





University of Potsdam  
Faculty of Human Sciences

# Bachelor Thesis

in Computational Linguistics

Submitted in Fulfillment of the Requirements for the Degree of  
Bachelor of Science

<b>Topic:</b>	Investigating the ability of RNN architectures to learn context-free grammars by example of Dyck(2)
<b>Author:</b>	Fynn Dobler <fynndobler@gmail.com> Matr.-Nr. 775710
<b>Version of:</b>	March 4, 2020
<b>1. Supervisor:</b>	Dr. Thomas Hanneforth
<b>2. Supervisor:</b>	Dr. Uladzimir Sidarenka

---

**Summary**

**Abstract**

---

## Zusammenfassung

---

# Contents

<b>Index of Figures</b>	<b>5</b>
<b>Index of Tables</b>	<b>5</b>
<b>List of Abbreviations</b>	<b>6</b>
<b>1 Introduction</b>	<b>7</b>
<b>2 Theoretical Background</b>	<b>8</b>
2.1 Formal Languages and Formal Grammars . . . . .	8
2.2 Formal Grammars and Natural Language . . . . .	9
2.2.1 Natural Language as supra-regular . . . . .	9
2.2.2 Natural Language as supra-context-free . . . . .	10
2.3 Dyck Languages . . . . .	11
2.4 Neural Network Architectures . . . . .	12
2.4.1 Simple RNN . . . . .	12
2.4.2 LSTM . . . . .	13
2.4.3 GRU . . . . .	14
2.5 Related Works . . . . .	14
<b>3 Experiment Setup</b>	<b>17</b>
3.1 Models . . . . .	17
3.2 Corpus Construction . . . . .	17
3.3 Evaluation . . . . .	19
3.4 Experiment 1: Long-Range Dependency . . . . .	19
3.5 Experiment 2: Deeper Nesting . . . . .	20
<b>4 Results</b>	<b>20</b>
4.1 Experiment 1: Long-Range Dependency . . . . .	20
4.2 Experiment 2: Deeper Nesting . . . . .	20
<b>5 Discussion</b>	<b>20</b>
<b>6 Conclusion</b>	<b>20</b>
<b>Bibliography</b>	<b>21</b>
<b>Appendix</b>	<b>24</b>
<b>Eidesstattliche Erklärung</b>	<b>25</b>

## List of Figures

1	Chomsky Hierarchy . . . . .	8
2	Unfolded RNN . . . . .	13
3	LSTM Memory Cell . . . . .	14
4	Illustration of a GRU . . . . .	15

## List of Tables

1	Formal grammar properties. . . . .	9
2	Corpus sizes in current works . . . . .	15
3	Reported values for performance in previous works . . . . .	16
4	Overview of investigated models . . . . .	16

# 1 Introduction

In the 2010s, seemingly every major technology company developed its own "personal assistant" system, a program that allows the end-user to interact with the company's services more intuitively by interpreting spoken natural language commands. Apple's Siri, Amazon's Alexa and Google's succinctly named Assistant have been irrevocably ingrained in day-to-day life. While the ethical and data security concerns raised by this development are still a point of contention, it is clear that Natural Language Processing (NLP) applications have boomed from a niche field to a rapidly growing multi-million dollar industry<sup>1</sup>. Despite state-of-the-art performance on NLP tasks such as machine translation, text classification, sentiment analysis and speech recognition having made leaps and bounds in the past decade, these systems are still far from acquiring a perfect understanding of natural language. Recently, many new Recurrent Neural Network (RNN) model ideas have been experimented with, like the Clockwork RNN (Koutník et al. (2014)) or the Recurrent Unit with a Stack-like State (RUSS) (Bernardy (2018)), often designed to excel at a specific task. To showcase the new model's superiority, its performance on a task is usually compared to that of a more well-known architecture, such as the Simple RNN (SRNN), the Long Short Term Memory (LSTM) or the Gated Recurrent Unit (GRU).

What is missing from the current state of literature is, however, a robust comparison of these three architectures on a task that adequately showcases their respective ability to perform well on natural language data. I seek to fill that gap with my work by trying to answer the following questions:

1. Can an SRNN, LSTM or GRU architecture learn the Dyck language with two pairs of brackets ( $D_2$ )?
2. If they cannot, what poses the highest difficulty in doing so?
3. What influence, if any, does corpus construction have on model performance, specifically generalizability?

The following work consists of five main parts, each of which will be briefly summarized hereunder:

In Chapter 2, I introduce the core concepts relevant for this thesis: formal languages, the complexity of Natural Language, Dyck languages, three neural network architectures and an overview of related works to contextualize my work within the current state of research. I describe model design and training, corpus construction and the two experiments I conduct in this work in Chapter 3. I report the results for each of the experiments in Chapter 4 and discuss them in detail in Chapter 5. Finally,

---

<sup>1</sup><https://www.tractica.com/newsroom/press-releases/natural-language-processing-market-to-reach-22-3-billion-by-2025/>

I seek to answer the research questions posed above with my experimental results in Chapter 6 and suggest further avenues of research on this topic.

## 2 Theoretical Background

### 2.1 Formal Languages and Formal Grammars

A formal language  $L(G)$  is defined as a subset of all words  $\Sigma^*$  over an alphabet  $\Sigma$ , where all words need to comply with the formal grammar  $G$ . As per Jurafsky and Martin (2009), the definition of a formal grammar is  $G = \{N, \Sigma, R, S\}$ , where  $N$  is a set of non-terminal symbols,  $\Sigma$  is a set of terminal symbols (alphabet),  $R$  is a set of rules of the form  $\alpha \rightarrow \beta$  (where  $\alpha$  and  $\beta$  are strings of symbols from  $(\Sigma \cup N)^*$ ) and  $S$  is a designated start symbol. Following the two definitions,  $L(G)$  consists of all strings  $w$  that can be derived from the start symbol  $S$  in a finite number of steps, formally  $\{w \in \Sigma^* | S \xRightarrow[G]{*} w\}$ . As such, a word  $w \in \Sigma^*$  that cannot be derived from  $S$  in a finite number of steps is not part of  $L(G)$ .

Formal grammars differ in terms of complexity and can be described in a hierarchical manner. Grammars of higher complexity have a greater generative power than grammars of lower complexity. The most commonly used hierarchy of grammars is the Chomsky hierarchy (Chomsky (1959)). In this hierarchy, formal grammars are classified into four types, sorted from most powerful to least powerful: Turing equivalent (Type 0), Context Sensitive (Type 1), Context Free (Type 2) and Regular (Type 3). The difference in generative power and complexity stems from increasing restrictions imposed on the rules of the grammar - a Type 3 grammar is more restrictive than a Type 0 grammar. As such, every grammar of a higher type is a subset of the previous type of grammar. A visual representation of this property can be found in Figure 1.

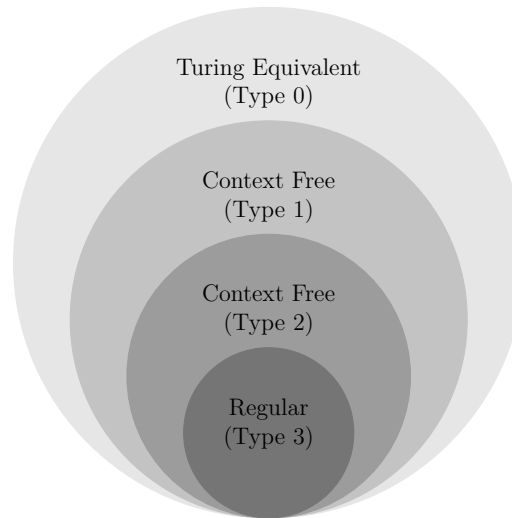


Figure 1: A visual representation of the Chomsky Hierarchy.

The four types of formal grammars can be defined by the form their rules can



take. An overview over these rules as per Jurafsky and Martin (2009) can be found in Table 1, where  $A$  is a single non-terminal,  $\alpha, \beta, \gamma$  are strings of terminal and non-terminal symbols, and  $x$  is a string of terminal symbols.  $\alpha, \beta$  and  $\gamma$  may be empty unless specifically disallowed. The table is supplemented with a column describing the corresponding automaton capable of accepting or recognizing the grammar.

## 2.2 Formal Grammars and Natural Language

The correspondence of formal grammars to automata (i.e. Kleene’s Theorem for regular languages and finite automata) and Computational Complexity Theory lends itself to consider natural languages under the same lense. While formal grammars constitute powerful tools with which phenomena in natural language can be modelled, assessing the precise complexity of Natural Language is the subject of ongoing investigation (Fitch et al. (2012), Petersson and Hagoort (2012), Newmeyer and Preston (2014)). Arguments answering that question usually seek to establish lower bounds: If there is a phenomenon in a natural language that cannot be described with a given type of grammar, natural language must be - however slightly - more complex than that type allows. Such arguments increase in credibility the more frequently they can be replicated for phenomena in multiple languages. The arguments establishing natural languages as supra-context-free (i.e. more complex than CFGs) as well as contrary evidence from empiric research shall be presented here.

### 2.2.1 Natural Language as supra-regular

English, as well as several other languages (Hagège (1976)) allow for center embedding, the embedding of a phrase into another phrase of the same type.

- (1) The man eats.
- (2) The man the boss fired eats.
- (3) The man the boss the investor distrusted fired eats.

<i>Type</i>	<i>Name</i>	<i>Rule Skeleton</i>	<i>Automaton</i>
0	Turing Equivalent	$\alpha \rightarrow \beta$ , s.t. $\alpha \neq \epsilon$	Turing Machine (recognized)
1	Context Sensitive	$\alpha A \beta \rightarrow \alpha \gamma \beta$ , s.t. $\gamma \neq \epsilon$	Linear Bound Automata (accepted)
2	Context Free	$A \rightarrow \gamma$	Push Down Automata (accepted)
3	Regular	$A \rightarrow xB$ or $A \rightarrow x$	Finite-State Automata (accepted)

Table 1: Overview of formal grammar properties according to Jurafsky and Martin (2009), augmented with corresponding automata.

- (4) The man the boss the investor the police investigated distrusted fired eats.

Let the set  $E$  contain all grammatical sentences of English, and let the noun phrases and transitive verbs constitute following sets:

$$A = \{\text{the boss, the investors, the police, } \dots\}$$

$$B = \{\text{fired, distrusted, investigated, } \dots\}$$

Then the following two sets can be defined.

$$E' = \{\text{the man } a^n b^n \text{ eats} \mid n \geq 0\}$$

$$R = \{\text{the man } a^* b^* \text{ eats}\}$$

$a^n$  and  $b^n$  are finite sequences of size  $n$  of elements of sets  $A$  and  $B$ , respectively.  $E'$  describes a subset of  $E$ , namely  $E \cap R$ . Since regular languages are closed under intersection and  $E'$  is not regular,  $E$  is not regular.<sup>2</sup>

While this proof is correct under the framework of Formal Language Theory, the validity of claiming that it shows natural language to be supra-regular is debatable. Research in psycholinguistics shows that native speakers have faced severe problems processing center embeddings of depth two or higher, yielding long processing times, an incomplete understanding of the presented sentence or leading the participants to judge the sentence as ungrammatical (Hamilton and Deese (1971), Frank et al. (2016)). Furthermore, the corpus-driven analysis by Karlsson (2007) suggests an upper limit of center embedding depth three in the seven investigated languages.

### 2.2.2 Natural Language as supra-context-free

Similarly to the proof given in Section 2.2.1, an argument characterizing natural language as supra-context-free can be brought forth. It is based on embedded infinitival verb phrases found in Swiss German (Shieber (1987)).

- (5) Jan säit das mer em Hans es huus haend wele hälfe  
 Jan said that we the Hans-DAT the house-ACC have wanted help  
 aastriiche.  
 paint  
 'Jan said that we have wanted to help Hans paint the house.'
- (6) Jan säit das mer d'chind em Hans es huus haend  
 Jan said that we the children-ACC the Hans-DAT the house-ACC  
 wele laa hälfe aastriiche.  
 have wanted let help paint  
 'Jan said that we have wanted to let the children help Hans paint the house.'

---

<sup>2</sup>The proofs for regular languages being closed under intersection and  $E'$  not being regular can be found in Hopcroft et al. (2006) and Sipser (2013).

Four finite sets can be constructed from these examples: accusative noun phrases ( $A = \{\text{d'chind}, \dots\}$ ), dative noun phrases ( $B = \{\text{em Hans}, \dots\}$ ), verbs taking accusative objects ( $C = \{\text{laa}, \dots\}$ ) and verbs taking dative objects ( $D = \{\text{hälfe}, \dots\}$ ). Let the set  $S$  then be the set of all grammatical sentences of Swiss German. Again, the two following sets can be defined:

$$\begin{aligned} S' &= \{\text{Jan säit das mer } a^n b^m \text{ es huus haend wele } c^n d^m \text{ aastriiche} \mid n, m \geq 0\} \\ R &= \{\text{Jan säit das mer } a^* b^* \text{ es huus haend wele } c^* d^* \text{ aastriiche}\} \end{aligned}$$

$S'$  is not context-free and results from  $S \cap R$ . Since context-free sets are closed under intersection with regular sets,  $G$  cannot be context-free.<sup>3</sup>

Curiously enough, empirical research into the matter of processing similar cross-serial dependencies in Dutch suggests them to be generally easier to process than nested dependencies (i.e. the ones used to prove natural language to be supra-regular) (Bach et al. (1986)).

### 2.3 Dyck Languages

Whether natural language is regular, context-free, supra-regular or supra-context-free is a distinction of only tangential relevance for this work. The first two cases are fully covered by CFGs, while the other two leave room for some natural language productions outside of the scope of CFGs. The characteristics of supra-context-free examples in natural language show a *weak* non-context-freeness, making CFGs sufficient for covering the vast majority of natural language productions. With this assumption, an appropriate CFG for a model to learn must be found. The most important property of this grammar is that model performance on its language must allow for strong conclusions about the learnability of any other CFG. In doing so, one can make reasoned assumptions about potential model performance on natural language data.

One such grammar is the Dyck Grammar, which can produce an array of Dyck Languages. Let  $D_n = \{N, \Sigma, R, S\}$  with

$$\begin{aligned} N &= \{S\} \\ \Sigma &= \{\epsilon, O_1, O_2, \dots, O_n, C_1, C_2, \dots, C_n\} \\ R &= \{ \\ &\quad S \rightarrow \epsilon \\ &\quad S \rightarrow SS \\ &\quad S \rightarrow O_n S C_n \}, \end{aligned}$$

---

<sup>3</sup>The respective proofs for  $S'$  not being context-free and context-free sets being closed under intersection with regular sets can be found in Hopcroft et al. (2006) and Sipser (2013).

where  $O_n$  represents an opening parenthesis,  $C_n$  represents a closing parenthesis and  $n$  denotes the number of distinct pairs of parentheses.  $D_1$ , then, denotes the Dyck Language with  $\Sigma = \{\epsilon, (, )\}$ ,  $D_2$  the Dyck Language with  $\Sigma = \{\epsilon, (, [, ], )\}$ , et cetera.

Within the family of Dyck Languages,  $D_2$  is of particular interest. According to the Chomsky-Schützenberger Representation Theorem (Chomsky and Schützenberger, 1963), for every context-free language  $L$  there exists a positive integer  $n$ , a regular language  $R$ , and a homomorphism  $h$  so that  $L = h(D_n \cap R)$ . Following the proof in Autebert et al. (1997), a homomorphism  $g_n$  can be constructed so that  $D_n = g_n^{-1}(D_2)$ . It follows that every context-free language can be represented as  $L = h(g_n^{-1}(D_2) \cap R)$ . As such, every CFL could be represented via homomorphisms on  $D_2$  and intersections with a regular language. Assuming natural languages to be context-free and bearing in mind that using a formal language is a choice of abstraction which allows for precise control over corpus composition, this makes  $D_2$  the language of choice when comparing neural network performance.

## 2.4 Neural Network Architectures

### 2.4.1 Simple RNN

Recurrent Neural Networks (RNNs) (Elman (1990)) are a neural network architecture particularly suited to processing sequential information by design: the RNN's output at a time step  $t$  is fed back as its input at the following time step  $t + 1$ . Not only does this enable RNNs to process sequences of arbitrary length, it also makes every output dependent on the previous computation as well as the current input. This property equips RNNs with a "memory" for previous inputs, allowing them to capture context dependencies a context-agnostic model cannot adequately learn.

Within the frame of this work, the specific case of the Simple RNN (SRNN) is considered. It is a three layer networks, consisting of an input layer, a hidden layer and an output layer. The hidden state  $h_t$  at time step  $t$  given the input vector  $x_t$  and the output vector  $y_t$  are calculated as per the following equations:

$$h_t = f(\mathbf{W}_{xh}x_t + \mathbf{W}_{hh}h_{t-1}) \quad (1)$$

$$y_t = \mathbf{W}_{hy}h_t \quad (2)$$

The function  $f$  constitutes a non-linear transformation, like tanh or ReLU.  $\mathbf{W}_{xh}$ ,  $\mathbf{W}_{hh}$ ,  $\mathbf{W}_{hy}$  are matrices of the weights connecting the input layer to the hidden layer, the hidden layer to itself and the hidden layer to the output layer, respectively.

When training RNNs, it is beneficial to think of the network as unfolding into an architecture with one layer per time step. A visualisation is provided in Figure 2. These conceptual layers share their parameters - if any weight changes at time step  $t$ , the weight also changes at  $t+1, t+2, \dots, t+n$ . Isolated changes are not possible. A popular

training algorithm for RNNs is Backpropagation Through Time (BPTT) (Williams and Zipser (1998)), a gradient based algorithm designed for recurrent rather than feedforward networks. However, as Bengio et al. (1994) and Hochreiter (1998) show, RNNs suffer from a fundamental flaw: the aptly named vanishing gradient problem, in which the training gradient diminishes to zero throughout the layers.

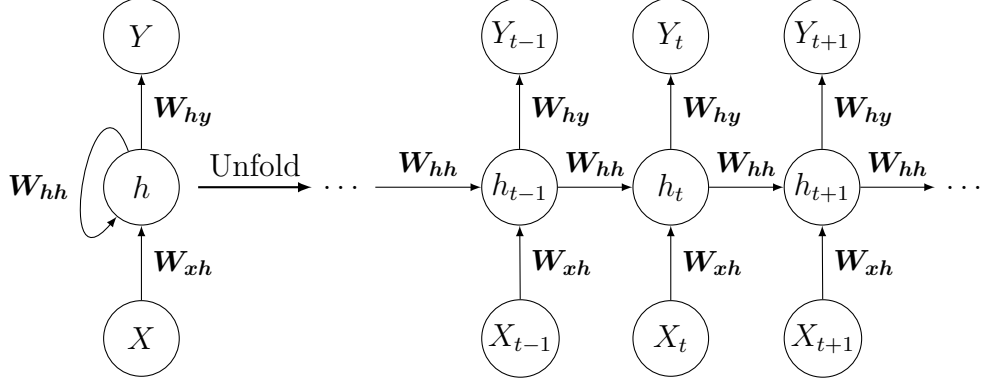


Figure 2: An RNN, unfolded through time.

### 2.4.2 LSTM

Long-Short Term Memory networks (LSTM) were designed by Hochreiter and Schmidhuber (1997) as an RNN architecture which preserves the RNN capabilities of processing sequential data of arbitrary length and capturing context dependencies, while circumventing the vanishing gradient problem.

LSTMs are based on self-connected linear units which are regulated by three gates consisting of a sigmoid layer  $\sigma$  each: input (in), output (out) and forget (forget). At every time step, the concatenated vector of the previous hidden state  $h_{t-1}$  and the current input  $x_t$  are received by all three gates. The sigmoid layer transforms every value in the concatenated vector to a value in range  $[0, \dots, 1]$  - a 0 translates to forgetting the information, while a 1 passes it through completely. Thus, the output of the gates determines what information is let through the input gate, passed through the output gate or forgotten by the self-connected linear unit.

$$\begin{aligned} \text{in}_t &= \sigma_{\text{in}}(\mathbf{W}_{\text{in}} \cdot [h_{t-1}, x_t] + b_{\text{in}}) \\ \text{out}_t &= \sigma_{\text{out}}(\mathbf{W}_{\text{out}} \cdot [h_{t-1}, x_t] + b_{\text{out}}) \\ \text{forget}_t &= \sigma_{\text{forget}}(\mathbf{W}_{\text{forget}} \cdot [h_{t-1}, x_t] + b_{\text{forget}}) \end{aligned}$$

Finally, the cell state  $C_{t-1}$  is updated to  $C_t$  and  $h_t$  is set.

$$\begin{aligned} C_t &= \text{forget}_t \odot C_{t-1} + \text{in}_t \odot \tanh(\mathbf{W}_C \cdot [h_{t-1}, x_t] + b_C) \\ h_t &= \text{out}_t \odot \tanh(C_t) \end{aligned}$$

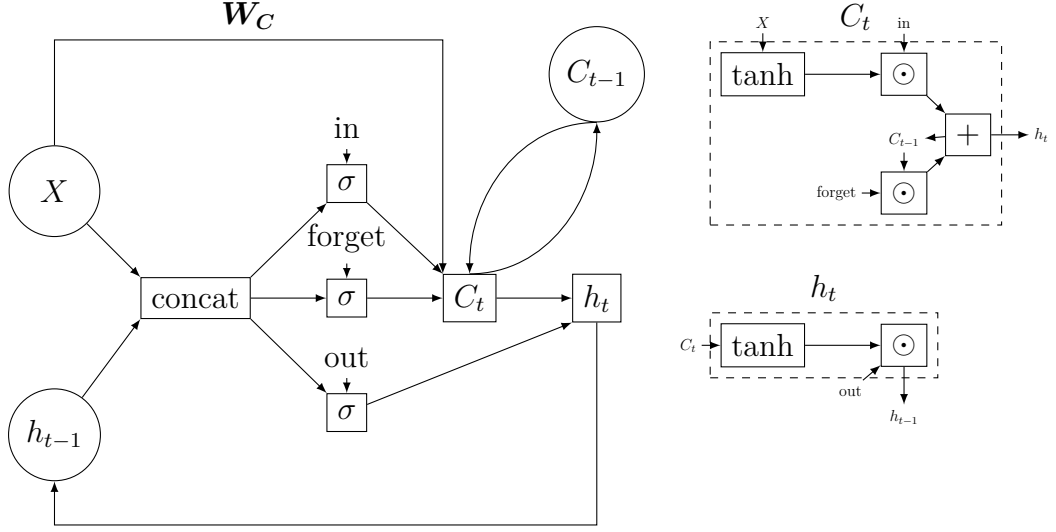


Figure 3: An LSTM memory cell.

### 2.4.3 GRU

A less complex alternative to LSTMs, the Gated Recurrent Unit (GRU) was developed by Cho et al. (2014). The information flow within the GRU is handled by just two gates: reset ( $r$ ) and update ( $z$ ). The update gate determines how much information from previous time steps is passed along for further time steps, while the reset gate enables the model to drop irrelevant information and only consider the current input rather than the previous hidden state, as described in the equations below, where  $j$  is the  $j$ -th hidden unit,  $\sigma$  is the squashing sigmoid function,  $\mathbf{W}$  and  $\mathbf{U}$  are learned gate-dependent weight matrices and  $\phi$  is a non-linear function.

$$\begin{aligned}
 r_j &= \sigma([\mathbf{W}_r x]_j + [\mathbf{U}_r h_{t-1}]_j) \\
 z_j &= \sigma([\mathbf{W}_z x]_j + [\mathbf{U}_z h_{t-1}]_j) \\
 h_j^t &= z_j h_j^{t-1} + (1 - z_j) \tilde{h}_j^t \\
 \tilde{h}_j^t &= \phi([\mathbf{W} x]_j + [\mathbf{U}(r \odot h_{t-1})]_j)
 \end{aligned}$$

## 2.5 Related Works

Currently, NLP models are trained and tested on vast datasets, such as the CoNLL Shared Tasks. By evaluating a variety of different models and approaches on the same data, it is possible to easily assess which one poses the current state-of-the-art for any given NLP task.

The earliest research on neural network architectures was done before such datasets were widely available and easy to process (Cleeremans et al. (1989), Elman (1990),

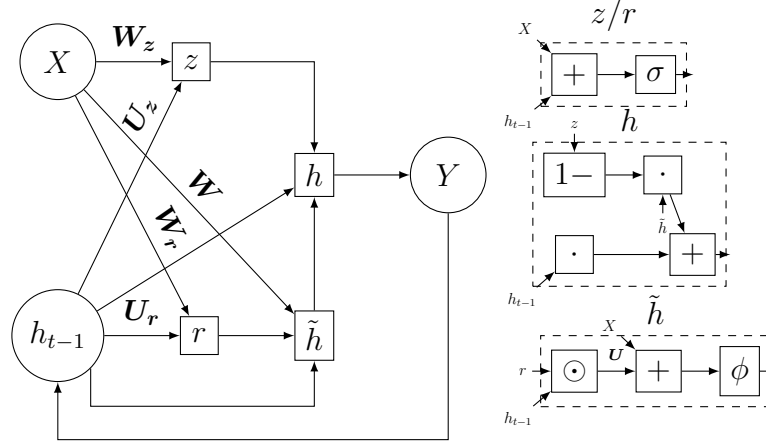


Figure 4: Illustration of a GRU.

<i>Paper</i>	$D_n$	<i>Grammar Probability</i>	<i>Training Corpus Size</i>
Deleu and Dureau (2016)	1	equal	unclear
Bernardy (2018)	1, 5	equal	102,400
Li et al. (2018)	1	modified	200 - 20,000
Skachkova et al. (2018)	1-5	modified	131,072
Sennhauser and Berwick (2018)	2	modified	1,000,000
Suzgun et al. (2019)	1-2	modified	10,000
Yu et al. (2019)	2	modified	1,000,000

Table 2: Overview of corpus sizes in current works.

Zeng et al. (1994), Hochreiter and Schmidhuber (1997), Rodriguez and Wiles (1998), Gers and Schmidhuber (2001)). During that time, novel architectures and algorithms were mostly scored on formal language datasets, with the test set containing longer words than the training set to assess learning success. However, evaluating on formal languages comes with its own advantages and challenges.

Primarily, it is undeniably cheaper than scoring on a natural language dataset. By deriving words from the grammar, datasets of arbitrary length with arbitrary properties can be generated. However, performance on a formal grammar dataset should always be understood as a simplified benchmark. As mentioned in Sections 2.2.1 and 2.2.2, the formal complexity of natural language is debatable, limiting the significance of formal benchmark performance for NLP tasks. Nevertheless, formal language datasets are still used to evaluate the performance of novel architectures to this day, as done by Joulin and Mikolov (2015), Bernardy (2018), Deleu and Dureau (2016), Li et al. (2018) and Yu et al. (2019).

In addition to the exploration of new architectures, formal languages are also still used to investigate particular behaviours of well-established architectures, such as LSTMs (Sennhauser and Berwick (2018)), or to compare several established models on a specific set of tasks (Skachkova et al. (2018), Suzgun et al. (2019)).

<i>Paper</i>	<i>Accuracy</i>	<i>Perplexity</i>	<i>Cell State</i>	<i>AUC</i>	<i>Error Rate</i>
Deleu and Dureau (2016)	No	No	No	Yes	No
Bernardy (2018)	Yes	No	No	No	No
Skachkova et al. (2018)	Yes	Yes	No	No	No
Sennhauser and Berwick (2018)	No	No	Yes	No	Yes
Suzgun et al. (2019)	Yes	No	Yes	No	No
Yu et al. (2019)	No	Yes	No	No	Yes

Table 3: Overview of reported values for performance. Cell State Analysis does not refer to a unified method, it merely means the paper investigates cell states at all. AUC refers to the area under the curve for an increasing length of Dyck words the model was able to generalize.

<i>Paper</i>	<i>Architectures</i>
Deleu and Dureau (2016)	Neural Turing Machine, LSTM
Bernardy (2018)	GRU, LSTM, RUSS
Skachkova et al. (2018)	SRNN, GRU, LSTM
Sennhauser and Berwick (2018)	LSTM
Suzgun et al. (2019)	SRNN, GRU, LSTM
Yu et al. (2019)	seq2seq

Table 4: Overview of investigated models.

While all these papers use a formal language to evaluate models, several factors prevent them from forming a solid basis upon which to compare their respective results: First, there is neither a benchmark train/dev/test set for Dyck languages (as is standard for most machine learning tasks) (Table 2), nor a set of measures that is reported consistently throughout the literature (Table 3). Additionally, only two papers compare the three well-established architectures (SRNN, GRU, LSTM) directly (Table 4). Finally, the employed training and test measures are not unified - in Bernardy (2018), for example, the models are trained to predict the next letter in words of variable length at any given time step, while in Suzgun et al. (2019), they predict the final letter of a word. Further model details, such as the inclusion of an Embedding and/or a Dropout layer or the number of hidden units also vary.

In conclusion, despite formal languages having been used to assess neural network model performance for decades, there is little to no comparative studies of SRNN, LSTM and GRU models performing on  $D_2$ . Any research comparing new architectures to any of these models does so with both varying training and testing methods, as well as with vastly differing corpora and as such cannot be directly compared with each other. This allows for no conclusive statement on the relative performance of these popular RNN architectures based on the current literature.



## 3 Experiment Setup

### 3.1 Models

For the following experiments, the three RNN architectures described in Section 2.4 have been used. All models consist of an embedding layer, a single layer of size  $n = \{2^1, 2^2, \dots, 2^n\}$  and a dense layer of size 1 with sigmoid activation. The activation value of the neuron in the dense layer acts as the output of the model: a value  $\geq 0.5$  means the model classified the input as a correct word. The models were implemented in Tensorflow 2.0.<sup>4</sup>

All models were trained with the same parameters. The training data was received one word at a time, in batches of 512. The loss was computed by binary cross-entropy, as is current standard for binary classification tasks. Furthermore, the Adam optimizer (Kingma and Ba (2014)) was applied with a learning rate of 0.0001. At the end of a training epoch, the models were evaluated for loss and accuracy on a validation set of 120,000 words. The models were trained until their loss on the validation set did not lower by more than 0.0001 for three consecutive epochs or for at most 100 epochs. The models with the lowest validation loss were used for all experiments.

The models were trained on the same training data for both experiments. To answer the question of training data influence on model performance, three distinct sets of training data were used, yielding a total of  $9 \times 3 \times 3 = 81$  (number of different hidden units  $\times$  number of training corpora  $\times$  number of different architectures) evaluated models.

### 3.2 Corpus Construction

To investigate the influence of corpus composition on model performance, three corpora were created: a baseline corpus which is directly sampled from a subset of  $D_2$ , as well as two modifications of the baseline corpus: one impoverishing the training data from long-range dependencies (Low LRD) and one enriching the training data with more long-range dependencies (High LRD). The sampling and modification processes will be explained later in this section.

The experiments were explicitly designed to test the models' abilities to generalize based on the training data they encounter. As such, it is prudent to give consideration to which properties the training data might possess to facilitate or inhibit generalizability - properties such as length, maximum nesting depth (ND) and the maximum distance between a pair of opening and closing brackets (BD). ND is, in this case, defined as the highest number of unresolved open brackets preceding an open bracket in a given word (i.e. in the word  $\{ \{ \{ \} \} \}$ , the square open bracket is at ND= 1, and the curly open bracket is at ND= 2, making the maximum ND of the word 2). Maximum

---

<sup>4</sup>The source code can be found at <https://github.com/FyDob/BSc-Thesis>.

BD, then, is the highest number of characters between a pair of brackets in a word. In the previous example word, the maximum BD would be 4. These measures are reported in Table ?? in terms of averages and variance.

Furthermore, the training corpora were chosen to be a small slice of a comparatively large subset of  $D_2$ . To facilitate generalization, the training corpora consist of words of varying length. As discussed in Section 2.5, previous works largely utilized similarly small language subsets and achieved encouraging results. For a discussion of Experiment 1 and 2 on a training corpus consisting of a majority of the target language, see Bernardy (2018).

In determining an eligible maximum length, a known fact about the size of  $D_n$  subsets was utilized: a Dyck language  $D_n$  contains  $n^m C_m$  words of length  $2m$ , where  $C_m$  is the  $m$ -th Catalan number (Skachkova et al. (2018)). It follows that a maximum length limit of  $2m$  produces a set of size  $\sum_{i=2}^{2m} n^i C_i$ . For example, a maximum length of 20 in  $D_2$  ( $D_2^{\leq 20}$ ) yields 20,119,506 words, which is a sufficiently large subset to sample from. The words were generated following the probabilistic grammar set forth by Sennhauser and Berwick (2018).

$$\begin{aligned} S &\rightarrow Z S \mid Z \\ Z &\rightarrow B \mid T \\ B &\rightarrow [ S ] \mid \{ S \} \\ T &\rightarrow [ ] \mid \{ \} \end{aligned}$$

The production  $Z \rightarrow B$  branches, whereas  $S \rightarrow Z S$  concatenates two smaller Dyck words. This representation provides a good intuition for understanding the merit of Experiment 1. The probabilities with which the rules were applied are calculated as follows, with alternative rules of course being applied with the complementary probability:

$$\begin{aligned} P_{\text{branch}} &= r_{\text{branch}} \cdot s(l) \quad \text{with } r_{\text{branch}} \sim U(0.7, 1.0) \\ P_{\text{concat}} &= r_{\text{concat}} \cdot s(l) \quad \text{with } r_{\text{concat}} \sim U(0.7, 1.0) \\ s(l) &= \min(1, -3 \cdot \frac{l}{n} + 3) \end{aligned}$$

with  $l$  being the number of already generated non-terminal characters and  $n$  the maximally desired length of the word.  $r_{\text{branch}}$ ,  $r_{\text{concat}}$  and  $l$  were sampled at every step of word generation.

Following this process, 500,000 words in  $D_2^{\leq 20}$  were generated. These words served as the basis for creating the three corpora. To create the Low LRD corpus, all words with a maximum bracket distance higher than 10 were modified<sup>5</sup> by first identifying

---

<sup>5</sup>This cut-off point was chosen as it significantly reduces the average maximum bracket distance

the bracket pair with the highest bracket distance, then simply moving the opening bracket from its original position to the position right before the closing bracket. (i.e.  $\{\{\{\}\}\}$  becomes  $\{\{\}\}\{\}$ ). This has the largest impact on bracket distance throughout the corpus, while ensuring grammaticality of the resulting word. The resulting set of long-range impoverished words was merged with all unmodified words, deleting all duplicates.

The High LRD corpus was created in a similar way: First, all words with a bracket distance lower than 19 were identified.<sup>6</sup> Then, the first pair of neighbouring closing brackets is found and deleted. The remaining word is wrapped in a randomly chosen pair of brackets, creating the longest possible bracket distance between the two (i.e.  $\{\{\{\}\}\}$  becomes  $\{\{\{\}\}\}$ ). The resulting set was merged with the unmodified words the same way as the Low LRD set.

Finally, the corpora were filled with 500,000 non-words obtained by corrupting the correct words in  $D_2^{\leq 20}$ . For one half of the words, a random opening bracket was replaced with a random closing bracket, while a random closing bracket was replaced with a random opening bracket for the other half.

In total, all corpora consist of 1,000,000 samples, of which 50% are incorrect.

### 3.3 Evaluation

Generating the words for the two experiments follow the procedure described by Bernardy (2018). Whereas the research put forth in his paper scrutinized the generative abilities of RNNs, I am investigating the models in a classification task.

### 3.4 Experiment 1: Long-Range Dependency

$$w_{\text{LRD}} = O_n w_1 w_2 C_n$$

While  $w_1$  and  $w_2$  might have been seen in training, the resulting word most certainly has not been observed. Neither could the model possibly have encountered a long-range dependency spanning 36 characters between the opening and closing bracket. As such, a high accuracy on the final character of the word serves as a strong indication for the model having not just learned the necessity of closing open brackets to form an acceptable Dyck word, but also for the model having retained information on the first opening bracket during all processing steps. For further insight in model behaviour, the accuracy is reported for every processed character.

---

without creating too many duplicates.

<sup>6</sup>The same considerations as for the Low LRD corpus cut-off apply.

### 3.5 Experiment 2: Deeper Nesting

To investigate how well a model performs on predicting brackets on a nesting level deeper than anything included in training, another test set was constructed. Since Experiment 1 already investigates Long-Range Dependency (LRD), this corpus was designed so its results are influenced as little as possible by LRD performance.

$$w_{\text{DN}} = O_n O_n O_n O_n O_n w C_n C_n C_n C_n C_n$$

This process still has the model extrapolate beyond the length of the training words, while increasing all present nesting depths by 5. Indeed, when analysing the results, the main focus lies on the infix word, since all its nesting depths are now increased by 5, possibly affecting model performance.

## 4 Results

### 4.1 Experiment 1: Long-Range Dependency

### 4.2 Experiment 2: Deeper Nesting

## 5 Discussion

## 6 Conclusion

This work has set out to answer three questions, as posed in Chapter 1. Related literature has been consulted to choose a proper approach. However, current literature contains neither a benchmark dataset to train and test models on, nor a unified set of tasks and measures to do so. Due to these facts, most results in current literature discussing model performance on  $D_2$  are incomparable to each other.

## Bibliography

- Autebert, J.-M., Berstel, J., and Boasson, L. (1997). *Context-Free Languages and Pushdown Automata*, pages 111–174. Springer Berlin Heidelberg, Berlin, Heidelberg.
- Bach, E., Brown, C., and Marslen-Wilson, W. (1986). Crossed and nested dependencies in german and dutch: A psycholinguistic study. *Language and Cognitive Processes*, 1:249–262.
- Bengio, Y., Simard, P., and Frasconi, P. (1994). Learning long-term dependencies with gradient descent is difficult. *IEEE Transactions on Neural Networks*, 5(2):157–166.
- Bernardy, J.-P. (2018). Can recurrent neural networks learn nested recursion? volume 16.
- Cho, K., van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., and Bengio, Y. (2014). Learning phrase representations using RNN encoder–decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1724–1734, Doha, Qatar. Association for Computational Linguistics.
- Chomsky, N. (1959). On certain formal properties of grammars. *Information and Control*, 2(2):137 – 167.
- Chomsky, N. and Schützenberger, M. (1963). The algebraic theory of context-free languages\*. In Braffort, P. and Hirschberg, D., editors, *Computer Programming and Formal Systems*, volume 35 of *Studies in Logic and the Foundations of Mathematics*, pages 118 – 161. Elsevier.
- Cleeremans, A., Servan-Schreiber, D., and McClelland, J. (1989). Finite state automata and simple recurrent networks. *Neural Computation - NECO*, 1:372–381.
- Deleu, T. and Dureau, J. (2016). Learning operations on a stack with neural turing machines. *CoRR*, abs/1612.00827.
- Elman, J. L. (1990). Finding structure in time. *Cognitive Science*, 14(2):179–211.
- Fitch, W. T., Friederici, A. D., and Hagoort, P. (2012). Pattern perception and computational complexity: introduction to the special issue. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 367(1598):1925–1932.
- Frank, S. L., Trompenaars, T., and Vasishth, S. (2016). Cross-linguistic differences in processing double-embedded relative clauses: Working-memory constraints or language statistics? *Cognitive Science*, 40(3):554–578.
- Gers, F. A. and Schmidhuber, E. (2001). Lstm recurrent networks learn simple context-free and context-sensitive languages. *IEEE Transactions on Neural Networks*, 12(6):1333–1340.
- Hagège, C. (1976). Relative clause, center-embedding, and comprehensibility. *Linguistic Inquiry*, 7(1):198–201.
- Hamilton, H. W. and Deese, J. (1971). Comprehensibility and subject-verb relations in complex sentences. *Journal of Verbal Learning and Verbal Behavior*, 10(2):163 – 170.

- Hochreiter, S. (1998). The vanishing gradient problem during learning recurrent neural nets and problem solutions. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 6:107–116.
- Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9:1735–80.
- Hopcroft, J. E., Motwani, R., and Ullman, J. D. (2006). *Introduction to Automata Theory, Languages, and Computation (3rd Edition)*. Addison-Wesley Longman Publishing Co., Inc., USA.
- Joulin, A. and Mikolov, T. (2015). Inferring algorithmic patterns with stack-augmented recurrent nets. *CoRR*, abs/1503.01007.
- Jurafsky, D. and Martin, J. H. (2009). *Speech and Language Processing (2Nd Edition)*. Prentice-Hall, Inc., Upper Saddle River, NJ, USA.
- Karlsson, F. (2007). Constraints on multiple center-embedding of clauses. *Journal of Linguistics*, 43(2):365–392.
- Kingma, D. and Ba, J. (2014). Adam: A method for stochastic optimization. *International Conference on Learning Representations*.
- Koutník, J., Greff, K., Gomez, F., and Schmidhuber, J. (2014). A clockwork rnn.
- Li, T., Rabusseau, G., and Precup, D. (2018). Nonlinear weighted finite automata. In Storkey, A. and Perez-Cruz, F., editors, *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pages 679–688, Playa Blanca, Lanzarote, Canary Islands. PMLR.
- Newmeyer, F. and Preston, L. (2014). *Measuring Grammatical Complexity*. Oxford linguistics. Oxford University Press.
- Petersson, K. M. and Hagoort, P. (2012). The neurobiology of syntax: Beyond string sets. *Philosophical transactions of the Royal Society of London. Series B, Biological sciences*, 367:1971–83.
- Rodriguez, P. and Wiles, J. (1998). Recurrent neural networks can learn to implement symbol-sensitive counting. In *Proceedings of the 1997 Conference on Advances in Neural Information Processing Systems 10, NIPS '97*, pages 87–93, Cambridge, MA, USA. MIT Press.
- Sennhauser, L. and Berwick, R. C. (2018). Evaluating the ability of lstms to learn context-free grammars. *CoRR*, abs/1811.02611.
- Shieber, S. M. (1987). *Evidence Against the Context-Freeness of Natural Language*, pages 320–334. Springer Netherlands, Dordrecht.
- Sipser, M. (2013). *Introduction to the Theory of Computation*. Course Technology, Boston, MA, third edition.
- Skachkova, N., Trost, T., and Klakow, D. (2018). Closing brackets with recurrent neural networks. pages 232–239.

- Suzgun, M., Gehrmann, S., Belinkov, Y., and Shieber, S. M. (2019). LSTM networks can perform dynamic counting. *CoRR*, abs/1906.03648.
- Williams, R. and Zipser, D. (1998). Gradient-based learning algorithms for recurrent networks and their computational complexity.
- Yu, X., Vu, N. T., and Kuhn, J. (2019). Learning the Dyck language with attention-based Seq2Seq models. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 138–146, Florence, Italy. Association for Computational Linguistics.
- Zeng, Z., Goodman, R. M., and Smyth, P. (1994). Discrete recurrent neural networks for grammatical inference. *IEEE transactions on neural networks*, 5 2:320–30.

## Appendix



## Eidesstattliche Erklärung

### Eidesstattliche Erklärung zur <-Arbeit>

Ich versichere, die von mir vorgelegte Arbeit selbstständig verfasst zu haben. Alle Stellen, die wörtlich oder sinngemäß aus veröffentlichten oder nicht veröffentlichten Arbeiten anderer entnommen sind, habe ich als entnommen kenntlich gemacht. Sämtliche Quellen und Hilfsmittel, die ich für die Arbeit benutzt habe, sind angegeben. Die Arbeit hat mit gleichem Inhalt bzw. in wesentlichen Teilen noch keiner anderen Prüfungsbehörde vorgelegen.

*Unterschrift :*                      *Ort, Datum :*

