

Categorial Grammar Induction as a Compositionality Measure for Understanding the Structure of Emergent Languages

Anonymous ACL submission

Abstract

This paper proposes a method for investigating the syntactic structure of emergent languages using categorial grammar induction. Although the structural property of emergent languages is an important topic, little has been done on syntax and its relation to semantics. Inspired by previous work on CCG induction for natural languages, we propose to induce categorial grammars from the sentence-meaning pairs of emergent languages. Since an emergent language born in a common environment called signaling game is represented as pairs of a message and a meaning, it is straightforward to extract sentence-meaning pairs to feed to categorial grammar induction. We also propose two compositionality measures that are based on the information obtained from induced grammars. Our experimental results reveal that our measures can recognize compositionality. While correlating with existing measure TopSim, our measures can gain more insights on the compositional structure of emergent languages from induced grammars.

1 Introduction

Communication among artificial agents born in an environment is called *emergent communication* and its protocols are *emergent languages* (Lazaridou and Baroni, 2020). Major motivations in this area are (1) to develop interactive AI (Foerster et al., 2016; Mordatch and Abbeel, 2018; Lazaridou et al., 2020), (2) to study language evolution (Kirby, 2001; Graesser et al., 2019; Dagan et al., 2021), and (3) to understand emergent languages or compare them with humans’ (Kottur et al., 2017; Chaabouni et al., 2019a; Kharitonov et al., 2020). (1) and (2) are important from the engineering or scientific points of view. In fact, (3) is fundamental since the first two are not achievable without recognizing and filling the gap between emergent and human languages. Despite its importance, few methods have been established to evaluate the struc-

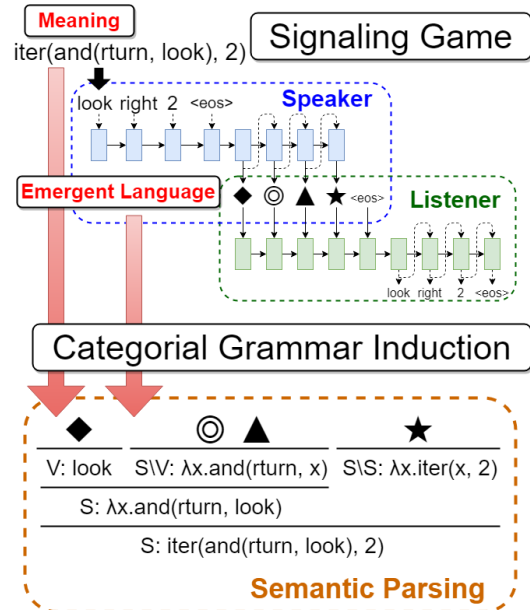


Figure 1: Illustration of a signaling game and categorial grammar induction (CGI). We first generate message-meaning pairs in the game, and then feed them to CGI.

ture of emergent languages with respect to *syntax* and *semantics*. Previous work frequently employs a *signaling game* (Lewis, 1969) or its variant, where agents are either a function from a *meaning space* to a *message space* or its inverse. The problem is that little has been analyzed on how syntax combines messages to yield semantics or meanings. Such a structural property is known as *compositionality*.

To analyze the syntax of emergent languages, we focus on categorial grammar induction (CGI, e.g., Zettlemoyer and Collins, 2005) and propose to apply it to emergent languages. Figure 1 illustrates the relationship between a signaling game and CGI. Since CGI derives an explicit lexicon and a semantic parser given sentence-meaning pairs, it is suitable for the syntactic analysis of a language emerging as message-meaning pairs in a signaling game. We also propose compositionality measures

061 built on the F1-score for unseen data and the lexi- 109
062 con size of CGI parsers. It is based on the intuition 110
063 that a compositional language is expected to be 111
064 generalized and described by a minimal lexicon. 112

065 Compositionality measures for emergent lan- 113
066 guages have been proposed, such as topographic 114
067 similarity (TopSim, Brighton and Kirby, 2006), 115
068 tree reconstruction error (TRE, Andreas, 2019), po- 116
069 sitional disentanglement (PosDis, Chaabouni et al., 117
070 2020), and bag-of-symbols disentanglement (Bos- 118
071 Dis, Chaabouni et al., 2020). We choose TopSim 119
072 and TRE to compare with ours, since TopSim is 120
073 most popular (e.g., Lazaridou et al., 2018) and TRE 121
074 is similar to ours in the sense that it assumes struc- 122
075 tured meaning representations. Note that they do 123
076 not consider the structure between a message and a 124
077 meaning space, whereas our approach is aware of 125
078 it with an explicit lexicon and a parser. 126

079 Pioneering and suggestive work by van der Wal 127
080 et al. (2020) on the syntax of emergent languages 128
081 proposes to apply unsupervised grammar induction 129
082 (UGI) originally developed for natural languages: 130
083 CCL (Seginer, 2007) and DIORA (Drozdo et al., 131
084 2019). UGI is reasonable if neither gold derivations 132
085 nor meanings are available¹. Note that UGI esti- 133
086 mates the structure of emergent languages given 134
087 only messages, whereas ours is intended to derive 135
088 not only the structure but also the systematic com- 136
089 position of messages to meanings given message- 137
090 meaning pairs. 138

091 Our contributions are (1) to propose to apply 139
092 categorial grammar induction (CGI) to emergent 140
093 languages for understanding their structure, (2) to 141
094 propose two CGI-based compositionality measures 142
095 that are more syntax-aware than existing composi- 143
096 tionality measures, and (3) to show they can indeed 144
097 measure compositionality. 145

098 2 Signaling Game in General 146

099 Most studies on emergent communication employ 147
100 *Lewis signaling game* (Lewis, 1969) or its vari- 148
101 ant as an environment for agents to communicate. 149
102 A signaling game contains a tuple (I, M, S, L) , 150
103 where I is an *input space*, M is a *message space*, a 151
104 mapping $S : I \rightarrow M$ is a *speaker*, and a map- 152
105 ping $L : M \rightarrow I$ is a *listener*. The goal is 153
106 $i = L(S(i))$ for a sampled input $i \in I$. Agents 154
107 S, L are trained to achieve the goal given I, M . 155
108 On the other hand, in a variant called *referential* 156

¹For example, if agents describe image data (e.g., Lazaridou et al., 2018), the meaning representations are unclear. 157

game or *discrimination game*, a listener is defined 109
as $L : M \times \mathcal{P}(I) \rightarrow I$, where $\mathcal{P}(I)$ is the power 110
set of I . The goal is to distinguish i from other dis- 111
tractors: $i = L(S(i), C)$ for candidates $C \in \mathcal{P}(I)$ 112
s.t. $i \in C$. An input space is typically a set of 113
image data (Havrylov and Titov, 2017; Lazaridou 114
et al., 2018; Bouchacourt and Baroni, 2018), se- 115
quential data (Li et al., 2020; Słowik et al., 2021), 116
or attribute-value objects (Li and Bowling, 2019; 117
Chaabouni et al., 2020; Ren et al., 2020). Besides, 118
a message space is a set of discrete sequences in 119
most studies. 120

Agent Architecture Each agent is typically rep- 121
resented as a neural network, in particular, an 122
encoder-decoder model. The speaker decoder and 123
listener encoder are often recurrent neural networks. 124
The speaker encoder can be a convolutional neural 125
network, recurrent neural network, or perceptron, 126
according to I . The listener decoder can be ei- 127
ther the same as the speaker encoder or a classifier, 128
depending on the goal. 129

Optimization Methods The speaker-listener pair 130
is trained in an End-to-End manner, regarded as a 131
single neural network. Previous work uses REIN- 132
FORCE (Williams, 1992) and/or Gumbel-Softmax 133
trick (Jang et al., 2017; Maddison et al., 2017), 134
since the standard backpropagation is not applica- 135
ble to discrete messages. 136

137 2.1 Existing Compositionality Measures 137

138 Compositionality is popular among those who are 139
interested in the structural similarity between emer- 140
gent and human languages. In the experiments, 141
we compare our measures with TopSim (Brighton 142
and Kirby, 2006) and TRE (Andreas, 2019). Let 143
 (I, M, S, L) be a signaling game defined above. 144

TopSim Let d_I, d_M be distance functions in I 145
and M . TopSim is defined as Spearman correla- 146
tion between $d_I(x, y)$ and $d_M(S(x), S(y))$ for all 147
 $(x, y) \in I \times I$. This score requires only d_I, d_M as 148
structural information for I, M . 149

TRE The intuition of TRE is that if an emergent 150
language is compositional, it should be approxi- 151
mated by another explicitly compositional function 152
 $f : I \rightarrow M$. Note that each $i \in I$ has to be a *binary* 153
tree t in which a node n is binary node denoted as 154
 $n = (n', n'')$, unary node denoted as $n = (n')$, or a 155
leaf node denoted as $n = l$. Besides, each $m \in M$ 156
has to be a sequence of a fixed length k over a fi- 157
nite alphabet A . The calculation of TRE involves

a distance δ and a composition κ_η with a trainable parameter $\eta = (V, W, \{E_l\}_l)$. δ is defined as L1 distance between $k \times |A|$ matrices. κ_η is defined as a mapping from a binary tree $t \in I$ to a $k \times |A|$ matrix:

$$\begin{aligned}\kappa_\eta(l) &= E_l \\ \kappa_\eta((n)) &= \kappa_\eta(n) \\ \kappa_\eta((n, n')) &= V\kappa_\eta(n) + W\kappa_\eta(n')\end{aligned}$$

where E_l is a $k \times |A|$ matrix for each leaf node l , and V, W are $k \times k$ matrices. Define $\text{one_hot}(m)$ as a $k \times |A|$ matrix, the r -th row of which is the one-hot vector of the r -th symbol in $m \in M$. Then, TRE is computed with stochastic gradient descent as follows:

$$\text{TRE} = \min_{\eta} \frac{1}{|I|} \sum_{i \in I} \delta(\kappa_\eta(i), \text{one_hot}(S(i)))$$

Note that the lower TRE is, the higher compositionality is judged. TRE is similar to ours in the sense that inputs are assumed to be tree-structured.

3 Categorial Grammar Induction

In this section, we introduce categorial grammar (CG) and review its induction (CGI) for natural languages. CGI is also eligible for the analysis of emergent languages in signaling games, as it derives a lexicon and a parser from message-meaning pairs. Although previous work is on combinatory categorial grammar (CCG), we restrict it to CG².

3.1 Categorial Grammar

The formalism for our semantic parsing is *categorial grammar* (CG, Steedman, 1996, 2000). Context-free grammars are described largely with rules, whereas CGs are described largely with *lexical entries* and their rules are simple. A lexical entry $w \vdash X : \psi$ is a triple of a word w , a category X (defined below), and a logical form ψ . Consider the following example pair of a message and its logical form:

“look left 1”
iter(and(1turn, look), 1)

²One might wonder why we do not use CCG. This is because the input spaces for our signaling games are described by context-free grammars, whose expressive power is known to be equal to that of CG. Nevertheless, it is interesting to speculate whether emergent languages can have complex rules like composition or type-raising. It is left for future work.

Their lexical entries can be described as follows:

look \vdash V : look 196
left \vdash S \ V : $\lambda x.$ and(1turn, x) 197
1 \vdash S \ S : $\lambda x.$ iter(x, 1) 198

Symbols like V, S \ V, and S \ S represent syntactic types or *categories*. A category is either an atomic category of the form N, V, or S, or a complex category of the form X/Y or $X \setminus Y$ where X, Y are categories. The atomic categories N, V, and S stand for the linguistic notions of noun, intransitive verb, and sentence respectively³.

In addition, CGs have *application rules* to describe the way to combine adjacent categories.

Application rules (with semantics):

$X/Y : f \quad Y : a \Rightarrow X : f(a) \quad (>)$ 211
 $Y : a \quad X \setminus Y : f \Rightarrow X : f(a) \quad (<)$ 212

where X, Y are categories. The first rule named “>” is called the *forward application rule*, while the second rule named “<” is called the *backward application rule*. Rule > (resp. <) means that a predicate f of category X/Y (resp. $X \setminus Y$) can take an argument a of category Y to yield $f(a)$ of category X .

With the lexical entries and the application rules, we can construct a derivation tree of “look left 1” as follows:

look	left	1	
V	S \ V	S \ S	
: look : $\lambda x.$ and(1turn, x)	: $\lambda x.$ iter(x, 1)		
S : and(1turn, look)			
S : iter(and(1turn, look), 1)			

3.2 Log-linear Probabilistic CGs

Given a lexicon Λ , a set of lexical entries, there might be multiple derivations for each message. Following previous work on CG induction (e.g., Zettlemoyer and Collins, 2005), we choose the most likely derivation by using a log-linear model, which contains a feature vector ϕ and a parameter vector θ . Given a message m , the joint probability of a logical form ψ and a derivation τ is defined as:

$$P(\tau, \psi \mid m; \theta, \Lambda) = \frac{e^{\theta \cdot \phi(m, \tau, \psi)}}{\sum_{(\tau', \psi')} e^{\theta \cdot \phi(m, \tau', \psi')}}.$$

³The category of intransitive verbs is usually S/N (S/NP) or S \ N (S \ NP), but we regard V as a atomic category. This is because the languages and logical forms we define in Section 5.1 take an imperative form without any subject.

Then, the parsing problem is to find the most likely logical form $\hat{\psi}$ given m :

$$\begin{aligned}\hat{\psi} &= \arg \max_{\psi} p(\psi \mid m; \theta, \Lambda) \\ &= \arg \max_{\psi} \sum_{\tau} P(\tau, \psi \mid m; \theta, \Lambda).\end{aligned}$$

3.3 CG Induction Algorithm

Algorithm 1 Common Structure of CG Induction

Input: A dataset $\mathcal{E} = \{(m^j, \psi^j)\}_{j=1}^N$ of message-meaning pairs, a seed lexicon Λ_{seed} , the number of iterations T , and a learning rate γ .

Output: Lexicon Λ and parameter vector θ

- 1: $\Lambda_0 \leftarrow \text{INITLEX}(\mathcal{E}, \Lambda_{\text{seed}})$
 - 2: $\theta_0 \leftarrow \text{INITPARAM}(\mathcal{E}, \Lambda_{\text{seed}})$
 - ▷ Step 0: Initialize lexicon and parameter
 - 3: **for** $t \in \{1, \dots, T\}$ **do**
 - 4: $\Lambda_t^+ \leftarrow \text{UPDATELEX}(\mathcal{E}, \theta_{t-1}, \Lambda_{t-1}, \Lambda_0)$
 - ▷ Step 1: Update Lexicon
 - 5: $\theta_t \leftarrow \text{UPDATEPARAM}(\mathcal{E}, \theta_{t-1}, \Lambda_t^+, \gamma)$
 - ▷ Step 2: Update Parameter
 - 6: $\Lambda_t \leftarrow \text{PRUNEXLEX}(\mathcal{E}, \theta_{t-1}, \Lambda_t^+)$
 - ▷ Step 3: Prune Lexicon (optional)
 - 7: **end for**
 - 8: **return** Λ_T and θ_T
-

Several CG induction (CGI) algorithms have been proposed. Algorithm 1 shows their common structure as a pseudo code. Generally, the inputs to CGI are a training data $\mathcal{E} = \{(m^j, \psi^j)\}_{j=1}^N$ of message-meaning pairs, a seed lexicon Λ_{seed} , the number of iterations T , and a learning rate γ , while the outputs are a lexicon Λ and a parameter θ . CGI involves four procedures: (1) lexicon and parameter initialization (INITLEX, INITPARAM) that helps learning in early iterations, (2) lexicon update (UPDATELEX) that introduces a new potential lexicon, (3) parameter update (UPDATEPARAM) with gradient descent, and optionally (4) lexicon pruning (PRUNEXLEX) that discards a lexicon no longer in use. ZC05 (Zettlemoyer and Collins, 2005) is the first paper formalizing CGI. ZC07 (Zettlemoyer and Collins, 2007) is its improved version. In ZC05/07, INITLEX is simply $\Lambda_0 = \Lambda_{\text{seed}}$ and UPDATELEX relies on hand-crafted templates to add a new lexicon. KZGS10/11 (Kwiatkowski et al., 2010, 2011) modified UPDATELEX so that it can create a new lexicon by automatically merging and splitting the existing entries in use. In KZGS10/11, INITLEX returns \mathcal{E} themselves with category S in addition to Λ_{seed} :

$$\Lambda_0 \leftarrow \Lambda_{\text{seed}} \cup \{m^j \vdash \mathcal{S} : \psi^j \mid j = 1, \dots, N\}$$

Then, the lexical entries are split or merged during the iteration, seeking an appropriate segmentation. A problem in KZGS10/11 is that the lexicon size increases monotonically over iterations. ADP14 (Artzi et al., 2014) addressed this issue by adding a lexicon pruning process (PRUNEXLEX), which discards the lexical entries no longer in use⁴.

4 CGI as a Compositionality Measure

We propose two compositionality measures CGF and CGL, which are based on an induced categorical grammar. Let $\mathcal{E}_{\text{train}}, \mathcal{E}_{\text{test}}$ be a training and test data for CGI. We train a log-linear model with $\mathcal{E}_{\text{train}}$ to derive a lexicon Λ and a parameter θ and test it with $\mathcal{E}_{\text{test}}$ to calculate the F1-score for semantic parsing:

$$\begin{aligned}\text{F1-score} &= 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}} \\ \text{precision} &= \frac{\# \text{ correctly parsed}}{\# \text{ parsed}} \\ \text{recall} &= \frac{\# \text{ correctly parsed}}{|\mathcal{E}_{\text{test}}|}\end{aligned}$$

following previous work (Zettlemoyer and Collins, 2005)⁵. Then, CGF and CGL are defined as:

$$\text{CGF} = \text{F1-score}, \quad \text{CGL} = |\Lambda|$$

Note that the higher CGF (resp. lower CGL) is, the more compositional a language is judged, since a compositional language should be generalized and described by a minimal lexicon.

4.1 Difference from Existing Measures

Although existing compositionality measures such as TopSim and TRE are also mappings from message-meaning pairs to a real number, neither they clarify the structure of a message space M nor they derive any compositional function from M to an input space I .

Remember that TopSim only involves distance functions d_I, d_M , the choice of which is left to humans, and it does not clarify the structure of M . On the other hand, our approach can derive the structure of M by deriving a lexicon. TRE induces a composition $\kappa_{\eta} : I \rightarrow M$, but not the inverse. As Andreas (2019) is aware, it causes a language with identical messages for all meanings to be judged

⁴ADP14 also has improvements in UPDATELEX, but we do not go into them in this paper.

⁵If Λ does not have sufficient lexical entries, the model fails to parse messages regardless of correctness.

compositional, contrary to our intuition. Again, ours would not regard it as compositional since a CGI parser is a function $M \rightarrow I$.

Therefore, what differentiates us from the existing measures is that our approach can derive an explicit lexicon and a semantic parser, whereas the existing measures cannot ⁶.

5 Experimental Setup

This section introduces a signaling game, optimization method, CGI algorithm, and evaluation metrics specific to our experiments. The overall experimental procedure is as follows:

1. Split an input space I in half: $\mathcal{D}_{\text{train}}, \mathcal{D}_{\text{test}}$.
2. Train a speaker S and a listener L with $\mathcal{D}_{\text{train}}$. Validate and test them with $\mathcal{D}_{\text{test}}$.
3. Given trained S , make datasets for CGI by pairing each message with its logical form:

$$\mathcal{E}_x = \{(S(i), \langle i \rangle) \mid i \in \mathcal{D}_x\}$$

where $x \in \{\text{train}, \text{test}\}$ and $\langle i \rangle$ is the logical form of i , to which CGI is applicable.

4. Train a CG parser with $\mathcal{E}_{\text{train}}$, test it with $\mathcal{E}_{\text{test}}$, and calculate CGF with $\mathcal{E}_{\text{test}}$ and CGF with a derived lexicon Λ .
5. Calculate TopSim and TRE with $\mathcal{D}_{\text{train}}$.

5.1 Input Space for Signaling Game

We define two input spaces for our signaling game: *Lang-attval* and *Lang-conj* ⁷. *Lang-attval* is the same as attribute-value inputs in previous work (e.g., Kottur et al., 2017), while *Lang-conj* is more complex. Moreover, we define logical forms for each input.

Lang-attval *Lang-attval* is defined as the set of sequences derived from the following context-free grammar with a start symbol S :

$$\begin{aligned} S &\rightarrow V' R \\ V' &\rightarrow V D \\ V &\rightarrow \text{look} \mid \text{jump} \mid \text{walk} \mid \text{run} \\ D &\rightarrow \text{left} \mid \text{right} \mid \text{up} \mid \text{down} \\ R &\rightarrow 1 \mid 2 \mid 3 \mid 4 \end{aligned}$$

⁶TopSim and TRE are still reasonable if our purpose is to distinguish partially (but insufficiently) compositional languages from the ones not compositional at all.

⁷They are inspired by the commands of Chaabouni et al. (2019b) or SCAN (Lake and Baroni, 2018).

Lang-attval is regarded as attribute-value objects (Kottur et al., 2017; Andreas, 2019; Li and Bowling, 2019; Ren et al., 2020). In our case, attributes are *verb*, *direction*, and *repetition*, each of which has 4 values (e.g., *look*, *jump*, *walk*, and *run* for *verb*).

Lang-conj Let S'' be a start symbol. Then, *Lang-conj* is the set of sequences derived from the above context-free grammar in addition to the following rules:

$$\begin{aligned} S'' &\rightarrow S \mid S S' \\ S' &\rightarrow \text{and } S \end{aligned}$$

Each element in *Lang-conj* is either an element in *Lang-attval* or a conjunction of two elements in *Lang-attval*.

Logical Form We define the logical form of each element in *Lang-attval/conj*, to which CGI is simply applicable. We temporarily denote elements *parenthetically* to clarify their derivation trees (e.g., “ $S(V'(V(\text{jump}), D(\text{left})), R(2))$ ” for “jump left 2”). Then, the logical form $\langle i \rangle$ of a derivation i is defined inductively as follows:

$$\begin{aligned} \langle S''(S(x)) \rangle &= \langle S(x) \rangle \\ \langle S''(S(x), S'(y)) \rangle &= \text{and}(\langle S(x) \rangle, \langle S'(y) \rangle) \\ \langle S'(\text{and}, S(x)) \rangle &= \langle S(x) \rangle \\ \langle S(V'(x), R(y)) \rangle &= \text{iter}(\langle V'(x) \rangle, \langle R(y) \rangle) \\ \langle V'(V(x), D(y)) \rangle &= \text{and}(\langle D(y) \rangle, \langle V(x) \rangle) \\ \langle X(x) \rangle &= \langle x \rangle \quad (X \in \{V, D, R\}), \end{aligned}$$

and for terminal symbols, $\langle \text{look} \rangle = \text{look}$, $\langle \text{left} \rangle = \text{lturn}$, $\langle 1 \rangle = 1$, and so forth.

Examples Here are some examples:

$$\begin{aligned} i &= \text{“jump left 2”} \in \text{Lang-attval} \cap \text{Lang-conj} \\ \langle i \rangle &= \text{iter}(\text{and}(\text{lturn}, \text{jump}), 2). \end{aligned}$$

Also,

$$\begin{aligned} i' &= \text{“jump left 2 and walk up 3”} \in \text{Lang-conj} \\ \langle i' \rangle &= \text{and}(\text{iter}(\text{and}(\text{lturn}, \text{jump}), 2), \\ &\quad \text{iter}(\text{and}(\text{uturn}, \text{walk}), 3)). \end{aligned}$$

5.2 Signaling Game for Sequential Data

Agent architectures and game procedure have to be adapted to the sequential inputs defined above. Hence, our signaling game takes a sequence-to-sequence procedure.

Architecture Speaker and listener agents are represented as a seq2seq model based on single-layer LSTMs (Hochreiter and Schmidhuber, 1997) with standard attention mechanisms (Dong and Lapata, 2016), similarly to Chaabouni et al. (2019b).

Game Procedure A sequential signaling game consists of a tuple $(I, A, k, \text{eos}, S, L)$, where I is an input space, A is a finite alphabet s.t. $\text{eos} \notin A$, k is a message length, and eos is a special symbol for end-of-sentence. A message space M is defined as the set of sequences of length k over A , $S : I \rightarrow M$ is a speaker, and $L : M \rightarrow I$ is a listener. Note that $x + \text{eos}$ denotes a sequence x attached with eos . The goal of the game is to minimize

$$\Delta(i + \text{eos}, L(S(i + \text{eos}) + \text{eos}))$$

for a uniformly sampled $i \in I$, where Δ is the humming distance.

5.3 Optimization for Agents

As Δ is indifferentiable, we use REINFORCE (Williams, 1992), which gives the following differentiable loss:

$$\begin{aligned} & \mathbb{E}[\{\Delta(i + \text{eos}, o) - b\} \log P_S(m|i + \text{eos})] \\ & + \mathbb{E}[\{\Delta(i + \text{eos}, o) - b\} \log P_L(o|m + \text{eos})] \\ & + \mathbb{E}[\lambda_S \mathcal{H}(P_S) + \lambda_L \mathcal{H}(P_L)] \end{aligned}$$

where P_S (resp. P_L) is the output distribution of speaker (resp. listener) over a message m (resp. output o) given an input i (resp. message m), b is a mean baseline, \mathcal{H} denotes entropy, and λ_S, λ_L are nonnegative hyper-parameters. The last term is an entropy regularizer (Williams and Peng, 1991).

5.4 CGI for Emergent Languages

We apply CGI to emergent languages. As there is no prior knowledge on them, CGI should avoid ad hoc methods, considering the following:

- (1) *Features in a log-linear model have to be as simple as possible.*
- (2) *Lexical entries have to be generated automatically without any manual templates.*
- (3) *Lexicon size has to be minimal; otherwise it is hard to interpret results, e.g., to measure compositionality with CGL.*

There is no existing method satisfying all of them simultaneously. We combine three methods. For

(1), we follow ZC05 (Zettlemoyer and Collins, 2005): each feature is the count of times that each lexical entry is used in a derivation. However, ZC05 generates lexical entries with manual templates, contrary to (2). Instead, we follow KZGS10 (Kwiatkowski et al., 2010) that creates a new lexicon by merging and splitting the existing entries in use. The problem in KZGS10 is that the lexicon size increases monotonically during iterations, which is against (3). Thus, we follow ADP14 (Artzi et al., 2014) to discard the entries no longer in use. Other modifications are detailed in Appendix A.

5.5 Other Languages for Comparison

To evaluate the effectiveness of our measures, we need more and less compositional languages as well as emergent languages to apply CGI. To this end, we use Lang-attval/conj and AdjSwap- x ($x \in \{1, 2\}$). AdjSwap- x is made by applying x -times random adjacent swaps to each message in emergent languages. As Lang-attval and Lang-conj are fully compositional by definition, they should be judged more compositional than emergent languages. On the other hand, AdjSwap- x should be judged less compositional. van der Wal et al. (2020) adopted three languages for the same purpose: fully-structured, random, and shuffled emergent languages. The fully-structured corresponds to Lang-attval/conj in our case. We use AdjSwap- x as instances of less-compositional languages rather than random and shuffled emergent languages. This is because preliminary experiments revealed that CGI totally fails for these languages (see Appendix C). While this is an expected behavior, we additionally employ AdjSwap- x as a language supposed to be more compositional than random and shuffled emergent languages, for obtaining more insights.

5.6 Evaluation Metrics for Compositionality

We use CGF/L as well as TopSim and TRE. When clarifying the target language, we write the metrics as (measure)-(language), e.g., TopSim-Emergent, CGF-AdjSwap-1, and CGL-Lang-attval.

6 Experiments

We show the experimental results in this section. Let $(I, A, k, \text{eos}, S, L)$ be a sequential signaling game as defined in Section 5.2.

For (hyper-)parameter settings, see Appendix B.

	1,1,1	16,13	25,1,1
	S/S	V	S\V
	: $\lambda x.\text{and}(x, \text{iter}(\text{and}(\text{rturn}, \text{walk}), 2))$: run : $\lambda x.\text{iter}(\text{and}(\text{rturn}, x), 3)$		
	S: $\text{iter}(\text{and}(\text{rturn}, \text{run}), 3)$		
	S: $\text{and}(\text{iter}(\text{and}(\text{rturn}, \text{run}), 3), \text{iter}(\text{and}(\text{rturn}, \text{walk}), 2))$		

Figure 2: Example correct derivation tree of a message 1, 1, 1, 16, 13, 25, 1, 1 when $(I, k, |A|) = (\text{Lang-conj}, 8, 31)$.

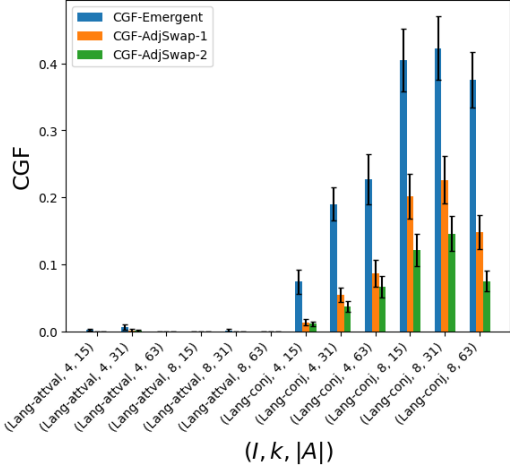


Figure 3: CGF plotted under various $(I, k, |A|)$. The error bars represent one standard error of mean.

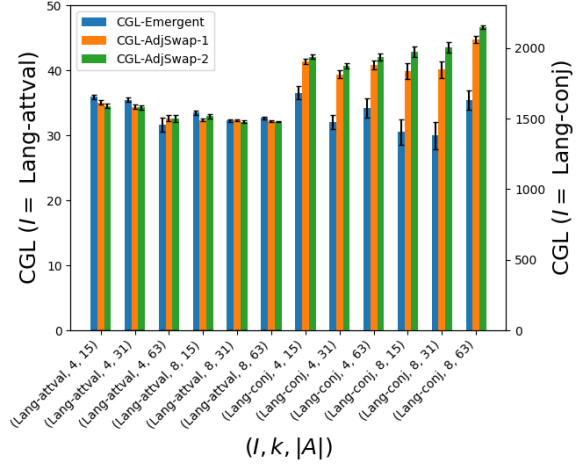


Figure 4: CGL plotted under various $(I, k, |A|)$. The error bars represent one standard error of mean.

6.1 Compositionality of Emergent Languages

We investigate whether CGF/L works as a measure of compositionality. If CGF works, the following inequality should hold: $\text{CGF-Lang-attval/conj} > \text{CGF-Emergent} > \text{CGF-AdjSwap-1} > \text{CGF-AdjSwap-2}$. Likewise, if CGL works, $\text{CGL-Lang-attval/conj} < \text{CGL-Emergent} < \text{CGL-AdjSwap-1} < \text{CGL-AdjSwap-2}$. First, we report that CGF-Lang-attval is 0.984 (± 0.0463), CGL-Lang-attval is 12.3 (± 0.852), CGF-Lang-conj is 0.868 (± 0.1173), and CGL-Lang-conj is 23.8 (± 17.59), where (\pm) denotes a standard error of mean⁸. For the rest, Figure 3 (resp. Figure 4) shows CGF (resp. CGL) under various $(I, k, |A|)$.

For $I = \text{Lang-attval}$, Figure 3 shows surprisingly that CGI fails: CGF-Emergent is near or equal to 0. Besides, CGL-Emergent and CGL-AdjSwap- x in Figure 4 do not show clear differences. Hence, neither CGF nor CGL does not recognize the compositionality of emergent languages. CGF is almost 0 (Figure 3) and CGL concentrates around the size of training data 32 (Figure 4), which means the models overfit the training data. There are two possible reasons for it: emergent languages

⁸We train models 32 times for Lang-attval and Lang-conj respectively.

are not compositional or the training data for CGI is insufficient. We suppose the former is true since CGF-Lang-attval is near perfect (0.984) and CGL-Lang-attval is almost minimal (12.3) with the same size of training data.

For $I = \text{Lang-conj}$, Figure 3 reveals that CGF exactly shows the order of compositionality as expected: $\text{CGF-Lang-conj} > \text{CGF-Emergent} > \text{CGF-AdjSwap-1} > \text{CGF-AdjSwap-2}$. Likewise, CGL in Figure 4 shows the expected order: $\text{CGL-Lang-conj} < \text{CGL-Emergent} < \text{CGL-AdjSwap-1} < \text{CGL-AdjSwap-2}$. Hence, CGF and CGL recognize the compositionality of emergent languages. Nevertheless, CGF-Emergent is less than half of CGF-Lang-conj and CGL-Emergent is over 50 times larger than CGL-Lang-conj. It suggests that emergent languages are not fully compositional.

6.2 Comparison with Existing Measures

Next, we check the relationships among CGF/L, TopSim, and TRE. We show the results for $I = \text{Lang-conj}$, where CGF/L recognizes the compositionality of emergent languages. Figure 5 shows the scatter plot of TopSim and CGF. It shows a correlation with Pearson $\rho = 0.644$ ($p = 8.77 \times 10^{-24} \ll 0.01$). We also note that TopSim and CGL show a correlation with Pearson $\rho = -0.689$

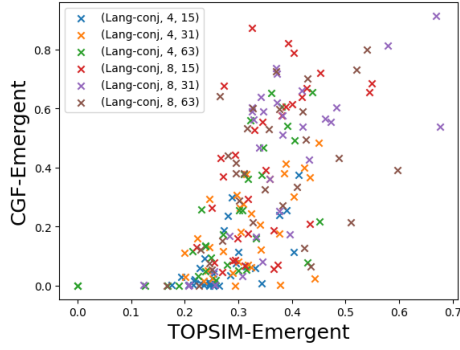


Figure 5: Scatter plot of CGF-Emergent and TopSim-Emergent, when $I = \text{Lang-conj}$. Pearson correlation is $\rho = 0.644$ ($p = 8.77 \times 10^{-24} \ll 0.01$).

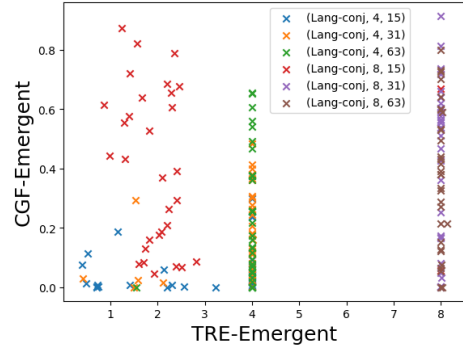


Figure 6: Scatter plot of CGF-Emergent and TRE-Emergent, when $I = \text{Lang-conj}$. Unnatural concentration around $k \in \{4, 8\}$ is observed.

($p = 2.88 \times 10^{-28} \ll 0.01$). Although p -values are considerably small, ρ s are moderate. Besides, Figure 5 shows several data points with high TopSim but low CGF. It suggests that TopSim tends to judge partially compositional languages more compositional than CGF.

Figure 6 shows the scatter plot of TRE and CGF. Astonishingly, it shows no correlation because of the unnatural concentration of TRE around $k \in \{4, 8\}$ if $|A| \in \{31, 63\}$. It means that a composition κ_η fails to learn so that its outputs are trapped between 0 and $1/|A|$. We speculate that the definition of κ_η or δ in Section 2.1 should have involved any nonlinear function. The scatter plots for CGLs are listed in Appendix D.

6.3 Example Derivation Tree of Emergent Language

Finally, Figure 2 exemplifies a derivation tree in an emergent language that CGI judges highly compositional ($\text{CGF} = 0.914$, $\text{CGL} = 423$). We can see how the message is combined to yield the meaning, which is a striking feature of CGI that the existing compositionality measures do not have. In this example, 16,13 means “run,” 25,1,1 means “__ right 3,” and 1,1,1 means “__ and walk right 2.” Interestingly, it suggests message and meaning segmentation does not necessarily match the intuitive segmentation as shown in Section 3.1.

7 Discussion

The experimental results show that CGF and CGL work as a compositionality measure for emergent languages. Note that the observations on Lang-conj are consistent with those of van der Wal et al. (2020) in a sense that fully structured languages are judged the most syntactical, the emergent lan-

guages are judged the second, and lower structured broken languages are the least. However, neither CGF nor CGL recognizes the compositionality when an input space is a small set of attribute-value objects. It casts doubt on attribute-value settings for studying structural similarities between emergent and human languages. We found a moderate correlation between CGF/L and TopSim which suggests that CGI is not as sensitive to partial compositionality as TopSim. On the other hand, TRE does not work if the alphabet size is too large, probably due to the choice of δ or κ_η in Section 2.1. Finally, we can directly observe the systematic composition of a message to a meaning, which is a salient feature of CGI that previous work does not have. We hope that it brings deeper insights on the syntax and semantics of emergent languages.

8 Conclusion

This paper introduces categorial grammar induction (CGI) as a new compositionality measure for the structure of emergent languages. We proposed to apply CGI to emergent languages and define two compositionality measures CGF and CGL. Our experiments revealed that CGF/L can measure compositionality as we expected. Unlike existing measures, our approach meets compositionality in a traditional sense, allowing us to analyze emergent languages with a lexicon and derivation trees. For future work, it would be interesting to study the structure of the derivations of emergent languages. Besides, we speculate that *situated CCGs* (Artzi and Zettlemoyer, 2013) are applicable, which induce CGs considering an external world. Hence, CGI may be applicable to visual referential games as well as 2D-grid world communication.

References

- Jacob Andreas. 2019. [Measuring compositionality in representation learning](#). In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.
- Yoav Artzi, Dipanjan Das, and Slav Petrov. 2014. [Learning compact lexicons for CCG semantic parsing](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 1273–1283. ACL.
- Yoav Artzi and Luke Zettlemoyer. 2013. [Weakly supervised learning of semantic parsers for mapping instructions to actions](#). *Trans. Assoc. Comput. Linguistics*, 1:49–62.
- Diane Bouchacourt and Marco Baroni. 2018. [How agents see things: On visual representations in an emergent language game](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, pages 981–985. Association for Computational Linguistics.
- Henry Brighton and Simon Kirby. 2006. [Understanding linguistic evolution by visualizing the emergence of topographic mappings](#). *Artif. Life*, 12(2):229–242.
- Rahma Chaabouni, Eugene Kharitonov, Diane Bouchacourt, Emmanuel Dupoux, and Marco Baroni. 2020. [Compositionality and generalization in emergent languages](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 4427–4442. Association for Computational Linguistics.
- Rahma Chaabouni, Eugene Kharitonov, Emmanuel Dupoux, and Marco Baroni. 2019a. [Anti-efficient encoding in emergent communication](#). In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 6290–6300.
- Rahma Chaabouni, Eugene Kharitonov, Alessandro Lazaric, Emmanuel Dupoux, and Marco Baroni. 2019b. [Word-order biases in deep-agent emergent communication](#). In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 5166–5175. Association for Computational Linguistics.
- Gautier Dagan, Dieuwke Hupkes, and Elia Bruni. 2021. [Co-evolution of language and agents in referential games](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, EACL 2021, Online, April 19 - 23, 2021*, pages 2993–3004. Association for Computational Linguistics.
- Li Dong and Mirella Lapata. 2016. [Language to logical form with neural attention](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, August 7-12, 2016, Berlin, Germany, Volume 1: Long Papers*. The Association for Computer Linguistics.
- Andrew Drozdo, Patrick Verga, Mohit Yadav, Mohit Iyyer, and Andrew McCallum. 2019. [Unsupervised latent tree induction with deep inside-outside recursive auto-encoders](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 1129–1141. Association for Computational Linguistics.
- Jakob N. Foerster, Yannis M. Assael, Nando de Freitas, and Shimon Whiteson. 2016. [Learning to communicate with deep multi-agent reinforcement learning](#). In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 2137–2145.
- Laura Graesser, Kyunghyun Cho, and Douwe Kiela. 2019. [Emergent linguistic phenomena in multi-agent communication games](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, pages 3698–3708. Association for Computational Linguistics.
- Serhii Havrylov and Ivan Titov. 2017. [Emergence of language with multi-agent games: Learning to communicate with sequences of symbols](#). In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 2149–2159.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. [Long short-term memory](#). *Neural Comput.*, 9(8):1735–1780.
- Eric Jang, Shixiang Gu, and Ben Poole. 2017. [Categorical reparameterization with gumbel-softmax](#). In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net.
- Eugene Kharitonov, Rahma Chaabouni, Diane Bouchacourt, and Marco Baroni. 2020. [Entropy minimization in emergent languages](#). In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 5220–5230. PMLR.
- Simon Kirby. 2001. [Spontaneous evolution of linguistic structure-an iterated learning model of the emergence of regularity and irregularity](#). *IEEE Trans. Evol. Comput.*, 5(2):102–110.

708	Satwik Kottur, José M. F. Moura, Stefan Lee, and Dhruv	<i>Information Processing Systems 2019, NeurIPS 2019,</i>	765
709	Batra. 2017. Natural language does not emerge 'nat-	<i>December 8-14, 2019, Vancouver, BC, Canada,</i>	766
710	urally' in multi-agent dialog. In <i>Proceedings of the</i>	15825–15835.	767
711	<i>2017 Conference on Empirical Methods in Natural</i>		
712	<i>Language Processing, EMNLP 2017, Copenhagen,</i>	Yaoyiran Li, Edoardo Maria Ponti, Ivan Vulic, and Anna	768
713	<i>Denmark, September 9-11, 2017,</i> pages 2962–2967.	Korhonen. 2020. Emergent communication pretrain-	769
714	Association for Computational Linguistics.	ing for few-shot machine translation. In <i>Proceedings</i>	770
		<i>of the 28th International Conference on Computa-</i>	771
715	Tom Kwiatkowski, Luke S. Zettlemoyer, Sharon Gold-	<i>tional Linguistics, COLING 2020, Barcelona, Spain</i>	772
716	water, and Mark Steedman. 2010. Inducing proba-	<i>(Online), December 8-13, 2020,</i> pages 4716–4731.	773
717	bilistic CCG grammars from logical form with higher-	International Committee on Computational Linguistics.	774
718	order unification. In <i>Proceedings of the 2010 Con-</i>		775
719	<i>ference on Empirical Methods in Natural Language</i>	Chris J. Maddison, Andriy Mnih, and Yee Whye Teh.	776
720	<i>Processing, EMNLP 2010, 9-11 October 2010, MIT</i>	2017. The concrete distribution: A continuous re-	777
721	<i>Stata Center, Massachusetts, USA, A meeting of SIG-</i>	laxation of discrete random variables. In <i>5th Inter-</i>	778
722	<i>DAT, a Special Interest Group of the ACL,</i> pages	<i>national Conference on Learning Representations,</i>	779
723	1223–1233. ACL.	<i>ICLR 2017, Toulon, France, April 24-26, 2017, Con-</i>	780
		<i>ference Track Proceedings.</i> OpenReview.net.	781
724	Tom Kwiatkowski, Luke S. Zettlemoyer, Sharon Gold-	Igor Mordatch and Pieter Abbeel. 2018. Emergence	782
725	water, and Mark Steedman. 2011. Lexical generaliza-	of grounded compositional language in multi-agent	783
726	tion in CCG grammar induction for semantic parsing.	populations. In <i>Proceedings of the Thirty-Second</i>	784
727	In <i>Proceedings of the 2011 Conference on Empirical</i>	<i>AAAI Conference on Artificial Intelligence, (AAAI-</i>	785
728	<i>Methods in Natural Language Processing, EMNLP</i>	<i>18), the 30th innovative Applications of Artificial</i>	786
729	<i>2011, 27-31 July 2011, John McIntyre Conference</i>	<i>Intelligence (IAAI-18), and the 8th AAAI Symposium</i>	787
730	<i>Centre, Edinburgh, UK, A meeting of SIGDAT, a Spe-</i>	<i>on Educational Advances in Artificial Intelligence</i>	788
731	<i>cial Interest Group of the ACL,</i> pages 1512–1523.	<i>(EAAI-18), New Orleans, Louisiana, USA, February</i>	789
732	ACL.	2-7, 2018,	790
		pages 1495–1502. AAAI Press.	
733	Brenden M. Lake and Marco Baroni. 2018. General-	Franz Josef Och and Hermann Ney. 2003. A systematic	791
734	ization without systematicity: On the compositional	comparison of various statistical alignment models.	792
735	skills of sequence-to-sequence recurrent networks.	<i>Comput. Linguistics,</i> 29(1):19–51.	793
736	In <i>Proceedings of the 35th International Conference on</i>		
737	<i>Machine Learning, ICML 2018, Stockholmsmässan,</i>	Yi Ren, Shangmin Guo, Matthieu Labeau, Shay B. Co-	794
738	<i>Stockholm, Sweden, July 10-15, 2018,</i> volume 80 of	hen, and Simon Kirby. 2020. Compositional lan-	795
739	<i>Proceedings of Machine Learning Research,</i> pages	guages emerge in a neural iterated learning model.	796
740	2879–2888. PMLR.	In <i>8th International Conference on Learning Repre-</i>	797
		<i>sentations, ICLR 2020, Addis Ababa, Ethiopia, April</i>	798
741	Angeliki Lazaridou and Marco Baroni. 2020. Emergent	26-30, 2020.	799
742	multi-agent communication in the deep learning era.	OpenReview.net.	
743	<i>CoRR,</i> abs/2006.02419.		
744	Angeliki Lazaridou, Karl Moritz Hermann, Karl Tuyls,	Yoav Seginer. 2007. Fast unsupervised incremental pars-	800
745	and Stephen Clark. 2018. Emergence of linguistic	ing. In <i>ACL 2007, Proceedings of the 45th Annual</i>	801
746	communication from referential games with sym-	<i>Meeting of the Association for Computational Lin-</i>	802
747	bolic and pixel input. In <i>6th International Conference</i>	<i>guistics, June 23-30, 2007, Prague, Czech Republic.</i>	803
748	<i>on Learning Representations, ICLR 2018, Vancouver,</i>	The Association for Computational Linguistics.	804
749	<i>BC, Canada, April 30 - May 3, 2018, Conference</i>		
750	<i>Track Proceedings.</i> OpenReview.net.	Agnieszka Słowik, Abhinav Gupta, William L. Hamil-	805
		ton, Mateja Jamnik, Sean B. Holden, and Christo-	806
751	Angeliki Lazaridou, Anna Potapenko, and Olivier Tiele-	pher J. Pal. 2021. Structural inductive biases in	807
752	man. 2020. Multi-agent communication meets nat-	emergent communication. In <i>Proceedings of The</i>	808
753	ural language: Synergies between functional and	<i>43rd Annual Meeting of the Cognitive Science Soci-</i>	809
754	structural language learning. In <i>Proceedings of the</i>	<i>ety, CogSci 2021.</i>	810
755	<i>58th Annual Meeting of the Association for Computa-</i>		
756	<i>tional Linguistics, ACL 2020, Online, July 5-10,</i>	Mark Steedman. 1996. Surface structure and interpre-	811
757	2020,	tation, volume 30 of <i>Linguistic inquiry.</i> MIT Press.	812
758	pages 7663–7674. Association for Computa-		
	tional Linguistics.	Mark Steedman. 2000. The syntactic process. Lan-	813
		guage, speech, and communication. MIT Press.	814
759	David K. Lewis. 1969. Convention: A Philosophical	Oskar van der Wal, Silvan de Boer, Elia Bruni, and	815
760	Study. Wiley-Blackwell.	Dieuwke Hupkes. 2020. The grammar of emergent	816
		languages. In <i>Proceedings of the 2020 Conference on</i>	817
761	Fushan Li and Michael Bowling. 2019. Ease-of-	<i>Empirical Methods in Natural Language Processing,</i>	818
762	teaching and language structure from emergent com-	<i>EMNLP 2020, Online, November 16-20, 2020,</i> pages	819
763	munication. In <i>Advances in Neural Information Pro-</i>	3339–3359. Association for Computational Linguistics.	820
764	<i>cessing Systems 32: Annual Conference on Neural</i>		821

822 Ronald J. Williams. 1992. [Simple statistical gradient-](#)
823 [following algorithms for connectionist reinforcement](#)
824 [learning](#). *Mach. Learn.*, 8:229–256.

825 Ronald J. Williams and Jing Peng. 1991. Function opti-
826 mization using connectionist reinforcement learning
827 algorithms. *Connection Science*, 3:241–268.

828 Luke S. Zettlemoyer and Michael Collins. 2005. [Learn-](#)
829 [ing to map sentences to logical form: Structured clas-](#)
830 [sification with probabilistic categorial grammars](#). In
831 *UAI '05, Proceedings of the 21st Conference in Un-*
832 *certainty in Artificial Intelligence, Edinburgh, Scot-*
833 *land, July 26-29, 2005*, pages 658–666. AUAI Press.

834 Luke S. Zettlemoyer and Michael Collins. 2007. [Online](#)
835 [learning of relaxed CCG grammars for parsing to](#)
836 [logical form](#). In *EMNLP-CoNLL 2007, Proceedings*
837 *of the 2007 Joint Conference on Empirical Meth-*
838 *ods in Natural Language Processing and Computa-*
839 *tional Natural Language Learning, June 28-30, 2007,*
840 *Prague, Czech Republic*, pages 678–687. ACL.

841 A Modifications of CGI

842 **INITLEX** We set $\Lambda_{\text{seed}} = \emptyset$, as we do not have
843 any prior knowledge on emergent languages.

844 **UPDATELEX** In KZGS10, UPDATELEX in-
845 cludes part of a potential new lexicon pruning the
846 rest, while ours includes all of them. This is be-
847 cause the PRUNEX of ADP14 would implicitly
848 do the same thing. Moreover, the original UP-
849 DATELEX splits lexical entries as a higher-order
850 unification problem to find f and g s.t. $h = f(g)$
851 or $h = f \circ g$, given a logical form h . On the other
852 hand, ours splits the entries as a problem only to
853 find $h = f(g)$, ensuring that $f \neq \lambda x.x$. and g is
854 not a function.

855 **INITPARAM** Since the algorithm can only search
856 limited space in practice, a reasonable parameter
857 initialization is required. KZGS10 used a statisti-
858 cal translation method⁹, while we simply compute
859 mean pointwise mutual information (pmi) between
860 n-grams and logical constants. Formally, given a
861 feature, i.e., a lexical entry $m \vdash X : \psi$, its initial
862 parameter is defined as:

$$863 \frac{1}{|\text{Cnst}(\psi)|} \sum_{c \in \text{Cnst}(\psi)} \text{pmi}(m, c)$$

864 if $|\text{Cnst}(\psi)| > 0$ otherwise 0. $\text{Cnst}(\psi)$ enumer-
865 ates the logical constants (e.g. look, left, or 1)
866 occurring in ψ .

⁹Giza++ Model 1 (Och and Ney, 2003).

867 B (Hyper-)parameters

868 **Agents** For agent architecture, the hidden state
869 size is 100. For agent optimization, the number of
870 mini-batches per epoch is 100, the size of mini-
871 batches is 1000, and the learning rate is 0.001.
872 Agents train either for 200 epochs or until loss
873 \mathcal{L} for a validation dataset reaches 0. Besides, the
874 weight of speaker’s (resp. listener’s) entropy regu-
875 larizer $\lambda_S = 0.1$ (resp. $\lambda_L = 1$). These parameters
876 are determined according to our preliminary exper-
877 iments.

878 **Signaling Game** For signaling games, an input
879 space $I \in \{\text{Lang-attval}, \text{Lang-conj}\}$, the size $|A|$
880 of an alphabet A is in $\{15, 31, 63\}$, and a message
881 length $k \in \{4, 8\}$.

882 **CGI** For CGI, the number of iterations $T = 10$,
883 a learning rate $\gamma = 0.1$, and a beam size for CKY
884 parsing is 10, referring to Artzi et al. (2014) and
885 our preliminary experiments.

886 **TRE** For TRE, a learning rate is 0.01 and the
887 number of steps is 1000 following the implementa-
888 tion of Andreas (2019).

889 C Shuffled Emergent Language and 890 Random Sequence

891 [Figure 7](#) and [Figure 9](#) show the compari-
892 son among CGF/L-Emergent, CGF/L-Shuffled,
893 CGF/L-Random.

894 D Other Experimental Results

895 [Figure 8](#) shows the scatter plot of TopSim and CGL
896 when $I = \text{Lang-conj}$. [Figure 10](#) shows the scatter
897 plot of TRE and CGL when $I = \text{Lang-conj}$.

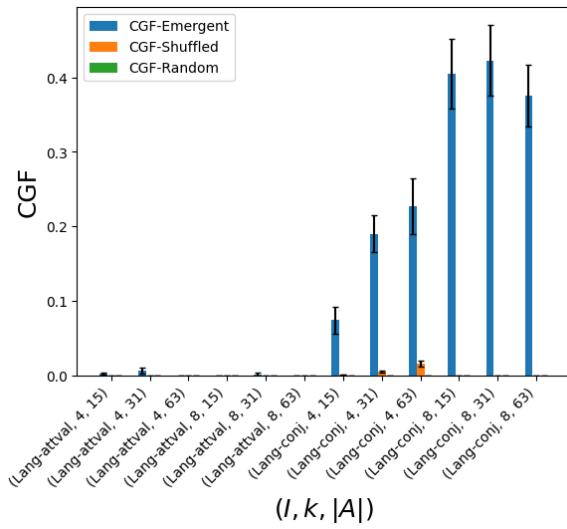


Figure 7: CGF plotted under various $(I, k, |A|)$. The error bars represent one standard error of mean.

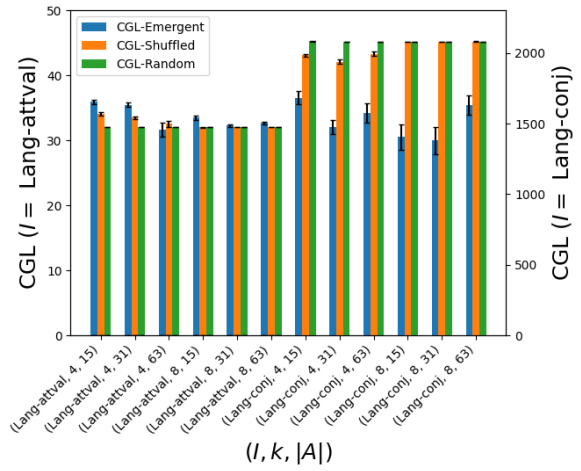


Figure 9: CGL plotted under various $(I, k, |A|)$. The error bars represent one standard error of mean.

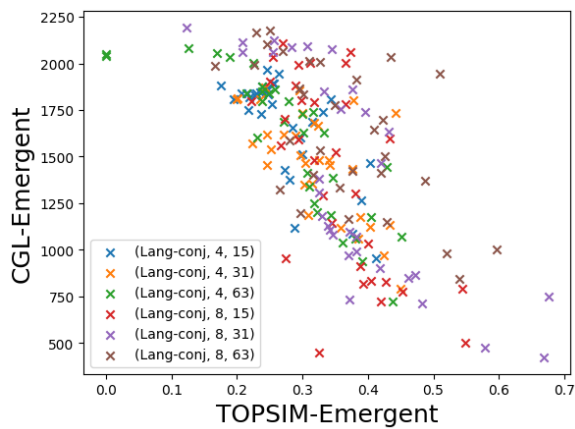


Figure 8: Scatter plot of CGL-Emergent and TopSim-Emergent.

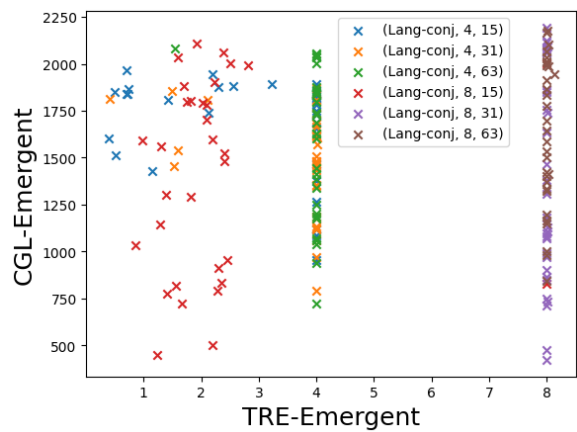


Figure 10: Scatter plot of CGL-Emergent and TRE-Emergent.