

# Gesture Recognition with Recurrent Neural Networks

**Justin Fu**

University of California, Berkeley  
justinfu@berkeley.edu

**Siddhartho Bhattacharya**

University of California, Berkeley  
siddhartho\_b@hotmail.com

## Abstract

Gesture recognition is typically done through images, but it requires extensive effort to process image data into a form useful for gesture recognition. Our project uses accelerometer data from smartphones, so no special hardware is required. Recurrent neural networks provide a natural way to model time sequences, and the LSTM architecture in particular provides a way to model long-range dependencies that are prevalent in the gestures we tested. TODO: We implemented S gestures, and achieve an accuracy of X percent on the test set we generated.

## 1 Introduction

TODO: Intro goes here. Talk about why this problem is hard/interesting, and possibly some recent work.

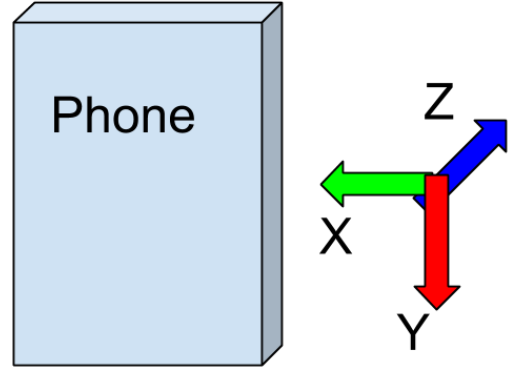
Potential applications of this project include making phones more intuitive to use. For example, imagine phones that automatically unlock and lock when you take them in and out of your pocket, picking up calls by just bringing it up to your ear, etc.

## 2 Data

We collected accelerometer data streamed from a smartphone. Specifically, we only considered values in the X and Y directions as seen from the perspective of the phone.

We implemented 6 different gestures representing some letters of the alphabet (M, Z, O, L, J, and S). Our training set contains 26 instances of each letter and our test set contains 10, for a total of 216 examples.

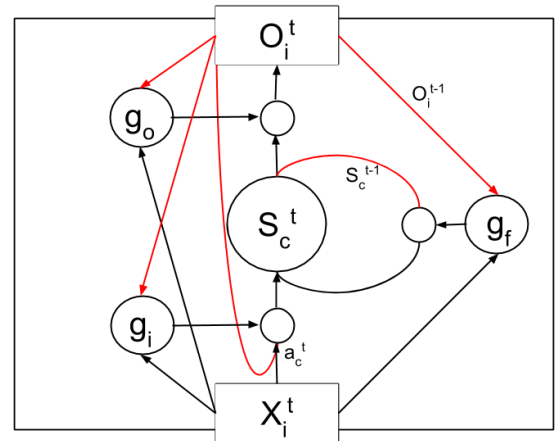
Figure 1: Directions of each axis on the phone



## 3 Model

We implemented an LSTM (long short term memory) network based on the architecture proposed by Graves [1]. Specifically, our network contains the forget, input, and output gates, but lacks the peephole connections that connect the internal cell state to these gates (as done in Vinyals et al. [2]).

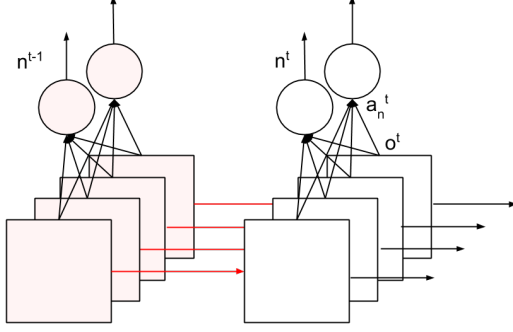
Figure 2: A single LSTM cell. Red arrows are connections forward in time.



We then interleave layers of LSTM units with

normal network layers to create a multilayer network. The output of the final layer is then fed through a softmax function.

Figure 3: Diagram of one interleaved layer of our network. LSTM cells are represented as squares and normal neurons are represented as circles.



### 3.1 Equations

We used cross-entropy as our objective function.  $O$  is the network output after the softmax and  $Y$  is the training label.  $n$  sums over each training example,  $t$  sums over time, and  $i$  sums over the dimensions of each training example at a given timestep. Note that  $t$  varies with each training example.

$$L(Y, O) = \sum_n \sum_t \sum_i y_i^{n,t} * \ln(o_i^{n,t})$$

Most of the equations on the forward and backward passes are the same as those in Graves [1], but some were missing from his paper and we dropped the peepholes so we specify all for completeness.

### 3.2 Forward Dynamics

The forward dynamics of the network are specified as follows. Equations are specified for a whole layer (rather than per unit). Bolded values represent vectors. Note that  $*$  represents a pointwise multiply, and the nonlinearities are applied element-wise to the input vector.

Gates

$$\mathbf{a}_i^t = W_{i,x} \mathbf{x}^t + W_{i,o} \mathbf{o}^{t-1}$$

$$\mathbf{g}_i^t = \sigma(\mathbf{a}_i^t)$$

$$\mathbf{a}_f^t = W_{f,x} \mathbf{x}^t + W_{f,o} \mathbf{o}^{t-1}$$

$$\mathbf{g}_f^t = \sigma(\mathbf{a}_f^t)$$

$$\mathbf{a}_o^t = W_{o,x} \mathbf{x}^t + W_{o,o} \mathbf{o}^{t-1}$$

$$\mathbf{g}_o^t = \sigma(\mathbf{a}_o^t)$$

Cell State

$$\mathbf{a}_c^t = W_{c,x} \mathbf{x}^t + W_{c,o} \mathbf{o}^{t-1}$$

$$\mathbf{s}_c^t = \mathbf{g}_i^t * \sigma(\mathbf{a}_c^t) + \mathbf{g}_f^t * \mathbf{s}_c^{t-1}$$

Cell Output (Hidden State)

$$\mathbf{o}^t = \mathbf{g}_i^t * \tanh(\mathbf{s}_c^t)$$

Normal Network Layer Output

$$\mathbf{a}_n^t = W_n \mathbf{o}^t$$

$$\mathbf{n}^t = \sigma(\mathbf{a}_n^t)$$

### 3.3 Backpropagation

$\delta_k^{t,l+1}$  represents the error backpropagating from the above layer in the current time step. As with before,  $*$  represents a pointwise multiplication. The  $\delta$  variables represent vectors for the whole layer.

Cell Output

$$\delta_n^t = \delta_k^{t,l+1} * \tanh'(\mathbf{a}_n^t)$$

$$\delta_h^t = W_{i,o}^T \delta_i^{t+1} + W_{f,o}^T \delta_f^{t+1} + W_{o,o}^T \delta_o^{t+1} + W_{c,o}^T \delta_s^{t+1}$$

$$\frac{\partial L}{\partial \mathbf{o}^t} = W_n^T \delta_n^t + \delta_h^t$$

Cell State

$$\delta_c^t = \mathbf{g}_o^t * \tanh'(\mathbf{s}_c^t) * \delta_h^t * \delta_c^{t+1} * \mathbf{g}_f^{t+1}$$

$$\delta_s^t = \mathbf{g}_i^t * \sigma'(\mathbf{a}_c^t) * \delta_c^t$$

Gates

$$\delta_o^t = \sigma'(\mathbf{a}_o^t) * \tanh(\mathbf{s}_c^t) * \frac{\partial L}{\partial \mathbf{o}^t}$$

$$\delta_f^t = \sigma'(\mathbf{a}_f^t) * \mathbf{s}_c^{t-1} * \delta_c^t$$

$$\delta_i^t = \sigma'(\mathbf{a}_i^t) * \sigma'(\mathbf{a}_c^t) * \delta_c^t$$

Input

$$\delta_k^{t,l} = W_{i,x}^T \delta_i^t + W_{f,x}^T \delta_f^t + W_{o,x}^T \delta_o^t + W_{c,x}^T \delta_s^t$$

### 3.4 Gradient

The gradient of a matrix  $W$  is calculated by taking the outer product of its input and the delta of its output.

For example, here is the derivative of the matrix of the normal layer output  $W_n$ . The input to this matrix are the cell outputs  $\mathbf{o}^t$  and its output,  $\mathbf{a}_n$  has the delta  $\delta_n^t$

$$\frac{\partial L}{\partial W_n} = \text{outer}(\mathbf{o}^t, \delta_n^t)$$

These gradients are summed over all time steps and all training examples per weight update.

### 3.5 Training

We trained our network using BPTT (backpropagation through time) and gradient descent with momentum, using the equations described in the previous section. All weights are randomly initialized in the range  $[-0.1, 0.1]$ .

Training inputs (acceleration values) are directly presented to the network on the lowest layer. We present training labels to the network as a one-hot vector representing the gesture to the network at all time steps.

### 3.6 Decoding

To retrieve a single gesture prediction from the network, we run a forward pass and collect the network outputs. Then, we average the network outputs after a softmax across time and select the class  $c$  with the maximum average response.

$$c = \arg \max_i \left( \sum_t o_i^t \right)$$

A problem with this approach is that our network does not explicitly model the summation over outputs. An approach to try is to backpropagate the error from this decision rule, but the network worked fine without it.

## 4 Results

As we collected data ourselves, there were no available datasets online to compare our architecture against.

TODO: Talk about training and test set sizes.

TODO: Talk about specific errors - ex. are there any gestures that are commonly confused?

## 5 Conclusion

TODO: How successful was the project? TODO: Future directions.

Gesture	Number Correct	Error %
Z	blah	10.23%
O	blah	5.32%
M	blah	5.32%
L	blah	5.32%
J	blah	5.32%
S	blah	5.32%

Table 1: Results

## 6 Example L<sup>A</sup>T<sub>E</sub>Xstuff

- Item Example 1
- Item Example 2

Column	Col	Col
row blah	blah	blah
row blah	blah	blah

Table 2: Example table!

Noindent

### 6.1 subsection

```
\usepackage{times}
\usepackage{latexsym}
```

“(Graves, 2008) Quote”

### 6.2 Sections

**Citations:** Citations within the text appear in parentheses as (?) or, if the author’s name appears in the text itself, as Gusfield (?). Append lowercase letters to the year in cases of ambiguity. Treat double authors as in (?), but write as in (?) when more than two authors are involved. Collapse multiple citations as in (??). Also refrain from using full citations as sentence constituents. We suggest that instead of

you use

“Gusfield (?) showed that ...”

If you are using the provided L<sup>A</sup>T<sub>E</sub>X and BibT<sub>E</sub>X style files, you can use the command `\newcite` to get “author (year)” citations.

As reviewing will be double-blind, the submitted version of the papers should not include the authors’ names and affiliations. Furthermore, self-references that reveal the author’s identity, e.g.,

“We previously showed (?) ...”

should be avoided. Instead, use citations such as

“Gusfield (?) previously showed ... ”

### 6.3 Footnotes

**Footnotes:** Put footnotes at the bottom of the page and use 9 points text. They may be numbered or referred to by asterisks or other symbols.<sup>1</sup> Footnotes should be separated from the text by a line.<sup>2</sup>

### Additional Materials

Our code and datasets are hosted on <https://github.com/justinjfu/lstm>

### Acknowledgments

Brian is our savior.

### References

- [1] Alex Graves. 2008. *Supervised Sequence Labelling with Recurrent Neural Networks*. PhD Thesis.
- [2] Oriol Vinyals, Alexander Toshev, Samy Bengio, Dumitru Erhan. 2014. *Show and Tell: A Neural Image Caption Generator*.

---

<sup>1</sup>This is how a footnote should appear.

<sup>2</sup>Note the line separating the footnotes from the text.