

مختن پاکبزهادا * ریاض سیر هم تبر خیرالبشر ده علمیه الصلوہ والسلام * مختل نشان
اڑاولوب * با غیبان طبع چن بروی ناو دان خامه عنبرت شاردن اساله جو بیار زلان
ایله بر حدیقه هطر از بیهای بیش در کد هر صفحه زنگی بر تخته هشکروفه نامدار *
و هر فقره نگذینی بر دسته ناز بوی نقده تاردر * سلسله سطوری صفحه نفشه راردن
دنجوی تر * هر کل سر مخ سخنلری کل دسته به شدن خوشبوی تر در * تعبیرات
جهان پسندی ز ب دسته اشتهر * و تلمیحات دل کشایی ممتاز و مسلم اصحاب

Natural Language Processing

Info 159/259

Lecture 9: Embeddings 2 (Sept 20, 2018)

David Bamman, UC Berkeley

ANDREW PIPER

Corpus Poetics: Thinking the Writer's Career with Data

5:30 pm - 7:00 pm (today!)

Geballe Room, Townsend Center, 220 Stephens Hall

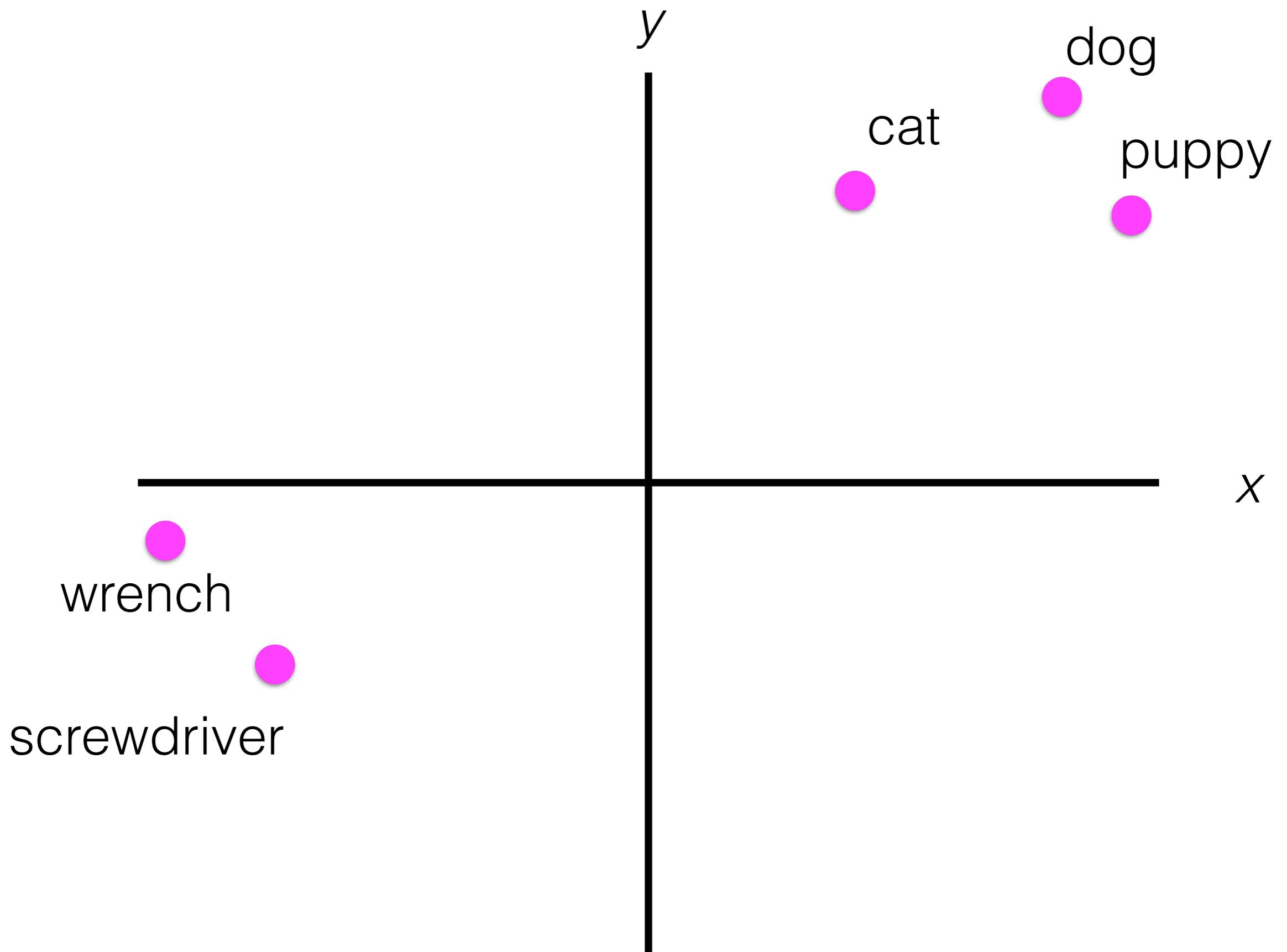
Distributed representation

- Vector representation that encodes information about the **distribution** of contexts a word appears in
- Words that appear in similar contexts have similar representations (and similar meanings, by the **distributional hypothesis**).

Dense vectors from prediction

- Learning low-dimensional representations of words by framing a predicting task: using context to predict words in a surrounding window
- Transform this into a **supervised** prediction problem; similar to language modeling but we're ignoring order within the context window

	1	2	3	4	...	50
the	0.418	0.24968	-0.41242	0.1217	...	-0.17862
,	0.013441	0.23682	-0.16899	0.40951	...	-0.55641
.	0.15164	0.30177	-0.16763	0.17684	...	-0.31086
of	0.70853	0.57088	-0.4716	0.18048	...	-0.52393
to	0.68047	-0.039263	0.30186	-0.17792	...	0.13228
...
chanty	0.23204	0.025672	-0.70699	-0.04547	...	0.34108
kronik	-0.60921	-0.67218	0.23521	-0.11195	...	0.85632
rolonda	-0.51181	0.058706	1.0913	-0.55163	...	0.079711
zsombor	-0.75898	-0.47426	0.4737	0.7725	...	0.84014
sandberger	0.072617	-0.51393	0.4728	-0.52202	...	0.23096



Word embeddings

- Pre-trained word embeddings great for words that appear frequently in data
- Unseen words are treated as UNKs and assigned zero or random vectors; everything unseen is assigned **the same representation**.

- supercalifragilisticexpialidocious
 - super, superior, supernatural
 - adventurous, fabulous, infamous
- supercalifragilisticexpialidociously
 - quickly, sadly, perfectly

Agglutinative languages

Muvaffakiyetsizleş(-mek)	(To) become unsuccessful
Muvaffakiyetsizleştir(-mek)	(To) make one unsuccessful
Muvaffakiyetsizleştirici	Maker of unsuccessful ones
Muvaffakiyetsizleştiricileş(-mek)	(To) become a maker of unsuccessful ones
Muvaffakiyetsizleştiricileştir(-mek)	(To) make one a maker of unsuccessful ones
Muvaffakiyetsizleştiricileştiriver(-mek)	(To) easily/quickly make one a maker of unsuccessful ones
Muvaffakiyetsizleştiricileştiriverebil(-mek)	(To) be able to make one easily/quickly a maker of unsuccessful ones
Muvaffakiyetsizleştiricileştiriveremeyebileceklerimi zdenmişsinizcesine	As though you happen to have been from among those whom we will not be able to easily/quickly make a maker of unsuccessful ones

Infinitive danser	Imperfect je dansais tu dansais il/elle dansait nous dansions vous dansiez ils/elles dansaient	Pluperfect j'avais dansé tu avais dansé il/elle avait dansé nous avions dansé vous aviez dansé ils/elles avaient dansé	Present Subjunctive je danse tu dances il/elle danse nous dansions vous dansiez ils/elles dansent
Past Participle dansé			
Gerund dansant			
Imperative danse (tu) dansons (nous) dansez (vous)	Future je danserai tu danseras il/elle dansera nous danserons vous danserez ils/elles danseront	Future Perfect j'aurai dansé tu auras dansé il/elle aura dansé nous aurons dansé vous aurez dansé ils/elles auront dansé	Imperfect Subjunctive je dansasse tu dansasses il/elle dansât nous dansassions vous dansassiez ils/elles dansassent
Present je danse tu dances il/elle danse nous dansons vous dansez ils/elles dansent	Conditional je danserais tu danserais il/elle danserait nous danserions vous danseriez ils/elles danseraient	Past Anterior j'eus dansé tu eus dansé il/elle eut dansé nous eûmes dansé vous eûtes dansé ils/elles eurent dansé	Present Perfect Subjunctive j'aie dansé tu aies dansé il/elle ait dansé nous ayons dansé vous ayez dansé ils/elles aient dansé
Present Perfect j'ai dansé tu as dansé il/elle a dansé nous avons dansé vous avez dansé ils/elles ont dansé	Past Historic je dansai tu dansas il/elle dansa nous dansâmes vous dansâtes ils/elles dansèrent	Conditional Perfect j'aurais dansé tu aurais dansé il/elle aurait dansé nous aurions dansé vous auriez dansé ils/elles auraient dansé	Pluperfect Subjunctive j'eusse dansé tu eusses dansé il/elle eût dansé nous eussions dansé vous eussiez dansé ils/elles eussent dansé

Shared structure

Even in languages like English that are not agglutinative and aren't highly inflected, words **share important structure.**

Even if we never see the word “unfriendly” in our data, we should be able to reason about it as:
un + friend + ly



Subword models

- Rather than learning a single representation for each word type w , learn representations z for the set of ngrams \mathcal{G}_w that comprise it [Bojanowski et al. 2017]
- The word itself is included among the ngrams (no matter its length).
- A word representation is the sum of those ngrams

$$w = \sum_{g \in \mathcal{G}_w} z_g$$

FastText

$$e(\text{where}) =$$

$e(<\text{wh}>)$
+ $e(\text{whe})$
+ $e(\text{her})$
+ $e(\text{ere})$
+ $e(\text{re}>)$

3-grams

+ $e(<\text{whe}>)$
+ $e(\text{wher})$
+ $e(\text{here})$
+ $e(\text{ere}>)$

4-grams

+ $e(<\text{wher}>)$
+ $e(\text{where})$
+ $e(\text{here}>)$

5-grams

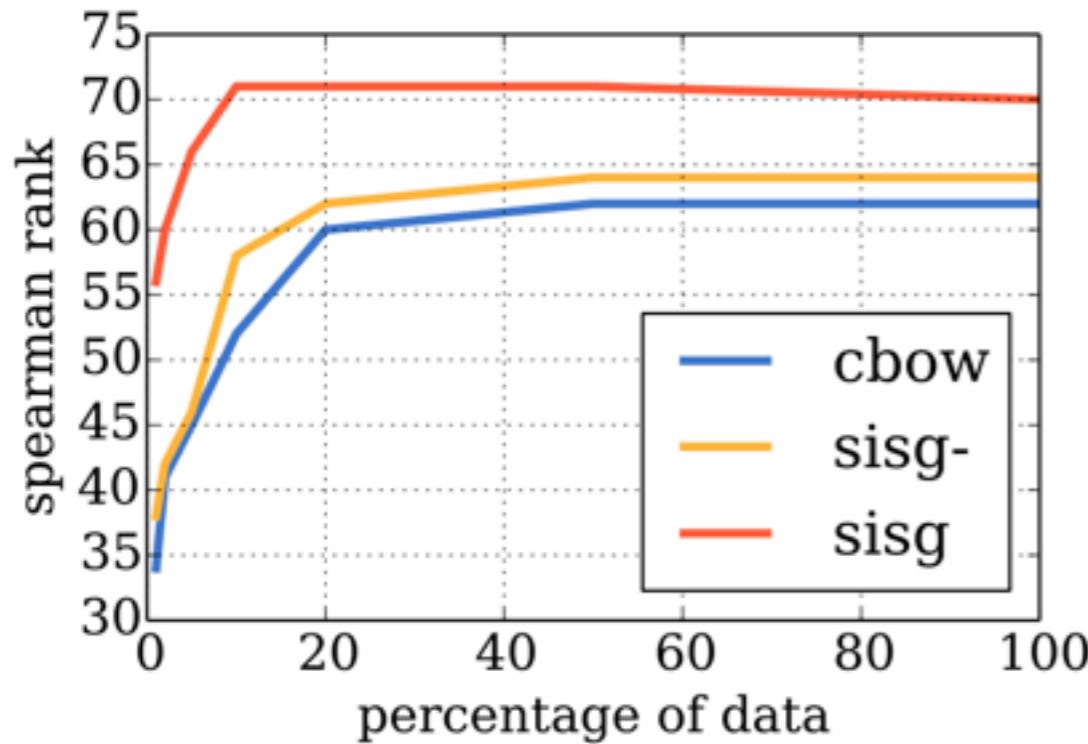
+ $e(<\text{where}>)$
+ $e(\text{where}>)$

6-grams

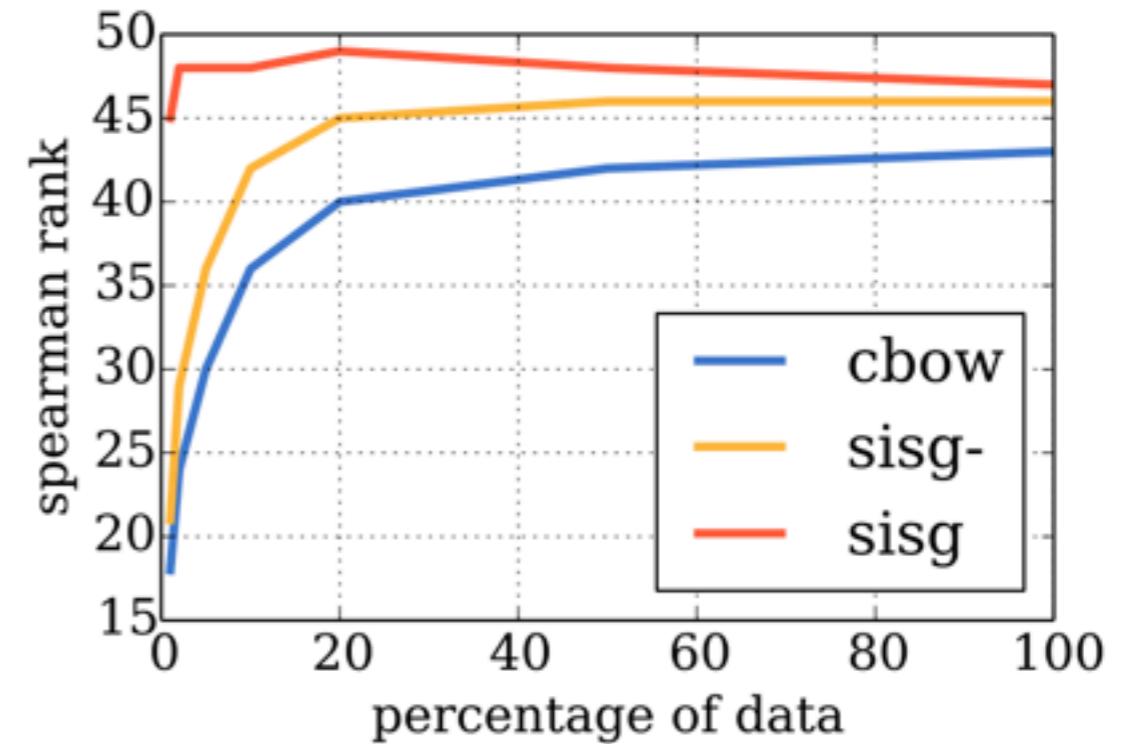
+ $e(<\text{where}>)$

word

$e(*)$ = embedding for *



(a) DE-GUR350



(b) EN-RW

1% =
~20M tokens

100% =
~1B tokens

- Subword models need less data to get comparable performance.

Low-dimensional distributed representations

- Low-dimensional, dense word representations are extraordinarily powerful (and are arguably responsible for much of gains that neural network models have in NLP).
- Lets your representation of the input share statistical strength with words that behave similarly in terms of their distributional properties (often **synonyms** or words that belong to the same **class**).

Using dense vectors

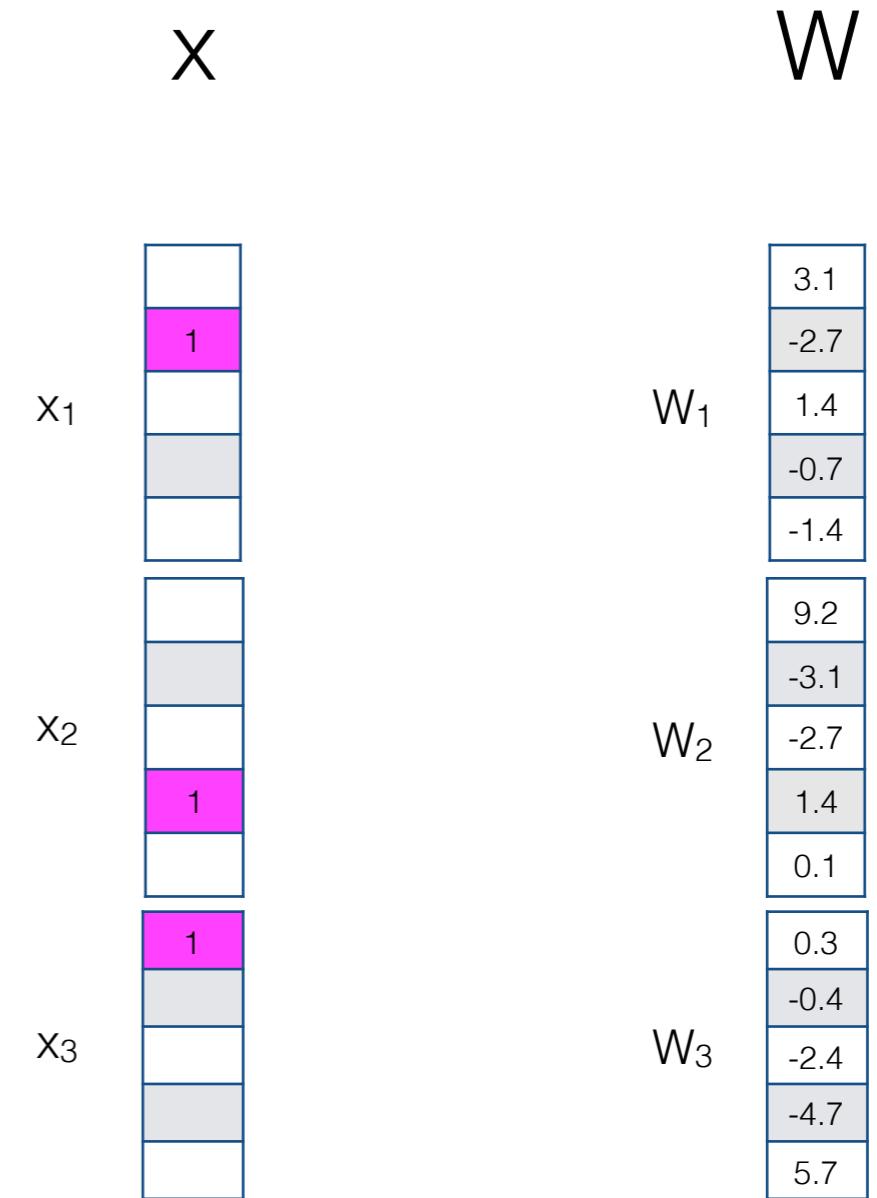
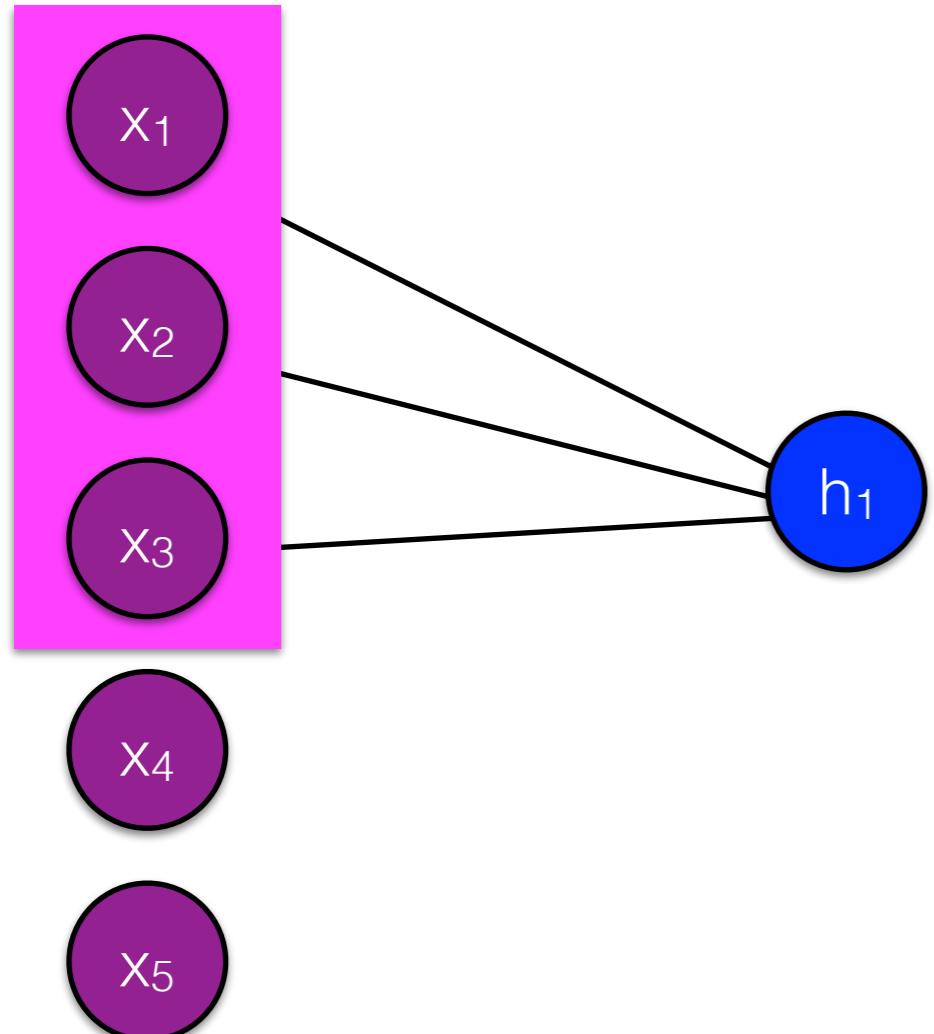
1. Trained word embeddings on a large collection (T tokens) of unlabeled text (Wikipedia, news, Twitter, books), preferably in the domain you want to use them for.
2. Use those pre-trained embeddings in a predictive model with S labeled examples

$$T \gg S$$

Using dense vectors

- In neural models (CNNs, RNNs, LM), replace the V -dimensional sparse vector with the much smaller K -dimensional dense one.
- Can also take the derivative of the loss function with respect to those representations to optimize for a particular task.

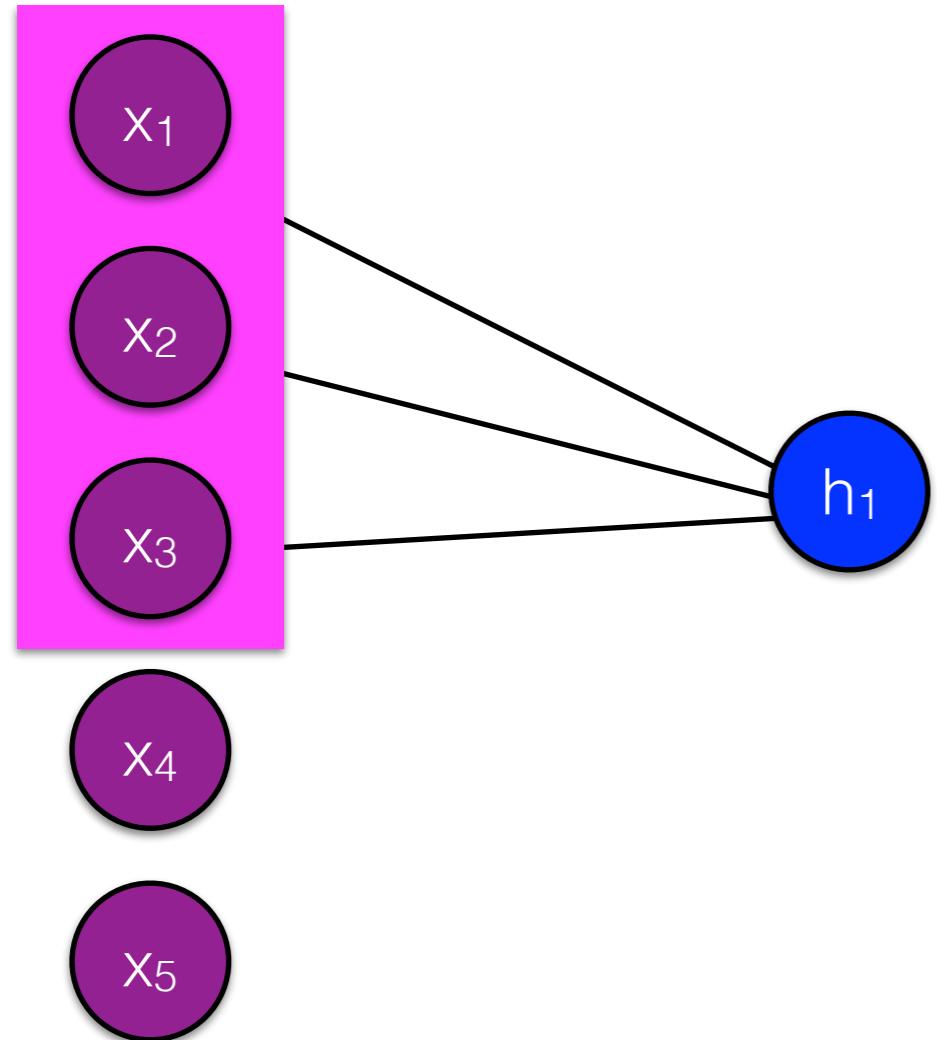
CNNs



$$h_1 = \sigma(x_1 W_1 + x_2 W_2 + x_3 W_3)$$

In a CNN, we just replace the input one-hot representation with the embedding

CNNs



	x	w					
x_1	<table border="1"><tr><td>0.4</td></tr><tr><td>0.8</td></tr><tr><td>1.2</td></tr><tr><td>-1.3</td></tr><tr><td>0.4</td></tr></table>	0.4	0.8	1.2	-1.3	0.4	w_1
0.4							
0.8							
1.2							
-1.3							
0.4							
x_2	<table border="1"><tr><td>0.2</td></tr><tr><td>-5.3</td></tr><tr><td>-1.2</td></tr><tr><td>5.3</td></tr><tr><td>0.4</td></tr></table>	0.2	-5.3	-1.2	5.3	0.4	w_2
0.2							
-5.3							
-1.2							
5.3							
0.4							
x_3	<table border="1"><tr><td>2.6</td></tr><tr><td>2.7</td></tr><tr><td>-3.2</td></tr><tr><td>6.2</td></tr><tr><td>1.9</td></tr></table>	2.6	2.7	-3.2	6.2	1.9	w_3
2.6							
2.7							
-3.2							
6.2							
1.9							

For dense input vectors (e.g., embeddings), full dot product

In a CNN, we just replace the input one-hot representation with the embedding

$$h_1 = \sigma(x_1 W_1 + x_2 W_2 + x_3 W_3)$$

Recurrent neural network

$$s_i = R(x_i, s_{i-1})$$

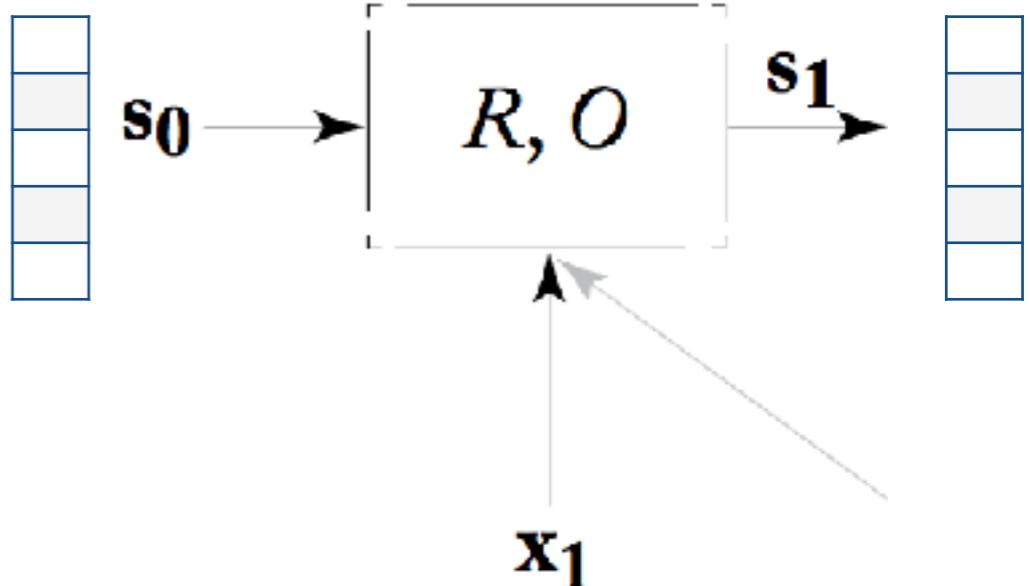
R is some function of the current input and previous state

$$y_i = O(s_i)$$

O is some function of the current state



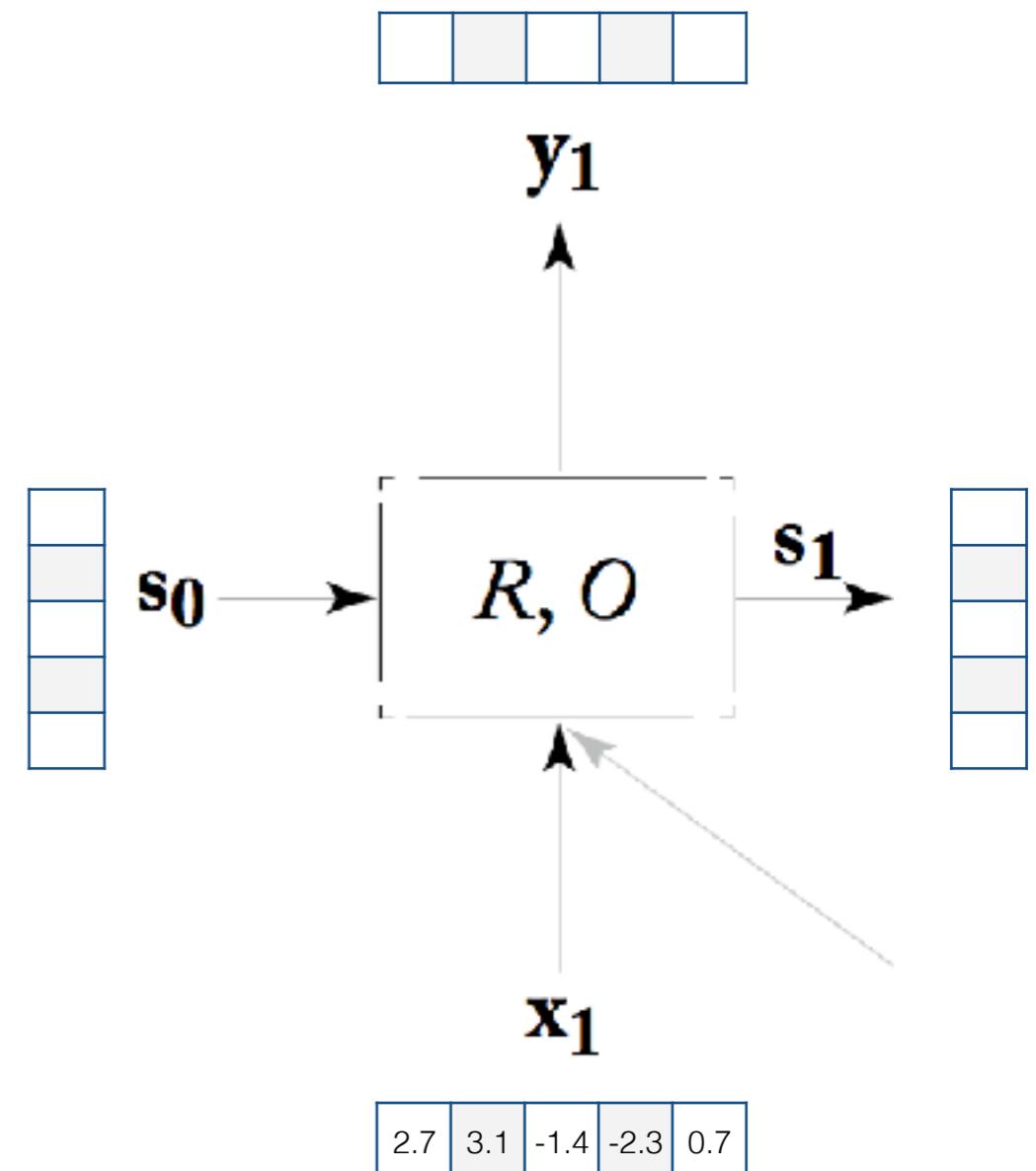
y₁



Recurrent neural network

$$s_i = R(x_i, s_{i-1})$$

$$y_i = O(s_i)$$



Equivalently

			1										
--	--	--	---	--	--	--	--	--	--	--	--	--	--

sparse one-hot vector

$$x \in \mathcal{R}^V$$

2.7	3.1	-1.4	-2.3	0.7

embeddings matrix

$$W \in \mathcal{R}^{V \times H}$$

2.7	3.1	-1.4	-2.3	0.7
-----	-----	------	------	-----

word embedding

$$xW \in \mathcal{R}^H$$

Recurrent neural network

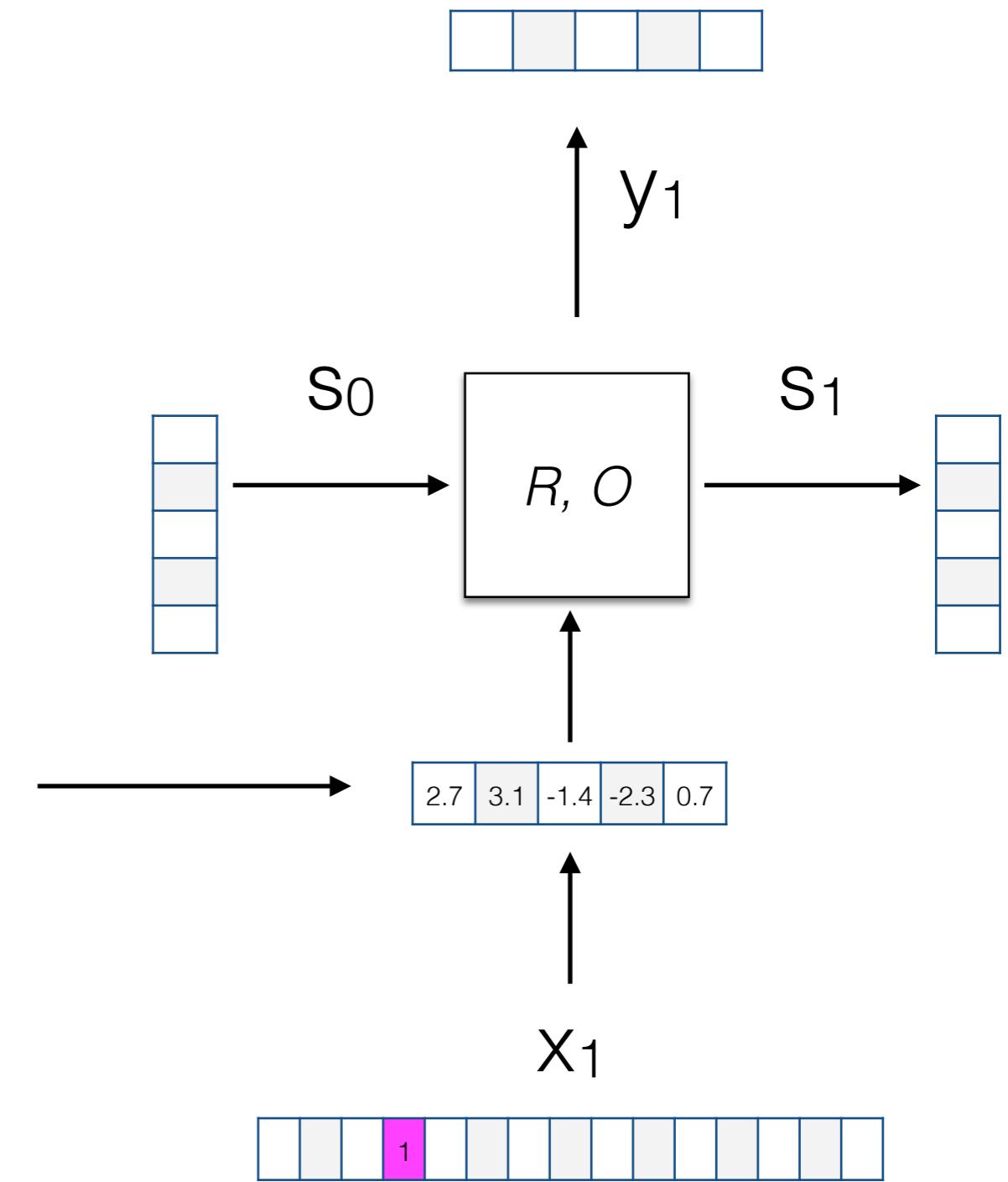
$$s_i = R(x_i, s_{i-1})$$

$$y_i = O(s_i)$$

$W_{\text{embedding}}$

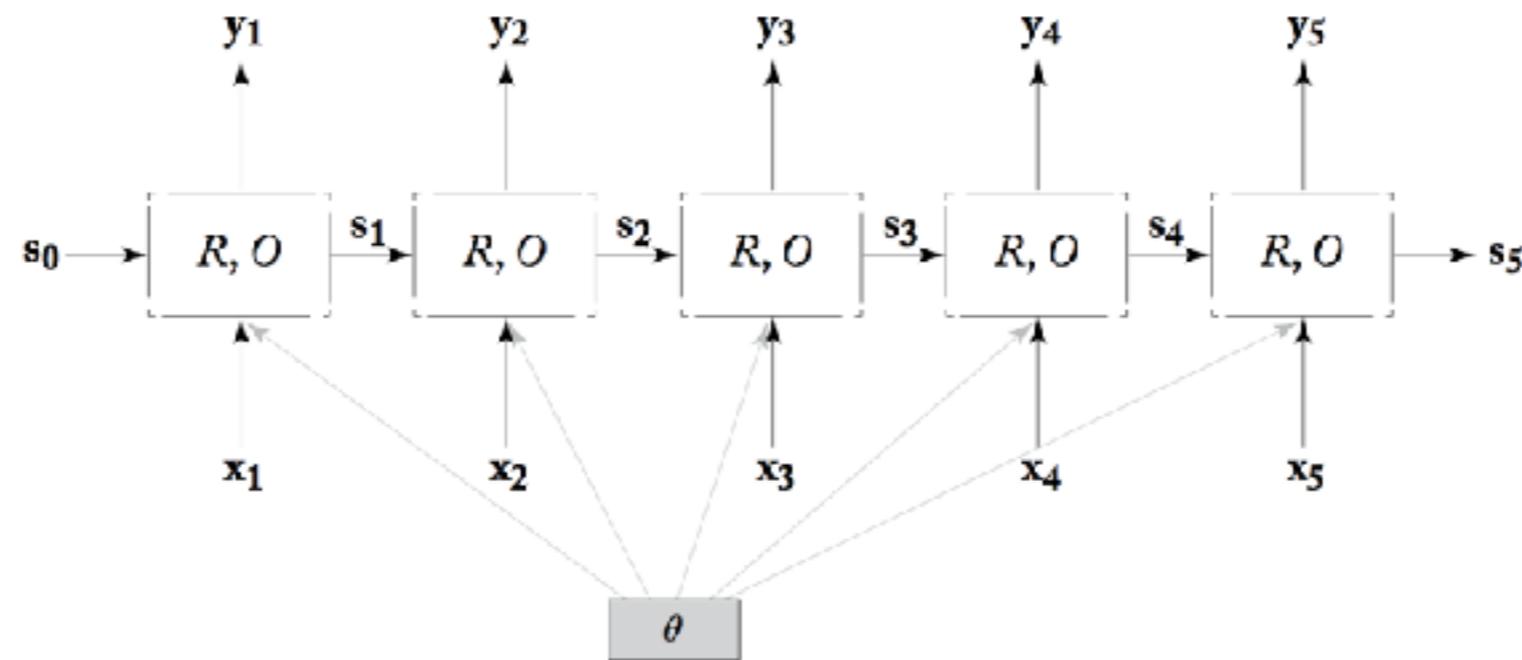
Embeddings are
parameters that can be
trained!

2.7	3.1	-1.4	-2.3	0.7	

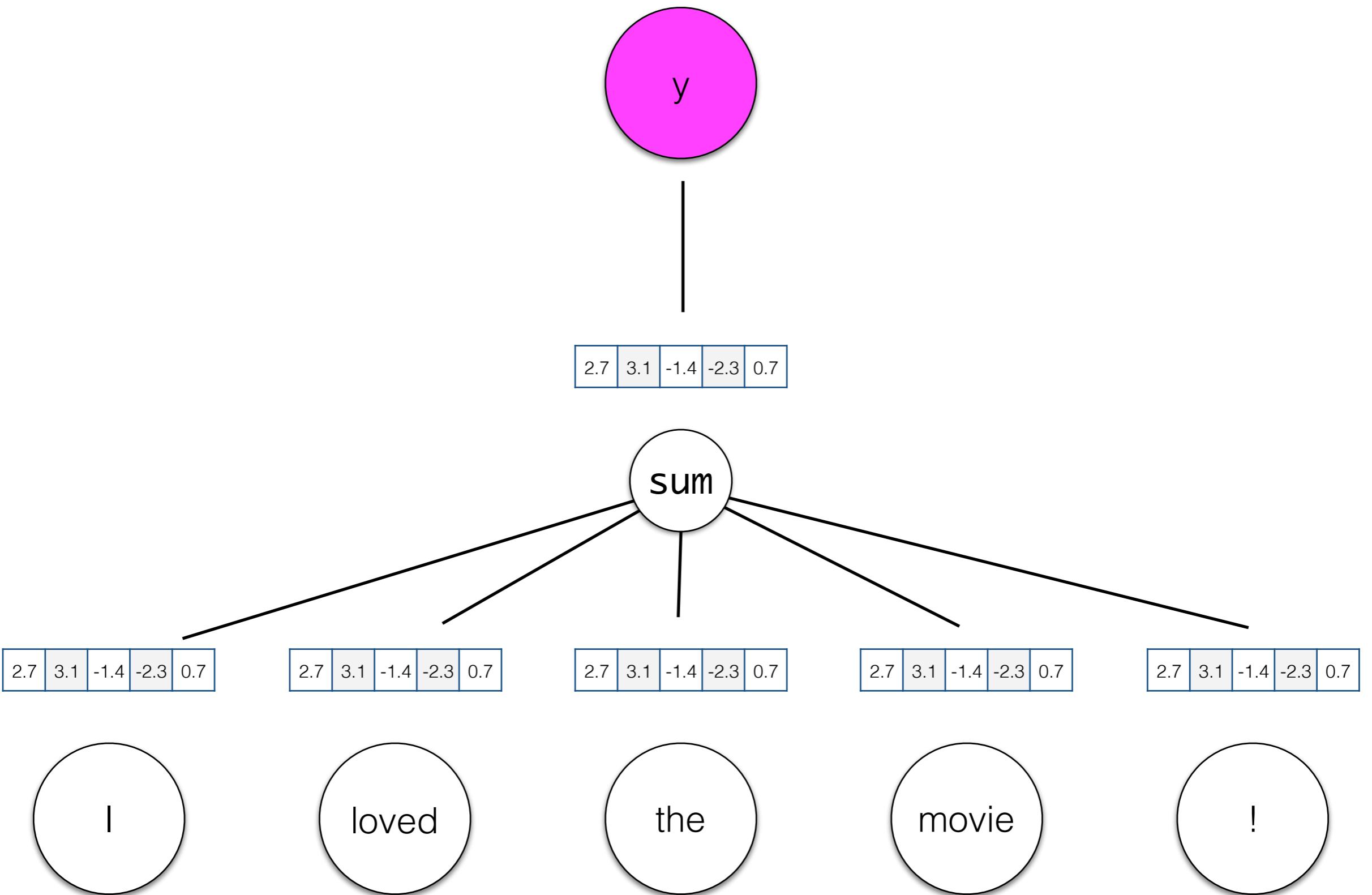


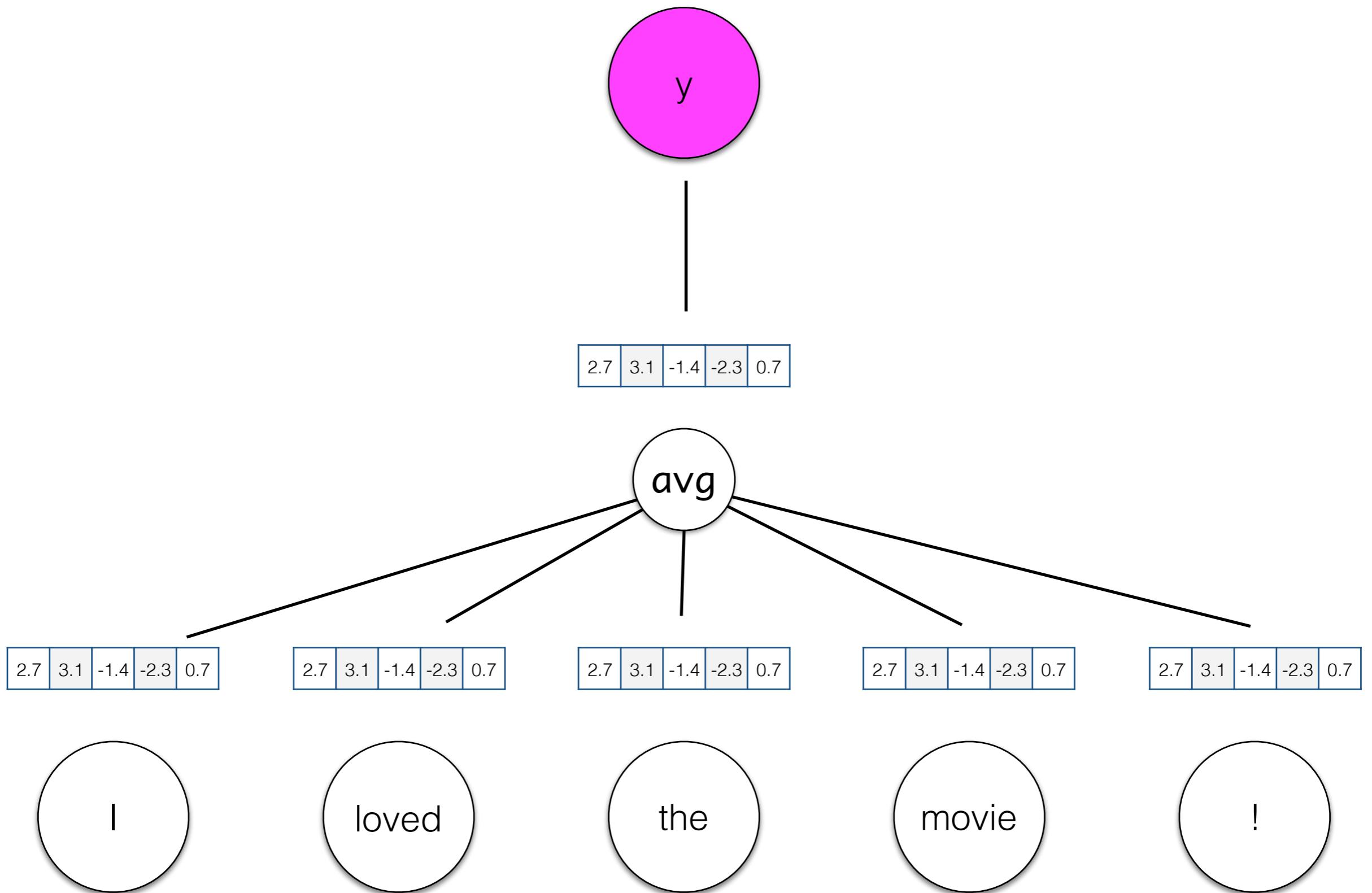
$$\frac{\partial L(\theta)_{y_1}}{\partial W_{embedding}}$$

we tried to prepare residents

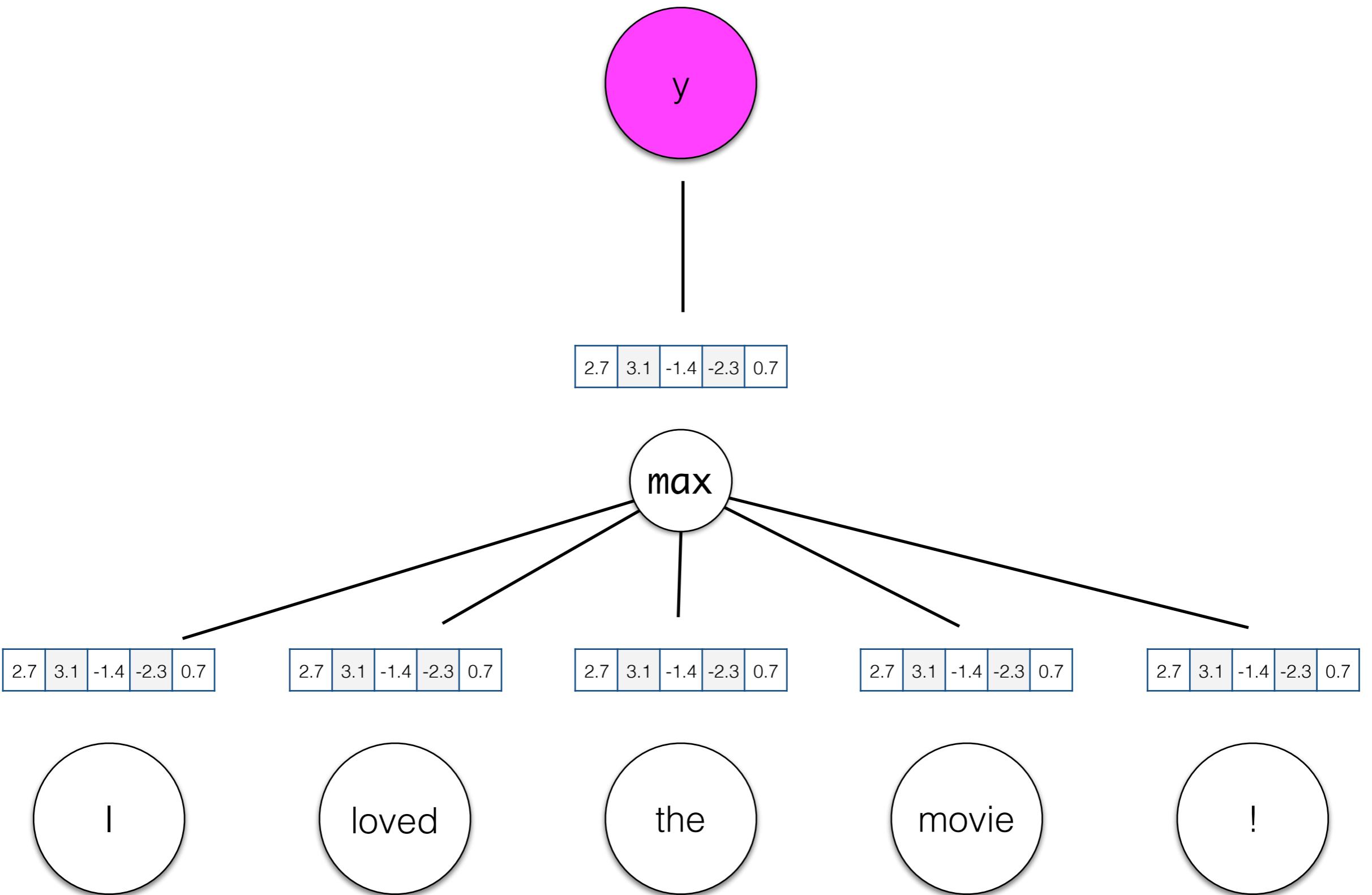


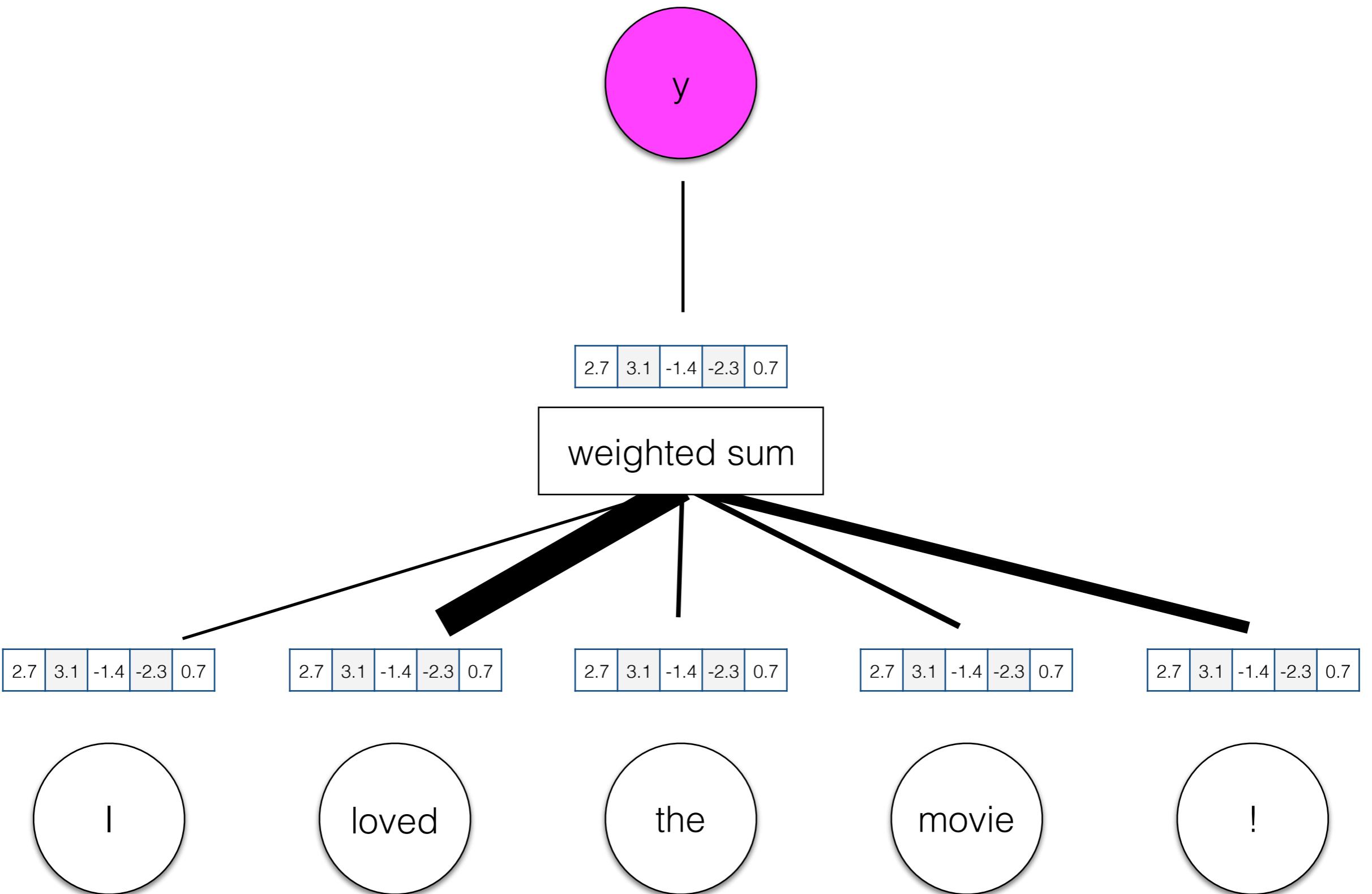
- We can optimize word embeddings for a specific task using by updating them using backpropagation as well.





Iyyer et al. (2015), “Deep Unordered Composition Rivals Syntactic Methods for Text Classification” (ACL)





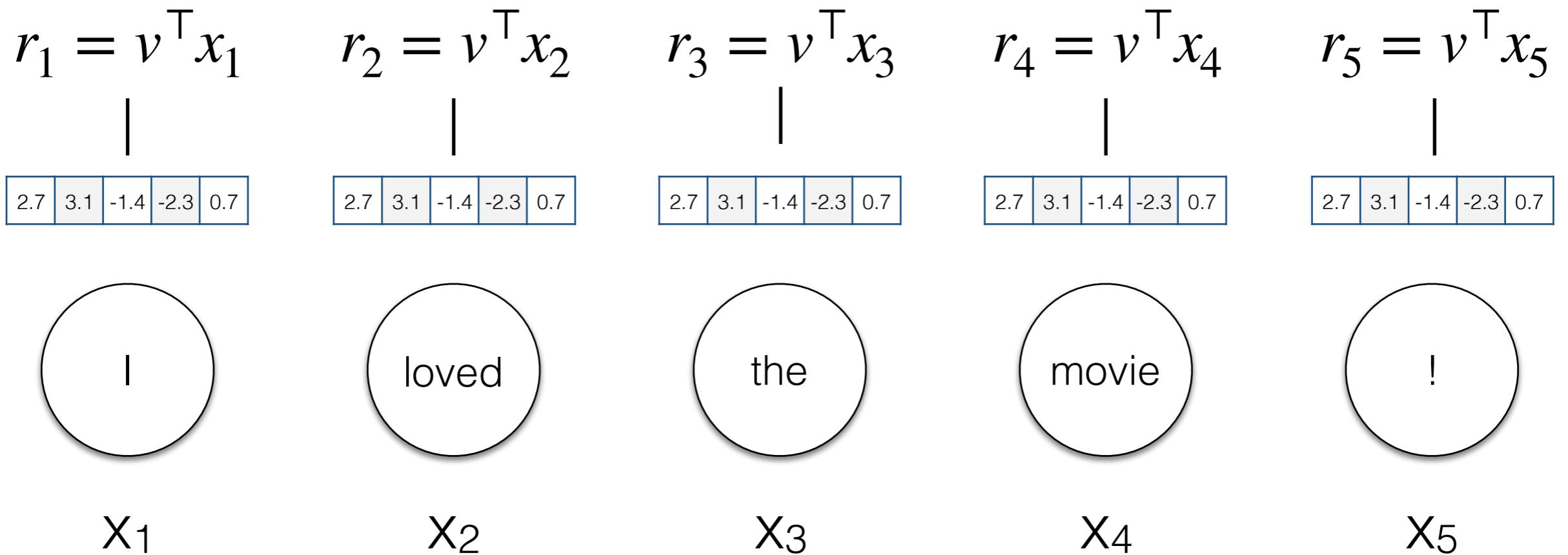
Attention

- Let's incorporate structure (and parameters) into a network that captures which elements in the input we should be **attending** to (and which we can ignore).

$v \in \mathcal{R}^H$

2.7	3.1	-1.4	-2.3	0.7
-----	-----	------	------	-----

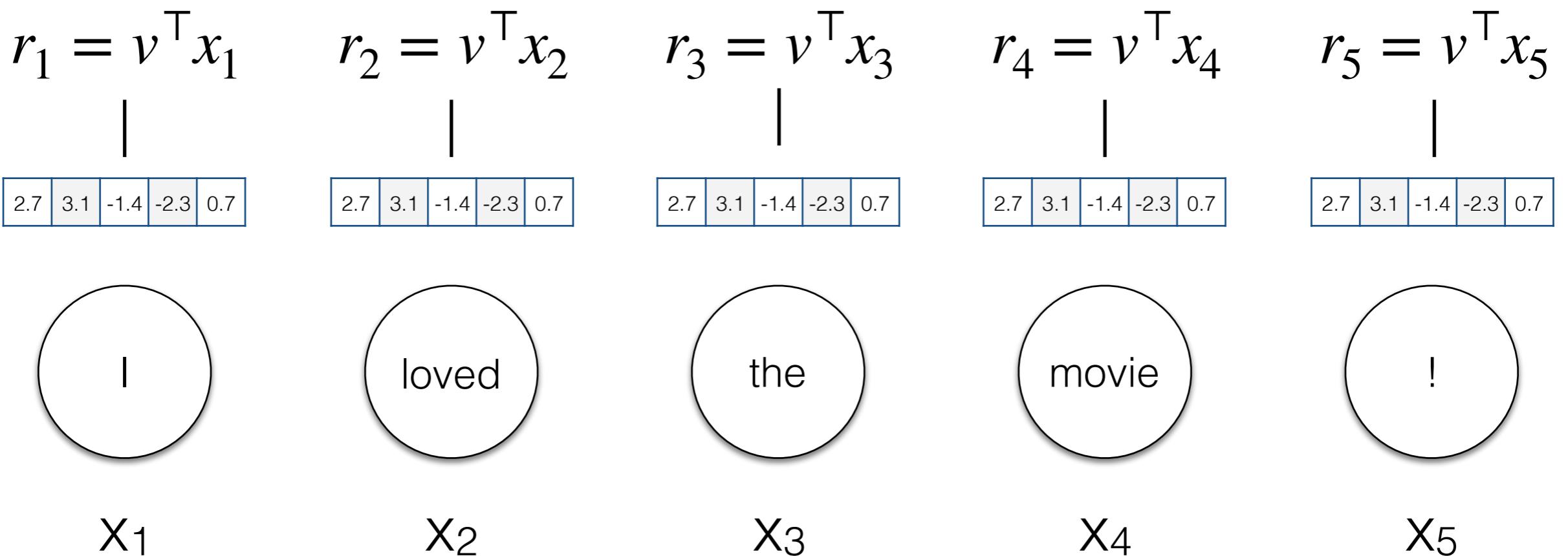
Define v to be a vector to be learned; think of it as an “important word” vector. The dot product here measures how similar each input vector is to that “important word” vector

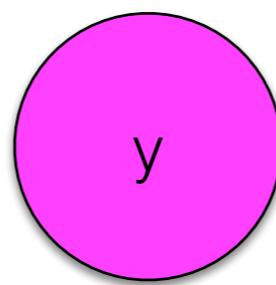


We convert r into a vector of normalized weights that sum to 1.

$$a = \text{softmax}(r)$$

$$r = [r_1, \dots, r_5]$$





2.7	3.1	-1.4	-2.3	0.7
-----	-----	------	------	-----

$$x_1a_1 + x_2a_2 + x_3a_3 + x_4a_4 + x_5a_5$$

weighted sum

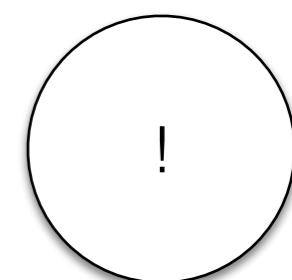
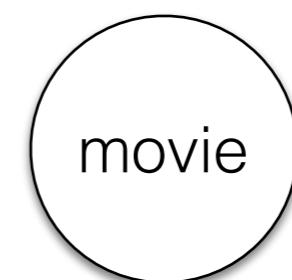
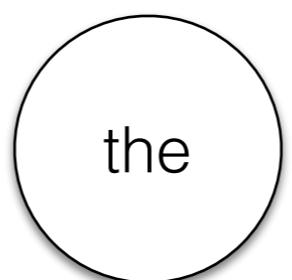
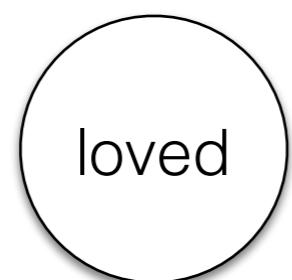
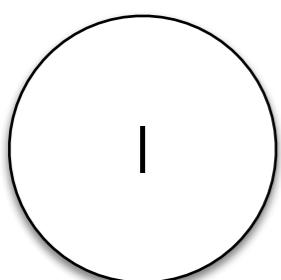
2.7	3.1	-1.4	-2.3	0.7
-----	-----	------	------	-----

2.7	3.1	-1.4	-2.3	0.7
-----	-----	------	------	-----

2.7	3.1	-1.4	-2.3	0.7
-----	-----	------	------	-----

2.7	3.1	-1.4	-2.3	0.7
-----	-----	------	------	-----

2.7	3.1	-1.4	-2.3	0.7
-----	-----	------	------	-----

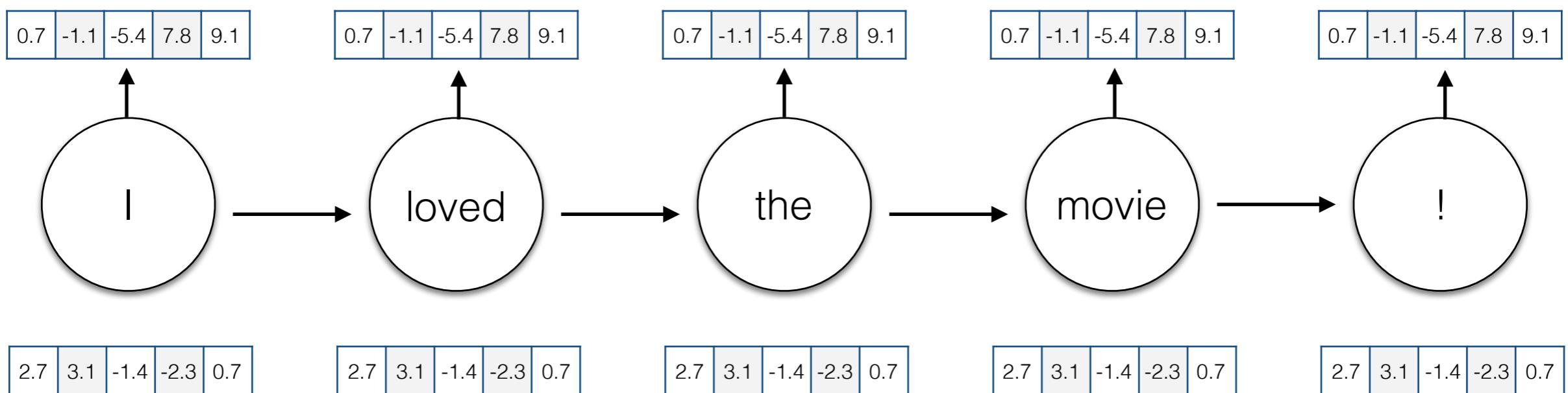


Attention

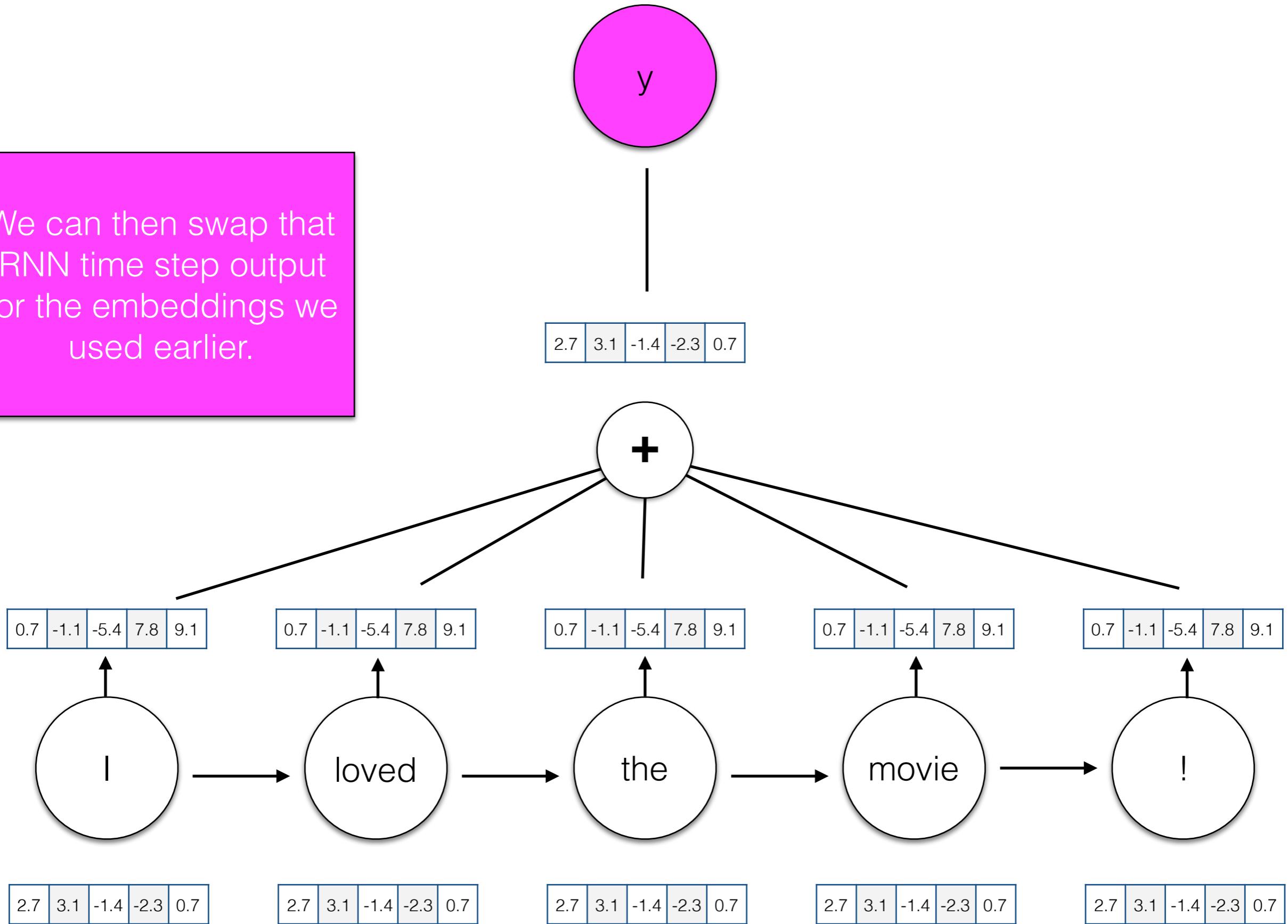
- Lots of variations on attention:
 - Linear transformation of x into before dotting with v
 - Non-linearities after each operation.
 - “Multi-head attention”: multiple v vectors to capture different phenomena that can be attended to in the input.
 - Hierarchical attention (sentence representation with attention over words + document representation with attention over sentences).

RNN

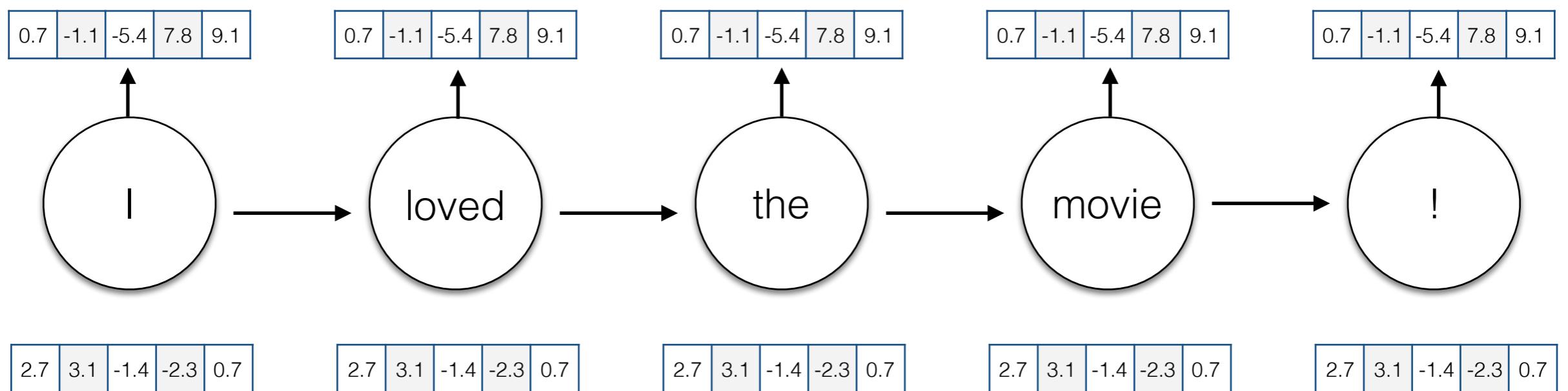
- With an RNN, we can generate a representation of the sequence as **seen through time t**.
- This encodes a representation of meaning specific to the local context a word is used in.



We can then swap that RNN time step output for the embeddings we used earlier.



What about the **future** context?

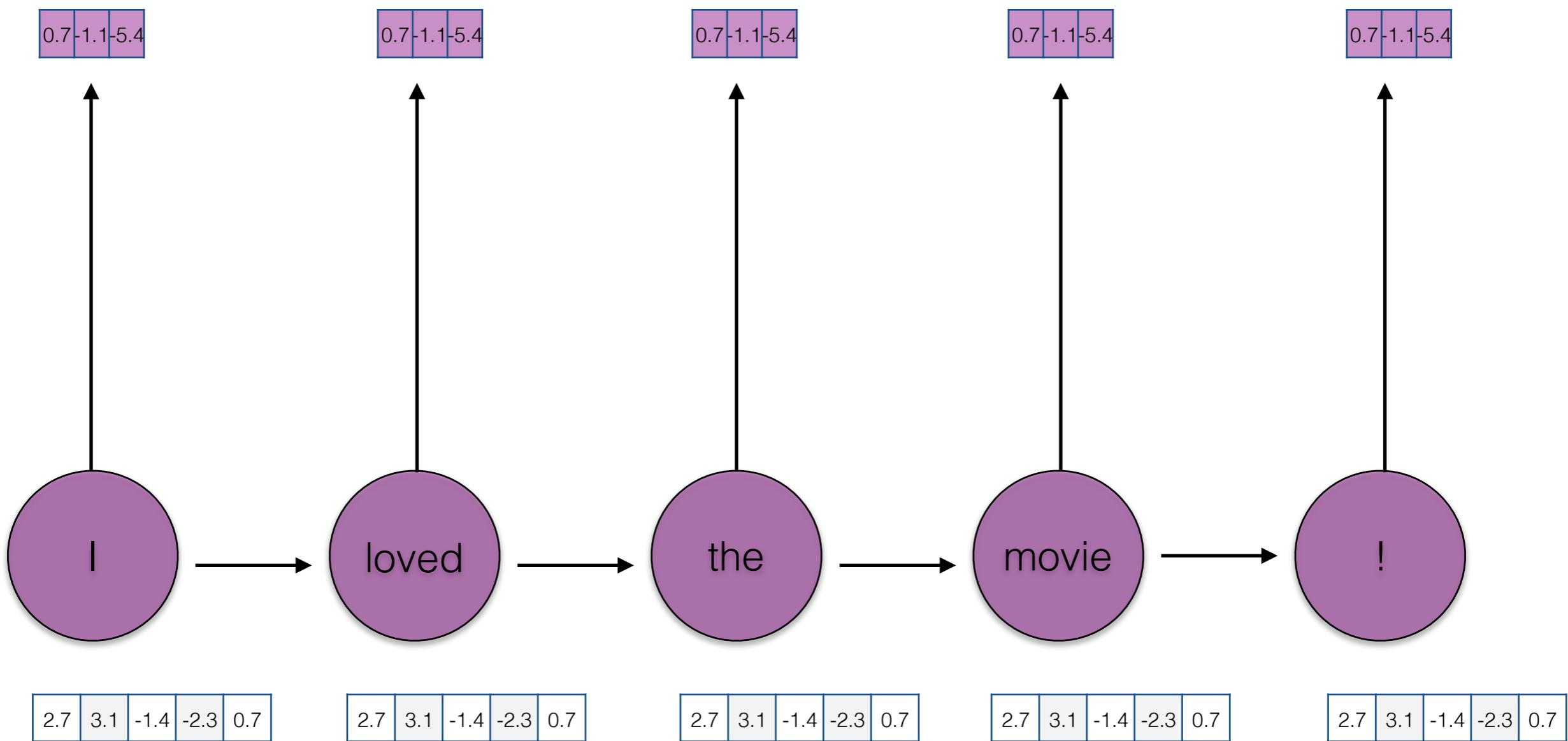


Bidirectional RNN

- A powerful alternative is make predictions conditioning both on the **past** and the **future**.
- Two RNNs
 - One running left-to-right
 - One right-to-left
- Each produces an output vector at each time step, which we concatenate

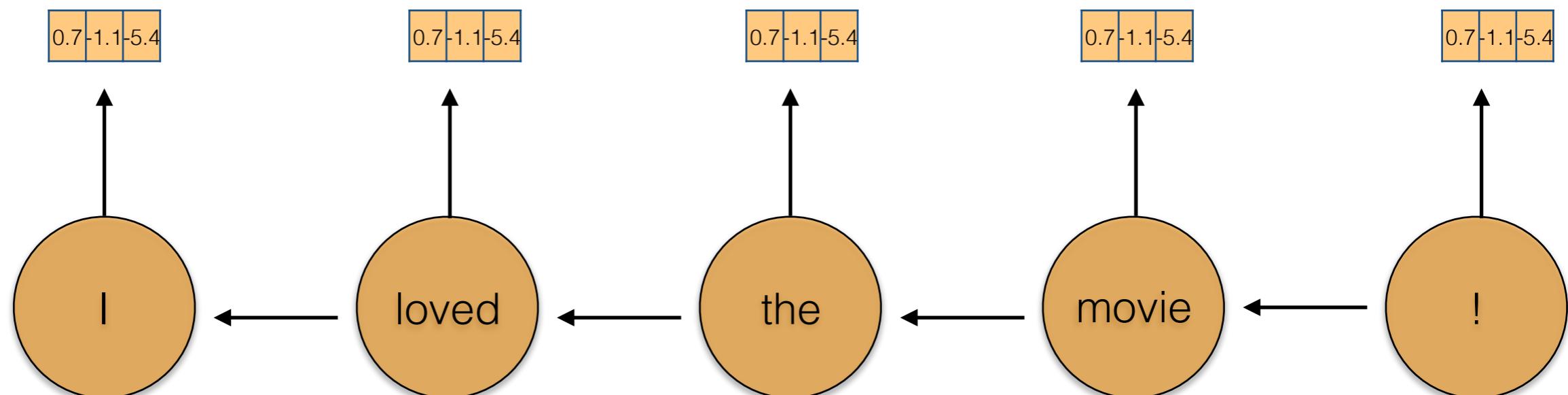
Bidirectional RNN

forward RNN



Bidirectional RNN

backward RNN



2.7	3.1	-1.4	-2.3	0.7
-----	-----	------	------	-----

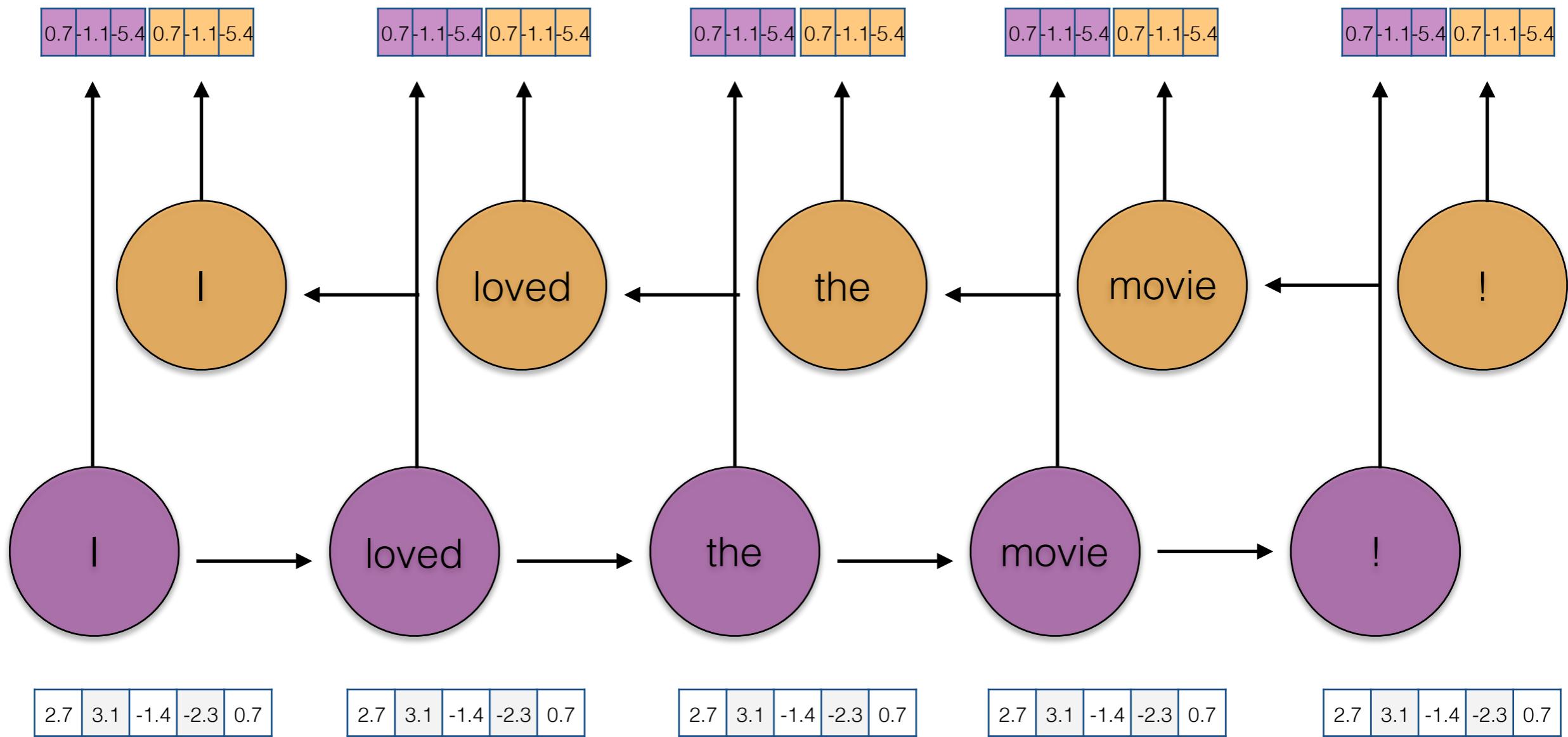
2.7	3.1	-1.4	-2.3	0.7
-----	-----	------	------	-----

2.7	3.1	-1.4	-2.3	0.7
-----	-----	------	------	-----

2.7	3.1	-1.4	-2.3	0.7
-----	-----	------	------	-----

2.7	3.1	-1.4	-2.3	0.7
-----	-----	------	------	-----

Bidirectional RNN



Bidirectional RNN

- The forward RNN and backward RNN each output a vector of size H at each time step, which we concatenate into a vector of size $2H$.
- The forward and backward RNN each have **separate parameters** to be learned during training.

Training BiRNNs

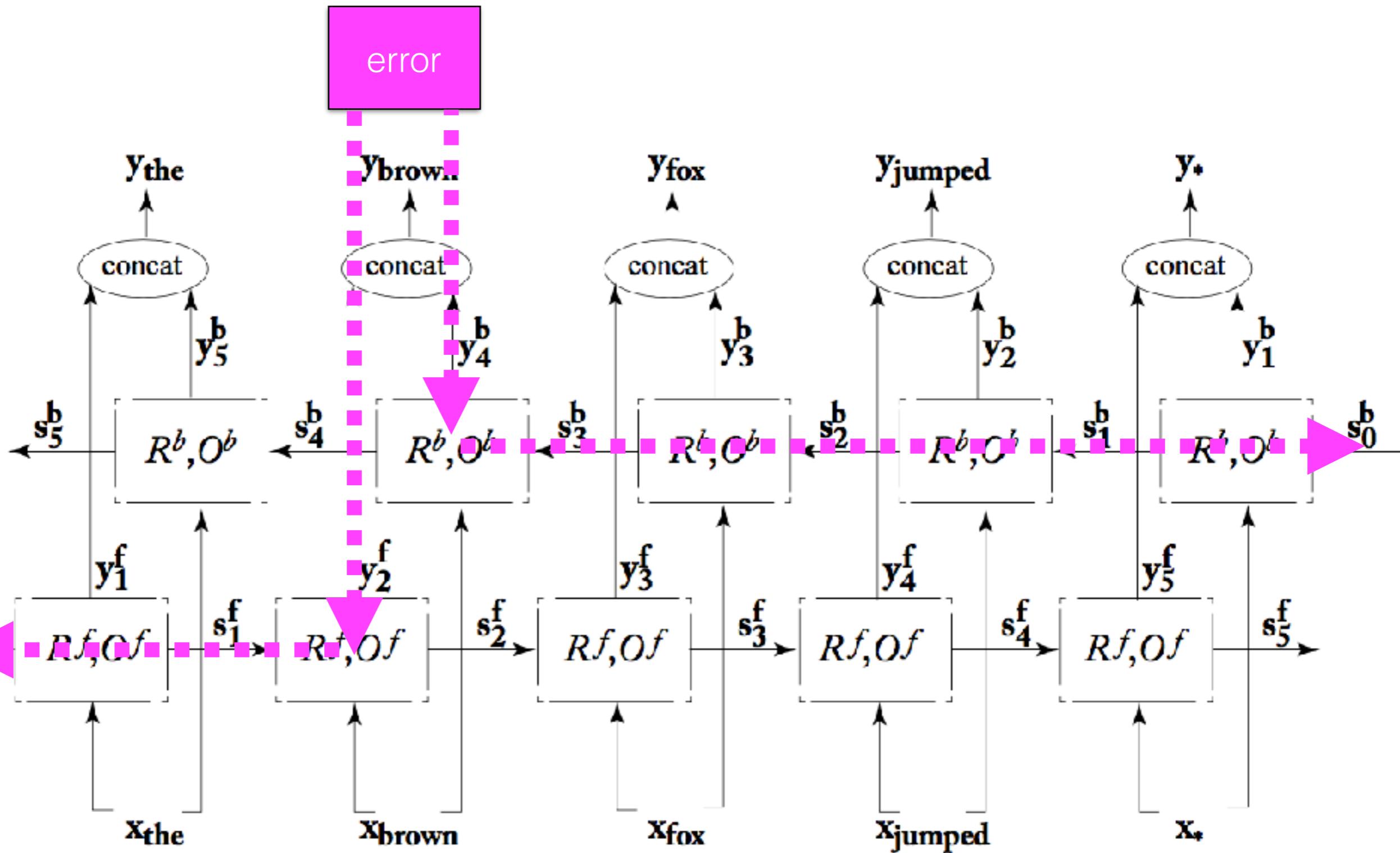
- Given this definition of an BiRNN:

$$s_b^i = R_b(x^i, s_b^{i+1}) = g(s_b^{i+1} \mathbf{W}_b^s + x^i \mathbf{W}_b^x + \mathbf{b}_b)$$

$$s_f^i = R_f(x^i, s_f^{i-1}) = g(s_f^{i-1} \mathbf{W}_f^s + x^i \mathbf{W}_f^x + \mathbf{b}_f)$$

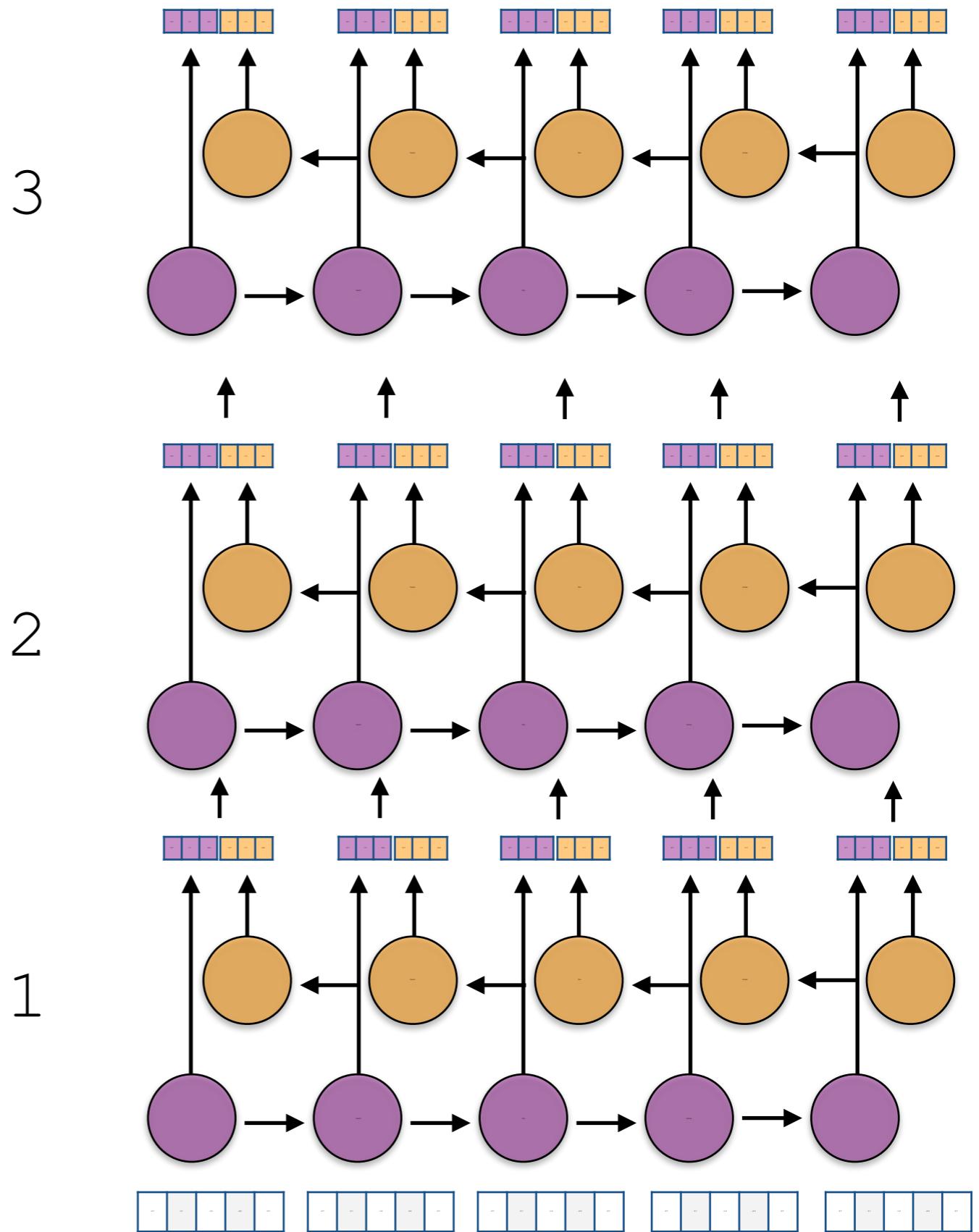
$$y_i = \text{softmax} \left([s_f^i; s_b^i] \mathbf{W}^o + \mathbf{b}^o \right)$$

- We have 8 sets of parameters to learn (3 for each RNN + 2 for the final layer)



Stacked RNN

- Multiple RNNs, where the output of one layer becomes the input to the next.



ELMo

- Peters et al. (2018), “Deep Contextualized Word Representations” (NAACL)
- Big idea: transform the representation of a word (e.g., from a static word embedding) to be sensitive to its local context in a sentence and optimized for a specific NLP task.
- Output = word representations that can be plugged into just about any architecture a word embedding can be used.

ELMo

- Peters et al. (2018), “Deep Contextualized Word Representations” (NAACL)
- Train a bidirectional RNN language model with L layers on a bunch of text.
- Learn parameters to combine the RNN output across all layers for each word in a sentence for a specific task (NER, semantic role labeling, question answering etc.). Large improvements over SOTA for lots of NLP problems.

ELMo

TASK	PREVIOUS SOTA		OUR BASELINE	ELMO + BASELINE	INCREASE (ABSOLUTE/ RELATIVE)
SQuAD	Liu et al. (2017)	84.4	81.1	85.8	4.7 / 24.9%
SNLI	Chen et al. (2017)	88.6	88.0	88.7 ± 0.17	0.7 / 5.8%
SRL	He et al. (2017)	81.7	81.4	84.6	3.2 / 17.2%
Coref	Lcc et al. (2017)	67.2	67.2	70.4	3.2 / 9.8%
NER	Peters et al. (2017)	91.93 ± 0.19	90.15	92.22 ± 0.10	2.06 / 21%
SST-5	McCann et al. (2017)	53.7	51.4	54.7 ± 0.5	3.3 / 6.8%

Table 1: Test set comparison of ELMo enhanced neural models with state-of-the-art single model baselines across six benchmark NLP tasks. The performance metric varies across tasks – accuracy for SNLI and SST-5; F₁ for SQuAD, SRL and NER; average F₁ for Coref. Due to the small test sizes for NER and SST-5, we report the mean and standard deviation across five runs with different random seeds. The “increase” column lists both the absolute and relative improvements over our baseline.

- Word embeddings can be substituted for one-hot encodings in many models (MLP, CNN, RNN, logistic regression).
- Subword embeddings allow you to create embeddings for word not present in training data; require much less data to train.
- Attention gives us a mechanism to learn which parts of a sequence to pay attention to more in forming a representation of it.
- BiLSTMs can transform word embeddings to be sensitive to their use **in context**.