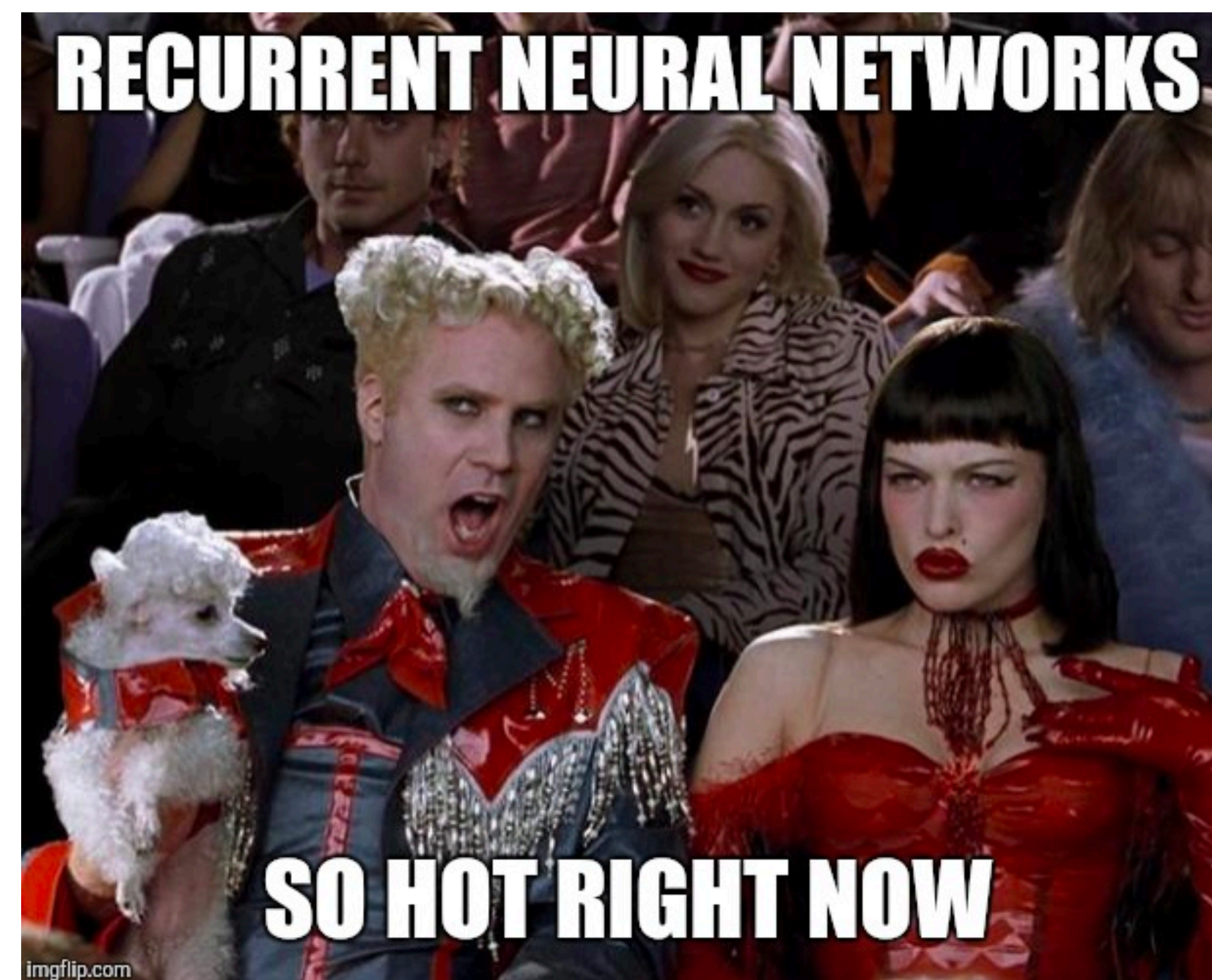


CS388: Natural Language Processing

Lecture 8: RNNs

Greg Durrett



Credit: Chelsea Voss csvoss.com



Administrivia

- ▶ Mini 1 back today
- ▶ Project 1 due tonight
- ▶ Mini 2 out tonight



Recall: Word Vectors

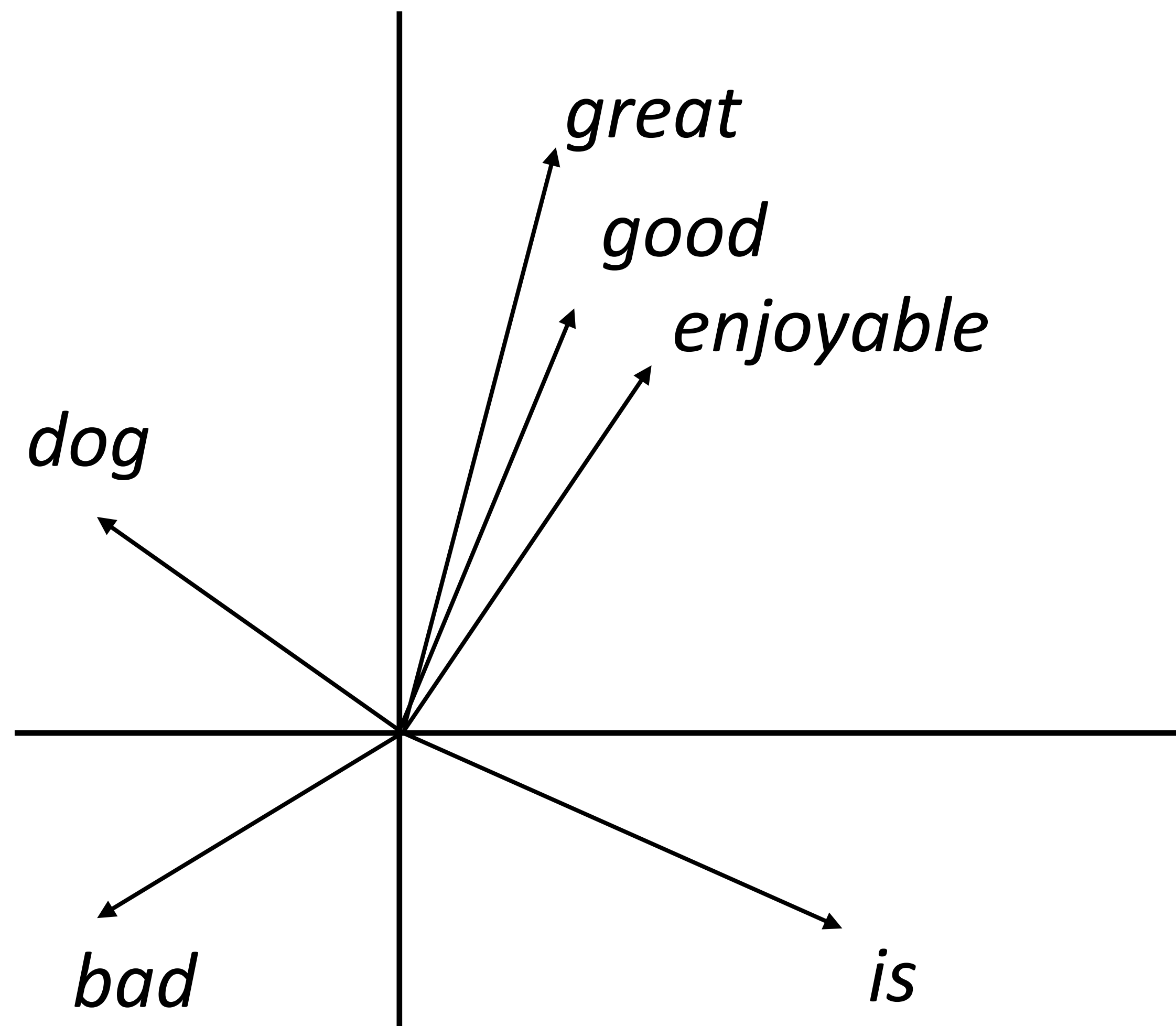
◆ *the president said that the downturn was over* ◆

<i>president</i>	<i>the</i> __ <i>of</i>
<i>president</i>	<i>the</i> __ <i>said</i> ←
<i>governor</i>	<i>the</i> __ <i>of</i>
<i>governor</i>	<i>the</i> __ <i>appointed</i>
<i>said</i>	<i>sources</i> __ ◆
<i>said</i>	<i>president</i> __ <i>that</i>
<i>reported</i>	<i>sources</i> __ ◆

president
governor

said
reported

the
a



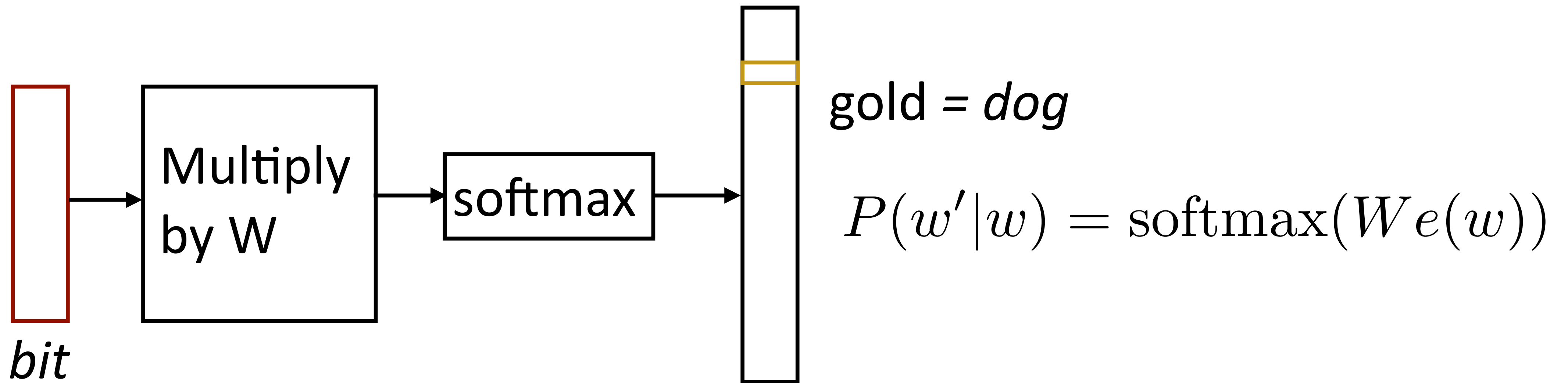
[Finch and Chater 92, Shuetze 93, many others]



Recall: Skip-Gram

- Predict one word of context from word

the *dog* *bit* *the* *man*



- Another training example: *bit* \rightarrow *the*
- Parameters: $d \times |V|$ **vectors**, $|V| \times d$ output parameters (W) (also usable as vectors!)



This Lecture

- ▶ Evaluating word embeddings
- ▶ Recurrent neural networks: basics, issues
- ▶ LSTMs / GRUs
- ▶ Applications / visualizations

Evaluating Word Embeddings



Evaluating Word Embeddings

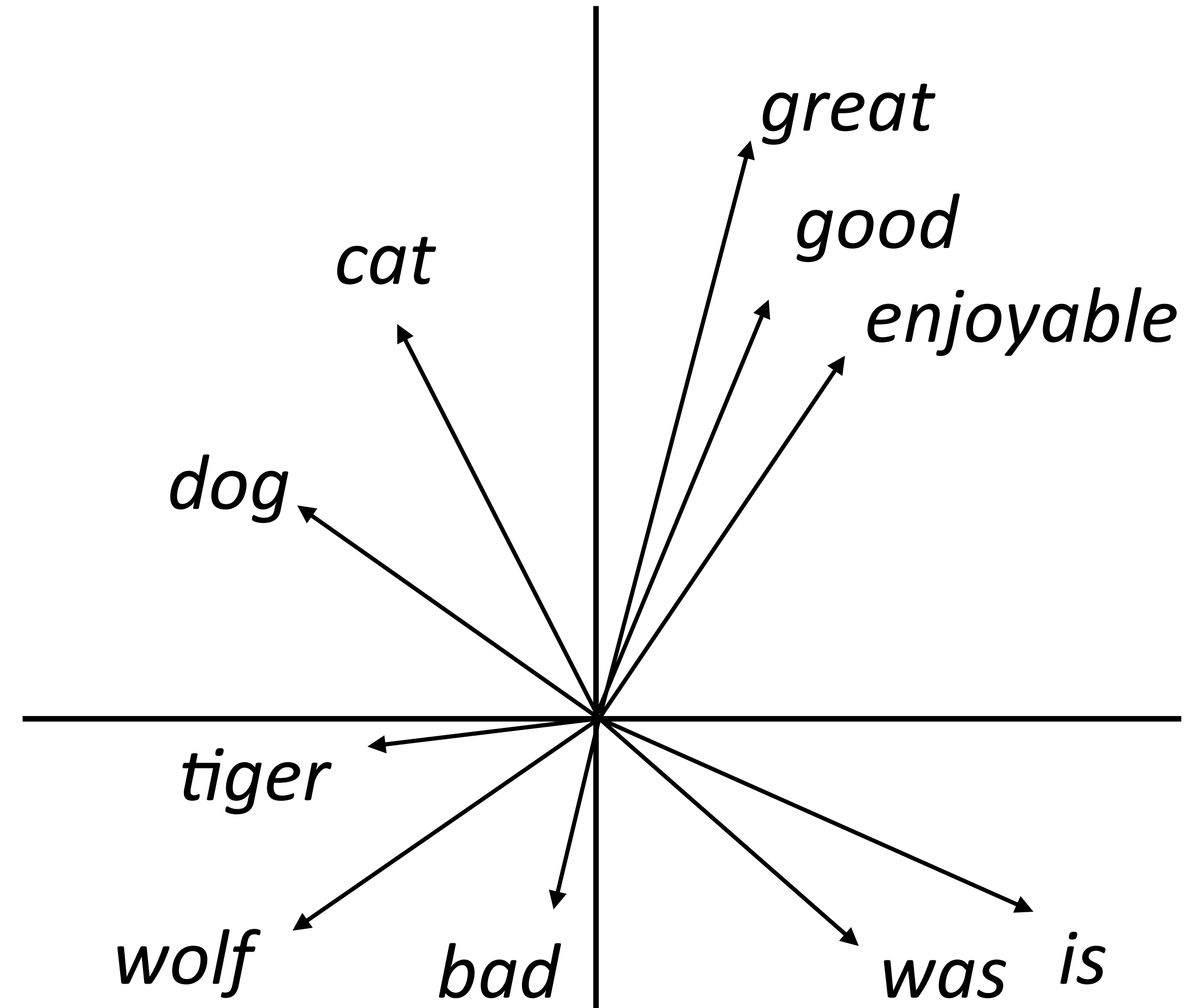
► What properties of language should word embeddings capture?

► Similarity: similar words are close to each other

► Analogy:

good is to best as smart is to ???

Paris is to France as Tokyo is to ???





Similarity

Method	WordSim Similarity	WordSim Relatedness	Bruni et al. MEN	Radinsky et al. M. Turk	Luong et al. Rare Words	Hill et al. SimLex
PPMI	.755	.697	.745	.686	.462	.393
SVD	.793	.691	.778	.666	.514	.432
SGNS	.793	.685	.774	.693	.470	.438
GloVe	.725	.604	.729	.632	.403	.398

- ▶ SVD = singular value decomposition on PMI matrix
- ▶ GloVe does not appear to be the best when experiments are carefully controlled, but it depends on hyperparameters + these distinctions don't matter in practice

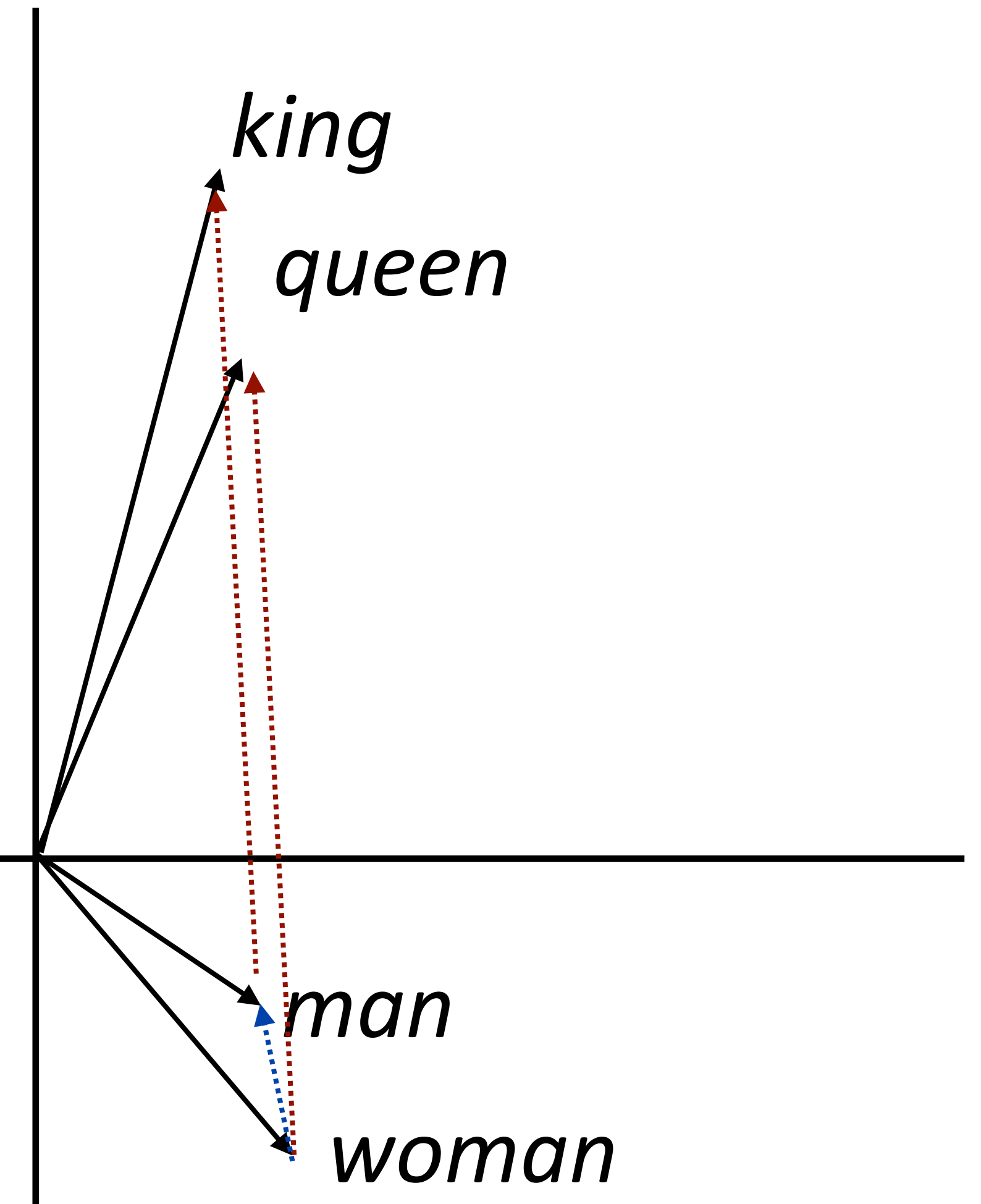


Analogies

$(king - man) + woman = queen$

$king + (woman - man) = queen$

- ▶ Why would this be?
- ▶ woman - man captures the difference in the contexts that these occur in
- ▶ Dominant change: more “he” with man and “she” with woman — similar to difference between king and queen
- ▶ Can evaluate on this as well





What can go wrong with word embeddings?

- ▶ What's wrong with learning a word's "meaning" from its usage?
- ▶ What data are we learning from?
- ▶ What are we going to learn from this data?



What do we mean by bias?

- Identify *she* - *he* axis in word vector space, project words onto this axis

Extreme *she* occupations

- | | | |
|-----------------|-----------------------|------------------------|
| 1. homemaker | 2. nurse | 3. receptionist |
| 4. librarian | 5. socialite | 6. hairdresser |
| 7. nanny | 8. bookkeeper | 9. stylist |
| 10. housekeeper | 11. interior designer | 12. guidance counselor |

Extreme *he* occupations

- | | | |
|----------------|-------------------|----------------|
| 1. maestro | 2. skipper | 3. protege |
| 4. philosopher | 5. captain | 6. architect |
| 7. financier | 8. warrior | 9. broadcaster |
| 10. magician | 11. fighter pilot | 12. boss |

Bolukbasi et al. (2016)

- Nearest neighbor of (b - a + c)

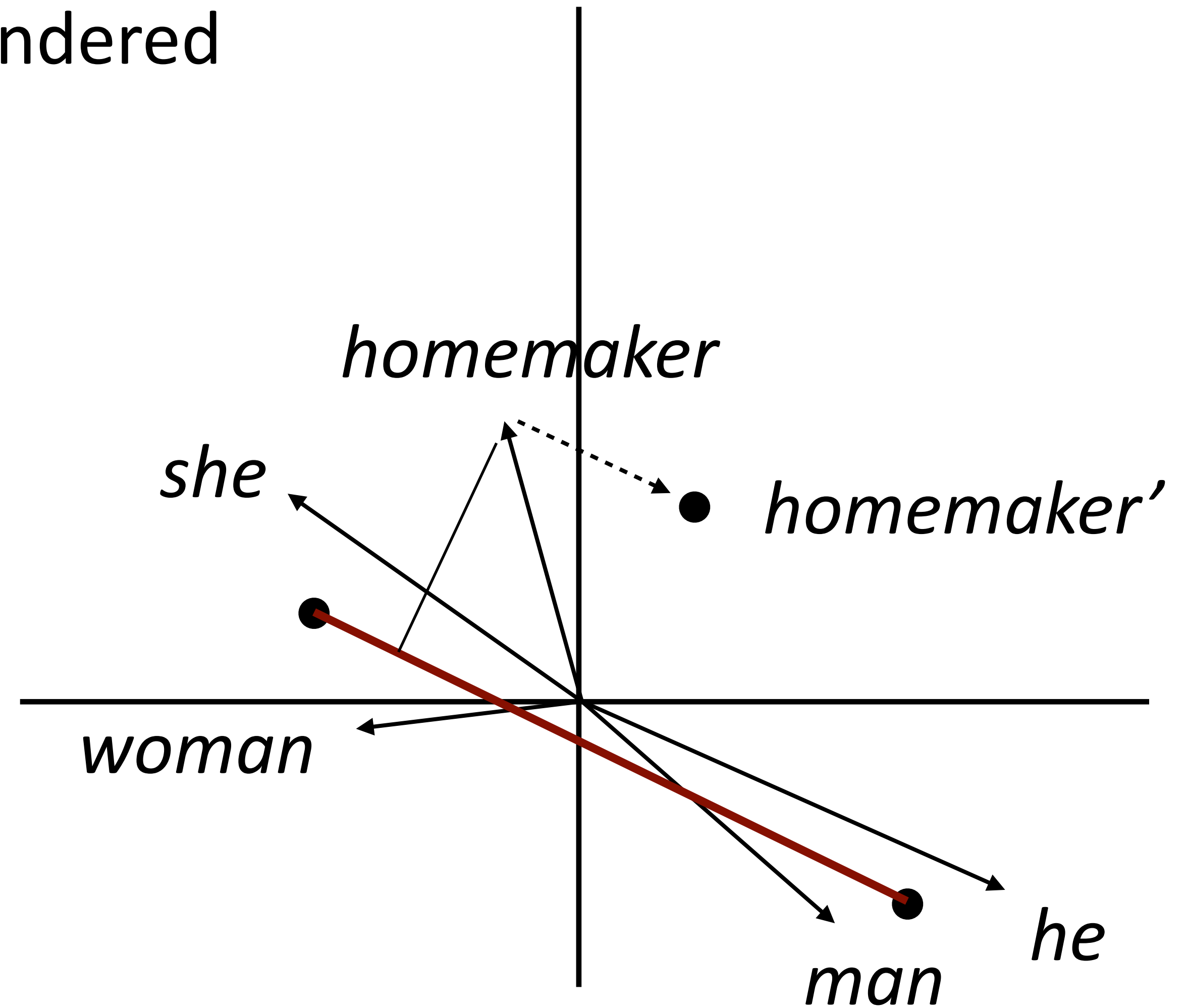
Racial Analogies	
black → homeless	caucasian → servicemen
caucasian → hillbilly	asian → suburban
asian → laborer	black → landowner
Religious Analogies	
jew → greedy	muslim → powerless
christian → familial	muslim → warzone
muslim → uneducated	christian → intellectually

Manzini et al. (2019)



Debiasing

- ▶ Identify gender subspace with gendered words
- ▶ Project words onto this subspace
- ▶ Subtract those projections from the original word

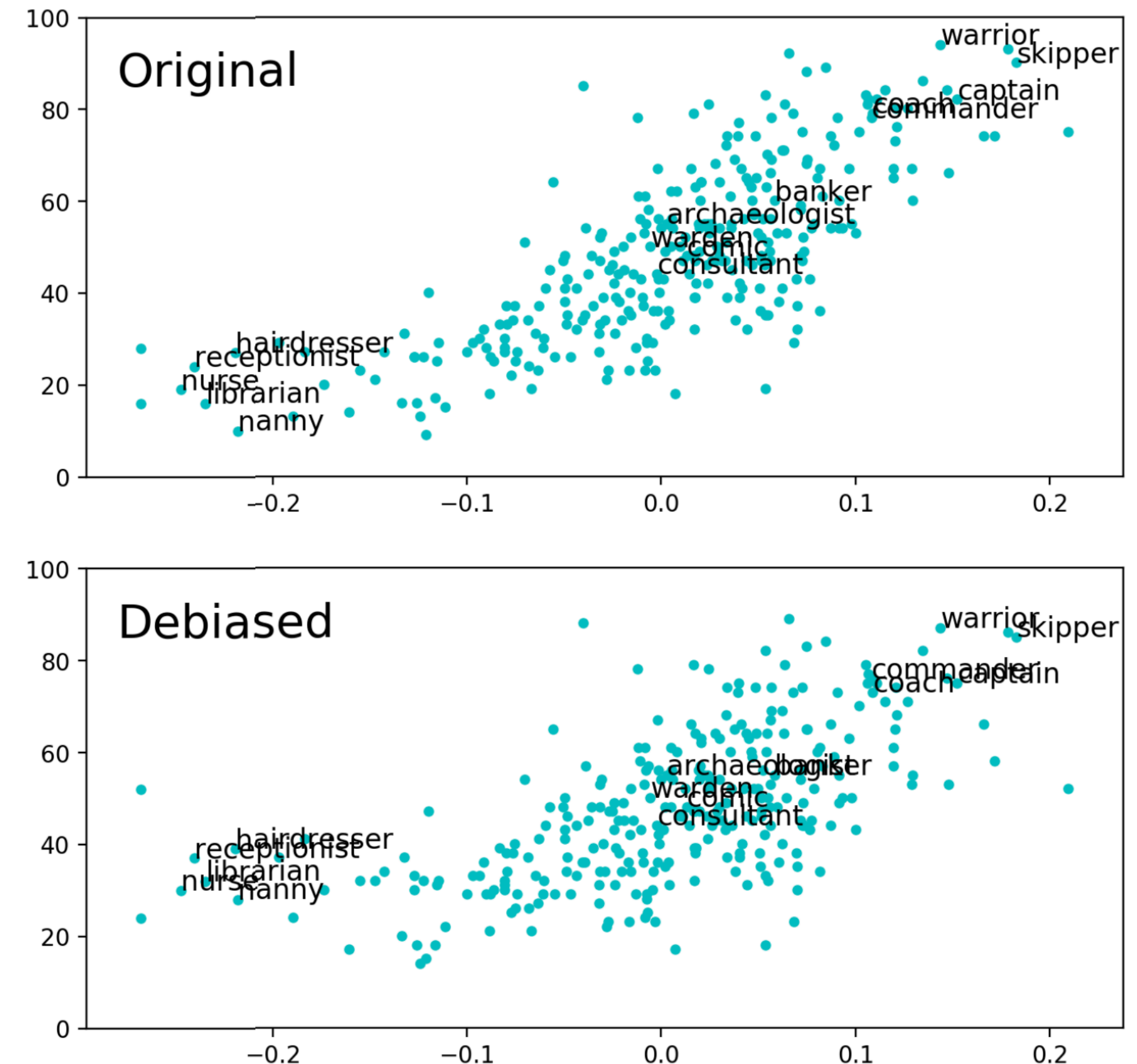


Bolukbasi et al. (2016)



Hardness of Debiasing

- ▶ Not that effective...and the male and female words are still clustered together
- ▶ Bias pervades the word embedding space and isn't just a local property of a few words



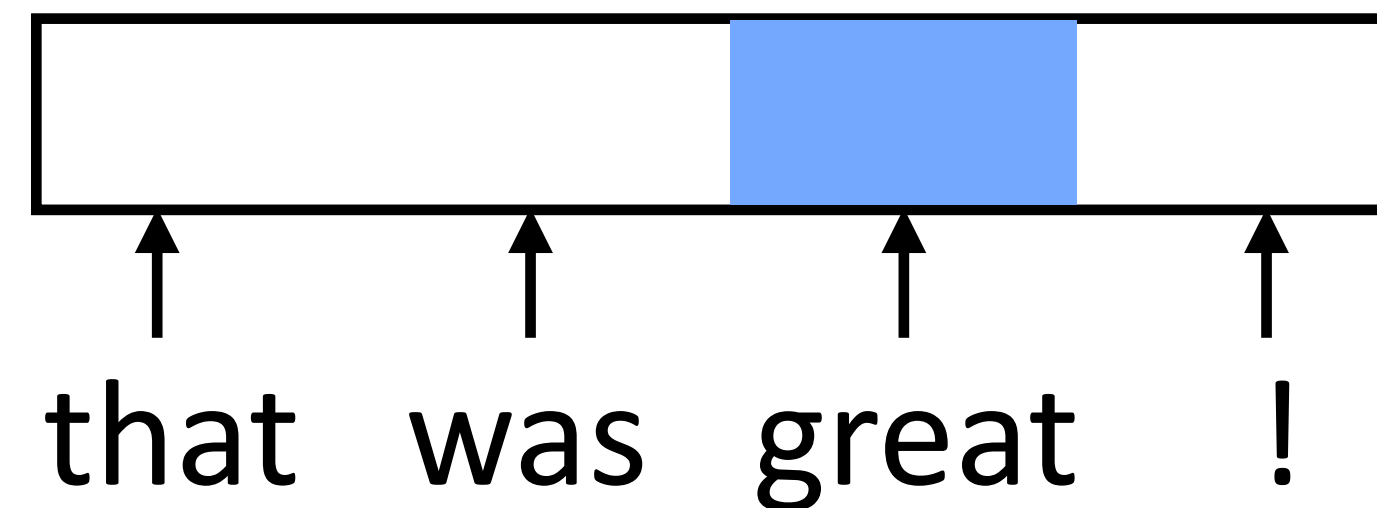
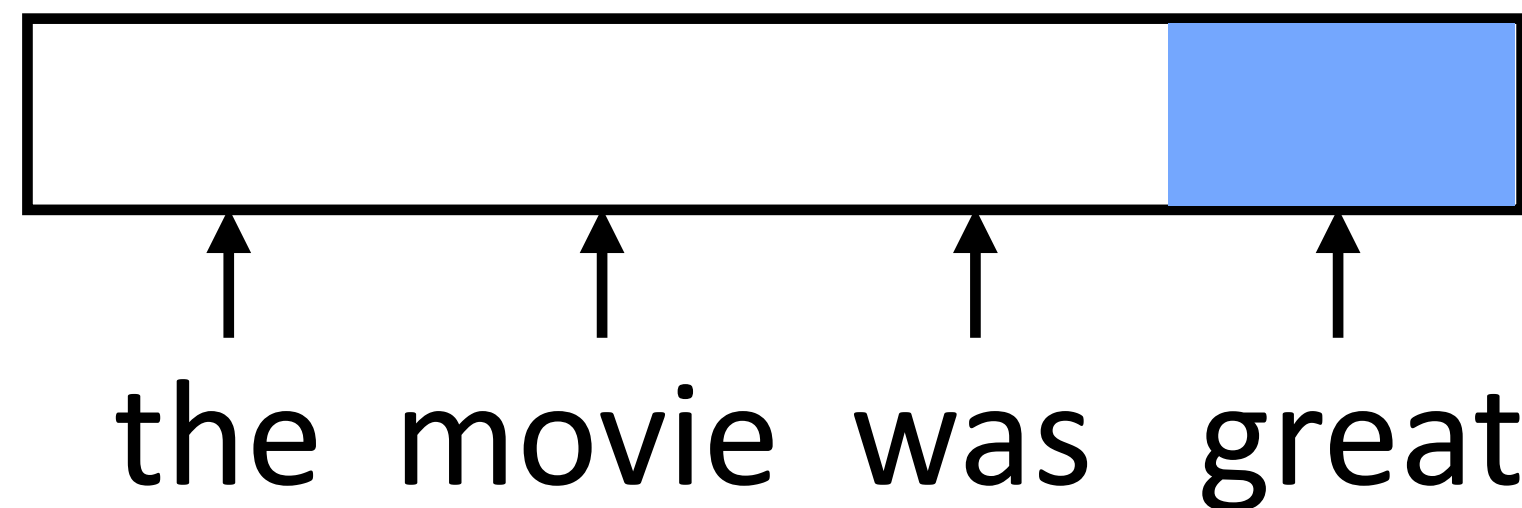
(a) The plots for HARD-DEBIASED embedding, before (top) and after (bottom) debiasing.

RNN Basics



RNN Motivation

- ▶ Feedforward NNs can't handle variable length input: each position in the feature vector has fixed semantics

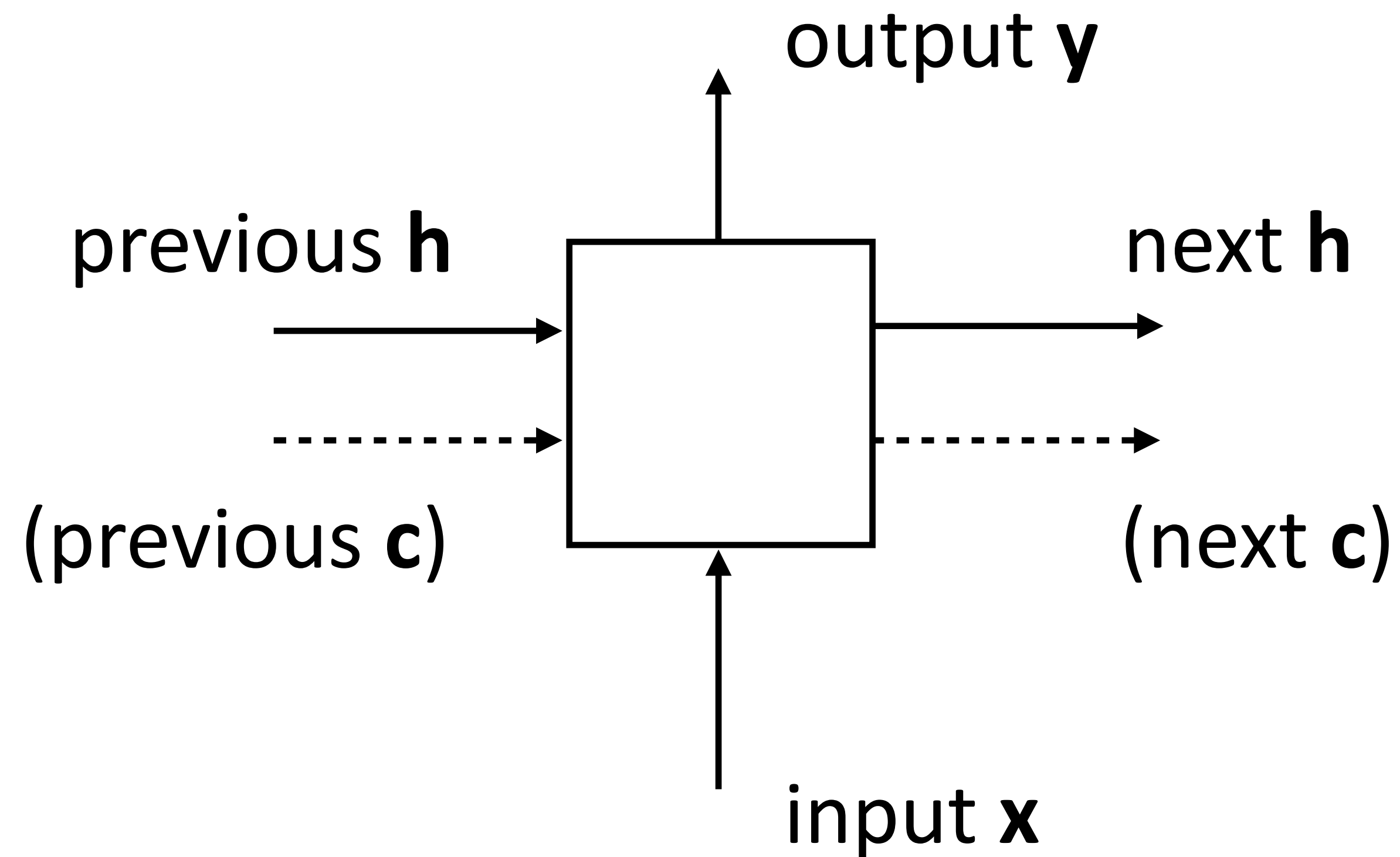


- ▶ These don't look related (*great* is in two different orthogonal subspaces)
- ▶ Instead, we need to:
 - 1) Process each word in a uniform way
 - 2) ...while still exploiting the context that that token occurs in



RNN Abstraction

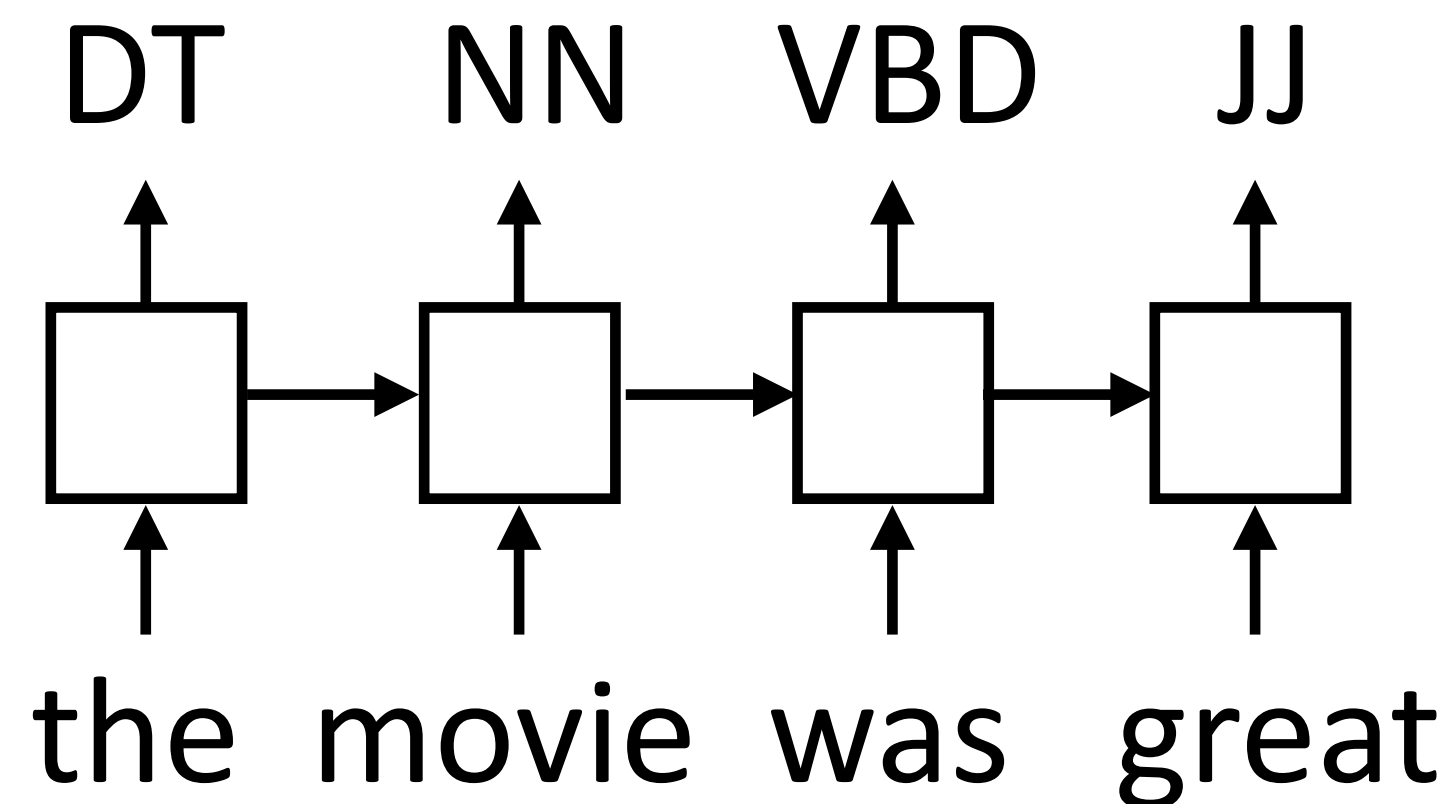
- ▶ Cell that takes some input \mathbf{x} , has some hidden state \mathbf{h} , and updates that hidden state and produces output \mathbf{y} (all vector-valued)





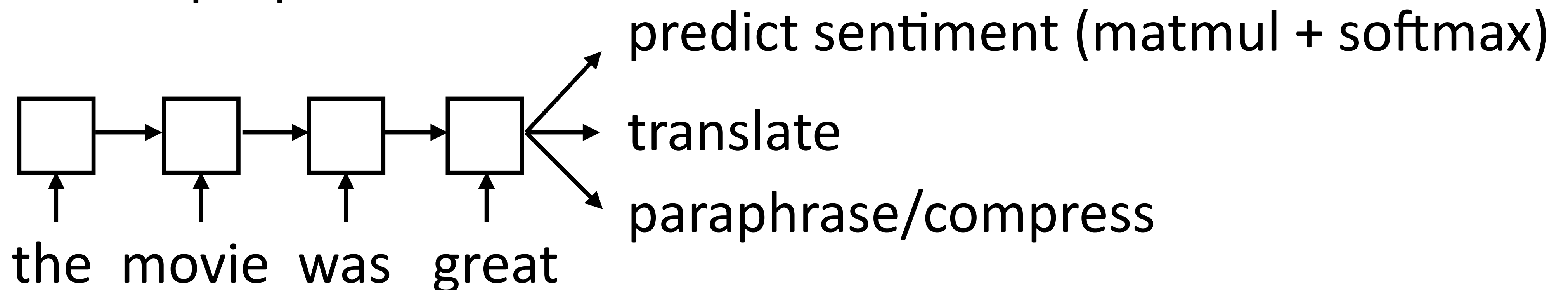
RNN Uses

- ▶ Transducer: make some prediction for each element in a sequence



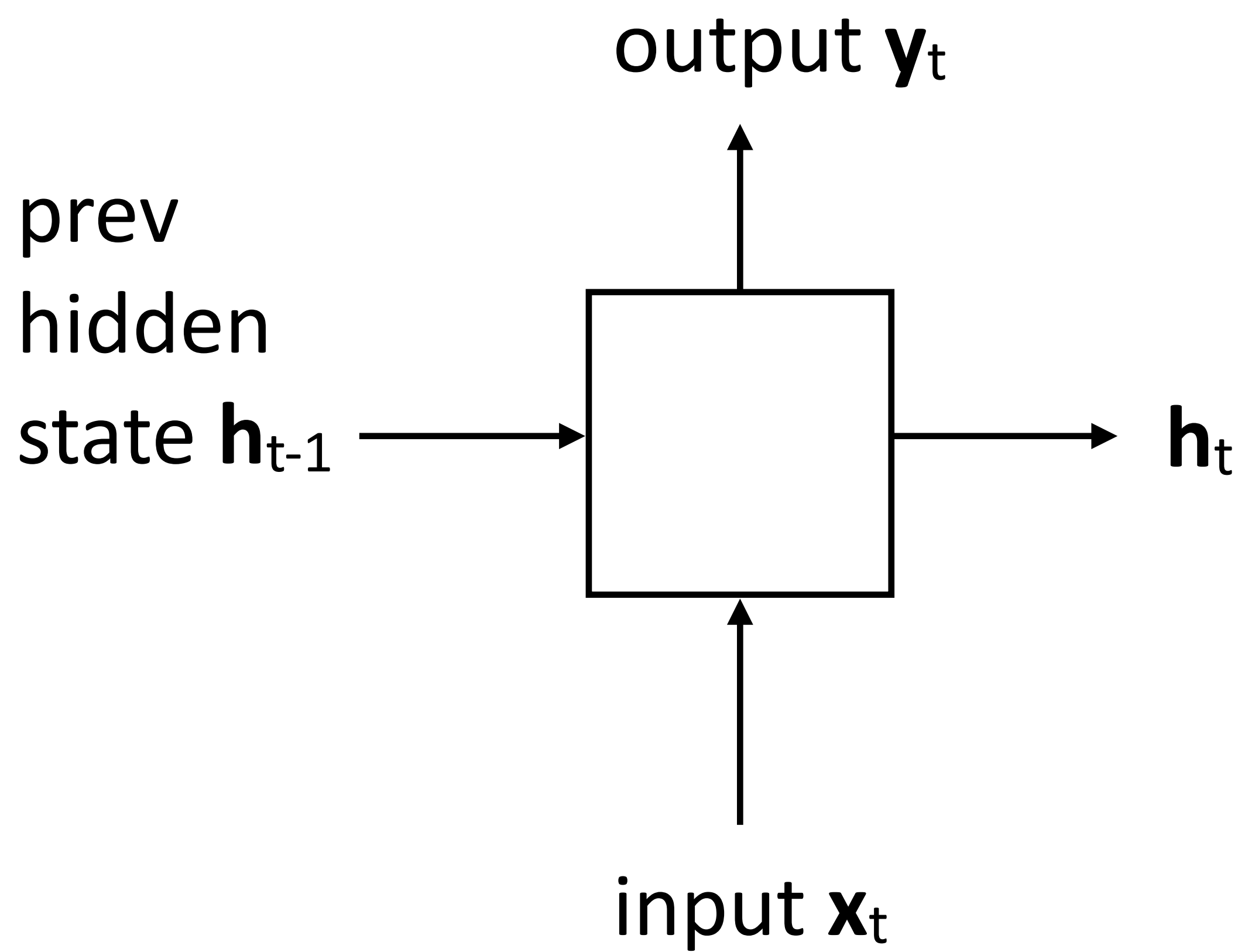
output \mathbf{y} = score for each tag, then softmax

- ▶ Acceptor/encoder: encode a sequence into a fixed-sized vector and use that for some purpose





Elman Networks



$$\mathbf{h}_t = \tanh(W\mathbf{x}_t + V\mathbf{h}_{t-1} + \mathbf{b}_h)$$

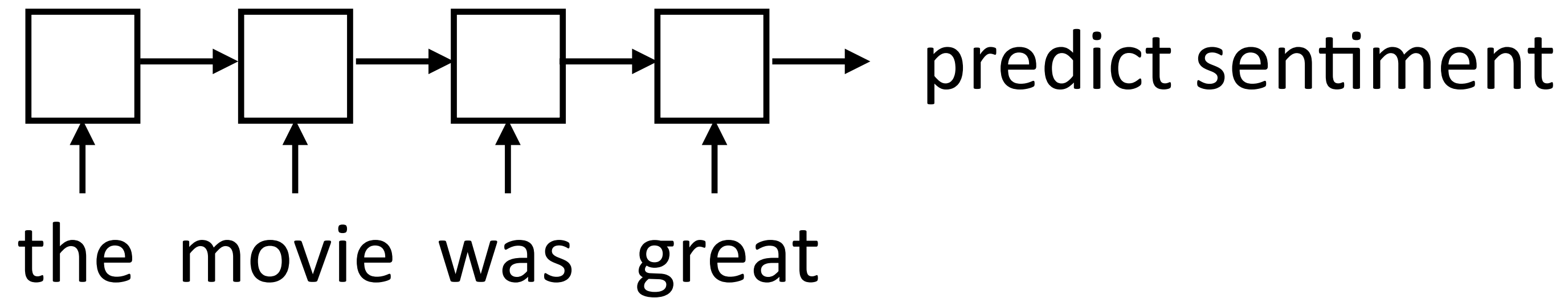
- Updates hidden state based on input and current hidden state

$$\mathbf{y}_t = \tanh(U\mathbf{h}_t + \mathbf{b}_y)$$

- Computes output from hidden state



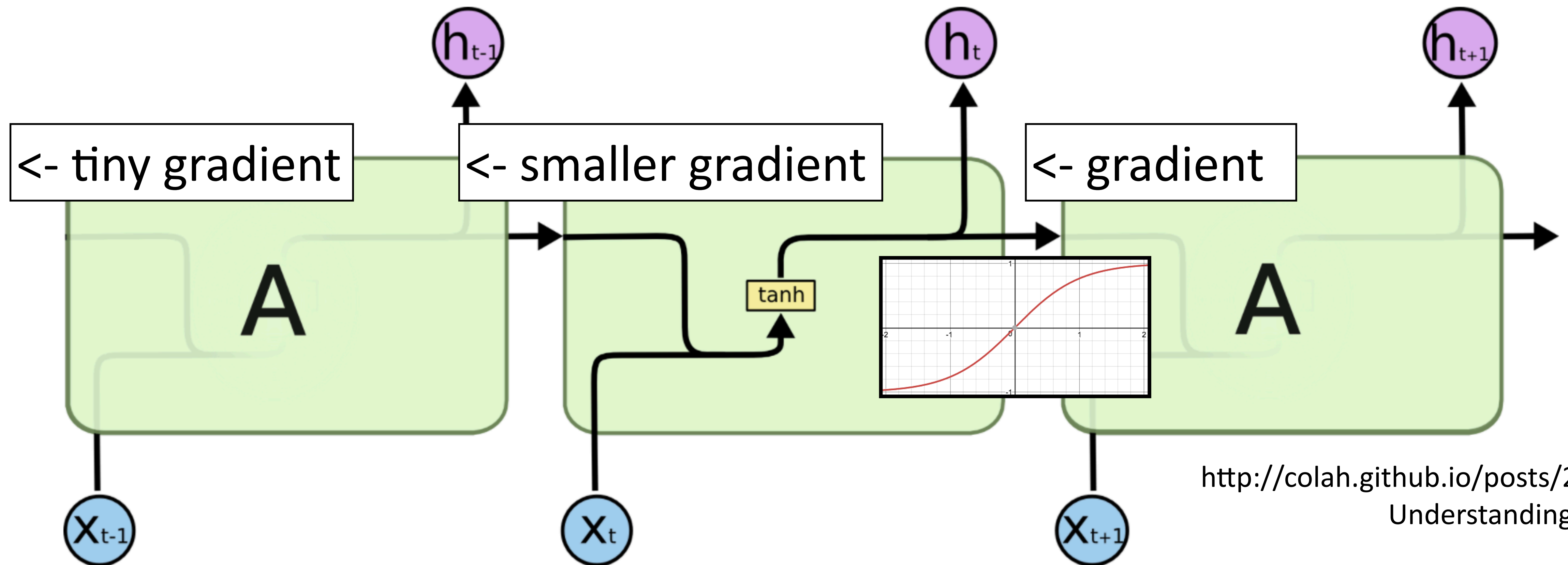
Training Elman Networks



- ▶ “Backpropagation through time”: build the network as one big computation graph, some parameters are shared
 - ▶ RNN potentially needs to learn how to “remember” information for a long time!
- it was my **favorite** movie of 2016, though it wasn't without **problems** -> +
- ▶ “Correct” parameter update is to do a better job of remembering the sentiment of *favorite*



Vanishing Gradient



- ▶ Gradient diminishes going through tanh; if not in $[-2, 2]$, gradient is almost 0
- ▶ Repeated multiplication by V causes problems $\mathbf{h}_t = \tanh(W\mathbf{x}_t + V\mathbf{h}_{t-1} + \mathbf{b}_h)$

LSTMs/GRUs



Gated Connections

- ▶ Designed to fix “vanishing gradient” problem using *gates*

$$\mathbf{h}_t = \mathbf{h}_{t-1} \odot \mathbf{f} + \text{func}(\mathbf{x}_t)$$

gated

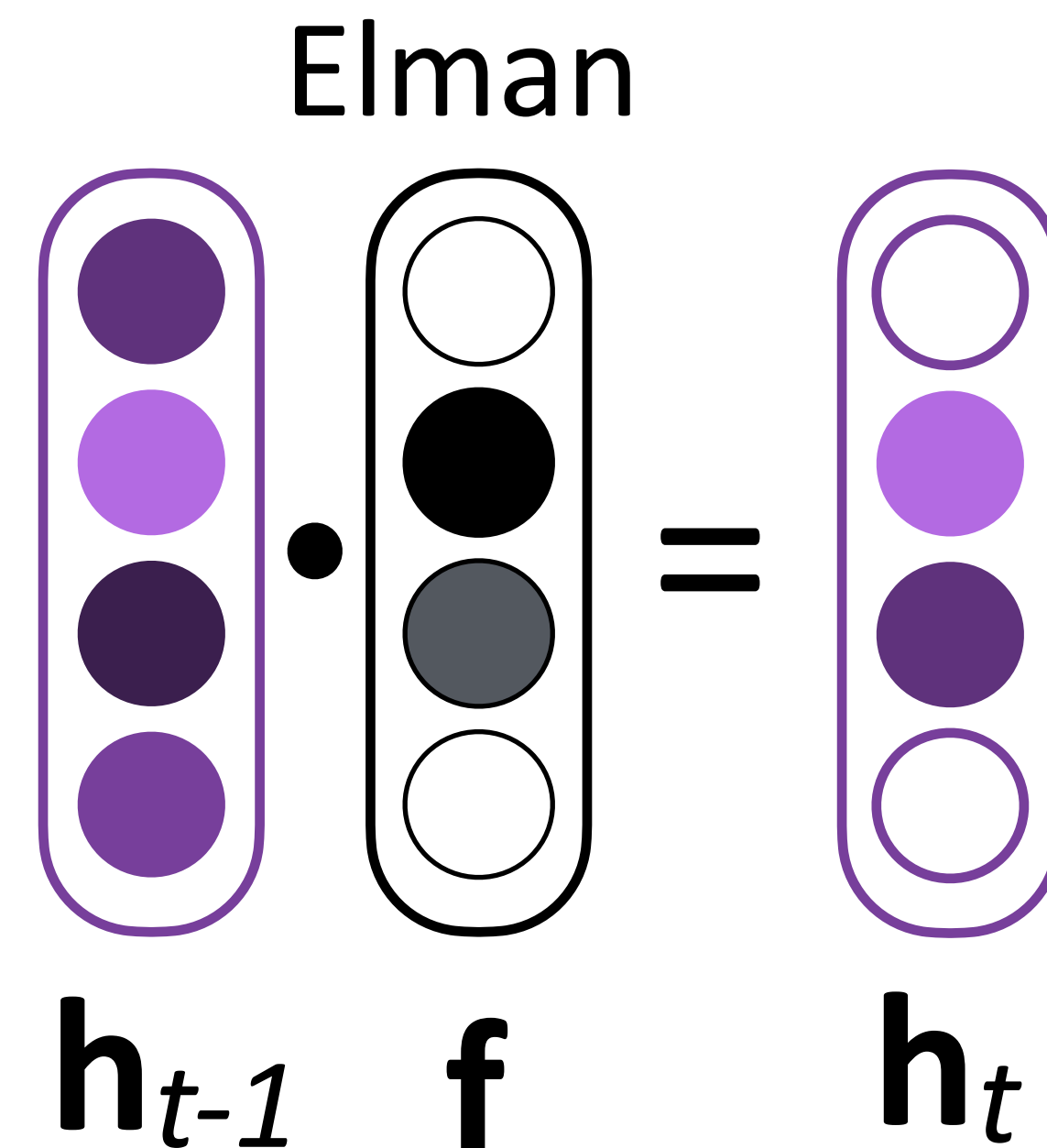
$$\mathbf{h}_t = \tanh(W\mathbf{x}_t + V\mathbf{h}_{t-1} + \mathbf{b}_h)$$

- ▶ Vector-valued “forget gate” \mathbf{f} computed based on input and previous hidden state

$$\mathbf{f} = \sigma(W^{xf}\mathbf{x}_t + W^{hf}\mathbf{h}_{t-1})$$

- ▶ Sigmoid: elements of \mathbf{f} are in $(0, 1)$

- ▶ If $\mathbf{f} \approx \mathbf{1}$, we simply sum up a function of all inputs — gradient doesn’t vanish! More stable without matrix multiply (V) as well





LSTMs

- ▶ “Long short-term memory” network: hidden state is a “short-term” memory

- ▶ “Cell” \mathbf{c} in addition to hidden state \mathbf{h}

$$\mathbf{c}_t = \mathbf{c}_{t-1} \odot \mathbf{f} + \text{func}(\mathbf{x}_t, \mathbf{h}_{t-1})$$

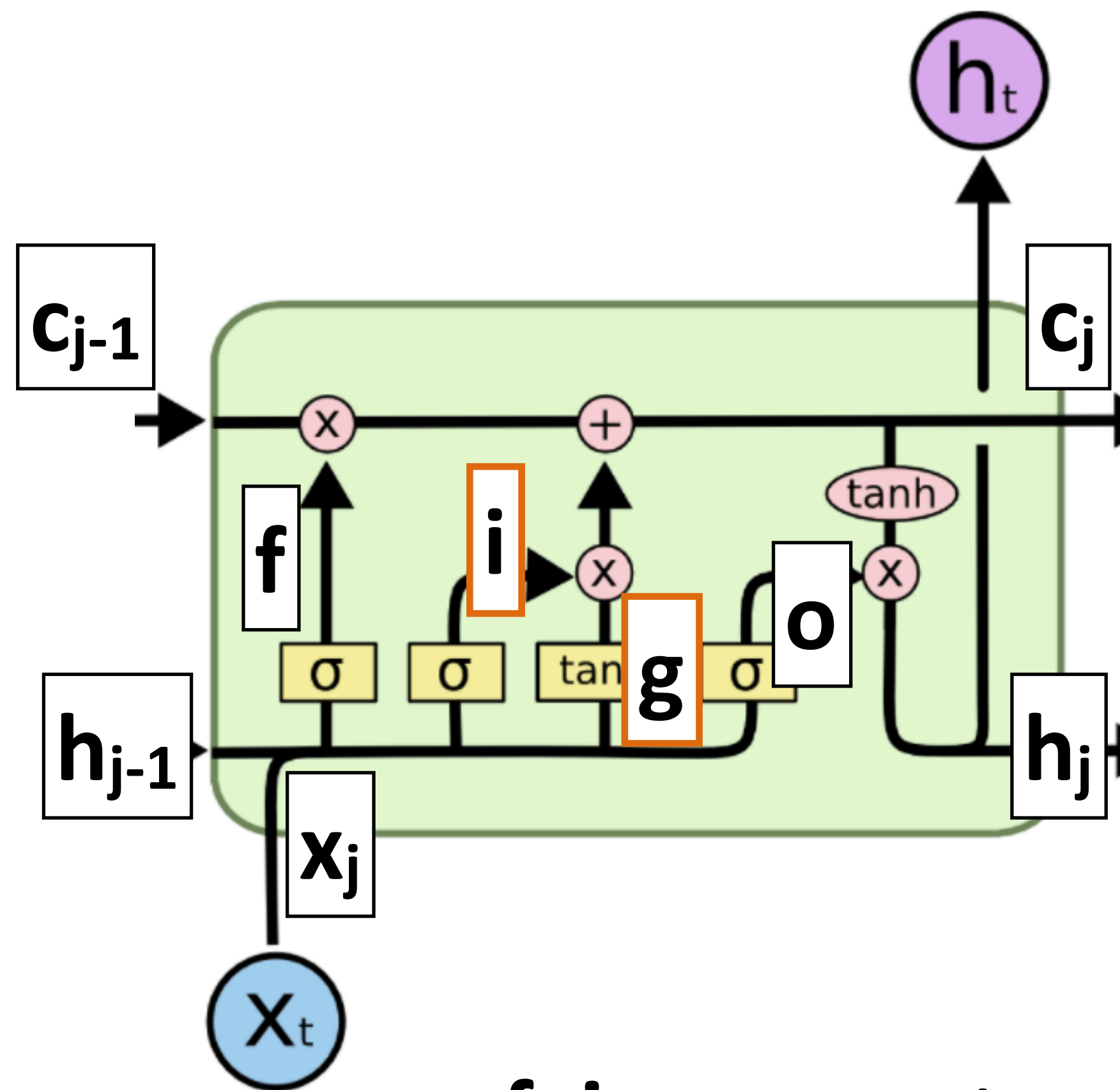
- ▶ Vector-valued forget gate \mathbf{f} depends on the \mathbf{h} hidden state

$$\mathbf{f} = \sigma(W^{xf}\mathbf{x}_t + W^{hf}\mathbf{h}_{t-1})$$

- ▶ Basic communication flow: $\mathbf{x} \rightarrow \mathbf{c} \rightarrow \mathbf{h} \rightarrow \text{output}$, each step of this process is gated in addition to gates from previous timesteps



LSTMs



$$\mathbf{c}_j = \mathbf{c}_{j-1} \odot \mathbf{f} + \mathbf{g} \odot \mathbf{i}$$

$$\mathbf{f} = \sigma(\mathbf{x}_j \mathbf{W}^{\mathbf{x}\mathbf{f}} + \mathbf{h}_{j-1} \mathbf{W}^{\mathbf{h}\mathbf{f}})$$

$$\mathbf{g} = \tanh(\mathbf{x}_j \mathbf{W}^{\mathbf{x}\mathbf{g}} + \mathbf{h}_{j-1} \mathbf{W}^{\mathbf{h}\mathbf{g}})$$

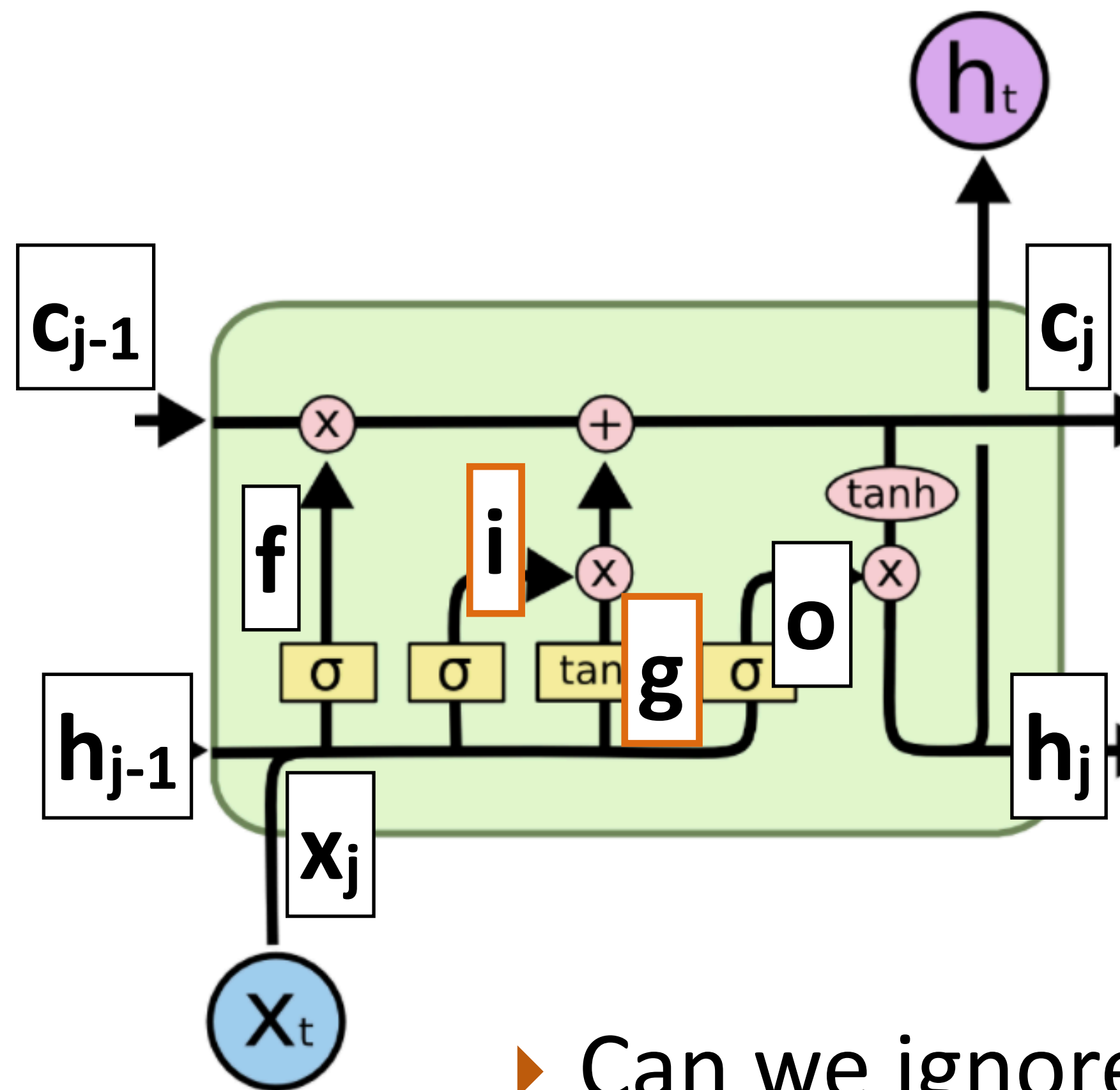
$$\mathbf{i} = \sigma(\mathbf{x}_j \mathbf{W}^{\mathbf{x}\mathbf{i}} + \mathbf{h}_{j-1} \mathbf{W}^{\mathbf{h}\mathbf{i}})$$

$$\mathbf{h}_j = \tanh(\mathbf{c}_j) \odot \mathbf{o}$$

$$\mathbf{o} = \sigma(\mathbf{x}_j \mathbf{W}^{\mathbf{x}\mathbf{o}} + \mathbf{h}_{j-1} \mathbf{W}^{\mathbf{h}\mathbf{o}})$$

- ▶ \mathbf{f} , \mathbf{i} , \mathbf{o} are gates that control information flow
- ▶ \mathbf{g} reflects the main computation of the cell

LSTMs



$$\mathbf{c}_j = \mathbf{c}_{j-1} \odot \mathbf{f} + \mathbf{g} \odot \mathbf{i}$$

$$\mathbf{f} = \sigma(\mathbf{x}_j \mathbf{W}^{\mathbf{x}\mathbf{f}} + \mathbf{h}_{j-1} \mathbf{W}^{\mathbf{h}\mathbf{f}})$$

$$\mathbf{g} = \tanh(\mathbf{x}_j \mathbf{W}^{\mathbf{x}\mathbf{g}} + \mathbf{h}_{j-1} \mathbf{W}^{\mathbf{h}\mathbf{g}})$$

$$\mathbf{i} = \sigma(\mathbf{x}_j \mathbf{W}^{\mathbf{x}\mathbf{i}} + \mathbf{h}_{j-1} \mathbf{W}^{\mathbf{h}\mathbf{i}})$$

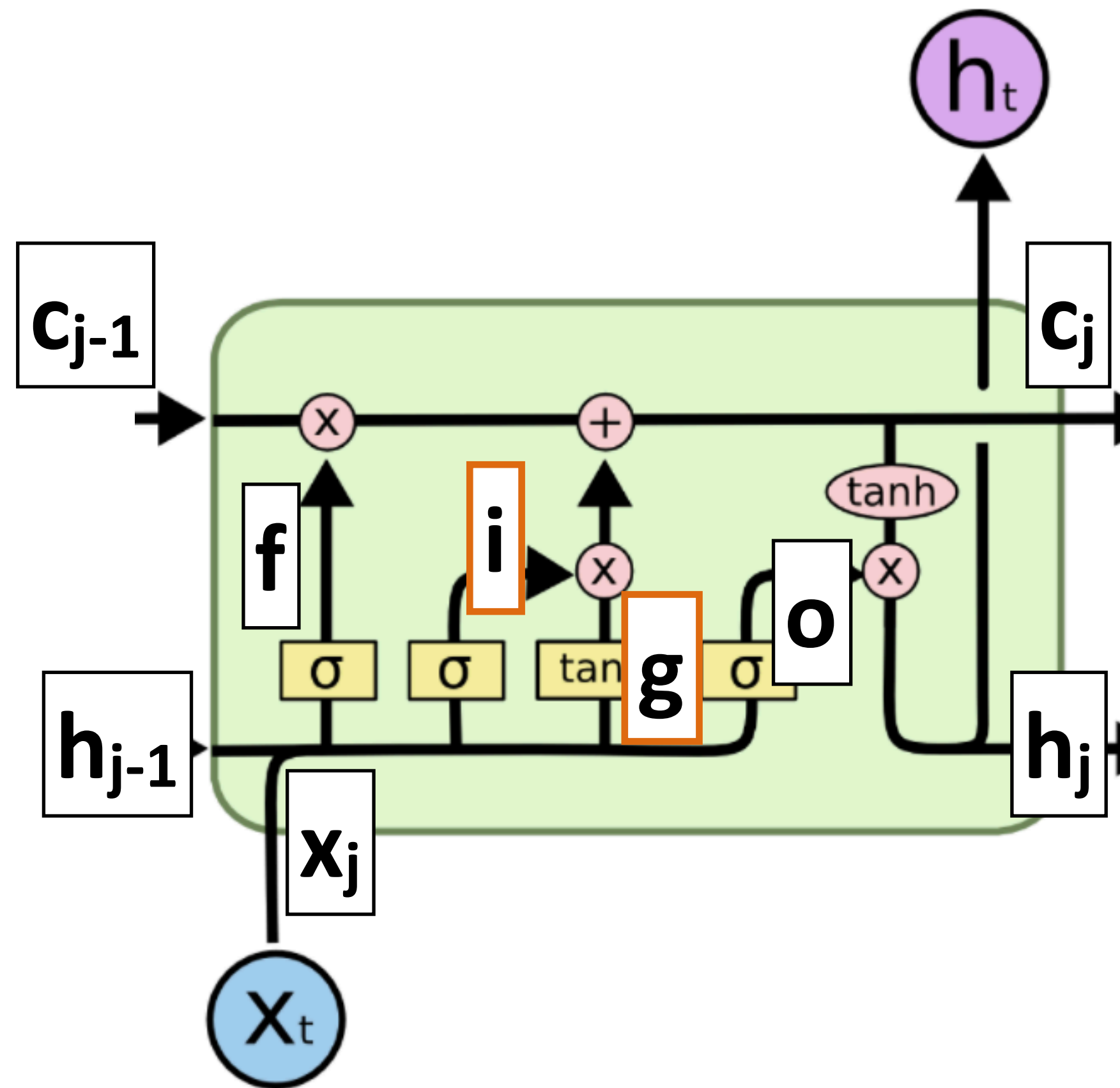
$$\mathbf{h}_j = \tanh(\mathbf{c}_j) \odot \mathbf{o}$$

$$\mathbf{o} = \sigma(\mathbf{x}_j \mathbf{W}^{\mathbf{x}\mathbf{o}} + \mathbf{h}_{j-1} \mathbf{W}^{\mathbf{h}\mathbf{o}})$$

- ▶ Can we ignore the old value of \mathbf{c} for this timestep?
- ▶ Can an LSTM sum up its inputs \mathbf{x} ?
- ▶ Can we ignore a particular input \mathbf{x} ?
- ▶ Can we output something without changing \mathbf{c} ?



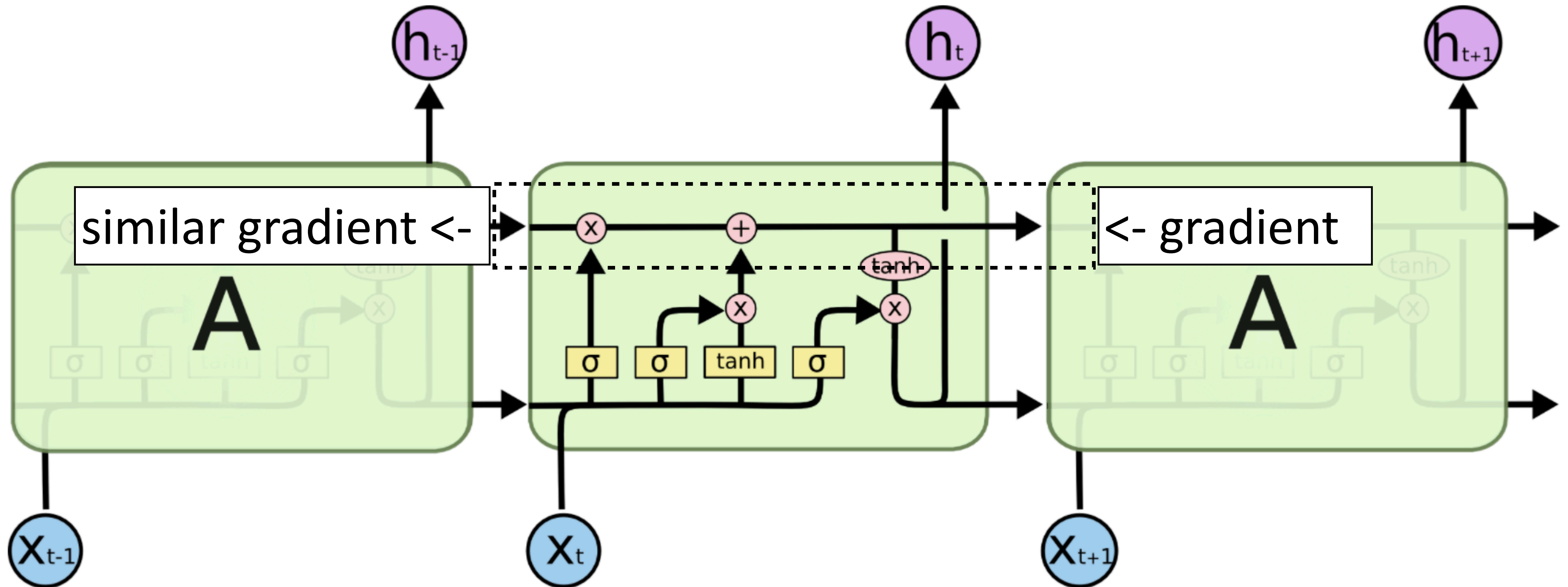
LSTMs



- ▶ Ignoring recurrent state entirely:
 - ▶ Lets us get feedforward layer over token
- ▶ Ignoring input:
 - ▶ Lets us discard stopwords
- ▶ Summing inputs:
 - ▶ Lets us compute a bag-of-words representation



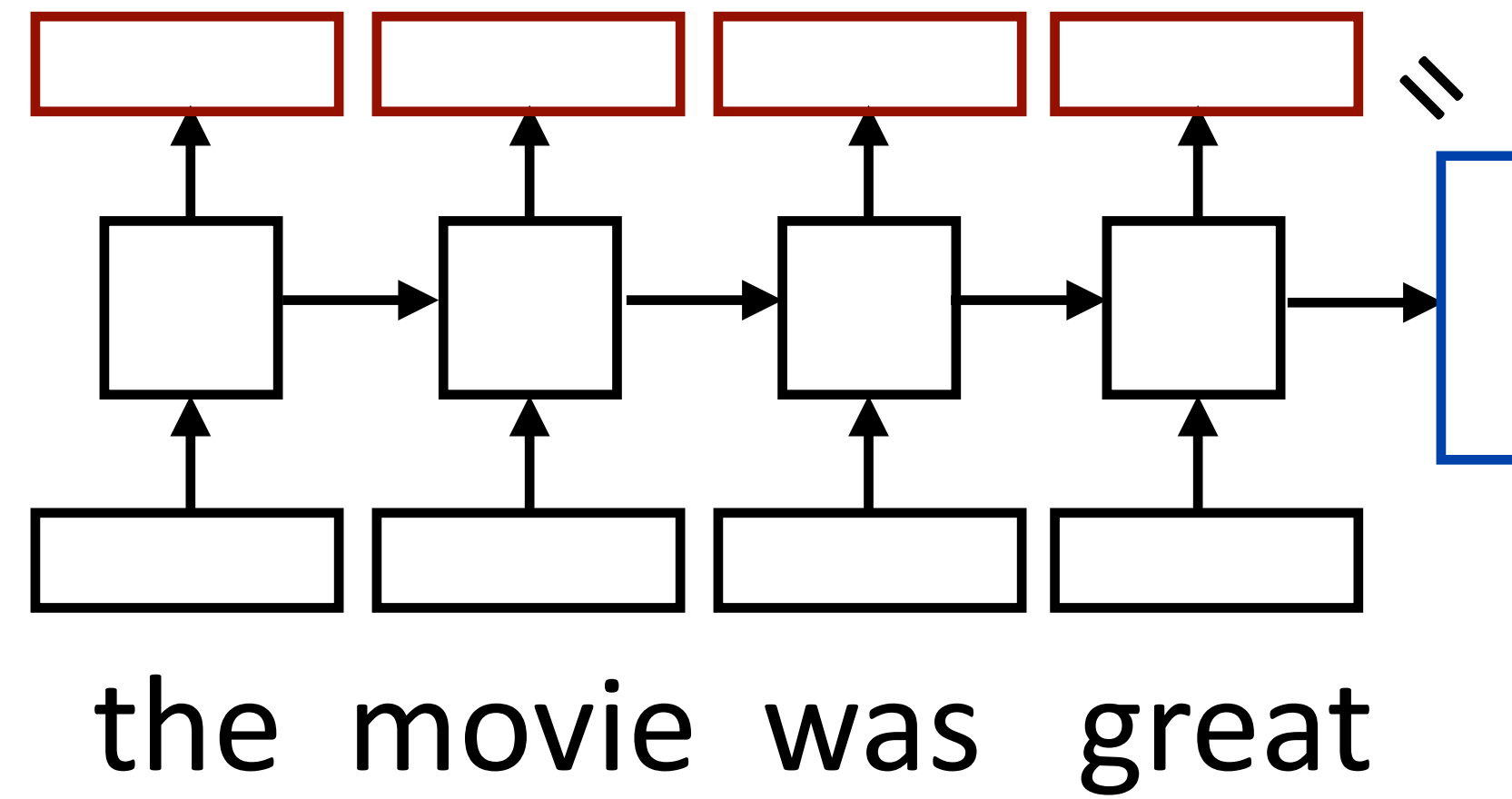
LSTMs



- ▶ Gradient still diminishes, but in a controlled way and generally by less — sometimes initialize forget gate = 1 to remember everything to start



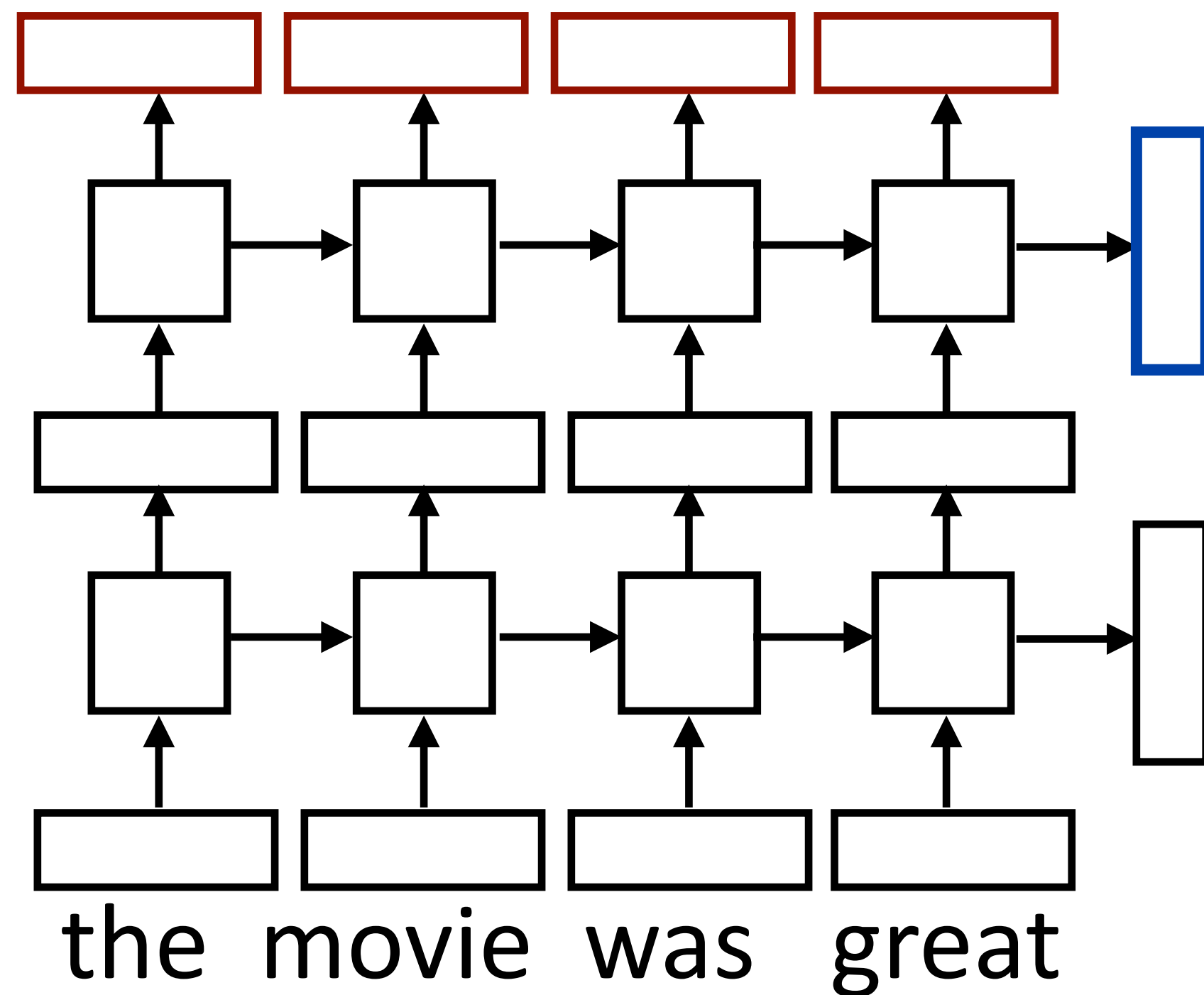
What do RNNs produce?



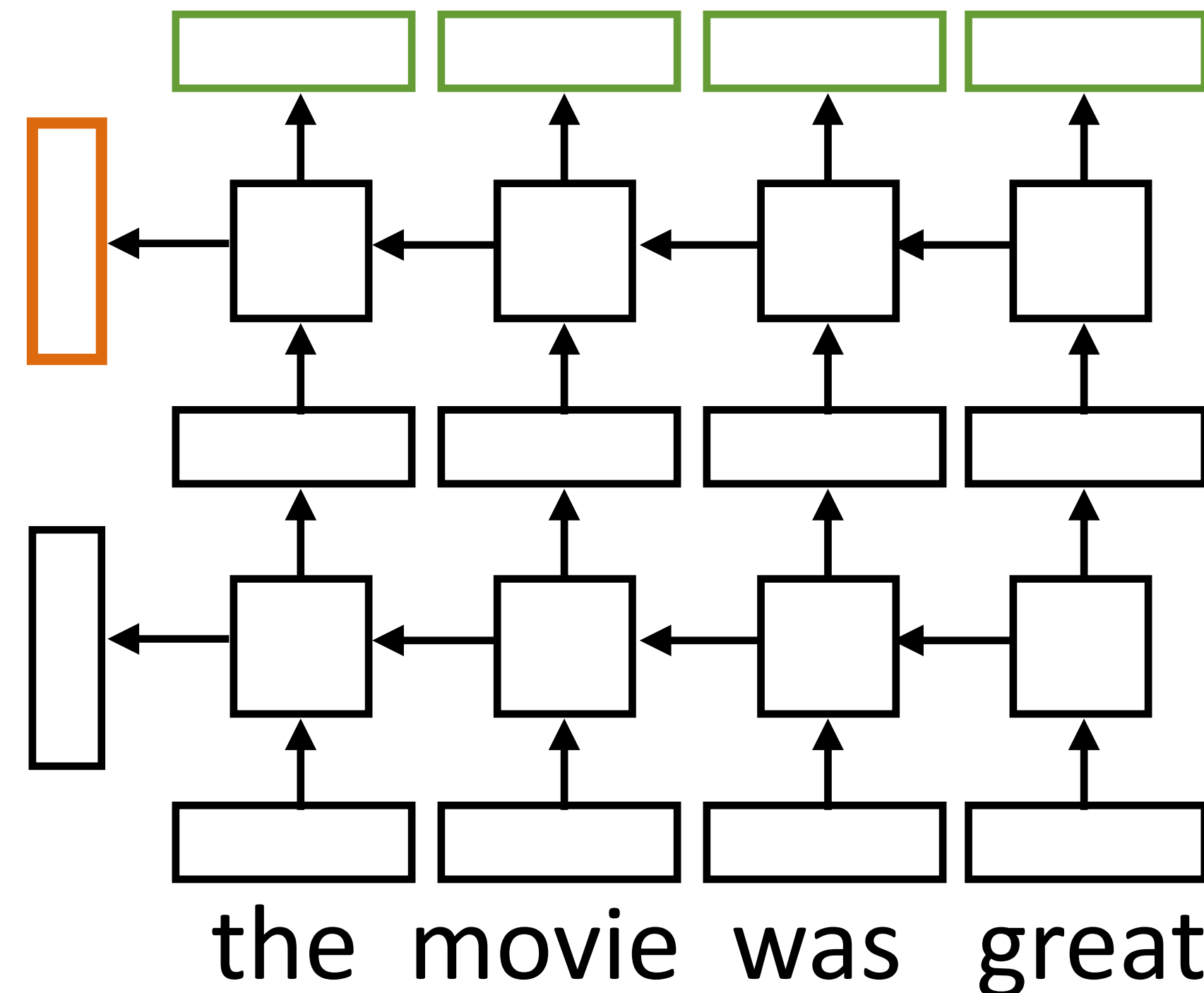
- ▶ **Encoding of the sentence** — can pass this a decoder or make a classification decision about the sentence
- ▶ **Encoding of each word** — can pass this to another layer to make a prediction (can also pool these to get a different sentence encoding)
- ▶ RNN can be viewed as a transformation of a sequence of vectors into a sequence of context-dependent vectors



Multilayer Bidirectional RNN



- Sentence classification based on concatenation of both final outputs

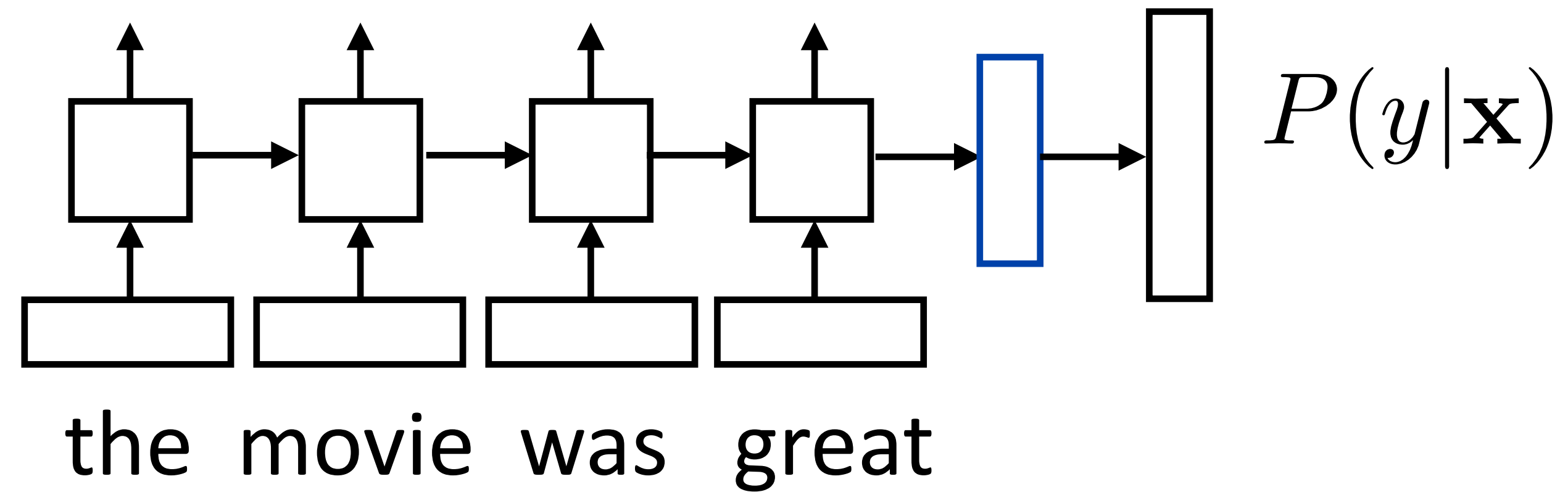


- Token classification based on concatenation of both directions' token representations





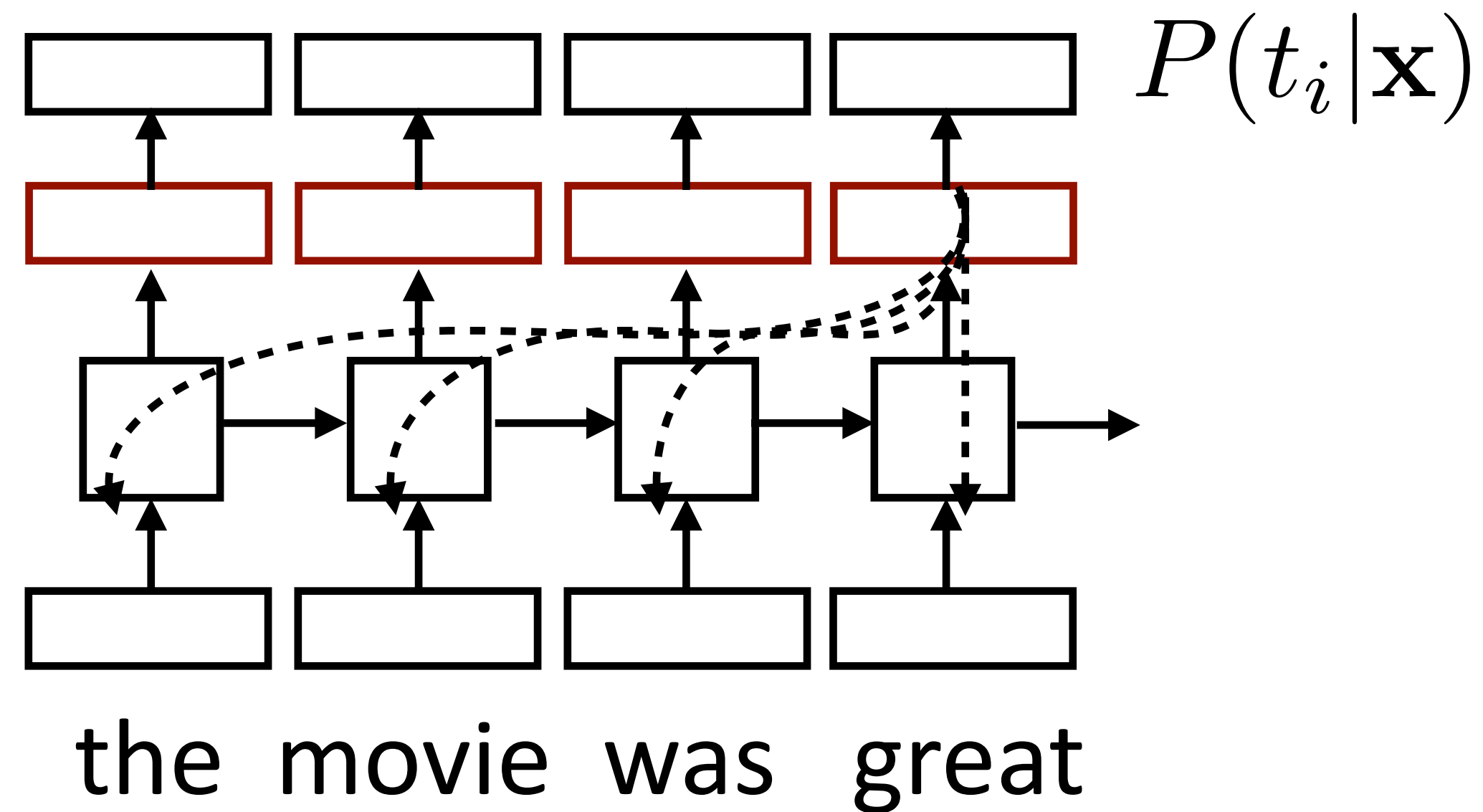
Training RNNs



- ▶ Loss = negative log likelihood of probability of gold label (or use SVM or other loss)
- ▶ Backpropagate through entire network
- ▶ Example: sentiment analysis



Training RNNs



- ▶ Loss = negative log likelihood of probability of gold predictions, summed over the tags
- ▶ Loss terms filter back through network
- ▶ Example: language modeling (predict next word given context) or POS tagging

Applications



What can LSTMs model?

- ▶ Sentiment
 - ▶ Encode one sentence, predict
- ▶ Language models
 - ▶ Move left-to-right, per-token prediction
- ▶ Translation
 - ▶ Encode sentence + then decode, use token predictions for attention weights (later in the course)



Visualizing LSTMs

- ▶ Train *character* LSTM language model (predict next character based on history) over two datasets: War and Peace and Linux kernel source code
- ▶ Visualize activations of specific cells (components of **c**) to understand them
- ▶ Counter: know when to generate \n

The sole importance of the crossing of the Berezina lies in the fact that it plainly and indubitably proved the fallacy of all the plans for cutting off the enemy's retreat and the soundness of the only possible line of action--the one Kutuzov and the general mass of the army demanded--namely, simply to follow the enemy up. The French crowd fled at a continually increasing speed and all its energy was directed to reaching its goal. It fled like a wounded animal and it was impossible to block its path. This was shown not so much by the arrangements it made for crossing as by what took place at the bridges. When the bridges broke down, unarmed soldiers, people from Moscow and women with children who were with the French transport, all--carried on by vis inertiae--pressed forward into boats and into the ice-covered water and did not, surrender.



Visualizing LSTMs

- ▶ Train *character* LSTM language model (predict next character based on history) over two datasets: War and Peace and Linux kernel source code
- ▶ Visualize activations of specific cells to see what they track
- ▶ Binary switch: tells us if we're in a quote or not

"You mean to imply that I have nothing to eat out of.... On the contrary, I can supply you with everything even if you want to give dinner parties," warmly replied Chichagov, who tried by every word he spoke to prove his own rectitude and therefore imagined Kutuzov to be animated by the same desire.

Kutuzov, shrugging his shoulders, replied with his subtle penetrating smile: "I meant merely to say what I said."



Visualizing LSTMs

- ▶ Train *character* LSTM language model (predict next character based on history) over two datasets: War and Peace and Linux kernel source code
- ▶ Visualize activations of specific cells to see what they track
- ▶ Stack: activation based on indentation

```
#ifdef CONFIG_AUDITSYSCALL
static inline int audit_match_class_bits(int class, u32 *mask)
{
    int i;
    if (classes[class]) {
        for (i = 0; i < AUDIT_BITMASK_SIZE; i++)
            if (mask[i] & classes[class][i])
                return 0;
    }
    return 1;
}
```




Visualizing LSTMs

- ▶ Train *character* LSTM language model (predict next character based on history) over two datasets: War and Peace and Linux kernel source code
- ▶ Visualize activations of specific cells to see what they track
- ▶ Uninterpretable: probably doing double-duty, or only makes sense in the context of another activation

```
/* Unpack a filter field's string representation from user-space
 * buffer. */
char *audit_unpack_string(void **bufp, size_t *remain, size_t len)
{
    char *str;
    if (!*bufp || (len == 0) || (len > *remain))
        return ERR_PTR(-EINVAL);
    /* Of the currently implemented string fields, PATH_MAX
     * defines the longest valid length.
     */
}
```



What can LSTMs model?

- ▶ Sentiment
 - ▶ Encode one sentence, predict
- ▶ Language models
 - ▶ Move left-to-right, per-token prediction
- ▶ Translation
 - ▶ Encode sentence + then decode, use token predictions for attention weights (later in the course)
- ▶ Textual entailment
 - ▶ Encode two sentences, predict



Sentiment Analysis

- ▶ Semi-supervised method: initialize the language model by training to reproduce the document in a seq2seq fashion (a type of pre-training called a sequential autoencoder)

Model	Test error rate
LSTM with tuning and dropout	13.50%
LSTM initialized with word2vec embeddings	10.00%
LM-LSTM (see Section 2)	7.64%
SA-LSTM (see Figure 1)	7.24%
Full+Unlabeled+BoW [21]	11.11%
WRRBM + BoW (bnc) [21]	10.77%
NBSVM-bi (Naïve Bayes SVM with bigrams) [35]	8.78%
seq2-bow _n -CNN (ConvNet with dynamic pooling) [11]	7.67%
Paragraph Vectors [18]	7.42%



Natural Language Inference

Premise

Hypothesis

A boy plays in the snow

entails

A boy is outside

A man inspects the uniform of a figure

contradicts

The man is sleeping

An older and younger man smiling

neutral

Two men are smiling and
laughing at cats playing

- ▶ Long history of this task: “Recognizing Textual Entailment” challenge in 2006 (Dagan, Glickman, Magnini)
- ▶ Early datasets: small (hundreds of pairs), very ambitious (lots of world knowledge, temporal reasoning, etc.)



SNLI Dataset

- ▶ Show people captions for (unseen) images and solicit entailed / neural / contradictory statements

- ▶ >500,000 sentence pairs

- ▶ Encode each sentence and process

100D LSTM: 78% accuracy

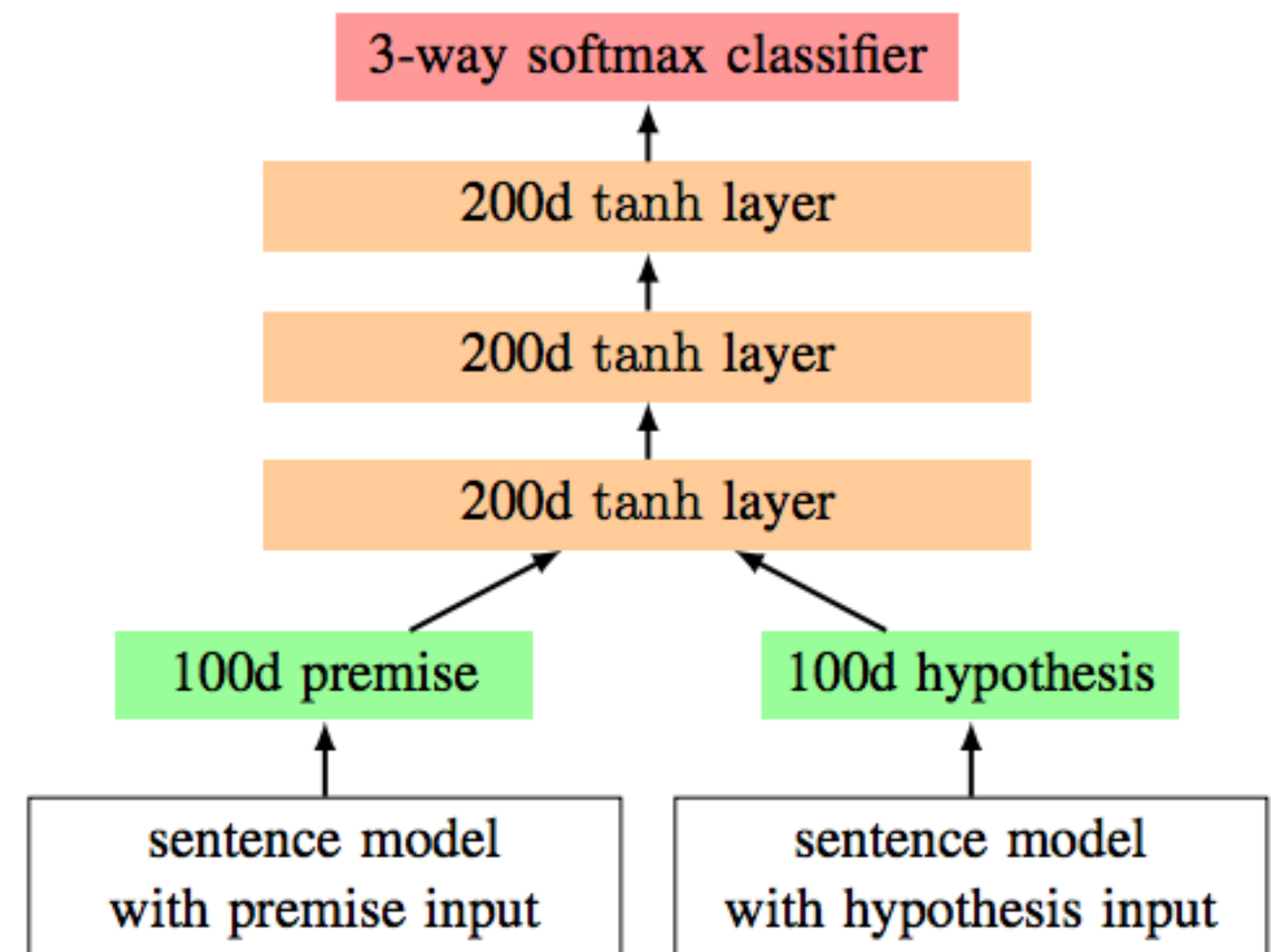
300D LSTM: 80% accuracy

(Bowman et al., 2016)

300D BiLSTM: 83% accuracy

(Liu et al., 2016)

- ▶ Later: better models for this



Bowman et al. (2015)



Takeaways

- ▶ RNNs can transduce inputs (produce one output for each input) or compress the whole input into a vector
- ▶ Useful for a range of tasks with sequential input: sentiment analysis, language modeling, natural language inference, machine translation