



HAL
open science

Constructing Languages to Explore Theoretical Principles

Guillaume Enguehard, Xiaoling Luo, Nicola Lampitelli

► **To cite this version:**

Guillaume Enguehard, Xiaoling Luo, Nicola Lampitelli. Constructing Languages to Explore Theoretical Principles. *RiCognizioni*, 2022, 9 (18), <https://www.ojs.unito.it/index.php/ricognizioni/article/view/7098>. 10.13135/2384-8987/7098 . halshs-03937995

HAL Id: halshs-03937995

<https://shs.hal.science/halshs-03937995v1>

Submitted on 11 May 2023

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

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



Distributed under a Creative Commons Attribution 4.0 International License

CONSTRUCTING LANGUAGES TO EXPLORE THEORETICAL PRINCIPLES

Guillaume ENGUEHARD, Xiaoliang LUO, Nicola LAMPITELLI

ABSTRACT • The construction of languages has always been related to linguistics. Most of these initiatives address real scientific questions but from a non-academic point of view. The fact that Ferdinand de Saussure's own brother, René de Saussure, wrote a theoretical essay on the construction of the Esperanto word (de Saussure, 1914) is an amusing illustration of this. In this paper, we propose a method inspired by experimental archaeology. The experiment consists in trying to obtain an artefact similar to the one observed using this or that construction method. An equivalent approach in linguistics would be the generation of linguistic systems based on explicitly formulated principles. Trying to generate similar systems pushes the linguist to explicitly define the principles that are needed and to explore all their consequences. In this context, we show that the use of notions induced by the observation of natural languages leads to a certain degree of circularity and that it is therefore more interesting to explore *a priori* principles based on very general assumptions.

KEYWORDS • Natural Languages; Natlang; Linguistic Theory; Constructed Languages; Conlangs.

Introduction

This paper aims to introduce an original approach to formal linguistics, which could be called 'simulationist linguistics', and which basically relies on two objectives:

1. to define a set of *a priori* principles concerning language.
2. to test the plausibility of those principles in constructing languages.

This approach differs strongly from the classical approach which consists in observing particular phenomena to induce general principles. In contrast, our goal is to consciously abstract as much as possible from particular phenomena in order to build a model whose value is strictly determined by its ability to generate data close to real language data, not to its ability to offer a good description of a particular language, to be realistic with regards to acquisition or to mimic cognitive mechanisms. Thus, the key word here is '*a priori* principles'. No matter these principles are of a structuralist, functionalist or generativist nature. The only thing that guides their selection is, again, their ability to produce naturalistic facts out of a minimum of assumptions. The present paper does not even try to define these principles, but to show the interest of the method.

In our first section, we show how constructions and discoveries are interrelated. We specifically focus on a parallel between experimental archeology and language construction. Second, we address important contrasts between modeling and simulation. We show that the former necessarily suffers from an inherent degree of circularity that only the second approach can overcome. Finally, we will discuss some examples of *a priori* principles for deriving phonological data in our third section.

1. Construction and discovery

1.1. Constructed languages and linguistic research

Constructing languages may be seen as a way to experiment or to exploit linguistic notions. The most famous constructed language, Esperanto, is in itself an application of theoretical principles in order to facilitate mutual understanding: Zamenhof (1887) based his work on phonemic principle, paradigm uniformity, agglutination.

In the same way, Ogden's Basic English - a controlled language proposed as an alternative to Esperanto - is founded on an interesting and profound consideration for semantic primes (Ogden, 1930). Though some aspects of Ogden's work are criticized for their lack of objectivity or motivation, they have proven to be effective in the field of English language teaching.

Some constructed languages have fewer practical purposes. For instance, Loglan was created by Brown (1960) in order to test the Sapir-Worth hypothesis and Toki Pona was designed by Lang (2014) to take full advantage of the possibilities of polysemy in a minimal language.

1.2. Experimental archaeology

Experimental approaches based on artificial constructions have been explored in other fields, namely when:

1. there is no way to physically observe the cause of a given phenomenon.
2. the cause can hardly be deduced from the observation of the result.

In this case, causes are investigated through speculative methods. In archaeology, for instance, the causes of an artifact are regularly outside our empirical field and it is not possible to make abstract generalizations as we do in linguistics. Researchers therefore sometimes use speculation and validation of hypotheses through direct experience (Bordes, 1947). The experiment consists in seeing if one obtains an artifact similar to the one observed using this or that construction method. One key contribution from experimental archeology is to understand how items were produced in prehistoric times, especially in the lithic industry. We cannot observe the process of carving and the result of this process offers only very partial indications of their production method. The only way to solve the mystery of their origin is to reproduce the process through trial and error. Though it never gives definite answers, it makes it possible to (in)validate hypotheses on the scale of a complete system.

1.3. Experimentation and simulation

An equivalent approach applied to linguistics would be the generation of linguistic systems based on explicitly formulated principles. This approach is not strictly speaking

‘experimental’ linguistics. In linguistics, experimental methods are often intended to produce speech by controlling the conditions of its occurrence. The experiment is not in itself a hypothetical cause being tested, but an exclusion of confounding variables (Gillioz and Zufferey, 2020, p. 8). Thus, it aims to produce data that are difficult to observe in a natural context, but the conclusion refers to something broader than the data themselves or the conditions of the experiment. In contrast, experimental archeology aims to test a hypothesis that directly refers to the conditions of the experiment. We can therefore speak of indirect experience in the first case and direct experience in the second, which can be illustrated as follows:

1. Experimental linguistics: if A involves B, then C - which contains A - is true (indirect experience)
2. Experimental archeology: if A involves B, then A is true (direct experience)

Since experimental linguistics is not the exact equivalent of experimental archeology, we will speak here of ‘simulationist linguistics’ to refer to direct experience. We aim at testing directly potential causes of language systems. Our goal is to manipulate directly those parameters leading to a linguistic grammatical form rather than inferring them.

The origin of these parameters does not really matter. They may be the result of data observation, cognitive, behavioral, functional or even *a priori* speculation. Their value is ultimately determined only by the degree of closeness between the results obtained and the natural languages as a whole. Trying to generate similar systems pushes the linguist to explicitly define the principles that are needed and to explore all the consequences of these principles.

In the remainder of this paper, we are going to show how to reproduce a linguistic system and to make discoveries without inference.

2. Modeling and simulation

2.1. Language generators

There are many online generators for producing artificial linguistic data. Language generators such as Vulgarlang (<https://www.vulgarlang.com/>) stem from well-established observations of real linguistic facts. They thus manipulate *a posteriori* parameters: They combine current linguistic categories such as onsets, nuclei, codas, syllables, feet and words that result from a careful typological observation.

It is also the way in which current formal theories proceed. A presumably universal model is constructed from natural data, which is then declined according to parameters that do not seem to be reducible to a universal principle (Prince and Smolensky, 1993). This is what we call modeling.

2.2. Modeling vs Simulation

2.2.1. Modeling

What we call modeling is a two-step process. First, formal linguistics builds categories induced from the observation of natural languages in a bottom-up way, then it deduces possible linguistic systems through the application of such categories in a top-down way.

In phonology, for instance, Element Theory - henceforth ET - (Kaye et al., 1985; Backley, 2011) claims that sound inventory of any language is built from a limited number of universal, primitive units called Elements, these are induced from the observation of a large range of natural languages. A simplified version of the Elements is shown in Table 1.

| Element | Phonological interpretation when associated to vocalic position | Phonological interpretation when associated to consonantal position |
|----------------|---|---|
| A | lowness | liquid |
| I | palatality | palatality |
| U | velarity, roundedness | velarity |
| ʔ | | stopness |
| h | high tone | fricative |
| L | low tone | voicing, nasality |
| v ⁰ | centralness | |

Table 1: Element Theory

Elements are then combined according to rules in order to derived balanced vowel systems. For instance, a mid-vowel [e] is the result of the following combination: AI.

We can take ET as a starting point to construct a language. A balanced vowel or consonant system can be considered as a structured application of combination rules, with a random variable handling the parameters of these combination rules. Table 2 shows two plausible vowel systems generated in a spreadsheet.

| | |
|---|-------|
| i | i y u |
| e | e ø o |
| a | a |

(a) No rounded palatal vowels (b) Rounded palatal vowels

Table 2: Example of generated vowel systems

Without going into details of the calculation of the functions in the spreadsheet, let us take 2a as an example. Table 2a allows the combination of I or U with A, resulting in three degrees of aperture. However, I and U cannot combine. In 2b, in contrast, the combination of I and U is allowed, which gives /y/ and /ø/ in the inventory.

The same principles can be applied to consonants. In Table 3 are shown two plausible consonant systems.

| | |
|--|--|
| p ^h t ^h k ^h | p ^h t ^h k ^h |
| p t k | b ^h d ^h g ^h |
| f s ʃ | p t k |
| m n ŋ | b d g |
| l | f s ʃ |
| r | v z ʒ |
| | m n ŋ |
| | l |
| | r |
| (a) No voice contrast | (b) Voice contrast |

Table 3: Example of generated consonant systems

In 3a, the marked consonant series is the aspirated one, there is no voiced consonant. In 3b, plosives can be marked by voicing, aspirated and both. The contrast between the two systems is due to the ability of L to combine with the element h found in obstruents.

As we can see in the above examples, one necessarily expects results derived from ET to be in conformity with natural languages since the former is based on generalisations from the latter. Generating constructed languages through modeling is thus efficient but tells us nothing new because of this circularity.

Moreover, *a posteriori* principles are partially biased because they are based on a limited sample of existing languages. Some languages are well described while others are not. These principles also hardly take into account areal explanations, and they cannot access prehistorical or future languages.

By the concept of ‘simulation’, we propose the reverse path: we define *a priori* principles from which we generate categories like onsets, nuclei, codas, and everything that makes linguistic data look natural.

2.2.2. Simulation

In contrast to empirical principles based on a sample of languages, *a priori* principles do not rely on the observation of languages. Their motivation can be diverse and ultimately matter very little. Only their ability to produce realistic data matters. By definition, such principles are universal and thus avoid the aforementioned biases.

Moreover, as modeling follows principles that are *a posteriori* derived from natural languages, generating a possible language on such a basis does nothing more than repeat an already known generalization. In the examples of Subsection 2.2.1, we managed to derive balanced systems with parametric variations only because we first assumed that phonological systems are balanced and subject to some specific parametric variations already implemented in our theory.

Conversely, if we can generate a possible language without assumptions drawn from natural languages, then it can tell us something new about the hidden mechanisms underlying the structure of natural languages. This is what we call ‘simulation’, as opposed to ‘modeling’.

2.3. How to define language naturalness?

It is difficult to define exactly what is natural in a language. To have the answer to this question, we need to know what is possible and what is impossible, and thus to have an exact model of how languages work.

In the meantime, we propose a minimalist definition according to which the naturalness of languages is a set of parametric restrictions affecting the following aspects of languages:

1. The inventory of units (= what)
2. The qualitative combination of units (= how)
3. The quantitative combination of units (= how much)

Modules and derivation

In this section, we present three different modules dedicated to deriving the phonological inventory, syllabic constraints, and word size respectively. These modules were implemented in a simple spreadsheet (LibreOffice Calc) with formula that everyone can reproduce.¹ In each case, we first present the *a priori* principle on which the module is based, and then we illustrate our results using a chosen simulation example as well as a randomly selected simulation example.

All the principles retained in our presentation remain debatable. They only serve to illustrate our point about the language simulation research method.

3.1. Module 1: Inventory

3.1.1. Assumption and implementation

Our first module is based on a functionalist assumption that the members of an inventory should be as far apart as possible in order to maximize the contrast between the different phonemes. This principle recalls Martinet's diachronic phonology (Martinet, 1955) and the Adaptive Dispersion Theory (Liljencrantz and Lindblom, 1972).

In practice, this postulate imposes to define a continuous space of phonetic values from which the inventory is built by selection. We define such a space in Table 4. For instance, this table does not show any dichotomy between consonants and vowels. Of course, this example is far from satisfactory. It is limited by a two-dimensional representation that does not account for the proximity between labials and gutturals, and the use of 'weird' symbols only loosely simulates a continuous space. For these two reasons, it is futile to present our category choices in more detail. But this example is sufficient, for the moment, to illustrate our point.

¹ The main formula can be found in the appendix.

| | | | | | | | | | | | |
|---|----|----|---|------|------|----|----|----|----|----|----|
| p | p̄ | t̄ | t | ĉp | ĉp̄ | t̄ | t̄ | c | k | q | ʔ |
| b | b̄ | d̄ | d | ĵb | ĵb̄ | d̄ | d̄ | ʃ | g | ɠ | |
| ϕ | f | θ | s | çϕ | xϕ | ʃ | ʃ̄ | ç | x | χ | h |
| β | v | ð | z | ĵβ | ŷβ | ʒ | z̄ | ʝ | ʎ | ʙ | fi |
| m | ɱ | n̄ | n | m̄n̄ | m̄n̄ | n̄ | n̄ | ɲ | ŋ | ɴ | |
| | | l̄ | l | | | l̄ | l̄ | ʎ | ʎ | | |
| | | r̄ | r | | | r̄ | r̄ | | | ʀ | |
| | | f̄ | f | | | f̄ | f̄ | | | | |
| | u | ø | ɹ | ɥ | w | ɹ̄ | ɹ̄ | ʝ | ɥ | ɥ̄ | ɹ̄ |
| | | | | y | u | ɹ̄ | ɹ̄ | i | ɥ | ɥ̄ | ī |
| | | | | ʏ | ʊ | ɹ̄ | ɹ̄ | ɹ̄ | ɥ̄ | ɥ̄ | ɹ̄ |
| | | | | ø | o | ē | ē | e | ʏ | ʏ̄ | ə |
| | | | | œ | ɔ | ɛ̄ | ɛ̄ | ɛ | ʌ | ʌ̄ | ɜ |
| | | | | œ | ɔ | ǣ | ǣ | æ | ɑ | ɑ̄ | ɑ |

Table 4: Module 1 - Suggestion for a phonic continuum

The definition of the phonological inventory is done by applying a categorial mesh to this space. This mesh follows the principles mentioned above, i.e. the fact that it exploits the entire phonological continuum and seeks to maximize the distance between its nodes.

In order to allow some variation in the realization of these principles, we introduce two random parameters following a normal distribution law: the first one manages the tightness of the mesh, as in Figures 1a and 1b; and the second one manages the regularity of its nodes distribution, as in Figures 1b and 1c.

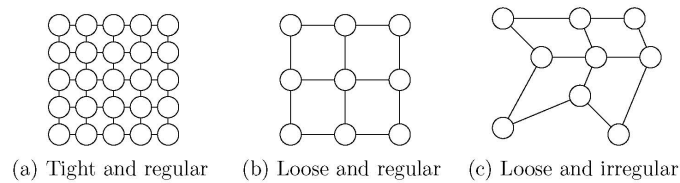


Figure 1: Mesh types

Each phonetic value falling under a node of the mesh is then selected in the phonological inventory.

3.1.3. Results

In the chosen example in Table 5, the implementation of our assumption simulates a system with 7 obstruants, 4 sonorants and 3 vowels. To the reader put off by the presence of unusual sounds, note that the selected symbols still represent specific phonetic values, not phonemes gathering a whole range of phonetic values.

A more phonological reading of this system reveals a system with three cardinal vowels /i,a,u/.

| | | | | |
|---|----|----|---|----|
| p | t̥ | | t | q |
| β | | ⱱβ | ʒ | n |
| | f | | ʈ | |
| | | | ĩ | ũ̃ |
| | | | a | |

Table 5: Module 1 - Chosen example (after 9 generations)

In the random example in Table 6, the vowel system is typologically distinct: it is more akin to a 5-vowel system like /i,e,a,o,u/. But in both cases, we observe universal patterns such as the presence of obstruents, sonorants and vowels, or the existence of several place features.

| | | | | |
|---|----|-----|---|----|
| p | t̥ | k̂p | t | q |
| | | | ʃ | x |
| v | | | | |
| | l | | ɭ | |
| | ǒ̃ | | ɹ | ɰ̃ |
| | | ʊ | ĩ | ũ̃ |
| | | | ɜ | ʌ |
| | | ɒ | | |

Table 6: Module 1 - Random example

Thus, a very simple *a priori* principle succeeds in simulating various types of phonological systems without help of induced typological conjectures.

3.2. Module 2: Qualitative combination

3.2.1. Assumption and implementation

Now that we have defined our inventory, we need to define the way all these values can combine. The use of concepts such as syllable, nucleus, onset, coda, foot, etc. is totally excluded, as they are based on an observation of facts.

Our *a priori* assumption for the second module is based on the same functionalist mechanism as above. Just as the units must be maximally contrastive on the paradigmatic axis, they must be maximally contrastive on the syntagmatic axis. This also recalls notions explored in Hjelmslev (1943) or OCP (Goldsmith, 1976; McCarthy, 1979). Thus, sequences of segments are ruled by the distance between two values in the phonetic space.

We calculate the distance between two values of a given inventory with tables like the one in 7. Each number represents the number of cells between two phonetic values.

| | p | t | k̂p | ʈ | q | θ | β | ĵβ | j | l | ʎ | ʊ | î | ũ | œ | ɛ | ʌ | æ |
|-----|----|----|-----|----|----|----|----|-----|----|----|----|----|----|----|----|----|----|----|
| p | 0 | 2 | 5 | 7 | 10 | 4 | 3 | 7 | 11 | 7 | 12 | 15 | 17 | 20 | 16 | 20 | 22 | 19 |
| t | 2 | 0 | 3 | 5 | 8 | 2 | 5 | 5 | 9 | 5 | 10 | 13 | 15 | 18 | 14 | 18 | 20 | 17 |
| k̂p | 5 | 3 | 0 | 2 | 5 | 5 | 8 | 4 | 6 | 8 | 7 | 10 | 12 | 15 | 13 | 15 | 17 | 14 |
| ʈ | 7 | 5 | 2 | 0 | 3 | 7 | 10 | 6 | 4 | 10 | 5 | 12 | 10 | 13 | 15 | 13 | 15 | 14 |
| q | 10 | 8 | 5 | 3 | 0 | 10 | 13 | 9 | 5 | 13 | 8 | 15 | 13 | 10 | 18 | 14 | 12 | 17 |
| θ | 4 | 2 | 5 | 7 | 10 | 0 | 3 | 3 | 7 | 3 | 8 | 11 | 13 | 16 | 12 | 16 | 18 | 15 |
| β | 3 | 5 | 8 | 10 | 13 | 3 | 0 | 4 | 8 | 4 | 9 | 12 | 14 | 17 | 13 | 17 | 19 | 16 |
| ĵβ | 7 | 5 | 4 | 6 | 9 | 3 | 4 | 0 | 4 | 4 | 5 | 8 | 10 | 13 | 9 | 13 | 15 | 12 |
| j | 11 | 9 | 6 | 4 | 5 | 7 | 8 | 4 | 0 | 8 | 3 | 10 | 8 | 9 | 13 | 9 | 11 | 12 |
| l | 7 | 5 | 8 | 10 | 13 | 3 | 4 | 4 | 8 | 0 | 5 | 8 | 10 | 13 | 9 | 13 | 15 | 12 |
| ʎ | 12 | 10 | 7 | 5 | 8 | 8 | 9 | 5 | 3 | 5 | 0 | 7 | 5 | 8 | 10 | 8 | 10 | 9 |
| ʊ | 15 | 13 | 10 | 12 | 15 | 11 | 12 | 8 | 10 | 8 | 7 | 0 | 2 | 5 | 3 | 5 | 7 | 4 |
| î | 17 | 15 | 12 | 10 | 13 | 13 | 14 | 10 | 8 | 10 | 5 | 2 | 0 | 3 | 5 | 3 | 5 | 4 |
| ũ | 20 | 18 | 15 | 13 | 10 | 16 | 17 | 13 | 9 | 13 | 8 | 5 | 3 | 0 | 8 | 4 | 2 | 7 |
| œ | 16 | 14 | 13 | 15 | 18 | 12 | 13 | 9 | 13 | 9 | 10 | 3 | 5 | 8 | 0 | 4 | 6 | 3 |
| ɛ | 20 | 18 | 15 | 13 | 14 | 16 | 17 | 13 | 9 | 13 | 8 | 5 | 3 | 4 | 4 | 0 | 2 | 3 |
| ʌ | 22 | 20 | 17 | 15 | 12 | 18 | 19 | 15 | 11 | 15 | 10 | 7 | 5 | 2 | 6 | 2 | 0 | 5 |
| æ | 19 | 17 | 14 | 14 | 17 | 15 | 16 | 12 | 12 | 12 | 9 | 4 | 4 | 7 | 3 | 3 | 5 | 0 |

Table 7: Module 2 - Proximity calculation

Then we apply a random parameter which, for each segment of a word, selects the next one among the most distant values following a normal distribution law.

3.2.2. Results

In the chosen example in Table 8, our implementation simulates pseudo-words with mostly CV, CVC, V, VC syllables. All the words contain a vowel. Only some of them, in bold, contain a consonant cluster, which is always [lp].

| | | | | | | | | | |
|------|------|---------|------------|------|--------|-------|-------------|------------|---------|
| œqæp | ũpΔp | ĩpΔpΔ | k̂pΔp | pΔpΔ | k̂pΔpΔ | βΔp | βΔp | lpΔ | ΔpΔ |
| θΔp | æpΔp | œqæp | pΔpΔ | ʎΔ | ɛpΔ | ɛpΔpΔ | ũp | k̂pΔp | θΔpΔ |
| ʎΔpΔ | æpΔp | ĵβΔpΔp | jæpΔ | ɛpΔ | qœqæ | pΔp | lpΔp | qœqæ | ʊpΔpΔ |
| ʊqæp | βΔpΔ | qœqæ | k̂pΔp | ʎp | ʊpΔ | qæp | œqæpΔ | βΔpΔp | θΔ |
| jæpΔ | ɛpΔ | Δp | jæp | jœq | k̂pΔp | ΔpΔp | ĵβΔ | βΔp | k̂pΔpΔp |
| jæ | œqæ | Δ | lpΔ | ʎæpΔ | qœ | ũpΔ | æp | jœq | ĵβΔp |
| æpΔp | ʎΔp | θΔpΔ | θΔp | œq | ɛpΔp | œqæ | ʊqœq | ũpΔ | ʎΔpΔ |
| θΔpΔ | ʊpΔp | βΔ | pΔp | ʎΔpΔ | ɛpΔ | ΔpΔ | ĩp | jæ | k̂pΔp |
| ʎΔp | ʎΔ | ũpΔ | æpΔp | æpΔp | ΔpΔ | ũp | βΔp | ʊp | œqœq |
| œq | βΔpΔ | βΔpΔ | ĩp | ʊp | æpΔ | æpΔp | k̂pΔp | k̂pΔp | ĩp |

Table 8: Module 2 - Chosen example (after 8 generations)

In the random example in Table 9, the simulated words are more complex. Many of them, in bold, contains consonant clusters, and some have no vowel. This absence of vowel can be intimidating, as in [ʎmʎpʎ] but it should be noted that most words contain a vowel.

Our implementing system says nothing about consonant syllabicity, but we could admit that it is derived from a consonant surrounded by more obstruent consonants. In this case, the following pseudo-language resembles Berber.

| | | | | | | | | | |
|--------------|----------------|--------------|------------|-------------|------------|--------------|---------------|-------------|------------|
| iĉp | θuuθ? | χo | ιφνφ | ɾv | ɬuθ | vx̄φ | ɾvχɬ | ηθuu | v̄x̄φv̄ɬ |
| ʔ | ɾmuuθ | x̄φv̄x̄φ | əθɣ | φʔ | θʔə | p | ɾɬi | θa | yɬiɬ |
| ĉp̄v̄ĉp̄u | iĉp̄ə | əχə | əp̄uɪφ | ɣ | v̄x̄φ̄aχ | kmɾəφ | mɾʔ | p̄v̄ɬ | v̄ĉp̄iĉp̄u |
| moɬ | ɣmə | yĉp̄ə | ɣĉp̄u | θvʔ | yəm | θʔ | əχə | yɬɣ | ɬχ |
| v̄φi | θɾʔp̄v̄ | v̄x̄φ | ηaχ | ηmə | imχ | φəɬo | ɣφə | ikv | ĉp̄ək |
| ɬmʔφɬ | θv̄ɬa | mʔ | ɬoɬv̄ | kimʔ | muʔu | əθa | oɬi | uʔv̄ɬ | ɬχp |
| əθɣ | χpɣ | u | ypə | akv̄x̄φ | χp | it | məɬ | miĉp̄i | axoʔ |
| iχo | əɬ | ɾĉp̄ɣ | ɾp | yφuɪp | ɬəp | iφv̄ | ikyk | ɬək | φv̄θ |
| ɣpəp | aĉp̄ | pə | ypim | əχv̄p | m | uɬɣθ | ʔpɣĉp̄ | ηip | x̄φu |
| pum | χi | ηɬ | ɾθi | θəφɬ | mky | θa | ɾθʔ | θɣmχ | əɬ |

Table 9: Module 2 - Random example

Once again, an *a priori* principle is sufficient to simulate different types of syllabic constraints without introducing the notion of syllable or without referring to the patterns of natural languages. The results obviously deserve to be improved, but they already show that we can explain a lot with very little.

3.2. Module 3: Quantitative combination

3.3.1. Assumption and implementation

Our last module aims to constrain the size of the simulated words. Still adopting a functionalist viewpoint, we assume that words must be as short as possible to preserve articulatory energy, but they must be large enough to allow the creation of a lexicon composed of distinct elements. In sum, the average word size must be strictly sufficient to allow a number of distinct combinations that corresponds to a universal lexicon size. This postulate is still in line with the logic of maximizing contrasts.

To implement this, we define a lexicon size arbitrarily (beyond a certain size, the difference doesn't make much difference) and we count the number of units in the phonological inventory. Then we calculate the average word size that is needed to derive a number of distinct combinations equal to the lexicon size, as in the example below:

Words: 200 000
 Segments: 34
 Average word size: 3,461379675

As before, we introduce the possibility of variation by applying a random parameter that defines the variation margin of each word from a normal distribution law.

3.3.2. Results

In the chosen example in Table 10, we obtain mostly 1-syllable words. Only a very small number of the pseudo-words, in bold, contain two syllables.

| | | | | | | | | | |
|-------------|------|------|-------------|------|-------------|------------|--------------|-------------|-------------|
| l | pΔj | ζφΔm | θχελ | jθ | θNθ | ϕρεmυ | υφενυ | εθ | ρυφε |
| ζφ | θχ | ινp | εcθ | jϕpΔ | nyυς | ρεζφ | ρε | ρε | em |
| υφ | πευ | ιR | Δςυp | pι | υ | ινυ | εϕρες | mum | œnf |
| γl | ςε | Nφep | yc | ρqι | ЈεNθ | tυp | ρεεζφ | jφΔ | ιqv |
| φυrΔ | θl | tœN | ϕpim | εθt | ru | cΔ | χι | ηΔιq | icυ |
| ciθq | εtε | ϕpυt | qupι | quφ | tym | ιtυ | φεν | tυt | χœ |
| γnœ | ιN | χεθχ | rΔ | ϕpι | εt | rq | ιqε | νιθe | Јγpυ |
| œχ | ιnqθ | ηt | ιmηΔ | υςm | θc | cθ | ene | ρεεq | ςγ |
| θςœt | pιθ | ςφυt | rΔj | ρυφ | r | υr | ηε | npυ | η |
| jme | ρευ | ιχœη | ιqε | qεN | ηθtι | ηq | εq | υφι | ιφ |

Table 10: Module 3 - Chosen example (after 15 generations)

In the random example in Table 11, we can see that most simulated words have more than one syllable, and some of them can go up to three syllables, as [qzũtũ].

| | | | | | | | | | |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| υpυ | ze | zyq | γbat | tyqv | zaq | γbũtũ | taqap | γqazũ | piq |
| apap | tatũ | zũva | vbĩ | rbata | ataz | tyrĩ | ĩpaq | ata | vazũ |
| ũvũ | vqt | ava | γbaqv | ĩ | zata | zaq | γpĩqa | avava | tĩpav |
| qzũtũ | taq | γqrũ | rav | pazũv | tũ | zũvũ | ap | qvʃat | ĩvaz |
| rqe | bpũ | γpĩ | zavũ | γt | γbγp | γbav | zav | tĩqũt | taz |
| tanũ | tiv | trq | ĩqvba | ataz | zapa | azav | pa | rqvũ | bpiv |
| pav | tĩ | tũita | ava | ĩvaqv | piv | ĩqĩ | γ | vatĩ | tĩt |
| γbavũ | zava | bpava | γbava | zap | γbavũ | rtap | taza | v | ũt |
| rqyta | qvʃa | vap | btat | zapat | bativ | ĩpaza | trqav | rqvũv | zpi |
| γpũ | tũvũt | zavq | zũt | γbaz | qĩva | zρ | ĩtaz | zũtat | vatĩ |

Table 11: Module 3 - Random example

Thus, our *a priori* assumption can simulate languages with very various word sizes. Of course, morphology has an obvious role in defining the word size and a complete simulation must take it into account. But it is interesting to note that phonological principles may also be sufficient to explain the variation of languages on this issue.

Conclusion

In this paper, we showed that modeling and simulation are two different things and that simulating plausible languages helps to understand what categories are primitive and what categories can be derived from more general principles.

Do our simulated languages look like natural languages? They do in the sense that our approach derives plausible languages displaying various types of restrictions on the inventory and the combination of units.

Our approach is very distinct from those based on modeling, which aim to define prototypical languages. Indeed, the aforementioned parametric restrictions do not emerge from categories induced by the observation of natural languages: they emerge from nothing but *a priori* principles.

The reader should keep in mind that the usual definition of a prototypical model is itself biased by the nature of the most observed languages, European languages. But many languages surprise us every day with their ‘unnatural’ aspects. Thus, though some of our results may seem odd or unnatural - and it is sure that they should be improved - this should not be used to reject without nuance this method. Experimental archaeologists faced similar difficulties in the early days of the discipline:

I don't claim to have the "orthodox" method of carving. It is very possible that the "Mousterian" technique exposed is not the one that was used by the Mousterians, or by all the Mousterians. These experiments are still in progress. They have not yet allowed me to completely reproduce the magnificent "en echarpe" of Egyptian flints.

(Bordes, 1947, p. 1-2)

The most important is that we can simulate some of the aspects found in natural languages with a minimum of assumptions. For instance, we found a pseudo-language with initial stress and vowel reduction without setting the category of stress. This can be observed in Table 12 where closed vowels, in bold, can be found only in initial syllables.² The other vowels can be found in initial or non-initial syllable. This reduction of the inventory in some syllables is what we usually call a vowel reduction, and this indicates here the presence of stress in initial syllable.

| | | | | | | | | | |
|--------------|--------------|--------------|--------------|------------|------------|--------------|-------------|-------------|-------------|
| ʃəʔəʔ | jʔə | ʔəʔə | wʔə | məpəp | əʔə | həʔə | pəpə | jʔ | həʔə |
| məpəp | ʊʔə | nəpəp | ɪʔəʔ | ʃəʔəʔ | ʔəʔə | ʃəp | ɪpəp | məpə | məpə |
| nəʔə | əpə | həʔə | əʔə | əpəp | pəpə | ʔə | ʃəʔə | əpə | ʔəʔə |
| həʔ | məpəp | jʔ | həʔ | pəp | ɪpə | ʊʔəʔ | pəpə | ɪpə | fəpə |
| ʊʔəʔə | wʔə | məp | ɪpəpə | pəp | pəpəp | əʔ | əʔəʔ | wʔəʔ | ɪʔə |
| əpə | həʔə | əpəp | ɪpə | pəpə | jʔə | ɪpəp | əʔə | əp | nəpə |
| wʔə | ʊʔəʔə | ʊʔəʔə | ʊʔ | əpəpə | əʔ | ʔəʔə | nəpə | nəʔəʔ | ʃəʔə |
| ʊʔə | ɪpəp | pəpəp | ə | fəp | pəpəp | ɪpəpə | əʔ | ɪpəp | nəʔə |
| ɪʔ | ɪpə | əpəp | əpəp | ɪpə | əʔəʔ | nəʔə | ɪpə | nəpə | ɪpəp |
| jʔ | mə | ʊʔəʔə | ɪʔəʔ | jʔə | fəp | hə | jʔə | ɪpəp | jʔə |

Table 12: Example of initial stress and vowel reduction

This implies to ask whether the notion of stress is ultimately necessary when we manage to simulate, without it, a process generally attributed to prosody.

Next, we should adapt our principles to morphology, syntax and semantics. We should also improve our results in phonology, for example by defining a more satisfying phonetic continuum. But for the time being, we showed to what extent constructing languages is an interesting way to do science.

REFERENCES

- Backley, P. (2011). *An Introduction to Element Theory*. Edinburgh University Press.
- Bordes, F. (1947). Etude comparative des différentes techniques de taille du silex et des roches dures. *L'anthropologie*, 51:1-29.
- Brown, J. C. (1960). Loglan. *Scientific American*, 202(6):53-63.
- de Saussure, R. (1914). *La vort-teorio en Esperanto*. Universala Esperantia Librejo.
- Gillioz, C. and Zufferey, S. (2020). *Introduction to Experimental Linguistics*. ISTE, London.
- Goldsmith, J. A. (1976). *Autosegmental Phonology*. PhD dissertation [ms], MIT.
- Hjelmslev, L. (1943). *Prolegomena to a Theory of Language*. University of Wisconsin Press, Madison, [1969] edition.
- Kaye, J., Lowenstamm, J., and Vergnaud, J.-R. (1985). The internal structure of phonological representations: a theory of charm and government. *Phonology Yearbook*, 2:305-328.
- Lang, S. (2014). *Toki Pona: The Language of Good*. CreateSpace, Charleston.

² It should be mentioned that these simulated words were obtained with a phonetic continuum slightly different from the one presented in this article

- Liljencrantz, J. and Lindblom, B. (1972). Numerical simulation of Vowel Quality Systems: The Role of Perceptual Contrasts. *Language*, 48(4):839- 862.
- Martinet, A. (1955). *Economie des changements phonétiques: traité de phonologie diachronique*. A. Francke, Berne.
- McCarthy, J. J. (1979). *Formal problems in Semitic phonology and morphology*. PhD dissertation [ms], MIT Cambridge, Mass.
- Ogden, C. K. (1930). *Basic English: A General Introduction with Rules and Grammar*. Paul Treber & Co., London.
- Prince, A. S. and Smolensky, P. (1993). *Optimality Theory: Constraint interaction in generative grammar*. Rutgers Center for Cognitive Science, New Brunswick, [2002] edition.
- Zamenhof, L.-L. (1887). *Mezhdunarodnyj Jazyk*. Kh. Kel'ter, Varsovia.

Appendix: formula

Module 1

```
=(SIN((2*PI()/V!$C$3)*((ROW(A1)-2)-(V!$C$3/4)))+1)+(SIN((2*PI()/V!$C$4)*((COLUMN(A1)-2)-(V!$C$4/4)))+1)+SQRT(NORM.INV(RAND();0;25)^2)*(SIN((2*PI()/V!$C$3)*((ROW(A1)-2)-(V!$C$3/4)))+1)+(SIN((2*PI()/V!$C$4)*((COLUMN(A1)-2)-(V!$C$4/4)))+1)/100
```

```
V!$C$3=ArrayFormula(LARGE(IF(P!$A$1:$Z$1000<>"";COLUMN(P!$A$1:$Z$1000));1))/((50+NORM.INV(RAND();0;12,5))*ArrayFormula(LARGE(IF(P!$A$1:$Z$1000<>"";COLUMN(P!$A$1:$Z$1000));1))/100)
```

```
V!$C$4=ArrayFormula(LARGE(IF(P!$A$1:$Z$1000<>"";ROW(P!$A$1:$Z$1000));1))/((50+NORM.INV(RAND();0;12,5))*ArrayFormula(LARGE(IF(P!$A$1:$Z$1000<>"";ROW(P!$A$1:$Z$1000));1))/100)
```

P!\$A\$1:\$Z\$1000= [phonetic space]

Modules 2 and 3

```
=IF(B1="";";IF(COLUMN(C1)-1>V!$B$6+INT(NORM.INV(RAND();0;25))*V!$B$6/100;";DECALER(INDIRECT(ADDRESS(MATCH(B1;T!$A:$A;0);ArrayFormula(SMALL(IF(INDIRECT(ADDRESS(MATCH(B1;T!$A:$A;0);2;;; "T")):INDIRECT(ADDRESS(MATCH(B1;T!$A:$A;0);ArrayFormula(SMALL(IF(T!$1:$1="";COLUMN(T!$1:$1);2))-1;;; "T"))>=INT(INDIRECT(ADDRESS(MATCH(B1;T!$A:$A;0);ArrayFormula(SMALL(IF(T!$1:$1="MAX";COLUMN(T!$1:$1);1)))); "T"))-(V!$B$2*INDIRECT(ADDRESS(MATCH(B1;T!$A:$A;0);ArrayFormula(SMALL(IF(T!$1:$1="MAX";COLUMN(T!$1:$1);1)))); "T"))/100));COLUMN(INDIRECT(ADDRESS(MATCH(B1;T!$A:$A;0);2;;; "T")):INDIRECT(ADDRESS(MATCH(B1;T!$A:$A;0);ArrayFormula(SMALL(IF(T!$1:$1="";COLUMN(T!$1:$1);2))-1;;; "T"))));INT(RAND()*NB.IF(INDIRECT(ADDRESS(MATCH(B1;T!$A:$A;0);2;;; "T")):INDIRECT(ADDRESS(MATCH(B1;T!$A:$A;0);ArrayFormula(SMALL(IF(T!$1:$1="";COLUMN(T!$1:$1);2))-1;;; "T")));">=INT(INDIRECT(ADDRESS(MATCH(B1;T!$A:$A;0);ArrayFormula(SMALL(IF(T!$1:$1="MAX";COLUMN(T!$1:$1);1)))); "T"))-(V!$B$2*INDIRECT(ADDRESS(MATCH(B1;T!$A:$A;0);ArrayFormula(SMALL(IF(T!$1:$1="MAX";COLUMN(T!$1:$1);1)))); "T"))/100))+1));;"T"));1-ROW(INDIRECT(ADDRESS(MATCH(B1;T!$A:$A;0);1;;; "T")));0))
```

V!\$B\$1=200,000

V!\$B\$2=SQRT(INT(NORM.INV(RAND();0;25))^2)+1

V!\$B\$6=log(\$B\$1;LIGNES(T!\$A:\$A)-NB.IF(T!\$A:\$A;""))

B1= [preceding segment]

T!= [table of proximity]