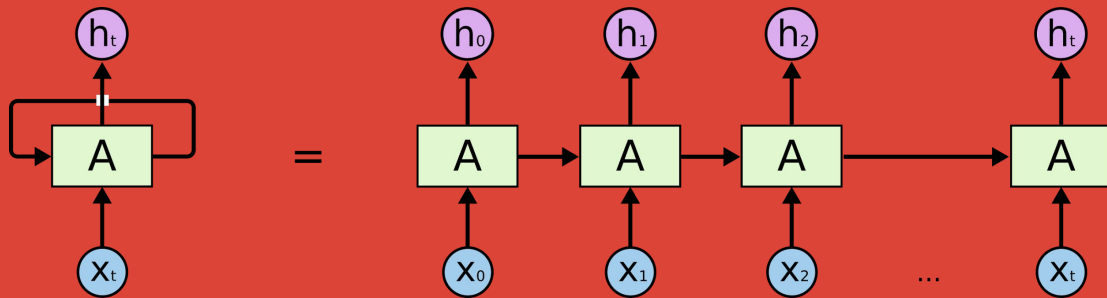# LSTMs

Learning to forget

Devansh Kukreja and Eric Nie

# Background/Motivation

- State-of-the-art in sequence learning
- Address the issue of remembering information
- Useful for analyzing sequential data
- Cells within RNNs and LSTMs maintain an internal state

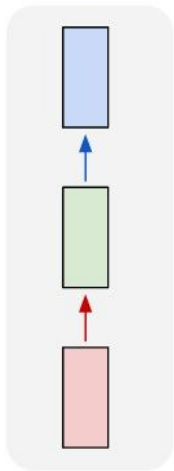# Recurrent Neural Networks (A Brief Introduction)
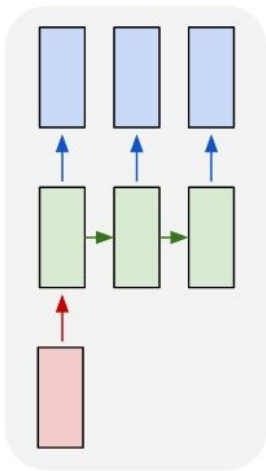
# Intro



Recurrent Neural Networks

- Internal states
- Sequences of inputs and outputs
- Share weights across timesteps

# More RNNs
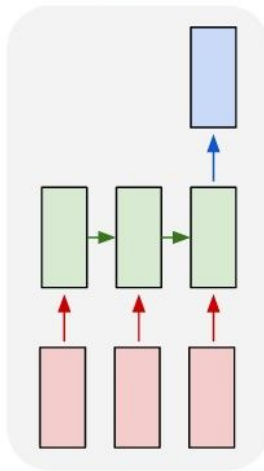


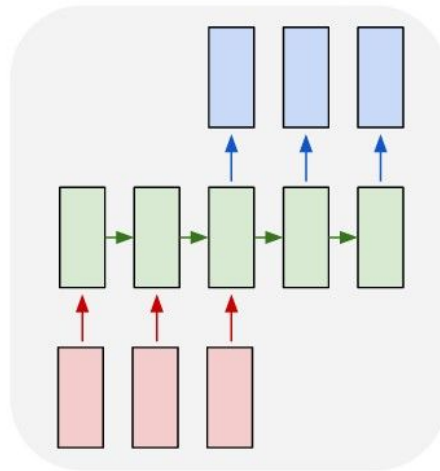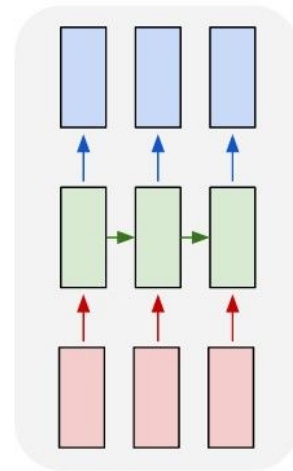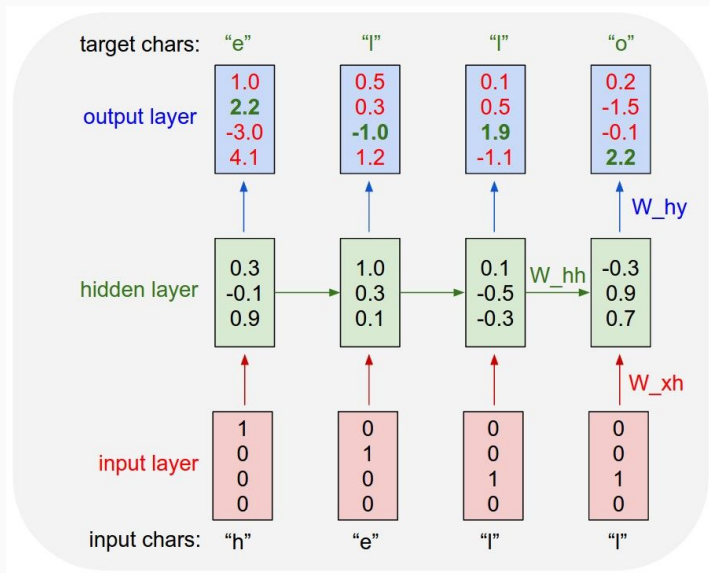one to one · one to many · many to one · many to many · many to many

# Character-level RNN



https://karpathy.github.io/2015/05/21/rnn-effectiveness/

- Predict the next character given input
- One-hot encoding

# What can an RNN do?

## Generate Text: Linux source code

```c
/*
 * Increment the size file of the new incorrect UI_FILTER group information
 * of the size generatively.
 */
static int indicate_policy(void)
{
    int error;
    if (fd == MARN_EPT) {
        /*
         * The kernel blank will coeld it to userspace.
         */
        if (ss->segment < mem_total)
            unblock_graph_and_set_blocked();
        else
            ret = 1;
        goto bail;
    }
    segaddr = in_SB(in.addr);
    selector = seg / 16;
    setup_works = true;
    for (i = 0; i < blocks; i++) {
        seq = buf[i++];
        bpf = bd->bd.next + i * search;
        if (fd) {
            current = blocked;
        }
    }
    rw->name = "Getjbbregs";
    bprm_self_clearl(&iv->version);
    regs->new = blocks[(BPF_STATS << info->historidac)] | PFMR_CLOBATHINC_SECONDS << 12;
    return segtable;
}
```

## Produce Research Papers:

# Training a RNN

Backpropagation through time:



- Minibatches of subsequences
- Update weights after some number of inputs

$$\begin{aligned}
\boldsymbol{a}^{(t)} &= \boldsymbol{b} + \boldsymbol{W}\boldsymbol{h}^{(t-1)} + \boldsymbol{U}\boldsymbol{x}^{(t)} & (10.8)\\
\boldsymbol{h}^{(t)} &= \tanh(\boldsymbol{a}^{(t)}) & (10.9)\\
\boldsymbol{o}^{(t)} &= \boldsymbol{c} + \boldsymbol{V}\boldsymbol{h}^{(t)} & (10.10)\\
\hat{\boldsymbol{y}}^{(t)} &= \text{softmax}(\boldsymbol{o}^{(t)}) & (10.11)
\end{aligned}$$

Equations for forward pass

# Adding Binary Numbers

# Issues with RNNs

- Training RNNs uses backpropagation through time
- Exploding and vanishing gradient problem similar to that of deep neural networks

# Extension to RNNs

Bidirectional RNNs



Deep Bidirectional RNNs

# Final Word on RNNs

- Turing complete
- Useful for learning sequential data
- Trouble capturing long-term dependencies in practice

**LSTMs** are a variant RNN that uses "**gating units**" to control information flow adaptively to combat "**gradient vanishing**" *

# Long Short Term Memory Networks

# Long Short Term Memory Networks

# Long Short Term Memory Networks

1. **Forget** Gate
2. **Input** Gate
3. **Output** Gate

# **Forget** Gate

1. Takes in $h_{t-1}$ concat $x_t$ and outputs a number between 0 and 1 for each number in the cell state.

2. That output filters the cell state, selectively choosing how much of each node to forget.

$$f_t = \sigma\left(W_f \cdot [h_{t-1}, x_t] + b_f\right)$$

# **Input** Gate

1. Sigmoid Input gate decides what new information can flow in.

2. Tanh layer creates a vector of new candidate values that could be added.

$$i_t = \sigma\left(W_i \cdot [h_{t-1}, x_t] \; + \; b_i\right)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] \; + \; b_C)$$

# Putting it together

1. Apply forget gate
2. Compute new candidate cell state
3. Filter in new cell state



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

# **Output** Gate

1. Run new cell state through tanh to put values between -1 and 1.

2. Multiply it by the output of the third sigmoid gate to decide what to output.

$$o_t = \sigma \left( W_o \left[ h_{t-1}, x_t \right] + b_o \right)$$

$$h_t = o_t * \tanh \left( C_t \right)$$

# **LSTM** Demo

Row 1: Characters as they're recognized

Row 2: States of some memory cells

Row 3: Writing as it's analyzed

Row 4: Gradient backpropagated to inputs from the most active character

https://arxiv.org/pdf/1308.0850.pdf

# Variations on LSTM

# Peephole Connections

1. Pass the gate layers the current cell state.

2. Helps the system learn from the size of the time lags.



$$f_t = \sigma\left(W_f \cdot [\boldsymbol{C_{t-1}}, h_{t-1}, x_t] + b_f\right)$$

$$i_t = \sigma\left(W_i \cdot [\boldsymbol{C_{t-1}}, h_{t-1}, x_t] + b_i\right)$$

$$o_t = \sigma\left(W_o \cdot [\boldsymbol{C_t}, h_{t-1}, x_t] + b_o\right)$$

ftp://ftp.idsia.ch/pub/juergen/TimeCount-IJCNN2000.pdf

http://colah.github.io/posts/2015-08-Understanding-LSTMs/

# Couple Forget/Input Gates

1. Pass the gate layers the current cell state.

2. Simplifies the cell but reduced performance slightly.



$$C_t = f_t * C_{t-1} + (1 - f_t) * \tilde{C}_t$$

# Gated Recurrent Unit (GRU)

1.  Combine forget and input gates into an **"update gate"**

2.  In general, equal to LSTMs in performance and faster in training



$$z_t = \sigma \left( W_z \cdot [h_{t-1}, x_t] \right)$$

$$r_t = \sigma \left( W_r \cdot [h_{t-1}, x_t] \right)$$

$$\tilde{h}_t = \tanh \left( W \cdot [r_t * h_{t-1}, x_t] \right)$$
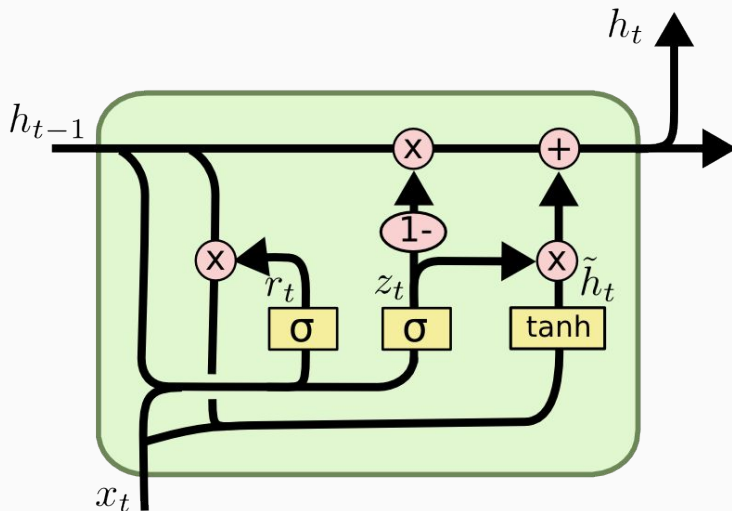
$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

# Depth Gated RNNs

1. Adds a fourth gate - the Depth Gate
2. Creates a gated linear connection between lower and upper layer memory cells
3. Relates to highway networks and Grid LSTM

Using the depth gate, a DGLSTM unit can be written as

$$i_t^{L+1} = \sigma(W_{xi}^{L+1}x_t + W_{hi}^{L+1}h_{t-1}^{L+1} + W_{ci}^{L+1}c_{t-1}^{L+1}) \tag{13}$$

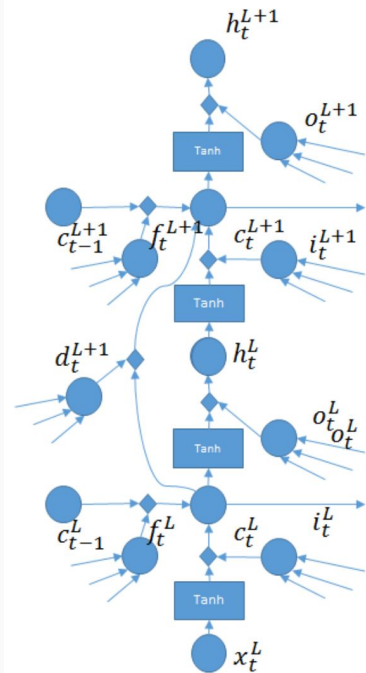$$f_t^{L+1} = \sigma(W_{xf}^{L+1}x_t + W_{hf}^{L+1}h_{t-1}^{L+1} + W_{cf}c_{t-1})^{L+1} \tag{14}$$

$$d_t^{L+1} = \sigma(b_d^{L+1} + W_{xd}^{L+1}x_t^{L+1} + W_{cd}^{L+1} \odot c_{t-1}^{L+1} + W_{ld}^{L+1} \odot c_t^L) \tag{15}$$

$$c_t^{L+1} = d_t^{L+1}c_t^L + f_t^{L+1} \odot c_{t-1} + i_t^{L+1} \odot tanh(W_{xc}x_t + W_{hc}h_{t-1}^{L+1}) \tag{16}$$

$$o_t^{L+1} = \sigma(W_{xo}^{L+1}x_t + W_{ho}^{L+1}h_{t-1} + W_{co}^{L+1}c_t^{L+1}) \tag{17}$$

$$h_t^{L+1} = o_t^{L+1} \odot tanh(c_t^{L+1}) \tag{18}$$

where $i_t^{L+1}$, $f_t^{L+1}$ $o_t^{L+1}$, and $d_t^{L+1}$ are the input gate, forget gate, output gate and the depth gate.
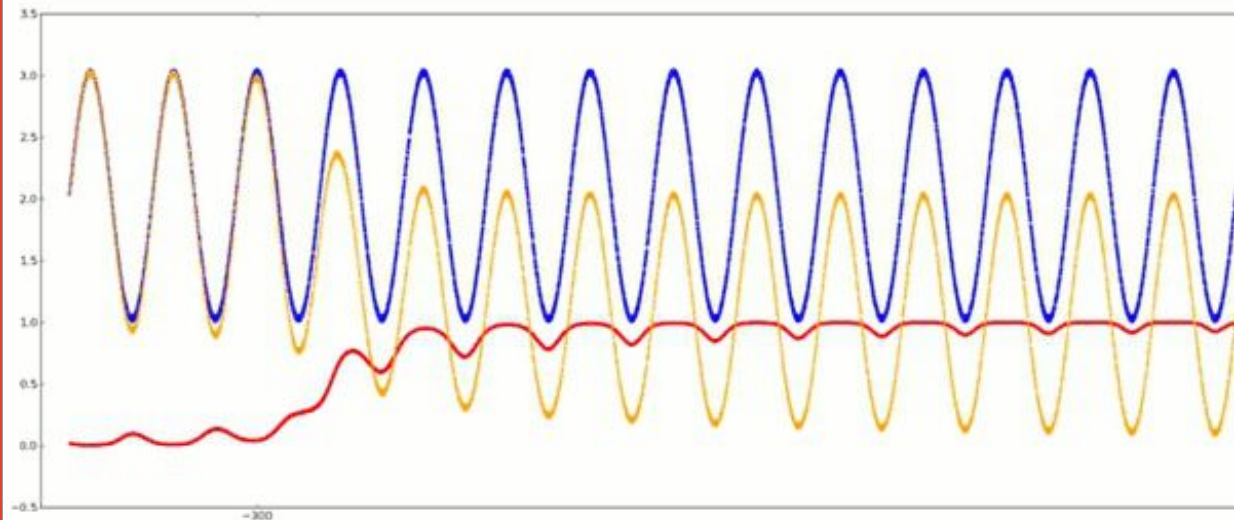
# Questions

Eric Nie
Devansh Kukreja

LSTMs



LSTM learning to predict next noised sinus value.

**Training Signal**
**LSTM Prediction**
**Error**