
LSTM : A Search Space Odyssey

IEEE Transactions on Neural Networks and Learning Systems, 2017

K. Greff et. al.

Presentation by

GIST College Physics Concentration Hanse Kim

Basic Information

Authors : K. Greff, et. al. 4

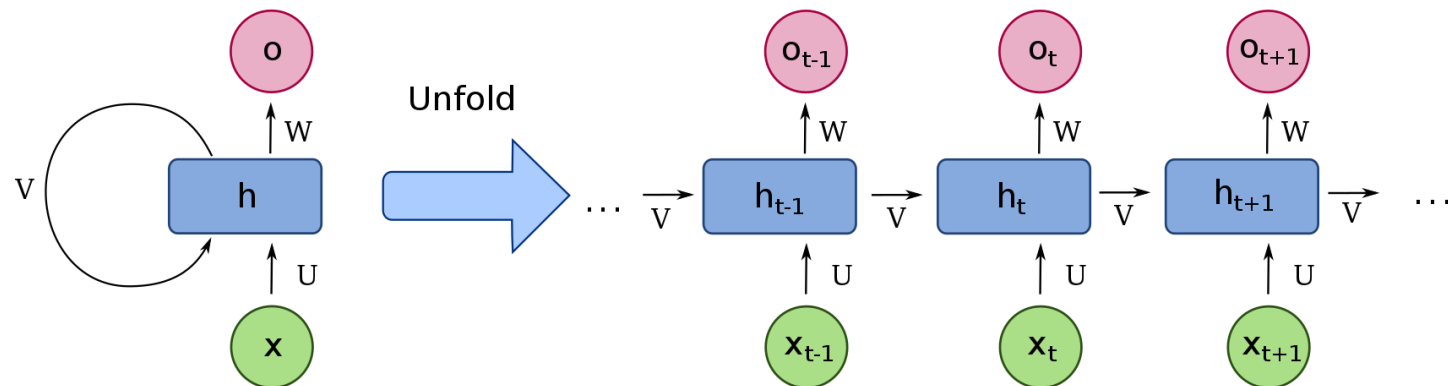
Journal : IEEE Transactions on Neural Networks and Learning Systems, 2017

Citations : 2982

- RNNs and LSTM
- Variants of LSTM
- Analysis of LSTM

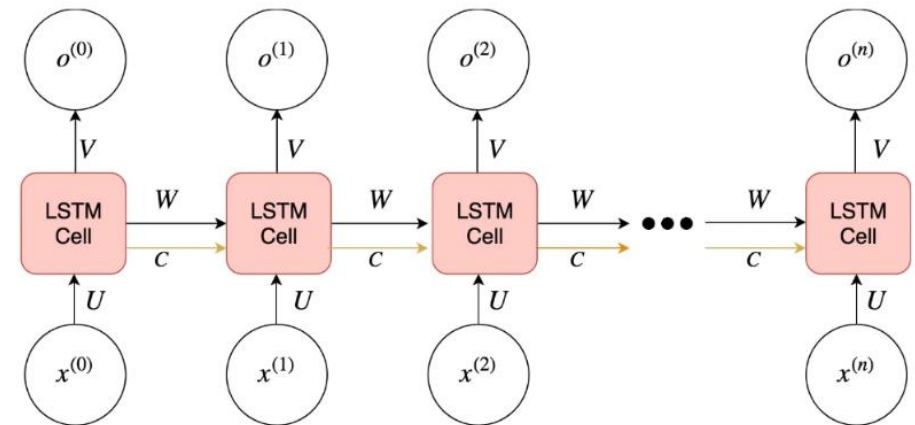
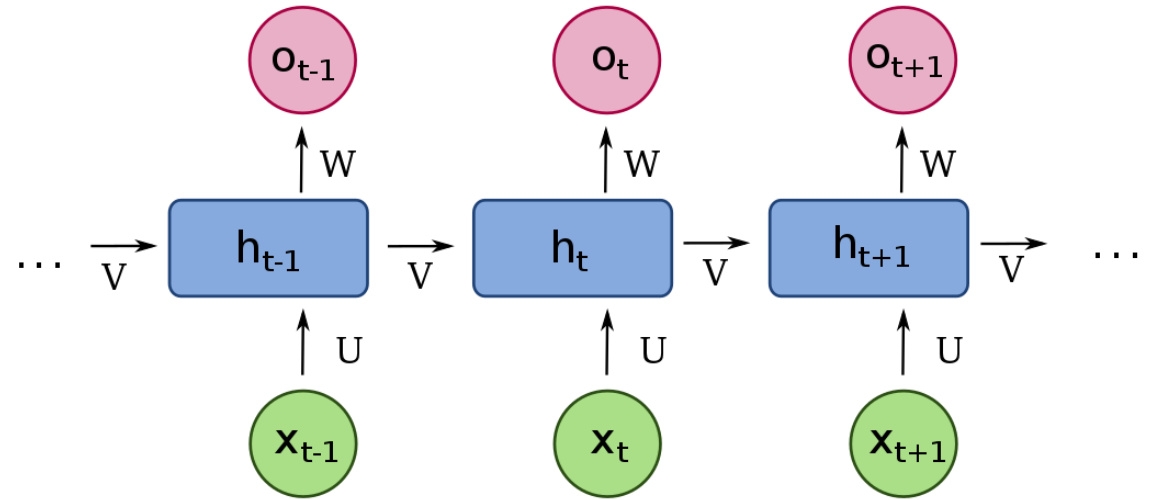
Introduction : RNNs

- RNN : Recurrent Neural Networks
 - NN that can input/output sequences of vectors of variable length
 - Internal state of previous events (memory) used in processing
 - FFNN hidden nodes connected along 'temporal sequence'
- Applications of RNN
 - Time series : Stock market predictions, cryptocurrency
 - NLP : Translation, sentiment analysis



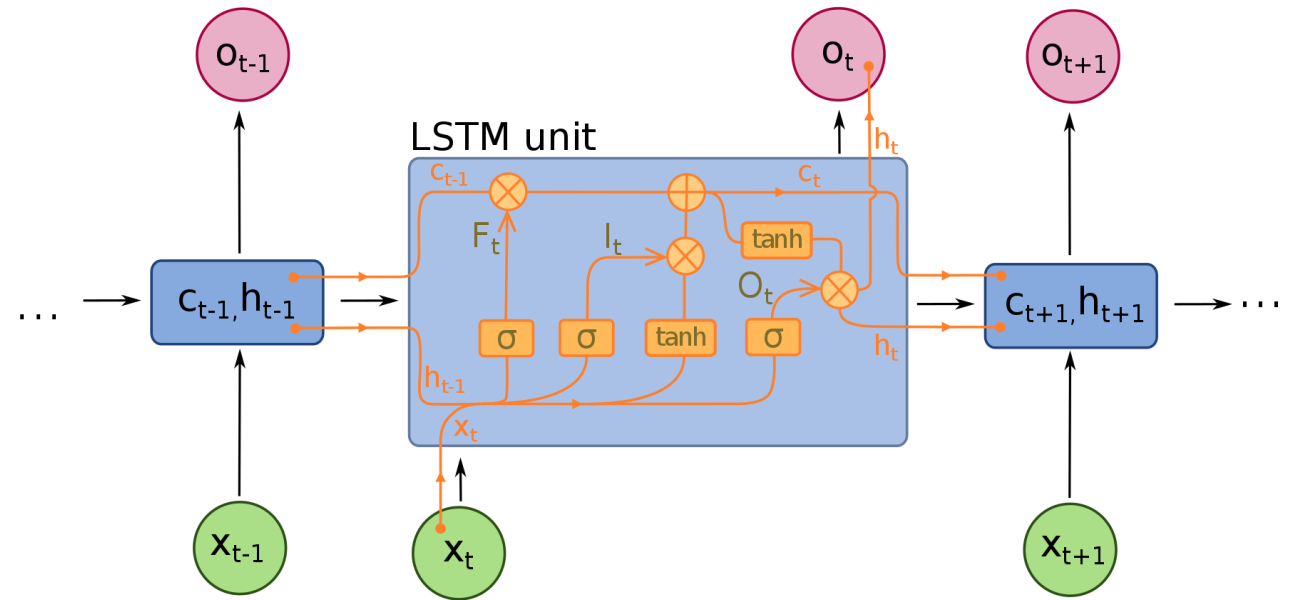
Introduction : Problems of RNNs

- Vanishing Gradient Problem
 - Problem more prominent than in DNNs; same weight V used over all hidden layers
- Skip Connections
 - Reduced rate of parameter vanishing
 - Layers that influence each other are independent from others; acts like DNN
- Gated Recurrent Networks
 - Leaky Recurrent Parameters
 - Set parameter for each time step; new parameter for network to design
 - GRU, LSTM etc.



Introduction : LSTM

- LSTM
 - Most popular variant of GRN
 - Memory cell & Gating Units
- LSTM Unit, Cell State
 - Value passed between LSTM units
 - Each cell can decide to reset it, write to it, or read from it
 - Explicitly expressed in forms of 'gates'; Forget, Input, Output



Introduction : Vanilla LSTM

- Vanilla LSTM
 - Peephole connections
 - Full BPTT

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

$$z^t = g(\bar{z}^t)$$

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

$$i^t = \sigma(\bar{i}^t)$$

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

$$f^t = \sigma(\bar{f}^t)$$

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

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

$$o^t = \sigma(\bar{o}^t)$$

$$y^t = h(c^t) \odot o^t$$

block input

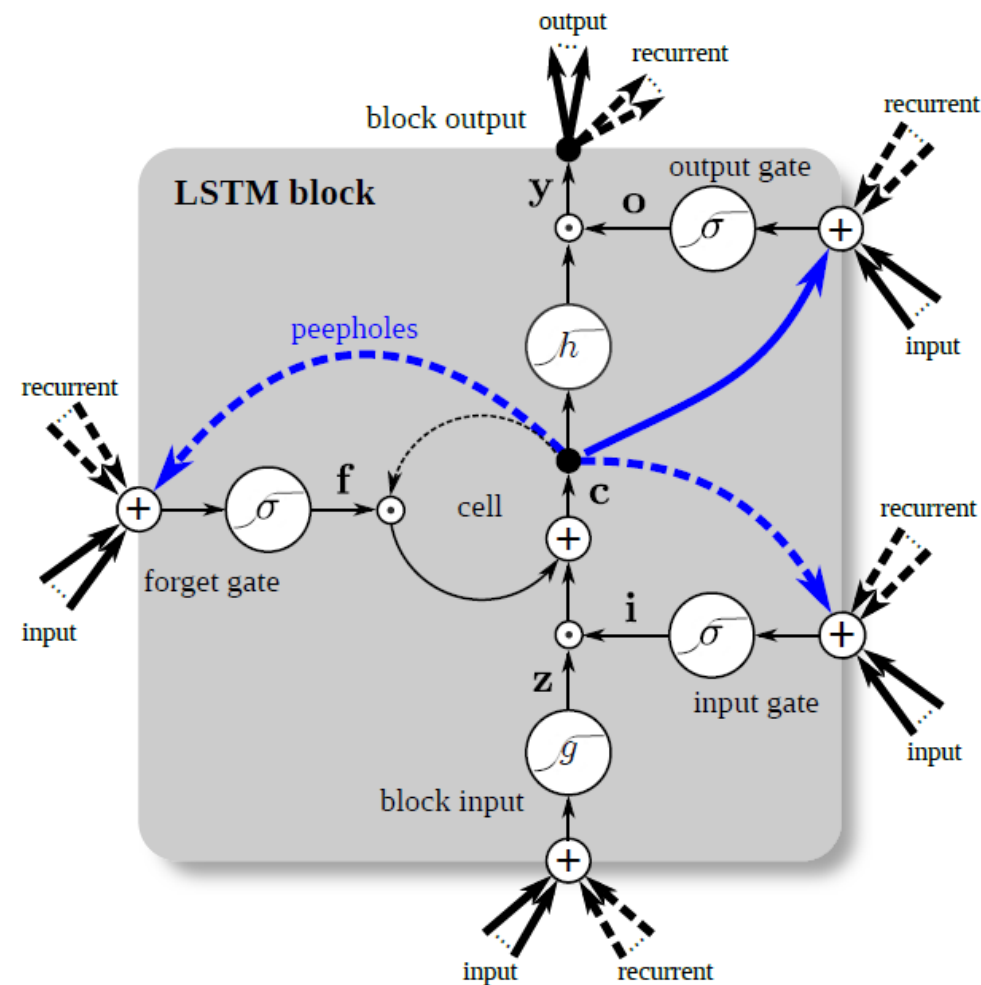
input gate

forget gate

cell

output gate

block output



LSTM Variants

- One aspect each tuned
 - No peepholes
 - Full Gate Recurrence
- Recurrent connections between all gates

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

NIG: No Input Gate: $\mathbf{i}^t = 1$
NFG: No Forget Gate: $\mathbf{f}^t = 1$
NOG: No Output Gate: $\mathbf{o}^t = 1$
NIAF: No Input Activation Function: $g(\mathbf{x}) = \mathbf{x}$
NOAF: No Output Activation Function: $h(\mathbf{x}) = \mathbf{x}$
CIFG: Coupled Input and Forget Gate: $\mathbf{f}^t = 1 - \mathbf{i}^t$
NP: No Peepholes:

$$\begin{aligned}\bar{\mathbf{i}}^t &= \mathbf{W}_i \mathbf{x}^t + \mathbf{R}_i \mathbf{y}^{t-1} + \mathbf{b}_i \\ \bar{\mathbf{f}}^t &= \mathbf{W}_f \mathbf{x}^t + \mathbf{R}_f \mathbf{y}^{t-1} + \mathbf{b}_f \\ \bar{\mathbf{o}}^t &= \mathbf{W}_o \mathbf{x}^t + \mathbf{R}_o \mathbf{y}^{t-1} + \mathbf{b}_o\end{aligned}$$

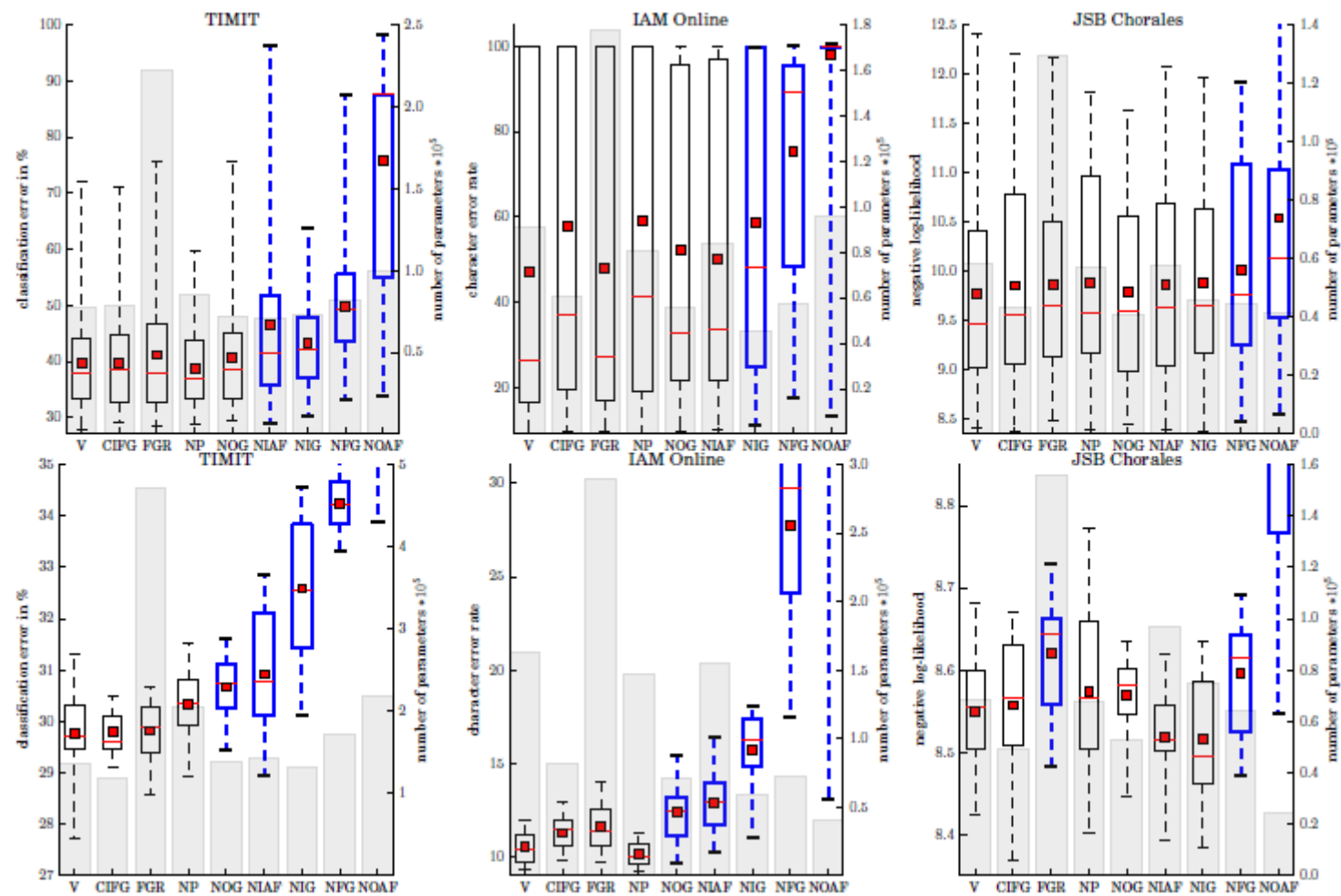
FGR: Full Gate Recurrence:

$$\begin{aligned}\bar{\mathbf{i}}^t &= \mathbf{W}_i \mathbf{x}^t + \mathbf{R}_i \mathbf{y}^{t-1} + \mathbf{p}_i \odot \mathbf{c}^{t-1} + \mathbf{b}_i \\ &\quad + \mathbf{R}_{ii} \mathbf{i}^{t-1} + \mathbf{R}_{fi} \mathbf{f}^{t-1} + \mathbf{R}_{oi} \mathbf{o}^{t-1} \\ \bar{\mathbf{f}}^t &= \mathbf{W}_f \mathbf{x}^t + \mathbf{R}_f \mathbf{y}^{t-1} + \mathbf{p}_f \odot \mathbf{c}^{t-1} + \mathbf{b}_f \\ &\quad + \mathbf{R}_{if} \mathbf{i}^{t-1} + \mathbf{R}_{ff} \mathbf{f}^{t-1} + \mathbf{R}_{of} \mathbf{o}^{t-1} \\ \bar{\mathbf{o}}^t &= \mathbf{W}_o \mathbf{x}^t + \mathbf{R}_o \mathbf{y}^{t-1} + \mathbf{p}_o \odot \mathbf{c}^{t-1} + \mathbf{b}_o \\ &\quad + \mathbf{R}_{io} \mathbf{i}^{t-1} + \mathbf{R}_{fo} \mathbf{f}^{t-1} + \mathbf{R}_{oo} \mathbf{o}^{t-1}\end{aligned}$$

Evaluation

- Datasets
 - TIMIT : Speech corpus, acoustic modelling benchmark
 - IAM Online : Handwriting database, time series of pen movement
 - JSB Chorales : Next-step prediction for music
- Network Architecture
 - JSB Chorales : Single-layer LSTM
 - TIMIT, IAM Online : Bi-directional LSTM
- Hyperparametres evaluated by random search

Results & Discussion



Conclusion

- LSTM attempts to improve upon RNN
- Vanishing gradient solved by non-linear output activation function, forget gate; ability to 'memorise' and 'forget'
- Empirical analysis backs the assertion

References

- Websites
 - https://en.wikipedia.org/wiki/Recurrent_neural_network