

SCHOOL OF COMPUTATION,
INFORMATION AND TECHNOLOGY —
INFORMATICS

TECHNISCHE UNIVERSITÄT MÜNCHEN

Master's Thesis in Robotics Cognition Intelligence

**Optimizing a minimal language using
Pre-trained Language Models**

Ajay Narayanan

SCHOOL OF COMPUTATION,
INFORMATION AND TECHNOLOGY —
INFORMATICS

TECHNISCHE UNIVERSITÄT MÜNCHEN

Master's Thesis in Robotics Cognition Intelligence

**Optimizing a minimal language using
Pre-trained Language Models**

**Optimierung einer Minimalsprache mithilfe
vorab trainierter Sprachmodelle**

Author:	Ajay Narayanan
Examiner:	Vincent Fortuin
Supervisor:	Advisor
Submission Date:	15-05-2025

I confirm that this master's thesis is my own work and I have documented all sources and material used.

Munich, 15-05-2025

Ajay Narayanan

Acknowledgments

I would like to thank Vincent Fortuin, my supervisor for his guidance and support throughout this thesis. I would also like to thank Sara Visconti, who worked with me to develop the language pipeline. I would also like to thank my family and friends for their continued support and encouragement during my work.

Abstract

Contents

Acknowledgments	iii
Abstract	iv
1. Introduction	1
1.1. Motivation	1
2. Background	3
2.1. Linguistics	3
2.1.1. Phonetics and Phonology	3
2.1.2. Morphology	4
2.1.3. Grammar	4
2.1.4. Grammatical Features	5
2.1.5. Syntax	7
2.2. Constructed Languages	8
2.2.1. The Process of Language Creation	8
2.3. Evaluation	8
2.3.1. Evaluation of Machine Translations	8
2.4. Automatic Evaluation Metrics	9
2.4.1. BLEU: Bilingual Evaluation Understudy	9
2.4.2. ROUGE: Recall-Oriented Understudy for Gisting Evaluation . .	9
2.4.3. METEOR: Metric for Evaluation of Translation with Explicit OR- dering	9
2.4.4. Information Theoretic Evaluation of Languages	10
2.5. Language Models and Embeddings	10
2.5.1. Clustering of Embeddings	10
3. Setting up the Pipeline	11
3.1. Pipeline Overview	11
3.1.1. Phonetics Modules	11
3.1.2. Phonotactics Modules	12
3.1.3. Syllable Builder Module	12
3.1.4. Grammar Modules	12

Contents

3.1.5. Vocabulary Modules	12
3.1.6. Translation Modules	12
3.1.7. Evaluation Modules	12
A. Appendix	13
Abbreviations	15
List of Figures	16
List of Tables	17
Bibliography	18

1. Introduction

This opening chapter will provide an overview of the research topic, reviewing the study's motivation, the research questions, and the methodology. The chapter will also provide a brief outline of the structure of the thesis.

1.1. Motivation

The goal of human language, much like any other human activity, is to aid our survival and propagation. It helps us to collaborate, compete and influence others [Lev22]. Efficiency is another hallmark of all living beings, as a product of biological evolution [Ha10]. Thus, languages form from a trade-off between expressiveness and efficiency. They tend to follow rules like Zipf's Law of Abbreviations CITE, that words that are more frequently used tend to be shorter than sparsely used words.

This is, of course, an idealization. In the real world, speakers must talk in noisy environments, which means that languages often have redundancy built in. This redundancy helps in error correction and understanding in such noisy environments. In addition, natural languages are not an optimal code in terms of information theory. This is due to the fact that human languages are constrained by the limitations of incremental language processing, which enforces constraints which force systematicity, compositionality and concatenation as means of combination [FH22].

Constructed Languages (ConLangs) are languages constructed by an individual or group of individuals rather than having evolved naturally [Sch21]. ConLangs are used for various purposes, such as to serve as an international auxiliary language (e.g. Esperanto), to create a fictional world (Dothraki, Klingon, or Quenya), to be logically rigorous, (e.g. Lojban), or simple and easy to learn (e.g., Toki Pona). Some, like Ithkuil [24] are designed to express more profound thoughts briefly but clearly.

Naturally, the next question that arises is whether we could design a language that is more optimal or efficient than natural languages. This question itself is more challenging than it seems. What does it mean for a language to be optimal? Is it efficient in terms of information theory? Is it ease of learning? Is it expressiveness? Is it the ability to convey complex thoughts? Is it the ability to be understood in even the most noisy environments? In addition, how do we measure these properties? Therefore, the first motivation of this thesis is to explore what it means for a language to be optimal

and whether we can design a language that is more optimal than natural languages. To develop this language, one must define its phonology, orthography, morphology, syntax, and vocabulary.

The second motivation of this thesis is to explore whether Large Language Models (LLMs) can provide insights into this problem. In the process of learning, Large Language Models and Machine learning models in general, learn to encode various features. These features potentially encode information about the structure of language itself, and facts about our world. Trained on a large and varied corpus, they could encode information about nearly everything humans communicate about. Thus, it is possible that LLMs could provide insights that could help us design a more optimal language.

To obtain relevant information from these models, we must use techniques that make the working of these models more interpretable. Given the large number of parameters in these models, it is difficult to interpret the learned features. In addition, neurons in the model often activate in response to multiple contexts, making it even harder to create human-understandable explanations [Elh+22]. This is known as *polysemanticity*, and the hypothesized reason is *superposition* of features, where networks represent more features than directions in the parameter space. Thus, we must use tools like Sparse Autoencoders [Cun+23] to disentangle these dense features and use techniques like Autointerpretability [Bil+] to provide explanations at scale.

To test these hypotheses, we setup a modular pipeline for generating constructed languages. The code for this pipeline is available at <https://github.com/fortuinlab/conlang>.

2. Background

2.1. Linguistics

Linguistics, the scientific study of languages is a broad and complex field encompassing various subfields. Although a comprehensive summary of Linguistics is beyond the scope of this thesis, we will briefly discuss some of the subfields and key concepts that are relevant to our work. [TS07] provides a glossary of linguistic terms that can be useful for readers interested in a more detailed overview of the field.

2.1.1. Phonetics and Phonology

Phonetics is the study of the physical sounds of human speech, their production, transmission and reception. [TS07]. The International Phonetic Alphabet (IPA) is a standardized system of phonetic notation that represents the sounds of spoken language. The system is based on the assumption that speech can be represented partly as a sequence of discrete sounds or *segments* [99]. In addition, the IPA also includes symbols for suprasegmental features such as stress and intonation. The full IPA Chart (reproduced here in A.1) shows all the symbols and diacritics used to represent sounds in the IPA. Sounds and words can be transcribed in IPA using [] brackets. For example, the sounds for the word *this* can be transcribed as [ðɪs]. The IPA helps linguistics transcribe sounds in a language-agnostic way, allowing them to compare sounds across languages. The IPA Handbook [99] provides a comprehensive guide to the use of the IPA.

Phonology is the study of the sound systems of languages, including the patterns of sounds and the rules that govern their distribution. [TS07]. The key difference in the disciplines is driven by the concept of a *phoneme*. A phoneme is an abstract unit of sound that can distinguished by a native speaker of a language. Phonemes and Phonemic transcriptions are represented using slashes / /. The key points about phonemes are:

1. Letters do not necessarily correspond to phonemes. For example, the English word *this* has four letters but 3 phonemes (/ðɪs/).
2. Phonemes can be realized as different sounds in different contexts. For example,

the English phoneme /p/ can be realized as [p^h] in the word *pin*([p^hm]) and [p] in the word *spin*([spɪn]). i.e. in English, the sounds [p] and [p^h] are *allophones* of the phoneme /p/.

3. Two sounds are considered different phonemes if changing them can change the meaning of a word. e.g. [dɛn] *den* and [ðɛn] *then* are distinct words in English.

Phonotactics defines the rules that govern the permissible sound sequences in a language [TS07]. For example, in English, the sequence /bl/ is permissible at the beginning of a word (e.g. *bled*) but not the sequence /bn/. Languages usually modify loanwords to fit their own phonotactic constraints. For example, the English word *beer* is borrowed into Japanese as *biru*.

2.1.2. Morphology

Morphology is the study of the structure of words and the rules that govern the formation of words in a language [TS07]. Most studies of morphology focus on the concept of a *morpheme*, the smallest unit of meaning in a language. For example, the word *unhappiness* can be broken down into three morphemes: *un-*, *happy* and *-ness*. Morphemes can be free or bound. Free morphemes can stand alone as words (e.g. *happy*) while bound morphemes must be attached to other morphemes (e.g. *-ness*).

Morphology can be further divided into **inflectional** and **derivational** morphology. **Inflectional morphology** involves the modification of a word for grammatical purposes such as tense, aspect, mood, number, e.g. the English verb *walk* can be inflected to *walked*, *walks*, *walking*, etc. **Derivational morphology** involves the creation of new words from existing words. For example, the English noun *happiness* can be derived from the adjective *happy* by adding the suffix *-ness*.

Lexicon

The **lexicon** of a language is the vocabulary of a language, i.e. the total set of words available for a speaker. It is better to consider the lexicon, not as a list of word, but a set of lexical resources including morphemes, and processes to construct words from these resources [TS07].

2.1.3. Grammar

Grammar is the set of rules that govern the structure of sentences in a language. Traditional Grammar describes certain terms for basic grammatical components such as *article*, *adjective*, *noun*, etc known as *parts of speech* [Yul20].

1. **Nouns** are words that refer to people, places, things, or abstract ideas, as if they were objects. For example, *cat*, *house*, and *happiness* are all nouns.
2. **Verbs** are words that express the actions or states of nouns. For example, *run*, *is*, and *happen* are all verbs.
3. **Adjectives** are words that describe or modify nouns. For example, *happy*, *red*, and *tall* are all adjectives.
4. **Adverbs** are words that describe or modify verbs, adjectives, or other adverbs. For example, *really*, *very*, and *well* are all adverbs.
5. **Articles** are words that define a noun as specific or unspecific. For example, *the* is a definite article, while *a* and *an* are indefinite articles.
6. **Pronouns** are words that take the place of noun phrases, typically when they are already known. For example, *he*, *she*, and *they* are all pronouns.
7. **Prepositions** are words that show the relationship between a noun or pronoun and other words in a sentence. For example, *in*, *on*, and *at* are all prepositions.
8. **Conjunctions** are words that connect words, phrases, or clauses, and indicate the relationship between them. For example, *and*, *but*, and *or* are all conjunctions.

Sometimes, parts of speech exhibit multiple forms, used in different grammatical circumstances. Each of these forms indicate a certain *grammatical category* or *feature*. For example, in English, verbs can be inflected for tense, aspect, mood, person, number, etc. This is known as *agreement*.

A language may choose to explicitly mark these features, i.e. *grammaticalize* them [Ros10]. All features can be expressed in any language (perhaps by adding explicit information), but every language chooses to express only a subset of these features grammatically.

2.1.4. Grammatical Features

Grammatical features or categories provide some extra information about the sentence, and different parts of speech may exhibit different forms to indicate these features. we will briefly discuss some of the most common grammatical features.

Grammatical Gender

In many languages, nouns are often classified into different classes, and different parts of speech will often agree with the specific class. These classes are called Grammatical Gender, and it often need not have anything to do with sex or gender. Languages with gender may only have 2 gender classes, but can also have many more. The gender assignment of many objects are often arbitrary.

Grammatical Number

Grammatical Number is a grammatical category that expresses count distinctions. The most common distinction is between singular(one) and plural(many), but some languages also have dual(two), trial(three), paucal(few), and other forms. For example, in English, the noun *cat* is singular, while *cats* is plural.

Grammatical Case

The **Grammatical Case** indicates one or more functions of a noun or noun phrase in a sentence. Many different cases have been identified in the worlds languages, such as

1. **Nominative** case, which indicates the subject of a verb.
2. **Accusative** case, which indicates the direct object of a verb.
3. **Dative** case, which indicates the indirect object of a verb.
4. **Genitive** case, which indicates possession.
5. **Locative** case, which indicates location.
6. **Instrumental** case, which indicates the means by which an action is performed.

These descriptions are not exact, and precise distinctions can heavily depend on the specific language.

Tense, Aspect and Mood

Tense, aspect and modality all provide some kind of information that is temporal in nature, or tell us about the status of the action or verb. They are often grouped together as **TAM** (Tense, Aspect, Modality).

Tense is a grammatical category that indicates the time at which an action takes place. The most common tenses are past, present, and future. For example, in English, the verb *walk* can be inflected to *walked* (past), *walks* (present), and *will walk* (future).

Aspect is a grammatical category that indicates the temporal structure of the action or event described by a verb. It indicates for example, whether the action is bounded, and unitary, or continuous or habitual. For example, in English, the sentences *She danced* and *She was dancing* have different aspects. They are both in the past tense, but the first sentence indicates a *perfective*, or completed aspect, while the second sentence indicates an *continuous*, or *progressive* aspect.

Modality is a grammatical category that indicates the speaker's attitude towards the action or event described by a verb. Modern linguists usually associate it with the expression of obligation, permission, prohibition, necessity, possibility and ability [TS07]. In English, modality is primarily expressed using Auxiliary verbs, such as *can*, *may*, *must*, etc. For example, the sentence *He can dance* indicates ability, while the sentence *He must dance* indicates obligation.

Grammatical Person

Grammatical Person is a grammatical category that indicates the different relationships between the speaker, the listener, and others in the discourse. Languages typically indicate this relationship using pronouns. The most common distinctions are between first person (the speaker), second person (the listener), and third person (others). For example, in English, the pronouns *I* (first person), *you* (second person), and *he/she/they/it* (third person) indicate the grammatical person.

Some languages also have a distinction in *clusivity* for the first person plural pronoun. The inclusive form includes the listener, while the exclusive form does not. For example in Malayalam, the pronoun *nammal* includes the listener, while the pronoun *njangal* does not.

2.1.5. Syntax

Syntax is the study of the structure of sentences and the rules that govern the formation of sentences in a language [TS07]. The goal of Syntactic Analysis is to have a finite set of rules that could be used to generate potentially infinite sentences. This set of rules is known as a **Generative Grammar** [Yul20]. We move from the concepts of Nouns to Noun Phrases, and Verbs to Verb Phrases, and so on.

A Noun Phrase is a phrase that is interchangeable with a noun. Take for example the sentence:

_____ bought a new car.

The cloze in the sentence could be filled with a noun, like "John", but also by say , "The young man".

With these definitions in mind, we could define a sentence as a Noun Phrase followed by a Verb Phrase. We can also have production rules for Noun Phrases, Verb Phrases, and so on. For example, a Noun Phrase could be defined as a determiner followed by an adjective followed by a noun. With these rules, known as **Phrase Structure Rules**, we can generate a tree structure for a sentence, known as a **Parse Tree** [JM25].

2.2. Constructed Languages

Constructed Languages are languages that have not naturally evolved, but were artificially constructed. Some conlangs are created for fictional word-building, like *Quenya* and *Sindarin* from J.R.R. Tolkien's Middle-Earth, or *Dothraki* and *High Valyrian* from George R.R. Martin's A Song of Ice and Fire.

2.2.1. The Process of Language Creation

2.3. Evaluation

2.3.1. Evaluation of Machine Translations

Evaluation of machine translations is a complex task, and is essential for assessing the accuracy and fluency of the translations. Human evaluations are often expensive and time-consuming [Pap+02], and are not always feasible for large datasets. As a result, many researchers have developed automatic evaluation metrics to assess the quality of machine translations. Classic methods like BLEU [Pap+02] measures the similarity between the machine translation and a reference translation by comparing n-grams. ROUGE [linROUGEEvaluationAutomatic2004] is another popular metric that measures recall as opposed to precision. It is often used for evaluating the quality of summaries, but can also be used for machine translation evaluation. METEOR [BL05] improves upon such methods by for example, considering synonyms and stemming.

We use these machine translation evaluation metrics by comparing the detranslated text with the original text. Although the actual values for these metrics would be dependent on the model and its parameters, it would still be useful to compare the performance between different generated conlangs.

2.4. Automatic Evaluation Metrics

2.4.1. BLEU: Bilingual Evaluation Understudy

BLEU (Bilingual Evaluation Understudy) [Pap+02] is one of the most widely used automatic metrics for evaluating machine translations. It measures the similarity between a machine translations and reference translations by analyzing the n-gram overlap. The BLEU score is given by:

$$BLEU = BP \cdot \exp \left(\sum_{n=1}^N w_n \log p_n \right) \quad (2.1)$$

where p_n represents the precision of n-grams, w_n are weights assigned to different n-grams, and BP (the brevity penalty) penalizes short translations to prevent artificially high scores. BLEU focuses primarily on precision but does not consider recall or semantic meaning, making it less reliable for evaluating translations with synonyms or paraphrasing.

2.4.2. ROUGE: Recall-Oriented Understudy for Gisting Evaluation

ROUGE (Recall-Oriented Understudy for Gisting Evaluation) [linROUGEPackageAutomatic2004] is a set of metrics primarily used for text summarization but also applicable to machine translation. Unlike BLEU, which emphasizes precision, ROUGE focuses on recall—how much of the reference translation appears in the generated text. The most common variant, ROUGE-N, computes the recall of n-gram matches:

$$ROUGE - N = \frac{\sum_{s \in ref} \sum_{n-gram \in s} count_{match}(n - gram)}{\sum_{s \in ref} \sum_{n-gram \in s} count(n - gram)} \quad (2.2)$$

Another popular variant, ROUGE-L, measures the longest common subsequence (LCS) between the reference and candidate translations, capturing fluency better. ROUGE is especially useful for evaluating translations with different word orders but similar meanings.

2.4.3. METEOR: Metric for Evaluation of Translation with Explicit ORdering

METEOR (Metric for Evaluation of Translation with Explicit ORdering) was designed to address some of BLEU's limitations by incorporating recall, stemming, synonym matching, and word order penalties. METEOR aligns words between candidate and reference translations using exact matches, stemmed matches, and synonym matches. The metric is computed as:

$$METEOR = F_{mean} \cdot (1 - Penalty) \quad (2.3)$$

where F_{mean} is the harmonic mean of precision and recall, and the penalty reduces the score for word order mismatches. Due to its ability to consider semantic variations, METEOR often correlates better with human judgments than BLEU.

These three metrics, though widely used, each have their strengths and weaknesses. BLEU remains a benchmark due to its efficiency, ROUGE excels in recall-based evaluation, and METEOR provides a more comprehensive assessment by considering meaning and word order.

2.4.4. Information Theoretic Evaluation of Languages

2.5. Language Models and Embeddings

2.5.1. Clustering of Embeddings

3. Setting up the Pipeline

In this chapter, we will describe the pipeline we have set up to generate constructed languages. The codebase is built in Python, and we used Poetry for Dependency Management. The pipeline is modular, to allow us to perform ablation studies.

3.1. Pipeline Overview

In order to generate constructed languages, we setup a modular pipeline. Although the codebase is flexible enough that we can reorder modules, as long as the input to any module has all the features required for its execution. This is facilitated by the `LanguageDescription` class, which contains all the features of a language. This class is piped through the modules, and each module can add or modify features of the language. Each module also generates a results dict which is stored as a JSON file after the run. If the module requires a certain element of the description to be generated beforehand, it check for those features and can throw an error if they are not present.

A pipeline can be setup by subclassing the `Pipeline` class, which takes care of executing the modules and saving the results. Each module of the pipeline is a subclass of the `Module` class, which has an `execute` method that takes a `LanguageDescription` object and other optional arguments, and returns the modified `LanguageDescription` object and other optional results.

During the course of development, we found ourselves using more or less similar modules in most of our experiments. We have identified the following modules that are common to most setups:

3.1.1. Phonetics Modules

The goal of a Phonetics Module is to generate the phonemic inventory of the language. The module configures the `PhonemeDataInventory` class, which contains a list of `PhonemeData`. The phoneme segments for this class are based on PHOIBLE [MM19] segments, with a `GlyphID` corresponding to their database. For our purposes, we also implemented an `alphabet` attribute, to have a simpler representation for phonemes that are hard to read or write. This was useful for debugging and visualization purposes.

3.1.2. Phonotactics Modules

The goal of a Phonotactics Module is to generate the phonotactic rules of the language. The module configures the `PhonotacticData` class, which specifies the rules for syllable structure and phonotactic constraints for word beginnings and endings.

3.1.3. Syllable Builder Module

The Syllable builder modules combines the phonemes generated by the Phonetics Module and the phonotactic rules generated by the Phonotactics Module to generate all the possible syllables in the language.

3.1.4. Grammar Modules

3.1.5. Vocabulary Modules

The goal of the Vocabulary Module is to generate the vocabulary of the language. The module configures the `VocabDictionary` class, which holds the list of `VocabularyEntry` objects. Each `VocabularyEntry` object contains the word in the constructed language, its translation, and its definition in english. A source words list is also part of each entry, to facilitate easy search for translation purposes. The class can also hold embeddings for each word, which could be also used for downstream tasks.

3.1.6. Translation Modules

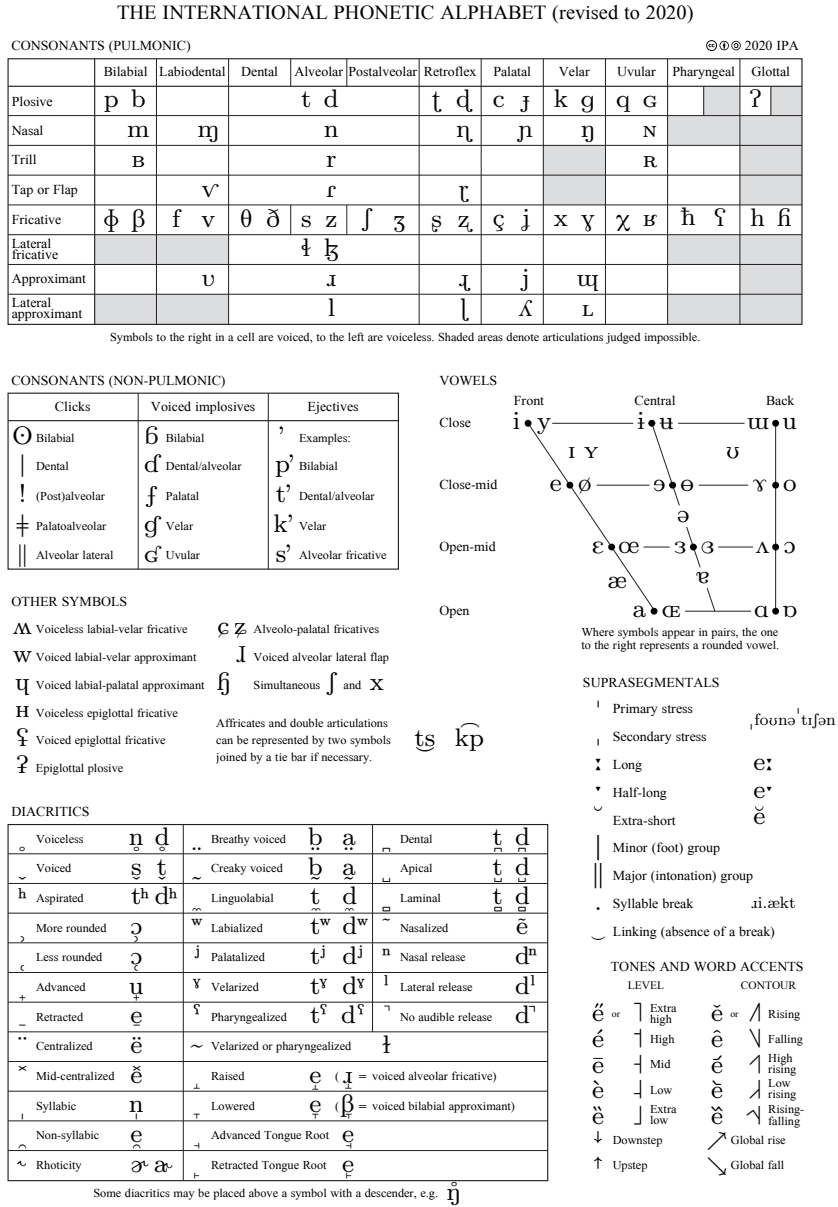
Once the language description is generated, we use a translation module to translate some text corpora represented by `AbstractSourceText`. The module translates the source text paragraph by paragraph, using an LLM model.

3.1.7. Evaluation Modules

The Evaluators are a separate class that does not inherit from the abstract `Module` class. They inherit the `Evaluator` class, which creates an evaluations folder inside the results folder, to store the results of the evaluations.

A. Appendix

A. Appendix



Abbreviations

List of Figures

- A.1. IPA Chart, available under a Creative Commons Attribution-Sharealike 3.0 Unported License. Copyright © 2018 International Phonetic Association. 14

List of Tables

Bibliography

- [24] “Ithkuil.” In: *Wikipedia* (Aug. 2024).
- [99] *Handbook of the International Phonetic Association : A Guide to the Use of the International Phonetic Alphabet*. 1. publ. Cambridge u.a.: Cambridge Univ. Press, 1999.
- [Bil+] S. Bills, N. Cammarata, D. Mossing, H. Tillman, L. Gao, G. Goh, I. Sutskevar, J. Leike, J. Wu, and W. Saunders. *Language Models Can Explain Neurons in Language Models*. <https://openaipublic.blob.core.windows.net/neuron-explainer/paper/index.html>.
- [BL05] S. Banerjee and A. Lavie. “METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments.” In: *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*. Ed. by J. Goldstein, A. Lavie, C.-Y. Lin, and C. Voss. Ann Arbor, Michigan: Association for Computational Linguistics, June 2005, pp. 65–72.
- [Cun+23] H. Cunningham, A. Ewart, L. Riggs, R. Huben, and L. Sharkey. *Sparse Autoencoders Find Highly Interpretable Features in Language Models*. Oct. 2023. doi: 10.48550/arXiv.2309.08600. arXiv: 2309.08600 [cs].
- [Elh+22] N. Elhage, T. Hume, C. Olsson, N. Schiefer, T. Henighan, S. Kravec, Z. Hatfield-Dodds, R. Lasenby, D. Drain, C. Chen, R. Grosse, S. McCandlish, J. Kaplan, D. Amodei, M. Wattenberg, and C. Olah. “Toy Models of Superposition.” In: *Transformer Circuits Thread* (2022).
- [FH22] R. Futrell and M. Hahn. “Information Theory as a Bridge between Language Function and Language Form.” In: *Frontiers in Communication* 7 (2022). ISSN: 2297-900X. doi: 10.3389/fcomm.2022.657725.
- [Ha10] R. Ha. “Cost-Benefit Analysis in Animal Behavior.” In: Jan. 2010, pp. 402–405.
- [JM25] D. Jurafsky and J. H. Martin. *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition with Language Models*. 3rd. 2025.

- [Lev22] N. Levshina. *Communicative Efficiency: Language Structure and Use*. Cambridge: Cambridge University Press, 2022.
- [MM19] S. Moran and D. McCloy, eds. *Phoible 2.0*. Jena: Max Planck Institute for the Science of Human History, 2019.
- [Pap+02] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu. “BLEU: A Method for Automatic Evaluation of Machine Translation.” In: *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*. ACL ’02. USA: Association for Computational Linguistics, July 2002, pp. 311–318. DOI: 10.3115/1073083.1073135.
- [Ros10] M. Rosenfelder. *The Language Construction Kit*. Yonagu Books, 2010. ISBN: 978-0-9844700-0-6.
- [Sch21] C. Schreyer. “Constructed Languages.” In: *Annual Review of Anthropology* 50. Volume 50, 2021 (2021), pp. 327–344. ISSN: 1545-4290. DOI: 10.1146/annurev-anthro-101819-110152.
- [TS07] R. Trask and P. Stockwell. *Language and Linguistics: The Key Concepts*. Key Concepts Series. Routledge, 2007. ISBN: 978-0-415-41359-6.
- [Yul20] G. Yule. *The Study of Language*. 7th ed. Cambridge: Cambridge University Press, 2020.