



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

Topic: Investigating the ability of RNN architectures to learn
context-free grammars by example of Dyck(2)

Author: Fynn Dobler <fynndobler@gmail.com>
Matr.-Nr. 775710

1. Supervisor: Dr. Thomas Hanneforth
2. Supervisor: Dr. Uladzimir Sidarenka

Summary

The capability of three major recurrent neural network (RNN) architectures - SRNN, LSTM and GRU - to learn the underlying structure of the Dyck(2) grammar has been investigated. To assess the influence of such factors as model complexity and training corpus composition, each architecture was instantiated in 27 models, with a model containing n hidden units ($n = \{2^1, 2^2, \dots, 2^9\}$) each being trained on a character level on one of three corpora. In total, 81 models were trained and tested. The corpora consisted of a Baseline corpus to compare to, as well as a corpus containing an increased frequency for long-range dependencies (High LRD) and one with a decreased frequency thereof (Low LRD). The models were evaluated on two classification experiments, designed to assess handling of long-range dependencies and increasing nesting levels. The results were reported in terms of accuracy, precision, recall and F1 score, as well as a closer look at the distribution of false positives. The findings show the models as likely to model Dyck(2) if they were trained on the Low LRD corpus, with LSTMs and GRUs achieving promising results.

Abstract

Formal grammars, specifically context-free grammars (CFGs), are powerful tools with which to model natural languages. In this thesis, the capability of several recurrent neural networks (RNNs) to learn CFGs by proxy of Dyck(2) was investigated. The impact of training corpus composition was assessed by training models on three different corpora of varying complexity. To assess whether Dyck(2) was learned, the networks classified words into belonging or not belonging to Dyck(2) in two experiments designed to test their ability to generalize to both extreme long-range dependencies (LRDs) and unseen nesting depths (NDs). Only few RNNs achieved above-chance accuracy. For the ones that did, low training data complexity facilitated generalization, while high complexity showed an inhibitive effect.

Zusammenfassung

Die Fähigkeit von drei bekannten RNN Architekturen - SRNN, LSTM und GRU - die unterliegende Struktur der Dyck(2)-Grammatik zu lernen wurde untersucht. Die Faktoren von Modellkomplexität und Korpuskomposition wurden berücksichtigt, indem zu jeder Architektur je 27 Modelle trainiert worden sind - je ein Modell mit n versteckten Einheiten ($n = \{2^1, 2^2, \dots, 2^9\}$) wurde auf einem von drei Korpora trainiert. Insgesamt wurden 81 Modelle trainiert und getestet. Die drei Korpora waren ein Baseline Korpus, zu dem die anderen verglichen wurden. Zusätzlich wurde ein Korpus kreiert, in dem sich Abhängigkeiten zwischen einzelnen Buchstaben häufig lange erstrecken (High LRD), sowie ein Korpus wo diese Abhängigkeiten häufig kurz sind (Low LRD). Die Modelle wurden auf zwei verschiedenen Klassifikations-Experimenten evaluiert, welche darauf ausgelegt waren, das Modellverhalten bei langen Abhängigkeiten und tiefer Rekursion zu untersuchen. Die Ergebnisse wurden in Form von Genauigkeit, Precision, Recall und F1 Score angegeben. Zusätzlich wurde die Verteilung von falsch als korrekt klassifizierten Wörtern betrachtet. Die Komplexität des Trainingskorpus' hat sich als großer Einfluss auf die Lernfähigkeit von Modellen erwiesen: Modelle, die auf dem Low LRD Korpus trainiert worden sind, waren erfolgreicher als die Baseline Modelle, während High LRD Modelle weniger erfolgreich waren.

Contents

List of Figures	4
List of Tables	5
List of Listings	5
1 Introduction	6
2 Theoretical Background	7
2.1 Formal Languages and Formal Grammars	7
2.2 Formal Grammars and Natural Language	8
2.2.1 Natural Language as supra-regular	8
2.2.2 Natural Language as supra-context-free	9
2.3 Dyck Languages	10
2.4 Neural Network Architectures	11
2.4.1 Simple RNN	11
2.4.2 LSTM	12
2.4.3 GRU	13
2.5 Related Works	13
3 Experiment Setup	16
3.1 Evaluation	16
3.2 Models	16
3.3 Corpus Construction	17
3.4 Experiment 1: Long-Range Dependency	19
3.5 Experiment 2: New Depths	19
4 Results	20
4.1 Architecture/Training Data	20
4.2 Experiment 1: Long-Range Dependency	21
4.3 Experiment 2: New Depths	21
4.4 Outliers: A Closer Look	22
5 Discussion	31
5.1 Learning D_2	31
5.2 Influence of Training Data	33
6 Conclusion	34
6.1 Further Research	36
Bibliography	37
Appendix	40
Eidesstattliche Erklärung	78

List of Figures

1	The Chomsky Hierarchy	7
2	Illustration of an unfolded RNN	12
3	Illustration of an LSTM Memory Cell	13
4	Illustration of a GRU	14

List of Tables

1	Formal grammar properties.	8
2	Corpus sizes in current works	14
3	Reported values for performance in previous works	15
4	Overview of investigated models	15
5	Properties of training corpora	17
6	Performance measures regardless of training corpus	23
7	Experiment 1: Base LRD network performance	24
8	Experiment 1: Low LRD network performance	25
9	Experiment 1: High LRD network performance	26
10	Experiment 2: Base LRD network performance	27
11	Experiment 2: Low LRD network performance	28
12	Experiment 2: High LRD network performance	29
13	Ratio of open/closed error categories misclassified as false positives. . .	30

Index of Listings

1	generate_raw_data.py	40
2	corpus_tools.py	44
3	RNN_classifier.py	58
4	concat_results.py	66
5	analyze_performance.py	70

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 industry¹. 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,

¹<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 Σ^* over an alphabet Σ , 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, Σ is a set of terminal symbols (alphabet), R is a set of rules of the form $\alpha \rightarrow \beta$ (where α and β 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.

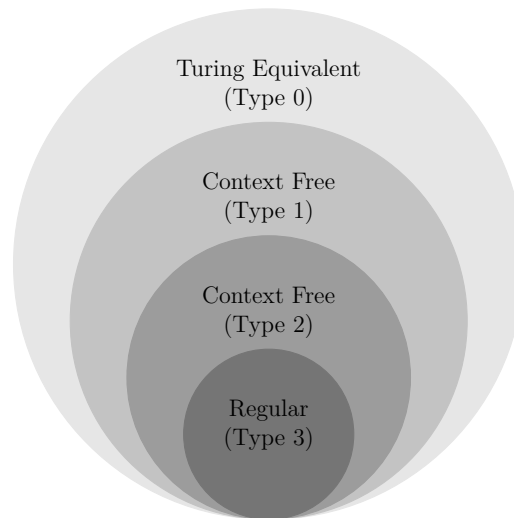


Figure 1: A visual representation of the Chomsky Hierarchy.

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, α, β, γ are strings of terminal and non-terminal symbols, and x is a string of terminal symbols. α, β and γ 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.²

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.'

²The 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.³

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

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

³The 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 (1995)), 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.

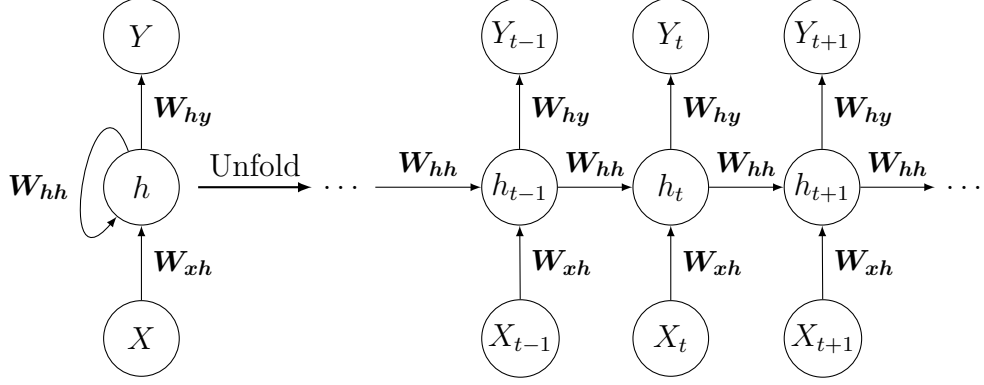


Figure 2: An RNN, unfolded through time.

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 σ 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.

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

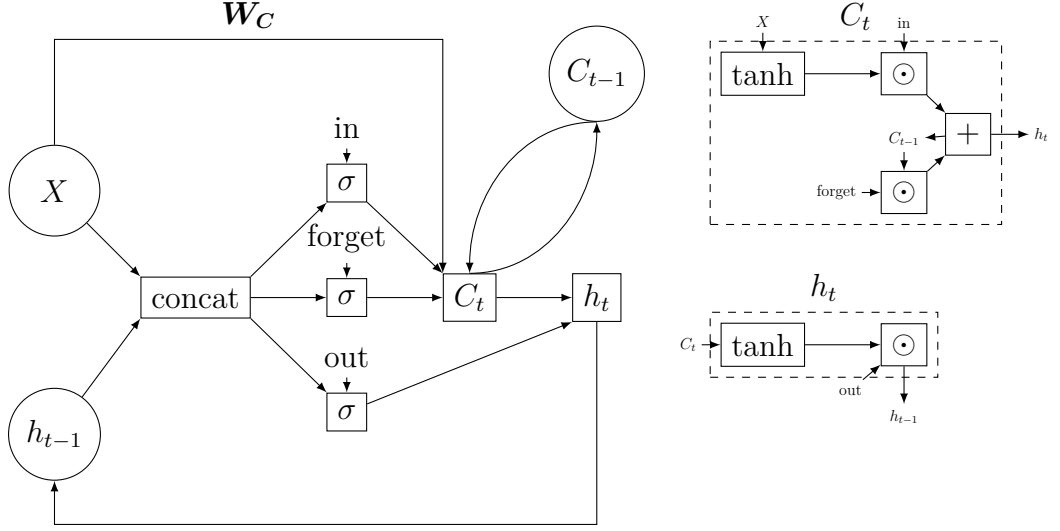


Figure 3: An LSTM memory cell.

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, σ is the squashing sigmoid function, \mathbf{W} and \mathbf{U} are learned gate-dependent weight matrices and ϕ 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),

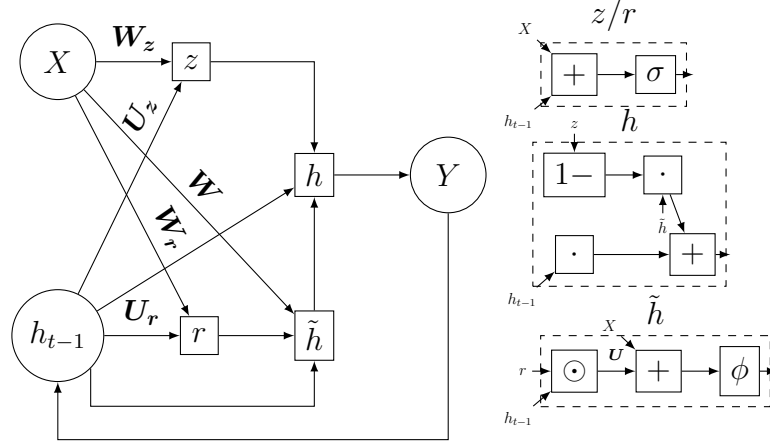


Figure 4: Illustration of a GRU.

<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 Evaluation

Generating the words for the two experiments follows the procedure described by Bernardy (2018) and will be explained in depth in the coming sections. Whereas the research put forth in his paper scrutinized the generative abilities of RNNs, I am investigating the models performing a classification task. As such, my training, test, validation and experiment data all consist of the same 1:1 ratio of correct-to-incorrect words. Within the incorrect words, a distinction between superfluous opening or closing brackets is made, also at a ratio of 1:1.

As such, a random guessing strategy would yield a baseline accuracy of 50%. A model is considered as having learned useful features from the training data if it scores above the baseline accuracy in the experiments. Furthermore, if the model has learned the underlying grammar of D_2 , there should be no difference in accuracy for the two classes of incorrect words, as both of them do not belong to D_2 , regardless of which bracket is replaced.

3.2 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 ≥ 0.5 means the model classified the input as a correct word. The models were implemented in Tensorflow 2.0. All code was implemented using Python 3.7.6, NumPy 1.18.1 (McKinney, 2011), pandas 0.25.3 (Oliphant, 2006) and scikit-learn 0.22.1 (Pedregosa et al., 2011).⁴

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

⁴The full source code can be found in the Appendix or under <https://github.com/FyDob/BSc-Thesis>.

<i>Corpus</i>	<i>Word Length</i>	<i>maxND</i>	<i>maxBD</i>
Baseline	18.37 (6.36)	4.31 (1.22)	13.00 (16.02)
High LRD	18.67 (4.75)	5.12 (1.04)	16.67 (4.75)
Low LRD	17.54 (8.04)	3.92 (0.98)	10.58 (9.02)

Table 5: Properties of the three corpora the models were trained on, reported in averages (variance in brackets).

units \times number of training corpora \times number of different architectures) evaluated models.

3.3 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 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 5 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_{branch} , r_{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 modified⁵ by first identifying 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.⁶ 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

⁵This cut-off point was chosen as it significantly reduces the average maximum bracket distance without creating too many duplicates.

⁶The same considerations as for the Low LRD corpus cut-off apply.

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.

3.4 Experiment 1: Long-Range Dependency

For this experiment, the test set consisted of 1,000,000 samples of length $1+18+18+1 = 38$, half of which were correct Dyck words. They were created by picking two random Dyck words $w_1, w_2 \in D_2^{=18}$ from the base corpus, concatenating them and wrapping the result in a randomly selected pair of matching brackets as follows:

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

To generate incorrect samples, the generated correct LRD words were corrupted in the same way as for the training corpora, yielding 250,000 incorrect LRD words with a superfluous opening or closing bracket each.

While w_1 and w_2 might have been seen in training (for models trained on the base corpus), 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 classification accuracy serves as a strong indication of the model having learned to generalize to longer, non-concatenated Dyck words.⁷ I report model performance on Experiment 1 in terms of accuracy, precision, recall and F1 score.

3.5 Experiment 2: New Depths

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 confounded as little as possible by LRD performance.

For this task, the test set consisted of 1,000,000 samples of length 30, half of which were correct Dyck words. First, 500,000 correct words were chosen at random from the base corpus. Then, they were wrapped by a prefix of five randomly chosen opening brackets and a suffix of the corresponding closing brackets as follows:

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

Generation of incorrect samples was done in accordance to Experiment 1 and corpus

⁷While the infixed sub-words are indeed concatenated, w_{LRD} cannot be created by concatenating two shorter words due to being wrapped by a matching bracket pair.

creation.

This process still has the model extrapolate beyond the length of the training words, while increasing all present nesting depths by 5. This is analogous to center embedding in natural language - processing increasing nesting levels is more complicated than processing a flat structure. A high classification accuracy in Experiment 2 indicates a capability to generalize to repeated application of grammar rules beyond what was seen in the training set. As such, it implies an understanding of the D_2 grammar. I report model performance on Experiment 2 in terms of accuracy, precision, recall and F1 score.

4 Results

During training, almost all 81 trained models achieved a validation accuracy near 100%, except for the SRNN-2 models trained on the base and low LRD corpus, which scored 75.0% and 50.4% respectively. I have included them in the experiments regardless of their low validation accuracy, since it was unclear whether validation accuracy would be a strong predictor for a network’s performance on the experiment data. I present my results with regard to three focus points: First, the overall performance of different architectures with respect to which corpus they were trained on, then the individual model performances on each of the two experiments, and finally a closer look at classifications made by outlier networks - networks which drastically over- or underperformed in either of the experiments - with regards to misclassified false positives.

4.1 Architecture/Training Data

As can be seen in Table 6, none of the architectures consistently achieved an accuracy far above the random guessing baseline of 50.0%. However, there was still a notable difference in performance between architectures: on average, the GRU networks scored the highest on accuracy and precision, while the SRNN networks achieved the best recall and F1 score. With 51.5%, LSTMs scored an average accuracy between SRNNs (50.0%) and GRUs (53.3%), but they underperformed in all other experiment measures.

Furthermore, the choice of training data had a notable effect on overall model performance: SRNNs and GRUs received a boost in performance in all measures when comparing the Base to the Low LRD models, elevating SRNNs from an accuracy below random guessing to 51.4%. While LSTMs lost 1.4% in terms of accuracy, all other performance measures improved significantly for the Low LRD models. Training on the High LRD corpus aided SRNNs in terms of accuracy, precision and F1 score, but worsened accuracy, recall and F1 score for LSTMs and GRUs.

4.2 Experiment 1: Long-Range Dependency

I report results for Experiment 1 in Tables 7, 8 and 9, which include the performance measures for all networks trained on the Base, Low and High LRD corpus respectively, as evaluated on Experiment 1.

When trained on the Base corpus, 8 of 27 models (29.6%) achieved an accuracy above random guessing. Among those 8, only 2 reached an accuracy above 60%: LSTM-8 and GRU-2. The vast majority of models - 20 in total - reached an accuracy of $50 \pm 5\%$. 4 performed even worse than that: the worst model (SRNN-16) only achieved 27.2% accuracy on the experiment data. There was no apparent correlation between validation accuracy and model performance in Experiment 1 - indeed, SRNN-2 with the lowest validation accuracy at 75.0% evaluated at below chance, but so did several models with a validation accuracy of 100%.

12 of 27 Low LRD-trained models (44.4%) scored an accuracy higher than the random guessing baseline. Among those, 3 models - SRNN-4, LSTM-16 and GRU-64 - reached an accuracy above 60%. While 22 models fell within the $50 \pm 5\%$ belt of accuracy, only 1 model performed significantly worse: LSTM-4 with 27.8%. Validation accuracy was entirely unrelated to model performance, with LSTM-4 having achieved a perfect score on the validation data, but completely failing at Experiment 1. On average, all measures have improved when compared to the Base models: accuracy improved by +2.7 percentage points (p.p.), precision by +11.0 p.p., recall by +9.5 p.p. and F1 score by +8.3 p.p.

9 of 27 High LRD-trained models (33.3%) achieved an accuracy above baseline, but only GRU-512 by a significant margin with 80%. 23 models performed within the $\pm 5\%$ margin around the baseline, and 3 models (LSTM-4, LSTM-8 and GRU-8) underperformed significantly. There was no relation between validation accuracy and performance on experiment data for High LRD models, either. Compared to the Base models, High LRD models almost consistently scored worse: accuracy, recall and F1 score went down by -1.0 , -2.2 and -1.8 p.p. respectively, but precision was improved by +5.6 p.p.

4.3 Experiment 2: New Depths

All results for Experiment 2 can be found in Tables 10, 11 and 12. When comparing mean performances across all models, they largely scored higher on Experiment 2 than Experiment 1. This suggests that Experiment 2 was easier regardless of training data. As with Experiment 1, a model’s validation accuracy did not correlate with its performance on the experiment data.

Among the Base models, 8 of 27 (29.6%) performed above guessing baseline in Experiment 2 - the same ratio as for Experiment 1, though there was minimal overlap in the best performers. Only GRU-64 and SRNN-64 performed above 50% accuracy

for both experiments. 20 of 27 models stayed in the $\pm 5\%$ margin of the baseline, with only 2 (SRNN-128 and SRNN-256) dropping below that. The best performing network - LSTM-128 - scored the highest accuracy across all models and all experiments with 99.3%.

The same number of Low LRD models performed above chance for Experiment 2 as for Experiment 1 (44.4%), with 6 models (SRNN-2, SRNN-8, LSTM-16, LSTM-512, GRU-4 and GRU-64) occurring in both groups. LSTM-512 and GRU-64 have both achieved an accuracy $> 90\%$. 19 models performed within the $\pm 5\%$ margin around the baseline, only 2 dropped below that. Compared to the Base models, low LRD models also performed better on Experiment 2: mean accuracy is up by +0.4 p.p., precision by +20.6, recall by +3.5 and F1 score by 6.5 p.p.

While the highest number of High LRD models have achieved an accuracy above random guessing - 10 of 27, or 37.0% - only SRNN-16 crossed the 60% threshold at all. Indeed, aside from SRNN-16, only one other model lay outside of the $\pm 5\%$ margin around the baseline - LSTM-64, with an accuracy of 28.5%. When comparing with the Base models, all measures except for precision, which improved by +10.7 p.p. Accuracy went down by -2.7 p.p., recall by -7.5 and F1 score by -5.0 p.p.

4.4 Outliers: A Closer Look

As is evident from the results described so far, a vast majority of models fell within a $\pm 5\%$ margin around random guessing in terms of accuracy. I consider these models to not have extracted any useful grammar information from the training data and discard them for further investigation. As such, I will only discuss models with an accuracy either $> 55\%$ or $< 45\%$, to both gather information on what caused models to succeed, and what caused them to fail.

	<i>Accuracy</i>	<i>Precision</i>	<i>Recall</i>	<i>F1 Score</i>	<i>Val Acc</i>
SRNN					
Base					
Mean	<i>0.476</i>	<i>0.278</i>	<i>0.181</i>	<i>0.193</i>	0.972
Variance	0.064	0.235	0.257	0.220	0.081
SRNN					
Low LRD					
Mean	0.514	0.443	0.288	0.316	<i>0.932</i>
Variance	0.055	0.176	0.259	0.213	0.158
SRNN					
High LRD					
Mean	0.511	0.415	0.153	0.200	0.981
Variance	0.043	0.211	0.171	0.187	0.054
SRNN					
Complete					
Mean	<i>0.500</i>	0.379	0.208	0.236	<i>0.962</i>
Variance	0.056	0.217	0.236	0.211	0.107
LSTM					
Base					
Mean	0.543	<i>0.219</i>	0.148	0.154	0.999
Variance	0.173	0.365	0.335	0.341	0.003
LSTM					
Low LRD					
Mean	0.529	0.391	0.176	0.195	0.999
Variance	0.158	0.295	0.311	0.302	0.002
LSTM					
High LRD					
Mean	<i>0.472</i>	0.258	<i>0.036</i>	<i>0.059</i>	1.000
Variance	0.075	0.278	0.057	0.090	0.000
LSTM					
Complete					
Mean	0.515	<i>0.289</i>	<i>0.120</i>	<i>0.136</i>	0.999
Variance	0.143	0.318	0.268	0.269	0.002
GRU					
Base					
Mean	0.531	<i>0.371</i>	0.147	0.185	0.999
Variance	0.097	0.297	0.226	0.247	0.001
GRU					
Low LRD					
Mean	0.554	0.507	0.206	0.243	0.999
Variance	0.133	0.198	0.302	0.285	0.001
GRU					
High LRD					
Mean	<i>0.514</i>	0.439	<i>0.126</i>	<i>0.170</i>	<i>0.972</i>
Variance	0.082	0.223	0.198	0.193	0.055
GRU					
Complete					
Mean	0.533	0.439	0.160	0.200	0.990
Variance	0.106	0.245	0.244	0.242	0.034

Table 6: Performance measures of all architectures across both experiments depending on which corpus they were trained on, as well as the compounded measures for all networks regardless of training data.

<i>Network</i>	<i>Accuracy</i>	<i>Precision</i>	<i>Recall</i>	<i>F1 Score</i>	<i>Val Acc</i>
GRU-2	0.890	0.982	0.794	0.878	0.995
GRU-4	0.487	0.413	0.060	0.105	1.000
GRU-8	0.488	0.331	0.023	0.043	1.000
GRU-16	0.550	0.713	0.168	0.272	1.000
GRU-32	0.434	0.267	0.075	0.118	1.000
GRU-64	0.510	0.537	0.149	0.234	1.000
GRU-128	0.553	0.611	0.293	0.396	0.999
GRU-256	0.497	0.364	0.009	0.018	1.000
GRU-512	0.500	0.487	0.007	0.015	1.000
LSTM-2	0.500	0.000	0.000	0.000	0.999
LSTM-4	0.451	0.000	0.000	0.000	1.000
LSTM-8	0.910	0.959	0.857	0.905	1.000
LSTM-16	0.343	0.001	0.000	0.000	1.000
LSTM-32	0.500	0.000	0.000	0.000	1.000
LSTM-64	0.505	0.596	0.030	0.057	1.000
LSTM-128	0.347	0.001	0.000	0.000	1.000
LSTM-256	0.455	0.128	0.016	0.028	1.000
LSTM-512	0.499	0.468	0.009	0.017	0.991
SRNN-2	0.465	0.357	0.087	0.140	<i>0.750</i>
SRNN-4	0.488	0.334	0.023	0.043	1.000
SRNN-8	0.461	0.037	0.003	0.006	1.000
SRNN-16	<i>0.272</i>	0.021	0.010	0.014	1.000
SRNN-32	0.498	0.492	0.141	0.219	1.000
SRNN-64	0.504	0.505	0.366	0.424	1.000
SRNN-128	0.484	0.040	0.001	0.003	1.000
SRNN-256	0.503	0.503	0.455	0.478	1.000
SRNN-512	0.484	0.017	0.001	0.001	1.000
Mean	0.503	0.339	0.132	0.163	0.990
Variance	0.130	0.294	0.233	0.253	0.048

Table 7: Performance measures for Experiment 1 of all networks that were trained on the Base LRD corpus.

<i>Network</i>	<i>Accuracy</i>	<i>Precision</i>	<i>Recall</i>	<i>F1 Score</i>	<i>Val Acc</i>
GRU-2	0.583	0.582	0.589	0.585	1.000
GRU-4	0.505	0.518	0.135	0.214	1.000
GRU-8	0.498	0.491	0.106	0.175	0.997
GRU-16	0.503	0.528	0.060	0.107	1.000
GRU-32	0.505	0.525	0.100	0.168	1.000
GRU-64	0.890	0.869	0.918	0.893	1.000
GRU-128	0.484	0.324	0.030	0.054	1.000
GRU-256	0.500	0.496	0.039	0.072	1.000
GRU-512	0.500	0.435	0.002	0.004	0.997
LSTM-2	0.500	0.500	0.475	0.487	0.996
LSTM-4	<i>0.278</i>	<i>0.000</i>	<i>0.000</i>	<i>0.000</i>	<i>1.000</i>
LSTM-8	0.500	0.000	0.000	0.000	1.000
LSTM-16	0.881	0.853	0.920	0.885	1.000
LSTM-32	0.491	0.092	0.002	0.004	1.000
LSTM-64	0.500	0.505	0.010	0.020	1.000
LSTM-128	0.496	0.374	0.011	0.021	0.998
LSTM-256	0.500	0.535	0.007	0.014	1.000
LSTM-512	0.507	0.566	0.058	0.106	1.000
SRNN-2	0.501	0.501	0.366	0.423	0.504
SRNN-4	0.708	0.646	0.922	0.760	0.930
SRNN-8	0.508	0.527	0.150	0.234	1.000
SRNN-16	0.500	0.000	0.000	0.000	1.000
SRNN-32	0.486	0.388	0.050	0.088	1.000
SRNN-64	0.487	0.375	0.040	0.072	1.000
SRNN-128	0.492	0.484	0.236	0.318	0.953
SRNN-256	0.501	0.501	0.412	0.452	1.000
SRNN-512	0.503	0.503	0.497	0.500	1.000
Mean	0.530	0.449	0.227	0.246	0.977
Variance	0.120	0.217	0.304	0.280	0.096

Table 8: Performance measures for Experiment 1 of all networks that were trained on the Low LRD corpus.

<i>Network</i>	<i>Accuracy</i>	<i>Precision</i>	<i>Recall</i>	<i>F1 Score</i>	<i>Val Acc</i>
GRU-2	0.494	0.480	0.136	0.212	1.000
GRU-4	0.499	0.482	0.036	0.066	1.000
GRU-8	0.342	0.000	0.000	0.000	1.000
GRU-16	0.523	0.590	0.153	0.242	1.000
GRU-32	0.504	0.521	0.093	0.158	1.000
GRU-64	0.504	0.522	0.088	0.150	1.000
GRU-128	0.489	0.419	0.056	0.098	1.000
GRU-256	0.503	0.549	0.035	0.066	0.902
GRU-512	0.800	0.768	0.861	0.812	0.849
LSTM-2	0.483	0.390	0.060	0.104	1.000
LSTM-4	<i>0.282</i>	<i>0.000</i>	<i>0.000</i>	<i>0.000</i>	1.000
LSTM-8	0.379	0.000	0.000	0.000	1.000
LSTM-16	0.502	0.640	0.007	0.015	1.000
LSTM-32	0.498	0.475	0.038	0.070	1.000
LSTM-64	0.500	0.000	0.000	0.000	1.000
LSTM-128	0.492	0.365	0.022	0.041	1.000
LSTM-256	0.497	0.367	0.008	0.016	1.000
LSTM-512	0.500	0.000	0.000	0.000	0.999
SRNN-2	0.500	0.499	0.042	0.078	<i>0.833</i>
SRNN-4	0.502	0.523	0.046	0.085	0.996
SRNN-8	0.539	0.622	0.201	0.304	1.000
SRNN-16	0.500	0.000	0.000	0.000	0.999
SRNN-32	0.500	0.497	0.037	0.070	1.000
SRNN-64	0.503	0.504	0.360	0.420	1.000
SRNN-128	0.497	0.494	0.235	0.318	1.000
SRNN-256	0.496	0.495	0.340	0.403	1.000
SRNN-512	0.490	0.459	0.112	0.180	1.000
Mean	0.493	0.395	0.110	0.145	0.984
Variance	0.083	0.230	0.180	0.182	0.045

Table 9: Performance measures for Experiment 1 of all networks that were trained on the High LRD corpus.

<i>Network</i>	<i>Accuracy</i>	<i>Precision</i>	<i>Recall</i>	<i>F1 Score</i>	<i>Val Acc</i>
GRU-2	0.500	0.000	0.000	0.000	0.995
GRU-4	0.563	0.655	0.265	0.377	1.000
GRU-8	0.472	0.164	0.014	0.025	1.000
GRU-16	0.500	0.000	0.000	0.000	1.000
GRU-32	0.500	0.000	0.000	0.000	1.000
GRU-64	0.515	0.550	0.160	0.248	1.000
GRU-128	0.496	0.000	0.000	0.000	0.999
GRU-256	0.500	0.000	0.000	0.000	1.000
GRU-512	0.601	0.597	0.620	0.608	1.000
LSTM-2	0.500	0.000	0.000	0.000	0.999
LSTM-4	0.500	0.000	0.000	0.000	1.000
LSTM-8	0.500	0.000	0.000	0.000	1.000
LSTM-16	0.500	0.000	0.000	0.000	1.000
LSTM-32	0.778	0.792	0.755	0.773	1.000
LSTM-64	0.500	0.000	0.000	0.000	1.000
LSTM-128	0.993	0.991	0.995	0.993	1.000
LSTM-256	0.500	0.000	0.000	0.000	1.000
LSTM-512	0.499	0.000	0.000	0.000	0.991
SRNN-2	0.580	0.548	0.917	0.686	<i>0.750</i>
SRNN-4	0.488	0.000	0.000	0.000	1.000
SRNN-8	0.500	0.000	0.000	0.000	1.000
SRNN-16	0.500	0.000	0.000	0.000	1.000
SRNN-32	0.515	0.564	0.136	0.220	1.000
SRNN-64	0.511	0.529	0.192	0.281	1.000
SRNN-128	0.441	0.409	0.266	0.322	1.000
SRNN-256	<i>0.382</i>	0.156	0.053	0.080	1.000
SRNN-512	0.497	0.497	0.615	0.550	1.000
Mean	0.531	0.239	0.185	0.191	0.990
Variance	0.113	0.311	0.308	0.290	0.048

Table 10: Performance measures for Experiment 2 of all networks that were trained on the Base LRD corpus.

<i>Network</i>	<i>Accuracy</i>	<i>Precision</i>	<i>Recall</i>	<i>F1 Score</i>	<i>Val Acc</i>
GRU-2	0.494	0.488	0.222	0.306	1.000
GRU-4	0.561	0.651	0.262	0.374	1.000
GRU-8	0.499	0.492	0.056	0.101	0.997
GRU-16	0.497	0.465	0.039	0.072	1.000
GRU-32	0.500	0.000	0.000	0.000	1.000
GRU-64	0.935	0.910	0.966	0.937	1.000
GRU-128	0.516	0.560	0.144	0.229	1.000
GRU-256	0.499	0.486	0.043	0.078	1.000
GRU-512	0.497	0.305	0.004	0.008	0.997
LSTM-2	0.499	0.000	0.000	0.000	0.996
LSTM-4	0.480	0.328	0.037	0.067	1.000
LSTM-8	0.555	0.660	0.226	0.337	1.000
LSTM-16	0.566	0.662	0.270	0.384	1.000
LSTM-32	0.327	0.207	0.122	0.154	1.000
LSTM-64	0.491	0.301	0.014	0.026	1.000
LSTM-128	0.500	0.000	0.000	0.000	0.998
LSTM-256	0.502	0.538	0.027	0.051	1.000
LSTM-512	0.953	0.917	0.996	0.955	1.000
SRNN-2	0.501	0.501	0.585	0.540	0.504
SRNN-4	0.442	0.361	0.149	0.211	0.930
SRNN-8	0.529	0.582	0.206	0.304	1.000
SRNN-16	0.498	0.000	0.000	0.000	1.000
SRNN-32	0.580	0.561	0.738	0.637	1.000
SRNN-64	0.514	0.546	0.170	0.260	1.000
SRNN-128	0.498	0.495	0.178	0.262	0.953
SRNN-256	0.500	0.499	0.348	0.410	1.000
SRNN-512	0.501	0.505	0.136	0.214	1.000
Mean	0.535	0.445	0.220	0.256	0.977
Variance	0.126	0.246	0.281	0.263	0.096

Table 11: Performance measures for Experiment 2 of all networks that were trained on the Low LRD corpus.

<i>Network</i>	<i>Accuracy</i>	<i>Precision</i>	<i>Recall</i>	<i>F1 Score</i>	<i>Val Acc</i>
GRU-2	0.535	0.580	0.253	0.353	1.000
GRU-4	0.502	0.518	0.070	0.123	1.000
GRU-8	0.500	0.000	0.000	0.000	1.000
GRU-16	0.492	0.269	0.010	0.019	1.000
GRU-32	0.511	0.552	0.114	0.189	1.000
GRU-64	0.502	0.512	0.081	0.140	1.000
GRU-128	0.541	0.601	0.244	0.347	1.000
GRU-256	0.503	0.530	0.048	0.089	0.902
GRU-512	0.500	0.000	0.000	0.000	0.849
LSTM-2	0.531	0.591	0.200	0.299	1.000
LSTM-4	0.500	0.000	0.000	0.000	1.000
LSTM-8	0.500	0.000	0.000	0.000	1.000
LSTM-16	0.517	0.658	0.070	0.127	1.000
LSTM-32	0.511	0.554	0.112	0.187	1.000
LSTM-64	0.285	0.000	0.000	0.000	1.000
LSTM-128	0.500	0.000	0.000	0.000	1.000
LSTM-256	0.521	0.608	0.121	0.202	1.000
LSTM-512	0.500	0.000	0.000	0.000	0.999
SRNN-2	0.491	0.333	0.017	0.033	0.833
SRNN-4	0.500	0.000	0.000	0.000	0.996
SRNN-8	0.500	0.000	0.000	0.000	1.000
SRNN-16	0.675	0.691	0.633	0.661	0.999
SRNN-32	0.507	0.557	0.072	0.128	1.000
SRNN-64	0.492	0.473	0.139	0.215	1.000
SRNN-128	0.485	0.281	0.019	0.036	1.000
SRNN-256	0.501	0.502	0.322	0.392	1.000
SRNN-512	0.512	0.534	0.183	0.273	1.000
Mean	0.504	0.346	0.100	0.141	0.984
Variance	0.056	0.267	0.141	0.164	0.045

Table 12: Performance measures for Experiment 2 of all networks that were trained on the High LRD corpus.

<i>Architecture</i>	<i>Experiment</i>	<i>Training Corpus</i>	<i>Open</i>	<i>Closed</i>	<i>Ratio</i>	<i>Total</i>
High Accuracy Models						
LSTM	LRD	base	18,161	0	inf	18,161
GRU	LRD	base	55,256	45,207	1.222	100,463
SRNN	LRD	low	101,557	150,905	0.673	252,462
LSTM	LRD	low	39,584	39,559	1.001	79,143
GRU	LRD	low	165,820	115,236	1.439	281,056
GRU	LRD	high	43,497	86,753	0.501	130,250
SRNN	ND	base	138,673	240,098	0.578	378,771
LSTM	ND	base	103,816	0	inf	103,816
GRU	ND	base	244,265	34,843	7.010	279,108
SRNN	ND	low	94,602	194,386	0.487	288,988
LSTM	ND	low	97,527	74,625	1.307	172,152
GRU	ND	low	82,933	34,981	2.371	117,914
SRNN	ND	high	90,850	50,458	1.801	141,308
Low Accuracy Models						
SRNN	LRD	base	228,451	4,387	52.075	232,838
LSTM	LRD	base	298,547	11,358	26.285	309,905
GRU	LRD	base	29,908	73,284	0.408	103,192
LSTM	LRD	low	222,045	0	inf	222,045
LSTM	LRD	high	338,852	387	875.587	339,239
GRU	LRD	high	0	157,850	0.000	157,850
SRNN	ND	base	307,497	29,619	10.382	337,116
SRNN	ND	low	116,613	15,213	7.665	131,826
LSTM	ND	low	234,224	0	inf	234,224
LSTM	ND	high	214,980	0	inf	214,980

Table 13: Ratio of open/closed error categories misclassified as false positives.

First, I report the number of the two error categories (superfluous open bracket, superfluous closed bracket) per network architecture in Tables 13 and. Hypothetically, since both categories were balanced in training and experiment data, there should be no significant difference between misclassifying one or the other - unless one category proves to be more complex to an architecture. The ratio of open-to-closed bracket misclassifications serves as a simple indication of whether the model extracted the correct information - that the amount of brackets needs to be balanced throughout a valid word - at all. A ratio of 1 implies no model difference between the two categories. Conversely, any skewing above or below 1 shows the model disproportionately struggling with one of the two categories. Only one model achieved a near perfect ratio: LSTM-16 trained on the Low LRD corpus, evaluating the Experiment 1 data.

Wrong words with superfluous open brackets have shown to be the most difficult to reject: 61.54% of the High Accuracy Models and 80% of the Low Accuracy Models skew towards a ratio > 1 . Furthermore, a model's inability to treat both error categories the same as indicated by the ratio correlates to model accuracy: will most High Accuracy Models err reasonably closely to 1, the Low Accuracy Models show a much more extreme distribution, frequently misclassifying almost all words in one error category, while handling the other one perfectly. There is no trend towards one architecture overall skewing towards 1. However, among the High Accuracy Models, the models trained on the Low LRD corpus have the lowest deviation from 1, with 0.256 on Experiment 1, and 0.652 on Experiment 2. The lowest deviation from 1 among the Low Accuracy Models comes from the models trained on the Base corpus, with 25.65 and 9.382, respectively.

5 Discussion

The results presented in Chapter 4 provide an in-depth look at model performance, reporting on established measures like accuracy and F1 score, but also on custom measures for this task, like misclassification ratio for the two categories of incorrect words. All of this was done to answer the question of whether RNNs can learn the underlying structure of D_2 , and whether specific properties of the training data can facilitate or inhibit that ability. I will now interpret the results from the perspective of these questions.

5.1 Learning D_2

Both experiments were designed to test the two most important aspects of what it means to generalize from a small language subset to more complex data: interpreting extreme long-range dependencies on long, unseen words, and handling unseen nesting depths. The former showcases the ability of generalizing to much greater length, while

the latter provides insight on how well the model handles deeper, more complex parse trees.

Judging from the average accuracy across all models reported in Tables 7, 8, 9, 10, 11 and 12, Experiment 1 was a harder task than Experiment 2, regardless of training corpus. By and large, models were struggling to correctly classify extreme long-range dependencies, while performing well on a deeper nesting task with shorter dependencies. Assuming that nesting depth, as it creates deeper and more complicated underlying parse trees than low-nesting depth long-range dependencies, would pose a higher difficulty to models that have learned to build a structural representation and as such, learned D_2 , this discrepancy suggests that most models did not achieve such a representation.

In conclusion, the experiments posed in this thesis were hard tasks. Not many models learned helpful information from the training data, some extracted the wrong kind of information, but most learned nothing to help them improve from random guessing. Successful model hidden unit counts range from 2 to 512, with no general trend connecting number of hidden units and model performance, making it impossible to assess lower and upper limits of model complexity required to learn D_2 .

However, 20 models in total have performed above chance on either of the two experiments. Accuracy ranges from 80.0% (High GRU-512) to 91.0% (Base LSTM-8) for Experiment 1, and 67.5% (High SRNN-16) to 99.3% (Base LSTM-128) for Experiment 2. Regardless of experiment and training data, GRUs slightly outperformed LSTMs and SRNNs with an average accuracy of 53.5% (see Table 6).

Individual model performance provides another intriguing implication: Several models that performed well in Experiment 1 barely achieved the baseline in Experiment 2 and vice-versa. This points to the tasks, despite both making a case for generalizability, posing different requirements to the representations the models learned. In fact, only two models performed above baseline for both experiments: Low LSTM-16 (88.1% and 56.66% for Experiment 1 and 2, respectively) and Low GRU-64 (89.0% and 93.5%). Considering there were 20 successful models in total, this number is fairly low, further pointing to learning D_2 being a difficult undertaking. Both of these models have been trained on the Low LRD corpus, indicating that the corpus facilitates all-purpose generalization.

To assess whether a High Accuracy model learned a valid representation of D_2 , incorrect words were split in two equal sized classes: words with superfluous open or closed brackets. A good model, then, should not make a difference between the two classes. While the High Accuracy models did not always succeed at that, they came a lot closer to treating both classes equally than the Low Accuracy models. Furthermore, the imbalance in false positive misclassification was skewed towards words with superfluous open brackets. This might be owed to the fact that an extra closed bracket at any position immediately renders a word of any length ungrammatical: it

resolves a dependency that does not exist. An extra open bracket, however, opens a dependency that might be resolved at a later point in the word. Indeed, all incorrect open bracket words are substrings of longer correct D_2 words. The same cannot be said for incorrect closed bracket words.⁸

In conclusion, the combination of task difficulty and volatile model performance makes it difficult to conclusively compare the three architectures. While GRUs achieved a high accuracy more consistently, LSTMs have produced the highest scoring models. Despite SRNNs scoring the worst in terms of overall accuracy, they - on average - outperformed the other two architectures in terms of F1 score. However, given the measure of misclassification ratio, the record high accuracy and fairly low number of hidden units in well-performing models, LSTMs have shown great promise to be capable of learning D_2 in this study.

5.2 Influence of Training Data

Whether comparing accuracy, precision, recall and F1 score across both experiments (see Table 6, individually comparing model instances based on what corpus they were trained on (i.e. Tables 7, 8 and 9 for Experiment 1) or looking for a trend in open/-closed misclassification ratio (see Table 13): the effect of training corpus complexity is consistent. Disregarding experiments, SRNN and GRU performance improved in all measures when trained on the Low LRD corpus. For LSTMs, it improved every measure but accuracy. When distinguishing between both experiments, models trained on the Low LRD corpus performed the best on average. Furthermore, High Accuracy Models showed the lowest deviation from the ideal 1 to 1 open/closed misclassification ratio if they were trained on the Low LRD corpus. Finally, Low LRD trained models constitute the majority class of High Accuracy Models. On the other hand, models trained on the High LRD corpus on average underperformed compared to the Base models, regardless of model architecture and experiment.

Overall, there is strong evidence for the Low LRD training corpus leading to the most consistently good results, while the High LRD corpus tends to worsen models. This is surprising: considering how the High LRD corpus is constructed, a sizeable portion of training data is similar to the experiment data, featuring at least one long-range dependency spanning the whole word. Frequently encountering this pattern in training as well as having the model encounter more complex structures seemed likely to boost model performance. Instead, the opposite is true: learning from a less complex corpus enhanced robustness and the ability to generalize, while complex training data inhibited these processes.

⁸It must be noted that this property makes the incorrect open words in no way more valid than the incorrect closed words. The experimental stimuli are of a fixed length, and their end is - like in the training data - signified by an end-of-word symbol. There is no reason for a model to anticipate an extension to resolve the open dependency.

These results imply not only that RNNs generalize from limited data to longer, more complex examples, they also do so by extracting generative rules from an underlying structure. This process is facilitated by giving the RNNs simple examples: once the connection is made that an open bracket must eventually be closed by its corresponding closing bracket, but not before more deeply nested pairs have been closed, the RNNs do not have to explicitly learn that the principle holds true at any nesting depth and at any character distance. Conversely, training RNNs mostly on complex words leads to an *overspecification* of the learned rules: the model rarely encounters the underlying principle of D_2 in a simple form and might assume that it only holds true for specific nesting depths or distances, leading to a needlessly - and inaccurately - complex rule set. This interpretation is similar to what Zeng et al. (1993) have found when experimenting with incrementally increasing the length of strings in training: Analog RNNs have been shown to learn 'soft' solutions that are then incrementally hardened as more restrictions are necessary.

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.

To assess model performance, I have adapted the two experiments proposed by Bernardy (2018) for a classification task. They were explicitly designed to investigate model performance on long-range dependencies, as well as a model's ability to generalize to deeper nesting depths. I have reported accuracy, precision, recall and the compound measure of F1 score across models and experiments, as well as providing a closer look at common error sources for the models.

In addition to assessing three different architectures on two experiments, I have investigated the impact of hidden unit number and training corpus composition on model performance. For the former, each architecture was implemented in 9 different models with hidden unit number $n \in \{2^1, 2^2, \dots, 2^9\}$. All models have been trained with the same hyperparameters. To achieve the latter, I have constructed three training corpora: a baseline corpus (Base), a corpus containing words with a high nesting depth and maximum bracket distance (High LRD) and a corpus containing words with a low nesting depth and maximum bracket distance (Low LRD). All corpora contained 1,000,000 words, 500,000 of which were correct. The incorrect words consisted of 250,000 words with extra open and 250,000 words with extra closing brackets. In total, 81 models were trained, and each model was evaluated on both experiment data sets.

The results of both experiments show learning D_2 to be a task of not trivial and

highly volatile difficulty: the number of hidden units in models with an accuracy well above random guessing ranges from 2 to 512. A vast majority of models failed to extract any useful information from the training data, either staying near the random decision baseline or vastly underperforming. Successful models mostly exhibit a behaviour indicating they have extracted a basic approximation of D_2 from the linear training data by way of barely differentiating between the two error categories.

In total, processing extreme long-range dependencies spanning the whole word was more difficult for the models than processing extreme unseen nesting depths. This might be owed to the fact that most models did not learn a valid representation of D_2 : the unseen nesting depths still featured fairly short long-range dependencies compared to Experiment 1, which can easily be resolved in memory without needing to understand recursion.

As expected, LSTMs and GRUs outperformed SRNNs. While successful LSTM models achieved the highest individual accuracy scores (topping out at 99.3%) and came closest to a perfect open/closed misclassification ratio of 1, GRUs were more consistently successful, both producing the most models performing better than chance and achieving the highest average accuracy. It stands to reason that, while not succeeding under the circumstances of this study, LSTMs show the highest capability of perfectly learning D_2 . Taking the Chomsky-Schützenberger Representation Theorem (Chomsky and Schützenberger (1963)) into account, it seems possible that these architectures are capable of learning the vast class of Type-2 languages - which likely contains natural language.

I have found corpus complexity to have a significant impact on model performance. Models trained on the Low LRD corpus largely outperformed the Base corpus models, while High LRD corpus models underperformed. While models trained on complex words failing to generalize to longer, more complex words seems paradoxical at first, these results suggest that whatever rules RNNs extract from the input data, rule extraction becomes harder the more complex the data is. Taking into account the open/closed bracket misclassification ratio, I suggested that RNNs can extract a small number of simple rules from simple data and generally apply them in a more complex context. The more complex training data encourages the models to instead learn a larger set of overly specified rules - and then fail to generalize them. If a Low LRD trained model learns that an open bracket must always be closed by its corresponding closing bracket, the High LRD trained model might learn that an open bracket may be closed by its corresponding closing bracket only after a certain number of characters, or only after a certain number of nesting levels have been resolved, leading to a bloated, overspecified and misleading ruleset.

6.1 Further Research

The most compelling result of this research is the effect of training corpus complexity on generalizability. While it holds true in this case, formal language data is by definition rigorously structured and, for D_2 , rather simple and limited. Natural language data features a far bigger alphabet, complex syntactic, semantic and morphological dependencies and irregularities. Nonetheless, RNNs for NLP tasks improving with structurally simplified training data poses an intriguing and possibly fruitful avenue of future research. Furthermore, while LSTMs emerged as the most promising architecture to learn D_2 , they have proven to be more volatile than GRUs. Whether there is an inherent difference between the architectures remains unclear and may be explored with more established and sophisticated methods of internal state analysis.

Bibliography

- Autebert, J.-M., Berstel, J., and Boasson, L. (1997). Context-Free Languages and Pushdown Automata. In *Handbook of Formal Languages: Volume 1 Word, Language, Grammar*, 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? *LiLT (Linguistic Issues in Language Technology)*, 16(1).
- 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 *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 Vasisht, 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. *31st International Conference on Machine Learning, ICML 2014*, 5.
- 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.
- McKinney, W. (2011). pandas: a foundational Python library for data analysis and statistics. *Python for High Performance and Scientific Computing*, 14.
- Newmeyer, F. and Preston, L. (2014). *Measuring Grammatical Complexity*. Oxford linguistics. Oxford University Press.
- Oliphant, T. E. (2006). *A guide to NumPy*, volume 1. Trelgol Publishing USA.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E. (2011). Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.
- 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. In *The Formal Complexity 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. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 232–239, Brussels, Belgium. Association for Computational Linguistics.
- Suzgun, M., Gehrmann, S., Belinkov, Y., and Shieber, S. M. (2019). LSTM Networks Can Perform Dynamic Counting. *CoRR*, abs/1906.03648.
- Williams, R. J. and Zipser, D. (1995). Gradient-based learning algorithms for recurrent networks and their computational complexity. In *Backpropagation: Theory, Architectures, and Applications*, page 433–486. L. Erlbaum Associates Inc., USA.
- 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. (1993). Learning Finite State Machines With Self-Clustering Recurrent Networks. *Neural Computation*, 5:976–990.
- 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.

S

Appendix

```
1 # generate_raw_data.py
2 # Script to sample words from Dyck(2) according to Sennhauser &
   Berwick 2018.
3 #
4 # Grammar:
5 # S -> Z S | Z
6 # Z -> B | T
7 # B -> [ S ] | { S }
8 # T -> [] | {}
9 #
10 # branch = Z -> B
11 # concat = S -> Z S
12 # s(l) = min(1, -3 * (1/n) + 3)
13 # P_branch = r_b * s(l); r_b ~ U(0.4, 0.8)
14 # P_concat = r_c * s(l); r_c ~ U(0.4, 0.8)
15 # l = number of already generated characters in the sentence, n =
   total goal length
16 # https://en.wikipedia.org/wiki/Uniform\_distribution\_\(continuous\) == U
   (0.4, 0.8)
17
18 import sys
19 import numpy as np
20 import os
21
22 LIMIT = int(sys.argv[1])
23 GOAL_CORPUS_SIZE = int(sys.argv[2])
24
25 def countCharacters(word):
26     '''Counts the number of generated terminal symbols in a string.
27     args:
28         word: string, string of terminal and non-terminal symbols.
29     returns:
30         generated_characters: int, number of generated terminal
31         symbols in word.'''
32     generated_characters = 0
33     for character in word:
34         if character in {"[", "]", "{", "}"}:
35             generated_characters += 1
36
37     # The beginning 'S' is counted as 2 generated terminal symbols here
38     # - if the word contains an 'S', the minimum amount of terminal
39     # symbols in the fully derived word is 2
40     return generated_characters + 2
```



```
38
39 def findNonterminals(word):
40     '''Assesses whether the input string contains any non-terminal
41         symbols.
42         args:
43             word: string, string of terminal and non-terminal symbols.
44         returns:
45             bool'''
46     if "S" in word:
47         return True
48     elif "Z" in word:
49         return True
50     elif "T" in word:
51         return True
52     elif "V" in word:
53         return True
54     elif "B" in word:
55         return True
56     else:
57         return False
58
59 def expand(word, limit):
60     '''Expands the input word by applying probabilistic grammar rules.
61         The generated word cannot be longer than the given limit.
62         args:
63             word: string, string of terminal and non-terminal symbols.
64             limit: int, maximum length the word is allowed to have.
65         returns:
66             word: string, string of terminal and non-terminal symbols.'''
67     # Sampling probabilities as per Sennhauser & Berwick 2018, once per
68     # word
69     prob_b = np.random.uniform(0.6, 0.9)
70     prob_c = np.random.uniform(0.6, 0.9)
71     for i in range(len(word)):
72         # Once word length has exceeded the limit, the candidate is
73         # invalid and will not be processed further
74         if len(word) > limit:
75             return word
76         else:
77
78             s_l = min(1, -3 * countCharacters(word)/float(limit) + 3)
79             P_branch = prob_b * s_l
80             P_concat = prob_c * s_l
81             random = np.random.uniform(0.0, 1.0)
82
83             # Applying probabilistic Dyck2 grammar by replacing left-hand
84             # non-terminals with their right-hand production
85             if word[i] == 'S':
```

```

81         if P_concat <= 0.0:
82             word = word[:i] + 'Z' + word[i+1:]
83         elif random < P_concat:
84             word = word[:i] + 'ZS' + word[i+1:]
85         else:
86             word = word[:i] + 'Z' + word[i+1:]
87     elif word[i] == 'Z':
88         if P_branch <= 0.0:
89             word = word[:i] + 'T' + word[i+1:]
90         elif random < P_branch:
91             word = word[:i] + 'B' + word[i+1:]
92         else:
93             word = word[:i] + 'T' + word[i+1:]
94     elif word[i] == 'B':
95         if random < 0.5:
96             word = word[:i] + '[S]' + word[i+1:]
97         else:
98             word = word[:i] + '{S}' + word[i+1:]
99     elif word[i] == 'T':
100         if random < 0.5:
101             word = word[:i] + '[]' + word[i+1:]
102         else:
103             word = word[:i] + '{} ' + word[i+1:]
104     else:
105         continue
106
107     return word
108
109 def createCorpus(limit, goal_corpus_size):
110     '''Generates a corpus of size goal_corpus_size of words with len <=
111         limit. Saves the result in steps of 10000 generated words as a txt
112         file. Saves a final corpus as a txt file.
113
114     args:
115         limit: int, maximum length for a word in the corpus.
116         goal_corpus_size: int, number of words in the full corpus.
117
118     returns:
119         none'''
120     counter = 0
121     corpus = set()
122     # Populates the corpus
123     while len(corpus) < goal_corpus_size:
124         prev_save_len = 0
125         counter += 1
126         word = 'S'
127         # Expands the word until it either is fully derived or exceeds the
128         # desired length
129         while findNonterminals(word) and len(word) <= limit:
130             word = expand(word, limit)

```

```

126     # If the word is fully derived and fulfills the length requirement
127     , it is added to the corpus
128     if not findNonterminals(word):
129         if len(word) <= limit:
130             old_length = len(corpus)
131             corpus.add(word)
132
133     # Print occasional update on the generation process to the console
134     if counter % 50000 == 0:
135         print("Corp Size: {} \t Gen words: {}".format(len(corpus), counter
136         ))
137
138     # Save interim results
139     if len(corpus) % 10000 == 0 and len(corpus) != prev_save_len:
140         print("Saving corpus, length {} ...".format(len(corpus)))
141         filename = os.path.join '..', 'corpus', "cumlen{}_{}.txt".format(
142             limit, len(corpus))
143         outfile = open(filename, 'w')
144         for item in corpus:
145             outfile.write(item)
146             outfile.write('$')
147         outfile.close()
148
149     # Save full corpus
150     filename = os.path.join '..', 'corpus', "cumlen{}_{}.txt".format(
151         limit, len(corpus))
152     outfile = open(filename, 'w')
153     for item in corpus:
154         outfile.write(item)
155         outfile.write('$')
156     outfile.close()
157
158 createCorpus(LIMIT, GOAL_CORPUS_SIZE)

```

Listing 1: Script to generate the base data from which the corpora are built. Sampling is done according to Sennhauser and Berwick (2018).

```
1 # corpus_tools.py
2 # Works with a base file of correct Dyck words to create
3 # Generates classification datasets with a 50/50 split on correct/
   incorrect words each from a basic sampling of correct Dyck(2)
   words.
4 # Creates following datasets:
5 # TRAINING
6 #   - base
7 #   - high LRD
8 #   - low LRD
9 # EXPERIMENTS
10 # - LRD
11 # - ND
12
13 import os
14 import sys
15 import random
16 import pickle
17 import numpy as np
18
19 INPUT_FILE = 'cumlen20_1180000.txt'
20 INPUT_PATH = os.path.join '..', 'corpus', INPUT_FILE
21 SIZE = 1000000 # Training Corpus Size Target
22 LENGTH = int(INPUT_FILE[6:8]) # max length of word in set
23 POSITIVE_RATIO = 0.5
24 NEGATIVE_RATIO = 1-POSITIVE_RATIO
25 BD_CUTOFF = 19
26 MAX_BD_CUTOFF = 10
27 OUTPUT_TRAINING = os.path.join '..', 'training', 'base.csv')
28 OUTPUT_HIGH_LRD = os.path.join '..', 'training', 'high.csv')
29 OUTPUT_LOW_LRD = os.path.join '..', 'training', 'low.csv')
30 OUTPUT_LRD = os.path.join '..', 'experiment', 'LRD.csv')
31 OUTPUT_ND = os.path.join '..', 'experiment', 'ND.csv')
32
33 def samePair(char1, char2):
34     '''Helper function to check whether the two characters are a valid
       bracket pair.
35     args:
36         char1: string, character in word
37         char2: string, character in word
38     returns:
39         bool'''
40     if char1 == '{' and char2 == '}':
41         return True
42     elif char1 == '[' and char2 == ']':
43         return True
44     else:
```

```
45     return False
46
47 def maxNestingDepth(word):
48     '''Calculates the maximum nesting depth within a D_2 word.
49     args:
50         word: string, Dyck word consisting of [, {, }, ] as brackets.
51     returns:
52         max_depth: int'''
53     # Remove EOW symbol when working with processed corpus
54     if word[-1] == '$':
55         word = word[:-1]
56
57     max_depth = 0
58     depth = 0
59     for character in word:
60         if character == "[" or character == "{":
61             depth += 1
62         else: # Any other character must be a closing bracket and thus
63             # reduce depth
64             depth -= 1
65             if depth < 0:
66                 return -1 # Indicates a corrupted word with a superfluous
67                     # closing bracket in analysis
68
69         if depth > max_depth:
70             max_depth = depth
71
72     return max_depth
73
74 def maxValidNestingDepth(word):
75     '''Calculates the maximum valid nesting depth within a word -- if
76     the word was corrupted, negative nesting depth might occur.
77     This function disregards that.
78     args:
79         word: string, Dyck word consisting of [, {, }, ] as brackets.
80     returns:
81         max_depth: int'''
82     # Remove EOW symbol when working with processed corpus
83     if word[-1] == '$':
84         word = word[:-1]
85
86     max_depth = 0
87     depth = 0
88     for character in word:
89         if character == "[" or character == "{":
90             depth += 1
91         elif character == "]" or character == "}":
92             depth -= 1
```

```

90     else:
91         continue
92     if depth < 0:
93         continue
94
95     if depth > max_depth:
96         max_depth = depth
97
98     return max_depth
99
100 def nestingDepthAtPosition(word, position):
101     '''Calculates the nesting depth of the character at word[position].
102     args:
103         word: string, Dyck word consisting of [, {, }, ] as brackets.
104         position: int, index for word.
105     returns:
106         depth: int'''
107     depth = 0
108     for i in range(position):
109         if word[i] == "[" or word[i] == "{":
110             depth += 1
111         else:
112             depth -= 1
113
114     return depth
115
116 def maxBracketDistance(word):
117     '''returns maximum distance in characters between an opening and its
118     corresponding closing bracket, i.e. maxBracketDistance('[{[]}]')
119     is 4.
120     args:
121         word: string, Dyck word consisting of [, {, }, ] as brackets.
122     returns:
123         max(max_distance_square,max_distance_curly): int, maximum
124         distance for either of the two bracket pair types = maximum
125         distance in the word.'''
126     # Remove EOW symbol when working with processed corpus.
127     if word[-1] == '$':
128         word = word[:-1]
129     distance_square = 0
130     max_distance_square = 0
131     distance_curly = 0
132     max_distance_curly = 0
133     square_stack = []
134     curly_stack = []
135
136     # Counts 'seen' characters between two corresponding brackets
137     for character in word:

```

```
134     # Fill/empty the two stacks to keep count of unclosed brackets
135     if character == "[":
136         square_stack.append(character)
137     elif character == "{":
138         curly_stack.append(character)
139     elif character == "]":
140         square_stack.pop()
141     else:
142         curly_stack.pop()
143
144     # Check if stack is empty - means all opening brackets have been
    closed
145     if not square_stack:
146         # Save new max_distance_square and reset 'seen' characters
    counter
147         if distance_square > max_distance_square:
148             max_distance_square = distance_square - 1 # Adjust for first
    character being counted
149             distance_square = 0
150         else: # Increment 'seen' characters counter as long as there is
    still an unclosed bracket on the stack
151             distance_square += 1
152
153     # Same process as for square_stack
154     if not curly_stack:
155         if distance_curly > max_distance_curly:
156             max_distance_curly = distance_curly - 1
157             distance_curly = 0
158         else:
159             distance_curly += 1
160
161     return max(max_distance_square, max_distance_curly)
162
163 def bracketDistanceAtPosition(word, position):
164     '''returns distance from closing bracket at a given position to
165     its corresponding opening bracket.
166     args:
167         word: string, Dyck word consisting of [, {, }, ] as brackets.
168         position: position: int, index for word.
169     returns:
170         distance: int, number of characters between word[position] and
    the corresponding opening bracket.'''
171     distance = 0
172     stack = []
173     character = word[position]
174
175     # Make sure a closing bracket is at the position. Otherwise, it's an
    opening bracket, for which the measure does not make sense
```

```
176     if character == "]":
177         opener = "["
178     elif character == "}":
179         opener = "{"
180     else:
181         return 0
182
183     # Go backwards through the word, starting at position. Same logic
184     # applies as in maxBracketDistance(word)
185     for i in range(position, -1, -1):
186         if word[i] == opener:
187             stack.pop()
188         elif word[i] == character:
189             stack.append(character)
190
191         if not stack:
192             return distance-1
193         else:
194             distance += 1
195
196 def determineError(word):
197     '''Given a corrupted D_2 word, this function determines the kind of
198     error - too many opening or too many closing brackets.
199     args:
200         word: corrupted Dyck-2 word
201     returns:
202         error: type of error'''
203     open = ('[', '{')
204     closed = (']', '}')
205     o = 0
206     c = 0
207     for character in word:
208         if character in open:
209             o+=1
210         elif character in closed:
211             c+=1
212         else:
213             continue
214     error = o-c
215     if error < 0:
216         return 'closed'
217     elif error > 0:
218         return 'open'
219     else:
220         return 'none'
221
222 def findErrorPosition(word):
223     '''Given a corrupted D_2 word, this function determines the position
```



```

    of the corrupted bracket.
222     args:
223         word: corrupted Dyck-2 word
224     returns:
225         position: character position of corrupted bracket'''
226 word = word.lstrip('0')
227 # Stacks are filled with character positions. If a stack is not
    empty by the time the word is over,
228 # the remaining position is the error position.
229 # If there is an attempt to pop an empty stack before the word is
    over, the current position is the error position.
230 square_stack = []
231 curly_stack = []
232
233 for i in range(len(word)):
234     if word[i] == '[':
235         square_stack.append(i)
236     elif word[i] == '{':
237         curly_stack.append(i)
238     elif word[i] == ']':
239         if square_stack:
240             square_stack.pop()
241         else:
242             return i
243     elif word[i] == '}':
244         if curly_stack:
245             curly_stack.pop()
246         else:
247             return i
248     else:
249         continue
250 if square_stack:
251     return square_stack[0]
252 elif curly_stack:
253     return curly_stack[0]
254 else:
255     raise ValueError("Encountered unknown corruption in word\n{}".
        format(word))
256
257 def measureLength(corpus):
258     '''Calculates average length and variance thereof for all words in
    the corpus.
259     args:
260         corpus: list of strings, list of Dyck words.
261     returns:
262         avg: float, average length of words in corpus.
263         var: float, variance of length of words in corpus.'''
264     if len(corpus) == 0:

```

```
265     return 0., 0.
266 total = []
267 for word in corpus:
268     total.append(len(word))
269 avg = sum(total)/float(len(total))
270 var = sum((length - avg)**2 for length in total)/float(len(total))
271
272 return avg, var
273
274 def measureMaxNestingDepth(corpus):
275     '''Calculates average maximum nesting depth and variance thereof for
276         all words in the corpus.
277     args:
278         corpus: list of strings, list of Dyck words.
279     returns:
280         avg: float, average maximum nesting depth of words in corpus.
281         var: float, variance of maximum nesting depth of words in corpus
282     '''
283     if len(corpus) == 0:
284         return 0., 0.
285     total = []
286     for word in corpus:
287         total.append(maxNestingDepth(word))
288     avg = sum(total)/float(len(total))
289     var = sum((depth - avg)**2 for depth in total)/float(len(total))
290
291     return avg, var
292
293 def measureMaxBracketDistance(corpus):
294     '''Calculates average maximum bracket distance and variance thereof
295         for all words in the corpus.
296     args:
297         corpus: list of strings, list of Dyck words.
298     returns:
299         avg: float, average maximum bracket distance of words in corpus.
300         var: float, variance of maximum bracket distance of words in
301             corpus.'''
302     if len(corpus) == 0:
303         return 0., 0.
304     total = []
305     for word in corpus:
306         total.append(maxBracketDistance(word))
307     avg = sum(total)/float(len(total))
308     var = sum((dist - avg)**2 for dist in total)/float(len(total))
309
310     return avg, var
311
312 def evaluateCorpus(corpus):
```

```

309 '''Calculates average and variance for three measures over all words
    in a corpus: word length, maximum nesting depth and maximum
    bracket distance.
310     args:
311         corpus:
312     returns:
313         3 x (avg, var): 3 tuples of floats, average and variance for the
    respective measure.'''
314 corpus = [entry[0][:-1] for entry in corpus if entry[1]] # Removing
    EOW
315
316 return measureLength(corpus), measureMaxNestingDepth(corpus),
    measureMaxBracketDistance(corpus)
317
318 def printStats(size, avgLen, varLen, avgMaxND, varMaxND, avgMaxBD,
    varMaxBD):
319     '''Prints a table with all calculated corpus stats to the console.
    Table is for copy-pasting into .tex file.
320     args:
321         size: int, number of words in the corpus
322         avgLen: float, average length of word in the corpus
323         varLen: float, variation of length of words in the corpus
324         avgMaxND: float, average maximum nesting depth of word in the
    corpus
325         varMaxND: float, variation of maximum nesting depth of words in
    the corpus
326         avgMaxBD: float, average maximum bracket depth of word in the
    corpus
327         varMaxBD: float, variation of maximum bracket depth of words in
    the corpus
328     returns:
329         none'''
330     print("Size\tAvg Length \t Avg MaxNestDepth \t Avg MaxBrackDist")
331     print("{}\t${:3.2f}$ ({:3.2f}$) & {:3.2f}$ ({:3.2f}$) & {:3.2f}$
    ({:3.2f}$)".format(size, avgLen, varLen, avgMaxND, varMaxND,
    avgMaxBD, varMaxBD))
332
333 def largerBD(corpus):
334     '''Increases average bracket distance for a corpus by finding words
    with a low maximum bracket distance, deleting the lowest distance
    pair from the word and then wrapping the word in a matching pair
    of brackets.
335     args:
336         corpus: list of strings, list of Dyck words.
337     returns:
338         corpus: list of strings, list of Dyck words.'''
339     # 'Word collectors' are initialized as lists to allow iteration
    through them

```

```

340 small_BD = []
341 big_BD = []
342
343 # Find words with low maximum bracket distance
344 for entry in corpus:
345     word = entry[0][:-1]
346     correct = entry[1]
347     if correct:
348         maxBD = maxBracketDistance(word)
349         if maxBD < BD_CUTOFF:
350             small_BD.append(word)
351         elif maxBD >= BD_CUTOFF:
352             big_BD.append(word)
353
354 # Modify low maximum bracket distance words
355 for word in small_BD:
356     prev_ND = 0 # Nesting depth to compare to
357     for i in range(len(word)):
358         # Calculate nesting depth at each position of the word. Once it
359         # decreases, a closing bracket has been found
360         ND = nestingDepthAtPosition(word, i)
361         if ND < prev_ND:
362             # Check if this position belongs to a bracket pair eligible
363             # for deletion - only {} and [] are eligible, since they have the
364             # shortest possible bracket distance
365             char = word[i-1]
366             prev_char = word[i-2]
367             if samePair(prev_char, char):
368                 bracket = random.randint(0, 1)
369                 if bracket:
370                     word = '[' + word[:i-2] + word[i:] + ']'
371                     break
372                 else:
373                     word = '{' + word[:i-2] + word[i:] + '}'
374                     break
375             prev_ND = ND # Update nesting depth
376             big_BD.append(word) # Populate modified list
377
378 BD_set = set(big_BD) # Fast deletion of duplicates.
379 BD_list = list(BD_set)
380 incorrect = [entry for entry in corpus if entry[1]==0]
381 BD = [[word+'$',1] for word in BD_list] # Complete newly created and
382     # old correct words with EOW and class
383 BD = BD + incorrect
384 random.shuffle(BD)
385
386 return BD

```

```

384 def smallerBD(corpus):
385     '''Decreases average bracket distance for a corpus by finding words
        with a high maximum bracket distance. In those words, the pair
        with the highest maximum bracket distance is found. The opening
        bracket is then moved right in front of the closing bracket.
386     args:
387         corpus: list of strings, list of Dyck words.
388     returns:
389         corpus: list of strings, list of Dyck words.'''
390     # 'Word collectors' are initialized as lists to allow iteration
        through them
391     big_BD = []
392     small_BD = []
393
394     # Find words with high maximum bracket distance
395     for entry in corpus:
396         word = entry[0][:-1]
397         correct = entry[1]
398         if correct:
399             maxBD = maxBracketDistance(word)
400             if maxBD > MAX_BD_CUTOFF:
401                 big_BD.append(word)
402             elif maxBD <= MAX_BD_CUTOFF:
403                 small_BD.append(word) # big_BD only features word entries,
        since all big_BD entries are correct
404
405     for word in big_BD:
406         max_BD = 0
407         max_pos = 0
408         for i in range(len(word)):
409             # Calculate bracket distance at each position of the word. Once
        the maximum bracket distance has been found, the position of the
        closing bracket is recorded
410             BD = bracketDistanceAtPosition(word, i)
411             if BD > max_BD:
412                 max_pos = i
413                 max_BD = BD
414             closer = max_pos # Position of longest distance closing bracket
415             opener = max_pos - max_BD - 1 # Fix off by one return of bd@pos
416             # New word is created by deleting the opener from its original
        position and moving it right in front of the closing bracket
417             # This ensures grammaticality and reduces maxBracketDistance
418             new_word = word[:opener] + word[opener+1:closer] + word[opener] +
        word[closer:]
419             small_BD.append(new_word)
420
421     BD_list = list(set(small_BD))
422     incorrect = [entry for entry in corpus if entry[1]==0]

```

```

423 BD = [[word+'$',1] for word in BD_list] # Complete newly created and
      old correct words with EOW and class
424 BD = BD + incorrect
425
426 return BD
427
428 def corrupt_words(correct_corpus, mode):
429     '''Turns correct Dyck_2 words into incorrect ones by replacing a
      random opening bracket with a random closing one or vice versa,
      depending on the mode.
430     args:
431         correct_corpus: list of correct word-class tuples
432         mode: string. "open"/"close" - determines brackets being changed
433     returns:
434         incorrect: list of incorrect word-class tuples'''
435     incorrect = set() # set ensures uniqueness
436     if mode == "open":
437         replace = ("{","[")
438         replacement = ["]","}"]
439     else:
440         replace = ("}","]")
441         replacement = ["{","["]
442
443     for entry in correct_corpus:
444         word = entry[0]
445         changed = 0
446         c = 0
447         while not changed:
448             c += 1
449             id = random.randint(0, len(word)-1)
450             if word[id] in replace:
451                 closing = random.choice(replacement)
452                 new_word = word[:id] + closing + word[id+1:]
453                 assert len(new_word)==len(word)
454                 incorrect.add(new_word)
455                 changed = 1
456             elif c >= 1000000:
457                 break
458     incorrect = [[word,0] for word in incorrect]
459
460     return incorrect
461
462 def create_LRD(base, subword_length):
463     '''Creates the dataset for the extreme long range dependency (LRD)
      experiment/experiment 1.
464     args:
465         base: list of correct word-class tuples
466     returns:

```

```

467     LRD: list of correct LRD word-class tuples'''
468     LRD = set()
469     base = [entry[0] for entry in base if len(entry[0])==subword_length]
470     for i in range(len(base)*5):
471         bracket = random.randint(0,1)
472         w1 = random.sample(base, 1)[0][: -1]
473         w2 = random.sample(base, 1)[0][: -1]
474         if bracket:
475             new_word = '[' + w1 + w2 + ']'
476         else:
477             new_word = '{' + w1 + w2 + '}'
478         LRD.add(new_word)
479     LRD = [[word,1] for word in LRD]
480
481     return LRD
482
483 def create_ND(base, infix_length):
484     '''Creates the dataset for the extreme new depths (ND) experiment/
485     experiment 2.
486     args:
487         base: list of correct word-class tuples
488     returns:
489         ND: list of correct LRD word-class tuples'''
490     ND = set()
491     base = [entry[0] for entry in base if len(entry[0])==infix_length]
492     for i in range(len(base)*5):
493         new_word = random.sample(base, 1)[0][: -1]
494         for i in range(5):
495             bracket = random.randint(0,1)
496             if bracket:
497                 new_word = '[' + new_word + ']'
498             else:
499                 new_word = '{' + new_word + '}'
500         ND.add(new_word + '$')
501     ND = [[word,1] for word in ND]
502
503     return ND
504
505 def create_corpus(data, type):
506     '''Creates full corpora for a classification task. Increases or
507     lowers bracket distance for the high/low LRD training corpora.
508     Prints corpus stats to the console, with the values in LaTeX
509     formatting.
510     args:
511         data: list of [word,correct-bool] pairs, filled with correct
512         generated D_2 words
513         type: string, high/low/[misc], triggers increasing/decreasing/
514         disregarding average bracket distance across the corpus

```

```

510         returns:
511         corpus: list of [word,correct-bool] pairs, filled with a
1:0.5:0.5 ratio of correct:incorrect_open:incorrect_closed words
,,,
512 if type == 'high':
513     corpus_correct = largerBD(data)
514 elif type == 'low':
515     corpus_correct = smallerBD(data)
516 else:
517     corpus_correct = data
518 corpus_incorrect_open = corrupt_words(corpus_correct, 'open')
519 corpus_incorrect_closed = corrupt_words(corpus_correct, 'closed')
520 # Debug print
521 #print(" === {} ===\nCorrect\tIncorrect O\tIncorrect C\n{}\t{}\t{}".
format(type.upper(), len(data), len(corpus_incorrect_open), len(
corpus_incorrect_closed)))
522 corpus = data[:int(POSITIVE_RATIO*SIZE)] + corpus_incorrect_open[:
int(NEGATIVE_RATIO/2.*SIZE)] + corpus_incorrect_closed[:int(
NEGATIVE_RATIO/2.*SIZE)]
523 random.shuffle(corpus)
524 (avLen, varLen), (avMaxND, varMaxND), (avMaxBD, varMaxBD) =
evaluateCorpus(corpus)
525 printStats(len(corpus), avLen, varLen, avMaxND, varMaxND, avMaxBD,
varMaxBD)
526
527 return corpus
528
529
530 def save2file(outpath, corpus):
531     '''Saves the corpus as a .csv file to the specified path.
532     args:
533         outpath: filepath of the output file
534         corpus: list of [word,correct-bool] pairs, to be used for
training/experiments on RNNs
535     returns:
536         none'''
537     outfile = open(outpath, 'w')
538     outfile.write('word,value\n')
539     for entry in corpus:
540         outfile.write('{}{}\n'.format(entry[0], entry[1]))
541
542 def create_corpora():
543     '''Creates all datasets needed for the classification tasks from the
input file.
544     args:
545         none
546     returns:
547         none'''

```



```
548 file = open(INPUT_PATH, 'r')
549 EOW = '$'
550 raw_text = file.read()
551
552 print("Creating Base...")
553 raw_classified = [[word+EOW,1] for word in raw_text.split(EOW)]
554
555 print("Corrupting Base...")
556 base = create_corpus(raw_classified, 'base')
557 save2file(OUTPUT_TRAINING, base)
558
559 print("Corrupting High LRD...")
560 highLRD = create_corpus(raw_classified, 'high')
561 save2file(OUTPUT_HIGH_LRD, highLRD)
562
563 print("Corrupting Low LRD...")
564 lowLRD = create_corpus(raw_classified, 'low')
565 save2file(OUTPUT_LOW_LRD, lowLRD)
566
567 print("Creating LRD...")
568 LRD_correct = create_LRD(raw_classified, LENGTH-2+1)
569 LRD = create_corpus(LRD_correct, 'LRD')
570 save2file(OUTPUT_LRD, LRD)
571
572 print("Creating ND...")
573 ND_correct = create_ND(raw_classified, LENGTH+1)
574 ND = create_corpus(ND_correct, 'ND')
575 save2file(OUTPUT_ND, ND)
```

Listing 2: Collection of functions to measure corpus and word properties. Used to generate five classification datasets: three different training sets and the two experiment sets.

```

1 # RNN_classifier.py
2 # Training, testing and experimenting on RNN models
3
4 import os
5 import sys
6 import pandas as pd
7 import numpy as np
8 from sklearn.model_selection import train_test_split
9 from sklearn.preprocessing import LabelEncoder
10 import tensorflow as tf
11 tf.keras.backend.clear_session() # Reduces memory leak during training
12     - seems to be issue with TF2.0
13 from tensorflow.keras import layers
14 from tensorflow.keras import losses
15 from tensorflow.keras.models import Sequential
16 from tensorflow.keras.layers import LSTM, Activation, Dense, Embedding
17     , SimpleRNN, GRU
18 from tensorflow.keras.optimizers import Adam
19 from tensorflow.keras.preprocessing.text import Tokenizer
20 from tensorflow.keras.preprocessing import sequence
21 from tensorflow.keras.utils import to_categorical
22 from tensorflow.keras.callbacks import EarlyStopping, ModelCheckpoint
23
24 MODE = sys.argv[1]
25 NETWORK = sys.argv[2]
26 HIDDEN_UNITS = int(sys.argv[3])
27 CORPUS = sys.argv[4]
28 EXPERIMENT = sys.argv[5]
29 EPOCHS = 100
30 TEST_SIZE = 0.2
31 VOCAB_SIZE = 7 # 6 characters + 1 do combat off by 1 error in
32     Embedding Layer
33 MAX_LEN = ((20-2)*2)+1+2 # Maximum length of words in all corpora and
34     experiments -- LRD length formula
35 BATCH = 512
36 CHECKPOINT_DIR = os.path.join '..', 'saved_models', CORPUS, NETWORK,
37     str(HIDDEN_UNITS))
38 CORPUS_DF = pd.read_csv('../training/{}.csv'.format(CORPUS), delimiter=
39     ',', encoding='latin-1')
40 EXP_DF = pd.read_csv('../experiment/{}.csv'.format(EXPERIMENT))
41
42 # ===== PROCESSING INPUT DATA =====
43
44 def filter_length(df, max_len_train):
45     '''Filters words of length > max_len_train out of a dataframe.
46     args:
47         df: pd.DataFrame of words

```

```

42     max_len_train: Maximum length words in df should have
43     returns:
44     df: pd.DataFrame of words'''
45     mask = (df['word'].str.len() <= max_len_train)
46     df = df.loc[mask]
47     df.head()
48     df.info()
49
50     return df
51
52 def prepare_training_data(df, max_word_length=MAX_LEN):
53     '''Preprocesses training data by splitting it into train/test and
54     transforming the strings into sequences of numbers.
55     args:
56     df: pd.DataFrame of words
57     max_word_length: Maximum wanted length of words
58     returns:
59     X_train: Preprocessed words for the train split
60     X_test: Preprocessed words for the test split
61     Y_train: Target values for the train split
62     Y_test: Target values for the test split
63     tok: Tokenizer fitted on training corpus, to be used on
64     experiment corpus'''
65     df = filter_length(df, max_word_length)
66     X = df.word
67     Y = df.value
68     # Encode target values
69     le = LabelEncoder()
70     Y = le.fit_transform(Y)
71     Y = Y.reshape(-1,1)
72
73     X_train_temp, X_test_temp, Y_train, Y_test = train_test_split(X, Y,
74     test_size=TEST_SIZE)
75
76     # Preprocess words into sequences of numbers corresponding to the
77     characters
78     tok = Tokenizer(char_level=True)
79     tok.fit_on_texts(X_train_temp)
80     sequences = tok.texts_to_sequences(X_train_temp)
81     X_train = sequence.pad_sequences(sequences, maxlen=MAX_LEN)
82     test_sequences = tok.texts_to_sequences(X_test_temp)
83     X_test = sequence.pad_sequences(test_sequences, maxlen=MAX_LEN)
84
85     return X_train, X_test, Y_train, Y_test, tok
86
87 def prepare_experiment_data(dataframe, tok, max_word_length=MAX_LEN):
88     '''Preprocesses experiment data by splitting it into train/test and
89     transforming the strings into sequences of numbers.

```

```

85     args:
86         df: pd.DataFrame of words
87         max_word_length: Maximum wanted length of words
88         tok: Tokenizer as fitted on the training corpus
89     returns:
90         X: Preprocessed words
91         Y: Target values'''
92     df = filter_length(dataframe, max_word_length)
93     X_temp = df.word
94     Y = df.value
95     # Encode target values
96     le = LabelEncoder()
97     Y = le.fit_transform(Y)
98     Y = Y.reshape(-1,1)
99
100    # Preprocess words into sequences of numbers corresponding to the
101    # characters
102    # The same tokenizer is used for training and experiment data (
103    # training tokenizer is called in this function)
104    sequences = tok.texts_to_sequences(X_temp)
105    X = sequence.pad_sequences(sequences,maxlen=MAX_LEN)
106
107    return X, Y
108
109    # ===== MODELS =====
110
111    def build_model(layer_size=HIDDEN_UNITS, network=NETWORK, vocabulary=
112    VOCAB_SIZE, max_len=MAX_LEN):
113    '''Creates a LSTM, GRU or SRNN based model with a specified number
114    of hidden units.
115    args:
116        layer_size: Number of hidden units in the LSTM, GRU or SRNN
117        layer
118        network: Name of network architecture
119        vocabulary: Number of units in the Embedding layer
120        max_len: Maximum length of a single input sequence
121    returns:
122        model: Trainable model'''
123    model = Sequential()
124    model.add(Embedding(vocabulary,vocabulary,input_length=max_len)) #
125    Embeds the input sequence in the same dimensionality - simplifies
126    the tf dataflow
127    if network == 'LSTM':
128        model.add(LSTM(layer_size))
129    elif network == 'GRU':
130        model.add(GRU(layer_size))
131    elif network == 'SRNN':
132        model.add(SimpleRNN(layer_size))

```

```

126     else:
127         raise ValueError("{} is not a valid network name. Valid network
128             names are 'SRNN', 'LSTM' and 'GRU'.".format(NETWORK))
129         return 0
130     model.add(Dense(1, activation='sigmoid'))
131     model.compile(loss='binary_crossentropy', optimizer=tf.optimizers.
132         Adam(learning_rate=0.0001), metrics=['accuracy'])
133     model.summary()
134
135     return model
136
137 def train_model(X, Y, checkpoint_dir=CHECKPOINT_DIR):
138     '''Trains the model on provided training data and saves the model
139     weights for every epoch.
140     Training stops either after the global maximum number of EPOCHS or
141     if the the loss on the validation set val_loss does not improve
142     for 3 epochs in a row, in which case the previously best model is
143     used for the rest of the session.
144     args:
145         X: Preprocessed words
146         Y: Target values
147         hidden_units: Number of hidden units in the LSTM, GRU or SRNN
148         layer
149         corpus: Name of the training corpus
150     returns:
151         model: Trained model
152     '''
153     # Saving checkpoints: Directory, filenames, callback
154     # Name of the checkpoint files, saving the values of loss, accuracy,
155     # val_loss and val_accuracy
156     checkpoint_prefix = os.path.join(checkpoint_dir, "ckpt_{epoch:02d}-{
157         loss:.4f}-{accuracy:.4f}-{val_loss:.4f}-{val_accuracy:.4f}")
158     #checkpoint_prefix = os.path.join(checkpoint_dir, "ckpt_{epoch:02d
159         }-{loss:.4f}-{accuracy:.4f}") # for len8 experiments
160     checkpoint_callback = tf.keras.callbacks.ModelCheckpoint(filepath=
161         checkpoint_prefix, save_weights_only=True)
162
163     # Create and train model
164     model = build_model()
165     model.fit(X,Y,batch_size=BATCH,epochs=EPOCHS,validation_split=0.15,
166         callbacks=[EarlyStopping(monitor='val_loss', min_delta=0.0001,
167             patience=3, restore_best_weights=True),
168             checkpoint_callback])
169
170     return model
171
172 # ===== TRAINING/TESTING/EXPERIMENTS =====
173 def set_best_checkpoint(checkpoint_dir=CHECKPOINT_DIR):

```

```

162 '''tf.train.latest_checkpoint() loads the latest saved model weights
    by referring to the file 'checkpoint' in the respective
    checkpoint directory.
163 However, the latest model is not always the best.
164 To circumvent this issue, this function overwrites the file '
    checkpoint' to refer
165 to the best (= lowest val_loss) model weight configuration, so tf.
    train.latest_checkpoint() instead retrieves the best model.
166     args:
167         checkpoint_dir: Directory whose checkpoint file is rewritten
168     returns:
169         None'''
170 files = os.listdir(checkpoint_dir)
171 weight_files = files[1:] # Ignore checkpoint file
172 filenames = [filename[:-6] for filename in weight_files if filename.
    endswith('.index')]
173
174 # Determine index of model weights with minimal val_loss
175 val_losses = []
176 for filename in filenames:
177     prefix, train_loss, train_acc, val_loss, val_acc = filename.split(
        '-')
178     val_losses.append(float(val_loss))
179 idx = val_losses.index(min(val_losses))
180 best_weights = filenames[idx]
181
182 # Overwrite checkpoint
183 outfile_name = os.path.join(checkpoint_dir, 'checkpoint')
184 outfile = open(outfile_name, 'w')
185 outfile.write('model_checkpoint_path: "{0}"\n
    nall_model_checkpoint_paths: "{0}"'.format(best_weights))
186 outfile.close()
187
188 return None
189
190 def train_test():
191     '''Trains and tests the model with a train_test_split. Saves the
        accuracy and loss on the test data in a file with name '{NETWORK}_
        {HIDDEN_UNITS}.csv'.
192     Overwrites 'checkpoint' file TF2.0 uses to store the name of the
        latest saved weights file with the name of the lowest val_loss
        weights file.
193     args:
194         none
195     returns:
196         none'''
197     # Preprocessing
198     X_train, X_test, Y_train, Y_test, tokenizer = prepare_training_data(

```

```

CORPUS_DF)
199
200 # Training model on train_split
201 model = train_model(X_train, Y_train)
202
203 # Testing on test_split
204 test_accr = model.evaluate(X_test, Y_test) # test_accr = [loss,
      accuracy] on test data
205
206 # Creating results file for test_split results
207 test_dir = os.path.join '..', 'test_results', CORPUS)
208 try:
209     os.makedirs(test_dir)
210 except FileExistsError:
211     print("Directory {} already exists, proceeding with saving test
      results.".format(test_dir))
212 test_out_filename = '{}_{}.csv'.format(NETWORK, HIDDEN_UNITS)
213 test_out_path = os.path.join(test_dir, test_out_filename)
214
215 # Saving test_split results
216 test_out = open(test_out_path, 'w')
217 test_out.write('Loss,Accuracy\n{:0.3f},{:0.3f}'.format(test_accr[0],
      test_accr[1]))
218 test_out.close()
219
220 # Set checkpoint file to refer to the best weights rather than the
      latest ones
221 set_best_checkpoint()
222
223 def run_experiment(max_word_length=MAX_LEN, checkpoint_dir=
      CHECKPOINT_DIR):
224     '''Makes the model classify experiment data specified in EXP_DF.
      Saves loss and accuracy in a '{NETWORK}_{HIDDEN_UNITS}.csv' file.
225     Saves predictions per word in a 'detail_{NETWORK}_{HIDDEN_UNITS}.csv
      file in the shape of [preprocessed word], [predicted_value], [
      correct_value] for
226     further analysis.
227     args:
228         max_word_length: Maximum wanted length of words
229         corpus: Name of the training corpus
230         network: Name of network architecture
231         hidden_units: Number of hidden units in the LSTM, GRU or SRNN
      layer
232     returns:
233         None'''
234     # Building model and assigning the weights of the model specified in
      its checkpoint file: model with the lowest val_loss
235     model = build_model()

```

```

236 model.load_weights(tf.train.latest_checkpoint(checkpoint_dir))
237
238 # Creating tokenizer from training corpus to apply to experiment
    corpus
239 discard_matrix, discard_matrix, discard_Y_train, discard_Y_test,
    tokenizer = prepare_training_data(CORPUS_DF, max_word_length)
240 # Preprocessing experiment data
241 exp_sequences_matrix, Y = prepare_experiment_data(EXP_DF, tokenizer)
242
243 # Complete evaluation
244 exp_accr = model.evaluate(exp_sequences_matrix, Y, batch_size=BATCH,
    verbose=0) # exp_accr = [loss, accuracy] on experiment data
245 # Classifying experiment words on a word by word basis for detailed
    results
246 exp = model.predict(exp_sequences_matrix)
247
248 # Creating results file
249 exp_dir = os.path.join '..', 'exp_results', CORPUS, EXPERIMENT)
250 try:
251     os.makedirs(exp_dir)
252 except FileExistsError:
253     print("Directory {} already exists, proceeding with experiment."
        format(exp_dir))
254
255 exp_out_filename = '{}_{}.csv'.format(NETWORK, HIDDEN_UNITS)
256 exp_out_path = os.path.join(exp_dir, exp_out_filename)
257
258 # Saving results
259 exp_out = open(exp_out_path, 'w')
260 exp_out.write('Loss, Accuracy\n{:0.3f},{:0.3f}'.format(exp_accr[0],
    exp_accr[1]))
261 exp_out.close()
262
263 # Creating detailed results file
264 detail_out_filename = 'detail_{}_{}.csv'.format(NETWORK,
    HIDDEN_UNITS)
265 detail_out_path = os.path.join(exp_dir, detail_out_filename)
266
267 # Saving detailed results
268 detail_out = open(detail_out_path, 'w')
269 for i in range(len(exp)):
270     detail_out.write('{} {}, {} \n'.format(exp_sequences_matrix[i], exp[
        i][0], Y[i][0]))
271 detail_out.close()
272
273 return None
274
275 if MODE == 'train':

```



```
276     train_test()
277 elif MODE == 'exp':
278     run_experiment()
279 else:
280     raise ValueError("{} is not an appropriate mode. Use 'train' or '
        test' instead.".format(MODE))
```

Listing 3: Preprocessing input data, creating and training of RNN models. The trained models are then loaded to classify experiment data.

```

1 # concat_results.py
2 # Creates a .csv overview over all model performances.
3 # Collecting calculated performance measures for all RNN models
4 import os
5 import pandas as pd
6
7 CORPORA = ['high', 'base', 'low']
8 EXPERIMENTS = ['ND', 'LRD']
9 OUTPATH_CSV = os.path.join '..', 'results'
10 OUTPATH_LATEX = os.path.join '..', 'latex', 'tab'
11 PERFORMANCE_COLUMNS = ['network', 'accuracy', 'precision', 'recall', '
    f1_score', 'val_acc']
12
13 def read_checkpoint(ckpt_file):
14     '''Reads 'checkpoint' file for epoch, train_loss, train_acc,
15         val_loss and val_acc of the lowest val_loss model.
16     args:
17         ckpt_file: Path to the 'checkpoint' file
18     returns:
19         epoch: Epoch of the lowest val_loss model
20         train_loss: Loss on training data of the lowest val_loss model
21         train_acc: Accuracy on training data of the lowest val_loss
22         model
23         val_loss: Loss on validation data of the lowest val_loss model
24         val_acc: Accuracy on validation data of the lowest val_loss
25         model'''
26     in_file = open(ckpt_file, 'r')
27     lines = in_file.readlines()
28     # Relevant field is in line 0 of checkpoint, second word and in
29     # quotation marks
30     # Accessing relevant field and stripping quotation marks
31     ckpt_string = lines[0].split()[1][1:-1] # "ckpt_epoch-train_loss-
32     train_acc-val_loss-val_acc" without quotation marks
33     ckpt_values = ckpt_string.split('_')[1]
34     epoch = int(ckpt_values.split('-')[0])
35     train_loss, train_acc, val_loss, val_acc = [float(value) for value
36         in ckpt_values.split('-')[1:]]
37
38     return epoch, train_loss, train_acc, val_loss, val_acc
39
40 def create_results(corpora=CORPORA, experiments=EXPERIMENTS):
41     '''Collect all results across all corpora and experiments in a list
42         of pd.DataFrames to create a dataframe containing all results.
43         Used for data exploration and to create tables for the thesis.
44     args:
45         corpora: list of corpora to include in the big dataframe
46         experiments: list of experiments to include in the big

```

```

dataframe
    returns:
39         results: dataframe with all specified results'''
40
41 single_frames = []
42 for corpus in corpora:
43     for experiment in experiments:
44         directory = os.path.join '..', 'exp_results', corpus, experiment
45         )
46         dirs = os.listdir(directory)
47         for filename in dirs: # Going through results for all
48             hidden_units configurations
49             if not filename.startswith('detail_'): # Disregarding the
50                 detailed results
51                 network, hidden_units = filename[:-4].split('_')
52                 model_directory = os.path.join '..', 'saved_models', corpus,
53                 network, hidden_units)
54                 models = os.listdir(model_directory)
55                 ckpt = os.path.join(model_directory, models[0])
56                 epoch, train_loss, train_acc, val_loss, val_acc =
57                 read_checkpoint(ckpt)
58
59                 df = pd.read_csv(os.path.join(directory, filename), delimiter
60                 = ',')
61                 df = df.loc[:, ~df.columns.str.contains('^Unnamed')]
62                 df = df.rename(str.lower, axis='columns')
63
64                 # Putting all measured values for this model together
65                 df['val_acc'], df['val_loss'], df['train_acc'], df['
66                 train_loss'], df['epoch'], df['network'], df['experiment'], df['
67                 corpus'], df['hidden_units'] = val_acc, val_loss, train_acc,
68                 train_loss, epoch, network, experiment, corpus, hidden_units
69                 col_order = ['network', 'hidden_units', 'corpus', '
70                 experiment', 'accuracy', 'precision', 'recall', 'f1_score', 'epoch
71                 ', 'train_acc', 'train_loss', 'val_acc', 'val_loss']
72                 df = df[col_order]
73                 single_frames.append(df)
74
75 # Join all small dfs into the complete overview of all results
76 results = pd.concat(single_frames, ignore_index=True)
77 results = results.sort_values(by='accuracy', ascending=False)
78
79 return results
80
81 def save2file(dataframe, experiment, corpus, mode, outpath):
82     '''Saves a given dataframe as either a .csv or .tex table. Latex
83     formatting was then manually modified.
84     args:
85         dataframe: dataframe containing measurements
86         experiment: string, name of the current experiment

```

```

74     corpus: string, name of the current corpus
75     mode: string, csv/latex, determines which file to generate
76     outpath: string, determines where the file is saved_models
77     returns:
78     none'''
79     try:
80         os.makedirs(outpath)
81     except:
82         print("Directory {} exists, proceeding.".format(outpath))
83     if mode=='csv':
84         dataframe.to_csv(os.path.join(outpath, 'results_{}_{}.csv'.format(
85             experiment, corpus)), index=False)
86     elif mode=='latex': # Latex output has been used as a basis for
87         # tables in the thesis, heavy modifications were made
88         dataframe.to_latex(os.path.join(outpath, 'results_{}_{}.tex'.
89             format(experiment, corpus)), columns=PERFORMANCE_COLUMNS,
90             float_format='{:0.3f}'.format, index=False)
91
92 def create_all_tables(df, mode):
93     '''Creates tables collecting all results given the supplied mode.
94     args:
95         df: dataframe containing measurements
96         mode: string, LRD/ND/SRNN/LSTM/GRU, determines what to filter
97         the dataframe for
98     returns:
99     none'''
100     corpora = ['base', 'low', 'high']
101     measures = []
102     if mode in ('SRNN', 'LSTM', 'GRU'):
103         data = df.loc[df.network == mode]
104         for corpus in corpora:
105             data_corpus = data.loc[data.corpus == corpus]
106             save2file(data_corpus, mode, corpus, 'csv', OUTPATH_CSV)
107             print(mode, corpus)
108             mean = data_corpus.mean()
109             std = data_corpus.std()
110             total = pd.concat([mean, std], axis=1)
111             print(total.T)
112             measures.append(total.T)
113         overview = pd.concat(measures)
114         print(overview)
115     elif mode in ('LRD', 'ND'):
116         data = df.loc[df.experiment == mode]
117         for corpus in corpora:
118             data_corpus = data.loc[data.corpus == corpus]
119             save2file(data_corpus, mode, corpus, 'csv', OUTPATH_CSV)
120
121 results = create_results()

```

```
117 results['hidden_units'] = results['hidden_units'].astype(int)
118 results = results.sort_values(by=['network', 'hidden_units'],
    ascending=True)
119
120 create_all_tables(results, 'LRD')
121 create_all_tables(results, 'ND')
122 create_all_tables(results, 'SRNN')
123 create_all_tables(results, 'LSTM')
124 create_all_tables(results, 'GRU')
```

Listing 4: Script to create the basis of the tables featured in Section 4.1.

```
1 # analyze_performance.py
2 # Functions to prepare network predictions for further analysis.
3 # Creates dataframes collecting the predictions of all outlier
   networks for error analysis.
4 # Only false positives were discussed in this thesis.
5 # To be used as 'python analyze_performance.py > some_outfile.csv'
6
7 import os
8 import pandas as pd
9 import corpus_tools
10 import tensorflow as tf
11 from tensorflow.keras.preprocessing.text import Tokenizer
12 from sklearn.preprocessing import LabelEncoder
13 from sklearn.model_selection import train_test_split
14 import numpy as np
15
16 def numbers_to_words(corpus, experiment, network, hidden_units):
17     '''Cleans up the saved detailed RNN predictions by removing
   additional characters and
18     translating character indices back to legible strings.
19     args:
20         corpus: string, name of the corpus
21         experiment: string, name of the experiment
22         network: string, name of the network architecture
23         hidden_units: int, number of hidden units
24     returns:
25         data_matrix: list of clean data, ready to be turned into a pd.
   DataFrame'''
26     detail_file = os.path.join '..', 'exp_results', corpus, experiment,
   'detail_{}_{}.csv'.format(network, hidden_units)
27     training_df = pd.read_csv('../training/{}.csv'.format(corpus),
   delimiter=',', encoding='latin-1')
28
29     table = str.maketrans(dict.fromkeys(' []')) # Remove superfluous
   characters
30     lines = []
31     with open(detail_file, 'r') as f:
32         for line in f:
33             line = line.translate(table).strip().split(',')
34             lines.append(line)
35     prev_line = []
36     mod_lines = []
37     # Merge lines split by \n in saving
38     for i in range(2, len(lines), 2):
39         line = lines[i]
40         prev_line = lines[i-1]
41         prev_line.extend(line)
```

```

42     mod_lines.append(prev_line)
43     # Decode numbers to characters
44     if experiment == 'LRD':
45         num_to_char = {'1': ']', '2': '[', '3': '{', '4': '}', '5': '$'}
46     else:
47         num_to_char = {'1': ']', '2': '{', '3': '[', '4': '}', '5': '$'}
48     # Prepare matrix to be turned into a df
49     data_matrix = [[line[0].translate(line[0].maketrans(num_to_char))+
50                     line[1].translate(line[1].maketrans(num_to_char)), float(line[2]),
51                     int(line[3])] for line in mod_lines]
52
53     return data_matrix
54
55 def prepare_dataframe(corpus, experiment, network, hidden_units,
56                       GLOBAL_WORDS):
57     '''Processes the data_matrix from numbers_to_words into a pd.
58     DataFrame for further analysis.
59     args:
60         detail_file: Path to a prediction-per-word file
61     returns:
62         df: pd.DataFrame with all columns relevant for analysis'''
63     data_matrix = numbers_to_words(corpus, experiment, network,
64                                    hidden_units)
65     df = pd.DataFrame(data_matrix, columns=['words', 'predictions', '
66         golds'])
67     df['words'] = GLOBAL_WORDS
68     # Append columns determining the
69     # - kind of error in the word if the word is wrong
70     # - position of the corrupted character if it can be determined
71     # - nesting depth at error position
72     # - running bracket distance at error position
73     df['max_valid_nesting_depth'] = df.words.map(lambda word:
74         corpus_tools.maxValidNestingDepth(word))
75     df['error'] = df.words.map(lambda word: corpus_tools.determineError(
76         word))
77     df['error_pos'] = df.words.loc[df.error != 'none'].map(lambda word:
78         corpus_tools.findErrorPosition(word))
79     df['error_depth'] = df.words.loc[df.error != 'none'].map(lambda word
80         : corpus_tools.nestingDepthAtPosition(word, corpus_tools.
81         findErrorPosition(word)))
82     df['error_distance'] = df.words.loc[df.error != 'none'].map(lambda
83         word: corpus_tools.bracketDistanceAtPosition(word, corpus_tools.
84         findErrorPosition(word)))
85     #print(df.head())
86     return df
87
88 def measure_performance(results):
89     '''Calculates precision, recall and F1 score of a dataframe

```

```

    containing predictions and gold labels.
77     args:
78         results: Dataframe containing predictions and gold labels
79     returns:
80         precision: TPs/(TPs+FPs)
81         recall: TPs/(TPs+FNs)
82         f1_score: 2*((precision*recall)/(precision+recall))'''
83     true_pos = len(results.loc[results.golds == 1].loc[results.
        predictions >= 0.5])
84     false_pos = len(results.loc[results.golds == 0].loc[results.
        predictions >= 0.5])
85     true_neg = len(results.loc[results.golds == 0].loc[results.
        predictions < 0.5])
86     false_neg = len(results.loc[results.golds == 1].loc[results.
        predictions < 0.5])
87     if true_pos+false_pos == 0:
88         precision = 0.
89         recall = 0.
90         f1_score = 0.
91     else:
92         precision = float(true_pos/(true_pos+false_pos))
93         recall = float(true_pos/(true_pos+false_neg))
94         if precision+recall == 0.:
95             f1_score = 0.
96         else:
97             f1_score = 2.*((precision*recall)/(precision+recall))
98     return precision, recall, f1_score
99
100 def extend_performance_measures():
101     '''Iterates through all files containing experiment results (
        accuracy, loss) and their corresponding detail files, which
        contain every single word, prediction and gold label. From that,
        precision, recall and F1 score are calculated and added to the
        experiment results file.
102     args:
103         none
104     returns:
105         none'''
106     corpora = ['base', 'high', 'low']
107     experiments = ['LRD', 'ND']
108     networks = ['SRNN', 'GRU', 'LSTM']
109     for corpus in corpora:
110         for experiment in experiments:
111             for network in networks:
112                 for hidden_units in [2**i for i in range(1,10)]:
113                     # Appending precision, recall, f1_score to sparse_file
114                     detail_file = os.path.join '..', 'exp_results', corpus,
                        experiment, 'detail_{}_{}.csv'.format(network, str(hidden_units)))

```



```

115         sparse_file = os.path.join('..', 'exp_results', corpus,
experiment, '{}_{}.csv'.format(network, str(hidden_units)))
116         sparse = pd.read_csv(sparse_file)
117         # Processing values in the details file to calculate
additional performance measures
118         details = numbers_to_words(detail_file)
119         precision, recall, f1_score = measure_performance(details)
120         sparse['precision'], sparse['recall'], sparse['f1_score'] =
precision, recall, f1_score
121         sparse.to_csv(sparse_file, index=False)
122         # Saving the extended values
123         prepend_header(detail_file)
124
125 def breakdownPredictions(df, corpus, experiment, network, hidden_units
, performance):
126     '''Transforms a full dataframe of word-prediction-gold_label data
into 4 dataframes specific to a single model: true positives,
false positives, true negatives and false negatives.
127     args:
128         df: Complete dataframe of word-prediction-gold_labels
129         corpus, experiment, network, hidden_units, performance: Strings
filtering the df for the specific model
130     returns:
131         true_positives, false_positives, true_negatives, false_negatives
: Dataframes containing all TPs, FPs, TNs and FNs of a specific
model'''
132     generals = pd.DataFrame()
133     generals['corpus'], generals['experiment'], generals['network'],
generals['hidden_units'], generals['performance'] = pd.Series(
corpus), experiment, network, hidden_units, performance
134     # Determine all predictions in their own series to analyze.
135     true_positives = df.loc[df.golds == 1].loc[df.predictions >= 0.5]
136     false_positives = df.loc[df.golds == 0].loc[df.predictions >= 0.5]
137     true_negatives = df.loc[df.golds == 0].loc[df.predictions <= 0.5]
138     false_negatives = df.loc[df.golds == 1].loc[df.predictions <= 0.5]
139
140     true_positives = true_positives.join(generals)
141     false_positives = false_positives.join(generals)
142     true_negatives = true_negatives.join(generals)
143     false_negatives = false_negatives.join(generals)
144
145     for column in generals.columns:
146         true_positives[column] = true_positives[column].fillna(generals[
column][0])
147         false_positives[column] = false_positives[column].fillna(generals[
column][0])
148         true_negatives[column] = true_negatives[column].fillna(generals[
column][0])

```

```

149     false_negatives[column] = false_negatives[column].fillna(generals[
150         column][0])
151
152     return true_positives, false_positives, true_negatives,
153         false_negatives
154
155 def prepend_header(detail_file):
156     '''Includes a descriptive header in a file containing individual
157         words, predictions and their gold label.
158         args:
159             detail_file: Path to a file containing individual words,
160                 predictions and their gold label
161         returns:
162             none'''
163     f = open(detail_file, 'r')
164     temp = f.read()
165     f.close()
166
167     f = open(detail_file, 'w')
168     f.write('words,predictions,golds\n')
169
170     f.write(temp)
171     f.close()
172
173 def create_mega_df():
174     '''Creates and saves four dataframes containing every single word,
175         prediction and gold label for every outlier model (accuracy
176         >55%/<45%), split into true positives, false positives, true
177         negatives and false negatives.'''
178     global_words_LRD = numbers_to_words('base', 'LRD', 'LSTM', '8')
179     df = pd.DataFrame(global_words_LRD, columns=['words', 'predictions',
180         'golds'])
181     GLOBAL_WORDS_LRD = df['words']
182     global_words_ND = numbers_to_words('base', 'ND', 'SRNN', '2')
183     df = pd.DataFrame(global_words_ND, columns=['words', 'predictions',
184         'golds'])
185     GLOBAL_WORDS_ND = df['words']
186
187     true_positives_frames = []
188     false_positives_frames = []
189     true_negatives_frames = []
190     false_negatives_frames = []
191     valid_nesting_depth_frames = []
192     error_pos_frames = []
193     error_depth_frames = []
194     networks = ['base LRD LSTM 8 good', 'base LRD GRU 2 good', 'base LRD
195         GRU 128 good',
196         'base LRD GRU 32 bad', 'base LRD LSTM 128 bad', 'base LRD LSTM 16

```

```

187     'bad', 'base LRD SRNN 16 bad',
188     'low LRD GRU 2 good', 'low LRD SRNN 4 good', 'low LRD LSTM 16 good',
189     'low LRD GRU 64 good',
190     'low LRD LSTM 4 bad',
191     'high LRD GRU 512 good',
192     'high LRD LSTM 8 bad', 'high LRD GRU 8 bad', 'high LRD LSTM 4 bad',
193     'base ND LSTM 128 good', 'base ND LSTM 32 good', 'base ND GRU 512 good',
194     'base ND SRNN 2 good', 'base ND GRU 4 good',
195     'base ND SRNN 128 bad', 'base ND SRNN 256 bad',
196     'low ND LSTM 512 good', 'low ND GRU 64 good', 'low ND SRNN 32 good',
197     'low ND LSTM 16 good', 'low ND GRU 4 good', 'low ND LSTM 8 good',
198     'low ND SRNN 4 bad', 'low ND LSTM 32 bad',
199     'high ND SRNN 16 good',
200     'high ND LSTM 64 bad']
201
202 for entry in networks:
203     corpus, experiment, network, hidden_units, performance = entry.split()
204     # Set translation table for numbers_to_words
205     if experiment == 'LRD':
206         GLOBAL_WORDS = GLOBAL_WORDS_LRD
207     else:
208         GLOBAL_WORDS = GLOBAL_WORDS_ND
209     details = prepare_dataframe(corpus, experiment, network,
210                                hidden_units, GLOBAL_WORDS)
211     true_positives, false_positives, true_negatives, false_negatives =
212         breakdownPredictions(details, corpus, experiment, network,
213                               hidden_units, performance)
214     true_positives_frames.append(true_positives)
215     false_positives_frames.append(false_positives)
216     true_negatives_frames.append(true_negatives)
217     false_negatives_frames.append(false_negatives)
218
219 # Concatenate all results
220 mega_true_positives = pd.concat(true_positives_frames)
221 mega_false_positives = pd.concat(false_positives_frames)
222 mega_true_negatives = pd.concat(true_negatives_frames)
223 mega_false_negatives = pd.concat(false_negatives_frames)
224
225 # Save results
226 tp_out = os.path.join('.', 'results', 'mega_true_positives.csv')
227 fp_out = os.path.join('.', 'results', 'mega_false_positives.csv')
228 tn_out = os.path.join('.', 'results', 'mega_true_negatives.csv')
229 fn_out = os.path.join('.', 'results', 'mega_false_negatives.csv')
230 mega_true_positives.to_csv(tp_out)
231 mega_false_positives.to_csv(fp_out)
232 mega_true_negatives.to_csv(tn_out)

```

```

225 mega_false_negatives.to_csv(fn_out)
226
227 def count_error_types(fp, experiment, network, corpus, performance):
228     '''Determines number of open/closed bracket error words in a
229         dataframe.
230         args:
231             fp: dataframe containing either false positives or true
232                 negatives (since other categories do not have incorrect words)
233             experiment: string, name of the experiment
234             network: string, name of the network architecture
235             hidden_units: int, number of hidden units
236         returns:
237             none'''
238     open = fp.loc[fp.experiment == experiment].loc[fp.network == network
239         ].loc[fp.corpus == corpus].loc[fp.performance == performance].loc[
240             fp.error == 'open'].count()[1]
241     closed = fp.loc[fp.experiment == experiment].loc[fp.network ==
242         network].loc[fp.corpus == corpus].loc[fp.performance ==
243         performance].loc[fp.error == 'closed'].count()[1]
244     total = open + closed
245     if total:
246         if closed:
247             ratio = open/closed
248         else:
249             ratio = np.inf
250
251     print("{}{},{},{},{},{},{},{},{}".format(network, experiment, corpus
252         .capitalize(), performance, open, closed, total, ratio))
253
254 def print_bracket_ratio_table():
255     '''Creates a table splitting false positives into open/closed error
256         types for each outlier network, sorted by 'good' and 'bad'
257         outliers.'''
258     fp_in = os.path.join '..', 'results', 'mega_false_positives.csv'
259     fp = pd.read_csv(fp_in)
260     experiments = ['LRD', 'ND']
261     networks = ['SRNN', 'LSTM', 'GRU']
262     corpora = ['base', 'low', 'high']
263     print("network,experiment,corpus,performance,open,closed,total,ratio
264         ")
265     for experiment in experiments:
266         for corpus in corpora:
267             for network in networks:
268                 count_error_types(fp, experiment, network, corpus, 'good')
269     for experiment in experiments:
270         for corpus in corpora:
271             for network in networks:
272                 count_error_types(fp, experiment, network, corpus, 'bad')

```

```
263  
264 create_mega_df()  
265 print_bracket_ratio_table()
```

Listing 5: Script to create the basis of the tables featured in Section 4.4.

Eidesstattliche Erklärung

Eidesstattliche Erklärung zur Bachelorarbeit

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