

Vocab Diet: Reshaping the Vocabulary of LLMs with Vector Arithmetic

Anonymous ACL submission

Abstract

Large language models (LLMs) encode word-form variation (e.g., “walk” vs. “walked”) as linear directions in embedding space. However, standard tokenization algorithms treat these variations as distinct tokens—filling the size-capped vocabulary with surface form variants (e.g., “walk”, “walking”, “Walk”), at the expense of less frequent words and multilingual coverage. We show that many of these variations can be captured by *transformation vectors*—additive offsets that yield the appropriate word’s representation when applied to the *base form* word embedding—in both the input and output spaces. Building on this, we propose a compact reshaping of the vocabulary: rather than assigning unique tokens to each surface form, we compose them from shared *base form* and *transformation* vectors (e.g., “walked”=“walk”+*past tense*). We apply our approach to multiple LLMs and across five languages, removing up to 10% of vocabulary entries—thereby freeing space to allocate new, more diverse tokens. Importantly, we do so while also expanding vocabulary coverage to out-of-vocabulary words, with minimal impact on downstream performance, while keeping the pretrained backbone frozen and only training lightweight adaptation modules. Our findings motivate a foundational rethinking of vocabulary design, moving from string enumeration to a compositional vocabulary that leverages the underlying structure of language.¹

1 Introduction

Modern large language models (LLMs) typically rely on subword tokenization algorithms like byte-pair encoding (BPE; Sennrich et al., 2016). Such methods allocate tokens to frequent words and split less frequent ones into sequences of subword tokens—minimizing the number of tokens needed to represent typical textual data. To fur-

¹We will release our code upon publication.

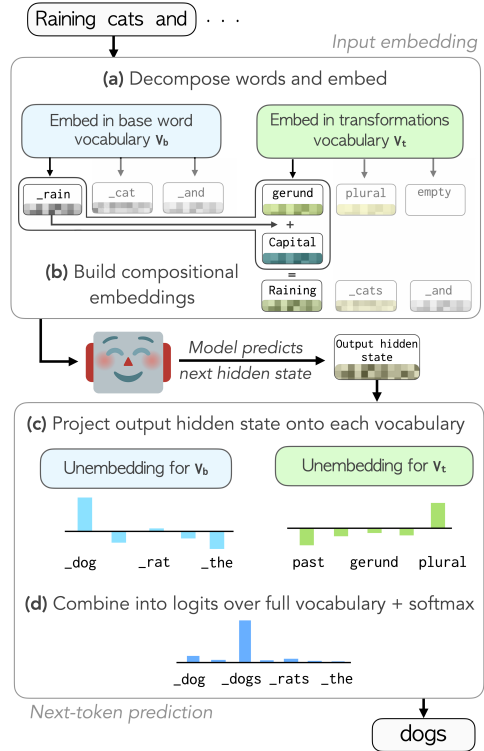


Figure 1: **Compositional vocabulary for LLMs.** **Top:** Input tokens are represented by (a) decomposing them into base words (\mathcal{V}_b) and transformations (\mathcal{V}_t), and (b) feeding the composite embeddings to the model. For example, “cats” becomes “cat” + *plural*. **Bottom:** The next token is predicted by (c) computing logits independently over base words and transformations, and (d) combining them into next-token probabilities. Our approach works seamlessly with pretrained LLMs with lightweight parameter updates, creating a more compact vocabulary that supports a wider array of words.

ther reduce inference costs, as well as broaden domain and multilingual coverage, recent models use ever-larger vocabularies, often exceeding 100k tokens (Grattafiori et al., 2024; OpenAI, 2024; Yang et al., 2024). While recent work calls for scaling up the vocabulary even further (Tao et al., 2024; Huang et al., 2025), the computational cost of supporting large vocabularies forces developers to cap

its size (Dagan et al., 2024; Wijmans et al., 2025).

Standard tokenization, while effective, often leads to a disproportionate allocation of the vocabulary (§3). Common words occupy multiple token slots for their various forms (e.g., “walk”, “walks”, “walking”, etc.), leaving less room for uncommon words and multilingual coverage—ultimately hurting both performance and inference costs (Petrov et al., 2023; Ahia et al., 2023; Ali et al., 2024). More fundamentally, it ignores a striking property of LLMs: their tendency to encode relationships between words as simple *linear directions* (Park et al., 2024; Marks and Tegmark, 2024). But can we harness this structure to build more compact vocabularies, without sacrificing expressivity?

We begin by investigating how LLMs represent word form variation. Building on the idea of vector arithmetic in embedding space (Mikolov et al., 2013b), we examine whether common word-form transformations—including morphological inflection (“walked”), derivation (“walker”) and capitalization (“Walk”)—can be captured as consistent *transformation vectors* added to a *base form* word embedding (§4). Focusing on five morphologically diverse languages, we use the UniMorph word-form database (Batsuren et al., 2022) to identify token pairs of base- and surface-form words exemplifying the same relation. We then compute the average offset vector for each relation, and use these as *transformation* vectors. Our results show that adding these vectors to *base form* embeddings yields representations that the model interprets similarly to the expected surface form (Ghandeharioun et al., 2024). Interestingly, this holds even when the target word is not represented as a single token in the vocabulary,² indicating that LLMs process and interpret word forms compositionally (§5).

Building on these insights, we propose a compact reshaping of the vocabulary, building word embeddings from shared components (Figure 1): a *base form* vector for the core lexical item and a *transformation* vector for encoding word-form variation. Rather than assigning a unique token embedding to each surface-form, we truncate the model’s embedding tables to remove any inflected forms, and introduce a small set of *transformation* embeddings—enabling us to represent the discarded words compositionally (e.g., “walked” as “walk”+*past tense*) in both input and output. Importantly, we only update lightweight adaptation

parameters: we fine-tune the *transformation* embeddings and train a LoRA adapter on the final $k = 6$ transformer blocks, leaving all other parameters frozen. In experiments across five models and five languages, our method removes up to 10% of the vocabulary tokens while maintaining performance over a suite of downstream tasks when representing words compositionally (§6). We further run a pretraining proof-of-concept experiment (§7), showing that a compositional model can be trained from scratch while removing 41.6% of vocabulary entries with comparable accuracy (38.3 vs. 38.1), and with an training throughput increase of 11.5%.

In summary, we introduce compositional structure into language model vocabularies, enabling efficient representation of linguistic diversity through shared building blocks. Our approach reduces redundancy in token allocation while expanding lexical coverage—with only lightweight adaptation. Our experiments demonstrate that LLMs can naturally operate with these representations, and establish compositional vocabularies as a competitive alternative to standard surface-form tokenization for future language models.

2 Background: Token Allocation in Language Model Vocabularies

Tokenization bridges natural language and model representations: it decomposes text into sequences of tokens from a fixed vocabulary, where each token is an atomic string unit for which the model learns specialized, single-vector embeddings. These vocabularies are almost universally built using byte-pair encoding (BPE; Sennrich et al., 2016), which iteratively merges the most frequent token pairs—from characters to subwords to words—in attempt to optimally compress the text using a predetermined vocabulary size.

As LLM vocabularies grow larger (e.g., Gemma Team, 2024; Aryabumi et al., 2024), there is growing recognition that vocabulary resources can be better allocated. Recent studies point to stark imbalances in token allocations across languages, negatively impacting both model cost (Petrov et al., 2023; Ahia et al., 2023) and performance (Ali et al., 2024; Limisiewicz et al., 2023; Toraman et al., 2022), motivating techniques for post-hoc vocabulary expansion to reduce costs for a specific language or domain (Han et al., 2025; Nakash et al., 2025; Liu et al., 2024b; Minixhofer et al., 2024).

Another line of research advocates for scaling

²E.g., a word like “walkable” is split into [_walk, ab1e].

up the vocabulary together with model size to unlock performance gains in the model’s main language (Tao et al., 2024; Huang et al., 2025; Liu et al., 2025). Still, expansion is ultimately bounded by memory and compute constraints (Dagan et al., 2024; Wijmans et al., 2025), underscoring the importance of carefully reconsidering how the token vocabulary is allocated.

3 Word Structure and Redundancy in Vocabulary Design

One underexplored source of inefficiency in current vocabulary design is the treatment of morphologically related word forms as independent tokens. In high-resource languages like English, this often results in large clusters of surface variants—*walk*, *walks*, *walking*, *walked*—each assigned a separate token, despite their shared meaning and structure.

To quantify this redundancy, we examine the English whole-word tokens in the GPT-4 tokenizer (OpenAI, 2024)—the base tokenizer for many recent LLMs (Grattafiori et al., 2024; Yang et al., 2024; OLMo et al., 2024). We use UniMorph’s English lexicon (Batsuren et al., 2022) to identify tokens that are English words,³ finding 24.6k such tokens (Figure 2, left side).⁴ Ignoring case (e.g., equating “walk” with “Walk”) reduces this to 17.7k unique types. Further accounting for inflectional and derivational relations reduces this to just 14.3k *base forms*, a total of 42% reduction.

Rather than assigning each word form with a distinct, independently-learned token, what if we could model these processes as *transformations* applied to a compact set of base words? Our analysis shows that, beyond reconstructing every in-vocabulary word, these tokens can further represent 98k out-of-vocabulary words (Fig. 2, right), which are currently represented using multiple tokens.

Altogether, this motivates a structured vocabulary design that composes word forms from shared blocks, yielding vocabularies that are simultaneously more compact and more expressive while scaling effectively across domains and languages.

³We only consider tokens that start with a leading space as whole word tokens; tokens without it can sometimes occur mid-word (like “ask” in “task”, compared to “_ask”).

⁴Out of 100k tokens, there are 41.3k tokens with a leading space in the vocabulary that are composed of English letters. Roughly 60% are identified as valid English words. The rest are either code-related terms, sub-words, or proper nouns. The other 60k tokens are either sub-words or non-English tokens.

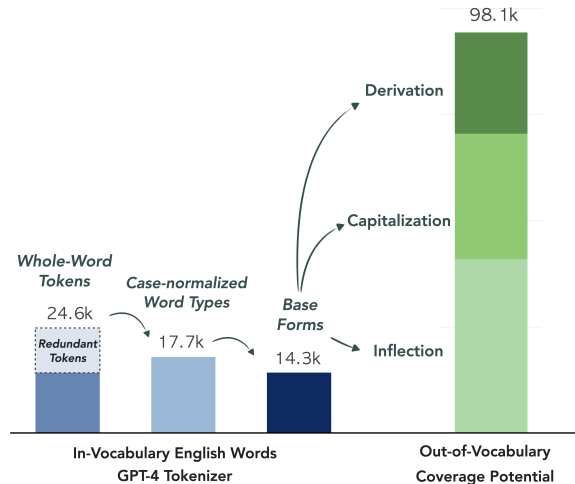


Figure 2: **Structure in LLM vocabularies and potential for compositional design.** **Left:** Many in-vocabulary English word tokens in the GPT-4 tokenizer are surface variants of other tokens—differing only by case, inflection, or derivation—reducing from 24k tokens to just 14k base form words. **Right:** The existing set of base forms and transformations can be used to compose over 98k currently out-of-vocabulary words, highlighting the inefficiencies of current vocabularies and the potential of a compositional design.

4 Composing Words from Base Forms and Transformations

We propose a compositional representation approach in which each surface form is constructed from a base word and a set of transformation vectors. Formally, let $\mathcal{V}_{\text{orig}}$ denote the model’s original token vocabulary. We define a subset $\mathcal{V}_b \subset \mathcal{V}_{\text{orig}}$ as the base-word vocabulary, consisting of canonical lexical forms (e.g., *walk*) and any auxiliary tokens (e.g., punctuation, sub-words, code segments, words in non-target languages). We also introduce a transformation vocabulary \mathcal{V}_t , which consists of a small number of vectors corresponding to morphological operations such as inflection or derivation, or other word-level processes like capitalization.

In our scheme, a word w is represented by a base $b_w \in \mathcal{V}_b$ and a set of transformations $T(w) \subset \mathcal{V}_t$:

$$\mathbf{e}_w = \mathbf{e}_{b_w} + \sum_{t_i \in T(w)} \mathbf{e}_{t_i} \quad (1)$$

where \mathbf{e}_{b_w} and \mathbf{e}_{t_i} are rows from embedding matrices E_b and E_t , respectively. For base words and auxiliary tokens, $T(w) = \emptyset$.

This decomposition applies both at input and output: At input, we replace direct lookup with Eq. 1. At output, we replace the model’s large unembedding matrix U with two separate, smaller matrices

for *base forms* and *transformations*: U_b and U_t . Given an output state \mathbf{h} , we score each candidate next-token w by separately projecting \mathbf{h} onto U_b and U_t and summing the relevant dot-products:

$$\text{logit}(w) = \mathbf{h} \cdot \mathbf{u}_{b_w} + \sum_{t_i \in T(w)} \mathbf{h} \cdot \mathbf{u}_{t_i} \quad (2)$$

where \mathbf{u}_{b_w} and \mathbf{u}_{t_i} are the corresponding columns of U_b and U_t for w 's components. Importantly, our method is agnostic to whether w is in-vocabulary (IV) or out-of-vocabulary (OOV), as long as its base form is IV. We define the compositional vocabulary \mathcal{V} as all words that can be constructed from $(b_w, T(w))$ combinations.

Vocabulary decomposition map. To apply this framework, we construct a mapping $w \mapsto (b_w, T(w))$ from surface forms to their base forms and matching transformations. We use UniMorph (Batsuren et al., 2022), a multilingual word form database, to identify base forms and their inflected and derivated forms. Transformation labels are drawn from UniMorph’s standardized tags (e.g., V; PST) with added rules for capitalization. Then, to build a decomposition map for a given tokenizer’s vocabulary $\mathcal{V}_{\text{orig}}$, we iterate over its tokens, identify base forms, and map all related surface forms—whether in-vocabulary or not—to their base and transformation sets.

Notably, the decomposition map could also be built using sources other than gold morphological annotations. Plausible alternatives include unsupervised morphological segmentation (Creutz and Lagus, 2005; Smit et al., 2014; Abdelali et al., 2016), statistical morphology learning (Creutz and Lagus, 2007), or bootstrapping morphological annotations via LLM-generated analyses. In this work, we focus on whether these compositions are interpreted and used correctly by models, and therefore use high-quality UniMorph annotations; future work could operationalize unsupervised alternatives.

Computing the transformation vectors. To define the transformation vectors themselves (i.e., the entries in E_t and U_t) we revisit the idea of vector arithmetic in embedding space (Mikolov et al., 2013b). Let O be an embedding matrix of $\mathcal{V}_{\text{orig}}$, and let $b(w) : \mathcal{V}_{\text{orig}} \mapsto \mathcal{V}_b$ be a function that maps a word to its base form. For each transformation t , we extract the set $R(t) = \{(w, b(w)) \mid t \in T(w)\}$ of word pairs in $\mathcal{V}_{\text{orig}}$ that exemplify t (e.g., *walk*

and *walked* for $t = \text{past tense}$).⁵ We then compute the average offset of their respective embeddings:

$$\mathbf{o}_t = \frac{1}{|R(t)|} \sum_{w \in R(t)} (\mathbf{o}_w - \mathbf{o}_{b(w)}) \quad (3)$$

We compute this separately for all $t \in \mathcal{V}_t$ in both the embedding and unembedding spaces, yielding *transformation* vectors for input and output. While prior work analyzed such linearity in the embeddings of LLMs (Park et al., 2024, 2025), to the best of our knowledge, our work is the first to leverage this for end-to-end language modeling.

5 Do LLMs Understand Compositional Word Representations?

We now turn to our core question: can LLMs that were pretrained with standard vocabularies interpret our compositional embeddings—sums of *base form* and *transformation* vectors—as intended?

Recent work has shown LLMs build-up and resolve the meanings of input tokens across their early layers, a process referred to as *detokenization* (Kaplan et al., 2025; Feucht et al., 2024; Gurnee et al., 2023). This was particularly observed for multi-token words or in-vocabulary words split into multiple tokens (e.g., due to typos). Building on this, we feed models with compositional inputs and inspect whether the embedding and early layer representations have successfully resolved into the intended surface form meanings. To interpret these internal representations, we follow Kaplan et al. and use Patchscopes (Ghandeharoun et al., 2024), a prompting method to probe the contents of a hidden state using natural language.

Languages and models. We experiment with five morphologically-diverse languages: English, Arabic, German, Russian and Spanish. For English, we use three LLMs: LLaMA-3-8B (Grattafiori et al., 2024), Qwen2.5-7B (Yang et al., 2024), and OLMo-2-7B (OLMo et al., 2024). As coverage of whole-word tokens in these models’ vocabularies for other languages is narrow,⁶ we use models with dedicated tokenizers for them: ALLaM-7B for Arabic (Bari et al., 2025) and EuroLLM-9B for the three other languages (Martins et al., 2025).⁷ In experiments for a specific model and language pair, we

⁵To obtain a “clean” signal for *transformations*, we only use w that demonstrate a *single* transformation ($|T(w)| = 1$).

⁶This restricts both the base-word lexicon, and the number of existing *base-inflected* pairs for extracting transformations.

⁷All models have vocabularies of 100k or more tokens, except ALLaM with 64k (but roughly 32k are for Arabic).

construct the vocabulary decomposition and transformation vectors (§4) only for that language, ignoring words in other languages.

Examining word representations. For each model and language pair, we iterate over all words w that could be composed from the *base forms* and *transformations* extracted from its vocabulary (§4). Next, given a surface form w , we replace the token embedding for w with its compositional representation e_w (Eq. 1), and feed it to Patchscopes to generate its textual description.⁸ We then evaluate whether the Patchscopes interpretation of the compositional embedding e_w matches the target word w (*embed*). We also examine whether the model successfully *detokenizes* compositional embeddings in its early layers: we feed e_w to the model without any context, extract the resulting hidden states at the first $k = 10$ layers, and report whether the Patchscopes interpretation matches the target word w in at least one layer (*detok*).

English results. We begin by examining English words that exist as single tokens in Llama-3-8B’s original vocabulary $\mathcal{V}_{\text{orig}}$ (Table 1, *in-vocab*). We observe that most inflectional transformations—such as verb tense (past, present participle) and number (plural)—as well as capitalization, are often correctly resolved by the model already at the embedding layer (*embed*), and almost always at early internal layers (*detok*). For example, $e_{\text{walk}} + e_{\text{past}}$ is interpreted by Patchscopes as “walked”. In contrast, derivations (e.g., “walk”→“walkable”), which rarely occur as single-tokens in the vocabulary, are seldom recognized by the model and often resolve as the base word instead. This suggests that models learn weaker linear structure for rare relations, or that *transformation* vectors built using small sample sizes show weaker generalization.

We next examine out-of-vocabulary words, i.e., English words that can be composed using the *base forms* and *transformations* but are *not* found as a single token in the original vocabulary (Table 1, *out-of-vocab*). Using our decomposition map, we construct single-vector representations for these words and feed them to the model. Surprisingly, many of these are resolved as the intended word form already at the embedding layer, with Patchscopes

⁸Following Kaplan et al. (2025), we use the Patchscopes prompt “[X], [X], [X], [X],”, where we replace the placeholder token ([X]) with a hidden state \mathbf{h} and let Patchscopes generate text. We expect Patchscopes to generate the intended word form if \mathbf{h} indeed captures it. For languages other than English, we add the prefix “In {language_name}”.

Transformation	In-vocab.			Out-of-vocab.		
	<i>embed</i>	<i>detok</i>	<i>N</i>	<i>embed</i>	<i>detok</i>	<i>N</i>
Inflection						
Plural (N)	92%	96%	0.8k	30%	56%	3.4k
Plural (N) & Present Singular (V)	87%	91%	1.6k	43%	75%	2.1k
Present Singular (V)	90%	91%	0.1k	64%	82%	0.3k
Past (V)	71%	81%	0.6k	9%	29%	2.9k
Past Participle (V)	64%	93%	14	14%	38%	21
Gerund (V)	83%	93%	0.2k	17%	34%	3.2k
Superlative (ADJ)	71%	94%	31	5%	29%	0.4k
Comparative (ADJ)	40%	83%	30	3%	36%	0.4k
Capitalization						
	80%	89%	6.0k	72%	85%	8.4k
Derivation						
-y	24%	47%	17	2%	12%	1.5k
-er	8%	17%	12	0%	6%	2.6k
-al	25%	25%	8	0%	9%	0.7k
un-	0%	33%	3	0%	2%	3.3k
re-	67%	67%	3	0%	10%	1.8k
-ic	100%	100%	2	4%	21%	0.4k
All derivatives	31%	45%	51	0%	3%	31.4k

Table 1: Accuracy of Patchscopes interpretations for compositional input representations (i.e., *base form* + *transformation* embeddings) of in-vocabulary and out-of-vocabulary English words in Llama-3.1-8B. We report successful resolution both at the embedding layer (*embed*), and after detokenization in early layers (*detok*). N indicates the number of surface forms evaluated per category. Compositional embeddings of capitalization and inflectional forms are very often resolved correctly—even for many out-of-vocabulary words, which never occur as single input vectors during pretraining. Derivatives remain challenging—likely due to rarely appearing as single tokens in the vocabulary.

generating the full, multi-token word, especially for inflections and capitalization. Similarly to in-vocabulary results, we observe higher successful resolution rates for early-layer detokenization, and representing out-of-vocabulary derivations using compositions generally fails. We observe similar results for English in other models (Appendix A).

Multilingual results. We repeat the same experiment on each of the other languages. Since each language has different types and number of inflectional and derivational processes,⁹ we aggregate results over five categories: adjective inflection, verb inflection, noun inflection, derivation and capitalization. Our results (Table 2) show that LLMs can correctly interpret compositional word representations across diverse languages and morphological structures. Surprisingly, some *transformation* vector types (e.g., adjective or verb inflections) work better for out-of-vocabulary representation than in English, hinting that models learn stronger linear encodings of morphological structure when the token vocabulary is more limited—a phenomenon

⁹We treat each UniMorph tag as its own *transformation*.

Language	Capitalization		Noun Inflection		Adjective Inflection		Verb Inflection		Derivation		
	<i>In-Vocab.</i>	<i>Out-Vocab.</i>	<i>In-Vocab.</i>	<i>Out-Vocab.</i>	<i>In-Vocab.</i>	<i>Out-Vocab.</i>	<i>In-Vocab.</i>	<i>Out-Vocab.</i>	<i>In-Vocab.</i>	<i>Out-Vocab.</i>	
<i>ALLaM</i>	Arabic	—	—	77% (1.8k)	14% (3.6k)	69% (0.5k)	23% (1.0k)	41% (1.0k)	14% (2.7k)	—	—
<i>EuroLLM</i>	German	95% (0.2k)	74% (0.4k)	—	—	21% (0.3k)	7% (1.3k)	82% (0.3k)	36% (1.2k)	—	—
	Russian	97% (66)	88% (0.7k)	63% (0.6k)	21% (4.2k)	100% (50)	89% (94)	83% (6)	30% (10)	—	—
	Spanish	97% (1.0k)	90% (2.8k)	76% (0.7k)	46% (1.9k)	79% (0.5k)	60% (1.1k)	67% (0.8k)	35% (6.9k)	37% (65)	14% (0.4k)
<i>Llama-3</i>	English	80% (6.0k)	72% (8.4k)	89% (2.4k)	35% (5.6k)	56% (61)	4% (0.9k)	76% (0.9k)	16% (6.4k)	20% (41)	0% (12.8k)

Table 2: Accuracy of Patchscopes interpretations for compositional input embeddings across languages. Numbers in parentheses indicate sample sizes. "—" indicates cases where no suitable *base-inflection* pairs found in the vocabulary or where there are no UniMorph entries for that category. For detokenization results, see Appendix B.3.

we further analyze in Appendix B.1. Overall, our results show that LLMs can naturally interpret compositional word embeddings across languages.

Analysis of transformation failures. Across languages and models, we observe a consistent gap between inflectional transformations (often resolved) and derivational transformations (rarely resolved). To characterize these failures, we analyze whether the number of in-vocabulary exemplar pairs used to estimate each transformation vector (Eq. 3) helps explain transformation failures. Qualitatively, we find that the number of exemplars mainly matters for *generalization*: transformations estimated from many pairs are much more likely to resolve multi-token surface forms, while in-vocabulary success is comparatively insensitive to exemplar count once a usable signal is available (see Appendix B.2).

Vocabulary size effects. We analyze the extent to which morphology is encoded *linearly and compositionally* as vocabulary size grows. We repeat the English Patchscopes experiment across models with various vocabulary sizes. Surprisingly, we find an inverse relationship between English whole-word vocabulary size and the linearity of morphological transformations; this suggests that larger vocabularies allow models to represent inflected forms of the same type as individual lexical units, rather than through a shared linear translation of their base forms (see Appendix B.1 for details).

6 Compositional Language Modeling

We have shown that *transformation* vectors capture meaningful operations in the input space of LLMs, and that these can be successfully composed with base word embeddings. We next investigate whether models can use compositional vocabularies effectively in end-to-end language modeling.

6.1 Experimental Setting and Implementation

Given a model’s vocabulary decomposition map (§4), we apply our compositional vocabulary framework and restructure the input and output embedding matrices. We replace the model’s input embedding of any surface form w with compositions of the corresponding *base form* and *transformation* embeddings (Eq. 1). For next-token prediction, we compute logits through summation of *base form* and *transformation* logits (Eq. 2). Importantly, any word not in the decomposition map maintains its original embedding and unembedding throughout training and inference, without modifications.

Fine-tuning the transformation vectors After initialization (Eq. 3), we train the *transformation* vectors jointly within the model: we treat the *transformation* embeddings and unembeddings matrices E_t and U_t as trainable weights (introducing fewer than 0.001% additional parameters), and freeze all other model parameters, including the embeddings and unembeddings of *base forms*. We use knowledge distillation loss (Hinton et al., 2015) to fine-tune the *transformation* vectors using two-stage distillation: We first freeze the output unembeddings and only train the *input transformations*, using the predictions of the original, unmodified model as targets. Next, we freeze the input embeddings and only train the *output transformations*, this time using the (frozen) model resulting from the first stage as the distillation target—ignoring all words $w \notin \mathcal{V}_{\text{orig}}$ in the loss. In both stages, we train on a fixed, small sample of the FineWeb-Edu corpus (Penedo et al., 2024).¹⁰ See Appendix C.

Lightweight LoRA adaptation. To allow lightweight adaptation to the reshaped output vocabulary, we add and train LoRA adapters to the final $k = 8$ model layers, keeping all other internal layers frozen. We use LoRA $r = \alpha = 256$.

¹⁰We use a sequence length of 256 and train on $\sim 5\text{M}$ tokens.

Filtering the decomposition map Our results in §5 indicate some out-of-vocabulary surface forms fail to be interpreted by the model as their intended word when given as compositions. We therefore filter out surface words with failed detokenization from the decomposition map, and fall back to using their original tokenization and embeddings in both input and output. We also exclude all derivational transformations due to their weak resolution rates. See analysis in Appendix B.4.

Downstream tasks We evaluate our compositional vocabulary models on a diverse suite of standard benchmarks. As a baseline, we compare performance to the original, unmodified models. For English, the benchmarks cover language understanding, knowledge, and commonsense: MMLU (Hendrycks et al., 2021), ARC (Clark et al., 2018), HellaSwag (Zellers et al., 2019), Winogrande (Sakaguchi et al., 2021), TriviaQA (Joshi et al., 2017), SQuAD (Rajpurkar et al., 2016), BoolQ (Clark et al., 2019), PIQA (Bisk et al., 2020) and COPA (Kavumba et al., 2019). We use the more efficient TinyBenchmark subsets when available (Maia Polo et al., 2024). For other languages, we use XNLI (Conneau et al., 2018), XQuAD (Artetxe et al., 2020) and Global MMLU (Singh et al., 2025). See Appendix E.

6.2 Results

We report our results for English on Llama-3-8B in Table 3,¹¹ and results for other languages in Table 4. Our compositional language modeling approach results in comparable average performance to the baseline models across languages, indicating that LLMs can effectively leverage compositional vocabularies with only lightweight adaptation.

We next inspect reductions in vocabulary size after applying our framework. Our approach removes roughly 10k surface form tokens from Llama3 and OLMo2 each, and 7.8k from Qwen2.5.¹² While seemingly small, recent work has shown that adding as few as several hundred dedicated tokens to the vocabulary can greatly improve tokenization efficiency and model performance on a language or expert domain (Ahia et al., 2023; Liu et al., 2024a; Nakash et al., 2025), emphasizing the importance of an efficient allocation of the vocabulary. These freed-up slots can be reallocated to new tokens

¹¹Results on other English models are found in Appendix A.

¹²Reductions in the multilingual models are more minimal due to the smaller initial token vocabularies in those languages, and vary from 0.6k to 3k, depending on language.

Category	Task	Baseline	End-to-end	Δ
Knowledge	TinyMMLU (Acc.)	53.0	53.8	+0.8
	TinyARC (Acc.)	46.1	45.6	-0.5
Reading Comprehension	BoolQ (Acc.)	83.2	83.2	+0.1
	TriviaQA (EM)	66.5	63.3	-3.3
	SQuAD (EM)	22.1	20.0	-2.1
Commonsense	TinyHellaswag (Acc.)	61.5	64.4	+2.9
	TinyWinogrande (Acc.)	60.3	58.9	-1.4
	PIQA (Acc.)	80.4	79.1	-1.3
	COPA (Acc.)	93.0	92.0	-1.0
Average		62.9	62.3	-0.6

Table 3: Downstream performance of English compositional-vocabulary models (*End-to-end*) and their original, unmodified version (*Baseline*) for Llama-3.1-8B. Our framework performs on-par with the baseline model, despite extensive changes to the model’s input and output representation mechanisms—highlighting the intrinsic ability of LLMs to process and predict word representations compositionally.

		XNLI Δ	XQuAD Δ	GMMLU Δ
<i>ALLaM</i>	Arabic	44.1 -0.3	42.7 -3.2	59.9 +0.2
<i>EuroLLM</i>	German	46.5 +0.6	51.3 -1.6	54.6 -0.7
	Russian	40.1 -4.5	37.4 -3.6	54.4 -0.3
	Spanish	43.3 -0.6	48.3 -4.1	55.2 -0.9

Table 4: Multilingual downstream performance of compositional vocabulary models, along with absolute performance difference from the baseline model (Δ).

via post-hoc expansion methods (Han et al., 2025; Minixhofer et al., 2024). Finally, our method has a marginal effect on decoding speed—only a 0.8% reduction compared to standard prediction (see Appendix B.5). We next show that vocabularies can also be built compositionally from the outset, with even greater vocabulary allocation efficiency, by pretraining a model with a compositional vocabulary from scratch (§7).

7 Compositional Vocabulary Pretraining

To validate that compositional vocabularies can also serve as a *design choice* for training new language models, we reshape English and Spanish BPE vocabularies into *base forms* and *transformations* and pretrain small baseline and compositional models from scratch. For English, we reshape the 50k-token GPT-2 tokenizer (Radford et al., 2019), while restricting the compositional model to predict exactly the same surface-form vocabulary as the BPE baseline (i.e., we do not extend to out-of-vocabulary words). For Spanish, we train a 32k-token BPE vocabulary¹³ and then reshape it, this

¹³We train the Spanish tokenizer on 10B bytes from the Spanish subset of FineWeb-2 (Penedo et al., 2025)

time allowing the compositional model to generate out-of-vocabulary surface forms via compositions.

We pretrain nanoGPT-124M models (Jordan et al., 2024) for 1B tokens in each language,¹⁴ comparing a baseline model against an otherwise-identical compositional model. In contrast to our post-hoc setup, the compositional model predicts directly in a factorized space (a base-form distribution and *transformation* predictions). We further include a space-prefix *transformation* (e.g., “_walking” vs. “walking”).¹⁵ We evaluate on a held-out validation set and report next-token accuracy for English (since both models share the same surface-form vocabulary). For Spanish, as the models have different vocabularies, we report Bits-Per-Bytes (BPB).¹⁶ We further examine tokenization efficiency in Spanish using average bytes per token (higher is better).

Reshaping the GPT-2 tokenizer yields a 41.6% reduction in the number of vocabulary entries used to represent the original BPE vocabulary (removing 19.8k tokens) while preserving exact coverage. Despite this reduction, the compositional model achieves comparable next-token accuracy on held-out validation of 38.10 (vs. 38.27). For Spanish, we observe a similar token vocabulary reduction of 41.8% with competitive BPB compared to the baseline (1.110 vs. 0.996). Allowing compositional generation improves tokenization efficiency, increasing average bytes per token (4.773→4.920) on held-out Spanish text. Finally, we measure training throughput and find that the compositional model trains faster due to the smaller vocabulary size, improving throughput by 11.5% (471k→525k tokens/sec).¹⁷

Together, these results show that compositional vocabularies can be trained from scratch effectively, offering compact vocabularies, improved tokenization efficiency, and practical training-speed benefits for future language models.

8 Related Work

Incorporating morphology into representations

A longstanding goal in NLP has been to integrate morphological knowledge into models. Early work on Transformer language models explored inject-

ing linguistic features post-hoc (Hofmann et al., 2021; Gan et al., 2022) or during pretraining (Park et al., 2021; Cui et al., 2022; Matthews et al., 2018; Blevins and Zettlemoyer, 2019; Hofmann et al., 2020; Seker et al., 2022; Peng et al., 2019), but such approaches are absent in modern LLMs. Recent work examined word segmentation effects on performance (Marco and Fraser, 2024; Lerner and Yvon, 2025), including morphology-aware tokenization to better reflect word structure (Bauwens and Delobelle, 2024; Asgari et al., 2025). Rather than injecting linguistic structure, we leverage compositional representations already present in LLMs.

Vector arithmetic of word representations

Linear structure in word representations was first observed in Word2Vec (Mikolov et al., 2013a,b; Levy and Goldberg, 2014; Vylomova et al., 2015). Recent work found similar structures in LLMs across the unembedding layer (Park et al., 2024, 2025), residual stream (Merullo et al., 2023; Hendel et al., 2023; Todd et al., 2024), and in behavior-steering directions (Subramani et al., 2022; Hernandez et al., 2024). We further show that such structure is usable for end-to-end language modeling. Beyond morphology, *transformation* vectors could capture semantic relations (e.g., country–nationality; Gladkova et al., 2016) or tie word embeddings across languages (Schut et al., 2025).

Post-hoc vocabulary modification

Recent work has proposed methods to expand or modify token vocabulary by training new embeddings and fine-tuning internal model layers (Kim et al., 2024; Takase et al., 2024; Han et al., 2025; Minixhofer et al., 2024; Ben-Artzy and Schwartz, 2025; Dobler and de Melo, 2023). We avoid fine-tuning model weights, and represent new forms compositionally using existing linguistic knowledge.

9 Conclusion

We have shown that word representations in LLMs are inherently compositional, and leveraged this property to introduce compositional vocabularies. Our framework enables token vocabularies that are more compact in the vocabulary size needed for linguistic coverage, while being more expressive. Our results highlight that language models can naturally operate with compositional vocabularies and that with integration into pretraining future models, LLMs could cover more words, languages, and domains—without sacrificing performance.

¹⁴We use FineWeb (English) and FineWeb-2 (Spanish).
¹⁵Modern BPE vocabularies include prefix whitespace characters when merging tokens, creating many near-duplicates.
¹⁶BPB allows comparison across models with different vocabularies; it normalizes negative log-likelihood by the number of UTF-8 bytes in the evaluation text. See Appendix D.
¹⁷Throughput is measured when training on 4 L40S GPUs.

611 Limitations

612 Our framework employs external morphological
613 resources to define transformation pairs. While
614 this allows for clean experimental control, it lim-
615 its applicability to languages or domains lacking
616 annotated morphological data. In future work, we
617 hope to explore whether similar transformation vec-
618 tors can be induced directly from data, using unsu-
619 pervised learning or joint training objectives that
620 encourage compositionality. Models that learn to
621 discover structure, rather than rely on it, would
622 offer broader generalization and better alignment
623 with language acquisition in humans.

624 Our vocabulary reshaping approach assumes a
625 one-to-one decomposition of each surface form
626 into a base word and a set of transformation vec-
627 tors. While effective for many cases, this simpli-
628 fication does not account for certain words which
629 could result from several distinct morphological
630 processes. Still, these problems are also encoun-
631 tered with standard tokenization approaches, with
632 models learning to disambiguate such words into
633 their intended meanings.

634 References

635 Ahmed Abdelali, Kareem Darwish, Nadir Durrani, and
636 Hamdy Mubarak. 2016. [Farasa: A fast and furious
637 segmenter for Arabic](#). In *Proceedings of the 2016
638 Conference of the North American Chapter of the
639 Association for Computational Linguistics: Demon-
640 strations*, pages 11–16, San Diego, California. Asso-
641 ciation for Computational Linguistics.

642 Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien
643 Bubeck, Ronen Eldan, Suriya Gunasekar, Michael
644 Harrison, Russell J Hewett, Mojan Javaheripi, Piero
645 Kauffmann, and 1 others. 2024. Phi-4 technical re-
646 port. *arXiv preprint arXiv:2412.08905*.

647 Orevaoghene Ahia, Sachin Kumar, Hila Gonen, Jungo
648 Kasai, David Mortensen, Noah Smith, and Yulia
649 Tsvetkov. 2023. [Do all languages cost the same?
650 tokenization in the era of commercial language mod-
651 els](#). In *Proceedings of the 2023 Conference on Em-
652 pirical Methods in Natural Language Processing*,
653 pages 9904–9923, Singapore. Association for Com-
654 putational Linguistics.

655 Mehdi Ali, Michael Fromm, Klaudia Thellmann,
656 Richard Rutmann, Max Lübbering, Johannes Lev-
657 eling, Katrin Klug, Jan Ebert, Niclas Doll, Jasper
658 Buschhoff, Charvi Jain, Alexander Weber, Lena Ju-
659 rkschat, Hammam Abdelwahab, Chelsea John, Pedro
660 Ortiz Suarez, Malte Ostendorff, Samuel Weinbach,
661 Rafet Sifa, and 2 others. 2024. [Tokenizer choice](#)

[for LLM training: Negligible or crucial?](#) In *Find-
662 ings of the Association for Computational Linguis-
663 tics: NAACL 2024*, pages 3907–3924, Mexico City,
664 Mexico. Association for Computational Linguistics.
665

Mohamed Taher Alrefaie, Nour Eldin Morsy, and Nada
666 Samir. 2024. Exploring tokenization strategies and
667 vocabulary sizes for enhanced arabic language mod-
668 els. *arXiv preprint arXiv:2403.11130*.
669

Zaid Alyafeai, Maged S. Al-Shaibani, Mustafa Ghaleb,
670 and Irfan Ahmad. 2021. [Evaluating various tokeniz-
671 ers for arabic text classification](#). *Neural Processing
672 Letters*, 55:2911–2933.
673

Mikel Artetxe, Sebastian Ruder, and Dani Yogatama.
674 2020. [On the cross-lingual transferability of mono-
675 lingual representations](#). In *Proceedings of the 58th
676 Annual Meeting of the Association for Computational
677 Linguistics*, pages 4623–4637, Online. Association
678 for Computational Linguistics.
679

Viraat Aryabumi, John Dang, Dwarak Talupuru,
680 Saurabh Dash, David Cairuz, Hangyu Lin, Bharat
681 Venkitesh, Madeline Smith, Jon Ander Campos,
682 Yi Chern Tan, and 1 others. 2024. Aya 23: Open
683 weight releases to further multilingual progress.
684 *arXiv preprint arXiv:2405.15032*.
685

Ehsaneddin Asgari, Yassine El Kheir, and Mohammad
686 Ali Sadraei Javaheri. 2025. [MorphBPE: A morpho-
687 aware tokenizer bridging linguistic complexity for
688 efficient llm training across morphologies](#). *Preprint*,
689 arXiv:2502.00894.
690

M Saiful Bari, Yazeed Alnumay, Norah A. Alzahrani,
691 Nouf M. Alotaibi, Hisham Abdullah Alyahya, Sultan
692 AlRashed, Faisal Abdulrahman Mirza, Shaykhah Z.
693 Alsubaie, Hassan A. Alahmed, Ghadah Alabduljab-
694 bar, Raghad Alkhathran, Yousef Almushayqih, Ra-
695 neem Alnajim, Salman Alsubaihi, Maryam Al Man-
696 sour, Saad Amin Hassan, Dr. Majed Alrubaian, Ali
697 Alammari, Zaki Alawami, and 7 others. 2025. [AL-
698 Lam: Large language models for arabic and english](#).
699 In *The Thirteenth International Conference on Learn-
700 ing Representations*.
701

Khuyagbaatar Batsuren, Omer Goldman, Salam Khal-
702 ifa, Nizar Habash, Witold Kieraś, Gábor Bella,
703 Brian Leonard, Garrett Nicolai, Kyle Gorman, Yusti-
704 nus Ghanggo Ate, Maria Ryskina, Sabrina Mielke,
705 Elena Budianskaya, Charbel El-Khaissi, Tiago Pi-
706 mentel, Michael Gasser, William Abbott Lane, Mo-
707 hit Raj, Matt Coler, and 76 others. 2022. [UniMorph
708 4.0: Universal Morphology](#). In *Proceedings of the
709 Thirteenth Language Resources and Evaluation Con-
710 ference*, pages 840–855, Marseille, France. European
711 Language Resources Association.
712

Thomas Bauwens and Pieter Delobelle. 2024. [BPE-
713 knockout: Pruning pre-existing BPE tokenisers
714 with backwards-compatible morphological semi-
715 supervision](#). In *Proceedings of the 2024 Conference
716 of the North American Chapter of the Association for
717 Computational Linguistics: Human Language Tech-
718 nologies (Volume 1: Long Papers)*, pages 5810–5832,
719

720	Mexico City, Mexico. Association for Computational Linguistics.	
721		
722	Amit Ben-Artzy and Roy Schwartz. 2025. Spellm: Character-level multi-head decoding . <i>Preprint</i> , arXiv:2507.16323.	
723		
724		
725	Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. 2020. Piqa: Reasoning about physical commonsense in natural language. In <i>Thirty-Fourth AAAI Conference on Artificial Intelligence</i> .	
726		
727		
728		
729	Terra Blevins and Luke Zettlemoyer. 2019. Better character language modeling through morphology. <i>arXiv preprint arXiv:1906.01037</i> .	
730		
731		
732	Danqi Chen, Jason Bolton, and Christopher D. Manning. 2016. A thorough examination of the CNN/Daily Mail reading comprehension task . In <i>Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2358–2367, Berlin, Germany. Association for Computational Linguistics.	
733		
734		
735		
736		
737		
738		
739	Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. BoolQ: Exploring the surprising difficulty of natural yes/no questions . In <i>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)</i> , pages 2924–2936, Minneapolis, Minnesota. Association for Computational Linguistics.	
740		
741		
742		
743		
744		
745		
746		
747		
748	Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge . <i>ArXiv</i> , abs/1803.05457.	
749		
750		
751		
752		
753	Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina Williams, Samuel Bowman, Holger Schwenk, and Veselin Stoyanov. 2018. XNLI: Evaluating cross-lingual sentence representations . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 2475–2485, Brussels, Belgium. Association for Computational Linguistics.	
754		
755		
756		
757		
758		
759		
760		
761	Mathias Creutz and Krista Lagus. 2005. Unsupervised morpheme segmentation and morphology induction from text corpora using morfessor 1.0. In <i>Unsupervised morpheme segmentation and morphology induction from text corpora using Morfessor 1.0</i> . Helsinki University of Technology.	
762		
763		
764		
765		
766		
767	Mathias Creutz and Krista Lagus. 2007. Unsupervised models for morpheme segmentation and morphology learning. <i>ACM Transactions on Speech and Language Processing (TSLP)</i> , 4(1):1–34.	
768		
769		
770		
771	Yiming Cui, Wanxiang Che, Shijin Wang, and Ting Liu. 2022. Lert: A linguistically-motivated pre-trained language model. <i>arXiv preprint arXiv:2211.05344</i> .	
772		
773		
	Gautier Dagan, Gabriele Synnaeve, and Baptiste Rozière. 2024. Getting the most out of your tokenizer for pre-training and domain adaptation. <i>arXiv preprint arXiv:2402.01035</i> .	774 775 776 777
	John Dang, Shivalika Singh, Daniel D’souza, Arash Ahmadian, Alejandro Salamanca, Madeline Smith, Aidan Peppin, Sungjin Hong, Manoj Govindassamy, Terrence Zhao, and 1 others. 2024. Aya expand: Combining research breakthroughs for a new multi-lingual frontier. <i>arXiv preprint arXiv:2412.04261</i> .	778 779 780 781 782 783
	Konstantin Dobler and Gerard de Melo. 2023. FOCUS: Effective embedding initialization for monolingual specialization of multilingual models . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 13440–13454, Singapore. Association for Computational Linguistics.	784 785 786 787 788 789 790
	Sheridan Feucht, David Atkinson, Byron C Wallace, and David Bau. 2024. Token erasure as a footprint of implicit vocabulary items in LLMs . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 9727–9739, Miami, Florida, USA. Association for Computational Linguistics.	791 792 793 794 795 796 797
	Guobing Gan, Peng Zhang, Sunzhu Li, Xiuqing Lu, and Benyou Wang. 2022. MorphTE: Injecting morphology in tensorized embeddings. <i>Advances in Neural Information Processing Systems</i> , 35:33186–33200.	798 799 800 801
	Bar Gazit, Shaltiel Shmidman, Avi Shmidman, and Yuval Pinter. 2025. Splintering nonconcatenative languages for better tokenization . In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 22405–22417, Vienna, Austria. Association for Computational Linguistics.	802 803 804 805 806 807
	Asma Ghandeharioun, Avi Caciularu, Adam Pearce, Lucas Dixon, and Mor Geva. 2024. Patchscopes: A unifying framework for inspecting hidden representations of language models. <i>arXiv preprint arXiv:2401.06102</i> .	808 809 810 811 812
	Anna Gladkova, Aleksandr Drozd, and Satoshi Matsuo. 2016. Analogy-based detection of morphological and semantic relations with word embeddings: what works and what doesn’t. In <i>Proceedings of the NAACL Student Research Workshop</i> , pages 8–15.	813 814 815 816 817
	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, and 1 others. 2024. The llama 3 herd of models. <i>arXiv preprint arXiv:2407.21783</i> .	818 819 820 821 822
	Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii, and Dimitris Bertsimas. 2023. Finding neurons in a haystack: Case studies with sparse probing . <i>Transactions on Machine Learning Research</i> .	823 824 825 826 827
	HyoJung Han, Akiko Eriguchi, Haoran Xu, Hieu Hoang, Marine Carpuat, and Huda Khayrallah. 2025.	828 829

830	Adapters for altering LLM vocabularies: What languages benefit the most? In <i>The Thirteenth International Conference on Learning Representations</i> .	886
831		887
832		888
833	Roe Hendel, Mor Geva, and Amir Globerson. 2023.	
834	In-context learning creates task vectors . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 9318–9333, Singapore. Association for Computational Linguistics.	890
835		891
836		892
837		893
838		894
839	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. <i>Proceedings of the International Conference on Learning Representations (ICLR)</i> .	895
840		896
841		897
842		898
843	Evan Hernandez, Belinda Z. Li, and Jacob Andreas. 2024. Inspecting and editing knowledge representations in language models . In <i>First Conference on Language Modeling</i> .	899
844		900
845		901
846		
847	Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. <i>arXiv preprint arXiv:1503.02531</i> .	902
848		903
849		904
850		905
851	Valentin Hofmann, Janet Pierrehumbert, and Hinrich Schütze. 2020. DagoBERT: Generating derivational morphology with a pretrained language model . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 3848–3861, Online. Association for Computational Linguistics.	906
852		907
853		908
854		909
855		910
856		911
857	Valentin Hofmann, Janet Pierrehumbert, and Hinrich Schütze. 2021. Superbizarre is not superb: Derivational morphology improves BERT’s interpretation of complex words . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 3594–3608, Online. Association for Computational Linguistics.	912
858		913
859		914
860		915
861		916
862		
863		
864		
865		
866	Valentin Hofmann, Hinrich Schuetze, and Janet Pierrehumbert. 2022. An embarrassingly simple method to mitigate undesirable properties of pretrained language model tokenizers . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 385–393, Dublin, Ireland. Association for Computational Linguistics.	917
867		918
868		919
869		920
870		921
871		922
872		923
873		
874	Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text de-generation . In <i>International Conference on Learning Representations</i> .	924
875		925
876		926
877		927
878		928
879		929
880		
881		
882		
883		
884		
885		
	Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b . <i>Preprint, arXiv:2310.06825</i> .	886
		887
		888
		889
	Keller Jordan, Jeremy Bernstein, Brendan Rappazzo, @fernbear.bsky.social, Boza Vlado, You Jiacheng, Franz Cesista, Braden Koszarsky, and @Grad62304977. 2024. modded-nanogpt: Speedrunning the nanogpt baseline .	890
		891
		892
		893
		894
	Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension . In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1601–1611, Vancouver, Canada. Association for Computational Linguistics.	895
		896
		897
		898
		899
		900
		901
	Guy Kaplan, Matanel Oren, Yuval Reif, and Roy Schwartz. 2025. From tokens to words: On the inner lexicon of LLMs . In <i>The Thirteenth International Conference on Learning Representations</i> .	902
		903
		904
		905
	Pride Kavumba, Naoya Inoue, Benjamin Heinzerling, Keshav Singh, Paul Reisert, and Kentaro Inui. 2019. When choosing plausible alternatives, clever hans can be clever . In <i>Proceedings of the First Workshop on Commonsense Inference in Natural Language Processing</i> , pages 33–42, Hong Kong, China. Association for Computational Linguistics.	906
		907
		908
		909
		910
		911
		912
	Seungduk Kim, Seungtaek Choi, and Myeongho Jeong. 2024. Efficient and effective vocabulary expansion towards multilingual large language models. <i>arXiv preprint arXiv:2402.14714</i> .	913
		914
		915
		916
	Stav Klein and Reut Tsarfaty. 2020. Getting the ##life out of living: How adequate are word-pieces for modelling complex morphology? In <i>Proceedings of the 17th SIGMORPHON Workshop on Computational Research in Phonetics, Phonology, and Morphology</i> , pages 204–209, Online. Association for Computational Linguistics.	917
		918
		919
		920
		921
		922
		923
	Paul Lerner and François Yvon. 2025. Unlike “likely”, “unlikely” is unlikely: BPE-based segmentation hurts morphological derivations in LLMs . In <i>Proceedings of the 31st International Conference on Computational Linguistics</i> , pages 5181–5190, Abu Dhabi, UAE. Association for Computational Linguistics.	924
		925
		926
		927
		928
		929
	Omer Levy and Yoav Goldberg. 2014. Linguistic regularities in sparse and explicit word representations . In <i>Conference on Computational Natural Language Learning</i> .	930
		931
		932
		933
	Tomasz Limisiewicz, Jivr’i Balhar, and David Marevcek. 2023. Tokenization impacts multilingual language modeling: Assessing vocabulary allocation and overlap across languages . In <i>Annual Meeting of the Association for Computational Linguistics</i> .	934
		935
		936
		937
		938
	Alisa Liu, Jonathan Hayase, Valentin Hofmann, Sewoong Oh, Noah A. Smith, and Yejin Choi. 2025. SuperBPE: Space travel for language models . <i>ArXiv, abs/2503.13423</i> .	939
		940
		941
		942

943	Chengyuan Liu, Shihang Wang, Lizhi Qing, Kun Kuang, Yangyang Kang, Changlong Sun, and Fei Wu. 2024a. Gold panning in vocabulary: An adaptive method for vocabulary expansion of domain-specific LLMs . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 7442–7459, Miami, Florida, USA. Association for Computational Linguistics.	Tomas Mikolov, Ilya Sutskever, Kai Chen, Gregory S. Corrado, and Jeffrey Dean. 2013a. Distributed representations of words and phrases and their compositionality . In <i>Neural Information Processing Systems</i> .	998 999 1000 1001
944		Tomas Mikolov, Wen-tau Yih, and Geoffrey Zweig. 2013b. Linguistic regularities in continuous space word representations . In <i>Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 746–751, Atlanta, Georgia. Association for Computational Linguistics.	1002 1003 1004 1005 1006 1007 1008
945		Benjamin Minixhofer, Edoardo Ponti, and Ivan Vulić. 2024. Zero-shot tokenizer transfer . In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	1009 1010 1011 1012
946		Itay Nakash, Nitay Calderon, Eyal Ben-David, Elad Hoffer, and Roi Reichart. 2025. Adaptivocab: Enhancing LLM efficiency in focused domains through lightweight vocabulary adaptation . In <i>Second Conference on Language Modeling</i> .	1013 1014 1015 1016 1017
947		Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, Nathan Lambert, Dustin Schwenk, Oyvind Tafjord, Taira Anderson, David Atkinson, Faeze Brahman, Christopher Clark, Pradeep Dasigi, Nouha Dziri, and 21 others. 2024. 2 olmo 2 furious . <i>ArXiv</i> , abs/2501.00656.	1018 1019 1020 1021 1022 1023 1024
948		OpenAI. 2024. tiktoken: A fast BPE tokeniser for use with openai’s models .	1025 1026
949		Hyunji Hayley Park, Katherine J. Zhang, Coleman Haley, Kenneth Steimel, Han Liu, and Lane Schwartz. 2021. Morphology matters: A multilingual language modeling analysis . <i>Transactions of the Association for Computational Linguistics</i> , 9:261–276.	1027 1028 1029 1030 1031
950		Kiho Park, Yo Joong Choe, Yibo Jiang, and Victor Veitch. 2025. The geometry of categorical and hierarchical concepts in large language models . In <i>The Thirteenth International Conference on Learning Representations</i> .	1032 1033 1034 1035 1036
951	Yihong Liu, Peiqin Lin, Mingyang Wang, and Hinrich Schuetze. 2024b. OFA: A framework of initializing unseen subword embeddings for efficient large-scale multilingual continued pretraining . In <i>Findings of the Association for Computational Linguistics: NAACL 2024</i> , pages 1067–1097, Mexico City, Mexico. Association for Computational Linguistics.	Kiho Park, Yo Joong Choe, and Victor Veitch. 2024. The linear representation hypothesis and the geometry of large language models . In <i>Proceedings of the 41st International Conference on Machine Learning</i> , volume 235 of <i>Proceedings of Machine Learning Research</i> , pages 39643–39666. PMLR.	1037 1038 1039 1040 1041 1042
952		Guilherme Penedo, Hynek Kydlíček, Anton Lozhkov, Margaret Mitchell, Colin A Raffel, Leandro Von Werra, Thomas Wolf, and 1 others. 2024. The fineweb datasets: Decanting the web for the finest text data at scale . <i>Advances in Neural Information Processing Systems</i> , 37:30811–30849.	1043 1044 1045 1046 1047 1048
953		Guilherme Penedo, Hynek Kydlíček, Vinko Sabolčec, Bettina Messmer, Negar Foroutan, Amir Hossein Kargaran, Colin Raffel, Martin Jaggi, Leandro Von Werra, and Thomas Wolf. 2025. Fineweb2: One pipeline to scale them all — adapting pre-training	1049 1050 1051 1052 1053
954			
955			
956			
957			
958	Manuel Mager, Arturo Oncevay, Elisabeth Mager, Katharina Kann, and Thang Vu. 2022. BPE vs. morphological segmentation: A case study on machine translation of four polysynthetic languages . In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 961–971, Dublin, Ireland. Association for Computational Linguistics.		
959			
960			
961			
962			
963			
964			
965	Felipe Maia Polo, Lucas Weber, Leshem Choshen, Yuekai Sun, Gongjun Xu, and Mikhail Yurochkin. 2024. tinyBenchmarks: evaluating LLMs with fewer examples . In <i>Proceedings of the 41st International Conference on Machine Learning</i> , volume 235 of <i>Proceedings of Machine Learning Research</i> , pages 34303–34326. PMLR.		
966			
967			
968			
969			
970			
971			
972	Marion Di Marco and Alexander Fraser. 2024. Subword segmentation in LLMs: Looking at inflection and consistency . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 12050–12060, Miami, Florida, USA. Association for Computational Linguistics.		
973			
974			
975			
976			
977			
978	Samuel Marks and Max Tegmark. 2024. The geometry of truth: Emergent linear structure in large language model representations of true/false datasets . In <i>First Conference on Language Modeling</i> .		
979			
980			
981			
982	Pedro Henrique Martins, Patrick Fernandes, João Alves, Nuno M Guerreiro, Ricardo Rei, Duarte M Alves, José Pombal, Amin Farajian, Manuel Faysse, Mateusz Klimaszewski, and 1 others. 2025. Euollm: Multilingual language models for europe . <i>Procedia Computer Science</i> , 255:53–62.		
983			
984			
985			
986			
987			
988	Austin Matthews, Graham Neubig, and Chris Dyer. 2018. Using morphological knowledge in open-vocabulary neural language models . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 1435–1445.		
989			
990			
991			
992			
993			
994			
995	Jack Merullo, Carsten Eickhoff, and Ellie Pavlick. 2023. Language models implement simple word2vec-style vector arithmetic . <i>arXiv preprint arXiv:2305.16130</i> .		
996			
997			

1054	data processing to every language . In <i>Second Conference on Language Modeling</i> .	<i>the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 18761–18799, Vienna, Austria. Association for Computational Linguistics.	1110 1111 1112
1056	Hao Peng, Roy Schwartz, and Noah A. Smith. 2019. PaLM: A hybrid parser and language model . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3644–3651, Hong Kong, China. Association for Computational Linguistics.	Peter Smit, Sami Virpioja, Stig-Arne Grönroos, and Mikko Kurimo. 2014. Morfessor 2.0: Toolkit for statistical morphological segmentation . In <i>Proceedings of the Demonstrations at the 14th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 21–24, Gothenburg, Sweden. Association for Computational Linguistics.	1113 1114 1115 1116 1117 1118 1119
1064	Aleksandar Petrov, Emanuele La Malfa, Philip Torr, and Adel Bibi. 2023. Language model tokenizers introduce unfairness between languages. <i>Advances in neural information processing systems</i> , 36:36963–36990.	Nishant Subramani, Nivedita Suresh, and Matthew Peters. 2022. Extracting latent steering vectors from pretrained language models . In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 566–581, Dublin, Ireland. Association for Computational Linguistics.	1120 1121 1122 1123 1124 1125
1069	Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.	Sho Takase, Ryokan Ri, Shun Kiyono, and Takuya Kato. 2024. Large vocabulary size improves large language models. <i>arXiv preprint arXiv:2406.16508</i> .	1126 1127 1128
1072	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 2383–2392, Austin, Texas. Association for Computational Linguistics.	Chaofan Tao, Qian Liu, Longxu Dou, Niklas Muenighoff, Zhongwei Wan, Ping Luo, Min Lin, and Ngai Wong. 2024. Scaling laws with vocabulary: Larger models deserve larger vocabularies . In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	1129 1130 1131 1132 1133 1134
1078	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. Winogrande: An adversarial winograd schema challenge at scale. <i>Communications of the ACM</i> , 64(9):99–106.	Falcon-LLM Team. 2024. The falcon 3 family of open models .	1135 1136
1082	Lisa Schut, Yarin Gal, and Sebastian Farquhar. 2025. Do multilingual llms think in english? <i>arXiv preprint arXiv:2502.15603</i> .	Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, and 1 others. 2024a. Gemma: Open models based on gemini research and technology. <i>arXiv preprint arXiv:2403.08295</i> .	1137 1138 1139 1140 1141 1142
1085	Amit Seker, Elron Bandel, Dan Bareket, Idan Brusilovsky, Refael Greenfeld, and Reut Tsarfaty. 2022. AlephBERT: Language model pre-training and evaluation from sub-word to sentence level . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 46–56, Dublin, Ireland. Association for Computational Linguistics.	Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, and 1 others. 2024b. Gemma 2: Improving open language models at a practical size. <i>arXiv preprint arXiv:2408.00118</i> .	1143 1144 1145 1146 1147 1148
1093	Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units . In <i>Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.	Eric Todd, Millicent Li, Arnab Sen Sharma, Aaron Mueller, Byron C Wallace, and David Bau. 2024. Function vectors in large language models . In <i>The Twelfth International Conference on Learning Representations</i> .	1149 1150 1151 1152 1153
1100	Shivalika Singh, Angelika Romanou, Clémentine Fourrier, David Ifeoluwa Adelani, Jian Gang Ngui, Daniel Vila-Suero, Peerat Limkonchotiwat, Kelly Marchisio, Wei Qi Leong, Yosephine Susanto, Raymond Ng, Shayne Longpre, Sebastian Ruder, Wei-Yin Ko, Antoine Bosselut, Alice Oh, Andre Martins, Leshem Choshen, Daphne Ippolito, and 4 others. 2025. Global MMLU: Understanding and addressing cultural and linguistic biases in multilingual evaluation . In <i>Proceedings of the 63rd Annual Meeting of</i>	Cagri Toraman, Eyup Halit Yilmaz, Furkan Şahinuç, and Oguzhan Ozcelik. 2022. Impact of tokenization on language models: An analysis for turkish . <i>ACM Transactions on Asian and Low-Resource Language Information Processing</i> , 22:1 – 21.	1154 1155 1156 1157 1158
1109		Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutit Bhosale, and 1 others. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .	1159 1160 1161 1162 1163 1164

- 1165 Reut Tsarfaty, Shoval Sadde, Stav Klein, and Amit
1166 Seker. 2019. [What’s wrong with Hebrew NLP?](#)
1167 [and how to make it right](#). In *Proceedings of the*
1168 *2019 Conference on Empirical Methods in Natural*
1169 *Language Processing and the 9th International*
1170 *Joint Conference on Natural Language Processing*
1171 *(EMNLP-IJCNLP): System Demonstrations*, pages
1172 259–264, Hong Kong, China. Association for Com-
1173 putational Linguistics.
- 1174 Ekaterina Vylomova, Laura Rimell, Trevor Cohn, and
1175 Timothy Baldwin. 2015. [Take and took, gaggle and](#)
1176 [goose, book and read: Evaluating the utility of vec-](#)
1177 [tor differences for lexical relation learning](#). *ArXiv*,
1178 [abs/1509.01692](#).
- 1179 Erik Wijmans, Brody Huval, Alexander Hertzberg,
1180 Vladlen Koltun, and Philipp Kraehenbuehl. 2025.
1181 [Cut your losses in large-vocabulary language mod-](#)
1182 [els](#). In *The Thirteenth International Conference on*
1183 *Learning Representations*.
- 1184 An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui,
1185 Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu,
1186 Fei Huang, Haoran Wei, and 1 others. 2024. Qwen2.
1187 5 technical report. *arXiv preprint arXiv:2412.15115*.
- 1188 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali
1189 Farhadi, and Yejin Choi. 2019. [HellaSwag: Can a ma-](#)
1190 [chine really finish your sentence?](#) In *Proceedings of*
1191 *the 57th Annual Meeting of the Association for Com-*
1192 *putational Linguistics*, pages 4791–4800, Florence,
1193 Italy. Association for Computational Linguistics.

A Supplementary Results

For the results on other models for the experiment in §5, see Table 5, Table 6. For downstream results (§6), see Table 8 and Table 9.

B Additional Analysis

B.1 Morphology in Embedding Space Scales Inversely with Vocabulary Size

Having established that models implicitly learn to represent words compositionally, and with recent calls to scale vocabularies even further, a natural question emerges: how does vocabulary size affect the way models encode linguistic structure?

To study this question, we evaluate the extent of compositional word representations across models with varying vocabulary sizes. For each model, we decompose its vocabulary and measure the average Patchscopes interpretation accuracy for each *transformation* vector we extract (see §5). We also separate models by their embedding architecture (*untied* vs. *tied*). We track each model’s English vocabulary size (the subset of tokens present in English UniMorph), and plot the results in order of increasing vocabulary size.

For regular embedding models (top panel in Figure 3), we use Llama2-7B (Touvron et al., 2023), Mistral-7B (Jiang et al., 2023), OLMo2-7B (OLMo et al., 2024), Phi4-14B (Abdin et al., 2024), Llama3-8B (Grattafiori et al., 2024), Qwen2.5-7B (Yang et al., 2024), and Falcon3-7B (Team, 2024). For tied input-output embedding models (bottom panel), where input and output embeddings share parameters, we analyze Llama3-3B (Grattafiori et al., 2024), Qwen2.5-3B (Yang et al., 2024), Phi4-Mini-Instruct-4B (Abdin et al., 2024), Aya-Expanse-8B (Dang et al., 2024), Gemma2-2B and Gemma2-9B (Team et al., 2024b). These models span English vocabulary sizes from 8k to 44k tokens and total vocabulary sizes from 32k to 256k tokens, providing coverage across different scales of vocabulary design.

Our results (Figure 3) reveal a general inverse relationship: models with compact English vocabularies (8-10k words, e.g., Llama2, Mistral) tend to encode morphology through consistent vector offsets that generalize across words. In contrast, large-vocabulary models (~40k words, e.g., Falcon3, Gemma2-9B) tend to represent inflected forms of the same type as individual lexical units, rather than through a shared linear translation of their base

forms, with weight tying further amplifying this trend. Overall, these results suggest that vocabulary scaling trades morphological compositionality in embedding space for lexical memorization.¹⁸

B.2 Exemplar Count Predicts Transformation OOV Generalization

We quantify whether the number of in-vocabulary exemplar pairs available to estimate each transformation vector (Eq. 3) predicts when the composed embedding is interpreted as the intended surface form. Concretely, for each transformation t , we use the number of single-token, in-vocabulary (IV) base/surface pairs as a proxy for the effective exemplar set size, and compute Spearman correlations with additive success measured separately on IV targets and on out-of-vocabulary (multi-token) targets.

For Llama-3.1-8B, across individual transformations ($n = 24$ with at least one IV exemplar), IV additive success is essentially independent of exemplar count (Spearman’s $\rho = 0.04$, $p = 0.87$). In contrast, the same IV exemplar count strongly predicts out-of-vocabulary performance ($n = 23$ transformations with both IV exemplars and OOV test cases; $\rho = 0.77$, $p = 1.5 \times 10^{-5}$). This dissociation suggests that exemplar richness is not the bottleneck for resolving IV targets (which largely saturate once a direction is available), but it is a prerequisite for the transformation vector to generalize to multi-token, out-of-vocabulary (OOV) surface forms.

When restricting the correlation to transformations within each coarse class, we observe the same qualitative pattern but with weaker statistical power due to small n . Among inflectional transformations ($n = 8$), IV success correlates moderately with the number of IV exemplars ($\rho = 0.50$, $p = 0.21$), whereas among derivational transformations ($n = 14$) IV success trends negative ($\rho = -0.44$, $p = 0.12$), consistent with the observation that derivations remain difficult even when many IV exemplars exist. Within each class, OOV success remains positively correlated with IV exemplar count (inflection: $\rho = 0.38$, $p = 0.35$; derivation: $\rho = 0.38$, $p = 0.18$), reinforcing that generalization to OOV targets depends on having sufficiently many IV pairs, but that this dependence does not by itself close the gap.

¹⁸Importantly, this does not imply that large-vocabulary models lack morphological knowledge, only that they rely less on linear encoding of morphology in embedding space.

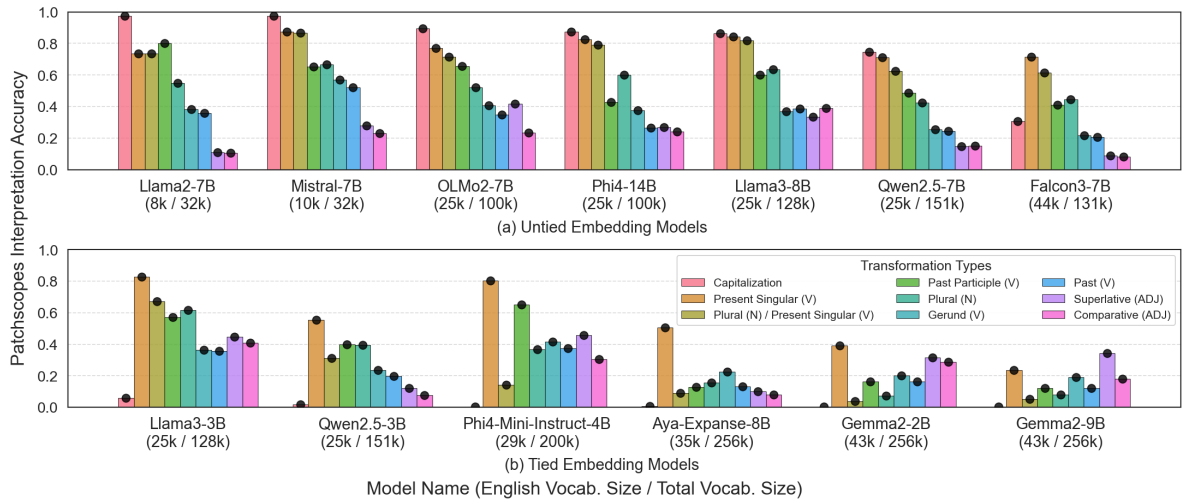


Figure 3: **Linear representation of morphology in embeddings weakens as vocabulary size increases.** Accuracy of Patchscopes interpretations of compositional word representations across models, in order of increasing English vocabulary size (English tokens present in UniMorph), separated by embedding architecture. Scaling vocabulary size leads models to represent inflection as individual lexical units, rather than through consistent vector offsets.

B.3 Multilingual Patchscopes Results

For the results on multilingual patchscopes interpretations of compositional input embeddings after detokenization, see Table 7.

B.4 Filtering the Vocabulary Decomposition for Failed Surface Forms

In §5 we have seen that, even though the compositional embeddings work well for many in- and out-of-vocabulary words, there are also failure cases where we are cannot be certain that the model interprets the compositional representation correctly. Intuitively, this means that using these representations in end-to-end language modeling might hurt model performance; Indeed, when we remove the surface forms corresponding to these failures from the decomposition map (and after fine-tuning the transformation vectors as usual), we observe an average 1.6 points improvement across downstream benchmarks, compared to no filtering. We note that for input-only restructuring, we observe no effect, likely because the model has more error-correction opportunities across its layers. We therefore apply this filtering in experiments in §6.

B.5 Decoding Speed Analysis

Our compositional language modeling approach introduces some additional complexity into next-token prediction: to compute token scores over the full, extended vocabulary, we map and sum up logit contributions from the *base form* and *transformation* vocabularies (Eq. 2). To validate that this does

not introduce meaningful overhead, we let both the baseline and compositional Llama-3-8B models generate texts in response to prompts from the CNN-DailyMail dataset (Chen et al., 2016), and measure the average number of tokens generated per second.¹⁹ Our approach introduces only a 0.8% drop in decoding speed (39.6 vs. 39.9 tokens/sec).

Still, since our compositional next-token prediction approach occurs in two-stages—first deciding on likely candidates for *base forms* and *transformations*—it naturally allows for optimizations like pruning base-form candidates before computing the logits over the full vocabulary (Holtzman et al., 2020), which could further decrease runtime.

C Experimental Details

For fine-tuning, we use a learning rate of $5e-5$, a warmup ratio of 0.03, a weight decay of 0.0, and a sequence length of $m = 256$. We train on $20k$ examples for 1 epoch. We run fine-tuning and inference on A10 GPUs, with fine-tuning taking roughly 30 minutes, and inference taking up to 1 hour.

D Pretraining Details

This appendix provides implementation and evaluation details for the pretraining experiment in §7.

Compositional tokenizer and coverage: We start from a 50k-token GPT-2 tokenizer and re-

¹⁹We use 50 random prompts and let models generate up to 256 tokens, on an L40S GPU.

1349 move any token that can be expressed as a base
1350 form plus *transformations*, including a whitespace-
1351 prefix *transformation* to capture pairs like “ walk-
1352 ing” vs. “walking”. Unlike our post-hoc setup, we
1353 do not filter out compositions based on Patchscopes
1354 interpretation failures (Appendix B.4).

Factorized prediction and base-conditioned transformations: The compositional model predicts a distribution over *base forms*, then predicts *transformations* conditioned on the chosen (or teacher-forced) base form by additionally providing its unembedding vector to the *transformation* head. At inference time, we decode by selecting a base form and then composing it with the predicted *transformations*.

Bits-Per-Byte (BPB): For both baseline and compositional models, we report BPB, computed as the average negative log-likelihood divided by the number of UTF-8 bytes in the evaluation text (lower is better). For the compositional model, we use the teacher-forced joint likelihood under the factorized distribution.

Tokenization efficiency: To test a morphologically richer setting, we train a Spanish 32k-token BPE tokenizer on 10B bytes from the Spanish subset of FineWeb-2 and apply our compositional reshaping, allowing compositions to represent out-of-vocabulary words up to three tokens long. On held-out Spanish text, the reshaped vocabulary improves tokenization efficiency, reducing token fertility (1.3028→1.265) and increasing average bytes per token (4.773→4.92).

1381 E Downstream Evaluation

1382 We include 5 in-context examples for every task.
1383 For each dataset, we use 5,000 examples (or the
1384 maximum available as some datasets have fewer
1385 available samples).

1386 **ARC** features 4-option multiple-choice science
1387 questions from grades 3 through 9. It has two sub-
1388 sets: ARC-Easy, focused on basic science knowl-
1389 edge, and ARC-Challenge, which involves more
1390 complex, procedural reasoning (Clark et al., 2018).

1391 **BoolQ** comprises naturally occurring yes/no
1392 questions accompanied by passages that support
1393 the answer (Clark et al., 2019).

1394 **COPA** offers binary multiple-choice questions
1395 centered around causal and consequential reason-
1396 ing (Kavumba et al., 2019).

HellaSwag includes 4-option multiple-choice
1397 questions where the task is to select the most plau-
1398 sible continuation of a given context (Zellers et al.,
1399 2019). 1400

MMLU presents 4-option multiple-choice ques-
1401 tions across 57 subject areas, testing both factual
1402 knowledge and reasoning skills (Hendrycks et al.,
1403 2021). 1404

PIQA provides multiple-choice questions de-
1405 signed to evaluate physical commonsense under-
1406 standing (Bisk et al., 2020). 1407

SQuAD pairs reading passages with related ques-
1408 tions, where the correct answer is always a text span
1409 from the passage itself (Rajpurkar et al., 2016). 1410

TriviaQA features open-domain questions aimed
1411 at assessing general world knowledge (Joshi et al.,
1412 2017). 1413

Winogrande contains questions modeled after
1414 the Winograd schema but scaled up in size and
1415 difficulty (Sakaguchi et al., 2021). 1416

XNLI provides natural language inference exam-
1417 ples in 15 languages, where the task is to determine
1418 whether a hypothesis is entailed by, contradicts, or
1419 is neutral with respect to a given premise (Conneau
1420 et al., 2018). 1421

XQuAD is a cross-lingual question answering
1422 dataset that pairs reading passages with related
1423 questions in 11 languages, where the correct an-
1424 swer is always a text span from the passage itself
1425 (Artetxe et al., 2020). 1426

Global MMLU extends the original MMLU
1427 benchmark to assess multilingual capabilities, fea-
1428 turing 4-option multiple-choice questions across
1429 57 subject areas in 42 languages including low-
1430 resource languages, testing both factual knowledge
1431 and reasoning skills in diverse linguistic contexts
1432 (Singh et al., 2025). 1433

1434 F Additional Related Work

**Tokenization for morphologically-rich lan-
1435 guages** Standard BPE tokenization often strug-
1436 gles to capture morphologically complex lan-
1437 guages (Klein and Tsarfaty, 2020; Park et al., 2021;
1438 Mager et al., 2022; Hofmann et al., 2022). Ara-
1439 bic inflection, for instance, uses non-concatenative
1440 morphology that breaks standard subword reusabil-
1441 ity (Alyafeai et al., 2021; Alrefaie et al., 2024; Tsar-
1442 faty et al., 2019; Gazit et al., 2025). Compositional
1443

Transformation	In-vocab.			Out-of-vocab.		
	<i>embed</i>	<i>detok</i>	<i>N</i>	<i>embed</i>	<i>detok</i>	<i>N</i>
Inflection						
Plural (N)	92%	92%	0.8k	24%	31%	3.4k
Plural (N) & Present Singular (V)	86%	87%	1.6k	35%	44%	2.1k
Present Singular (V)	91%	91%	0.1k	54%	64%	0.3k
Past (V)	65%	68%	0.6k	10%	15%	2.9k
Past Participle (V)	79%	79%	14	24%	29%	21
Gerund (V)	83%	84%	0.2k	17%	22%	3.2k
Superlative (ADJ)	87%	87%	31	3%	10%	0.4k
Comparative (ADJ)	47%	67%	30	4%	12%	0.4k
Capitalization	72%	73%	6.0k	74%	76%	8.3k
Derivation						
-y	17%	22%	18	2%	6%	1.5k
-er	25%	25%	12	1%	3%	2.6k
-al	62%	62%	8	1%	2%	0.7k
un-	0%	33%	3	0%	1%	3.3k
re-	67%	67%	3	0%	1%	1.8k
-ic	100%	100%	2	5%	7%	0.4k
All derivatives	40%	44%	52	0%	1%	31.4k

Table 5: Accuracy of Patchscopes interpretations for Qwen-2.5-7B.

Transformation	In-vocab.			Out-of-vocab.		
	<i>embed</i>	<i>detok</i>	<i>N</i>	<i>embed</i>	<i>detok</i>	<i>N</i>
Inflection						
Plural (N)	93%	94%	0.8k	34%	42%	3.4k
Plural (N) & Present Singular (V)	86%	90%	1.6k	41%	58%	2.1k
Present Singular (V)	90%	91%	0.1k	60%	71%	0.3k
Past (V)	74%	85%	0.6k	12%	24%	2.9k
Past Participle (V)	100%	100%	14	24%	43%	21
Gerund (V)	93%	97%	0.2k	26%	38%	3.2k
Superlative (ADJ)	97%	97%	31	20%	38%	0.4k
Comparative (ADJ)	87%	90%	30	7%	18%	0.4k
Capitalization	80%	96%	6.0k	50%	85%	8.3k
Derivation						
-y	65%	65%	17	13%	19%	1.5k
-er	25%	33%	12	6%	19%	2.6k
-al	75%	88%	8	4%	11%	0.7k
un-	33%	33%	3	1%	6%	3.3k
re-	100%	100%	3	1%	17%	1.8k
-ic	100%	100%	2	10%	15%	0.4k
All derivatives	63%	67%	51	2%	6%	31.4k

Table 6: Accuracy of Patchscopes interpretations for OLMo-2-7B.

1444 vocabularies can bypass such limitations by rep-
1445 resenting surface forms as transformations over
1446 lexical roots, enabling reuse of base forms even
1447 when their surface realizations use diverging token
1448 sequences.

Language	Capitalization		Noun Inflection		Adjective Inflection		Verb Inflection		Derivation		
	<i>In-Vocab.</i>	<i>Out-Vocab.</i>	<i>In-Vocab.</i>	<i>Out-Vocab.</i>	<i>In-Vocab.</i>	<i>Out-Vocab.</i>	<i>In-Vocab.</i>	<i>Out-Vocab.</i>	<i>In-Vocab.</i>	<i>Out-Vocab.</i>	
<i>ALLaM</i>	Arabic	—	—	78% (1.8k)	16% (3.6k)	69% (0.5k)	25% (1.0k)	43% (1.0k)	15% (2.7k)	—	—
<i>EuroLLM</i>	German	100% (0.2k)	89% (0.4k)	—	—	27% (0.3k)	11% (1.3k)	88% (0.3k)	44% (1.2k)	—	—
	Russian	98% (66)	96% (0.7k)	72% (0.6k)	28% (4.2k)	100% (50)	93% (94)	100% (6)	50% (10)	—	—
	Spanish	100% (1.0k)	97% (2.8k)	83% (0.7k)	59% (1.9k)	82% (0.5k)	67% (1.1k)	72% (0.8k)	42% (6.9k)	46% (65)	20% (0.4k)
<i>Llama-3</i>	English	89% (6.0k)	85% (8.4k)	93% (2.4k)	63% (5.6k)	89% (61)	32% (0.9k)	85% (0.9k)	34% (6.4k)	34% (41)	4% (12.8k)

Table 7: Accuracy of Patchscopes *detokenization* interpretations for compositional input embeddings across languages.

Category	Task	Baseline	End-to-end	Δ
Knowledge	TinyMMLU _(Acc.)	62.9	62.5	-0.4
	TinyARC _(Acc.)	51.8	48.8	-3.0
Reading Comprehension	BoolQ _(Acc.)	87.6	87.6	+0.1
	TriviaQA _(EM)	58.3	52.9	-5.4
	SQuAD _(EM)	37.3	31.9	-5.5
Commonsense	TinyHellaswag _(Acc.)	52.1	54.3	+2.2
	TinyWinogrande _(Acc.)	62.3	64.3	+2.0
	PIQA _(Acc.)	79.5	78.7	-0.8
	COPA _(Acc.)	91.0	90.0	-1.0
Average		64.7	63.4	-1.3

Table 8: Downstream performance of English compositional-vocabulary models (*End-to-end*) and their original, unmodified version (*Baseline*) for Qwen2.5-7B.

Category	Task	Baseline	End-to-end	Δ
Knowledge	TinyMMLU _(Acc.)	51.2	51.9	+0.7
	TinyARC _(Acc.)	51.9	49.1	-2.8
Reading Comprehension	BoolQ _(Acc.)	84.5	84.8	+0.3
	TriviaQA _(EM)	65.4	56.8	-8.6
	SQuAD _(EM)	40.0	30.6	-9.3
Commonsense	TinyHellaswag _(Acc.)	62.6	63.4	+0.8
	TinyWinogrande _(Acc.)	63.5	63.9	+0.4
	PIQA _(Acc.)	80.4	79.5	-0.9
	COPA _(Acc.)	91.0	89.0	-2.0
Average		65.6	63.2	-2.4

Table 9: Downstream performance for OLMo-2-7B.