# Visualizing and Understanding Recurrent Networks

**Andrej Karpathy**[*]        **Justin Johnson**[*]        **Li Fei-Fei**
Department of Computer Science, Stanford University
{karpathy,jcjohns,feifeili}@cs.stanford.edu

## Abstract

Recurrent Neural Networks (RNNs), and specifically a variant with Long Short-Term Memory (LSTM), are enjoying renewed interest as a result of successful applications in a wide range of machine learning problems that involve sequential data. However, while LSTMs provide exceptional results in practice, the source of their performance and their limitations remain rather poorly understood. Using character-level language models as an interpretable testbed, we aim to bridge this gap by providing a comprehensive analysis of their representations, predictions and error types. In particular, our experiments reveal the existence of interpretable cells that keep track of long-range dependencies such as line lengths, quotes and brackets. Moreover, an extensive analysis with finite horizon $n$-gram models suggest that these dependencies are actively discovered and utilized by the networks. Finally, we provide detailed error analysis that suggests areas for further study.

## 1   Introduction

Recurrent Neural Networks, and specifically a variant with Long Short-Term Memory (LSTM) (14), have recently emerged as an effective model in a wide variety of applications that involve sequential data. These include language modeling (22), handwriting recognition and generation (9), machine translation (28; 1), speech recognition (10), video analysis (7) and image captioning (30; 19).

However, the source of their impressive performance remains poorly understood. This raises concerns of interpretability and limits our ability design better architectures. A few recent ablation studies analyzed the effects on performance as they removed or modified various gates and connections (12; 5). While this analysis illuminates the performance-critical pieces of the architecture, it is still limited to measuring the test set perplexity alone.

An often cited advantage of the LSTM architecture is that it can store and retrieve information over long time scales using its gating mechanisms. However, it is not immediately clear that good solutions can be discovered with stochastic gradient descent and truncated backpropagation through time, and that these mechanisms are effectively utilized on real-world data. In this work we use character-level language models as an interpretable testbed for illuminating the long-range dependencies learned by the LSTM architecture. Our analysis reveals the existence of cells that robustly identify interpretable, high-level patterns such as line lengths, brackets and quotes. We further quantify the LSTM predictions with comprehensive comparison to $n$-gram models, where we find that LSTMs perform significantly better on characters that require long-range reasoning. Finally, we conduct a detailed error analysis in which we *"peel the onion"* of errors with a sequence of oracles. These results allow us to quantify the extent of remaining errors in several categories and to suggest specific areas for further study.

## 2   Related Work

**Recurrent Networks**. Recurrent Neural Networks (RNNs) have a long history of applications in various sequence learning tasks (31; 26; 25). Despite their early successes, the difficulty of training

---

[*]Both authors contributed equally to this work.

1

simple recurrent networks (2; 24) has encouraged various proposals for improvements to their basic architecture. Among the most successful variants are the Long Short Term Memory networks (14), which can in principle store and retrieve information over long time periods with explicit gating mechanisms and a built-in constant error carousel. In the recent years there has been a renewed interest in further improving on the basic architecture by modifying the functional form as seen with Gated Recurrent Units (4), incorporating content-based soft attention mechanisms (1; 32), push-pop stacks (18), or more generally external memory arrays with both content-based and relative addressing mechanisms (11). In this work we focus the majority of our analysis on the LSTM due to its widespread popularity and proven track record.

**Understanding Recurrent Networks**. While there is an abundance of work that modifies or extends the basic LSTM architecture, relatively little attention has been paid to understanding its computational properties or the source of its performance. Greff et al. (12) recently conducted a comprehensive study of different LSTM components and concluded that the forget gates are its most critical components. Chung et al. evaluated GRU compared to LSTMs (5). Pascanu et al. examined the effects of depth (23). These approaches study recurrent network architectures based on variations in the final performance, while our approach dives deep into the statistical patterns in their activations and predictions. We also conduct detailed error analysis that breaks down the the final performance into interpretable categories.

## 3 Experimental Setup

We first describe three variants of a deep recurrent network (RNN, LSTM and the GRU models), then explain how they are used in sequence learning and finally describe the optimization.

### 3.1 Recurrent Neural Network Models

The simplest instantiation of a deep recurrent network arranges hidden state vectors $h_t^l$ in a two-dimensional grid, where $t = 1 \dots T$ is thought of as time and $l = 1 \dots L$ is the depth. The bottom row of vectors $h_t^0 = x_t$ at depth zero holds the input vectors $x_t$ and each vector in the top row $\{h_t^L\}$ is used to predict an output vector $y_t$. All intermediate vectors $h_t^l$ are computed as a function of $h_{t-1}^l$ and $h_t^{l-1}$. Through these hidden vectors, each output $y_t$ at some particular time step $t$ becomes a function of all input vectors up to that time, $\{x_1, \dots, x_t\}$. The precise mathematical form of the recurrence $(h_{t-1}^l, h_t^{l-1}) \to h_t^l$ varies from model to model and we describe these details next.

**Vanilla Recurrent Neural Network** (RNN) has a recurrence of the form

$$h_t^l = \tanh W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

where we assume that all $h \in \mathbb{R}^n$. The parameter matrix $W^l$ on each layer has dimensions $[n \times 2n]$ and $\tanh$ is applied elementwise. Note that $W$ varies between layers but is shared through time. We omit the bias vectors for brevity. Interpreting the equation above, the inputs from the layer below $h_t^{l-1}$ and before in time $h_{t-1}^l$ are transformed and interact through additive interaction before the non-linearity. This is known to be a weak form of coupling (27). Both the LSTM and the GRU (discussed next) include more powerful multiplicative interactions.

**Long Short-Term Memory** (LSTM) (14) was designed to address the difficulties of training RNNs (2). In particular, it was observed that the backgropagation dynamics led the gradients in an RNN to either vanish or explode. It was later found that the exploding gradient concern can be alleviated with a heuristic of clipping the gradients at some maximum value (24). On the other hand, LSTMs were designed to mitigate the vanishing gradient problem. In addition to a hidden state vector $h_t^l$, LSTMs also maintain a memory vector $c_t^l$. At each time step the LSTM can choose to read from, write to, or reset the cell using explicit gating mechanisms. The precise form of the update is as follows:

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \mathrm{sigm} \\ \mathrm{sigm} \\ \mathrm{sigm} \\ \tanh \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix} \qquad \begin{aligned} c_t^l &= f \odot c_{t-1}^l + i \odot g \\ h_t^l &= o \odot \tanh(c_t^l) \end{aligned}$$

Here, the sigmoid function $\mathrm{sigm}$ and $\tanh$ are applied element-wise, and $W^l$ is a $[4n \times 2n]$ matrix. The three vectors $i, f, o \in \mathbb{R}^n$ are each thought of as binary gates that control whether each memory cell is updated, whether it is reset to zero, and whether its local state is revealed in the hidden vector,

respectively. The activations of these gates are based on the sigmoid function and hence allowed to range smoothly between zero and one to keep the model differentiable. The vector $g \in \mathbb{R}^n$ ranges between -1 and 1 and is used to additively modify the memory contents. This additive interaction is a critical feature of the LSTM's design, because during backpropagation a sum operation merely distributes gradients. This allows gradients on the memory cells $c$ to flow backwards through time uninterrupted for long time periods, or at least until the flow is disrupted with the multiplicative interaction of an active forget gate. Lastly, note that an implementation of the LSTM requires one to maintain both $h_t^l$ and $c_t^l$ at every point in the network.

**Gated Recurrent Unit** (GRU) (4) was recently proposed as a simpler alternative to the LSTM:

$$\begin{pmatrix} r \\ z \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \end{pmatrix} W_r^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix} \qquad \begin{aligned} \tilde{h}_t^l &= \tanh(W_x^l h_t^{l-1} + W_g^l(r \odot h_{t-1}^l)) \\ h_t^l &= (1-z) \odot h_{t-1}^l + z \odot \tilde{h}_t^l \end{aligned}$$

Here, $W_r^l$ are $[2n \times 2n]$ and $W_g^l$ and $W_x^l$ are $[n \times n]$. The GRU has the interpretation of computing a *candidate* hidden vector $\tilde{h}_t^l$ and then smoothly interpolating towards it, as gated by $z$.

### 3.2 Character-level Language Modeling

We use character-level language modeling as an interpretable testbed for sequence learning. In this setting, the input to the network is a sequence of characters and the network is trained to predict the next character in the sequence with a Softmax classifier at each time step. Concretely, assuming a fixed vocabulary of $K$ characters we encode all characters with $K$-dimensional 1-of-$K$ vectors $\{x_t\}, t = 1, \ldots, T$, and feed these to the recurrent network to obtain a sequence of $D$-dimensional hidden vectors at the last layer of the network $\{h_t^L\}, t = 1, \ldots, T$. To obtain predictions for the next character in the sequence we project this top layer of activations to a sequence of vectors $\{y_t\}$, where $y_t = W_y h_t^L$ and $W_y$ is a $[K \times D]$ parameter matrix. These vectors are interpreted as holding the (unnormalized) log probability of the next character in the sequence and the objective is to minimize the average cross-entropy loss over all targets.

### 3.3 Optimization

Following previous work (28) we initialize all parameters uniformly in range $[-0.08, 0.08]$. We use mini-batch stochastic gradient descent with batch size 100 and RMSProp (6) per-parameter adaptive update with base learning rate $2 \times 10^{-3}$ and decay 0.95. These settings work robustly with all of our models. The network is unrolled for 100 time steps. We train each model for 50 epochs and decay the learning rate after 10 epochs by multiplying it with a factor of 0.95 after each additional epoch. We use early stopping based on validation performance and cross-validate the amount of dropout for each model individually. Our supplementary material contains additional details regarding efficient implementation and further experiments with other initialization ranges and updates.

## 4 Experiments

**Datasets**. Two datasets previously used in the context of character-level language models are the Penn Treebank dataset (20) and the Hutter Prize 100MB of Wikipedia dataset (16) . However, both datasets contain a mix of common language and special markup. Our goal is not to compete with previous work but rather to study recurrent networks in a controlled setting and on both ends on the spectrum of degree of structure. Therefore, we chose to use Leo Tolstoy's *War and Peace* (WP) novel, which consists of 3,258,246 characters of almost entirely English text with minimal markup, and at the other end of the spectrum the source code of the *Linux Kernel* (LK). We shuffled all header and source files randomly and concatenated them into a single file to form the 6,206,996 character long dataset. We split the data into train/val/test splits as 80/10/10 for WP and 90/5/5 for LK. Therefore, there are approximately 300,000 characters in the validation/test splits in each case. The total number of characters in the vocabulary is 87 for WP and 101 for LK.

### 4.1 Comparing Recurrent Networks

We first train several recurrent network models to support further analysis and to compare their performance in a controlled setting. In particular, we train models in the cross product of type (LSTM/RNN/GRU), number of layers (1/2/3), number of parameters (4 settings), and both datasets (WP/KL). The 4 parameter sizes were chosen to be in units of 1-layer LSTMs with 64, 128, 256, and 512 cells. With our character vocabulary sizes this gives approximately 50K, 130K, 400K, and 1.3M parameters respectively. The sizes of hidden layers in the other models were always chosen to match the closest prototype size as close as possible.

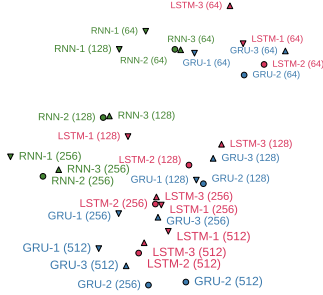| Layers | LSTM | | | RNN | | | GRU | | |
|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 1 | 2 | 3 | 1 | 2 | 3 |
| Size | War and Peace Dataset | | | | | | | | |
| 64 | 1.449 | 1.442 | 1.540 | 1.446 | 1.401 | 1.396 | 1.398 | **1.373** | 1.472 |
| 128 | 1.277 | **1.227** | 1.279 | 1.417 | 1.286 | 1.277 | 1.230 | **1.226** | 1.253 |
| 256 | 1.189 | **1.137** | 1.141 | 1.342 | 1.256 | 1.239 | 1.198 | 1.164 | **1.138** |
| 512 | 1.161 | 1.092 | 1.082 | - | - | - | 1.170 | 1.201 | **1.077** |
| | Linux Kernel Dataset | | | | | | | | |
| 64 | 1.355 | **1.331** | 1.366 | 1.407 | 1.371 | 1.383 | 1.335 | 1.298 | 1.357 |
| 128 | 1.149 | **1.128** | 1.177 | 1.241 | 1.120 | 1.220 | 1.154 | 1.125 | 1.150 |
| 256 | 1.026 | **0.972** | 0.998 | 1.171 | 1.116 | 1.116 | 1.039 | 0.991 | 1.026 |
| 512 | 0.952 | 0.840 | 0.846 | - | - | - | 0.943 | 0.861 | **0.829** |



Figure 1: **Left:** The **test set cross-entropy loss** for all models and datasets (low is good). Models in each row have nearly equal number of parameters. The test set has 300,000 characters. The standard deviation, estimated with 100 bootstrap samples, is less than $4 \times 10^{-3}$ in all cases. **Right:** A t-SNE embedding based on the probabilities assigned to test set characters by each model on War and Peace. The color, size, and marker correspond to model type, model size, and number of layers.

The test set results are shown in Figure 1. Our consistent finding is that depth of at least two is beneficial. However, between two and three layers our results are mixed. Additionally, the results are mixed between the LSTM and the GRU, but both significantly outperform the RNN. We also computed the fraction of times that each pair of models agree on the most likely character and use it to render a t-SNE (29) embedding (we found this more stable and robust than the KL divergence). The plot (Figure 1, right) further supports the claim that the LSTM and the GRU make similar predictions while the RNNs form their own cluster.

## 4.2 Internal Mechanisms of an LSTM

**Interpretable, long-range LSTM cells.** LSTMs can in principle use their memory cells to remember long-range information and keep track of various attributes of text they are in. For instance, it is a simple exercise to write down toy cell weights that would allow the cell to keep track of whether it is inside a quoted string. However, to our knowledge, the existence of such cells has never been experimentally demonstrated on real-world data. In particular, it could be argued that even if the LSTM is in principle capable of using these operations, practical optimization challenges (i.e. SGD dynamics, or approximate gradients due to truncated backpropagation through time) might prevent it from discovering these solutions. In this experiment we verify that multiple interpretable cells do in fact exist in these networks. We show several examples in Figure 2. Note that truncated backpropagation prevents the gradient signal from noticing dependencies longer than 100 characters, but we still observe cells that reliably keep track of quotes or comment blocks for periods much longer than 100 characters. We hypothesize that these cells first develop on patterns shorter than 100 characters but then generalize to longer sequences.

**Gate activation statistics**. We also recorded and examined the statistics of gate activations in the networks while evaluating the test set. We were particularly interested in studying the distributions of saturation regimes in the networks, where we define a gate to be left or right-saturated if its activation is less than 0.1 or more than 0.9, respectively. We then compute the fraction of times that each gate's activation falls into either category and show these statistics in Figure 3. The amount of right-saturated forget gate activations in an LSTM is particularly interesting, since this corresponds to cells that remember their value from the previous iteration. For instance, note that there are multiple cells that are almost always right-saturated and hence function as perfect integrators. Conversely, there are no cells that function in purely feed-forward fashion, as these would show as consistently left-saturated forget gates. The output gate statistics also reveal that there are no cells that get consistently revealed to the hidden state. Lastly, a surprising fact is that unlike the other two layers, the activations in the first layer are diffuse. This is a finding that we struggle to explain but it is present across all of our models, including GRUs.

## 4.3 Understanding Long-Range Interactions

Good performance of LSTMs is frequently attributed to their ability to store long-range information. In this section we test this hypothesis by comparing an LSTM with baseline models that can only utilize information from a fixed number of previous steps. In particular, we consider two baselines:

Figure 2: Several examples of cells with interpretable activations discovered in our best Linux Kernel and War and Peace LSTMs. Text color corresponds to $tanh(c)$, where -1 is red and +1 is blue.

Figure 3: **Left three:** Saturation plots for an LSTM. Each circle is a gate in the LSTM and its position is determined by the fraction of time it is left and right-saturated. These fractions must add to at most one (indicated by the diagonal line). **Right two:** Saturation plot for a 3-layer GRU model.

*1. n-NN*: A fully-connected neural network with one hidden layer and $tanh$ nonlinearities. The input to the network is a sparse binary vector of dimension $nK$ that concatenates the one-of-$K$ encodings of $n$ consecutive characters. We optimize the model as described in Section 3.3 and cross-validate the size of the hidden layer.

*2. n-gram*: An unpruned $(n + 1)$-gram character-level language model using modified Kneser-Ney smoothing (3). This is a standard smoothing method for language models (15). All models were trained using the popular KenLM software package (13).

**Performance comparisons.** The performance of both $n$-gram models is shown in Table 2. The $n$-gram and $n$-NN models perform nearly identically for small values of $n$, but for larger values the $n$-NN models start to overfit and the $n$-gram model performs better. Moreover, we see that on both datasets our best recurrent network outperforms the 20-gram model (1.077 vs. 1.195 on WP and 0.84 vs.0.889). It is difficult to make a direct model size comparison, but the 20-gram model file has 3GB, while our largest checkpoints are 11MB. However, the assumptions encoded in the Kneser-Ney smoothing model are intended for word-level modeling of natural language and may not be optimal for character-level data. Despite this concern, these results provide some evidence that the recurrent networks are effectively utilizing information beyond 20 characters.

**Error Analysis.** It is instructive to delve deeper into the errors made by both recurrent networks and $n$-gram models. In particular, we define a character to be an error if the probability assigned to it by a model is below 0.5. Figure 4 (left) shows the overlap between the test-set errors for the 3-layer LSTM, and the best $n$-NN and $n$-gram models. We see that the majority of errors are shared by all three models, but each model also has its own unique errors.

| $n$<br>Model | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 20 |
|---|---|---|---|---|---|---|---|---|---|---|
| | War and Peace Dataset | | | | | | | | | |
| $n$-gram | 2.399 | 1.928 | 1.521 | 1.314 | 1.232 | 1.203 | **1.194** | 1.194 | 1.194 | 1.195 |
| $n$-NN | 2.399 | 1.931 | 1.553 | 1.451 | 1.339 | **1.321** | - | - | - | - |
| | Linux Kernel Dataset | | | | | | | | | |
| $n$-gram | 2.702 | 1.954 | 1.440 | 1.213 | 1.097 | 1.027 | 0.982 | 0.953 | 0.933 | **0.889** |
| $n$-NN | 2.707 | 1.974 | 1.505 | 1.395 | **1.256** | 1.376 | - | - | - | - |

Table 2: The **test set cross-entropy loss** on both datasets for $n$-gram models (low is good). The standard deviation estimate using 100 bootstrap samples is below $4 \times 10^{-3}$ in all cases.



Figure 4: **Left:** Overlap between test-set errors between our best 3-layer LSTM and the $n$-gram models (low area is good). **Middle/Right:** Mean probabilities assigned to a correct character (higher is better), broken down by the character, and then sorted by the difference between two models. "<s>" is the space character. LSTM (red) outperforms the 20-gram model (blue) on special characters that require long-range reasoning. Middle: LK dataset, Right: WP dataset.
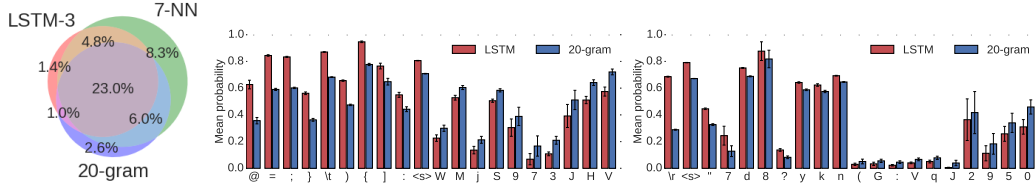
To gain deeper insight into the errors that are unique to the LSTM or the 20-gram model, we compute the mean probability assigned to each character in the vocabulary across the test set. In Figure 4 (middle,right) we display the 10 characters where each model has the largest advantage over the other. On the Linux Kernel dataset, the LSTM displays a large advantage on special characters that are used to structure C programs, including whitespace and brackets. The War and Peace dataset features an interesting long-term dependency with the carriage return, which occurs approximately every 70 characters. Figure 4 (right) shows that the LSTM has a distinct advantage on this character. To accurately predict the presence of the carriage return the model likely needs to keep track of its distance since the last carriage return. The cell example we've highlighted in Figure 2 (top, left) seems particularly well-tuned for this specific task. Similarly, to predict a closing bracket or quotation mark, the model must be aware of the corresponding open bracket, which may have appeared many time steps ago. The fact that the LSTM performs significantly better than the 20-gram model on these characters provides strong evidence that the model is capable of effectively keeping track of long-range interactions.

**Case study: closing brace**. Of these structural characters, the one that requires the longest-term reasoning is the closing brace ("}") on the Linux Kernel dataset. Braces are used to denote blocks of code, and may be nested; as such, the distance between an opening brace and its corresponding closing brace can range from tens to hundreds of characters. This feature makes the closing brace an ideal test case for studying the ability of the LSTM to reason over various time scales. We group closing brace characters on the test set by the distance to their corresponding open brace and compute the mean probability assigned by the LSTM and the 20-gram model to closing braces within each group. The results are shown in Figure 5 (left). We see not only that the LSTM consistently outperforms the 20-gram model at predicting closing braces, but that the performance difference increases up to a distance of 60. After this point the performance difference between the models remains relatively constant, suggesting that no additional time horizon ranges are being captured and that perhaps the LSTM has trouble keeping track of braces longer than this amount.

**Training dynamics.** We can gain further insight into the training dynamics of the LSTM by comparing it with trained $n$-NN models during training using the (symmetric) KL divergence between the predictive distributions on the test set. We visualize the divergence and the difference in the mean loss in Figure 5 (right). Notably, we see that in the first few iterations the LSTM behaves like the 1-NN model but then diverges from it soon after. The LSTM then behaves most like the 2-NN, 3-NN, and 4-NN models in turn. This experiment suggests that the LSTM "grows" its competence over increasingly longer dependencies. We believe that this insight is related to why Sutskever et al. (28)
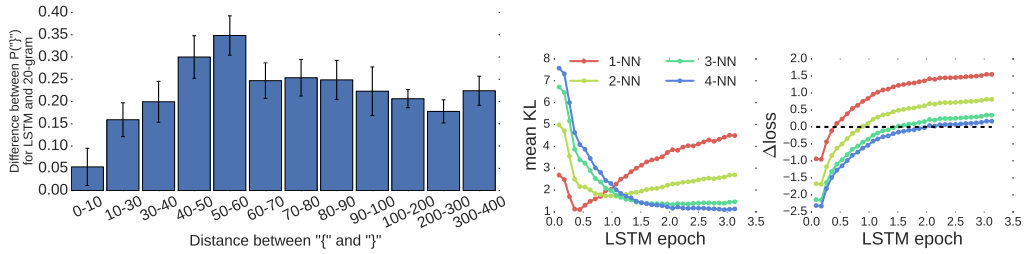
Figure 5: **Left**: Mean difference between the probabilities that the LSTM and 20-gram model assign to the "}" character, bucketed by the distance to the matching "{". **Right**: Comparison of the similarity between 3-layer LSTM and the $n$-NN baselines over the first 3 epochs of training, as measured by the symmetric KL-divergence (middle) and the test set loss (right). Low KL indicates similar predictions, and positive $\Delta$loss indicates that the LSTM outperforms the baseline.

had to reverse the source sentences in their encoder-decoder architecture for machine translation. Without the inversion it is as if the model was immediately presented with $n$-grams of very high $n$. This prevents the RNN from slowly "growing" its competence from shorter to longer dependencies.

## 4.4 Error Analysis: Breaking Down the Failure Cases

In this section we break down LSTM's errors into categories to understand its limitations, the relative severity of each error type, and to suggest areas for further study. We focus on the War and Peace dataset where it is easier to analyze the errors. Our approach is to *"peel the onion"* by iteratively removing the errors with a series of constructed oracles. We discuss these first:

$n$-**gram oracle.** First, we construct optimistic $n$-gram oracles that eliminate errors that might be fixed with better modeling of the previous $n$ characters. In particular, we evaluate our $n$-gram model ($n = 1, \ldots, 9$) and remove a character error if it is correctly classified by any of these models.

**Dynamic $n$-long memory oracle.** Consider the string *"Jon yelled at Mary but Mary couldn't hear him."* One interesting and consistent failure mode is that if the LSTM fails to recognize the first occurrence of *"Mary"* then it will almost always also fail to recognize the second, with an identical pattern of errors. In principle the first mention should make the second easier: the LSTM could conceivably store a summary of previously seen characters in the data and fall back on this memory when it is uncertain. However, this does not appear to take place in practice. This observation is related to the improvements seen in "dynamic evaluation" (21; 17) of recurrent language models, where an RNN is allowed to train on the test set characters during evaluation as long as it sees them only once. This oracle optimistically removes errors in all words (starting with the second character) that can also be found as a substring in the last $n$ characters (we use $n \in \{100, 500, 1000, 50000\}$).

**Rare words oracle.** Next, we construct an oracle that eliminates errors for rare words that occur only up to $n$ times in the training data ($n = 0, \ldots, 5$). Note that a character-level model could in principle hypothesize the spelling of unseen words by recognizing common letter patterns. Errors corrected by this oracle indicate its failure to do so.

**Word model oracle.** We noticed that a large portion of the errors occur on the first character of each word. Intuitively, the task of selecting the next word in the sequence is much harder than completing the last few characters of a known word. We constructed an oracle that eliminated all errors after a space, quote or a newline. Interestingly, a high portion of errors can be found after a newline character, since the models have to learn that new line is identical to a space.

**Punctuation oracle.** This oracle removes errors on all punctuation.

**Boost oracles.** The remaining errors that do not show salient structures or patterns are removed by an oracle that boosts the probability of the correct letter by a fixed amount. These oracles allow us to understand the distribution of the difficulty of the remaining errors.

We now subject two LSTM models to the error analysis: First, our best LSTM model and second, the best LSTM model in the smallest model category (50K parameters). The smaller second model will allow us to quantify the reduction in each error category as we scale up our models. For the following analysis we define a character to be an error when its assigned probability is less than 0.5. We then iterate over oracles in the order presented above, remove each error type in turn and visualize the fractions of errors attributed to each type in Figure 6.

7

**LSTM-3 (512)**

*Pie chart labels (counter-clockwise from top):* 1-gram, 2-gram, 3-gram, 4-gram, 5-gram, 6-gram, 7-gram, 8-gram, 9-gram, 100 memory, 500 memory, 1000 memory, 5000 memory, Word 0-train, Word 1-train, Word 2-train, Word 3-train, Word 4-train, Word 5-train, After space or quote, After newline, Punctuation, 0.1 boost, 0.2 boost, 0.3 boost, 0.4 boost, 0.5 boost

**LSTM-2 (64)**

*Pie chart labels (counter-clockwise from top):* 1-gram, 2-gram, 3-gram, 4-gram, 5-gram, 6-gram, 7-gram, 8-gram, 9-gram, n memory, Word n-train, After space or quote, After newline, Punctuation, 0.1 boost, 0.2 boost, 0.3 boost, 0.4 boost, 0.5 boost

**1-gram to 5-gram**

"My wife," continued Prince Andrew, "is an excellent woman, one of those rare women with whom a man's honor is safe; but, O God, what would I not give now to be unmarried! You are the first and only one to whom I mention this, because I like you."

**Up to 500 memory**

circular, memorandum, or report, skillfully, pointedly, and elegantly. Bilibin's services were valued not only for what he wrote, but also for his skill in dealing and conversing with those in the highest spheres.

Bilibin liked conversation as he liked work, only when it could be made elegantly witty. In society he always awaited an opportunity to say something striking and took part in a conversation only when that was possible. His conversation was always sprinkled with wittily original,

**Less than 3 training examples of word**

Nicholas and Sonya, the niece. Sonya was a slender little brunette with a tender look in her eyes which were veiled by long lashes, thick black plaits coiling twice round her head, and a tawny tint in her complexion and especially in the color of her slender but graceful and muscular arms and neck. By the grace of her movements, by the softness and flexibility of her small limbs, and by a certain coyness and reserve of

**After space or quote**

"No, impossible!" said Prince Andrew, laughing and pressing Pierre's hand to show that there was no need to ask the question. He wished to

**After newline**

Anna Pavlovna smiled and promised to take Pierre in hand. She knew his father to be a connection of Prince Vasili's. The elderly lady who had been sitting with the old aunt rose hurriedly and overtook Prince Vasili

**Punctuation**

"There now!... So you, too, are in the great world?" said he to Pierre.

**0.4 to 0.5 boost**

"Educate this bear for me! He has been staying with me a whole month and this is the first time I have seen him in society. Nothing is so necessary for a young man as the society of clever women."

Figure 6: **Left:** LSTM errors removed one by one with oracles, starting from top of the pie chart and going counter-clockwise. The area of each slice corresponds to fraction of errors contributed. "$n$-memory" refers to dynamic memory oracle with context of $n$ previous characters. "Word $t$-train" refers to the rare words oracle with word count threshold of $t$. **Right:** Concrete examples of text from the test set for each error type. Blue color highlights the relevant characters with the associated error. For the memory category we also highlight the repeated substrings with red bounding rectangles.

**The error breakdown.** The best LSTM model makes a total of 140K errors out of 330K test set characters (42%). Of these, the $n$-gram oracle eliminates 18%, suggesting that the model is not taking full advantage of the last 9 characters. The dynamic memory oracle eliminates 6% of the errors. In principle, a dynamic evaluation scheme could be used to mitigate this type of error, but we believe that a more principled approach could involve a model similar to Memory Networks (32), where the model is allowed to attend to a recent history of the sequence while making its next prediction. Finally, the rare words oracle accounts for 9% of the errors. This error type could be mitigated with transfer learning, or by increasing the size of the training set. The majority of the remaining errors (37%) follow a space, a quote, or a newline. This suggests that substantially longer time dependencies, and possibly hierarchical context models, could eliminate a large fraction of the errors. See Figure 6 for examples of each error type.

**Errors eliminated by scaling up**. The smaller LSTM model makes a total of 184K errors (or 56% of the test set), approximately 44K more than the large model. Surprisingly, 36K of these errors (81%) are $n$-gram errors, 5K come from the boost category, and the remaining 3K are distributed across the other categories relatively evenly. That is, scaling the model up by a factor 26 in the number of parameters has almost entirely provided gains in the local, $n$-gram error rate and has left the other error categories untouched in comparison. This analysis suggests that in order to fully remove all errors it might be necessary to develop new architectural improvements instead of simply scaling up the basic model. For full details of this analysis refer to our supplementary material.

## 5  Conclusion

We have presented a comprehensive analysis of Recurrent Neural Networks and their representations, predictions and error types. In particular, qualitative visualization experiments, cell activation statistics and in-depth comparisons to finite horizon $n$-gram models demonstrate that these networks learn powerful, long-range interactions. We have also conducted a detailed error analysis that illuminates the limitations of recurrent networks and allows us to suggest specific areas for further study. In particular, $n$-gram errors can be significantly reduced by scaling up the models and rare words could be addressed with bigger datasets. However, further architectural innovations may be needed to eliminate the remaining errors.

# References

[1] D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.

[2] Y. Bengio, P. Simard, and P. Frasconi. Learning long-term dependencies with gradient descent is difficult. *Neural Networks, IEEE Transactions on*, 5(2):157–166, 1994.

[3] S. F. Chen and J. Goodman. An empirical study of smoothing techniques for language modeling. *Computer Speech & Language*, 13(4):359–393, 1999.

[4] K. Cho, B. van Merriënboer, D. Bahdanau, and Y. Bengio. On the properties of neural machine translation: Encoder-decoder approaches. *arXiv preprint arXiv:1409.1259*, 2014.

[5] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014.

[6] Y. N. Dauphin, H. de Vries, J. Chung, and Y. Bengio. Rmsprop and equilibrated adaptive learning rates for non-convex optimization. *arXiv preprint arXiv:1502.04390*, 2015.

[7] J. Donahue, L. A. Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, and T. Darrell. Long-term recurrent convolutional networks for visual recognition and description. *CVPR*, 2015.

[8] J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *The Journal of Machine Learning Research*, 12:2121–2159, 2011.

[9] A. Graves. Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850*, 2013.

[10] A. Graves, A.-R. Mohamed, and G. Hinton. Speech recognition with deep recurrent neural networks. In *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*, pages 6645–6649. IEEE, 2013.

[11] A. Graves, G. Wayne, and I. Danihelka. Neural turing machines. *arXiv preprint arXiv:1410.5401*, 2014.

[12] K. Greff, R. K. Srivastava, J. Koutník, B. R. Steunebrink, and J. Schmidhuber. LSTM: A search space odyssey. *CoRR*, abs/1503.04069, 2015.

[13] K. Heafield, I. Pouzyrevsky, J. H. Clark, and P. Koehn. Scalable modified Kneser-Ney language model estimation. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics*, pages 690–696, Sofia, Bulgaria, August 2013.

[14] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.

[15] X. Huang, A. Acero, H.-W. Hon, and R. Foreword By-Reddy. *Spoken language processing: A guide to theory, algorithm, and system development*. Prentice Hall PTR, 2001.

[16] M. Hutter. The human knowledge compression contest. 2012.

[17] F. Jelinek, B. Merialdo, S. Roukos, and M. Strauss. A dynamic language model for speech recognition. In *HLT*, volume 91, pages 293–295, 1991.

[18] A. Joulin and T. Mikolov. Inferring algorithmic patterns with stack-augmented recurrent nets. *CoRR*, abs/1503.01007, 2015.

[19] A. Karpathy and L. Fei-Fei. Deep visual-semantic alignments for generating image descriptions. *CVPR*, 2015.

[20] M. P. Marcus, M. A. Marcinkiewicz, and B. Santorini. Building a large annotated corpus of english: The penn treebank. *Computational linguistics*, 19(2):313–330, 1993.

[21] T. Mikolov. Statistical language models based on neural networks. *Presentation at Google, Mountain View, 2nd April*, 2012.

[22] T. Mikolov, M. Karafiát, L. Burget, J. Cernockỳ, and S. Khudanpur. Recurrent neural network based language model. In *INTERSPEECH 2010, 11th Annual Conference of the International Speech Communication Association, Makuhari, Chiba, Japan, September 26-30, 2010*, pages 1045–1048, 2010.

[23] R. Pascanu, Ç. Gülçehre, K. Cho, and Y. Bengio. How to construct deep recurrent neural networks. *CoRR*, abs/1312.6026, 2013.

[24] R. Pascanu, T. Mikolov, and Y. Bengio. On the difficulty of training recurrent neural networks. *arXiv preprint arXiv:1211.5063*, 2012.

[25] D. E. Rumelhart, G. E. Hinton, and R. J. Williams. Learning internal representations by error propagation. Technical report, DTIC Document, 1985.

[26] J. Schmidhuber. Deep learning in neural networks: An overview. *Neural Networks*, 61:85–117, 2015. Published online 2014; based on TR arXiv:1404.7828 [cs.NE].

[27] I. Sutskever, J. Martens, and G. E. Hinton. Generating text with recurrent neural networks. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pages 1017–1024, 2011.

[28] I. Sutskever, O. Vinyals, and Q. V. Le. Sequence to sequence learning with neural networks. In *Advances in Neural Information Processing Systems*, pages 3104–3112, 2014.

[29] L. Van der Maaten and G. Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(2579-2605):85, 2008.

[30] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan. Show and tell: A neural image caption generator. *CVPR*, 2015.

[31] P. J. Werbos. Generalization of backpropagation with application to a recurrent gas market model. *Neural Networks*, 1(4):339–356, 1988.

[32] J. Weston, S. Chopra, and A. Bordes. Memory networks. *CoRR*, abs/1410.3916, 2014.

## Supplementary Material

**Exact Model Sizes**

| Layers | LSTM | | | RNN | | | GRU | | |
|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 1 | 2 | 3 | 1 | 2 | 3 |
| Size | | | | War and Peace Dataset | | | | | |
| 64 | 64 | 45 | 37 | 141 | 96 | 78 | 77 | 53 | 44 |
| 128 | 128 | 89 | 68 | 273 | 174 | 139 | 151 | 98 | 79 |
| 256 | 256 | 159 | 126 | 531 | 325 | 257 | 300 | 185 | 146 |
| 512 | 512 | 308 | 241 | 1045 | 623 | 487 | 596 | 357 | 280 |
| | | | | Linux Kernel Dataset | | | | | |
| 64 | 64 | 46 | 38 | 142 | 98 | 81 | 78 | 54 | 45 |
| 128 | 128 | 85 | 69 | 274 | 178 | 143 | 152 | 100 | 81 |
| 256 | 256 | 161 | 128 | 534 | 329 | 260 | 300 | 187 | 149 |
| 512 | 512 | 310 | 243 | 1048 | 628 | 492 | 596 | 359 | 282 |

Table 3: Exact sizes of the models used in all experiments. The models in each row have approximately the same number of parameters.


**Training Details**

In this section we provide additional details regarding training.

**Efficient implementation**. It is impractical to compute the loss function over the entire training sequence, which can in practice be be millions of characters long. Instead, a good strategy is to take the input sequence of characters are reshape it into one data matrix $X$ of size $[B \times N]$, where $B$ is the batch size (we use $B = 100$). We then march from left to right in chunks of length $S$ (we use $S = 100$), and every non-overlapping block of $[B \times S]$ is fed as input to the network. The ground truth labels matrix $Y$ is the same as $X$, but shifted by 1 columnwise (i.e. next character in the sequence). It is important to keep track of the hidden state activations $h_T^l$ at the last time step in the previous batch and copy them into the vectors $h_0^l$ in the next batch. We evaluate the loss and backpropagate only on the $[B \times S]$ chunk of data. In particular, note that the backpropagated signal is truncated $S$ steps backwards in time, limiting the extent of long term interactions that the networks can learn. In practice the networks do learn much longer dependencies (as supported by our qualitative cell visualization experiments), but this can liekly only happen if the same dependencies also occur on time scales less than $S$ and then happen to generalize to longer sequences.

**Initialization**. Following previous work (28) we initialize all parameters uniformly in the range $[-0.08, 0.08]$. We experimented with multiple settings of this hyperparameter on multiple architectures and found its sensitivity to be relatively low. In particular, It is safer to err on the side of making this range smaller, since some models start to diverge when it is greater than 1, but the performance degrades gracefully, even down to $10^{-4}$ and less.

**Optimization**. We experimented with multiple per-parameter learning rate updates and consistently found RMSProp (6) to provide the best results, and with the widest tolerance to the learning rate setting. Adagrad (8) also provided a wide tolerance but its performance did not match that of RMSProp. SGD with momentum was able to nearly match the performance of RMSProp but required a precise setting of the learning rate. For example, for a 3-layer LSTM trained with SGD converged only in the range $[10^{-0.4}, 10^{-0.8}]$ and diverged otherwise, while RMSProp diverged outside of $[10^{-2.2}, 10^{-3.2}]$ interval. We clipped the gradients elementwise in range $[-5, 5]$ and found $[-1, 1]$ to be too aggressive, at least in our setting (these settings correspond to case where we average the loss over both the batch size and the sequence length). We did not use L2 regularization because cross-validation always preferred very small coefficients. We used early stopping based on the validation performance. We decay the learning rate after 10 epochs by multiplying it with a factor of 0.95 after each additional epoch. We train each model for 50 epochs.

**Activation Distributions**

During our analysis we noticed a distinct peak of the $\tanh c$ (Figure 7, bottom) at $\pm 0.76$, which is exactly $\tanh(1)$. This is caused by the fact that the largest amount that the LSTM can add or subtract to a cell is 1, since $g$ is the output of $\tanh$. This raises a concern because the hidden state $h$ is computed as $h = o \odot \tanh(c)$. Therefore if a cell happens to be zero because it was reset in the previous iteration, the LSTM does not have the capacity to fully saturate the cell in a single time step. We investigated a simple fix to this problem by changing the update to the form $c_t^l = f \odot c_{t-1}^l + \alpha i \odot g$ where $\alpha$ is a constant. We observed significant significant reduction in performance for $\alpha < 0.25$. Using $\alpha = 2$ allows the LSTM to saturate the cell in one step. In this case we observed a consistently faster convergence (e.g. 700 vs. 900 iterations to reach loss of 2.0) across all architectures. However, the final performance at epoch 50 remains unchanged.
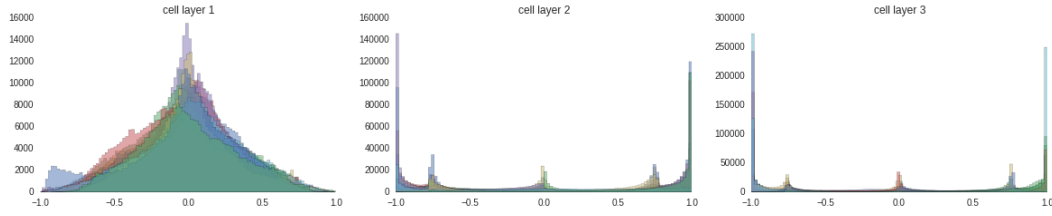
Figure 7: Histograms of $\tanh(c)$ for 20 randomly chosen cells (each in different color) across three layers. Distinct peaks show at $\tanh(1)$, which are caused when cell memory is reset with a forget gate and then maimally added to with the input gate.
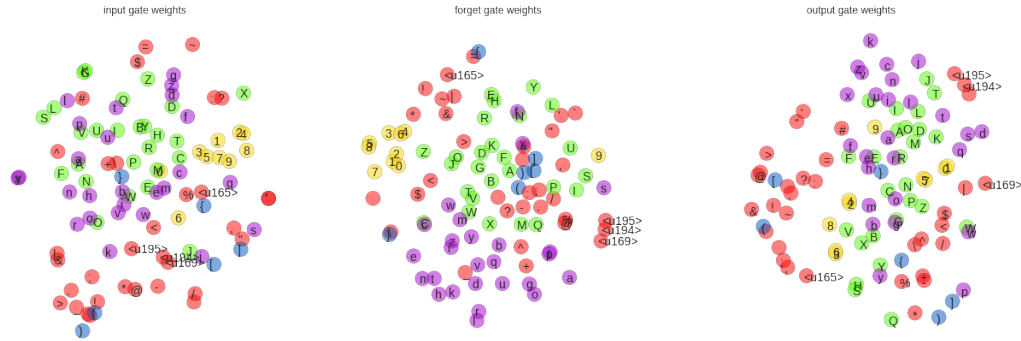


Figure 8: t-SNE of the weights connecting to the individual characters. Each circle corresponds to an input/forget/output gate. For example, numbers cluster together, indicating that their influence on the input/forget gates are very similar. Additionally, the brackets (except for curly brackets) are all clustered in the forget gate indicating similar influence, but they are spread out in the other gates.

**Weights**

We examined the weights of a 1-layer LSTM to gain an understanding of the effect of the inputs on the LSTM's representation. First, we compute the total strength of influence that each character exhibits on the representation, which we quantify as the sum of absolute values of its weights to all gates. For an LSTM trained on the Linux Kernel dataset, the characters with the strongest influence turn out to be special characters (\;} (' {v.w+gb+), while the characters that exhibit weakest influence include rare unicode symbols, followed by `$789^5. This suggests that the LSTMs learns to some extent ignore numerals.

We wanted to gain a more holistic understanding of the influence of each character on the activations of the network. To this effect, we embed the rows of the LSTM's weight matrix connecting the characters to the gates with t-SNE (29) (Figure 8). This figure shows, for instance, that numerals, which are clustered, give rise to very similar write/forget behavior in the LSTM.

**Error Analysis**

| Oracle | LSTM-3 (512) | | LSTM-2 (64) | |
|---|---|---|---|---|
| | # Errors | Fraction of Errors | # Errors | Fraction of Errors |
| 1-gram | 1581 | 0.011 | 1697 | 0.009 |
| 2-gram | 4001 | 0.029 | 4637 | 0.025 |
| 3-gram | 2959 | 0.021 | 8200 | 0.045 |
| 4-gram | 3425 | 0.025 | 14053 | 0.076 |
| 5-gram | 3183 | 0.023 | 11602 | 0.063 |
| 6-gram | 2974 | 0.021 | 8257 | 0.045 |
| 7-gram | 2825 | 0.020 | 5859 | 0.032 |
| 8-gram | 2406 | 0.017 | 4268 | 0.023 |
| 9-gram | 1635 | 0.012 | 2634 | 0.014 |
| 100 memory | 989 | 0.007 | 1074 | 0.006 |
| 500 memory | 2714 | 0.019 | 2985 | 0.016 |
| 1000 memory | 1407 | 0.010 | 1566 | 0.009 |
| 5000 memory | 3409 | 0.024 | 3828 | 0.021 |
| Word 0-train | 3937 | 0.028 | 4026 | 0.022 |
| Word 1-train | 2220 | 0.016 | 2302 | 0.013 |
| Word 2-train | 1687 | 0.012 | 1775 | 0.010 |
| Word 3-train | 1557 | 0.011 | 1636 | 0.009 |
| Word 4-train | 1330 | 0.001 | 1377 | 0.008 |
| Word 5-train | 1197 | 0.009 | 1244 | 0.007 |
| after space or quote | 45114 | 0.323 | 45834 | 0.250 |
| After newline | 6843 | 0.049 | 7277 | 0.040 |
| Punctuation | 8544 | 0.061 | 8728 | 0.048 |
| 0.1 boost | 3060 | 0.022 | 2445 | 0.013 |
| 0.2 boost | 4384 | 0.031 | 4288 | 0.023 |
| 0.3 boost | 6113 | 0.044 | 7069 | 0.038 |
| 0.4 boost | 8337 | 0.060 | 10091 | 0.055 |
| 0.5 boost | 11986 | 0.09 | 14993 | 0.082 |

Table 4: Full breakdown of errors corrected by all oracles for Figure 6 from main paper.

| Oracle Group | LSTM-3 (512) | | LSTM-2 (64) | |
|---|---|---|---|---|
| | # Errors | Fraction of Errors | # Errors | Fraction of Errors |
| $n$-gram | 24989 | 0.179 | 61207 | 0.333 |
| memory | 8519 | 0.061 | 9453 | 0.051 |
| rare words | 11928 | 0.085 | 12360 | 0.067 |
| spacing | 60501 | 0.433 | 61839 | 0.337 |
| boost | 33880 | 0.242 | 38886 | 0.212 |

Table 5: Breakdown of errors correct by oracles, grouped by oracle type.
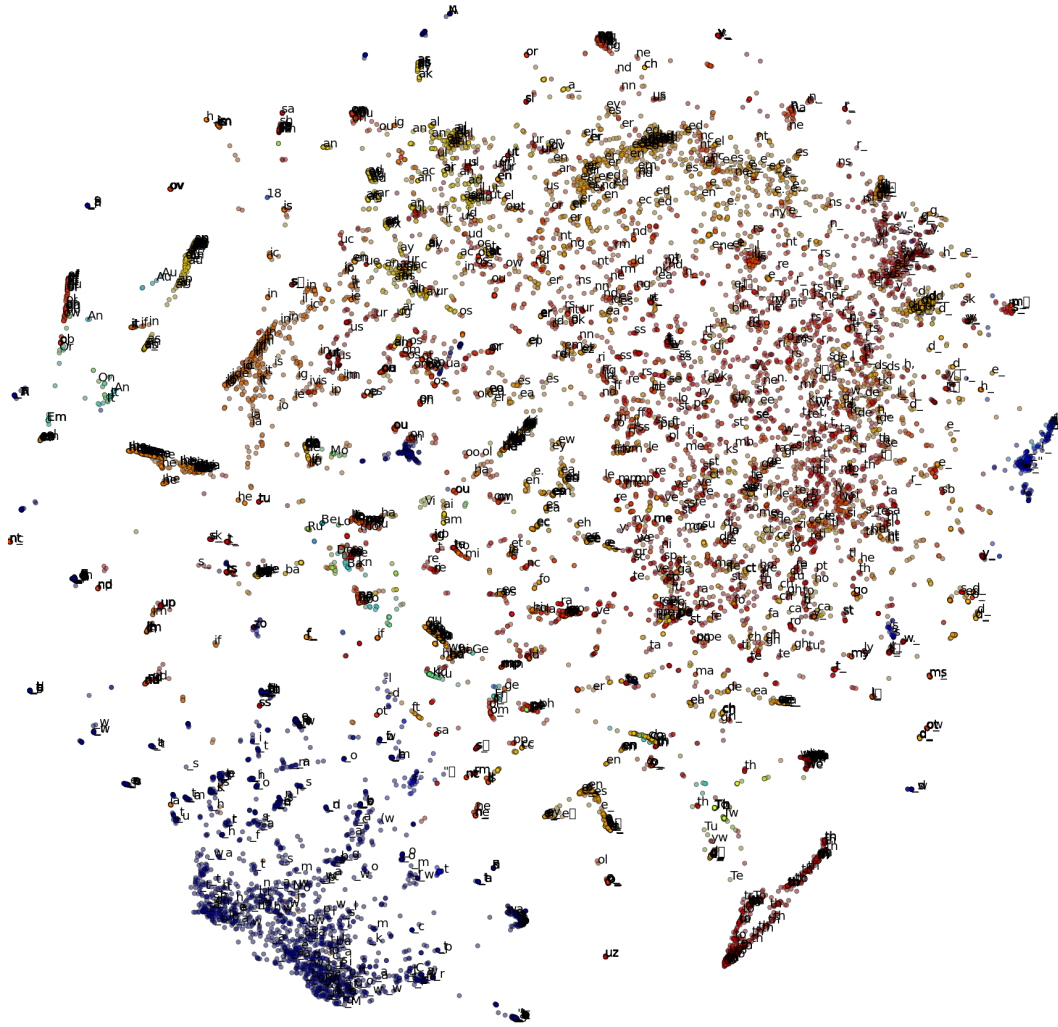
Figure 9: A t-SNE visualization of LSTM cell state representations on the War and Peace dataset. Each point is colored based on the character most recently read by the LSTM, and is captioned with both the previous and next character. Spaces are shown as underscores.

**Cell state representation**

In Figure 9 we use t-SNE to visualize the (gated) cell state $\tanh(c)$ of an LSTM, allowing us to gain insight into the type of information it encodes. In the lower left we see a cluster of cell states that occur immediately after a space, and in the upper right we see a cluster of cell states where a space is about to occur. This demonstrates that the LSTM uses its cell state both to record characters that have been recently seen and to hypothesize about characters that may appear next.