

Plausibility Processing in Transformer Language Models: Focusing on the Role of Attention Heads in GPT2

Anonymous ACL submission

Abstract

The goal of this paper is to enhance our understanding of how Transformer language models process semantic knowledge, especially regarding the plausibility of noun-verb relations. First, I demonstrate GPT2 exhibits a higher degree of similarity with humans in plausibility processing compared to other Transformer language models. Next, I delve into how knowledge of plausibility is contained within attention heads of GPT2 and how these heads causally contribute to GPT2’s plausibility processing ability. Through several experiments, it was found that: i) GPT2 has a number of attention heads that detect plausible relationships between nouns and verbs; ii) these heads collectively contribute to the Transformer’s ability to process plausibility, albeit to varying degrees; and iii) attention heads’ individual performance in detecting plausible noun does not necessarily build a causal relation with GPT2’s plausibility processing ability.

1 Introduction

Transformers are attention-based neural network models (Vaswani et al., 2017), and they have brought breakthroughs in the field of Natural Language Processing achieving state-of-the-art performance in diverse downstream tasks such as machine translation, sentiment analysis, and text summarization, to name a few. Such great performance is thought to be attributed to Transformers’ ability to build dependencies even between long-distant words which attention heads are developed for (Merx and Frank, 2020). To be specific, unlike previous neural network language models (e.g., Simple Neural Networks or Recurrent Neural Networks) that have issues retaining linguistic information coming from distant tokens, attention heads in Transformers enable to represent the meaning of tokens by integrating their contextual information without losing information from distant tokens (Bahdanau et al., 2014).

Provided that Transformer language models consist of multiple attention heads that serve different roles, previous studies examined functions that individual attention heads serve and how language processing work is divided inside Transformers (Clark et al., 2019; Voita et al., 2019; Vig, 2019; Jo and Myaeng, 2020). However, previous studies mostly focused on finding attention heads that process linguistic knowledge intrinsic to language systems such as morphosyntactic rules, and little attention has been paid to semantic knowledge, which requires much of world knowledge going beyond rules in language systems.

Consequently, we only have limited knowledge of how attention heads contribute to Transformers’ general ability to process semantic knowledge. A number of studies (Bhatia et al., 2019; Bhatia and Richie, 2022; Ettinger, 2020; Han et al., 2022; Misra et al., 2020, 2021; Pedinotti et al., 2021; Peng et al., 2022; Ralethe and Buys, 2022) examined how Transformers process semantic knowledge in comparison with humans, but their focus was mostly on the models’ performance from the final hidden state without answering where the specific type of knowledge is preserved or processed in Transformer models. A few studies started investigating how world knowledge is stored in Transformers (e.g., Meng et al. (2022) examined how GPT stores factual associations). However, the previous findings are yet generalizable to all types of semantic knowledge, and thus more studies are needed to understand how Transformers process other types of semantic knowledge. In this regard, the present study aims to advance our knowledge of Transformer language models’ semantic processing by closely investigating the models’ ability to process the plausibility of the relation between nouns and verbs.

As shown in sentences in (1) from Cunnings and Sturt (2018), the semantic plausibility of the relationship between nouns and verbs can be de-

042
043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082

083 terminated by the degree to which semantic features
084 of nouns and verbs match. For instance, in (1a),
085 the syntactic dependent (*plate*) of the verb (*shat-*
086 *tered*) has a feature [+shatterable], which builds a
087 plausible relation with the verb (*shattered*). In (1b),
088 however, the syntactic dependent *letter* does not
089 have a feature [+shatterable], and thus it is semanti-
090 cally implausible dependent of the verb (*shattered*).

- 091 (1) a. Sue remembered the **plate** that the but-
092 **ter shattered** ...
093 b. Sue remembered the **letter** that acci-
094 **dentally shattered** ...

095 In order to examine how such knowledge is pre-
096 served and processed inside Transformer-based lan-
097 guage models, this paper answers the following
098 questions: (i) How similar are Transformer’s plu-
099 sibility processing patterns to humans’?; (ii) How
100 sensitive is each of the attention heads in Trans-
101 formers to plausibility relation?; and (iii) How do
102 these heads make causal effects on Transformers’
103 ability to process semantic plausibility?

104 After comparing patterns in plausibility process-
105 ing between a group of Transformer-based lan-
106 guage models and humans, it was found that GPT2
107 tends to process the plausibility between nouns
108 and verbs in a way that is more similar to humans
109 than other language model types. Several follow-
110 up experiments that especially focus on GPT2 an-
111 swered the last two questions. Specifically, it was
112 uncovered that GPT2 has a set of attention heads
113 that detect semantic plausibility, which are rela-
114 tively diffusely distributed from the bottom layers
115 to the top layers and that they exert causal effects
116 on Transformers’ semantic plausibility processing
117 ability. GPT2’s plausibility processing ability al-
118 most disappeared when the plausibility-processing
119 attention heads are pruned, but the effects of remov-
120 ing a plausibility-processing attention head was not
121 balanced nor proportional to the attention heads’
122 performance in detecting plausible nouns. Rather,
123 it was found that a single attention head accounts
124 for most of plausibility processing ability of GPT2.

125 2 Background

126 **What roles do attention heads serve?** There
127 have been a lot of studies that attempted to ex-
128 plain the language processing mechanism in Trans-
129 formers with analyzing functions distinct attention
130 heads serve (Voita et al., 2019; Vig, 2019; Clark
131 et al., 2019; Jo and Myaeng, 2020). Specifically,

132 Voita et al. (2019) found attention heads specialized
133 for a position, syntactic relation, rare words detec-
134 tion; Vig (2019) found attention heads specialized
135 in part-of-speech and syntactic dependency; Clark
136 et al. (2019) found attention heads specialized in
137 coreference resolution; and Jo and Myaeng (2020)
138 examined how linguistic properties at the sentence
139 level (e.g., length of sentence, depth of syntactic
140 trees and etc.) are processed in attention heads.

141 Despite numerous attempts in examining the
142 roles of attention heads, the focus has been mostly
143 on linguistic knowledge intrinsic to language sys-
144 tems which does not require much world knowl-
145 edge that is indispensable for semantic processing.
146 Thus, it needs to be closely examined how Trans-
147 formers preserve and process such knowledge that
148 facilitates sentence processing.

149 **How do we learn attention heads are specialized**
150 **for certain linguistic knowledge?** In previous
151 studies, attention heads are considered to be able
152 to process a certain type of linguistic knowledge
153 if attention distribution patterns in the attention
154 heads are consistent with the linguistic knowledge
155 (Voita et al., 2019; Vig and Belinkov, 2019; Ryu
156 and Lewis, 2021). However, such regional analysis
157 does not explain how much contribution attention
158 heads make to Transformers’ ability to process lin-
159 guistic knowledge because such information from
160 the attention heads may fade away or be lumped
161 along with the information flows - from bottom
162 layers to top layers - eventually making little con-
163 tribution to Transformers’ ability to process the lin-
164 guistic knowledge. Thus, to rigorously confirm the
165 role of attention heads in processing a certain type
166 of knowledge, it is crucial to analyze the causal
167 effects that they make on Transformer’s ability to
168 process linguistic information (Belinkov and Glass,
169 2019; Meng et al., 2022; Vig et al., 2020).

170 In this sense, this paper will not only examine
171 which attention heads can form attention distribu-
172 tions that are consistent with semantic plausibility
173 knowledge, but also examine how much influence
174 the attention heads can exert on Transformers’ gen-
175 eral ability to process plausibility.

176 3 Comparison between humans and 177 Transformer language models in 178 plausibility processing patterns

179 This section examines how a set of Transformer
180 language models process plausibility of noun-verb
181 relations in comparison with human data.

3.1 Data

In [Cunnings and Sturt \(2018\)](#), it was investigated how the degree of noun-verb plausibility affects the way humans process sentences. There are 32 sets of sentences with varying not only the plausibility of dependent-verb relations but also the plausibility distractor-verb relations¹.

(2)

- a. *plausible - plausible*
... that the **plate** that the butler with the cup accidentally **shattered** ...
- b. *plausible - implausible*
... that the **plate** that the butler with the tie accidentally **shattered** ...
- c. *implausible - plausible*
... that the **letter** that the butler with the cup accidentally **shattered** ...
- d. *implausible - implausible*
... that the **letter** that the butler with the tie accidentally **shattered** ...

3.2 Method

[Cunnings and Sturt \(2018\)](#) measured the degree of difficulty that people have when processing a certain noun-verb pair with reading times that are measured at verb² (*shattered* in (2)). To compare humans' responses with Transformer language models, I compute surprisals ([Hale, 2001](#); [Levy, 2008](#)), also measured at verbs, as a metric that represents processing difficulty of the model, given a large set of evidence manifesting that surprisals computed from neural network language models can simulate human sentence processing patterns ([Futrell et al., 2019](#); [Michaelov and Bergen, 2020](#); [Van Schijndel and Linzen, 2021](#); [Wilcox et al., 2020](#)).

Surprisal is a term that estimates the degree of the unexpectedness of tokens given their preceding context, which is computed by taking the negative log probability of a token conditioned on its preceding words (See Equation (A)). In neural network language models, the surprisal of a word is computed using the softmax-activated hidden state before consuming the word ([Wilcox et al., 2018](#)).

$$\text{Surprisal}(w) = -\log_2 P(w|h) \quad (\text{A})$$

¹In experiments with language models, I removed sets of sentences whose tokens of interest are not recognized as a single token by the tokenizer.

²The original paper also talks about the spillover region following the verbs of interest, but this study focuses on the reading times (total viewing times) measured at the verb.

where h is the softmax-activated hidden state of the sentence before encountering the current word.

Both reading times and surprisals measured at verbs are expected to be greater in sentences with implausible nouns than in ones with plausible nouns since it is less likely to anticipate a certain verb after encountering a noun in an implausible relationship with the verb.

A set of Transformer language models to be tested includes ALBERT ([Lan et al., 2019](#)), RoBERTa ([Liu et al., 2019](#)), BERT ([Kenton and Toutanova, 2019](#)), and GPT2 ([Radford et al., 2019](#)). The versions of models that are tested have 144 attention heads, which are spread across 12 layers with 12 attention heads each. Models are accessed through Huggingface ([Wolf et al., 2019](#)).

3.3 Results

As shown in [Figure 1](#), GPT2 exhibits the highest level of similarity to humans in processing the plausibility of noun-verb pairs, in comparison to other Transformer-based language models.

In addition, further statistical analysis using regression models supports GPT2's similarity with humans in plausibility processing. First, significantly lower processing difficulties are observed when syntactic dependents are in a plausible relationship with the verb than when they are in an implausible relation for both human (estimate = .11, SE = .01, $t = 9.26$, $p < .001$) and GPT2 (estimate = 4.81, SE = .84, $t = 4.86$, $p < .001$).

Also, GPT2 showed marginally significant plausibility effects even with distractors that do not form a dependency relation with the verb (estimate = 1.57, SE = .84, $t = 1.87$, $p = .06$) (i.e., processing difficulties are greater in (b) and (d) than in (a) and (c)), similar to the human data where significant plausibility effects from distractors are found (estimate = .04, SE = .13, $t = 2.85$, $p < .05$)³.

Being inconsistent with the human reading time data that show the interaction effects of dependent-plausibility and distractor-plausibility (estimate = .02, SE = .01, $t = 2.29$, $p < .05$), GPT2 data do not show significant interaction effects (estimate = .89, SE = 1.19, $t = .75$, $p = .46$). This absence of evidence for interaction effects in GPT2 may be

³Plausibility effects observed for distractors in GPT2 and humans are due to the illusion of plausibility ([Cunnings and Sturt, 2018](#)): even distractors that cannot build syntactic dependency with cues (verbs) can be illusorily considered as the syntactic dependents, causing moderate plausibility effects while sentence processing.

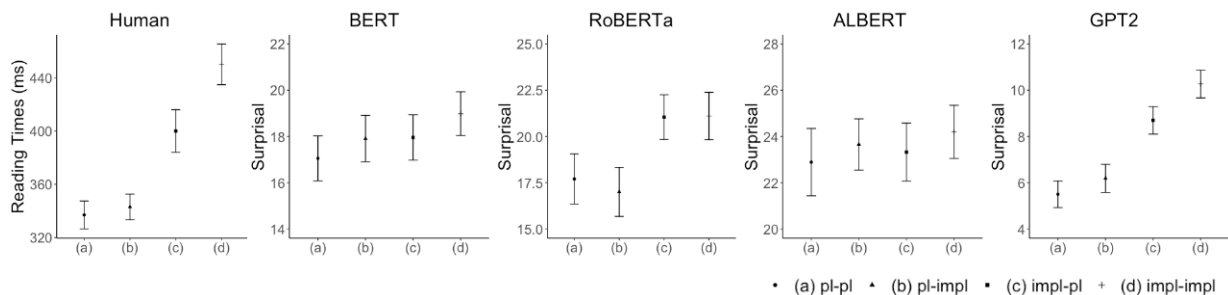


Figure 1: Surprisals computed from Transformer language models and reaction times from human subjects for processing different types of noun-verb pairs. Human reading times are from [Cunnings and Sturt \(2018\)](#). Shapes at the center and intervals for each condition represent means and standard errors.

due to the difference in sample sizes, which can impact the level of statistical significance. It would be possible to observe the interaction effects with the increased data size especially given a trend of interaction in GPT2: the surprisal difference between (a) and (b) is smaller than the surprisal difference between (c) and (d), consistent with human data. For the statistical results from other Transformer-based language models, see [Appendix A](#)

3.4 Discussion

Compared to other language models, GPT2 is found to process plausibility between nouns and verbs in a similar way as humans do. While more rigorous study is required to explain the origin of GPT2’s superior performance in simulating human plausibility processing patterns, I assume that the GPT2’s similarity to humans arises from the psychological plausibility of its decoder-only architecture. In particular, it processes sentences incrementally much like the way humans process sentences (i.e., it constructs the meaning of a certain word only given its prefix, without any influence from the ‘unseen’ next coming words), unlike other types of language models that are tested exploit bidirectional processing (i.e., it process each word of sentences not incrementally, but integrating both preceding and following words.)

Given that GPT2 shows the most similar patterns as humans in processing plausibility of noun-verb relations, the following sections will examine the role that attention heads in plausibility processing, focusing on the GPT2 model.

4 Plausibility processing attention heads in GPT2

This section will examine whether GPT2 has a specific set of attention heads that can sensitively

detect plausibility of noun-verb relations, irrespective of syntactic dependency relation. Experimental stimuli were the same as previous experiment.

4.1 Method

In GPT2’s attention heads, each token allocates different amounts of attention to previous tokens depending on the relevance of the two tokens⁴.

With such a property of Transformers, the capacity of attention heads in detecting plausibility is measured in terms of *accuracy* that indicates how likely the plausible noun is to get higher attention than the implausible noun in a certain attention head (See Equation (B)).

$$Accuracy_{lh} = \frac{\sum_{j=1}^k [Attn(pl_j, v_j) > Attn(impl_j, v_j)]}{k} \quad (\text{B})$$

, where lh refers to the location of attention heads (h for the h th head in the l th layer), j refers to the sentence id, pl_j and $impl_j$ refer to the plausible and implausible nouns to be compared in the j th sentence set, v_j refers to the verb in the j th sentence, and k is the number of sentence sets.

In order to ensure that the heads do not particularly work for tokens that form syntactic dependency but work for semantically related tokens, I measured the accuracy not only using pairs of syntactic dependents (*plate* vs. *letter* in (2)), but using pairs of distractors (*cup* vs. *tie* in (2)). Considering both of noun types enabled to find attention heads that can judge the plausibility between nouns and verbs regardless of syntactic compatibility between them. Thus, there are four comparisons between

⁴The relevance can be defined in terms of functions that attention heads serve. For instance, if an attention head is specialized for detecting *subject-verb* dependency relation, the amount of attention can reflect how likely two tokens are in the *subject-verb* relationship ([Voita et al., 2019](#))

336 *plausible* and *implausible* conditions for each set
337 of sentences: (pl-pl vs. pl-impl), (impl-pl vs. impl-
338 impl), (pl-pl vs. impl-pl), (pl-impl vs. impl-impl),
339 where the first and the second corresponds to syn-
340 tactic dependents and distractors, respectively.

341 4.2 Results

342 I consider attention heads are able to process plau-
343 sible relationships between nouns and verbs when
344 their accuracy in identifying appropriate nouns sur-
345 passes the chance level, having the cutoff as 70%
346 at my discretion. To select attention heads that can
347 process the semantic plausibility regardless of the
348 syntactic dependency relation between the noun
349 and the verb, I consider attention heads whose ac-
350 curacies are greater than 70% in both noun types.

351 With such criteria, eighteen attention heads are
352 recognized to be able to process plausibility: [(0,
353 1), (0, 5), (0, 10), (1, 5), (1, 6), (1, 11), (3, 0), (4,
354 3), (4, 4), (4, 10), (5, 10), (5, 11), (6, 6), (7, 1), (7,
355 9), (8, 3), (8, 10), (9, 4), (10, 7)], where the first
356 numbers refer to indexes of layers and the second
357 refer to indexes of heads (i.e., (i, j) refers to the j th
358 head in the i th layer.) Among the attention heads
359 that are found to process semantic plausibility, two
360 attention heads - (1, 6) and (5, 10) - especially
361 show noteworthy performance in detecting plau-
362 sible, achieving 95% of accuracy. Please refer to
363 Appendix B to see the values from each head.

364 4.3 Discussion

365 This section showed that a set of attention heads
366 are particularly good at processing semantic plausi-
367 bility between nouns and verbs. Such plausibility
368 processing ability seems independent of their abil-
369 ity to process syntactic dependencies since their
370 ability to process plausibility is not limited to pro-
371 cessing syntactic dependents of verbs, but it is also
372 applicable to distractors that do not form any syn-
373 tactic dependencies with verbs.

374 Unlike attention heads specialized for processing
375 a certain syntactic relation and superficial linguistic
376 information such as word position or word rarity is
377 clustered in a relatively small region (Voita et al.,
378 2019), it seems that the components that process se-
379 mantic plausibility are relatively evenly distributed
380 across twelve layers and take up an even greater re-
381 gion: 18 attention heads out of 144 attention heads
382 in the GPT2-small model. In the next section, it
383 will be discussed how these plausibility-processing
384 attention heads collectively exert causal effects on
385 GPT2's plausibility-processing ability.

5 Causal effects of plausibility-processing attention heads on GPT2's plausibility sensitivity

386 In the previous experiment, attention heads capa-
387 ble of detecting plausible relations between nouns
388 and verbs are found. The present section will ex-
389 amine how such attention heads make causal influ-
390 ence on GPT2's sensitivity to plausibility between
391 nouns and verbs. In particular, I attempt to an-
392 swer two questions: (i) How GPT2's responses
393 to plausible/implausible verb-noun pairs change
394 when plausibility-processing attention heads are
395 removed? and (ii) How does GPT2's plausibility-
396 sensitivity change as attention heads are gradually
397 pruned?
398
399
400

5.1 Influence of a set of plausibility-processing heads to plausibility sensitivity

401 In this study, I examine how GPT2's responses
402 to plausible and implausible noun-verb relations
403 change when the plausibility-processing heads are
404 removed.
405
406

5.1.1 Method

407 Surprisals are computed from two models: i) GPT2
408 without plausibility-processing heads and ii) GPT2
409 after removing the same number of attention heads
410 as i), but the heads to prune selected randomly. I
411 included the random-removal model to see whether
412 the disappearance of the plausibility sensitivity in
413 GPT2 is simply attributed to taking away some
414 part of the information in GPT2, or it is caused
415 by specifically removing plausibility processors.
416 In order for reliability, we used 100 different ran-
417 dom attention head sets for ii), and computed the
418 average of surprisals from the 100 models.
419

420 Attention heads were pruned by replacing at-
421 tention values with zeros, following Michel et al.
422 (2019).

5.1.2 Results

423 When removing the plausibility processing atten-
424 tion heads (left in Figure 2), no plausibility effects
425 are found for syntactic dependents (estimate = .77,
426 SE = .53, $t = 1.43$, $p = .15$) and for distractors (esti-
427 mate = .71, SE = .54, $t = 1.32$, $p = .19$). Also, no
428 interaction effects are found (estimate = 0.06, SE =
429 0.76, $t = 0.08$, $p = 0.94$)
430

431 Importantly, such a decrease is not the effect that
432 is caused by simply removing some random compo-
433 nents in GPT2. When randomly selected eight-teen
434 attention heads are pruned (right in Figure 2), the

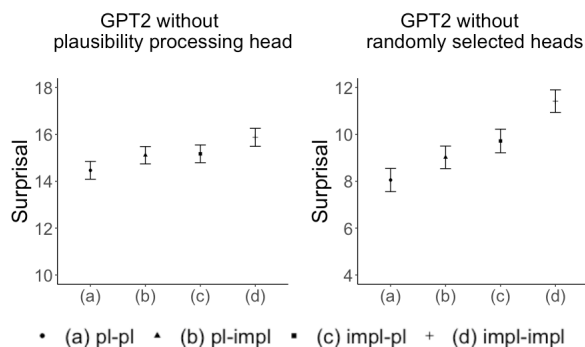


Figure 2: Surprisals computed from GPT2s after removing different sets of attention heads and reaction times from human subjects for processing different types of noun-verb pairs.

GPT2 model better simulates human responses in processing plausibility. In this case, the significant plausibility effects are observed both in syntactic dependents (estimate = 2.40, SE = .69, $t = 3.46$, $p < .001$) and in distractors (estimate = 1.70, SE = .69, $t = 2.45$, $p < .05$), although interaction effects are not found as well (estimate = 0.73, SE = 0.98, $t = 0.75$, $p = 0.46$).

5.2 Gradual changes in GPT2’s plausibility sensitivity as attention heads are pruned

The previous section examined how the set of plausibility-processing attention heads influences GPT2’s responses to plausible or implausible noun-verb relations. Though it was shown that plausibility processing attention heads collectively contribute to GPT2’s ability to process plausibility unlike other sets of attention heads, it is unanswered how individual attention heads contribute to GPT2’s plausibility-processing ability. Do they have balanced contributions to GPT2’s ability to process plausibility? Or, is it that only a small set of plausibility-processing attention heads account for most of the plausibility-processing ability of GPT2? In order to answer these questions, the following experiment investigates how GPT2’s general sensitivity to plausibility gradually changes as attention heads are pruned one by one.

5.2.1 Method

This study operationalizes GPT2’s plausibility sensitivity as the difference in *surprisals* measured at the verbs of interest (*‘shattered’* in (2)) in sentences with plausible nouns and in ones with implausible

nouns as shown in Equation (C).

$$Plausibility\ Sensitivity = surprisal_{impl}(verb) - surprisal_{pl}(verb) \quad (C)$$

, where $surprisal_{pl}(verb)$ and $surprisal_{impl}(verb)$ refer to surprisals measured at the verb in a sentence with a plausible noun and in a sentence with an implausible noun, respectively.

I computed two plausibility sensitivities: one that compares surprisals at verbs when having plausible syntactic dependents of verbs in sentences and having implausible syntactic dependents ($\{(c)+(d)\} - \{(a)+(b)\}$) and the other that compares surprisals when having plausible distractors of verbs and implausible distractors ($\{(b)+(d)\} - \{(a)+(c)\}$).

Both types of plausibility sensitivities are measured at each point after gradually removing a plausibility processing attention head one by one.

Attention heads were pruned in decreasing order of their accuracies⁵ in detecting plausible nouns over implausible nouns.

5.2.2 Results

Figure 3 plots how the plausibility sensitivities for both types of noun-verb relations change as plausibility-processing attention heads are removed gradually.

When it comes to the plausibility sensitivity for distractors, the changes seem to be continuous. Such patterns suggest that the set of plausibility processing attention heads make a collective contribution to plausibility effects for distractors. Such collective contribution that plausibility processing attention heads make is especially supported by the fact that the gradual decrease in plausibility sensitivity over the course of removing 18 attention heads eventually led to the elimination of the statistically significant plausibility effects for distractors as observed in Section 5.1.

In contrast, the sensitivity to plausibility for the relation between syntactic dependents and verbs shows a drastic decrease upon the removal of the attention head (0, 10). The effect from the removal of the head (0, 10) shows that this particular head exerts a huge amount of causal effects on GPT2’s general sensitivity to plausible relations between syntactic subjects and verbs⁶. Figure 4 confirms that the head (0, 10) causes a huge amount of causal

⁵I used the average values of accuracies for dependents and for distractors that were computed in Section 3.

⁶The drastic drop after the removal of the head (0, 10) was also found when attention heads are removed in random order.

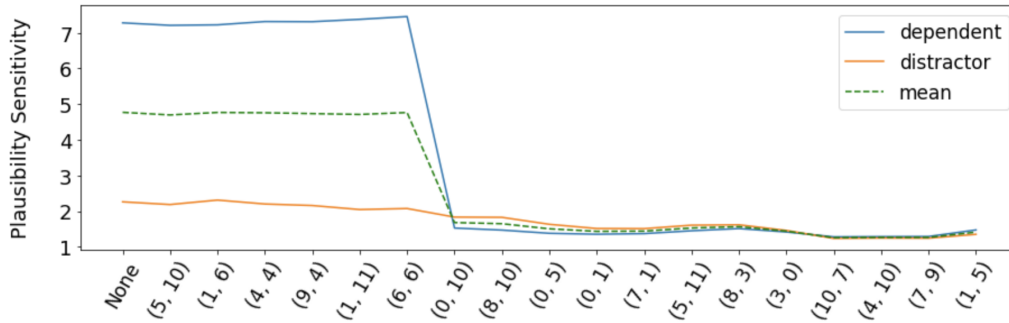


Figure 3: Changes in plausibility sensitivity by noun types as attention heads are gradually pruned. X-axis indicates plausibility-processing attention heads that are pruned at a certain point.

contribution on GPT2’s plausibility processing ability since it reduces the difference in surprisals between plausible conditions and implausible conditions, though it does not alone eliminate the significance in plausible effects for syntactic dependents (estimate = 1.29, SE = 0.61, $t = 2.10$, $p < .05$) or for distractors (estimate = 1.40, SE = 0.61, $t = 2.29$, $p < .05$).

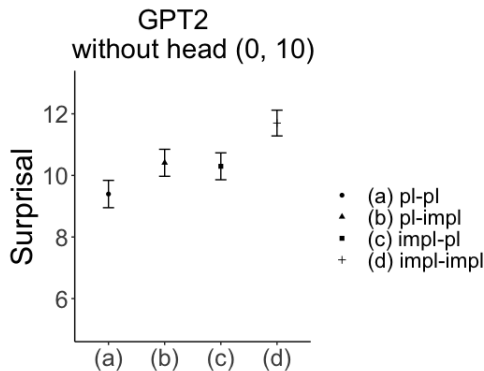


Figure 4: Surprisals by conditions computed with the GPT2 without a single attention head (0, 10)

One additional interesting finding is that the general level of surprisals upon the removal of the attention head (0, 10) increases considerably regardless of the condition. For instance, the removal of the single attention head (0, 10) increases surprisals by 2.79 bits on average across the four conditions, which seems to be huge given that the randomly selected 18 attention heads only led to the 1.89 bits of increase. Such trends indicate one possible explanation of the role of the head (0, 10): it contributes to GPT2’s general ability to predict the next word, and such impact arises in any sentence, not only in the sentences that require plausibility-processing. In the next section, further analysis on the role of the attention head (0, 10) will be

provided to address such a possibility.

5.3 Further analysis on the role of the attention head (0, 10)

To better understand the origin of GPT2’s plausibility processing ability, the present study aims to further examine the role of (0, 10) that make great contribution to plausibility sensitivity in GPT2. In particular, I examine whether the (0, 10) is only specialized for semantic plausibility or is responsible for predicting next words in general sentences which leads to influence plausibility processing.

5.3.1 Method

Perplexity in Equation (D) is the average value of surprisals computed from every tokens in corpus, which can be used to estimate the predictive power of language models in predicting next words given preceding context (Goodkind and Bicknell, 2018).

$$Perplexity(LM) = \frac{1}{M} \sum_{i=1}^m \log_2 P(w_i|h) \quad (D)$$

, where i is the index of words, m is the number of words in corpus, and h refers to the softmax-activated hidden state of the preceding context.

To examine how the general predictive power gets affected by the removal of the head (0, 10) in comparison with the removal of other heads, I computed the perplexities of GPT2 after removing each of 144 attention heads and compared those values. Andersen (1855)’s “The Money Box” story which has 41 sentences was used to compute perplexities.

5.3.2 Results

The perplexity of GPT2 with the entire set of attention heads was 5.47. In most of the cases, the removal of a single head does not seem to considerably affect GPT2’s perplexity, since the perplexity

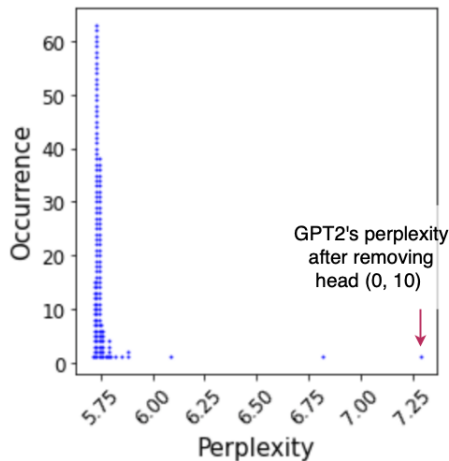


Figure 5: Histogram of 144 perplexities of GPT2 after removing single attention head

remains to be in a similar range after the removal as shown in Figure 5⁷. However, it is clear that the removal of the head (0, 10) seriously harms the general predictive power of GPT2 because the perplexity becomes 7.27 after removing it, which is much greater compared to the most of other attention heads. This suggests that the head having the greatest influence on GPT2’s plausibility processing ability is not specifically specialized for plausibility processing, but rather the attention head contributes to the general predictive power of any kind of sentence.

5.4 Discussion

Results of this section suggest plausibility processing in GPT2 requires a collective contribution from a large set of plausibility processing attention heads, given that plausibility sensitivity decreases continuously as attention heads are gradually pruned.

At the same time, however, it was also shown that the amount of causal effects that each attention head makes are highly imbalanced because the attention head (0, 10), which contributes to GPT2’s general predictive power, leads to a much more drastic decrease in plausibility sensitivity for dependents than other heads. Taken together, although a single attention head can account for a great portion of the plausibility effects, other plausibility-processing attention heads make an additional contribution to GPT2’s plausibility-processing ability.

Interestingly, the head (0, 10) did not achieve noteworthy performance in detecting plausible

⁷For 95% of attention heads, the perplexities change by less than 0.1 bit after the removal.

nouns over implausible nouns in Section 4. This suggests that analyzing the causal effects each attention head makes is indispensable to understanding the role that attention heads serve, provided that the performance that each attention head shows in processing particular linguistic information does not necessarily lead to the eventual contribution to the model’s performance in processing the specific information.

In addition, how the plausibility-processing attention heads affect Transformers’ general ability needs to be investigated in relation to other attention heads. This is especially the case given the results showing that the way plausibility sensitivity decreases as attention heads are pruned depending on the relation types that nouns have (i.e., syntactic dependents or distractors).

6 Conclusion & Limitations

The present study has shown how semantic plausibility is processed in Transformer language models, especially focusing on the role of attention heads. First, I demonstrated that GPT2, whose decoder-only architecture is more aligned with the way humans process sentences, shows greater similarity to humans in plausibility processing compared to other models. Then, a set of experiments showed a number of attention heads, which are diffusely distributed across 12 layers in GPT2, contribute to the model’s sensitivity to plausible relations between nouns and verbs. Moreover, it was observed that they make imbalanced but collective causal contributions to GPT2’ plausibility-processing ability, which establishes the importance of causal effect analysis in attention-head-probing studies.

Although the results provide a window into how Transformers process semantic knowledge of plausibility, this study has a few limitations. First, the scope of the study is restricted to the plausibility of noun-verb relations although there exist many different types of semantic knowledge. For generalizability, the scope of the study needs to be extended. Also, it does not explain how attention heads interact with other components, such as hidden states in different layers or multi-layer perceptrons, in plausibility processing. Such information would be crucial to deepen our understanding of the roles of plausibility processing attention heads since it could explain the mechanism of how the attention heads contribute to Transformers’ ability to process plausibility.

References

Hans Andersen. 1855. *Hans Andersen's Fairy Tales: The money box.*

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.

Yonatan Belinkov and James Glass. 2019. Analysis methods in neural language processing: A survey. *Transactions of the Association for Computational Linguistics*, 7:49–72.

Sudeep Bhatia and Russell Richie. 2022. Transformer networks of human conceptual knowledge. *Psychological Review*.

Sudeep Bhatia, Russell Richie, and Wanling Zou. 2019. Distributed semantic representations for modeling human judgment. *Current Opinion in Behavioral Sciences*, 29:31–36.

Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D Manning. 2019. What does bert look at? an analysis of bert's attention. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 276–286.

Ian Cunnings and Patrick Sturt. 2018. Retrieval interference and semantic interpretation. *Journal of Memory and Language*, 102:16–27.

Allyson Ettinger. 2020. What bert is not: Lessons from a new suite of psycholinguistic diagnostics for language models. *Transactions of the Association for Computational Linguistics*, 8:34–48.

Richard Futrell, Ethan Wilcox, Takashi Morita, Peng Qian, Miguel Ballesteros, and Roger Levy. 2019. Neural language models as psycholinguistic subjects: Representations of syntactic state. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 32–42.

Adam Goodkind and Klinton Bicknell. 2018. Predictive power of word surprisal for reading times is a linear function of language model quality. In *Proceedings of the 8th workshop on cognitive modeling and computational linguistics (CMCL 2018)*, pages 10–18.

John Hale. 2001. A probabilistic earley parser as a psycholinguistic model. In *Second meeting of the north american chapter of the association for computational linguistics*.

Simon Jerome Han, Keith Ransom, Andrew Perfors, and Charles Kemp. 2022. Human-like property induction is a challenge for large language models.

Jae-young Jo and Sung-Hyon Myaeng. 2020. Roles and utilization of attention heads in transformer-based neural language models. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3404–3417. 701
702
703
704
705

Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT*, pages 4171–4186. 706
707
708
709

Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2019. Albert: A lite bert for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942*. 710
711
712
713
714

Roger Levy. 2008. Expectation-based syntactic comprehension. *Cognition*, 106(3):1126–1177. 715
716

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*. 717
718
719
720
721

Kevin Meng, David Bau, Alex J Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in gpt. In *Advances in Neural Information Processing Systems*. 722
723
724
725

Danny Merckx and Stefan L Frank. 2020. Human sentence processing: Recurrence or attention? *arXiv preprint arXiv:2005.09471*. 726
727
728

James A Michaelov and Benjamin K Bergen. 2020. How well does surprisal explain n400 amplitude under different experimental conditions? *arXiv preprint arXiv:2010.04844*. 729
730
731
732

Paul Michel, Omer Levy, and Graham Neubig. 2019. Are sixteen heads really better than one? *Advances in neural information processing systems*, 32. 733
734
735

Kanishka Misra, Allyson Ettinger, and Julia Rayz. 2020. Exploring bert's sensitivity to lexical cues using tests from semantic priming. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4625–4635. 736
737
738
739
740

Kanishka Misra, Allyson Ettinger, and Julia Rayz. 2021. Do language models learn typicality judgments from text? In *Proceedings of the Annual Meeting of the Cognitive Science Society*, volume 43. 741
742
743
744

Paolo Pedinotti, Giulia Rambelli, Emmanuele Chersoni, Enrico Santus, Alessandro Lenci, and Philippe Blache. 2021. Did the cat drink the coffee? challenging transformers with generalized event knowledge. *arXiv preprint arXiv:2107.10922*. 745
746
747
748
749

Hao Peng, Xiaozhi Wang, Shengding Hu, Hailong Jin, Lei Hou, Juanzi Li, Zhiyuan Liu, and Qun Liu. 2022. Copen: Probing conceptual knowledge in pre-trained language models. *arXiv e-prints*, pages arXiv–2211. 750
751
752
753

754	Alec Radford, Jeff Wu, Rewon Child, David Luan,	et al. 2019. Huggingface’s transformers: State-of-	809
755	Dario Amodei, and Ilya Sutskever. 2019. Language	the-art natural language processing. <i>arXiv preprint</i>	810
756	models are unsupervised multitask learners.	<i>arXiv:1910.03771</i> .	811
757	Sello Ralethe and Jan Buys. 2022. Generic overgen-		
758	eralization in pre-trained language models. In <i>Pro-</i>		
759	<i>ceedings of the 29th International Conference on</i>		
760	<i>Computational Linguistics</i> , pages 3187–3196.		
761	Soo Hyun Ryu and Richard L Lewis. 2021. Accounting		
762	for agreement phenomena in sentence comprehen-		
763	sion with transformer language models: Effects of		
764	similarity-based interference on surprisal and atten-		
765	tion. <i>arXiv preprint arXiv:2104.12874</i> .		
766	Marten Van Schijndel and Tal Linzen. 2021. Single-		
767	stage prediction models do not explain the magnitude		
768	of syntactic disambiguation difficulty. <i>Cognitive sci-</i>		
769	<i>ence</i> , 45(6):e12988.		
770	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob		
771	Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz		
772	Kaiser, and Illia Polosukhin. 2017. Attention is all		
773	you need. <i>Advances in neural information processing</i>		
774	<i>systems</i> , 30.		
775	Jesse Vig. 2019. A multiscale visualization of attention		
776	in the transformer model. In <i>Proceedings of the 57th</i>		
777	<i>Annual Meeting of the Association for Computational</i>		
778	<i>Linguistics: System Demonstrations</i> , pages 37–42.		
779	Jesse Vig and Yonatan Belinkov. 2019. Analyzing		
780	the structure of attention in a transformer language		
781	model. In <i>Proceedings of the 2019 ACL Workshop</i>		
782	<i>BlackboxNLP: Analyzing and Interpreting Neural</i>		
783	<i>Networks for NLP</i> , pages 63–76.		
784	Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov,		
785	Sharon Qian, Daniel Nevo, Simas Sakenis, Jason		
786	Huang, Yaron Singer, and Stuart Shieber. 2020.		
787	Causal mediation analysis for interpreting neural		
788	nlp: The case of gender bias. <i>arXiv preprint</i>		
789	<i>arXiv:2004.12265</i> .		
790	Elena Voita, David Talbot, Fedor Moiseev, Rico Sen-		
791	nrich, and Ivan Titov. 2019. Analyzing multi-		
792	head self-attention: Specialized heads do the heavy		
793	lifting, the rest can be pruned. <i>arXiv preprint</i>		
794	<i>arXiv:1905.09418</i> .		
795	Ethan Wilcox, Roger Levy, Takashi Morita, and Richard		
796	Futrell. 2018. What do rnn language models learn		
797	about filler–gap dependencies? In <i>Proceedings of the</i>		
798	<i>2018 EMNLP Workshop BlackboxNLP: Analyzing</i>		
799	<i>and Interpreting Neural Networks for NLP</i> , pages		
800	211–221.		
801	Ethan Gotlieb Wilcox, Jon Gauthier, Jennifer Hu,		
802	Peng Qian, and Roger Levy. 2020. On the predic-		
803	tive power of neural language models for human		
804	real-time comprehension behavior. <i>arXiv preprint</i>		
805	<i>arXiv:2006.01912</i> .		
806	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien		
807	Chaumond, Clement Delangue, Anthony Moi, Pier-		
808	ric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz,		

Table 1: Statistical analysis on plausibility effects in human and Transformer-based language models.

		Human	BERT	RoBERTa	ALBERT	GPT2
Difficulty measurement		reading times	surprisals			
Plausibility effects (syntactic dependents)	estimate	.11	1.10	4.11	.55	4.81
	SE	.01	1.38	1.83	1.77	.84
	t	9.26	.78	2.24	.31	4.86
	p	<.001	.44	<.05	.76	<.001
Plausibility effects (distractors)	estimate	.04	1.03	.06	.87	1.57
	SE	.13	1.38	1.83	1.77	.84
	t	2.85	.75	.03	.49	1.87
	p	<.05	.46	.97	.62	<.10
Interaction effects (dependents \times distractors)	estimate	.02	.17	.76	.11	.89
	SE	.01	1.95	2.59	2.50	1.19
	t	2.29	.09	.29	.04	.75
	p	<.05	.93	.77	.96	.46

A Statistical analysis on plausibility effects

In order for quantitative analysis on how well Transformer language models simulate plausibility effects found in human data (Cunnings and Sturt, 2018), linear regression models for language model data were fit with the following equation: $surprisal \sim$

$subject_plausibility * distractor_plausibility$.

The results are shown in Table 1. Results for human data are from Cunnings and Sturt (2018).

B Scores for detecting the plausible noun-verb relations by attention heads

The performance of attention heads in selecting the plausible nouns in relation with verbs over the implausible ones was measured in terms of *accuracy* in the main text. The details of the method are provided in Section 4.

In addition to accuracy, I also computed attention differences which indicate how much more attention values plausible nouns get compared to implausible nouns (See Equation (E)). The attention differences obtained from all attention heads are shown in Figure 6.

$Attention\ Difference_{lh} =$

$$\sum_{j=1}^k [Attn(pl_j, v_j) - Attn(impl_j, v_j)] \quad (E)$$

,where lh refers to the location of attention heads (hth head in the l th layer), j refers to the sentence id,

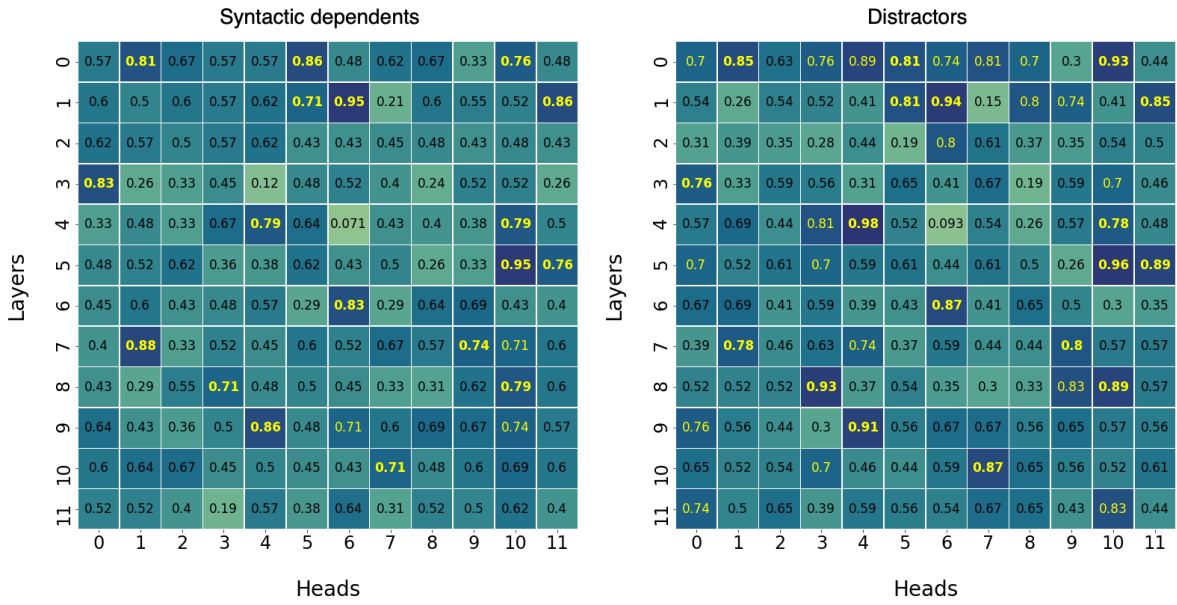
pl_j and $impl_j$ refer to the plausible and implausible nouns to be compared in the j th sentence set, v_j refers to the verb in the j th sentence, and k is the number of sentence sets.

Metrics were computed two times: one by comparing plausible syntactic dependents and implausible syntactic dependents, and the other by comparing plausible distractors and implausible distractors.

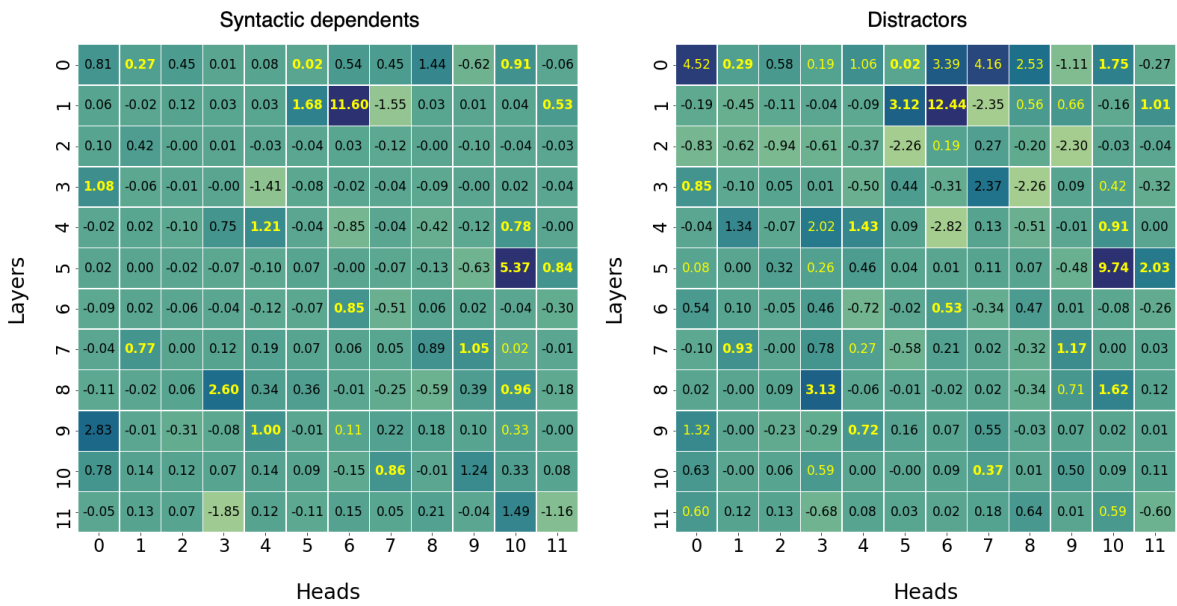
C Changes in surprisal values as attention heads are gradually pruned

In Section 5.2, it was observed how the plausibility sensitivity changes as the plausibility-processing attention heads are gradually pruned. To provide additional information, this section shows how the surprisals for each condition change along with the gradual head-pruning process.

Surprisals were computed at the verb for each sentence in Cunnings and Sturt (2018)’s experimental data. The metrics were computed multiple times after removing one of the plausibility-processing attention heads. The computed surprisal values were then averaged by conditions. The plot that shows how surprisal values change by conditions is given in Figure 7.



(a) Accuracy



(b) Attention Difference

Figure 6: Accuracy and attention difference by attention heads. Attention heads annotated with bold-yellow showed accuracy greater than 0.70 in both subjects-comparison and distractors-comparison and thus considered to be specialized for plausibility processing; Attention heads annotated with non-bold-yellow are the ones that showed accuracy greater than 0.70 only for the corresponding condition; Attention heads annotated with black are found to be insensitive to plausibility (accuracies are less than 0.7 for both noun types).

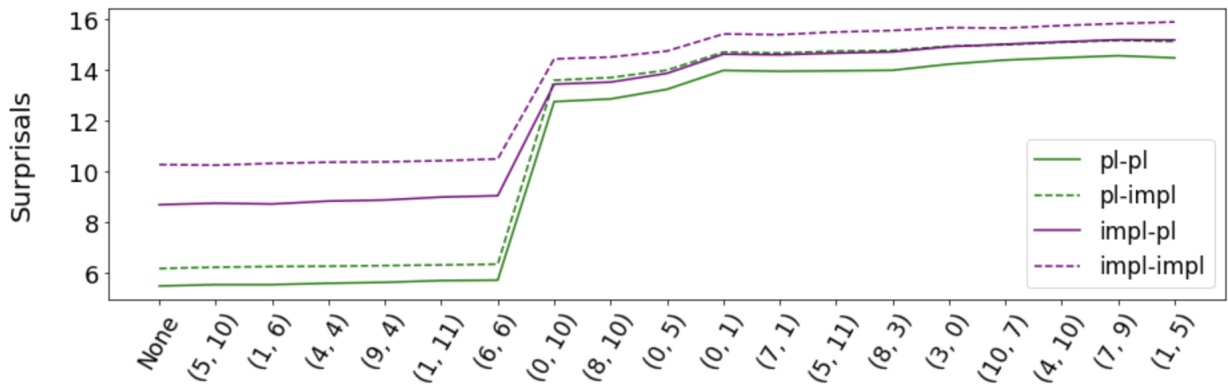


Figure 7: Changes in surprisals by conditions as attention heads are gradually pruned. X-axis indicates plausibility processing attention heads that are pruned at a certain point. Attention heads were removed in decreasing order of accuracies in selecting plausible nouns over implausible nouns.