



HAL
open science

Towards the Machine Translation of Scientific Neologisms

Paul Lerner, François Yvon

► **To cite this version:**

Paul Lerner, François Yvon. Towards the Machine Translation of Scientific Neologisms. Rapport D2-3.1, ISIR, Université Pierre et Marie Curie UMR CNRS 7222. 2025. hal-04852293

HAL Id: hal-04852293

<https://hal.science/hal-04852293v1>

Submitted on 20 Dec 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - ShareAlike 4.0 International License

Towards the Machine Translation of Scientific Neologisms

Paul Lerner and François Yvon

January 2025

MaTOS — Livrable D2-3.1

Machine Translation for Open Science - ANR-22-CE23-0033



Towards the Machine Translation of Scientific Neologisms

Paul Lerner and François Yvon

January 2025

Abstract

Scientific research continually discovers and invents new concepts, which are then referred to by new terms, neologisms, or *neonyms* in this context. As the vast majority of publications are written in English, disseminating this new knowledge to the general public often requires translating these terms. However, by definition, no parallel data exist to provide such translations. Therefore, we propose to leverage term definitions as a useful source of information for the translation process. As we discuss, Large Language Models are well suited for this task and can benefit from in-context learning with co-hyponyms and terms sharing the same derivation paradigm. These models, however, are sensitive to the superficial and morphological similarity between source and target terms. Their predictions are also impacted by subword tokenization, especially for prefixed terms.

We also extended experiments on segmentation into sub-lexical units with a controlled corpus, with negative prefixation and adverbial suffixation of adjectival bases or pseudowords. Our results confirm the previous ones: language models struggle to generate prefixations due to sub-optimal segmentation, which can be resolved through morphological segmentation. We enrich these results with analyses of the alignment between subword embeddings.

Résumé

La recherche scientifique découvre et invente continuellement de nouveaux concepts qui sont alors désignés par de nouveaux termes, des néologismes ou néonymes dans ce contexte. Puisque les publications se font très majoritairement en anglais, il convient de traduire fidèlement ces termes dans d'autres langues, comme le français, tout en évitant une multiplication d'anglicismes. Toutefois, il n'existe par définition pas de données parallèles où trouver des néologismes. Nous proposons donc d'exploiter la définition du terme afin de le traduire plus fidèlement. Pour ce faire, nous explorons les capacités de modèles de langues multilingues, qui parviennent à traduire des néologismes scientifiques dans une certaine mesure. Nous montrons notamment qu'ils utilisent souvent des procédés morphosyntaxiques appropriés mais

sont limités par la segmentation en unités sous-lexicales, particulièrement pour la préfixation, et biaisés par la fréquence d'occurrences des termes ainsi que par des similarités de surface entre l'anglais et le français. Afin de pallier ces limites, nous proposons de sélectionner des exemples (*in-context learning*) co-hyponymes du terme ou issus du même paradigme dérivationnel.

Nous avons également approfondi les expériences sur la segmentation en unités sous-lexicales avec un corpus contrôlé, avec une préfixation négative et une suffixation adverbiale par base adjectivale ou pseudo-mot. Nos résultats confirment les précédents: les modèles de langues peinent à générer des préfixations en raison d'une segmentation sous-optimale, ce qui peut être résolu grâce à une segmentation morphologique. Nous enrichissons ces résultats par des analyses sur l'alignement entre les plongements des sous-mots.

Contents

1. Introduction	6
2. Related Work	7
3. Discussion	8
4. Neological and Morphological Processes	9
5. Translation Methods	10
5.1. Implementation	11
5.2. Evaluation	12
6. Translation Results	13
6.1. Datasets	13
6.2. Definition-augmented Translation	14
6.3. In-Context Learning	14
6.4. Frequency Bias and Semantic Change	15
6.5. Morphosyntactic Analysis	16
6.6. Morphosyntactic Divergences	17
6.7. Translation or Orthographic Conversion?	17
6.8. Prefixation, Fertility, and BPE	18
7. Derivational Morphology	19
8. Derivation Methods	21
8.1. Definition to Word Generation	21
8.2. In-Context Learning	21
8.3. Large Multilingual Models	21
8.4. Segmentation	22
8.5. Controlled Datasets	22
9. Derivation Results	23
9.1. Prefixations vs. Suffixations	23
9.2. Initial- vs. Intra-word Alignment	24
10. Conclusion	25
11. Limitations	26
11.1. On Translation	26
11.2. On Morphology	27
Bibliography	28
A. Machine-Translated Definitions	41

B. Frequency and Neology	42
C. Morphosyntactic Classification	42
C.1. Datasets	43
C.2. Implementation	44
C.3. Results	44
D. Implementation Details	45
D.1. LLM Implementation	45
D.2. mBART Fine-tuning on SciPar	45
D.3. Corpus frequency	45
E. Rules of Morphotactics	45
F. Complete Results	47

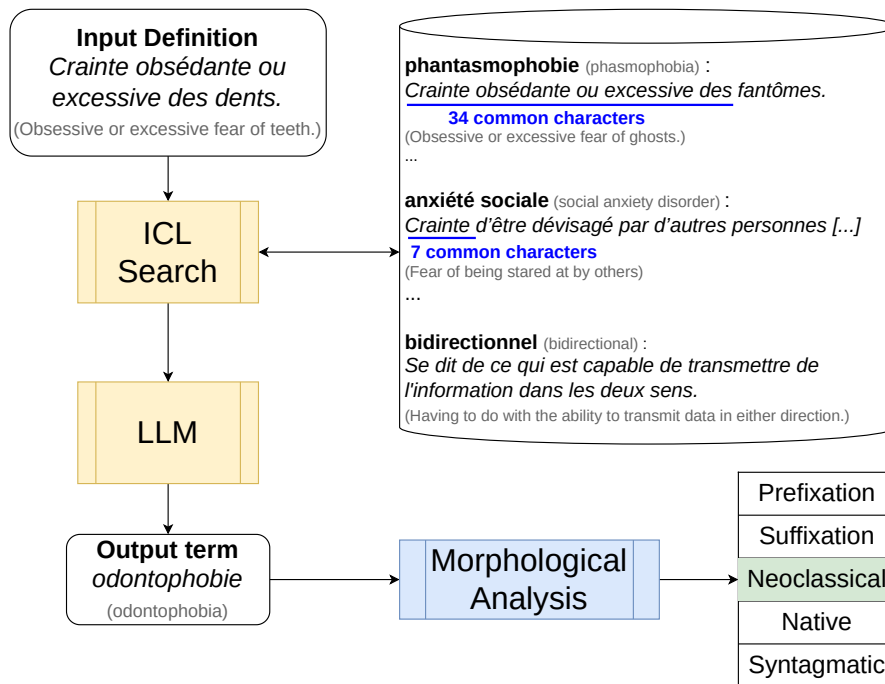


Figure 1: Overview of our experiments: in DEF setting, given a definition, we study how to retrieve relevant ICL examples, here co-hyponyms. An LLM is then tasked to generate a term matching the definition. We also perform several analyses, including a morphological analysis of the output term. See text for details.

1. Introduction

New concepts are constantly introduced by researchers around the world, which leads to a profusion of neologisms. These are also known as *neonyms* [Cabré, 1999], as opposed to neologisms of everyday language [Cartier et al., 2018]. Because most of this research is published in English [Gordin, 2015, Larivière and Riddles, 2021],¹ communicating in another language, such as French, requires translating these terms to facilitate scientific dissemination.² For example, a teacher wanting to instruct their French students about “Large Language Models” would be hardly understandable if they directly borrowed every term from English, e.g.:

EN: large language models are self-supervised

?? les *large language models* sont *self-supervised*

FR: les grands modèles de langue sont auto-supervisés

Quoting Liu et al. [2021]: “Precisely defining the terminology is the first step in scientific communication”.

Translating scientific neologisms is a fundamental problem for traditional Machine Translation

¹In French-speaking countries, a significant part of research in humanities and social sciences is still disseminated in French. The same holds for other major linguistic areas.

²See, e.g., <https://www.helsinki-initiative.org/>.

(MT) systems that rely on parallel data, which, by definition, can not contain such new words.³ Therefore, we propose to leverage definitions of terms as a way to translate them more accurately. We study how to take this information into account and, in particular, how to select relevant examples for in-context learning, in a linguistically motivated manner. We conduct extensive experiments on two thesauri covering 13 diverse domains, from Humanities to Computer Science and find our methods to be domain-agnostic. As we focus on translation from English into French, we rely on the fact that neologisms are mostly formed through five non-exclusive morphological processes (prefixation, suffixation, and neoclassical, native, or syntagmatic compounding), and study (i) how morphological divergences between the source and target impact translation; (ii) whether systems outputs conform to attested morphological patterns (see Figure 1).

Terminology remains a major source of critical errors for MT [Haque et al., 2020], which is often tackled by augmenting MT systems with domain-specific resources and dedicated (pre-)processing modules [Semenov et al., 2023]. Our work could benefit such approaches by enriching said thesaurus or providing on-the-fly translations by extracting definitions from source documents [Jin et al., 2013, Head et al., 2021, August et al., 2022, Huang et al., 2022].

We tackle Neologism Translation with Large Multilingual Language Models (mLLMs), which are effective for many MT and NLP tasks [Xu et al., 2024]. We show that these models are able, to some extent, to translate terms from English to French, to generate a term from its (French) definition, and also to combine both sources of information. We also show that LLMs benefit from in-context learning examples that are co-hyponyms or belong to the same derivational paradigm as the source term/definition (see Figure 1). However, we also highlight several limitations of these models: (i) their tokenizer, based on crude heuristics such as BPE [Gage, 1994], tends to over-segment prefixed terms, which is detrimental to translation quality; (ii) they perform much better if the source and target term are superficially similar (likely cognates or loanwords), which makes the task closer to orthographic conversion than translation (e.g. *exocytosis* → *exocytose*); (iii) their performance correlates with terms frequency in a large corpus, which may be used as a proxy of their degree of lexicalization.

This work opens up new challenges for MT and more broadly NLP, on an important topic for knowledge dissemination. It also sheds light on the somewhat overlooked issue of morphological processing in LLMs. We propose several avenues for future work to address the limitations outlined above. This work has been published in two conferences papers [Lerner and Yvon, 2025b,a]. Our code and data are freely available.⁴

2. Related Work

While we rely on definitions to generate neologisms, some work has been done in the opposite direction, to generate the definition of a given word [Noraset et al., 2017]. Interestingly, like us, they leverage the structure of definitions in *genus* and *differentiae* [Chodorow et al., 1985, Montemagni and Vanderwende, 1992]. The *genus* is a hypernym of the input term (see Figure 1, *phasmophobia* is a kind of *fear*). We will find that terms sharing the same hypernym prove to be useful examples for In-Context Learning.

³At least, not with their new intended meaning.

⁴<https://github.com/PaulLerner/neott>

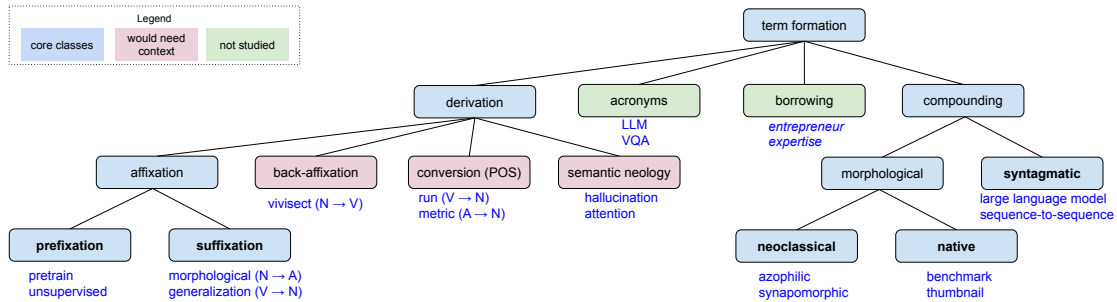


Figure 2: Overview of the studied neological processes. Adapted from [Daille, 2017].

Neologism Translation is related to Multilingual Term Extraction [Laroche and Langlais, 2010, Delpech et al., 2012, Rigouts Terryn et al., 2020], except that, importantly, we do not assume that the target term exists *anywhere*. Indeed, we will see that a significant part of the terms in our test data do not appear even a single time in a large corpus such as OSCAR [Abadji et al., 2022].

Our framing of Neologism Translation somewhat resembles the *Reverse Dictionary* task [Hill et al., 2016, Pilehvar, 2019]. However, Reverse Dictionary is an Information Retrieval task that consists of mapping the representation of a definition to an existing word embedding of a *known word*. On the contrary, we design here a fully *generative* task for *unknown words*.

The study of Zhang et al. [2020] comes closest to our work but is restricted to a monolingual setting in the very specific domain of genetics, where a term is linked to several genes according to its molecular function, biological process, and cellular component.

Hofmann et al. [2021] and Truong et al. [2024] are also interested in derivational morphology and LLMs but consider binary classification tasks that assess whether the LLM “understands” words, while we are interested in the actual generation of new forms. Our results are consistent with theirs: notably, Hofmann et al. [2021] also find that LLMs are unable to process prefixations, compared to suffixations, for the same reason. They also find that enforcing morphological segmentation improves performance. Hofmann et al. [2020] is similar to our work but always relies on morphological segmentation, except in their preliminary experiment.

Oh and Schuler [2024] and Pimentel and Meister [2024] discuss another effect of BPE marking the beginning of words: the miscomputation of word probabilities, an indicator of word surprise used in psycholinguistic studies. Both propose a simple rescaling method to recover the correct values.

3. Discussion

BPE is ubiquitous in NLP as virtually all LLMs depend on it. However, marking strings in the beginning of words leads to caveats that are overlooked. We show that this faulty tokenization limits the ability of LLMs to generate prefixations, a morphological process that is however productive in many languages. Such defects in morphological abilities may partly explain the recurrent difficulties of LLMs to generate a sufficiently large number of new lexemes, as attested by low Type-to-Token Ratio scores in generated texts [Muñoz-Ortiz et al., 2024]. We also show that an accessible solution is morphological segmentation, which enables even “small” models

(of a few hundred million parameters) to reach near-perfect generation accuracy.

4. Neological and Morphological Processes

Typology Our typology of neologisms is adapted from Lieber [2010] and Daille [2017], and relies on morphosyntactic features that can easily be detected automatically. Complementary typologies, which vary according to the studied phenomena, have also been proposed, see, e.g. [Lombard and Huyghe, 2020]. We retain the following five constructions that cover the largest part of our corpus, both in English and French:

(i) **Prefixation**, where an affix is concatenated at the beginning of a word to form a new one (e.g., *pre+train* = *pretrain*).

(ii) **Suffixation**, where affixation is performed at the word’s end (e.g. *generalize+ation* = *generalization*).

(iii) **Native compounding**, which compounds two independent words. This process is more regular in English (e.g. *bench+mark* = *benchmark*) than in French [Arnaud, 2003].

(iv) **Neoclassical compounding**, which compounds only *bound morphemes*, i.e. morphemes that cannot act as independent words (e.g., *azo+philic* = *azophilic*). Like native English but unlike native French, the head of neoclassical compounds is always located at the rightmost position, in both languages: e.g. *azophilic* means “attracted to azote”, not “azote is attracted” [Namer, 2003, Amiot and Dal, 2008].

(v) **Syntagmatic compounding**, where syntagms that follow syntactic rules of the language are lexicalized into terms, thereby losing the compositionality of meaning. Therefore, they often cannot be translated by a composition of translations of its constituents [Daille and Morin, 2005], e.g. “zero-shot learning” translates to “*apprentissage sans exemple*” in French, literally “learning without example”.

Note that for (i), (ii), and (iv), derivation is often accompanied by a phonological or graphemic change at the junction between morphemes. Finally, note that these processes are not exclusive but can be combined, e.g. *bidirectional* is a prefixation (*bi-*) of a suffixation (*-al*).⁵ All studied morphological processes are illustrated in Figure 2.

Figure 2 also includes rarer processes that would require a disambiguating context and are therefore not handled by the morphological classifier introduced below: (i) **Semantic neology**, where a lexical unit is associated with a new concept through a metaphoric transfer between two domains, resulting in a homonym. (ii) **Conversion**, where the part-of-speech (POS) of a word changes without affixation, resulting again in a homonym [Tribout, 2010]. (iii) **Back-affixation**, which requires a diachronic perspective to recognize it among other affixations (e.g. *vivisect* is formed by removing *tion* from *vivisection*, and not the other way around).

We finally do not study the following processes, although they are frequent in both English and French: (i) **Borrowing**, because we precisely seek to avoid it (e.g. *entrepreneur* is borrowed as is from French). (ii) **Acronyms**, which cannot be translated without their expanded form.

⁵It could also be interpreted as the suffixation of the noun **bidirection* although it is unattested [Corbin, 2012a]. See also Copot and Bonami [2024] for a “baseless” approach to derivation where both *directional* and **bidirection* could interact with *bidirectional*.

Setting	Prompt template
TERM	Le terme anglais {src_term} peut se traduire en français par : <i>“The English term {src_term} can be translated in French as :”</i>
DEF	{def} définit le terme : <i>“{def} defines the term :”</i>
DEF+TERM	{def} définit le terme anglais {src_term} qui peut se traduire en français par : <i>“{def} defines the English term {src_term} which can be translated to French as :”</i>

Figure 3: Prompt templates used with LLMs corresponding to our three settings, with English translations

The reader should refer to Dal [2003a], Lieber [2010], or Corbin [2012a] for a more complete introduction to morphology,⁶ going beyond English and French, and therefore, the above processes (e.g. templates in Semitic languages). Finally note that we are not interested in inflections (e.g. singular/plural), which do not form new lexemes.

Morphosyntactic Classification We build two multi-label classifiers, one per language, to identify the morphosyntactic processes described above. They rely on character n-gram features and are trained on Wiktionary in the FastText framework [Joulin et al., 2017]. They are very accurate with 92.5 F1 in English and 95.8 F1 in French, see Appendix C for details. This classifier is used below to analyze the morphological processes used to coin new terms (see Figure 1), to evaluate English-French congruences and divergences and how they impact the performance of the models.

5. Translation Methods

We study the translation of neologisms in three settings, always in the EN-FR direction, which is our main application scenario (see Section 6.1):⁷ (i) **TERM**: translate the contextless source term. This is our baseline condition. (ii) **DEF**: generate the target term from its definition in the same language, one of the main novelties of our work (see Figure 1); (iii) **DEF+TERM**: translate the source term given its definition, combining the two sources of information. Both input terms and definitions are extracted from public thesauri (see Section 6.1).

We cast these three subtasks in a text-to-text generation framework, where an LLM is tasked to complete a prompt [Brown et al., 2020, Raffel et al., 2020]. Because of the mixed language input in setting DEF+TERM, we use mLLMs. The prompt may contain several examples to enable in-context learning (ICL). We study four ways to select these examples: the first two serve as baselines, while the last two are linguistically motivated:

- (i) **Random**: sampling from the set of examples for ICL.

⁶See also Aronoff [1976] and Fradin [2015] for a lexematic approach to morphology and Dal [2003b] and Mattiello [2017] on analogy.

⁷Moreover, as most neonyms are first formed in English, then translated to French, studying the reverse direction (FR-EN) would be plagued by translationese, which is known to lead to overoptimistic results [Zhang and Toral, 2019].

(ii) **Domain**: similar to *Random*, additionally requiring ICL examples to belong to the same domain as the target term (“oracle” condition).

(iii) **Co-hyponyms**: terms sharing the same hypernym are often formed in the same way. To find co-hyponyms, we simply rely on the longest common string with the beginning of the *input definition* (see Figure 1). Therefore, this method does not apply to the *TERM* setting, which does not have access to definitions. For instance, definitions starting with “*Crainte obsédante ou excessive des*”⁸ identify several *phobias*, e.g. *traumatophobie* (traumatophobia) or *odontophobie* (odontophobia). With “*Opération consistant à*”,⁹ we find deverbals in *-ation* or *-age*, e.g. *dénaturation* (denaturation), *quantification* (quantizing), or *tricotage* (knitting).

(iv) **Derivation paradigms**: as hinted at above, terms stemming from the same derivational paradigm, i.e. sharing a base, prefix, or suffix, may serve as analogical context to form new terms. For example, *pretraining* was likely formed on the model of *preprocessing*; likewise for *underfitting* modeled after *overfitting*. Like for co-hyponyms, we rely on the longest common string, but this time between *source terms*, either at the beginning or the term ending. Therefore, this method does not apply to the *DEF* setting, which does not have access to the source term. Note that this method is not limited to morphological affixes but can also find whole words in common between syntagms. For example, “*air gap*” and “*air flotation*” share the word *air* in their initial and “*unmoderated newsgroup*” and “*unmerchantable*” share the prefix *un-*.

The last two methods can be both combined in the *DEF+TERM* setting by concatenating their top results, while keeping the total number of examples to five. The hyperparameters for this fusion are set through grid search on the validation set.¹⁰ We limit the number of examples to five to keep a reasonable input length and as we found the performance to quickly saturate, consistently with prior work (e.g. [Bawden and Yvon, 2023](#)).

5.1. Implementation

We experiment with two mLLMs: BLOOM [[BigScience et al., 2023](#)] and CroissantLLM [[Faysse et al., 2024](#)]. BLOOM was the first open-source mLLM to scale up to billions of parameters. It is highly multilingual, trained on 46 natural languages, including EN and FR. We experiment with both 1.1B and 7.1B parameters versions. CroissantLLM is an EN-FR bilingual model, trained on an equally large amount of data in the two languages. With only 1.3B parameters, it was designed to be efficient at inference time, to make up for its costly pretraining, following [Liu et al. \[2019\]](#) and [Hoffmann et al. \[2022\]](#).

Each of our three prompt templates (see Figure 3) correspond to one settings presented above. We experimented with a few different wordings but found that the prompt content hardly mattered because of ICL examples, consistently with prior work (e.g. [Zebaze et al., 2024](#)). ICL examples use the same prompt template, but include both the instruction and the target term. Different examples are separated by the three characters *###*, which serves as end-of-sequence signal.

Apart from LLMs, we use mBART as a standard sequence-to-sequence baseline for the *TERM* setting (standard MT). More precisely, we fine-tune mBART50-One-to-Many, a 610M parameter

⁸“Obsessive or excessive fear of”.

⁹“Operation consisting of”.

¹⁰The optimum for *Derivation paradigms* is three prefixes and two suffixes. When fusing with *Co-hyponyms* the optimum is one co-hyponym from the definition, three prefixes, and one suffix.

model [Tang et al., 2021], on 1.1M EN-FR parallel sentences from SciPar [Roussis et al., 2022]. This process ensures that the model is robust to scientific vocabulary. Still, mBART only translates from EN to FR and is not suited for the conditions DEF and DEF+TERM. This model achieves 37.3 BLEU on a held-out test set of 3K sentences [Papineni et al., 2002]. See Peng et al. [2024] and Appendix D for additional details.

5.2. Evaluation

We draw inspiration from standard Question Answering metrics (e.g. Rajpurkar et al., 2016) and considered: (i) Exact Match (EM) between the target and output strings;¹¹ (ii) token-level F1 score after standard preprocessing (case insensitive, stop-words and punctuation filtering). At a time when LLM-based metrics flourish, one might criticize these metrics for being overly strict and not modeling semantic similarity. However, we argue that evaluating terminological equivalence is mostly not a semantic matter: the meaning of the terms is highly dependent on the domain and words that would otherwise be synonymous often cannot be used interchangeably. For instance “*big language model” is an incorrect variant of “large language model”, although *big* and *large* are synonyms (i.e. semantically close, even with a non-neural metric like METEOR; Banerjee and Lavie, 2005). Moreover, LLM-based metrics are known to bias towards models with the same architecture or training data [He et al., 2023, Panickssery et al., 2024], while EM is equally strict for all models.

In addition to EM and F1, we also assess whether our models generate terms with the same morphological processes as the reference, as described in Section 4 (see Figure 1).

Model	Setting	FranceTerme		TERMIUM	
		EM	F1	EM	F1
mBART	TERM	<u>26.3</u>	41.3	<u>31.1</u>	49.7
CroissantLLM	TERM	25.6	<u>42.2</u>	30.3	<u>50.3</u>
CroissantLLM	DEF	4.6	19.8	3.8	22.7
CroissantLLM	DEF+TERM	25.3	42.9	30.2	51.5
BLOOM-1.1B	TERM	15.9	31.3	17.1	37.1
BLOOM-1.1B	DEF	1.1	11.3	1.4	15.4
BLOOM-1.1B	DEF+TERM	17.8	34.9	20.0	41.2
BLOOM-7.1B	TERM	23.7	40.3	27.5	47.7
BLOOM-7.1B	DEF	<u>10.0</u>	<u>24.7</u>	<u>7.5</u>	<u>26.6</u>
BLOOM-7.1B	DEF+TERM	27.1	44.6	32.1	53.5

Table 1: Definition-augmented Translation results on the test sets of FranceTerme and TERMIUM, with 5 randomly selected ICL examples for LLMs. Best overall results are bolded while best results in settings TERM and DEF are underlined.

¹¹EM is also used to evaluate morphological reinflection in the SIGMORPHON Shared Task, where it is referred to as “accuracy” [Cotterell et al., 2016].

6. Translation Results

6.1. Datasets

We experiment with two EN-FR bilingual thesauri in this work: FranceTerme¹² and TERMIUM,¹³ which are curated by the French and Canadian governments, respectively. Both of these thesauri are well-studied in the neology literature [Pecman, 2012, Tonti, 2023, Holeš, 2024]. We filter loanwords (cf. Section 4) by removing terms that are identical in EN and FR (case insensitive; 2.9% of FranceTerme, 4.6% of TERMIUM). To filter acronyms, we discard terms with two consecutive upper-case letters (1.8% of FranceTerme, 2.3% of TERMIUM). We also filter entries with missing data to only keep triples of (EN term, FR term, FR definition).¹⁴ FranceTerme finally amounts to 6,623 terms equally and randomly split into validation and test sets. When testing, the validation set will serve for ICL and vice-versa. TERMIUM is much larger so we randomly keep 5,000 terms for validation, 5,000 for testing, and the remaining 194,992 for ICL. TERMIUM broadly covers 13 coarse-grained domains (listed in Table 3), which are balanced enough so that we can confidently compute statistics for each of them (from 83 samples in Metal. to 895 in MPS in the test set). On the other hand, FranceTerme covers ≈ 70 very imbalanced domains (some containing just one sample) so we only consider it as a whole.

Setting	ICL	FranceTerme		TERMIUM	
		EM	F1	EM	F1
TERM	Random	23.7	40.3	27.5	47.7
TERM	Domain	26.3	42.6	29.6	49.7
TERM	Paradigm	<u>27.0</u>	<u>43.8</u>	<u>36.3</u>	<u>55.4</u>
DEF	Random	10.0	24.7	7.5	26.6
DEF	Domain	10.1	25.1	8.6	27.5
DEF	Co-hyponyms	<u>10.7</u>	<u>25.8</u>	<u>10.5</u>	<u>30.0</u>
DEF+TERM	Random	27.1	44.6	32.1	53.5
DEF+TERM	Domain	28.5	46.0	32.5	54.2
DEF+TERM	Fusion	31.2	48.2	40.7	60.0

Table 2: Results of BLOOM-7.1B on the test sets of FranceTerme and TERMIUM according to our ICL selection strategy: (i) random (baseline); (ii) domain (baseline); (iii) derivation paradigm (not applicable to DEF); (iv) co-hyponyms (not applicable to TERM); (v) fusion of the latter two. Best overall results are bolded while best results in settings TERM and DEF are underlined.

¹²<https://www.culture.fr/franceterme>, open license compatible with CC-BY 2.0, version of November 17 2023.

¹³<https://www.btb.termiumplus.gc.ca/> Open Government Licence - Canada, version of February 6 2023.

¹⁴FranceTerme definitions are only available in FR, the target language. TERMIUM provides both EN and FR definitions, so we provide additional results in Appendix A with machine-translated definitions. We find our results to be consistent with both reference and machine-translated French definitions.

Setting	ICL	Agr.	CS	Indus.	MPS	Mech.	Med.	Hum.	Env.	Tele.	Jus.	Eco.	Elec.	Metal.
TERM	Random	20.5	36.2	16.5	33.0	18.9	31.3	27.2	32.9	32.7	<u>33.1</u>	26.5	24.4	19.3
TERM	Paradigm	<u>22.6</u>	<u>44.1</u>	<u>22.7</u>	<u>46.8</u>	<u>28.1</u>	<u>50.0</u>	<u>34.5</u>	<u>39.1</u>	<u>38.1</u>	30.5	<u>31.9</u>	<u>31.1</u>	<u>24.1</u>
DEF	Random	<u>5.6</u>	6.0	5.9	6.9	5.1	11.2	10.2	9.1	5.4	5.9	8.4	6.7	4.8
DEF	Co-hyponyms	<u>5.6</u>	<u>8.6</u>	<u>9.7</u>	<u>11.7</u>	<u>9.7</u>	<u>15.3</u>	<u>11.5</u>	<u>13.2</u>	<u>6.5</u>	<u>7.6</u>	<u>8.6</u>	<u>7.7</u>	<u>10.8</u>
DEF+TERM	Random	29.2	40.5	21.2	36.9	21.9	37.1	34.3	37.4	32.7	32.2	31.3	25.4	28.9
DEF+TERM	Fusion	28.7	44.8	28.4	48.5	31.1	52.9	40.9	46.0	44.6	38.1	37.3	34.9	33.7

Table 3: Exact Match of BLOOM-7.1B on the 13 domains of TERMIUM according to our ICL selection strategy: Agriculture (Agr.), Electronic and Computer Science (CS), Industries (Indus.), Maths Physics and Natural Sciences (MPS), Mechanics (Mech.), Medicine (Med.), Humanities (Hum.), Environmental Sciences (Env.), Telecommunications (Tele.), Law and Justice (Jus.), Economy (Eco.), Electricity (Elec.), and Metallurgy (Metal.). Best overall results are bolded while best results in settings TERM and DEF are underlined.

6.2. Definition-augmented Translation

We now explore the three settings of Neologisms Translation with our four models, keeping ICL selection random (see Table 1). We find that TERM, translating the contextless source term, is much easier than DEF, where the input is the FR definition. However, the performance of models in setting TERM are limited, with mBART, BLOOM-7.1B, and CroissantLLM all reaching similar performance. We find that BLOOM-7.1B is able to combine information from source term and definition in setting DEF+TERM, significantly outperforming TERM. Model size is particularly important in this setting, as we observe that BLOOM-1.1B and CroissantLLM, which are roughly the same size, barely outperform or even deteriorate TERM when using the additional definition. Therefore, we focus on BLOOM-7.1B in the following experiments. BLOOM-7.1B DEF+TERM is so effective that it outperforms an oracle late fusion of TERM and DEF, suggesting an interaction between the two sources of information. For instance, BLOOM-7.1B DEF+TERM correctly predicts *capteur de mission* for *mission sensor* “capteur réalisant des mesures qui font partie de l’objet de la mission d’un engin spatial”,¹⁵ unlike TERM which predicts *mission de reconnaissance* and DEF which predicts *instrument de mesure* (“measuring instrument”).

6.3. In-Context Learning

Results according to our different ICL strategies are in Table 2. We find that our strategies consistently improve over random and domain selection, even though the latter accesses the ground-truth domain through an oracle. The performance gains are especially high for TERMIUM, where the set of examples for ICL is much larger. Furthermore, we show in Table 3 that our methods are domain-agnostic, with significant improvements in 12 out of 13 domains of TERMIUM, from Humanities to Computer Science. In the rest of this section, we will focus on FranceTerme for the sake of space, but our results are consistent on both datasets.

¹⁵“sensor performing measurements that are part of the mission of a spacecraft”.

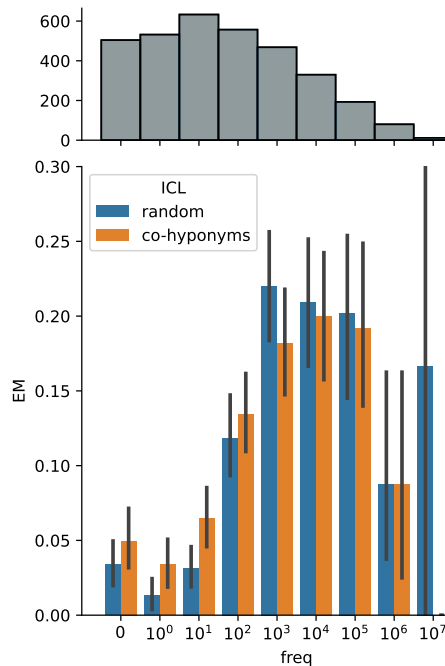


Figure 4: Exact Match (EM) of BLOOM-7.1B (DEF) w.r.t. term’s corpus frequency, comparing random and co-hyponym ICL selection, on FranceTerme’s test set. The upper part shows the number of examples in each bin. Note the logarithmic scale of the x -axis.

6.4. Frequency Bias and Semantic Change

Our main research interest lies in Neologism Translation. However, assessing whether a term is neological or lexicalized is a subjective matter [Lombard and Huyghe, 2020]. Therefore, we choose a continuous scale of neology based on the term’s frequency in large corpora, namely ROOTS-fr-open [Laurençon et al., 2022] and OSCAR-fr 22.01 [Abadji et al., 2022]. ROOTS-fr-open is a French CC-licensed subset of ROOTS, the dataset used to train BLOOM. It consists of 4 billion words (20 GB), mostly from Wikimedia. OSCAR-fr 22.01 is a French cleaned subset of Common Crawl, which was also partly used to train BLOOM. It consists of 42 billion words (382 GB).

Figure 4 shows that 15.8% of FranceTerme target (French) terms do not appear even a single time in this huge corpus, and most appear less than 100 times (i.e. the frequency of monolexical terms is less than 2×10^{-9}). See Appendix B for examples of each decile. We find that the neological feeling [Lombard and Huyghe, 2020] is weaker after 1,000 occurrences (e.g. *effet de rebond* “rebound effect”). It is not a coincidence that BLOOM (DEF) predicts terms much more accurately above this 1,000 occurrences threshold (Figure 4). However, the bulk of the distribution lies before 1,000, where we find our co-hyponym ICL selection method to significantly and consistently improve results. For example, given “*Enzyme qui déphosphoryle les résidus sérine*,

thréonine ou tyrosine préalablement phosphorylés, présents dans les protéines”,¹⁶ BLOOM, with random ICL, fails to generate *protéine-phosphatase* (“protein phosphatase”, 0 occurrences), while our co-hyponym selection strategy succeeds because of relevant ICL examples such as *protéine-kinase* (“protein-kinase”): “*Enzyme qui phosphoryle les résidus sérine, thréonine ou tyrosine présents dans les protéines.*”¹⁷

On the other hand, we observe that most frequent terms are indeed *semantic neologisms*, i.e. terms transferred from one domain to another, with a meaning change. We find that BLOOM is unable to generate semantic neologisms, as its performance drops after 10^6 occurrences (Figure 4). For example, for *pression* “*marquage serré de l’adversaire en possession du ballon*”¹⁸, which metaphorically transfers the concept of *pressure* from physics to sports, the model generates the literal syntagm *marquage individuel* (“individual marking”).

Setting	ICL	Pre.	Suff.	Neo.	Native	Synt.
TERM	Random	71.5	86.2	<u>61.1</u>	14.8	87.7
TERM	Paradigm	<u>73.4</u>	<u>87.4</u>	59.8	<u>24.4</u>	<u>88.0</u>
DEF	Random	59.2	82.0	<u>39.5</u>	15.8	77.6
DEF	Co-hyponyms	<u>59.7</u>	<u>82.5</u>	36.7	<u>18.7</u>	<u>79.5</u>
DEF+TERM	Random	71.8	86.9	63.3	17.9	87.6
DEF+TERM	Fusion	74.8	88.4	65.3	26.5	88.8

Table 4: F1 scores of morphosyntactic processes prediction by BLOOM-7.1B on FranceTerme test set. The best overall results are in boldface while the best results in settings TERM and DEF are underlined.

6.5. Morphosyntactic Analysis

The multi-label classifier described in Section 4 allows us to analyze the morphological processes used to coin new terms. We compare the morphological processes of the models’ outputs with the corresponding reference (see Table 4). We find that, even when the output term is incorrect, the morphological analysis of the output term agrees mostly with the reference. For example, while *énantiomère* (“enantiomer”) does not match the reference *distomère* (“distomer”), both are neoclassical compounds. The only exception is for native compounds, which are rare in French: only 2.8% of EN native compounds are translated as native compounds into FR. Overall, these performance are in line with previous results (Table 2): our ICL selection strategies consistently improves the scores.

¹⁶“Enzyme that dephosphorylates previously phosphorylated serine, threonine or tyrosine residues in proteins”

¹⁷“Enzyme that phosphorylates serine, threonine or tyrosine residues present in proteins.”

¹⁸“close marking of opponents in possession of the ball”

6.6. Morphosyntactic Divergences

The multi-label classifier also enables us to evaluate the divergences between English source terms and their reference French counterparts. We study here how this divergence impacts the performance of the models. Given E and F , the sets of EN and FR morphosyntactic processes involved in the generation of a given term, respectively, we rely on the symmetric difference between these two sets to define a distance metric: $\Delta = |(E \cup F) \setminus (F \cap E)|$. We find that model performance is negatively correlated with this distance, especially when relying on the EN source term, see Figure 5. For example, the TERM model translates the syntagm of suffixation “homing head” using the same processes, resulting in *tête de guidage*, not matching the reference prefixation *autodirecteur*.

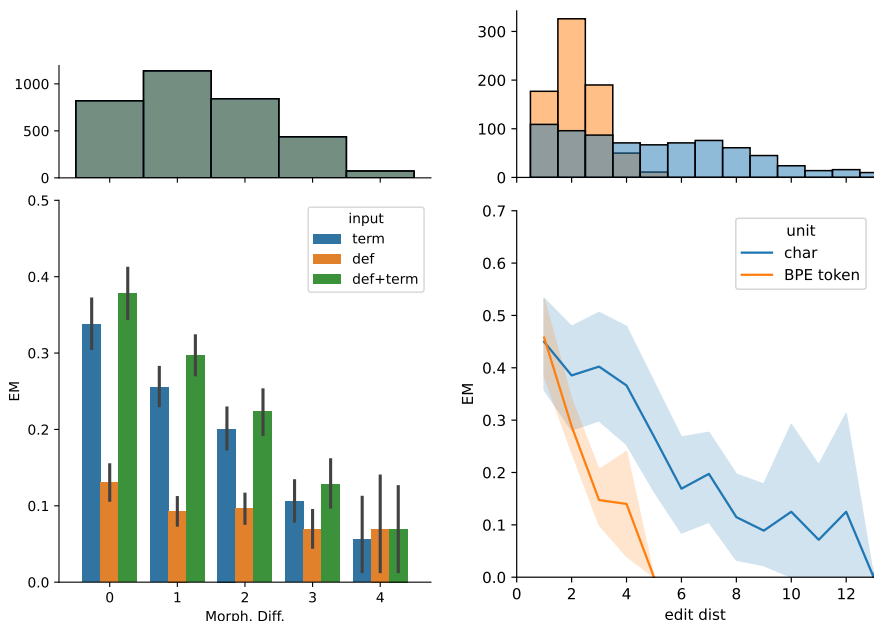


Figure 5: Exact Match (EM) of BLOOM-7.1B outputs w.r.t. morphosyntactic difference Δ between EN and FR processes, in the three usual settings with randomly selected ICL examples on FranceTerme’s test set (left). The upper part shows the number of examples for $\Delta \in [0, 4]$. Exact Match (EM) of BLOOM-7.1B (TERM) outputs w.r.t. edit distance between EN and FR monolexical terms, with randomly selected ICL examples on FranceTerme’s test set (right). The upper part displays the number of examples in each bin. Edit distance is at least 1 because loanwords were filtered out.

6.7. Translation or Orthographic Conversion?

We saw in Section 6.2 that setting TERM was much easier than DEF. We show that this is due to frequent surface similarities between EN and FR, which makes the translation akin to an orthographic conversion. We quantify this by computing the edit distance between EN and

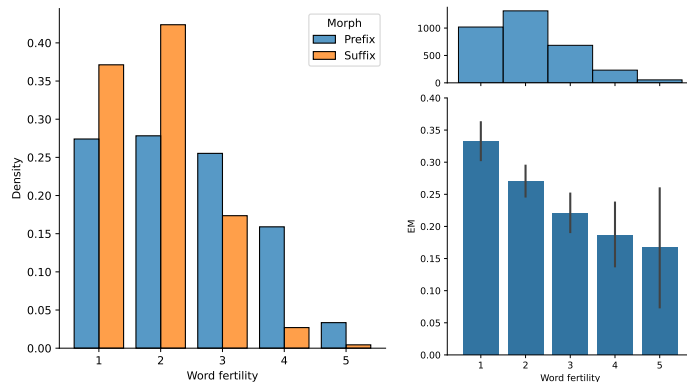


Figure 6: Distribution of word fertilities for prefixed and suffixed terms on FranceTerme’s test set (left). Density is normalized separately for prefixes and suffixes to ease visualization. Exact Match (EM) of BLOOM-7.1B (DEF+TERM) outputs w.r.t. word fertilities (right). The upper part shows the number of examples in each bin.

FR monolexical terms.¹⁹ Figure 5 shows that the performance in setting TERM is negatively correlated with the edit distance, while DEF does not suffer from this bias. For example, the model correctly predicts the following terms with an edit distance of 3 or less: *mycotoxin* → *mycotoxine*, *exocytosis* → *exocytose*, *iconomatic* → *iconomatique*. This result holds for both character-level and token-level edit distance. For token-level distance, we may assume that the model directly copies tokens from source to target. The examples above actually share the following tokens: “_my c oto”, “_ex”, and “_ic onom”, respectively.

6.8. Prefixation, Fertility, and BPE

BLOOM, as mBART and CroissantLLM, relies on BPE tokenization, like most LLMs [Gage, 1994, Sennrich et al., 2016]. While BPE circumvents out-of-vocabulary (OOVs) issues by splitting rare words into subwords, it only relies on character n-grams co-occurrences and rarely generates morphologically sound segmentations [Church, 2020]. When pre-tokenizing text on whitespace, tokens beginning a word bear a special mark “_”; without pre-tokenization, a whitespace will occur before each word start [Kudo and Richardson, 2018, Wolf et al., 2020]. This means that prefixations and suffixations are not treated equally, with two issues for prefixations: (i) even if segmented correctly, the base and derivation will not share any representation (e.g. “_collision” vs. “_pré collision”; Hofmann et al., 2021); (ii) most likely, the derived term will be over-segmented, as the occurrences of the base in word-internal position are too rare to warrant a dedicated vocabulary entry (e.g. “collision”). For our running example, *précollision* is split as “_préc oll ision”²⁰. Unlike suffixations which are often reasonably well segmented and share representations with their base (e.g. “_collision neur”).

Figure 6 shows that prefixed terms suffer from this BPE tokenization more than suffixed forms

¹⁹Doing so for polylexical terms would require more caution, because of syntactic divergences between EN and FR.

²⁰Note that these three tokens are not meaningful morphemes in French.

and have a much higher word fertility.²¹ Furthermore, in the same figure, we show that word fertility is negatively correlated with EM. For example, BLOOM fails to predict *téléconsultation* (segmented as “_tél éc ons ult ation”, although “_consultation” has a dedicated token).

7. Derivational Morphology

Large Language Models (LLMs) constitute a workhorse of modern Natural Language Processing applications, owing to their unprecedented ability to generate syntactically correct, semantically coherent, and pragmatically relevant utterances, responses to a wide array of queries, in a growing number of languages. As recent studies have shown, during their training process, LLMs also acquire some sort of morphological abilities, e.g., to generate inflected forms for known and unknown lemmas [Weissweiler et al., 2023] – at least when they follow regular morphological patterns (see also Hofmann et al., 2020, Mortensen et al., 2024). These abilities extend even to previously unknown languages, given that some examples of the targeted patterns are provided in the prompt [Tanzer et al., 2024, Zhang et al., 2024]. Such morphological knowledge is essential to achieve good performance in constrained (e.g., Machine Translation) as well as unconstrained text generation applications. The ability to manipulate and recombine substrings and to handle unknown word forms can be attributed to the use of subword vocabularies, e.g., relying on Byte Pair Encoding (BPE; Gage, 1994, Sennrich et al., 2016).

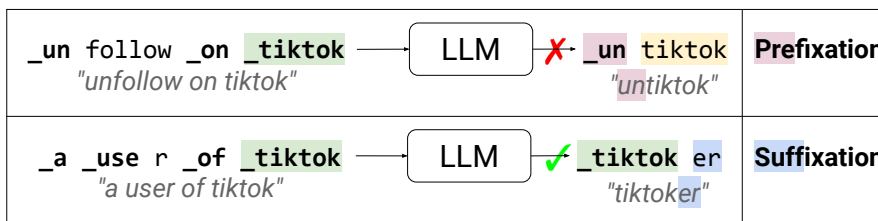


Figure 7: In BPE tokenization, marking word-initial tokens with “_” hinders the generation of prefixed forms (e.g., “_un tiktok”), as they do not share any token with their base (e.g., “_tiktok”). Identical tokens are highlighted in the same color.

In this contribution, we show that if BPE-based tokenizers enable morphological generalization, they do not handle all morphological processes equally well. The reason, we claim, is that BPE marks word-initial substrings with a special character “_”, to make tokenization reversible [Kudo and Richardson, 2018].²² Therefore, suffixed and prefixed forms are handled differently: once tokenized, suffixed forms such as “*tiktok er*” may share a subword with their base “*tiktok*”, implying also some semantic similarity. Crucially, this cannot happen with prefixed forms like

²¹Fertility is the number of tokens in a given form; for polylexical terms, we define word fertility as the maximum fertility over words occurring in the term.

²²Equivalently, Sennrich et al. [2016] marked intra-words with “@@”, e.g., “tiktok @@er”. Marking the end of words instead of their start would hinder suffixations. The issue is identical for all subword tokenizers that we are aware of: BPE, Unigram [Kudo, 2018], and WordPiece [Wu et al., 2016], as they all mark word-initial substrings to make tokenization reversible. A similar case is made by Hofmann et al. [2021] about WordPiece, as discussed in Section 2. We simply focus on BPE as it has now been widely adopted by all modern LLMs (e.g., GPT-4, [OpenAI, 2023]; Llama-3, [Llama Team, 2024]; and Gemma, [Gemma Team, 2024]).

“*untiktok*”, which, even assuming a morphologically plausible tokenization “*_un tiktok*”, are represented using a distinct token “*tiktok*”, unrelated to the base “*_tiktok*” (see Figure 7; Hofmann et al., 2020).²³

We present here experiments that highlight this problem in a very controlled setting. For this, we consider two regular affixation processes in English and French: negative prefixations (e.g., EN “*un-*”, FR “*in-*”) and adverbial suffixations (e.g., EN “*-ly*”, FR “*-ment*”). As both processes apply to adjectives, we can compare on a fair basis the capacities of LLMs to generate prefixations and suffixations of the same set of lexemes. Our experiments include both attested adjectival bases and nonce words. We find, across several LLM families and sizes, that (i) LLMs often fail to derive new words via prefixation compared to suffixation; (ii) the cases where prefixation is successful may be explained by the alignment between word-initial and word-internal embeddings of the same string (e.g., when “*_tiktok*” \approx “*tiktok*” in the embedding space), which is dependent on the model size and amount of pretraining data; (iii) this tendency can be mitigated with in-context learning (ICL), especially with a consistent selection of ICL examples; (iv) the issue disappears when using a morphological segmentation, which leads to near-perfect accuracy, for both prefixations and suffixations.

Derivational Morphology is central to the structure of the lexicon, so as to move away from the arbitrariness of the sign [De Saussure, 1916, Lieber, 2010, Corbin, 2012b]. Affixation is cross-linguistically the most common process that human languages use to derive new lexemes [Štekauer, 2012, Goethem, 2020]. For example, Turkish’s *-li* attaches to nouns to make personal nouns (e.g., *şehir* ‘city’ \rightarrow *şehirli* ‘city dweller’); Chinese’s *-xué* attaches to nouns to make nouns meaning ‘the study of X’ (e.g., *dòngwù* ‘animal’ \rightarrow *dòngwùxué* ‘zoology’); Samoan’s *fa’a-* attaches to nouns to make verbs meaning ‘make X’ (e.g., *goto* ‘sink’ \rightarrow *fa’agoto* ‘make sink’, Lieber, 2010). In this paper, we study English and French as mere examples to motivate our finding about BPE, which formally applies to any text in any language. Regular affixation processes are routinely used to form *neologisms*, new lexemes or terms, either in everyday conversations or in specific domains [Daille, 2017, Cartier et al., 2018].

Formally, *prefixation* operates at the beginning of a lexeme (e.g., “*untiktok*”), whereas *suffixation* applies at lexeme’ ends (e.g., “*tiktoker*”). This implies, as discussed above, that the two types of derivative will be handled differently by subword tokenizers. Affixation may additionally cause phonological or graphemic change(s), resulting in variation (*allomorphy*) in the surface realization of some lexemes [Lieber, 2010]. This is another cause of possible divergence between the tokenization of a base and a derived lexeme. In our experiments, we make sure to only consider cases of purely concatenative affixations,²⁴ as in the above examples, to isolate the tokenization challenge from other segmentation issues. In English and French, other differences, not developed here, between prefixation and suffixation are that the latter tends to play a more syntactic role (e.g., converting adjectives to adverbs with “*-ly*”) while the former holds a more semantic role (e.g., negating adjectives with “*un-*”).

²³Word-internal tokens only make their way to the tokenizer’s vocabulary if supported by enough prefixed forms in the training corpus; otherwise, such derivatives are over-segmented, e.g., “*_un tik t ok*”.

²⁴This means that valid morphological segmentations will always be either “<prefix> <base>” for prefixations or “<base> <suffix>” for suffixations.

8. Derivation Methods

Lang.	Affix	Definition
EN	un-	Not <base>
FR	<i>in-</i>	<i>Qui n'est pas <base></i>
EN	-ly	In a <base> manner
FR	<i>-ment</i>	<i>D'une manière <base></i>

Table 5: Affixations and associated definition templates

8.1. Definition to Word Generation

To measure differences in the way suffixations and prefixations are handled by LLMs, we consider the simple morphological task of generating a lexeme from its definition, framed here as a text-to-text problem [Brown et al., 2020, Raffel et al., 2020]. Given the definition of a lexeme (e.g., “*a user of tiktok*”), an LLM is prompted to generate the derivative (e.g., “*tiktoker*”), cf. Figure 7. Models are prompted in the same language as the definition (<def>) and the target lexeme, i.e. (i) EN: “<def> defines the term :”; (ii) FR: “<def> définit le terme :”. The expected continuation is the derived form. Definitions always include the base lexeme and unambiguously correspond to either a prefixed or a suffixed derivative (Table 5).

8.2. In-Context Learning

LLMs can further generalize to such tasks by leveraging In-Context Learning. Our early results suggested that LLMs were not too sensitive to the exact prompt formulation, but mostly leveraged ICL examples, consistently with prior work (e.g., Zebaze et al., 2024). We thus use five ICL examples in each prompt, formatted as above, separated by the three characters ###, which serve as end-of-sequence signal. Here is an example from the ADJ-EN dataset, using a single ICL example: “*Not pluvial defines the term : unpluvial ### Not lightfast defines the term :*” (the model should generate “*unlightfast*”).

We limit the number of ICL examples to five to keep a reasonable input length.²⁵ We compare two ICL selection methods: (i) Random sampling, examples can be either a prefixation or a suffixation; (ii) Morphological: sampling only prefix (resp. suffix) for prefix (resp. suffix) generation tasks.

8.3. Large Multilingual Models

We conduct experiments with three model families: BLOOM [BigScience et al., 2023], CroissantLLM-1.3B [Faysse et al., 2024], and Llama-2-7B [Touvron et al., 2023], including various sizes for BLOOM, ranging from 560M to 7.1B parameters. All models are multilingual and cover EN

²⁵Early experiments suggest that the difference between prefixes and suffixes is only stronger with fewer examples.

and FR to different degrees: BLOOM is highly multilingual, trained on 46 natural languages; CroissantLLM is bilingual, trained on an equal share of EN and FR; Llama-2 is mostly trained on EN (89.70%) but does include some FR (0.16%).

8.4. Segmentation

We compare two segmentation strategies:

(i) BPE, used by all studied LLMs. Keeping the same example as above, the base and derived word are tokenized as follow by BPE (for BLOOM but beginning of words are always marked by

BPE, regardless of the LLM): *pluvial* **_pluv** ial
unpluvial **_un pl uv** ial

Notice how the derived word does not include the tokens of its base.

(ii) Morphological segmentation, where we enforce that derived words in the ICL samples share tokens with their base by adding an extra space to the affix. In that case, the output is expected to be also space separated. For example:

un pluvial **_un _pluv** ial

8.5. Controlled Datasets

Dataset	Base	Prefixation	Suffixation
ADJ-FR	démontable	indémontable	démontablement
ADJ-EN	lightfast	unlightfast	lightfastly
PSEUDO-FR	géniable	ingéniable	généablement
PSEUDO-EN	orionful	unorionful	orionfully

Table 6: Examples of a base and its derivatives for each dataset

We perform controlled experiments, where each base has one derived prefixation and suffixation. We study two regular affixations that apply to adjectival bases: (i) negative prefixations (EN: “*un-*”, FR: “*in-*”); (ii) adverbial suffixations (EN: “*-ly*”, FR: “*-ment*”), paired with the definitions listed in Table 5, e.g.: (i) “*Not lightfast*” → “*unlightfast*”; (ii) “*In a lightfast manner*” → “*lightfastly*”.

We experiment with two sets of bases, in each language: (i) ADJ, attested adjectives from MorphyNet [Batsuren et al., 2021], which is built upon Wiktionary; (ii) PSEUDO, pseudo-words generated with UniPseudo [New et al., 2024] (see examples for each dataset in Table 6). As explained above, we restrict ourselves to purely concatenative affixation and avoid allomorphy phenomena using “morphotactic” rules described in Appendix E. MorphyNet has fewer samples in FR than EN, and FR morphotactics are more strict, so ADJ-FR contains 2,313 adjectival bases, i.e. 4,626 derived words (one prefixation and suffixation per base), while ADJ-EN contains 14,455 bases. Pseudo-adjectives are generated with UniPseudo, using a character n-gram model trained with attested adjectives. In each language, we first generate 5,000 nonce words of L letters, for $L \in \llbracket 6, 12 \rrbracket$. After filtering these with morphotactic rules, we obtain PSEUDO-FR (comprising

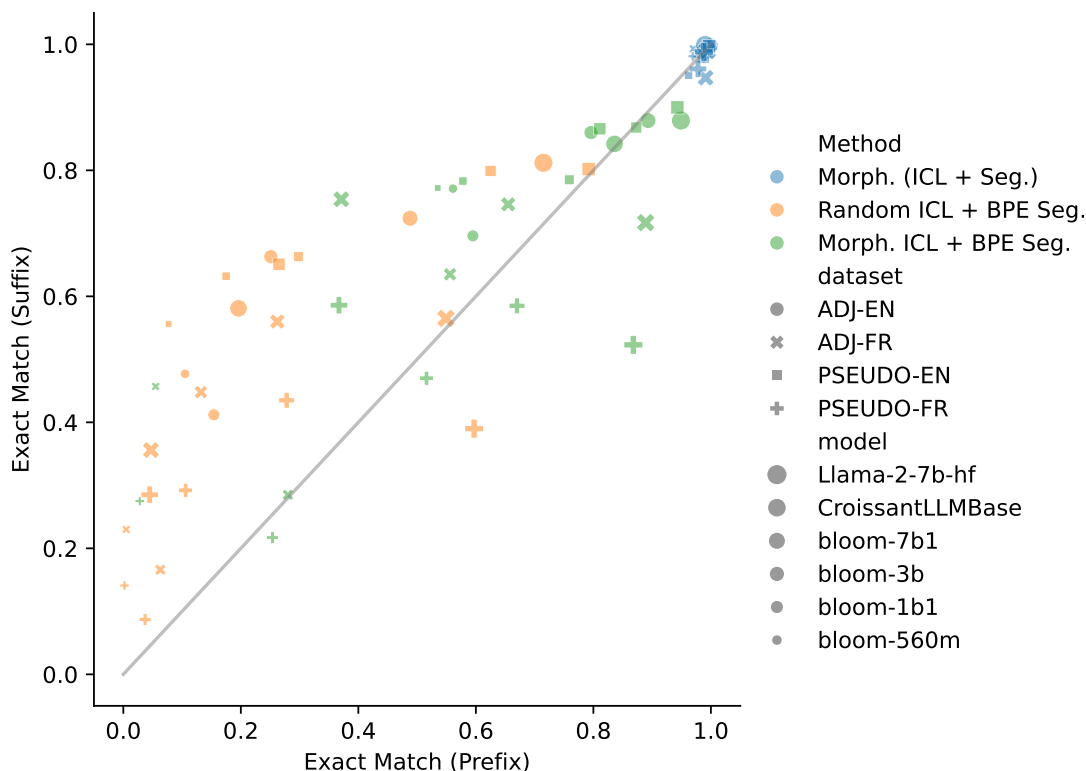


Figure 8: Exact match scores for prefixes vs. suffixes for four datasets (attested adjectival bases and pseudo-words, EN and FR; plotted with different shapes), three model families and four BLOOM model sizes ranging from 560M to 7.1B parameters (plotted with different sizes), according to ICL examples and segmentation method (different colors). Most points are above the line $y = x$, because suffixes are better generated than prefixes.

8,507 bases) and PSEUDO-EN (29,177 bases). The datasets are equally and randomly split in ICL-test splits, without overlap between bases.²⁶

9. Derivation Results

9.1. Prefixations vs. Suffixations

Figure 8 displays our main results on the four datasets with three different model families: BLOOM, CroissantLLM-1.3B, and Llama-2-7B (detailed scores are in Table 13 in Appendix F). We use Exact Match (EM), also known as *accuracy* to evaluate generation [Cotterell et al., 2016]. Clearly, with standard BPE segmentation, suffix generation is overall far superior to prefix generation (e.g., 26.2 EM for prefixes vs. 56.0 for suffixes, with BLOOM-7.1B on FR attested adjectives; above the $y = x$ line). Errors in prefixations also include cases of morphotactically

²⁶See Appendix D for implementation details and github.com/PaulLerner/neott for code and data.

incorrect forms, with the generated prefixes containing extraneous letters, dashes, or spaces (e.g., “*incgrandiose*”, “*in_onirique*”, or “*in_cognitive*”, with BLOOM-7.1B on FR attested adjectives). We also find that prefixes are sensitive to the choice of ICL examples: selecting only prefixes (resp. suffixes) for prefix (resp. suffix) prediction helps to reduce the gap (green vs. orange dots). This finding is consistent with Hofmann et al. [2024] who find that LLMs generalize through analogies rather than rules. We finally observe that Llama outperforms BLOOM and CroissantLLM for this task.

Morphological segmentation, on the other hand, solves the initial- vs. intra-word tokenization issue and yields near-perfect accuracy, for both prefixes and suffixes and all models (blue vs. green dots).²⁷ Figure 8 shows that even small versions of BLOOM, with 560M or 1.1B parameters (small dots), achieve near-perfect accuracy for prefixes with a morphological segmentation, when the corresponding Exact Match score was close to zero with BPE segmentation. BLOOM-7.1B is still able to correctly generate some prefixed form, probably due to its larger number of parameters. These results are consistent on the four datasets, i.e., for both attested adjectival bases and pseudo-words, in EN and FR.

9.2. Initial- vs. Intra-word Alignment

In this section we ask: how is it possible at all for BPE-based models to generate prefixations? We argue that, like for suffixations where the model simply needs to copy the base (e.g., “*_tiktok*”) and append a suffix (e.g., “*er*”),²⁸ for prefixations the model first needs to generate the prefix (e.g., “*_un*”) then an intra-word token whose representation is close to that of the base (e.g., “*tiktok*”), therefore to model the similarity between the two tokens (e.g., “*_tiktok*” \approx “*tiktok*”).

We find that, when a string has dedicated embeddings respectively covering word-initial and word-internal occurrences (e.g. “*_like*” and “*like*” or “*_vraisemblable*” and “*vraisemblable*”), both are often aligned, i.e., close in the embedding space. To evaluate this, for each pair pairs of vocabulary units of the form (*_x, x*) made of a word-initial and a matched word-internal vocabulary entry, we compute the cosine similarity of *_x* with all existing word-internal entries and measure the ability to retrieve the matched entry *x* with Precision@1 (P@1). Depending on the model, we find P@1 values ranging from 71.9 to 83.0, reported in Table 7. These values are well correlated with the EM scores for prefixes reported above (for the four BPE-based BLOOM models, we find Pearson $r = 0.639$, $p < 0.01$, across the four datasets).

This finding is consistent with Itzhak and Levy [2022], who find that word embeddings encode the string of characters that compose it; and Tytgat et al. [2024] who find that word embeddings are sensitive to surface similarities (e.g. edit distance).

Figure 9 shows that alignment increases with the number of tokens seen in training: for CroissantLLM, P@1 increases from 55.2 (after 300B tokens) up to 71.9 after 3T (again correlated with EM scores of prefixes of BPE-based models with Pearson $r = 0.338$, $p < 0.05$, across the four datasets). Therefore, gigantic amounts of data are used to implicitly learn an alignment that

²⁷Note that, while suffixations can always be tokenized by BPE as “<base> <suffix>” (e.g., “*_lightfast ly*”), the optimal tokenization (according to BPE) may not necessarily preserve the base (e.g., “*_light fastly*”). This explains why morphological segmentation also improves suffixation results.

²⁸Empirically, we find across all models and datasets that BPE-based models tend to copy the base tokens at a 63% rate in average when generating suffixations.

could be made explicit using morphological segmentation.

Model	# Pairs	# Intra	P@1
CroissantLLM-1.3B	3,771	14,296	71.9
BLOOM-7.1B	13,365*	111,326	76.3
Llama-2-7B	5,272	15,590	83.0

Table 7: Alignment between embeddings of word-initial types and the corresponding word-internal variant, for three models. *BLOOM’s vocabulary contains a lot of noise so we evaluate only on fully Latin strings (matching [A-Za-z]), otherwise P@1 would drop to 65.4.

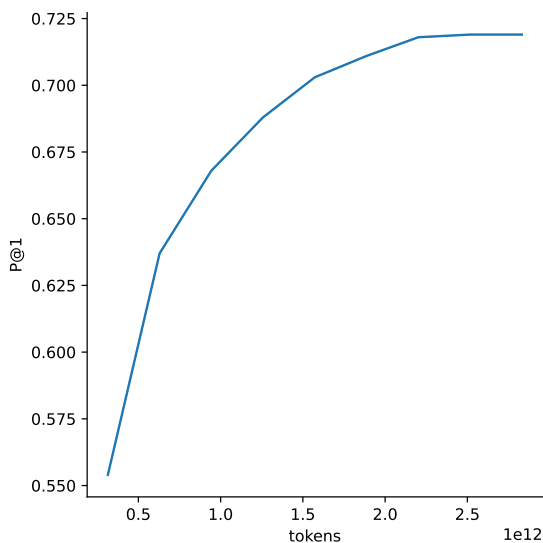


Figure 9: Alignment between embeddings of word-initial types and the corresponding word-internal variant for various checkpoints of CroissantLLM, according to the number of tokens used in training (in trillions).

10. Conclusion

Neologism translation is a challenge for standard MT systems that rely on parallel data. We propose a first effort to leverage definitions to accurately translate neologisms with Large Language Models. We found that LLMs were, to some extent, able to generate terms from their definition. Moreover, they can also combine the definition with the source term to translate it more accurately. As these models rely on In-Context Learning, we proposed to retrieve co-hyponyms or terms from the same derivation paradigm as the source term, which consistently improved results over

two datasets covering 13 diverse domains. The more terms are neological, which we assess from their corpus frequency, the more co-hyponyms retrieval improves performance.

However, we also pinpoint several limitations of these models: (i) they are sensitive to the similarity of source and target terms, either superficial or morphological; (ii) they rely on BPE tokenization, which is not morphologically sound and therefore impacts performance, especially for prefixations. This first limit is likely to be persistent but should be controlled in future work. The second limit, however, may be tackled using morphological segmentation [Smit et al., 2014, Batsuren et al., 2022] or character-based models [Cherry et al., 2018, Wang et al., 2024].

Our models may prove useful to enrich thesauri (e.g., providing suggestions to FranceTerme’s translators and lexicographers). Another obvious application is terminology-constrained MT [Semenov et al., 2023], with challenging research questions, especially for document-level MT, where one must find the right balance between terminological consistency and variation. Finally, in our future work we would also like to study the translation of terms in a more dynamic settings, considering new derivatives or complex noun phrases as they are coined or proposed to denote novel concepts in emerging research works. The latter category, which generalize our “syntagmatic compounds”, in particular, is likely to pose difficult translation problems, due to the opaqueness of the semantic relationships between their subparts.

11. Limitations

11.1. On Translation

Our study is limited to a single language pair, namely EN-FR, which, however, is highly demanding of such technology.²⁹ Moreover, French has a strong tradition of scientific writings as well as scientific terminology, as a large body of literature was published in French until a decline in the mid-20th century [Bacaër, 2019, Larivière and Riddles, 2021] and higher education is given in French. This is not the case for many low-resource languages due to a general tendency, observed in many countries, to use English for higher education, or for which scientific terminology simply does not exist [Gordin, 2015].

Furthermore, we conduct extensive experiments on EN-FR and find our results to be consistent across two datasets and 13 diverse domains. Our method could be extended to other languages with a tradition of scientific writing, e.g., Russian, Chinese, or German [Céspedes et al., 2024]. In the latter case, we could leverage multilingual thesauri such as IATE [Zorrilla-Agut and Fontenelle, 2019]). It would be particularly interesting to study other morphosyntactic processes than those of Section 4. We also plan to study the FR-EN direction, which is especially relevant for humanities and social sciences, where a large body of work is still published in French. However, many concepts in humanities are culture-dependent and challenging to translate.

As a first step to study definition-to-term generation, we assume that the definition of the term is available. In future work, we plan to extract definitions on the fly from source documents [Jin et al., 2013, Head et al., 2021, August et al., 2022, Huang et al., 2022]. Because of FranceTerme, experiments of Sections 6.2 and 6.3 were conducted with definitions in French (the target

²⁹Both France and Québec are pushing to disseminate scientific findings in multiple languages. See, e.g., Second French Plan for Open Science [Vidal, 2021].

language). However, we provide additional results in Appendix A with TERMIUM definitions machine-translated from English. Our findings of Sections 6.2 and 6.3 are consistent with these machine-translated definitions.

Studying neologisms is necessarily a race against the clock. We find that some terms in FranceTerme and TERMIUM already appear in large corpora such as OSCAR (cf. Section 6.4). However, most terms of FranceTerme appear less than 100 times in a 46 billion words corpus (i.e. 2×10^{-9} frequency). We recommend future work to conduct a similar analysis and focus on the performance on these rare terms. Our ICL method significantly improves performance on low-frequency terms. Also note that terms recorded in a thesauri show institutionalization, which is a step towards lexicalization [Hohenhaus, 2005]. Finally, we find that very frequent terms are indeed neologisms but have gone through semantic change. We plan to better assess this latter phenomenon by studying diachronic corpora [Ryskina et al., 2020].

11.2. On Morphology

We study only two languages: English and French. However, we focus on a formal issue of the BPE method, which would be identical for any text and therefore any language. We assume that this caveat would affect only more strongly less-resourced languages.

We are limited to one prefixation and one suffixation per language. This restriction was inevitable to allow stratified data generation (Section 8.5): the chosen negative prefixations and adverbial suffixations are very regular in English and French, both can be applied to any adjective. However, formally, the affixation process is identical regardless of the actual affix, be it *-ly*, *-ation*, or *-ical*.

Hofmann et al. [2021] had already pointed out the issue of marking beginning of words with WordPiece (instead of BPE), and also proposed to fix it by leveraging morphological segmentation. However, we propose a new framework (generation instead of classification) and provide additional analysis to understand the phenomenon through in-context learning (Figure 8), alignment of initial- and intra-word embeddings (Table 7), and amount of pretraining data (Figure 9). Additionally, we conduct extensive experiments on three different LLM families, while Hofmann et al. [2020, 2021] only use BERT [Devlin et al., 2019].

We propose to use morphological segmentation to solve the issue with the BPE tokenizer. This, however, is easier said than done: BPE has the advantage of being language-agnostic and therefore allows transfer learning between languages within a multilingual language model. In contrast, we are not aware of a morphological segmentation method that could be applied to all languages. It would most likely require a language identification pipeline followed by language-specific segmentation.

Acknowledgments

We thank Natalie Kübler, Mathilde Huguin, Alexandra Mestivier, and Lichao Zhu for their helpful feedback on an initial draft of this report. We also thank Ziqian Peng for providing mBART results and Felix Herron for his initial work on this topic. This work was performed using HPC resources from GENCI-IDRIS (Grant 2023-AD011014881).

Bibliography

- Julien Abadji, Pedro Ortiz Suarez, Laurent Romary, and Benoît Sagot. Towards a Cleaner Document-Oriented Multilingual Crawled Corpus. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 4344–4355, Marseille, France, June 2022. European Language Resources Association. URL <https://aclanthology.org/2022.lrec-1.463>. [Cited on pages 8 and 15.]
- Alfred V. Aho and Margaret J. Corasick. Efficient string matching: an aid to bibliographic search. *Commun. ACM*, 18(6):333–340, jun 1975. ISSN 0001-0782. doi: 10.1145/360825.360855. URL <https://doi.org/10.1145/360825.360855>. [Cited on page 45.]
- Duarte M. Alves, José Pombal, Nuno M. Guerreiro, Pedro H. Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, Pierre Colombo, José G. C. de Souza, and André F. T. Martins. Tower: An Open Multilingual Large Language Model for Translation-Related Tasks. In *COLM 2024*, 2024. URL <https://openreview.net/pdf?id=EHPns3hVkj>. [Cited on page 41.]
- Dany Amiot and Georgette Dal. La composition néoclassique en français et l’ordre des constituants. *La composition dans une perspective typologique*. Arras: Artois Presses Université, pages 89–113, 2008. URL https://www.academia.edu/download/36426597/Dal___Amiot_2008_Composition_neoclassique_Ordre_des_constituants.pdf. [Cited on page 9.]
- Pierre JL Arnaud. *Les composés timbre-poste*. Presses Universitaires Lyon, 2003. [Cited on page 9.]
- Mark Aronoff. Word formation in generative grammar. *Linguistic Inquiry Monographs Cambridge, Mass*, (1):1–134, 1976. [Cited on page 10.]
- Tal August, Katharina Reinecke, and Noah A. Smith. Generating Scientific Definitions with Controllable Complexity. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8298–8317, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.569. URL <https://aclanthology.org/2022.acl-long.569>. [Cited on pages 7 and 26.]
- Nicolas Bacaër. Quelques aspects de la disparition du français dans la recherche scientifique. *FIU Francophonie et innovation à l’université*, 1:16–27, 2019. URL <https://hal.science/hal-02268776>. [Cited on page 26.]
- Satanjeev Banerjee and Alon Lavie. METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments. In Jade Goldstein, Alon Lavie, Chin-Yew Lin, and Clare Voss, editors, *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan, June 2005. Association for Computational Linguistics. URL <https://aclanthology.org/W05-0909>. [Cited on page 12.]

- Khuyagbaatar Batsuren, Gábor Bella, and Fausto Giunchiglia. Morphynet: a large multilingual database of derivational and inflectional morphology. In *Proceedings of the 18th sigmorphon workshop on computational research in phonetics, phonology, and morphology*, pages 39–48, 2021. URL <https://aclanthology.org/2021.sigmorphon-1.5/>. [Cited on pages 22 and 43.]
- Khuyagbaatar Batsuren, Gábor Bella, Aryaman Arora, Viktor Martinovic, Kyle Gorman, Zdeněk Žabokrtský, Amarsanaa Ganbold, Sárka Dohnalová, Magda Sevcíková, Kateřina Pelegrinová, Fausto Giunchiglia, Ryan Cotterell, and Ekaterina Vylomova. The SIGMORPHON 2022 Shared Task on Morpheme Segmentation. In Garrett Nicolai and Eleanor Chodroff, editors, *Proceedings of the 19th SIGMORPHON Workshop on Computational Research in Phonetics, Phonology, and Morphology*, pages 103–116, Seattle, Washington, July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.sigmorphon-1.11. URL <https://aclanthology.org/2022.sigmorphon-1.11>. [Cited on pages 26 and 43.]
- Rachel Bawden and François Yvon. Investigating the translation performance of a large multilingual language model: the case of BLOOM. In Mary Nurminen, Judith Brenner, Maarit Koponen, Sirkku Latomaa, Mikhail Mikhailov, Frederike Schierl, Tharindu Ranasinghe, Eva Vanmassenhove, Sergi Alvarez Vidal, Nora Aranberri, Mara Nunziatini, Carla Parra Escartín, Mikel Forcada, Maja Popovic, Carolina Scarton, and Helena Moniz, editors, *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 157–170, Tampere, Finland, June 2023. European Association for Machine Translation. URL <https://aclanthology.org/2023.eamt-1.16>. [Cited on page 11.]
- BigScience, Teven Le Scao, and et al. BLOOM: A 176B-Parameter Open-Access Multilingual Language Model, June 2023. URL <http://arxiv.org/abs/2211.05100>. arXiv:2211.05100 [cs]. [Cited on pages 11 and 21.]
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language Models are Few-Shot Learners. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf>. [Cited on pages 10 and 21.]
- Maria Teresa Cabré. *Terminology: Theory, methods, and applications*, volume 1. John Benjamins Publishing, 1999. [Cited on page 6.]
- Emmanuel Cartier, Jean-François Sablayrolles, Najet Boutmgharine, John Humbley, Massimo Bertocci, Christine Jacquet-Pfau, Natalie Kübler, and Giovanni Tallarico. Détection automatique, description linguistique et suivi des néologismes en corpus: point d’étape sur les tendances du français contemporain. In *6e Congrès Mondial de Linguistique Française-Université*

- de Mons, Belgique, 9-13 juillet 2018*, volume 46, pages 1–20. EDP Sciences, 2018. URL <https://iris.univr.it/handle/11562/983330>. [Cited on pages 6 and 20.]
- Colin Cherry, George Foster, Ankur Bapna, Orhan Firat, and Wolfgang Macherey. Revisiting Character-Based Neural Machine Translation with Capacity and Compression. In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun’ichi Tsujii, editors, *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4295–4305, Brussels, Belgium, October 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1461. URL <https://aclanthology.org/D18-1461>. [Cited on page 26.]
- Martin S. Chodorow, Roy J. Byrd, and George E. Heidorn. Extracting Semantic Hierarchies From a Large On-Line Dictionary. In *23rd Annual Meeting of the Association for Computational Linguistics*, pages 299–304, Chicago, Illinois, USA, July 1985. Association for Computational Linguistics. doi: 10.3115/981210.981247. URL <https://aclanthology.org/P85-1037>. [Cited on page 7.]
- Kenneth Ward Church. Emerging trends: Subwords, seriously? *Natural Language Engineering*, 26(3):375–382, May 2020. ISSN 1351-3249, 1469-8110. doi: 10.1017/S1351324920000145. URL <https://www.cambridge.org/core/journals/natural-language-engineering/article/emerging-trends-subwords-seriously/619F28526E833A6B623E6D2009F37B82>. Publisher: Cambridge University Press. [Cited on page 18.]
- Maria Copot and Olivier Bonami. Baseless derivation: the behavioural reality of derivational paradigms. *Cognitive Linguistics*, 35(2):221–250, May 2024. ISSN 1613-3641. doi: 10.1515/cog-2023-0018. URL <https://www.degruyter.com/document/doi/10.1515/cog-2023-0018/html>. Publisher: De Gruyter Mouton. [Cited on page 9.]
- Danielle Corbin. *Morphologie dérivationnelle et structuration du lexique*, volume 193. Walter de Gruyter, 2012a. [Cited on pages 9 and 10.]
- Danielle Corbin. *Morphologie dérivationnelle et structuration du lexique*. Walter de Gruyter, October 2012b. ISBN 978-3-11-135838-3. Google-Books-ID: AYwjAAAQBAJ. [Cited on page 20.]
- Ryan Cotterell, Christo Kirov, John Sylak-Glassman, David Yarowsky, Jason Eisner, and Mans Hulden. The SIGMORPHON 2016 shared Task—Morphological reinflection. In Micha Elsner and Sandra Kuebler, editors, *Proceedings of the 14th SIGMORPHON Workshop on Computational Research in Phonetics, Phonology, and Morphology*, pages 10–22, Berlin, Germany, August 2016. Association for Computational Linguistics. doi: 10.18653/v1/W16-2002. URL <https://aclanthology.org/W16-2002>. [Cited on pages 12 and 23.]
- Lucía Céspedes, Diego Kozłowski, Carolina Pradier, Maxime Holmberg Sainte-Marie, Natsumi Solange Shokida, Pierre Benz, Constance Poitras, Anton Boudreau Ninkov, Saeideh Ebrahimi, Philips Ayeni, Sarra Filali, Bing Li, and Vincent Larivière. Evaluating the Linguistic Coverage of OpenAlex: An Assessment of Metadata Accuracy and Completeness, September 2024. URL <http://arxiv.org/abs/2409.10633>. arXiv:2409.10633. [Cited on page 26.]

- Béatrice Daille and Emmanuel Morin. French-English terminology extraction from comparable corpora. In *Second International Joint Conference on Natural Language Processing: Full Papers*, 2005. doi: 10.1007/11562214_62. URL <https://aclanthology.org/I05-1062>. [Cited on page 9.]
- Béatrice Daille. *Term Variation in Specialised Corpora: Characterisation, automatic discovery and applications*, volume 19 of *Terminology and Lexicography Research and Practice*. John Benjamins Publishing Company, Amsterdam, August 2017. ISBN 978-90-272-2343-2 978-90-272-6535-7. doi: 10.1075/tlrp.19. URL <http://www.jbe-platform.com/content/books/9789027265357>. [Cited on pages 8, 9, and 20.]
- Georgette Dal. Productivité morphologique: définitions et notions connexes. *Langue française*, pages 3–23, 2003a. [Cited on page 10.]
- Georgette Dal. Analogie et lexique construit: quelles preuves? 2003b. URL <https://lilloa.univ-lille.fr/handle/20.500.12210/65173>. Publisher: Toulouse: Université de Toulouse-le-Mirail, 1979-2006. [Cited on page 10.]
- Ferdinand De Saussure. *Cours de linguistique générale*, volume 1. Otto Harrassowitz Verlag (1989 reedition), 1916. [Cited on page 20.]
- Estelle Delpech, Béatrice Daille, Emmanuel Morin, and Claire Lemaire. Extraction of Domain-Specific Bilingual Lexicon from Comparable Corpora: Compositional Translation and Ranking. In Martin Kay and Christian Boitet, editors, *Proceedings of COLING 2012*, pages 745–762, Mumbai, India, December 2012. The COLING 2012 Organizing Committee. URL <https://aclanthology.org/C12-1046>. [Cited on page 8.]
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Tamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://aclanthology.org/N19-1423>. [Cited on page 27.]
- Manuel Faysse, Patrick Fernandes, Nuno Guerreiro, António Loison, Duarte Alves, Caio Corro, Nicolas Boizard, João Alves, Ricardo Rei, Pedro Martins, Antoni Bigata Casademunt, François Yvon, André Martins, Gautier Viaud, Céline Hudelot, and Pierre Colombo. CroissantLLM: A Truly Bilingual French-English Language Model, February 2024. URL <http://arxiv.org/abs/2402.00786>. arXiv:2402.00786 [cs]. [Cited on pages 11 and 21.]
- Bernard Fradin. *Nouvelles approches en morphologie*. PUF, 2015. URL <https://books.google.com/books?hl=en&lr=&id=Hd0JCwAAQBAJ&oi=fnd&pg=PP4&dq=info:zgdXhQrVrh0J:scholar.google.com&ots=KtybHbgDvc&sig=tKFm8iA8aSeGy7Q80V7WjkgjN8g>. [Cited on page 10.]

- Philip Gage. A New Algorithm for Data Compression. *Computer Users Journal*, 12(2):23–38, February 1994. ISSN 0898-9788. URL https://www.derczynski.com/papers/archive/BPE_Gage.pdf. Place: USA Publisher: R & D Publications, Inc. [Cited on pages 7, 18, and 19.]
- Google Gemma Team. Gemma 2: Improving Open Language Models at a Practical Size, October 2024. URL <http://arxiv.org/abs/2408.00118>. arXiv:2408.00118. [Cited on page 19.]
- Kristel Van Goethem. Affixation in morphology, 07 2020. URL <https://oxfordre.com/linguistics/view/10.1093/acrefore/9780199384655.001.0001/acrefore-9780199384655-e-678>. [Cited on page 20.]
- Michael D. Gordin. *Scientific Babel: How Science Was Done Before and After Global English*. University of Chicago Press, April 2015. ISBN 978-0-226-00032-9. Google-Books-ID: UrnnBgAAQBAJ. [Cited on pages 6 and 26.]
- Rejwanul Haque, Mohammed Hasanuzzaman, and Andy Way. Analysing terminology translation errors in statistical and neural machine translation. *Machine Translation*, 34(2-3):149–195, September 2020. ISSN 0922-6567, 1573-0573. doi: 10.1007/s10590-020-09251-z. URL <https://link.springer.com/10.1007/s10590-020-09251-z>. [Cited on page 7.]
- Tianxing He, Jingyu Zhang, Tianle Wang, Sachin Kumar, Kyunghyun Cho, James Glass, and Yulia Tsvetkov. On the blind spots of model-based evaluation metrics for text generation. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12067–12097, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.674. URL <https://aclanthology.org/2023.acl-long.674>. [Cited on page 12.]
- Andrew Head, Kyle Lo, Dongyeop Kang, Raymond Fok, Sam Skjonsberg, Daniel S. Weld, and Marti A. Hearst. Augmenting Scientific Papers with Just-in-Time, Position-Sensitive Definitions of Terms and Symbols. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, CHI '21, pages 1–18, New York, NY, USA, May 2021. Association for Computing Machinery. ISBN 978-1-4503-8096-6. doi: 10.1145/3411764.3445648. URL <https://dl.acm.org/doi/10.1145/3411764.3445648>. [Cited on pages 7 and 26.]
- Felix Hill, Kyunghyun Cho, Anna Korhonen, and Yoshua Bengio. Learning to Understand Phrases by Embedding the Dictionary. *Transactions of the Association for Computational Linguistics*, 4:17–30, 02 2016. ISSN 2307-387X. doi: 10.1162/tacl_a_00080. URL https://doi.org/10.1162/tacl_a_00080. [Cited on page 8.]
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. Training Compute-Optimal Large Language Models, March 2022. URL <http://arxiv.org/abs/2203.15556>. arXiv:2203.15556 [cs]. [Cited on page 11.]

- Valentin Hofmann, Janet Pierrehumbert, and Hinrich Schütze. DagoBERT: Generating derivational morphology with a pretrained language model. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3848–3861, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.316. URL <https://aclanthology.org/2020.emnlp-main.316>. [Cited on pages 8, 19, 20, and 27.]
- Valentin Hofmann, Janet Pierrehumbert, and Hinrich Schütze. Superbizarre Is Not Superb: Derivational Morphology Improves BERT’s Interpretation of Complex Words. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *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)*, pages 3594–3608, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.279. URL <https://aclanthology.org/2021.acl-long.279>. [Cited on pages 8, 18, 19, and 27.]
- Valentin Hofmann, Leonie Weissweiler, David Mortensen, Hinrich Schütze, and Janet Pierrehumbert. Derivational Morphology Reveals Analogical Generalization in Large Language Models, November 2024. URL <http://arxiv.org/abs/2411.07990>. arXiv:2411.07990. [Cited on page 24.]
- Peter Hohenhaus. Lexicalization and institutionalization. In *Handbook of word-formation*, pages 353–373. Springer, 2005. [Cited on page 27.]
- Jan Holeš. Quels termes pour communiquer ? Autour des néologismes officiels dans le domaine de la communication sur FranceTerme. *Çédille: Revista de Estudios Franceses*, (25):423–441, 2024. ISSN 1699-4949. URL <https://dialnet.unirioja.es/servlet/articulo?codigo=9568687>. Publisher: Asociación de Francesistas de la Universidad Español Section: Çédille: Revista de Estudios Franceses. [Cited on page 13.]
- Jie Huang, Hanyin Shao, Kevin Chen-Chuan Chang, Jinjun Xiong, and Wen-mei Hwu. Understanding Jargon: Combining Extraction and Generation for Definition Modeling. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3994–4004, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.266. URL <https://aclanthology.org/2022.emnlp-main.266>. [Cited on pages 7 and 26.]
- Itay Itzhak and Omer Levy. Models In a Spelling Bee: Language Models Implicitly Learn the Character Composition of Tokens. In Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz, editors, *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5061–5068, Seattle, United States, July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.373. URL <https://aclanthology.org/2022.naacl-main.373>. [Cited on page 24.]
- Yiping Jin, Min-Yen Kan, Jun Ping Ng, and Xiangnan He. Mining scientific terms and their definitions: A study of the ACL anthology. In *Proceedings of the 2013 Conference on*

- Empirical Methods in Natural Language Processing*, pages 780–790, 2013. URL <https://aclanthology.org/D13-1073.pdf>. [Cited on pages 7 and 26.]
- Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. Bag of Tricks for Efficient Text Classification. In Mirella Lapata, Phil Blunsom, and Alexander Koller, editors, *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 427–431, Valencia, Spain, April 2017. Association for Computational Linguistics. URL <https://aclanthology.org/E17-2068>. [Cited on pages 10 and 42.]
- Taku Kudo. Subword Regularization: Improving Neural Network Translation Models with Multiple Subword Candidates. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 66–75, Melbourne, Australia, July 2018. Association for Computational Linguistics. doi: 10.18653/v1/P18-1007. URL <https://aclanthology.org/P18-1007>. [Cited on page 19.]
- Taku Kudo and John Richardson. SentencePiece: A simple and language independent subword tokenizer and detokenizer for Neural Text Processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium, November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-2012. URL <https://aclanthology.org/D18-2012>. [Cited on pages 18 and 19.]
- Vincent Larivière and Amanda Riddles. Langues de diffusion des connaissances: quelle place reste-t-il pour le français. *Magazine de l’Acfas*, 2021. [Cited on pages 6 and 26.]
- Audrey Laroche and Philippe Langlais. Revisiting context-based projection methods for term-translation spotting in comparable corpora. In *Proceedings of the 23rd international conference on computational linguistics (Coling 2010)*, pages 617–625, 2010. URL <https://aclanthology.org/C10-1070.pdf>. [Cited on page 8.]
- Hugo Laurençon, Lucile Saulnier, Thomas Wang, Christopher Akiki, Albert Villanova del Moral, Teven Le Scao, Leandro Von Werra, Chenghao Mou, Eduardo González Ponferrada, and Huu Nguyen. The BigScience ROOTS corpus: A 1.6 tb composite multilingual dataset. *Advances in Neural Information Processing Systems*, 35: 31809–31826, 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/hash/ce9e92e3de2372a4b93353eb7f3dc0bd-Abstract-Datasets_and_Benchmarks.html. [Cited on page 15.]
- Paul Lerner and François Yvon. Unlike “Likely”, “Unlike” is Unlikely: BPE-based Segmentation hurts Morphological Derivations in LLMs. In *Proceedings of the 31st International Conference on Computational Linguistics*. International Committee on Computational Linguistics, 2025a. [Cited on page 7.]
- Paul Lerner and François Yvon. Towards the Machine Translation of Scientific Neologisms. In *Proceedings of the 31st International Conference on Computational Linguistics*. International Committee on Computational Linguistics, 2025b. [Cited on page 7.]

- Rochelle Lieber. *Introducing morphology*. Cambridge University Press, Cambridge, 2010. ISBN 978-0-511-77018-0. OCLC: 650278652. [Cited on pages 9, 10, and 20.]
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. RoBERTa: A Robustly Optimized BERT Pretraining Approach. September 2019. URL <https://openreview.net/forum?id=SyxSOT4tvS>. [Cited on page 11.]
- Zequn Liu, Shukai Wang, Yiyang Gu, Ruiyi Zhang, Ming Zhang, and Sheng Wang. Graphine: A Dataset for Graph-aware Terminology Definition Generation. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3453–3463, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.278. URL <https://aclanthology.org/2021.emnlp-main.278>. [Cited on page 6.]
- Meta Llama Team. The Llama 3 Herd of Models, 2024. [Cited on page 19.]
- Alizée Lombard and Richard Huyghe. Catégorisation comme néologisme et sentiment des locuteurs. *Langue française*, 207(3):123–138, 2020. ISSN 0023-8368. doi: 10.3917/lf.207.0123. URL <https://www.cairn.info/revue-langue-francaise-2020-3-page-123.htm>. Place: Paris Publisher: Armand Colin. [Cited on pages 9 and 15.]
- Elisa Mattiello. *Analogy in word-formation: A study of English neologisms and occasionalisms*, volume 309. Walter de Gruyter GmbH & Co KG, 2017. URL <https://books.google.com/books?hl=fr&lr=&id=RYkIDwAAQBAJ&oi=fnd&pg=PT4&ots=pg685jtz-b&sig=fb4RPrGtQOzZqPASXTX2bi41en0>. [Cited on page 10.]
- Simonetta Montemagni and Lucy Vanderwende. Structural patterns vs. string patterns for extracting semantic information from dictionaries. In *COLING 1992 Volume 2: The 14th International Conference on Computational Linguistics*, 1992. URL <https://aclanthology.org/C92-2083>. [Cited on page 7.]
- David R. Mortensen, Valentina Izrailevitch, Yunze Xiao, Hinrich Schütze, and Leonie Weissweiler. Verbing weirds language (models): Evaluation of English zero-derivation in five LLMs. In Nicoletta Calzolari, Min-Yen Kan, Veronique Hoste, Alessandro Lenci, Sakriani Sakti, and Nianwen Xue, editors, *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 17359–17364, Torino, Italia, May 2024. ELRA and ICCL. URL <https://aclanthology.org/2024.lrec-main.1508>. [Cited on page 19.]
- Alberto Muñoz-Ortiz, Carlos Gómez-Rodríguez, and David Vilares. Contrasting Linguistic Patterns in Human and LLM-Generated News Text. *Artificial Intelligence Review*, 57(10):265, August 2024. ISSN 1573-7462. doi: 10.1007/s10462-024-10903-2. URL <https://link.springer.com/10.1007/s10462-024-10903-2>. [Cited on page 8.]

- Fiammetta Namer. Automatiser l'analyse morpho-sémantique non affixale: le système DériF. *Cahiers de grammaire*, 28:31–48, 2003. URL <http://w3.erss.univ-tlse2.fr/publications/CDG/28/CG28-3-Namer.pdf>. [Cited on page 9.]
- Boris New, Jessica Bourgin, Julien Barra, and Christophe Pallier. UniPseudo: A universal pseudoword generator. *Quarterly Journal of Experimental Psychology*, 77(2):278–286, February 2024. ISSN 1747-0218. doi: 10.1177/17470218231164373. URL <https://doi.org/10.1177/17470218231164373>. Publisher: SAGE Publications. [Cited on page 22.]
- Thanapon Noraset, Chen Liang, Larry Birnbaum, and Doug Downey. Definition Modeling: Learning to Define Word Embeddings in Natural Language. *Proceedings of the AAAI Conference on Artificial Intelligence*, 31(1), February 2017. ISSN 2374-3468. doi: 10.1609/aaai.v31i1.10996. URL <https://ojs.aaai.org/index.php/AAAI/article/view/10996>. Number: 1. [Cited on page 7.]
- Byung-Doh Oh and William Schuler. Leading whitespaces of language models' subword vocabulary pose a confound for calculating word probabilities. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 3464–3472, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.202. URL <https://aclanthology.org/2024.emnlp-main.202>. [Cited on page 8.]
- OpenAI. GPT-4 Technical Report, March 2023. URL <http://arxiv.org/abs/2303.08774>. arXiv:2303.08774 [cs]. [Cited on page 19.]
- Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. fairseq: A fast, extensible toolkit for sequence modeling. In Waleed Ammar, Annie Louis, and Nasrin Mostafazadeh, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 48–53, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-4009. URL <https://aclanthology.org/N19-4009>. [Cited on page 45.]
- Arjun Panickssery, Samuel R. Bowman, and Shi Feng. LLM Evaluators Recognize and Favor Their Own Generations, April 2024. URL <http://arxiv.org/abs/2404.13076>. arXiv:2404.13076 [cs]. [Cited on page 12.]
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a Method for Automatic Evaluation of Machine Translation. In Pierre Isabelle, Eugene Charniak, and Dekang Lin, editors, *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA, July 2002. Association for Computational Linguistics. doi: 10.3115/1073083.1073135. URL <https://aclanthology.org/P02-1040>. [Cited on page 12.]
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An

- Imperative Style, High-Performance Deep Learning Library. *Advances in Neural Information Processing Systems*, 32, 2019. URL <https://papers.nips.cc/paper/2019/hash/dbbca288fee7f92f2bfa9f7012727740-Abstract.html>. [Cited on page 45.]
- Mojca Pecman. Tentativeness in term formation: A study of neology as a rhetorical device in scientific papers. *Terminology. International Journal of Theoretical and Applied Issues in Specialized Communication*, 18(1):27–58, January 2012. ISSN 0929-9971, 1569-9994. doi: 10.1075/term.18.1.03pec. URL <https://www.jbe-platform.com/content/journals/10.1075/term.18.1.03pec>. Publisher: John Benjamins. [Cited on page 13.]
- Ziqian Peng, Rachel Bawden, and François Yvon. À propos des difficultés de traduire automatiquement de longs documents. In *Actes de la conférence conjointe JEP-TALN-RECITAL 2024*, Toulouse, France, 2024. [Cited on pages 12 and 45.]
- Mohammad Taher Pilehvar. On the importance of distinguishing word meaning representations: A case study on reverse dictionary mapping. In Jill Burstein, Christy Doran, and Tamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2151–2156, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1222. URL <https://aclanthology.org/N19-1222>. [Cited on page 8.]
- Tiago Pimentel and Clara Meister. How to compute the probability of a word. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 18358–18375, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.1020. URL <https://aclanthology.org/2024.emnlp-main.1020>. [Cited on page 8.]
- Matt Post. A Call for Clarity in Reporting BLEU Scores. In Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Christof Monz, Matteo Negri, Aurélie Névél, Mariana Neves, Matt Post, Lucia Specia, Marco Turchi, and Karin Verspoor, editors, *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium, October 2018. Association for Computational Linguistics. doi: 10.18653/v1/W18-6319. URL <https://aclanthology.org/W18-6319>. [Cited on page 45.]
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *Journal of Machine Learning Research*, 21:1–67, 2020. [Cited on pages 10 and 21.]
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, 2016. [Cited on page 12.]
- Ayla Rigouts Terryn, Véronique Hoste, and Els Lefever. In no uncertain terms: a dataset for monolingual and multilingual automatic term extraction from comparable corpora. *Language*

- Resources and Evaluation*, 54(2):385–418, June 2020. ISSN 1574-0218. doi: 10.1007/s10579-019-09453-9. URL <https://doi.org/10.1007/s10579-019-09453-9>. [Cited on page 8.]
- Dimitrios Roussis, Vassilis Papavassiliou, Prokopis Prokopidis, Stelios Piperidis, and Vassilis Katsouros. SciPar: A collection of parallel corpora from scientific abstracts. In Nicoletta Calzolari, Frédéric Béchet, Philippe Blache, Khalid Choukri, Christopher Cieri, Thierry Declerck, Sara Goggi, Hitoshi Isahara, Bente Maegaard, Joseph Mariani, H  l  ne Mazo, Jan Odijk, and Stelios Piperidis, editors, *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 2652–2657, Marseille, France, June 2022. European Language Resources Association. URL <https://aclanthology.org/2022.lrec-1.284>. [Cited on page 12.]
- Maria Ryskina, Ella Rabinovich, Taylor Berg-Kirkpatrick, David R. Mortensen, and Yulia Tsvetkov. Where New Words Are Born: Distributional Semantic Analysis of Neologisms and Their Semantic Neighborhoods. In *Proceedings of the Society for Computation in Linguistics 2020*, pages 367–376, 2020. URL <https://aclanthology.org/2020.scil-1.43.pdf>. [Cited on page 27.]
- Kirill Semenov, Vil  m Zouhar, Tom Kocmi, Dongdong Zhang, Wangchunshu Zhou, and Yuchen Eleanor Jiang. Findings of the WMT 2023 Shared Task on Machine Translation with Terminologies. In *Proceedings of the Eighth Conference on Machine Translation*, pages 663–671, Singapore, 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.wmt-1.54. URL <https://aclanthology.org/2023.wmt-1.54>. [Cited on pages 7 and 26.]
- Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural Machine Translation of Rare Words with Subword Units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany, August 2016. Association for Computational Linguistics. doi: 10.18653/v1/P16-1162. URL <https://aclanthology.org/P16-1162>. [Cited on pages 18 and 19.]
- Peter Smit, Sami Virpioja, Stig-Arne Gr  nroos, and Mikko Kurimo. Morfessor 2.0: Toolkit for statistical morphological segmentation. In Shuly Wintner, Marko Tadi  c, and Bogdan Babych, editors, *Proceedings of the Demonstrations at the 14th Conference of the European Chapter of the Association for Computational Linguistics*, pages 21–24, Gothenburg, Sweden, April 2014. Association for Computational Linguistics. doi: 10.3115/v1/E14-2006. URL <https://aclanthology.org/E14-2006>. [Cited on page 26.]
- Pavol   tekauer. *Word-formation in the World’s Languages: A Typological Survey*. Cambridge University Press, 2012. [Cited on page 20.]
- Yuqing Tang, Chau Tran, Xian Li, Peng-Jen Chen, Naman Goyal, Vishrav Chaudhary, Jiatao Gu, and Angela Fan. Multilingual translation from denoising pre-training. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3450–3466, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.findings-acl.304. URL <https://aclanthology.org/2021.findings-acl.304>. [Cited on page 12.]

Garrett Tanzer, Mirac Suzgun, Eline Visser, Dan Jurafsky, and Luke Melas-Kyriazi. A Benchmark for Learning to Translate a New Language from One Grammar Book. In *ICLR 2024*, October 2024. URL <https://openreview.net/forum?id=tbVWug9f2h>. [Cited on page 19.]

Michela Tonti. *Le phraséotérme à la confluence de la langue naturelle, de la langue de spécialité et des néoformations*. « Ajustement cosmétique », « injonction de diversité », « fonction contrôle gestion garde-fou » et bien d'autres, pages 151–188. De Gruyter, Berlin, Boston, 2023. ISBN 9783110749854. doi: doi:10.1515/9783110749854-009. URL <https://doi.org/10.1515/9783110749854-009>. [Cited on page 13.]

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open Foundation and Fine-Tuned Chat Models, July 2023. URL <http://arxiv.org/abs/2307.09288>. arXiv:2307.09288 [cs]. [Cited on page 21.]

Delphine Tribout. *Les conversions de nom à verbe et de verbe à nom en français*. PhD thesis, Université Paris Diderot (Paris 7), 2010. [Cited on page 9.]

Thinh Truong, Yulia Otmakhova, Karin Verspoor, Trevor Cohn, and Timothy Baldwin. Revisiting subword tokenization: A case study on affixal negation in large language models. In Kevin Duh, Helena Gomez, and Steven Bethard, editors, *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 5082–5095, Mexico City, Mexico, June 2024. Association for Computational Linguistics. URL <https://aclanthology.org/2024.naacl-long.284>. [Cited on page 8.]

Julie Tytgat, Guillaume Wisniewski, and Adrien Betrancourt. Évaluation de la Similarité Textuelle : Entre Sémantique et Surface dans les Représentations Neuronales. In BALAGUER, Mathieu, BENDAHDAN, Nihed, HO-DAC, Lydia-Mai, MAUCLAIR, Julie, MORENO, Jose G, PINQUIER, and Julien, editors, *35èmes Journées d'Études sur la Parole (JEP 2024) 31ème Conférence sur le Traitement Automatique des Langues Naturelles (TALN 2024) 26ème Rencontre des Étudiants Chercheurs en Informatique pour le Traitement Automatique des Langues (RECITAL 2024)*, pages 85–96, Toulouse, France, July 2024. ATALA & AFPC. URL <https://inria.hal.science/hal-04623010>. [Cited on page 24.]

- Frédérique Vidal. Second French Plan for Open Science, 2021. URL <https://www.enseignementsup-recherche.gouv.fr/sites/default/files/2021-10/second-french-plan-for-open-science-13715.pdf>. [Cited on page 26.]
- Junxiong Wang, Tushaar Gangavarapu, Jing Nathan Yan, and Alexander M. Rush. MambaByte: Token-free Selective State Space Model, January 2024. URL <http://arxiv.org/abs/2401.13660>. arXiv:2401.13660 [cs]. [Cited on page 26.]
- Leonie Weissweiler, Valentin Hofmann, Anjali Kantharuban, Anna Cai, Ritam Dutt, Amey Hengle, Anubha Kabra, Atharva Kulkarni, Abhishek Vijayakumar, Haofei Yu, Hinrich Schuetze, Kemal Oflazer, and David Mortensen. Counting the Bugs in ChatGPT’s Wugs: A Multilingual Investigation into the Morphological Capabilities of a Large Language Model. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 6508–6524, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.401. URL <https://aclanthology.org/2023.emnlp-main.401>. [Cited on page 19.]
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. Transformers: State-of-the-art natural language processing. In Qun Liu and David Schlangen, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online, October 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-demos.6. URL <https://aclanthology.org/2020.emnlp-demos.6>. [Cited on pages 18 and 45.]
- Stephen T Wu, Hongfang Liu, Dingcheng Li, Cui Tao, Mark A Musen, Christopher G Chute, and Nigam H Shah. Unified Medical Language System term occurrences in clinical notes: a large-scale corpus analysis. *Journal of the American Medical Informatics Association*, 19(e1):e149–e156, 04 2012. ISSN 1067-5027. doi: 10.1136/amiajnl-2011-000744. URL <https://doi.org/10.1136/amiajnl-2011-000744>. [Cited on page 45.]
- Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Łukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, and Jeffrey Dean. Google’s Neural Machine Translation System: Bridging the Gap between Human and Machine Translation. *arXiv:1609.08144 [cs]*, October 2016. URL <http://arxiv.org/abs/1609.08144>. arXiv: 1609.08144. [Cited on page 19.]
- Yuemei Xu, Ling Hu, Jiayi Zhao, Zihan Qiu, Yuqi Ye, and Hanwen Gu. A Survey on Multilingual Large Language Models: Corpora, Alignment, and Bias, June 2024. URL <http://arxiv.org/abs/2404.00929>. arXiv:2404.00929. [Cited on page 7.]

Armel Zebaze, Benoît Sagot, and Rachel Bawden. In-Context Example Selection via Similarity Search Improves Low-Resource Machine Translation, August 2024. URL <http://arxiv.org/abs/2408.00397>. arXiv:2408.00397. [Cited on pages 11 and 21.]

Kexun Zhang, Yee Man Choi, Zhenqiao Song, Taiqi He, William Yang Wang, and Lei Li. Hire a Linguist!: Learning Endangered Languages with In-Context Linguistic Descriptions, February 2024. URL <http://arxiv.org/abs/2402.18025>. arXiv:2402.18025 [cs]. [Cited on page 19.]

Mike Zhang and Antonio Toral. The Effect of Translationese in Machine Translation Test Sets. In Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, André Martins, Christof Monz, Matteo Negri, Aurélie Névél, Mariana Neves, Matt Post, Marco Turchi, and Karin Verspoor, editors, *Proceedings of the Fourth Conference on Machine Translation (Volume 1: Research Papers)*, pages 73–81, Florence, Italy, August 2019. Association for Computational Linguistics. doi: 10.18653/v1/W19-5208. URL <https://aclanthology.org/W19-5208>. [Cited on page 10.]

Yanjian Zhang, Qin Chen, Yiteng Zhang, Zhongyu Wei, Yixu Gao, Jiajie Peng, Zengfeng Huang, Weijian Sun, and Xuan-Jing Huang. Automatic term name generation for gene ontology: task and dataset. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4705–4710, 2020. URL <https://aclanthology.org/2020.findings-emnlp.422/>. [Cited on page 8.]

Paula Zorrilla-Agut and Thierry Fontenelle. Iate 2: Modernising the eu’s iate terminological database to respond to the challenges of today’s translation world and beyond. *Terminology*, 25(2):146–174, 2019. [Cited on page 26.]

A. Machine-Translated Definitions

In the main text, experiments of definition-augmented translation (Sections 6.2 and 6.3) were conducted with French definitions, the target language, as it is the only language available in FranceTerme. We provide here additional results for TERMIUM, which includes both French and English definitions. This enables us to study a more general setting, where we do not assume that a French definition exists.

For this, we automatically translate English definitions into French using TowerInstruct-7B-v0.2 [Alves et al., 2024], and reproduce the experiments of Section 5 with these machine-translated definitions.³⁰

We find the results of Sections 6.2 and 6.3 to be consistent with these machine-translated definitions, as reported in Table 8: (i) definition-augmented translation (DEF+TERM) improves term translation (TERM); (ii) the co-hyponym and derivation paradigm strategies improve over random sampling and domain strategies.

³⁰Using the tower_instruct_0_shot configuration as instructed in <https://github.com/deep-spin/tower-eval>.

Setting	ICL	EM	F1
TERM	Random	27.5	47.7
TERM	Domain	29.6	49.7
TERM	Paradigm	<u>36.3</u>	<u>55.4</u>
DEF	Random	6.2	23.4
DEF	Domain	6.4	23.2
DEF	Co-hyponyms	<u>8.2</u>	<u>25.7</u>
DEF+TERM	Random	29.2	50.4
DEF+TERM	Domain	29.8	50.8
DEF+TERM	Fusion	36.6	57.0

Table 8: Results of BLOOM-7.1B on the TERMIUM test set with machine-translated definitions. Results are broken down by ICL selection strategy, like in Table 2: (i) random (baseline); (ii) domain (baseline); (iii) derivation paradigm (not applicable to DEF); (iv) co-hyponyms (not applicable to TERM); (v) fusion of the latter two. Best overall results are bolded while best results in settings TERM and DEF are underlined. Results in setting TERM are copied from Table 2.

B. Frequency and Neology

In addition to the analysis of Section 6.4, Table 9 displays random examples of terms for each decile, which accurately reflects the feeling of neology. After the 7th decile, i.e. 1,000 occurrences, the neological feeling is weaker. Note that *pas*, the most frequent term, is a semantic neologism from the electronics domain and relates to the distance between two adjacent interconnection lines in an integrated circuit. However, *pas* has many different meanings, including as negation adverb “not”, which covers most of its occurrences.

C. Morphosyntactic Classification

We build a multi-label classifier for four of the five classes defined in section 4: prefixation, suffixation, neoclassical or native compounding. For the fifth (syntagmatic compounding), we rely on a simple heuristic: the number of words segmented by spaCy. If there are several words, we consider the term to be a syntagm.

To detect these four morphological processes, we use FastText’s architecture [Joulin et al., 2017], which provides a linear classifier for character sequences, represented by the set of words and character n-grams found in them. This classifier is trained in a one-versus-all fashion, equivalent to a binary classifier for each class.

In this section, we describe in more detail the data used to train and evaluate this classifier.

Decile	Term	Occurrences
min	<i>classification semi-dirigée</i> “semi-supervised classification”	0
0.1	<i>moment d’exécution</i> “timing”	0
0.2	<i>stellarateur</i> “stellarator”	2
0.3	<i>horloge à fontaine atomique</i> “atomic fountain clock”	7
0.4	<i>sondage au limbe</i> “limb sounding”	22
0.5	<i>sauvetage côtier sportif</i> “surf life saving”	74
0.6	<i>planche nautique</i> “aquatic board”	273
0.7	<i>effet de rebond</i> “rebound effect”	1,052
0.8	<i>embarquée</i> “nudging”	4,327
0.9	<i>clonage</i> “cloning”	45,680
max	<i>pas</i> “pitch”	232,506,256

Table 9: Random examples of terms from FranceTerme according to their frequency in a large corpus, one per decile

C.1. Datasets

We build a training and evaluation set from the MorphyNet etymological databases [Batsuren et al., 2021] and the one used for the SIGMORPHON 2022 shared task [Batsuren et al., 2022], both extracted from English Wiktionary.³¹ We combine the two databases because they contain complementary information: SIGMORPHON contains native compoundings but only provides morphological segmentation, while MorphyNet provides the base of all words and differentiates between prefixation and suffixation.

However, these two databases share the same shortcoming: they do not consider neoclassical compounds, which are found mixed in with affixations. To differentiate between them, we use a simple heuristic: if all morphemes in a word are categorized as affixes within MorphyNet, then

³¹<https://en.wiktionary.org/>

none of them are free, so it is a neoclassical compound.

Our algorithm is recursive for decomposing complex words (with more than two morphemes). For example, *bidirectional* will be decomposed into *bi+directional* (prefixation) and *directional* will in turn be decomposed into *direction+al* (suffixation). *Bidirectional* will therefore inherit these two labels.

C.2. Implementation

Statistics from the English and French lexicons are in Table 10, which confirm that native compounds are much rarer in French. We also note that neoclassical compounds are less systematically annotated in French than in English, perhaps because MorphyNet and SIGMORPHON come from English Wiktionary. We also show how the different processes combine in Table 12. Derived terms are often prefixed and suffixed at the same time, which is impossible for neoclassical compounds, by construction.

These lexicons are randomly divided into training (80%), validation (10%), and test (10%) sets. We train one model for each language. Monomorphemes (inflected or not) are kept and serve as negative examples for all classes during training.

FastText hyperparameters are determined automatically on the validation set using the fastText Python library. For both languages, we find it optimal to use character n-grams for $n \in \llbracket 3, 6 \rrbracket$.

C.3. Results

Results on the test set are in Table 11. The classifier is very accurate and has very good recall, with the exception of native compounds in French which are under-represented, due to their rarity, and for which recall is modest. To a lesser extent, recall for neoclassical compounds is lower in French than in English, due to their under-representation in SIGMORPHON, as mentioned above.

Process	# EN	# FR
Native	45,463	2,854
Neoclassical	32,766	7,583
Prefixation	190,305	96,721
Suffixation	217,404	155,169

Table 10: Number of words in our English and French morphological classification corpora for each process independently

	English			French		
	P	R	F1	P	R	F1
Native	95.3	93.0	94.1	89.7	66.3	76.2
Neo.	93.4	91.4	92.4	92.2	87.2	89.6
Pre.	91.5	91.3	91.4	93.8	93.5	93.6
Suff.	93.2	93.3	93.2	97.4	98.0	97.7
Overall	92.7	92.4	92.5	95.9	95.7	95.8

Table 11: Precision (P), Recall (R) and F1 for multi-label morphological classification, in English and French

D. Implementation Details

D.1. LLM Implementation

LLMs are implemented in the transformers library [Wolf et al., 2020] itself based on pytorch [Paszke et al., 2019]. LLMs are quantized in 8 bits for effective inference on a single V100 GPU with 32GB of RAM. We use greedy decoding.

D.2. mBART Fine-tuning on SciPar

mBART is implemented with fairseq [Ott et al., 2019]. It is fine-tuned with a single NVIDIA RTX A6000 GPU with 48GB of RAM. It uses a batch size of 4,096 samples and accumulates gradients for 4 steps. Early stopping is done according to the validation BLEU score [Peng et al., 2024].³²

D.3. Corpus frequency

For the analysis of Section 6.4, we compute corpus frequency (case insensitive) using Aho-Corasick’s algorithm [Aho and Corasick, 1975, Wu et al., 2012], implemented in the pyahocorasick Python library.³³

E. Rules of Morphotactics

The following rules were used to create controlled datasets pairing a base (e.g., “*lightfast*”) with a prefixed derivative (e.g., “*unlightfast*”) and a suffixed derivative (e.g., “*lightfastly*”; see Section 8.5).

³²SacreBLEU signature [Post, 2018]:

nrefs:1lcase:mixedleff:noltok:13alsmooth:explversion:2.3.1

³³<https://pyahocorasick.readthedocs.io>

Native	Neo.	Pre.	Suff.	# EN	# FR
				207,074	118,811
			✓	109,353	90,646
		✓		91,115	35,646
		✓	✓	88,349	60,307
	✓			17,191	3,508
	✓		✓	9,677	3,640
	✓	✓		5,593	432
	✓	✓	✓	0	0
✓				34,425	2,162
✓			✓	5,552	353
✓		✓		808	115
✓		✓	✓	4,373	221
✓	✓			138	1
✓	✓		✓	100	2
✓	✓	✓		67	0
✓	✓	✓	✓	0	0

Table 12: Number of words in our English and French morphological classification corpora for each process combination

For English (i) The base should not start with “*un*” to avoid a double negation. (ii) The base should not end with:

- “*y*” because it would then have to be substituted by “*i*” (as in “*easy*” → “*easily*”);
- “*le*” because it would be deleted (as in “*noble*” → “*nobly*”);
- “*ll*” because the suffix would then be “*-y*” instead of “*-ly*” (as in “*full*” → “*fully*”);
- “*ic*” to avoid allomorphy with the suffix “*-ally*” (as in “*allergic*” → “*allergically*”).

For French (i) The base should not start with:

- “*i*” to avoid a double negation;
- “*b*”, “*l*”, “*m*”, “*n*”, “*p*”, or “*r*” to avoid allomorphy with the “*i-*” prefix (also respectively written “*il-*”, “*im-*”, or “*ir-*”), as in “*irréaliste*”.

(ii) The base *should* end with an “*e*” so that the “*-ment*” suffixation is morphotactic (e.g., avoid impossible words like “**absentment*”) and orthographic (e.g., adverbs are often formed on the feminine adjectival form that ends with an “*e*”: “*amicalement*” and not “**amicalment*”).

Model	Dataset	Random ICL + BPE Seg.		Morph. ICL + BPE Seg.		Morph. (ICL + Seg.)	
		Prefix	Suffix	Prefix	Suffix	Prefix	Suffix
CroissantLLM-1.3B	ADJ-EN	0.196	0.581	0.836	0.842	0.988	0.988
	ADJ-FR	0.047	0.356	0.371	0.754	0.991	0.947
	PSEUDO-EN	0.265	0.651	0.811	0.866	0.987	0.980
	PSEUDO-FR	0.045	0.285	0.367	0.586	0.978	0.961
Llama-2-7B	ADJ-EN	0.715	0.812	0.949	0.879	0.990	0.999
	ADJ-FR	0.549	0.565	0.889	0.717	0.994	0.990
	PSEUDO-EN	0.792	0.802	0.943	0.900	0.997	0.997
	PSEUDO-FR	0.597	0.390	0.868	0.523	0.988	0.988
BLOOM-560M	ADJ-EN	0.105	0.477	0.561	0.771	0.998	0.996
	ADJ-FR	0.005	0.230	0.055	0.457	0.970	0.993
	PSEUDO-EN	0.077	0.556	0.535	0.772	0.996	0.992
	PSEUDO-FR	0.002	0.141	0.028	0.275	0.969	0.981
BLOOM-1.1B	ADJ-EN	0.154	0.412	0.595	0.696	0.995	0.996
	ADJ-FR	0.063	0.166	0.280	0.285	0.978	0.985
	PSEUDO-EN	0.175	0.632	0.578	0.783	0.962	0.951
	PSEUDO-FR	0.037	0.087	0.254	0.217	0.981	0.981
BLOOM-3B	ADJ-EN	0.251	0.663	0.796	0.860	0.998	0.995
	ADJ-FR	0.132	0.448	0.556	0.635	0.994	0.987
	PSEUDO-EN	0.298	0.663	0.759	0.785	0.997	0.995
	PSEUDO-FR	0.106	0.292	0.516	0.470	0.988	0.981
BLOOM-7.1B	ADJ-EN	0.488	0.724	0.893	0.879	0.999	0.998
	ADJ-FR	0.262	0.560	0.655	0.746	0.995	0.996
	PSEUDO-EN	0.625	0.799	0.873	0.868	0.999	0.998
	PSEUDO-FR	0.278	0.435	0.670	0.585	0.998	0.994

Table 13: Numbers in Figure 8

F. Complete Results

Table 13 reports the scores that are plotted in Figure 8.