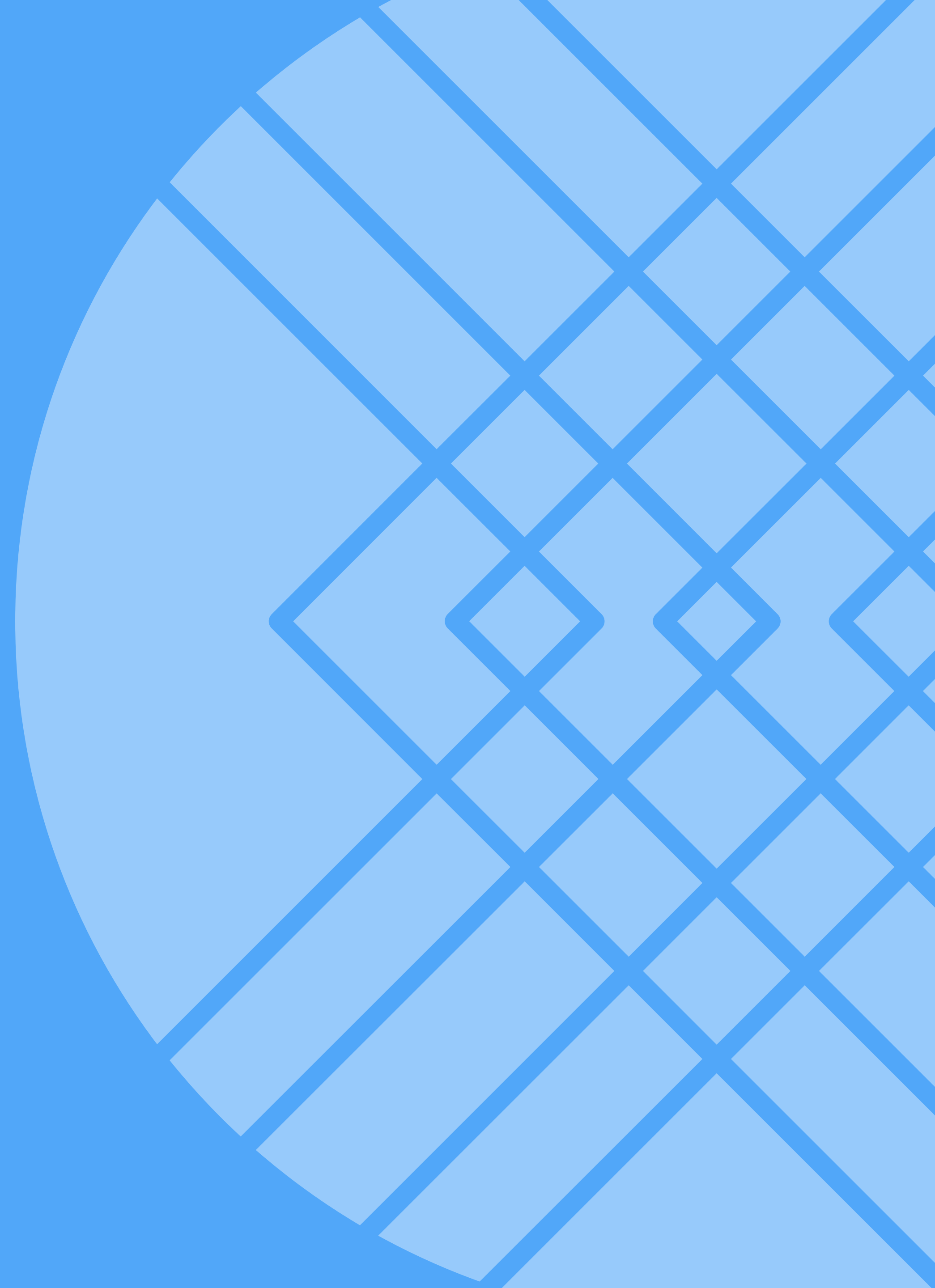


# **Word2Vec** and all the things

Data Labs @  
StitchFix



# About



@chrisemoody

Former astrophysicist, supercomputing

Data Labs at Stitch Fix

Fun stuff with Gaussian Processes,  
t-SNE & word2vec

# Credit

Large swathes of this talk are from previous presentations by:

- [Tomas Mikolov](#)
- [Christopher Olah](#)
- [Radim Rehurek](#)
- [Omer Levy & Yoav Goldberg](#)
- [Richard Socher](#)
- [Xin Rong](#)

**1**

**word2vec**

**2**

**Stitch Fix**

**3 Fun Stuff**

## word2vec

1. *king - man + woman = queen*
2. Learns from raw text
3. Huge splash in NLP world
4. Pretty simple algorithm
5. Comes pretrained

Predict current word given previous words

$$P(W) = \prod_i P(w_i | w_1 \dots w_{i-1})$$

What if we estimate  $P(w)$  empirically?

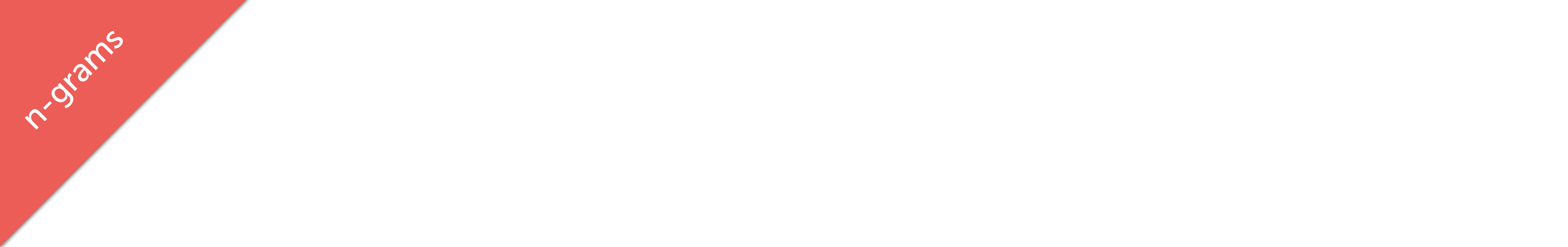
How far back should  $w_1 \dots w_{i-1}$  go?



n-grams

**Q**

$$\operatorname{argmax}[p(w \mid w_0 \dots w_k)]$$



n-grams

**Q**

$\text{argmax}[p(w \mid \text{'now', 'with', 'a', 'new', 'model'})]$



If context is too **long** you overfit....

**Q**

$\text{argmax}[p(w \mid \text{'now', 'with', 'a', 'new', 'model'})]$

**A**

'GM'

now with a new model GM

If context is too **long** you overfit....

**Q**

$\text{argmax}[p(w \mid \text{'with', 'a', 'new', 'model', 'GM'})]$

**A**

'is'

now with a new model GM is

If context is too **long** you overfit....

**Q**

$\text{argmax}[p(w \mid 'a', 'new', 'model', 'GM', 'is')]$

**A**

'trying'

now with a new model GM is trying

If context is too **long** you overfit....

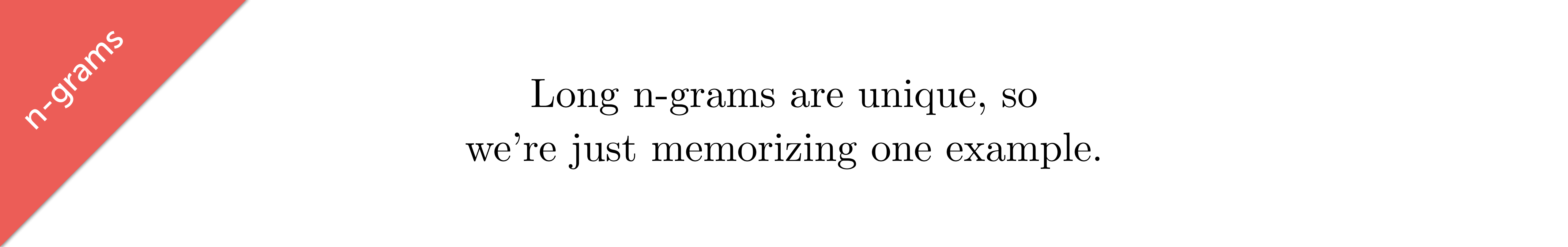
**Q**

$\text{argmax}[p(w \mid \text{'new', 'model', 'GM', 'is', 'trying'})]$

**A**

'to'

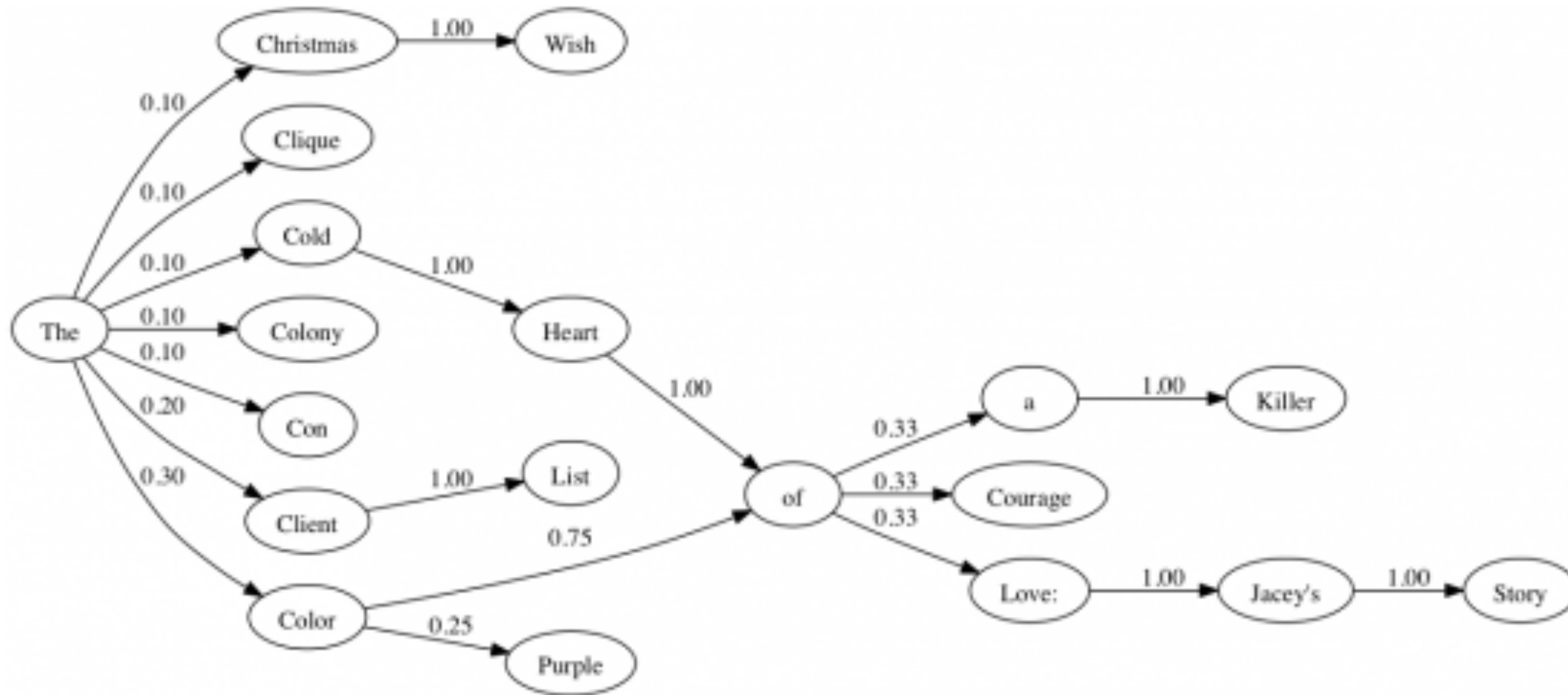
now with a new model GM is trying to



Long n-grams are unique, so we're just memorizing one example.

“prestige. Now, with a new model and a move, G.M. is trying to recapture the swagger.”

If context is too **short** your sentences don't make sense...



...just because the current word *only* depends on the previous word.  
(but it usually works *grammatically*)

# Garkov

by Josh Millard  
via Jim Davis



(this is how most internet chatbots work)

\*checkout [Garkov](#)

**Idea:** replace one-hot coded words with **dense vectors**\*

$$V_{\text{run}} = [0, 0, 1, 0, 0]$$

$$V_{\text{running}} = [0, 1, 0, 0, 0]$$



$$V_{\text{run}} = [0.2, 0.3, -.1, -.9, 0.0]$$

$$V_{\text{running}} = [0.7, 1, -.7, 0.5, 0.1]$$

\*or distributed representations Bengio et al.



**Idea:** replace one-hot coded words with **dense vectors**\*

$$V_{\text{cat}} = [0.2, 0.3, -.1, -.9, 0.0]$$

$$V_{\text{dog}} = [0.7, 1, -.7, 0.5, 0.1]$$

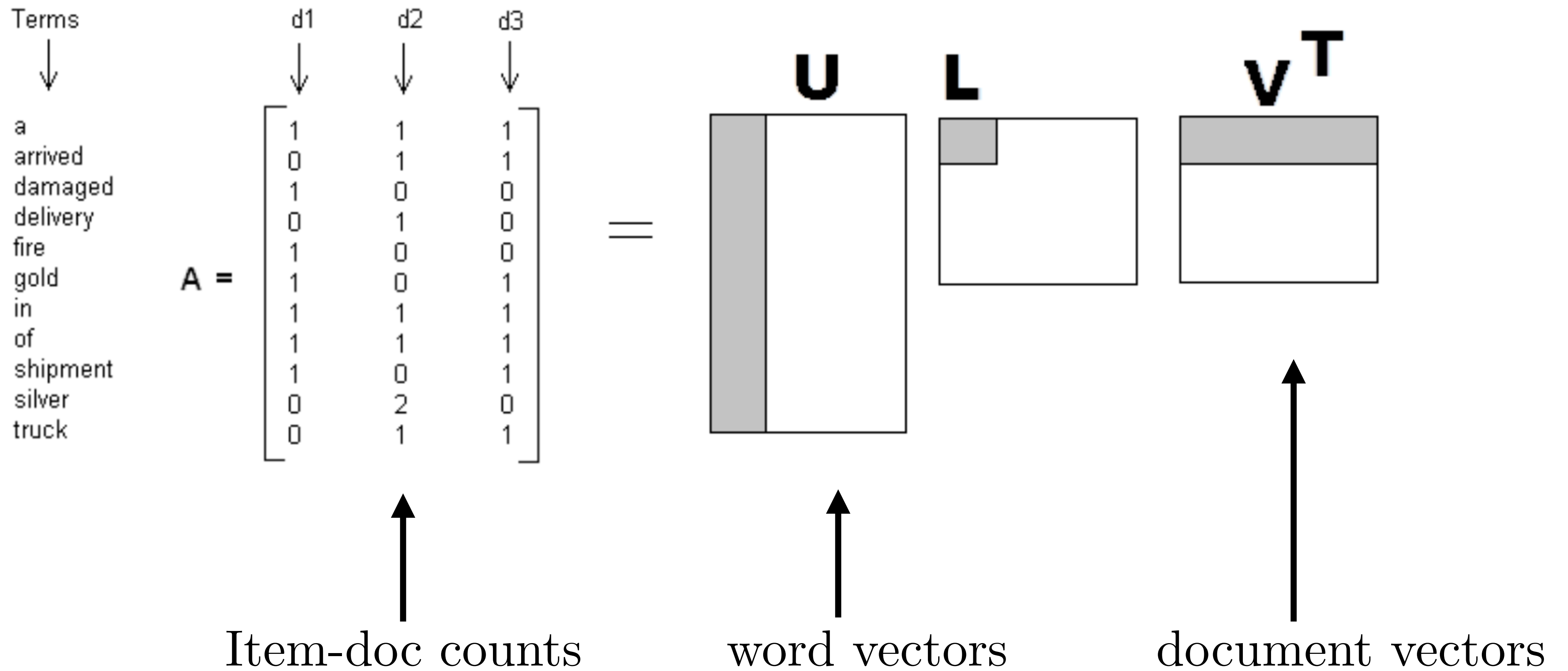
- *Locality* — properties change gradually from word to word  
*quickly ~ quick*
- *Regularity* — directions will be meaningful & consistent  
*king - man + woman = queen*

\*see more by Rumelhart, Collobert, Bengio

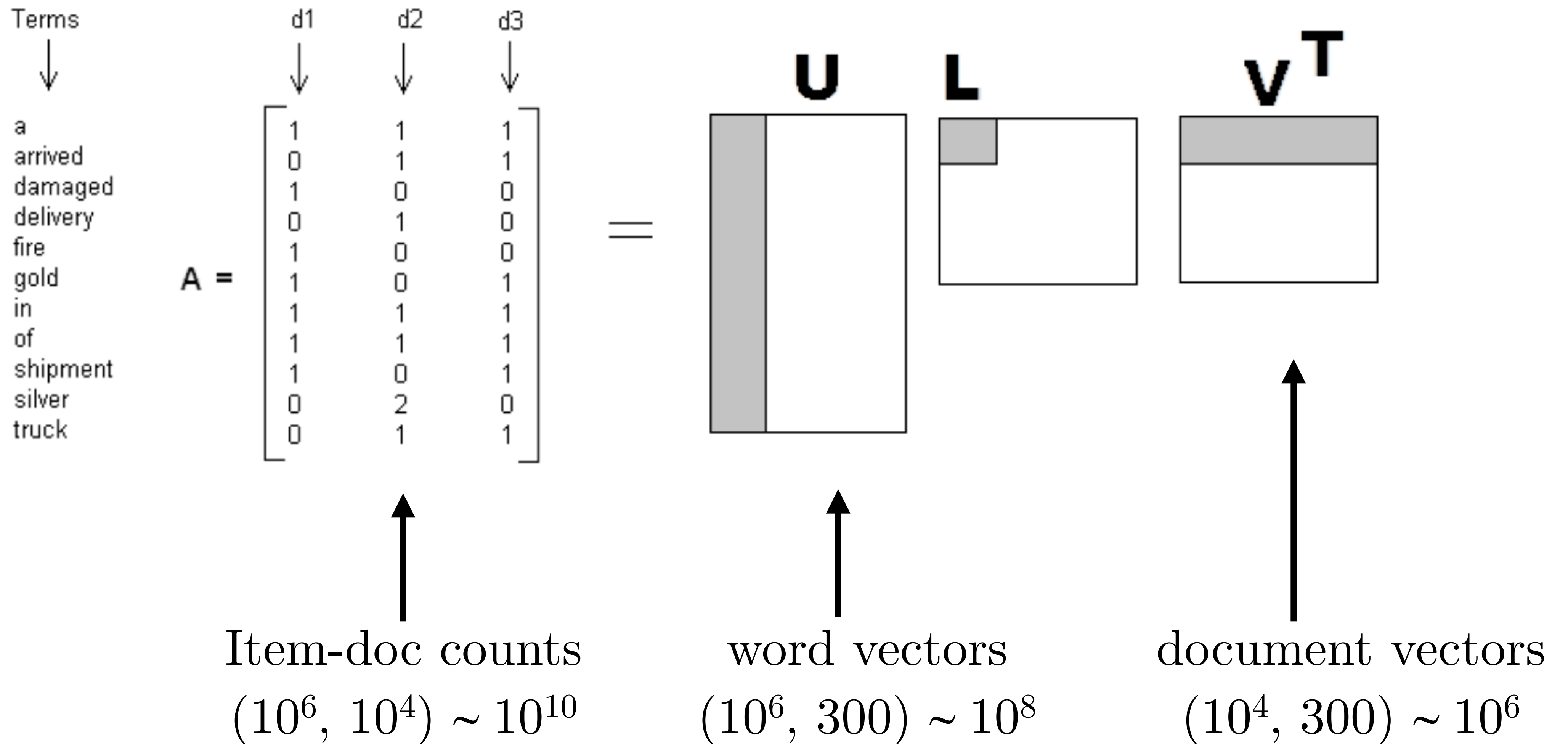
Build a co-occurrence matrix.

Terms	d1	d2	d3
a	1	1	1
arrived	0	1	1
damaged	1	0	0
delivery	0	1	0
fire	1	0	0
gold	1	0	1
in	1	1	1
of	1	1	1
shipment	1	0	1
silver	0	2	0
truck	0	1	1

Let's SVD the co-occurrence matrix.



Let's SVD the co-occurrence matrix.



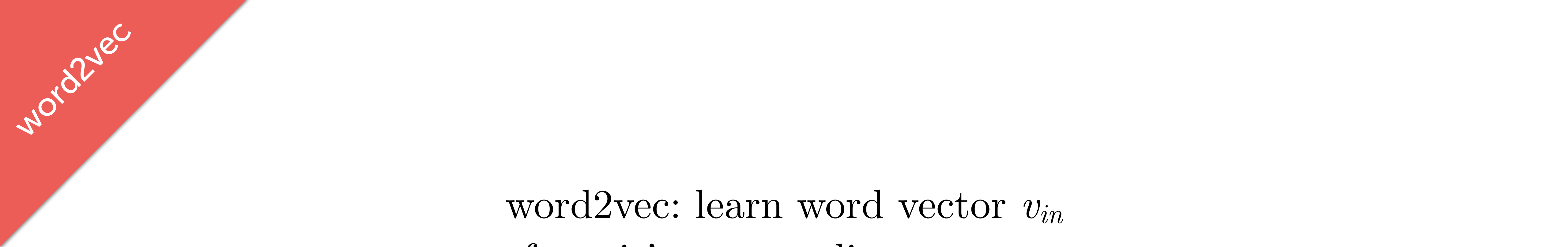
Doesn't scale.

$$\sim O(nm^2)$$

Hard to read more than a few millions of documents.

# word2vec

1. Set up an objective function
2. Randomly initialize vectors
3. Do gradient descent



word2vec

word2vec: learn word vector  $v_{in}$   
from it's surrounding context

$$v_{in}$$

“The fox jumped **over** the lazy dog”

Maximize the likelihood of seeing this **context** given the word *over*.

$$P(\textit{the}|\textit{over})$$

$$P(\textit{fox}|\textit{over})$$

$$P(\textit{jumped}|\textit{over})$$

$$P(\textit{the}|\textit{over})$$

$$P(\textit{lazy}|\textit{over})$$

$$P(\textit{dog}|\textit{over})$$

...instead of maximizing the likelihood of co-occurrence counts.



What should this be?

$$P(\text{fox}|\text{over})$$

Should depend on the word vectors.

$$P(\text{fox}|\text{over})$$



$$P(v_{\text{fox}}|v_{\text{over}})$$

Twist: we have *two* vectors for every word.  
Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

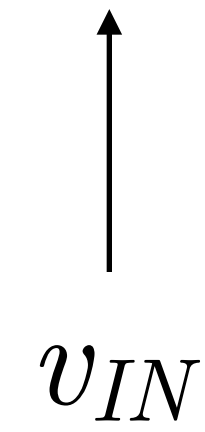
“The fox jumped **over** the lazy dog”

Twist: we have *two* vectors for every word.  
Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped **over** the lazy dog”



Twist: we have *two* vectors for every word.  
Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped **over** the lazy dog”



$v_{OUT}$



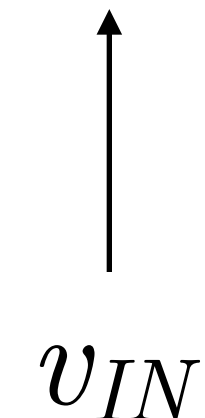
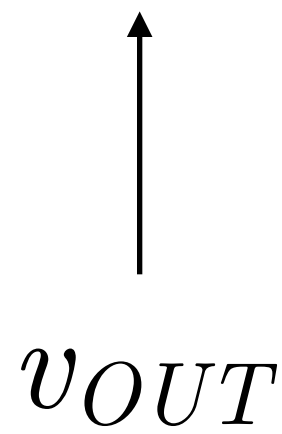
$v_{IN}$

Twist: we have *two* vectors for every word.  
Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped **over** the lazy dog”

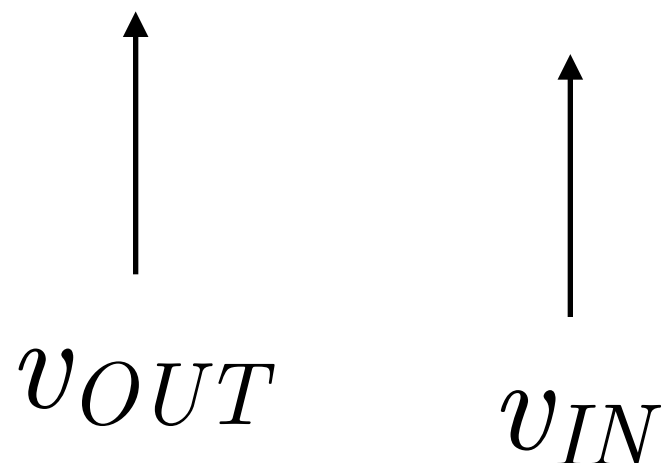


Twist: we have *two* vectors for every word.  
Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped **over** the lazy dog”

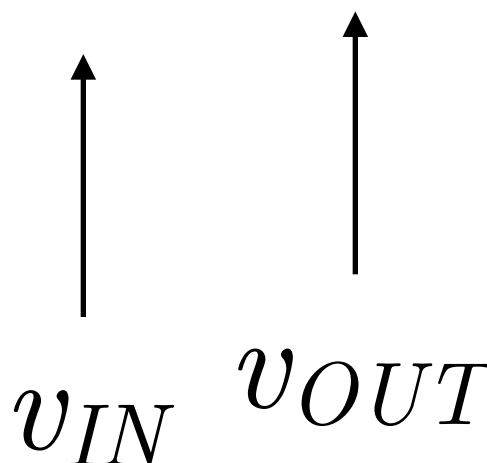


Twist: we have *two* vectors for every word.  
Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped **over** the lazy dog”



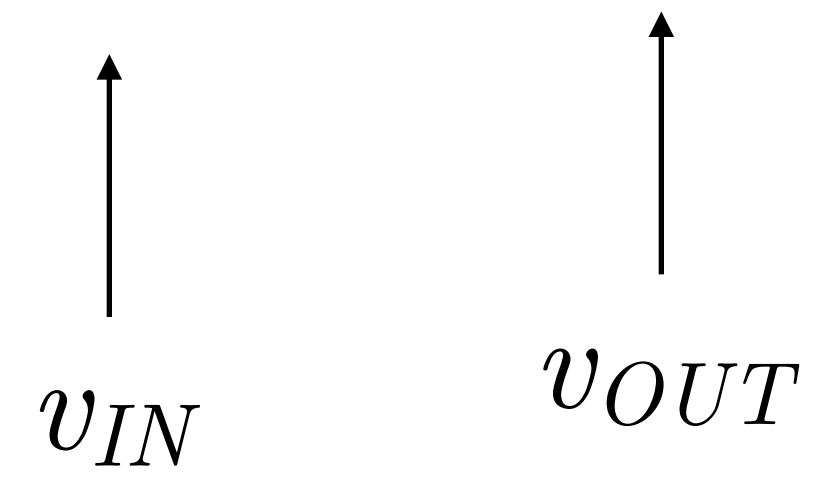


Twist: we have *two* vectors for every word.  
Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped **over** the lazy dog”

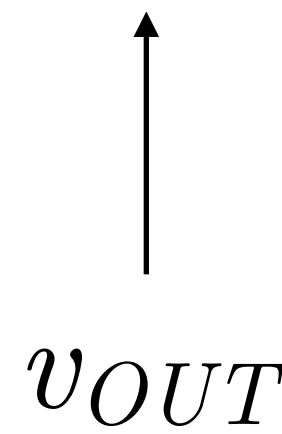
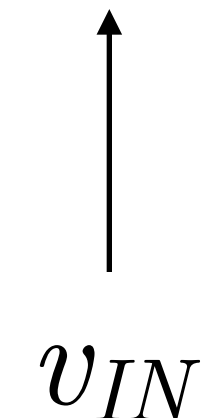


Twist: we have *two* vectors for every word.  
Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped **over** the lazy dog”

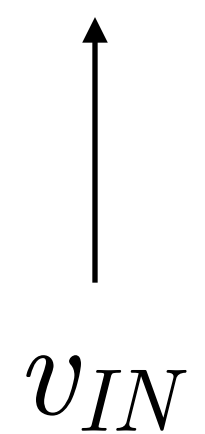


Twist: we have *two* vectors for every word.  
Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped over **the** lazy dog”

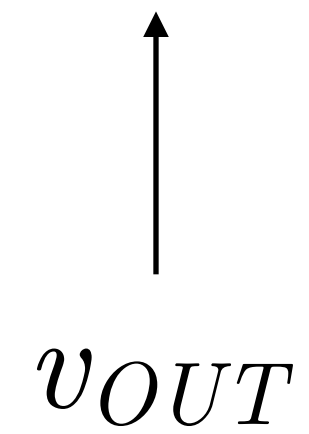


Twist: we have *two* vectors for every word.  
Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped over **the** lazy dog”

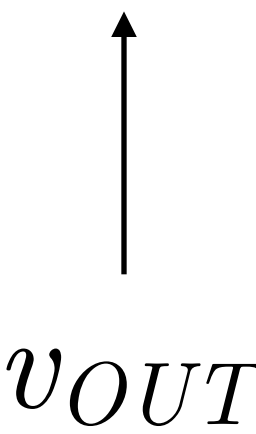


Twist: we have *two* vectors for every word.  
Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped over **the** lazy dog”



Twist: we have *two* vectors for every word.  
Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped over **the** lazy dog”

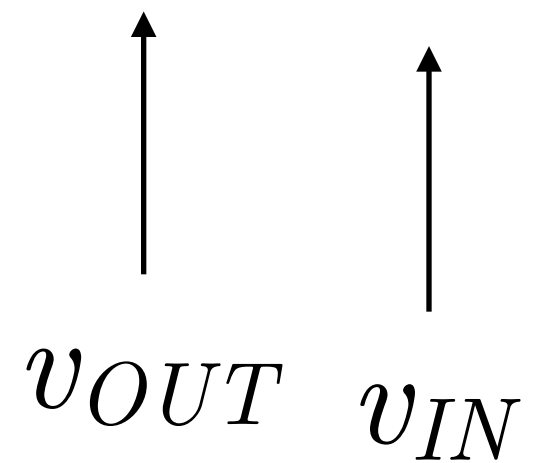


Twist: we have *two* vectors for every word.  
Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped over **the** lazy dog”

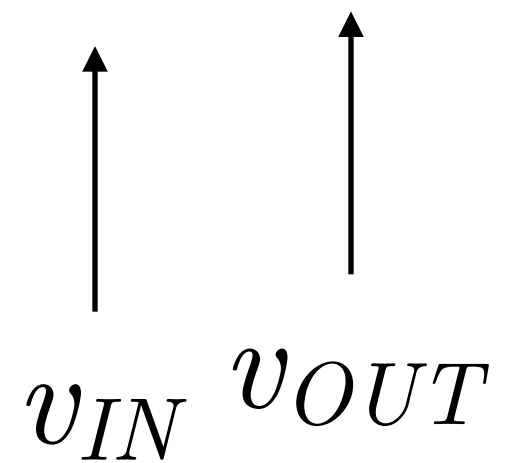


Twist: we have *two* vectors for every word.  
Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped over **the** lazy dog”



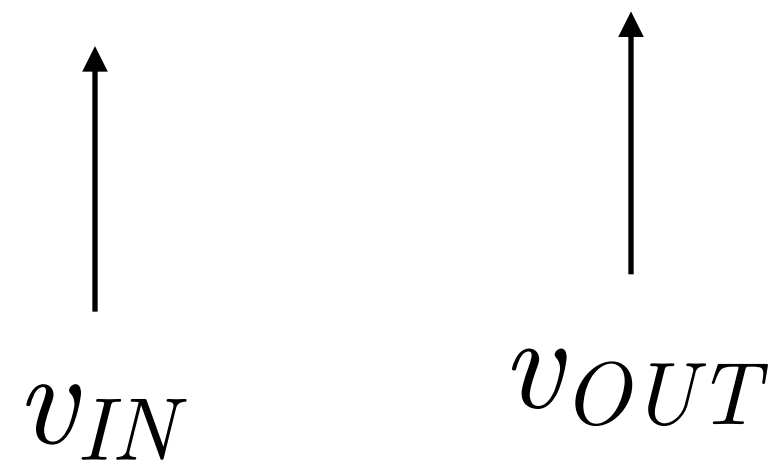


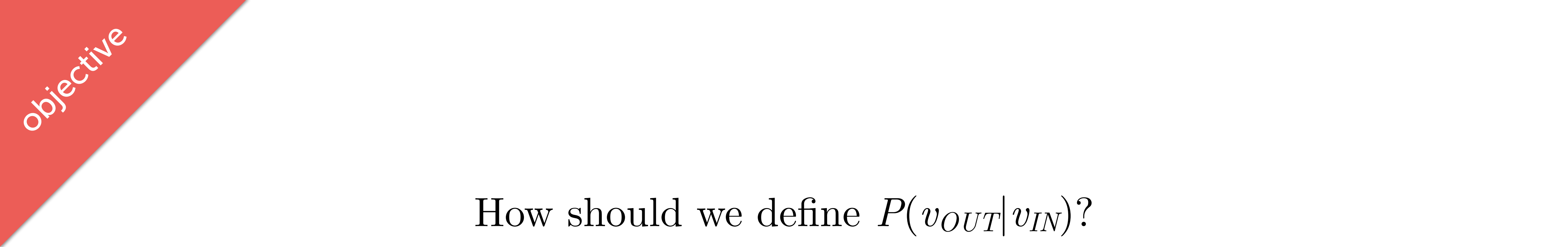
Twist: we have *two* vectors for every word.  
Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped over **the** lazy dog”



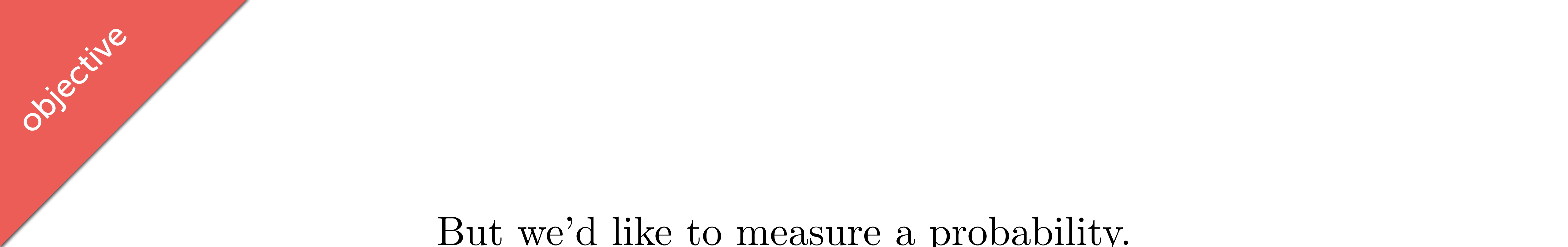


objective

How should we define  $P(v_{OUT}|v_{IN})$ ?

Measure loss between  
 $v_{IN}$  and  $v_{OUT}$ ?

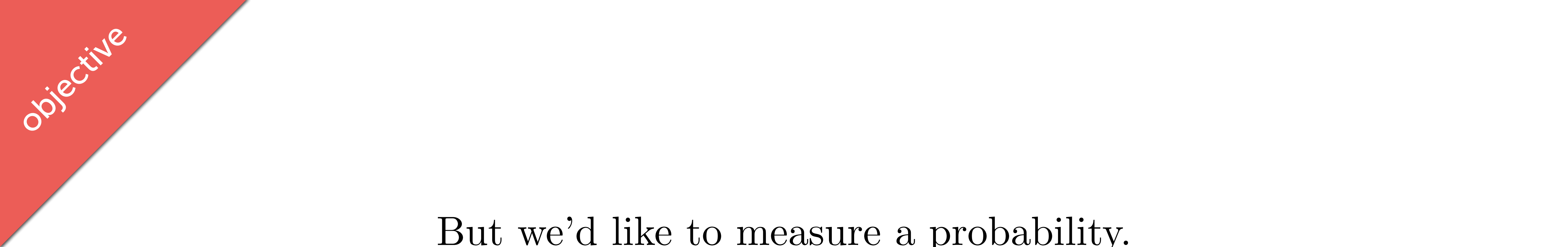
$$v_{in} \bullet v_{out}$$



objective

But we'd like to measure a probability.

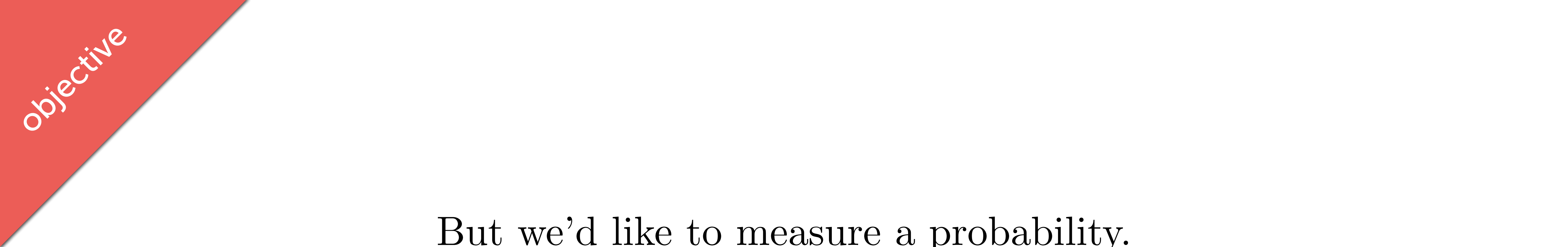
$$v_{in} \bullet v_{out} \in [-1,1]$$



objective

But we'd like to measure a probability.

$$\text{softmax}(v_{in} \cdot v_{out}) \in [0,1]$$



objective

But we'd like to measure a probability.

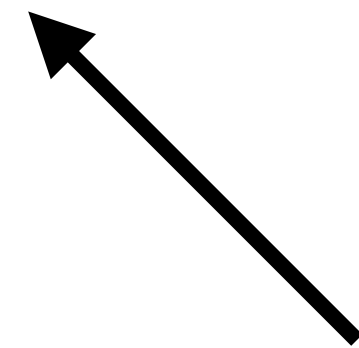
$$\textit{softmax}(v_{in} \cdot v_{out})$$

Probability of choosing 1 of N discrete items.

Mapping from vector space to a multinomial over words.

But we'd like to measure a probability.

$$\textit{softmax} = \frac{\exp(v_{in} \cdot v_{out})}{\sum_{k \in V} \exp(v_{in} \cdot v_k)}$$



Normalization term over all words

But we'd like to measure a probability.

$$\textit{softmax} = \frac{\exp(v_{in} \cdot v_{out})}{\sum_{k \in V} \exp(v_{in} \cdot v_k)} = P(v_{out} | v_{in})$$

Learn by gradient descent on the softmax prob.

For every example we see update  $v_{in}$

$$v_{in} := v_{in} + \frac{\partial}{\partial v_{in}} P(v_{out} | v_{in})$$

$$v_{out} := v_{out} + \frac{\partial}{\partial v_{out}} P(v_{out} | v_{in})$$



What's our performance?

$$\frac{\exp(v_{in} \cdot v_{out})}{\sum_{k \in V} \exp(v_{in} \cdot v_k)}$$

What's our performance?

$$\frac{\exp(v_{in} \cdot v_{out})}{\sum_{k \in V} \exp(v_{in} \cdot v_k)} \leftarrow O(V)$$

What's our performance?

$$\frac{\exp(v_{in} \cdot v_{out})}{\sum_{k \in V} \exp(v_{in} \cdot v_k)} \quad \leftarrow O(V)$$

$V$  operations for every update.

$VC$  operations per input word.

$VCN$  over the whole corpus.

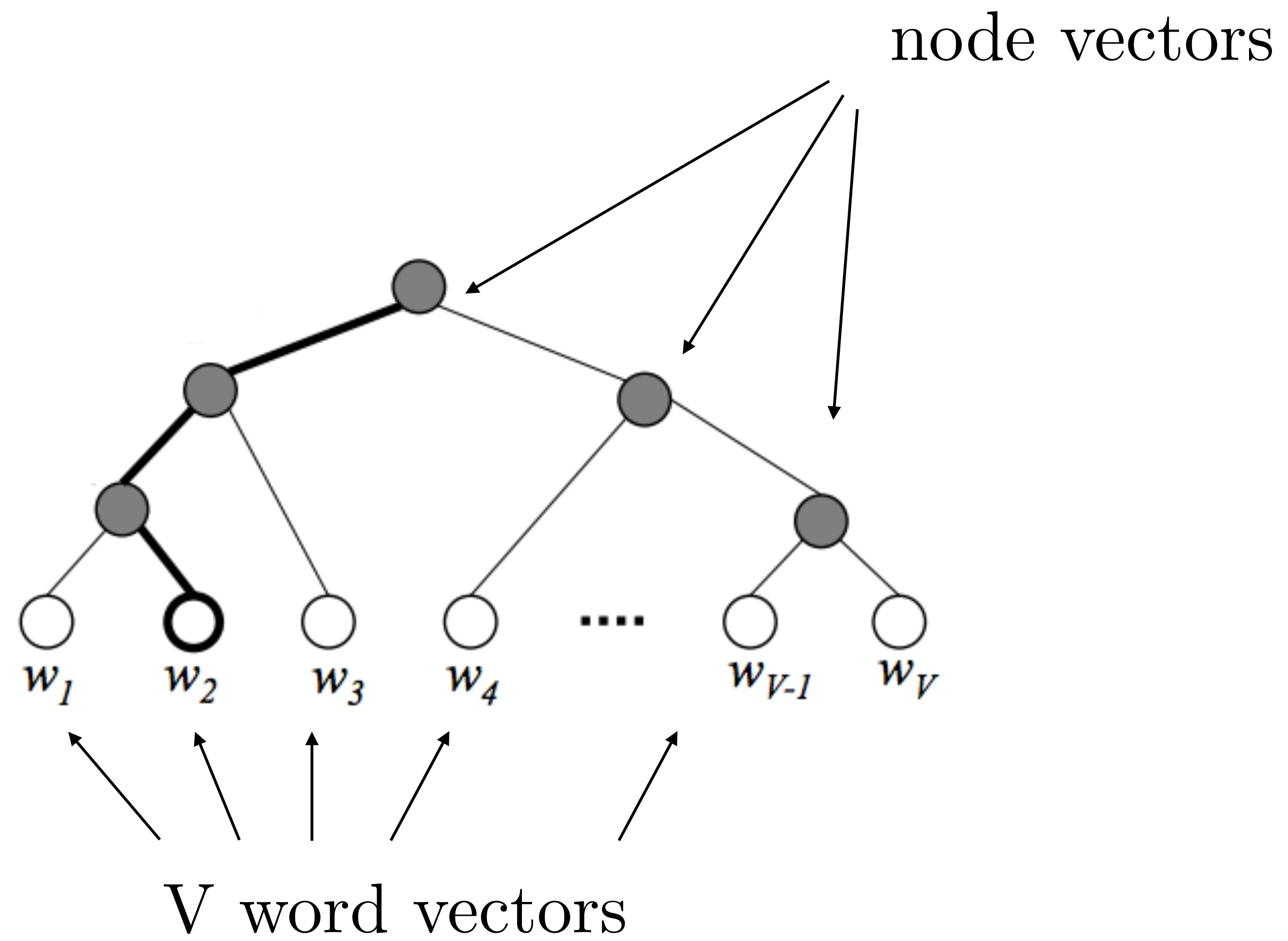
ಠ\_ಠ

How is  $O(VCN)$  supposed to be better than SVD that was  $O(NV^2)$ ?

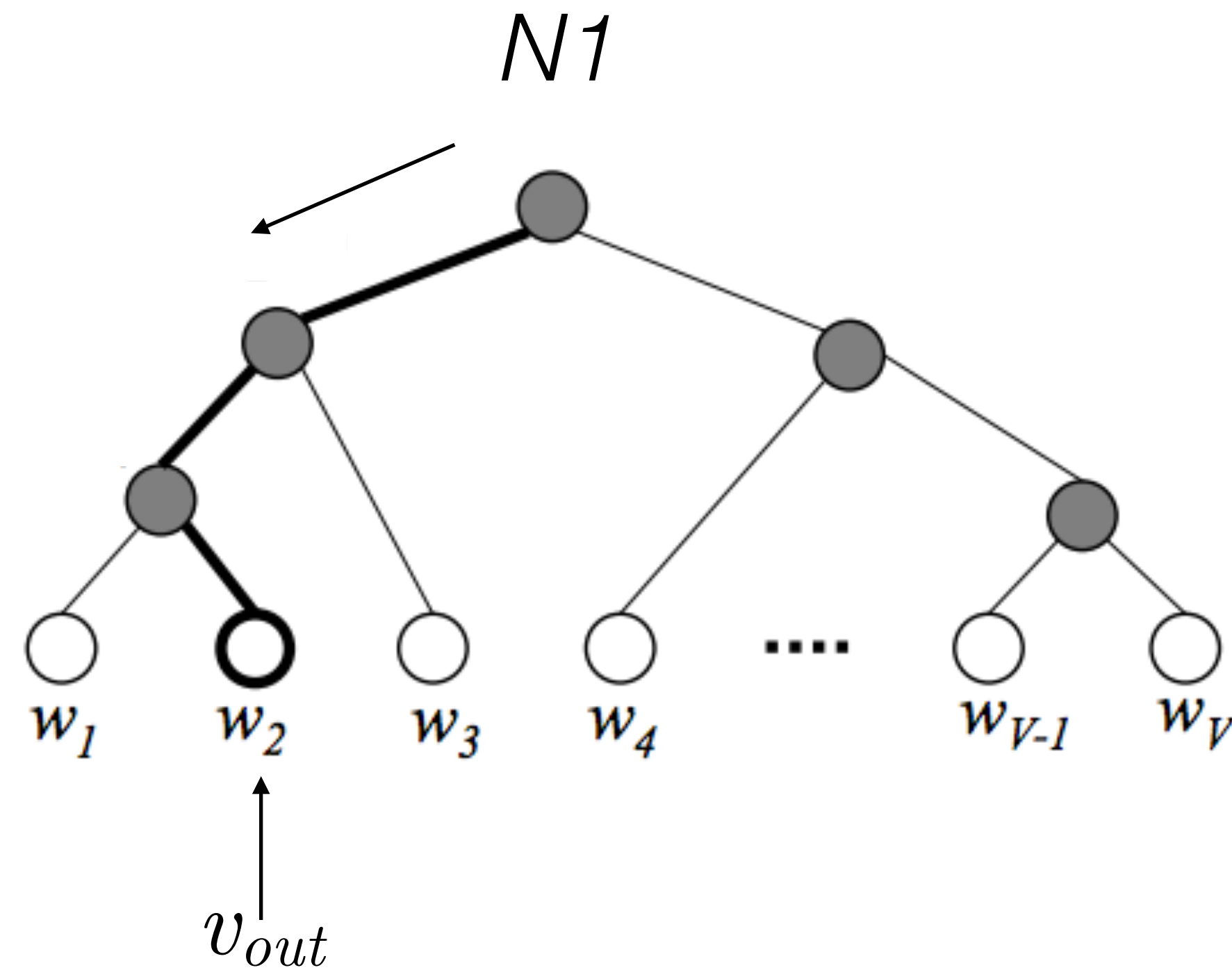
Have an  $O(V)$  problem?

Build a tree and get a  $O(\log V)$  problem!

# Hierarchical softmax

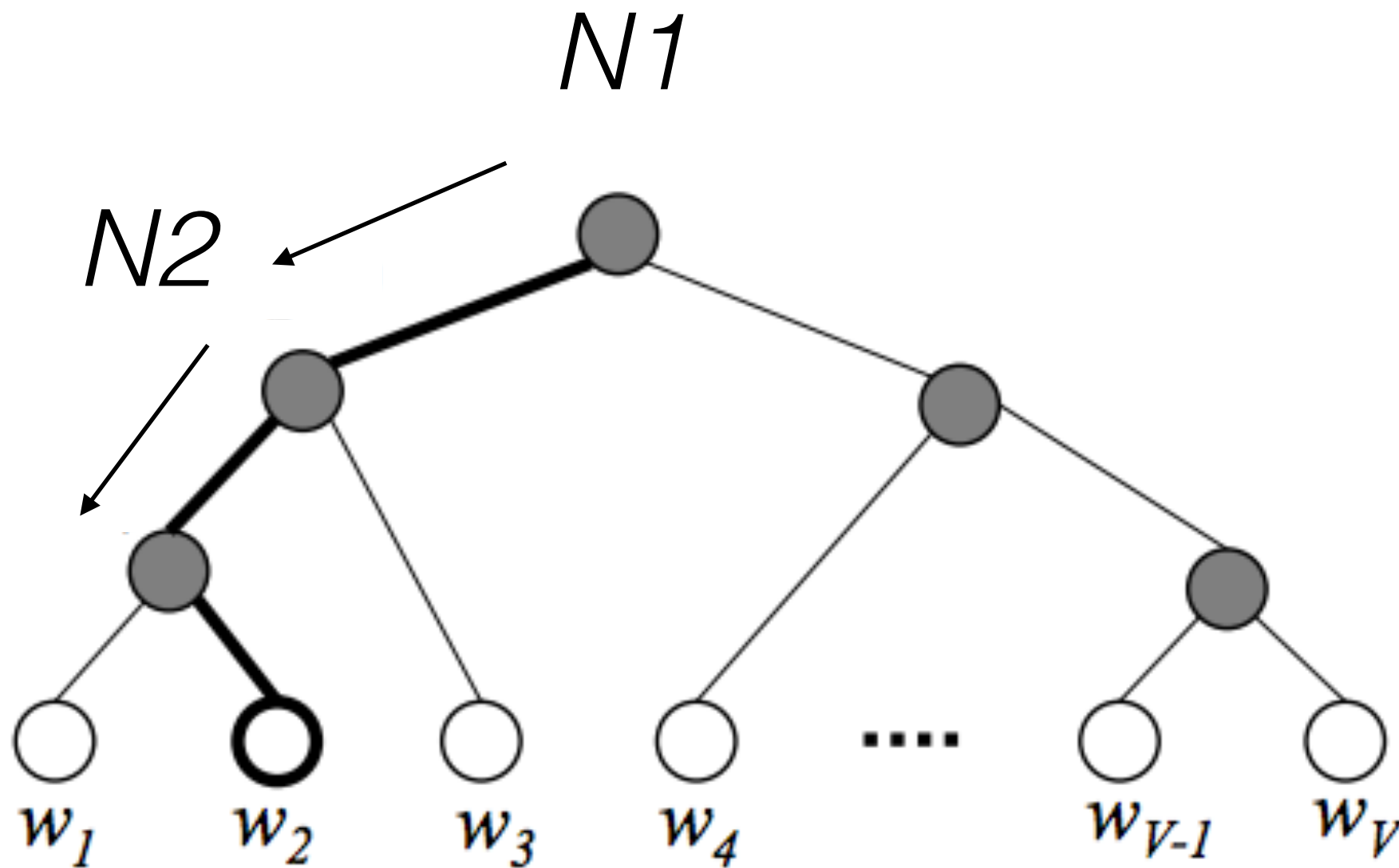


# Hierarchical softmax



$$P(v_{out}|v_{in}) = P(\text{going left at } N1 | v_{out})$$

# Hierarchical softmax

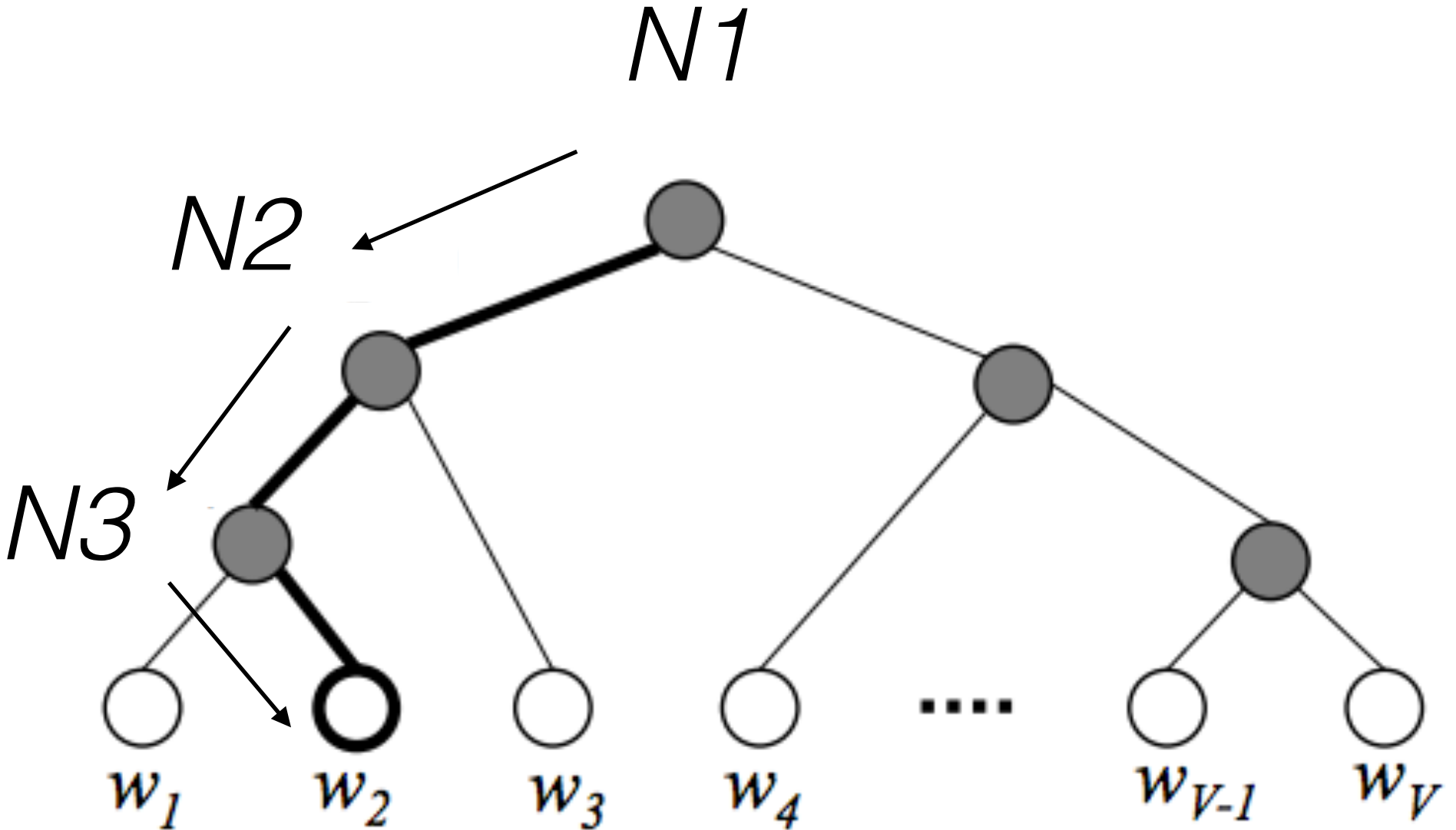


$$P(v_{out}|v_{in}) = P(\text{going left at } N1 | v_{in}) P(\text{going left at } N2 | v_{in})$$

\*Bengio 2005



# Hierarchical softmax

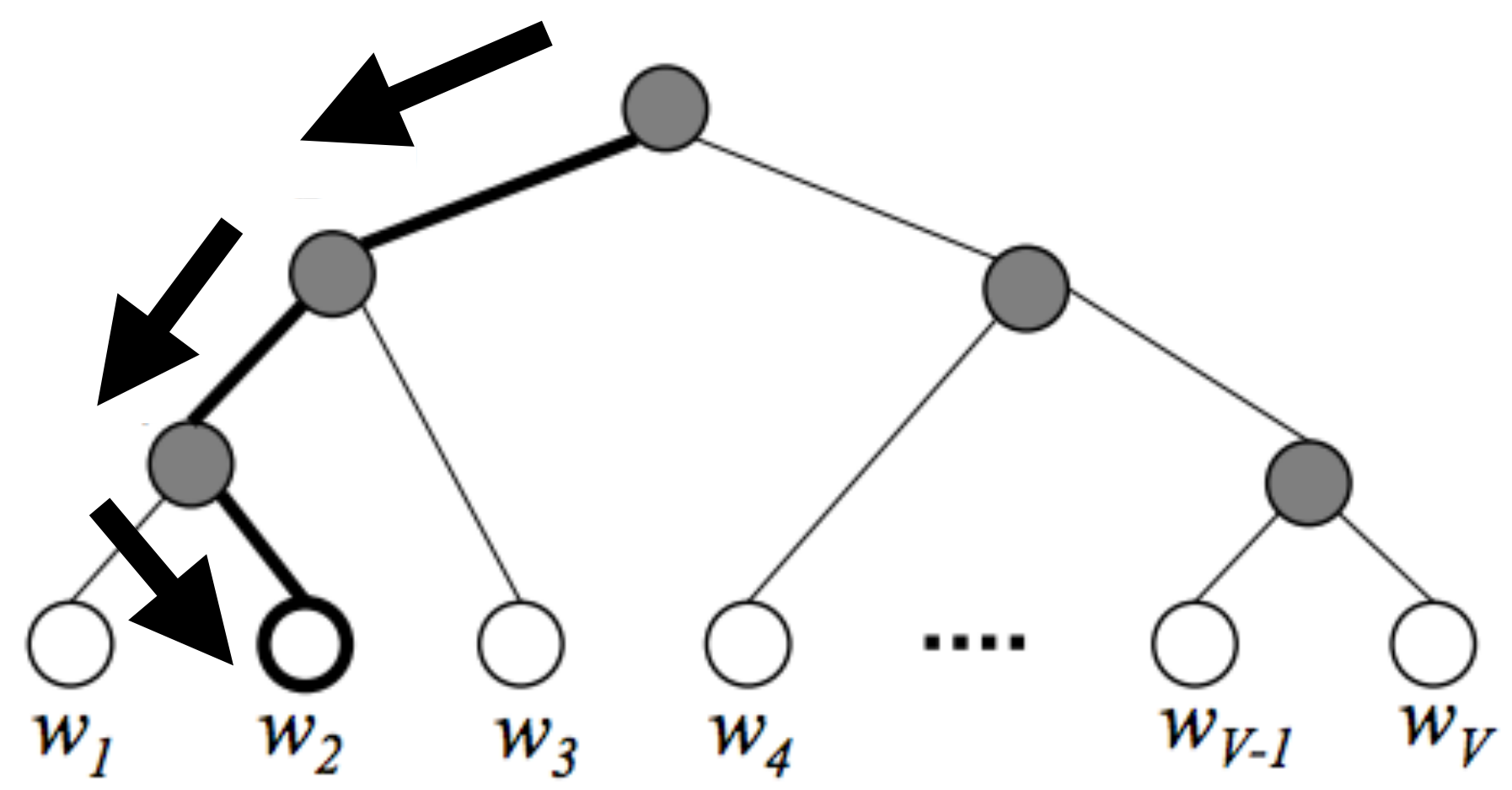


$$P(v_{out}|v_{in}) = P(\text{left at } N1 | v_{in}) P(\text{left at } N2 | v_{in}) P(\text{right at } N3 | v_{in})$$

\*Bengio 2005

# Hierarchical softmax

$O(\log N)$  steps



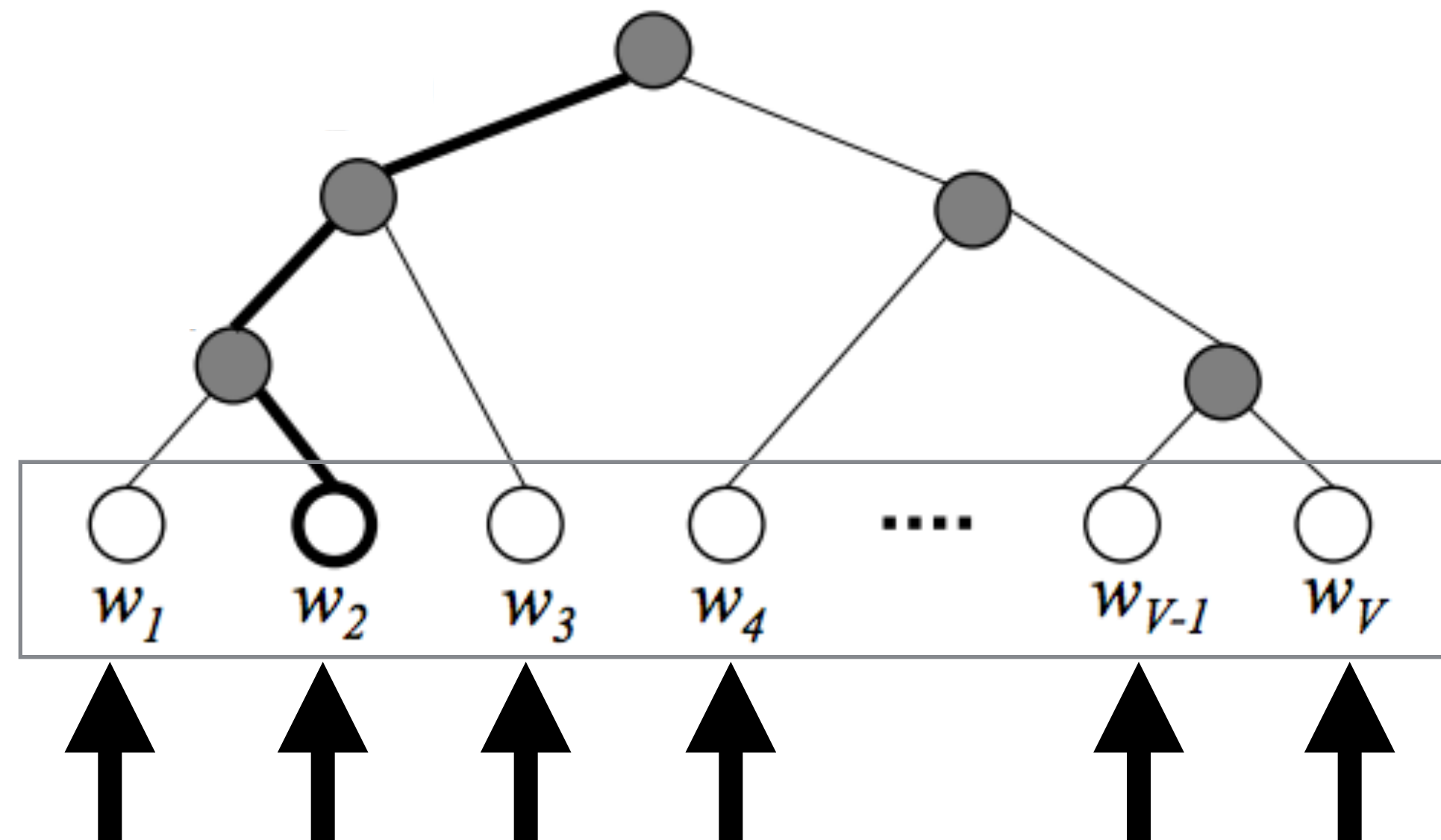
~10 comparisons

# Hierarchical softmax

$\sim O(\log N)$  steps

vs

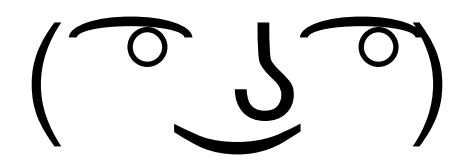
$O(N)$  steps



$\sim 10$  comparisons

$\sim 50k$  comparisons

$$\sum \exp(v_{in} \cdot v_{out})$$

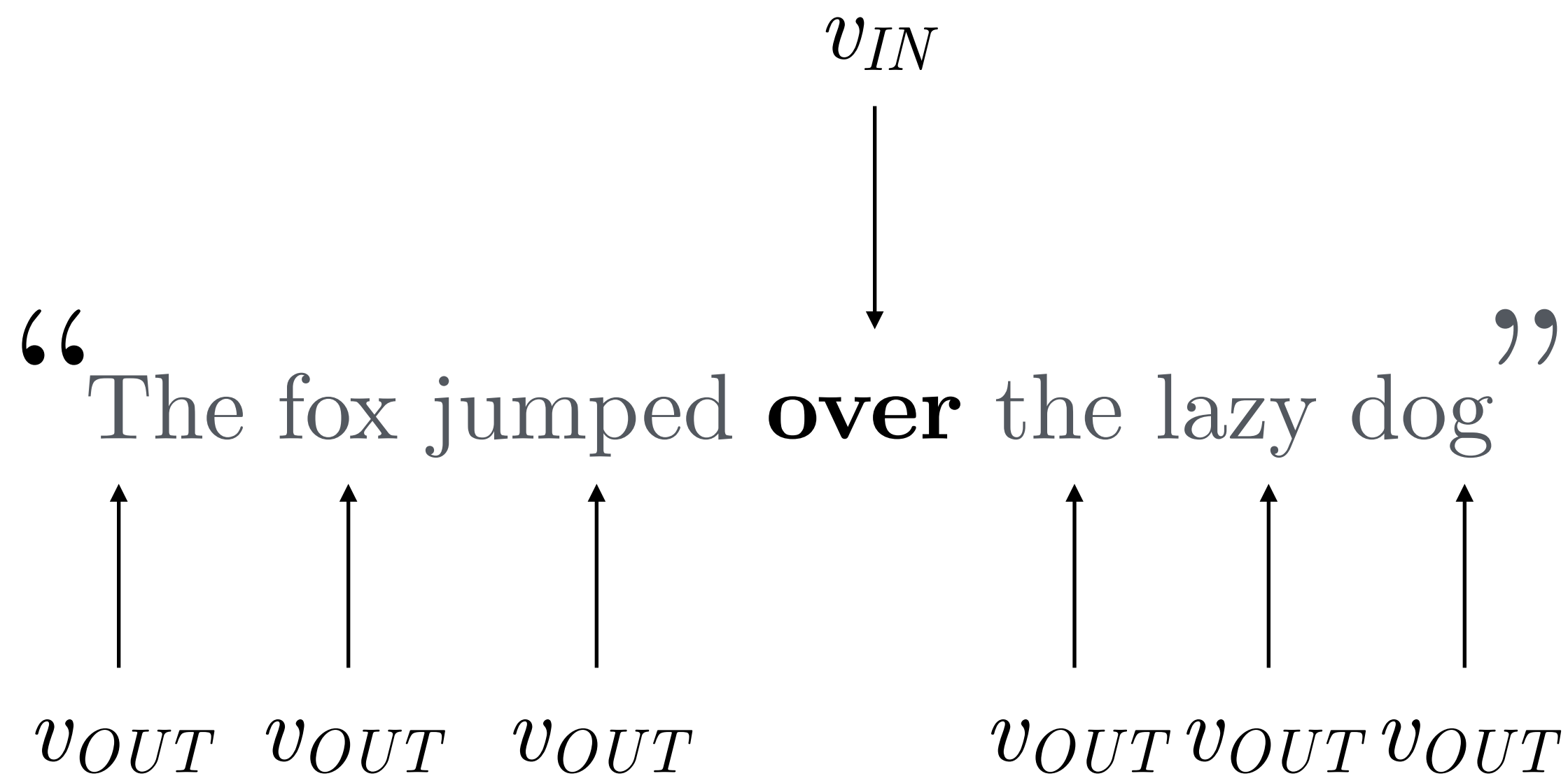


Now performance is  $O(N C \log V)$ !

Now we can scale to a 100 billion word corpus.

# SkipGram

Guess the context  
given the word

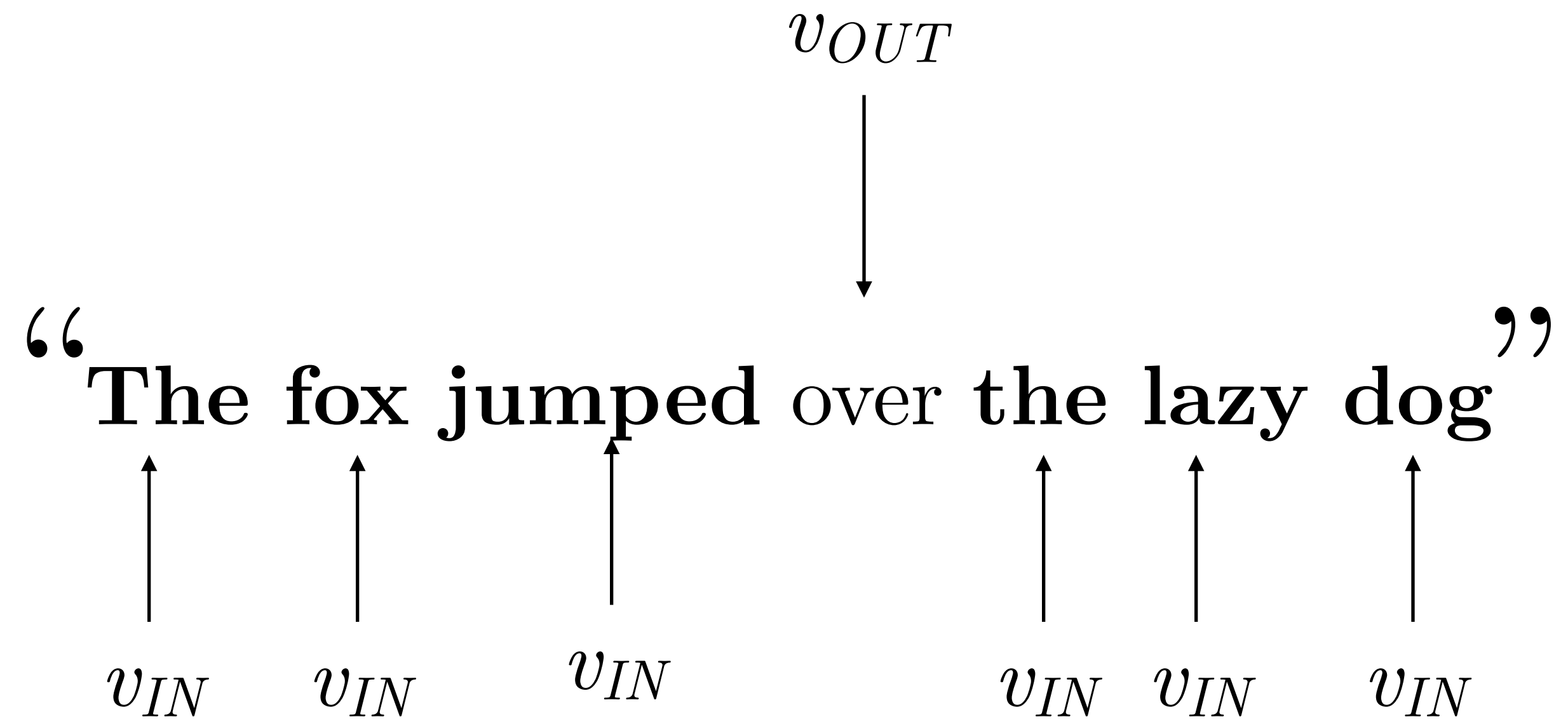


Better at syntax.

(this is the one we went over)

# CBOW

Guess the word  
given the context



~20x faster.

(this is the alternative.)

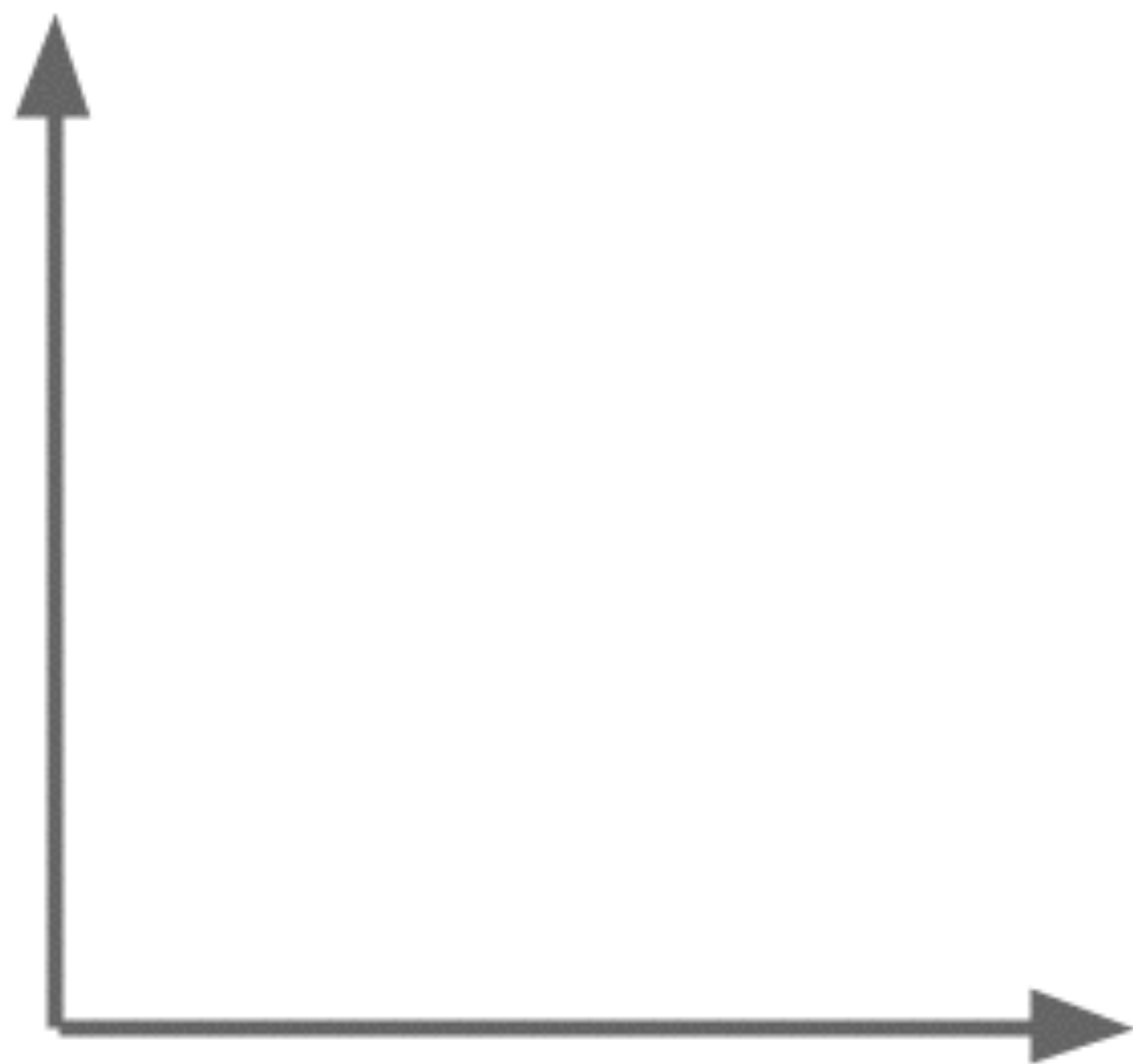
Model (training time)	Redmond	Havel	ninjutsu
Collobert (50d) (2 months)	conyers lubbock keene	plauen dzerzhinsky osterreich	reiki kohona karate
Turian (200d) (few weeks)	McCarthy Alston Cousins	Jewell Arzu Ovitz	- - -
Mnih (100d) (7 days)	Podhurst Harlang Agarwal	Pontiff Pinochet Rodionov	- - -
Skip-Phrase (1000d, 1 day)	Redmond Wash. Redmond Washington Microsoft	Vaclav Havel president Vaclav Havel Velvet Revolution	ninja martial arts swordsmanship

← word2vec

Model	Vector Dimensionality	Training words	Accuracy [%]		
			Semantic	Syntactic	Total
Collobert-Weston NNLM	50	660M	9.3	12.3	11.0
Turian NNLM	50	37M	1.4	2.6	2.1
Turian NNLM	200	37M	1.4	2.2	1.8
Mnih NNLM	50	37M	1.8	9.1	5.8
Mnih NNLM	100	37M	3.3	13.2	8.8
Mikolov RNNLM	80	320M	4.9	18.4	12.7
Mikolov RNNLM	640	320M	8.6	36.5	24.6
Huang NNLM	50	990M	13.3	11.6	12.3
CBOW	300	783M	15.5	53.1	36.1
Skip-gram	300	783M	<b>50.0</b>	55.9	<b>53.3</b>

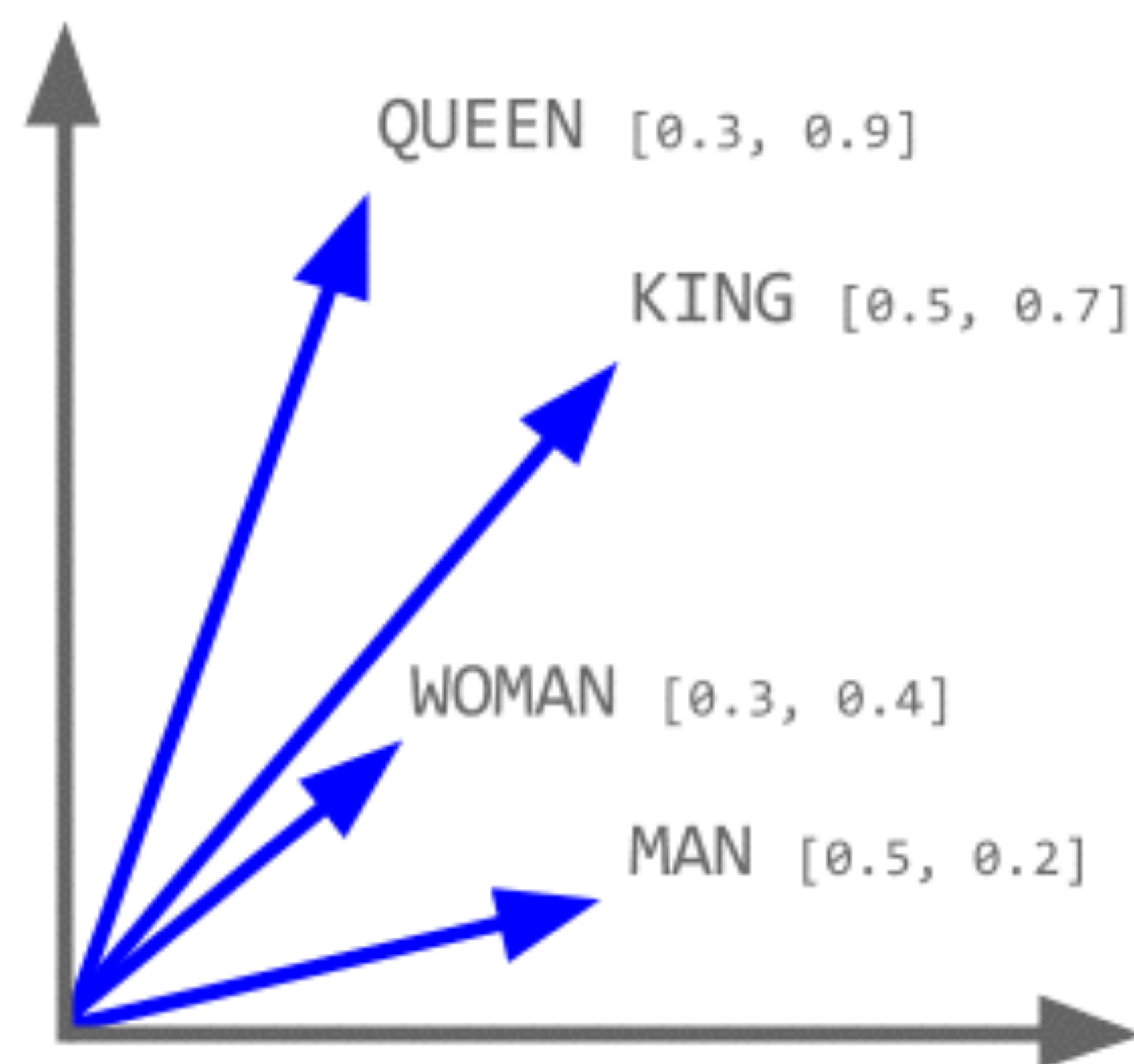
← word2vec

What is king + man - woman?

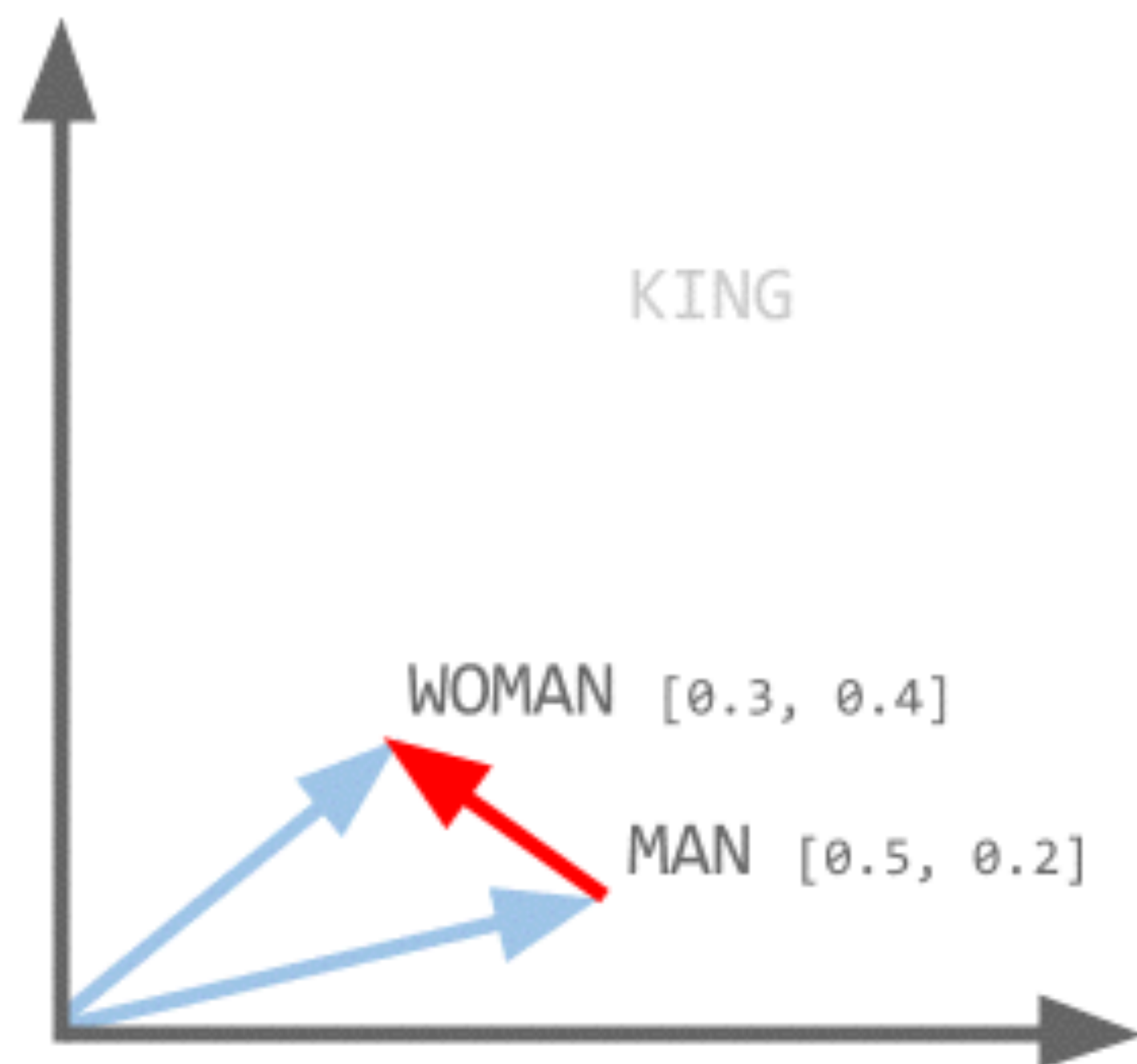




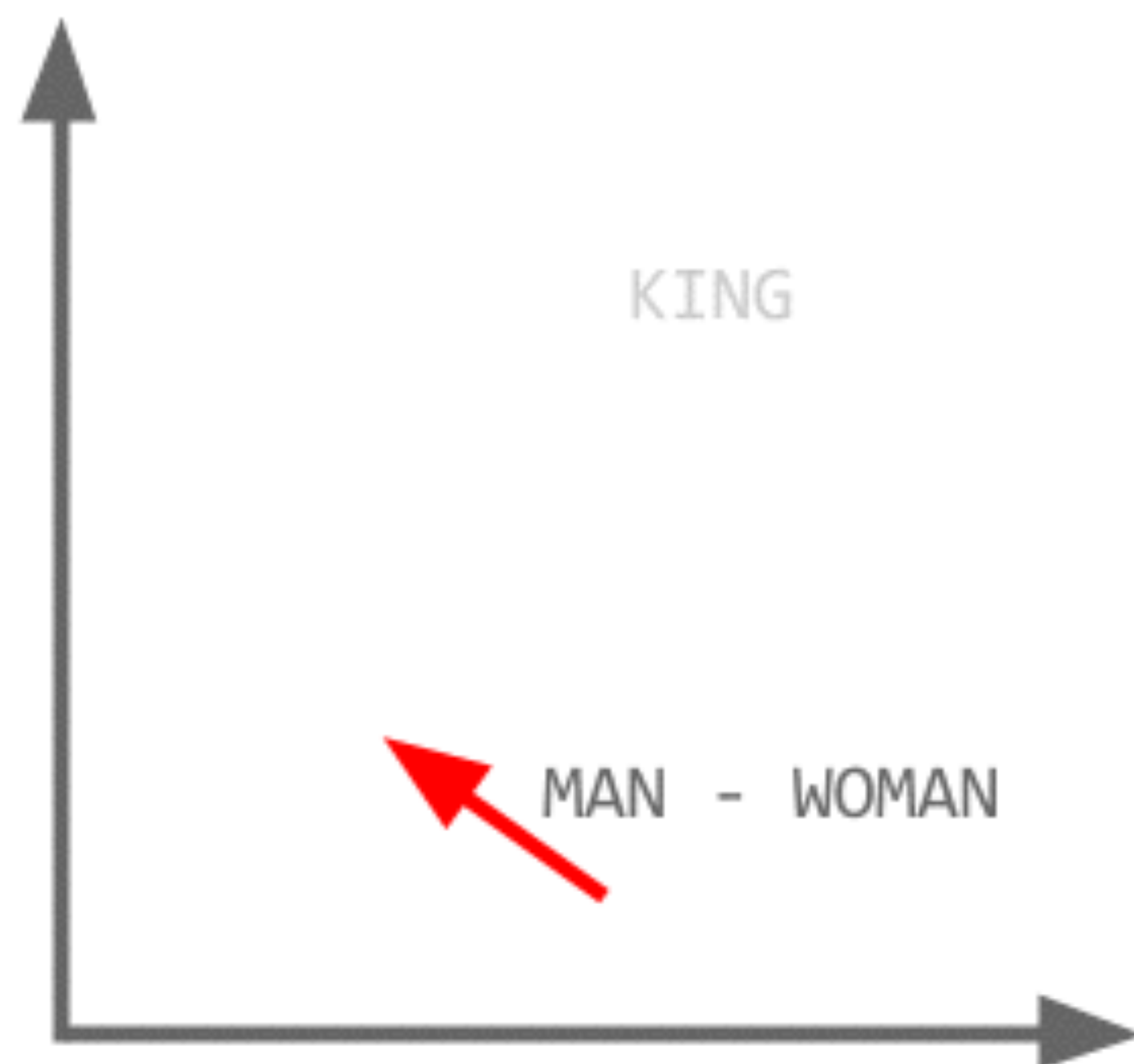
Load up the word vectors



Start with man - woman



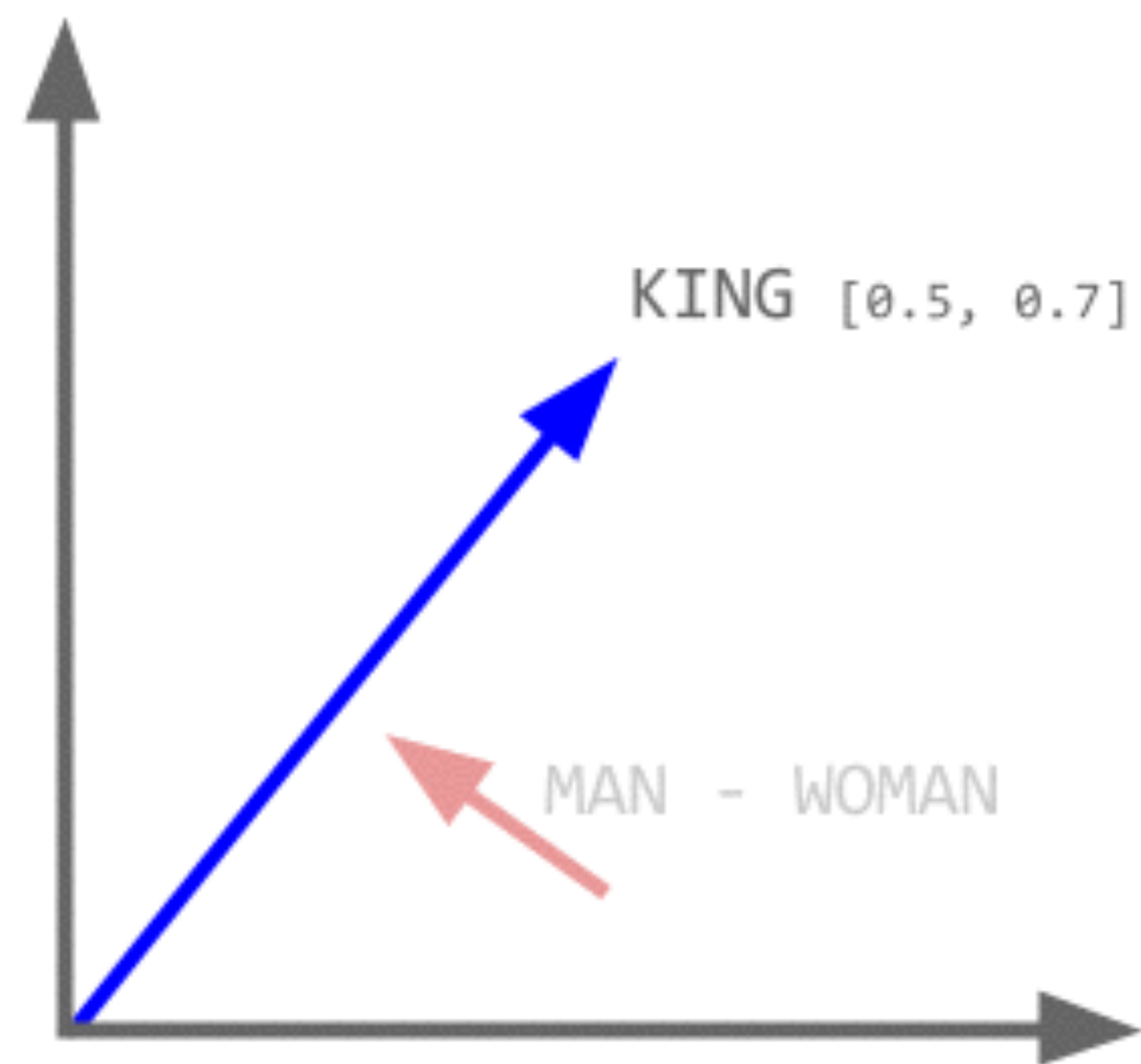
Start with man - woman



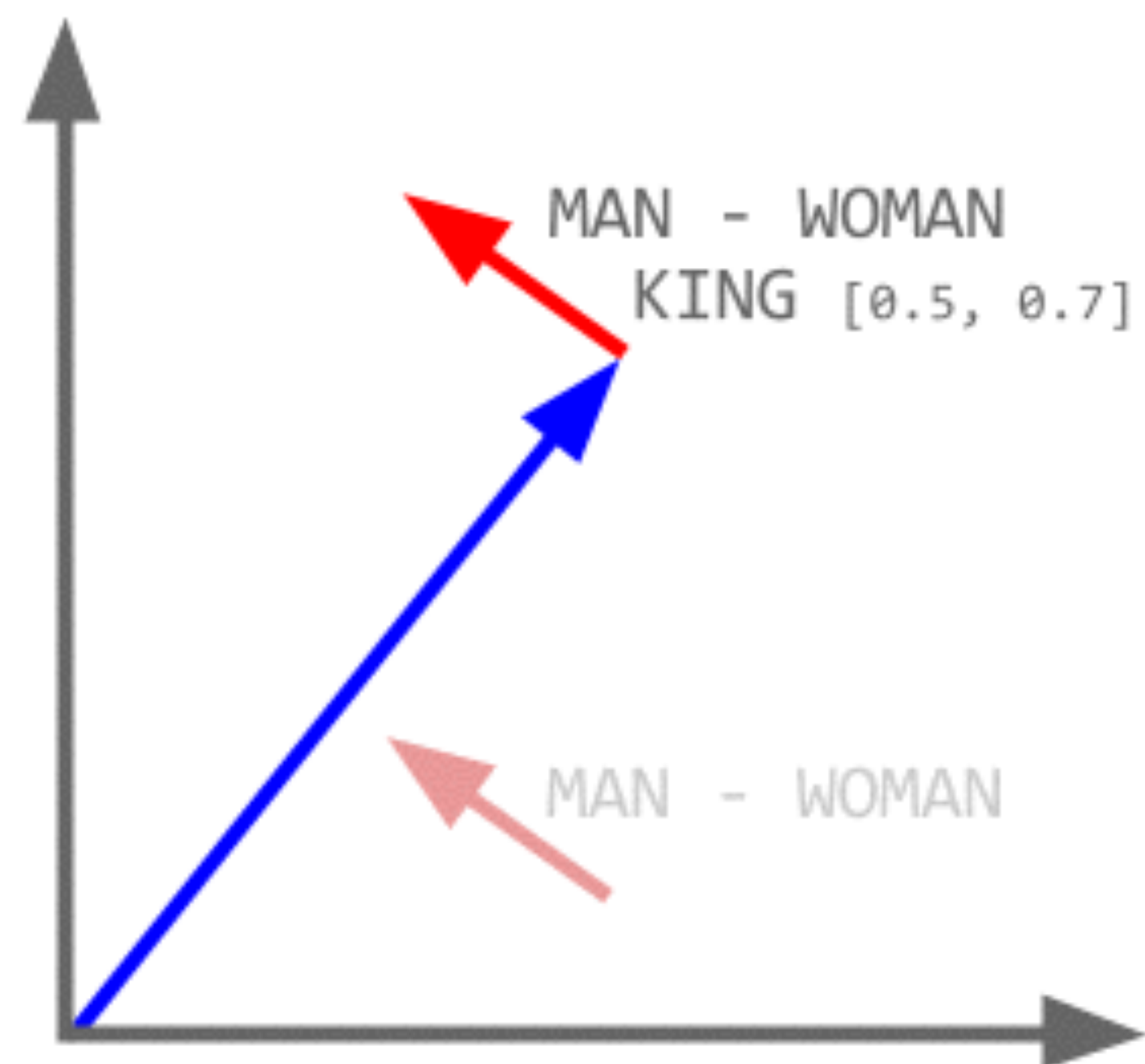
KING

MAN - WOMAN

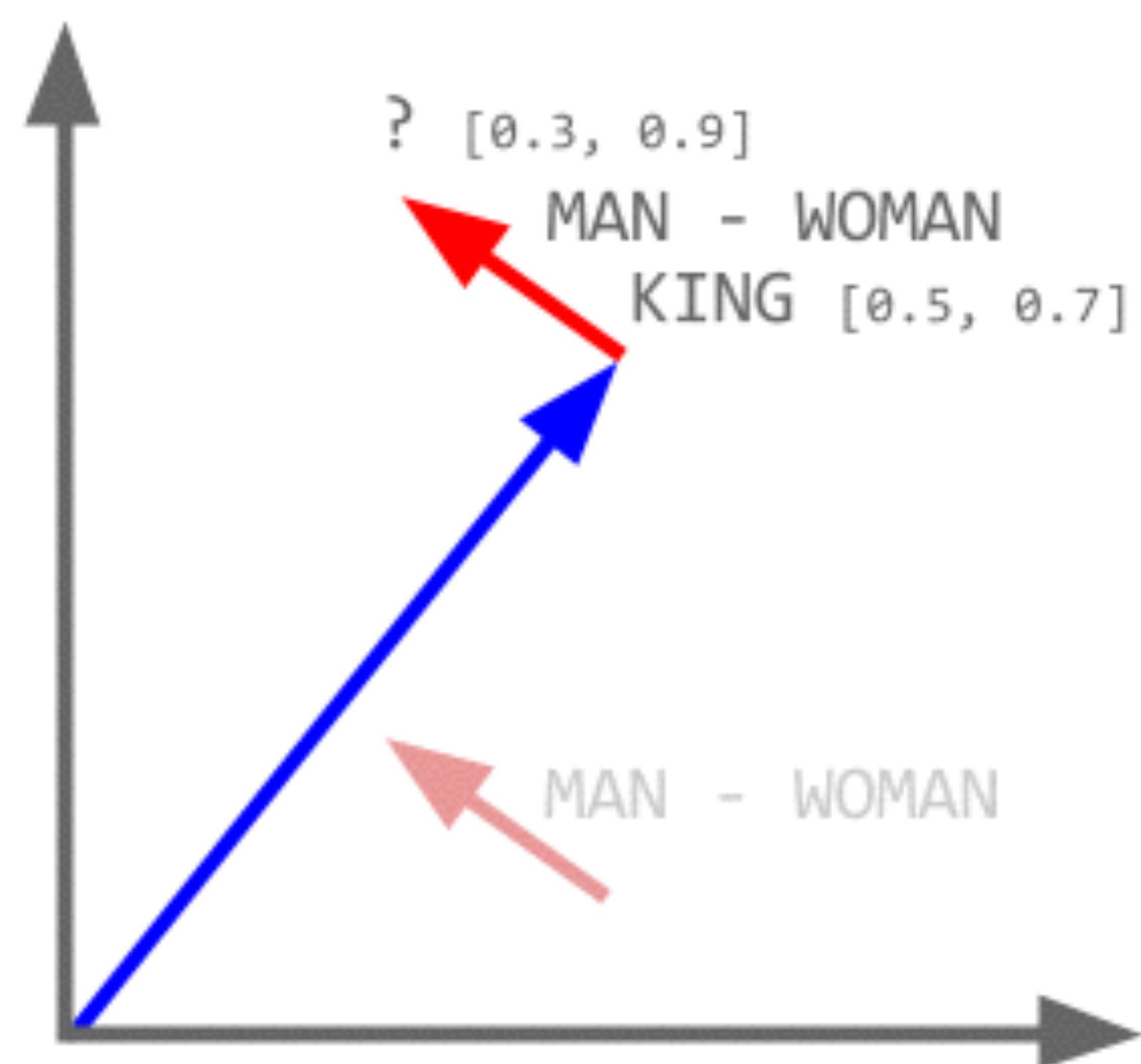
Then take king



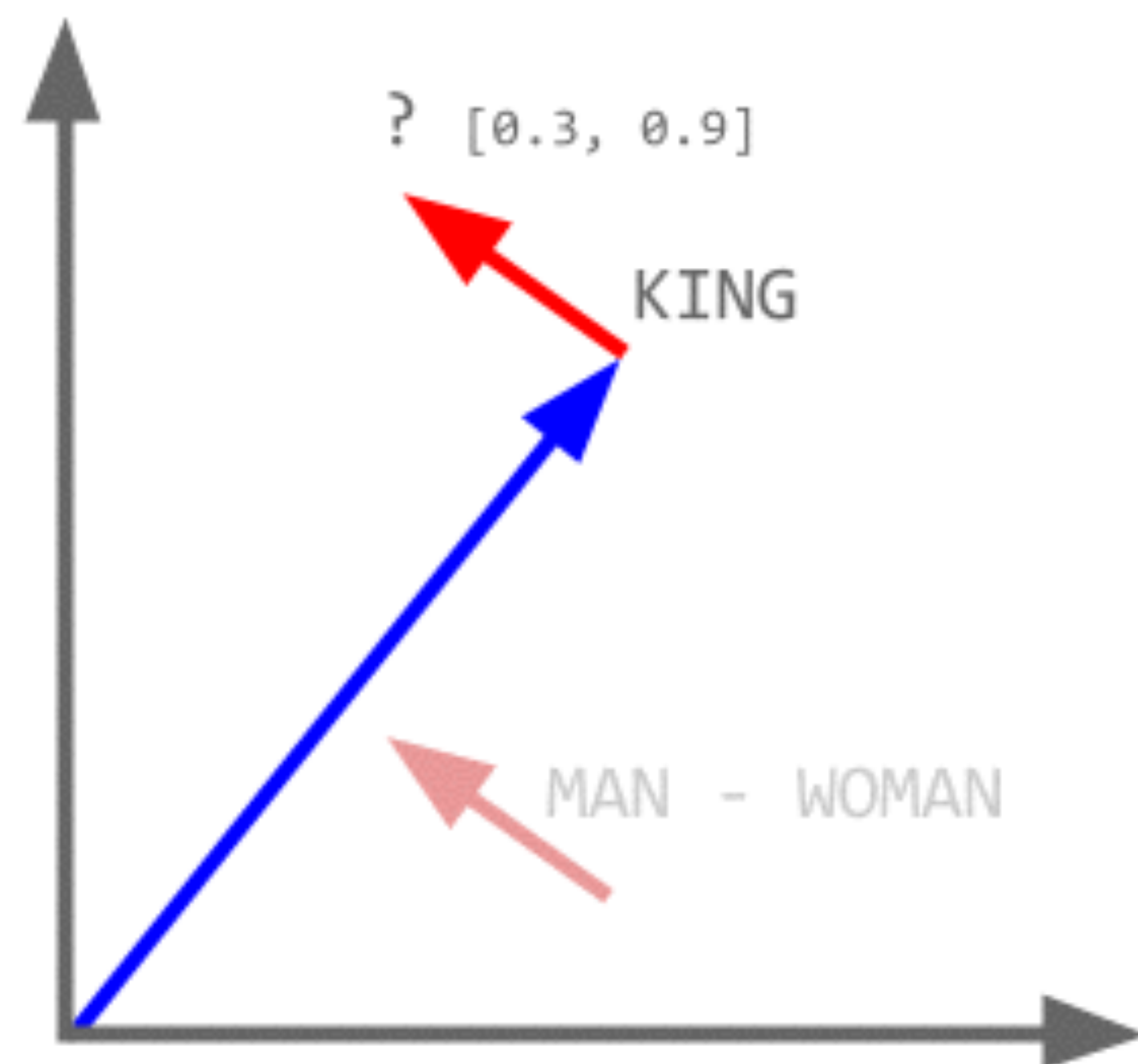
And add man - woman



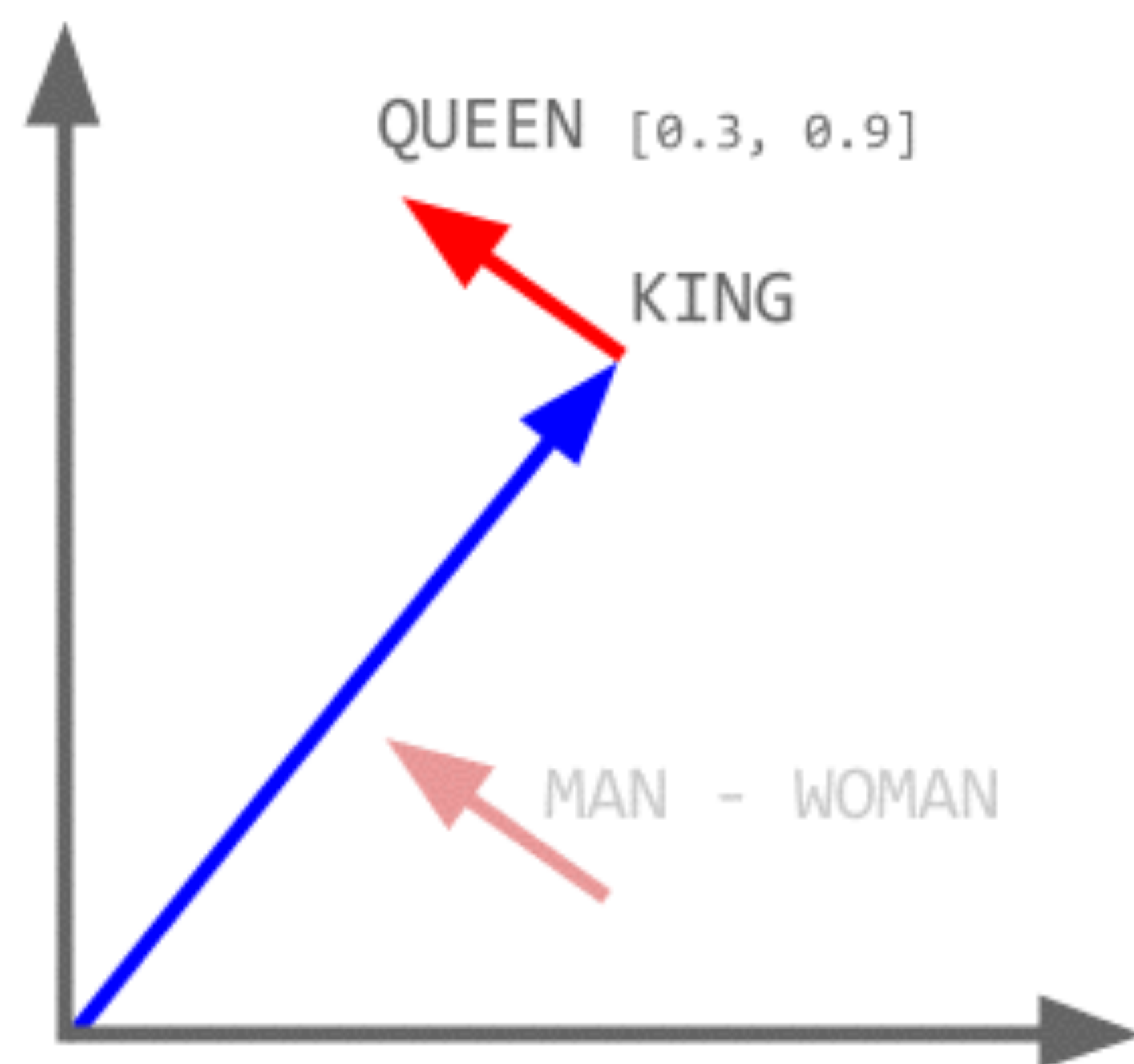
And add man - woman



Find nearest word to result

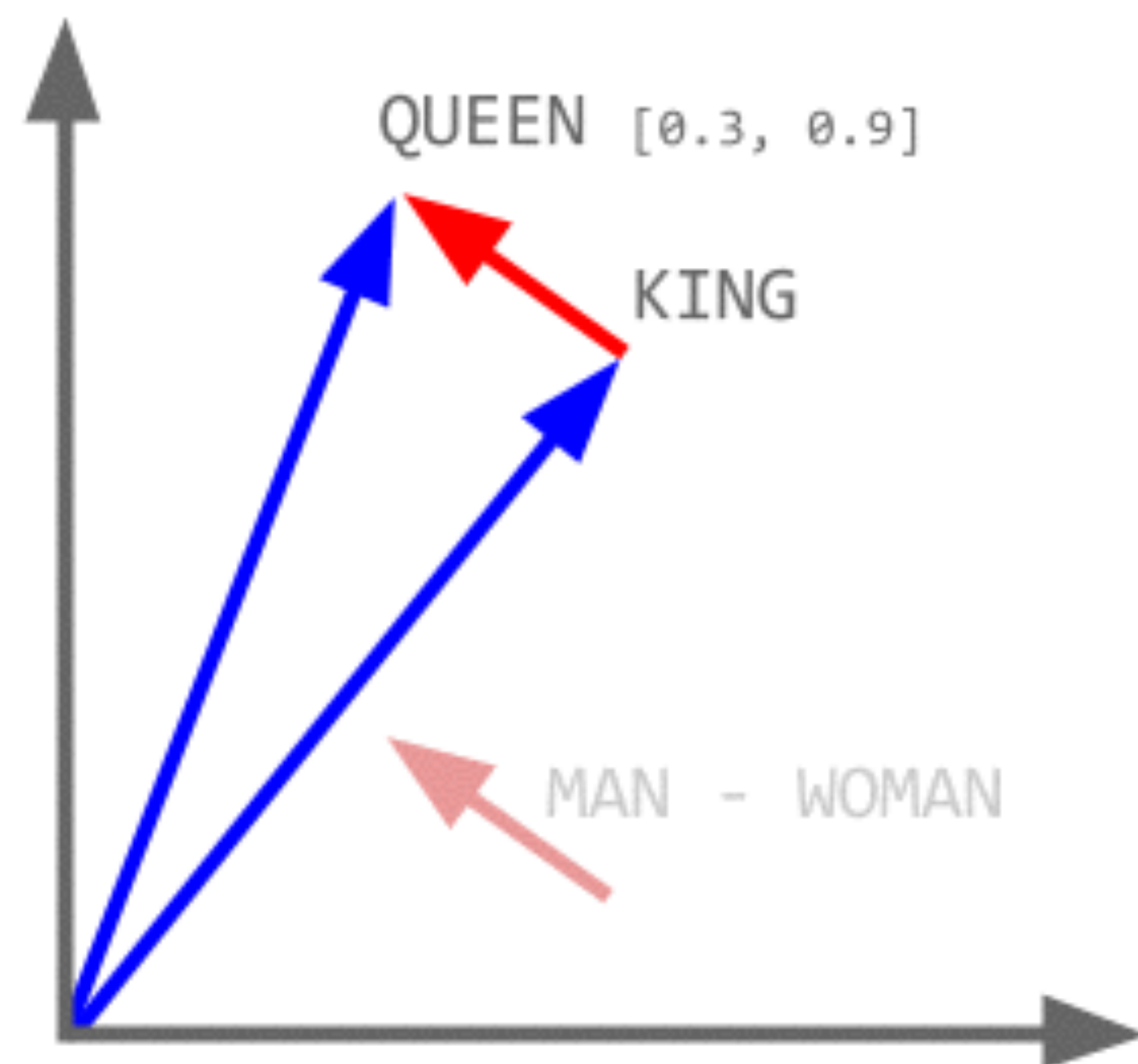


queen is closest to resulting vector

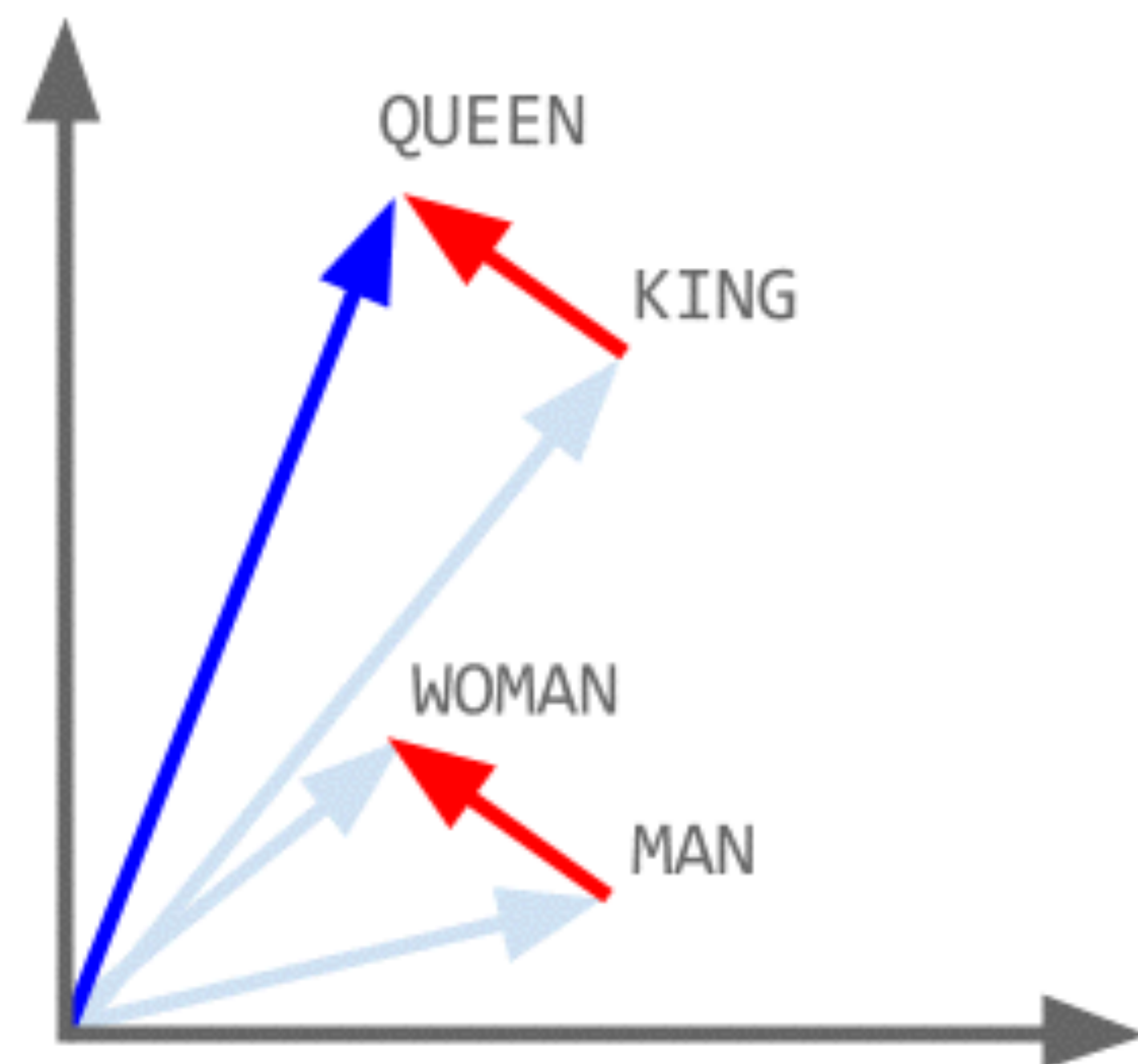




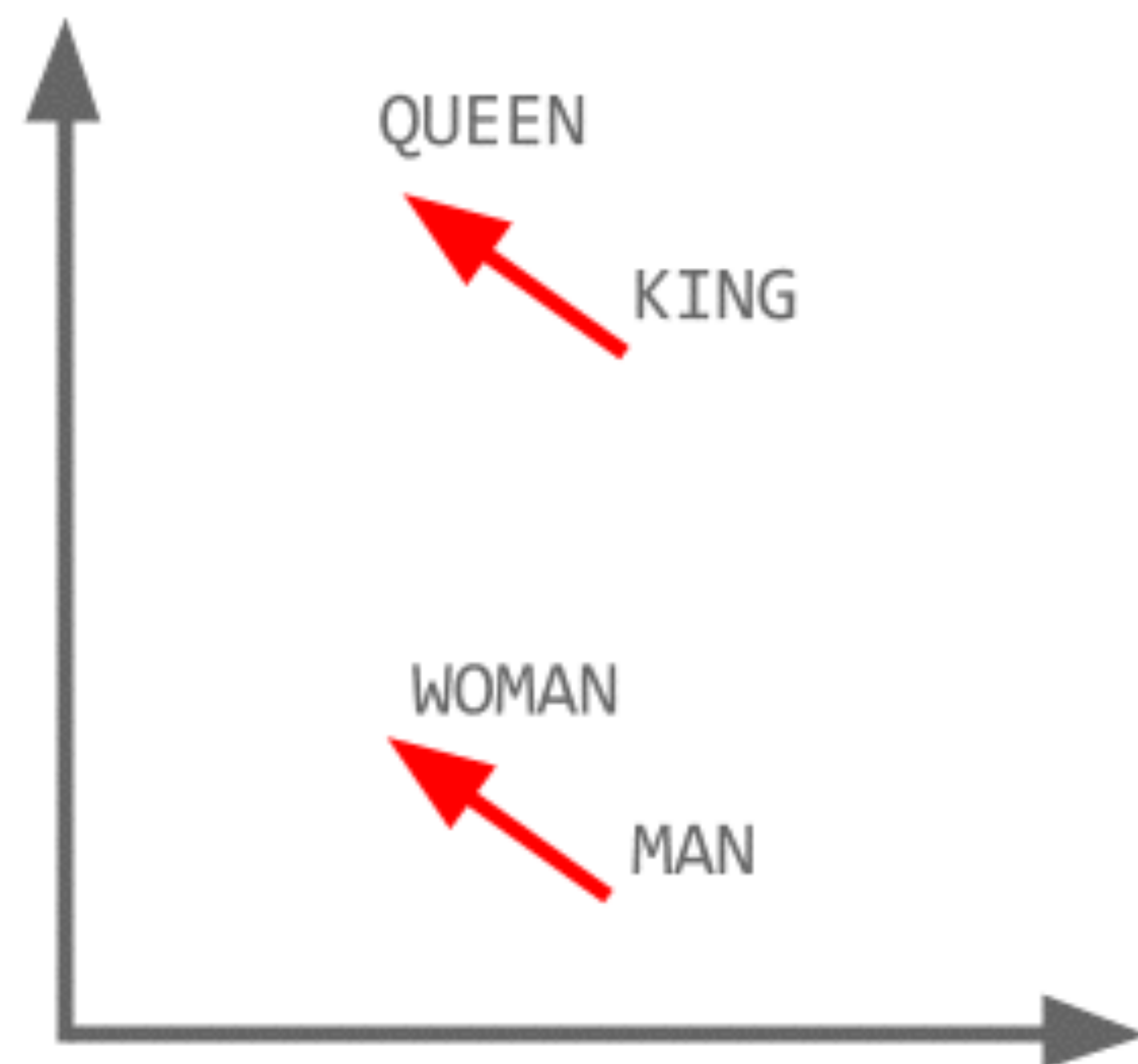
queen is closest to resulting vector



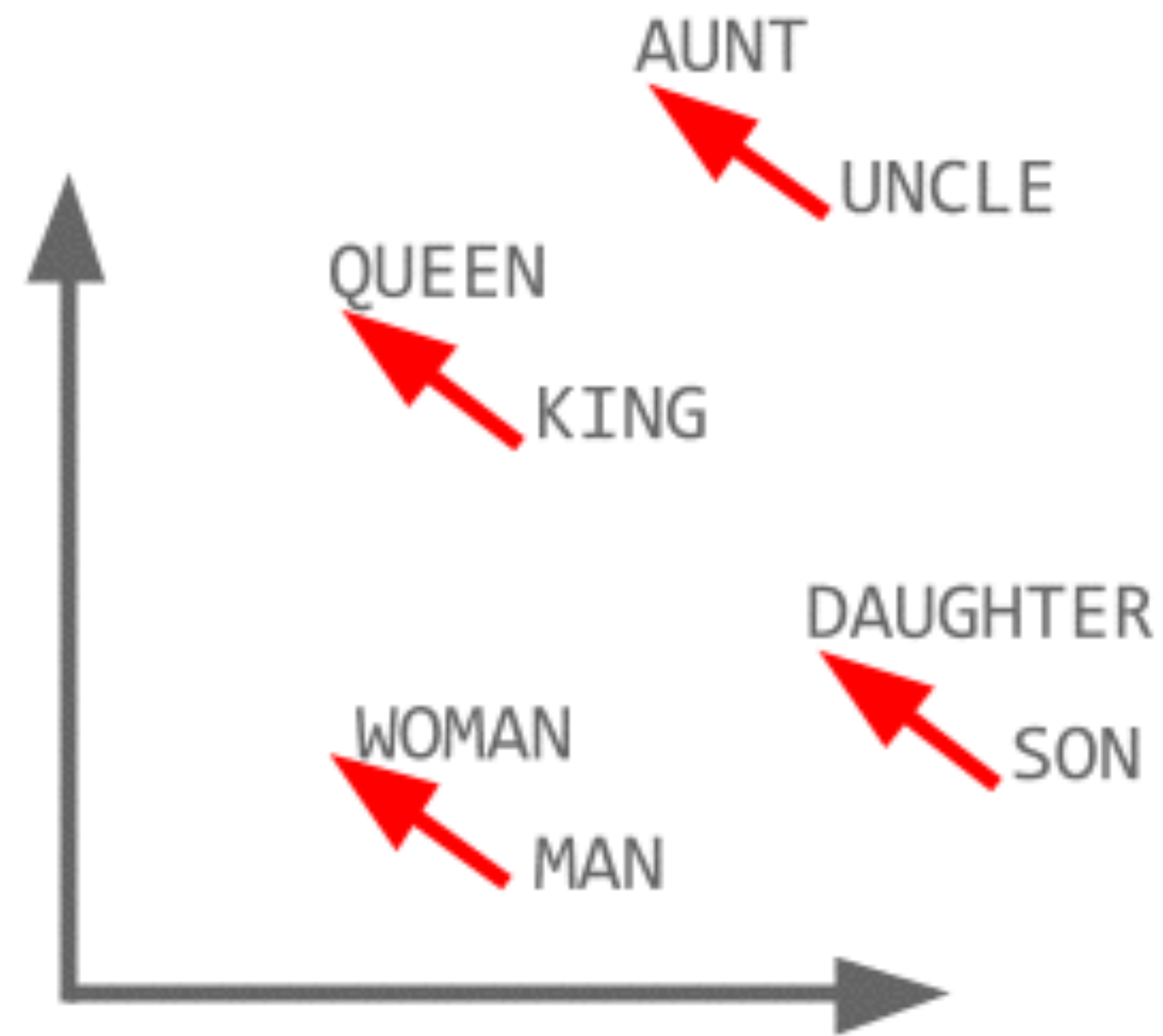
So  $\text{king} + \text{man} - \text{woman} = \text{queen!}$



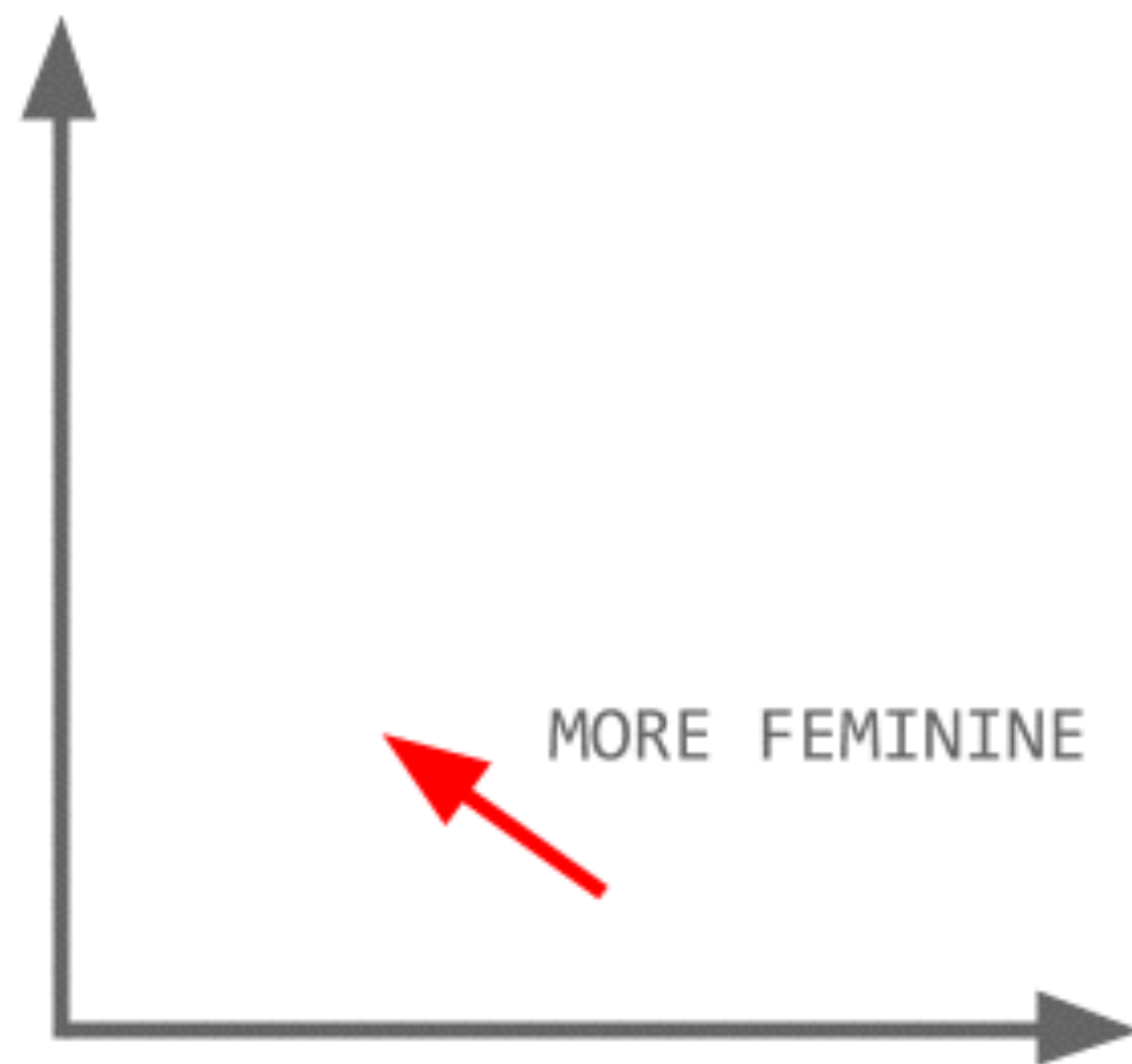
The **red direction** encodes gender



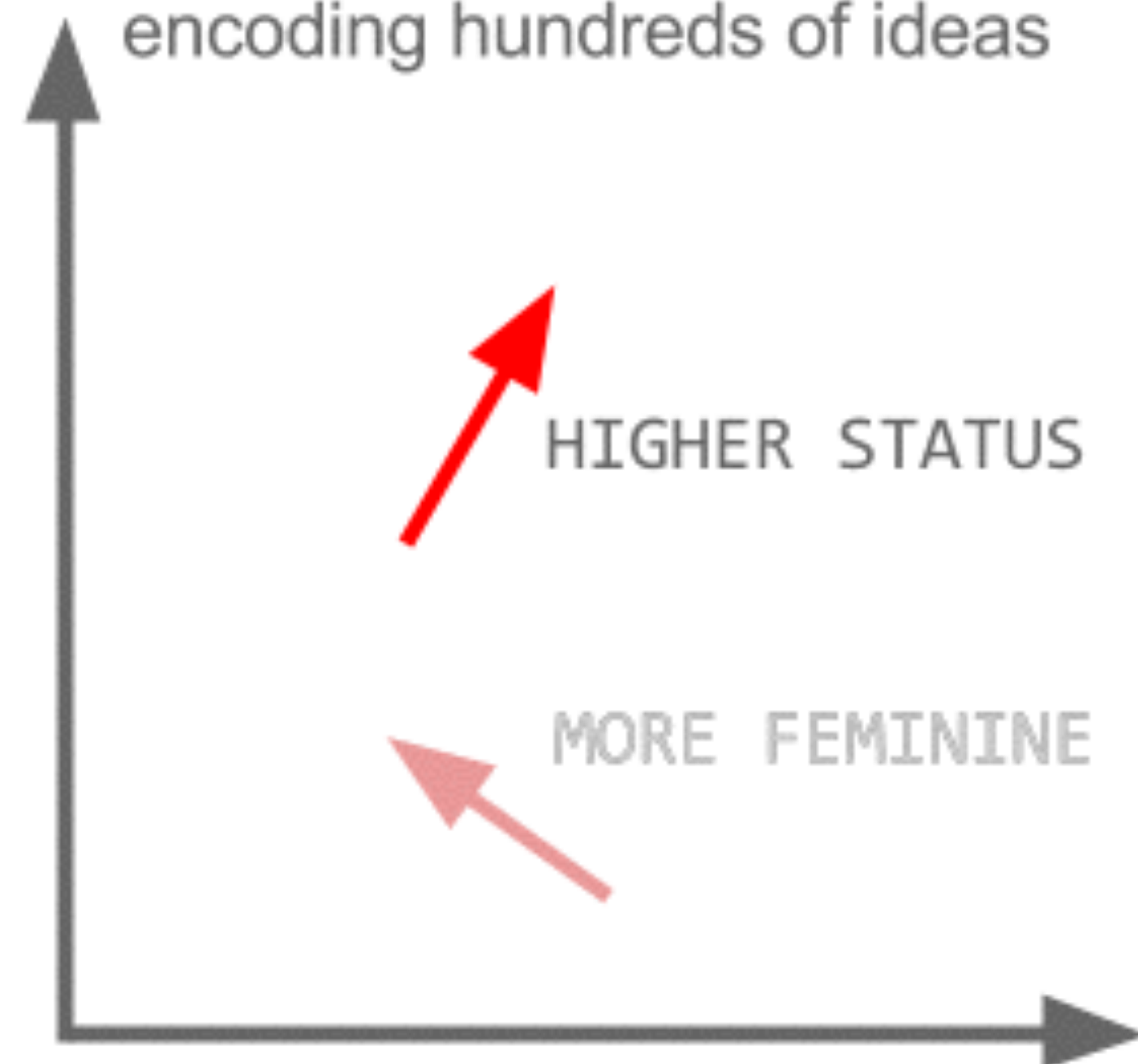
Which is consistent across all words



This **direction** always means **gender**



We have hundreds of **directions**  
encoding hundreds of ideas



Relationship	Example 1	Example 2	Example 3
France - Paris	Italy: Rome	Japan: Tokyo	Florida: Tallahassee
big - bigger	small: larger	cold: colder	quick: quicker
Miami - Florida	Baltimore: Maryland	Dallas: Texas	Kona: Hawaii
Einstein - scientist	Messi: midfielder	Mozart: violinist	Picasso: painter
Sarkozy - France	Berlusconi: Italy	Merkel: Germany	Koizumi: Japan
copper - Cu	zinc: Zn	gold: Au	uranium: plutonium
Berlusconi - Silvio	Sarkozy: Nicolas	Putin: Medvedev	Obama: Barack
Microsoft - Windows	Google: Android	IBM: Linux	Apple: iPhone
Microsoft - Ballmer	Google: Yahoo	IBM: McNealy	Apple: Jobs
Japan - sushi	Germany: bratwurst	France: tapas	USA: pizza

Czech + currency	Vietnam + capital	German + airlines	Russian + river	French + actress
koruna	Hanoi	airline Lufthansa	Moscow	Juliette Binoche
Check crown	Ho Chi Minh City	carrier Lufthansa	Volga River	Vanessa Paradis
Polish zolty	Viet Nam	flag carrier Lufthansa	upriver	Charlotte Gainsbourg
CTK	Vietnamese	Lufthansa	Russia	Cecile De







+ 'Pregnant'





# word2vec

Learns word vectors  
Learn doc vectors

Text feature generation is  
great for ML models

# LDA

Learns a topic distribution for every  
document

Each of the  $k$  topics is a great tag

and now for something **completely crazy**

All of the following ideas will change what  
'words' and 'context' represent.

What about summarizing documents?

On the day he took office, President Obama reached out to America's enemies, offering in his first inaugural address to **extend** a hand if you are willing to unclench your fist. More than six years later, he has arrived at a moment of truth in testing that



paragraph  
vector

IN



On the day he took office, President Obama reached out to America's enemies, offering in his first inaugural address to **extend** a hand if you are willing to unclench

your fist. More than six years later, he has arrived at a moment of truth in testing that



The framework nuclear agreement he reached with Iran on Thursday did not provide the definitive answer to whether Mr. Obama's audacious gamble will pay off. The fist

OUT

OUT

Iran has shaken at the so-called Great Satan since 1979 has not completely relaxed.

Normal skipgram extends  $C$  words before, and  $C$  words after.

paragraph  
vector

IN  
↓

**doc\_1347**

OUT

OUT



On the day he took office, President Obama reached out to America's enemies, offering in his first inaugural address to extend a hand if you are willing to unclench your fist. More than six years later, he has arrived at a moment of truth in testing that



The framework nuclear agreement he reached with Iran on Thursday did not provide the definitive answer to whether Mr. Obama's audacious gamble will pay off. The fist Iran has shaken at the so-called Great Satan since 1979 has not completely relaxed.

OUT

OUT

A document vector simply extends the context to the whole document.

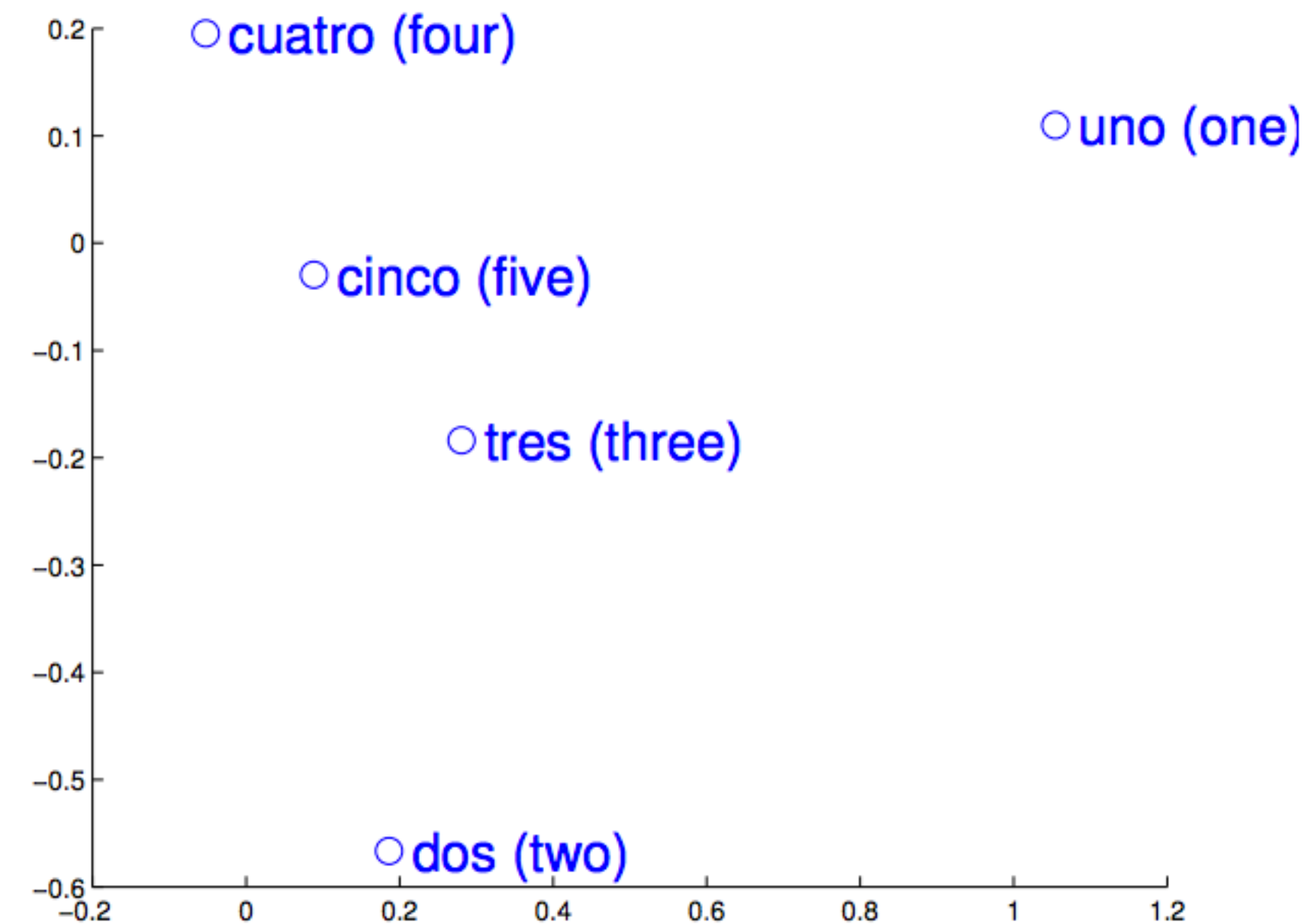
```
from gensim.models import Doc2Vec
fn = "item_document_vectors"
model = Doc2Vec.load(fn)
model.most_similar('pregnant')
matches = list(filter(lambda x: 'SENT_' in x[0], matches))

# ['...I am currently 23 weeks pregnant...',
#  '...I'm now 10 weeks pregnant...',
#  '...not showing too much yet...',
#  '...15 weeks now. Baby bump...',
#  '...6 weeks post partum!...',
#  '...12 weeks postpartum and am nursing...',
#  '...I have my baby shower that...',
#  '...am still breastfeeding...',
#  '...I would love an outfit for a baby shower...']
```

# translation

(using just a rotation matrix)

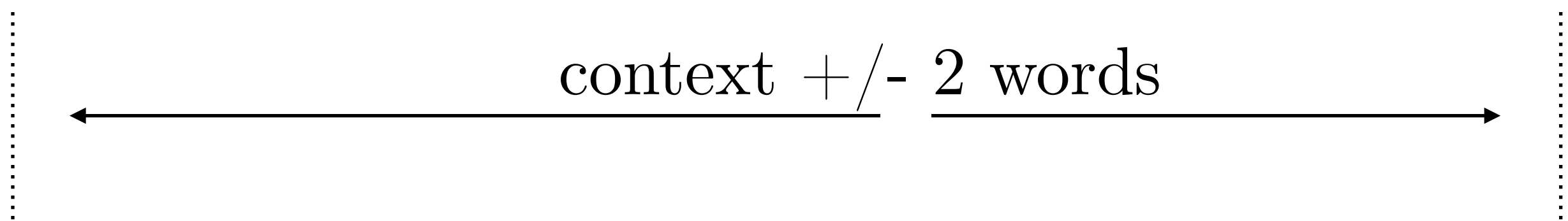
English  
↓  
Matrix  
Rotation  
↓  
Spanish



context  
dependent

Australian scientist **discovers** star with telescope

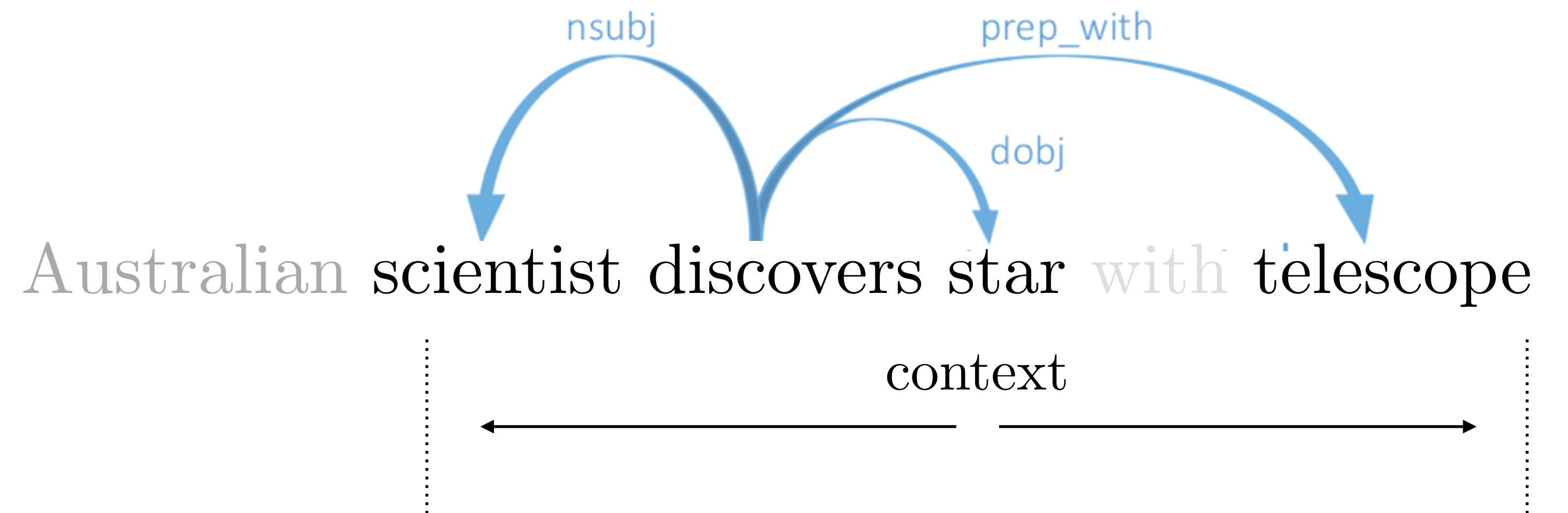
context +/- 2 words



**context  
dependent**



context  
dependent



context  
dependent

	BoW	DEPS
hogwarts	dumbledore hallows half-blood malfoy snape	sunnydale collinwood calarts greendale millfield
topically-similar	vs	'functionally' similar



# word2vec

learn word vectors from sentences

“The fox jumped over **the** lazy dog”

$v_{OUT}$   $v_{OUT}$   $v_{OUT}$   $v_{OUT}$   $v_{OUT}$   $v_{OUT}$

# deepwalk

‘words’ are graph vertices

‘sentences’ are random walks on the graph

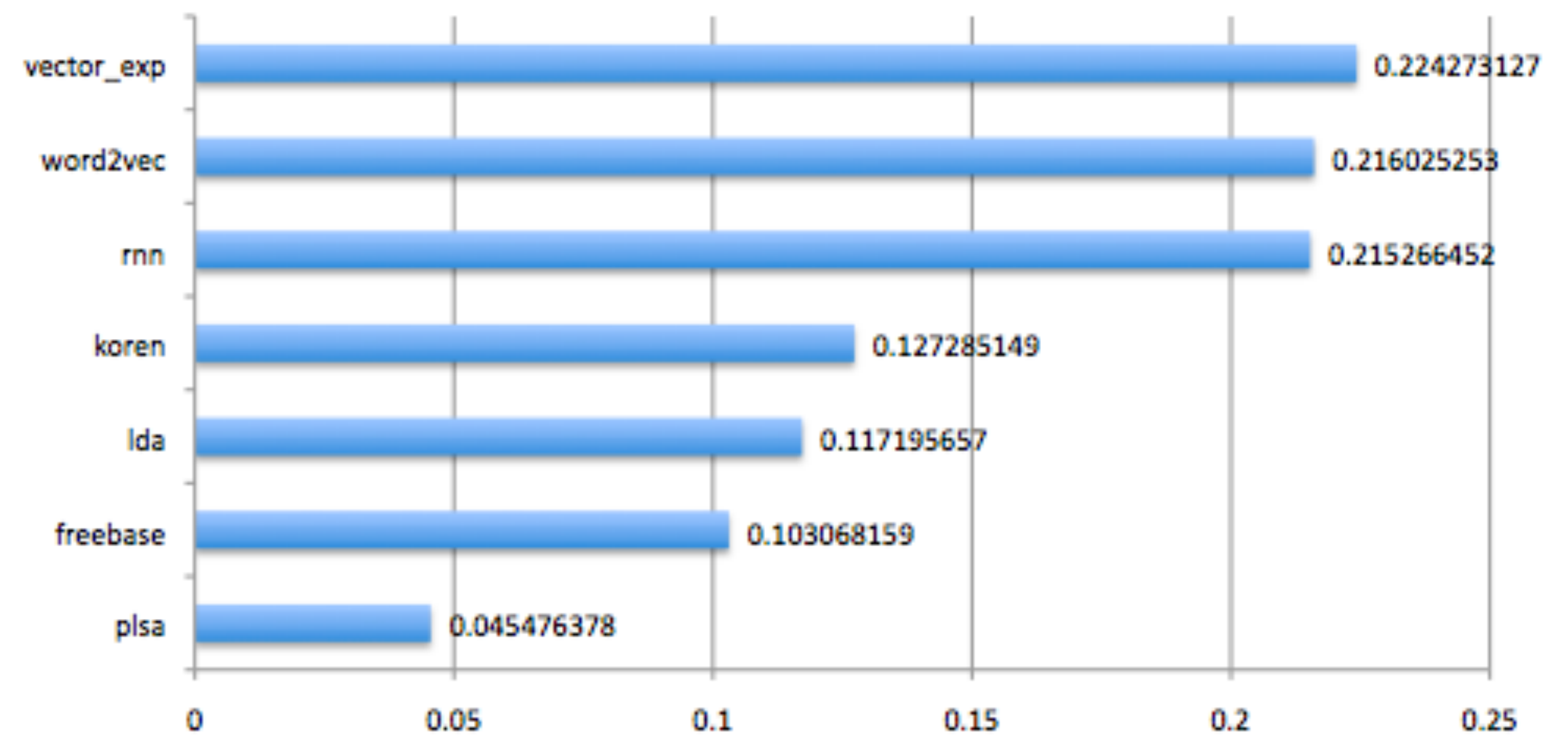
$v_{46} \rightarrow v_{45} \rightarrow v_{71} \rightarrow v_{24} \rightarrow v_5$

# Playlists at Spotify

‘words’ are songs  
‘sentences’ are playlists

# Playlists at Spotify

Great performance on 'related artists'



# Fixes at Stitch Fix?

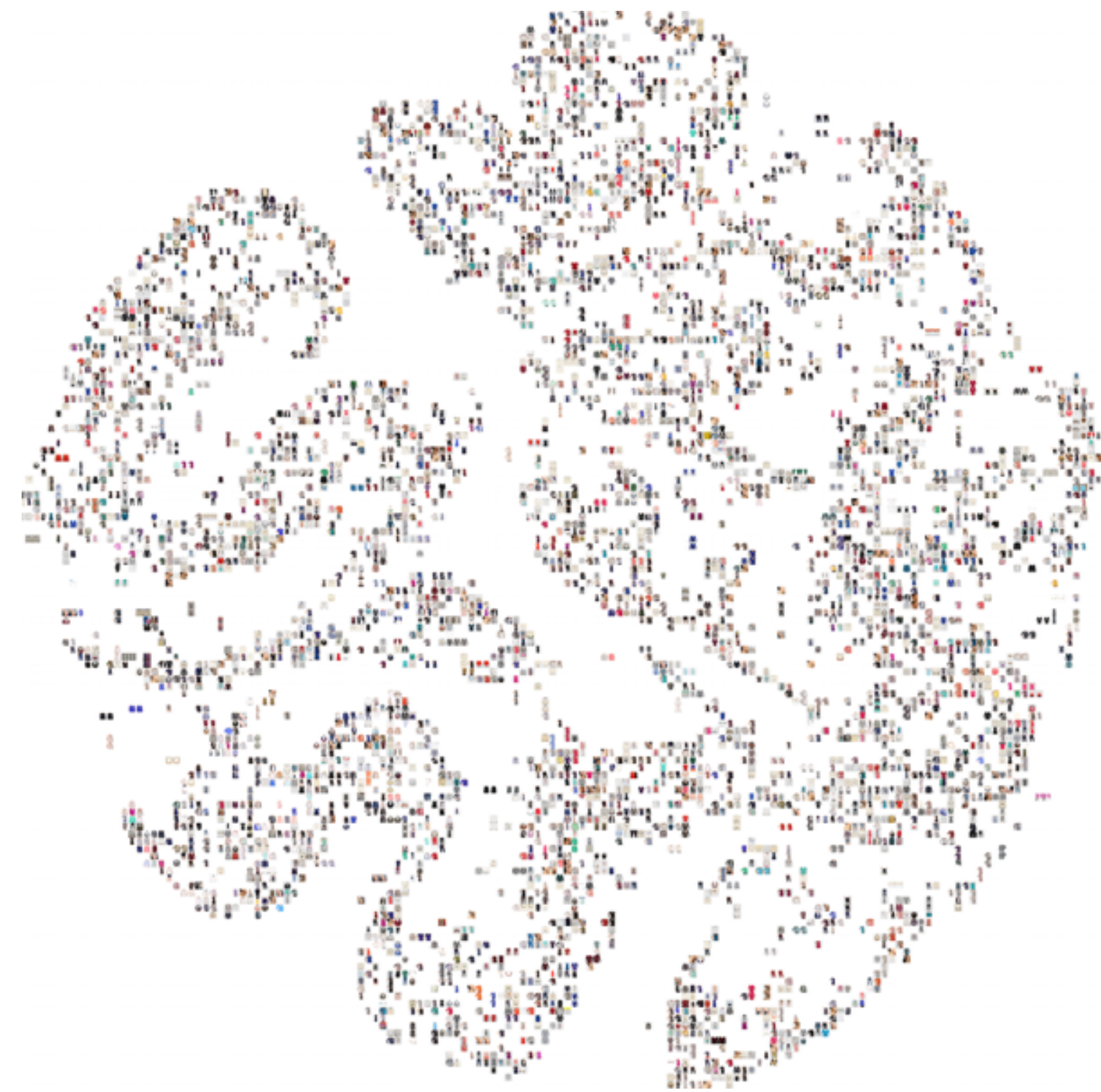
Let's try:  
'words' are styles  
'sentences' are fixes

# Fixes at Stitch Fix?

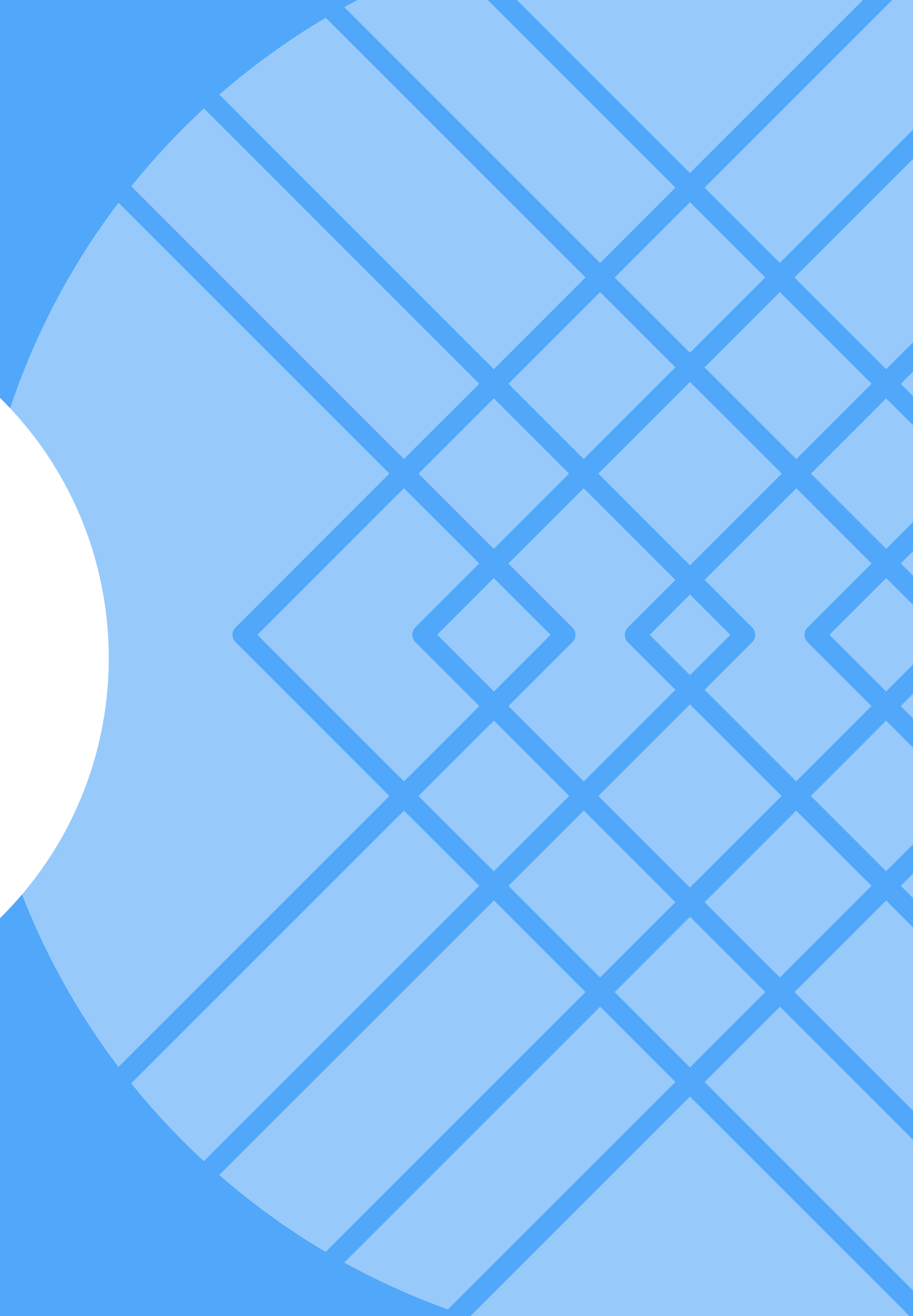
Learn similarity between styles  
because they co-occur

Learn 'coherent' styles

# Fixes at Stitch Fix?



Got lots of structure!



GloVe

sequence  
learning