

Language Emergence for Activities

Manas Agarwal

manas@cse.iitk.ac.in

Under the guidance of Dr. Amitabha Mukerjee

April 16, 2010

1 Abstract

Language development is an important aspect of human cognitive processes. Robotic and Artificial Intelligent models have been used to simulate language emergence situations [6]. The aim is to study the process of specification of the universal features of language and emergence of structures (symantic/syntactic) in human language through a model known as *Language Games*. The games consist of short two body interactions which are classified by game playing agents and utterances are produced. Our first aim is to develop large number of such interactions and correctly classify them. The second part consist of language learning and lexicon development through naming games. We see that our models produce consistent results.

2 Language emergence

There are many theories and researches trying to explain the formation of structures in human languages. Some researchers assume that compositional structures emerged from use of regularities found in expressions based on holophrases (single word expressions) [1]. Some other researchers have assumed that the ability to use syntax has evolved as a biological adaptation. Another important line of thinking is that the emergence of language structures is based on exploiting regularities found in the interaction with natural world. This hypothesis has been investigated using multi-agent models. A Language game model have been applied to study language emergence which mainly consisted of object naming and identification [8]. We intend to apply similar model to study emergence of compositional language consisting of objects and simple actions.

3 Language Game

The scientific aim of evolutionary linguistics is to study how modern languages have evolved from a stage prior to language. This is useful, because it provides a complimentary methodology that can help researchers to develop and test hypothesis on language

evolution in the simulated world. Vogt and Steels [7] developed a model which implements a scenario in which the population of one generation transmits their language to the next generation by engaging in language games. This is called a *naming game*. The game is played by two agents: the *speaker*, typically from older generation and a *hearer*, typically from a new generation. Different games can be played. The following basic steps are involved in the game:

1. *Making Contact* : Agents assume their respective roles and come close enough to have a shared context.
2. *Perception*: Each agent categorize the sensory experience of the objects and identify a distinctive feature set that distinguish the present context with any other one.
3. *Encoding*: The speaker identifies one distinctive feature set and encodes this into expression. It means finding smallest set of words that describe the feature set from present lexicons. If no word can be formed, a new word is created and added.
4. *Decoding*: Hearer looks up all the words in its lexicon and reassembles a feature set that cover all the words.
5. *Feedback*: The hearer compares the decoded feature with the distinctive feature set it was expecting. If sets matches, the game is successful and a positive feedback is given by hearer. Otherwise the game ends.

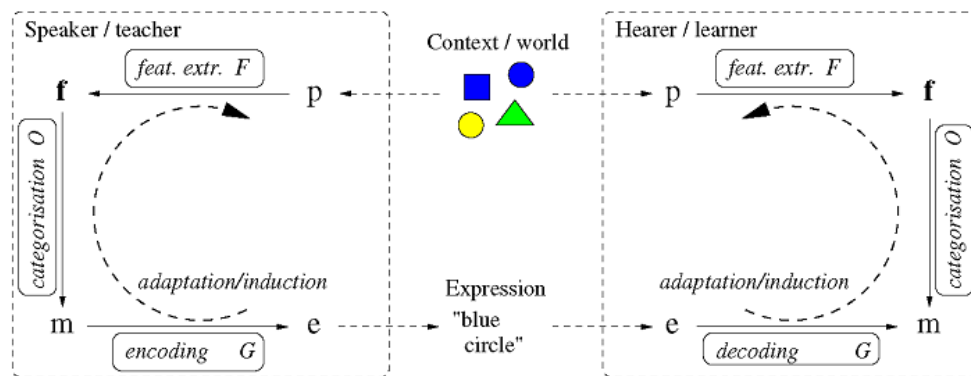


Figure 1: The Naming game processes used by Vogt. On the left are speaker's processes and on the right are hearer's processes. Between them are the feedback processes.

4 The Iterated Learning Model

The ILM is a general framework for modelling cultural transmission of language [2]. The model iterates over generations which population is divided into: adults and learners.

Each generations consists a population of agents which learn language from utterances produced by the previous generations. The learner starts his life as a novice user and acquire language by observing behaviour of adults, who are assumed to have mastered the language. A number of such language games are played. At the end of each iteration, adults leave and are replaced by novel learners. This cycle is repeated. ILM is a simplifid model of population dynamics and cultural transmission of language and is also used in the *naming game*.

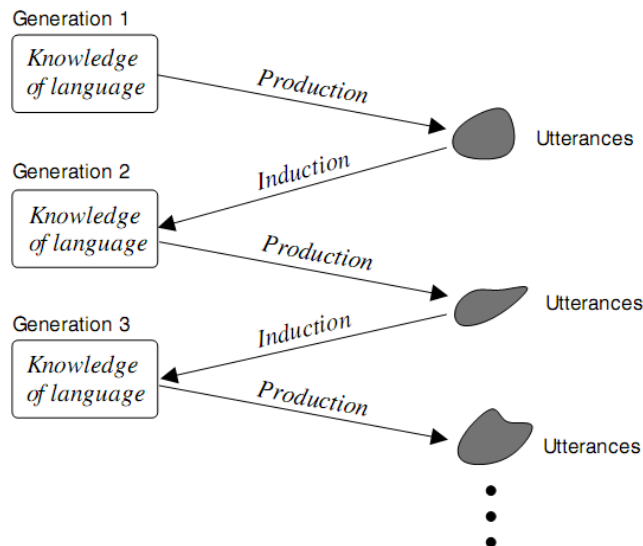


Figure 2: The schematic view of the iterative learning model. The utterances are used by the agents in the next generations to induce the knowledge of language.

5 The New model: Learning for Activity

Notably, Vogt and Steel’s work focussed on learning language structure only for objects, i.e. common nouns. The objects are described by color and shapes. Other works have also studied for proper nouns. But language emergence for activities have largely been untouched. In the formative years of human life, the concepts are largely acquired from dynamic activities rather than static sceneries [4]. Thus, study of formation of verbs during language evolution can provide a useful insight to the actual human language learning.

Our model uses the Vogt model of *naming game* to simulate learning of language composition when their are some agents and an activity is being performed by them. The objective of our game is to find a distinctive category for this context. Categorization means identifying a distinctive feature for the underlying action between the agents. Work on the construction of an activity template from visual inputs is mostly centered

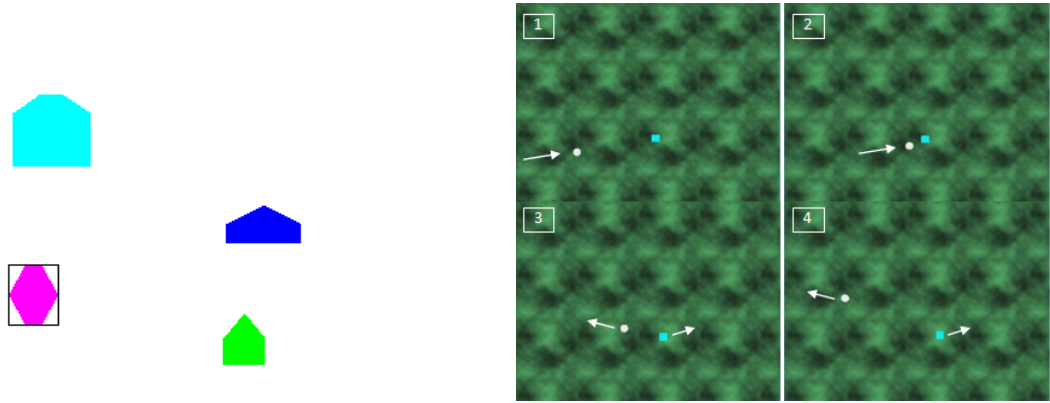


Figure 3: [Left] static scenes used by Vogt, [Right] A two body interaction (bounce) used by the proposed model

on single agent activities and that too require visual priors for recognition. However, Satish, Mukerjee[5] has provided an unsupervised learning technique, which can cluster the activities into basic groups using attentive focus. It claims that semantics of certain actions may be learned prior to language, that is, we can form a crude linguistic description of a scene just by the visual inputs only. These actions were prominently change of positions of agents in a multi-agent setups, and the result were semantics describing the relative motions amongst them. They also incorporated centre of attention feature since they claim that a commentator is more likely to talk about things that are in attentive focus. The results of the work [Fig. 4] show that by restricting the constituents participating in an action using computational model of visual attention, they have achieved better correlation with the human commentary provided.

	C_1	C_2	C_3	C_4	Total	%	TCA
CC	399	6	10	29	444	90	
MA	16	311	5	48	380	82	84
Chase	21	59	149	154	383	79	

HUMANS							
ALGO							
HUMANS							
ALGO							
HUMANS							
ALGO							

Figure 4: Clustering by Merge Neural Gas. The i^{th} row, j^{th} column gives the number of i^{th} action labels in j^{th} cluster. Below given is comparison of human and algorithm labelling.

6 The Model in detail

6.1 Interaction generation and classification

- Our world (The game) consist only of two agents (distinct) in a 2d space. We would like to simulate basic movements that can be categorized later. [3] describe some trajectories for such type of interaction that are shown to be classified as same action universally during learning by infants. These trajectories [Fig. 5] are used for action generations. We are finally using four actions : *Push*, *Hit*, *Harass* and *Bounce*.

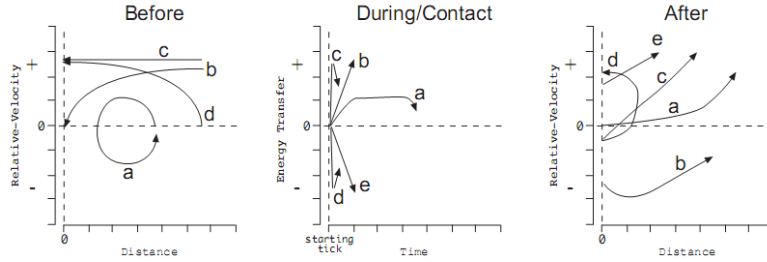


Figure 5: The labeled trajectories characterize the component phases of seven interaction types, as described by the verbs: push, shove, hit, harass, bounce, counter-shove and chase.

- A feature vector set is extracted from each action sequence which is used by Merge Neural Gas Algorithm to classify the interaction. The feature vector set includes the vector $[f_1, f_2]$ for each frame of the video, where $f_1 = (\vec{x}_B - \vec{x}_A) \cdot (\vec{v}_B - \vec{v}_A)$, $f_2 = (\vec{x}_B - \vec{x}_A) \cdot (\vec{v}_B + \vec{v}_A)$.
- Using these feature vectors, we would like to categorize the relationship between different balls into one of the four actions that we are expecting. In the above mentioned work [5], Merge Neural Gas Algorithm is used. Neural Gas Algorithm learns important topological relations in a given set of input vectors in an unsupervised manner by means of simple Hebb-like learning rule. Since input features are arriving from temporally connected data(from different time frames), a *context* information is also added to pass on temporal information. This changed algorithm is the Merge Neural Gas Algorithm.

6.2 Language Game

- Two agents are chosen from the population, one becomes Speaker S and another a hearer H .
- S plays a discrimination game to find a distinct category for an action sequence using Merge Neural Gas Algorithm.

- S produces an utterance $u = w_i$ where word w_i is from lexical element $l_i \in L_S$ for which $m_i = c_d$ and $\sigma_i \geq \sigma_j$ for all other elements $l_j \in L_S$ for which $m_j = c_d$. If no such element is found, a new element $l = \langle w, c_d, 0.01 \rangle$ is added to the lexicon L_S , where w is the new created word. This happens with a certain *word creation probability*.
- When H receives utterance u , it plays its own discrimination game to produce a category D . If $D \neq Null$, H tries to interpret the utterance by searching an element $l_i \in L_H$ with $w_j = u$ and $m_i \in D$ and $\sigma_i \geq \sigma_j$ for all other elements $l_j \in L_H$ for which $m_j = D$.
- If H finds a lexical element to cover u , the game is successful and both agents increase the association score of the used element l_i by $\sigma_i = \eta \cdot \sigma_i + 1 - \eta$ where $\eta \in \langle 0, 1 \rangle$ is a learning parameter (default value is 0.9). Also, association score of competing elements l_j are laterally inhibited by $\sigma_j = \eta \cdot \sigma_j$. If the game fails, H adopts u and adds the element $l = \langle u, c_d, 0.01 \rangle$ to L_H . Also, S lowers the association score of the used element l_i by $\sigma_i = \eta \cdot \sigma_i$.
- *Iterated Learning Model* A population of ILM contains a group of N adult agents and a group of N learner agents. Iterations are done in two steps:
 1. K language games are played.
 2. The adults are removed from the population and replaced by the learners. N new agents are placed in the learners group.

By default, *Speakers* are always selected from adult population and *Hearers* are always selected from the learner group, except for the first iteration where each agent is equally likely to be selected.

7 Results

7.1 Classification of Interaction videos

There were four types of videos: Push, Hit, Bounce, Harass
The percentage accuracy of classification is as shown:

	C1	C2	C3	C4	Total	% Accuracy
Bounce	1888	42	51	19	2000	94.4%
Push	20	1933	15	32	2000	96.6%
Hit	91	30	1841	32	2000	92.1%
Harass	12	27	34	1927	2000	96.4%

Figure 6: Clustering results

7.2 Language Game

Parameters: number of Iterations=3 , number of games= 1000

The table shows the associative scores of the total lexicon generated.

Words	Verb	Score of Last Speaker	Score of Last Hearer
"geda"	Bounce	0.98	0.96
"cabe"	Bounce	0.42	0.23
"defe"	Bounce	0.13	0.11
"dega"	Hit	0.96	0.92
"feged"	Hit	0.28	0.16
"dace"	Push	0.96	0.89
"gabe"	Push	0.34	0.26
"daceg"	Harass	0.92	0.92
"cebag"	Harass	0.23	0.17

Figure 7: Lexicons and their scores

The plot shows a statistic called *Discriminative Success*. It is calculated over past 100 games and is the average of the number of successful games during this period. It can be noted that at the end of the iteration, the probability that a game would be successful tends to one.

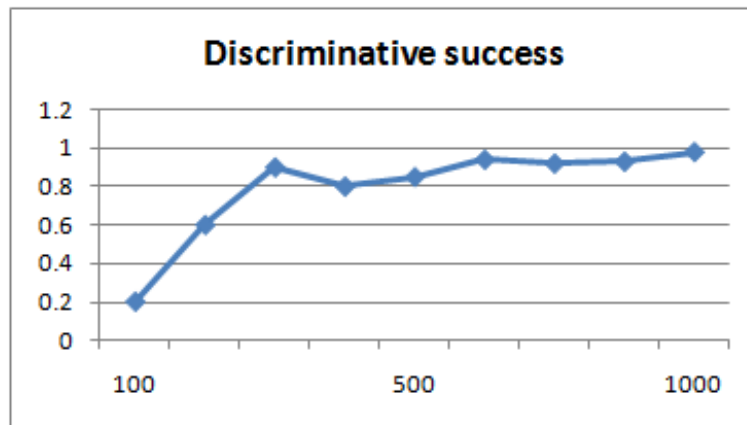


Figure 8: Discriminative Success plot

8 Conclusion

This work presents a novel model and its implementation that can be used by investigators who are interested in grounded lexicon emergence and formation for verbs. The results are as expected. The work is limited by the type of interaction involved. Future work may involve incorporation of syntax emergence and combination of object and verb learning.

References

- [1] BRIGHTON, H. Compositional syntax from cultural transmission. *Artificial Life* 8, 1 (2002), 25–54.
- [2] BRIGHTON, H., KIRBY, S., AND SMITH, K. Situated cognition and the role of multi-agent models in explaining language structure. In *Adaptive Agents and Multi-Agent Systems: Adaptation and Multi-Agent Learning* (2003), D. Kudenko, E. Alonso, and D. Kazakov, Eds., Springer, pp. 88–109.
- [3] COHEN, P. R., MORRISON, C. T., AND CANNON, E. Maps for verbs: The relation between interaction dynamics and verb use. In *Proc. IJCAI-05* (2005), pp. 1022–1027.
- [4] MANDLER, J. Actions organize the infant’s world. In *Action Meets Word: How Children Learn Verbs* (Hirsh-Pasek and Golinkoff, 2006). Oxford University Press (2006), 2006, p. 111133.
- [5] SATISH, G., AND MUKERJEE, A. Acquiring linguistic argument structure from multimodal input using attentive focus. In *Development and Learning, 2008. ICDL 2008. 7th IEEE International Conference on* (2008), pp. 43–48.
- [6] STEELS, L. A case study in the behaviour-oriented design of autonomous agents. In *3rd simulation of adaptive behaviour Conference, The MIT Press* (1994).
- [7] STEELS, L., AND VOGT, P. Grounding adaptive language games in robotic agents. In *Proceedings of the 4th European Conference on Artificial Life* (Cambridge, MA, 1997), I. Harvey and P. Husbands, Eds., The MIT Press, pp. 474–482.
- [8] VOGT, P. The emergence of compositional structures in perceptually grounded language games. *Artificial Intelligence* 167 (2005), 206–242.