

tomo suli pi kama sona

NIMI PI TOMO NI

tsbohc

**TOKI PONA: DISTRIBUTIONAL APPROACH TO
SEMANTIC ANALYSIS OF A CONSTRUCTED LANGUAGE**

tenpo ni la, lipu ni li pona ala

(spelling and general unenglishness will be corrected at a later date)

ma pona 2022

CONTENTS

Introduction	4
Distributional semantics and constructed languages	7
1 Natural language processing	7
2 Distributional semantics	7
2.1 Distributional representations	10
2.1.1 Context types	11
2.1.2 Semantic similarity metric	12
2.1.3 Curse of dimensionality	13
2.2 Notable implementations	13
2.2.1 Count vector model	13
2.2.2 Neural probabilistic language model	14
2.2.3 Recurrent neural net language model	15
2.2.4 Continuous bag-of-words model	15
2.2.5 Continuous skip-gram model	16
2.3 Applications	16
3 Artificial languages	18
3.1 The notion of a constructed language	18
3.2 Brief history	18
3.2.1 Antiquity	18
3.2.2 Linguistic mysticism	18
3.2.3 Early artistic languages	20
3.2.4 Early modern constructed languages	20
3.2.5 J. R. R. Tolkien	20

3.2.6 Modern constructed languages	20
3.3 Motivation	20
3.4 Classification	20
3.4.1 Traditional: structure and source material	20
3.4.2 Traditional: purpose	22
3.4.3 Blanke's functional classification	24
Language modelling and Toki Pona	25
1 Toki Pona	25
1.1 History	25
1.2 Phonology	25
1.3 Grammar	25
1.4 Vocabulary	25
1.5 Tokiponidos	25
2 Vector space model	25
2.1 Text tokenisation	25
2.2 Text normalisation	25
2.3 Model construction	25
2.4 Projection and visualisation	25
2.5 Observations	25
REFERENCES	26
Supplementary marterial	31
1 Vector space model	31
2 Dictionaries	31
2.1 nimi pu	31

INTRODUCTION

Toki Pona is the second most spoken constructed language in the world. Its core vocabulary consists of only 120-140 words, not including words that are rare and/or considered non-standard by the majority of speakers. Despite the small vocabulary size, Toki Pona can be used to convey a wide range of ideas of varying complexity [[Meulen 2021](#)].

Problem

While some cases are covered by the lessons present in the original book [[Lang 2014](#)], the definitions provided by the official publicly available dictionary do not fully reflect how the vocabulary is used today.

This could be amended by the way of inquiry of a large enough sample of proficient speakers, but this approach would prove highly subjective. Moreover, there are weak dialectical variations and idiolects present in the language.

A more neutral and less manual approach is to use distributional semantics to construct a model of the language. The data provided by the model can serve as the basis for identifying the undocumented word senses and ranking the existing ones.

Goals

The primary goal of this paper is to comment on the semantics of the core Toki Pona vocabulary, as well as to organise individual words into groups based on their semantic relatedness and the context they are most prevalently used in.

The secondary goal is to discuss the semantics of the core vocabulary of Toki Pona — as it is spoken by the majority of the community — in relation to

the official dictionary of the language.

- **Subject.** Semantic evaluation and classification of vocabulary.
- **Object.** Toki Pona, a constructed language.
- **Methodology.** Distributional semantics and natural language processing, namely language modelling (word embedding).

Objectives

1. Define natural language processing and distributional semantics, discuss modern implementations as well as other related concepts.
2. Define and classify constructed languages.
3. Describe Toki Pona, its philosophy, history, and unique features.
4. Obtain the necessary corpora.
5. Construct a vector space model of the language.
6. Make observations on the model.
7. Classify the words of the vocabulary based on the observed semantic relationships between them.

Relevance

Constructed languages are rapidly gaining popularity. Despite this, the only constructed language that has seen much representation in scientific writing is Esperanto.

The existing dictionaries or other resources concerned with teaching Toki Pona to new speakers could benefit from the findings of this research. New tools can be developed which will aid newcomers to the language.

The vector space model of Toki Pona developed in the course of this research can find further use in machine translation, topic modelling, text prediction, sentiment analysis, and many other areas.

DISTRIBUTIONAL SEMANTICS AND CONSTRUCTED LANGUAGES

1 Natural language processing

“Linguistics is concerned not only with language per se, but must also deal with how humans model the world. The study of semantics, for example, must relate language expressions to their meanings, which reside in the mental models possessed by humans. <...> Whereas computational linguistics, as a subfield of linguistics, is concerned with the formal or computational description of rules that languages follow [[Tsuji 2021](#)]”.

The aim of this research is to bridge the gap between the two disciplines, to use computational linguistics to build a semantic model of a constructed language. This model can then be used to explore the nuances of how humans speak the said language.

In turn, “Natural Language Processing is a field at the intersection of computer science, artificial intelligence, and linguistics [[Vajjala, Majumder 2020, p. 7](#)]”. “Natural language processing includes a range of algorithms, tasks, and problems that take human-produced text as an input and produce some useful information, such as labels, semantic representations, and so on, as an output [[Hagiwara 2021, p. 4](#)]”.

2 Distributional semantics

The core idea behind distributional semantics has roots in American structuralism (Harris) and British lexicology (Firth) and is known as the distributional hypothesis. In its simplest form, it states that “similarity in

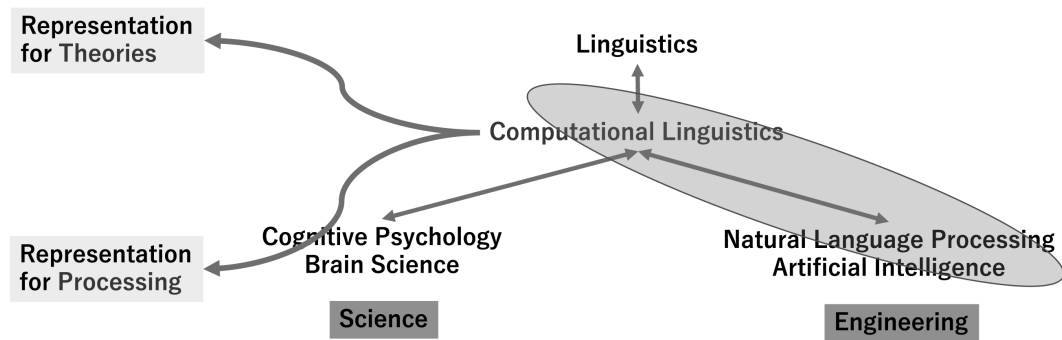


Figure .1: Language-related disciplines [Tsuji 2021]

meaning results in similarity of linguistic distribution [Harris 1954]”.

The reverse of this statement is also true. Meaning that “the statistical distribution of linguistic items in context plays a key role in characterizing their semantic behavior [Lenci 2018]”. The aim of distributional semantics is exactly that, to learn the meanings of linguistics units from a corpus of text.

Distributional semantics was popularised by Firth in the 1950s. In a 1957 publication he wrote, “the placing of a text as a constituent in a context of situation contributes to the statement of meaning since situations are set up to recognise use. <...> You shall know a word by the company it keeps! [Firth 1957, p. 11]”.

The ideas introduced by the distributional hypothesis have received attention in cognitive science [Mcdonald 2008] and language learning [Yarlett, Ramscar, Dye 2008].

Overview

Distributional semantics has become widespread with the adoption of information technology in the field of linguistic research.

Distributional semantics are most frequently applied by taking large amounts of text as input and pushing it through an abstraction algorithm to produce a distributional model as output [Emerson 2020].

Distributional models rely on context to produce semantic

representations. That is, distributional models characterise the meanings of words through the context in which they have been observed [Erk 2016].

Planets of the solar system
are orbiting the *sun*. The
moon is orbiting the earth.
It's his antique *typewriter*
clacking. <...>

→
algorithm

	dim1	dim2
sun	0.11023	0.53848
moon	0.21575	0.44034
typewriter	0.52834	0.05389

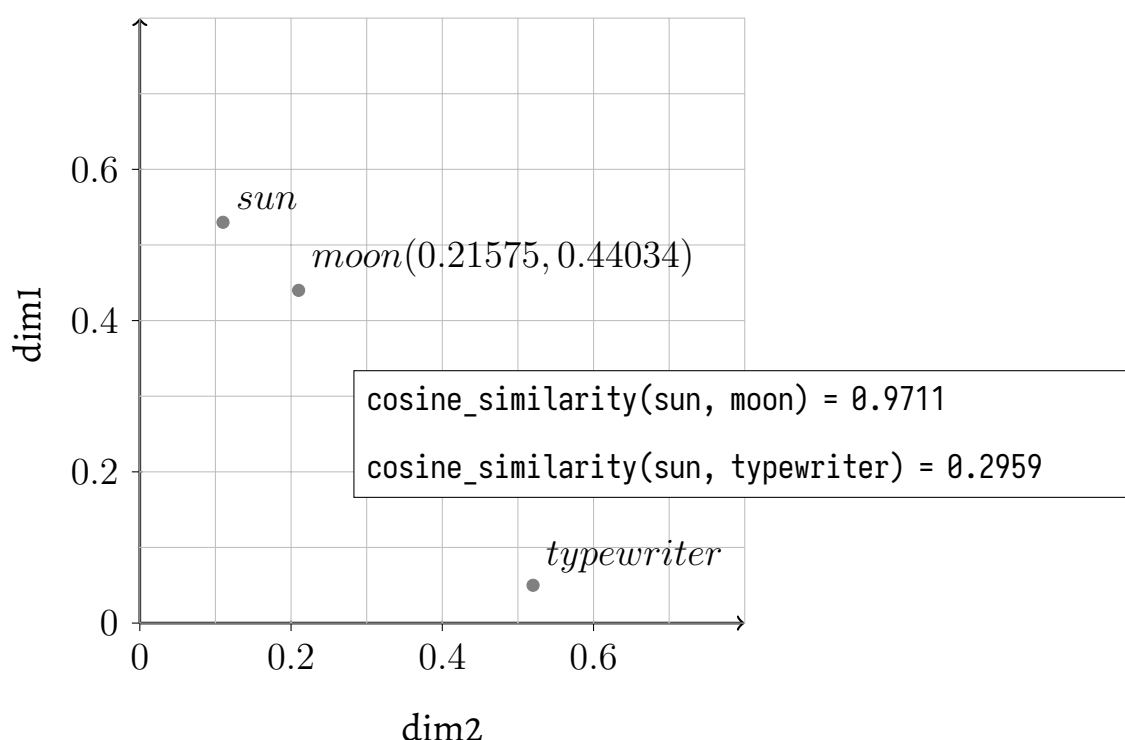


Figure .2: Distributional semantics, an illustrated overview

In a model, the semantic representations are stored in the form of vectors. Vectors are essentially lists of numbers that refer to points in a multi-dimensional space. These vectors are referred to as word vectors.

In the illustrated example, the model only has the dimensionality of two and thus can be mapped onto a two dimensional plane without any further processing.

If this is not the case, the multi-dimensionality of the word vector

encodings can be reduced to only two or three dimensions. The resulting dimensions can then be used to create a projection of the model which can be observed by the human eye.

All of the approaches to distributional semantics share the quality of learning semantic representations from a corpus in an unsupervised manner. Meaning that it is not required for the corpus to be preprocessed by hand.

2.1 Distributional representations

Distributional representations are mathematic encodings of the distributional properties of words. Typically, in the form of a sequence of numbers. This sequence of numbers can be viewed as a multi-dimensional vector for the purposes of applying to them principles derived from linear algebra.

“Word vectors represent words as multidimensional continuous floating point numbers where semantically similar words are mapped to proximate points in geometric space [[Ahire 2018](#)]”.

In simpler terms, a word vector is a numerical representation of a word in a corpus relative to every other word in that corpus.

“Vectors have geometrical interpretations: Vectors with n components define points (or arrows) in n -dimensional spaces. Therefore, distributional representations are geometrical representations of the lexicon in the form of a distributional vector space. The positions of lexemes in a distributional semantic space depend on their co-occurrences with linguistic contexts [[Lenci 2018](#)]”.

2.1.1 Context types

Distributional representations output by a distributional model differ with respect to how the linguistic context is defined.

The contexts can be of the following types [[Lenci 2018](#)]:

- **Undirected window-based collocate.** This context type includes words around the current word. No information as to whether the context words precede or follow after the current word is provided to the model. The window size typically ranges from 2 to 10.
- **Directed window-based collocate.** Unlike the previous context type, directed window-based contexts provide the direction in which the context word was seen relative to the current word.
- **Dependency-filtered syntactic collocate.** This context restricted the words which are analysed by the algorithm based on their syntactic roles. This information is however not provided to the model.
- **Dependency-typed syntactic collocate.** This context type provides the previously omitted syntactic type to the model.
- **Text region.** A text region context can represent any text sample that is uniquely identifiable: book chapters, web pages, or simply text portions of any fixed size.

The term window provides a physical analogy to a linguistic context. As the algorithm processes the corpus, the window of the context slides across the text, accounting for the words that can be seen through it.

2.1.2 Semantic similarity metric

The semantic similarity between two vectors is primarily measured in two ways: using cosine similarity or the Euclidean distance.

The primary advantage of using one of these two methods is that they can be calculated for vectors of any dimensionality.

Euclidean distance

The Euclidean distance between two points is the length of a line segment between the two points. It can also be defined as the shortest distance between two points in an n -dimensional space. For the purposes of calculating the Euclidean distance, the vectors are viewed as point coordinates [Oduntan, Adeyanju 2018].

$$d_{Euc}(p, q) = \sqrt{\sum_{i=1}^n (p_i - q_i)^2} = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + \dots + (p_n - q_n)^2}$$

Cosine similarity

Cosine similarity is a measurement of similarity between two sequences of numbers. When calculating cosine similarity, the two sequences of numbers are viewed as vectors. Cosine similarity is equal to the cosine of the angle between two vectors, that is, the dot product of the vectors divided by the product of their lengths [Oduntan, Adeyanju 2018].

Cosine similarity always falls into the interval $[-1, 1]$. Two parallel vectors have a cosine similarity of 1, two orthogonal (perpendicular to each other) vectors have a cosine similarity of 0, while two opposite vectors have a cosine similarity of -1 .

$$s_{cos}(A, B) := \cos(\theta) = \frac{A \cdot B}{\|A\| \cdot \|B\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$

This method was chosen to simplify the process of comparing similarities between vector pairs. Where the Euclidean distance provides an absolute value, the cosine similarity provides a fraction.

2.1.3 Curse of dimensionality

The curse of dimensionality refers to the phenomena that arise when organising data in high-dimensional spaces.

In the context of distributional models, dimensionality is determined by how many word relationships are accounted for by the model.

As the dimensionality of representations increases, the volume of the space they take up increases so fast that the available data becomes sparse. In other words, it becomes hard to make sense of the data as it becomes spread too thinly across the multi-dimensional space [Venkat 2018].

A common solution to this is dimensionality reduction.

2.2 Notable implementations

2.2.1 Count vector model

The simplest implementations of distributional models feature counting algorithms. “<...> these models just record other words that have been observed in the vicinity of a target word in large text corpora, and form some sort of aggregate over the recorded context items. They then estimate the semantic

similarity between words based on contextual similarity [Erk 2016]”. These models are referred to as count models.

“Context items are counted only if they appear close to the target word, that is, if they are within the relevant context [Erk 2016]”.

The count models operate on window-based context. The window size is typically narrow (2-4 words). The window can be allowed to cross the boundaries of sentences or not [Baroni, Dinu, Kruszewski 2014].

2.2.2 Neural probabilistic language model

In the recent years, the distributional model architecture has seen as notable shift to machine learning algorithms. With the improvements of hardware performance, the training of complex neural networks on corpora of larger sizes has become possible.

The earlier machine learning based models were plagued by the curse of dimensionality. This problem was solved in the model proposed by [Bengio, Ducharme, Vincent 2000]. The proposed neural network model learns distributional representations and the generalisation function at the same time. The generalisation function is based on the estimates of probability of a word appearing in the given context.

The architecture of this model “consists of input, projection, hidden and output layers. At the input layer, N previous words are encoded using 1-of- V coding, where V is size of the vocabulary. The input layer is then projected to a projection layer P that has dimensionality $N \times D$, using a shared projection matrix. As only N inputs are active at any given time, composition of the projection layer is a relatively cheap operation [Mikolov, Chen, Corrado 2013]”.

The training complexity of this model is

$$Q = N \times D + N \times D \times H + H \times V$$

where H is the size of the hidden layer.

2.2.3 Recurrent neural net language model

This model architecture contains recurrent neural networks, meaning that as the model learns from the input, it produces output that is fed back into the model as input. The recurrent matrix connects hidden layers to itself using time-delayed connections. “This allows the recurrent model to form some kind of short term memory, as information from the past can be represented by the hidden layer state that gets updated based on the current input and the state of the hidden layer in the previous time step [Mikolov, Chen, Corrado 2013]”.

This model architecture consists of only input, hidden, and output layers, thus allowing for a reduction of complexity when compared to the neural probabilistic language model [Mikolov, Chen, Corrado 2013].

The training complexity of this model is

$$Q = H \times H + H \times V$$

2.2.4 Continuous bag-of-words model

The first architecture of Word2vec proposed by Mikolov removes the non-linear hidden layer, further reducing complexity. The projection layer is shared for all words [Mikolov, Chen, Corrado 2013].

The continuous bag-of-words model is not influenced by history like the previous one. In a continuous bag-of-words model not only the words preceding the current word are used for context, but also the words that follow it.

This model attempts to predict the current word from the sum of the context vectors. This sum of vectors is referred to as a “bag of words”, giving the name to the model. If the prediction of the word is correct after comparing

it with the current word, its distributional representation is reinforced. If the prediction is wrong, the distributional representation is corrected.

The training complexity of this model is

$$Q = N \times D + D \times \log_2 V$$

Because this model architecture produces the prediction as output, the learned weights of the hidden layer is what represents the word vectors.

2.2.5 Continuous skip-gram model

The second architecture of Word2vec proposed by Mikolov has the opposite objective of the continuous bag-of-words model. The continuous skip-gram model predicts the surrounding context from the current word. Similar to the continuous bag-of-words model, when the continuous skip-gram model succeeds in predicting the context words, the semantic representation of the current word is reinforced. When it fails, it is corrected [Mikolov, Chen, Corrado 2013].

The training complexity of this model is

$$Q = N \times D + N \times D \times \log_2 V$$

While the complexity of this model is greater, the accuracy is also much greater [Mikolov, Chen, Corrado 2013]. Similar to a continuous bag-of-words model, the weights of the hidden layer are the distributional representations.

2.3 Applications

The data provided by the distributional model can be used directly to analyse the semantics of a language:

1. **Semantic similarity.** By definition, distributional models provide data that quantifies semantic relatedness between individual words or expressions.

This data can be interpreted by humans to draw conclusions about the meanings of words or used in other areas of natural language processing.

2. **Word clustering.** Semantic representations tend to form groups in the multi-dimensional space. Word clustering refers to the ways and means by which these groups can be extracted as formal clusters [Bekkerman, El-Yaniv, Tishby 2003].
3. **Automatic creation of thesauri.** The semantic similarity data can be further processed to produce lists of homonyms, synonyms, or even antonyms [Henestroza Anguiano, Denis 2011].
4. **Word sense disambiguation.** This refers to a problem in computational linguistics that is concerned with identifying which sense of a word is used in a particular sentence [Musto, Narducci, Basile 2011].
5. **Information retrieval.** Distributional models can be used to access semantically similar words to those of a query, expanding the retrieved results from exact word matching to semantically fuzzy matching [Silva, Maia 2019].
6. **Data mining.** In data mining, namely text mining, distributional models can provide means of identifying similar documents, thus narrowing the scope of a search [Dalianis 2018, p. 89].
7. **Paraphrasing.** The data provided by distributional models can supply paraphrasing algorithms with vocabulary, or aid in judging the relative semantic similarity between two paraphrases on a sentence level basis [El Desouki, Gomaa 2019].
8. **Sentiment analysis.** Given a small list of words manually tagged with emotive potentials, distributional models can propagate these potentials

through a corpus based on the semantic similarity between the tagged words [[Alshari, Azman, Doraisamy 2017](#)].

3 Artificial languages

3.1 The notion of a constructed language

3.2 Brief history

“The dream of a perfect language did not only obsess European culture. The story of the confusion of tongues, and of the attempt to redeem its loss through the rediscovery or invention of a language common to all humanity, can be found in every culture [[Eco 1995, p. 1](#)]”.

3.2.1 Antiquity

The concept of a constructed language was first mentioned in writing by Athenaeus of Naucratis in his *Deipnosophistae* (circa AD 230). The languages he presented were not full languages, but languages of rudimentary type known as a naming language. These languages were collections of neologisms that could be used to replace existing vocabulary or to refer to things that otherwise had no name. Athenaeus further writes of other people who invented their own words [[Sanders 2020](#)].

3.2.2 Linguistic mysticism

Irish myths of the seventh century described the origin of Gaelic. According to *Auraicept na n-Éces*, King Fénus Farsaid of Scythia traveled to the Tower of Babel after God fragmented human language. King Fénus and his many scholars studied the remains of the human language and combined its

best fragments into a new and more perfect language, Gaelic [Williams 2016].

The earliest attempt at language creation was *Lingua Ignota* by Hildegard of Bingen, a German Benedictine abbess and polymath, in the eleventh century. The underlying structure of the language was Latin, but the spelling was significantly altered. “She did compose one macaronic antiphon, “O orzchis Ecclesia,” in which Latin and *Lingua Ignota* vocabulary alternate <...> Alternatively, if the *Lingua* had a second use as a secret language (possibly in the presence of outsiders) for Hildegard and her nuns, as some have suggested, this reviewer submits that verbs are not always needed for the achievement of communication: “Enpholianz warinz nascutil” (bishop / wart / nose) provides, if not exactly a sentence, an entirely understandable lexical string [Straubhaar 2008]”.

“Hildegard also created *Litteraë Ignotæ* ‘unknown letters’, a constructed writing system or *neography*, which she used to represent *Lingua Ignota* [Sanders 2020]”. Despite the lack in grammatical depth by comparison to modern constructed languages, *Lingua Ignota* is widely praised as the first constructed language ever created.

From the belief in the connection between language creation and the divine arose the notion that languages inherited by humans from gods had become corrupted and now were causing confusion among the scholars and philosophers, impeding scientific progress. At the same time, by the 1600s, “Latin, which was the sole language of education and the only language taught, experienced a decline, and there was no other language to replace it anytime soon [Stria 2016, p. 51]”. This caused a search for a new language to replace Latin, and the subsequent creation of many philosophical constructed languages. These languages were centered around the idea of providing a top-down categorical view of the universe.

In his essay, Wilkins proposed one such universal philosophical language to replace Latin. This language was to be unambiguous and to encompass every

concept in the universe, to overcome the curse of Babel [[Wilkins 1668, ch. 2, p. 1](#)].

Wilkins constructed a table of 40 major genera, which he then divided further into 251 characteristic differences. From them he derived 2030 species, which appear in pairs. For example, “starting from the major genus of Beasts, after having divided them into viviparous and oviparous, and after having subdivided the viviparous ones into whole footed, cloven footed and clawed, Wilkins arrives at the species Dog/Wolf [[Eco 1995, p. 239](#)]”.

Despite the initial acclaim and the interest it received from the king, the language soon fell into obscurity [[Okrent 2010, p. 25](#)]. This style of language creation continued into the seventeenth century and then was abandoned.

3.2.3 Early artistic languages

3.2.4 Early modern constructed languages

3.2.5 J. R. R. Tolkien

3.2.6 Modern constructed languages

3.3 Motivation

3.4 Classification

3.4.1 Traditional: structure and source material

The classification proposed by Couturat and Leau groups constructed languages based on their relationship with source material or a lack thereof [[Couturat, Leau 2014](#)]:

- **A priori.** Constructed languages that are not based on the elements of natural languages.

- **A posteriori.** Constructed languages that borrow elements from natural languages.
- **Mixed.** A combination of the two.

Several linguists have adopted this classification, most notably Janton, who provided several additions [[Janton, Tonkin, Edwards 1993, p. 5](#)]:

1. **A priori.** Metalanguages, schematic languages. These constructed languages are “characterized by largely artificial, nonethnic word roots, schematic derivation, and fixed word categories (i.e., philosophical languages) [[Janton, Tonkin, Edwards 1993, p. 6](#)]”.
2. **A posteriori.** Naturalistic languages, pseudolanguages. “These languages consciously imitate, in varying degrees, natural languages [[Janton, Tonkin, Edwards 1993, p. 5](#)]”.
 - **Minimal languages.** Simplified natural languages, living or dead.
 - **Mixed languages.** Languages that use natural and non-natural roots.
 - Languages that were schematically derived with natural word roots in distorted form (Volapük, 1880) or with both artistic and natural word roots (Perio, 1904).
 - Languages that are partly schematic and partly naturalistic. Natural word roots in this group are rarely distorted (Esperanto, 1887).
 - **Naturalistic languages.**
 - Languages with some schematic traits (Unial, 1903; Novial, 1928)
 - Languages with natural derivation (Interlingua, 1940s).

This classification presents a scale of artificiality, with languages derived from natural languages on one end and deliberately designed languages on the

other. It is also primarily concerned with morphological and syntactic natures of derivation, not the intention with which the language was created.

3.4.2 Traditional: purpose

“Another type of classification categorises artificial languages according to the purpose of creation [Stria 2016, p. 93]”. Kennaway provides the following division: universal languages (from a strive for perfection), international languages, languages of fiction, languages as recreation [Kennaway 2013].

The primary concern with this way of classifying constructed languages this way lies in the fact that purpose is rarely binary. The prominent example are the languages which were constructed as recreation and later served as a tool in fictional worldbuilding.

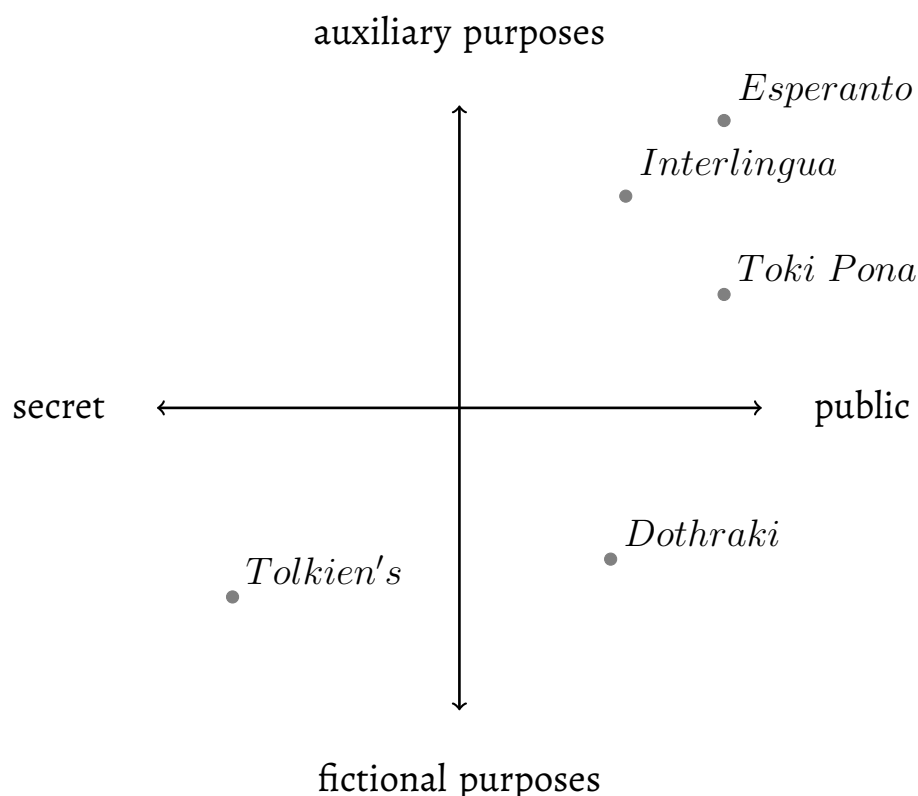


Figure .3: Federico Gobbo's coordinate system

What followed were classifications that took a far more granular

approach. “A detailed typology is proposed by Albani and Buonarotti (1994), where a division is made into sacred and non-sacred languages. Sacred languages are further divided into structured (Bālaibalan) and non-structured with six subdivisions. Nonsacred languages split into languages with communicative and expressive goals both with further detailed subdivisions [Stria 2016, p. 93]”.

The classifications that followed saw a decrease in complexity and a shift to a more visual approach. Gobbo proposed a coordinate system which placed constructed languages between the “secret” and “public” extremes on the *x axis* and “auxiliary purposes” and “fictional purposes” on the *y axis* [Gobbo 2012]. This however, is also not ideal, as the publicity of languages tends to shift greatly over time.

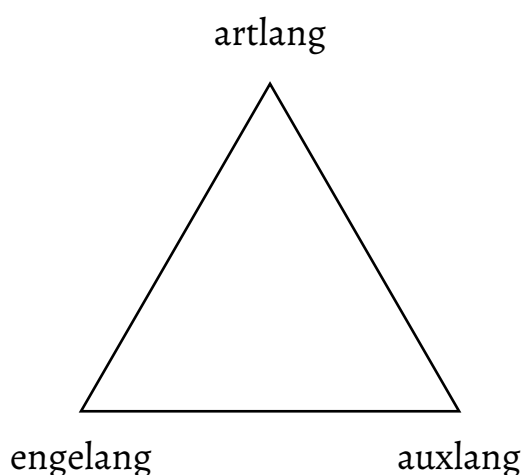


Figure .4: The modified Gnoli triangle

“One of the newest propositions widely spread on the Internet is the so-called Gnoli triangle. Claudio Gnoli, dissatisfied with the fact that his constructed language Liva was not easily classified, came up with an idea of a triangle whose vertices were labelled ‘artlang’ (artistic language), ‘auxlang’ (auxiliary language) and ‘loglang’ (logical language; the term ‘engelang’ was proposed later by And Rosta, apparently in 2001) [Stria 2016, p. 97]”.

- **Auxlang.** An international auxiliary language, that is, a language which

was devised with the intention of being a means of international communication. The majority of international auxiliary languages are meant to be second languages and not to replace native languages [Libert 2018].

- **Artlang.** A language created as a form of artistic expression or to fill an artistic role. Artistic languages often have irregular grammar systems, much like natural languages.
- **Engelang.** A language whose grammar or other feature is based on logic (a loglang) or a language which was created to as an experiment or to prove a hypothesis of how languages function.

3.4.3 Blanke's functional classification

LANGUAGE MODELLING AND TOKI PONA

1 Toki Pona

1.1 History

1.2 Phonology

1.3 Grammar

1.4 Vocabulary

1.5 Tokiponidos

2 Vector space model

2.1 Text tokenisation

2.2 Text normalisation

2.3 Model construction

2.4 Projection and visualisation

2.5 Observations

REFERENCES

1. *Ahire J. B.* Introduction to Word Vectors // Retrieved March. — 2018. — Vol. 12. — P. 2018.
2. *Alshari E. M., Azman A., Doraisamy S.* Improvement of sentiment analysis based on clustering of Word2Vec features // 2017 28th international workshop on database and expert systems applications (DEXA). — IEEE. 2017. — P. 123–126.
3. *Baroni M., Dinu G., Kruszewski G.* Don't count, predict! A systematic comparison of context-counting vs. context-predicting semantic vectors //. Vol. 1. — 06/2014. — P. 238–247. — DOI: [10.3115/v1/P14-1023](https://doi.org/10.3115/v1/P14-1023).
4. *Bekkerman R., El-Yaniv R., Tishby N.* Distributional Word Clusters vs. Words for Text Categorization // Journal of Machine Learning Research. — 2003. — Mar. — Vol. 3. — P. 1183–1208. — DOI: [10.1162/153244303322753625](https://doi.org/10.1162/153244303322753625).
5. *Bengio Y., Ducharme R., Vincent P.* A Neural Probabilistic Language Model //. Vol. 3. — 01/2000. — P. 932–938. — DOI: [10.1162/153244303322533223](https://doi.org/10.1162/153244303322533223).
6. *Couturat L., Leau L.* Histoire de La Langue Universelle (1903). — Literary Licensing, LLC, 2014. — ISBN 9781498147248. — URL: https://books.google.ru/books?id=s2%5C_poQEACAAJ.
7. *Dalianis H.* Clinical text mining: Secondary use of electronic patient records. — Springer Nature, 2018.
8. *Eco U.* The search for the perfect language. — Search, 1995.

9. *El Desouki M. I., Gomaa W. H.* Exploring the recent trends of paraphrase detection // International Journal of Computer Applications. — 2019. — Vol. 975, S 8887.
10. *Emerson G.* What are the Goals of Distributional Semantics? — 2020. — DOI: [10.48550/ARXIV.2005.02982](https://doi.org/10.48550/ARXIV.2005.02982). — URL: <https://arxiv.org/abs/2005.02982>.
11. *Erk K.* What do you know about an alligator when you know the company it keeps // Semantics and Pragmatics. — 2016. — Vol. 9.
12. *Firth J.* A Synopsis of Linguistic Theory, 1930-1955. — 1957. — URL: <https://books.google.ru/books?id=T8LDtgAACAAJ>.
13. *Gobbo F.* Alan Turing creator of Artificial Languages // InKoj. — 2012. — Sept. — Vol. 3. — DOI: [10.13130/2037-4550/2385](https://doi.org/10.13130/2037-4550/2385).
14. *Hagiwara M.* Real-World Natural Language Processing: Practical Applications with Deep Learning. — Manning, 2021. — ISBN 9781617296420. — URL: <https://books.google.ru/books?id=0k5NEAAQBAJ>.
15. *Harris Z. S.* Distributional Structure // WORD. — 1954. — Vol. 10, no. 2/3. — P. 146–162. — DOI: [10.1080/00437956.1954.11659520](https://doi.org/10.1080/00437956.1954.11659520). — eprint: <https://doi.org/10.1080/00437956.1954.11659520>. — URL: <https://doi.org/10.1080/00437956.1954.11659520>.
16. *Henestroza Anguiano E., Denis P.* FreDist: Automatic construction of distributional thesauri for French. — 2011. — June.
17. *Janton P., Tonkin H., Edwards J.* Esperanto: Language, Literature, and Community. — State University of New York Press, 1993. — ISBN 9781438407807. — URL: https://books.google.sk/books?id=R%5C_ZGfG2a1tIC.
18. *Kennaway R.* Artificial Languages //. — 05/2013. — ISBN 978-0-415-42432-5.
19. *Lang S.* Toki Pona: The Language of Good. — Sonja Lang, 2014. — ISBN 9780978292300. — URL: <https://books.google.ru/books?id=5P0ZjwEACAAJ>.

20. *Lenci A.* Distributional Models of Word Meaning // Annual Review of Linguistics. — 2018. — Feb. — Vol. 4. — DOI: [10.1146/annurev-linguistics-030514-125254](https://doi.org/10.1146/annurev-linguistics-030514-125254).
21. *Libert A. R.* Artificial Languages. — 06/2018. — DOI: [10.1093/acrefore/9780199384655.013.11](https://doi.org/10.1093/acrefore/9780199384655.013.11). — URL: <https://oxfordre.com/linguistics/view/10.1093/acrefore/9780199384655.001.0001/acrefore-9780199384655-e-11>.
22. *Mcdonald S.* Testing the distributional hypothesis: The influence of context on judgements of semantic similarity. — 2008. — July.
23. *Meulen S. v. d.* Request for New Language Code Element in ISO 639-3 // ISO 639-3 Registration Authority. — 2021. — Aug. — Vol. 2021.
24. *Mikolov T., Chen K., Corrado G.* Efficient Estimation of Word Representations in Vector Space // Proceedings of Workshop at ICLR. — 2013. — Jan. — Vol. 2013.
25. *Musto C., Narducci F., Basile P.* Cross-Language Information Filtering: Word Sense Disambiguation vs. Distributional Models //. Vol. 6934. — 09/2011. — P. 250–261. — ISBN 978-3-642-23953-3. — DOI: [10.1007/978-3-642-23954-0_24](https://doi.org/10.1007/978-3-642-23954-0_24).
26. *Oduntan O., Adeyanju O.* A Comparative Analysis of Euclidean Distance and Cosine Similarity Measure for Automated Essay-Type Grading // Journal of Engineering and Applied Sciences. — 2018. — July. — Vol. 13. — P. 4198–4204. — DOI: [10.3923/jeasci.2018.4198.4204](https://doi.org/10.3923/jeasci.2018.4198.4204).
27. *Okrent A.* In the Land of Invented Languages: Adventures in Linguistic Creativity, Madness, and Genius. — Random House Publishing Group, 2010. — (Spiegel & Grau trade paperbacks). — ISBN 9780812980899. — URL: <https://books.google.ru/books?id=3anWeY0G2moC>.

28. *Sanders N.* A primer on constructed languages //. — 08/2020. — P. 6–26. — ISBN 9780198829874. — DOI: [10.1093/oso/9780198829874.003.0002](https://doi.org/10.1093/oso/9780198829874.003.0002).
29. *Silva F. T. da, Maia J. E.* Query Expansion in Text Information Retrieval with Local Context and Distributional Model. // J. Digit. Inf. Manag. — 2019. — Vol. 17, no. 6. — P. 313.
30. *Straubhaar S. B.* Sarah Higley, Hildegard of Bingen's Unknown Language: An Edition, Translation and Discussion. (The New Middle Ages series.) Palgrave Macmillan, 2007 // Medieval Feminist Forum: A Journal of Gender and Sexuality. Vol. 44. — Society for Medieval Feminist Scholarship. 2008. — P. 158–161.
31. *Stria I.* Inventing languages, inventing worlds. Towards a linguistic worldview for artificial languages. — 01/2016. — ISBN 978-83-947609-1-5. — DOI: [10.14746/9788394760915](https://doi.org/10.14746/9788394760915).
32. *Tsujii J.* Natural Language Processing and Computational Linguistics // Computational Linguistics. — 2021. — Dec. — Vol. 47, no. 4. — P. 707–727. — ISSN 0891-2017. — DOI: [10.1162/coli_a_00420](https://doi.org/10.1162/coli_a_00420). — eprint: https://direct.mit.edu/coli/article-pdf/47/4/707/1979478/coli_a_00420.pdf. — URL: https://doi.org/10.1162/coli%5C_a%5C_00420.
33. *Vajjala S., Majumder B.* Practical Natural Language Processing: A Comprehensive Guide to Building Real-World NLP Systems. — O'Reilly Media, 2020. — ISBN 9781492054009. — URL: <https://books.google.ru/books?id=hPrrDwAAQBAJ>.
34. *Venkat N.* The Curse of Dimensionality: Inside Out. — 09/2018. — DOI: [10.13140/RG.2.2.29631.36006](https://doi.org/10.13140/RG.2.2.29631.36006).

35. *Wilkins J.* An Essay Towards a Real Character, and a Philosophical Language. — Sa: Gellibrand, for John Martyn printer to the Royal Society, 1668. — URL: <https://books.google.ru/books?id=BCctZjBtiEYC>.
36. *Williams M.* Ireland's Immortals: A History of the Gods of Irish Myth. — Princeton University Press, 2016. — ISBN 9781400883325. — URL: <https://books.google.ru/books?id=o5cTDAAQBAJ>.
37. *Yarlett D., Ramscar M., Dye M.* Language Learning Through Similarity-Based Generalization. — 2008. — Jan.

SUPPLEMENTARY MATERIAL

1 Vector space model

The two-dimensional projection of the semantic model constructed as a result of this research.

2 Dictionaries

2.1 nimi pu

The dictionary of Toki Pona as it appears in Toki Pona: The Language of Good [Lang 2014, p. 125–134]. This dictionary is licensed under public domain.

Word		Definition
a or kin	<i>particle</i>	(emphasis, emotion or confirmation)
akesi	<i>noun</i>	non-cute animal; reptile, amphibian
ala	<i>adjective</i>	no, not, zero
alasa	<i>verb</i>	to hunt, forage
ale or ali	<i>adjective</i>	all; abundant, countless, bountiful, every, plentiful
	<i>noun</i>	abundance, everything, life, universe
	<i>number</i>	100
anpa	<i>adjective</i>	bowing down, downward, humble, lowly, dependent
ante	<i>adjective</i>	different, altered, changed, other
anu	<i>particle</i>	or
awen	<i>adjective</i>	enduring, kept, protected, safe, waiting, staying

Table .1: nimi pu

	Word	Definition
	<i>pre-verb</i>	to continue to
e	<i>particle</i>	(before the direct object)
en	<i>particle</i>	(between multiple subjects)
esun	<i>noun</i>	market, shop, fair, bazaar, business transaction
ijo	<i>noun</i>	thing, phenomenon, object, matter
ike	<i>adjective</i>	bad, negative; non-essential, irrelevant
ilo	<i>noun</i>	tool, implement, machine, device
insa	<i>noun</i>	centre, content, inside, between; internal organ, stomach
jaki	<i>adjective</i>	disgusting, obscene, sickly, toxic, unclean, unsanitary
jan	<i>noun</i>	human being, person, somebody
jelo	<i>adjective</i>	yellow, yellowish
jo	<i>verb</i>	to have, carry, contain, hold
kala	<i>noun</i>	fish, marine animal, sea creature
kalama	<i>verb</i>	to produce a sound; recite, utter aloud
kama	<i>adjective</i>	arriving, coming, future, summoned
	<i>pre-verb</i>	to become, manage to, succeed in
kasi	<i>noun</i>	plant, vegetation; herb, leaf
ken	<i>pre-verb</i>	to be able to, be allowed to, can, may
	<i>adjective</i>	possible
kepeken	<i>preposition</i>	to use, with, by means of
kili	<i>noun</i>	fruit, vegetable, mushroom
kiwen	<i>noun</i>	hard object, metal, rock, stone

Table .1: nimi pu

	Word	Definition
ko	<i>noun</i>	clay, clinging form, dough, semi-solid, paste, powder
kon	<i>noun</i>	air, breath; essence, spirit; hidden reality, unseen agent
kule	<i>adjective</i>	colourful, pigmented, painted
kulupu	<i>noun</i>	community, company, group, nation, society, tribe
kute	<i>noun</i>	ear
	<i>verb</i>	to hear, listen; pay attention to, obey
la	<i>particle</i>	(between the context phrase and the main sentence)
lape	<i>adjective</i>	sleeping, resting
laso	<i>adjective</i>	blue, green
lawa	<i>noun</i>	head, mind
	<i>verb</i>	to control, direct, guide, lead, own, plan, regulate, rule
len	<i>noun</i>	cloth, clothing, fabric, textile; cover, layer of privacy
lete	<i>adjective</i>	cold, cool; uncooked, raw
li	<i>particle</i>	(between any subject except mi alone or sina alone and its verb; also to introduce a new verb for the same subject)
lili	<i>adjective</i>	little, small, short; few; a bit; young
linja	<i>noun</i>	long and flexible thing; cord, hair, rope, thread, yarn

Table .1: nimi pu

	Word	Definition
lipu	<i>noun</i>	flat object; book, document, card, paper, record, website
loje	<i>adjective</i>	red, reddish
lon	<i>preposition</i>	located at, present at, real, true, existing
luka	<i>noun</i>	arm, hand, tactile organ
	<i>number</i>	five
lukin or oko	<i>noun</i>	eye
	<i>verb</i>	to look at, see, examine, observe, read, watch
	<i>pre-verb</i>	to seek, look for, try to
lupa	<i>noun</i>	door, hole, orifice, window
ma	<i>noun</i>	earth, land; outdoors, world; country, territory; soil
mama	<i>noun</i>	parent, ancestor; creator, originator; caretaker, sustainer
mani	<i>noun</i>	money, cash, savings, wealth; large domesticated animal
meli	<i>noun</i>	woman, female, feminine person; wife
mi	<i>noun</i>	I, me, we, us
mije	<i>noun</i>	man, male, masculine person; husband
moku	<i>verb</i>	to eat, drink, consume, swallow, ingest
moli	<i>adjective</i>	dead, dying
monsi	<i>noun</i>	back, behind, rear
mu	<i>particle</i>	(animal noise or communication)
mun	<i>noun</i>	moon, night sky object, star
musi	<i>adjective</i>	artistic, entertaining, frivolous, playful, recreational

Table .1: nimi pu

	Word	Definition
mute	<i>adjective</i>	many, a lot, more, much, several, very
	<i>noun</i>	quantity
nanpa	<i>particle</i>	-th (ordinal number)
	<i>noun</i>	numbers
nasa	<i>adjective</i>	unusual, strange; foolish, crazy; drunk, intoxicated
nasin	<i>noun</i>	way, custom, doctrine, method, path, road
nenā	<i>noun</i>	bump, button, hill, mountain, nose, protuberance
ni	<i>adjective</i>	that, this
nimi	<i>noun</i>	name, word
noka	<i>noun</i>	foot, leg, organ of locomotion; bottom, lower part
o	<i>particle</i>	hey! O! (vocative or imperative)
olin	<i>verb</i>	to love, have compassion for, respect, show affection to
ona	<i>noun</i>	he, she, it, they
open	<i>verb</i>	to begin, start; open; turn on
pakala	<i>adjective</i>	botched, broken, damaged, harmed, messed up
pali	<i>verb</i>	to do, take action on, work on; build, make, prepare
palisa	<i>noun</i>	long hard thing; branch, rod, stick
pan	<i>noun</i>	cereal, grain; barley, corn, oat, rice, wheat; bread, pasta
pana	<i>verb</i>	to give, send, emit, provide, put, release

Table .1: nimi pu

	Word	Definition
pi	<i>particle</i>	of
pilin	<i>noun</i>	heart (physical or emotional)
	<i>adjective</i>	feeling (an emotion, a direct experience)
pimeja	<i>adjective</i>	black, dark, unlit
pini	<i>adjective</i>	ago, completed, ended, finished, past
pipi	<i>noun</i>	bug, insect, ant, spider
poka	<i>noun</i>	hip, side; next to, nearby, vicinity
poki	<i>noun</i>	container, bag, bowl, box, cup, cupboard, drawer, vessel
pona	<i>adjective</i>	good, positive, useful; friendly, peaceful; simple
pu	<i>adjective</i>	interacting with the official Toki Pona book
sama	<i>adjective</i>	same, similar; each other; sibling, peer, fellow
	<i>preposition</i>	as, like
seli	<i>adjective</i>	fire; cooking element, chemical reaction, heat source
selo	<i>noun</i>	outer form, outer layer; bark, peel, shell, skin; boundary
seme	<i>particle</i>	what? which?
sewi	<i>noun</i>	area above, highest part, something elevated
	<i>adjective</i>	awe-inspiring, divine, sacred, supernatural
sijelo	<i>noun</i>	body (of person or animal), physical state, torso
sike	<i>noun</i>	round or circular thing; ball, circle, cycle, sphere, wheel
	<i>adjective</i>	of one year

Table .1: nimi pu

Word		Definition
sin or namako	<i>adjective</i>	new, fresh; additional, another, extra
sina	<i>noun</i>	you
sinpin	<i>noun</i>	face, foremost, front, wall
sitelen	<i>noun</i>	image, picture, representation, symbol, mark, writing
sona	<i>verb</i>	to know, be skilled in, be wise about, have information on
	<i>pre-verb</i>	to know how to
soweli	<i>noun</i>	animal, beast, land mammal
suli	<i>adjective</i>	big, heavy, large, long, tall; important; adult
suno	<i>noun</i>	sun; light, brightness, glow, radiance, shine; light source
supa	<i>noun</i>	horizontal surface, thing to put or rest something on
suwi	<i>adjective</i>	sweet, fragrant; cute, innocent, adorable
tan	<i>preposition</i>	by, from, because of
taso	<i>particle</i>	but, however
	<i>adjective</i>	only
tawa	<i>preposition</i>	going to, toward; for; from the perspective of
	<i>adjective</i>	moving
telo	<i>noun</i>	water, liquid, fluid, wet substance; beverage
tenpo	<i>noun</i>	time, duration, moment, occasion, period, situation
toki	<i>verb</i>	to communicate, say, speak, say, talk, use language, think
tomo	<i>noun</i>	indoor space; building, home, house, room

Table .1: nimi pu

	Word	Definition
tu	<i>number</i>	two
unpa	<i>verb</i>	to have sexual or marital relations with
uta	<i>noun</i>	mouth, lips, oral cavity, jaw
utala	<i>verb</i>	to battle, challenge, compete against, struggle against
walo	<i>adjective</i>	white, whitish; light-coloured, pale
wan	<i>adjective</i>	unique, united
	<i>number</i>	one
waso	<i>noun</i>	bird, flying creature, winged animal
wawa	<i>adjective</i>	strong, powerful; confident, sure; energetic, intense
weka	<i>adjective</i>	absent, away, ignored
wile	<i>pre-verb</i>	must, need, require, should, want, wish