

Verbing Weirds Language (Models)

Evaluation of English Zero-Derivation in Five LLMs

David R. Mortensen, Valentina Izrailevitch, Yunze Xiao, Hinrich Schütze, and Leonie Weissweiler May 2, 2024

Carnegie Mellon University and LMU Munich





English speakers like to verb words!

Verbing ⊂ (Conversion = Zero-Derivation)

CONVERSION OF ZERO-DERIVATION is pervasive in English (and many other languages): a word with one (or more) prototypical parts of speech is used in a context that calls for **another** part of speech.

You can "verb" various parts of speech in English:

Adjective His hair has begun to gray.

Mass Noun If you don't want to *water* the plants, please *coffee* the graduate students instead.

Count Noun The fascist tried to *knife* me in the back.

Note: to gray, to knife are established usage; to coffee is not.

Are LLMs like **GPT-3.5**, **GPT-4**, **Llama2**, **Mistral**, and **Falcon** robust to <u>conversion</u> of words to <u>non-prototypical parts</u> of <u>speech</u> (as is seen when people English)?

We curated five word lists:

Category	Num	UniMorph Noun	UniMorph Verb
transitive verbs	42	×	1
intransitive verbs	42	×	1
mass nouns	51	✓	×
count nouns	79	 Image: A second s	×
nounce words	49	×	×

All lexical sets were manually curated by a native-speaker linguist.

We tested the ability of LLMs to generalize about zero-derivation by forcing them to answer questions that required construing the same orthographic word as having non-prototypical (and prototypical) parts of speech:

Prototypical: If I thrive daily, do I thrive every day? Non-prototypical: If I health daily, do I health every day? Nonce: If I volice daily, do I volice every day?

We were looking for differences between the non-prototypical condition, on the one hand, and the prototypical and nonce conditions, on the other, in how their responses matched our reference responses We tested the ability of LLMs to generalize about zero-derivation by forcing them to answer questions that required construing the same orthographic word as having non-prototypical (and prototypical) parts of speech:

Prototypical: If I thrive daily, do I thrive every day? Non-prototypical: If I health daily, do I health every day? Nonce: If I volice daily, do I volice every day?

We were looking for differences between the non-prototypical condition, on the one hand, and the prototypical and nonce conditions, on the other, in how their responses matched our reference responses (In the examples, "Yes").

- 1. LLMs answer less consistently with the reference in the non-prototypical condition than the prototypical condition
- 2. LLMs answer less consistently with the reference in the nonce condition than the non-prototypical condition
- 3. There is a correlation between performance on the prototypical conditions and the other two conditions
- 4. The difference in LLM performance can be explained primarily by the size of the models

Results



- GPT-4 nearly perfect across all frame types except INTRANSITIVE
- GPT-3.5 performs similarly
- Falcon performs very well with INTRANSITIVE frames but poorly with MASS NOUN frames
- Llama performs very well on MASS frames and Mistral performs very well on INTRANSITIVE frames but they are otherwise comparable
- Prototypical > Non-Prototypical, Nonce

- Model: Logistic Regression
- Factors:
 - Prototypical part of speech
 - Model type
 - Prototypicality of filler given frame
 - Answer ("yes", "no", "null")
- · Results:
 - All factors significant (p < 0.01)
 - ANSWER TYPE as strongest predictor.
 - PROTOTYPICAL PART OF SPEECH is also a strong predictor

- GPT-4 is the best, followed by GPT-3.5
- The best open model is Falcon, even though it is smaller than Llama2 70B.
- What drags Falcon down seems to be its reluctance to follow instructions (not generalization ability per se)
- Performance on these tasks is not a function of model size, but of other aspects of their training.

prototypical performance > non-prototypical performance

prototypical performance > non-prototypical performance

Supported

prototypical performance > non-prototypical performance non-prototypical performance > nonce performance Supported

prototypical performance > non-prototypical performance non-prototypical performance > nonce performance Supported Not supported prototypical performance > non-prototypical performanceSupportednon-prototypical performance > nonce performanceNot supportedCorrelation between prototypical, non-prototypical, nonce performanceNot supported

prototypical performance > non-prototypical performanceSupportednon-prototypical performance > nonce performanceNot supportedCorrelation between prototypical, non-prototypical, nonce performanceSupported

prototypical performance > non-prototypical performanceSupportednon-prototypical performance > nonce performanceNot supportedCorrelation between prototypical, non-prototypical, nonce performanceSupported

Difference between model size accounts for difference in performance

prototypical performance > non-prototypical performanceSupportednon-prototypical performance > nonce performanceNot supportedCorrelation between prototypical, non-prototypical, nonce performanceSupported

Difference between model size accounts for difference in performance Not supported

- GPT-3.5 and (especially GPT-4) are very good at the verbing task, in part—but not completely—because they follow instructions well
- The open models lag behind, but not in a way that can be explained by model size (Mistral-7B is roughly as good as Llama-70B and Falcon-40B is better than either if you factor out null responses).
- Unlike inflection, existing language models are able to perform this task well.

However—in the case of conversion—this appears to be a limitation of degree, and not in kind.

However—in the case of conversion—this appears to be a limitation of degree, and not in kind.

GPT-4 has near perfect performance on the verbing task, and approaches the human ceiling on the wug task even though GPT-3.5 does not.

However—in the case of conversion—this appears to be a limitation of degree, and not in kind.

GPT-4 has near perfect performance on the verbing task, and approaches the human ceiling on the wug task even though GPT-3.5 does not. However, the ultimate test of such models is whether they can achieve human-like performance with human-like levels of training data, a subject of future work.

Thanks to my Collaborators







Lorenzo

Xiao



Valentina Izrailevitch

Leonie Weissweiler

H Sc

Hinrich Schütze