# LSTM: A Search Space Odyssey

Authors: Klaus Greff, Rupesh K. Srivastava, Jan Koutník, Bas R. Steunebrink, Jürgen Schmidhuber
Presenter: Sidhartha Satapathy

## Scientific contributions of the paper:

- The paper aims at evaluating different elements of the most popular LSTM architecture.
- The paper shows the performance of various variants of the vanilla LSTM by making a single change which allows us to isolate the effect of each of these changes on the performance of the architecture.
- The paper also provide insights gained about hyperparameters and their interaction.

# Dataset 1: IAM Online Handwriting Database

- IAM Online Handwriting Database: The IAM Handwriting Database contains forms of handwritten English text which can be used to train and test handwritten text recognizers and to perform writer identification and verification experiments.

In mid-april Anglesey moved his family and entourage from Rome to Naples, there to await the arrival of

Each sequence or line in this case is made up of frames and the task at hand is to classify each of these frames into one of the 82 characters.

Here are the output characters:
abcdefghijklmnopqrstuvwxyz
ABCDEFGHIJKLMNOPQRSTUVWXYZ
0123456789 !"#&\'()*+,-./[]:;? And the empty symbol.

The performance in this case is the character error rate.

# Dataset 2: TIMIT

- TIMIT Speech corpus: **TIMIT** is a corpus of phonemically and lexically transcribed speech of American English speakers of different sexes and dialects.
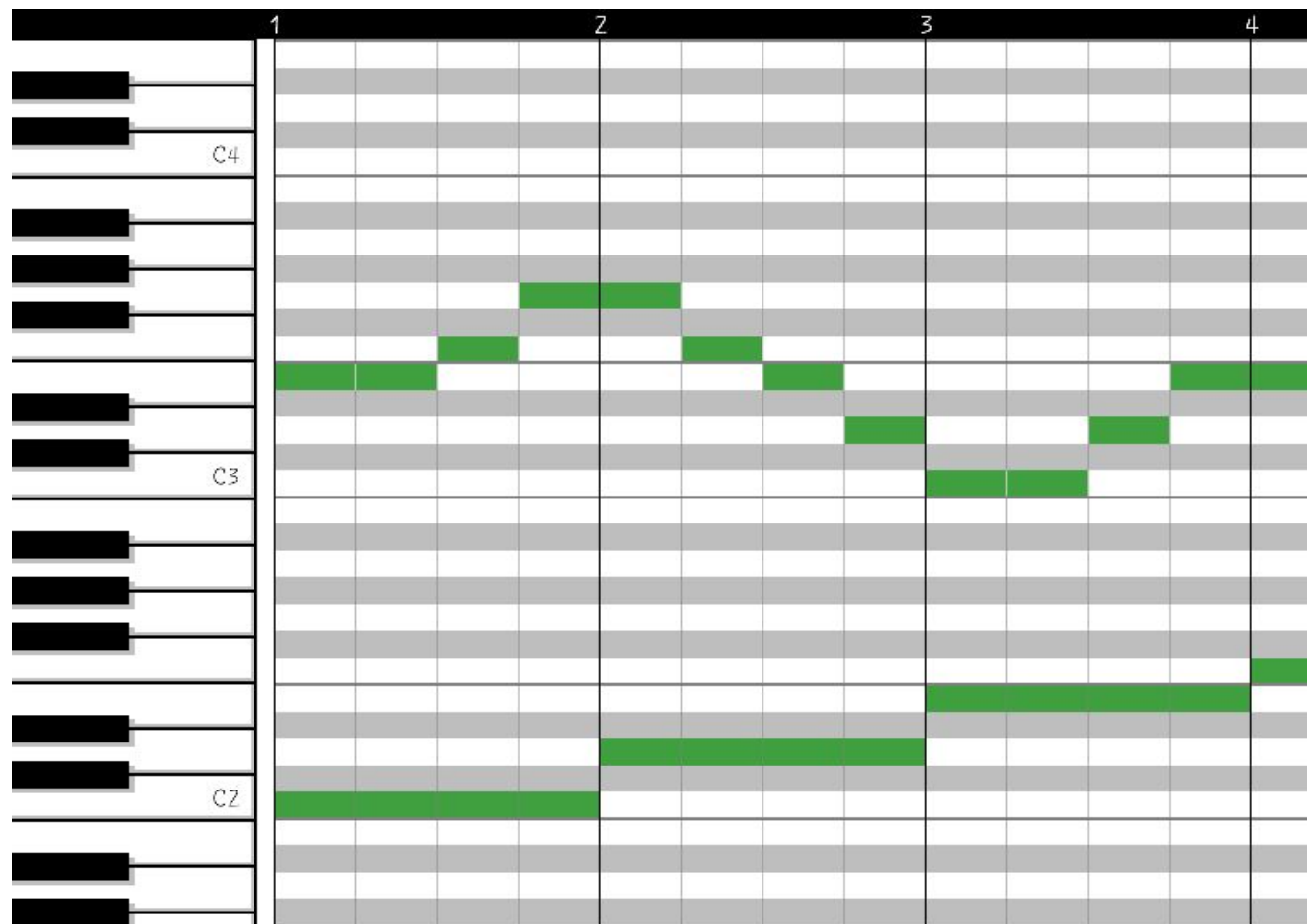
- Our experiments focus on the frame-wise classification task for this dataset, where the objective is to classify each audio-frame as one of 61 phones.
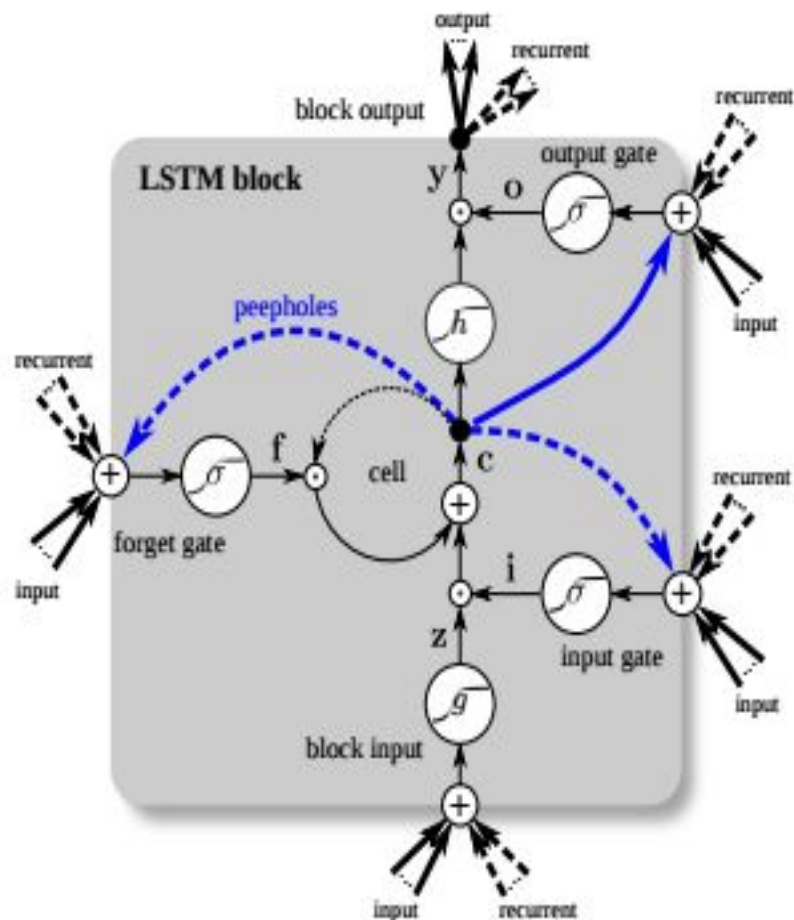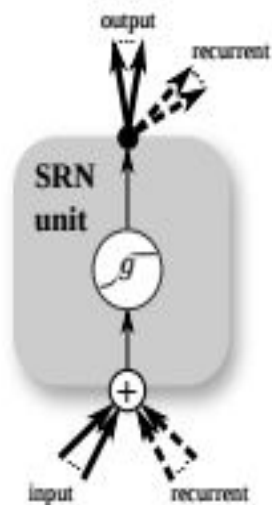- The performance in this case is the classification error rate.

| | Phone Label | Example | | Phone Label | Example | | Phone Label | Example |
|---|---|---|---|---|---|---|---|---|
| 1 | iy | b*ee*t | 22 | ch | *ch*oke | 43 | en | butto*n* |
| 2 | ih | b*i*t | 23 | b | *b*ee | 44 | eng | Washi*ng*ton |
| 3 | eh | b*e*t | 24 | d | *d*ay | 45 | l | *l*ay |
| 4 | ey | b*ai*t | 25 | g | *g*ay | 46 | r | *r*ay |
| 5 | ae | b*a*t | 26 | p | *p*ea | 47 | w | *w*ay |
| 6 | aa | b*o*b | 27 | t | *t*ea | 48 | y | *y*acht |
| 7 | aw | b*ou*t | 28 | k | *k*ey | 49 | hh | *h*ay |
| 8 | ay | b*i*te | 29 | dx | mu*dd*y | 50 | hv | a*h*ead |
| 9 | ah | b*u*t | 30 | s | *s*ea | 51 | el | bott*le* |
| 10 | ao | b*ou*ght | 31 | sh | *sh*e | 52 | bcl | b closure |
| 11 | oy | b*oy* | 32 | z | *z*one | 53 | dcl | d closure |
| 12 | ow | b*oa*t | 33 | zh | a*z*ure | 54 | gcl | g closure |
| 13 | uh | b*oo*k | 34 | f | *f*in | 55 | pcl | p closure |
| 14 | uw | b*oo*t | 35 | th | *th*in | 56 | tcl | t closure |
| 15 | ux | t*oo*t | 36 | v | *v*an | 57 | kcl | k closure |
| 16 | er | b*ir*d | 37 | dh | *th*en | 58 | q | glotal stop |
| 17 | ax | *a*bout | 38 | m | *m*om | 59 | pau | pause |
| 18 | ix | deb*i*t | 39 | n | *n*oon | 60 | epi | epenthetic silence |
| 19 | axr | butt*er* | 40 | ng | si*ng* | | | |
| 20 | ax-h | s*u*spect | 41 | em | bott*om* | 61 | h# | begin/end marker |
| 21 | jh | *j*oke | 42 | nx | wi*nn*er | | | |

Table 2. 61 TIMIT original phone set.

# Dataset 3: JSB Chorales

- JSB Chorales:  JSB Chorales is a collection of 382 four part harmonized chorales by J. S. Bach, the networks where trained to do next-step prediction.

**Legend**

- unweighted connection
- weighted connection
- connection with time-lag
- • branching point
- ⊙ mutliplication
- ⊕ sum over all inputs
- σ gate activation function (always sigmoid)
- g input activation function (usually tanh)
- h output activation function (usually tanh)

SRN unit

LSTM block

block output

output gate

peepholes

forget gate

cell

input gate

block input

$$\bar{\mathbf{z}}^t = \mathbf{W}_z \mathbf{x}^t + \mathbf{R}_z \mathbf{y}^{t-1} + \mathbf{b}_z$$

$$\mathbf{z}^t = g(\bar{\mathbf{z}}^t) \qquad\qquad\qquad block\ input$$

$$\bar{\mathbf{i}}^t = \mathbf{W}_i \mathbf{x}^t + \mathbf{R}_i \mathbf{y}^{t-1} + \mathbf{p}_i \odot \mathbf{c}^{t-1} + \mathbf{b}_i$$

$$\mathbf{i}^t = \sigma(\bar{\mathbf{i}}^t) \qquad\qquad\qquad input\ gate$$

$$\bar{\mathbf{f}}^t = \mathbf{W}_f \mathbf{x}^t + \mathbf{R}_f \mathbf{y}^{t-1} + \mathbf{p}_f \odot \mathbf{c}^{t-1} + \mathbf{b}_f$$

$$\mathbf{f}^t = \sigma(\bar{\mathbf{f}}^t) \qquad\qquad\qquad forget\ gate$$

$$\mathbf{c}^t = \mathbf{z}^t \odot \mathbf{i}^t + \mathbf{c}^{t-1} \odot \mathbf{f}^t \qquad\qquad\qquad cell$$

$$\bar{\mathbf{o}}^t = \mathbf{W}_o \mathbf{x}^t + \mathbf{R}_o \mathbf{y}^{t-1} + \mathbf{p}_o \odot \mathbf{c}^t + \mathbf{b}_o$$

$$\mathbf{o}^t = \sigma(\bar{\mathbf{o}}^t) \qquad\qquad\qquad output\ gate$$

$$\mathbf{y}^t = h(\mathbf{c}^t) \odot \mathbf{o}^t \qquad\qquad\qquad block\ output$$

# Variants of the LSTM Block:

- NIG: No Input Gate
- NFG: No Forget Gate
- NOG: No Output Gate
- NIAF: No Input Activation Function
- NOAF: No Output Activation Function
- CIFG: Coupled Input and Forget Gate
- NP: No Peepholes
- FGR: Full Gate Recurrence

# NIG: No Input Gate

**NIG:** No Input Gate: $\mathbf{i}^t = 1$

$$\bar{\mathbf{z}}^t = \mathbf{W}_z \mathbf{x}^t + \mathbf{R}_z \mathbf{y}^{t-1} + \mathbf{b}_z$$

$$\mathbf{z}^t = g(\bar{\mathbf{z}}^t) \qquad \text{\textit{block input}}$$

$$\bar{\mathbf{i}}^t = \mathbf{W}_i \mathbf{x}^t + \mathbf{R}_i \mathbf{y}^{t-1} + \mathbf{p}_i \odot \mathbf{c}^{t-1} + \mathbf{b}_i$$

$$\mathbf{i}^t = \sigma(\bar{\mathbf{i}}^t) \qquad \text{\textit{input gate}}$$

$$\bar{\mathbf{f}}^t = \mathbf{W}_f \mathbf{x}^t + \mathbf{R}_f \mathbf{y}^{t-1} + \mathbf{p}_f \odot \mathbf{c}^{t-1} + \mathbf{b}_f$$

$$\mathbf{f}^t = \sigma(\bar{\mathbf{f}}^t) \qquad \text{\textit{forget gate}}$$

$$\mathbf{c}^t = \mathbf{z}^t \odot \mathbf{i}^t + \mathbf{c}^{t-1} \odot \mathbf{f}^t \qquad \text{\textit{cell}}$$

$$\bar{\mathbf{o}}^t = \mathbf{W}_o \mathbf{x}^t + \mathbf{R}_o \mathbf{y}^{t-1} + \mathbf{p}_o \odot \mathbf{c}^t + \mathbf{b}_o$$

$$\mathbf{o}^t = \sigma(\bar{\mathbf{o}}^t) \qquad \text{\textit{output gate}}$$

$$\mathbf{y}^t = h(\mathbf{c}^t) \odot \mathbf{o}^t \qquad \text{\textit{block output}}$$

# NFG: No Forget Gate

**NFG: No Forget Gate: $\mathbf{f}^t = 1$**

$$\bar{\mathbf{z}}^t = \mathbf{W}_z \mathbf{x}^t + \mathbf{R}_z \mathbf{y}^{t-1} + \mathbf{b}_z$$
$$\mathbf{z}^t = g(\bar{\mathbf{z}}^t) \qquad \text{\textit{block input}}$$
$$\bar{\mathbf{i}}^t = \mathbf{W}_i \mathbf{x}^t + \mathbf{R}_i \mathbf{y}^{t-1} + \mathbf{p}_i \odot \mathbf{c}^{t-1} + \mathbf{b}_i$$
$$\mathbf{i}^t = \sigma(\bar{\mathbf{i}}^t) \qquad \text{\textit{input gate}}$$
$$\bar{\mathbf{f}}^t = \mathbf{W}_f \mathbf{x}^t + \mathbf{R}_f \mathbf{y}^{t-1} + \mathbf{p}_f \odot \mathbf{c}^{t-1} + \mathbf{b}_f$$
$$\mathbf{f}^t = \sigma(\bar{\mathbf{f}}^t) \qquad \text{\textit{forget gate}}$$
$$\mathbf{c}^t = \mathbf{z}^t \odot \mathbf{i}^t + \mathbf{c}^{t-1} \odot \mathbf{f}^t \qquad \text{\textit{cell}}$$
$$\bar{\mathbf{o}}^t = \mathbf{W}_o \mathbf{x}^t + \mathbf{R}_o \mathbf{y}^{t-1} + \mathbf{p}_o \odot \mathbf{c}^t + \mathbf{b}_o$$
$$\mathbf{o}^t = \sigma(\bar{\mathbf{o}}^t) \qquad \text{\textit{output gate}}$$
$$\mathbf{y}^t = h(\mathbf{c}^t) \odot \mathbf{o}^t \qquad \text{\textit{block output}}$$

# NOG: No Output Gate

**NOG:** No Output Gate: $o^t = 1$

$$\bar{z}^t = W_z x^t + R_z y^{t-1} + b_z$$

$$z^t = g(\bar{z}^t) \qquad \text{\textit{block input}}$$

$$\bar{i}^t = W_i x^t + R_i y^{t-1} + p_i \odot c^{t-1} + b_i$$

$$i^t = \sigma(\bar{i}^t) \qquad \text{\textit{input gate}}$$

$$\bar{f}^t = W_f x^t + R_f y^{t-1} + p_f \odot c^{t-1} + b_f$$

$$f^t = \sigma(\bar{f}^t) \qquad \text{\textit{forget gate}}$$

$$c^t = z^t \odot i^t + c^{t-1} \odot f^t \qquad \text{\textit{cell}}$$

$$\bar{o}^t = W_o x^t + R_o y^{t-1} + p_o \odot c^t + b_o$$

$$o^t = \sigma(\bar{o}^t) \qquad \text{\textit{output gate}}$$

$$y^t = h(c^t) \odot o^t \qquad \text{\textit{block output}}$$

# NIAF: No Input Activation Function

**NIAF:** No Input Activation Function: $g(\mathbf{x}) = \mathbf{x}$

$$\bar{\mathbf{z}}^t = \mathbf{W}_z \mathbf{x}^t + \mathbf{R}_z \mathbf{y}^{t-1} + \mathbf{b}_z$$
$$\mathbf{z}^t = g(\bar{\mathbf{z}}^t) \qquad \qquad \textit{block input}$$
$$\bar{\mathbf{i}}^t = \mathbf{W}_i \mathbf{x}^t + \mathbf{R}_i \mathbf{y}^{t-1} + \mathbf{p}_i \odot \mathbf{c}^{t-1} + \mathbf{b}_i$$
$$\mathbf{i}^t = \sigma(\bar{\mathbf{i}}^t) \qquad \qquad \textit{input gate}$$
$$\bar{\mathbf{f}}^t = \mathbf{W}_f \mathbf{x}^t + \mathbf{R}_f \mathbf{y}^{t-1} + \mathbf{p}_f \odot \mathbf{c}^{t-1} + \mathbf{b}_f$$
$$\mathbf{f}^t = \sigma(\bar{\mathbf{f}}^t) \qquad \qquad \textit{forget gate}$$
$$\mathbf{c}^t = \mathbf{z}^t \odot \mathbf{i}^t + \mathbf{c}^{t-1} \odot \mathbf{f}^t \qquad \qquad \textit{cell}$$
$$\bar{\mathbf{o}}^t = \mathbf{W}_o \mathbf{x}^t + \mathbf{R}_o \mathbf{y}^{t-1} + \mathbf{p}_o \odot \mathbf{c}^t + \mathbf{b}_o$$
$$\mathbf{o}^t = \sigma(\bar{\mathbf{o}}^t) \qquad \qquad \textit{output gate}$$
$$\mathbf{y}^t = h(\mathbf{c}^t) \odot \mathbf{o}^t \qquad \qquad \textit{block output}$$

# NOAF: No Output Activation Function

**NOAF:** No Output Activation Function: $h(\mathbf{x}) = \mathbf{x}$

$$\bar{\mathbf{z}}^t = \mathbf{W}_z \mathbf{x}^t + \mathbf{R}_z \mathbf{y}^{t-1} + \mathbf{b}_z$$
$$\mathbf{z}^t = g(\bar{\mathbf{z}}^t) \qquad \qquad \textit{block input}$$
$$\bar{\mathbf{i}}^t = \mathbf{W}_i \mathbf{x}^t + \mathbf{R}_i \mathbf{y}^{t-1} + \mathbf{p}_i \odot \mathbf{c}^{t-1} + \mathbf{b}_i$$
$$\mathbf{i}^t = \sigma(\bar{\mathbf{i}}^t) \qquad \qquad \textit{input gate}$$
$$\bar{\mathbf{f}}^t = \mathbf{W}_f \mathbf{x}^t + \mathbf{R}_f \mathbf{y}^{t-1} + \mathbf{p}_f \odot \mathbf{c}^{t-1} + \mathbf{b}_f$$
$$\mathbf{f}^t = \sigma(\bar{\mathbf{f}}^t) \qquad \qquad \textit{forget gate}$$
$$\mathbf{c}^t = \mathbf{z}^t \odot \mathbf{i}^t + \mathbf{c}^{t-1} \odot \mathbf{f}^t \qquad \qquad \textit{cell}$$
$$\bar{\mathbf{o}}^t = \mathbf{W}_o \mathbf{x}^t + \mathbf{R}_o \mathbf{y}^{t-1} + \mathbf{p}_o \odot \mathbf{c}^t + \mathbf{b}_o$$
$$\mathbf{o}^t = \sigma(\bar{\mathbf{o}}^t) \qquad \qquad \textit{output gate}$$
$$\mathbf{y}^t = h(\mathbf{c}^t) \odot \mathbf{o}^t \qquad \qquad \textit{block output}$$

# CIFG: Coupled Input and Forget Gate

**CIFG:** Coupled Input and Forget Gate: $\mathbf{f}^t = 1 - \mathbf{i}^t$

$$\bar{\mathbf{z}}^t = \mathbf{W}_z \mathbf{x}^t + \mathbf{R}_z \mathbf{y}^{t-1} + \mathbf{b}_z$$

$$\mathbf{z}^t = g(\bar{\mathbf{z}}^t) \qquad\qquad\qquad block\ input$$

$$\bar{\mathbf{i}}^t = \mathbf{W}_i \mathbf{x}^t + \mathbf{R}_i \mathbf{y}^{t-1} + \mathbf{p}_i \odot \mathbf{c}^{t-1} + \mathbf{b}_i$$

$$\mathbf{i}^t = \sigma(\bar{\mathbf{i}}^t) \qquad\qquad\qquad input\ gate$$

$$\bar{\mathbf{f}}^t = \mathbf{W}_f \mathbf{x}^t + \mathbf{R}_f \mathbf{y}^{t-1} + \mathbf{p}_f \odot \mathbf{c}^{t-1} + \mathbf{b}_f$$

$$\mathbf{f}^t = \sigma(\bar{\mathbf{f}}^t) \qquad\qquad\qquad forget\ gate$$

$$\mathbf{c}^t = \mathbf{z}^t \odot \mathbf{i}^t + \mathbf{c}^{t-1} \odot \mathbf{f}^t \qquad\qquad cell$$

$$\bar{\mathbf{o}}^t = \mathbf{W}_o \mathbf{x}^t + \mathbf{R}_o \mathbf{y}^{t-1} + \mathbf{p}_o \odot \mathbf{c}^t + \mathbf{b}_o$$

$$\mathbf{o}^t = \sigma(\bar{\mathbf{o}}^t) \qquad\qquad\qquad output\ gate$$

$$\mathbf{y}^t = h(\mathbf{c}^t) \odot \mathbf{o}^t \qquad\qquad\qquad block\ output$$

# NP: No Peepholes

**NP: No Peepholes:**

$$\bar{\mathbf{i}}^t = \mathbf{W}_i \mathbf{x}^t + \mathbf{R}_i \mathbf{y}^{t-1} + \mathbf{b}_i$$

$$\bar{\mathbf{f}}^t = \mathbf{W}_f \mathbf{x}^t + \mathbf{R}_f \mathbf{y}^{t-1} + \mathbf{b}_f$$

$$\bar{\mathbf{o}}^t = \mathbf{W}_o \mathbf{x}^t + \mathbf{R}_o \mathbf{y}^{t-1} + \mathbf{b}_o$$

# NP: No Peepholes

$$\bar{z}^t = W_z x^t + R_z y^{t-1} + b_z$$
$$z^t = g(\bar{z}^t) \qquad \qquad \qquad \text{\textit{block input}}$$
$$\bar{i}^t = W_i x^t + R_i y^{t-1} + p_i \odot c^{t-1} + b_i$$
$$i^t = \sigma(\bar{i}^t) \qquad \qquad \qquad \text{\textit{input gate}}$$
$$\bar{f}^t = W_f x^t + R_f y^{t-1} + p_f \odot c^{t-1} + b_f$$
$$f^t = \sigma(\bar{f}^t) \qquad \qquad \qquad \text{\textit{forget gate}}$$
$$c^t = z^t \odot i^t + c^{t-1} \odot f^t \qquad \qquad \qquad \text{\textit{cell}}$$
$$\bar{o}^t = W_o x^t + R_o y^{t-1} + p_o \odot c^t + b_o$$
$$o^t = \sigma(\bar{o}^t) \qquad \qquad \qquad \text{\textit{output gate}}$$
$$y^t = h(c^t) \odot o^t \qquad \qquad \qquad \text{\textit{block output}}$$

# FGR: Full Gate Recurrence

**FGR:** Full Gate Recurrence:

$$\bar{\mathbf{i}}^t = \mathbf{W}_i \mathbf{x}^t + \mathbf{R}_i \mathbf{y}^{t-1} + \mathbf{p}_i \odot \mathbf{c}^{t-1} + \mathbf{b}_i$$
$$+ \mathbf{R}_{ii} \mathbf{i}^{t-1} + \mathbf{R}_{fi} \mathbf{f}^{t-1} + \mathbf{R}_{oi} \mathbf{o}^{t-1}$$

$$\bar{\mathbf{f}}^t = \mathbf{W}_f \mathbf{x}^t + \mathbf{R}_f \mathbf{y}^{t-1} + \mathbf{p}_f \odot \mathbf{c}^{t-1} + \mathbf{b}_f$$
$$+ \mathbf{R}_{if} \mathbf{i}^{t-1} + \mathbf{R}_{ff} \mathbf{f}^{t-1} + \mathbf{R}_{of} \mathbf{o}^{t-1}$$

$$\bar{\mathbf{o}}^t = \mathbf{W}_o \mathbf{x}^t + \mathbf{R}_o \mathbf{y}^{t-1} + \mathbf{p}_o \odot \mathbf{c}^{t-1} + \mathbf{b}_o$$
$$+ \mathbf{R}_{io} \mathbf{i}^{t-1} + \mathbf{R}_{fo} \mathbf{f}^{t-1} + \mathbf{R}_{oo} \mathbf{o}^{t-1}$$

# FGR: Full Gate Recurrence

$$\bar{\mathbf{z}}^t = \mathbf{W}_z \mathbf{x}^t + \mathbf{R}_z \mathbf{y}^{t-1} + \mathbf{b}_z$$

$$\mathbf{z}^t = g(\bar{\mathbf{z}}^t) \qquad\qquad\qquad block\ input$$

$$\bar{\mathbf{i}}^t = \mathbf{W}_i \mathbf{x}^t + \mathbf{R}_i \mathbf{y}^{t-1} + \mathbf{p}_i \odot \mathbf{c}^{t-1} + \mathbf{b}_i$$

$$\mathbf{i}^t = \sigma(\bar{\mathbf{i}}^t) \qquad\qquad\qquad input\ gate$$

$$\bar{\mathbf{f}}^t = \mathbf{W}_f \mathbf{x}^t + \mathbf{R}_f \mathbf{y}^{t-1} + \mathbf{p}_f \odot \mathbf{c}^{t-1} + \mathbf{b}_f$$

$$\mathbf{f}^t = \sigma(\bar{\mathbf{f}}^t) \qquad\qquad\qquad forget\ gate$$

$$\mathbf{c}^t = \mathbf{z}^t \odot \mathbf{i}^t + \mathbf{c}^{t-1} \odot \mathbf{f}^t \qquad\qquad cell$$

$$\bar{\mathbf{o}}^t = \mathbf{W}_o \mathbf{x}^t + \mathbf{R}_o \mathbf{y}^{t-1} + \mathbf{p}_o \odot \mathbf{c}^t + \mathbf{b}_o$$

$$\mathbf{o}^t = \sigma(\bar{\mathbf{o}}^t) \qquad\qquad\qquad output\ gate$$

$$\mathbf{y}^t = h(\mathbf{c}^t) \odot \mathbf{o}^t \qquad\qquad block\ output$$

# Hyperparameter Search

- While there are other methods to efficiently search for good hyperparameters, this paper uses random search has several advantages for our setting:
  - it is easy to implement
  - trivial to parallelize
  - covers the search space more uniformly, thereby improving the follow-up analysis of hyperparameter importance.

- The paper shows 27 random searches (one for each combination of the nine variants and three datasets). Each random search encompasses 200 trials for a total of 5400 trials of randomly sampling the hyperparameters.

- The hyperparameters and ranges are:
  - hidden layer size: log-uniform samples from [20; 200]
  - learning rate: log-uniform samples from [$10^{-6}$; $10^{-2}$]
  - momentum: 1 - log-uniform samples from [0:01; 1:0]
  - standard deviation of Gaussian input noise: uniform samples from [0; 1].

# Results and Discussions:

| Datasets: | State of the art: | Best result: |
|---|---|---|
| IAM Online | 26.9%  (Best LSTM Result) | **9.26%** |
| TIMIT | **26.9%** | 29.6% |
| JSB Chorales | -5.56 | **-8.38** |

# Hyperparameter Analysis:

- Learning Rate: It is the most important hyperparameter and accounts for 67% of the variance on the test set performance.
- We observe there is a sweet-spot at the higher end of learning rate, where the performance is good and the training time is small.
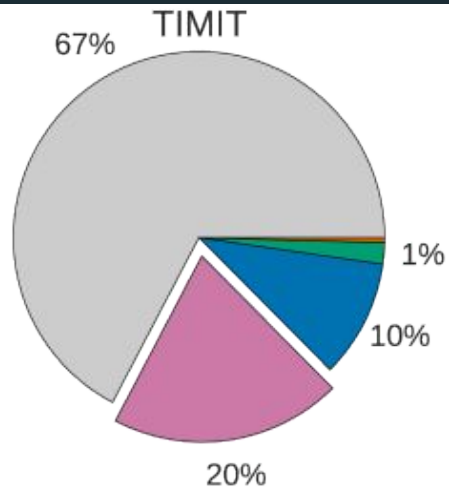
# Hyperparameter Analysis:

- Hidden Layer Size: Not surprisingly the hidden layer size is an important hyperparameter affecting the LSTM network performance. As expected, larger networks perform better.
- It can also be seen in the figure that the required training time increases with the network size.

# Hyperparameter Analysis:

- Input Noise: Additive Gaussian noise on the inputs, a traditional regularizer for neural networks, has been used for LSTM as well. However, we find that not only does it almost always hurt performance, it also slightly increases training times. The only exception is TIMIT, where a small dip in error for the range of [0:2; 0:5] is observed.

# Conclusion:

- We conclude that the most commonly used LSTM architecture (vanilla LSTM) performs reasonably well on various datasets.
- None of the eight investigated modifications significantly improves performance. However, certain modifications such as coupling the input and forget gates or removing peephole connections, simplified LSTMs in our experiments without significantly decreasing performance.

- The forget gate and the output activation function are the most critical components of the LSTM block. Removing any of them significantly impairs performance.
- The learning rate (range: log-uniform samples from [$10^{-6}$; $10^{-2}$]) is the most crucial hyperparameter, followed by the hidden layer size( range: log-uniform samples from [20; 200]).
- The analysis of hyperparameter interactions revealed no apparent structure.

THANK YOU