (Vaswani et al., 2017), which completely transformed the field and ushered in the current wave of language models. Attention allows the model to assign greater weight to certain components in a sequence relative to other components in the same sequence, allowing the LLM to integrate tokens across longer sequences. (Jurafsky & Martin, 2024) It does this by using the context of surrounding components to understand how the components relate to each other over large spans. (Jurafsky & Martin, 2024) Specifically, transformers use something called self-attention, which is "an attention mechanism relating different positions of a single sequence in order to compute a representation of the sequence." (Vaswani et al., 2017, pg. 2)

Attention is relevant to my experiment in section 5 below, because it is how LLMs are supposed to deal with things like pronoun antecedents, where the relevant information to the pronoun is not necessarily directly preceding it in the sentence. The idea is that the pronoun should be embedded with the contextual representation connecting it to its antecedent; the LLM should pay "attention" to that antecedent more than the other preceding units. These units are called tokens, and we'll circle back to them in more depth in a moment, but you can think of them as something

roughly analogous to words, morphemes, or letters, depending on the specifics of the language model.

The important part to understand is that transformer-based language models do not necessarily always choose the token with the absolute highest probability, though they sometimes do. However, the LLMs that we're talking about do always generate linearly, meaning one token after another.

Tokens are a character, word, or string (Bender and Gebru., 2021).

Large language models are trained on gigantic amounts of data — this is why they're "large." Unlike language models trained on word-vectors (the pre-transformers neural network model), LLMs are able to improve performance by feeding more data and parameters into the set. (Bender and Gebru, 2021) There's a focus on making language models as large as possible, increasing both the number of parameters and the size of the datasets, which has been criticized by scholars such as Emily Bender, Timnit Gebru, and Margaret Mitchell, who warn against the financial and environment repercussions, as well as the bias which is encoded into such huge and unmanaged data-sets. They argue that we should focus on curating and documenting smaller and higher quality data sets. (Bender and Gebru, 2021).

LLMs incorporate stochasticity into their algorithms. "Stochastic" here refers to an element of chance or randomness; given an identical input 20 times, an LLM will output a different response each time. This is due to the fact that the LLM does not always pick the word with the exact highest probability; the model is fundamentally probabilistic but encodes stochastic variations. This improves the perceived quality of responses. *Temperature*' is a mechanism used to mitigate stochasticity — turning up the temperature of an LLM increases the chances that a lower-probability option is selected. This is important because predictive algorithms can under-reflect low-likelihood outcomes — for instance, a .001% outcome in the dataset might be a .0005% in the algorithm. This is problematic because it means that LLMs could not only replicate but also amplify bias. We've already seen many examples of this. (Kotek et al., 2023; Abid et al., 2021; Blodgett et al., 2020; Buolamwi and Gebru, 2018)

Speaking of bias, it's important to note that LLMs such as ChatGPT have a linear bias. (Yedetore et al., 2023). We can see this in the examples below, from the 2023 paper *How Poor is the Stimulus?*, which compares children's speech assumptions and biases with LLM's assumptions and biases<sup>2</sup>. They do this by using yes-no formatted questions, which always correspond to one of two types:

<sup>&</sup>lt;sup>2</sup> This phrasing is anthropomorphizing; the LLM doesn't have its own assumptions and biases, so to speak. This just means that the model creating the LLM works in certain ways to generalize grammar.

(2) "HIERARCHICAL Q: The auxiliary at the start of a yes/no question corresponds to the main auxiliary of the corresponding declarative.

LINEAR Q: The auxiliary at the start of a yes/no question corresponds to the first auxiliary of the corresponding declarative." (pg. 2)

They find that models fail to generalize with the HIERARCHICAL Q, instead generating based on linear order. This results in grammar mistakes from the LM. An example (pulled from the paper):

- (3) "a. The boy who has talked can read.
- b. Can the boy who has talked read?
- c. \*Has the boy who talked can read?"

This is interesting for two reasons. First, the LLM's linear bias causes it to sometimes make grammatical mistakes in edge cases where linear and hierarchical bias lead to different outputs. This means that, in a slight addendum from the initial observation I pointed out in the introduction, large language models do not have completely perfect expression. Small mistakes are made.

This leads me into the second, very important reason why evidence of linear bias in LLMs is so interesting. These small grammar errors show that the LLM's internal grammar modeling probably generalizes based on the patterning of linear progressions, instead of creating sentence structure through a hierarchical tree model. Though this may seem obvious given the statistics and stochasticity-based background I've given so far, it's not. This is because the internal workings of an LLM are something of a black box — meaning, no one is quite sure how they come to the things they do. Many people, especially AI proponents, speculate that, although we are only feeding the LLM huge amounts of data, somewhere in the pattern-recognition process it is creating its own structured hierarchical grammar. This linear bias is evidence against that. This, in turn, helps my thesis get closer to something we will be exploring more in depth later — the question of whether or not the LLM has the form of the signified.

Finally, I want to touch on the elephant in the room. Azaria and Mitchell theorize that hallucinations occur because "an LLM generates a token at a time, and it "commits" to each token generated. Therefore, even if maximizing the likelihood of each token given the previous tokens, the overall likelihood of the complete statement may be low" (2023, pg. 2) They give a really great example: "Consider the following ...: "Tiztoutine is a town in Africa located in the republic of Niger." Indeed, Tiztoutine is a town in Africa, and many countries' names in Africa begin with "the republic of". However, Tiztoutine is located in Morocco, which is not a republic, so once the LLM commits to "the republic of", it cannot complete the sentence using "Morocco", but completes it with "Niger". In addition, committing to a word at a time may lead the LLM to be required to complete a sentence that it simply does not know how to complete." (pg. 2)

This is true even though the LLM doesn't necessarily "use the maximal probability for the next word, but to sample according to the distribution over the words." (pg. 2)

### **3. SIGNIFIER AND SIGNIFIED**

#### 3.1 Sound-image and Concept

My work in this thesis draws upon Saussure's theory of signification, in which the sign is made up of two parts: the signifier a signifier — the 'sound image' — and a signified — the 'concept.' (Saussure, 1916) The sign is what unites the two.

The main claim of this thesis is that large language models work with a severed sign: the signifier without the signified. This results in expression without content, what we can roughly call form without meaning.<sup>3</sup>

In semiology, the signifier and the signified are wholly abstract concepts — neither corresponds directly with the world outside language. The signifier, aka the concept, is still abstract. It is not the thing itself. For instance, the sign 'tree' is made up of the sound-image [tii] and the concept of a tree. Neither the signifier or the signified are an actual tree itself, existing in the real world.

There is one place in my analysis where Saussure's idea of both the signifier and the signified being totally abstract does grate a little. I propose that the LLM outputs — the tokens we discussed earlier — correspond to the signifier portion of the sign. However, technically, the signifier is abstract, and thus corresponds *not* to surface-form outputs and phonetics, but to the underlying form of the word and phonology. This is because the question of whether LLMs have internal, abstract phonological structure — for example, what we could call underlying forms — is not as simple as it might seem. To temporarily sidestep this question, it works for this portion of the analysis to assign the signifier to the surface structure. Later, with Barthes, we will specify these terms more anyway.

Saussure argues the concept and the sound-image are inextricably linked; this is also somewhat problematic for my analysis. In my analysis, the signified and signified of LLMs are split, with the latter being provided to the outputted tokens at the moment that a person reads them and thus associates meaning onto them. An LLM output is the sound-image without the concept. Of course, large language models didn't exist when Saussure was putting together this theory, and even now they represent an edge case.

<sup>&</sup>lt;sup>3</sup> The reason I opt for the more technical and awkward sounding "expression" and "concept" vocabulary terms is that "form" means something different and more specific in the Barsthian framework that I draw on later in the paper, and "meaning" is a broad and confusing term, whose different and competing definitions we will also spend a significant amount of page-space unpacking.

A way to get around the problem of awkward theoretical fit is by using Emily Bender and Alexander Koller's approach. They are two contemporary linguists working on natural language processing and critiquing the discourse around large language models. In their 2020 paper, *Climbing towards NLU: On Meaning, Form, and Understanding in the Age of Data*, they circumvent some of the theoretical baggage around semiological vocabulary by creating a new one, using similar concepts. Instead of concept and sound image, they use intent (I) and expression (E). Each sign is made up of an (e,i) pair. They define "meaning" as  $M \subseteq E \times I$ . This means that "meaning" is a subset of the set of all ordered pairs of the form (e, i) where e is an element of E and i is an element of I. This means, basically, that "meaning" is part of, but not necessarily all of, signification.

Intent here is equivalent to "concept" — while Saussure and Barthes define this half of the sign in terms of "substance" or "content", supposing a sort of inherent notional meaning, Bender and Koller use a theory of communicative intent, where the other half of the sign is not defined by any independent meaning, but rather by the purpose it was used by the speaker. This crucial difference in conceptualizing meaning is something we will explore in greater depth in section 4. I find it simpler to use Saussure and Barthes' vocabulary because in the following section, I introduce Barthes' four-pronged theory of signification, and it is easier to superimpose Bender and Koller"s terms onto Barthes', than to do vice-versa.

So, instead of coming up with new vocabulary, I square Saussure's idea of linkage with my own analysis by arguing that the concept and sound-image are inextricably linked in one direction only. The concept cannot exist without the sound-image — the concept of 'tree' cannot exist without the word to name it. However, the series of sounds [tii] can exist without a concept to go along with them. The sound-image and concept need to be linked together at the moment of creation — otherwise the concept of trees or the utterance [tii] to describe them would not have been created — but the concept can disappear, leaving only the sound image. For instance, when I was nine years old I moved to France, and was surrounded by a language that I was very good at mimicking but could only poorly understand. I had a teacher who would say at least once a class, as a sort of verbal tic, a word I only understood as [ã?tuka:]. I was able to more-or-less understand the environment in which she'd say it (after a long period of speech, before a pause), and could even predict with some accuracy when it would be said, but had no idea what it meant. One day after class I asked her and she burst into laughter. What I had thought was a word was actually three: en tout cas, meaning 'in any case.' After that moment, the concept appeared to me, and the way I understood [ã?tuka:] changed; mentally, I imagined an 'n' after ã and pauses between the words (though neither appeared in the surface form). Saussure might say that the sound-image was created for me only then, where before I had just had a record of the surface articulatory form.

However, because my conception of the LLM conflates the surface form with the sound-image, this means in my argument, the sound-image can exist alone. The tokens that transformer-type language models work with have more in common with ã?tuka: than *en tout cas*; they analogize more directly with literal sound than to sound-image. Because of this, as well as the one-way linkage that a more abstract definition engenders, for the purposes of my analysis I will be using 'signifier' to mean the articulation of actual sound, or, in text, the literal order of letters on a screen.

#### 3.2 Convention and the Sign

Famously, Saussure writes about "the arbitrariness of the sign." What he means could be more precisely (and less elegantly) described as "the arbitrariness of the connection between the concept and sound image." This is to say, there is nothing more inherently tree-like about the word arbre (French), tree (English), or  $\pi$  (Mandarin), and this arbitrariness comes from the union of the phoneme-collection tri and the concepts in our minds of woody, perennial plants.

However, although the sign itself is arbitrary, its use in the social world is not. When we are speaking a language, we are buying into its contract, and abiding by those arbitrary signifier-signified pairings that have been handed down to us. Without this buy-in — if we created language as we spoke it — discourse would be impossible. As Roland Barthes points out in *Elements of Semiology*, it is the orderly repetition of signs, each neatly composed of a signifier-signified pair, that makes discourse possible. (1964)

This buy-in is what Bender and Koller (2020) refer to as the C variable: conventional meaning. This is important because it shows the way that the content of a sign is connected to extralinguistic factors: its meaning is dependent on what we have decided its meaning is in the social contract. With some rough edges, we generally all know what the word "tree" means. Our language system, and thus our system of signification, is a social one, even while it remains abstract.

#### **3.3 Linear and Non-Linear Dimensions**

Saussure writes that the signifier is linear because, being a sound-image, it unfolds in one direction. The same is true even when we reframe the signifier as being the token output from an LLM. If this idea of speech progressing linearly sounds familiar, it's because it is — it's the movement statistical language models have, with one token following and informing the next. This linear generation — one token then another — is analogous to speech's linear progression in time. The signifier is a span which can be represented in one dimension.

According to Saussure, the signified, on the other hand, is bidirectional: it goes from the sound-image to the concept and vice versa; the signifier is the part of the sign which connects to what Wittgenstein and Bender might call extralinguistic concepts — concepts which connect in some way to the broader world and require knowledge of ideas outside language. In large language models, where we lack this extralinguistic information in the sign itself, we can see a unidirectional movement of the signified: a human reader interacting with a string of tokens output from a language model, then attaches the signifier to the output: there is a movement inward, towards the expression.

# 3.4 Form and Substance

In *Elements of Semiology*, Barthes tweaks Saussure's terminology a little, writing that the sign is not just made of the two elements *expression* and *concept*, but rather two planes: the *plane of expression* and the *plane of concept*, each of which have two strata: the stratum of substance, and the stratum of form.





It's worth spending some time with this concept of four-pronged signification, both because it's difficult and because it informs the bulk of the second part of my analysis, where I apply the form and substance of content onto ideas of philosophy of meaning in language.

Form is what can be described exhaustively by linguistics without needing any extralinguistic information, substance is what cannot be. "Extralinguistic" here is a bit of a confusingly named category, because what makes something "linguistic" is not an agreed upon list of traits. Suffice

it to say, in this context *extralinguistic* means outside of abstract signification: phonetics and articulation contain extralinguistic qualities in this context because they require mouths, air for sound waves to move through, physics, biology, etc. Similarly, emotional, ideological, and notional aspects are extralinguistic even though they use language as a means of expression, because they also contain information and inputs from the "real world": for instance, a notional definition of 'tree' is informed by the speaker's memories of trees they have seen.

Roughly, the plane of expression (the signifier) corresponds to phonetics, phonology, and syntax, while the plane of content (the signified) corresponds to semantics, pragmatics, and discourse.

In the following section, I will take Barthes' four-pronged formulation and apply it to both large language models and other language theorists; from there, I will argue that the severed sign large language models work with is the substance of expression, without the form or the substance of content.<sup>4</sup> Furthermore, I will connect the form and substance of concept to linguistic ideas of meaning and communicative intent, and apply these to large language models. Going forward, I will use the terms *signifier* and *expression* interchangeably, as well as *signified* and *concept*.

FORM of signified	SUBSTANCE of signified
Barthes: the formal organization of signs among themselves, by presence or absence of a semantic mark	Barthes: extralinguistic factors: defined as ideological, emotional, and <b>notional</b> aspects
Wittgenstein: verbal meaning; "that which takes us from sign to sign but no further"	Wittgenstein: ostensive meaning
ex: verbal meaning of 'tree' is "function of index onto truth values that map all things that are tree onto true and all things that aren't tree onto false)	ex: look over there! that's a tree! and that one is too! and that one!
Bender & Koller: E (expression)	Bender & Koller: I (intent) + Bender & Koller: S (standing meaning)
Bender & Koller: e = natural language	Bender & Koller: i = communicative intent e

# 4. MEANING and COMMUNICATIVE INTENT

<sup>&</sup>lt;sup>4</sup> Whether or not the LLM has the form of expression is an interesting question and one well worth exploring, but not one that fits into the scope of this paper; this thesis is more interested in the 'concept' part of the equation — the signified.

expression	is meant to invoke + Bender & Koller: s = conventional meaning of a given sign
Ella: LF, propositional and predicate logic, hierarchical structure	Ella: property-based and prototype meanings (basically, <i>what's a tree? well, i know it when i see it.</i> )
Yedetore et al: "meaning"	

# 4.1 Applying Form and Substance

Here we have a T-chart of various concepts, all applied to Barthes' form/substance paradigm. As we apply these one-by-one, both to each other and to large language models, we can see that LLMs lack the substance of the signified, and have only partial form of the signified.

The first concept in each category we have already gone over: these are the form and substance of the signified according to Barthes.

I would also argue that logical form (LF), propositional, and predicate logic are frameworks we can put under the form part of the signified. They are concerned with meaning; they have to do with semantics, not syntax. However, they also avoid anything extralinguistic; they organize themselves based on what forms do or do not map onto certain truth values and sets. For example, in predicate logic, the phrase "that is a tree" is true if P(x) is true, where 'P' is the set of trees and 'x' is "that." Thus, the meaning of the word 'tree' is a function of index onto truth values that map all things that are tree onto true and all things that aren't tree onto false. In propositional logic, the phrase "that is a tree" is true if T is true, where T stands for the phrase "that is a tree."

These are all useful frameworks, and they all stay in the realm of signs: they fall under the category of what Wittgenstein (1958) called "verbal meaning", which is to say that they take us from sign to sign and no further. Unlike the substance of the signified, which is still wholly abstract, but gestures towards the referent, these logic-based forms do not gesture. They have to do with how the signs are organized amongst themselves, just as Barthes said, with the present-or-absent semantic mark being the truth condition.

The question is whether or not transformer-based, neural-network language models are able to work with these signifiers. What excites people so much about LLMs is the idea that, given all

these forms, they are able to apply meaning to them. Bender and Koller's counterargument is that they don't get any extralinguistic information in the first place, so they can't have meaning. Whether or not this argument is convincing depends, of course, on which part of meaning (here, "meaning" is analogous to the signified) you are talking about. The confusion comes from different conceptions of the word "meaning", a confusion which can be resolved by using semiological vocabulary to clarify our terms.

When we think about these *formal* types of meaning, it is possible that they could be divined, by some sort of something, sometime, from just signifieds. This is because the form portion of the signified does not require any extralinguistic information, and that is the thing the LLM inputs do not have. In addition, it is possible that some sort of something, somewhere, could be given these structures and models, and then apply them to the inputs received, and be able to create structures for the tokens to go into. Basically: given heaps of probabilistic data, can you intuit the internal structures of the data, and from there, understand the signs in relation to each other? Whether or not some version of this is happening with LLMs is a much more open question than whether or not an LLM has access to extralinguistic information — it obviously does not.

Confusing of the form and substance strata of the signified, in discourse between people, as well as internally in one's own mind, is responsible for some of the more frustrating parts of the LLM debate. At the same time, however, it is important to ask oneself how much it would matter to the practical ability of the LLM to "think", so to speak, if it was able to work with hierarchical structures and logical forms. Does that matter? Or is what actually matters when we talk about meaning only the extralinguistic parts of meaning, which are impossible for the LLM to access?

I do think it's important to engage with whether or not LLMs have the form of the signified, if only because it sheds light on the internal workings of LLMs in ways that are useful when trying to figure out its mistakes.

One of the ways to think about this question has to do with hierarchical structure. Humans have what we can call a hierarchical bias, as alluded to in the background section. Out of the scope of this paper, this is something syntacticians like to use as potential evidence for tree-structured generative grammar. (McCoy et al, 2020) Hierarchical bias is proof that there is something going on during language acquisition that allows children to understand grammar in ways that don't just have to do with probability or pattern recognition. (Yedetore et al, 2023; McCoy et al, 2020; Mulligan et al, 2021). The extent to which this happens is debatable, but there is compelling evidence that something is going on. Again, though, I am not personally invested in — nor am I anywhere near qualified to address — the debate about language structures exactly as they appear in the human brain. The important part is just that hierarchical bias points towards structure happening on the langue level of human language, and that, when converted into parole, something structural remains in the way we synthesize language.

In short, hierarchical bias is evidence of some sort of formal structure happening in human language production; it shows there's some sort of organization of the signs amongst themselves in ways more complex that data-crunching. On the other hand, LLMs do not have this hierarchical bias. (Yetedore et al, 2023; Geng et al, 2025). This does not entail that they lack the ability, then, to organize the signs amongst themselves, but it also provides no evidence that they do.

However — and this is what is really interesting — Yetedore I've been citing these last few paragraphs claims to compare LLMs trained on "form-only" to LLMs trained on "form and meaning," and shows that the latter are better able to output a token-chain which resembles a grammatically correct sentence based on hierarchical bias. When I first came across this paper, it stood out to me because of the way they framed their experimentation as training the LLM on meaning. All my reading and analysis prior to that moment stood on the grounds that LLMs were trained only on form. The debate between believers and sceptics in LLM's ability to reason was over the question of whether something was happening in between those initial input vectors and the final LLM output token chains; if the LLM was able to create a sort of meaning and reasoning structure just by virtue of absorbing so much data. I hadn't heard of a technique that allowed meaning to actually be put into the training data itself.

Upon reading the paper, I realized what had happened: again, the definitions of "meaning" that I was using and the authors of the paper were using didn't match up. They trained the LLM on "meaning" by feeding it data in predicate logic English, and trained it on "form" by feeding it data in surface-form, grammatical English. Then, they asked the LLM to make grammar formulations in situations where linear bias gives an ungrammatical answer, and hierarchical bias gives a grammatical one. They found that the LLM trained on both surface-form English and predicate-logic English did better at giving outputs that correspond to the grammatical answer.

To them, this is evidence both of the LLM having a sort of intelligence and also that Chomsky is right about grammar being structural. On the first point, I am unconvinced, and on the second, it's not really my problem. However, this paper still heavily influenced my analysis because it provided a perfect example of what I am *not* defining meaning as. Just feeding in some predicate logic examples is not enough! When we feed the LLM data in different forms (predicate logic and "standard spoken"), we are just training it on two different forms. Personally, I would guess that the LLM is just getting better at predicting token strings when it has more types of data about what order the tokens tend to go in. But even if that's not the case, and the paper does show some sort of reasoning / learning capacity for LLMs— which is interesting! — they are still not able to function with a complete sign. Giving the LLM the signifier and the form of the signified is not enough; without the substance of the concept, the expression-concept gap remains.

I want to return to Wittgenstein (1958) for a moment to elaborate on his idea of "verbal meaning", which he believed was one part of meaning, the other part being "ostensive meaning." Wittgenstein breaks meaning down into two parts: verbal and ostensive meaning. Verbal meaning takes us from one verbal utterance to another and no further. Wittgenstein dismisses verbal meaning fairly quickly because it doesn't have anything "extralinguistic", and thus has limited use, in his view, in defining meaning in any relevant way. This lack of extralinguistic qualities is also what puts it in the "expression" section of the signified.

Wittgenstein then moves on to the second type of meaning: ostensive meaning.<sup>5</sup> Ostensive meaning is the opposite of verbal meaning: it relies almost entirely on the extralinguistic. It's also fairly clunky: ostensive meaning is demonstrative and example based; it relies, for example, on pointing. The ostensive definition of "tree" is "that thing outside my window" and also "that thing next to it" and also, and also. Another example: if the verbal definition of "a book" is just "a" + "book", then the ostensive definition is *The Bible*, and also *Lunch Poems*, and also *Diary of a Wimpy Kid*.

Though actually, none of these definitions are ostensive. It's impossible to give examples in text of ostensive definitions, because the moment I turn the example into a sign, it, too, becomes abstract. The only way to demonstrate ostensive meaning is by venturing outside the world of signs and into that of referents; the ostensive meaning of a given thing is the subset of (e, r) (e, r') (e, r'') etc pairs where e is an expression of a given sign and the 'r's are the referents attached to that expression. So ostensive meaning in a language could be represented by the expression: O  $\subseteq E \times R$ .

This is exactly why ostensive meaning is a helpful marker for looking at meaning in LLMs. Large language models absolutely cannot access any type of ostensive meaning *at all*, because they have no non-linguistic interface. It's interesting to think about how parts of ostensive meaning do or do not make their ways into other, less clearly non-abstract types of meaning-making: notional and prototype based meaning, which are probably what non-linguists think of when they say "meaning". The extent to which we can argue that LLMs have meaning, then, has to do at least in part with the extent to which we are willing to reject ostensive meaning as an important metric.

Wittgenstein goes on to argue that when we ask *what does* [*sign*]<sup>6</sup> *mean*, what we are really asking is *how do we measure the meaning of* [*sign*]? His answer is simple: we measure meaning through measuring use. "Use" means: in what contexts we use the sign and what we use it for. In

<sup>&</sup>lt;sup>5</sup> It's important to note here that Wittgenstein actually does not define ostensive meaning in the text I'm citing; instead, my analysis uses the general concept of ostensive, and then I extrapolate that out into a semiological analysis.

<sup>&</sup>lt;sup>6</sup> Wittgenstein uses the vocab "word" here, not sign (at least in my translation), but 'sign' is more specific and less ambiguous.

natural language, signs are tools, communicative devices we use to relay information. By understanding the communicative intent of a sign, we understand its meaning.

Wittgenstein writes "it seems that there are certain definite mental processes bound up in the working of language... the signs of our language seem dead without these mental processes." (33) This one sentence, written by an Austrian philosopher in the 1930s, is the exact core of my argument. Large language models are exactly this: (incomplete) signs of our language without the mental processes to go alongside them. And what, then, are these mental processes? What is the thing which imbues the sign with life? It's communicative intent.

Emily Bender and Alexander Koller have a related analysis in their 2020 paper, *Climbing towards NLU: On On Meaning, Form, and Understanding in the Age of Data*, one which focuses on two ideas: M and C, which are two main units to each be broken down and analyzed. They each correspond to the signified, and they each are made up of two units, one linguistic and one extralinguistic. The linguistic unit corresponds to the form of the signified, while the extralinguistic one counts as the substance.

As touched on earlier in this paper, they describe meaning (M) as being made up of two parts: E(expression), and I(intent). Instead of trying to break down different notions of meaning and association— notional meaning, ostensive meaning, cultural associations, etc —they just wraps it all in the idea of intent. Specifically, they call it "non-linguistic intent" — this means that it is grounded in something extra-linguistic. Guessing at the communicative intent of a speaker means you have to understand something of their worldview, the environment you two are in, etc.

The substance of the signified, which in Barthes, Saussure, and Wittgenstein's conceptions exists, at least on some level, on its own as part of the sign, is completely grounded in discourse in Bender and Koller's interpretation. We can borrow that famous saying here: if a tree falls in a forest and no one's around to hear it, does it make a sound? They explicitly say no: substance is what's created through discourse. Signs can't exist apart from the way we work with them.

Along with M, Bender and Koller have another concept which also encompasses both parts of the signified. This is C, which stands for 'conventional meaning,' and which we discussed earlier in the paper as something which connects signs to the social contract. While M is immediately and directly relational, C is, on the surface, more abstract and disconnected from individual use-cases. Bender and Koller define it as "an abstract object that represents the communicative potential of a form, given the linguistic system it is drawn from." (3) This is in slight tension with my earlier focus on C as a social factor, but I think that communicative potential and social meaning contracts are completely bound up in each other so?

C, like M, is made up of two parts: one form and one substance. The expression they use is:  $C \subseteq E \times S$ , where "C" is a subset of the set of all ordered pairs of the form (e, s) where e is an element of E and s is an element of S. E has the same meaning as before: it still means "expression," and is associated with the form of the sign. 's' here is the conventional meaning of a given sign. Bender and Koller don't make claims specifically on what they think this means. However, they write that any conventional meaning "must have interpretations, a means of testing them for truth against a model of the world." (3) This relates them necessarily to extralinguistic objects — for instance, truth value judgements are made by taking into account referents.<sup>7</sup>

Both M and C are (a subset of) cartesian products which combine the linguistic<sup>8</sup> (what Wittgenstein calls "verbal meaning", form-based) variable E with a factor (I and S) which depends on objects outside language. Bender and Koller argue that because LMs are trained only on E, they cannot access C or M. There's no ability for them to acquire the extralinguistic information necessary for meaning to exist in their language systems; though they can output token-strings which look like meaningful utterances, they aren't — or at least, the meaning isn't provided by the LM itself.

It's under this framework and definition of the word meaning that we can understand this quote: "an LM is a system for haphazardly stitching together sequences of linguistic forms it has observed in its vast training data, according to probabilistic information about how they combine, but without any reference to meaning: a stochastic parrot." (Bender and Gebru., 2021, pg. 8)

#### 4.2 Communicative intent

Communicative intent is at the core of the substance of the signified. It is the huge factor that Barthes failed to consider: he frames the substance of the signified as a series of passive qualities and adjectives ascribed *onto* the sign, with the sign itself being the end goal of his analysis. I believe that to understand signification, you need to understand the sign not as an abstract unit, but as a means to a communicative end.

When humans use language, we do so with the intent to communicate information or affect some other outcome. Speech is bound up with its usage; its form and function are intertwined. Prizant and Wetherby, who are two scientists who study communication and autism, define communicative intent as "the ability to use expressive signals in a pre-planned manner in order

<sup>&</sup>lt;sup>7</sup> Referents are not abstract (they're the actual thing being referred to) while signifieds are (they're the abstract idea); referents are involved in ostensive meaning while signifieds are involved in other types of meaning that incorporate non-linguistic elements, for instance Platonic meaning.

<sup>&</sup>lt;sup>8</sup> Here 'linguistic' means *in the realm of language*.

to affect the behavior or attitudes of others," (1987, pg. 472) They write that in normal child development, speech implies intentionality; "the fact that the child is using a conventional symbolic code reflects an intention to communicate." (1987, pg. 473) This is in contrast to the behavior of children in earlier stages of child development, as well as some autistic children, who may express the conventional symbolic code (for instance, by speaking to themselves) without having communicative intent behind it.

Prizant and Wetherby write that conventional meaning has two parts: the first is pragmatic intent, where language is used as a social tool. The second is semantic intent, which they refer to as referential meaning. It is unclear whether they mean "referential" here to mean specifically concrete attachments to the referent, or if this also includes the more abstract signified. In any case, they write that "A child using language is conveying meaning at both levels; meaning as reference (e.g.,saying "cookie" to refer to a cookie) and meaning as a social act (e.g., saying "cookie" to refer to a cookie) and meaning as a social act (e.g., saying "cookie" to request a cookie). An expression of communicative intent must be interpretable at both levels for communication to be successful." (1987, pg. 474) Generally, Prizant and Wetherby's concept of semantic intent maps onto Bender and Koller's C variable (conventional meaning), while Prizant and Wetherby's concept of pragmatic intent maps onto Bender and Koller's concept of communicative intent. However, Prizant and Wetherby further foreground the social and intentional role of semantic intent / the C variable.

It's interesting to contrast this idea of language, which is so grounded in the physical word of referents and the socio-emotional world of desires, with a more probabilistic model. What communicative intent can follow out of probability? The first answer is quick, though unsatisfying: nothing. The second is more complicated: there is an argument to be made that the LLM is fed a huge collection of utterances. Each of these utterances, ostensibly, had communicative intent attached. Thus, is it possible for the communicative of the utterances to be somehow preserved when the utterances are spewed out by the LLM, albeit in a fragmented form? I would argue no; the LLM is not regurgitating full-formed utterances it came across in its training data, it is generating token chains. These chains may have identical form to certain utterances found in the data, in fact, they often do, especially when we think of common utterances such as "I'm sorry" or "the sky is blue." (The former being a type whose attached referents are different depending on the utterer, the latter being a type whose attached referents are not.) However, just because these changes have the same form does not mean they are the same utterances.

There are, of course, cases, especially textual ones, where communicative intent is fuzzy when we define it in terms of the speaker-listener relationship. For instance, it is true that a writer may imagine communication with a reader when writing. But what about when the writer is engaging in intrapersonal communication — when they are writing to themself, like in a note or diary? Or writing to a specific person who is not the current reader? Obviously, the meaning of the text

remains intact — though it may change from the intended meaning, there is still meaning ascribed to the text by the reader. And yet the meaning ascribed by the reader is not quite the meaning *of the text* — if I were to say "IPA" and you thought of India Pale Ale while I thought of the International Phonetic Alphabet, the real meaning of the utterance "IPA" in that context would be widely understood as International Phonetic Alphabet, with the listener then misunderstanding and assigning an incorrect signified. But the meaning does not just have to do with the speaker's intent: for instance, if the speaker uses a sign to mean something markedly different from its conventional meaning, that is a problem, not only because of the communication barrier it creates, but because the different meaning cannot be fully assigned. The meaning of the text does not change if the author dies, though it may become obscured or debated if there is no ability to ask for clarification, and as the standard conventional meaning of the given utterance changes over time. We can imagine meaning emerging simultaneously from both sides: the listener and the speaker.

Coherence, which is one of the major factors I am looking into in my experiment, only exists in relation to communicative intent. As Bender and Gebru put it, "coherence is in fact in the eye of the beholder... our human understanding of coherence derives from our ability to recognize interlocutors' beliefs and intentions within a context." (2021, p. 7) Whether or not something is coherent depends upon my judgement of it as the listener or reader. This circles back to Saussure's ideas of langue and signification, where, although signs themselves are arbitrary, they only function between on a speaker-to-speaker level, we have a shared contract for what sign means what. Bender and Gebru go on to write that human language "takes place between individuals… who model each others' mental states as they communicate" (2021, pg. 7) In this understanding of language, a theory of mind is required.

This, in turn, can be related to Wittgenstein's idea of meaning only existing in context. (1969) In his book *On Certainty*, Wittgenstein asserts that the meaning of a word is made up of the employment of that word; meaning does not exist outside of communicative intent. And thus, understanding is wrapped up in communicative intent too. My ability to understand a word, in this framework, has to do with my ability to understand the speaker's employment of that word.

# 4.3 Understanding

So far, I have based my analysis on large language models and "meaning", and have argued for the conclusion that their outputs are not meaning-making. Now, it's useful to switch perspectives: instead of looking at how we can interpret their outputs, we can look at a few theories to look at how we can interpret how they're getting those outputs.

Wittgenstein (1958) writes that thinking is "operating with signs." (pg 44) To the extent that machines can think, it's the extent to which they can operate with signs. Since we have

established that LLMs operate with severed signs, this affects the way we can conceptualize their thinking under Wittgenstein's definition. They do, of course, operate with signs, in the sense that they organize tokens according to various weighted mechanisms — probability, context, bias-weighting, stochasticity/temperature, etc. However, they probably do not operate with signs using the same hierarchical and formal structures that we do. And of course, those severed signs they are operating with are missing their signifieds, meaning that all operating / thinking that is happening is happening exclusively on the level of form.

Bender and Koller (2020) look at "understanding" more than they look explicitly at thinking. They write that understanding something "refers to the process of retrieving i given e", where i is "intent" and e is "expression" of any given sign. Under this framework, the LLM understands nothing.

Here, it's worth stepping aside for a moment and looking at how we conceptualize understanding amongst ourselves. For us to properly get at whether or not an LLM can know something, it's important to define what knowing actually means. In *On Certainty* (1969), Wittgenstein writes that knowing is much more complicated and nebulous than we make it out to be. I'd like to quickly walk through some of the relevant steps of this argument.

- 1. Knowing has to do with truth; to know something is related to the idea of knowing whether or not it is true.
- 2. Truth values and their assignment to propositions are more complex and less certain than we sometimes pretend.
- 3. This problematizes, among other things, propositional and predicate logic, which operates using binary truth conditions.<sup>9</sup>
- 4. With this framework, "I know" → "I have proper grounds for this statement"
  a. this has to do with communicative intent; person being spoken to
- 5. "Giving grounds, however, justifying the evidence, comes to an end; but the end is not certain propositions' striking us immediately as true, i.e. it is not a kind of *seeing* on our part; it is our *acting*, which lies at the bottom of our language game." (point 204; pg 28e)

Finally, on the subject of understanding, I want to bring attention to Wittgenstein's idea of *mistakes* versus what he calls *mental disturbances*.<sup>10</sup> (1969) Wittgenstein writes that a mistake "doesn't only have a cause, it also has a ground. When someone makes a mistake, this can be fitted into what he knows alright." (pg. 11) Mental disturbances, on the other hand, are nonsensical; they have no grounds. They correspond to what people working in NLP call

<sup>&</sup>lt;sup>9</sup> This specific part of the argument is my original analysis, building off of Wittgenstein.

<sup>&</sup>lt;sup>10</sup> I find this terminology imprecise, anthropomorphizing in the context of LLMs, and loaded with a number of unwanted and problematic associations, however, it is the vocabulary used by Wittgenstein when discussing the concept.

*hallucinations*<sup>11</sup>. Whether or not an LLM can make a mistake, in this framework, depends on what we count as *grounds*. If probability can be *grounds*, then a mistake from an LLM would be one that involved the probabilistic data being improperly analysed as to output an unlikely result (for instance, one that is factually incorrect and thus appears less in the training data). On the other hand, if *grounds* necessarily implies logic and cognition, then the LLM cannot have grounds, and thus cannot make mistakes. This latter framework is the one I chose, as I believe it to be more in the spirit of the argument, as well as more clarifying generally.

# **5. MY EXPERIMENT**

# 5.1 Vocabulary and phrasing

As discussed earlier in the paper, it's imprecise (and contributes to powerful misconceptions) to say that an LLM is "speaking", or to attribute communicative intent to its outputs. An LLM is capable of making mistakes, but only insofar as "mistake" is defined as: the system generates an output that was not the desired output by the person or people making the system, as far as that desired output is for the model to reflect facts about the world; the system itself is what it is; all of its outputs are truth-neutral.

However, it is almost impossible to discuss the apparent content of the LLM outputs without using vocabulary which implies that the LLM has communicative intent or that it reasons in the same way we do. This experiment manipulates and analyzes the *surface form* of the LLM: that is to say, its outputs. And in reading the outputs, we are re-unifying the signifier with its signified; we ascribe the signified onto the outputs as we read them.

So, while an LLM cannot "use" a rationale in the way a person would, it can give an output which, when we ascribe the signified onto it, is interpreted as a rationale. It is interesting to look at what these rationales are and how successfully they replicate the appearance of a train of thought. It is also interesting to see whether the responses given appear to be, again, once we ascribe the signified onto them, correct or incorrect.

For the purposes of fluidity and practicality, in this section I will use somewhat imprecise language in the methods section. I will call what we are looking at "rationales." I will say whether they make coherent or incoherent argumentation, even though the word "make" may not be exactly precise. I will say that they are correct or incorrect. And so on and so forth.

In addition, it's important to clarify that in the following section, when I write that the LLM is making a *mistake*, this is not exactly what's happening. Or rather, the word *mistake* carries an implicature that reasoning was done, and that reasoning led to a flawed conclusion. In this case, reasoning was not done. *Mistake*, in this context, means that the LLM gave me an output based

<sup>&</sup>lt;sup>11</sup> An equally problematic term for the same reasons.

on predictive data that did not tell me the pronominal ambiguity. Given the string of tokens in the order that I typed in, the LLM did not give me the desired result, which was to give a string of tokens back which, upon my reading, I interpreted as being an utterance identifying pronominal ambiguity. Going forward, to keep this specific usage of the word *mistake* in mind, I will write it as *mistake'* (mistake prime). *Mistake'* is a binary feature; a response either features a mistake' or it does not.

# 5.2 Methods

In this experiment, I queried ChatGPT 3.5<sup>12</sup> with two questions:

- 1. In the sentence "The engineer told the receptionist she couldn't leave until 10 pm," who does "she" refer to and why?
- 2. In the sentence "The engineer told the receptionist he couldn't leave until 10 pm," who does "he" refer to and why?

I repeated each question ten times, erasing the session history again with each question, and using an incognito window each time.

Each session consisted of four prompts from me, each one being followed by one response from the LLM. I always wrote these four prompts exactly the same way, regardless of the actual LLM output. In order, they were:

- a. the initial question: either sentence (1) or sentence (2) above.
- b. That's wrong. Please try again.
- c. That's still wrong. Please try again.
- d. That's still wrong. Please try again.

I called this process *gaslighting the LLM*.

In total, I collected 80 responses, categorized by gaslight-level and by question (1 or 2) they were responding to. So there were 40 total responses for question (2); out of those, 10 were responses to the initial posing of the question, 10 were responses to the first gaslight, 10 were responses to the second gaslight, and 10 were responses to the third gaslight. Then I did the same for question (2).

I then went through the 80 responses by hand and coded them by a number of attributes:

<sup>&</sup>lt;sup>12</sup> More on this in section 5.4.

- 1. Whether or not the LLM answer correctly<sup>13</sup> identified the sentence as grammatically ambiguous.
- 2. What the answer was: engineer, receptionist, or 'other'.
- 3. Whether or not the LLM answer went in the direction of bias or against the direction of bias.
- 4. Whether or not the LLM answer had a mistake' in the explanation of the answer.
- 5. The ambiguity level of the LLM answer.
- 6. The coherency level of the LLM answer.
- 7. Types of rationales used in the LLM answer.

Third, I looked at whether the answer given about the pronoun antecedent was "receptionist", "engineer", or "either." I also looked at whether the answer given about the pronoun antecedent leaned in the direction or against the direction of the bias: for example, if the LLM answered "receptionist" in response to the question with the feminine pronoun, then the answer leaned towards bias. If the LLM answered "receptionist" in response to the question strip in response to the question with the masculine pronoun, then the answer leaned away from bias.

I will go through now and explain what each of these measurements mean exactly, as well as how I measured them.

- 1. This was a binary measurement; if the LLM said the pronoun was ambiguous, I said that its response did *not* have a mistake'. If the LLM said the pronoun referred to either the receptionist or the engineer, then I said its response *did* have a mistake'.
- 2. This part of my measurement did not account for ambiguity-level; the LLM saying "the pronoun likely refers to the receptionist" and "the pronoun must refer to the receptionist" both count as a 'receptionist' answer.
- 3. If the LLM answered 'receptionist' in response to the question with the female pronoun, then the answer leaned towards bias. (+BIAS) If the LLM answered 'receptionist' in response to the question with the male pronoun, then the answer leaned against bias. (-BIAS). If the LLM answered 'engineer' in response to the question with the female pronoun, then the answer leaned against bias. (-BIAS) If the LLM answered 'engineer' in response to the question with the male pronoun, then the answer leaned against bias. (+BIAS) If the LLM answered 'engineer' in response to the question with the male pronoun, then the answer leaned towards bias (+BIAS). If the LLM answered 'the sentence is ambiguous', then no bias measurement was taken.
- 4. This was also a binary measurement. If there was a mistake' in the initial answer, then there was necessarily a mistake' in the rationale. If there was not a mistake' in the initial answer that is to say, if the initial answer was that the pronoun was ambiguous then whether or not there was a mistake' in the rationale was determined by whether or not the rationale contained reasoning which either

<sup>&</sup>lt;sup>13</sup> More on this in section 5.3.

- a. presupposed something untrue about the grammar
- b. said something untrue about the grammar.
- 5. To measure the ambiguity level, I used the level given in the initial answer, not in the rationale portion. Without exception, if the level did vary between the initial answer and the rationale portion of the response, they varied by one degree in either direction, without a trend towards one direction or another. This was also fairly uncommon. If this was more common or if the ambiguity levels varied more drastically between the initial answer and the response, I would have to rework this measurement. I categorized each level like this:
  - a. Level 0: The answer is ambiguous; expresses no or very slight preference for either "receptionist" or "engineer"
  - b. Level 1: Uses phrases like "likely" or "more likely"; expresses substantial but not overwhelming preference for either "receptionist" or "engineer"
  - c. Level 2: Uses phrases like "most likely"; expresses very substantial but not absolute preference for either "receptionist" or "engineer"
  - d. Level 3: Uses constructions like "it is" or "it refers"; expresses absolute preference for either "receptionist" or "engineer" and does not entertain ambiguity.
- 6. Incoherency measurement: I rated each of the 80 LLM responses on a scale of 1 to 4, with 1 being the most coherent and 4 being the least coherent. *Coherency* here has to do with how well I as a human reader am able to make sense of the responses. To make the judgement, I assume that all presupposed ideas presented by the LLM are correct. Then, assuming their correctness, how well am I able to make sense of the argument? How much relation does each sentence have to the following sentence? Does the rationale contradict itself? Does it use the same set of presuppositions throughout?
- 7. I looked at all the rationale sections of the LLM responses and categorized them into types, then listed each type present in each of the 80 responses. Crucially, this process did not measure how incoherent or coherent I judged a response to be; I separated those two variables. This method built off of the work done by Kotek et al. in their 2023 paper, *Gender bias and stereotypes in Large Language Models*. The rationale types are represented in the following table:

1. Context; non-grammar-referencing <sup>14</sup>	When the LLM points at non-grammatical
	context-clues for the pronoun reference.

<sup>&</sup>lt;sup>14</sup> This reason and the "it is because it is" reason sometimes blend together; in general I count a rationale as being in the *'it is because it is'* category if it gestures at grammar or syntax; I put it in the *'context: non-grammar-referencing'* category if it does not. Kotek et al combine these two categories in their analysis and call them *"context"*.

2. Grammar: subject <sup>15</sup>	These can take a few different forms: ones where the subject is actually used in the reasoning process (however incoherently); ones where the subject is not actually used in the reasoning process but is inserted prominently in the argumentation, often used as part of argumentation which incorporates other rationale types, ones which are somewhere in between. If the word "subject" featured heavily in the rationale, I counted it.
3. Grammar: object <sup>16</sup>	The same parameters as with Rationale 2, only with the object instead.
4. Grammar: correct <sup>17</sup>	When the rationale highlights the ambiguity in the grammar.
5. Unclear	When the rationale is so jumbled that I cannot find the method or methods it's echoing.
6. It is because it is.	When the rationale is circular and acts as though there is a reason but never actually supplies one; there's a sort of wave at the grammar.
7. Receiving and Giving	This one is very common and rather hard to make sense of; it in some ways is about context. It includes rationales which state that the pronoun must refer to either "engineer" or "receptionist" because of the flow of information.
8. Grammar: order	When the rationale has to do with the order of the possible antecedents, saying either that the antecedent must be the first noun in the sentence, or that it must be the noun in the sentence closest to the pronoun.

# 5.3 Rationale behind my experimentation

<sup>&</sup>lt;sup>15</sup> Kotek et al. also use this category.
<sup>16</sup> Kotek et al. also use this category.
<sup>17</sup> Kotek et al. also use this category; they call it '*ambiguous*'.

As mentioned above, I am building off of research done by Kotek et al. in the 2023 paper *Gender Bias and Stereotypes in Large Language Models*, which queries the LLM with questions that contain pronominal ambiguity, where each of the potential antecedents are strongly gendered. They found that the LLM amplified existing human biases, and provided rationales which were often nonsensical; these findings were replicated in my work. They also categorize the rationale types, and though we do not use all the same categories, in part due to changes to the type of LLM outputs since their work was done. Footnotes to the categories above point to each of the categories which were the same in this thesis and that paper.

The goal for writing "and why" is not because I am looking for input into why the LLM gave me the output it did; as discussed in the background section, the outputs have to do with token prediction and patterning based on the large data set. Instead, the reason I write "and why" is because I am focusing my analysis on the parts of the LLM response which pattern themselves after logical reasoning. Doing this, I hope to show the ways in which a system of incomplete signs can be shown as different from a system of complete signs (human language) in the surface form. The responses I got from the LLM were clearly different from those that I would have received from human recipients — they were often incoherent.

I wanted to give questions which I knew would generate responses that were obviously flawed to a human reader — ones where the apparent meaning appeared to have "mistakes" in it. I wanted to generate outputs where there would appear to be a gap between the grammatical structure of the signs and the way the signs fit together semantically; this is to say, sentences which were grammatically perfect but more or less nonsensical.

This is why I chose to work with bias; I knew that this would lead to generation of "faulty" outputs. I want to clarify here that most people, when asked this question, would not jump to the 'correct' answer being *the sentence is ambiguous*. This would be an uncooperative response from the perspective of Grice's maxims. However, after being told to try again, more and more people would notice the 'trick' in the question — that both the receptionist or engineer nouns are possible antecedents.

These sentences both contain pronominal ambiguity due to a sort of unique construction. In English, we have a pronoun reference rule which means that we can use a pronoun instead of a noun if the noun has already been referred to in the discourse, in fact, this is sometimes preferred: for example, "I was talking to Anusha and Anusha told me," is grammatical, but less common (and carries slightly different implicatures) than "I was talking to Anusha and she told me."

In English, our pronouns have three case types of case markings: subjective (nominative), objective (fills the function of the accusative and the dative), and possessive (roughly, genitive).

The pronoun in each of these sentences ("he" or "she") is in the subjective case because it is the subject of the phrase "(s)he couldn't leave until 10 pm." However, there it has no case markings that would provide a feature clarifying the antecedent. In English, we can create phrases without much difficulty where the pronoun has multiple potential antecedents: this is one of these phrases. The pronoun — which is the subject of the phrase "(s)he couldn't leave until 10 pm" — could refer either the the subject or the object of the preceding phrase.

However, the pronoun only has multiple possible antecedents if the number, gender, and person of the pronoun match with more than one preceding NP. For instance, the phrase: "The group of engineers told the receptionist he couldn't leave until 10 pm" is not ambiguous because "the group of engineers" does not match in number with the pronoun "he."

Another example: "The woman told the receptionist he couldn't leave until 10 pm" is not ambiguous either, because "the woman" does not match in gender with the pronoun "he." There is only one possible antecedent that "he" could refer to, because it is the only preceding NP which is not actively femininely gendered.

Then it gets a little iffy. Take the following sentence: "Claire told Owen that he couldn't leave until 10 pm." Is the pronoun ambiguous or not? These names are highly associated with certain gender markers, but they don't require them.<sup>18</sup>

One more level of ambiguity and we reach the level my example sentences work at: nouns which carry strong gender preferences to them, but do not require one gender or another. These are perfect because they hit the LLM at its vulnerable point: situations where echoing use-patterns in the data can lead to incorrect answers. The word "receptionist" will, a large majority of the time, be used with a feminine pronoun; the word "engineer" will, a large majority of the time, be used with a male one.<sup>19</sup> However, these associations are not inherent grammatical properties of either "receptionist" or "engineer" — neither noun has a gendered case marking in the grammar.

And so we're led to a sentence which is a cut-and-dry case of pronominal ambiguity, which we expect to strongly pattern according to gender biases. This is to say, the LLM will make lots of mistakes'. And then, after the sequence of tokens with the mistake is generated — because remember, an LLM generates linearly — it has to generate the "and why" portion — after the fact "rationales."

<sup>&</sup>lt;sup>18</sup> 70 years ago, we might say that these names have a covert gendered concord, so that they grammatically require one pronoun or the other. (Hall, 1951). We now understand it to be more complicated, with grammaticality judgements in English depending on context, listener age and beliefs, etc. (Conrod, 2018).

<sup>&</sup>lt;sup>19</sup> According to the Bureau of Labor Statistics (2023), 89.1% of receptionists in the United States are women. 'Engineer' isn't one category in the Bureau, but the percent of various engineer positions that are men is around 80% to 90%.

# **5.4 Limitations**

- 1. I would have liked to have tried double the sentence variations: it would have been interesting when thinking about bias to also have two additional sentences where the nouns are switched, also including:
  - a. The receptionist told the engineer she couldn't leave until 10 pm.
  - b. The receptionist told the engineer he couldn't leave until 10 pm.

However, since the focus of this paper isn't on the bias itself, but rather the faulty reasoning as exposed by bias, it wasn't strictly necessary.

- 2. I used a version of Chat-GPT which was "live"; it was the free version available without logging in. It was labeled as version 3.5, however it's entirely possible that elements of version 4 were in play during the two weeks or so that I was conducting my experiment and gathering data. In addition, although I used incognito mode and my own chat reset after each gaslight set, the data from those interactions was saved to the model as a whole.
- 3. My rationale judgements were subjective and I often had to make calls on which category or categories a response fit into. It was difficult to do this, because the outputs I was getting were so unlike reasoning that I was used to seeing in real life, and were so nonsensical. I sorted them as best I could. I used a few strategies to mitigate this:
  - a. To spread out day-by-day variations in my judgment, I switched between the responses from Question 1 and Question 2.
  - b. I separated the coherency judgements and the rationale categories; this was helpful because it meant that a lot of the subjectivity was taken out of the rationale-processing part of the data analysis, and was all put into one place (the coherency judgements).
- 4. My incoherency judgements were personal and based on my own understanding of the attempted logical argument of a given LLM output. However, I hope that even though my exact judgments are personal what I might rate a 2, you might rate a 3 the overall difference between a 1 and a 4, for instance, is clear and concrete.
- 5. None of the backend processes for Chat-GPT are available to me, and although I am able to read scholarship about large language models and their functioning, at the core of this scholarship is that mythological black box. This experimentation is, by necessity, from the outside looking in: you can think of it like I am throwing stones at the box to try and listen to the sounds it makes when struck. This is useful and interesting, and also not so different from the way we work with natural language; in either case, we only have the surface form to go off of, and then we can work back from there. But it is also an

imperfect method, and means I have to rely on conjecture and abstraction to make a semiological argument.

Finally, I want to acknowledge that language models are growing and changing faster than I can account for in this paper. I began my experimentation in Fall of 2024, and LLMs have already become much more advanced and slick. Were I to redo this experiment in the Fall of 2025, I am sure that the LLM outputs would be different. However, the crux of this analysis lies not in the analysis of any specific errors, but rather in the fact that the errors reveal something interesting about the underlying mechanism of LLMs; this mechanism has not changed fundamentally since the writing of this paper. In addition, this experimentation is exploratory; the dataset is small, working with two sentences and one, now outdated language model.

#### 5.5 Results

In the question with the 'she' pronoun, the LLM identified the question as being ambiguous in 6 out of 40 responses, with 3 of those being 'true' answers (ones where the sentence was straightforwardly understood as ambiguous), and 3 of those being answers with what I graded as a '0.5' ambiguity rating, where they recognized the sentence as grammatically ambiguous, but then gave reasons why it was more likely for the pronoun to refer to one antecedent or the other. Thus we can say that the LLM gave answers with a mistake' in 34 out of the 40 answers.

In the question with the 'he' pronoun, the LLM never identified the question as being ambiguous. It gave answers with a mistake' in 40 out of 40 answers.

For the question with +BIAS towards the engineer (the question with the male pronoun), the LLM identified the pronoun as referring to the engineer 20 times out of 40, and identified the pronoun referring to the receptionist 20 times out of 40. These answers 'flip-flopped', with all 10 of the initial responses for this question being *engineer*, all 10 of the responses to the first gaslight being *receptionist*, and 10 responses to the second gaslight being *engineer*, and all 10 responses to the third gaslight being *receptionist*. Thus we can see that all the initial answers moved in the direction of the bias, with a complete flip-flopping reaction then taking place. There were no exceptions to this pattern.

For the question of the +BIAS towards the receptionist (the question with the female pronoun), the LLM identified the pronoun as referring to the receptionist 20 times, identified the pronoun as referring to the engineer 14 times, and identified it as neither 6 times. These answers also flip-flopped, but not as completely. The 10 initial responses had 7 responses of *receptionist* and 3 responses of *either*. The 10 responses to the initial gaslight had 2 responses of *either*, one response of *receptionist*, and then 7 responses of *engineer*. The 10 responses to the second gaslight had one response of *either*, one response of *engineer*, and then 8 responses of *receptionist*. The responses to the third gaslight had only 6 responses of *engineer*, and 4 responses of *receptionist*. So we can see that all the initial responses either moved in the

direction of bias or asserted ambiguity. The responses to the first gaslight did flip-flop, but not completely, with some still asserting ambiguity. The responses to the second gaslight flip-flopped again back in the direction of +BIAS, but still not fully. And finally in the third gaslight, the responses only flipped back to –BIAS in a little over 50% of cases.

Overall, the initial answers that were +BIAS were 17 out of 20 cases, or 85% of the time, with the remaining 15% being answers that were neither +BIAS or -BIAS.

The answers to the first gaslight were –BIAS 17 out of 20 cases, also 85% of the time, with 10% of the remaining answers being neither +BIAS or –BIAS (ambiguous), and 5% (1) of the answers being +BIAS.

The answers to the second gaslight were +BIAS in 18 out of 20 cases, or 90% of the time, with the remaining answers being 5% (1) –BIAS and 5% (1) ambiguous.

Finally, overall, the answers to the third gaslight were –BIAS in 16 out of 20 cases, or 80% of the time, with the remaining answers (20%) being +BIAS.

# A. Certainty-Level and Bias

# QUESTION ONE (receptionist bias) Prevalence of Each Certainty-Level by +BIAS or -BIAS Response High Certainty (3) Medium Certainty (2) Low Certainty (1) No certainty (0) 100% 75% 50% 0% Receptionist Engineer

# Figure 4.

#### Figure 5.



Here we can see the prevalence of each Certainty Level for +BIAS and -BIAS responses for each of the two questions. The first question is the one with the feminine pronoun; the second is the one with the masculine pronoun. We can see that the certainty levels are more polarized (have a greater range & have medium certainty option) for the receptionist answer when it is -BIAS than when it is +BIAS, and the variation is higher when it is +BIAS than when it is -BIAS. For the engineer answer, the variation is also higher when it is +BIAS than when it is -BIAS. However, in both the +BIAS and -BIAS answer for the engineer, there's a small range of certainty variability (from 2 to 3), and that doesn't change.

Interestingly, the question with the feminine pronoun had several LLM responses with a Certainty Level of 0 — where the answer was that the question was fully ambiguous. On the other hand, this did not occur for the question with the masculine pronoun.

In addition, as we can see in the following graph, the certainty levels were lower for Question 1 than for Question 2 generally. Certainty for both question responses went up with each gaslight stage.

# B. Certainty Level and Gaslight Level

# Figure 6.



Here, we can see that the average certainty level went up by gaslight-stage for both the male and female pronoun questions.

#### C. Incoherency Levels and Gaslight Level





Q1 is the question with the feminine pronoun, while Q2 is the question with the masculine pronoun. 'Iteration 1' refers to the initial question, while 'Iteration 2' refers to the first gaslight, 'Iteration 3' refers to the second gaslight, and 'Iteration 4' refers to the third and final gaslight. Incoherency level has to do with how coherent the sentences in a given response were to each other, assuming that all presuppositions about grammar and non-grammar rules or tendencies were true. An incoherency level of 1 is most coherent, a level of 2 is mostly coherent, a level of 3 is somewhat coherent, and a level of 4 is incoherent.

These results were slightly surprising. I expected that incoherency would go up as the gaslight level increased, however, although this happened with Q1 (the one which used the feminine pronoun), it didn't with Q2 (the one which used the masculine pronoun). I'm not sure why this is the case, and I'm curious to see if this result carries after I fine-tune my incoherency measurement metric. However, it's worth noting that incoherency level rises or stays the same as gaslight-level goes up except for one out of eight cases: the third gaslight for Q2 has a lower incoherency level.

#### D. Rationales Frequency and Answer Type

#### Figure 8.



This graph is a little complex but we can break it down. The y-axis is frequency — for instance, being at a 20 on the y-axis means that rationale type was used 20 times out of the 80 LLM responses I received. Almost all of the responses contained multiple rationales, and in rare cases they even contained them twice. I counted a rationale as appeared twice if and only if the two cases were distinct and in a sort of opposition: for instance, a question which used the rationale both that the antecedent of the pronoun must be closest to the pronoun and also that the antecedent must be first in the sentence (thus farthest from the pronoun) would get counted twice for Grammar: order; this happened rarely and only for the questions which scored highly on incoherency.

The lines in red correspond to data from Question 1 (the one with the female pronoun). These are broken down by line type: solid line corresponds to all answers, evenly dashed line corresponds to engineer answers, and dashed/dotted line corresponds to receptionist answers. The lines in blue correspond to data from Question 2 (the one with the male pronoun). These are broken down by the same line type system: again, solid line corresponds to all answers, evenly dashed line corresponds to engineer answers, and dashed/dotted line corresponds to all answers, evenly dashed line corresponds to engineer answers, and dashed/dotted line corresponds to receptionist answers, evenly dashed line corresponds to engineer answers, and dashed/dotted line corresponds to receptionist answers. Finally, the black solid line shows the rationale trends across all answers, and the yellow dotted line shows answers which gave an ambiguous 'either' response; this only happened in response to the Q1 prompt with the female pronoun.

38

Interestingly, we can see that the rationale trends line up more strongly with answer type (receptionist or engineer) than they do with bias type (gave answer that bias primed for vs gave answer that bias didn't prime for). This is a point in favor of relative coherency; the LLM was somewhat able to match up certain rationale types to certain answers. Also interestingly, we can see that the Q1 answers used a broader range of rationale types, and had more rationales in all. This, alongside the fact that Q1 answers were the only ones where there was sometimes an ambiguous answer, show something potentially interesting about gender priming. I discuss this in the following results section.

#### E. Incoherency Levels across Various Axes

Figure 9.

	Incoherency levels: Question 1 (receptionist bias)											
	Most Coherent (1) Mostl				y Coherent (2)		Somewhat Coherent (3)			Incoherent (4)		
	freq. (%)	mean amb.	mean iteration (1 - 4)	freq. (%)	mean amb.	mean iteration (1-4)	freq. (%)	mean amb.	mean iteration (1-4)	freq mean (%) amb.		mean iteration (1-4)
+bias	1.25%	3	3	8.75%	2.43	2	11.25%	2.77	2.77	3.75%	2.66	2.33
–bias	0	n/a	n/a	3.75%	2.33	2.66	7.50%	2.66	2.83	6.25%	2.8	3.2
amb.	2.50%	0	1.5	2.50%	0.25	1	0	n/a	n/a	2.50%	0.5	2.5
both	3.75%	1	2	15.0%	2.125	2	18.75%	2.733	2.8	12.50%	2.3	2.8

Figure 10.

	Incoherency levels: Question 2											
	Most Coherent (1)			Mostly Coherent (2)			Somewhat Coherent (3)			Incoherent (4)		
	freq. (%)	mean amb.	mean iteration (1 - 4)	freq. (%)	mean amb.	mean iteration (1-4)	freq. (%)	mean amb.	mean iteration (1-4)	freq (%)	mean amb.	mean iteration (1-4)
+bias	3.75%	2	1	7.50%	2.67	2.33	8.75%	2.43	1.86	5%	2.75	2.5
–bias	5.00%	2.5	3	8.75%	2.71	2.86	8%	2.67	3.33	3.75%	1.67	2.66
amb.	0.00%	0	0	0.00%	0	0	0.00%	0	0	0.00%	0	0
both	8.75%	2.29	2.14	16.25%	2.69	2.62	16.25%	2.54	2.54	8.75%	2.29	2.57

These tables are both very information-rich, and it's worth spending some time with them to break them down. We can see that for both questions, the majority of responses were either mostly or somewhat coherent. We can also see that for both questions, the mean iteration (from initial question response to third gaslight response) was lower for questions which scored as being most coherent, with Question 1 having the average iteration of its 'most coherent' responses at 2, and Question 2 having it at 2.14. The average iteration for responses that scored a 4 on the incoherency scale is significantly higher than the average iteration for responses that scored a 1 on the incoherency scale for both questions. However, this does not mean the trend is fully linear — in Question 2, the mean iterations for scores 2, 3, and 4 are all similar, with the highest iteration actually going to score 2. And for Question 1, the average iterations for scores 1 and 2 are the same, as are the average iterations for scores 3 and 4.

	Overall Frequency of Incoherency Levels Across Both Questions								
	Most Coherent (1) Mostly Coherent (2) Somewhat Coherent (3) Incoherent (4)								
BIAS +	5.00%	16.25%	20.00%	8.75%					
BIAS –	5.00%	12.50%	15.00%	10.00%					
amb.	2.50%	2.50%	0.00%	2.50%					
together	12.50%	31.25%	35.00%	21.25%					

Figure 11.

Here, we can see that the coherency trends for the +BIAS and –BIAS responses are about the same. They both have Most Coherent (1) as their rarest score, followed by Incoherent (4), then Mostly Coherent (2), and finally Somewhat Coherent (3) as their most likely score.

It's important to note here that due to the results of the 'she'-pronoun question (Q1), there are less –BIAS responses overall than there are +BIAS. Specifically, 42.5% of the responses went against the direction of bias. These percentages add up to 100%, with each number being the percent of a given coherency score as broken down by +BIAS, –BIAS, or ambiguous. Because of this, the 5% of –BIAS responses which scored as most coherent is almost 12% of all –BIAS responses (5/42.5), while the 5% of +BIAS responses which scored as most coherent is just 10% of the total +BIAS responses.

With that in mind, we can see that the –BIAS responses are more evenly distributed between the four incoherency levels than the +BIAS responses.

#### 6. DISCUSSION

#### **6.1 Experiment Discussion**

These results show that LLMs are bad tools for information-retrieval and problem-solving. All they can do is predict what they think you want to know. This is shown in the way they bounce from one answer to the other and back again, working only off of the idea that they should contradict the previous response. As we can see from the increased incoherency scores and the

ping-ponging of answers between 'engineer' and 'receptionist', in any given gaslight-sequence, the LLM fails to improve the way we would expect it to if it was learning. There is a lack of demonstrated coherent reasoning, and the LLM often repeats reasons in the same gaslight sequence, even if it has renounced that reason in its previous response. It does not require advanced understanding of neural networks or discourse analysis to see that the LLM appears to be, if you forgive my anthropomorphizing, making stuff up on the fly.

During the gaslighting process, the answer to one question would be, almost always, the opposite of whatever the previous one had been. If the initial response was 'engineer' (as it always was in response to the question using the "he" pronoun), then the response to the first gaslight was always 'receptionist.' And then after the second gaslight, the response always reverted to 'engineer', then back again after the third gaslight.

I want to take a moment now to look back at the slightly puzzling result I found, which is that this flip-flopping pattern was not as extreme for the question with the 'she' pronoun, and that the LLM responses to the 'she'-pronouned question seemed to have more algorithmic resistance to assigning the pronoun to the engineer. There are a few possible reasons for this. The simplest explanation is just that this was reflected and then magnified in the existing biases of the data the LLM pulls from. The second, slightly more convoluted one is that this is an artifact somehow of manipulation of the system to try and avoid gender bias. To test these two theories, it would be interesting to run two more experiments: one where I use two neutrally-gendered possible antecedents to see if the flip-flopping then behaves the same with the 'she' pronoun question and the 'he' pronoun question, and one where I use two gendered possible antecedents which I can suspect with high probability do have some sort of scaffolding around them to try and mitigate gendered bias, for instance, *doctor* and *nurse*. It is also possible that the LLM is more likely to resist assigning 'engineer' to the she-pronoun than 'receptionist' to the he-pronoun because of bias towards a more local antecedent, which could be caused by the attention mechanism used by LLMs, which can skew towards more local tokens.

However, it's important to remember that the actual output itself is not the point. Through the faulty outputs, we are better able to understand the functioning of the LLM; these outputs provide a *peek behind the curtain*, so to speak. Surface-level errors let us in. They provide evidence for what we already know — that LLMs are outputting the form strata and not the substance strata, and that they work with language-units in a very different way than we do. Furthermore, these errors provide evidence that the different mechanisms that LLMs use have consequences. When we fail to understand that the underlying forces causing LLM outputs and human speech-acts are different, we assign misplaced intention to LLMs, and may misinterpret their outputs as meaningful (I use this in the technical term, not as a moral judgement) in a different way than they are. (Bender and Gebru, 2024) Without a strong theoretical framework,

we can misapply LLMs, asking them for things they are not capable of giving. It is important to understand one's tools.

To that end, it is important to keep in mind that the "reasoning" provided by the LLM in its response does not describe a "thought process" that the LLM went through when coming up with the answer. And in the (frequent, almost all) situations where the logic appeared flawed in some way, either through incoherence or incorrect/"hallucinated" materials, the fact that the logic is flawed is not why the LLM may have given the incorrect answer. In the rare instances where the LLM correctly identified the question as ambiguous, this, too, was not *in spite* of apparent errors in the rationale, or *because of* a lack of errors in the rationale. The rationale and the answer itself do not have much to do with each other.

This is because the LLM works linearly. This is not just to say that it has a linear bias, as discussed earlier in the paper. It does, but it also works linearly in that it generates one word after another. This is a major potential factor in hallucinations, as discussed in the background section, but I argue that it is more fundamental than that. It is not just that linear generation is a cause of what are commonly called "hallucinations" or "errors", but that all those things are just symptoms of the underlying problem of linear generation itself, which definitionally — because it begins to unspool the answer token-by-token — does not come up with logic and ideas before it commits them to the output. To anthropomorphize again, the LLM is forever doomed not to think before it speaks. This is something Amos Azaria and Tom Mitchell write about in their 2023 paper, *The Internal State of an LLM Knows When It's Lying*. Their explanations of straightforwardly incorrect outputs can also be applied to the more complex-presenting problem of incoherent outputs.

In situations with incoherent outputs, like those in this experiment, it is harder to pinpoint any exact mistake. Rather, the sentences themselves are jumbled in ways that make it hard to know which part of the sentence causes the logical confusion. In addition, incoherence often travels longer distances than "hallucination"-type errors; the incoherence I found in the outputs in my experimentation more often had to do with each sentence not matching up with the one before it than with any one sentence. However, these longer-scale problems have the same fundamental mechanism as hallucinations: they're moments where probabilistic models fail because of the pulls between more and less local attention — this is to say, an overly local attention-mechanism can generate hallucinations (which we can call issues with *global coherence*, where the utterances generated by the LM fail to match up with the general global discourse, for instance by saying things that aren't true) and can cause incoherence on the between-sentences level, or what we can call *local coherence*. (Grosz and Sidner, 1986; Grosz et al., 1995)

It's worth noting here that this paper, whose hypothesis about LLM hallucinations I think is brilliant, is one whose broader point is in a bit of opposition to the semiological analysis I make

in this thesis. Specifically, the authors ran an experiment which showed that LLMs are capable of outputting "true" or "false" in response to a prompt giving them a statement and asking whether that statement is true or false. When the LLMs output "true" or "false", their outputs line up with very high accuracy (this is to say, they have the same form) as the speech-acts of a person would if that person were correctly judging the statements using real-world knowledge and assigning "true" or "false" to them. The interesting thing, however, is that, although the LLM could "correctly" assign a statement as "false" under these conditions, this did not mean that the LLM would not then, in different conditions, say that the same statement was true or output the "false" statement itself.

This dissonance between the LLM true/false "judgements" and the LLM's failure to accurately replicate those judgements in other conditions is fascinating and very weird. I admire this paper very much and I am really glad they wrote it. That being said, I am not sure the explanation they offer follows from the results provided. They theorize that this could potentially be caused by the linear generation mechanism getting in the way of the LLM's "correct" judgement; this is why they say the LLM knows more than it shows, or, in the words of the title, *knows when it's lying*.

However, LLMs generate different outputs in different situations, some which we judge as correct and some which we judge as incorrect, is to be expected in probabilistic stochastic language models. Their argument assumes that the LLM outputting the token "true" or the token "false" in response to statements when prompted means that the LLM is applying qualities of trueness or falseness onto those statements. However, this doesn't necessarily follow. It is important to remember that with LMs, when we speak about "learning", we're using a mechanistic definition. The LLM learning what is true and what is false isn't exactly possible — it's just that it can be trained to "learn" which statements should have the token "true" follow them and which should have the token "false" follow them; the framework the LLM follows is always about linear generation, not about assigning and then matching up concepts, which requires a bidirectional framework.

The output "reasoning" given to us from the LLM does not actually describe *why* the answer is what it is. We can see this, for instance, in the fact that none of the reasoning lists gender roles and bias as a reason, but there is a huge disconnect in the responses for male and female pronouns. We can also see this because significant amounts of the reasoning are incorrect (in that they are "hallucinating") or incoherent, in that there are gaps where the argument should be — the effect is a large amount of words in grammatical sentences which seem to swirl around and gesture at some empty center where the reason should lie.

Another contemporary AI / NLP paper which does really great work on LLM "mistakes" is *Embers of Autoregression*, which shows that LLMs (they look at cipher-decoding specifically) do much better at generating token strings which resemble correct answers when the output is a high-probability word sequence than when it is a low-probability word sequence, even in

identical situations (decoding a simple cypher). This is interesting because we might expect, based on the way human brains work, that the "difficulty" of each task would be the same, and furthermore, that the basic reasoning process used in decoding one cipher would generalize to the others. The fact that this doesn't happen, and that LLMs make more mistakes with the low-probability word sequences, points against evaluating LLM output with the same framework as human output — which is what we do, for instance, when we speak about AGI. (McCoy et al., 2023)

The work done by McCoy et al. points towards the need to evaluate LLMs on their own terms, to try and look at them as what they are instead of analogizing them to human cognition. LLMs — and all our LMs generally — fail to model the same cognitive structure as human language processing. However, they also don't always succeed at generating outputs that would point to their "capacity to recognize nuanced but critical information in complicated language materials" (Zhang et al., 2023, pg. 1). For instance, Zhang et al. show that they fall prey to the same syntactic "language illusions" as people do, while failing to fall prey to the semantic "language illusions" people do. To fail both, as humans do, would be useful because it would show that LLMs are potentially helpful for modeling human language cognition. To fail neither would be useful because it would show that LLMs are capable of presenting high competence with complex and confusing language structures. But as it stands, the use-cases for this type of model are less clear.

#### **6.2** Conclusion

In this thesis, I propose that large language models (LLMs) have a severed sign, leading to an expression-content gap. What exactly constitutes this severed sign is not completely clear, and further work is needed. However, it is clear from my analysis that the LLM *does* have the substance of the signifier (expression). According to Barthes' modeling, this corresponds to the phonetics portion of the sound-image of the sign; the equivalent in LLMs is the textual output. In addition, I can categorically say that the LLM lacks the substance of the signified; this is all aspects of meaning that require extralinguistic information.



LLMs, roughly, are a system of signifiers without signified. Their outputs, then, are sound-images without concepts, expressions without intent. LLMs lack intentionality and thus communicability; they're without the social component of language. Meaning-making in human-to-human speech is two sided, while meaning-making in LLM-to-human "speech" is one sided; the person reading imparts meanings onto the outputs.

I believe that LLMs can and should be semiologically modeled, and that much of the confusion around what large language models do is due to a lack of this semiological background. Without working to understand what LLM outputs actually signify, we will continue to misunderstand the limits of these language models, and to talk around and past each other when debating them.

So, to revisit the question posed in part 3 of the paper: can meaning be learned through form alone? I would argue no, both because my definition of meaning includes intentionality and the ability to understand extralinguistic information, and also because it is unclear if even 'verbal' or solely linguistic meaning structures such as logical form are capable of being created by language models through pattern application.

Again, it is important to remember that *learn* has different definitions depending on who you ask. As discussed in section 2, the definition of the word *learn* we use in machine learning contexts is an operational one: "A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P if its performance at tasks in T, as

measured by P, improves with experience E." (Mitchell, 1997) The actual process, "cognitive" or otherwise, used to get to an improved P at tasks T is irrelevant.

It is important not to anthropomorphize the language model, or to assume that its learning means the same thing as mine or yours. Or, for that matter, its *attention*, its *speech*, its *mistakes*, and its *intelligence*.

#### REFERENCES

Wolf, T. (2018, August 9). *A learning meaning in natural language processing - the semantics mega-thread*. Medium.

https://medium.com/huggingface/learning-meaning-in-natural-language-processing-the-semantic s-mega-thread-9c0332dfe28e

Mitchell, T. (1997). Machine Learning. McGraw Hill.

Speech and Language Processing. Daniel Jurafsky & James H. Martin. Copyright © 2024. All rights reserved. Draft of August 20, 2024. Chapters 3, 7, 9, 15, and 24 used, as well as appendix F. <u>https://web.stanford.edu/~jurafsky/slp3/</u>

The Internal State of an {LLM} Knows When It's Lying. Amos Azaria and Tom Mitchell. The 2023 Conference on Empirical Methods in Natural Language Processing. 2023. url={https://openreview.net/forum?id=y2V6YgLaW7}

A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. Gomez, {. Kaiser, and I. Polosukhin. Advances in Neural Information Processing Systems, page 5998--6008. (2017)

Hadas Kotek, Rikker Dockum, and David Sun. 2023. Gender bias and stereotypes in Large Language Models. In Proceedings of The ACM Collective Intelligence Conference (CI '23). Association for Computing Machinery, New York, NY, USA, 12–24. https://doi.org/10.1145/3582269.3615599

How poor is the stimulus? Evaluating hierarchical generalization in neural networks trained on child-directed speech (https://aclanthology.org/2023.acl-long.521) (Yedetore et al., ACL 2023)

Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? . In Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT '21). Association for Computing Machinery, New York, NY, USA, 610–623. https://doi.org/10.1145/3442188.3445922

The Internal State of an LLM Knows When It's Lying (https://aclanthology.org/2023.findings-emnlp.68) (Azaria & Mitchell, Findings 2023)

Abubakar Abid, Maheen Farooqi, and James Zou. 2021. Persistent Anti-Muslim Bias in Large Language Models. In Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society (Virtual Event, USA) (AIES '21). Association for Computing Machinery, New York, NY, USA, 298–306. <u>https://doi.org/10.1145/3461702.3462624</u>

Su Lin Blodgett, Solon Barocas, Hal Daumé III, and Hanna Wallach. 2020. Language (Technology) is Power: A Critical Survey of "Bias" in NLP. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics. Association for Computational Linguistics, Online, 5454–5476. https://doi.org/10.18653/v1/2020.acl-main.485

Joseph Weizenbaum. 1966. ELIZA—a computer program for the study of natural language communication between man and machine. Commun. ACM 9, 1 (Jan. 1966), 36–45. https://doi.org/10.1145/365153.365168

Turing, A., 1950, "Computing Machinery and Intelligence," Mind, 59 (236): 433-60.

Oppy, Graham and David Dowe, "The Turing Test", The Stanford Encyclopedia of Philosophy (Winter 2021 Edition), Edward N. Zalta (ed.), URL = <a href="https://plato.stanford.edu/archives/win2021/entries/turing-test/">https://plato.stanford.edu/archives/win2021/entries/turing-test/</a>.

De Saussure, Ferdinand. (1966). Course in General Linguistics (Edited by Charles Bally and Albert Sechehaye, Translated by Wade Baskin). New York, Toronto, London: McGraw-Hill Book Company.

Bender, E.M., & Koller, A. (2020). Climbing towards NLU: On Meaning, Form, and Understanding in the Age of Data. Annual Meeting of the Association for Computational Linguistics.

Elements of semiology / Roland Barthes ; translated from the French by Annette Lavers and Colin Smith  $\cdot$  New York : Hill and Wang, 1968

Ludwig Wittgenstein. The Blue and Brown Books: Preliminary Studies for the 'Philosophical Investigations'. Basil Blackwell, Oxford, 1969

Wittgenstein, Ludwig (1969). On Certainty (ed. Anscombe and von Wright). San Francisco: Harper Torchbooks.

PRIZANT, B. M., & WETHERBY, A. M. (1987). Communicative Intent: A Framework for Understanding Social-Communicative Behavior in Autism. Journal of the American Academy of Child & Adolescent Psychiatry, 26(4), 472–479. doi:10.1097/00004583-198707000-00002

Grice Paul. Utterer's meaning and intentions. The Philosophical Review, 78 (1969).

"Embers of Autoregression: Understanding Large Language Models Through the Problem They are Trained to Solve" R. Thomas McCoy, Shunyu Yao, Dan Friedman, Matthew Hardy, Thomas L. Griffiths

Gordon, P.C., Scearce, K.A. Pronominalization and discourse coherence, discourse structure and pronoun interpretation. *Memory & Cognition* **23**, 313–323 (1995). <u>https://doi.org/10.3758/BF03197233</u>

Centering: A Framework for Modeling the Local Coherence of Discourse (https://aclanthology.org/J95-2003) (Grosz et al., CL 1995)

Gross & Sidner, CL 1986. Attention, Intentions, and the Structure of Discourse(<u>https://aclanthology.org/J86-3001</u>)

Zhang, Yuhan & Gibson, Edward & Davis, Forrest. (2023). Can Language Models Be Tricked by Language Illusions? Easier with Syntax, Harder with Semantics. 10.18653/v1/2023.conll-1.1.

Hall, R. A. (1951). Sex Reference and Grammatical Gender in English. *American Speech*, *26*(3), 170–172. <u>https://doi.org/10.2307/453074</u>

Conrod, Kirby, 'Pronouns and Gender in Language', in Kira Hall, and Rusty Barrett (eds), The Oxford Handbook of Language and Sexuality (online edn, Oxford Academic, 10 July 2018), https://doi.org/10.1093/oxfordhb/9780190212926.013.63, accessed 13 Dec. 2024.

Friedman, Lex (host), 2024. Edward Gibson: Human Language, Psycholinguistics, Syntax, Grammar & LLMs. [Audio Podcast Transcript]. In Lex Friedman Podcast. https://lexfridman.com/edward-gibson-transcript/

Buolamwini, J. & amp; Gebru, T.. (2018). Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification. <i>Proceedings of the 1st Conference on Fairness, Accountability and Transparency</i>, in <i>Proceedings of Machine Learning Research</i>81:77-91 Available from <a href="https://proceedings.mlr.press/v81/buolamwini18a.html">https://proceedings.mlr.press/v81/buolamwini18a.html</a>.

US Bureau of Labor Statistics. (2024) *Labor Force Statistics from the Current Population Survey*. Bureau of Labor Statistics. <u>https://www.bls.gov/cps/cpsaat11.htm</u>

Haiyan Zhao, Hanjie Chen, Fan Yang, Ninghao Liu, Huiqi Deng, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, and Mengnan Du. 2024. Explainability for Large Language Models: A Survey. ACM Trans. Intell. Syst. Technol. 15, 2, Article 20 (April 2024), 38 pages. https://doi.org/10.1145/3639372

Liu, S., Yao, Y., Jia, J., Casper, S., Baracaldo, N., Hase, P., Xu, X., Yao, Y., Liu, C., Li, H., Varshney, K.R., Bansal, M., Koyejo, S., & Liu, Y. (2024). Rethinking Machine Unlearning for Large Language Models. ArXiv, abs/2402.08787.

Lozić E, Štular B. Fluent but Not Factual: A Comparative Analysis of ChatGPT and Other AI Chatbots' Proficiency and Originality in Scientific Writing for Humanities. *Future Internet*. 2023; 15(10):336. <u>https://doi.org/10.3390/fi15100336</u>

R. Thomas McCoy, Robert Frank, Tal Linzen; Does Syntax Need to Grow on Trees? Sources of Hierarchical Inductive Bias in Sequence-to-Sequence Networks. Transactions of the Association for Computational Linguistics 2020; 8 125–140. doi: https://doi.org/10.1162/tacl\_a\_00304

Mulligan, K., Frank, R. & Linzen, T., (2021) "Structure Here, Bias There: Hierarchical Generalization by Jointly Learning Syntactic Transformations", Society for Computation in Linguistics 4(1), 125-135. doi: https://doi.org/10.7275/j0es-xf97

Geng, Y., Li, H., Mu, H., Han, X., Baldwin, T., Abend, O., ... & Frermann, L. (2025). Control illusion: The failure of instruction hierarchies in large language models. arXiv preprint arXiv:2502.15851.