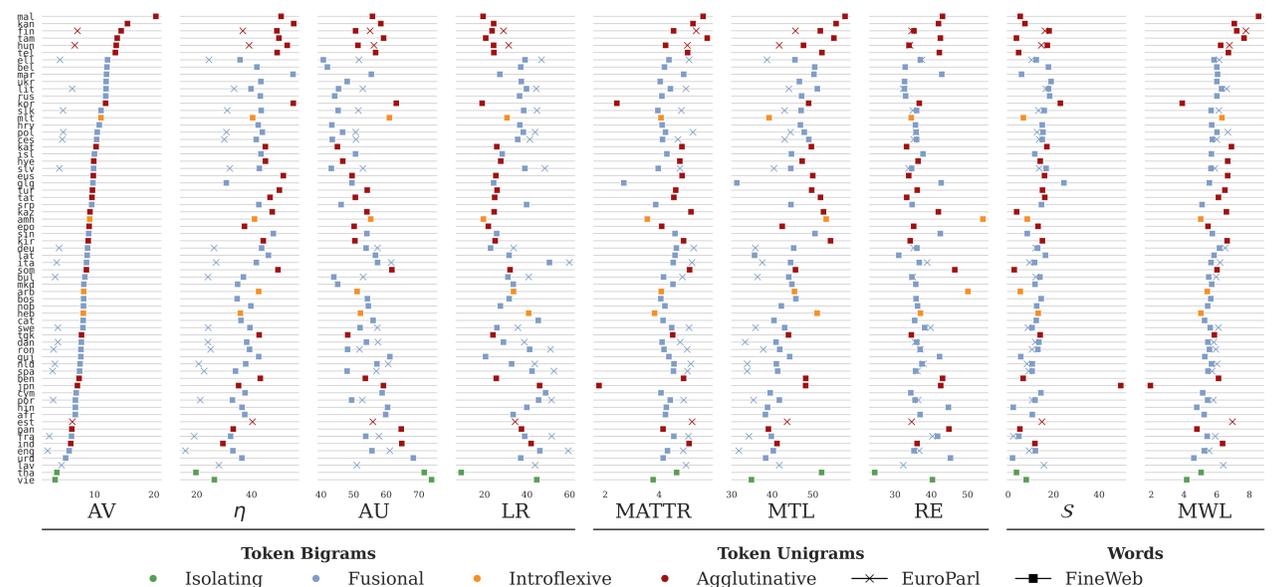
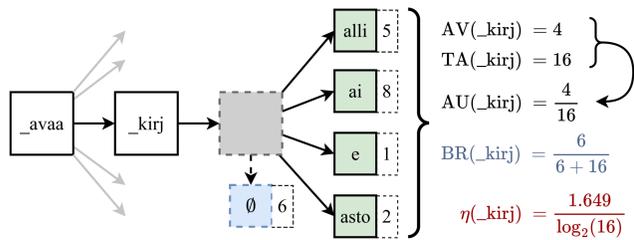


Relating Language Modeling Performance to Morphology

Wessel Poelman* Thomas Bauwens* & Miryam de Lhoneux
LAGoM-NLP, Department of Computer Science, KU Leuven
Empirical Methods in Natural Language Processing 2025



1. Background

- The relation between **intrinsic language differences** and **language modeling** remains an open question.
- Morphology** has been suggested both unimportant and crucial.
- Existing experimental setups contain **confounding factors** that make it hard to draw reliable conclusions.
- We identify these factors, suggest methodological improvements, and introduce corpus-based token bigram metrics as a gradient proxy for morphological complexity.

2. Hypotheses

Starting point: hypotheses and conclusions from award winning paper by Arnett & Bergen (2025):

- H1:** Subword tokenization is less morphologically aligned for agglutinative languages.
- H2:** Subword vocabularies are used more inefficiently for agglutinative languages (“worse quality”).
- H3:** Less training data is available for agglutinative languages.

- H1:** MorphScore; evaluate recall of stem-suffix boundaries.
- H2:** Relate PPLs of monolingual models to unigram metrics.
- H3:** Relate PPLs of monolingual models to data-size adjusted monolingual models based on the language’s UTF-8 encoding compared to English.

Ultimate conclusion: morphology is not a factor, just text encoding and dataset scaling.

3. Confounding Factors

- Languages:** What set of languages is under consideration?
- Grouping:** If results/languages are grouped, is there enough in-group agreement to justify this?
- Tokenization:** What subword tokenization algorithm is used? The vocabulary and segmentation are both crucial.
- Vocabulary size vs. data size:** How does the amount of subword types relate to the amount of training data?
- Corpus domain:** Are tokenizers and models trained on the same data? Are datasets comparable across languages?
- Performance indicator:** What metric is used to evaluate and compare tokenizers and models across languages?

Controlling for these factors can be seen as criteria for the “ideal” experiment. One has to work *backwards* from this to something feasible in terms of data and analysis.

4. Tokenizer Alignment

Segmentation	MS	F_1
<code>gathered</code> → <code>gather/ed</code>	1	1.0
<code>gathered</code> → <code>gathered</code>	0	0.0
<code>gathered</code> → <code>g/a/t/h/e/r/e/d</code>	1	0.25
<code>arabaları</code> → <code>araba/lar/ı</code>	1	1.0
<code>arabaları</code> → <code>araba/ları</code>	1	0.5
<code>arabaları</code> → <code>arabalar/ı</code>	0	0.5

MorphScore vs full alignment.

Consider a fusional word of the form `aaaaa/bc` and an agglutinative word of the form `wwwww/xx/yy/zz`.

- Both have a stem-suffix boundary, but the odds of the stem and suffix sticking together is lower in agglutinative languages, since there is a chance that the suffix morphs already form a bigger token: `wwwww/xxyy/zz`,
- Missing a stem-suffix boundary produces highly specific tokens for both: `aaaaab/c` or `aaaa/abc`, and `wwwwwx/x/yy/zz` or `wwwww/wxx/yy/zz`.
- Missing an agglutinative suffix-suffix boundary is potentially much worse: in `wwwww/xy/y/zz`, the `yy` morph has lost half its length, making it potentially meaningless.
- Misses can cascade: `wwwww/xy/yz/z` or `wwwww/xy/yz/z`; the three morphemes will be harder for a model to piece together from the embeddings of two tokens.

Morphological alignment is a bad predictor for language modeling.

5. Metrics

Morphological Complexity

Unigram-based

- Moving average type-token ratio (MATTR)
- Mean token length (MTL)
- Rényi Efficiency (RE)

Word-based

- Tokens per character averaged per word (S)
- Mean word length (MWL)

Language Modeling Performance

Model	Sequence	PPL
A	<code>Sabe_jugar_al_ajedrez</code>	20
B	<code>Do_you_know_how_to_play_chess</code>	22
B	<code>Can_you_play_chess</code>	18

Which model is “better” is arbitrary.

6. Accessor Variety

- Existing metrics are based on **unigrams** or **words**.
- Unigrams are not measuring complexity.
- Language models do not use words.

Bigrams! Idea stems from Harris (1955) and Feng et al. (2004). Results in gradient, corpus-based view of morphological complexity.

Proposed Metrics

- Accessor Variety (AV): tokens following each other per type
- Entropic efficiency of AV distribution (η)
- Accessor Uniqueness (AU): how many unique types?
- Lexicalization Ratio (LR): rare or memorized types; ignored

Findings

- AV and η provide a gradient view of morphological complexity that is closely related to language modeling.
- Difficulty is probably caused by having *more* and *more equally likely follow-up* options at each timestep. This is what AV and η measure.
- Coarse groupings like *fusional* or *agglutinative* hide information in these setups.

7. Conclusions

- Morphological alignment is a bad predictor for a performance gap between morphological systems.
- Unigram metrics do not measure language modeling difficulty.
- Dataset size alone cannot explain away morphology.
- Bigram metrics such as AV and η are better proxies for morphological complexity that do not require expert annotations.
- And finally...

Keep confounding factors in mind!

Contact & Acknowledgements



Email



LinkedIn



Publications

WP and TB are funded by a KU Leuven Bijzonder Onderzoeksfonds C1 project with reference C14/23/096.

