

Position Paper

Losing our Tail, Again: (Un)Natural Selection & Multilingual LLMs

Anonymous ACL submission

Abstract

Multilingual Large Language Models (LLMs) considerably changed how technologies can influence language. While previous technologies could mediate or assist humans, there is now a tendency to *offload* the task of writing itself to these technologies, enabling these models to change our languages more directly. While they provide us quick access to information and impressively fluent output, beneath their (apparent) sophistication lies a subtle, more insidious threat: the gradual decline and loss of linguistic diversity. In this position paper, I explore how model collapse, with a particular focus on translation technology, can lead to the loss of linguistic forms, grammatical features, and cultural nuance. Model collapse refers to the eventual consequence of self-consuming training loops, where automatically generated data enters the training data thereby reinforcing biases and losing linguistic diversity. Drawing on recent work in Computer Vision, Natural Language Processing and Machine Translation, I argue that the *tails* of our linguistic distributions are vanishing, and with them, the narratives and identities they carry. This paper is a call to resist linguistic flattening and to reimagine Natural Language Processing as a field that encourages, values and protects expressive multilingual diversity and creativity.

1 A Simple Evolutionary Trade-Off

A few million years ago, we lost our tails. Once useful (and presumably quite fun), they became obsolete, absorbed by evolutionary trade-offs. They disappeared through natural selection. Darwin already drew parallels between the evolution of languages and those of species stating that they follow ‘curiously parallel’ paths in *The Descent of Man* (Darwin, 1871). He went so far as to suggest that “*the survival of certain favoured words in the struggle for existence is natural selection*” (Darwin, 1871, p. 61), highlighting how language itself is

subject to various evolutionary forces. However, it is unlikely that even Darwin could have anticipated the current path of language evolution and the latest forces that came into play.

Today, it is not just ‘natural’¹ forces but technological ones that exert strong selective pressures on languages and influence the struggle for linguistic survival on various levels. Given that Multilingual LLMs are essentially next-word predictors over large datasets, they amplify statistically likely languages and linguistic forms, whether these are words, subwords, syntactic constructions – pruning away our rich, long statistical tails. Will this technology-driven artificial selection merely accelerate an already inevitable trajectory or does it mark a more disruptive turn? In other words, does losing our ‘language’ tails imply the disappearance of obsolete, fun and decorative elements or does it signal a more fundamental erosion of diversity and richness? In this paper I take a position on these questions arguing that the LLM-era marks a disruptive turn: unlike ‘natural’ evolution, technological selection in current mainstream LLMs (under prevailing data and training regimes) reduces rather than reshapes linguistic diversity.

A quick note on language tails Like many complex systems, language follows long-tailed distributions: a few elements occur very frequently, while most are rare.

2 Natural Language Selection

Before turning to how technology reshapes language, it is helpful to briefly consider how languages evolve ‘naturally’. Evidently, this is not

¹While the technological forces that shape language can certainly be seen as part of the ‘natural’ forces influencing our linguistic ecosystem, there is a crucial shift when we move from tools that help us record, polish, or disseminate language to models that actively *generate* it. In such systems, the technology becomes not just a medium but a core driver of language change. I will return to this later in the paper.

076 meant to be a comprehensive account of language
077 evolution since doing justice to that topic requires
078 several volumes.² The goal here is to highlight a
079 few key aspects of how humans learn and transmit
080 language, to lay the groundwork for drawing par-
081 allels with how Multilingual LLMs learn, transmit
082 and affect language. Broadly speaking, natural se-
083 lection in languages operates at two main, related
084 levels: *across languages* (S. 2.1), where some lan-
085 guages survive and thrive while others gradually
086 disappear, and *within languages* (S. 2.2), where cer-
087 tain words, expressions, or grammatical structures
088 persist over time while others fall out of use.

089 **2.1 Across Languages**

090 Historically, population shifts, cultural dominance
091 and environmental factors have repeatedly driven
092 language extinction (Bromham et al., 2022). Nearly
093 half of the world’s roughly 6000 languages³ are en-
094 dangered, many disappearing at an alarming rate.⁴

095 While language loss is often viewed as some-
096 thing tragic, it can stem from a shared communica-
097 tive goal. For instance, speakers might adopt a
098 dominant language to improve mutual understand-
099 ing. Evidently, even in these seemingly positive
100 scenarios, something linguistically and culturally
101 valuable is lost. Yet, more subtly and at the same
102 time more *universally*, we might also be witness-
103 ing a different type of loss: the loss of richness
104 and diversity **within** languages themselves. Such a
105 loss is less likely to be driven by collective needs
106 like mutual intelligibility⁵, but can most likely be
107 attributed to converging pressures and constraints
108 from language technologies.

109 **2.2 Within Languages**

110 Whereas the extinction of languages is often driven
111 by external forces, changes within a language are
112 also shaped by internal cognitive bottlenecks. The
113 acquisition of a language, for instance, happens un-
114 der rather constrained conditions: humans are ex-
115 posed to sparse, ambiguous, and often noisy and im-
116 perfect input. In addition to this so-called *poverty*

²For readers interested in a more broad, comprehensive,
nuanced history on, among others, the evolution of language, I
recommend “The Language Puzzle” (Mithen, 2024) and “How
Language Began” (Everett, 2017).

³Estimates vary between 3000 and 10000 (Crystal, 2002).

⁴<https://www.unesco.org/en/articles/multilingual-education-bet-preserve-indigenous-languages-and-justice>.

⁵Of course, in some specific cases or contexts, simplifi-
cation can help communication across a more diverse set of
speakers.

117 *of stimulus* (Chomsky, 1980), we face cognitive
118 limitations including a finite memory, a limited
119 attention span and bounded processing capacities.
120 Yet, despite these many hurdles (or perhaps because
121 of them, as suggested by DeCaro et al. (2008))
122 (dominant) languages continue to thrive.

123 In the following sections, I will distinguish be-
124 tween two types of language-internal changes: (i)
125 the emergence of structure, briefly revisiting some
126 foundational research on grammar and composi-
127 tionality (S. 2.2.1); and (ii) the evolution of our
128 the lexicon, that is, how words are lost, gained or
129 transformed over time (S. 2.2.2). While in reality,
130 these processes are to some extent intertwined, they
131 are shaped by different pressures.

132 **2.2.1 Towards Structure**

133 Iterated learning (Brighton et al., 2005; Smith et al.,
134 2003b; Ren et al., 2024a) models how linguistic
135 structure emerges through repeated cultural trans-
136 mission, where each generation of learners infers
137 patterns from the behaviour of the previous one.

138 To cope with limited and noisy input, humans
139 rely on (implicit) inductive biases such as simplic-
140 ity and compressibility to learn language (Kirby
141 et al., 2015). These structured preferences guide
142 learning by enabling pattern extraction and struc-
143 tural generalization from sparse data. Research
144 by Smith et al. (2003a) showed that when lan-
145 guage is transmitted under cognitive and commu-
146 nicative pressure, structure and compositionality
147 emerge through repeated cultural transmission cy-
148 cles—processes shaped by our cognitive biases that
149 favor more learnable and generalizable patterns.
150 Far from being detrimental, such biases are often
151 considered essential (DeCaro et al., 2008), allowing
152 generalization without linguistic impoverishment.

153 This also implies that our learning process in-
154 volves *convergence*. Convergence towards system-
155 atic forms, towards structure and patterns. This
156 convergence is not of the kind that leads to a loss of
157 expressivity. It is a structural, higher-level conver-
158 gence that supports more open-ended productivity,
159 one where nonsense syllables can become, for in-
160 stance, compositional morphemes.

161 While iterated learning and cultural transmis-
162 sion help explain how structure in language be-
163 comes more regular and learnable over time, not
164 all aspects of language evolve in such a way. Our
165 lexicon, the system of words and meanings, often
166 changes in more irregular and dynamic ways. Vo-
167 cabulary is more susceptible to social influence,

168 borrowing, innovation, and drift.

169 2.2.2 Lexical Changes

170 Despite our limited memory, languages are and re-
171 main remarkably rich. Words may disappear, fall
172 out of fashion or decrease in usage, their overall
173 vocabularies remain relatively stable, often even ex-
174 panding when languages flourish driven by social
175 needs, interactions, scientific progress or creativ-
176 ity. For instance, based on over five million digi-
177 tized books,⁶ Michel et al. (2011) estimated that
178 English grew by more than 70% over the past fifty
179 years, adding on average about 8500 new words per
180 year. A later study by Gerlach and Altmann (2013)
181 showed that vocabulary growth largely scales with
182 corpus size, with the rate of new-word introduc-
183 tion slowing but never vanishing as datasets ex-
184 pand. These findings can be contrasted however
185 with claims of convergence towards a maximum
186 vocabulary size (Bernhardsson et al., 2009).

187 The vocabulary growth reported for English by
188 Michel et al. (2011) and Gerlach and Altmann
189 (2013) aligns with findings that languages with
190 larger speaker populations tend to have larger vo-
191 cabularies (Reali et al., 2018). At the same time,
192 widely spoken languages often exhibit simpler mor-
193 phology (fewer cases and inflections) while smaller
194 communities may develop structurally more com-
195 plex systems (Lupyan and Dale, 2010).

196 Word survival is frequently influenced by posi-
197 tive frequency-dependent selection (Pagel et al.,
198 2019), yet frequency alone does not determine
199 a word’s fate (Altmann et al., 2011). Competi-
200 tion among words is also shaped by cognitive and
201 communicative pressures that favour efficiency and
202 clear categorization. Naming systems for colours
203 or emotions, for example, tend to converge across
204 cultures, reflecting shared communicative needs
205 rather than pure frequency effects (Petersen et al.,
206 2012). Even common words can disappear if they
207 no longer align with prevailing cognitive, social,
208 or technological forces. Consider the case of ‘ra-
209 diogram’, ‘Roentgenogram’, and ‘X-ray’, once
210 competing terms for the same concept (Petersen
211 et al., 2012). Although ‘Roentgenogram’ domi-
212 nated early scientific discourse, it was gradually re-
213 placed by the shorter and more efficient ‘X-ray’ (Pi-
214 antadosi et al., 2011), a shift reinforced by the
215 global rise of English as the scientific lingua franca.

⁶The 2010 analysis covered roughly 4% of all books ever published.

Vocabulary, in short, is shaped not just by usage
frequency but by communicative efficiency, cogni-
tive constraints, technological mediation, and so-
ciocultural forces, showing a dynamic and context-
sensitive interplay between convergence and diver-
gence. Yet as technological mediation increasingly
shapes communication, these dynamics may now
be governed by a different kind of selection.

3 A Force of *Technology*

Darwin argued that languages evolve through a pro-
cess of variation and selection, convergence and
divergence, similar to much like species. We have
seen how such dynamics give rise to compo-
sitionality and structure, as well as changes in our
lexicon. Broadly speaking, we indeed established
that language reflects a tension between conver-
gence and divergence. However, we set out to
explore what happens when the strongest selec-
tive forces are no longer human cognition or social
need, but instead arise from (large-scale) techno-
logical forces.

First, we briefly examine the relationship be-
tween language and cultural technologies that pre-
date the rise of (L)LMs (S. 3.1). This provides con-
text for the *Descent of Language* section (S. 3.3),
where I examine early sign of linguistic decline
in computational models and discuss more recent
work on LLMs.

3.1 Language & Cultural Technologies

It is important to note that it is not the first time
that technology plays a disruptive role in shaping
language. Language has never existed in isolation
and (cultural) technologies have continuously influ-
enced how language is transmitted. While it is hard
to imagine now, even *writing* was once considered
a disruptive technological innovation:

*They seem to talk to you as though they
were intelligent, but if you ask them any-
thing about what they say from a desire
to be instructed they go on telling just the
same thing forever.*

— Plato, *Phaedrus* 275d

As cognitive scientist Alison Gopnik pointed out
in her 2023 ACL keynote⁷, without any context,
this critique on *writing* could just as well apply to
LLMs today. There is a surface-level fluency that

⁷Gopnik, Alison (2023). Keynote at ACL 2023: *Large Language Models as Cultural Technologies: Imitation and Innovation in Children and Models*.

(over)confidently mimics understanding, yet lacks genuine responsiveness, real-world grounding or explanatory depth. While we no longer worry about the “dangers” of writing, history reminds us that each new “tool” reshapes how language is used, transmitted, and transformed. Some languages, the so-called “long-tailed” ones⁸, have been particularly impacted. Their marginalization has been reinforced by successive cultural technologies (from writing to the internet), and by their exclusion from the digital sphere. This phenomenon, described as *digital language death*, is said to affect roughly 95% of the world’s languages (Kornai, 2013).

Earlier technologies also marked turning points in the evolution of language. The printing press, for instance, helped standardize spelling and grammar, elevating dominant dialects while marginalizing others (Sasaki, 2017). More recently, automatic spell-checkers boosted the ‘reproductive fitness’ of recognized forms at the expense of other alternatives (Petersen et al., 2012), contributing, for example, to the rising dominance of ‘colour’ over ‘color’ (Petersen et al., 2012)⁹.

While the impact of technology on language and concerns regarding their disruption are not new, the current wave of generative models operate at an different level, scale, and speed. They represent a new kind of intervention, one where technology is not simply storing, mediating or ‘correcting’ language, but actively **generating** it across different domains, tasks and platforms. In a sense, the technology quite suddenly moved from the passenger’s to the main driver’s seat.

Concerns regarding the loss of diversity *across* languages due to recent technologies, such as LLMs, have been explicitly raised in recent, seminal, work within the field of Natural Language Processing (NLP) by, among others, Joshi et al. (2020) and Bender et al. (2021). Here, I shift the focus inwards, to the internal diversity *within languages*. This dimension received comparatively little attention, even though current multilingual LLMs exert significant influence on what is preserved, erased, or amplified within a language itself.

⁸Languages that receive limited localization attention or commercial investment; this does not always correspond to speaker numbers. For instance, Bengali is the 7th most spoken language, but falls outside the top 50 most localized languages (<https://www.lionbridge.com/blog/translation-localization/localizing-long-tail-languages/>)

⁹Amusingly, the tool I am writing in still nudges me away from ‘colour’.

3.2 Large Language Models

LLMs mark a significant shift in how language is produced, accessed, and reused; and perhaps more interestingly, how the cognitive labor of writing and idea formulation is now often *offloaded* to technology. Earlier Language technologies were built for (domain-) specific tasks; today’s foundation models, however, operate across domains and applications. Large-scale models are rapidly becoming an integral part of a broad range of our everyday activities, allowing them to shape and influence language more directly than before. Unlike earlier tools that assisted human writing, these models now drive change by producing language themselves. Emerging research, suggests that LLMs introduce subtle, yet cumulative, distortions (Shumailov et al., 2023). Albeit initially *imperceptible*, these small changes can accumulate across multiple training cycles leading to a phenomenon coined *model collapse*, where models trained on their own outputs progressively lose quality and diversity. In computer vision, such a collapse leads to *visible* artefacts in AI-generated images (Alemohammad et al., 2023), yet for language, the consequences remain underexplored. This despite the fact that language shapes and possibly constrains human thought, meaning that this gradual impoverishment or distortion of linguistic output could have profound and far-reaching implications.

LLMs will likely continue to substantially influence both the content we are exposed to (images, texts, audio...) and the systems that generate this content given that its output will increasingly re-enter the training cycle. Therefore, at this point we can assume that interactions between models are not hypothetical but inevitable (Martínez et al., 2023). These interactions can occur through (partial) training on output from another LLM or a model’s own output (Ren et al., 2024b). The subsequent feedback loops this creates, where models learn from their own output or that of others, accelerate concerns regarding the distortion of language.

So far, emphasis and research efforts have largely focused on sustaining the benefits of training from large-scale, human-generated data scraped from the Web, summarized in this recent article as: “*LLMs world is our word*”¹⁰. But perhaps more concern-

¹⁰https://www.theguardian.com/technology/article/2024/sep/07/if-journalism-is-going-up-in-smoke-i-might-as-well-get-high-off-the-fumes-confessions-of-a-chatbot-helper?utm_source=chatgpt.com

ingly, I fear that: “*Our world* might be turning into *their word*”.

3.3 The Descent of Language

Already before terms such as ‘model collapse’ (Shumailov et al., 2023), ‘Model Autophagy Disorder’ (MAD) (Alemohammad et al., 2023) or ‘Habsburg AI’¹¹ gained traction, earlier research empirically showed how statistical language models¹² indeed amplify dominant linguistic forms while forgetting or flattening rarer, low-probability ones.

3.3.1 Precursor: Statistical and Neural Translation Models

As a precursor to this line of work, studies in Machine Translation (MT) (Vanmassenhove et al., 2019, 2021) provided empirical evidence that both statistical and neural models systematically favor frequent lexical and morphological patterns, reducing linguistic diversity. This raised two concerns: (i) technically, frequency bias diminishes lexical richness and can eliminate infrequent but grammatically necessary forms; and (ii) sociolinguistically, machine-generated translationese may, over time, influence language itself (Kranich, 2014). For instance, Vanmassenhove et al. (2021) showed that MT systems disproportionately produce masculine *pr’esident* over feminine *pr’esidente* when translating English *president* into French, reflecting data imbalances. Under iterative training, rare forms may disappear entirely: in their experiments, the plural *pr’esidentes* vanishes. Such low-probability forms are not noise but carriers of grammatical precision, expressiveness, and social meaning—yet current models fail to distinguish these cases, discarding rare forms indiscriminately. Similarly, Luo et al. (2024) found MT outputs to be structurally closer to the source text than human translations, with fewer morphosyntactic divergences. Beam search biases toward “safe,” high-probability constructions, reinforcing convergence and reducing syntactic diversity.

3.3.2 Generative Models

In recent work on LLMs, concerns around **model collapse** have gained traction: when models are trained on data increasingly composed of their own

¹¹A term coined by Jathan Sadowski (<https://x.com/jathansadowski/status/1625245803211272194?lang=en>).

¹²Whether they are called statistical, neural or large language models, in the end, they are still all statistical models.

(or other models’) output, their behavior begins to *converge*, potentially degrading over time (Shumailov et al., 2023). In computer vision this de-generation has been made *visible*, with iterative training cycles producing recognizably distorted ‘Habsburg Jaw’ artefacts (Alemohammad et al., 2023). In language, such collapse is likely subtler and harder to detect, yet potentially more consequential, given the foundational role of language in human cognition and cultural evolution.

Empirical studies suggest early signs of such dynamics. McCoy et al. (2023) show that even on simple deterministic tasks, LLM behavior is sensitive to probability: GPT-4’s performance drops dramatically when the correct output sequence is low-probability, revealing what they call *em-bers of autoregression*. They argue for evaluations grounded in model training objectives and constraints — a more explicitly *teleological* perspective. Evidence from real-world text production points in a similar direction. Kobak et al. (2024) report sharp, sudden spikes in terms such as *delve*, *crucial*, and *significant* in PubMed abstracts following the release of public LLM tools, exceeding even pandemic-related shifts. This suggests subtle, distributional nudging of academic discourse toward statistically likely phrasing. Meanwhile, Shumailov et al. (2023) demonstrate how low-probability events disappear first in iterative training, with models gradually converging toward high-probability sequences — an effect reminiscent of earlier “poisoning” cycles in search engines, though with no obvious filtration mechanism available for LLMs.

Focusing specifically on linguistic structure, Guo et al. (2024) find that LLM-generated text exhibits reduced diversity compared to human language, particularly in creative tasks. Instruction tuning improves lexical variety but narrows syntactic and semantic flexibility overall, suggesting a redistribution rather than resolution of expressive constraints.

4 Thoughts & Discussions

Darwin noted the similarities between the evolution of species and that of languages. Of course, he was referring to mechanisms of gradual change through natural evolution. I could, however, not help to be amused by the rather unexpected parallel that can be drawn between the evolution of our species and the effect of recent technologies on languages. It seems that, once again, our species

448 and language are following a similar path: *We are*
449 *losing our tails*. While our physical appendages be-
450 came obsolete and even cumbersome, I argue that
451 its language equivalent is anything but that, cap-
452 turing the rich diversity and continuous evolution
453 of language, shaped both by the need to articulate
454 novel concepts and by our ongoing desire to signal
455 belonging and distinction within social groups.

456 While one could argue that across languages,
457 various natural forces, sometimes in combination
458 with technological ones, have contributed to a sig-
459 nificant loss of diversity across languages leading
460 even to *language death*, when we shift the focus to
461 the internal diversity within languages, natural evo-
462 lution has led to structure, compositionality, and a
463 growth in terms of vocabulary, at least for thriving
464 languages. This stands in contrast with what we
465 observe when language is *driven* (largely or solely)
466 by technological forces, whose internal pressures
467 seem to lead to a reduction of expressivity, creativ-
468 ity and a loss of overall lexical diversity regardless
469 of whether the language is flourishing or not.

470 In the paragraphs below, I set out my reflections
471 and formulate a position on what it means when
472 current technological forces become the main en-
473 gines of language change.

474 (L)LMs reduce linguistic richness and amplify

475 **biases** While current models' ability to gener-
476 alize over large amounts of data is one of their
477 biggest assets, their statistical nature and the pres-
478 sures that shape these models have drawbacks re-
479 garding diversity and creativity. Up until recently,
480 this might not have been a priority for our field,
481 given that these technologies were largely regarded
482 as domain-specific tools that were often still su-
483 pervised or post-edited by humans (e.g. chatbots,
484 MT...). However, the fairly recent public release
485 of large, general-purpose, language models raises
486 concerns about the potential implications for lan-
487 guage (and language technologies) in the longer
488 term.

489 From the literature, it seems that indeed, (L)LMs
490 exacerbate imbalances in the data by (i) forgetting
491 low-probability events/words, and (ii) overgener-
492 ating high(er) probability ones (Vanmassenhove
493 et al., 2021; Shumailov et al., 2023). This could
494 lead to self-reinforcing loop where less frequent
495 linguistic forms and expressions are underrepre-
496 sented and risk being entirely lost in translation.
497 Low(er)-probability events (words, subwords, etc.)
498 contribute to the complexity and richness across

499 and within language(s), and are important to ensure
500 that models do not converge toward oversimplified,
501 biased language. After all, languages are full of
502 improbable events.

503 Furthermore, the effects of model collapse are
504 likely not confined to obvious levels such as vocabu-
505 lary. They may also appear at the morphosyntactic
506 level (e.g., Guo et al. (2024)), at non-linguistically
507 motivated subword or character levels, in the fact
508 that (the) dominant or target language(s) struc-
509 ture(s) may "bleed through" (e.g., preferences for
510 Subject-Verb-Object constructions), or in the prop-
511 agation of specific ideologies (e.g., representations
512 of gender). Prior work (e.g. Cao et al. (2023) or
513 Dokic et al. (2025)) has demonstrated how stereo-
514 types encoded in English can "leak" into other lan-
515 guages in multilingual LLMs, with typologically
516 distant languages being particularly vulnerable.
517 This asymmetry is potentially overlooked given the
518 dominance of English in evaluation benchmarks,
519 which could lead to an underestimation of col-
520 lapse effects in other languages. Similarly, one
521 could hypothesize that multilingual models might
522 start confabulating words similar to English¹³ or
523 translate English idiomatic expressions literally,
524 even when translating into typologically distant lan-
525 guages, a well-known cause for errors (Karakanta
526 et al., 2025).

527 And last but not least, given the unbalanced na-
528 ture of language representation on the web, we
529 can furthermore assume that minority, long-tailed
530 and morphologically-richer languages are likely af-
531 fected disproportionately. As model outputs feed
532 back into future models, the result is a **compound-**
533 **ing distortion of language use, progressively**
534 **shifting further away from diversity observed**
535 **in the real world on the language-internal and**
536 **cross-linguistic level.**

537 Methodological Blind Spots in Measuring Lin-

538 **guistic Diversity** Current evaluation practices in
539 NLP often involve pairwise comparisons between
540 a single AI-generated text and a single human-
541 written text. Because the model has been trained
542 on massive datasets containing the writing styles,
543 vocabularies, and linguistic patterns of millions
544 of humans, its outputs often exhibit surface-level
545 lexical or syntactic variety that surpasses that of

¹³Castilho et al. (2025) showed how, just like humans, LLMs tend to resort to what they call "Lazy Gaelicisations" which involve the adaptation of English words towards the Irish orthography. This is, however, also a common strategy among Irish speakers.

546 an individual human writer/text (Reviriego et al.,
547 2024). This could lead to potentially misleading
548 conclusions, where one might start claiming that AI
549 outputs are *overall* more diverse or more lexically
550 rich than human texts. This way of comparing AI-
551 written vs human-written text, *overlooks people’s*
552 *individual variation*.

553 LLMs can adopt and mimic many different
554 styles, but at inference time, they tend to converge
555 on high-probability patterns. When many genera-
556 tions of outputs are compared over time, these out-
557 puts are likely increasingly uniform. This can be
558 contrasted against what we observe in human lan-
559 guage, which is inherently diverse *across* speakers,
560 contexts, and time. If we instead evaluate diversity
561 at scale (i.e., across many texts or over generations),
562 human-generated language maintains a rich vari-
563 ation through individual idiolects, sociolects, and
564 cultural registers shaped by our individual biases
565 rather than a common one. **Short-term superficial**
566 **lexical diversity does not guarantee long-**
567 **term linguistic sustainability.** Methodological
568 approaches that treat diversity as an isolated textual
569 property risk drawing premature conclusions with
570 respect to diversity on the longer term.

571 Aside from this potential methodological
572 blindspot regarding linguistic diversity, it is worth
573 highlighting once more that current evaluation
574 metrics for translation (BLEU (Papineni et al.,
575 2002), METEOR (Banerjee and Lavie, 2005),
576 TER (Snover et al., 2006), COMET (Rei et al.,
577 2020)) are not designed to capture loss of diver-
578 sity. Nor should they: diversity is not a property
579 of a single translation, nor even necessarily of a
580 single text. This does not mean, however, that we
581 should ignore it: diversity emerges at the level of
582 systems and over time, and understanding how it
583 evolves across models and generations is important
584 for assessing the broader impact of language tech-
585 nologies. More broadly, in line with McCoy et al.
586 (2023), we conclude that we should not evaluate
587 LLMs as if they are humans but should instead treat
588 them as a distinct type of system, one that has been
589 shaped by its own particular set of pressures. It is
590 thus important for metrics to be designed in order
591 to reveal *their* idiosyncratic weaknesses.

592 **Compositionality and Systemacity still largely**
593 **elude LLM capabilities.** Humans learn language
594 under limited memory capacities, there is a critical
595 period where we can easily learn languages, and
596 our exposure to language is (in some ways) much

597 more restricted than that of LLMs. Our bottlenecks,
598 however, seems to serve as pressure mechanisms
599 for the emergence of structure and compositionality,
600 which allows us, among others, to create new
601 words that can often immediately be understood
602 by others speaking the same language. We gen-
603 eralize and converge, but we do so towards a pro-
604 ductive system. Even though neural networks can
605 behave compositionally and systematically, it is
606 not straightforward for them. With (virtually) un-
607 constrained memory and faster convergence, these
608 models, in particular LLMs, can learn a language
609 in no time albeit with (overall) smaller and less
610 diverse vocabularies.

611 Regarding compositionality in NMT using
612 Transformers (Vaswani et al., 2017), recent work
613 by Yin et al. (2024) compares the performances of
614 different Transformer-based models (Transformer
615 trained from scratch, pre-trained decoder-only mod-
616 els (BLOOMZ-7b (Muennighoff et al., 2022) and
617 LLaMA2-13b (Touvron et al., 2023)) and a pre-
618 trained encoder-decoder model (mT5-large (Xue
619 et al., 2021))). They illustrated that all these models
620 still struggle when translating new or long com-
621 pounds. Additionally, lower perplexity source sen-
622 tences are more likely correctly translated into the
623 target and error rates go up, when the length of
624 the compounds increase. These findings are in line
625 with some of the phenomena discussed in the pre-
626 vious section. While they do find that fine-tuned
627 Pretrained LLMs outperform Transformer models
628 trained from scratch, they point out that this advan-
629 tage could be due to pretraining exposure rather
630 than true compositional generalizations.

631 In experiments similar to those conducted by
632 Kirby et al. (2014) illustrating the effect of cultural
633 transmission through an iterated learning frame-
634 work, Kouwenhoven et al. (2025b) and Kouwen-
635 hoven et al. (2025a) empirically evaluated and com-
636 pared human-human, LLM-LLM and human-LLM
637 (artificial) language learning to compare how arti-
638 ficial languages differ when optimized by LLMs
639 or humans’ inductive biases.¹⁴ Their comparisons
640 of language learning across the three different con-
641 ditions revealed that, while similar to human vo-
642 cabularies, LLM languages are subtly different.

¹⁴They do not focus on behavioural biases but on the im-
plicit inductive ones. For humans, these are biases such as pref-
erences for compressibility, simplicity or efficiency), while
for LLMs they focused on *increasingly apparent* biases of
the Transformer architecture (e.g. simplicity, structure, re-
cency)(Kouwenhoven et al., 2025a).

The LLM optimized languages showed less diversity and variation, making them more *degenerate* in comparison to those optimized for humans. These differences were alleviated when humans and LLMs collaborated, which underscores that to achieve successful interactions between humans and machines, it is essential to optimise for communicative success since the need to be expressive in human language can prevent convergence.

More generally, Zhou et al. (2023) looked at the link between the complexity of the dataset and the ability of models to generalize. More complex datasets provide: (i) more diversity in terms of the examples the model is exposed to but also, (ii) a reduced repetition preventing the model from *ungeneralizable* surface memorization. Yet, as pointed out by Dziri et al. (2023) but also by McCoy et al. (2023), these models still fail on sometimes surprisingly trivial problems and quickly decay once task complexity increases - indicating once more that these are symptoms of a more fundamental limitation.

Averaging Biases Humans are fundamentally biased. For decision-making we rely on frugal heuristics (Gigerenzer and Goldstein, 1996) which is efficient but obviously imperfect. Again, this relates to our cognitive constraints (limited attention, processing capacity, memory). In this context, it is important to highlight that our biases are **not monolithic**. While it is true that they are partially shaped by our society, environment and direct surroundings, we each develop a slightly different set of heuristics and biases over time. These are based on our unique experiences and contexts. Regardless of whether they are good or bad: *they are many*. Besides, some of us actively fight or question our own biases when we recognize that they could be harmful, unfair or simply undesirable.

The biases multilingual LLMs propagate, in contrast, are averaged, dominant ones that are present in an already biased sample of training data. Rather than capturing the diversity and heterogeneity of biases in human reasoning, by letting LLMs drive language change we risk exacerbating and normalizing dominant biases without having a critical self-reflection or self-correction component.

Invisible Gaps: The Missing Data

That which we ignore reveals more than what we give our attention to.

Finally, model collapse is as much about what is not generated as what is. Over-reliance on high-frequency or majority-culture content leads to blind spots in representation. The absence of specific linguistic structures, sociolinguistic variants, or cultural references in training data can render certain forms of expression invisible in model outputs. As models are increasingly retrained on AI-generated content, these omissions risk becoming permanent.

5 Conclusions

A few million years ago, we lost our tails. Once functional, they became obsolete and cumbersome. I set out in this paper with the question: Do the many statistical language tails face the same fate, and if so, would we merely be losing the obsolete, fun and decorative elements?

Based on recent LLM-related research, I argue that the long statistical tails of language may indeed face a similar fate. The artificial selection driven by LLMs marks a rather disruptive shift from language being shaped by generational, cultural transmission. Unlike the evolution of language driven by humans, which despite (or because of) our cognitive constraints shaped language into a structured, productive and compositional tool with a rich vocabulary; language shaped by models tends to collapse towards what is likely, driven by statistical biases. The interplay and balance between convergence and divergence, that characterizes human behaviour and communication on multiple levels risks being lost. Unlike humans, models are not intrinsically motivated to be creative, to express belonging or differences, or to innovate. There is no capacity for critical self-reflection, self-correction and a limited plurality of voices.

As these systems increasingly *ingest* their own or other models' outputs, the risk of flattening linguistic diversity grows, with rare words, less-resourced languages, and culturally significant variation most at risk. Preserving the long tails of language means rethinking how we evaluate and train these systems, not just for accuracy or fluency, but for the communicative richness that makes language human. **Without our tails, we risk losing the balancing act** by converging, collapsing and flattening the expressive multilingual linguistic, social and cultural diversity.

¹⁵<https://github.com/MimiOnuoha/missing-datasets>

740 **Limitations**

741 The arguments presented in this paper are intended
742 to provoke critical reflection on the trajectory of
743 language in the era of LLMs; however, they are
744 subject to several limitations regarding scope, em-
745 pirical generalization, and my own perspective.
746 While we draw on a synthesis of recent findings
747 in model collapse, iterated learning, and sociolin-
748 guistics, the long-term impact of LLMs on natural
749 human language and speech remains a develop-
750 ing phenomenon. The limitations and concerns
751 discussed reflect the current state-of-the-art. It is
752 possible that future architectures, perhaps those
753 incorporating neuro-symbolic reasoning or novel
754 inductive biases, can mitigate certain aspects of
755 model collapse. Our position is therefore a critique
756 of the current trajectory of generative AI rather than
757 an immutable law of artificial/machine intelligence.

758 We furthermore acknowledge an inherent selec-
759 tion bias in the literature reviewed and the exam-
760 ples provided. The author's perspective is situated
761 within an academic tradition that values linguistic
762 richness. We recognize that researchers from more
763 functionalist or engineering-driven backgrounds
764 might interpret the "flattening" of language as an
765 increase in communicative efficiency or standard-
766 ization rather than a loss of richness.

767 **References**

768 Sina Alemohammad, Josue Casco-Rodriguez, Lorenzo
769 Luzi, Ahmed Imtiaz Humayun, Hossein Babaei,
770 Daniel LeJeune, Ali Siahkoochi, and Richard G Bara-
771 niuk. 2023. Self-consuming generative models go
772 mad. *arXiv preprint arXiv:2307.01850*, 4:14.

773 Eduardo G Altmann, Janet B Pierrehumbert, and Adil-
774 son E Motter. 2011. Niche as a determinant of word
775 fate in online groups. *PloS one*, 6(5):e19009.

776 Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An
777 automatic metric for mt evaluation with improved cor-
778 relation with human judgments. In *Proceedings of*
779 *the acl workshop on intrinsic and extrinsic evaluation*
780 *measures for machine translation and/or summariza-*
781 *tion*, pages 65–72.

782 Emily M Bender, Timnit Gebru, Angelina McMillan-
783 Major, and Shmargaret Shmitchell. 2021. On the
784 dangers of stochastic parrots: Can language models
785 be too big? In *Proceedings of the 2021 ACM confer-*
786 *ence on fairness, accountability, and transparency*,
787 pages 610–623.

788 Sebastian Bernhardsson, Luis Enrique Correa da Rocha,
789 and Petter Minnhagen. 2009. The meta book and

size-dependent properties of written language. *New*
Journal of Physics, 11(12):123015.

Henry Brighton, Kenny Smith, and Simon Kirby. 2005.
Language as an evolutionary system. *Physics of Life*
Reviews, 2(3):177–226.

Lindell Bromham, Russell Dinnage, Hedvig Skirgård,
Andrew Ritchie, Marcel Cardillo, Felicity Meakins,
Simon Greenhill, and Xia Hua. 2022. Global pre-
dictors of language endangerment and the future
of linguistic diversity. *Nature ecology & evolution*,
6(2):163–173.

Yang Trista Cao, Anna Sotnikova, Jieyu Zhao, Linda X
Zou, Rachel Rudinger, and Hal Daume III. 2023.
Multilingual large language models leak human
stereotypes across language boundaries. *arXiv*
preprint arXiv:2312.07141.

Sheila Castilho, Zoe Fitzsimmons, Claire Holton, and
Aoife Mc Donagh. 2025. Synthetic fluency: Hallu-
cinations, confabulations, and the creation of irish
words in llm-generated translations. *Proceedings of*
Machine Translation Summit XVII: Research Track.

Noam Chomsky. 1980. Rules and representations. *Be-*
havioral and brain sciences, 3(1):1–15.

David Crystal. 2002. *Language death*. Cambridge
university press.

Charles Darwin. 1871. *The Descent of Man, and Se-*
lection in Relation to Sex, volume 1. John Murray,
London.

Marci S DeCaro, Robin D Thomas, and Sian L Beilock.
2008. Individual differences in category learning:
Sometimes less working memory capacity is better
than more. *Cognition*, 107(1):284–294.

Kristian Dokic, Barbara Pisker, and Bojan Radisic. 2025.
Mirroring cultural dominance: Disclosing large lan-
guage models social values, attitudes and stereotypes.
Societies, 15(5):142.

Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang Lorraine
Li, Liwei Jiang, Bill Yuchen Lin, Sean Welleck, Peter
West, Chandra Bhagavatula, Ronan Le Bras, Jena D.
Hwang, Soumya Sanyal, Xiang Ren, Allyson Et-
tinger, Zaid Harchaoui, and Yejin Choi. 2023. **Faith**
and fate: Limits of transformers on compositional-
ity. In *Advances in Neural Information Processing*
Systems 36: Annual Conference on Neural Informa-
tion Processing Systems 2023, NeurIPS 2023, New
Orleans, LA, USA, December 10 - 16, 2023.

Daniel L Everett. 2017. *How language began: The*
story of humanity's greatest invention. Liveright Pub-
lishing.

Martin Gerlach and Eduardo G Altmann. 2013. Stochas-
tic model for the vocabulary growth in natural lan-
guages. *Physical Review X*, 3(2):021006.

842	Gerd Gigerenzer and Daniel G Goldstein. 1996. Reasoning the fast and frugal way: models of bounded rationality. <i>Psychological review</i> , 103(4):650.	894
843		895
844		896
845	Yanzhu Guo, Guokan Shang, and Chloé Clavel. 2024. Benchmarking linguistic diversity of large language models. <i>arXiv preprint arXiv:2412.10271</i> .	897
846		898
847		899
848	Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. 2020. The state and fate of linguistic diversity and inclusion in the nlp world. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 6282–6293.	900
849		901
850		902
851		903
852		904
853		
854	Alina Karakanta, Mayra Nas, and Aletta G. Dorst. 2025. Metaphors in literary machine translation: Close but no cigar? <i>Proceedings of Machine Translation Summit XVII: Research Track</i> .	905
855		906
856		907
857		908
858	Simon Kirby, Tom Griffiths, and Kenny Smith. 2014. Iterated learning and the evolution of language. <i>Current opinion in neurobiology</i> , 28:108–114.	909
859		910
860		911
861	Simon Kirby, Monica Tamariz, Hannah Cornish, and Kenny Smith. 2015. Compression and communication in the cultural evolution of linguistic structure. <i>Cognition</i> , 141:87–102.	912
862		913
863		914
864		915
865	Dmitry Kobak, Rita González-Márquez, Emőke-Ágnes Horvát, and Jan Lause. 2024. Delving into chatgpt usage in academic writing through excess vocabulary. <i>arXiv preprint arXiv:2406.07016</i> .	916
866		917
867		918
868		919
869	András Kornai. 2013. Digital language death. <i>PloS one</i> , 8(10):e77056.	920
870		921
871	Tom Kouwenhoven, Max Peeperkorn, Roy de Kleijn, and Tessa Verhoef. 2025a. Shaping shared languages: Human and large language models’ inductive biases in emergent communication. <i>arXiv preprint arXiv:2503.04395</i> .	922
872		923
873		924
874		925
875		926
876	Tom Kouwenhoven, Max Peeperkorn, and Tessa Verhoef. 2025b. Searching for structure: Investigating emergent communication with large language models . In <i>Proceedings of the 31st International Conference on Computational Linguistics</i> , pages 9977–9991, Abu Dhabi, UAE. Association for Computational Linguistics.	927
877		928
878		929
879		930
880		931
881		932
882		933
883	Svenja Kranich. 2014. Translations as a locus of language contact. In <i>Translation: A multidisciplinary approach</i> , pages 96–115. Springer.	934
884		935
885		936
886	Jiaming Luo, Colin Cherry, and George Foster. 2024. To diverge or not to diverge: A morphosyntactic perspective on machine translation vs human translation. <i>Transactions of the Association for Computational Linguistics</i> , 12:355–371.	937
887		938
888		939
889		940
890		941
891	Gary Lupyan and Rick Dale. 2010. Language structure is partly determined by social structure. <i>PloS one</i> , 5(1):e8559.	942
892		943
893		944
	Gonzalo Martínez, Lauren Watson, Pedro Reviriego, José Alberto Hernández, Marc Juárez, and Rik Sarkar. 2023. Towards understanding the interplay of generative artificial intelligence and the internet. In <i>International Workshop on Epistemic Uncertainty in Artificial Intelligence</i> , pages 59–73. Springer.	945
		946
		947
		948
		949
	R Thomas McCoy, Shunyu Yao, Dan Friedman, Matthew Hardy, and Thomas L Griffiths. 2023. Embers of autoregression: Understanding large language models through the problem they are trained to solve. <i>arXiv preprint arXiv:2309.13638</i> .	950
		951
		952
	Jean-Baptiste Michel, Yuan Kui Shen, Aviva Presser Aiden, Adrian Veres, Matthew K Gray, Google Books Team, Joseph P Pickett, Dale Hoiberg, Dan Clancy, Peter Norvig, et al. 2011. Quantitative analysis of culture using millions of digitized books. <i>science</i> , 331(6014):176–182.	953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

950	Yi Ren, Shangmin Guo, Linlu Qiu, Bailin Wang, and Danica J Sutherland. 2024a. Bias amplification in language model evolution: An iterated learning perspective. <i>arXiv preprint arXiv:2404.04286</i> .	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in neural information processing systems</i> , 30.	1004
951			1005
952			1006
953			1007
954	Yi Ren, Shangmin Guo, Linlu Qiu, Bailin Wang, and Danica J Sutherland. 2024b. Bias amplification in language model evolution: An iterated learning perspective. <i>arXiv preprint arXiv:2404.04286</i> .	Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mt5: A massively multilingual pre-trained text-to-text transformer. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> . Association for Computational Linguistics.	1009
955			1010
956			1011
957			1012
958	Pedro Reviriego, Javier Conde, Elena Merino-Gómez, Gonzalo Martínez, and José Alberto Hernández. 2024. Playing with words: Comparing the vocabulary and lexical diversity of chatgpt and humans. <i>Machine Learning with Applications</i> , 18:100602.		1013
959			1014
960			1015
961			1016
962			
963	Yu Sasaki. 2017. Publishing nations: Technology acquisition and language standardization for european ethnic groups. <i>The Journal of Economic History</i> , 77(4):1007–1047.	Yongjing Yin, Lian Fu, Yafu Li, and Yue Zhang. 2024. On compositional generalization of transformer-based neural machine translation. <i>Information Fusion</i> , 111:102491.	1017
964			1018
965			1019
966			1020
967	Ilia Shumailov, Zakhar Shumaylov, Yiren Zhao, Yarin Gal, Nicolas Papernot, and Ross Anderson. 2023. The curse of recursion: Training on generated data makes models forget. <i>arXiv preprint arXiv:2305.17493</i> .	Xiang Zhou, Yichen Jiang, and Mohit Bansal. 2023. Data factors for better compositional generalization. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 14549–14566.	1021
968			1022
969			1023
970			1024
971			1025
972	Kenny Smith, Henry Brighton, and Simon Kirby. 2003a. Complex systems in language evolution: the cultural emergence of compositional structure. <i>Advances in complex systems</i> , 6(04):537–558.		
973			
974			
975			
976	Kenny Smith, Simon Kirby, and Henry Brighton. 2003b. Iterated learning: A framework for the emergence of language. <i>Artificial life</i> , 9(4):371–386.		
977			
978			
979	Matthew Snover, Bonnie Dorr, Richard Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In <i>Proceedings of the 7th Conference of the Association for Machine Translation in the Americas: Technical Papers</i> , pages 223–231.		
980			
981			
982			
983			
984			
985	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .		
986			
987			
988			
989			
990			
991	Eva Vanmassenhove, Dimitar Shterionov, and Matthew Gwilliam. 2021. Machine translationese: Effects of algorithmic bias on linguistic complexity in machine translation. In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 2203–2213.		
992			
993			
994			
995			
996			
997			
998	Eva Vanmassenhove, Dimitar Shterionov, and Andy Way. 2019. Lost in translation: Loss and decay of linguistic richness in machine translation. In <i>Proceedings of Machine Translation Summit XVII: Research Track</i> , pages 222–232, Dublin, Ireland. European Association for Machine Translation.		
999			
1000			
1001			
1002			
1003			