

Natural Language and Speech Processing

Lecture 10: Recurrent Neural Networks for Modeling Language

Tanel Alumäe

NLP and Sequential Data

- .NLP is full of sequential data
 - Words in sentences
 - Characters in words
 - Sentences in document
 -

Long-distance dependencies

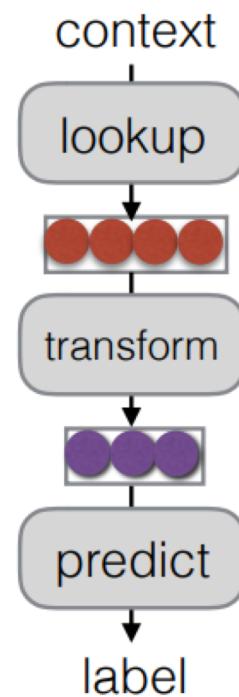
- .Agreement in number, gender, etc
 - He does not have very much confidence in himself.***
 - She does not have very much confidence in herself.***
- .Selectional preference
 - The reign has lasted as long as the life of the queen.***
 - The rain has lasted as long as the life of the clouds.***

Long-distance dependencies, II

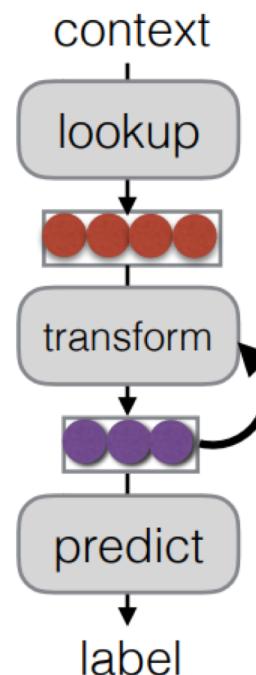
- .What does „it” refer to?
 - The trophy would not fit in the brown suitcase because it was too **big**.
Trophy
 - The trophy would not fit in the brown suitcase because it was too **small**.
Suitcase

Recurrent neural networks

.Feed-forward
neural network

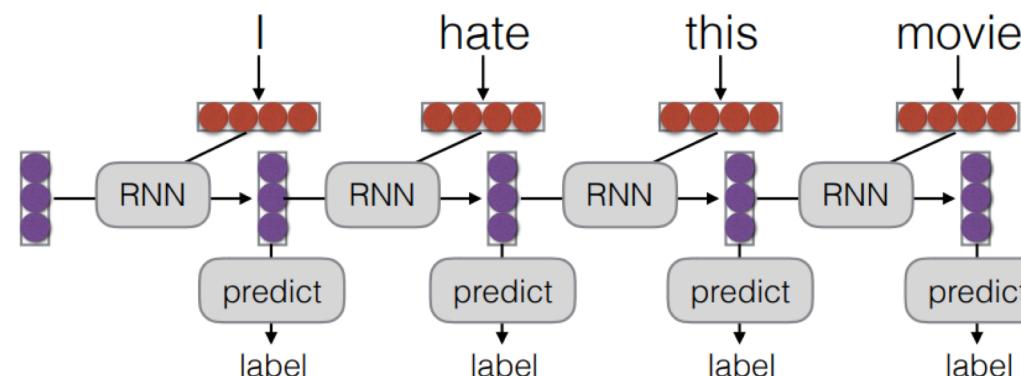


.Recurrent neural
network



Unrolling a neural network in time

- .Processing a sequence (e.g. a sentence) with a neural network
- .The current output of RNN depends on the hidden state at the prior state (a **vector**)

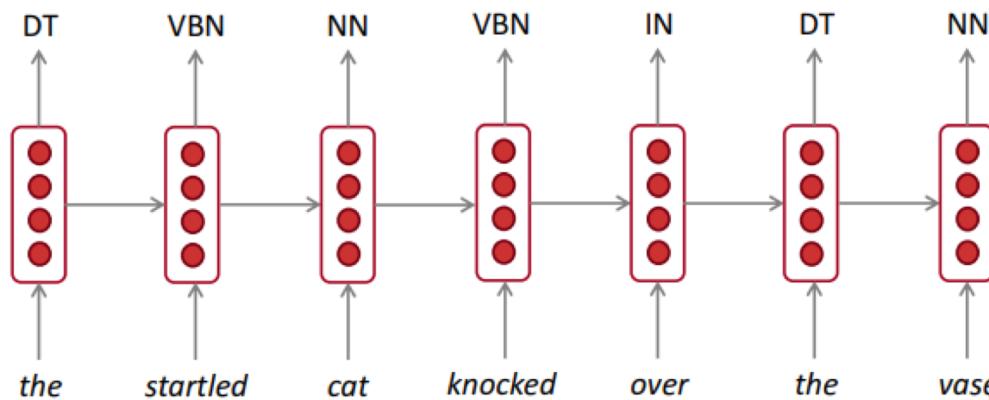


- .In NLP the current output depends on the current input word and the hidden layer output of the previous word

This is the reason RNNs are considered to

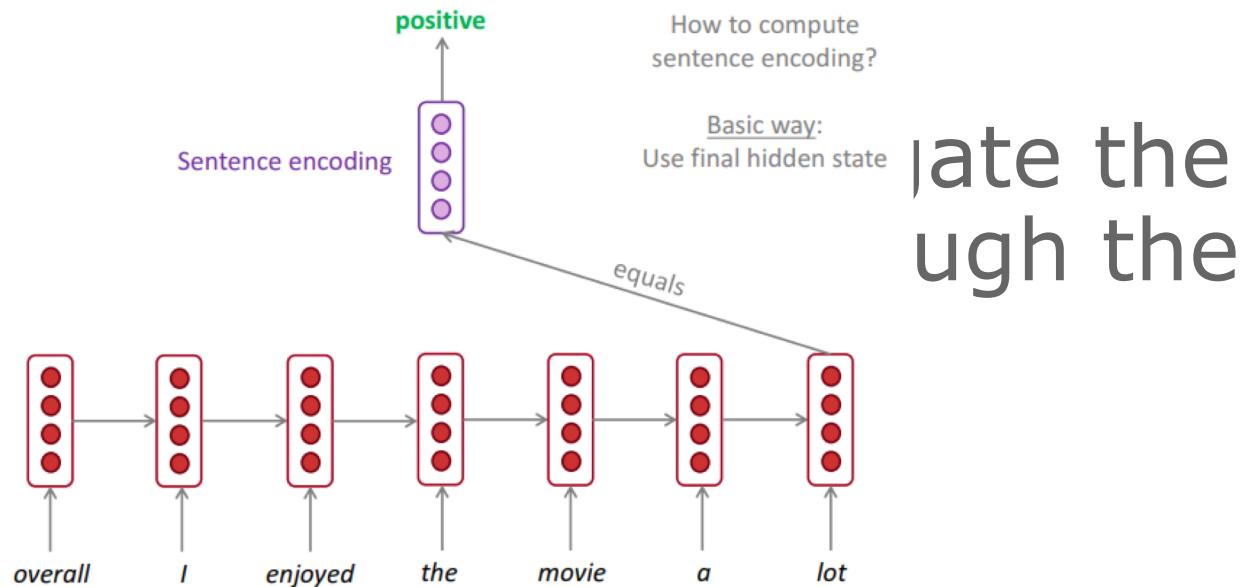
RNNs for word labeling

- .RNNs are very flexible and can be customized for different tasks
- .Below: RNN for word labeling
 - i.e., classify every element of a sequence



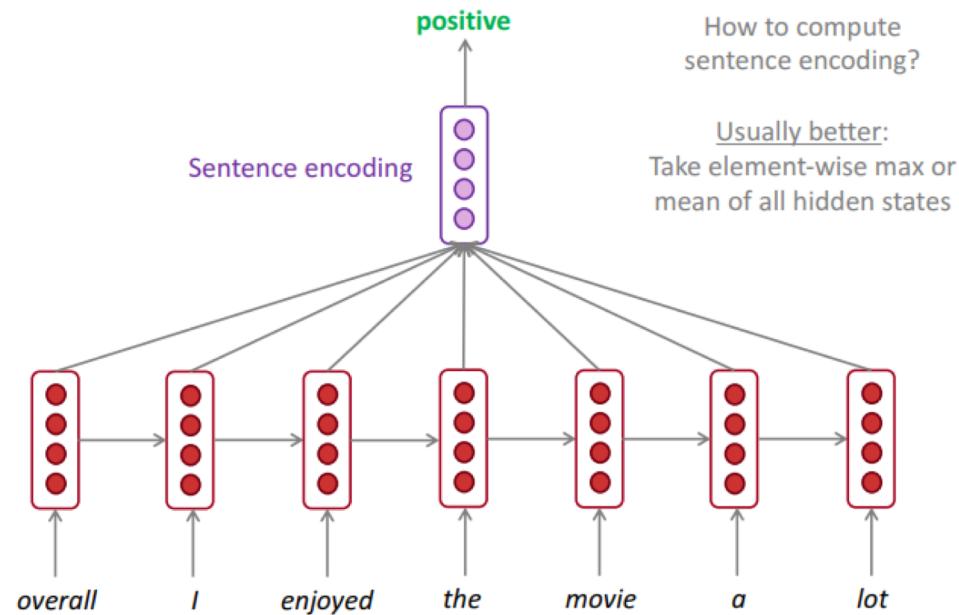
RNN for sequence labeling

- .RNNs can be used for labeling the whole sequence
- E.g., for classifying sentences/documents
- .Simplest: just use the final hidden state for classification
- The RNN most important sentence



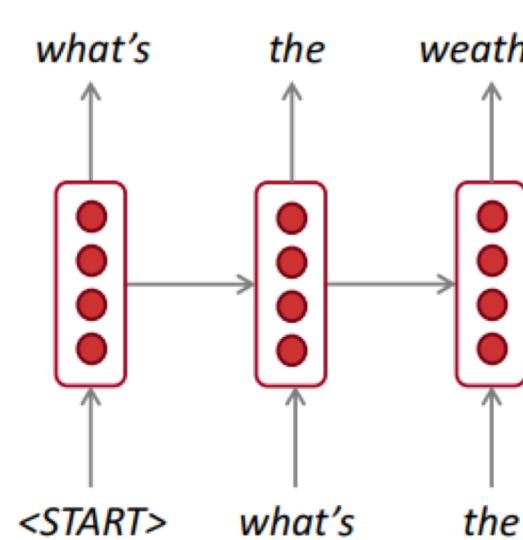
RNN for sequence labeling

- .Alternative: average or max the hidden states over the whole sequence
- Even better: attention mechanism – learn a function that weights the hidden states appropriately (e.g. downweight the function)



RNN for generating text

- .RNN can be used for generating text
 - Start with some random initial state, use a pseudo-word <START> as the initial input
 - Generate a word (or character)
Sample it from a distribution word, with the RNN, and
 - Feed it back into the RNN, and generate another word



RNN text generation samples I

.Shakespeare:

KING LEAR:

O, if you were a feeble sight, the courtesy of your law,

Your sight and several breath, will wear the gods

~~With his heads; and my hands are wonder'd at the Naturalism and decision for the majority of Arab~~
~~countries' capitalide was grounded by the Irish~~
~~so drop upon your lordship's head, and your opinion language by [[John Blair]] [[An Imperial Japanese~~
~~Shall be against your honour.~~

Revolt]], associated with Guangzham's sovereignty.

His generals were the powerful ruler of the Portugal in the [[Protestant Immineners]], which could be

RNN generated text, II

.PDF compiled from generated LaTeX

S
r

For $\bigoplus_{n=1,\dots,m} \mathcal{L}_{m,n} = 0$, hence we can find a closed subset \mathcal{H} in \mathcal{H} and any sets \mathcal{F} on X , U is a closed immersion of S , then $U \rightarrow T$ is a separated algebraic space.

Proof. Proof of (1). It also start we get

$$S = \text{Spec}(R) = U \times_X U \times_X U$$

and the comparicoly in the fibre product covering we have to prove the lemma generated by $\coprod Z \times_U U \rightarrow V$. Consider the maps M along the set of points Sch_{fppf} and $U \rightarrow U$ is the fibre category of S in U in Section, ?? and the fact that any U affine, see Morphisms, Lemma ???. Hence we obtain a scheme S and any open subset $W \subset U$ in $Sh(G)$ such that $\text{Spec}(R') \rightarrow S$ is smooth or an

$$U = \bigcup U_i \times_{S_i} U_i$$

which has a nonzero morphism we may assume that f_i is of finite presentation over S . We claim that $\mathcal{O}_{X,x}$ is a scheme where $x, x', s'' \in S'$ such that $\mathcal{O}_{X,x'} \rightarrow \mathcal{O}'_{X',x'}$ is separated. By Algebra, Lemma ?? we can define a map of complexes $GL_{S'}(x'/S'')$ and we win. \square

To prove study we see that $\mathcal{F}|_U$ is a covering of X' , and \mathcal{T}_i is an object of $\mathcal{F}_{X/S}$ for $i > 0$ and \mathcal{F}_p exists and let \mathcal{F}_i be a presheaf of \mathcal{O}_X -modules on C as a \mathcal{F} -module. In particular $\mathcal{F} = U/\mathcal{F}$ we have to show that

$$\widetilde{M}^\bullet = \mathcal{I}^\bullet \otimes_{\text{Spec}(k)} \mathcal{O}_{S,s} - i_X^{-1} \mathcal{F}$$

is a unique morphism of algebraic stacks. Note that

$$\text{Arrows} = (Sch/S)_{fppf}^{\text{opp}}, (Sch/S)_{fppf}$$

and

$$V = \Gamma(S, \mathcal{O}) \longmapsto (U, \text{Spec}(A))$$

is an open subset of X . Thus U is affine. This is a continuous map of X is the inverse, the groupoid scheme S .

Proof. See discussion of sheaves of sets. \square

The result for prove any open covering follows from the less of Example ???. It may replace S by $X_{\text{spaces},\text{étale}}$ which gives an open subspace of X and T equal to S_{Zar} , see Descent, Lemma ???. Namely, by Lemma ?? we see that R is geometrically regular over S .

Lemma 0.1. Assume (3) and (3) by the construction in the description.

Suppose $X = \lim |X|$ (by the formal open covering X and a single map $\underline{\text{Proj}}_X(\mathcal{A}) = \text{Spec}(B)$ over U compatible with the complex

$$\text{Set}(\mathcal{A}) = \Gamma(X, \mathcal{O}_{X,\mathcal{O}_X}).$$

When in this case of to show that $\mathcal{Q} \rightarrow \mathcal{C}_{Z/X}$ is stable under the following result in the second conditions of (1), and (3). This finishes the proof. By Definition ?? (without element is when the closed subschemes are catenary. If T is surjective we may assume that T is connected with residue fields of S . Moreover there exists a closed subspace $Z \subset X$ of X where U in X' is proper (some defining as a closed subset of the uniqueness it suffices to check the fact that the following theorem

(1) f is locally of finite type. Since $S = \text{Spec}(R)$ and $Y = \text{Spec}(R)$.

Proof. This is form all sheaves of sheaves on X . But given a scheme U and a surjective étale morphism $U \rightarrow X$. Let $U \cap U = \coprod_{i=1,\dots,n} U_i$ be the scheme X over S at the schemes $X_i \rightarrow X$ and $U = \lim_i X_i$. \square

The following lemma surjective restrocomposes of this implies that $\mathcal{F}_{x_0} = \mathcal{F}_{x_0} = \mathcal{F}_{x,\dots,0}$.

Lemma 0.2. Let X be a locally Noetherian scheme over S , $E = \mathcal{F}_{X/S}$. Set $\mathcal{I} = \mathcal{J}_1 \subset \mathcal{I}'_n$. Since $\mathcal{I}'_n \subset \mathcal{I}^n$ are nonzero over $i_0 \leq p$ is a subset of $\mathcal{J}_{n,0} \circ \mathcal{A}_2$ works.

Lemma 0.3. In Situation ???. Hence we may assume $q' = 0$.

Proof. We will use the property we see that p is the next functor (??). On the other hand, by Lemma ?? we see that

$$D(\mathcal{O}_{X'}) = \mathcal{O}_X(D)$$

where K is an F -algebra where δ_{n+1} is a scheme over S . \square

RNN generated text, III

.Generated baby names, based on 8000 actual names (only generated names not in the training data are listed):

*Rudi Levette
Berice Lussa Hany Mareanne Chrestina Carissy
Marylen Hammie Janye Marlise Jacacie Hendred
Romand Charienna Nenotto Ette Dorane Wallen
Marly Darine Salina Elvyn Ersia Maralena Minoria
Ellia Charmin Antley Nerille Chelon Walmor Evena
Jeryly Stachon Charisa Allisa Anatha Cathanie
Geetra Alexie Jerin Cassen Herbett Cossie Velen
Daurenge Robester Shermond Terisa Licia Roselen
Ferine Jayn Lusine Charvanne Sales Sanny Resa*

RNN for conditional text generation

- .As before, but the initial state encodes some input (e.g. an image)
 - E.g., use a convolutional neural network to reduce the image to a fixed length feature vector

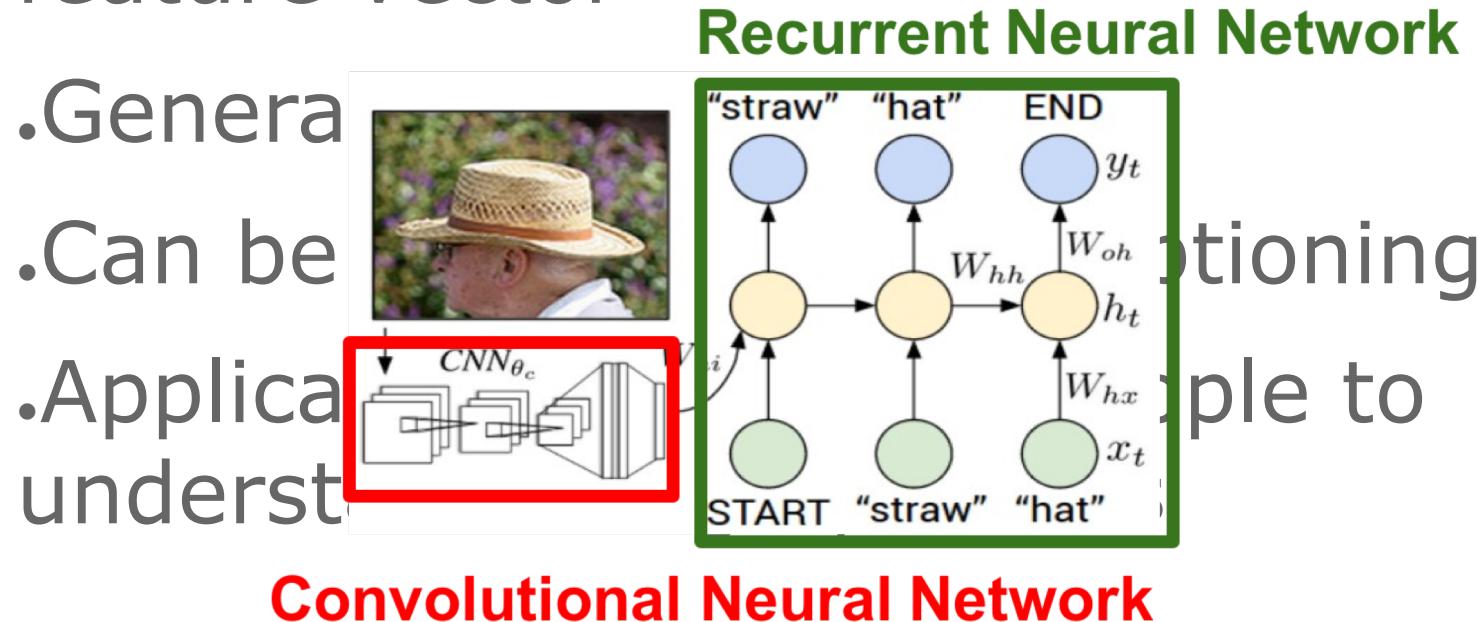


Image captioning examples



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."



"boy is doing backflip on wakeboard."



"a young boy is holding a baseball bat."



"a cat is sitting on a couch with a remote control."

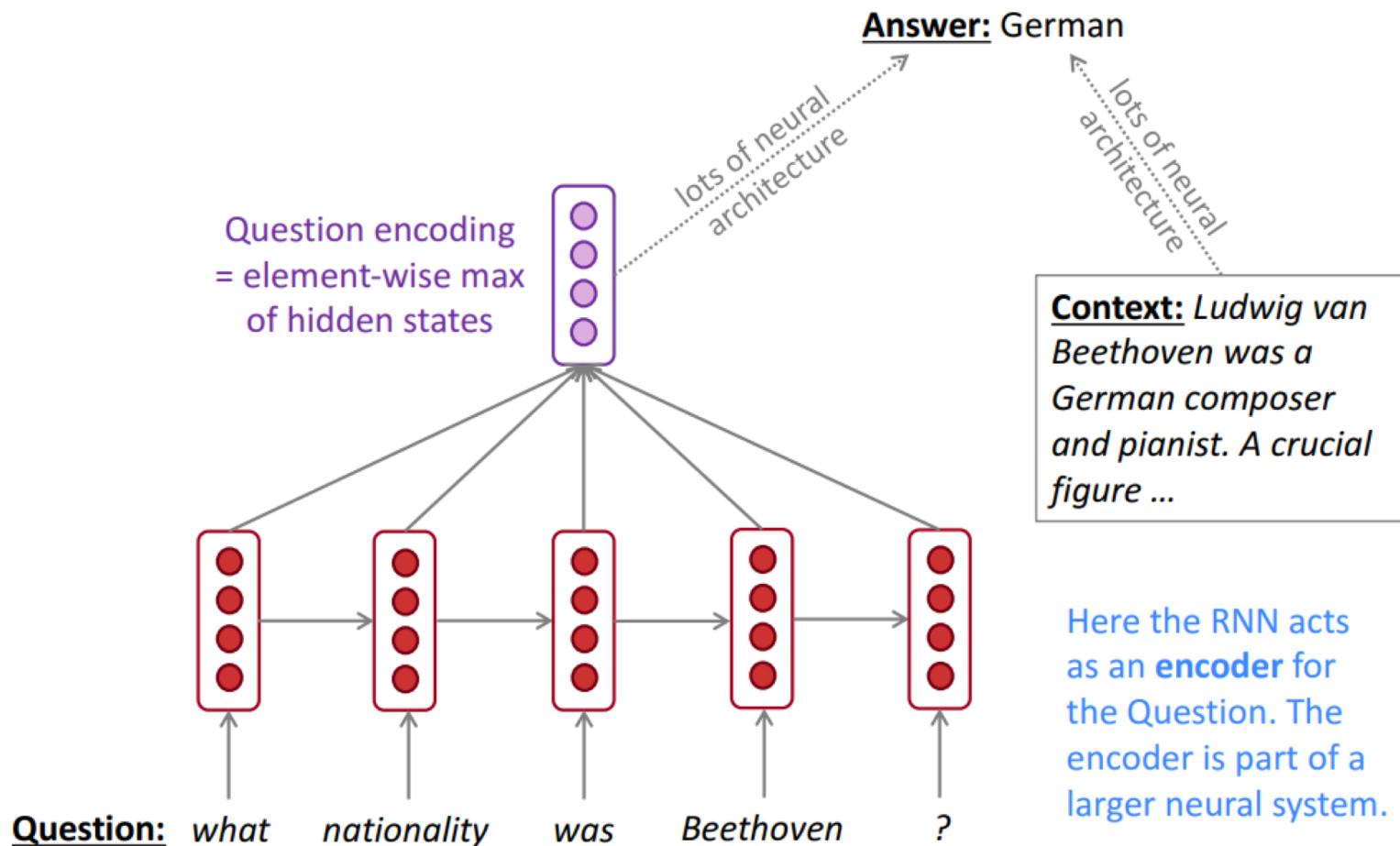


"a woman holding a teddy bear in front of a mirror."



"a horse is standing in the middle of a road."

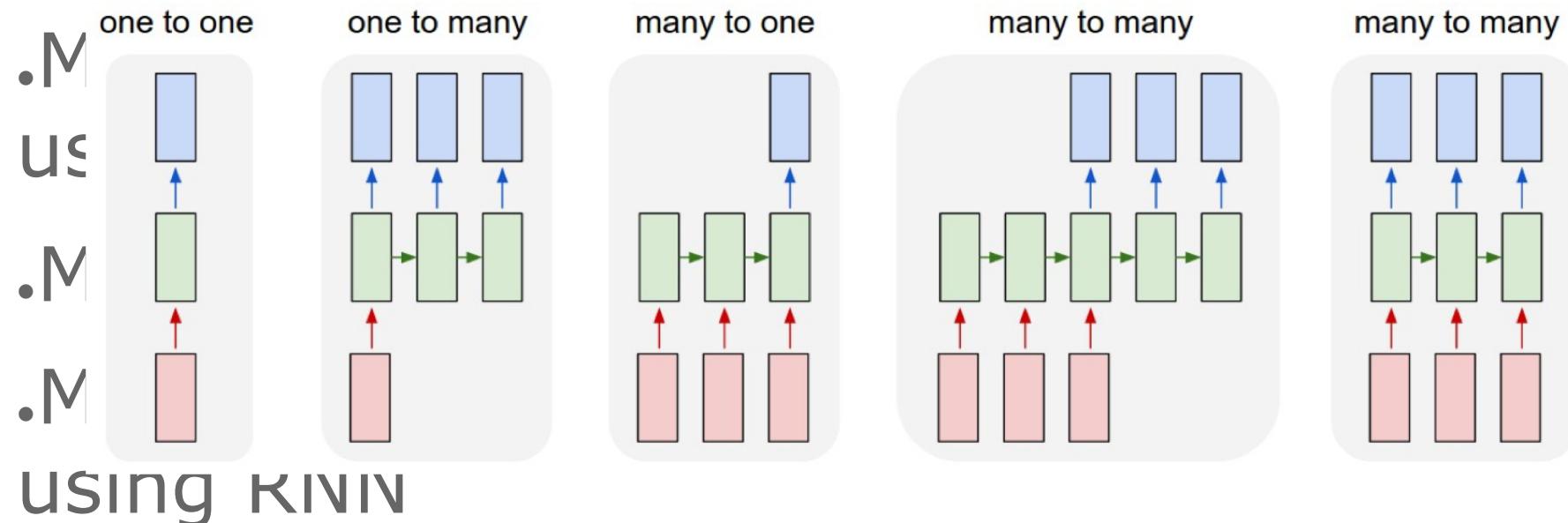
RNN as encoder module



Summary of neural network architectures

.One-to-one: feed-forward DNN we have studied previously (e.g. document classification with CNN, named entity recognition)

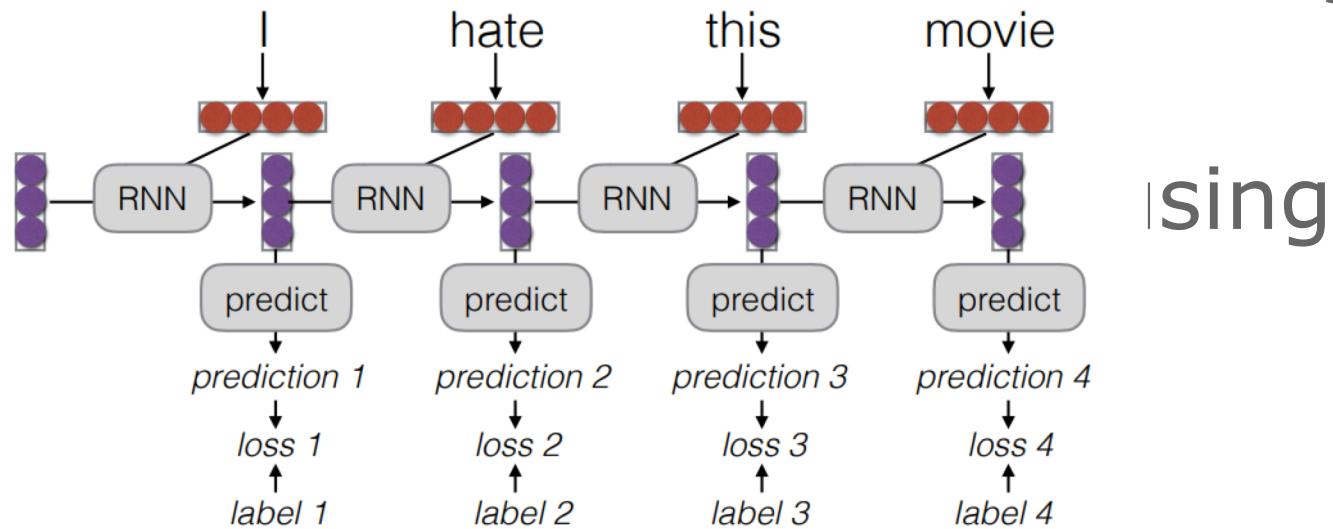
.One-to-many: image captioning



Training RNN

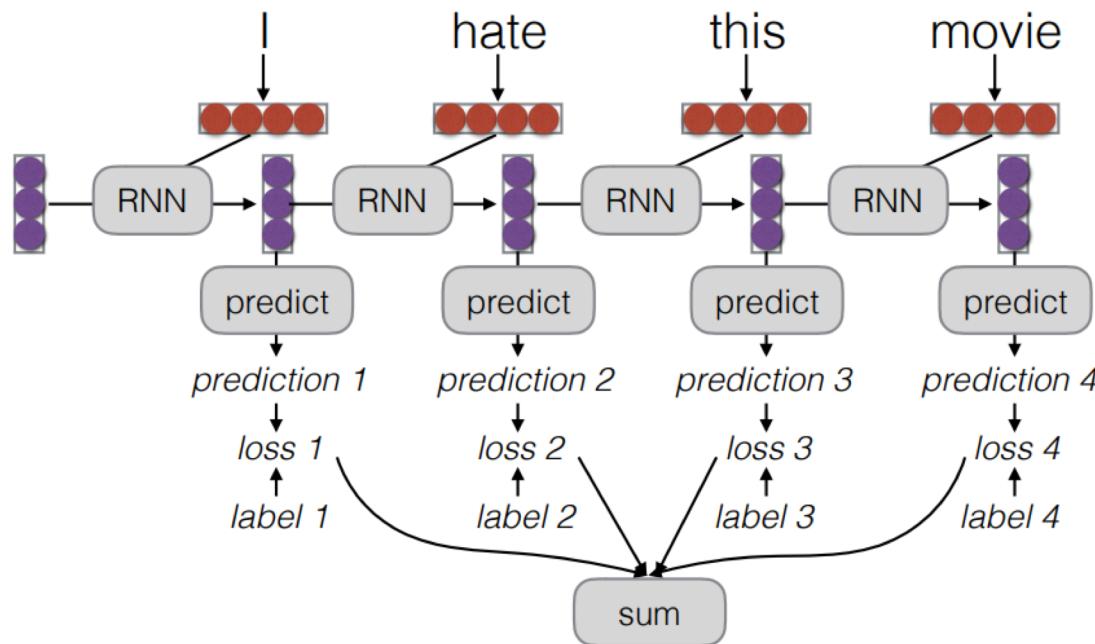
- .Training RNN can be done using backpropagation, as with feed-forward NNs
- .As before, we (or our toolkit) need to have a ~~computation graph~~ connecting the RNN

.Training
an unro



Training RNN

- Training RNN can be done using backpropagation, as with feed-forward NNs
- As before, we (or our toolkit) need to have a context vector c for the RNN.
- Training an unrolled RNN using backpropagation involves the following steps:
 - The RNN processes each word in the sentence sequentially.
 - For each word, the RNN produces a hidden state vector (represented by red dots).
 - The RNN then makes a prediction (represented by purple dots) based on the current hidden state.
 - The predicted word is compared against the actual word (label), and a loss is calculated for each prediction.
 - The losses are summed to provide a total training loss for the entire sentence.

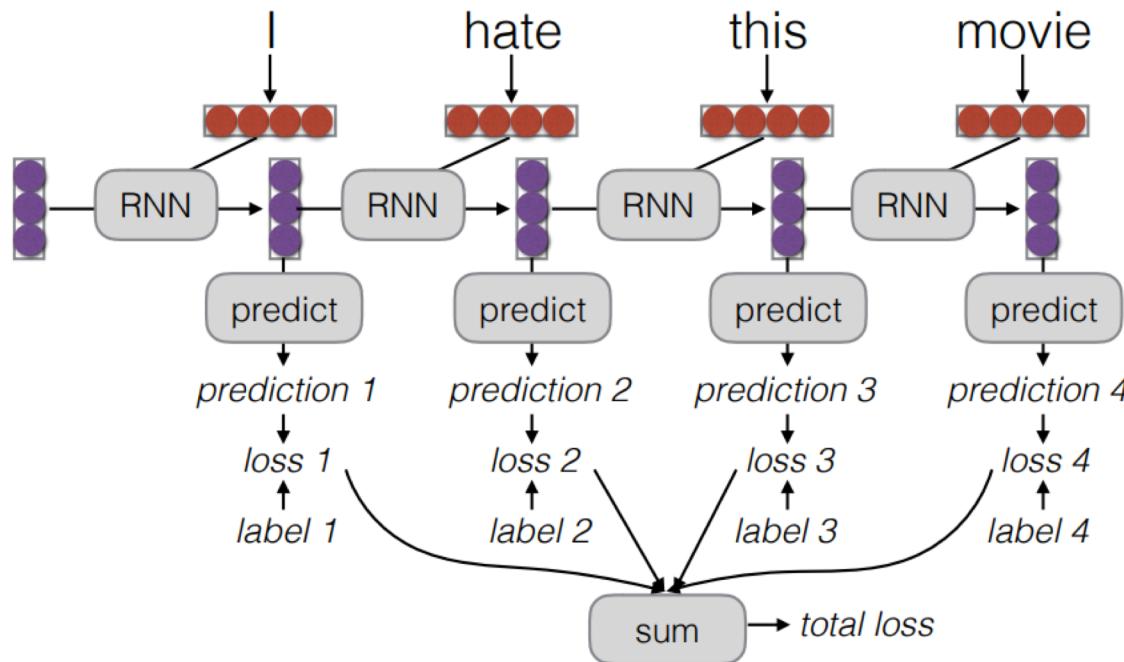


Training RNN

. Training RNN can be done using backpropagation, as with feed-forward NNs

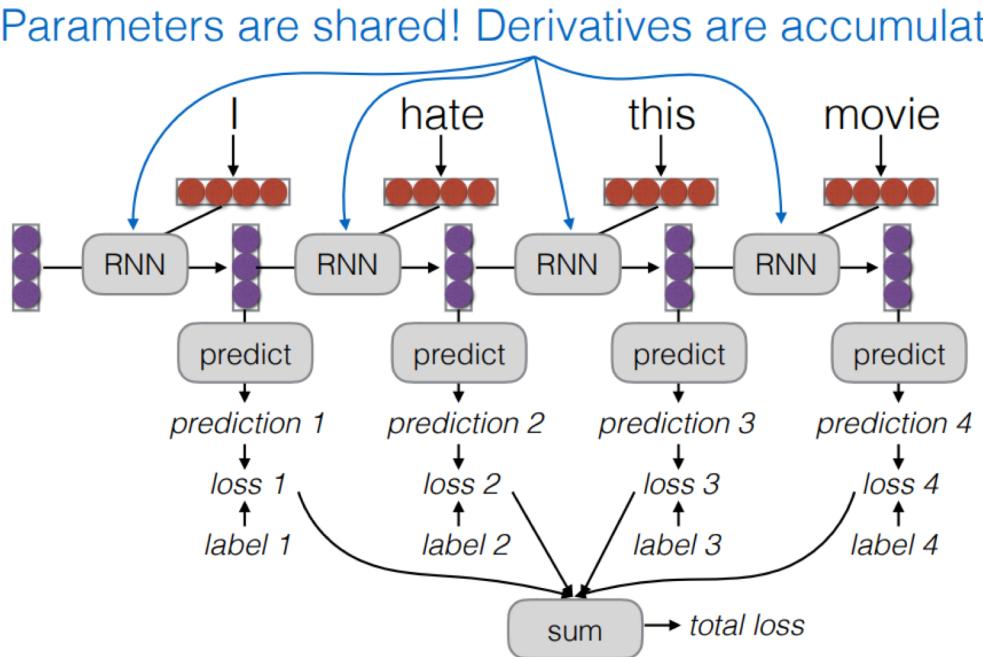
. And before we (or our toolkit) need to have a look at how to train the RNN

. Training an unrolled RNN



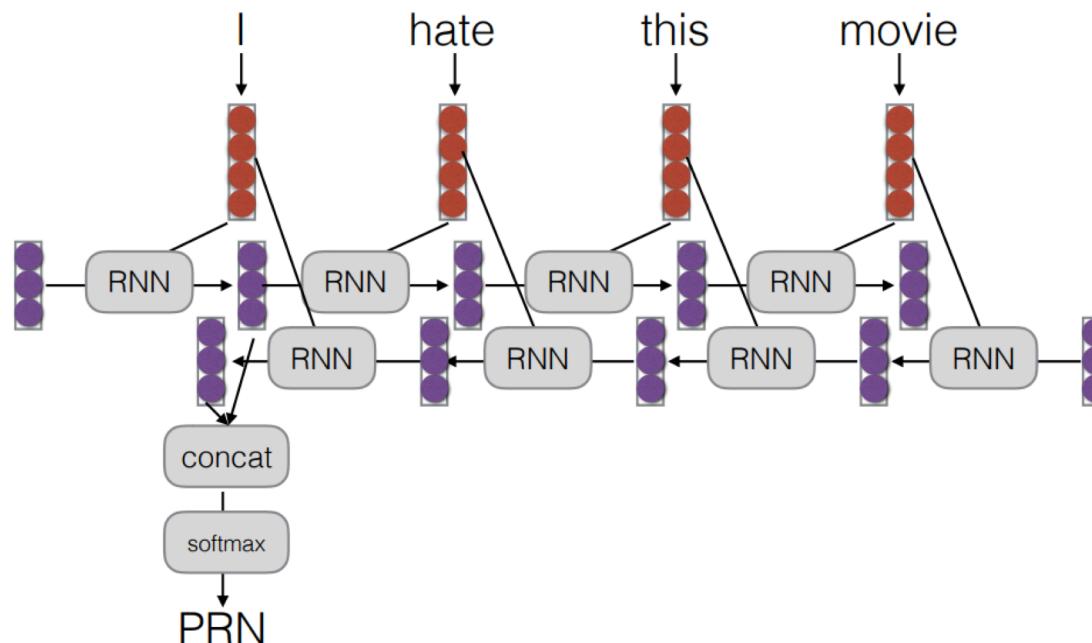
Training RNN

- . Training RNN can be done using backpropagation, as with feed-forward NNs
- . And before, we (or our toolkit) need to have a look at how the RNN is trained.
- . Training an unrolled RNN using backpropagation



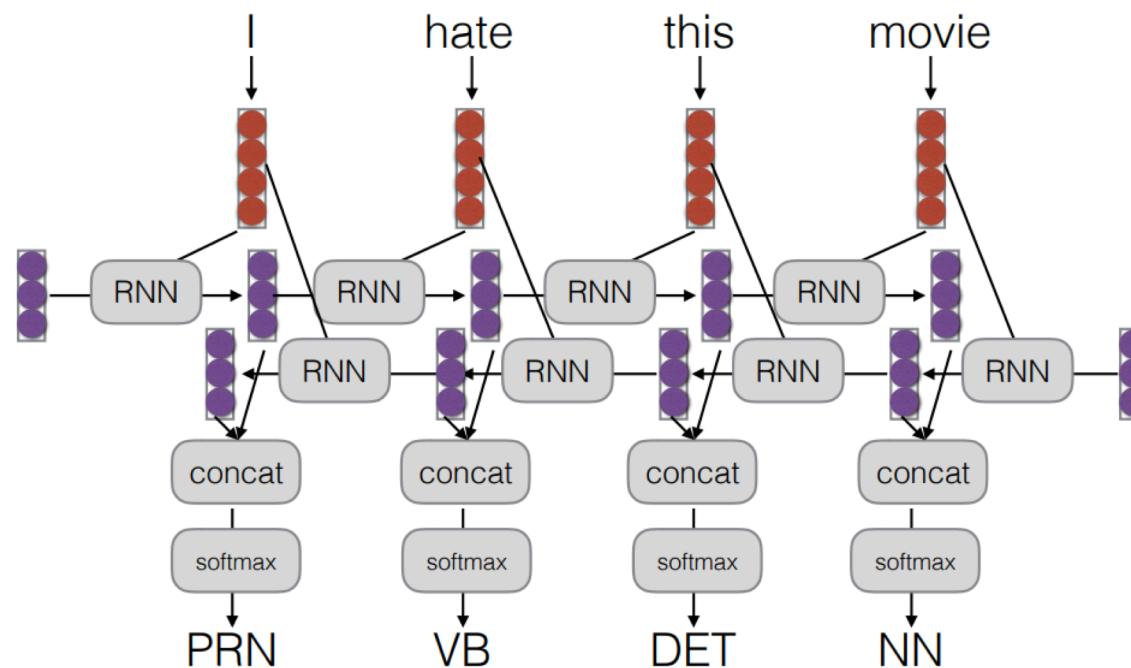
Bidirectional RNNs

- .Simple extension: have two RNNs, running in both directions
- .The hidden states from the two RNNs are simply concatenated



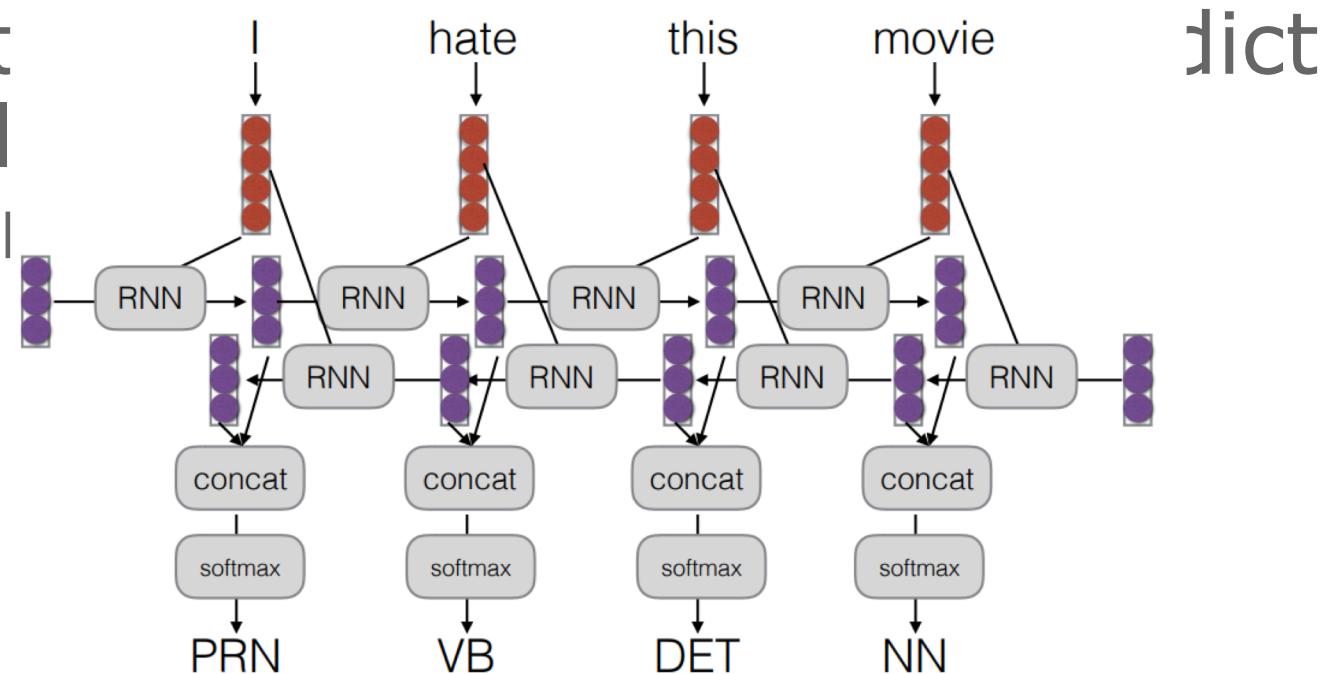
Bidirectional RNNs

- .Simple extension: have two RNNs, running in both directions
- .The hidden states from the two RNNs are simply concatenated



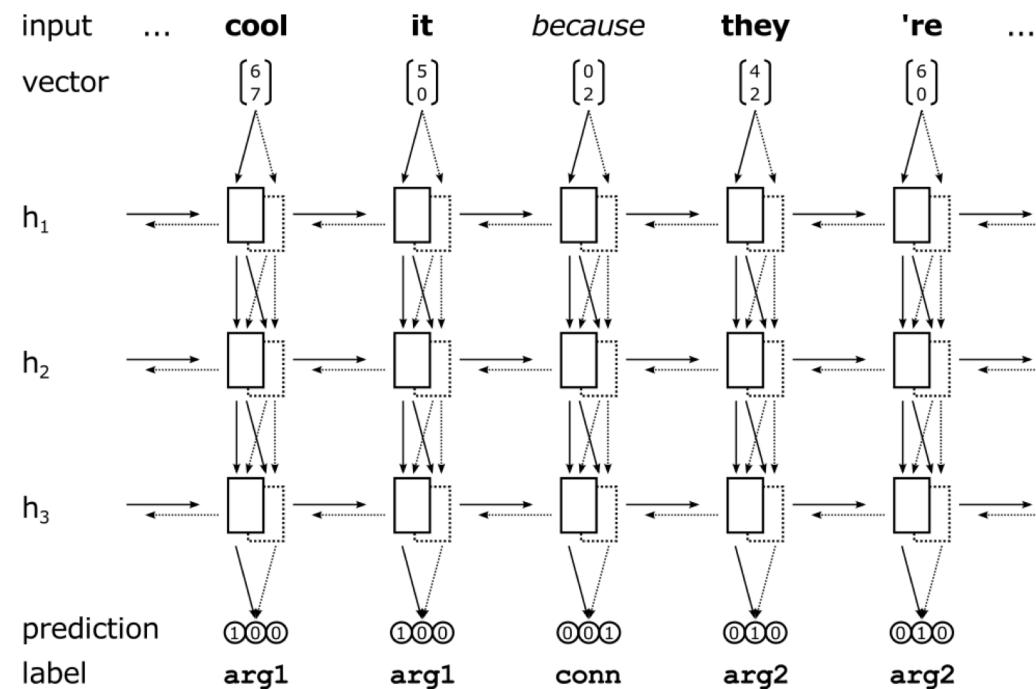
Bidirectional RNNs

- Why?
 - Word classification: enables to use information from both past and future words
 - Alternative to unidirectional RNNs: the label consumer



Deep RNNs

- .We can stack many recurrent layers
- .Each subsequent layer uses the output of the previous recurrent layer as input

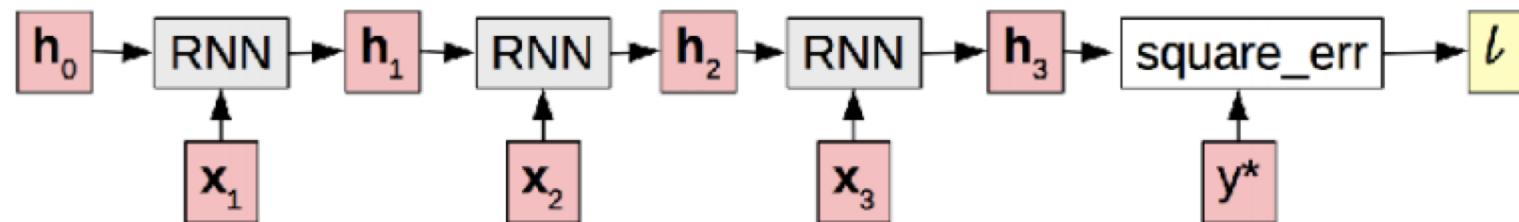


Vanishing Gradients

.Result: the RNN doesn't learn to use information that is **many time steps in the past**

.Or, (in the exploding case), the parameters of the RNN get quickly very large

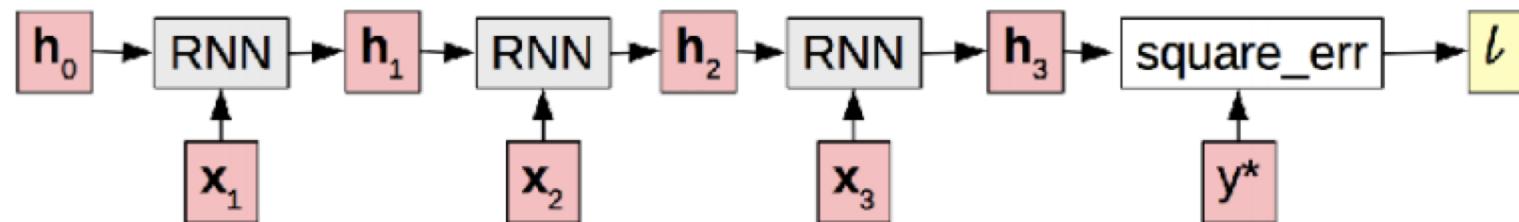
large $\frac{dl}{dh_0}$ =tiny $\frac{dl}{dh_1}$ =small $\frac{dl}{dh_2}$ =med. $\frac{dl}{dh_3}$ =large



Vanishing Gradients

- .Gradients decrease (vanish) or increase (explode) when they get pushed back through the RNN
- .They get „squashed” through non-linearities and small (or large) weight

$$\mathcal{M} \quad \frac{dl}{d_{h_0}} = \text{tiny} \quad \frac{dl}{d_{h_1}} = \text{small} \quad \frac{dl}{d_{h_2}} = \text{med.} \quad \frac{dl}{d_{h_3}} = \text{large}$$



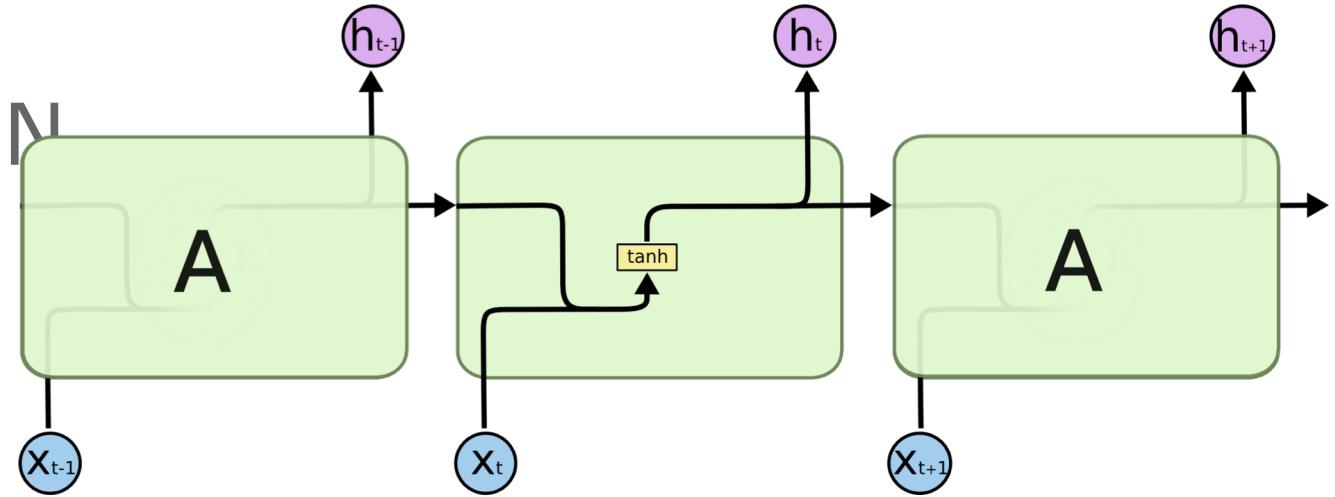
Solution to vanishing gradients: gated units

- .The main solution to the vanishing gradient problem is to use a more complex hidden unit in the recurrent cell
- .Gated Recurrent Unit (**GRU**) and Long Short-Term Memory (**LSTM**)
- .Main ideas:
 - Make additive connections between time steps (not multiplicative)
 - .Addition does not modify the gradient -- no vanishing

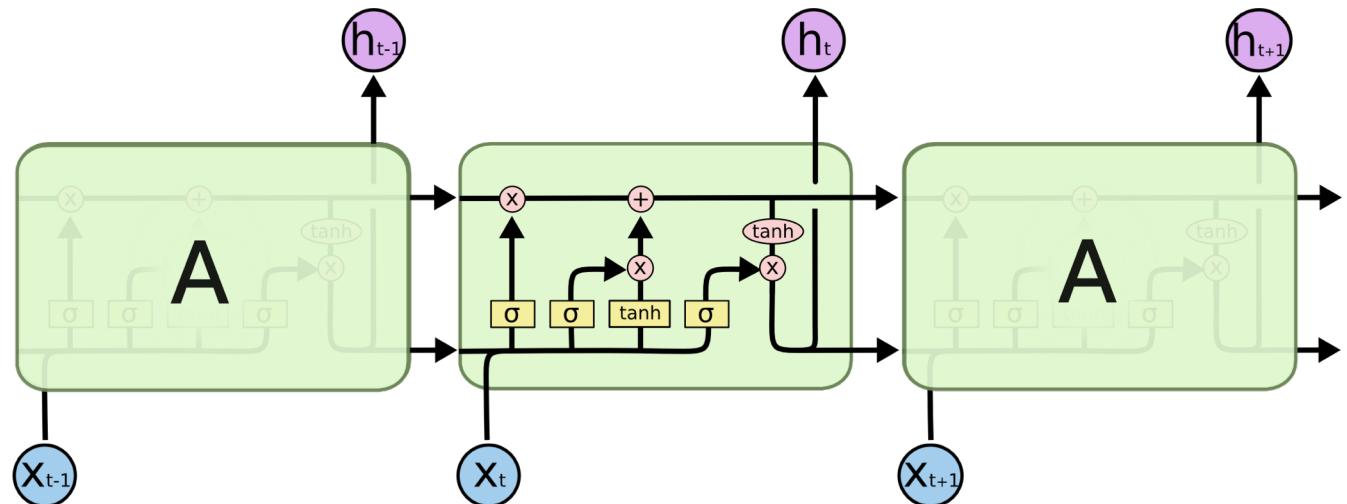
.Standard RNN

.

LSTM

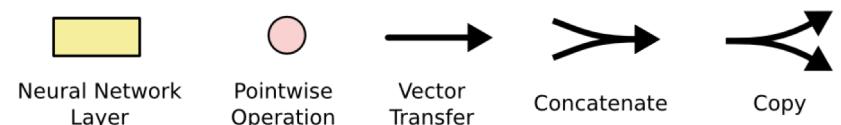


.LSTM



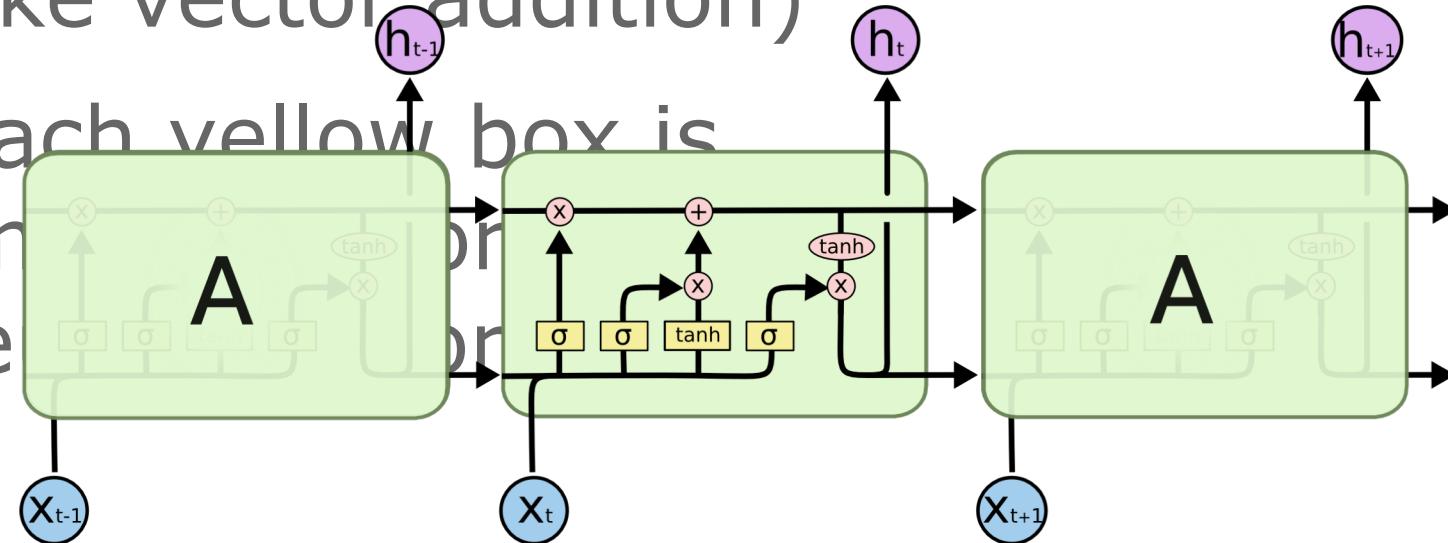
LSTM

.Each line carries a vector



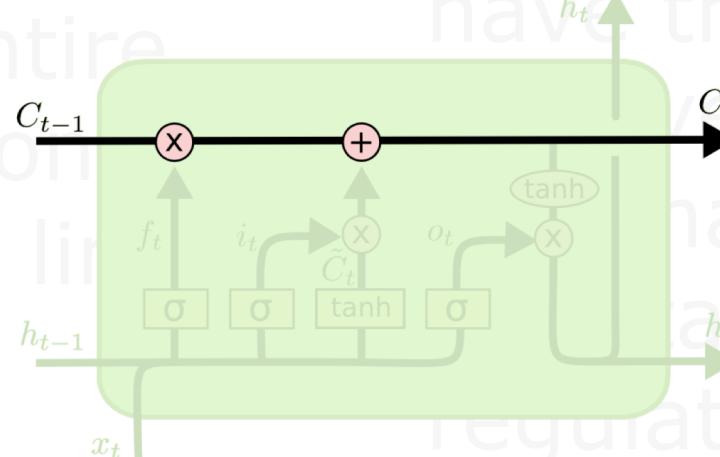
.Pink circle is a pointwise operation (like vector addition)

.Each yellow box is a small neural net



LSTM walkthrough

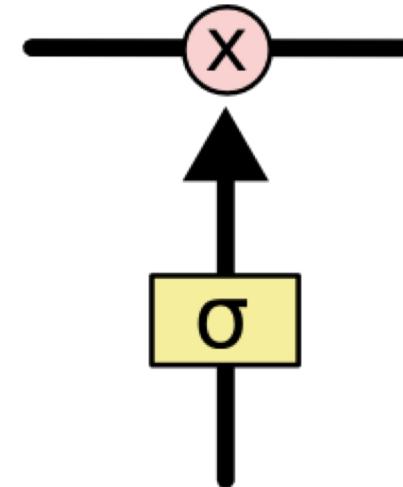
- .Each LSTM cell has a **state**
- .Cell state is kind of like a conveyor belt
- .It runs straight down the entire chain, with only some minor lateral interactions
- .It's very easy for information to just flow along it unchanged
- .The LSTM does have the ability to gate or add information to the state, carefully regulated by structures called



structures called

LSTM walkthrough: gates

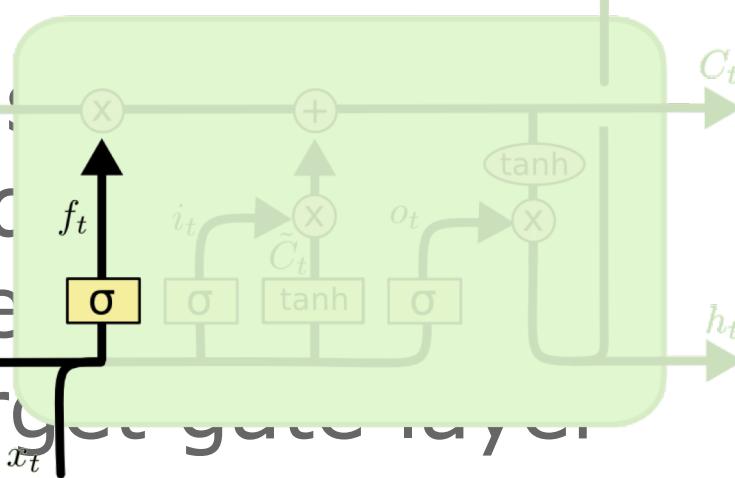
- .Gates are a way to optionally let information through
- .They are composed out of a sigmoid
- .The sigmoid layer outputs numbers between zero and one, describing how much of each component should be let through operation
- .A value of zero means “let nothing through,” while a value of one means “let everything through!”



LSTM walkthrough: forget gate

.The first step in our LSTM is to decide what information we're going to throw away from the cell state

.This makes a layer "forget gate layer".



.It looks at h_{t-1} and x_t , and outputs a number between 0 and 1

.0 – forget the cell state

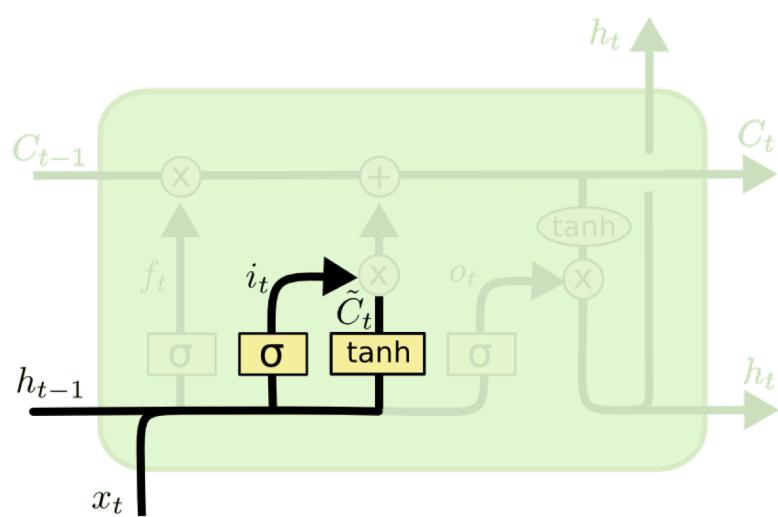
1 – keep the cell state

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

LSTM walkthrough: input gate layer

.The next step is to decide what new information we're going to store in the cell state

.First, a sigmoid layer called the "input gate layer" decides which values we'll update

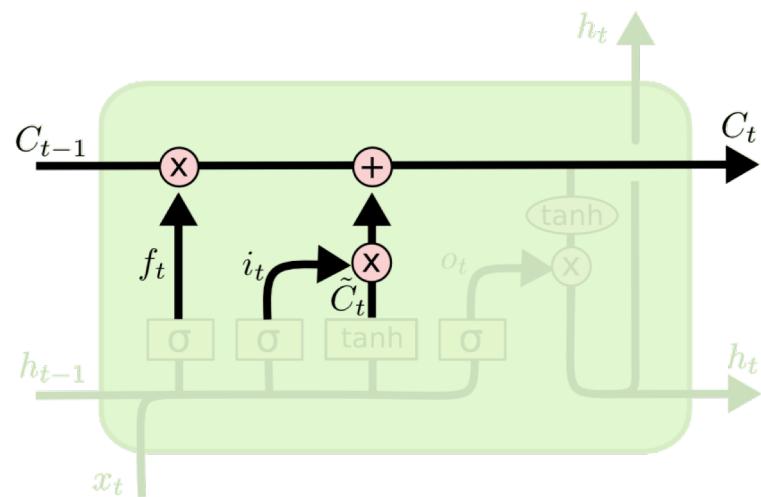


.Next, a tanh layer creates a vector of new candidate values \hat{C}_t that could be added to the state.

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$
$$\hat{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

LSTM walkthrough: update cell state

.We multiply the old state by f_t , forgetting the things we decided to forget earlier



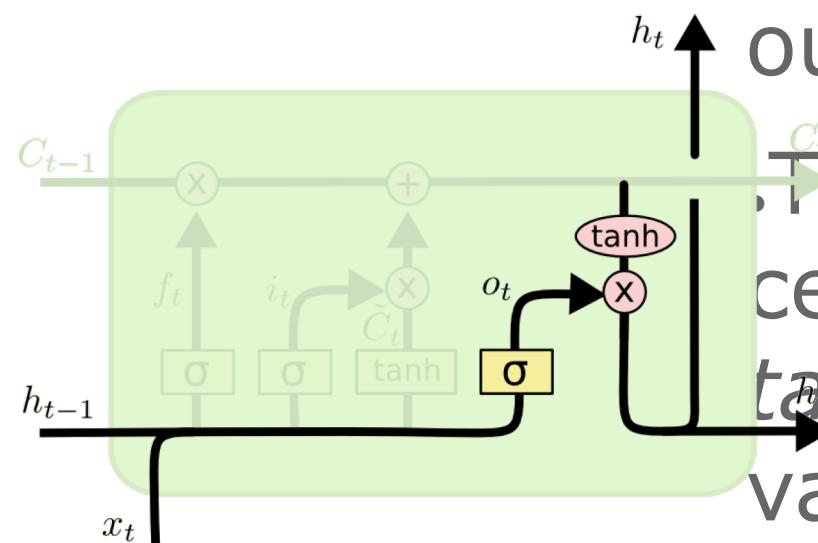
.Then we add $i_t * \hat{C}_t$
.This is the new candidate values, scaled by how much we decided to update each state value.

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

LSTM walkthrough: compute output

.Output will be based on our cell state, but will be a filtered version

.First, we run a sigmoid layer which decides what parts of the cell state we're going to output



.Then, we put the cell state through $tanh$ (to push the values to be between -1 and 1)
$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$
$$h_t = o_t * \tanh(C_t)$$

LSTM vs RNN

- .When we talk about remarkable results achieved with RNNs, we really usually talk about LSTMs
- .They really work better than raw RNNs for most things
- .GRU is similar to LSTM, but a bit simpler
- .GRU might work better than LSTM if your training data is rather limited (and it's faster)
- .LSTM formulas seem intimidating

RNN strengths and weaknesses

- .RNNs, particularly LSTMs have much more modeling power than feed-forward NNs
- .However, they can take time to tune
 - Often, initial experiments produce disappointing results
 - Sensitive to learning rate, model initialization, scale of the features, etc
 - Tuning them is even deeper black magic than for feed-forward NNs

They also require a lot of training data

Teacher-student training

. Sometimes, teacher-student training can help to learn a good LSTM

. Idea:

- Train a feed-forward NN (like CNN) for the task (e.g., for sentiment analysis)

- Now, train the LSTM, but instead of using real labels as targets, use **soft labels** produced by the feed-forward NN (i.e., the probability distribution predicted by the simpler model)

- In other words, train the RNN to **mimic**

Pre-training/Transfer

.Idea:

- First train the RNN for the simpler task that has a lot of training data
- Then, use the pretrained RNN as part of a new model for the more complex task

.Example: LM -> Sentence Classifier

- First train a RNN to do language modeling (predicting the next word) – lots of data available
- Then use the trained RNN as a starting point for training a sentence classification

Questions?