

# Emergence of sound systems through self-organisation

Bart de Boer  
Artificial Intelligence Laboratory  
Vrije Universiteit Brussel  
Pleinlaan 2, 1050 Brussel  
bartb@arti.vub.ac.be

## Abstract

This paper describes a model for explaining the emergence and the universal structural tendencies of vowel systems. Both are considered as the result of self-organisation in a population of language users. The language users try to imitate each other and to learn each other's vowel systems as well as possible under constraints of production and perception, while at the same time maximising the number of available speech sounds. It is shown through computer simulations that coherent and natural sound systems can indeed emerge in populations of artificial agents. It is also shown that the mechanism that is responsible for the emergence of sound systems can be used for learning existing sound systems as well. Finally, it is argued that the simulation of agents that can only produce isolated vowels is not enough. More complex utterances are needed for other interesting universals of sound systems and for explaining realistic sound change. Work in progress on implementing agents that can produce and perceive complex utterances is reported.

## Introduction

The research described in this paper tries to explain the emergence and structure of systems of speech sounds. It investigates how a coherent system of speech sounds can emerge in a population of agents and how the constraints under which the system emerges impose structure through self-organisation. If self-organisation can explain structure, then innate and biologically evolved mechanisms are not necessary. This effectively decreases the number of linguistic phenomena that have to be explained by biological evolution.

What are the phenomena that have to be explained by a theory of the emergence of speech sounds? The systems of speech sounds in the world's languages show remarkable regularities. First of all, certain sounds occur much more frequently than others. In the UPSID<sup>1</sup>, a database that contains the phoneme inventories of 451 languages, (the first version with 317 languages is described in Maddieson 1984) the vowels [i], [a] and [u] appear in 87%, 87% resp. 82% of the languages, while the vowels [y], [œ] and [ɯ] occur in only 5%, 2% resp. 9% of the languages. This holds even more for consonants. Some consonants, e.g. [m] (94%), [k] (89%) or [j] (84%) appear very frequently, while others, e.g. [ʀ] (1%), [ʃ'] (1%) and [ʔ] (1%) appear very rarely. The sound systems of languages also display a fair amount of symmetry. If a language has a front unrounded vowel of a given height, for example an [e] (occurring in 27% of the languages), it is

---

<sup>1</sup>UCLA Phonological Segment Inventory Database

quite likely that it also has the corresponding back rounded vowel [o] (which occurs in 29% of all languages, but in 85% of the languages with [e]). In the case of consonants, if a language has a voiced stop at a given place of articulation, e.g. [d] (27%) it usually also has a [t] (40% in whole sample vs. 83% in languages with [d]).

Sometimes these universal characteristics are explained by innate properties of the brain (Jakobson & Halle 1956, Chomsky & Halle 1968). However the question then becomes how these innate properties have evolved. Also, if there are innate constraints it is not clear why there is still such huge variation between different languages. It is clearly preferable to have an explanation that does not need innate mechanisms.

Functional explanations of the above mentioned phenomena are more satisfying. A number of articulatory, perceptual and cognitive criteria have been proposed (e.g. Carré et al. 1995, Liljencrants & Lindblom 1972, Lindblom 1992, Stevens 1972). Some of these have been tested with computers simulations. These criteria can be summarised as articulatory ease, acoustic distinctiveness and minimum effort of learning.

These functional explanations are not the full explanation, either. They assume that the systems of speech sounds one finds are the result of an optimisation of one or more of the proposed criteria. However, it is not clear who is doing the optimisation. Certainly children that learn a language do not do an optimisation of the system of speech sounds they learn. Rather, they try to imitate their parents (and peers) as accurately as possible. This explains the fact that people can speak the same language with different accents, from which one can identify their place of birth or their social group.

If none of the individual speakers does an explicit optimisation of their sound system, but still (near-) optimal sound systems are found more frequently than non-optimal ones, it is clear that the optimisation must be an emergent property of the interactions in the population. Therefore, if one wants to explain the sound systems that are found in populations of agents that imitate and learn each other's sounds under acoustic, articulatory and cognitive constraints.

A first attempt at building a computer model of a population of interacting agents for explaining the shape of vowel systems was undertaken by Glotin (Glotin 1995) later followed by Berrah (Berrah 1998). Both methods have the drawback that the population is subject to some genetic evolution and that the agents still do local optimising by pushing the vowels in their vowel systems away from each other. Also the number of vowels in every agent has to be fixed beforehand in these simulations.

In this paper a system is presented in which a population of agents that are each able to produce, perceive and learn vowels, develops a coherent system of vowel sounds that confirms to the tendencies of vowel systems in human languages. The number of vowels need not be fixed beforehand and there is no genetic evolution of the agents. Although the agents are able to change their repertoire of vowels in order to optimise the successfulness of imitation, they only do this in reaction to interactions with other agents. They also cannot change the positions of their vowels in any global way. The emerging vowel systems are therefore truly the result of the interactions between the agents. The research is based on Steels' (Steels 1996, 1997, 1998) ideas on the origin of language. Steels considers language as the result of a process of mainly cultural evolution, while the universal tendencies of language can be explained as the results of self-organisation under constraints of perception and production. Steels has applied his ideas mainly to lexicon and meaning formation, and is now working on syntax.

In the next two sections, the agents and their interactions are described in considerable detail. In section 3 some results of the simulations that were performed with this system are presented. Finally, in section 4 conclusions, future work and a discussion of the work are represented.

tions are described in considerable detail. In section 3 some results of the simulations that were performed with this system are presented. Finally, in section 4 conclusions, future work and a discussion of the work are represented.

## 1 The agents

The agents are equipped with an articulatory synthesiser for perception and a prototype list for storage of vowel constructed to be as human like as possible, in order to make to research in linguistics and in order to make it possible vowels.

or production, a model of human hearing. All the elements of the agent were the results of the research applicable to use the agents to learn *real* human

An agent consists of three parts ( $S, D, V$ ) where  $S$  is the synthesis function,  $D$  is the distance measure and  $V$  is the agent's set of vowels. The synthesiser function is a function  $S: A_r \rightarrow A_c$ , where  $A_r$  is the set of possible articulations and  $A_c$  is the set of possible acoustic signals. For the agents presented in this section the set of possible articulations is the set of articulatory vectors  $(p, h, r)$  where  $p, h, r$  are real numbers in the range  $[0, 1]$ . Parameters  $p, h$  and  $r$  are the major vowel features (Ladefoged and Maddieson, 1996: ch. 9) *position, height* and *rounding*. Position corresponds (roughly) to the position of the highest point of the tongue in the front to back dimension, height corresponds to the vertical distance between the highest part of the tongue and the roof of the mouth and rounding corresponds to the rounding of the lips. Position zero means most fronted, height zero means lowest and rounding zero means that the lips are maximally spread. The parameter values for the high, front, unrounded vowel [i], such as in "leap" are  $(0, 1, 0)$ . For the high, back rounded vowel [u], such as in "loop" they are  $(1, 1, 1)$ . For the low, back, unrounded vowel [a] such as in "father" they are  $(1, 0, 0)$ .

The set  $A_c$  of possible outputs of the synthesiser function consists of vectors  $(F_1, F_2, F_3, F_4)$  where  $F_1, F_2, F_3, F_4 \in \mathbf{R}$  are the first four formant frequencies of the generated vowel. These formant frequencies correspond to the peaks in the power spectrum of the vowel. When agents communicate among each other, they exchange only the formant values, not a real signal. This is done to reduce the amount of data. A certain amount of noise is added, however. This noise consists of a random shifting of the formant frequencies, according to the following formula:

$$1) F_i \leftarrow \left( 1 + \frac{Noise\%}{100} U(-0.5, 0.5) \right) F_i.$$

In which  $U(-0.5, 0.5)$  is a random number drawn from the uniform distribution between  $-0.5$  and  $0.5$ , *Noise%* is the noise percentage (a parameter of the system) and  $F_i$  represents the formants.

The formant frequencies are generated by a three dimensional quadratic interpolation between sixteen data points that have been generated by Maeda's articulatory synthesiser (Maeda, 1989, Vallée, 1994, pp. 162–164). The equations for calculating the synthesiser function are shown in. The formant values for [i] are  $(252, 2202, 3242, 3938)$ , for [u] :  $(276, 740, 2177, 3506)$  and for [a]:  $(703, 1074, 2356, 3486)$ . An important property of the synthesiser function is that it is easy to calculate the formant frequencies from the articulatory description, but that it is very hard to calculate the articulatory description from the acoustic description. With this synthesiser all basic vowels can be generated. It is therefore *language-independent*.

A vowel  $v$  consists of elements  $(ar, ac, s, u)$ , where  $ar \in A_r$  is the articulatory prototype,  $ac \in A_c$  is the corresponding acoustic prototype and  $s, u$  are the success and use scores, (which will be explained with the imitation game) respectively. The vowels are represented as prototypes as this

seemed to be both a realistic and computationally effective way to represent vowels. Research in human perception of speech sounds (e.g. Cooper et al. 1976) seems to indicate that humans perceive speech sounds in terms of prototypes. If human subjects are presented with acoustic signals that vary continuously from one speech sound to another, (i.e. from [ga] to [ba]) they tend to perceive these signals as either the one category [ba] or the other [ga], never a something “in between”. Perception suddenly switches somewhere in the middle.

An agent's vowels are stored in the set  $V$ , which we will call the vowel set. When an agent decides it has encountered a new vowel  $v_{new}$  (we will describe below how and when this is decided), it adds both the acoustic and the articulatory descriptions of  $v_{new}$  to  $V$ :  $V \leftarrow V \cup v_{new}$ . A sound  $A$  that the agent hears will be compared to the acoustic prototypes  $ac_v$  of the vowels  $v$  in its vowel set, and the distance between  $A$  and all  $ac_v$  ( $v \in V$ ) is calculated using the distance function  $D: Ac^2 \rightarrow \mathbf{R}$  (described below). It then considers that it has recognised the vowel  $v_{rec}$  that has minimal distance to  $A$ :  $\{v_{rec} | v_{rec} \in V \cap \neg \exists v_2 : (v_2 \in V \cap D(A, ac_{v_2}) < D(A, ac_{v_{rec}}))\}$ .

The distance between two vowels is determined by using a weighted distance in the  $F_1$ - $F_2'$  space, where  $F_1$  is the frequency of the first formant (expressed in Bark, a logarithmic frequency scale) and  $F_2'$  is the weighted average of the second, third and fourth formants (also expressed in Barks). This distance measure is based on the distance measure described by Mantakas et. al (1986). The distance measure is based on weighting formant peaks differently depending on their distance relative to a critical distance  $c$ , which is taken to be 3.5 Bark. In order to calculate  $F_2'$ , two weights have to be calculated:

$$2) w_1 = \frac{c - (F_3 - F_2)}{c}, w_2 = \frac{(F_4 - F_3) - (F_3 - F_2)}{F_4 - F_2}$$

Where  $w_1$  and  $w_2$  are the weights and  $F_1$ - $F_4$  are the formants in Bark. The value of  $F_2'$  can now be calculated as follows:

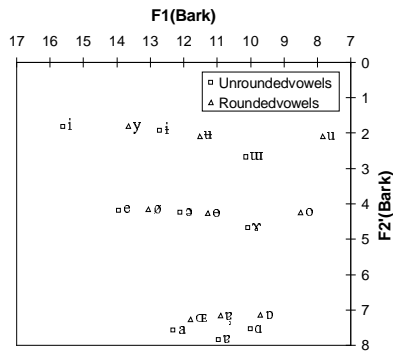
$$3) F_2' = \begin{cases} F_2, & \text{if } F_3 - F_2 > c \\ \frac{(2 - w_1)F_2 + w_1F_3}{2}, & \text{if } F_3 - F_2 \leq c \text{ and } F_4 - F_2 > c \\ \frac{w_2F_2 + (2 - w_2)F_3}{2} - 1, & \text{if } F_4 - F_2 \leq c \text{ and } F_3 - F_2 < F_4 - F_3 \\ \frac{(2 - w_2)F_3 + w_2F_4}{2} - 1, & \text{if } F_4 - F_2 \leq c \text{ and } F_3 - F_2 \geq F_4 - F_3 \end{cases}$$

The values of  $F_1$  and  $F_2'$  for a number of vowels is shown in figure 1. We can see from this figure that the distribution of the vowel through the acoustic space is quite natural. However, as it is a 2-dimensional projection of an essentially 4-dimensional space, not all distances between all phonemes can be represented accurately. The distance between two signals,  $a, b \in Acc$  can now be calculated using a weighted Euclidean distance:

$$4) D(a, b) = \sqrt{(F_1^a - F_1^b)^2 + \lambda(F_2'^a - F_2'^b)^2}$$

The value of the parameter  $\lambda$  is 0.5 for all experiments that will be described.

With the synthesis function and the distance measure that have been described in this section, the agents can produce and perceive speech sounds in a human-like way. The results that are gener-



ated with this system can therefore be compared with the results of research into human sound systems.

the results of research into human sound systems.

## 2 The imitation game

The imitation game was designed for allowing the agents to determine the vowels of the other agents and to develop a realistic vowel system. The imitation game is played in a population of agents (size 20 in all the experiments presented here). From this population two agents are picked at random: an *initiator* and an *imitator*. The initiator starts the imitation game by producing a sound that the imitator has to imitate. The imitator listens to the sound, and tries to analyse it in terms of the sound prototypes it already knows. It then produces the acoustic signal of the prototype it found. The initiator then listens to this signal and analyses it in terms of its prototypes. If the prototype it finds is the same as the one it used to produce the original sound, the game is considered *successful*. Otherwise it is a *failure*. This is communicated to the imitator. The exact steps of the imitation game are illustrated in table 1. Note non-verbal feedback is needed to indicate whether the game was a success or a failure. If observed, the non-verbal feedback can be compared to gesture or facial expression or the failure to achieve a communicative goal. Making the imitation game dependent on non-verbal communication might seem like introducing a very unrealistic element in the agents' learning. To human children it is hardly ever directly indicated whether the sounds they produce are right or wrong. However, there are more indirect ways of discovering that the right sound was not used, such as a failure to achieve the desired goal of the communication. But our imitation game abstracts from this and assumes that a feedback signal is somehow available. Depending on the outcome of the imitation game, the imitator can alter its vowel inventory. The way this is done is described in table 2, together with a number of other routines that are reused. However, if the imitation game was a failure and the vowel that was used has a low use-to-

to determine the vowels of the other agents and to develop a realistic vowel system. The imitation game is played in a population of agents (size 20 in all the experiments presented here). From this population two agents are picked at random: an *initiator* and an *imitator*. The initiator starts the imitation game by producing a sound that the imitator has to imitate. The imitator listens to the sound, and tries to analyse it in terms of the sound prototypes it already knows. It then produces the acoustic signal of the prototype it found. The initiator then listens to this signal and analyses it in terms of its prototypes. If the prototype it finds is the same as the one it used to produce the original sound, the game is considered *successful*. Otherwise it is a *failure*. This is communicated to the imitator. The exact steps of the imitation game are illustrated in table 1. Note non-verbal feedback is needed to indicate whether the game was a success or a failure. If observed, the non-verbal feedback can be compared to gesture or facial expression or the failure to achieve a communicative goal. Making the imitation game dependent on non-verbal communication might seem like introducing a very unrealistic element in the agents' learning. To human children it is hardly ever directly indicated whether the sounds they produce are right or wrong. However, there are more indirect ways of discovering that the right sound was not used, such as a failure to achieve the desired goal of the communication. But our imitation game abstracts from this and assumes that a feedback signal is somehow available. Depending on the outcome of the imitation game, the imitator can alter its vowel inventory. The way this is done is described in table 2, together with a number of other routines that are reused. However, if the imitation game was a failure and the vowel that was used has a low use-to-

**table 1: Basic organisation of the imitation game.**

<b>initiator</b>	<b>imitator</b>
<b>if</b> ( $V = \emptyset$ ) Add random vowel to $V$ Pick random vowel $v$ from $V$ $u_v := u_{v-1} + 1$ Produce signal $A_1 := ac_v$	
	Receives signal $A_1$ . <b>if</b> ( $V = \emptyset$ ) $v_{new} := \text{Findphoneme}(A_1)$ $V := V \cup v_{new}$ Calculate $v_{rec}$ : $v_{rec} \in V \wedge \neg \exists v_2 : (v_2 \in V \wedge D(A_1, ac_{v_2}) < D(A_1, ac_{v_{rec}}))$ Produce signal $A_2 := ac_{v_{rec}}$
Receives signal $A_2$ . Calculate $v_{rec}$ : $v_{rec} \in V \wedge \neg \exists v_2 : (v_2 \in V \wedge D(A_2, ac_{v_2}) < D(A_2, ac_{v_{rec}}))$ <b>if</b> ( $v_{rec} = v$ ) Send non-verbal feedback: <i>success</i> . $s_v := s_v + 1$ <b>else</b> Send non-verbal feedback: <i>failure</i> .	
Do other updates of $V$ .	Receive non-verbal feedback. Update $V$ according to feedback signal. Do other updates of $V$ .

**table2:Actionsperformedbytheagents**

<pre> Shiftcloser ( v, A); return v<sub>best</sub> {   v<sub>best</sub> := v   for(allsixneighbors v<sub>neigh</sub>of v)do:     if (D(ac<sub>vneigh</sub>, A) &lt; D( ac<sub>vrec</sub>, A))       v<sub>best</sub> := v<sub>neigh</sub> } </pre>	<pre> Findphoneme ( A); return v<sub>best</sub> {   vowel v:   ar<sub>v</sub>=(0.5,0.5,0.5)   ac<sub>v</sub>=S( ar<sub>v</sub>)   s<sub>v</sub>=0   u<sub>v</sub>=0   do     v<sub>best</sub> := v     v:=Shiftcloser( v<sub>best</sub>,A )   until( v= v<sub>best</sub>) } </pre>	<pre> Updateaccordingtofeedbacksignal {   u<sub>vrec</sub> := u<sub>vrec</sub>+1   if (feedbacksignal= success)     v<sub>rec</sub>:=Shiftcloser( v<sub>rec</sub>, A<sub>1</sub>)     s<sub>vrec</sub> := s<sub>vrec</sub>+1   else     if( u<sub>vrec</sub>/s<sub>vrec</sub>&gt; threshold)       v<sub>new</sub>:=Findphoneme( A<sub>1</sub>)       V:= V ∪ v<sub>new</sub>     else       v<sub>rec</sub>:=Shiftcloser( v<sub>rec</sub>, A<sub>1</sub>) } </pre>
--	--	---

success ratio, the vowel was probably not a good imitation shifted toward the signal that was heard in the hopeth. The phoneme is not thrown away. This is done in the These routines do three things: they throw away bad vowel minimum number of times (five times in all experiments presented). Their use-to-success ratio is less than a threshold (0.7 in all experiments presented). Also, vowels that are too close in articulatory and acoustic space can prevent a cluster of bad phonemes to emerge at a position required. This has been observed in experiments without merging. The articulatory threshold for merging is the minimal distance to a neighbouring prototype set. The acoustic threshold for merging is determined by the noise level. If two vowels are so close that they can be confused by the noise that is added to the formant frequencies, they are merged. The last change agents can make to their vowel inventory is adding a random new vowel. This is done with a low probability (0.01 in all experiments presented). The values for the articulatory and acoustic parameters of the new vowel are chosen randomly from a uniform distribution between 0 and 1. The imitation game contains all the elements that are necessary for the emergence of vowel systems and innovation: the noise, the imperfect imitations and the random insertions of vowels. Other mechanisms take care of (implicit) selection of good quality vowels: vowels are only retained if they exist in other agents as well, otherwise successful limitations are possible, and their success score will drop. Unsuccessful vowels will eventually be removed. The merging ensures that phonemes will stay apart, so that sufficiently spaced vowel systems emerge. Note that all the actions of the agents can be performed using local information only. The agents do not need to look at each other's vowel systems di-

ion of any other sound. It is therefore at it will become a better imitation. *other updates* routines, described in table 3. Vowels that have been tried at least a minimum number of times (five times in all experiments presented). Vowels are considered bad if their use-to-success ratio is less than a threshold (0.7 in all experiments presented). Also, vowels that are too close in articulatory and acoustic space can prevent a cluster of bad phonemes to emerge at a position where only one good vowel would be required. This has been observed in experiments without merging. The articulatory threshold for merging is the minimal distance to a neighbouring prototype set to be 0.03 in all experiments. The acoustic threshold for merging is determined by the noise level. If two vowels are so close that they can be confused by the noise that is added to the formant frequencies, they are merged. The last change agents can make to their vowel inventory is adding a random new vowel. This is done with a low probability (0.01 in all experiments presented). The values for the articulatory and acoustic parameters of the new vowel are chosen randomly from a uniform distribution between 0 and 1. The imitation game contains all the elements that are necessary for the emergence of vowel systems and innovation: the noise, the imperfect imitations and the random insertions of vowels. Other mechanisms take care of (implicit) selection of good quality vowels: vowels are only retained if they exist in other agents as well, otherwise successful limitations are possible, and their success score will drop. Unsuccessful vowels will eventually be removed. The merging ensures that phonemes will stay apart, so that sufficiently spaced vowel systems emerge. Note that all the actions of the agents can be performed using local information only. The agents do not need to look at each other's vowel systems di-

**table3:Otherupdatesoftheagents'vowelsystems**

<pre> Merge( v<sub>1</sub>, v<sub>2</sub>, V) {   if( s<sub>v1</sub>/u<sub>v1</sub>&lt; s<sub>v2</sub>/u<sub>v2</sub>)     s<sub>v2</sub> := s<sub>v2</sub>+ s<sub>v1</sub>     u<sub>v2</sub> := u<sub>v2</sub>+ u<sub>v1</sub>     V := V - v<sub>1</sub>   else     s<sub>v1</sub> := s<sub>v1</sub>+ s<sub>v2</sub>     u<sub>v1</sub> := u<sub>v1</sub>+ u<sub>v2</sub>     V := V - v<sub>2</sub> } </pre>	<pre> Dootherupdatesof V {   for( ∇ v ∈ V)//Removebadvowels     if( s<sub>v</sub>/u<sub>v</sub>&lt; throwawaythreshold ∧ u<sub>v</sub>&gt;min.uses )       V := V - v   for( ∇ v<sub>1</sub> ∈ V)//Mergingofvowels     for( ∇ v<sub>2</sub>:( v<sub>2</sub> ∈ V ∧ v<sub>2</sub> ≠ v<sub>1</sub>))       if(D( ac<sub>v1</sub>, ac<sub>v2</sub>)&lt; acousticmergethreshold )         Merge( v<sub>1</sub>, v<sub>2</sub>, V)       if(Euclidean distance between ar<sub>v1</sub> and ar<sub>v2</sub>&lt;         articulatorymergethreshold )         Merge( v<sub>1</sub>, v<sub>2</sub>, V)   Addnewvowelto Vwithsmallprobability. } </pre>
--	---

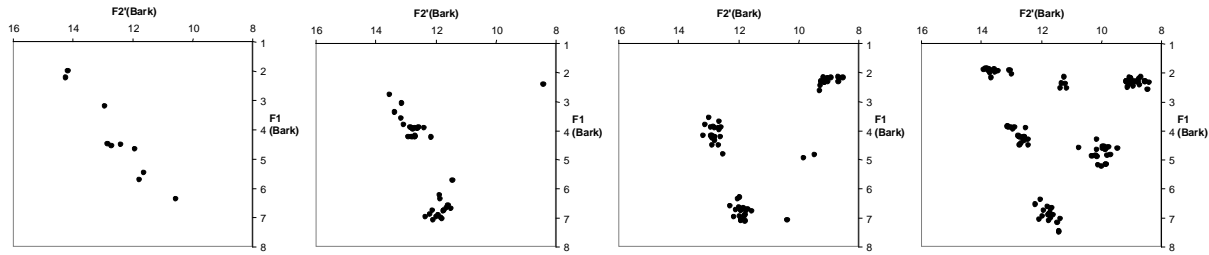


figure2: Vowelsystem after 20, 200, 1000 and 2000 games, 10% noise

rectly.

### 3 Vowel experiments

So far, only experiments with vowels have been done. The partly described in (de Boer 1997a, 1997b). The first aim of the coherent sound system can indeed emerge in a population of agents that learn such a sound system, but that do not have a sound system. It was to show that the system that is learnt has the same characteristics as human sound systems. Vowels were the signals of choice, as they are easy to represent, generate and perceive and because of the universal characteristics of human vowel systems more thoroughly described than those of other speech signals. A typical example of the emergence of a vowel system in a population of twenty agents with ten percent noise is illustrated in figure 2. In this figure the vowel systems of the agents in the population are shown after different numbers of imitation games. All vowels of all agents in the population are plotted on top of each other, consisting of the first formant  $F_1$  and the weighted sum of the second, third and fourth formants ( $F_2'$ ). The frequency of the formants is shown in the Bark frequency scale. Due to articulatory limitations the acoustic space that can be reached by the agents is roughly triangular with the apex at the bottom of the graph.

In the leftmost graph the agents' vowels after 20 imitation games are shown. One can see hardly any structure at all; the vowels are dispersed throughout the acoustic space (the apparent linear correlation is just an artefact). This is caused by the fact that initially vowels are mostly added at random. After 200 imitation games, clusters emerge. This happens because the agents try to imitate each other as closely as possible while at the same time there is a pressure of having a maximal number of vowels. Almost every agent in the population now has two vowels: one in each cluster.

After 1000 imitation games the available acoustic space starts to get full, and the clusters become tighter. Every agent in the population now has at least three vowels. Some agents have more (the isolated dots in the graph), other agents have not had the opportunity to copy these yet. Finally, after 2000 imitation games, the available acoustic space is completely covered. The system that emerges consists of tight clusters that are approximately equally spaced. The vowels that emerge are [i], [e]-[ø], [a], [o], [u] and [ɨ] which, except for the

se experiments have already been done. The first aim of the experiments was to show that agents that are in principle able to learn a sound system at the beginning. The second aim was to show that the system that is learnt has the same characteristics as human sound systems. Vowels were the signals of choice, as they are easy to represent, generate and perceive and because of the universal characteristics of human vowel systems more thoroughly described than those of other speech signals.

A typical example of the emergence of a vowel system in a population of twenty agents with ten percent noise is illustrated in figure 2. In this figure the vowel systems of the agents in the population are plotted on top of each other, consisting of the first formant  $F_1$  and the weighted sum of the second, third and fourth formants ( $F_2'$ ). The frequency of the formants is shown in the Bark frequency scale. Note that due to articulatory limitations the acoustic space that can be reached by the agents is roughly triangular with the apex at the bottom of the graph.

In the leftmost graph the agents' vowels after 20 imitation games are shown. One can see hardly any structure at all; the vowels are dispersed throughout the acoustic space (the apparent linear correlation is just an artefact). This is caused by the fact that initially vowels are mostly added at random. After 200 imitation games, clusters emerge. This happens because the agents try to imitate each other as closely as possible while at the same time there is a pressure of having a maximal number of vowels: one in each

cluster. After 1000 imitation games the available acoustic space starts to get full, and the clusters become tighter. Every agent in the population now has at least three vowels. Some agents have more (the

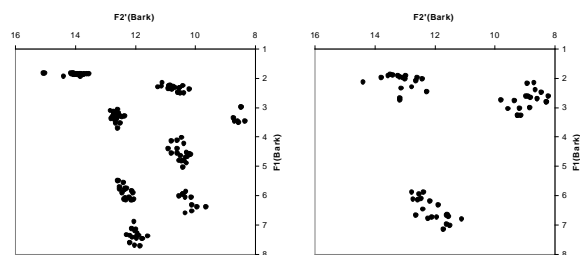


figure3: Systems with 10% and 25% noise

rounding of the front midsegment, is a possible six-vowel system (such as found, for example in the Saami language of Lapland).

The noise level determines the number and size of the clusters. The number of clusters will be lower and they will be more spread out. In a system with 10% noise, the clusters are still spread near-optimally throughout the available acoustic space. Both systems are also natural. The one with 10% noise has eight vowel clusters, while the one with 25% noise is the canonical three-vowel system, consisting of [i], [a] and [u]. Note that the vowel system that obtained under 10% noise in this simulation run is not the same as the one that obtained in figure 2. This is because the population does not converge to one optimal solution, rather it converges to a good system, which might, apparently, consist of 6 or 8 vowel clusters. Both systems, however, show similar characteristics of symmetry and spread of vowel clusters. These experiments show that a coherent sound system can emerge in a population of agents and that these sound systems show the same universal characteristics as sound systems from natural languages. However, there is no transfer from one generation of speakers to the next, yet. In real life, when speakers enter (they are born) and leave (they die or move away) the community constantly. Still, the language remains relatively stable. The simulation presented here can be used to test whether it is possible to transfer the sound system in a stable way from one generation to the next.

Succession of generations can be modelled by adding and removing agents from the population at random. These processes model birth and death of language users. After a sufficiently long period of time, all the original agents in the population will have learnt their sound system from the original population of new agents. The sound system in the population of new agents can then be compared with the original sound system. This is done in figure 4. The white squares represent the positions of the original agents' vowels and the black circles represent the positions of the vowels after 2000 imitation games. On average every 50 imitation games an agent was removed from or added to the population. The original population consisted of 20 agents, the final population consisted of 11 agents for the left graph and 14 agents for the right graph (the number of agents was not fixed, due to moving agents.) The noise level was a constant 10%.

In the simulation that resulted in the left graph, agents could change their vowel repertoire more easily when they were young than when they were old. Comparing the two graphs, it can be observed that both systems preserve the approximate positions of the clusters. However, in the left graph the clusters have become more dispersed, have merged and new clusters have emerged. In the right graph, the positions and number of clusters has hardly changed at all.

Apparently cultural transfer of sound systems is possible in both simulations. Extra stability is ensured when older agents can change their vowel systems less easily than younger agents. Apparently the older agents provide a stable target to which the younger agents can adapt their vowel systems.

usters. If the noise level is higher, the clusters are widely dispersed. This is shown in figure 3, where a system with 25% noise is compared with a system with 10% noise. Note however, that the clusters are still spread near-optimally throughout the available acoustic space. Both systems are also natural. The one with 10% noise has eight vowel clusters, while the one with 25% noise is the canonical three-vowel system, consisting of [i], [a] and [u]. Note that the vowel system that obtained under 10% noise in this simulation run is not the same as the one that obtained in figure 2. This is because the population does not converge to one optimal solution, rather it converges to a good system, which might, apparently, consist of 6 or 8 vowel clusters. Both systems, however, show similar characteristics of symmetry and spread of vowel clusters. These experiments show that a coherent sound system can emerge in a population of agents and that these sound systems show the same universal characteristics as sound systems from natural languages. However, there is no transfer from one generation of speakers to the next, yet. In real life, when speakers enter (they are born) and leave (they die or move away) the community constantly. Still, the language remains relatively stable. The simulation presented here can be used to test whether it is possible to transfer the sound system in a stable way from one generation to the next.

Succession of generations can be modelled by adding and removing agents from the population at random. These processes model birth and death of language users. After a sufficiently long period of time, all the original agents in the population will have been replaced and the new agents in the population. The sound system in the population of new agents can then be compared with the original sound system. This is done in figure 4. The white squares represent the positions of the original agents' vowels and the black circles represent the positions of the vowels after 2000 imitation games. On average every 50 imitation games an agent was removed from or added to the population. The original population consisted of 20 agents for the left graph and 14 agents for the right graph (the number of agents was not fixed, due to the independence of adding and removing agents.) The noise level was a constant 10%.

In the simulation that resulted in the left graph, agents could learn equally well, independent of whether they were young or old. Comparing the two graphs, it can be observed that both systems preserve the approximate positions of the clusters. However, in the left graph the clusters have become more dispersed, have merged and new clusters have emerged. In the right graph, the positions and number of clusters has hardly changed at all.

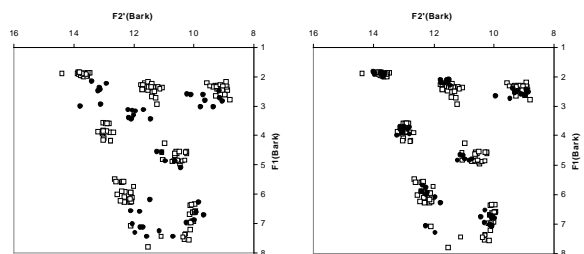


figure 4: Systems after population replacement.



## 4 Conclusions and discussion

The results of the simulations show clearly that coherent sound systems can emerge as the result of local interactions between the members of a population. They also show that the systems that emerge show characteristic tendencies similar to the ones that are found in human sound systems, such as more frequent use of certain vowels and symmetry of the system. This means that we do not need to look for evolutionary ways of explaining the universal tendencies of vowel systems. Apparently the characteristics emerge as the result of self-organisation under constraints of perception, production and learning. The systems that are found can be considered attractors of the dynamical system that consists of the agents and their interactions. Of course we still need an evolutionary account of the shape of the human vocal tract and of the performance of human perception, but we do not need any specific innate mechanisms for explaining the structure of the vowel systems that appear in human languages.

It has also been shown that the vowel systems can be transferred from one generation of agents to the next. For this, no change in the interactions and the behaviour of the agents has to be made, although the transfer from generation to generation is improved if older agents are made to learn less quickly than young agents. Apparently the same mechanism can be used to learn an existing vowel system as well as to produce a sound system in a population where no sound system existed previously. This lends support to Steels' (Steels 1997, 1998) thesis that the same mechanism that is responsible for the ability to learn language is responsible for the emergence of language in the first place. The use of computer simulations makes it easy for the researcher to perform experiments like these, and thus provides an extra means to test and fine-tune linguistic theories. The ability to explain the emergence, the learning and the universal structural tendencies of sound systems as the result of local interactions between agents that exist in a population is a remarkable result. It indicates that not all aspects of language need to be explained through biological evolution. This makes it easy to explain that language evolved in a relatively short time. It needs to be tested, however, whether these results also hold for more complex utterances than isolated vowels. Work is in progress (de Boer, 1998) on building agents that can produce and perceive complex utterances, using articulatory synthesizers, dynamically moving articulators, models of human perception and models of infant learning of speech. In any case, modelling aspects of language as the result of interactions in a population seems to be a promising way to learn more about the origins of language, especially so because it provides an extra mechanism next to biological evolution for explaining the complexity and structure of language.

## 5 Acknowledgements

This research was done at the artificial intelligence laboratory of the Vrije Universiteit Brussel. It is part of an ongoing research project into the origins of language and intelligence. Funding was provided by the GOA2 project of the Vrije Universiteit Brussel. Part of the work was done at the Sony CSL in Paris, France. I thank Luc Steels for valuable discussion of the ideas and the work presented here and for providing the research environment of the VUB AI-lab. I thank Edwin de Jong, Tony Belpaeme and Paul Vogt of the VUB AI-lab for their comments and suggestions and y. I also thank Björn Lindblom, Christine Ericsson and the other people of the phonetics laboratory of Stockholm University for the opportunity to present my work there and for their feedback and suggestions.

## References

- Berrah, A. R. (1998) Évolution artificielle d'une société d'agents de parole: Un modèle pour l'émergence du code phonétique, Thèse de l'Institut National Polytechnique de Grenoble, Spécialité Sciences Cognitives.
- Carré, R, M. Bordeau & J.-P. Tubach (1995) Vowel-vowel production: The distinctive region model (DRM) and vowel harmony, *Phonetica* **52**, pp.205–214
- Chomsky, N. & M. Halle (1968) *The sound pattern of English*, MIT Press, Cambridge, Mass.
- Cooper, F. S, P. C. Delattre, A. M. Liberman, J. M. Borsari & L. J. Gerstman (1976), Some experiments on the perception of synthetic speech sounds, in: D. B. Fry (ed.) *Acoustic phonetics*, Cambridge University Press, pp.258–283
- de Boer, B. (1997a) Generating vowels in a population of agents, in: P. Husbands & I. Harvey (eds.) *Fourth European Conference on Artificial Life*, MIT Press, pp.503–510
- de Boer, B. (1997b) Self organisation in vowel systems through imitation, in: J. Coleman (ed.) *Computational Phonology, Third Meeting of the ACL SIGPHON*, July 12, 1997, pp.19–25
- de Boer, B. (1998) *A realistic model of emergent phonology*, Vrije Universiteit Brussel AI-lab AI-memo98-04.
- Glotin, H. (1995) La Vie artificielle d'une société de robots parlants: émergence et changement du code phonétique. DE Sciences cognitives-Institut National Polytechnique de Grenoble
- Jakobson, R. & M. Halle (1956) *Fundamentals of language*, The Hague: Mouton & Co.
- Ladefoged, P. & I. Maddieson (1996) *The sounds of the world's languages*, Oxford: Blackwell.
- Liljencrants, L. & B. Lindblom (1972) Numerical simulations of vowel quality systems: The role of perceptual contrast, *Language* **48** pp.839–862.
- Lindblom, B. (1992) Phonological units as adaptive emergents of lexical development, in: C. A. Ferguson, L. Menn & C. Stoel-Gammon, *Phonological Development*, pp.131–163
- Maddieson, I. (1984) *Patterns of sounds*, Cambridge University Press.
- Maeda, S. (1989) Compensatory articulation during speech: Evidence from the analysis and synthesis of vocal tract shapes using an articulatory model, in: W. J. Hardcastle and A. Marchal (eds.) *Speech Production and Speech Modelling*, Kluwer, pp.131–149
- Mantakas, M, J. L. Schwartz & P. Escudier (1986) Modèle de prédiction du 'deuxième formant effectif' F2'—application à l'étude de la labialité des voyelles avant du français. In: *Proceedings of the 15th journées d'études sur la parole*. Société Française d'Acoustique, pp.157–161.
- Steels, L. (1996) The spontaneous self-organization of an adaptive language, in: S. Muggleton (ed.) *Machine Intelligence* **15**.
- Steels, L. (1997) The synthetic modelling of language origins, *Evolution of Communication* **1**(1): pp.1–34
- Steels, L. (1998) Synthesising the origins of language and meaning using co-evolution, self-organisation and level formation, in: J. R. Hurford, M. Studdert-Kennedy & C. Knight (eds.) *Approaches to the evolution of language*, Cambridge: Cambridge University Press pp.384–404
- Stevens, K. N. (1972). The quantal nature of speech: Evidence from articulatory-acoustic data. In: E. E. David, Jr. & P. B. Denes (Eds.) *Human communication: a unified view*. New York: McGraw-Hill.
- Vallée, N. (1994) Systèmes vocaliques: de la typologie aux prédictions, Thèse préparée au sein de l'Institut de la Communication Parlée (Grenoble-URAC.N.R.S .no368)