

SCC5809 – Redes Neurais

Avaliação Empírica de Unidades Recorrentes na Modelagem de Sequências

Leonardo Leite de Melo
Wagner W. Ávila Bombardelli

04 de outubro
2018

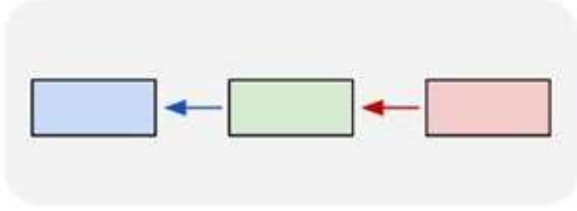
SUMÁRIO

1. Introdução
Aplicações com RNN / Resumo do Artigo
2. Contextualização
RNN / LTSM / GRU / Discussão
3. Experimentação
Tarefas e Conjunto de Dados / Modelos
4. Resultados e Análise
5. Conclusão

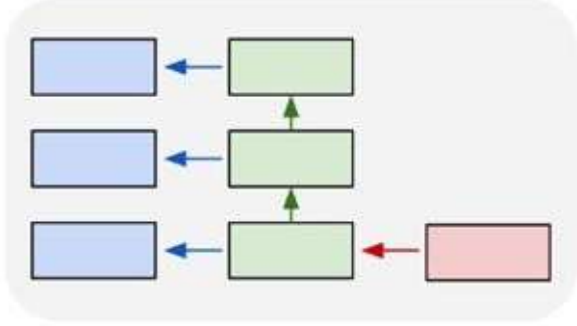
INTRODUÇÃO



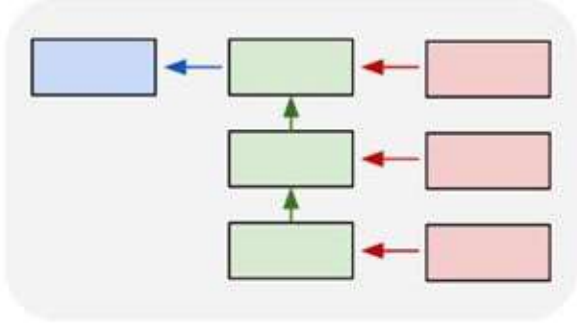
one to one



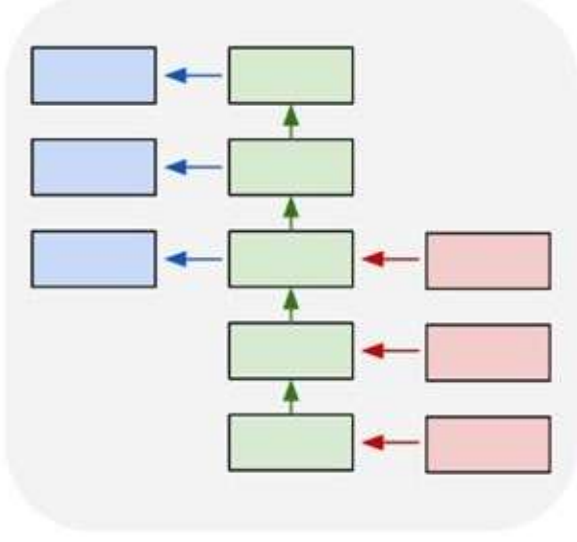
one to many



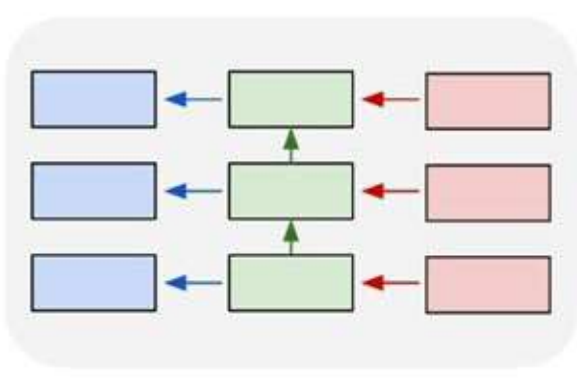
many to one



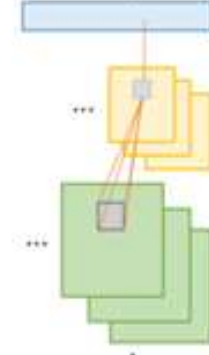
many to many



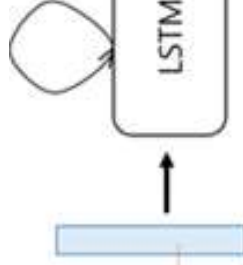
many to many



1. Input Image



2. Deep CNN



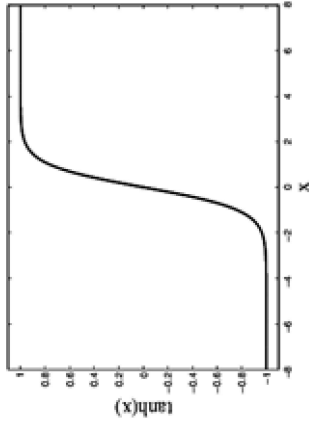
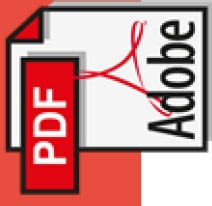
3. RNN

The man at bat
readies to swing
at the pitch

4. Caption



RESUMO DO ARTIGO



LSTM

GRU

On the Properties of Neural Machine Translation: Encoder–Decoder Approaches

Kyunghyun Cho Bart van Merriënboer
Université de Montréal

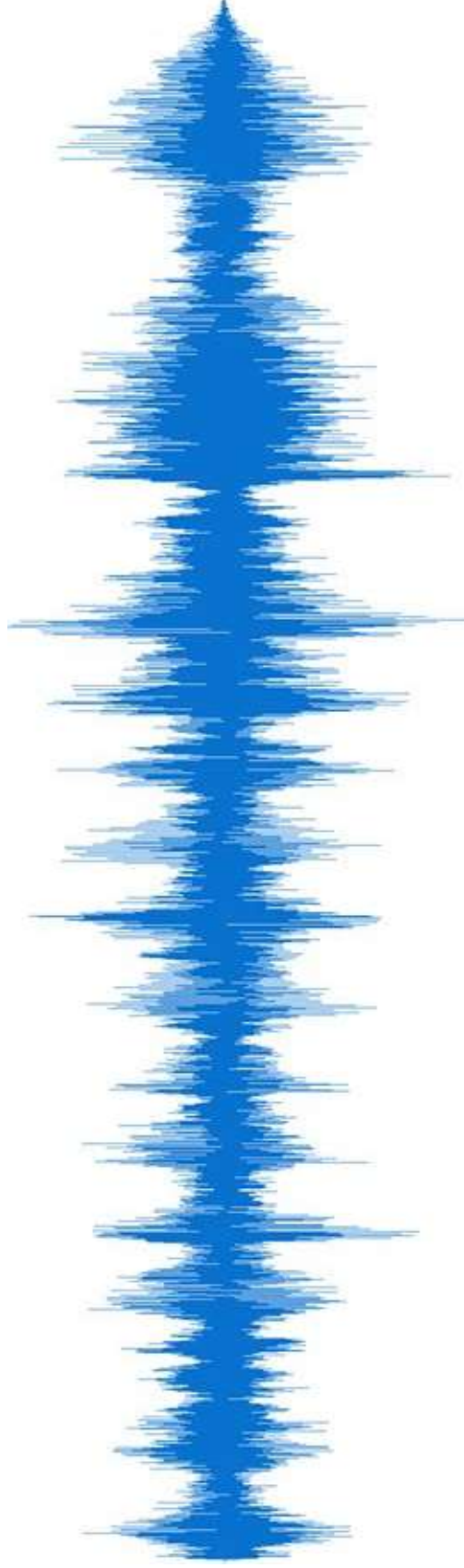
Dzmitry Bahdanau*
Jacobs University, Germany

Empirical Evaluation of
Gated Recurrent Neural Networks
on Sequence Modeling

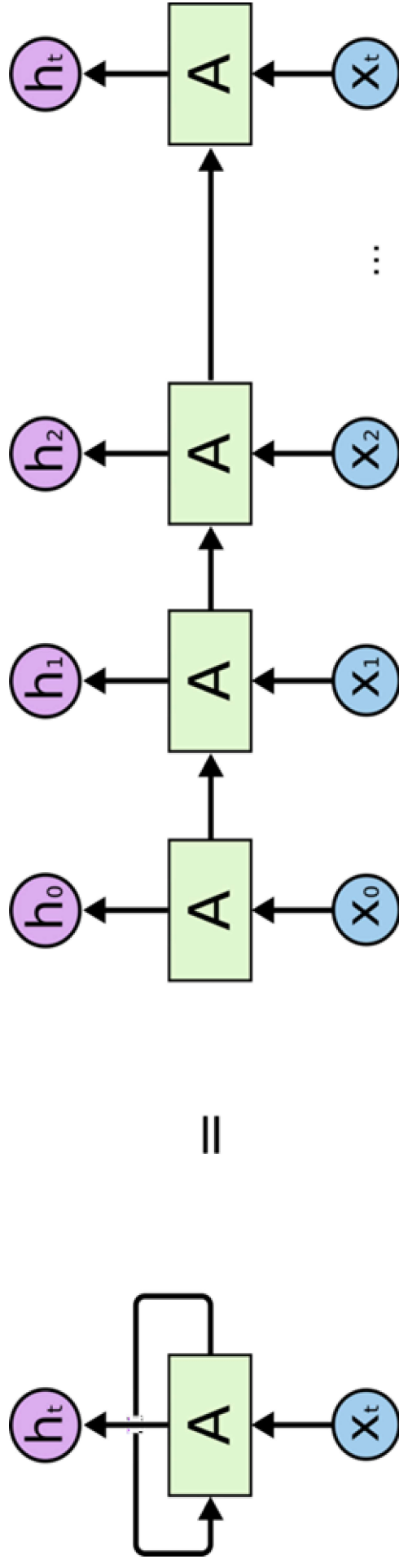
Junyoung Chung
Caglar Gulcehre

Kyunghyun Cho
Université de Montréal

Yoshua Bengio
Université de Montréal
CIFAR Senior Fellow



CONTEXTUALIZAÇÃO



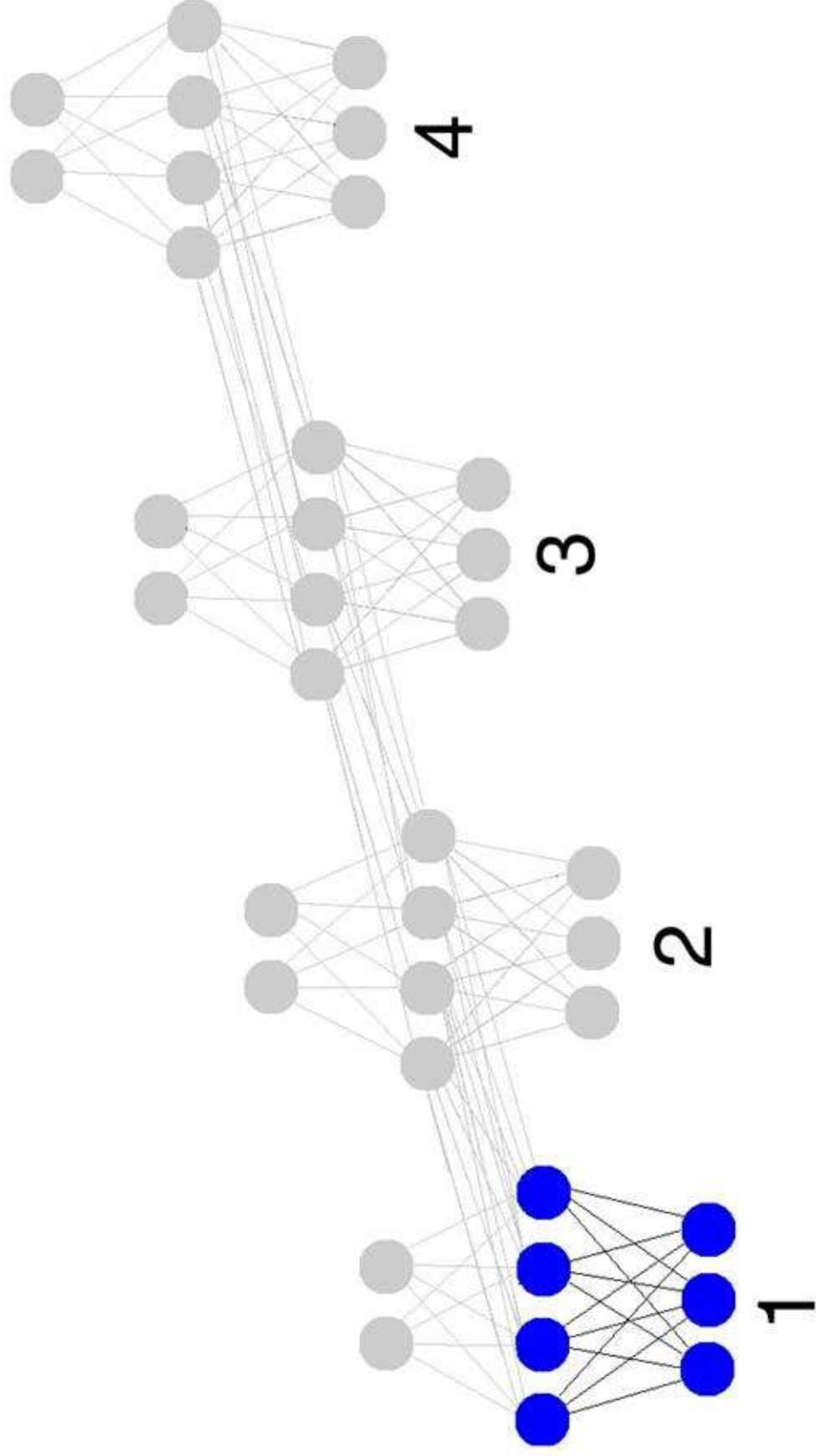
$$\mathbf{h}_t = \begin{cases} 0, & t=0 \\ \phi(\mathbf{h}_{t-1}, \mathbf{x}_t), & \text{otherwise} \end{cases}$$

$$\mathbf{h}_t = g(W\mathbf{x}_t + U\mathbf{h}_{t-1})$$

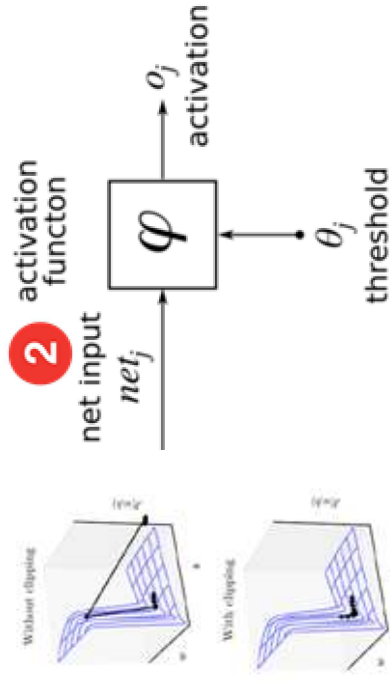
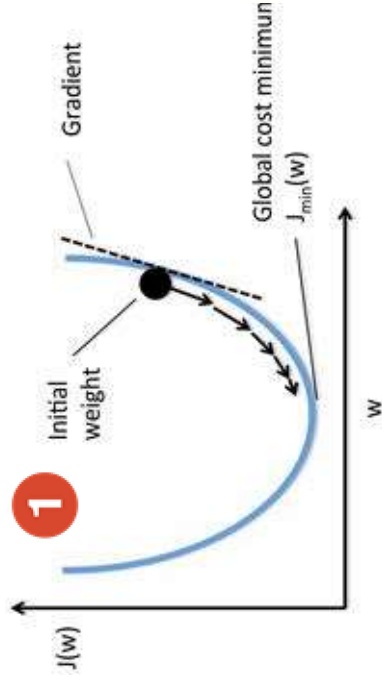
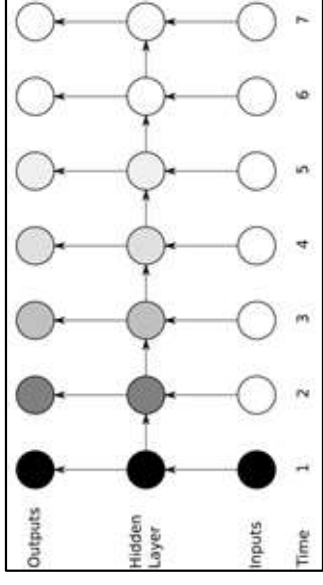
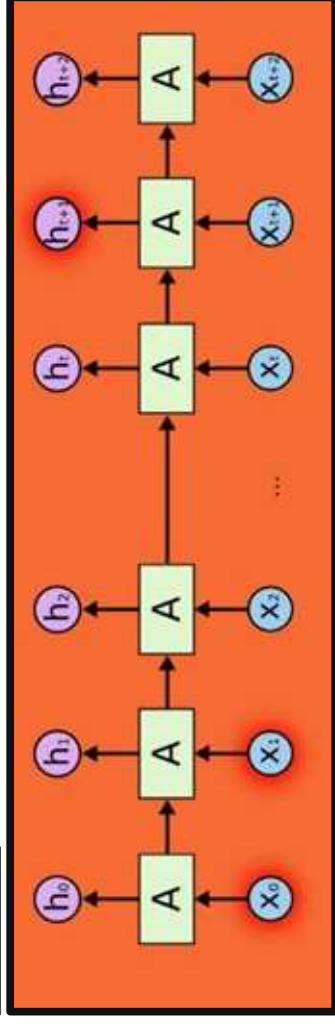
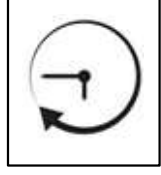
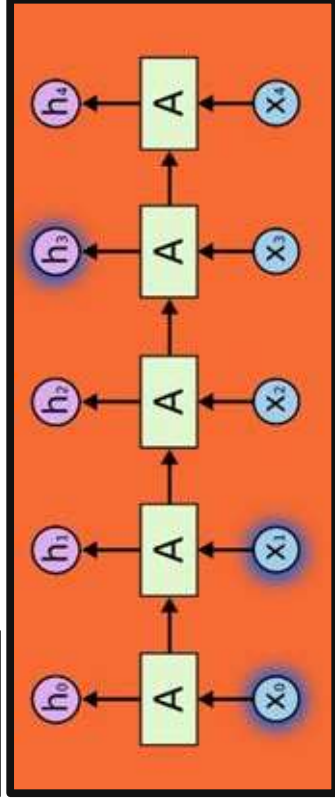
$$p(x_1, \dots, x_T) = p(x_1)p(x_2 | x_1)p(x_3 | x_1, x_2) \dots p(x_T | x_1, \dots, x_{T-1})$$

$$p(x_t | x_1, \dots, x_{t-1}) = g(h_t)$$

CONTEXTUALIZAÇÃO

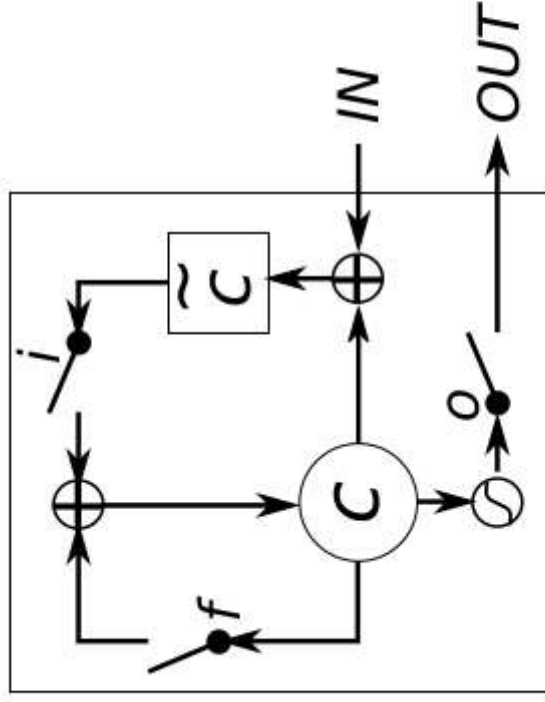
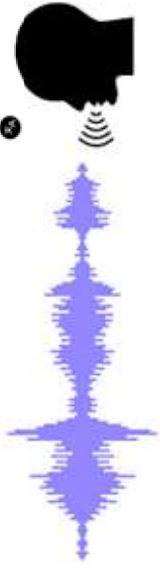


VANISHING GRADIENT PROBLEM



LSTM x GRU

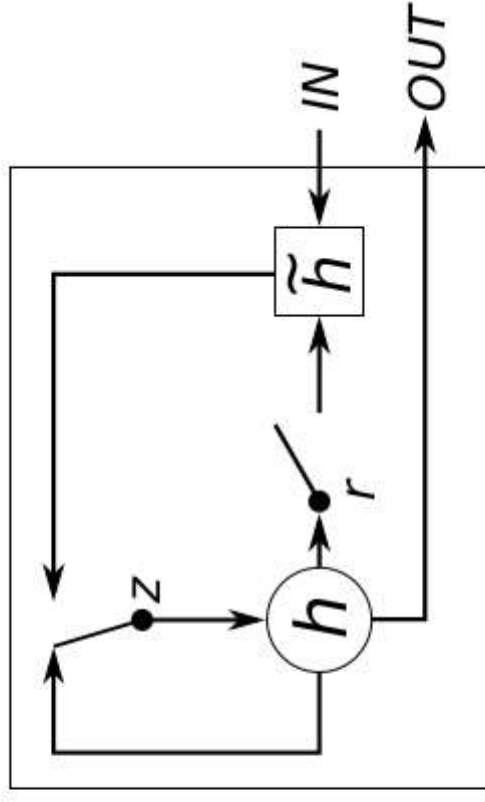
Hello你好



(a) Long Short-Term Memory



[Hochreiter and Schmidhuber, 1997]



(b) Gated Recurrent Unit

[Cho et al., 2014]

LSTM was the winner of the ICDAR (International Conference on Document Analysis and Recognition) 2009 - Handwriting Competition for the best known results in h

LSTM

$$h_t^j = o_t^j \tanh(c_t^j)$$

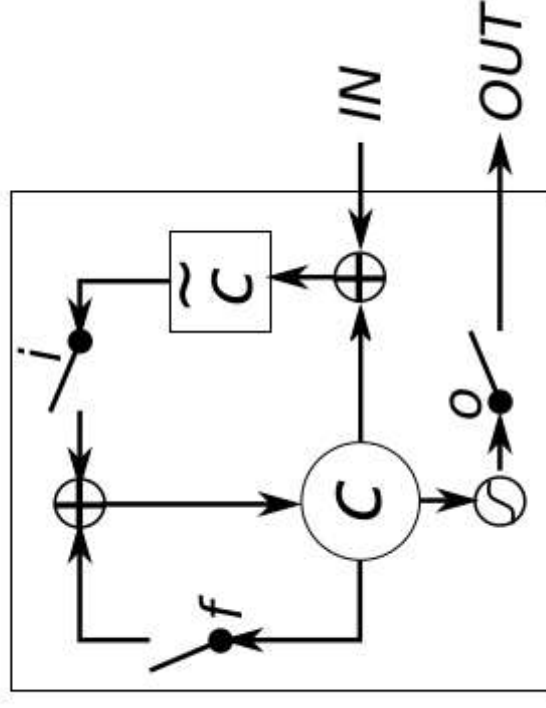
$$o_t^j = \sigma(W_o \mathbf{x}_t + U_o \mathbf{h}_{t-1} + V_o \mathbf{c}_t)^j$$

$$c_t^j = f_t^j c_{t-1}^j + i_t^j \tilde{c}_t^j$$

$$\tilde{c}_t^j = \tanh(W_c \mathbf{x}_t + U_c \mathbf{h}_{t-1})^j$$

$$f_t^j = \sigma(W_f \mathbf{x}_t + U_f \mathbf{h}_{t-1} + V_f \mathbf{c}_{t-1})^j$$

$$i_t^j = \sigma(W_i \mathbf{x}_t + U_i \mathbf{h}_{t-1} + V_i \mathbf{c}_{t-1})^j$$

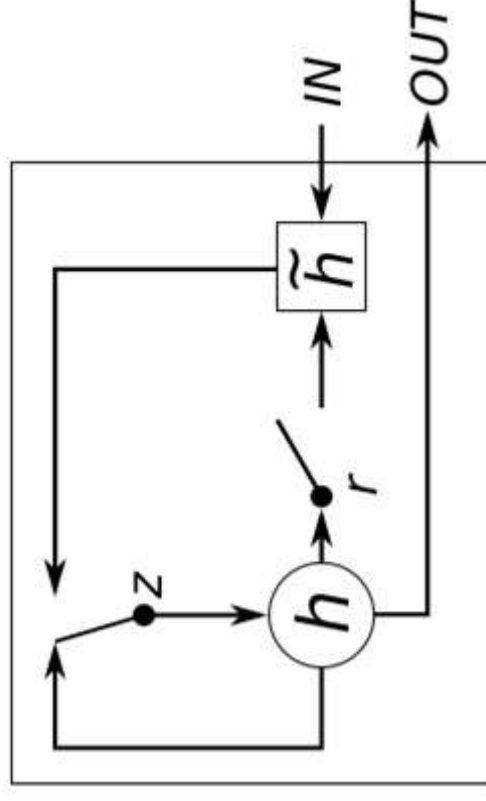


(a) Long Short-Term Memory

Introduziu o conceito de célula de memória.
A célula de memória pode manter seu valor por um tempo curto ou longo como uma função de suas entradas, o que permite que a célula se lembre daquilo que é

GRU

$$\begin{aligned}h_t^j &= (1 - z_t^j)h_{t-1}^j + z_t^j\tilde{h}_t^j \\z_t^j &= \sigma(W_z\mathbf{x}_t + U_z\mathbf{h}_{t-1})^j \\ \tilde{h}_t^j &= \tanh(W\mathbf{x}_t + U(\mathbf{r}_t \odot \mathbf{h}_{t-1}))^j \\r_t^j &= \sigma(W_r\mathbf{x}_t + U_r\mathbf{h}_{t-1})^j\end{aligned}$$



(b) Gated Recurrent Unit

A GRU é mais simples do que a LSTM, pode ser treinada com mais rapidez e pode ser mais eficiente em sua execução. No entanto, a LSTM pode ser mais expressiva.

DISCUSSÃO

SEMELHANÇA

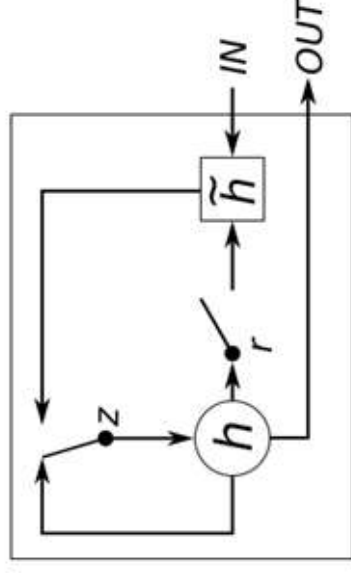
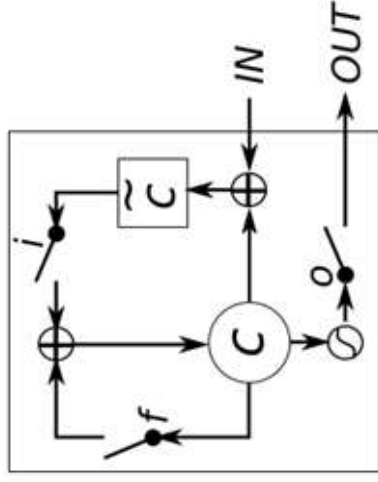


- (i) É fácil para cada unidade lembrar a existência de um recurso específico no fluxo de entrada
- (ii) Essa adição cria efetivamente caminhos de atalho que ignoram várias etapas temporais.

DISCUSSÃO

DIFERENÇAS

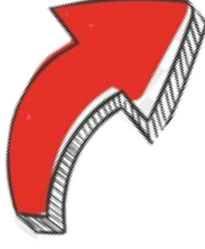
- (i) Na unidade LSTM, a quantidade de conteúdo da memória acessada por outras unidades n
- (ii) Localização do gate de *input* e gate de *reset* correspondente.



INFORMAÇÕES ADICIONAIS

Attention Is All You Need

Ashish Vaswani* Google Brain avaswani@google.com	Noam Shazeer* Google Brain noam@google.com	Niki Parmar* Google Research nikip@google.com	Jakob Uszkoreit* Google Research usz@google.com
Llion Jones* Google Research llion@google.com	Aidan N. Gomez* † University of Toronto aidan@cs.toronto.edu	Illia Polosukhin* ‡ illia.polosukhin@gmail.com	Lukasz Kaiser* Google Brain lukaszkaizer@google.com



Recurrent models typically factor computation along the symbol positions of the input and output sequences. Aligning the positions to steps in computation time, they generate a sequence of hidden states h_t , as a function of the previous hidden state h_{t-1} and the input for position t . This inherently sequential nature precludes parallelization within training examples, which becomes critical at longer sequence lengths, as memory constraints limit batching across examples. Recent work has achieved significant improvements in computational efficiency through factorization tricks [21] and conditional computation [22], while also improving model performance in case of the latter.

RNN (LSTM e GRU em particular) foram firmemente estabelecidas como abordagens de última geração em modelagem de sequência, como processamento de lin

Inúmeros esforços, desde então, continuaram a empurrar os limites de modelos de linguagem recorrente e arquiteturas e ncoder-decoder .

(VASWANI, Ashish et al. Attention is all you need. In: Advances in Neural Information Processing Systems. 2017. p. 5998-6008.)

CONFIGURAÇÃO DO EXPERIMENTO I

- Comparação de LSTM, GRU e tanh na tarefa de modelar uma sequência de dados.
- Aprender a probabilidade de distribuição em uma sequência de dados.

$$\max_{\theta} \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} \log p(x_t^n \mid x_1^n, \dots, x_{t-1}^n; \theta),$$

CONFIGURAÇÃO DO EXPERIMENTO II

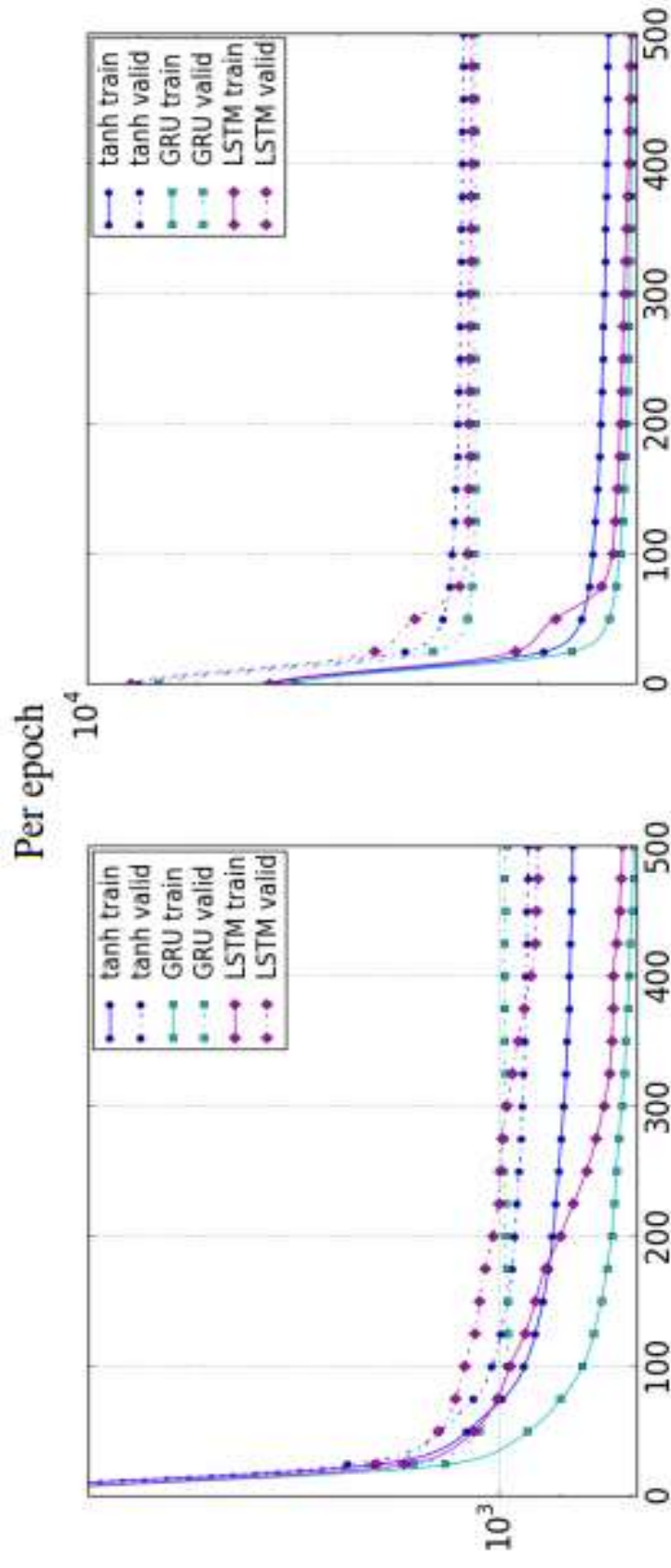
Unit	# of Units	# of Parameters
Polyphonic music modeling		
LSTM	36	$\approx 19.8 \times 10^3$
GRU	46	$\approx 20.2 \times 10^3$
tanh	100	$\approx 20.1 \times 10^3$
Speech signal modeling		
LSTM	195	$\approx 169.1 \times 10^3$
GRU	227	$\approx 168.9 \times 10^3$
tanh	400	$\approx 168.4 \times 10^3$

Table 1: The sizes of the models tested in the experiments.

		tanh		GRU		LSTM	
Music Datasets	Nottingham	train	test	3.22	2.79	3.08	3.08
		test		3.13	3.23	3.20	3.20
	JSB Chorales	train	test	8.82	6.94	8.15	8.15
		test		9.10	8.54	8.67	8.67
	MuseData	train	test	5.64	5.06	5.18	5.18
		test		6.23	5.99	6.23	6.23
Ubisoft Datasets	Piano-midi	train	test	5.64	4.93	6.49	6.49
		test		9.03	8.82	9.03	9.03
	Ubisoft dataset A	train	test	6.29	2.31	1.44	1.44
		test		6.44	3.59	2.70	2.70
	Ubisoft dataset B	train	test	7.61	0.38	0.80	0.80
		test		7.62	0.88	1.26	1.26

Table 2: The average negative log-probabilities of the training and test sets.

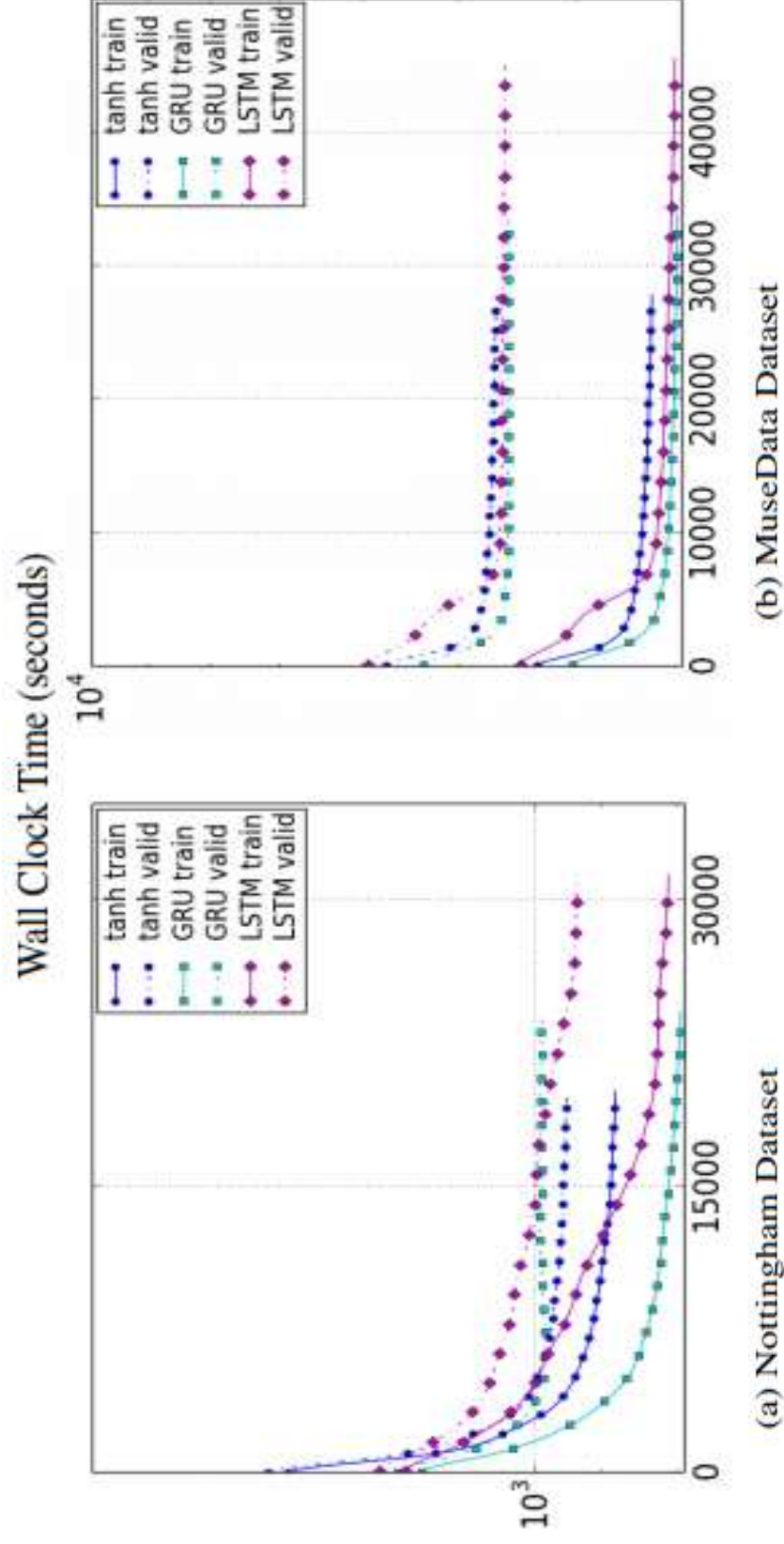
RESULTADOS E ANÁLISE I



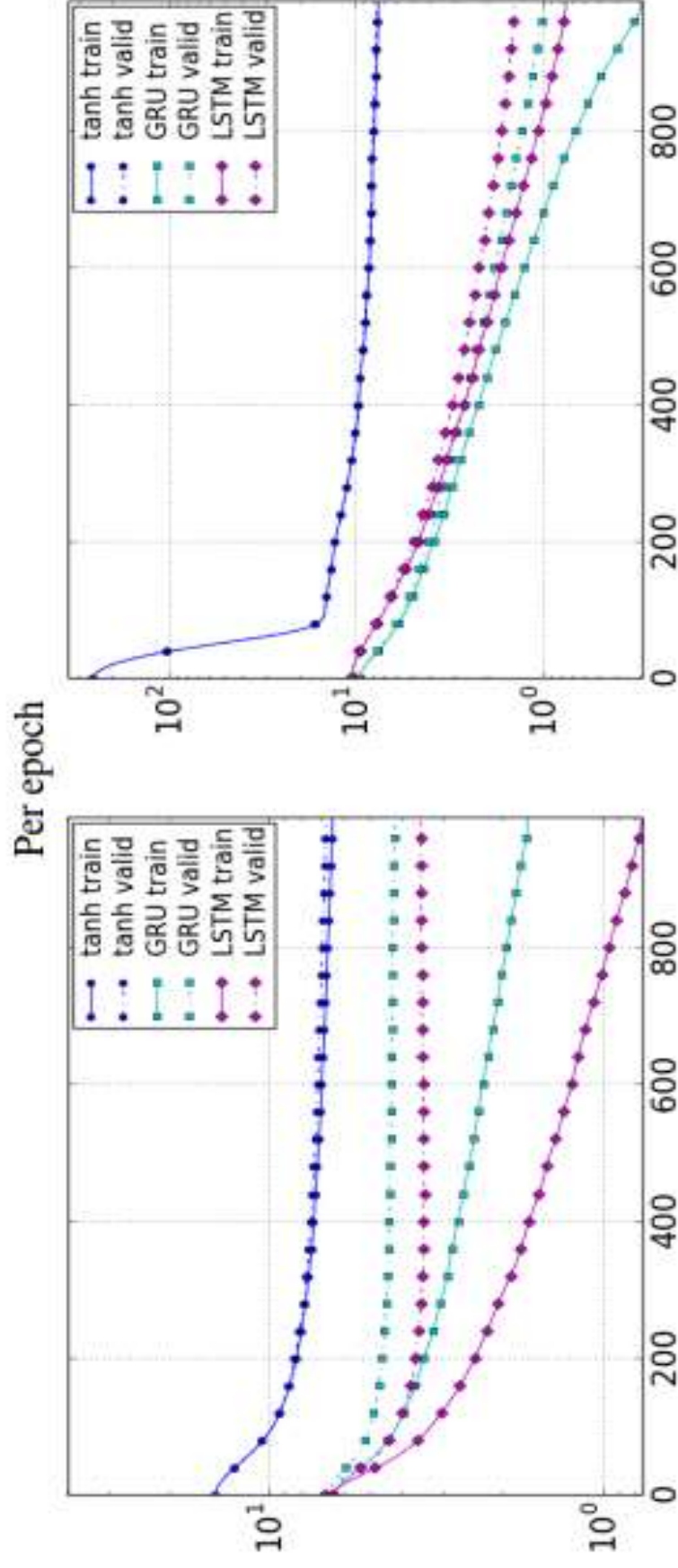
(a) Nottingham Dataset

(b) MuseData Dataset

RESULTADOS E ANÁLISE II



RESULTADOS E ANÁLISE III



RESULTADOS E ANÁLISE IV

