# Language Models Model Language

**Łukasz Borchmann**
Snowflake AI Research
`lukasz.borchmann@snowflake.com`

## Abstract

Linguistic commentary on LLMs, heavily influenced by the theoretical frameworks of de Saussure and Chomsky, is often speculative and unproductive. Critics challenge whether LLMs can legitimately model language, citing the need for "deep structure" or "grounding" to achieve an idealized linguistic "competence." We argue for a radical shift in perspective towards the empiricist principles of Witold Mańczak, a prominent general and historical linguist. He defines language not as a "system of signs" or a "computational system of the brain" but as the totality of all that is said and written. Above all, he identifies frequency of use of particular language elements as language's primary governing principle. Using his framework, we challenge prior critiques of LLMs and provide a constructive guide for designing, evaluating, and interpreting language models.

## 1 Introduction

For two thousand years, language scholars have produced myriad works without clearly defining a way to validate their claims (Mańczak, 1981, 1996c). This lack of rigor left the field uniquely unprepared for the empirical success of Large Language Models (LLMs).

When GPT or Claude synthesizes complex responses, most users of LLMs perceive this ability as a sign of linguistic competence. At the same time, orthodox linguists rush to claim that apparent "competence" is merely an illusion, because LLMs' architecture and training techniques do not fit established theories of language. These objections reveal more about current linguistic dogma than about the capabilities of LLMs.

Debate remains anchored in theoretical frameworks inherited from de Saussure's abstract "system of signs," Chomsky's postulated innate mental grammar, and even Plato's conception of language as a communication tool. When scholars criticize LLMs for failing to "explain the rules of English syntax" (Chomsky et al., 2023), point out that they don't distinguish "between correctness and likelihood" (Fox and Katzir, 2024), or question whether they can "have access to meaning" despite experiencing only form (Bender et al., 2021), these researchers apply criteria derived from specific—and contestable—theories. Such critiques function less as objective assessments of LLMs' capacities than as defenses of a linguistic paradigm that has been challenged by new empirical evidence.

We propose a radical theoretical reorientation and turn to Witold Mańczak (1924-2016), whose critique of abstract, dualistic linguistics offers an alternative foundation. While Bender et al. (2021) advocate constraining the LLMs' tendency to act like "stochastic parrots," we call for a new science of ornithology that is equipped to understand what has actually taken flight.

## 2 The Mańczakian Framework

The Mańczakian framework is best understood as a response to what its author saw as a foundational crisis in linguistics: the absence of scientific criteria for determining "truth." Mańczak argued that the field had long substituted the authority of famous scholars for the authority of verifiable evidence (Mańczak, 1980, 1981, 1982, 1988a, 1996c).

The central goal of his work was to reground linguistics as an inductive, quantitative science. He insisted that any valid claim must be subject to statistical or experimental verification and explained by observable phenomena rather than abstract speculation. Applying quantitative methods to vast corpora, he discovered that frequency of use of language elements is a primary force shaping language.

We will now unpack these core tenets, demonstrating how the empiricist lens of Mańczak re-

frames long-standing linguistic debates and provides a robust foundation for the science of language and language models.

## 2.1 The Map Is Not the Territory

Mańczak found it absurd that people would claim to study an intangible "system of signs" or mental "competence" when their actual work involved poring over concrete texts. He proposed a simple, naturalistic definition: language is the totality of all that is said and written.

**Language Definition**

The fundamental error of modern linguistics is the false equivalence between the product of a particular analysis and the object of the study itself (Mańczak, 1969a).

**The Primacy of Frequency.** Where structuralists and their successors saw dichotomies, Mańczak saw a continuum governed by statistics. For example, in his view, the dichotomy between grammar and lexicon is false, because grammar is "the quintessence, condensation, abbreviation, or generalization of lexicography" (Mańczak, 1996h). Grammar covers only high-frequency patterns that apply to large classes of words. Lexicography, in contrast, handles the information for all individual words, including the less frequent and more idiosyncratic.

Even the distinction between "rule" and "exception" is quantitative, not qualitative: high-frequency patterns are rules, while low-frequency patterns are exceptions. We can interpret Mańczak's perspective in the light of Jaynes' theory, according to which the absolute "True/False" statements of classic logic are seen as a special case of a more general, probabilistic logic (Jaynes, 2003). In this framing, a Mańczakian "rule" corresponds to a high-plausibility inference, while an "exception" represents a low-plausibility inference.

**Frequent Errors Become Norms**

Importantly, according to Mańczak (1969a) "the rules of grammar, once abstracted from texts, do not prevent the subsequent evolution of language. This evolution consists of errors which, if their frequency increases sufficiently, become new norms. Conversely, old norms, if their frequency diminishes, become errors."

Example 1 presents three examples from the evolution of Latin into Romance languages that illustrate how frequency drives linguistic change and, crucially, how grammatical rules themselves are merely abstractions from this ongoing evolution.

**Language Acquisition.** Mańczak's claim that frequency of use in language is a primary organizing principle is strongly supported by decades of research in cognitive science. Studies show that frequency affects nearly every level of language processing in humans, from the recognition of sounds and words to the processing of syntax (Saffran et al., 1996; Romberg and Saffran, 2010; Fló et al., 2025, *inter alia*). High-frequency patterns are processed more quickly and accurately, and learned earlier by children, confirming that the human mind is finely tuned to the frequency of exposure to particular language elements (Ellis, 2002).

This statistics-based view of language contrasts strikingly with dominant linguistic paradigms. As Mańczak noted, "the fact that, in the three hundred pages of [de Saussure's] *Course in General Linguistics* the term *frequency of use* does not appear even once, has serious implications" (Mańczak, 1969a). When frequency is ignored, linguists must invent increasingly elaborate theoretical constructs—deep structures, innate faculties, universal grammars— to explain what simple statistical patterns readily predict.

**Result of Ignoring Frequency**

## 2.2 Talkin' All About That Syntax When You Can't Even Generate a Haiku

Once invented, these constructs then take on a life of their own. Critics of LLMs now assert that *true* linguistic competence requires an innate "deep structure" or "internal form of language." This claim mistakenly interprets the linguist's abstract analysis as the functional prerequisite of language itself. As Mańczak (1996a) puts it:

> If a layman were to try to convince a chemist that water is composed not of two elements (i.e., oxygen and hydrogen) but of four, the chemist would shrug his shoulders and conclude that the third and fourth elements exist only in the layman's imagination, because chemists have succeeded countless times in synthesizing water from only two elements. Now, the principle that synthesis [reconstructing a coherent whole] is required to validate analysis [breaking things down into components] . . . should be applied to linguistics. Linguists who analyze text should look for only those components that are essential for its synthesis and must validate their analyses by means of synthesis.

He similarly criticized proponents of the generative and transformational grammars that dominated

Table 1: Comparison of dominant and Mańczakian linguistic paradigms.

| Aspect | Dominant Paradigms | Mańczakian Paradigm |
|---|---|---|
| **Definition of Language** | An abstract "system of signs" (de Saussure) or mental "competence" distinct from performance (Chomsky). | The material totality of all that is said and written—the actual corpus of texts. |
| **Organizing Principle** | Innate universal grammar; binary structural oppositions; algebraic rules for generating well-formed sentences. | Frequency of use. High-frequency patterns become rules; low-frequency patterns are exceptions on the same continuum. |
| **Status of Grammar** | A pre-existing mental mechanism in the mind. | A post-hoc descriptive abstraction created by linguists from observed patterns in texts. |
| **Validation Method** | Subjective sense of grammaticality; theoretical elegance; appeal to linguistic intuition and dogma. | "Synthesis validates analysis." A theory is valid only if it can be used to reconstruct what it claims to explain. |
| **View of Meaning** | Requires grounding in deep structure or real-world referents. Form and meaning are fundamentally separate. | Primarily relational, derived from network of connections between words. Most meanings exist within the axiomatic system of language. |
| **Assessment of LLMs** | Fundamentally flawed "stochastic parrots" lacking access to deep structure or grounded meaning. | A large-scale, empirical validation. Not broken in principle but incomplete in practice. |

---

**Three studies from Latin's evolution into Romance languages**

**(A) Analogical Regularization.** Classical Latin used a mix of rules to form words representing quantities. While 17 was additive (*septendecim* 'seven-ten'), 18 and 19 were subtractive (*duodeviginti* 'two-from-twenty,' *undeviginti* 'one-from-twenty'). In the evolution of Latin to the Romance languages, these irregular subtractive forms were almost universally abandoned. These were replaced by more regular, compound forms analogically derived from simpler additive rules (e.g., Italian *diciotto* 'ten-and-eight'). Additive rules served as the model for simplifying the rarer, more complex subtractive ones (Mańczak, 1958, 1978, 1996e).

**(B) Irregular Development due to Frequency.** The evolution of the Latin verb *ambulare* 'to walk' illustrates how usage frequency shapes language. The form of the word diverged based on its function. Everyday use resulted in significant irregular shortening, leading to French *aller*, Italian *andare*, and Spanish *andar*. Infrequent use in a specific context—e.g., 'to amble,' which began as a specialized equestrian term—retained a more regular and archaic form: French *ambler*, Italian *ambiare*, and Spanish *amblar*. The fact that high-frequency words are prone to irregular change, while less common words often preserve their regular, historical forms, is a common linguistic phenomenon (Mańczak, 1969b,c, 1987, 1988b, 1996d, *inter alia*).

**(C) Grammaticalization.** A fundamental shift in the history of Romance languages is the emergence of the 'have' + participle construction to express the perfect tense (e.g., French *il a chanté* 'he has sung'). This developed from Classical Latin, where the verb *habere* was a normal lexical verb meaning 'to possess' or 'to hold'. This new, constant grammatical use created a functional split. The word's original meaning has been mostly replaced by its modern role as a grammatical tense marker (Mańczak, 1996g).

Example 1: Examples from the evolution of Latin into the Romance languages, demonstrating the impact of frequency and temporal character of grammar: (A) analogical regularization of numerals, (B) divergent evolution depending on usage frequency, and (C) grammaticalization from lexical verb to perfect tense marker.

20th-century linguistics. Mańczak noted with irony that while champions of these theories filled entire volumes with abstract formalisms, their analysis was never validated by synthesis. More than half a century after the "generative" turn, its adherents had "not yet written a single generative or transformational grammar of a concrete language" (Mańczak, 1996b). The roar of generative theory

had produced not even a whisper of practical synthesis.

This failure becomes evident if we review the actual analyses. As Mańczak pointed out, "Chomsky, to analyze a sentence as simple as *sincerity may frighten the boy*, needed 10 pages in his book *Aspects of the Theory of Syntax*, whereas to reconstruct this sentence, it is enough to cite only five simple positional rules" as illustrated in Example 2 (Mańczak, 1996b). The principle that should guide any scientific grammar—"synthesis validates analysis"—was, as Mańczak observed, "completely unknown to Chomsky."

Yet linguists tolerated lack of validation for decades, mistaking Chomsky's theoretical complexity for demonstrated accuracy. At the same time, when they needed to assess grammaticality, they searched corpora or polled native speakers instead of consulting a generative rule system. Today, practitioners of linguistics already operate in Mańczak's world, even if linguistic theorists haven't caught up.

---

**Validation through synthesis**

**The analysis of a simple sentence.** Mańczak examined how linguists with competing methodologies analyze a simple English sentence: *Sincerity may frighten the boy*. He observed that creating such a sentence using generative grammar required a long, convoluted process.

In contrast, to synthesize the sentence, one needs only a few simple rules: "1) Every noun can be accompanied by a third-person verb that agrees with it in number. 2) Certain verbs, including *may*, can be accompanied by infinitives without the particle *to*. Since these verbs are rare, grammars enumerate them. 3) Transitive verbs (including *frighten*) are accompanied by a noun. Since these verbs are common, grammars refer to dictionaries for detailed information. 4) Every noun can be accompanied by an article under certain conditions. Grammars enumerate these conditions. 5) The article precedes the noun, the infinitive follows the verb it accompanies, etc." (Mańczak, 1996b).

These straightforward rules highlight the gap between the theoretical complexity of Chomsky's analysis and the practical utility of Mańczak's alternative approach.

---

Example 2: The theory requiring immense abstract machinery to analyze a simple sentence that can be generated with a few practical rules has failed the fundamental test of synthesis.

Stunningly effective LLMs have arrived as the answer to the unfulfilled promise of generative grammar. Their capacity to synthesize coherent language based on probabilistic analysis of text rather than abstract rules or postulated deep structures is the ultimate vindication of Mańczak's approach.

Empirically, LLM performance increases smoothly with the amount of pretraining data (Kaplan et al., 2020; Hoffmann et al., 2022). Pretraining minimizes expected next-token surprisal (*cross-entropy*), pushing the model's conditional predictions to match empirical next-token frequencies. Estimation of the language's frequency structure improves and sharpens with a larger training set, especially in the long tail. Interestingly, models naturally retain the useful token frequency information in embeddings (Zhou et al., 2021; Gong et al., 2018).

LLMs' success is not a mystery or a "stochastic parrot" trick. It is a large-scale validation of Mańczak's central thesis: language is text, and frequency is not a secondary, peripheral aspect, but its primary organizing force.[1]

## 2.3 An Obituary for the Language Organ

According to the Chomskyan paradigm, language acquisition requires a specialized "language organ" because children lack sufficient linguistic input to develop language from pattern recognition alone. Mańczak rejected this view, arguing that debates about hypothetical brain structures fell outside the proper naturalistic focus of linguistics (Grochowski, 2017; Mańczak, 1969a).

We see now that his skepticism was justified. The Chomskyan paradigm has been challenged by evidence from cross-linguistic research and developmental psychology, causing many experts to abandon it (Ibbotson and Tomasello, 2016; Pullum and Scholz, 2002; Christiansen and Chater, 2008; Snow, 1972; Tomasello, 2003; Bybee, 2006, *inter alia*). Instead, a "usage-based" alternative has emerged. Its proponents argue that children learn grammar from the ground up by applying general cognitive tools like pattern recognition and categorization to the sounds they hear. This evidence-based theory sees grammar as an emergent property of history and psychology—a set of templates discovered by observing that some sentences or word forms are built analogically. It aligns perfectly with

---

[1]For a complementary view, see Piantadosi (2024), who argues that LLMs refute Chomsky's theory.

Table 2: Architectural evolution and its impact on analogical capability.

| Model | Core Mechanism | Representation of Words | Analogical Capability |
|---|---|---|---|
| **N-gram** | Surface-level counting of token co-occurrence; $P(\text{word} \mid \text{previous } n-1 \text{ words})$. | Atomic, discrete symbols. No notion of similarity between words. Sentences "Anna likes cats" and "Lily loves dogs" are not similar. | **None.** Can reproduce frequent sequences but cannot grasp words or sentences similarity. |
| **CBOW** | Learning vector representation of meaning by predicting a word from its context. | Dense, low-dimensional vectors. Words occurring in similar contexts are mapped to nearby points in a geometric space. | **External to model.** Series of analogous word pairs can be demonstrated to share similar geometric relationship. |
| **LLMs** | Learning relationships between learned vector representations of meaning. | Dense, high-dimensional vectors that are dynamically adjusted based on the specific context of the entire sequence. | **Inherent, core mechanism.** The model operates on a series of embeddings and generalizes patterns. |

Mańczak's insistence on an empirical, text-first science of language.[2]

To criticize the ability of LLMs to generate plausible human speech as a mere application of learned patterns is to miss the point that use of analogy is an essential aspect of linguistic competence. LLMs can generate plausible real-time answers in part because they recognize—during previous training—the same frequency patterns that created human grammar in the first place.

**Key Leap: *Relations* Between Embeddings**

The success of modern LLMs stems from a crucial architectural innovation: the replacement of flat, surface-level n-gram counts with high-dimensional embeddings (Table 2). Earlier n-gram models, limited by their reliance on simple tables of memorized word combinations, failed to recognize that the sentences 'Anna likes cats' and 'Lily loves dogs' are analogous. While CBOW models could demonstrate that analogous word pairs share a similar geometric relationship, the breakthrough improvement was to operate on sequences of learned vector representations. When presented with a novel problem to solve, Transformer uses its vast internalized map of relationships to find and apply the closest learned analogy.

This ability to represent and manipulate relationships—the very essence of analogy—is the key to genuine linguistic generalization.

### 2.4 This Page is Intentionally Left Ungrounded

Perhaps the most persistent criticism of LLMs is that they are ungrounded: they manipulate symbols (*form*) without access to their real-world referents (*meaning*).[3] While Mańczak did not directly engage with modern semantic theories, his view aligns with a core assumption of componential and reductionist semantics: the necessity of undefinable primitives (Grochowski, 2017).

**Axiomatic Meaning**

Mańczak argued that attempting to describe language without external reference (ignoring meaning *entirely*) was a descent into a "nebulous darkness." He proposed a simple resolution (Mańczak, 1996h):

> Just as in mathematics, most but not all statements can be proved (with the help of other statements) … most words in a given language can be defined with the help of other words, with the unavoidable exception that [to avoid circularity] the meaning of certain words must be taken as self-evident (axiomatic).

Nevertheless, he recognized that the meaning of most words is relational, derived from an intricate web of connections between terms.[4]

An LLM correctly using a concept like "justice" requires only mastery of the vast, multidimensional web of relationships that connect that word to "fairness," "law," "equality," "crime," and thousands of other terms. Whether this represents a human-like *understanding* is irrelevant.

We do not demand that a calculator understands what "1+2" *truly means* to accept its utility. We do

---

[2]See Goldberg (2024) for a usage-based constructionist account explicitly linking constructions to LLMs.

[3]See Bender and Koller (2020) for a representative formulation of the form-vs-meaning objection.

[4]Concepts without referents—such as "perpetual motion machine" or "king of San Francisco"—further illustrate how meaning can exist purely through relational networks (Piantadosi and Hill, 2022).

not dismiss the results of theorem-proving software because it cannot understand the philosophical basis of Zermelo-Fraenkel set theory. While simple counting may be grounded in early childhood experiences, advanced mathematical reasoning involves manipulation of a highly formal system according to a set of rules. For most people, including professional mathematicians, the "meaning" of an axiom lies in its role within that formal system, rather than an intuitive or sensory experience.

Demanding that a language model meet a higher standard of "grounded meaning" is misguided. In the Mańczakian view, the relevant test for LLM quality is not whether the model has access to an outside world, but whether it has mastered the internal, relational logic of the textual world it was given.[5]

## 3  Conclusion

LLM users often feel amazed, disappointed, or a mixture of both. They're amazed that LLMs can produce human-like text about any topic. They're disappointed when models confidently present factually incorrect information.

Beyond these reactions, prominent linguists argue that LLMs are intrinsically flawed because of gaps in linguistic competence and inability to "understand" language. They call for models to demonstrate knowledge of "deep structure" or reliance on "grounded meaning"—concepts their theories treat as requirements of language.

We argue that by applying these theory-laden standards, such critics interfere with useful analysis. The "stochastic parrots" do not merely mimic language but in fact reveal what language has been all along. LLMs are imperfect tools not because they *fail* to model language but because they *only* model language.

The implications of our proposed shift in perspective extend beyond theoretical linguistics. The path to improving *linguistic* competence of LLMs is to design, evaluate, and deploy systems that have already proven mastery of language's relational logic. Satisfying abstract theoretical requirements for language competence—demanded by the vocal critics of LLMs—will neither improve linguistic

---

[5]This textual world might suffice. Mandelkern and Linzen (2024) argue that LLMs can use words to talk about real things because training texts already link those words to the physical world. In their view, the key issue is whether the model counts as part of our speech community.

competence nor bring us closer to Artificial General Intelligence.

Mańczak saw it clearly: language is the totality of texts, and frequency is its organizing principle. Decades later, engineers unknowingly built LLMs on this very foundation. Their models now draft our contracts, structure our arguments, and increasingly impact our futures—making Mańczak's case more powerfully than any argument could.

## 4  Limitations

We acknowledge that proposing a novel theoretical lens for understanding language and LLMs leaves some critical questions unaddressed.

**Your paper ignores urgent ethical dilemmas.** Carefully defining a technology is a prerequisite for ethical debate. We address the former to set the stage for the latter.

**A text-only definition of language excludes some linguistic subfields, such as pragmatics, psycholinguistics, etc.** To qualify as linguistic, a claim must generate falsifiable predictions about *distributions over texts* or about *observable properties of utterances*, all of which must be testable by corpus statistics or experiments. Pragmatics satisfies these conditions when it predicts context-conditioned choices among forms (e.g.,

---

### Witold Mańczak (1924–2016)

Witold Mańczak was a Polish linguist specializing in Indo-European, Romance, and Slavic studies. Over his career, he published 24 books and more than 960 other texts, developing a distinct theoretical framework for language analysis (Dębowiak, 2014, 2016; Grochowski, 2017).

A defining feature of Mańczak's work was his application of statistical methods to verify linguistic hypotheses. He argued that linguistic analysis should be based on quantifiable data, and proposed that language development is driven by three factors: regular phonetic development, analogical development, and irregular phonetic development caused by frequency of use. With his data-centric approach, Mańczak challenged several foundational linguistic theories (see Appendix A).

His life's work stands as a testament to a tireless pursuit of verifiable statements, establishing him as the creator of a coherent, statistics-based theory of language and one of the most significant linguists of his era.

hedging, politeness markers). Sociolinguistics does so when it predicts socially conditioned alternations in variants and styles across communities. Ethnolinguistics is linguistic when it predicts culturally conditioned textual patterns in a community corpus (e.g., motifs, metaphor families, collocational frames).

Claims that fail to meet these conditions may still be valuable—in sociology, psychology, or cultural studies—but they do not count as linguistic in the Mańczakian sense.

**Your focus on relational meaning seems to embrace a purely structuralist view of semantics.** We do not claim to offer a complete theory of meaning. Our central thesis is that the vast majority of the time "meaning" can be inferred (and in the case of LLMs is inferred) solely from the relational structure of the text. The success of LLMs demonstrates how much linguistic competence can be achieved by assigning a central role to distributional semantics, and a limited role—if any—to direct grounding.

**LLMs are trained on unrepresentative corpora, very distant from "the totality of all that is said and written."** This is true and aligns with our view that LLMs are not fundamentally flawed, but incomplete in practice. The Mańczakian framework offers a clear path forward, favoring "rational selection of texts" based on circulation and influence (Mańczak, 1996f, 1961). A truly Mańczakian approach is a principled, frequency-weighted corpus construction that reflects how language is actually used.

**LLMs lack a "world model."** The Mańczakian framework we adopt posits that language is not a map of the physical world but a self-contained universe of texts. From this perspective, the demand for grounding in physical reality is misplaced. Our model should be seen not as a flawed attempt to simulate a mind interacting with a physical world but as a successful and direct model of language itself.

**You are closing off valuable inquiry into the cognitive plausibility of LLMs.** Our goal is not to declare cognitive comparisons invalid but to argue that they should not be the *primary* benchmark. An LLM is, first and foremost, a direct model of textual corpus. This is not a metaphor; it is a description of its construction and function. We argue that any comparison to human cognition must come *after* this primary, non-metaphorical evaluation is established. The Mańczakian framework provides a necessary baseline.

**Many of Mańczak's core ideas are central to the well-established "usage-based linguistics."** While Mańczak's ideas do align with modern usage-based linguistics, he reached similar conclusions decades earlier through the empirical analysis of contemporary and historical texts.

First, Mańczak's radical, text-only simplicity provides a direct rebuttal to critiques of LLMs grounded in Saussurean or generative theories: his framework simply discards abstractions that cannot be found in the textual record. Second, his focus on the structure of the input text itself—rather than human cognitive processing—aligns directly with how text-trained models actually function.

The fact that Mańczak's textual analysis independently yielded insights that are now central to cognitivism makes his framework particularly compelling.

**If linguistic competence requires merely applying frequency-weighted patterns, how do you explain creativity?** High-frequency patterns don't mechanically reproduce themselves; they serve as templates for novel combinations. LLMs demonstrate this empirically. Creativity isn't the opposite of pattern utilization—it's pattern mastery.

## Acknowledgments

# A  Selected Contributions to Historical and Comparative Linguistics

This summary was initially published by Mańczak in French and served as a foreword to *Linguistique générale et linguistique indo-européenne*. It outlines his career-long pursuit against fundamental methodological flaws of modern linguistics. Some of these findings oppose prevailing views—and that is precisely why his methodological proposal cannot be ignored. I strongly encourage the reader to seek out the original publications.

## English Translation

The fundamental problem of linguistics is that of the criteria of truth. Unfortunately, this problem is taboo. Given that linguistics has existed for two thousand years and that the *Linguistic Bibliography* recorded 21,000 works for the year 2001, it follows that linguists have published, in total, several hundreds of thousands of works, and yet none of these has been devoted to the criteria of truth. Even the term "criteria of truth" is never used by linguists. This is an extraordinary thing, considering that linguists unanimously agree that linguistics is a science, and that science is a search for truth. Why do linguists keep the question of distinguishing true from false in their discipline secret?

This enigma has intrigued me for a long time. As linguistic works provide no information capable of resolving this question, I began to observe how linguists react when they learn an opinion previously unknown. To my great astonishment, I found that linguists never intend to verify the opinion in question, but are interested only in the question of who shares this opinion. If they learn that this opinion is shared by one or more authorities, they consider it true. If, on the contrary, they learn that this opinion comes from someone who does not have the reputation of being an authority, this view appears false to them. The criterion of truth used by linguists is the following: X has formulated an opinion, X is an authority, therefore this opinion is true; Y has formulated an opinion, Y is not an authority, consequently this opinion is false. Obviously, this criterion of truth is medieval and unscientific, which is why linguists prefer not to talk about it.

In this state of affairs, I reflected on the criteria of truth likely to be employed in linguistics. I came to the conclusion that linguists can resort to statistics (and, exceptionally, to experiment) and that, in the science of language, many opinions rely on faith in the infallibility of authorities and are invalidated by statistical data. Here are some examples:

1. In all languages, the form of words depends on three main factors, not only on regular phonetic development and analogical development, but also on what I call irregular phonetic development due to frequency (*Le développement phonétique des langues romanes et la fréquence*, Kraków, 1969; *Słowiańska fonetyka historyczna a frekwencja*, Kraków, 1977; *Frequenzbedingter unregelmäßiger Lautwandel in den germanischen Sprachen*, Wrocław, 1987). According to a professor of applied mathematics, the chance that the theory of irregular phonetic development due to frequency is erroneous is less than 1 in 10 million (*Etymologia przyimka dla a nieregularny rozwój fonetyczny spowodowany frekwencją*, Prace Filologiczne 60, 2011, p. 189–195).

2. Bartoli's "norm" according to which lateral areas are more archaic than central areas is invalidated by statistical data (*La Roumanie et l'Espagne sont-elles des territoires archaïques de la Romania?*, Limba română, limbă romanică. Omagiu acad. M. Sala la împlinirea a 75 de ani, Bucureşti, p. 313–317).

3. Since 1925, when Meillet introduced the notion of "empty slot" (*case vide*), it has been imagined that phonetic evolution consists of filling "empty slots" in phonological systems. But I have examined a large number of facts and have come to the conclusion that it is not symmetry, but asymmetry that characterizes languages, that it is possible to formulate a law according to which more frequently used linguistic elements are more differentiated than less used elements (*Do the "cases vides" exist?*, Linguistique générale et linguistique indo-européenne, Kraków, 2008, p. 59–62).

4. Laryngeal theory is invalidated by statistical data (*Critique de la théorie des laryngales*, Analecta Indoeuropaea Cracoviensia I. Safarewicz memoriae dicata, Cracoviae, 1995, p. 237–247; *Encore un argument contre la théorie des laryngales*, Lingua Posnaniensis 46, 2004, p. 41–44).

5. In my opinion, Verner's law requires revision (*La restriction de la règle de Verner à la position médiane et le sort du s final en germanique*, Historische Sprachforschung 103, 1990, p. 92–101; *La règle de Verner s'applique-t-elle à la position*

*finale?*, Historische Sprachforschung 109, 1996, p. 110–116).

6. In light of statistical data, Old Church Slavonic is a compromise between the Macedono-Bulgarian dialect and the Moravo-Pannonian speech (*Przedhistoryczne migracje Słowian i pochodzenie języka staro-cerkiewno-słowiańskiego*, Kraków, 2004; *Pochodzenie języka staro-cerkiewno-słowiańskiego a Kodeks zografski*, Warszawa, 2006).

7. In light of statistical data, the original homeland of the Indo-Europeans is identical with that of the Slavs (*De la préhistoire des peuples indo-européens*, Kraków, 1992; *L'habitat primitif des Indo-Européens se trouvait-il vraiment en Arménie?*, Folia Orientalia 33, 1997, p. 65–74).

8. The German orientalist Ludolf (17th c.) was the first to affirm that "die Sprachverwandtschaft offenbart sich nicht im Wörterbuch, sondern in der Grammatik" [language relationship reveals itself not in the dictionary, but in grammar]. But one can justify the division of Indo-European languages into Germanic, Slavic, Baltic, Romance, etc. only through lexical convergences, and not inflectional or phonetic ones (*La classification des langues romanes*, Kraków, 1991, p. 22–36).

9. The number one problem of Romance etymology is that of verbs meaning "to go": Fr. *aller*, It. *andare*, Sp. *andar*, Prov. *ana*, etc. Since the 16th century, about sixty etymologies have been proposed in total, which is a record, and not only for Romance etymology. Among researchers, there are adherents to monogenesis (affirming that all these forms come from, for example, *ambulare*) and advocates of polygenesis (claiming, for example, that *aller < *advehulare, andar < *am(bi)vehitare, ana < *amvehinare*). Probability calculus allows this question to be decided in favor of monogenesis (*Une étymologie romane controversée: aller, andar, etc.*, Revue roumaine de linguistique 19, 1974, p. 89–101; *Étymologie de fr. aller, esp. andar, etc. et calcul des probabilités*, Revue roumaine de linguistique 20, 1975, p. 735–739).

10. Since 1435, it has been affirmed that Romance languages derive from Vulgar Latin, but, in light of statistical data, they derive from Classical Latin (*Le problème de l'origine des langues romanes dans le livre de H. Lüdtke et celui de R. Kiesler*, Actes du XXVe Congrès International de Linguistique et de Philologie Romanes, t. VI, Berlin, 2010, p. 207–211).

11. Since Jordanes, that is, for 1400 years, it has been estimated that the original homeland of the Goths was in Scandinavia. But comparison of parallel texts in Gothic, High German, Middle German, Low German, Danish and Swedish revealed that the original homeland of the Goths was in the southernmost part of ancient Germania (*Le mythe de l'origine scandinave des Goths*, L'art de la philologie. Mélanges en l'honneur de L. Löfstedt, Helsinki, 2007, p. 137–145).

12. The division of words into stressed and unstressed (articles, pronouns, prepositions, etc.), which dates back to Antiquity, is the result of a false generalization. It is true that there are homonymies *le vent = levant*, *à voir = avoir*, and *moi = émoi* and that the syllables *le-, a-, é-* in *levant, avoir, émoi* are unstressed, but it is erroneous to conclude from this that *le, à,* and are unstressed because "stressed" words are treated in the same way. *Dix vers, vingt cœurs, va tôt*, pronounced without pauses, are homonymous with *divers, vainqueur, Watteau*, where the syllables *di-, vain-, Wa-* are unstressed. It is affirmed that *Long vient* = stressed word + stressed word, while *l'on vient* = proclitic + stressed word, but a very simple experiment proves that these expressions are homonymous (*La division des mots en toniques et atones est-elle justifiée?*, Lingua Posnaniensis 32–33, 1991, p. 181–185).

13. Since Antiquity, the question of what constitutes the difference between proper nouns and common nouns has been discussed. About ten definitions of the proper noun have been proposed so far, none of which applies to all proper nouns. In my opinion, the difference between proper nouns and common nouns consists in the fact that common nouns are, in the vast majority of cases, translated from one language to another, while proper nouns almost never are. For example, a common noun like *ville* is translated into Italian as *città*, into English as *town*, etc., whereas a proper noun like *Paris* is not, cf. It. *Parigi*, Eng. *Paris*, etc. Among all definitions of the proper noun, mine suffers the fewest exceptions (*La notion de nom propre*, Proceedings of 13th International Congress of Onomastic Sciences, Kraków, 1982, p. 101–106).

## French Original

Le problème fondamental de la linguistique est celui des critères de vérité. Malheureusement, cette question constitue un tabou. Étant donné que la linguistique existe depuis deux mille ans et que la Bibliographie linguistique a enregistré, pour l'année 2001, 21 000 travaux, il en résulte que les linguistes en ont publié, au total, plusieurs centaines de milliers, et pourtant aucun de ces derniers n'a été consacré aux critères de vérité. Même le terme « critères de vérité » n'est jamais employé par les linguistes. C'est une chose extrêmement étrange, si l'on considère que les linguistes sont unanimes pour dire que la linguistique est une science, et que la science n'est pas autre chose qu'une recherche de la vérité. Pourquoi donc les linguistes gardent-ils un secret sur la question de savoir comment ils distinguent le vrai du faux dans leur discipline?

Cette énigme m'a intrigué depuis longtemps. Comme les travaux linguistiques ne fournissent aucun renseignement susceptible de résoudre cette question, j'ai commencé à observer comment les linguistes réagissent quand ils apprennent une opinion qui leur était inconnue auparavent. A mon grand étonnement, j'ai constaté que les linguistes n'ont jamais l'intention de vérifier l'opinion en question, mais s'intéressent uniquement à la question de savoir qui partage cette opinion. S'ils apprennent que cette opinion est partagée par une ou plusieurs autorités, ils considèrent cette opinion comme vraie. Si, au contraire, ils apprennent que cette opinion provient de quelqu'un qui n'a pas la réputation d'être une autorité, cette vue leur paraît fausse. Il en résulte que le critère de vérité utilisé par les linguistes est le suivant: X a formulé une opinion, X est une autorité, par conséquent cette opinion est vraie; Y a formulé une opinion, Y n'est pas une autorité, par conséquent cette opinion est fausse. Évidemment, ce critère de vérité est médiéval, non scientifique, et c'est la raison pour laquelle les linguistes préfèrent ne pas en parler.

Dans cet état de choses, il m'est venu à l'esprit de réfléchir sur les critères de vérité susceptibles d'être employés en linguistique et je suis arrivé à la conclusion que les linguistes peuvent recourir à la statistique (et, exceptionnellement, à l'expérience) et que, dans la science du langage, il y a beaucoup d'opinions qui s'appuient sur la foi en l'infaillibilité des autorités et qui sont infirmées par des données statistiques. Voici quelques exemples.

1. Dans toutes les langues, la forme des mots dépend de trois facteurs principaux, non seulement du développement phonétique régulier et du développement analogique, mais aussi de ce que j'appelle un développement phonétique irrégulier dû à la fréquence (*Le développement phonétique des langues romanes et la fréquence*, Kraków, 1969; *Słowiańska fonetyka historyczna a frekwencja*, Kraków, 1977; *Frequenzbedingter unregelmäßiger Lautwandel in den germanischen Sprachen*, Wrocław, 1987). De l'avis d'un professeur de mathématiques appliquées, la chance que la théorie du développement phonétique irrégulier dû à la fréquence soit erronée, est moindre que 1 sur 10 millions (*Etymologia przyimka dla a nieregularny rozwój fonetyczny spowodowany frekwencją*, Prace Filologiczne 60, 2011, p. 189–195).

2. La « norme » de Bartoli d'après laquelle les aires latérales sont plus archaïques que les aires centrales, est infirmée par des données statistiques (*La Roumanie et l'Espagne sont-elles des territoires archaïques de la Romania?*, Limba română, limbă romanică. Omagiu acad. M. Sala la împlinirea a 75 de ani, Bucureşti, p. 313–317).

3. Depuis 1925, où Meillet a introduit la notion de « case vide », on imagine que l'évolution phonétique consiste à remplir des « cases vides » dans les systèmes phonologiques. Mais j'ai examiné un grand nombre de faits et suis arrivé à la conclusion que ce n'est pas la symétrie, mais l'asymétrie qui caractérise les langues, qu'il est possible de formuler une loi d'après laquelle les éléments linguistiques plus employés sont plus différenciés que les éléments moins utilisés (*Do the "cases vides" exist?*, Linguistique générale et linguistique indo-européenne, Kraków, 2008, p. 59–62).

4. La théorie des laryngales est infirmée par des données statistiques (*Critique de la théorie des laryngales*, Analecta Indoeuropaea Cracoviensia I. Safarewicz memoriae dicata, Cracoviae, 1995, p. 237–247; *Encore un argument contre la théorie des laryngales*, Lingua Posnaniensis 46, 2004, p. 41–44).

5. A mon avis, la règle de Verner exige une révision (*La restriction de la règle de Verner à la position médiane et le sort du s final en germanique*, Historische Sprachforschung 103, 1990, p. 92–101; *La règle de Verner s'applique-t-elle à la position finale?*, Historische Sprachforschung 109, 1996, p. 110–116).

6. A la lumière de données statistiques, le vieux slave est un compromis entre le dialecte macédo-bulgare et le parler moravo-pannonien (*Przedhistoryczne migracje Słowian i pochodzenie języka staro-cerkiewno-słowiańskiego*, Kraków, 2004; *Pochodzenie języka staro-cerkiewno-słowiańskiego a Kodeks zografski*, Warszawa, 2006).

7. A la lumière de données statistiques, l'habitat primitif des Indo-Européens est identique avec celui des Slaves (*De la préhistoire des peuples indo-européens*, Kraków, 1992; *L'habitat primitif des Indo-Européens se trouvait-il vraiment en Arménie?*; Folia Orientalia 33, 1997, p. 65–74).

8. L'orientaliste allemand Ludolf (17e s.) a été le premier à affirmer que « die Sprachverwandtschaft offenbart sich nicht im Wörterbuch, sondern in der Grammatik ». Mais on peut justifier la division des langues indo-européennes en germaniques, slaves, baltes, romanes, etc. uniquement par des convergences lexicales, et non flexionnelles ou phonétiques (*La classification des langues romanes*, Kraków, 1991, p. 22–36).

9. Le problème numéro un de l'étymologie romane est celui des verbes ayant pour sens « aller » : fr. *aller*, it. *andare*, esp. *andar*, prov. *ana*, etc. Depuis le 16e siècle, on a, au total, proposé une soixantaine d'étymologies, ce qui est un record, et cela non seulement pour l'étymologie romane. Parmi les chercheurs, il y a des adhérents à la monogenèse (affirmant que toutes ces formes proviennent, par exemple, de *ambulare*) et des adeptes de la polygenèse (prétendant, par exemple, que *aller < *advehulare, andar < *am(bi)vehitare, ana < *amvehinare*). Le calcul des probabilités permet de trancher cette question en faveur de la monogenèse (*Une étymologie romane controversée: aller, andar, etc.*, Revue roumaine de linguistique 19, 1974, p. 89–101 ; *Étymologie de fr. aller, esp. andar, etc. et calcul des probabilités*, Revue roumaine de linguistique 20, 1975, p. 735–739).

10. Depuis 1435, on affirme que les langues romanes proviennent du latin vulgaire, mais, à la lumière de données statistiques, elles sont issues du latin classique (*Le problème de l'origine des langues romanes dans le livre de H. Lüdtke et celui de R. Kiesler*, Actes du XXVe Congrès International de Linguistique et de Philologie Romanes, t. VI, Berlin, 2010, p. 207–211).

11. Depuis Jordanès, c'est-à-dire depuis 1400 ans, on estime que l'habitat primitif des Goths se trouvait en Scandinavie. Mais la comparaison de textes parallèles en gotique, allemand supérieur, moyen allemand, bas allemand, danois et suédois a révélé que l'habitat primitif des Goths se trouvait dans la partie la plus méridionale de la Germanie ancienne (*Le mythe de l'origine scandinave des Goths*, L'art de la philologie. Mélanges en l'honneur de L. Löfstedt, Helsinki, 2007, p. 137–145).

12. La division des mots en toniques et atones (articles, pronoms, prépositions, etc.), qui remonte à l'Antiquité, est le résultat d'une fausse généralisation. Il est vrai qu'il y a des homonymies *le vent = levant, à voir = avoir, et moi = émoi* et que les syllabes *le-, a-, é-* dans *levant, avoir, émoi* sont atones, mais il est erroné d'en conclure que le, à, et sont atones parce que les mots « toniques » sont traités de la même manière. *Dix vers, vingt cœurs, va tôt*, prononcés sans pauses, sont homonymes de *divers, vainqueur, Watteau*, où les syllabes *di-, vain-, Wa-* sont atones. On affirme que *Long vient* = mot tonique + mot tonique, alors que *l'on vient* = proclitique + mot tonique, mais une expérience bien simple prouve que ces expressions sont homonymes (*La division des mots en toniques et atones est-elle justifiée?*, Lingua Posnaniensis 32–33, 1991, p. 181–185).

13. Depuis l'Antiquité, on discute la question de savoir en quoi consiste la différence entre noms propres et noms communs. On a jusqu'ici proposé une dizaine de définitions du nom propre, dont aucune ne s'applique à tous les noms propres. A mon avis, la différence entre noms propres et noms communs consiste en ce que les noms communs sont, dans la grande majorité des cas, traduits d'une langue à l'autre, tandis que les noms propres ne le sont presque jamais. Par exemple, un nom commun comme *ville* est traduit en italien par *città*, en anglais par *town*, etc., alors qu'un nom propre comme *Paris* ne l'est pas, cf. it. *Parigi*, angl. *Paris*, etc. Parmi toutes les définitions du nom propre, la mienne souffre le moins d'exceptions (*La notion de nom propre*, Proceedings of 13th International Congress of Onomastic Sciences, Kraków, 1982, p. 101–106).

# References

Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '21, page 610–623, New York, NY, USA. Association for Computing Machinery.

Emily M. Bender and Alexander Koller. 2020. Climbing towards NLU: On meaning, form, and understanding in the age of data. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5185–5198, Online. Association for Computational Linguistics.

Joan L. Bybee. 2006. From usage to grammar: The mind's response to repetition. *Language*, 82(4):711–733.

Noam Chomsky, Ian Roberts, and Jeffrey Watumull. 2023. The false promise of ChatGPT. *The New York Times*.

M. H. Christiansen and N. Chater. 2008. Language as shaped by the brain. *Behavioral and Brain Sciences*, 31(5):489–558.

Przemysław Dębowiak. 2014. O dorobku naukowym profesora witolda mańczaka z okazji jubileuszu 90. urodzin. *Język Polski*, XCIV(3):194–199.

Przemysław Dębowiak. 2016. Witold Mańczak (12 VIII 1924–12 I 2016). *Onomastica*, (No 60):5–10.

Nick C. Ellis. 2002. Frequency effects in language processing: A review with implications for theories of implicit and explicit language acquisition. *Studies in Second Language Acquisition*, 24(2):143–188.

Ana Fló, Lucas Benjamin, Marie Palu, and Ghislaine Dehaene-Lambertz. 2025. Statistical learning beyond words in human neonates. *eLife*, 13:RP101802.

Danny Fox and Roni Katzir. 2024. Large language models and theoretical linguistics. *Theoretical Linguistics*, 50(1-2):71–76.

Adele E. Goldberg. 2024. Usage-based constructionist approaches and large language models. *Constructions and Frames*, 16(2):220–254.

ChengYue Gong, Di He, Xu Tan, Tao Qin, Liwei Wang, and Tie-Yan Liu. 2018. FRAGE: Frequency-agnostic word representation. *CoRR*, abs/1809.06858.

Maciej Grochowski. 2017. O poglądach profesora Witolda Mańczaka na paradygmaty językoznawstwa. *LingVaria*, 12(spec):19–28.

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, and 3 others. 2022. Training compute-optimal large language models. *Preprint*, arXiv:2203.15556.

Paul Ibbotson and Michael Tomasello. 2016. Evidence rebuts Chomsky's theory of language learning. *Scientific American*.

E.T. Jaynes. 2003. *Probability theory: The logic of science*. Cambridge University Press.

Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *Preprint*, arXiv:2001.08361.

Matthew Mandelkern and Tal Linzen. 2024. Do language models' words refer? *Computational Linguistics*, 50(3):1191–1200.

Witold Mańczak. 1958. Tendances générales des changements analogiques. *Lingua*, 7:298–325, 387–420.

Witold Mańczak. 1961. O racjonalny dobór haseł w słownikach. *Poradnik Językowy*, pages 471–476.

Witold Mańczak. 1969a. Critique du structuralisme. *Folia Linguistica*, 3(3-4):169–177.

Witold Mańczak. 1969b. Le développement phonétique des langues romanes et la fréquence. Zeszyty naukowe Uniwersytetu Jagiellońskiego 205, Kraków. Nakładem Uniwersytetu Jagiellónskiego.

Witold Mańczak. 1969c. Nieregularny rozwój fonetyczny spowodowany częstością użycia w prasłowiańskim. *Slavia*, 38:52–62.

Witold Mańczak. 1978. Les lois du développement analogique. *Linguistics*, 205:53–60.

Witold Mańczak. 1980. Critères de vérité dans la linguistique. *General Linguistics*, 20:140–145.

Witold Mańczak. 1981. Kryteria prawdy w językoznawstwie. *Biuletyn Polskiego Towarzystwa Językoznawczego*, 38:135–142.

Witold Mańczak. 1982. Linguistique et autres sciences. *Biuletyn Polskiego Towarzystwa Językoznawczego*, 39:147–152.

Witold Mańczak. 1987. *Frequenzbedingter unregelmäßiger Lautwandel in den germanischen Sprachen*. Ossolineum, Wrocław.

Witold Mańczak. 1988a. Critères de vérité. leurs conséquences pour la linguistique. *Langages*, 89:51–64.

Witold Mańczak. 1988b. O nieregularnym rozwoju fonetycznym spowodowanym frekwencją. *Biuletyn Polskiego Towarzystwa Językoznawczego*, 41:105–111.

Witold Mańczak. 1996a. Gramatyka opisowa. In *Problemy językoznawstwa ogólnego*, pages 139–146. Wrocław.

Witold Mańczak. 1996b. Gramatyka transformacyjno-generatywna. In *Problemy językoznawstwa ogólnego*, pages 183–190. Wrocław.

Witold Mańczak. 1996c. Największy problem językoznawstwa: kryteria prawdy. In *Problemy językoznawstwa ogólnego*, pages 13–19. Wrocław.

Witold Mańczak. 1996d. Nieregularny rozwój fonetyczny spowodowany frekwencją. In *Problemy językoznawstwa ogólnego*, pages 52–76. Wrocław.

Witold Mańczak. 1996e. Prawa rozwoju analogicznego. In *Problemy językoznawstwa ogólnego*, pages 81–97. Wrocław.

Witold Mańczak. 1996f. Racjonalny dobór haseł w słownikach. In *Problemy językoznawstwa ogólnego*, pages 147–149. Zakład Narodowy im. Ossolińskich - Wydawnictwo, Wrocław.

Witold Mańczak. 1996g. Rozwój semantyczny a frekwencja. In *Problemy językoznawstwa ogólnego*, pages 121–127. Wrocław.

Witold Mańczak. 1996h. Słownik a gramatyka. In *Problemy językoznawstwa ogólnego*, pages 128–132. Zakład Narodowy im. Ossolińskich - Wydawnictwo, Wrocław.

Steven T. Piantadosi. 2024. Modern language models refute chomsky's approach to language. In Edward Gibson and Moshe Poliak, editors, *From fieldwork to linguistic theory: A tribute to Dan Everett (Empirically Oriented Theoretical Morphology and Syntax 15)*, pages 353—414. Berlin: Language Science Press.

Steven T. Piantadosi and Felix Hill. 2022. Meaning without reference in large language models. *Preprint*, arXiv:2208.02957.

G. K. Pullum and B. C. Scholz. 2002. Empirical assessment of stimulus poverty arguments. *The Linguistic Review*, 19(1-2):9–50.

Alexa R Romberg and Jenny R Saffran. 2010. Statistical learning and language acquisition. *Wiley Interdisciplinary Reviews: Cognitive Science*, 1(6):906–914.

Jenny R. Saffran, Elissa L. Newport, and Richard N. Aslin. 1996. Word segmentation: The role of distributional cues. *Journal of Memory and Language*, 35(4):606–621.

Catherine E. Snow. 1972. Mothers' speech to children learning language. *Child Development*, 43(2):549–565.

Michael Tomasello. 2003. *Constructing a language: A usage-based theory of language Acquisition*. Harvard University Press.

Kaitlyn Zhou, Kawin Ethayarajh, and Dan Jurafsky. 2021. Frequency-based distortions in contextualized word embeddings. *Preprint*, arXiv:2104.08465.