

Modelos de Linguagem Neurais

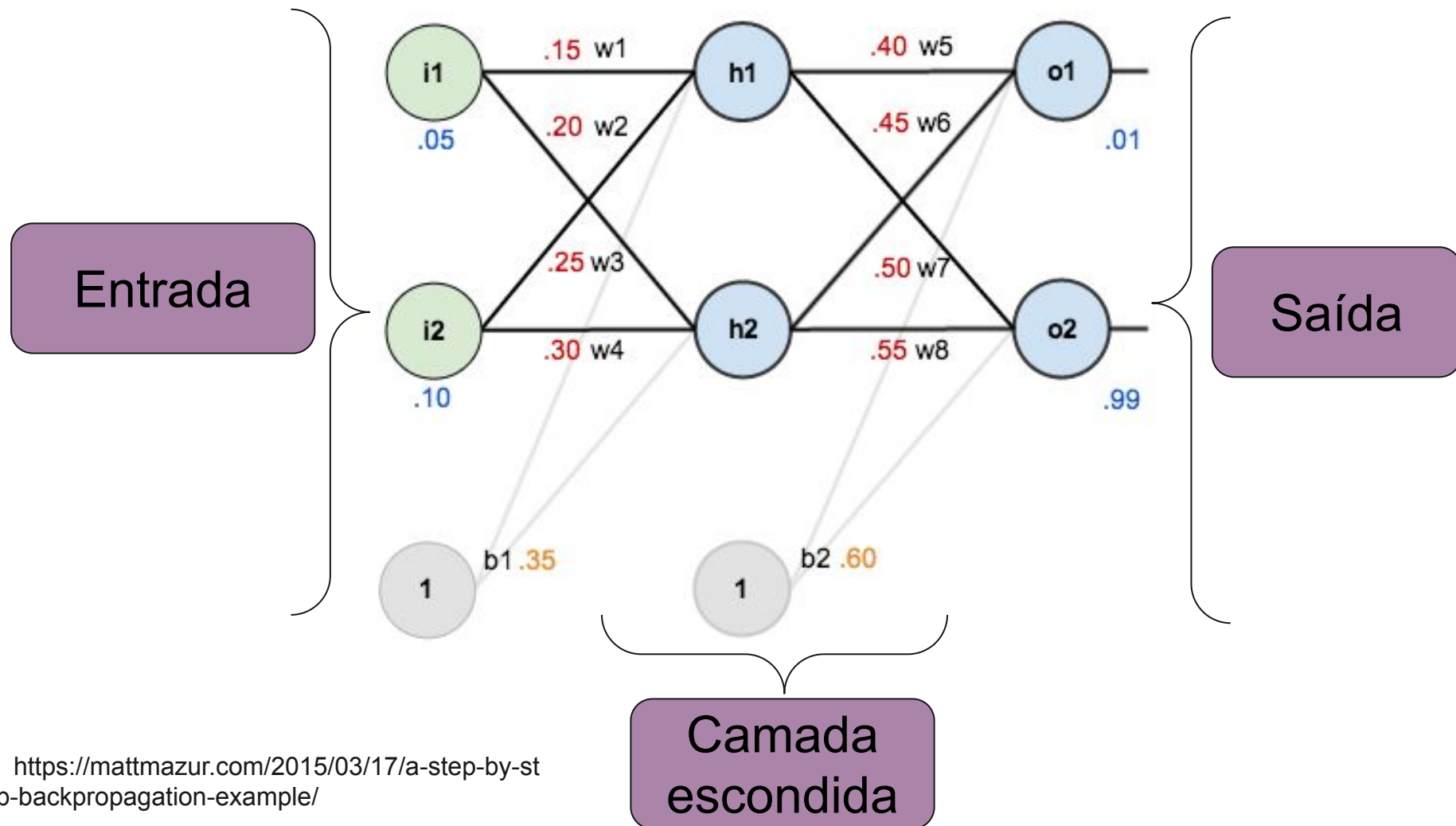
Com Redes Recorrentes

"Time will explain"

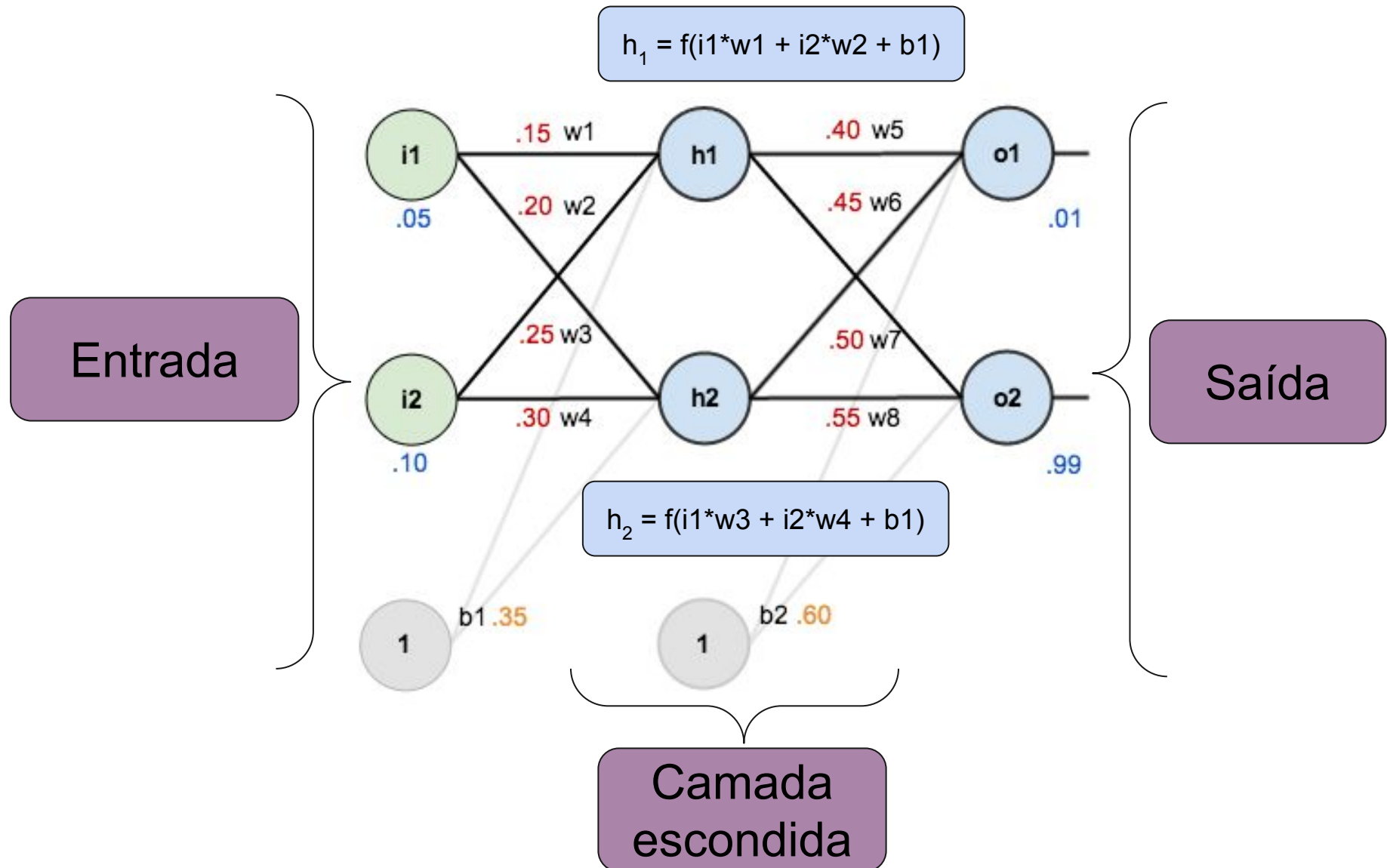


Profa Aline Paes
alinepaes@ic.uff.br

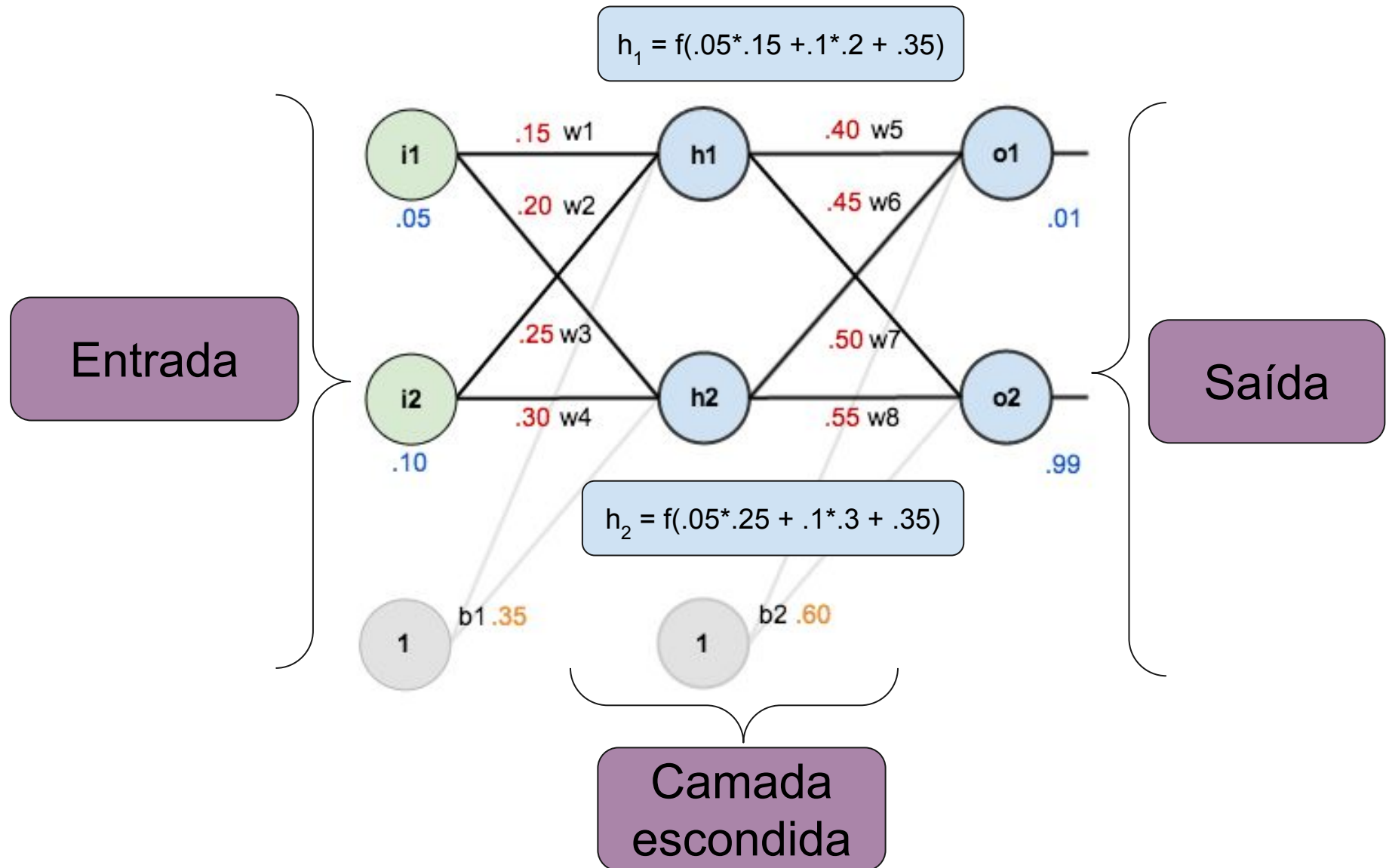
Rede neural feed-forward



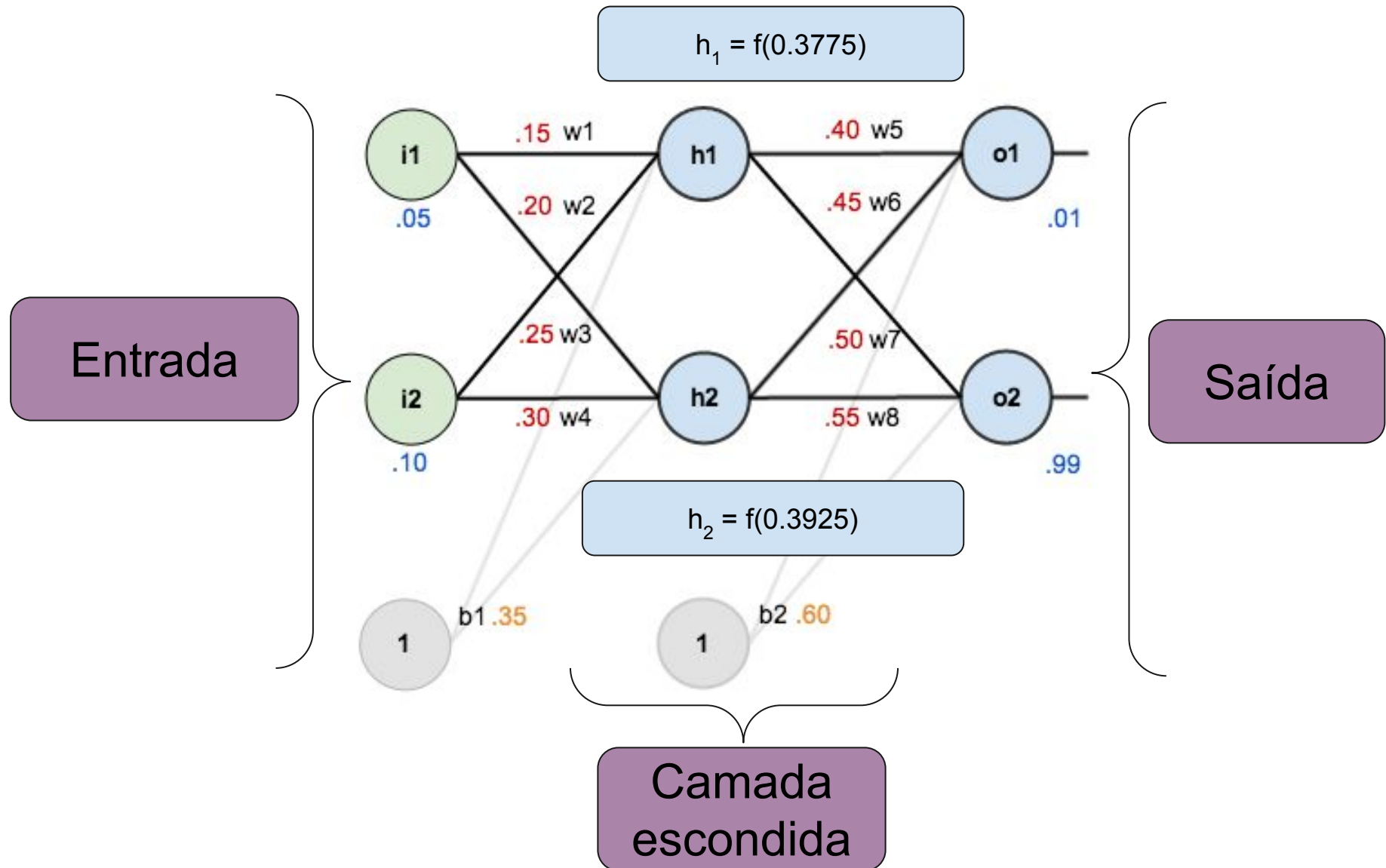
Rede neural feed-forward



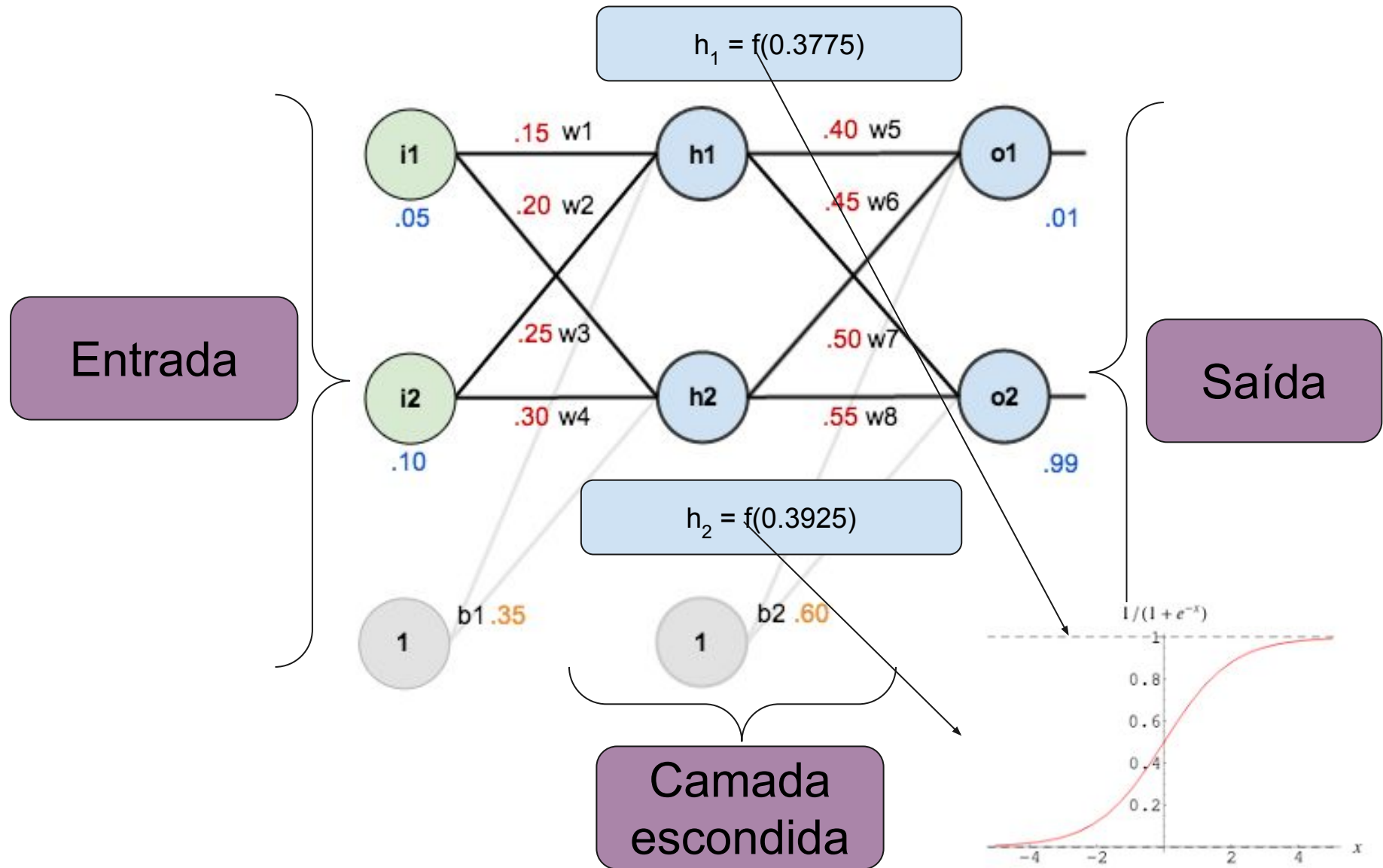
Rede neural feed-forward



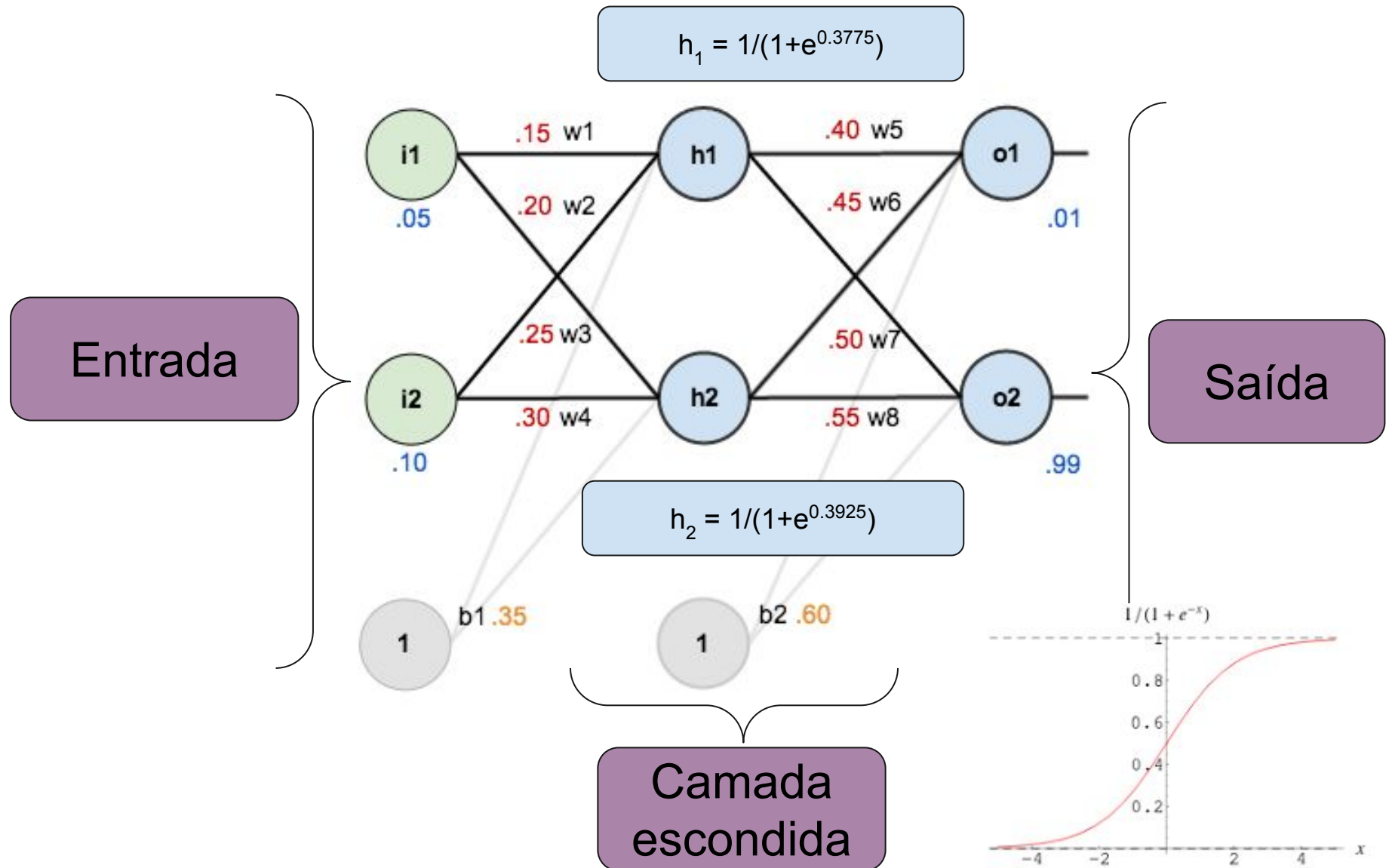
Rede neural feed-forward



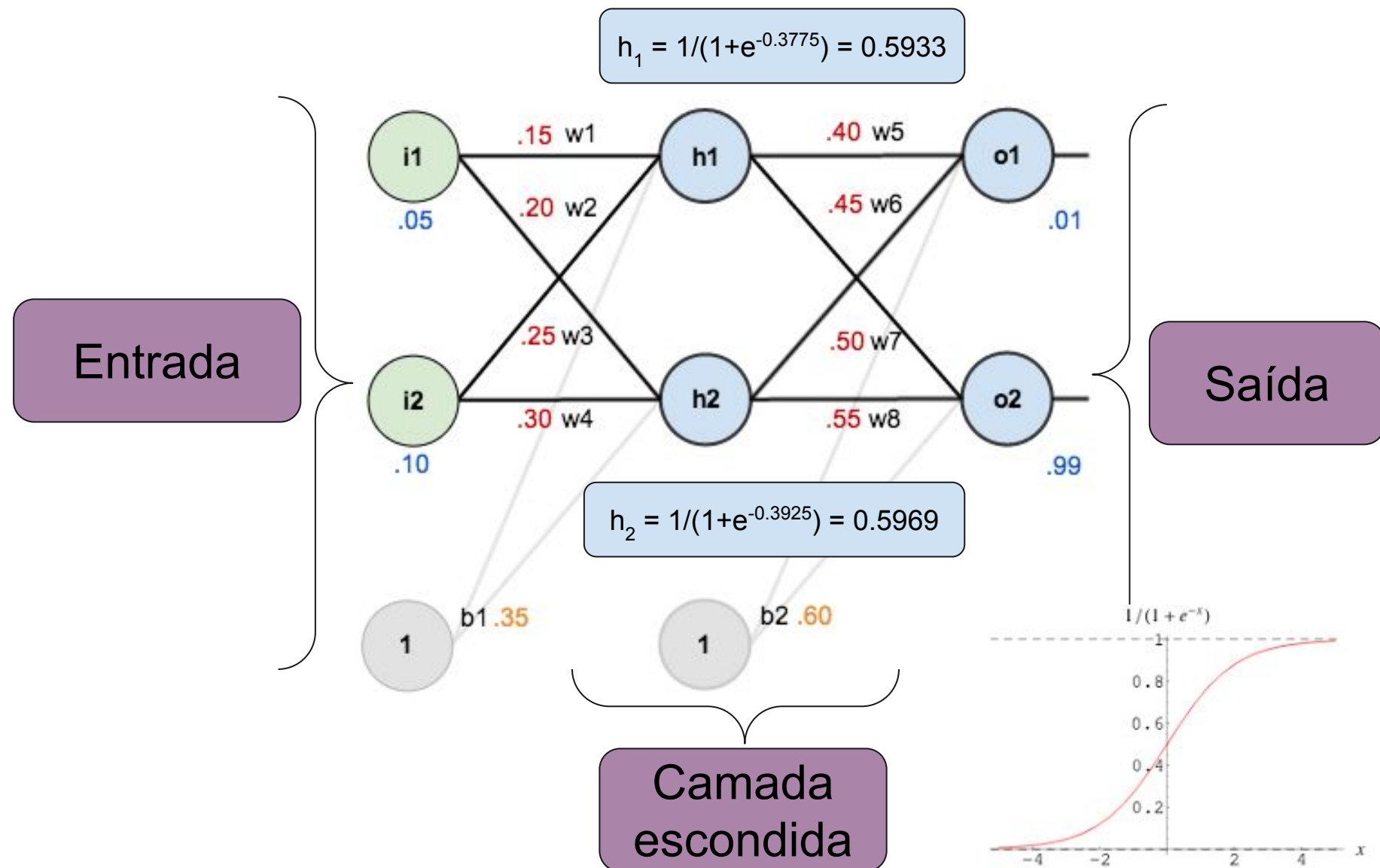
Rede neural feed-forward



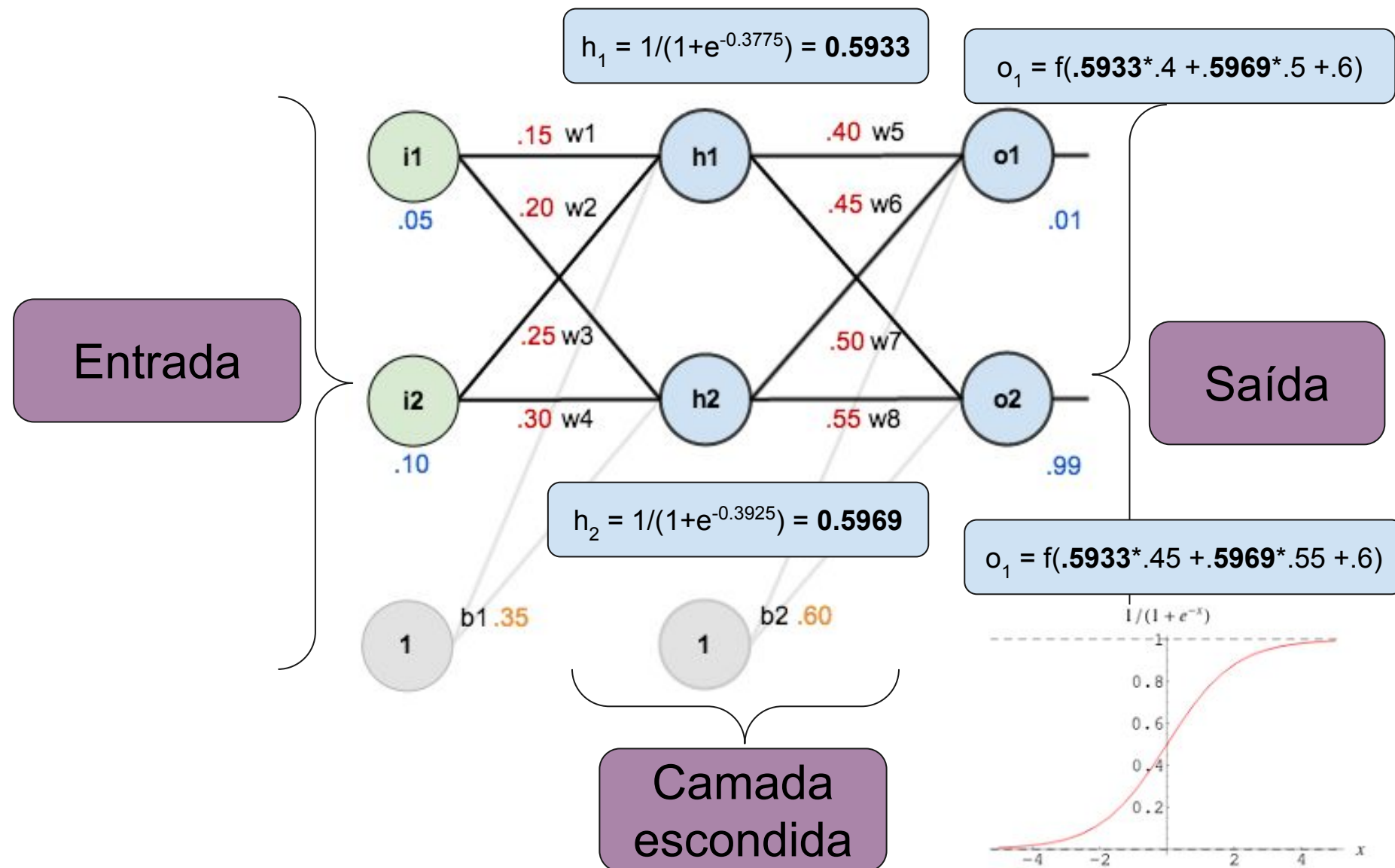
Rede neural feed-forward



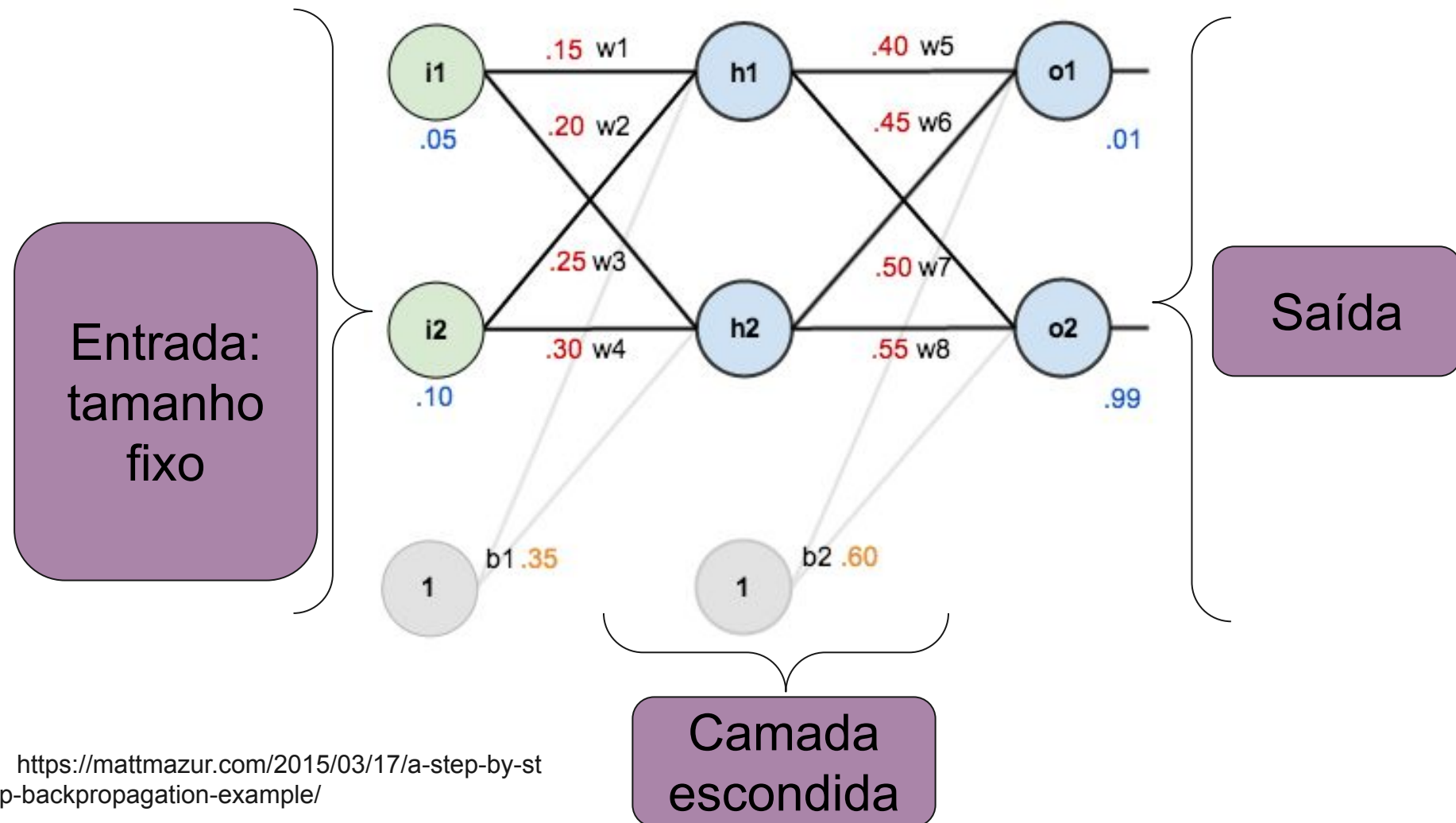
Rede neural feed-forward



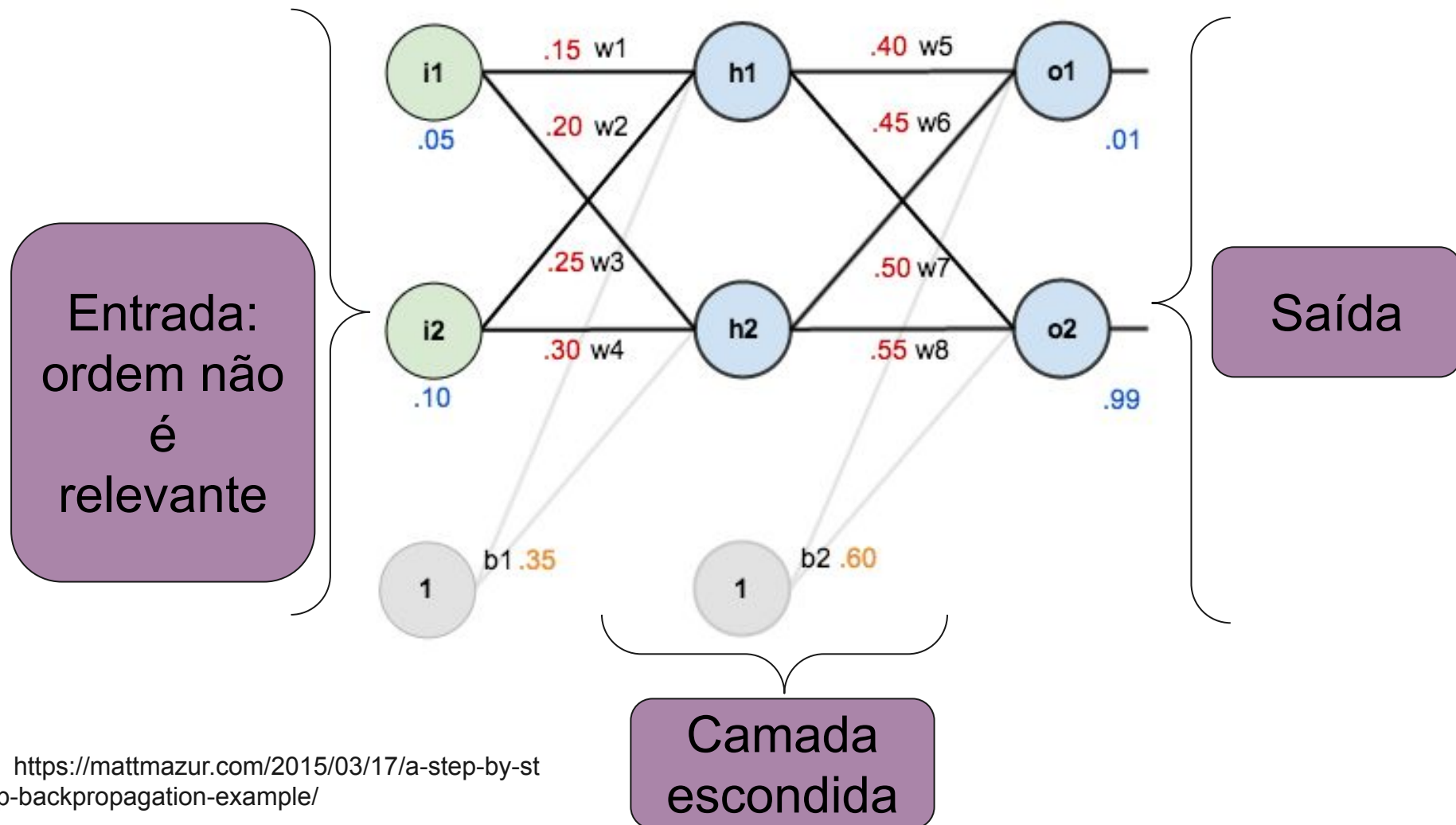
Rede neural feed-forward



Rede neural feed-forward



Rede neural feed-forward



Sequências em NLP

"De casa eu decidi antes que as coisas mais difíceis
ficassem mudar."

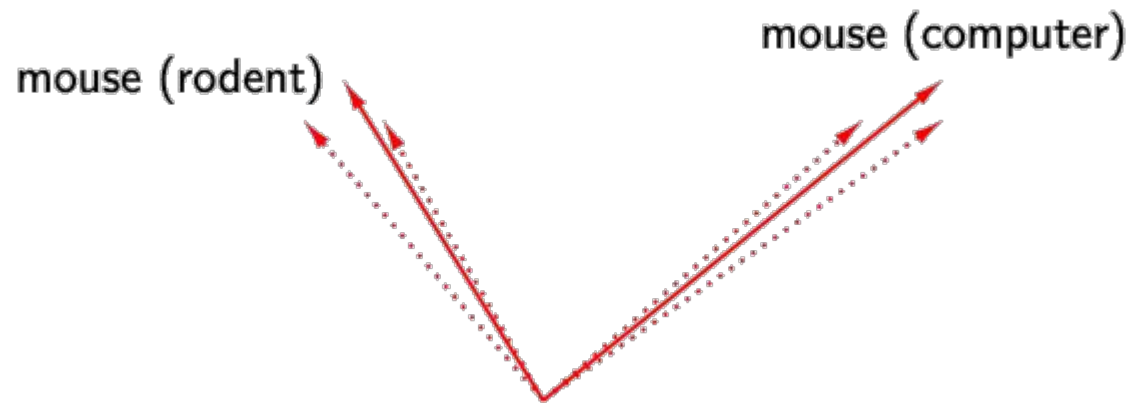
"De eu decidi casa antes que as coisas mais difíceis
ficassem mudar todo dia."

"Casa de antes decidi casa eu as que coisas mais
mudar ficassem difíceis."

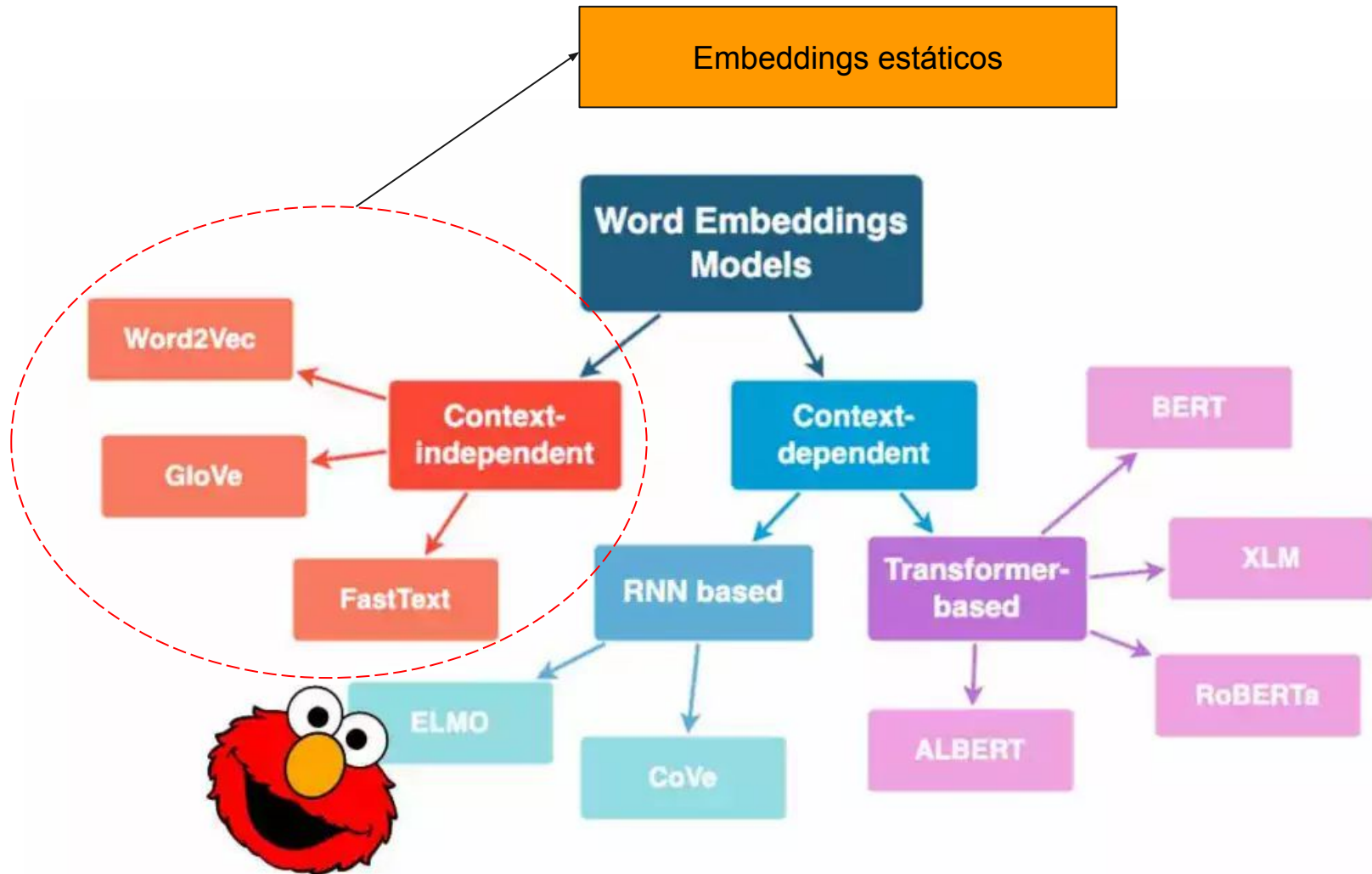
Sequências em NLP

- Diversos problemas precisam considerar as **dependências** entre termos
 - Co-referência
 - Apesar das *suas* obrigações familiares, *Wilma* consegue se dedicar aos estudos.
 - Concordância de número e gênero
 - Lula e FHC *foram* presidentes do Brasil.
 - Coesão de textos
 - A caixa não coube na mala pois ela era muito...
 - Grande:
 - Pequena:

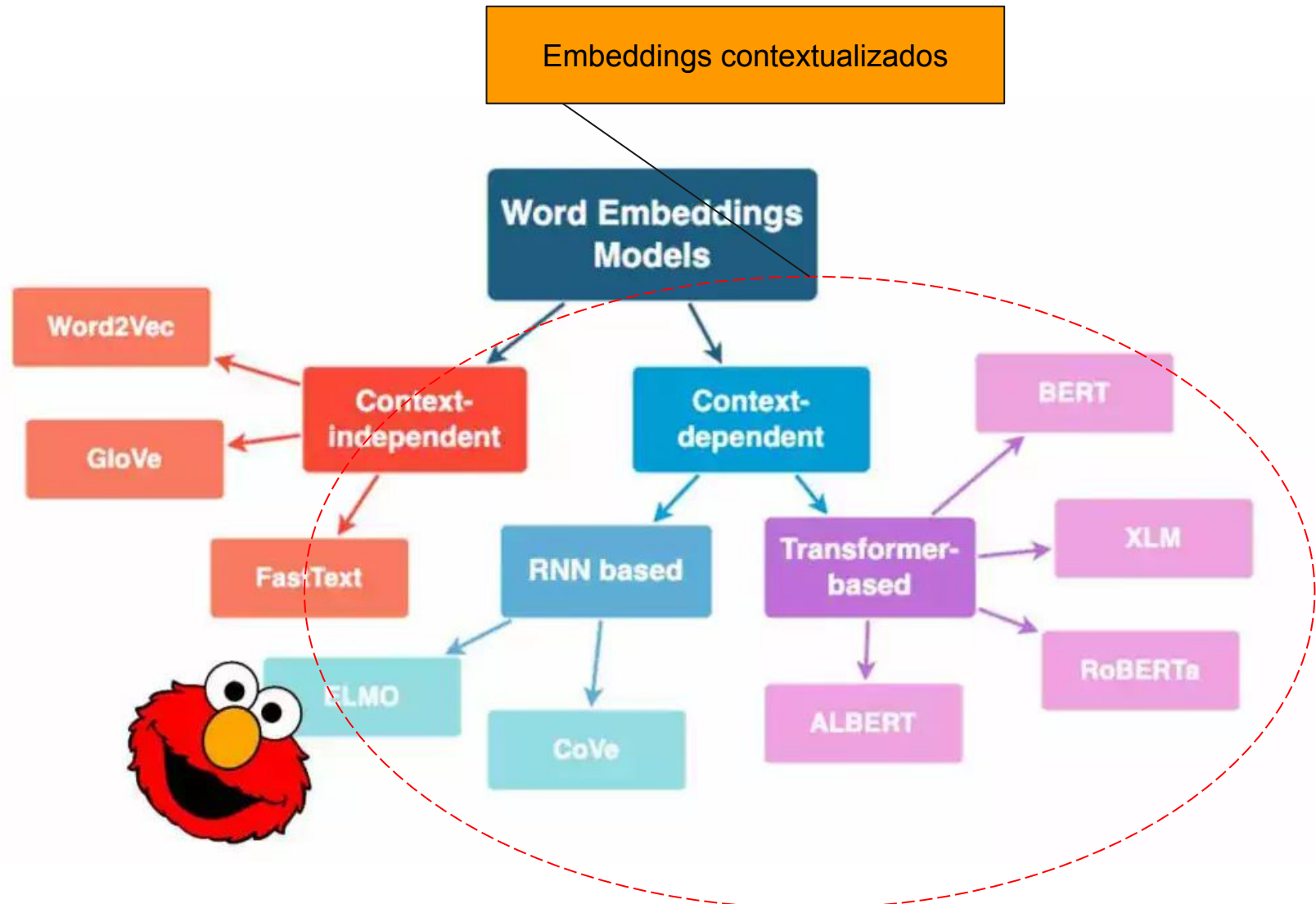
Contexto - de verdade



Contexto - de verdade



Contexto - "de verdade"



Embeddings contextualizados : Elmo (Peters et al., 2018)

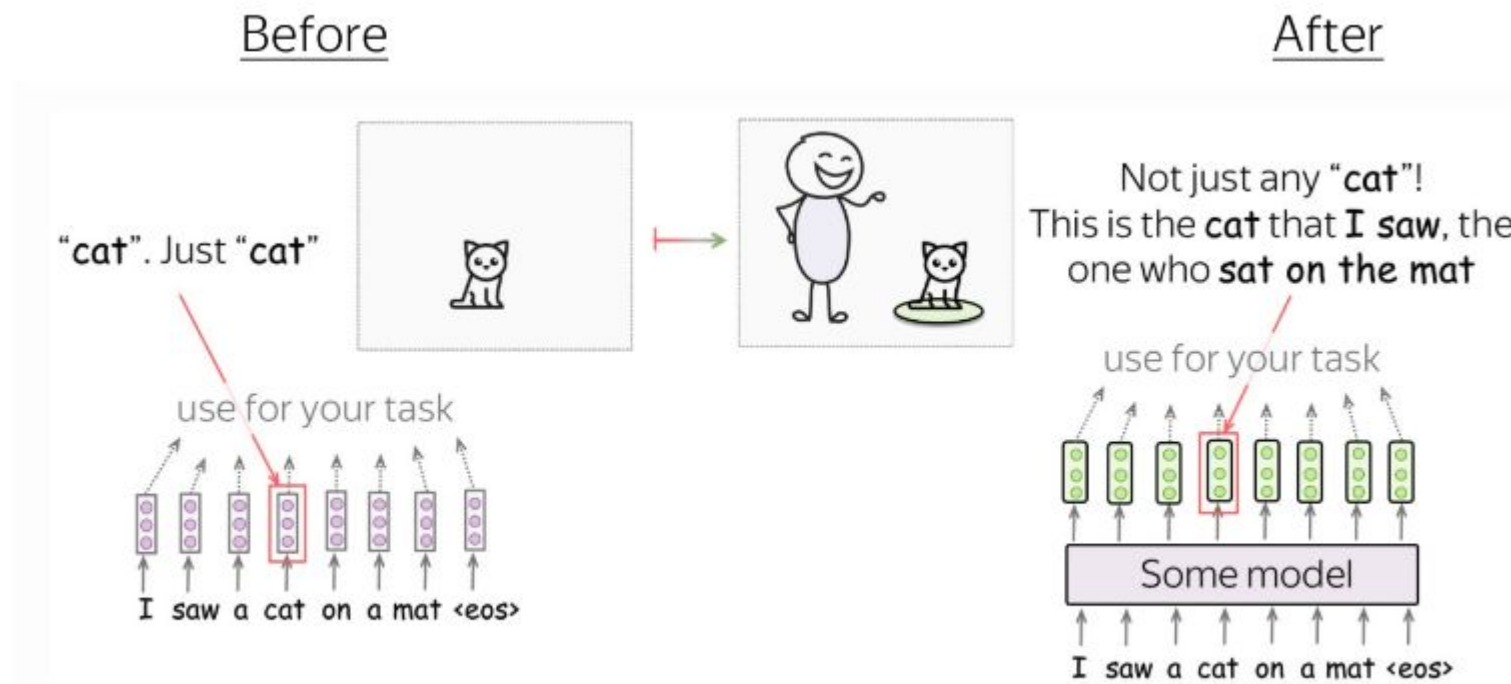


Contextualized word-embeddings can give words different embeddings based on the meaning they carry in the context of the sentence.

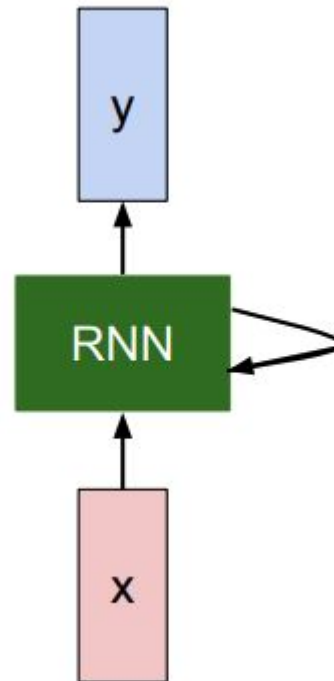
Also, RIP Robin Williams

* <http://jalammar.github.io/illustrated-bert/>

Embeddings contextualizados

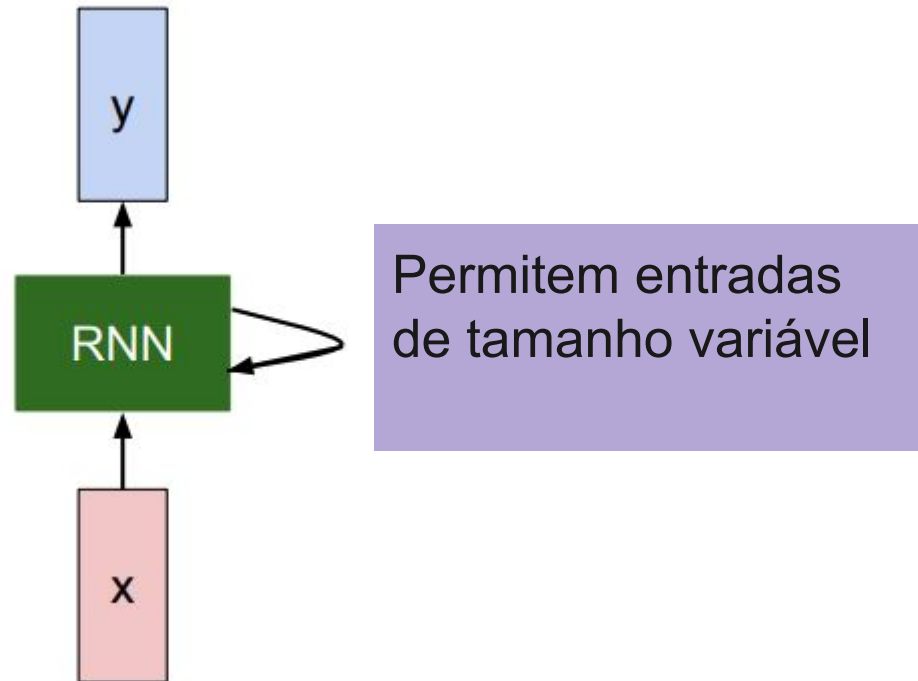


Rede neural recorrente (Elman, 1990)



Estado interno:
atualizado conforme
uma **sequência** é
processada

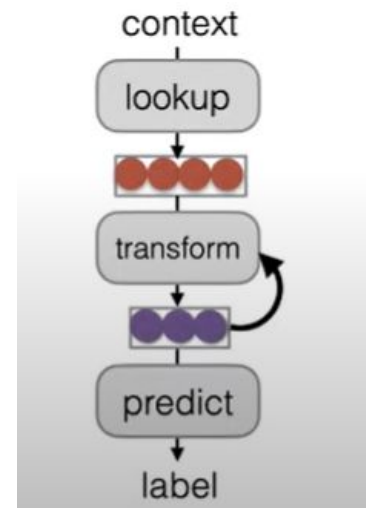
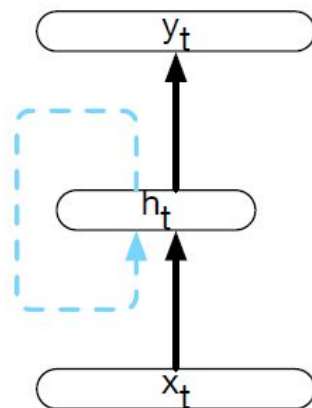
Rede neural recorrente (Elman, 1990)



Redes Neurais Recorrentes (Elman, 1990)

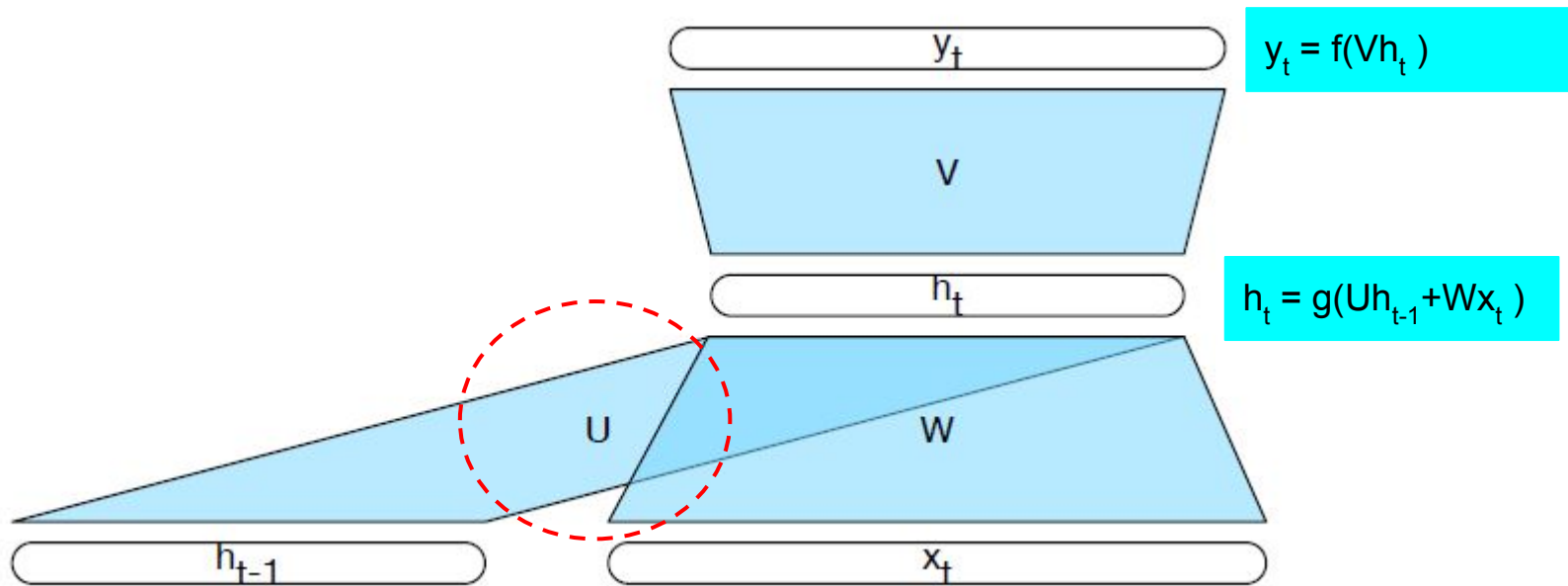
- Contém um ciclo em suas conexões
 - O valor de uma unidade é direta ou indiretamente dependente de uma saída anterior
- Simulam memória
- Permitem entradas de tamanho variável

*Speech and Language Processing

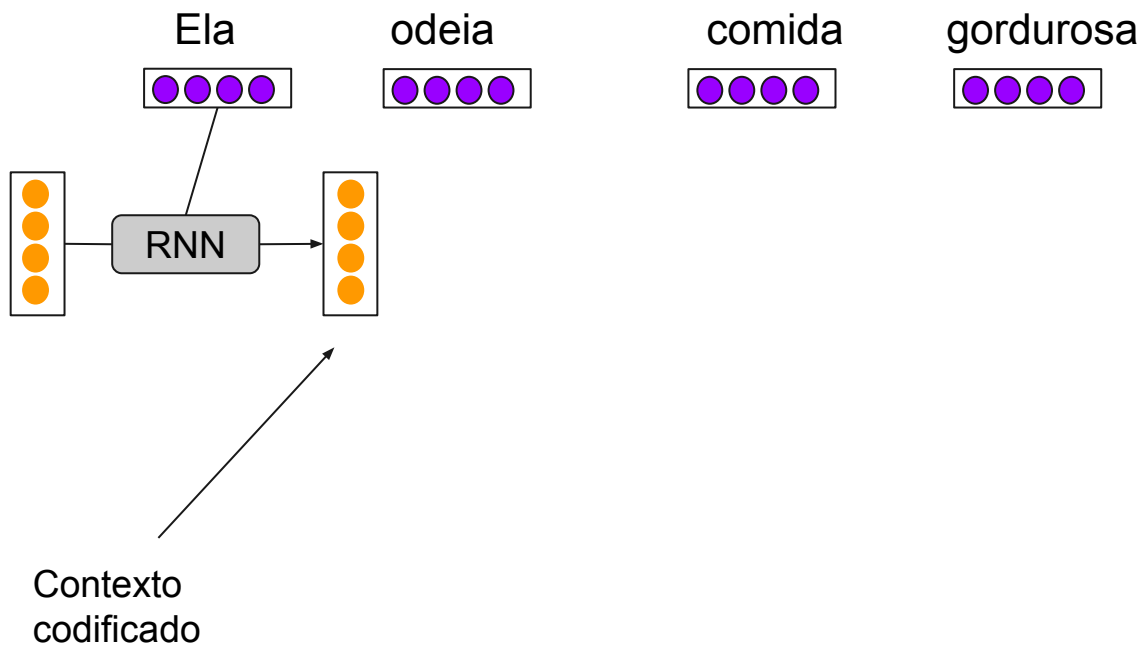


*CMU Neural Nets for NLP

Unfolding

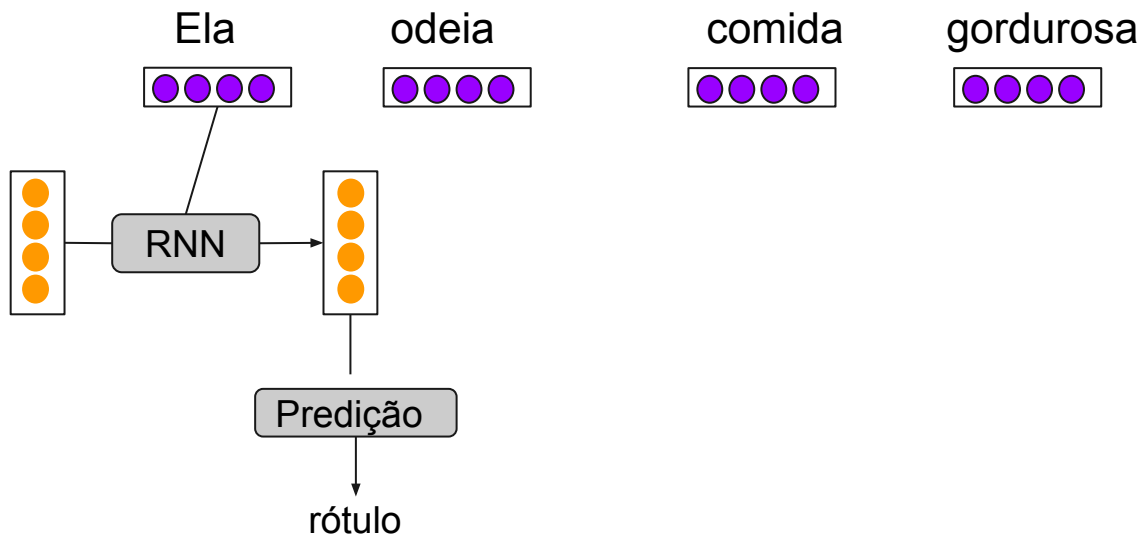


Unfolding

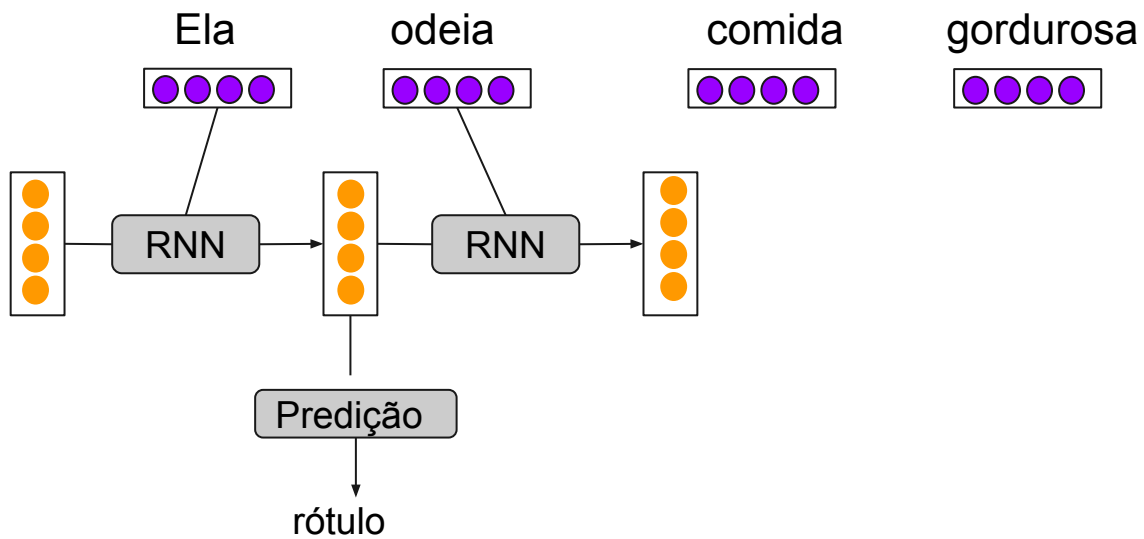


*CMU Neural nets for NLP

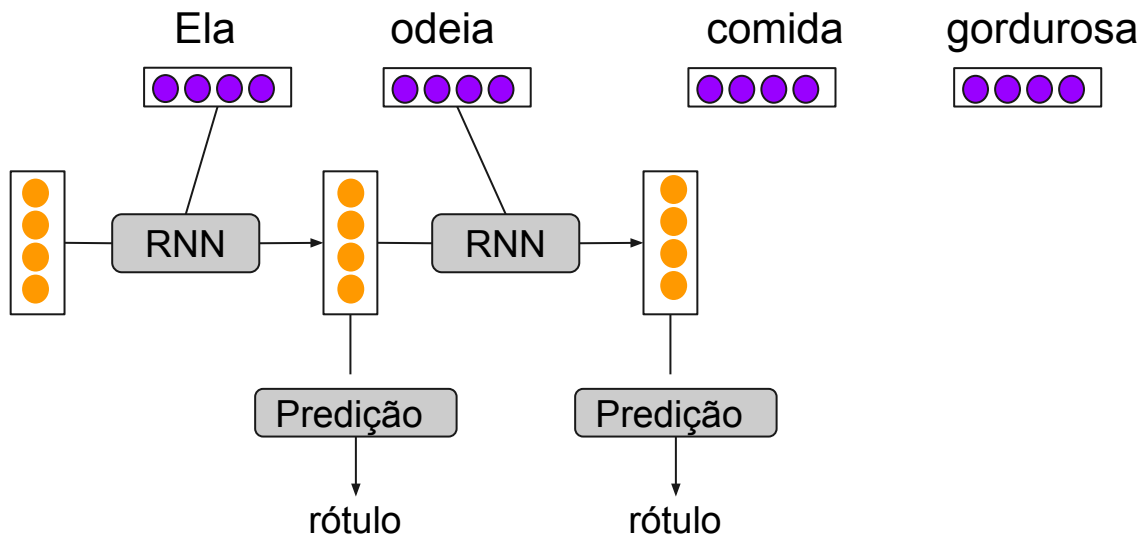
Unfolding



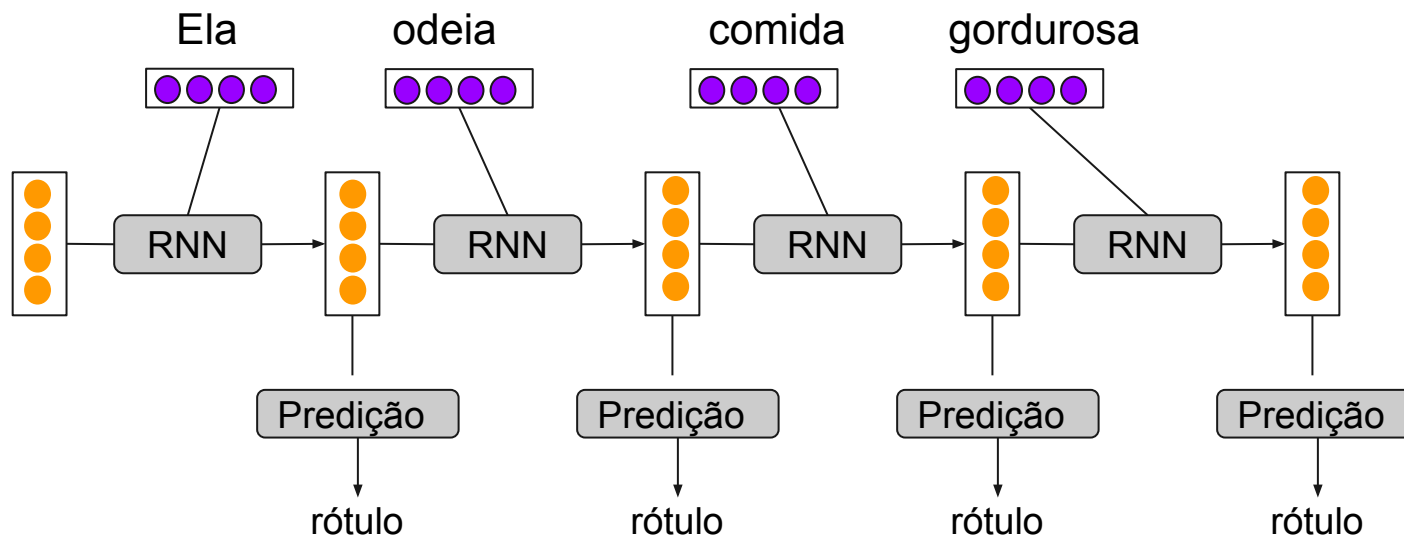
Unfolding



Unfolding

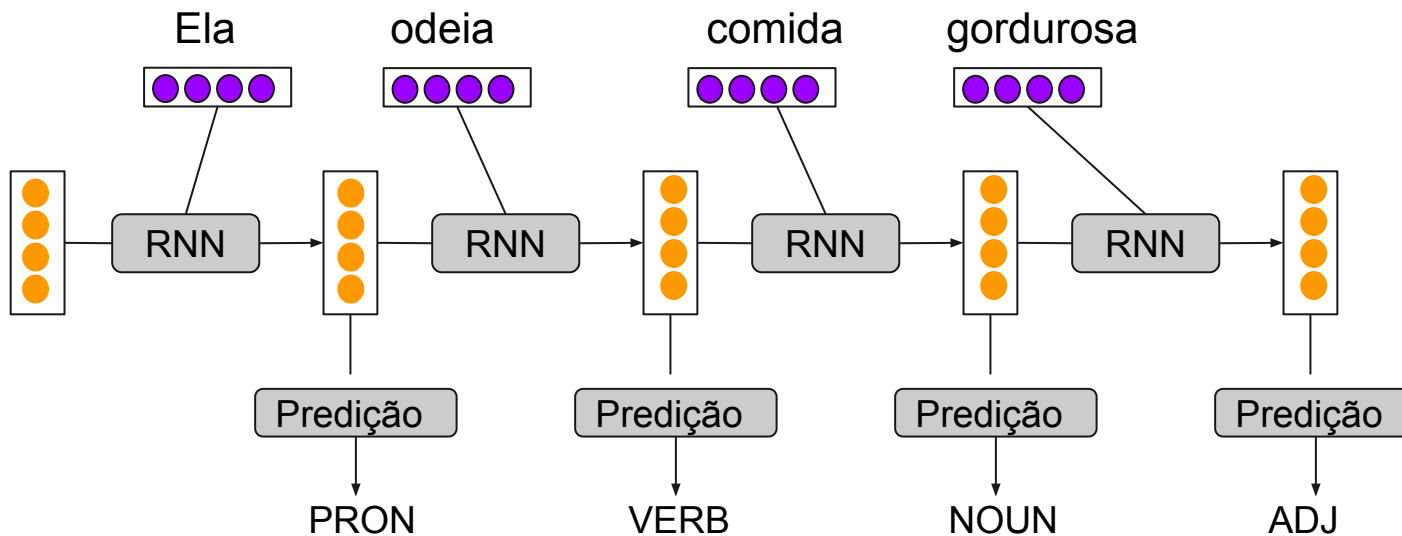


Unfolding



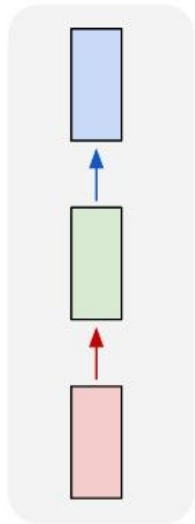
*CMU Neural nets for NLP

Unfolding - Tagging, parsing

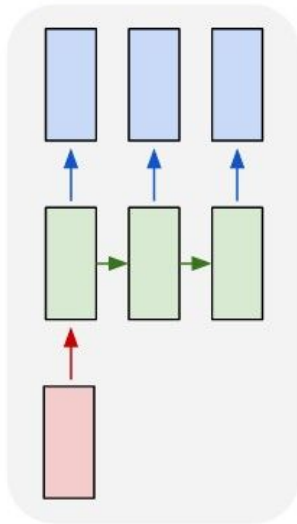


Resumo de tipos

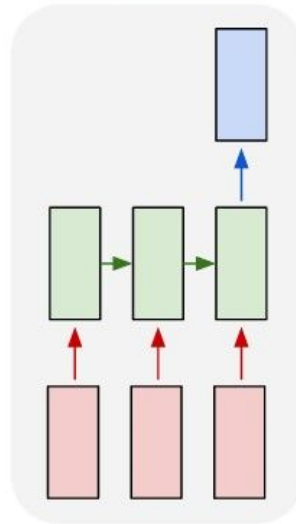
one to one



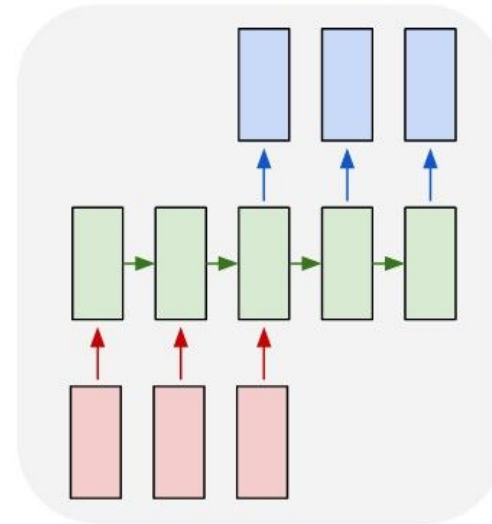
one to many



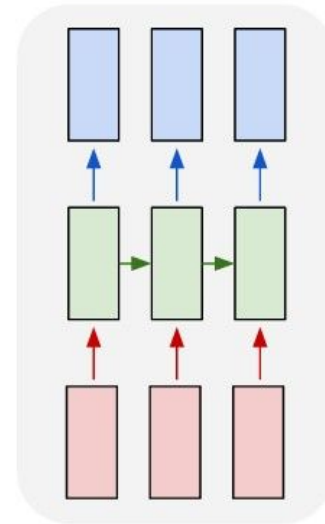
many to one



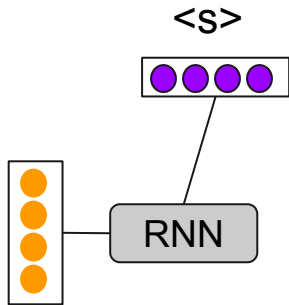
many to many



many to many

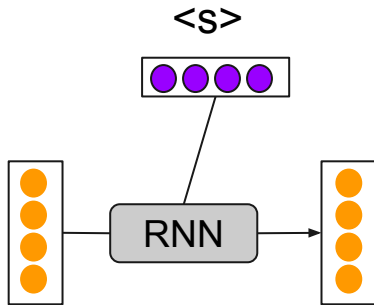


Inferência - geração de texto



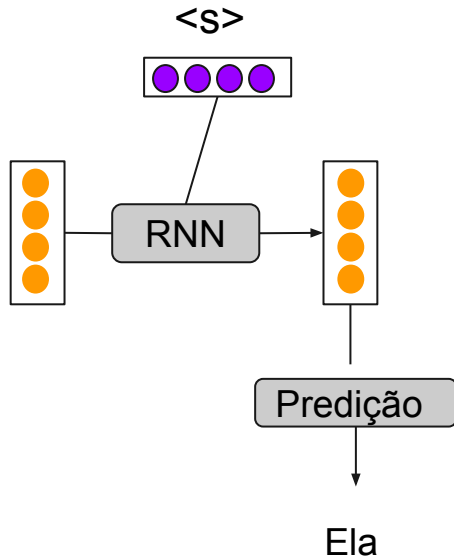
*CMU Neural nets for NLP

Inferência - geração de texto

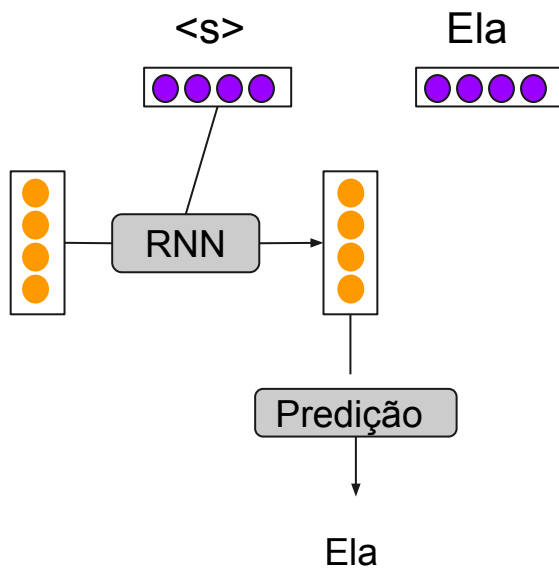


*CMU Neural nets for NLP

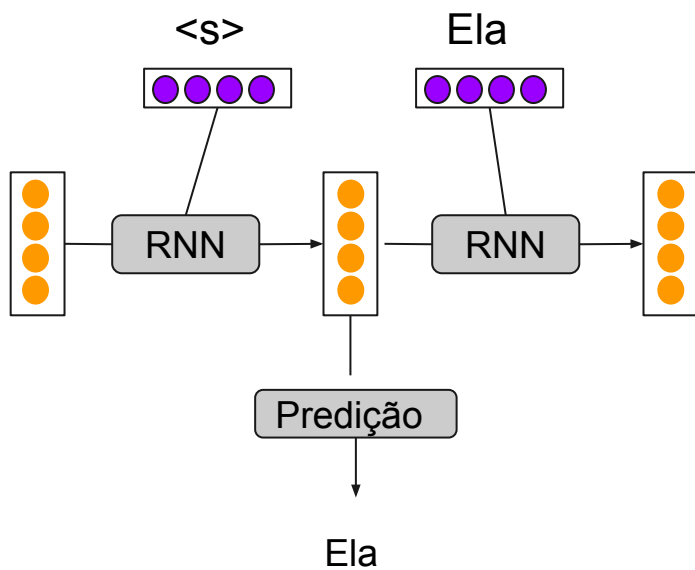
Inferência - geração de texto



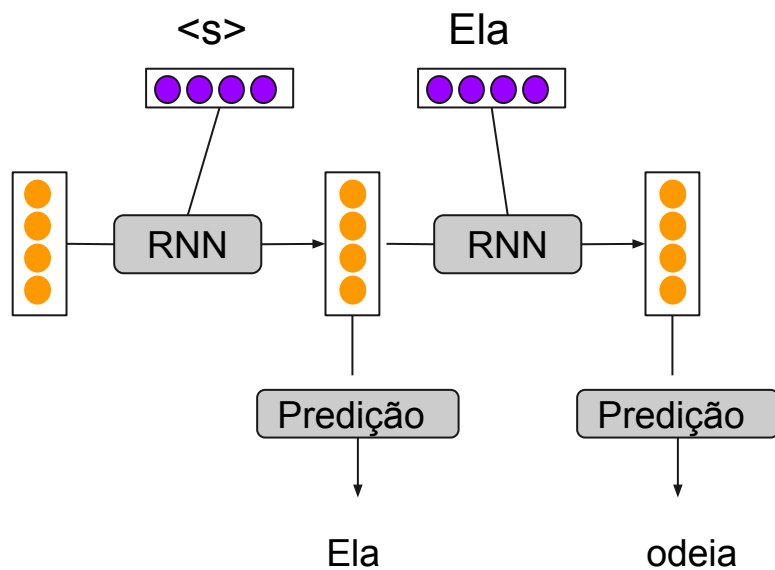
Inferência - geração de texto



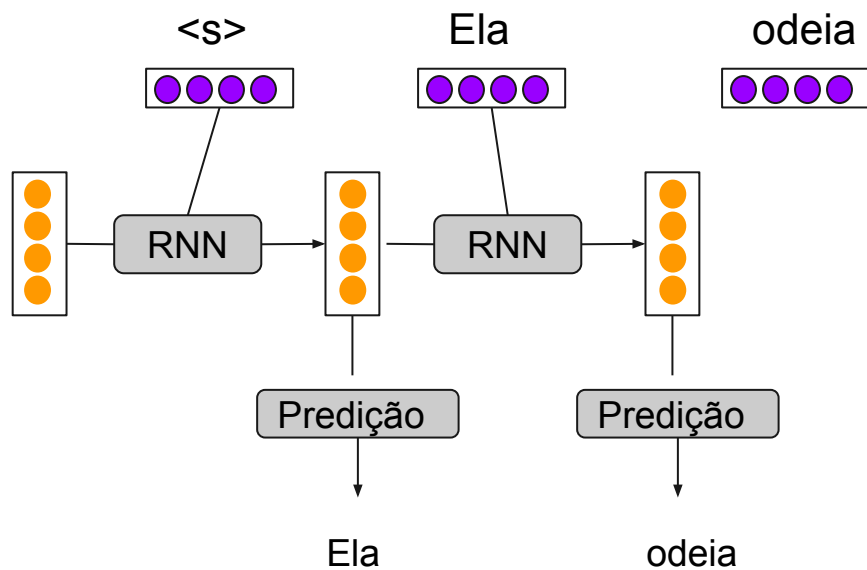
Inferência - geração de texto



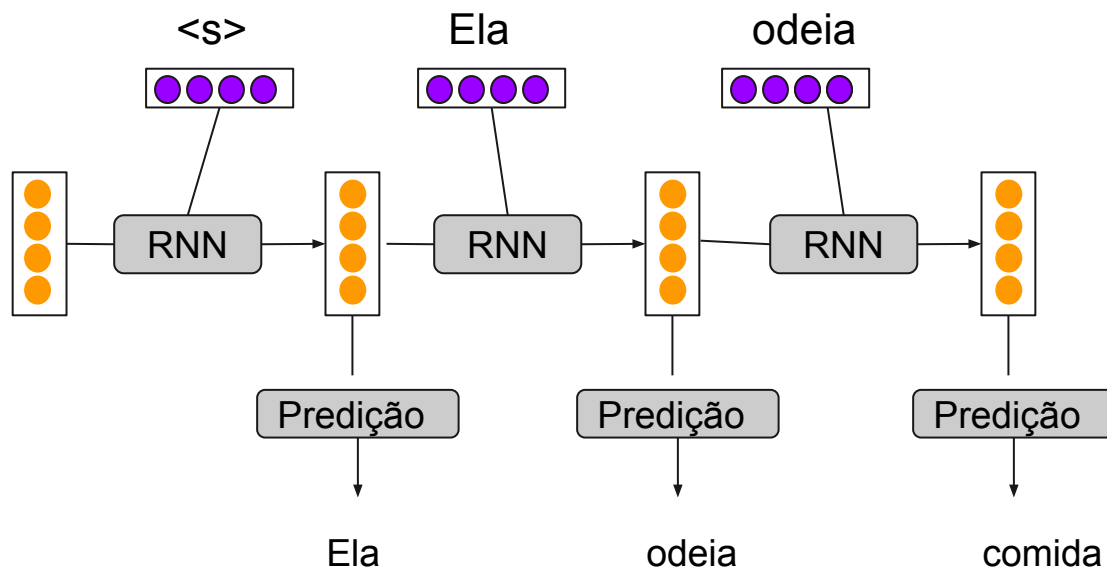
Inferência - geração de texto



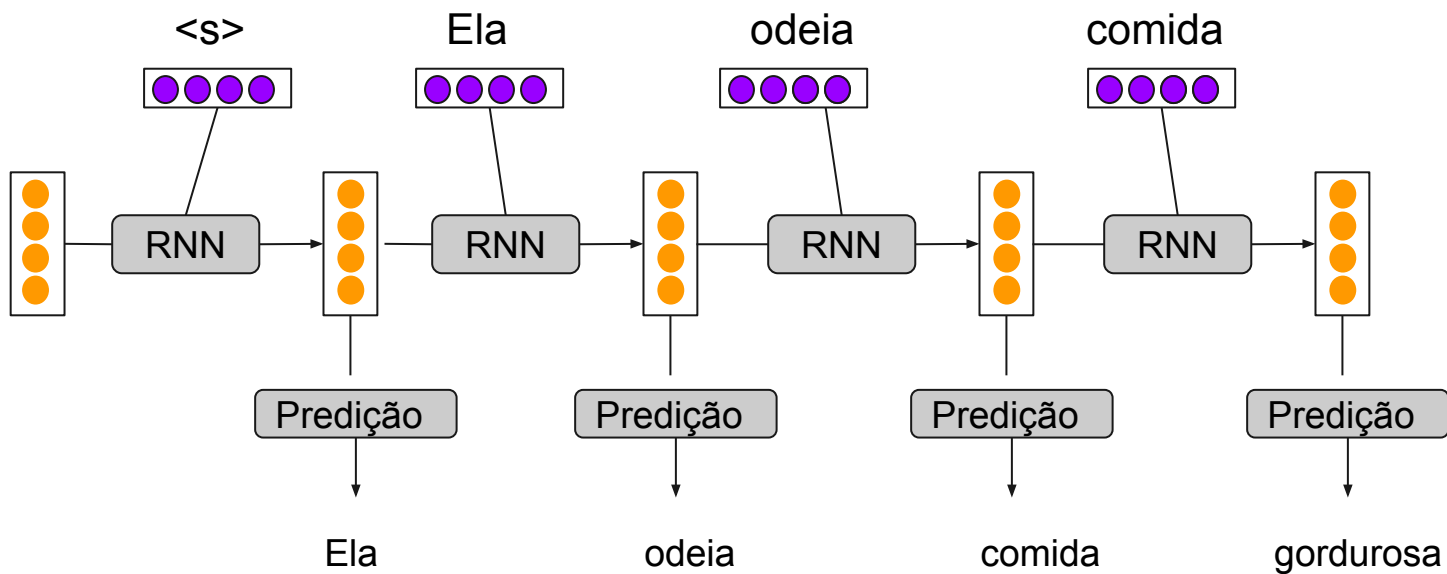
Inferência - geração de texto



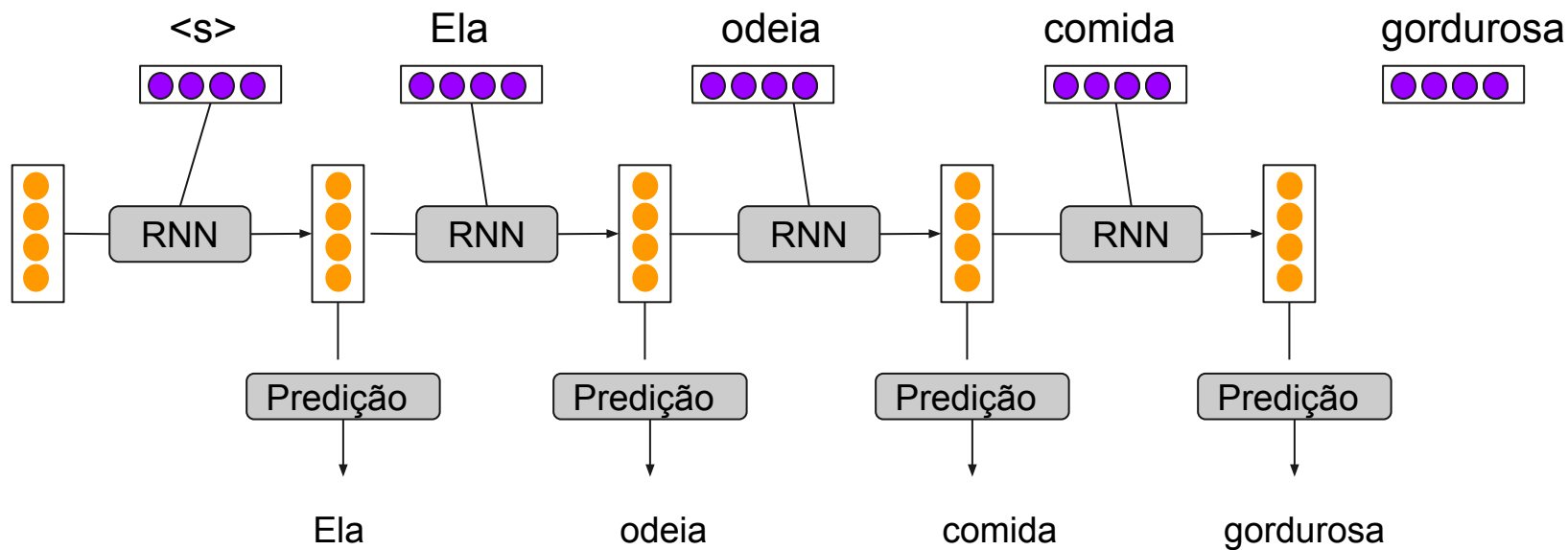
Inferência - geração de texto



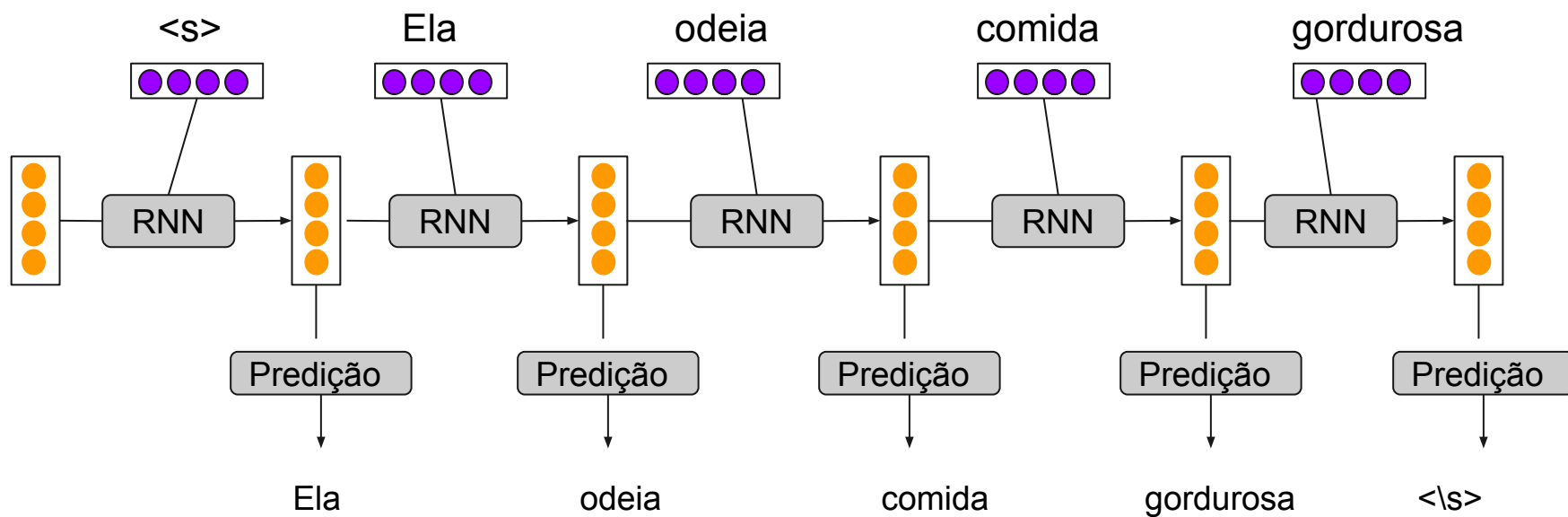
Inferência - geração de texto



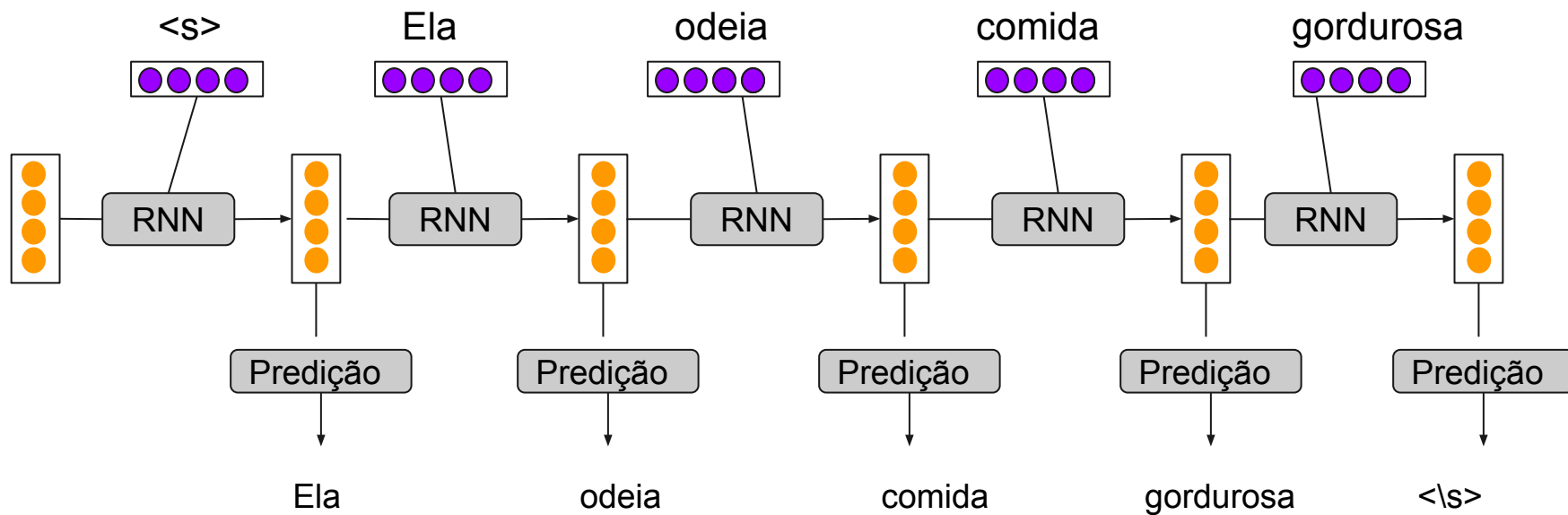
Inferência - geração de texto



Inferência - geração de texto

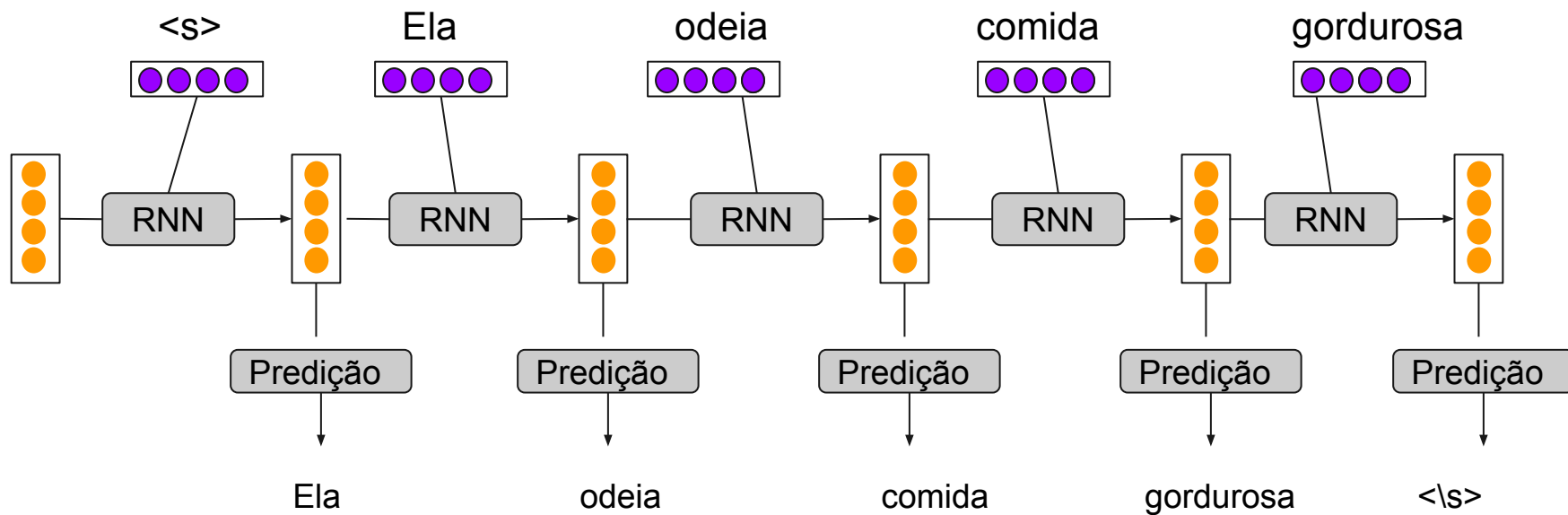


Inferência - geração de texto

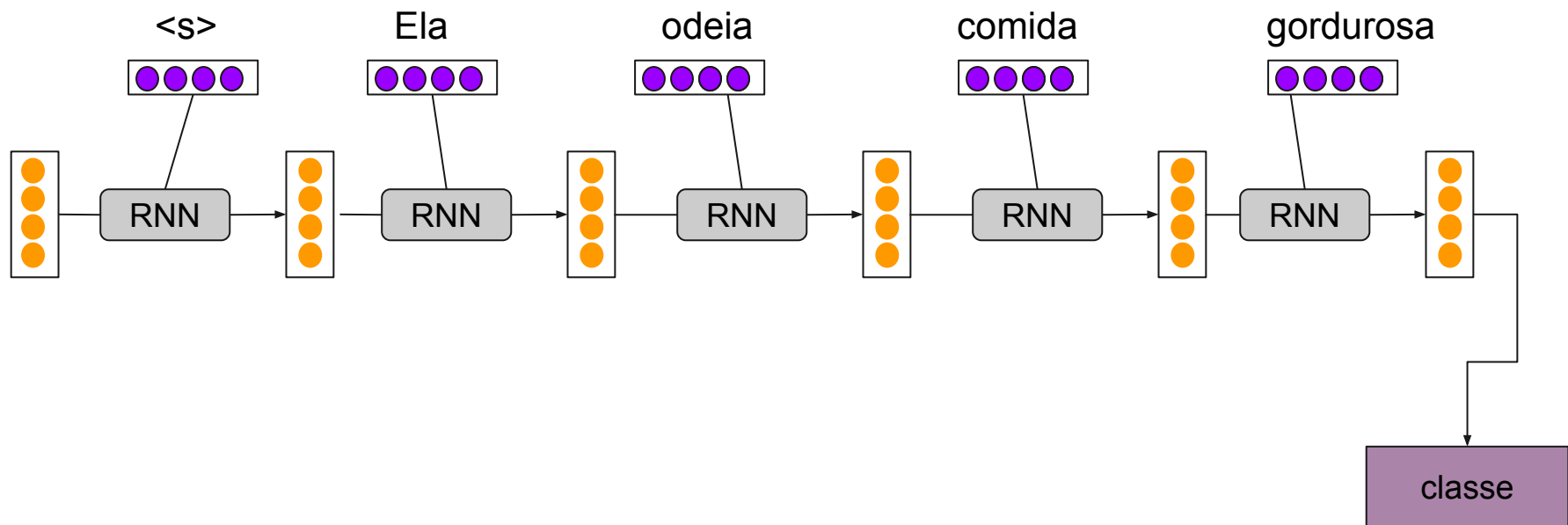


Geração autorregressiva

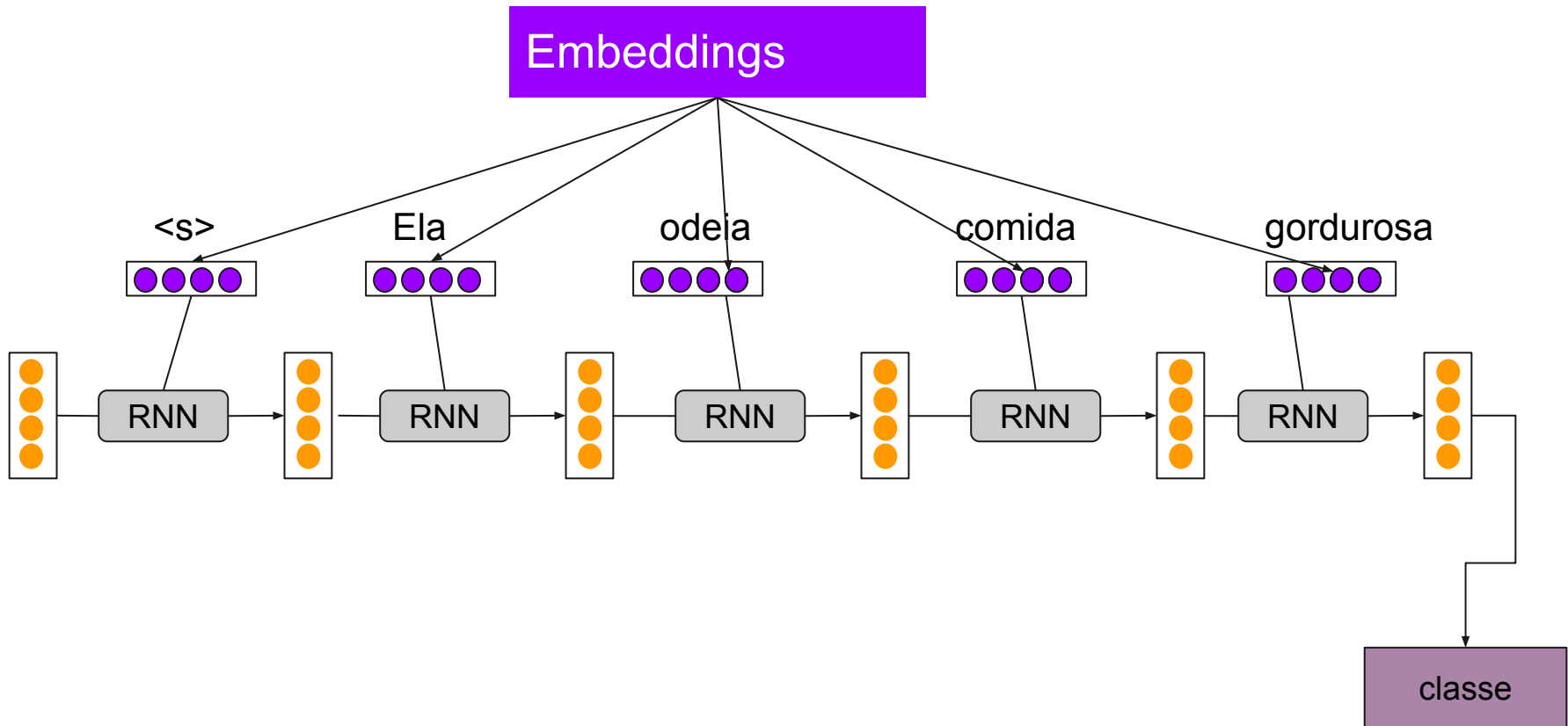
Inferência - classificação de texto?



Inferência - classificação de texto



Inferência - classificação de texto

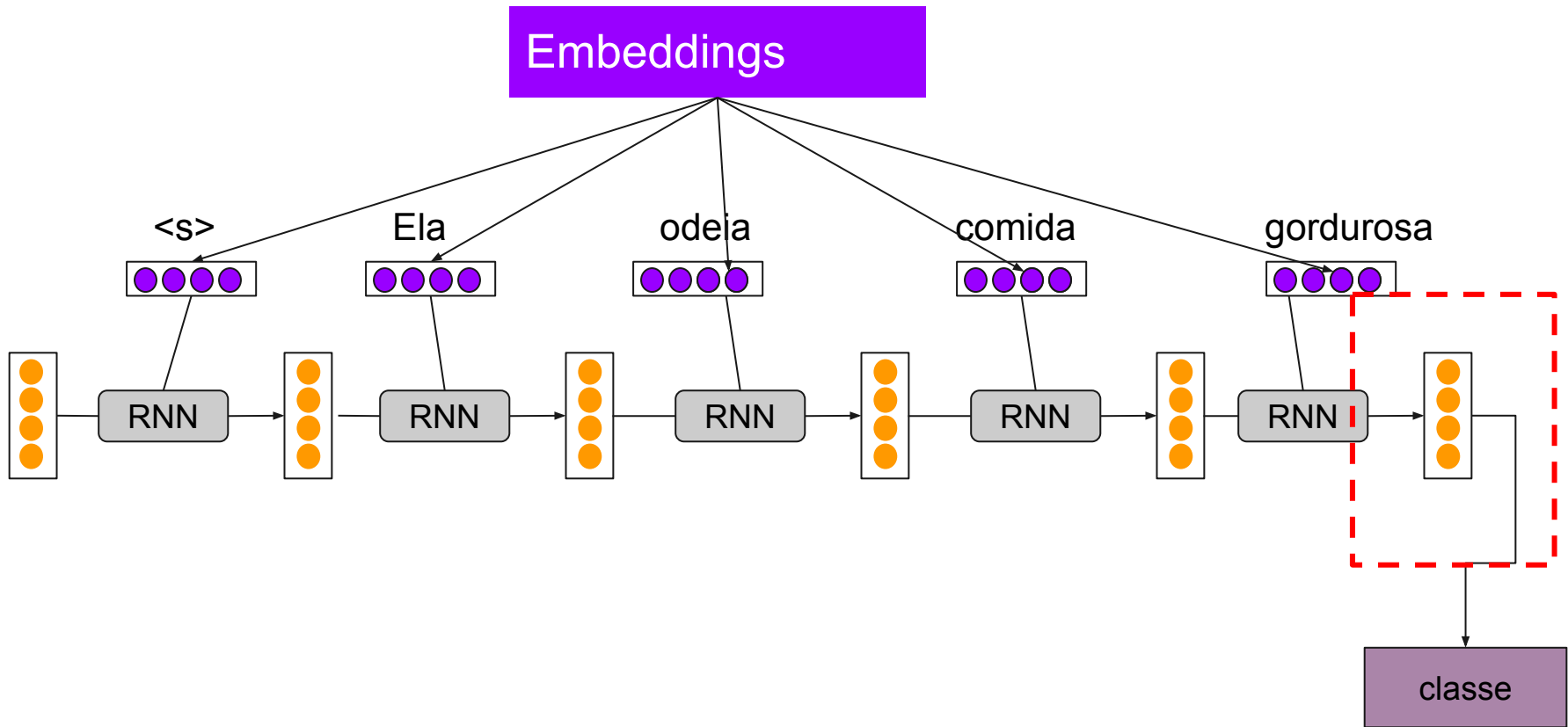


Pergunta

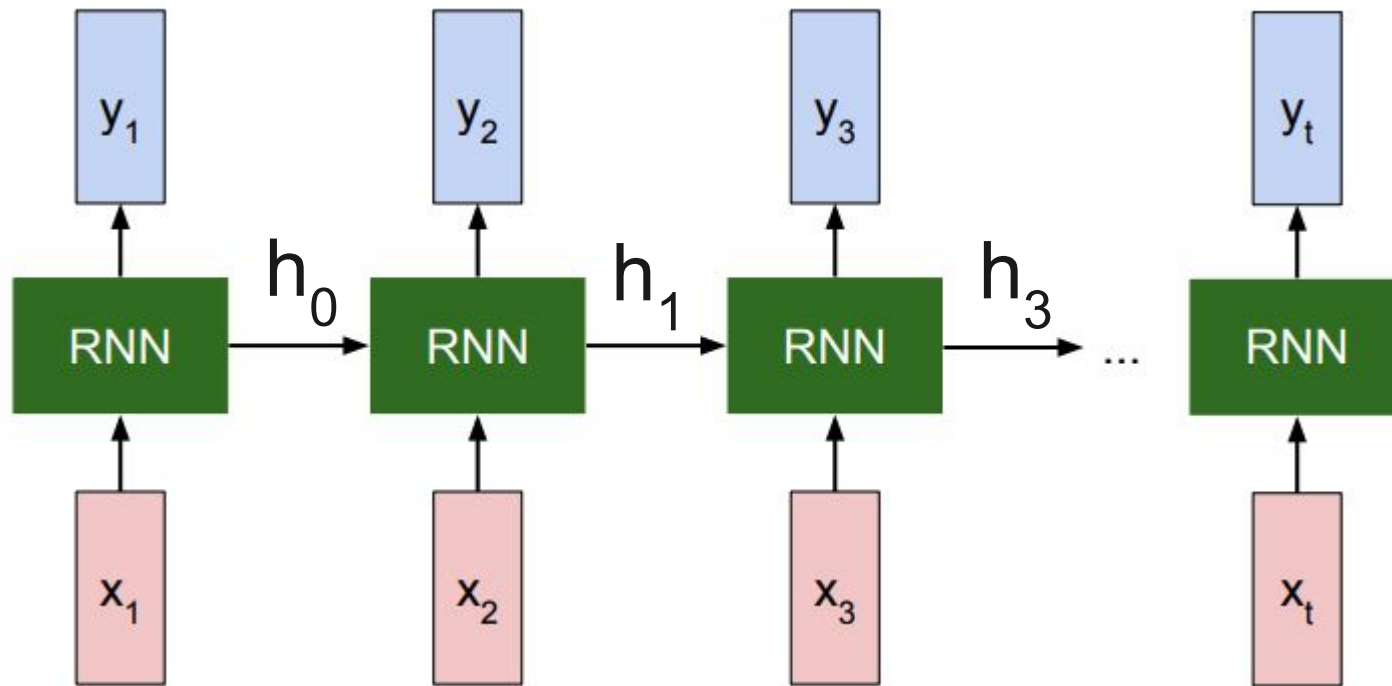
Como representar uma sentença com RNN?

- a) Média dos embeddings
- b) Saída no instante t
- c) Soma dos embeddings
- d) Concatenação dos embeddings

Inferência - classificação de texto



Rede neural recorrente

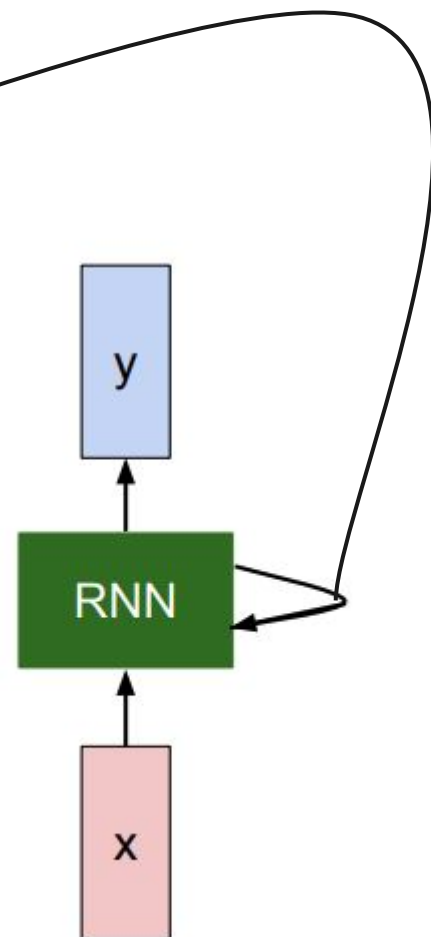


Rede neural recorrente (Elman, 1990)

Fórmula de recorrência

$$h_t = f_w(h_{t-1}, x_t)$$

Entrada no
instante t

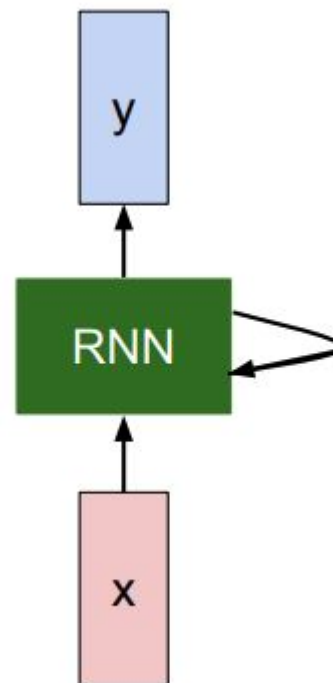


Rede neural recorrente (Elman, 1990)

Fórmula de recorrência

$$h_t = f_w(h_{t-1}, x_t)$$

Estado
anterior

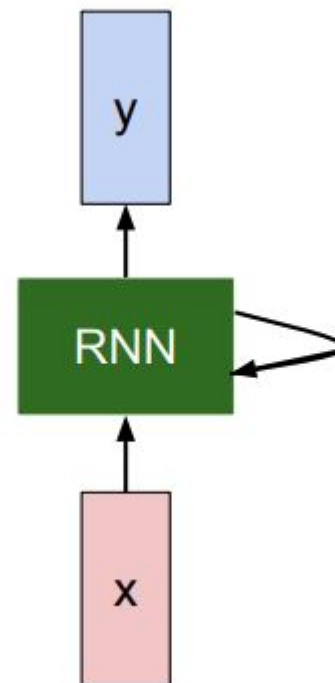


Rede neural recorrente (Elman, 1990)

Fórmula de recorrência

$$h_t = f_w(h_{t-1}, x_t)$$

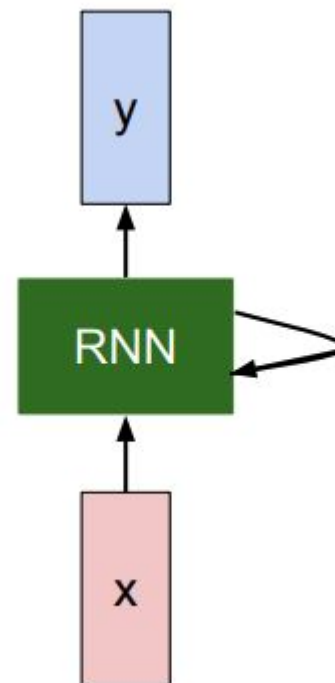
Novo
estado



Rede neural recorrente (Elman, 1990)

$$h_t = f_w(h_{t-1}, x_t)$$

Função parametrizada

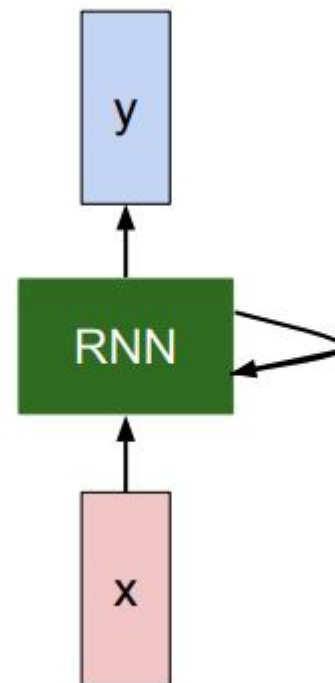


Rede neural recorrente (Elman, 1990)

Fórmula de recorrência

$$h_t = f_w(h_{t-1}, x_t)$$

os mesmos parâmetros para
todos os estados

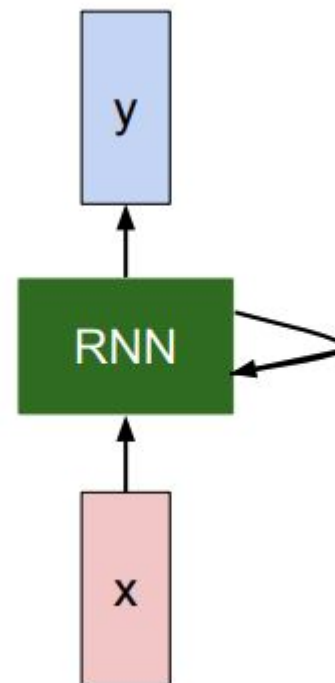


Rede neural recorrente (Elman, 1990)

Fórmula da saída

$$y_t = f_{w'}(h_t)$$

Estado novo

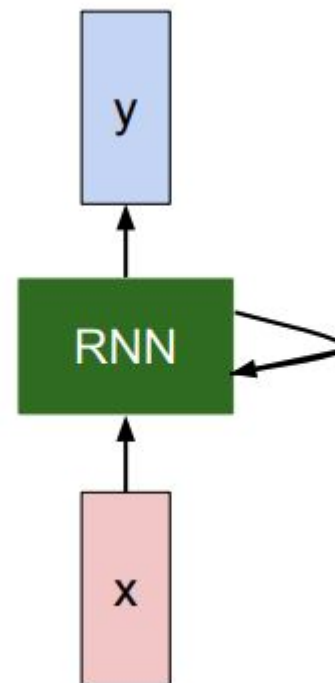


Rede neural recorrente (Elman, 1990)

Fórmula da saída

$$y_t = f_{w'}(h_t)$$

Saída

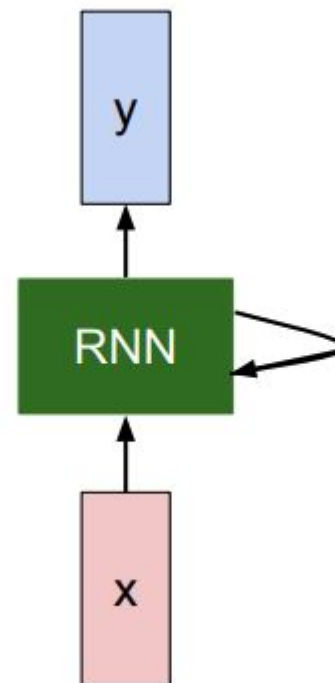


Rede neural recorrente (Elman, 1990)

Fórmula da saída

$$y_t = f_{w'}(h_t)$$

Outra função
parametrizada

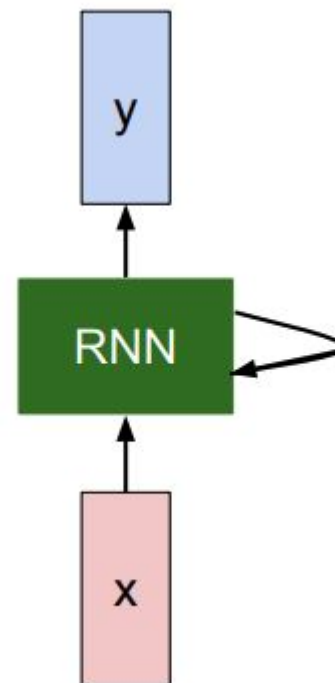


Rede neural recorrente (Elman, 1990)

Fórmula da saída

$$y_t = f_{w'}(h_t)$$

Conjunto diferente
de parâmetros



Rede neural recorrente (Elman, 1990)

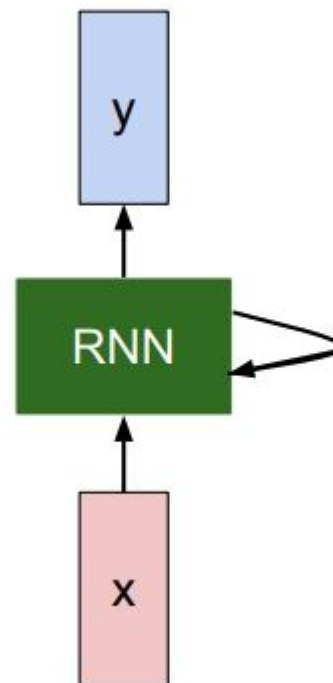
Vanilla RNN

$$h_t = f_W(h_{t-1}, x_t)$$



$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = W_{hy}h_t$$



Rede neural recorrente (Elman, 1990)

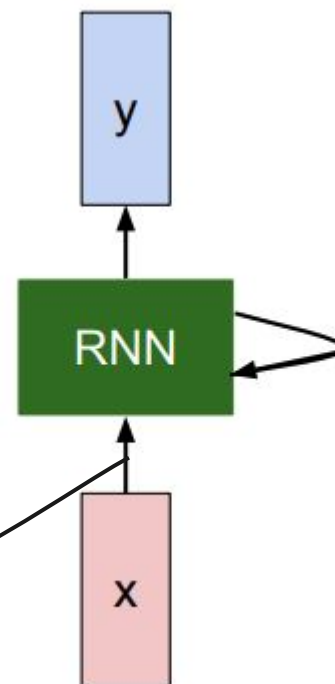
Vanilla RNN

$$h_t = f_W(h_{t-1}, x_t)$$



$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = W_{hy}h_t$$



Rede neural recorrente (Elman, 1990)

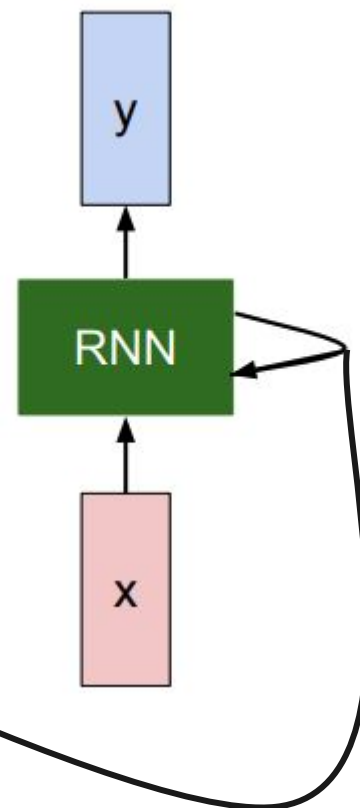
Vanilla RNN

$$h_t = f_W(h_{t-1}, x_t)$$



$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = W_{hy}h_t$$



Rede neural recorrente (Elman, 1990)

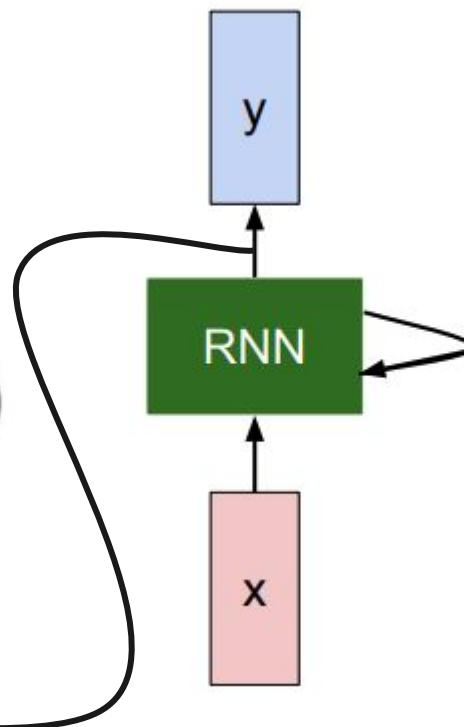
Vanilla RNN

$$h_t = f_W(h_{t-1}, x_t)$$

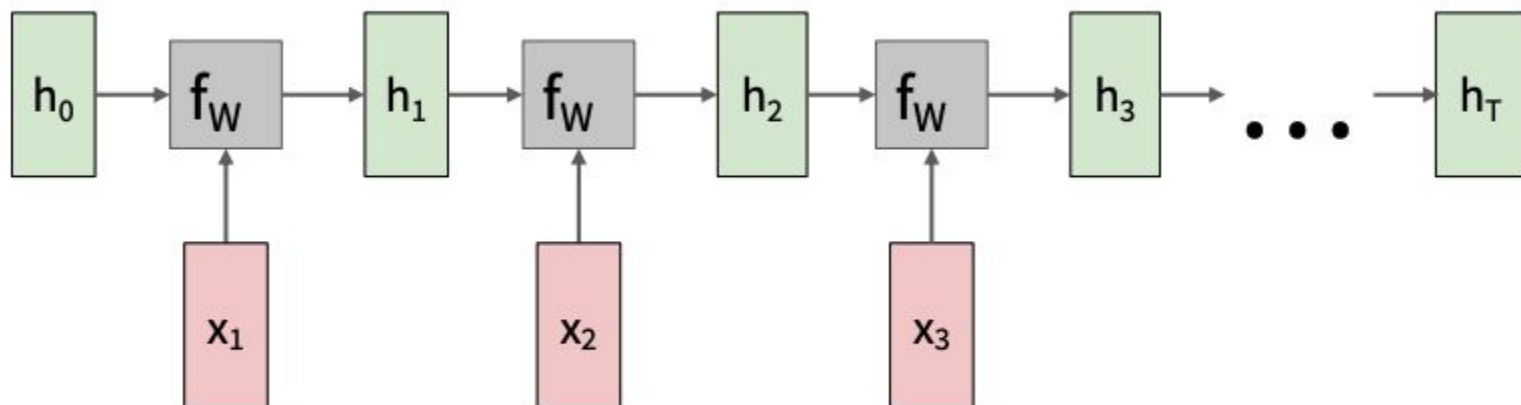


$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

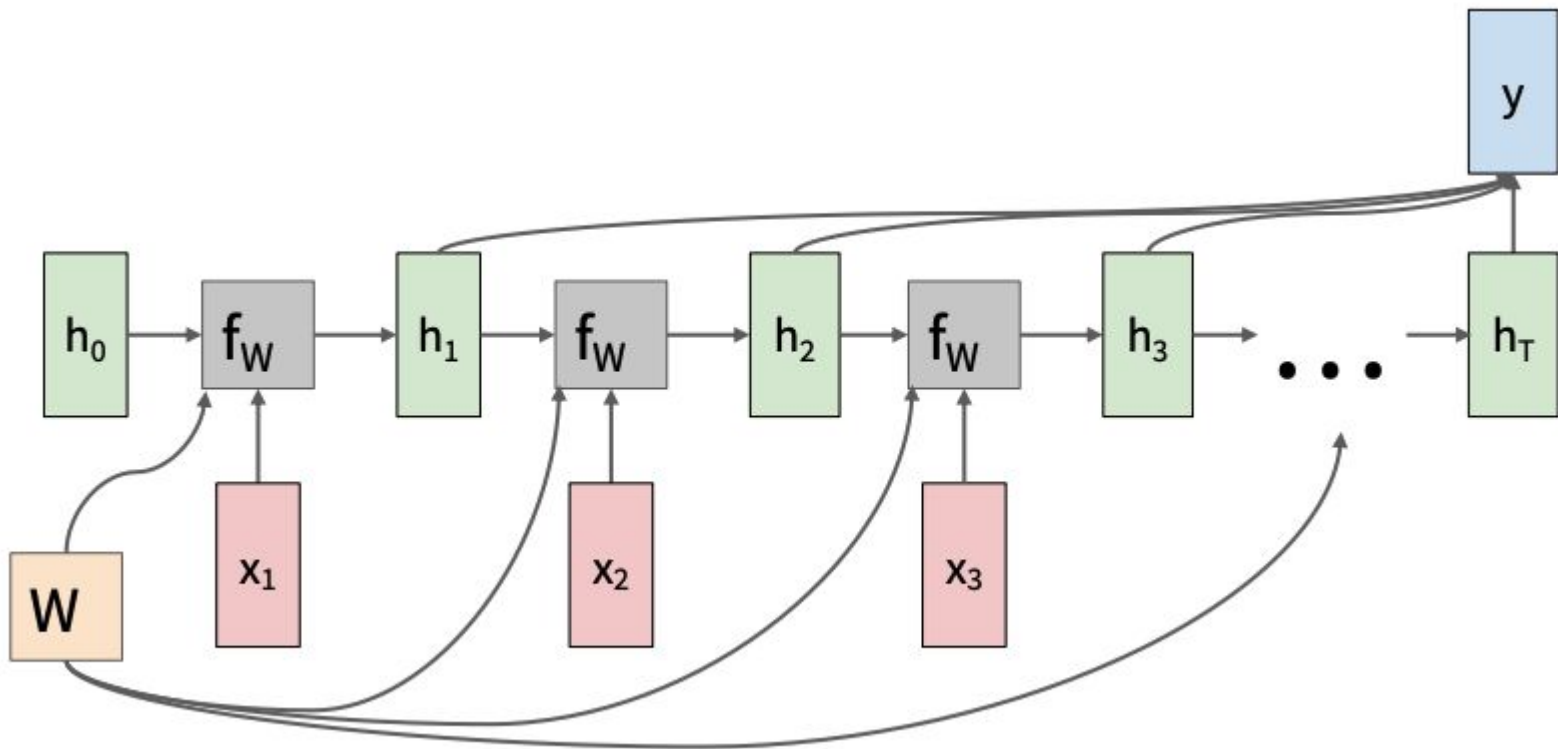
$$y_t = W_{hy}h_t$$



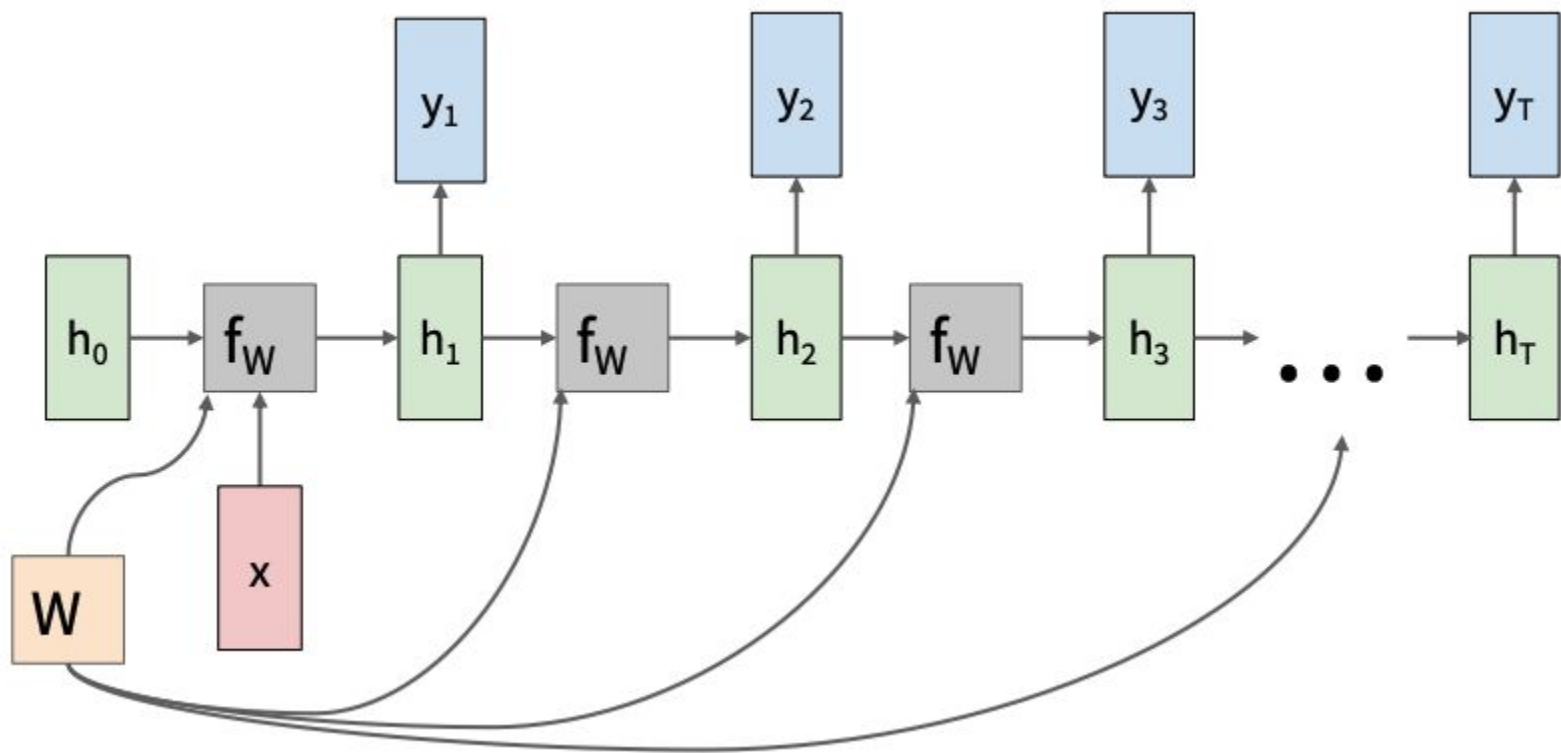
Grafo computational



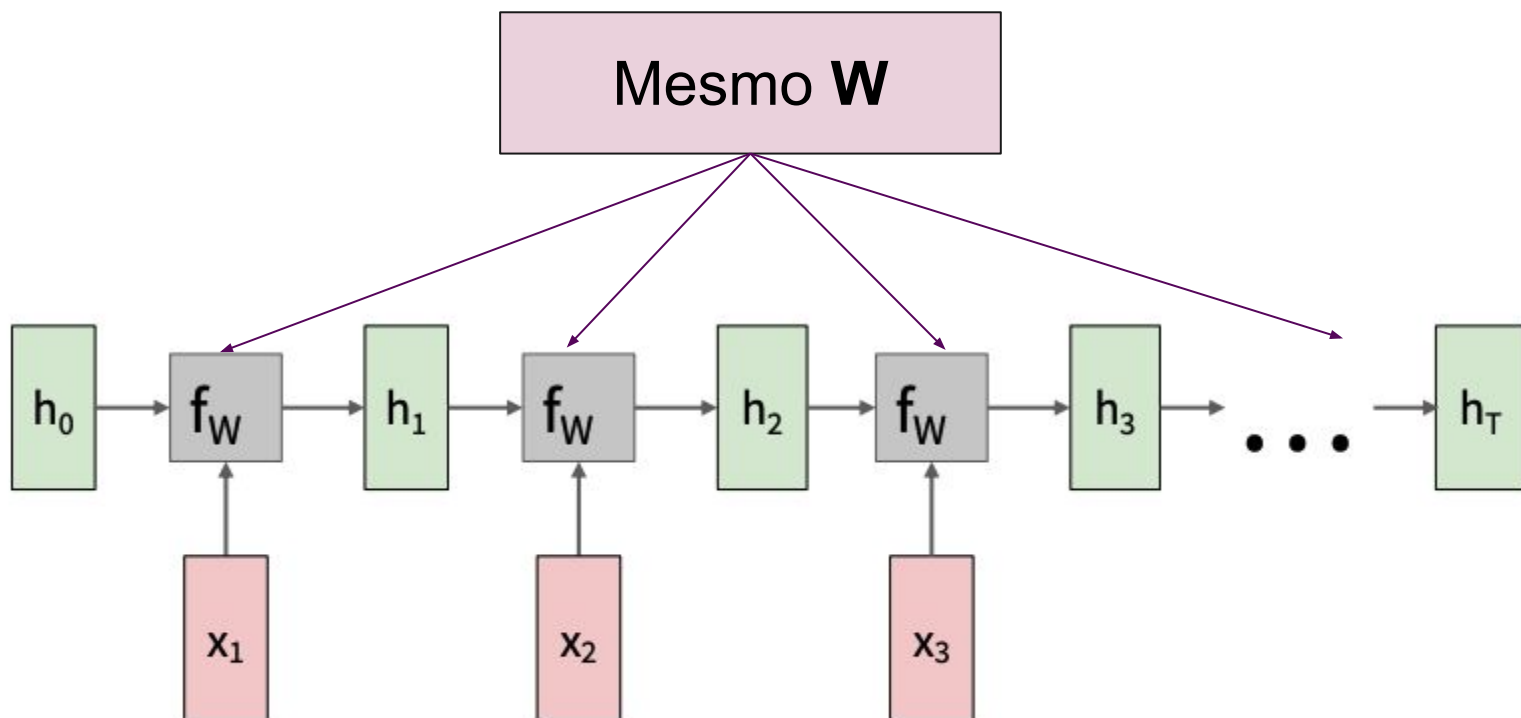
Muitos para um



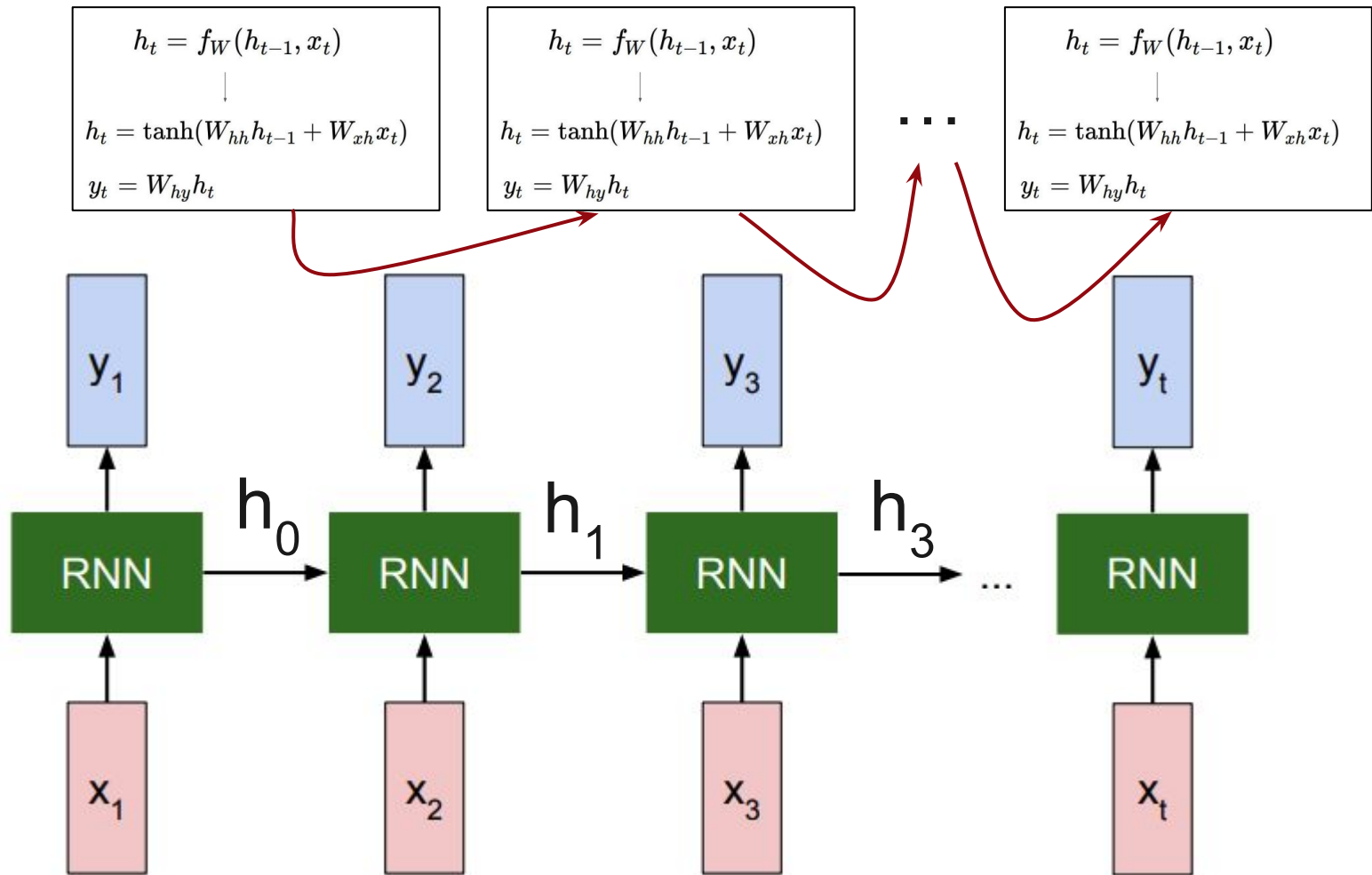
Um para muitos



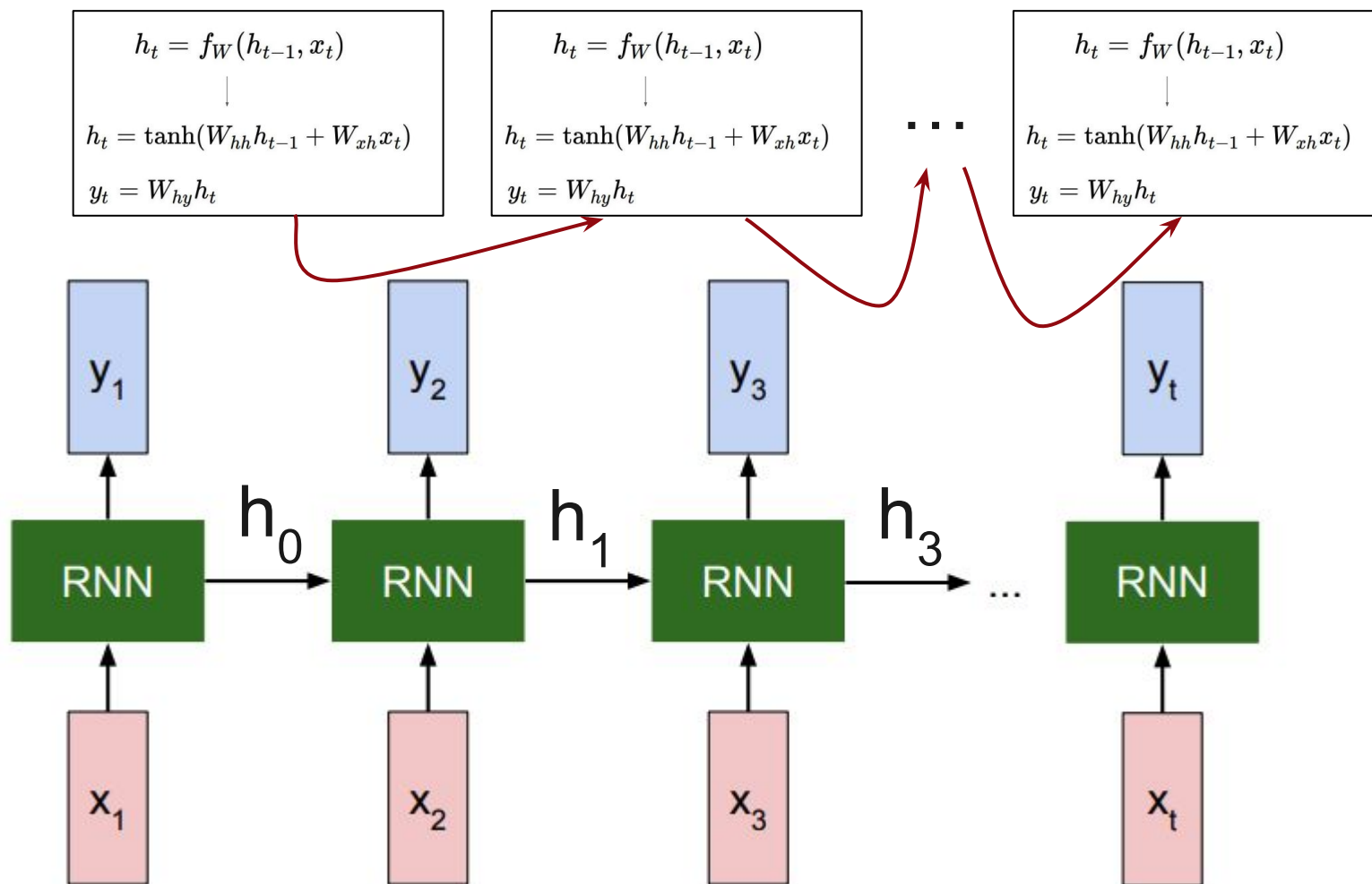
Grafo computational



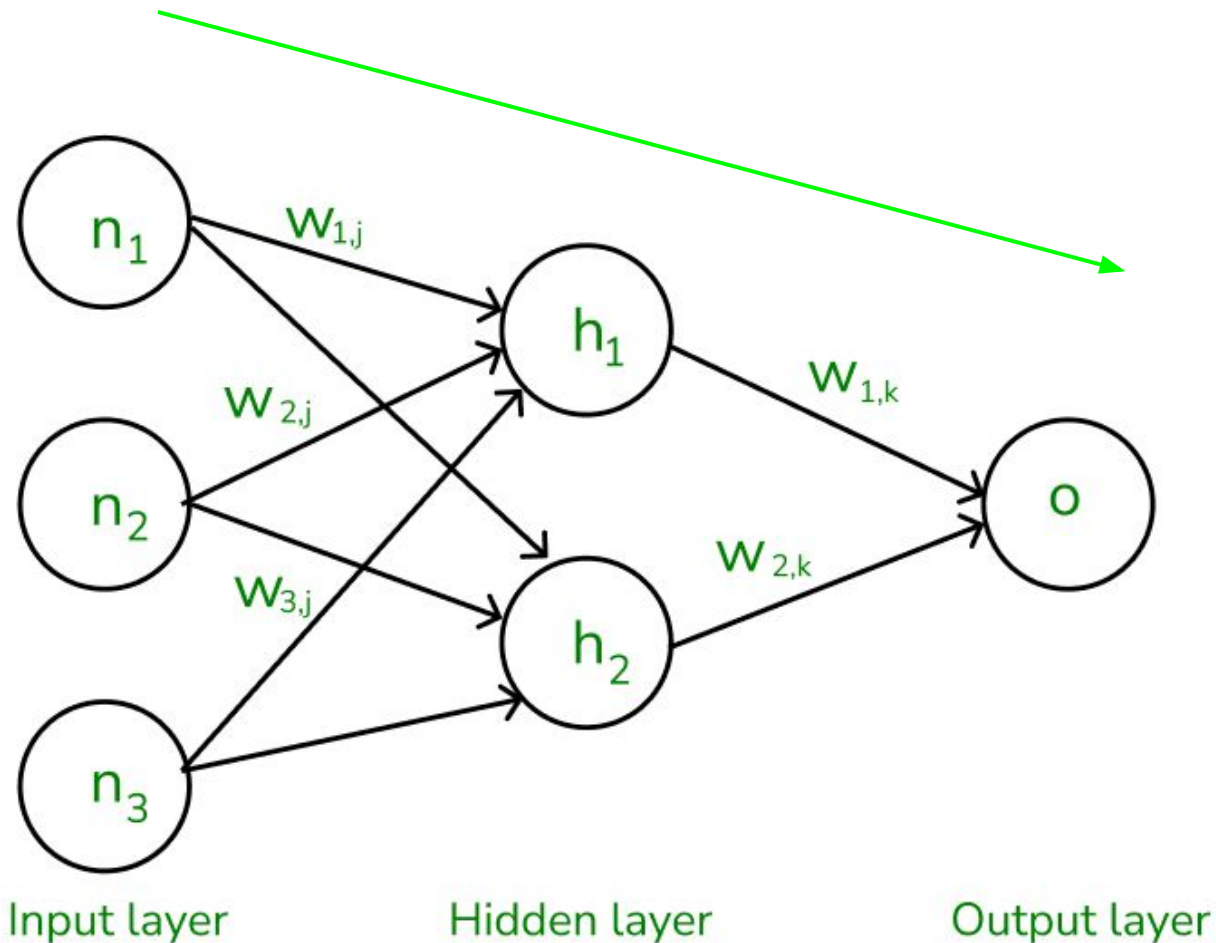
Rede neural recorrente



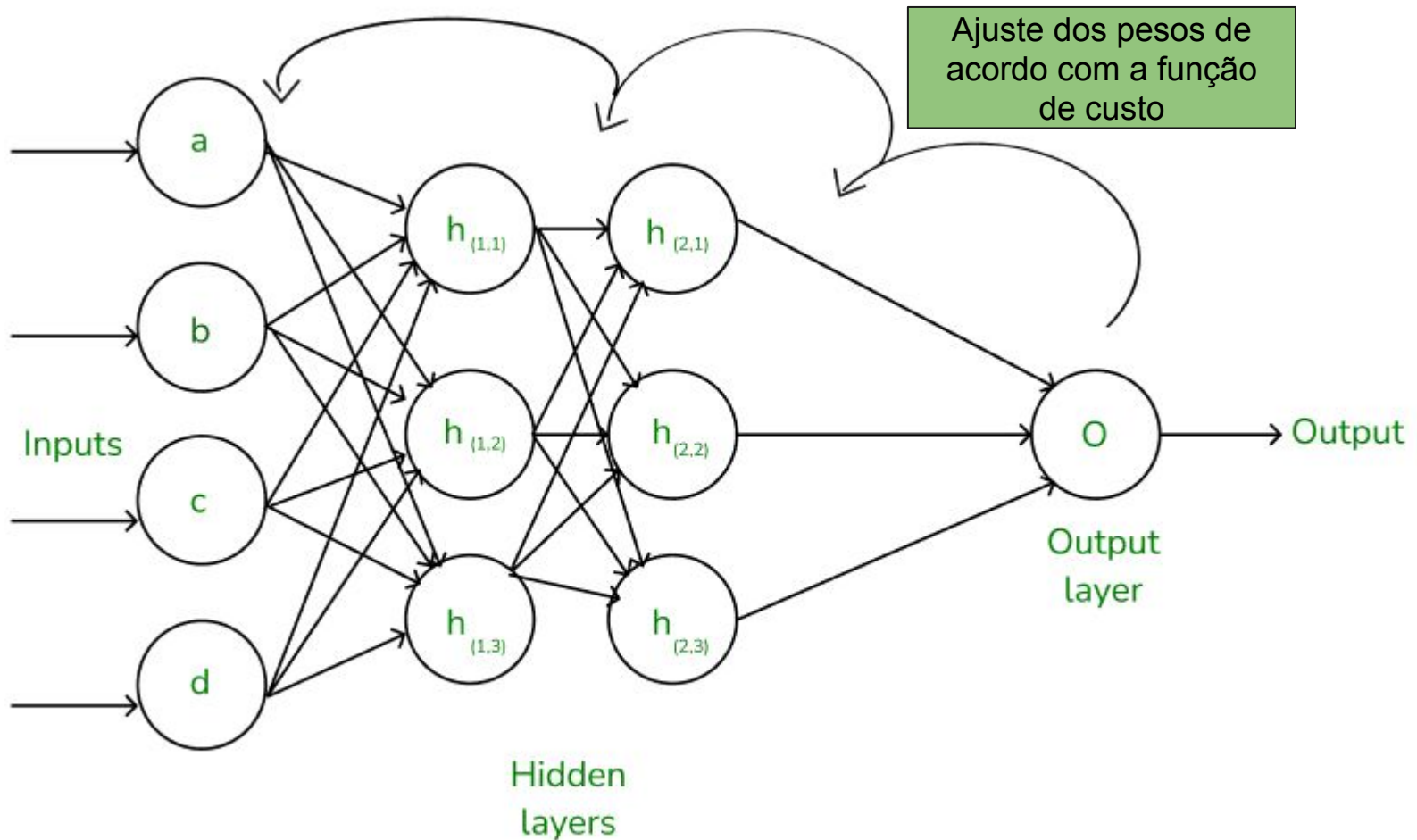
De onde vêm os pesos?



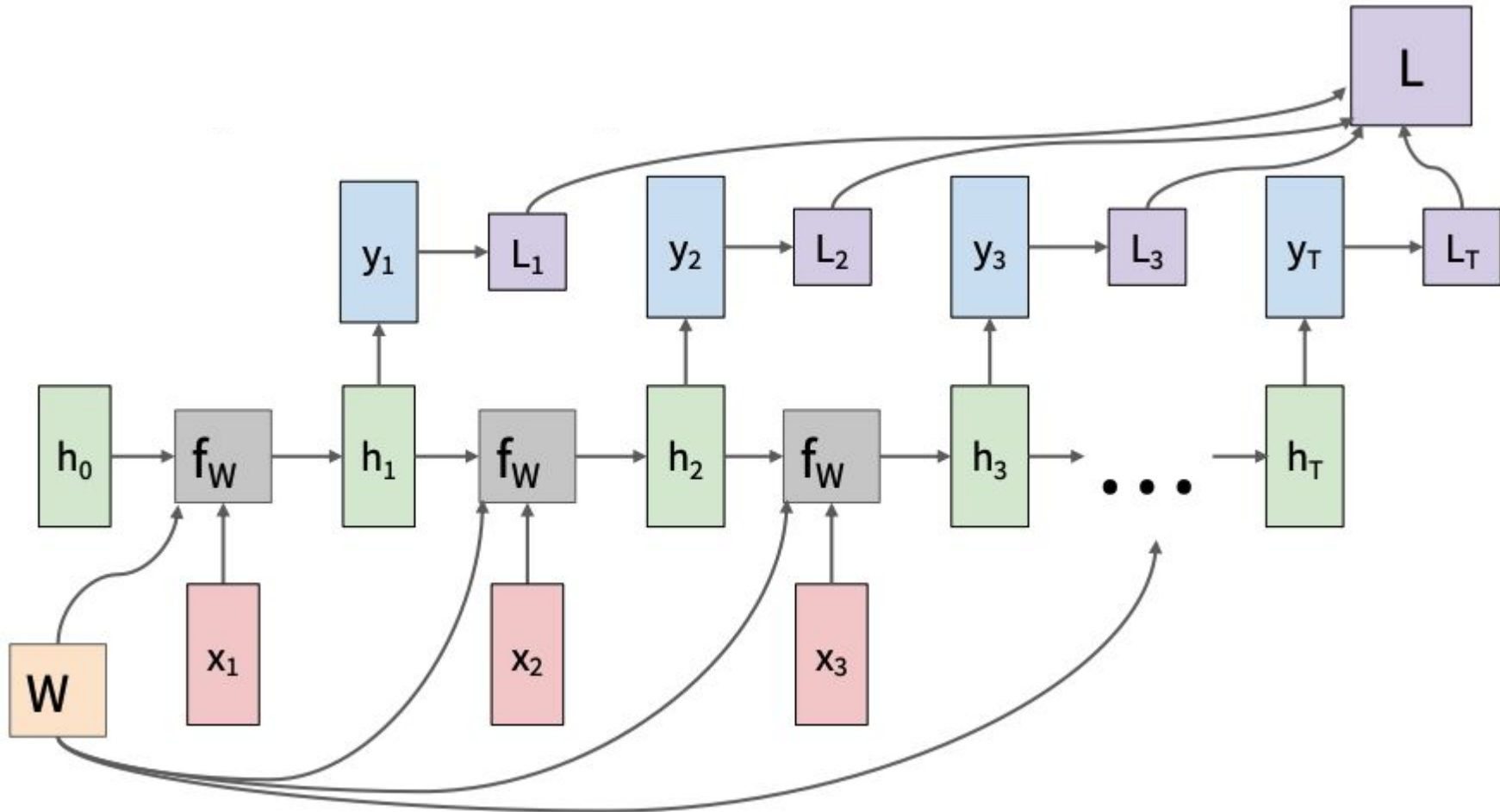
Forward para redes não recorrentes



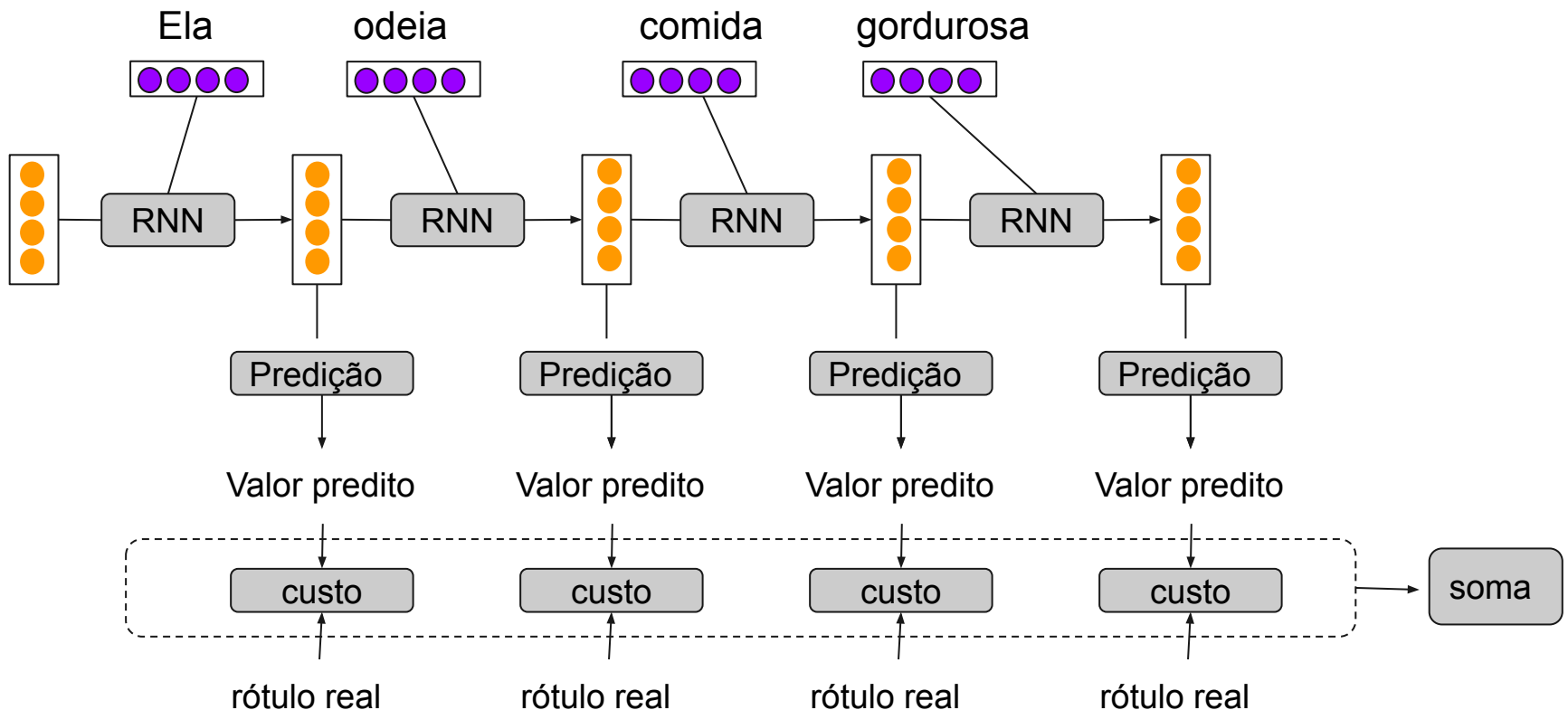
Backpropagation para redes não recorrentes



Muitos para muitos



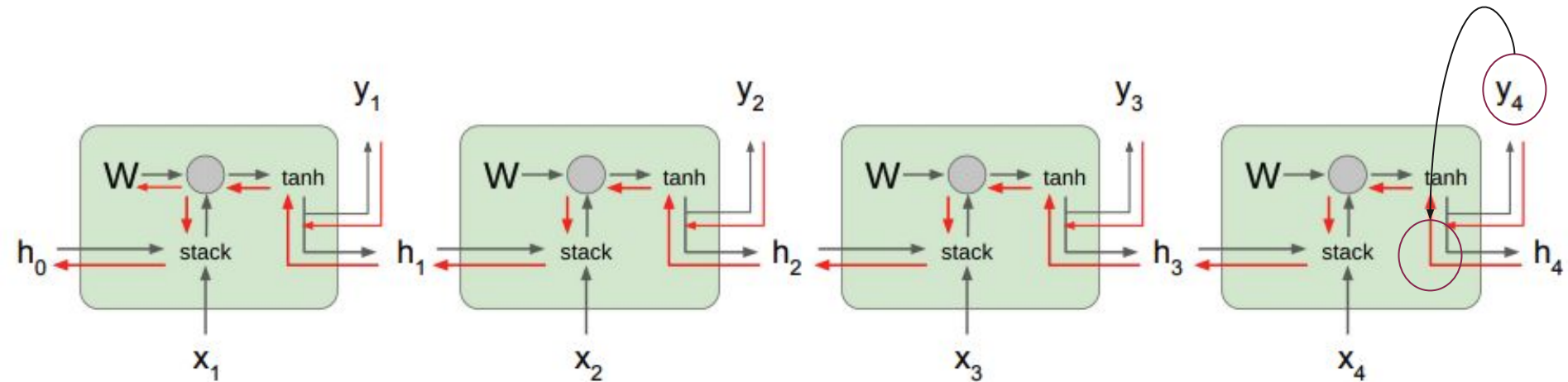
Treinamento



Backpropagation through time

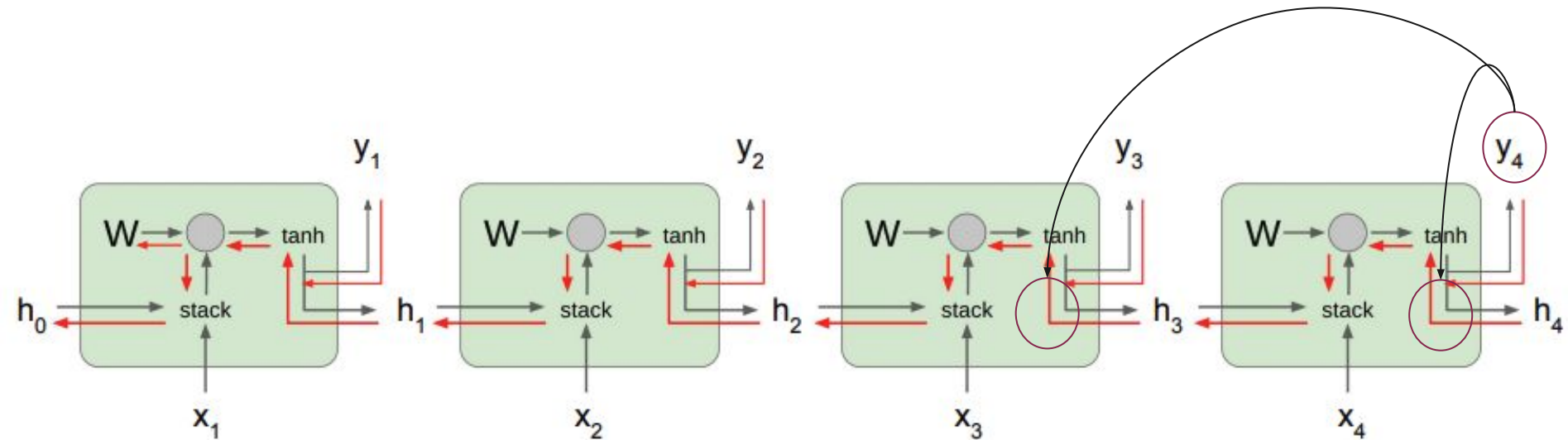
- O estado da camada escondida no instante de tempo t contribui para
 - A saída e seu erro associado no tempo t
 - A saída e o erro no instante de tempo $t+1$

Treinamento



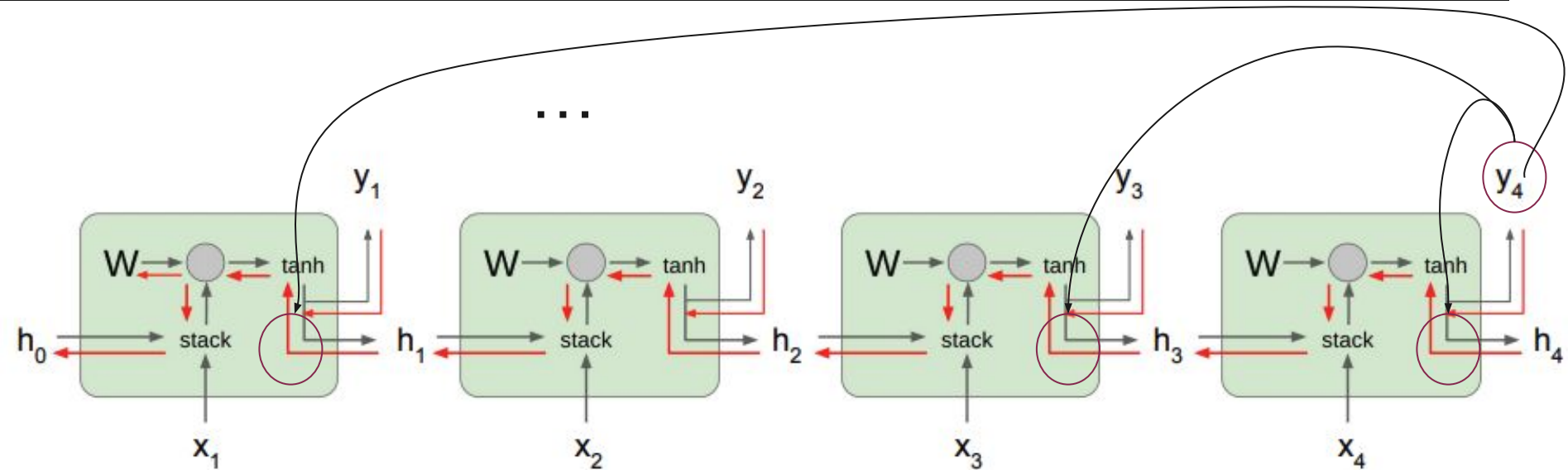
Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013
Cs231n. Stanford 2022

Treinamento



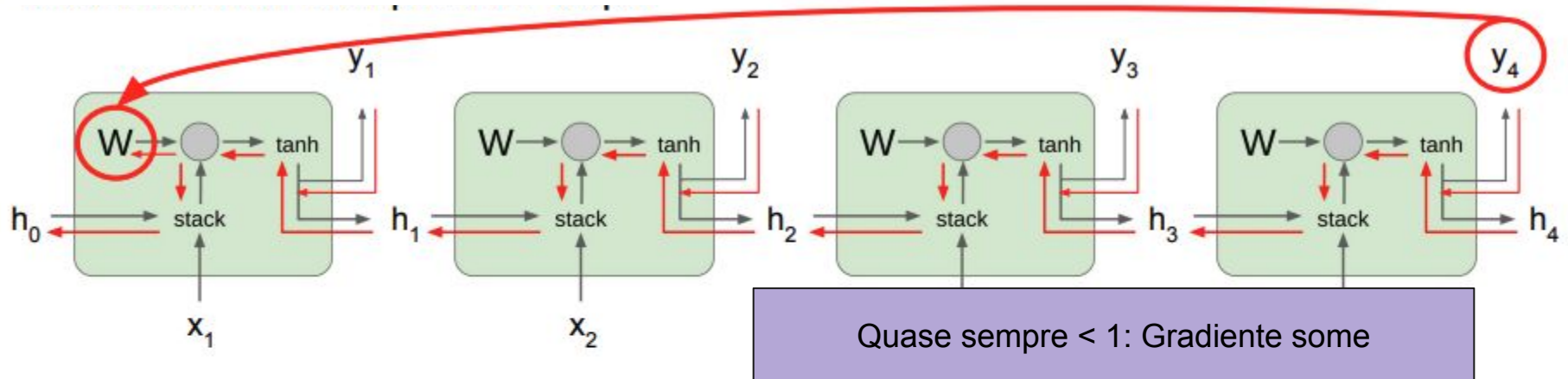
Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013
Cs231n. Stanford 2022

Treinamento



Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013
Cs231n. Stanford 2022

Treinamento



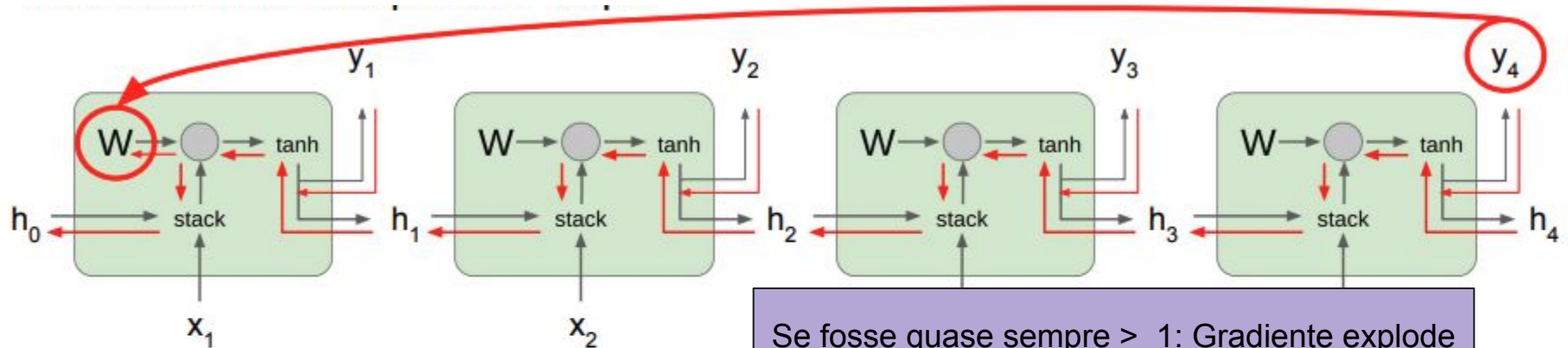
$$\frac{\partial L}{\partial W} = \sum_{t=1}^T \frac{\partial L_t}{\partial W}$$

$$\frac{\partial h_t}{\partial h_{t-1}} = \tanh'(W_{hh} h_{t-1} + W_{xh} x_t) W_{hh}$$

$$\frac{\partial L_T}{\partial W} = \frac{\partial L_T}{\partial h_T} \frac{\partial h_t}{\partial h_{t-1}} \cdots \frac{\partial h_1}{\partial W} = \frac{\partial L_T}{\partial h_T} \left(\prod_{t=2}^T \frac{\partial h_t}{\partial h_{t-1}} \right) \frac{\partial h_1}{\partial W}$$

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
 Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013
 Cs231n. Stanford 2022

Treinamento



$$\frac{\partial L}{\partial W} = \sum_{t=1}^T \frac{\partial L_t}{\partial W}$$

$$\frac{\partial h_t}{\partial h_{t-1}} = \tanh'(W_{hh} h_{t-1} + W_{xh} x_t) W_{hh}$$

$$\frac{\partial L_T}{\partial W} = \frac{\partial L_T}{\partial h_T} \frac{\partial h_t}{\partial h_{t-1}} \cdots \frac{\partial h_1}{\partial W} = \frac{\partial L_T}{\partial h_T} \left(\prod_{t=2}^T \frac{\partial h_t}{\partial h_{t-1}} \right) \frac{\partial h_1}{\partial W}$$

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
 Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013
 Cs231n. Stanford 2022

Mas e se eu precisar processar uma dependência de longo prazo?

FOLHA DE REDAÇÃO		029	
EXAME NACIONAL DO ENSINO MÉDIO - ENEM 2021			
Nome completo do Participante: EVELY APARECIDA SILVA LIMA		INSTRUÇÕES	
Nº de Inscrição: _____		1. Verifique se o seu CPF, o seu nome e a sua data de nascimento estão corretos e assine no local indicado.	
CPF: _____	Data de Nascimento: _____	2. Transcreva a sua redação com caneta esferográfica de tinta preta, fabricada em material transparente.	
Assinatura do Participante: <i>Evelyn Aparecida Silva Lima</i>		3. Não haverá substituição desta FOLHA DE REDAÇÃO por erro de preenchimento do Participante.	
		4. Escreva a sua redação com letra legível. No caso de erro, risque, com um traço simples, a palavra, a frase, o trecho ou o sinal gráfico e escreva, em seguida, o respectivo substitutivo.	
		5. Não será avaliado texto escrito em local indevido. Respeite rigorosamente as margens.	

1	A Constituição Federal de 1988, norma de maior hierarquia do sistema jurídico brasileiro,
2	possui em si o caráter de bem-estar da população. Entretanto, quando se observa a deficiência de
3	visibilidade social como forma de garantir o acesso à cidadania no Brasil, verifica-se que esse processo é
4	total e a falta de acesso é mais desfavoravelmente na prática. Dessa forma, essa realidade se deve, à
5	inequívoca realidade social.
6	Primeiramente, vale ressaltar que a debilidade da Poder Público, possui íntima relação com o
7	Estado de Direito, Thomas Hobbes, em seu livro "Leviatã" defende a obrigação de Estado um
8	que auxilia e protege de todos os cidadãos, todavia, não de todos os cidadãos. Hobbes, um
9	que possui um papel íntimo em relação à invisibilidade de pessoas com o registro civil e, por
10	consequência disso, a falta de uma segurança estabelecida pela Constituição, um dos
11	livros de pessoas não possuem a certeza de movimento, mostrando um alto nível de
12	insegurança e migração, incluindo a ausência de acesso ao registro civil, devido a
13	diversidade e falta de acesso a serviços públicos essenciais como hospitais, escolas,
14	que a invisibilidade a essas pessoas, seja decorrente da falta de dados governamentais,
15	Adicionalmente, uma grande parcela da população, os mestres chamados "Povoado de
16	é um livro escrito pelo jornalista Vladimir Stankovic que a importância a seguir a ética de
17	medir, ou seja, a possibilidade das pessoas terem as informações importantes para
18	na sociedade, para-se que a garantia de acesso à cidadania, encontra-se em
19	na realidade, ocorre porque, inicialmente, a realidade não se encontra em
20	uma problemática, pelo contrário, ela adquire uma perspectiva individualista, por
21	a falta de um registro civil, ou seja, a impossibilidade de acessar outros
22	dados, é essencial suprir esses pontos que afetam, sobretudo, um
23	Para isso, portanto, a necessidade de garantir o acesso à cidadania, para
24	o Estado, o Poder Público, responsável por administrar e
25	de Ministério da Cidadania, a partir de medidas governamentais
26	atribuição de serviços essenciais para cidadãos que não
27	um registro civil. Essa não será a única medida, mas
28	tanto, para que também a realidade não se torne uma
29	conjuntura de tais ações, os brasileiros terão o direito
30	garantido pela Constituição, como uma

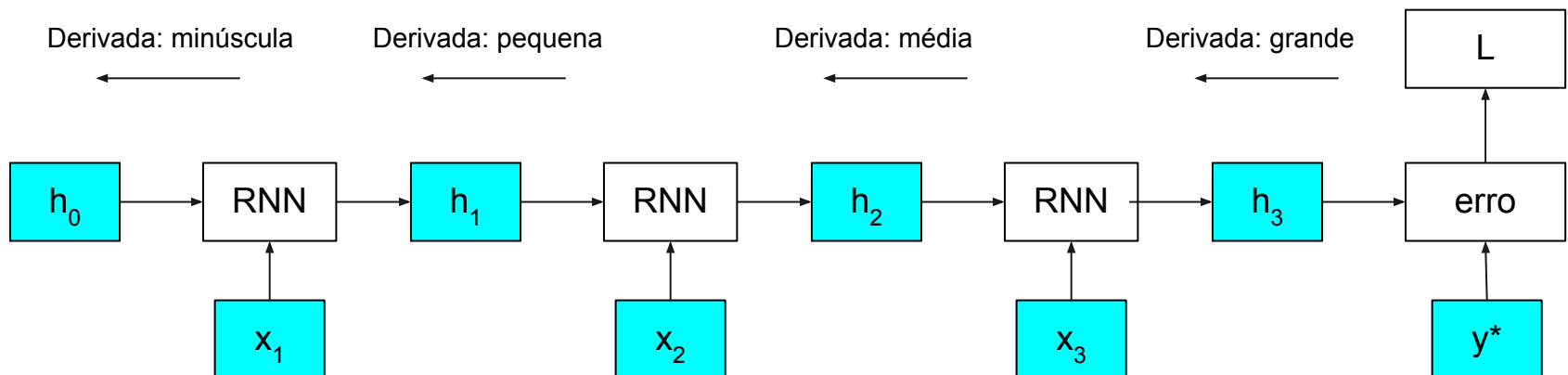
029221100245676500

Memória de longo prazo

- Redes recorrentes têm dificuldade em lidar com dependência de longa distância
 - Camadas intermediárias devem
 - Fornecer informação útil para o instante corrente
 - Atualizar e carregar informação de contexto para decisões futuras

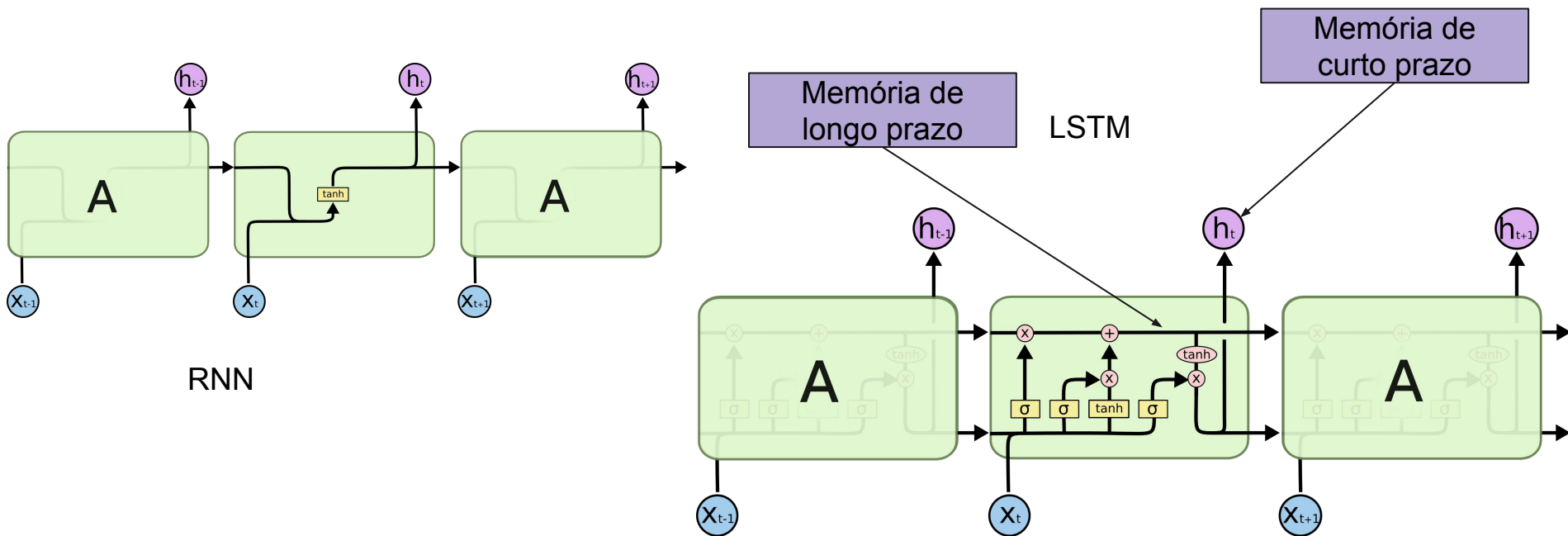
Vanishing gradient

- Redes recorrentes têm dificuldade em lidar com dependência de longa distância
 - *Backpropagação* do sinal do erro através do tempo
 - Camada escondida contribui para a perda do instante de tempo seguinte



Long Short-Term memory (LSTM) (Hochreiter and Schmidhuber, 1997)

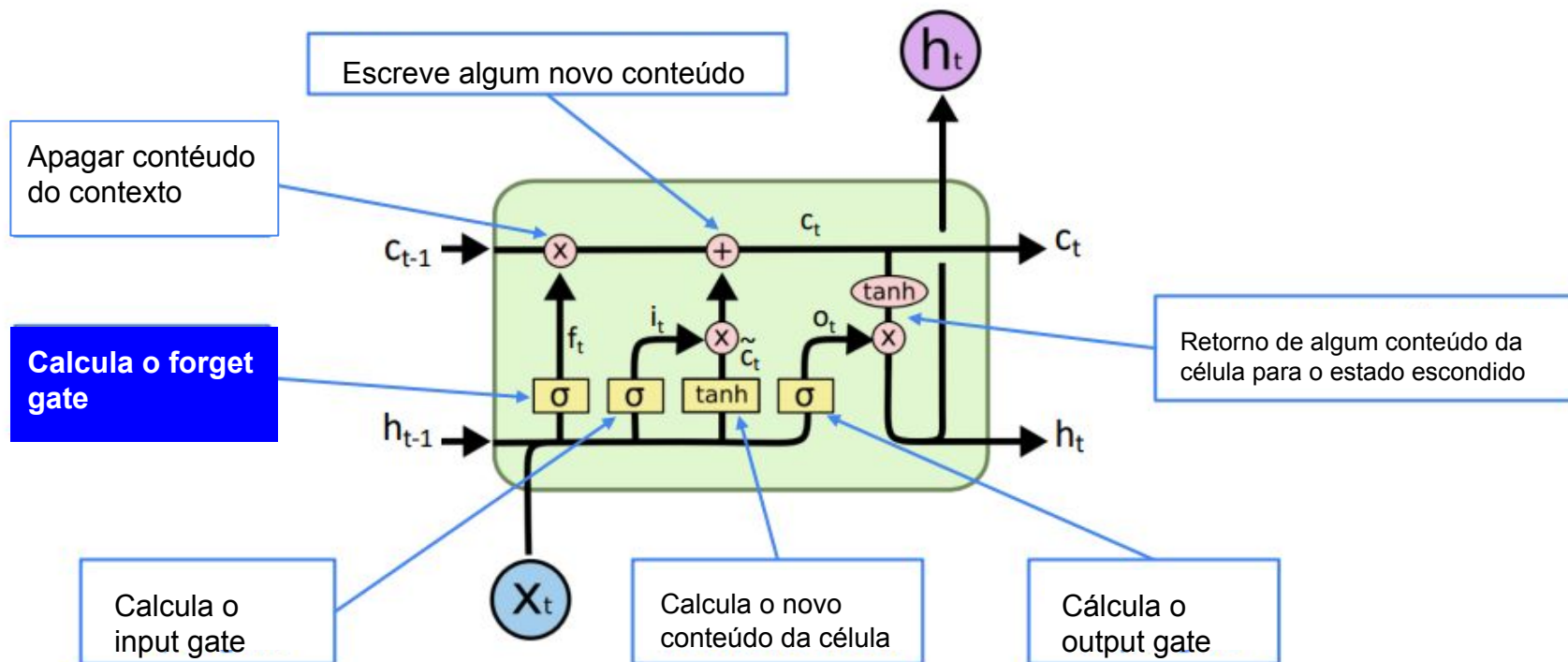
- RNN com uma estrutura de “memória”
- A cada passo, existe um estado escondido e uma célula de estado (vetores)
 - Contexto explícito
- A célula armazena informação de “longo termo”



Long Short-Term memory (LSTM) (Hochreiter and Schmidhuber, 1997)

- A rede pode apagar, escrever e ler informação da célula
 - A seleção de qual informação passará por cada operação é controlada por *gates* (vetores)
 - Conexões aditivas
 - Camada feedforward + sigmoid + multiplicação
 - A cada passo, as operações nos gates podem ser: abrir (1), fechar (0) ou algo no meio do caminho
 - Gates são dinâmicos: seu valor é calculado com base no contexto corrente

Long Short-Term Memory (LSTM)



Long Short-Term Memory

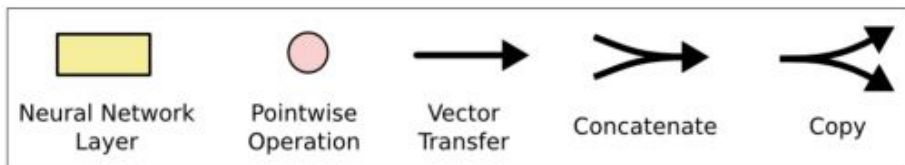
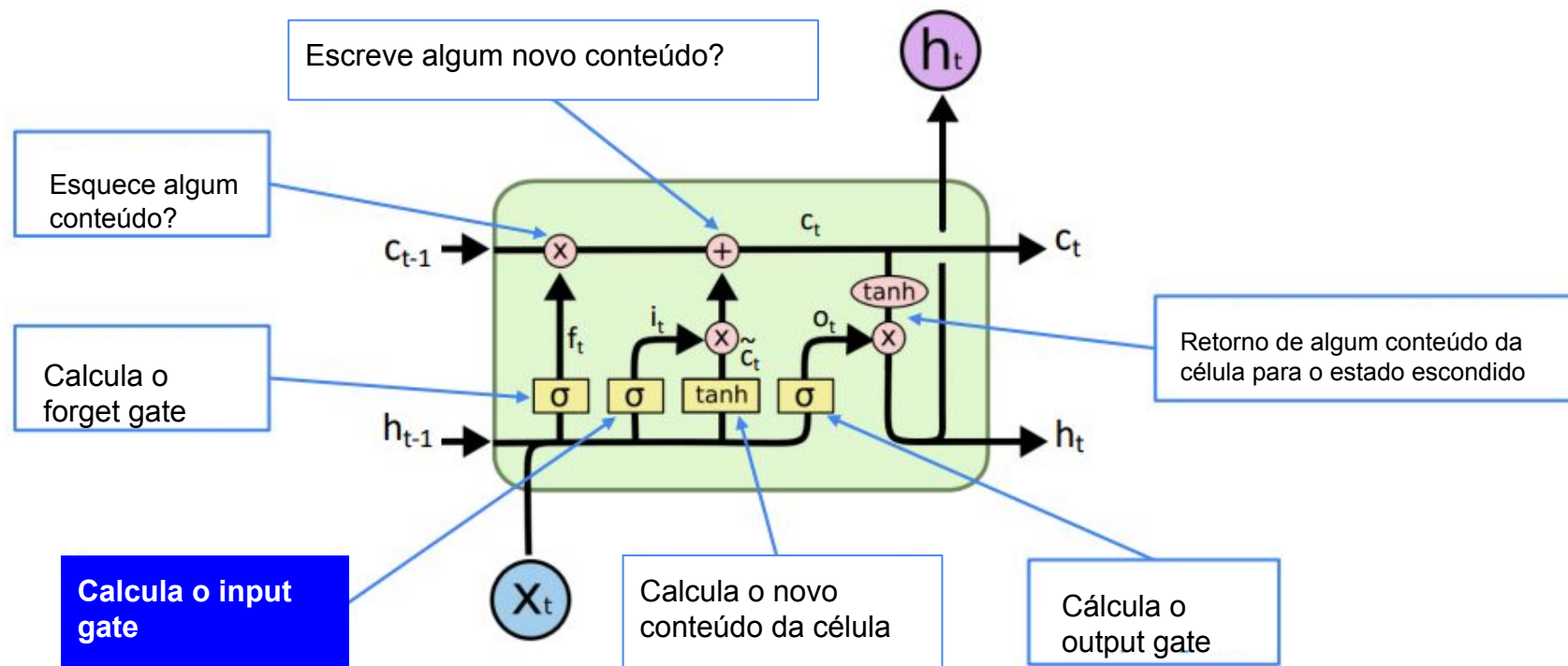
Forget gate: o que eu não devo esquecer
(ou o que eu devo lembrar do curto prazo)

Sigmoid: valores entre 0
e 1

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

concatenação

Long Short-Term Memory (LSTM)



Long Short-Term Memory

Forget gate: o que eu não devo esquecer

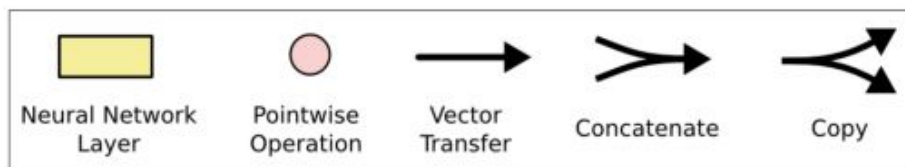
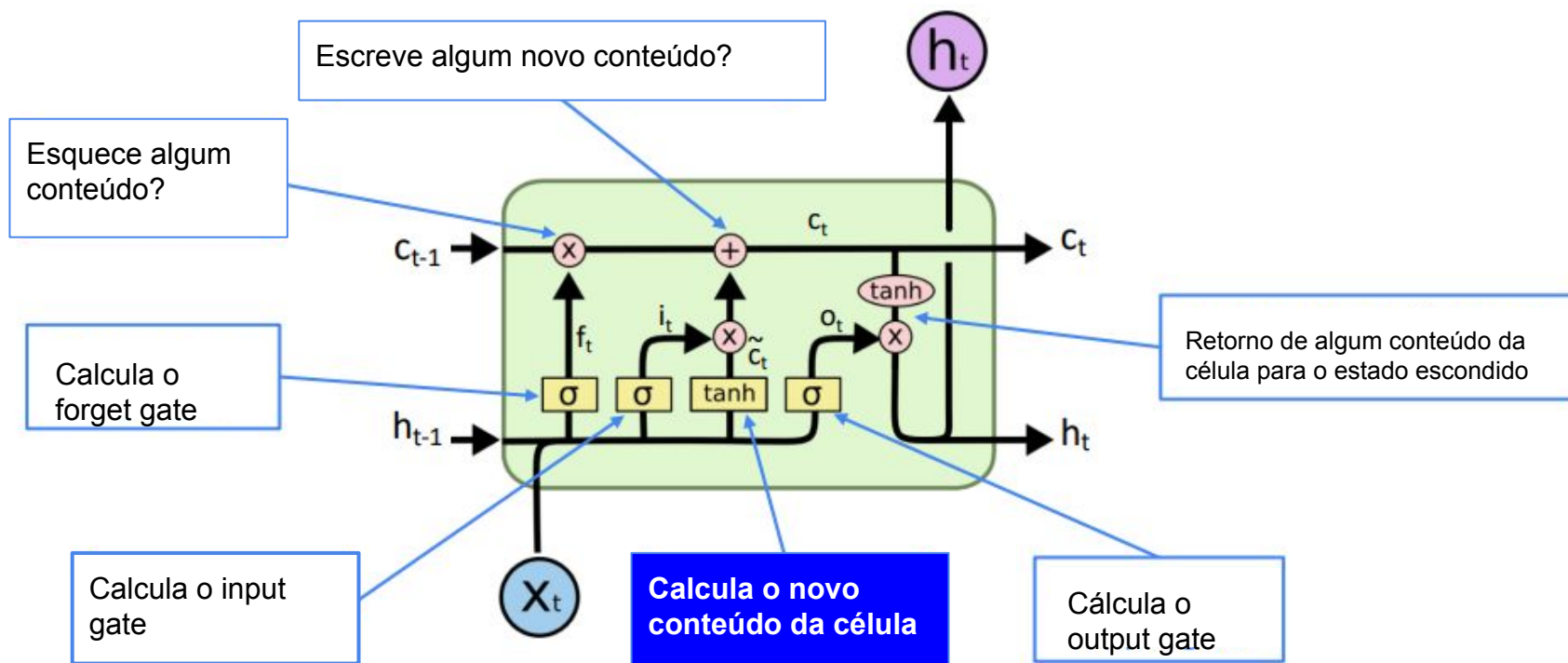
Input gate: Escreve na célula?
(escrita)

Sigmoid: valores entre 0
e 1

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

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

Long Short-Term Memory (LSTM)



Long Short-Term Memory

Forget gate: o que é armazenado vs o que é esquecido, a partir da célula anterior

Input gate: Escreve na célula?

Sigmoid: valores entre 0 e 1

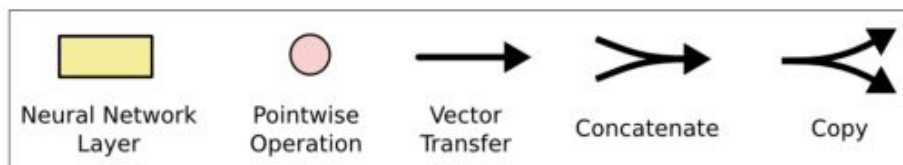
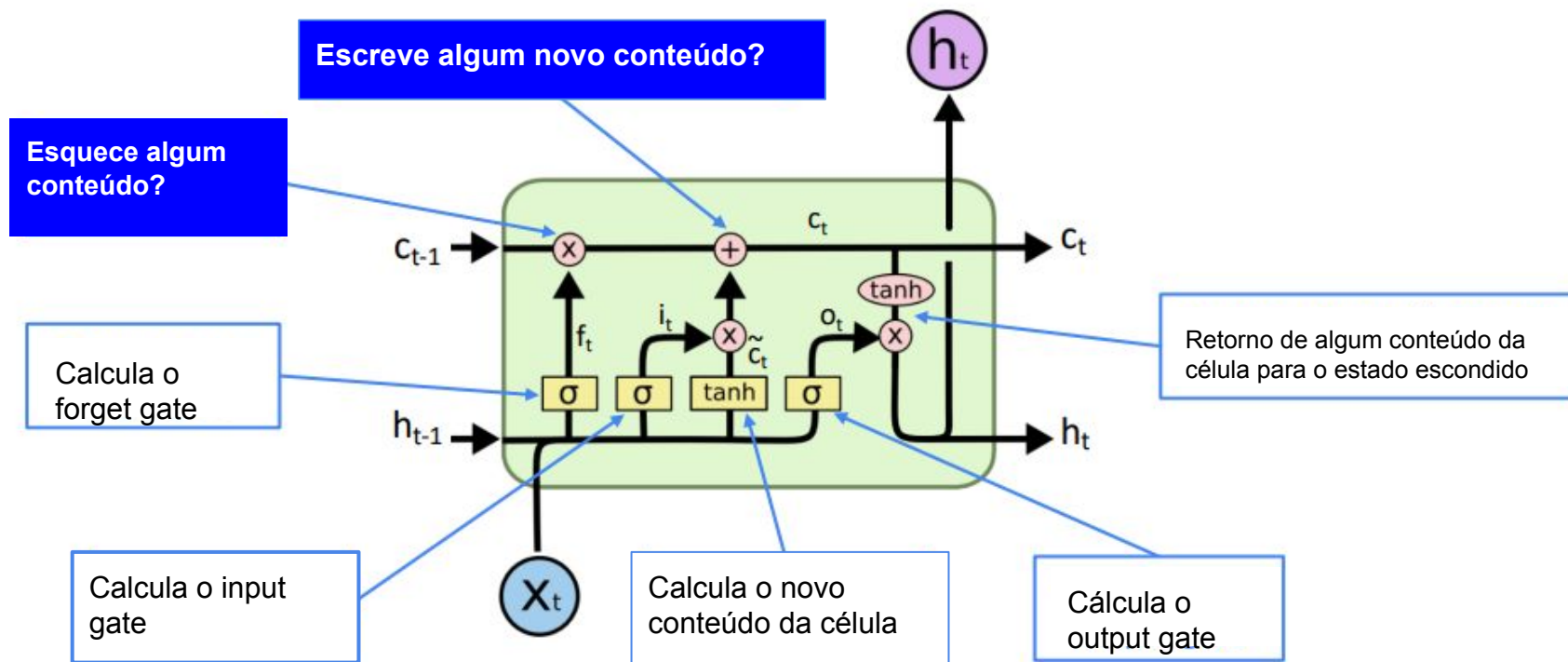
$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

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

Novo conteúdo a ser escrito na célula

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

Long Short-Term Memory (LSTM)



Long Short-Term Memory

Forget gate: o que é armazenado vs o que é esquecido, a partir da célula anterior

Input gate: ESCRITA

Sigmoid: valores entre 0 e 1

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

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

Novo conteúdo a ser escrito na célula

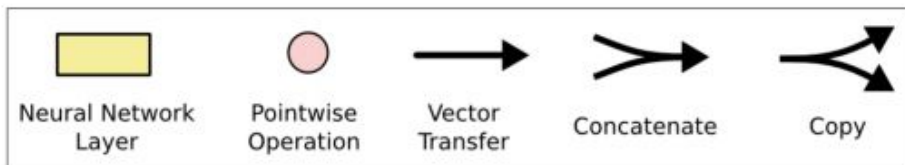
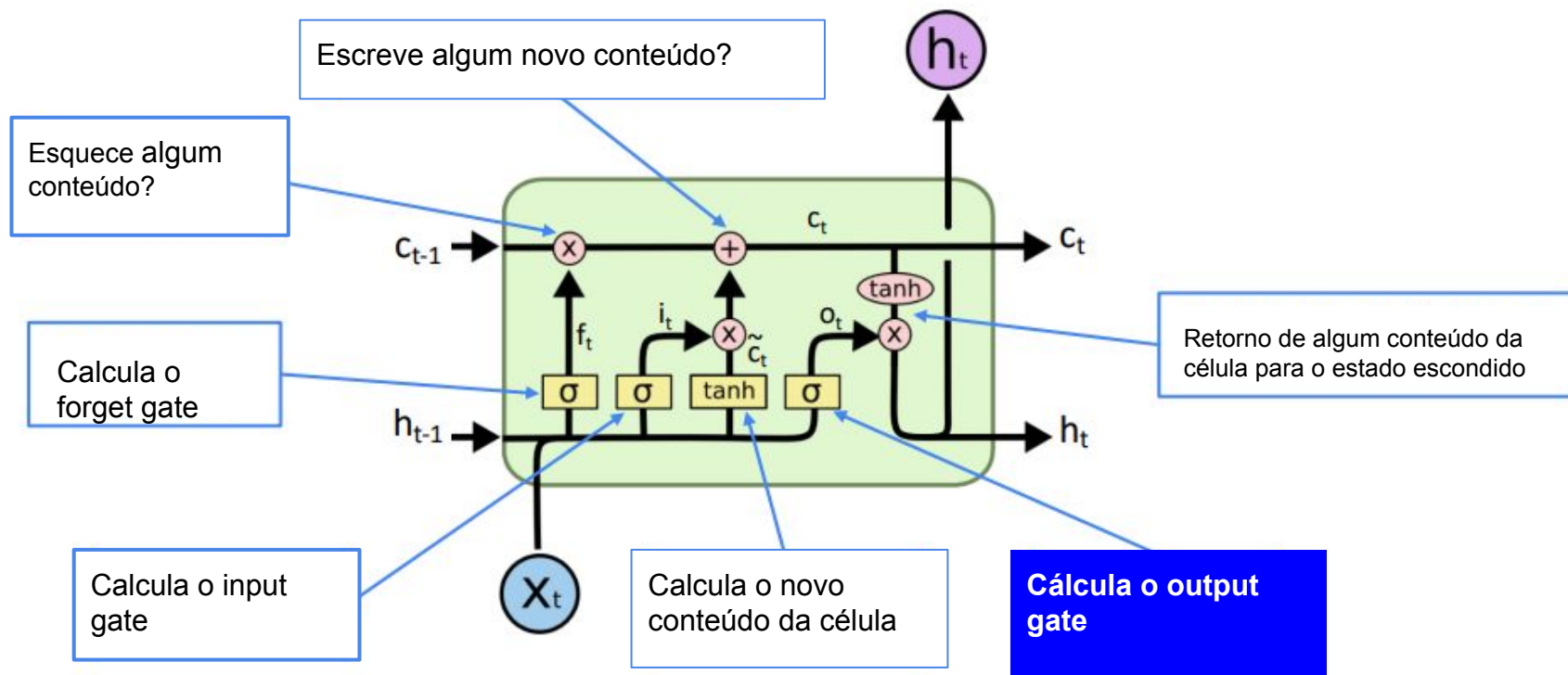
Estado da célula: “esquece” algum conteúdo do estado anterior da célula e escreve algum conteúdo novo

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

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

Produto de elemento a elemento

Long Short-Term Memory (LSTM)



Long Short-Term Memory

Forget gate: o que é armazenado vs o que é esquecido, a partir da célula anterior

Input gate: que partes do novo conteúdo são escritos para a célula (escrita)

Output gate: o que levar da célula? (leitura)

Novo conteúdo a ser escrito na célula

Estado da célula: “esquece” algum conteúdo do estado anterior da célula e escreve algum conteúdo novo

Sigmoid: valores entre 0 e 1

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

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

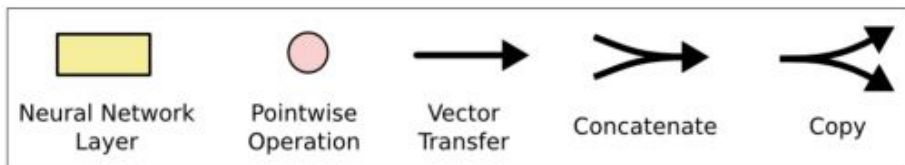
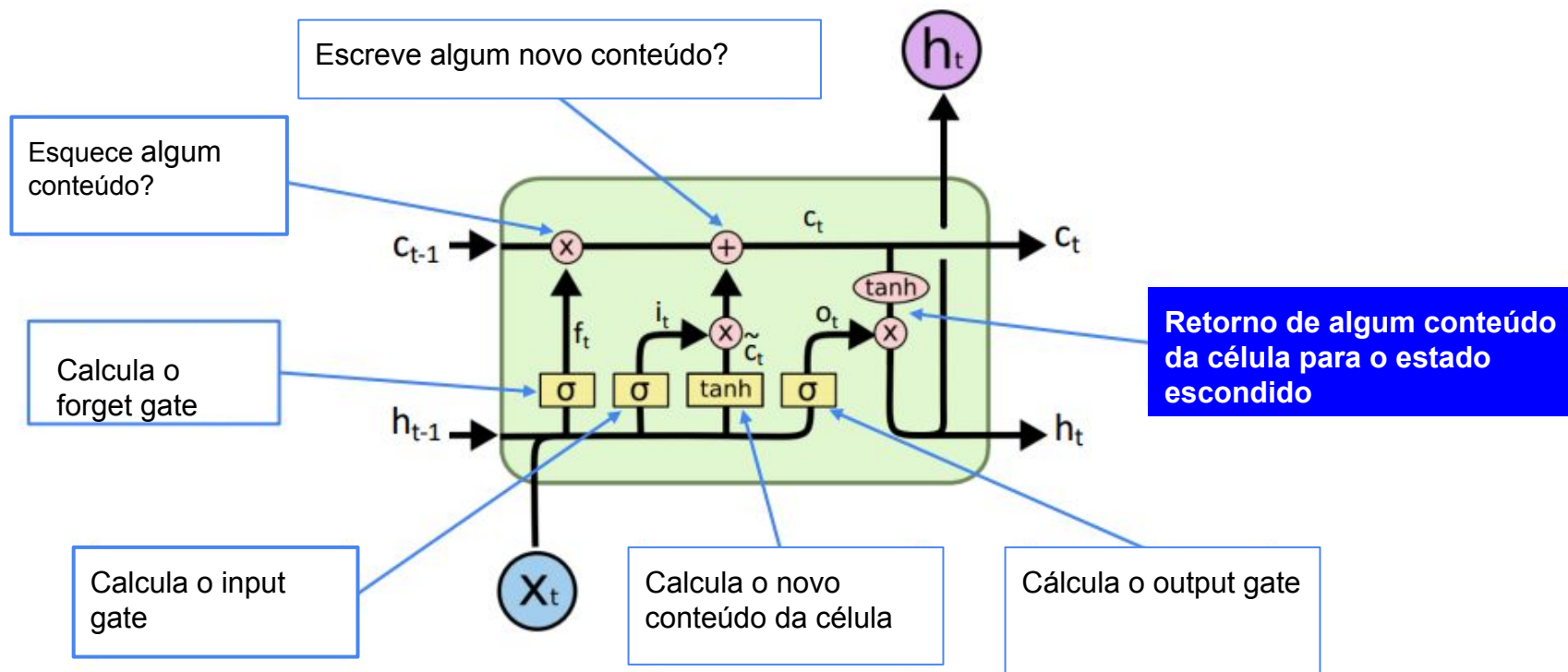
$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

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

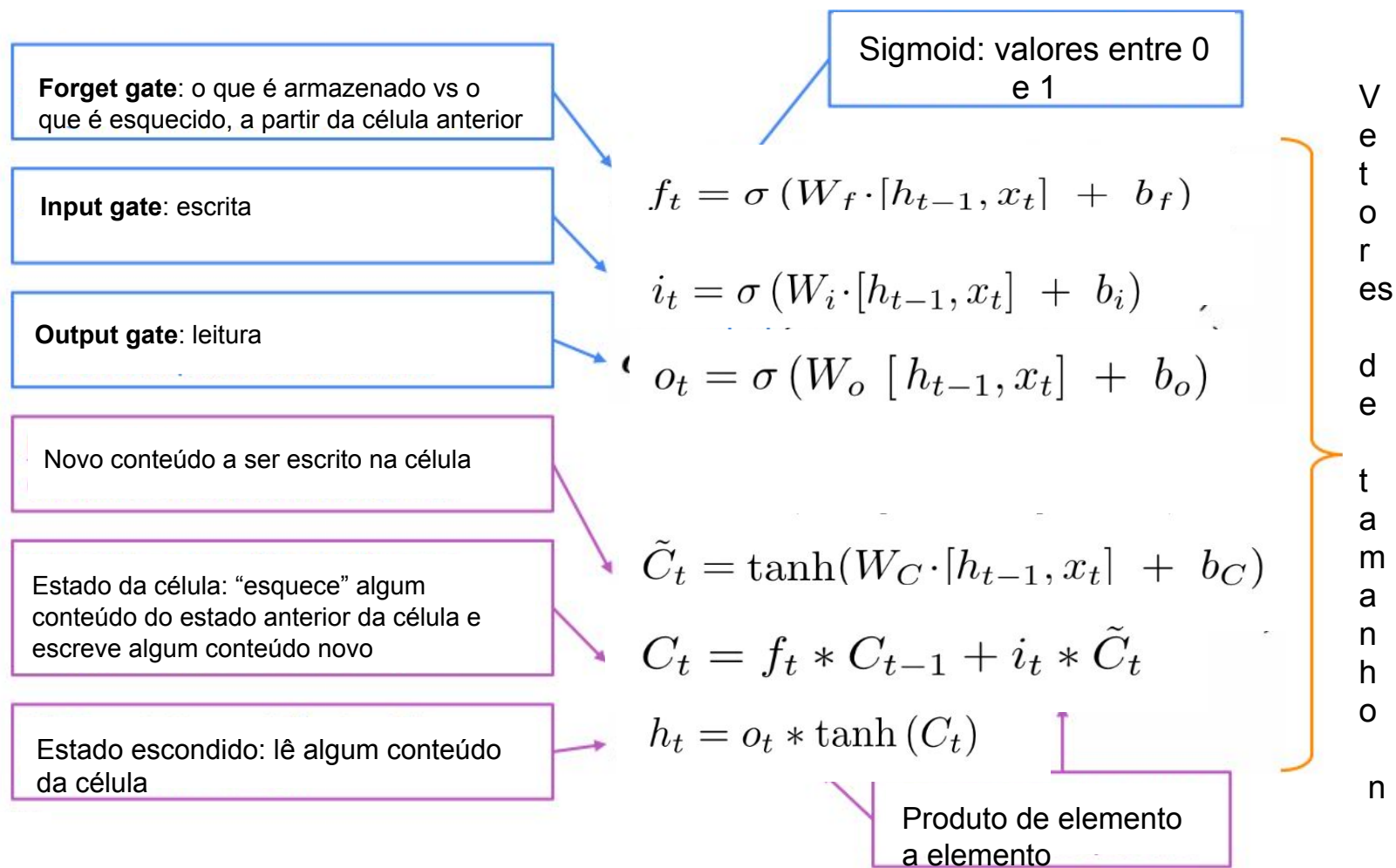
$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

Produto de elemento a elemento

Long Short-Term Memory (LSTM)



Long Short-Term Memory



Gated Recurrent Units (GRUs) (Cho et al., 2014)

- LSTM pode ser bem custosa para treinar
 - 8 matrizes de peso (duas para cada gate)
- GRUs
 - Dispensam o vetor de contexto (célula)
 - Usam apenas dois gates
 - Reset
 - O que é relevante no estado anterior e o que pode ser ignorado?

$$r_t = \sigma(W_r[h_{t-1}; x_t])$$

$$h'_t = \tanh(U(r_t \odot h_{t-1}) + Wx_t)$$

Gated Recurrent Units (GRUs) (Cho et al., 2014)

- LSTM pode ser bem custosa para treinar
 - 8 matrizes de peso (duas para cada gate)
- GRUs
 - Dispensam o vetor de contexto (célula)
 - Usam apenas dois gates

- Reset $r_t = \sigma(W_r[h_{t-1}; x_t])$

$$h'_t = \tanh(U(r_t \odot h_{t-1}) + Wx_t)$$

- Update

- O que de h'_t será usado diretamente no novo estado escondido h_t e o que precisa ser preservado de h_{t-1}

$$z_t = \sigma(U_z h_{t-1} + W_z x_t)$$

$$h_t = (1 - z_t) h_{t-1} + z_t h'_t$$

Embeddings contextualizados : Elmo (Peters et al., 2018)

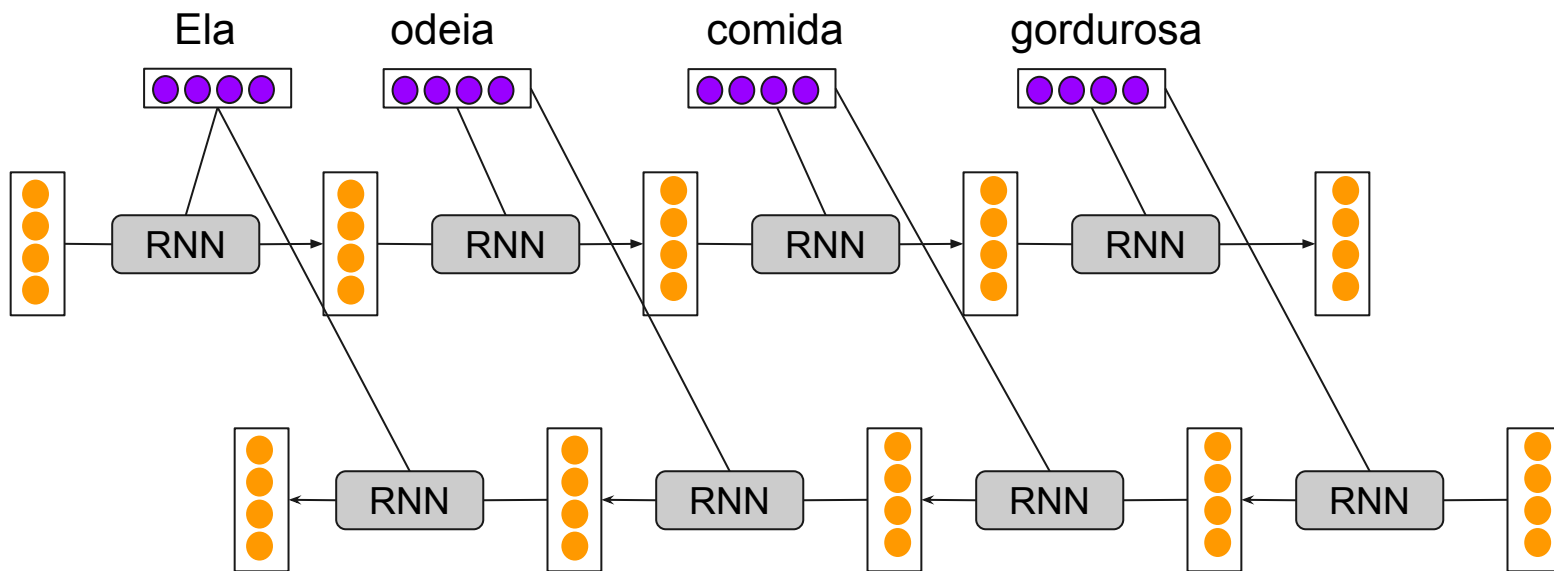


Contextualized word-embeddings can give words different embeddings based on the meaning they carry in the context of the sentence.

Also, RIP Robin Williams

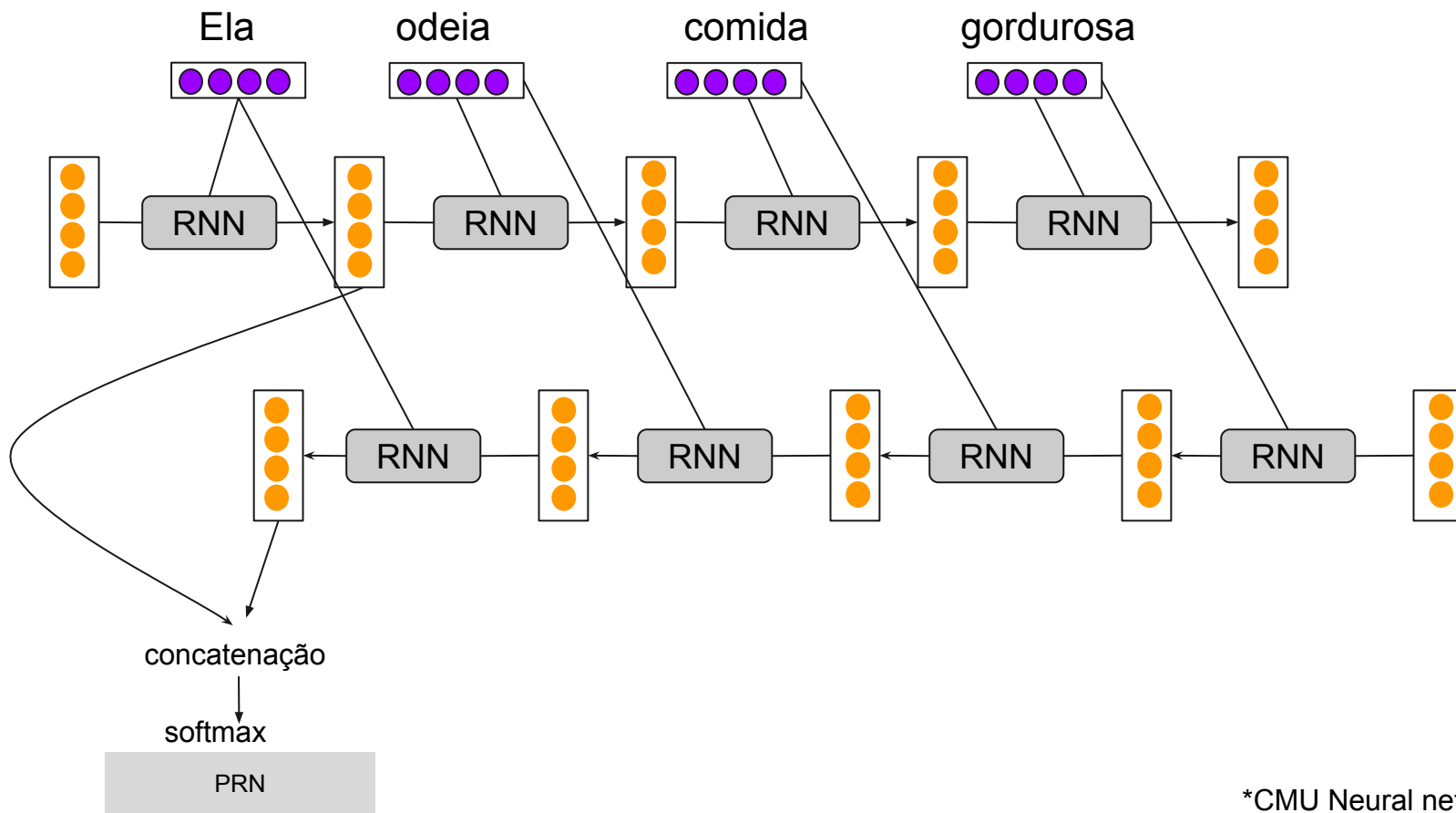
* <http://jalammar.github.io/illustrated-bert/>

RNN bidirecional

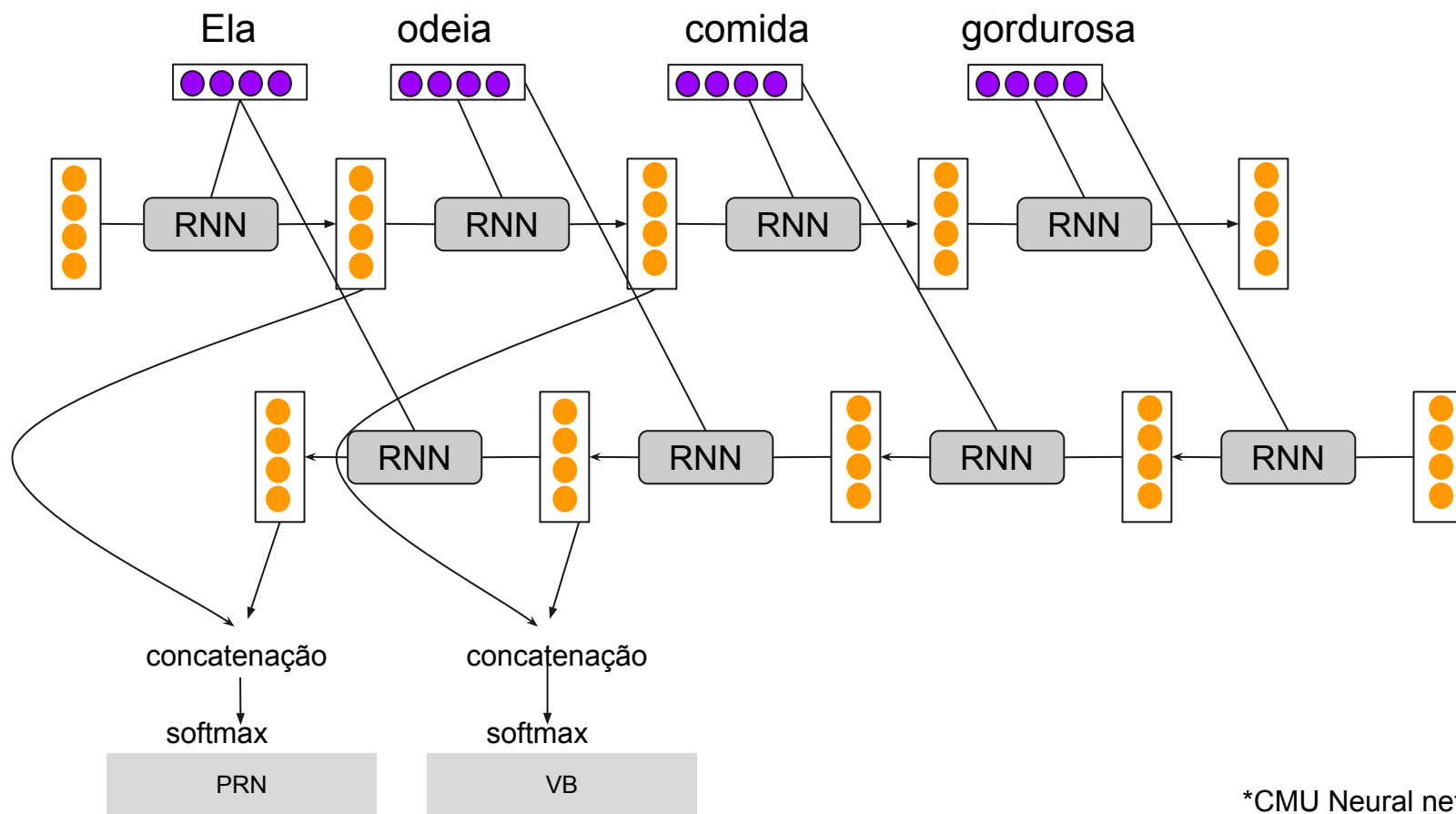


*CMU Neural nets for NLP

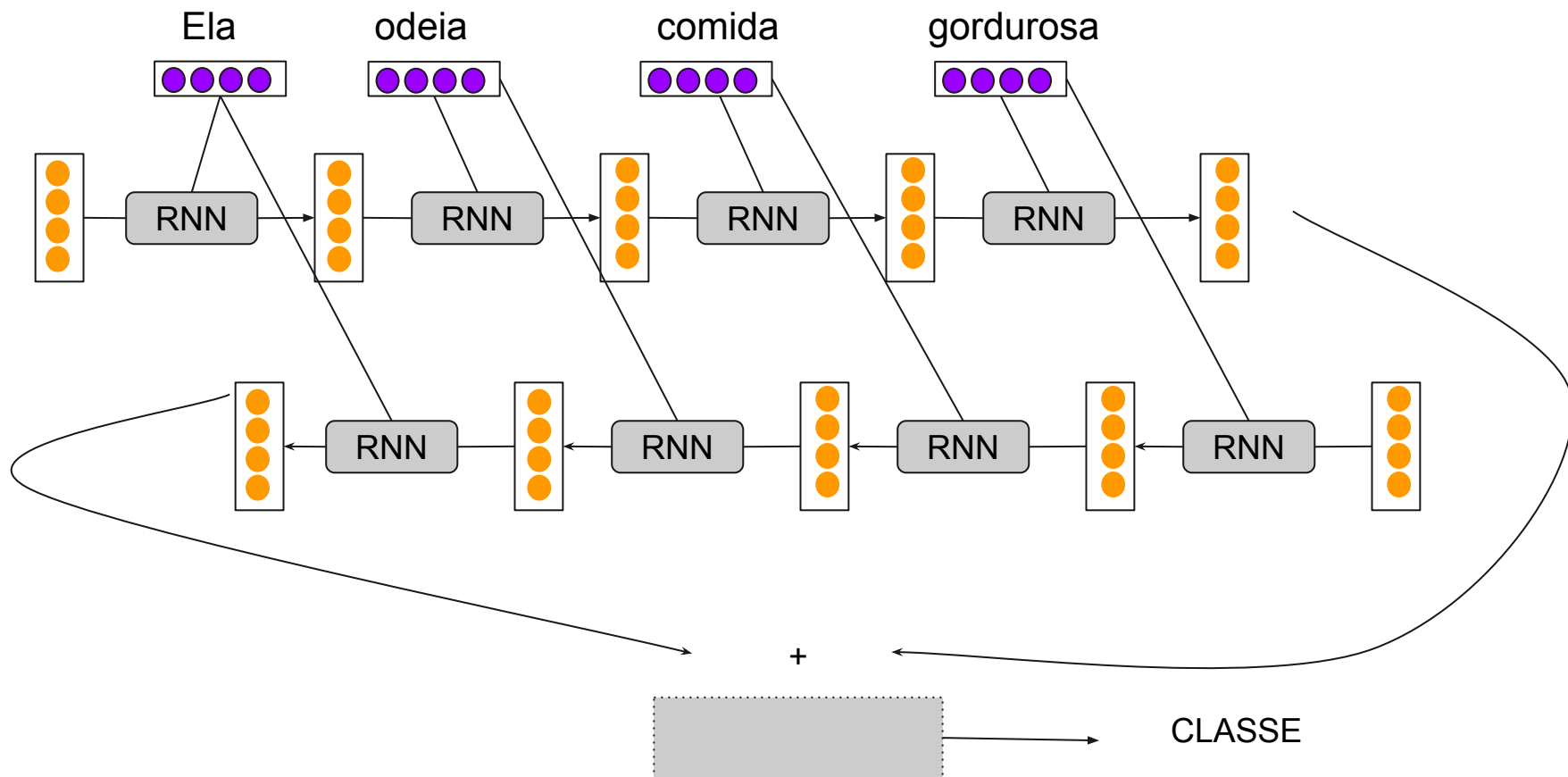
RNN bidirecional



RNN bidirecional



RNN bidirecional



Embeddings contextualizados : Elmo (Peters et al., 2018)

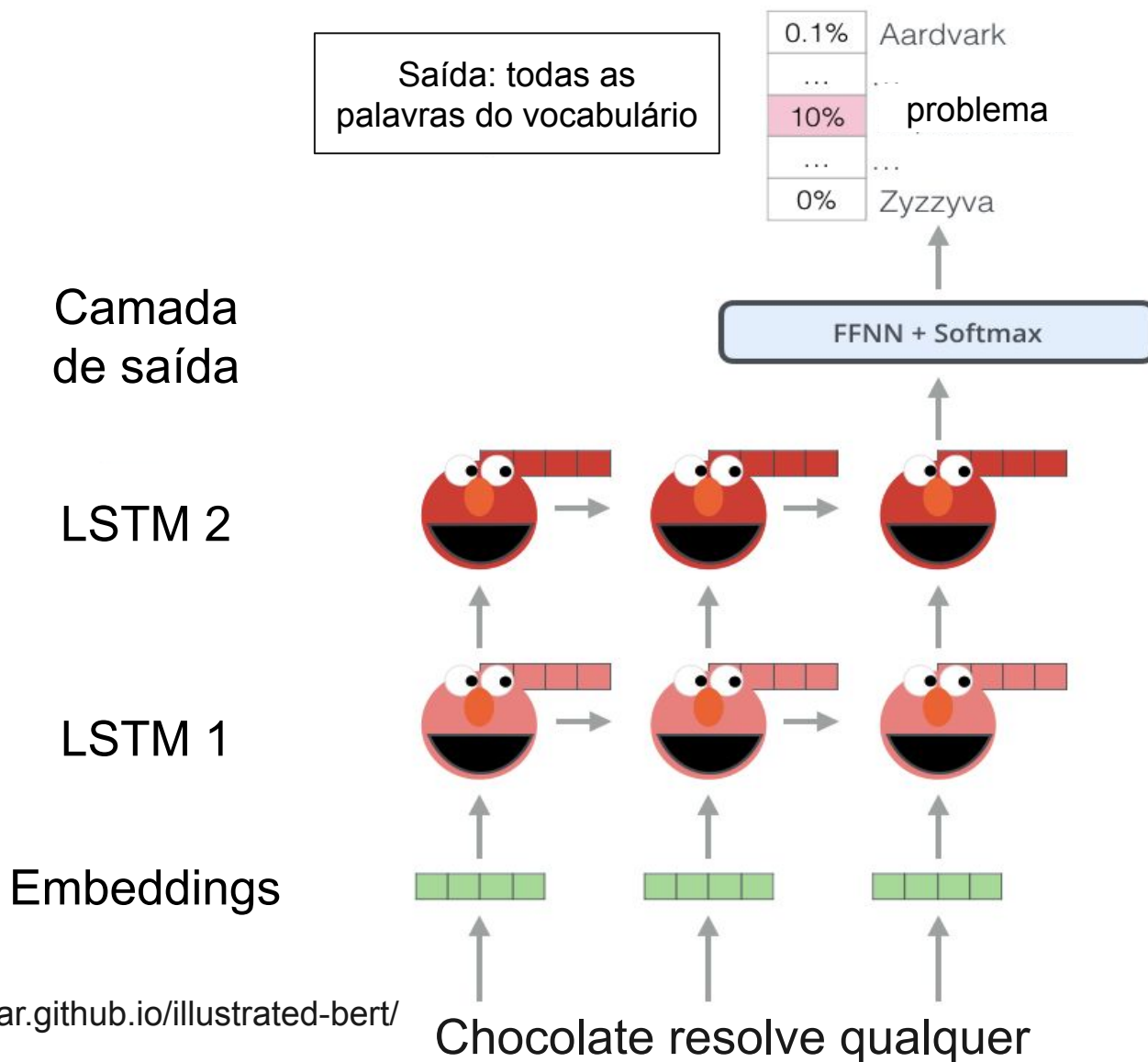
- Modelo de linguagem bidirecional

$$p(t_1, t_2, \dots, t_N) = \prod_{k=1}^N p(t_k \mid t_1, t_2, \dots, t_{k-1})$$

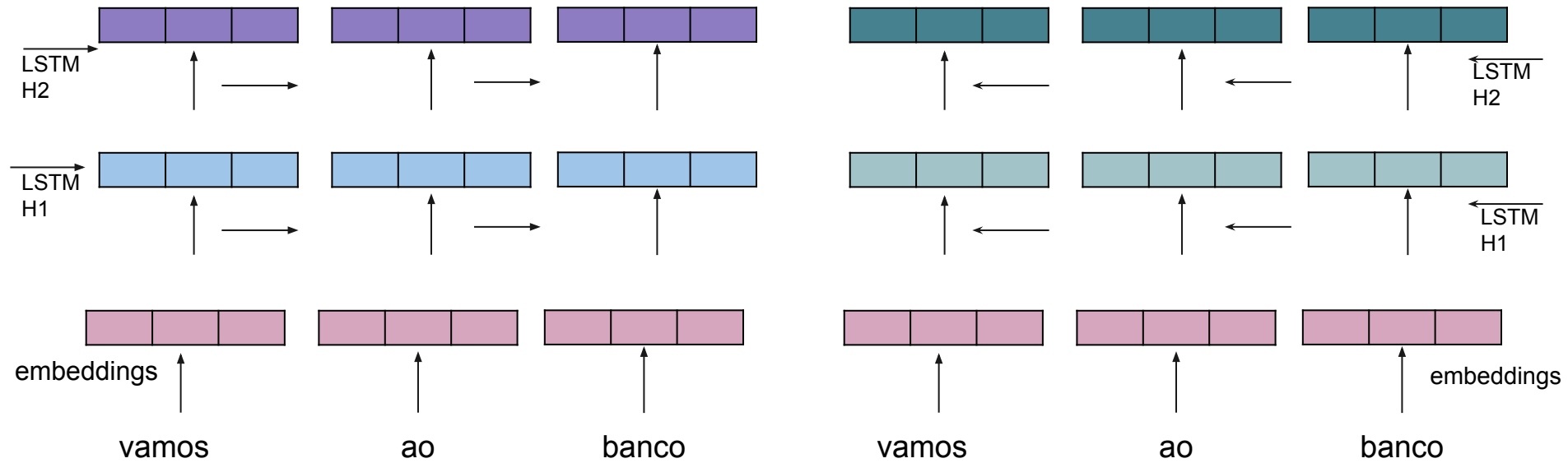
$$p(t_1, t_2, \dots, t_N) = \prod_{k=1}^N p(t_k \mid t_{k+1}, t_{k+2}, \dots, t_N)$$

- Embedding da palavra: combinação linear das camadas escondidas correspondentes

Treinamento genérico contextualizado



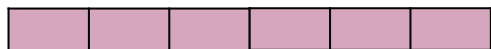
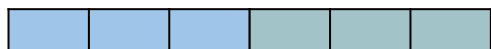
Elmo



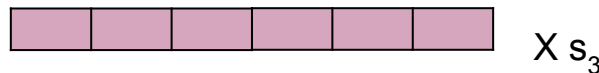
Elmo



Concatena pesos das camadas escondidas e embeddings



Cada um é multiplicado por um peso, baseado na tarefa

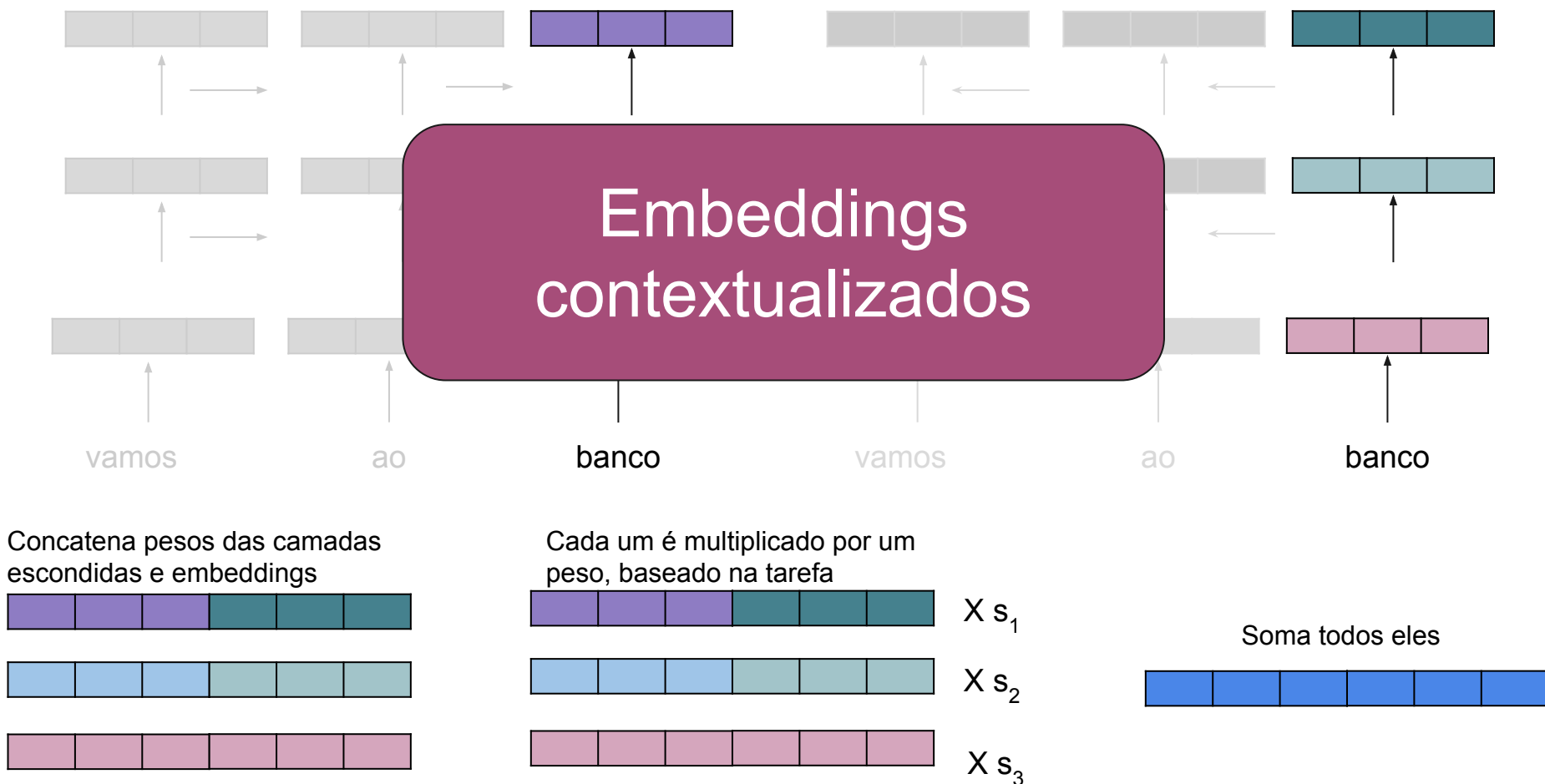


Soma todos eles



Elmo

Embeddings contextualizados



ELMo

