

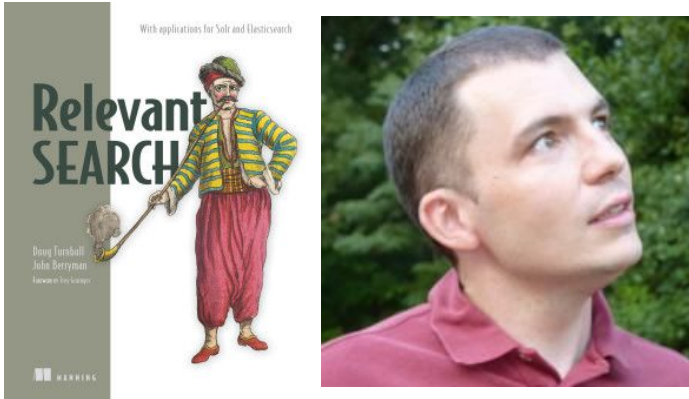


The Neural Search Frontier

How Deep Learning is poised to revolutionize the search relevance landscape



About us



Doug Turnbull, CTO - OpenSource Connections
(<http://o19s.com>)

Relevance raining July 10,11!
<https://opensourceconnections.com/blog/2018/07/10/think-like-relevance-training/>



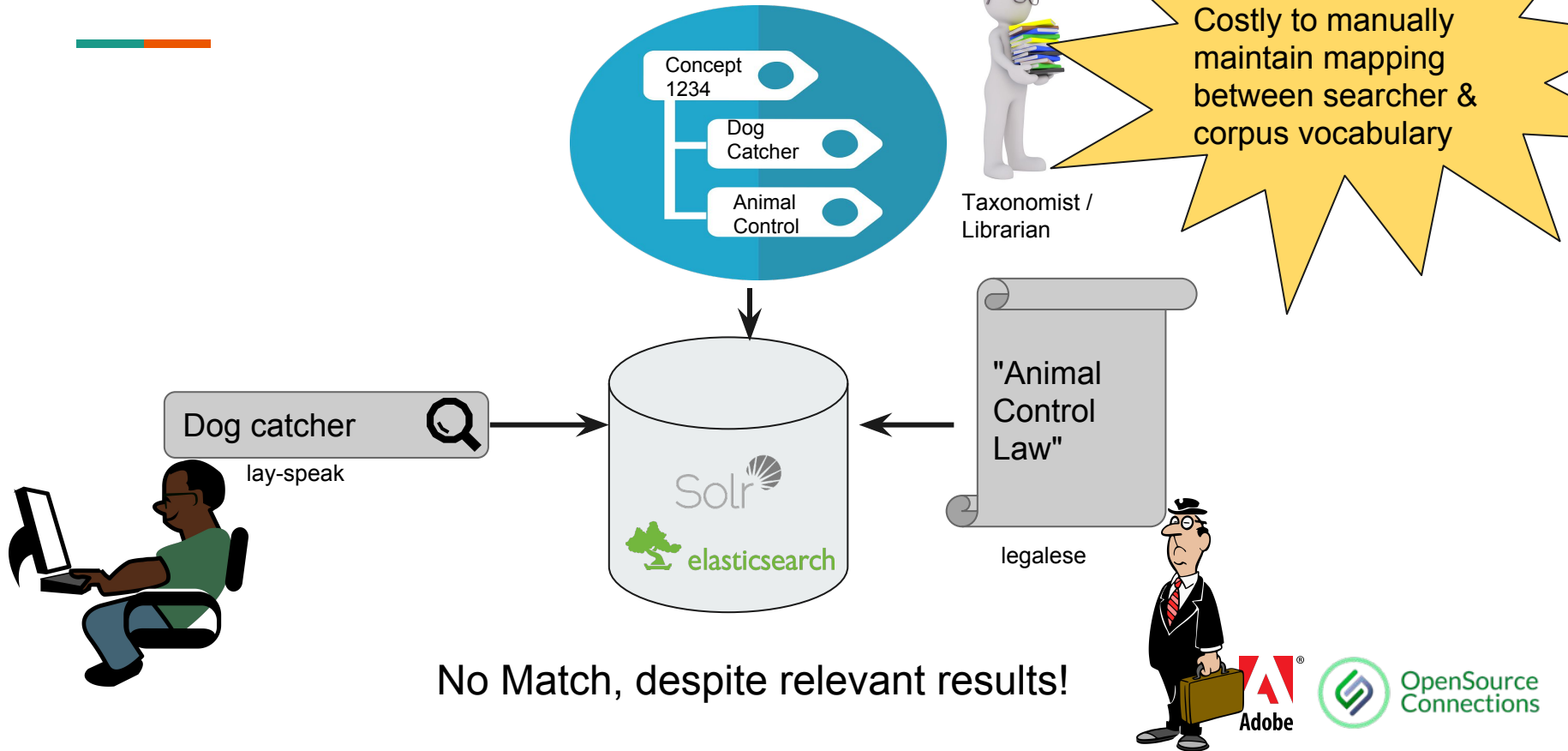
Tommaso Teofili

Computer Scientist, Adobe
ASF member

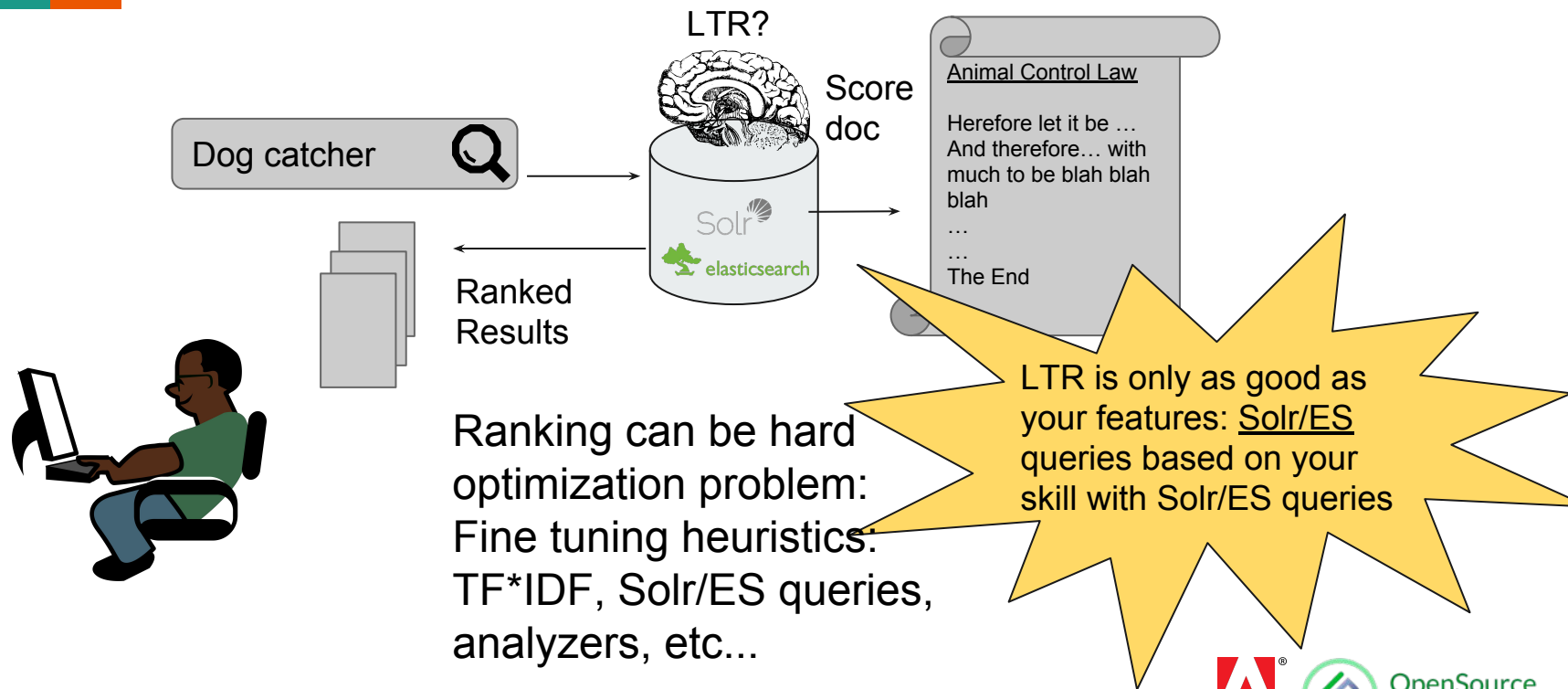


Search Relevance Challenges

Matching: vocab mismatch!



Ranking optimization hard!



Ranking can be hard optimization problem:
Fine tuning heuristics:
TF*IDF, Solr/ES queries,
analyzers, etc...

Embedding-Based Ranking

Embeddings

Word representations that ‘remember’ something about the context they tend appear in

Context: usually defined by the surrounding words

“Cryptocurrency is hyped?”

Appear in similar contexts

“Have you heard the hype about bitcoin?”

I would like to cancel my reservation

Appear in similar contexts

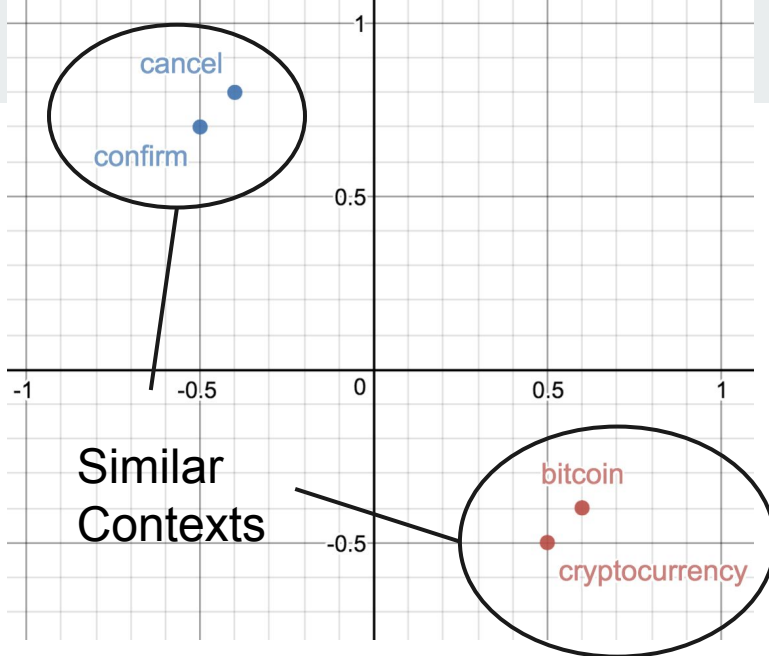
I would like to confirm my reservation



Embeddings are N-Dimensional Vectors that remember context

2D Vector:
[0.5, -0.5] cryptocurrency
[0.6, -0.4] bitcoin

[-0.5, 0.7] cancel
[-0.4, 0.8] confirm



“Cryptocurrency is hyped?”

Appear in similar contexts

“Have you heard the hype about bitcoin?”

I would like to cancel my reservation

Appear in similar contexts

I would like to confirm my reservation

Search is based on similar vector similarity

bitcoin regulation 🔍

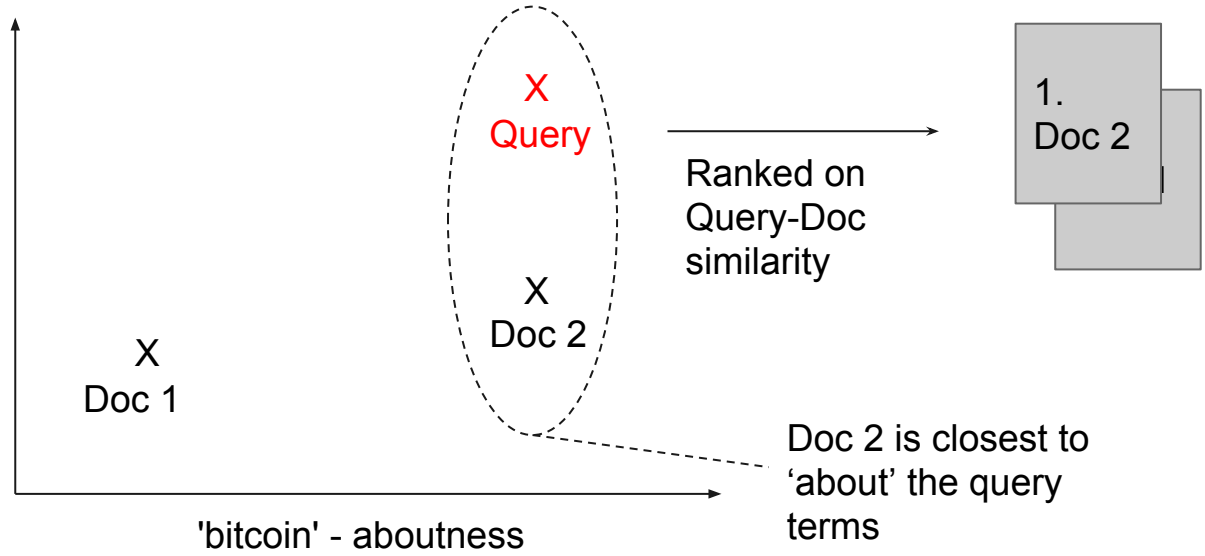
Doc 1
Bitcoin is fun. Gasoline is regulated ...

'regulation' aboutness

Doc 1

Doc 2
Bitcoin regulation guide. Bitcoin is a cryptocurrency. You'll like bitcoin!

Doc 2



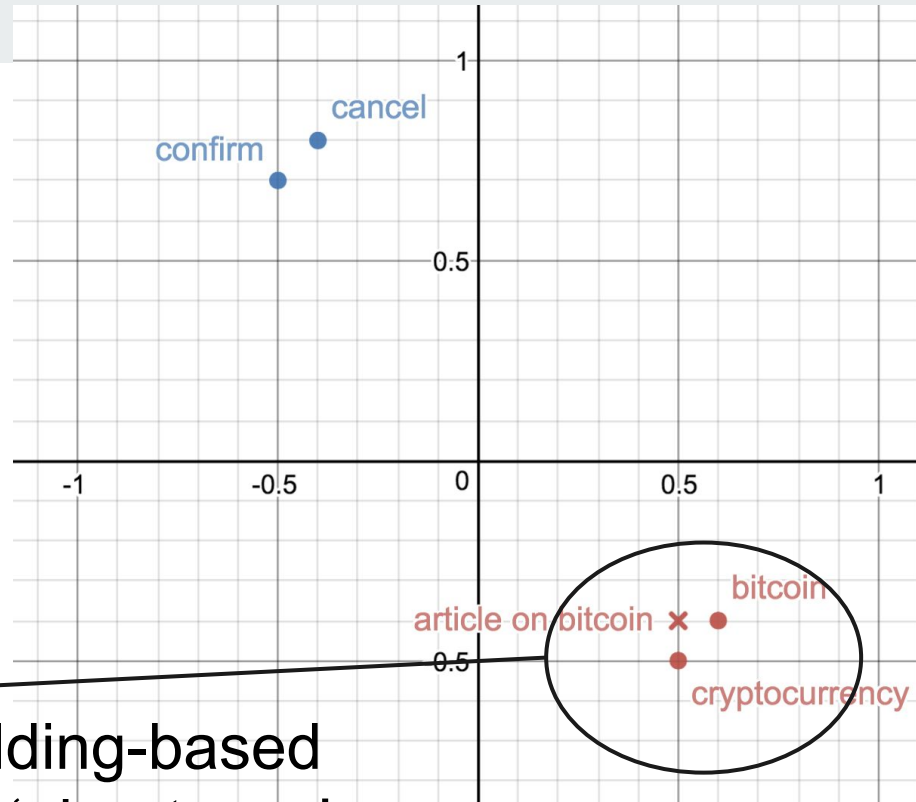
(Term 'aboutness' usually measured via TF*IDF based stats)



Documents also can be embeddings

2D Vector:

- $[0.5, -0.5]$ cryptocurrency
- $[0.6, -0.4]$ bitcoin
- $[0.5, -0.4]$ article on bitcoin



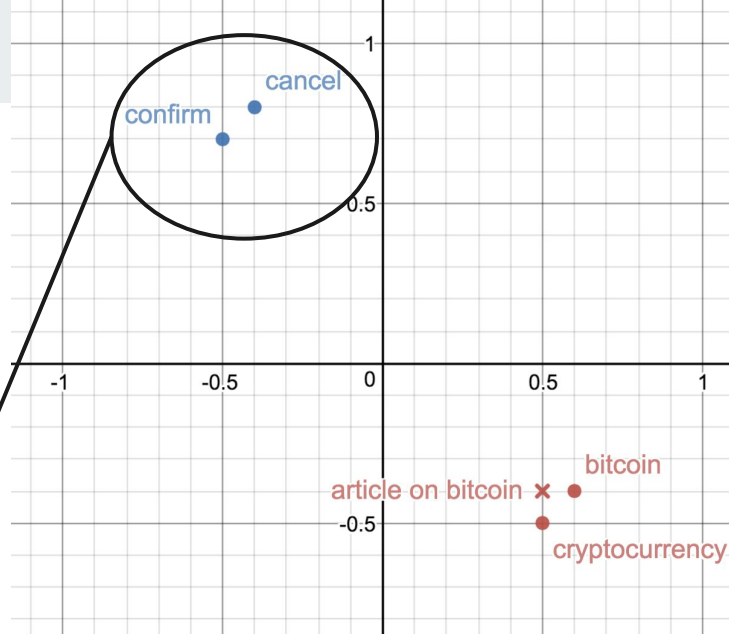
Hypothesis: embedding-based similarity related to 'aboutness' in relevance

(ie article is 'about' cryptocurrency even if it never says the term 'cryptocurrency')

Hypothesis doesn't always hold... (corpuses & search probs vary)

2D Vector:

- [0.5, -0.5] cryptocurrency
- [0.6, -0.4] bitcoin
- [0.5, -0.4] article on bitcoin

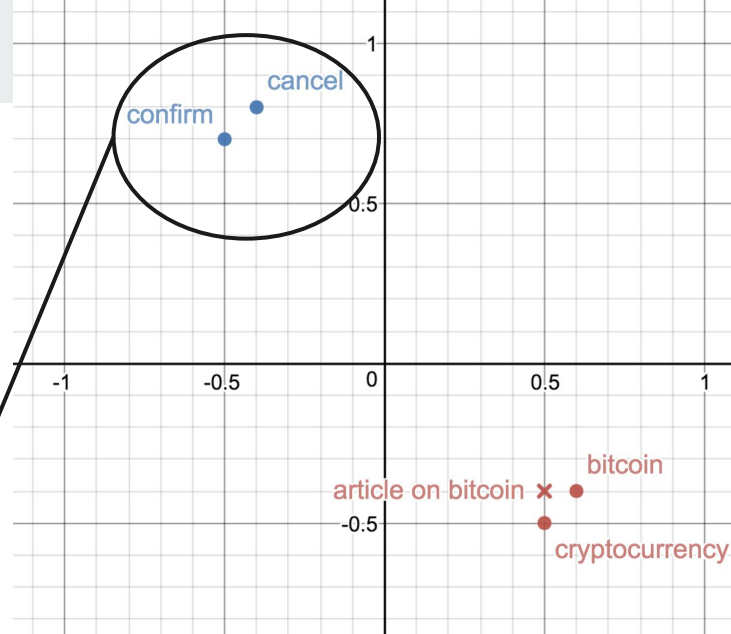


Hypothesis Fails: articles here might be similar in embedding space, but may not map to user notions of relevance

Context alone is not enough (sometimes)

2D Vector:

[0.5, -0.5] cryptocurrency
[0.6, -0.4] bitcoin
[0.5, -0.4] article on bitcoin



We can adjust them using:

- Synsets
- Sentiment
- e.g. <https://arxiv.org/abs/1805.07966>

Low hanging fruit #1

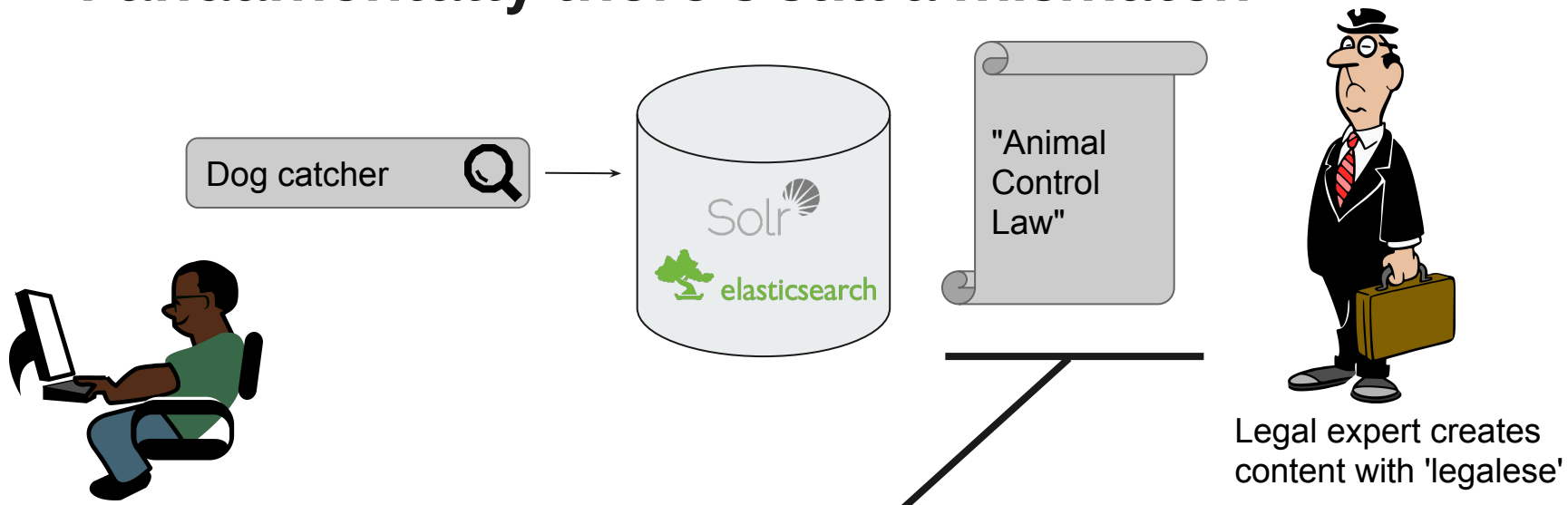


The probability distribution of words in the corpus can provide enough information to predict whether two words appear in similar contexts

Synonyms

[0.5, -0.5] cryptocurrency
[0.6, -0.4] bitcoin

Fundamentally there's still a mismatch



Average Citizen searches in lay-speak

Embeddings derived from corpus, not our searcher's language

Legal expert creates content with 'legalese'

Ranking-based embeddings

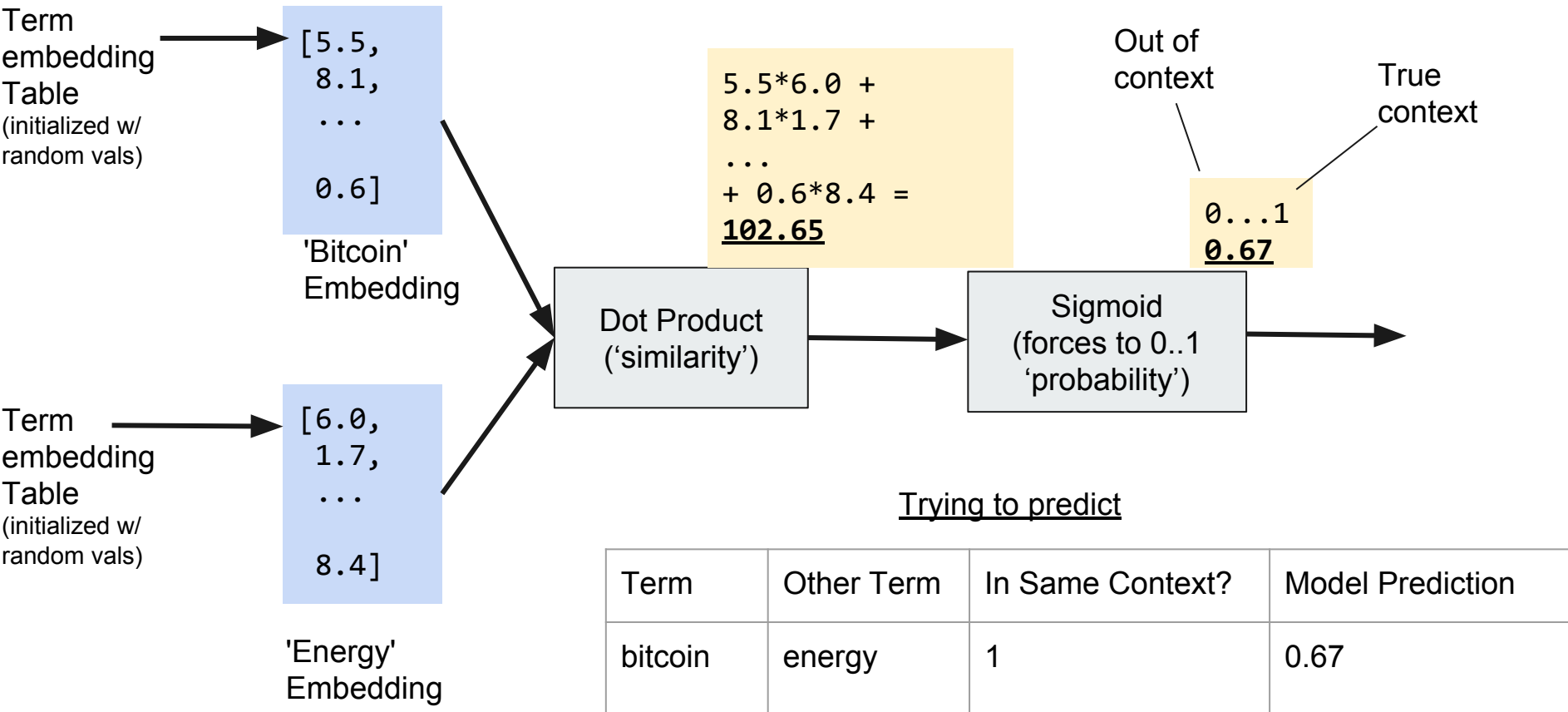
Embetterings - Neural Networks 101

Can we build embeddings that are closer to our searcher's notion of 'relevance'?

Maybe?

But to hack embeddings... we must first understand them

The word2vec Model



Tweak...

	Term	Other Term	True Context?	Prediction	
Terms sharing context	bitcoin	energy	1	0.67	↑
Random unrelated terms	bitcoin	dongles	0	0.78	↓
	bitcoin	relevance	0	0.56	↓
	bitcoin	cables	0	0.34	↓

Goal:
Tweak 'bitcoin' to get prediction closer to 1

Goal:
Tweak 'bitcoin' to get prediction closer to 0

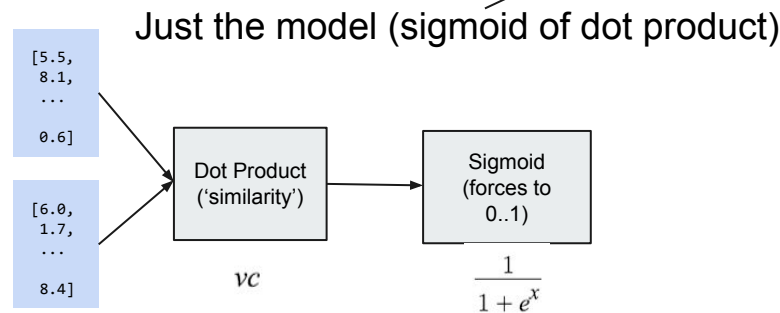
(showing negative sampling method)

Maximize Me!

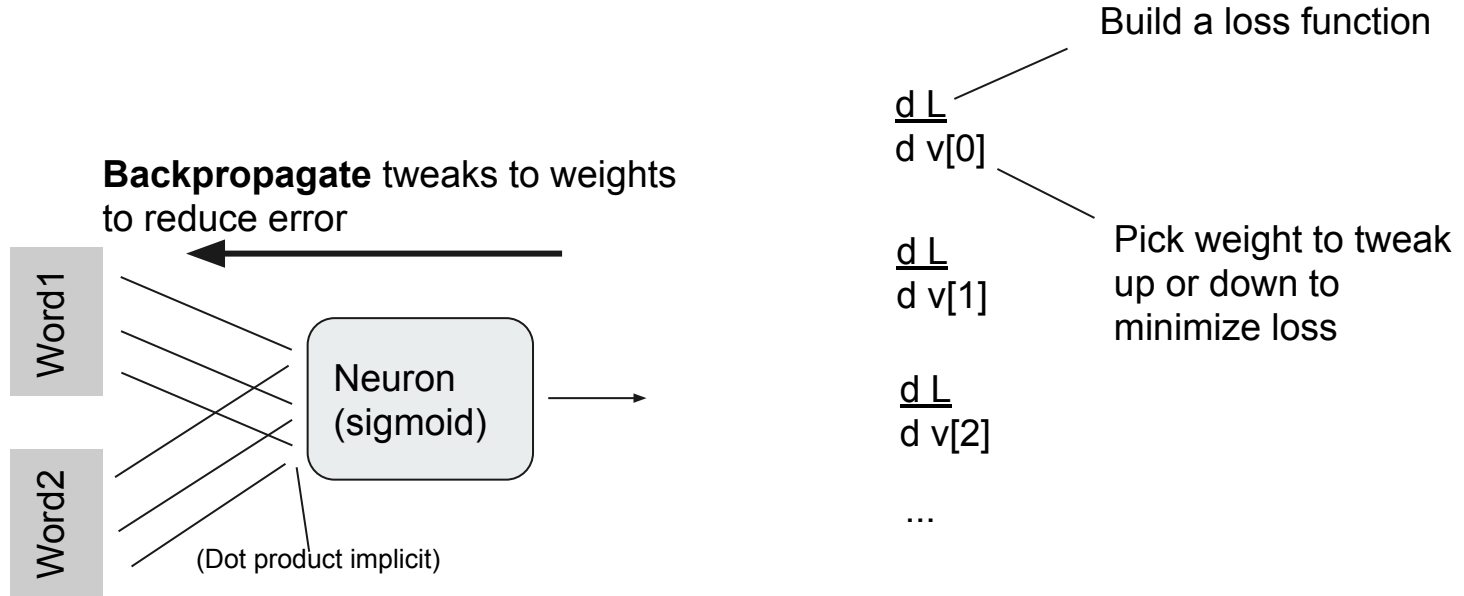
Maximize 'True' Contexts Minimize 'False' Contexts

$$L = \log\left(\frac{1}{1 + e^{vc}}\right) - \sum_{i=0}^k \log\left(\frac{1}{1 + e^{vc_i}}\right)$$

Objective /
'Loss' Function /
'Cost' Function



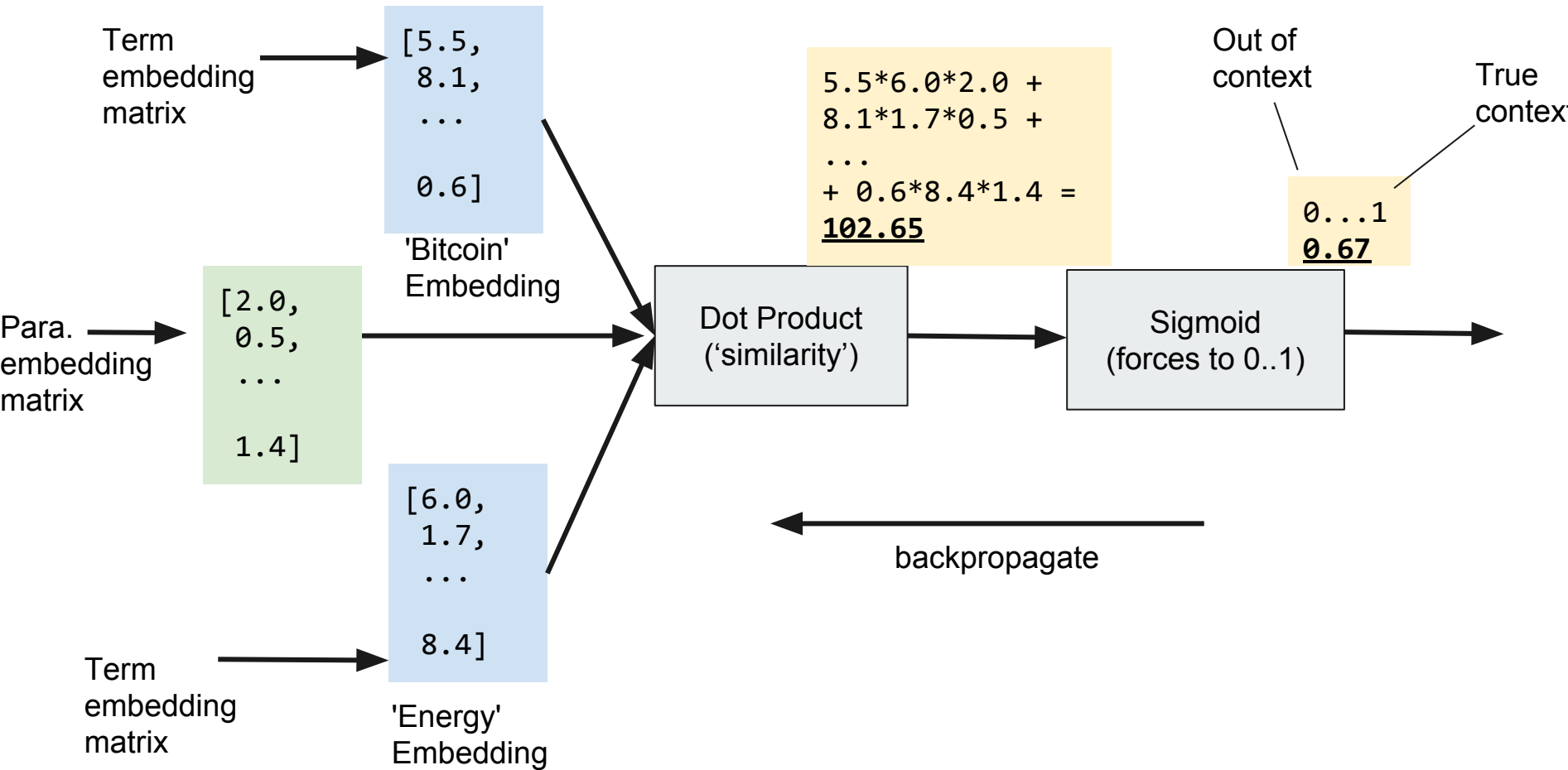
Did you just neural network?



More complex/'deep' neural nets, can **propagate** back error to earlier layers to learn weights



Doc2Vec, train paragraph vector w/ term vectors



Same neg sampling can apply here



Terms sharing context

Random unrelated docs

Term	Para	True Context?	Prediction
bitcoin	Doc 1234 Para 1	1	0.67
bitcoin	Doc 5678 Para 5	0	0.78
bitcoin	Doc 1234 Para 8	0	0.56
bitcoin	Doc 1537 Para 1	0	0.34



Tweak embeddings to get prediction closer to 1



Tweak embeddings to get prediction closer to 0



Low hanging fruit #2

Similar word embeddings lie close to one another

Similar document embeddings lie close to one another

Document embeddings lie close to word embeddings representing important words or topics for them

Retrieval

Did anyone say
terms???

Ranking

Did anyone say
frequencies ???

~~Low~~ hanging fruit #2

Corpus and search vocabularies can differ a lot

Combining traditional IR and neural models can help in filling the gaps

Retrieval

Terms and vectors

Ranking

Frequencies and distances

Neural Search Hack1: Relevance-based embeddings?



	Query Term	Doc	Relevant?	Relevance Score	
Relevant doc for term	bitcoin	Doc 1234	1	0.67	↑
Random unrelated docs	bitcoin	Doc 5678	0	0.78	↓
	bitcoin	Doc 1234	0	0.56	↓
	bitcoin	Doc 1537	0	0.34	↓

Tweak embeddings to get query-doc score closer to 1

Tweak embeddings to get prediction closer to 0

(derived from clicks, etc)

What about queries/docs we haven't seen?



Pretrain word2vec with corpus/sessions?

(a kind of prior)

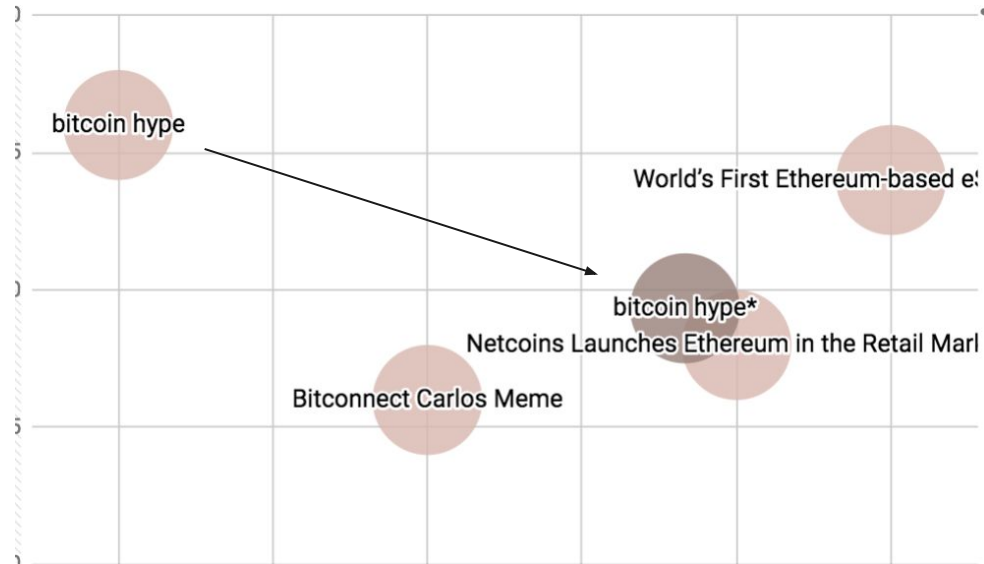
Query Term	Doc	Corpus Model	True Relevance	Tweaked Score
bitcoin	Doc 1234	0	1	0.01
bitcoin	Doc 5678	1	0	0.99
bitcoin	Doc 1234	0	0	0.56
bitcoin	Doc 1537	0	0	0.34



Takes a lot to overcome original unified model, and its a case-by-case basis

An online / evolutionary approach to converge on improved embeddings?

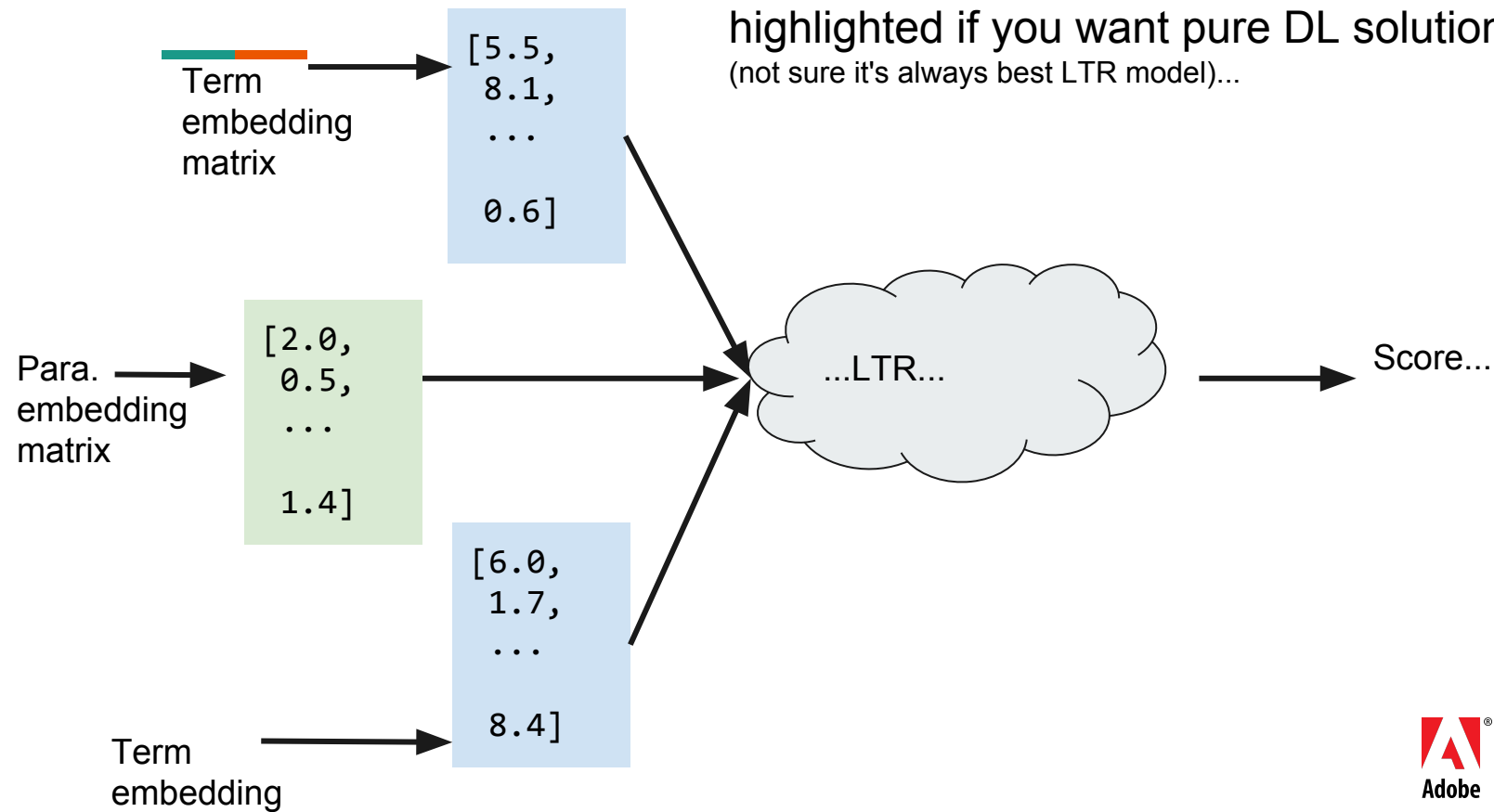
Geometrically adjust query embeddings



Averaging query vectors by means of some possibly relevant documents

LTR feature...

See today's mobile.de talk where **RankNet** is highlighted if you want pure DL solution (not sure it's always best LTR model)...

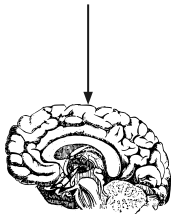


Neural Language Models

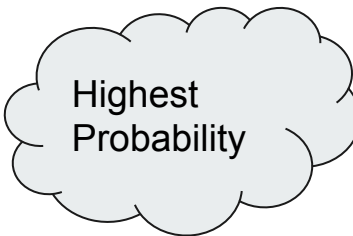
Language Model

Given what's been spoken thus far, predict the next word / character

“eat ?”



cat: 0.001
...
chair: 0.001
pizza: 0.02
...
nap: 0.001



Obvious applications: autosuggest, etc



Markov language model

From corpus we can build a **transition matrix** reflecting frequency of word transitions:

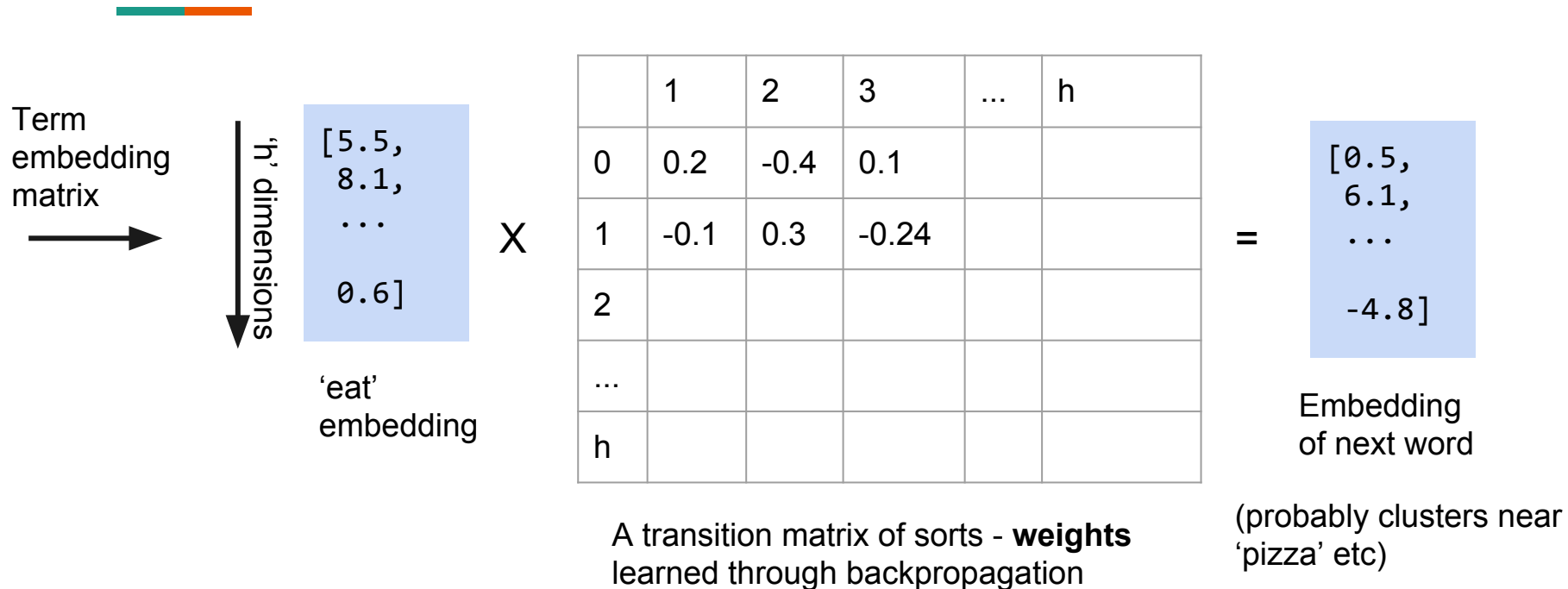
V words

V words

	pizza	chair	nap	...
eat	0.02	0.0002	0.0001	...
cat	0.0001	0.0001	0.02	...
...	

Probability of 'pizza' following 'eat' (as in "eat pizza")

Markov model for embeddings?



Accuracy requires context...

The race is on! eat pizza!

Context + Word => Next Word

(and new context)

Old Context + Word => New Context & Word

[-1.6,
7.3
...]
Output
Embedding

'Embedding' markov model from earlier

'Embedding' for next word

[1.6,
-5.4
...]

||

[0.5, 1.6 ...
6.1, -4.8 ...
...
]

X

[5.5,
8.1,
...]

Eat (embedding)

Context + new word => New Context/word

Prev Context -> New Context Transition

X [0.5, 1.6 ...
6.1, -4.8 ...
...
]

= [0.3,
4.5
...]

+

[-0.9,
-1.2
...]

New Context

||

[...]

X

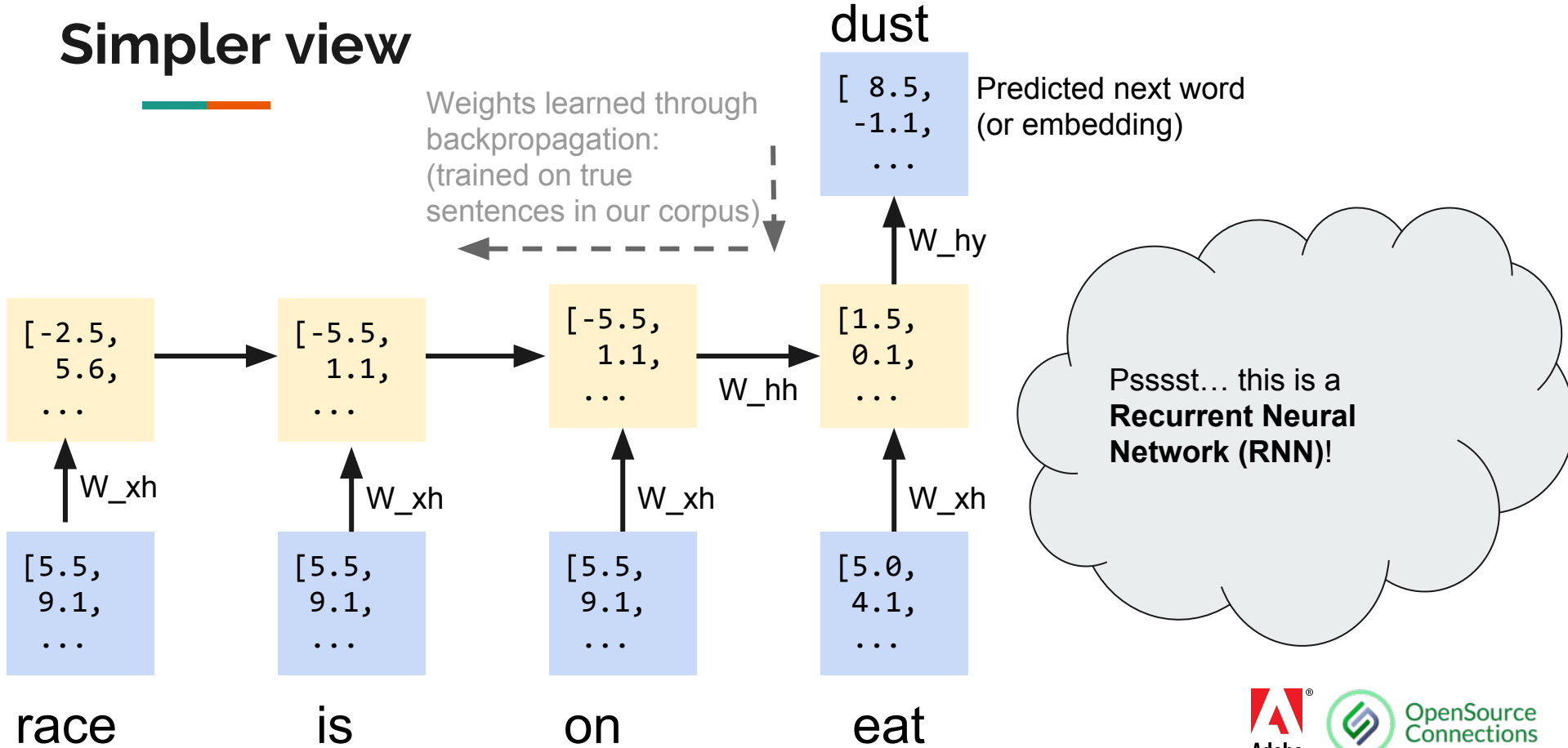
[5.5,
8.1,
...]

Pizza (embedding)

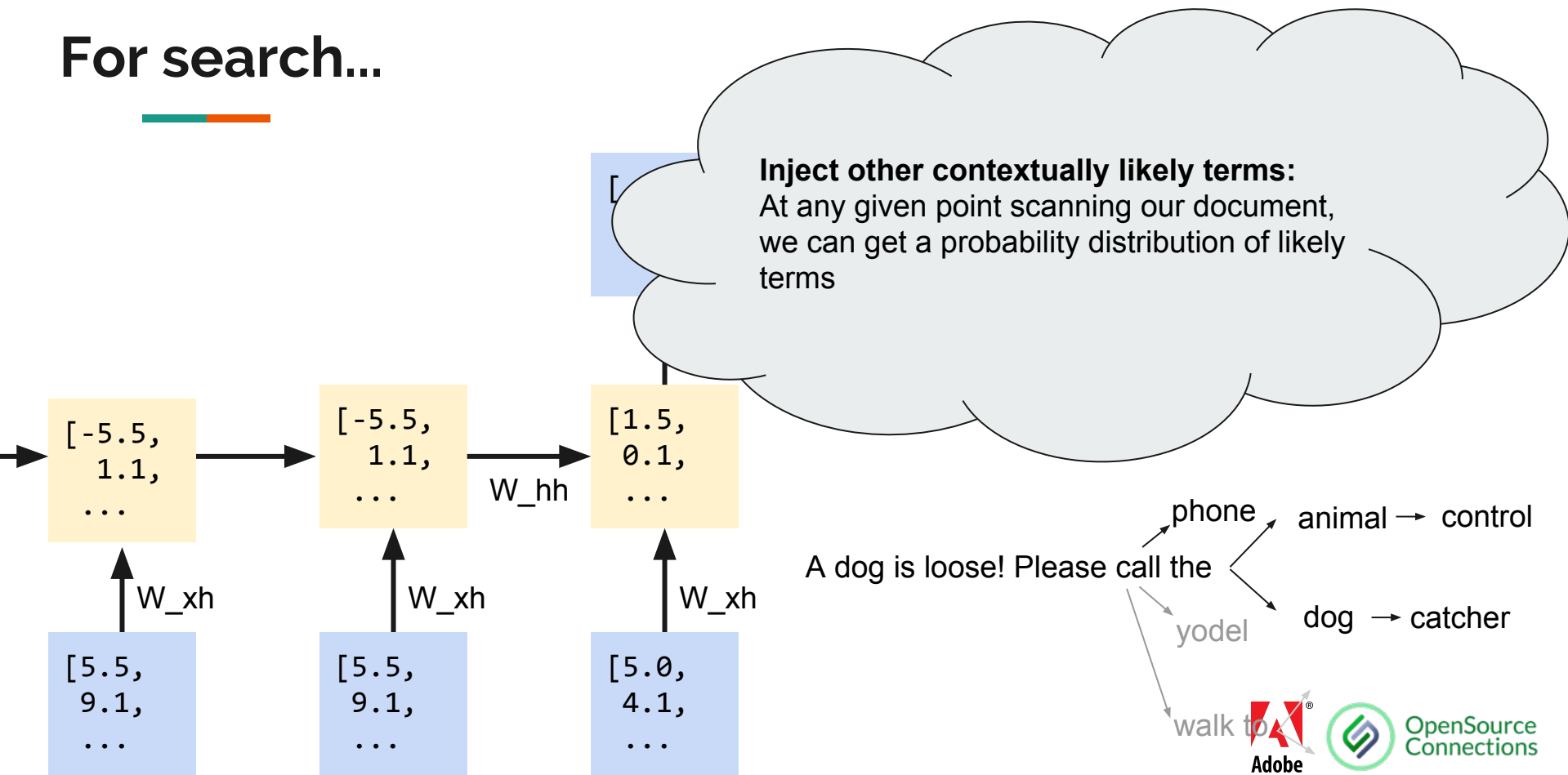
Context -> Output Transition

Next...

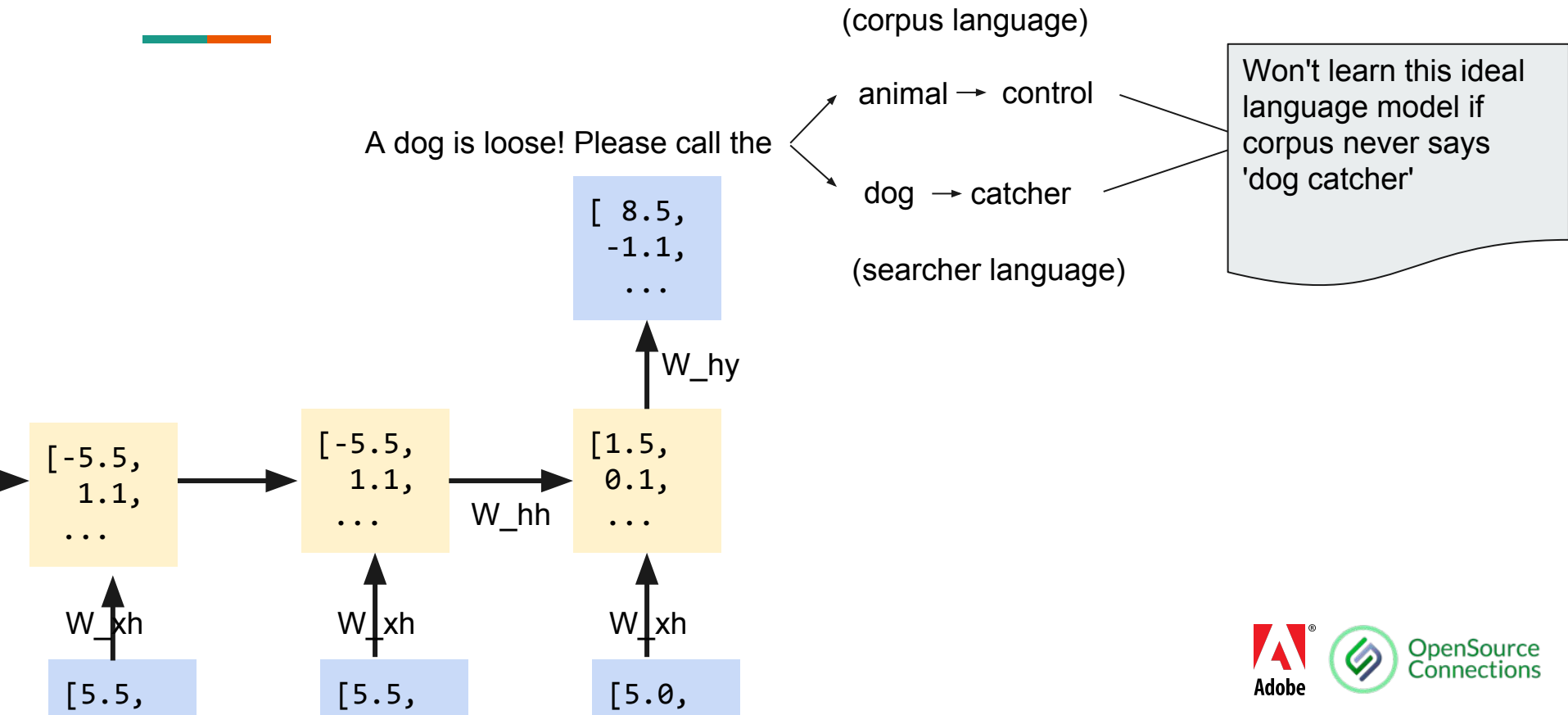
Simpler view



For search...



Not a silver bullet



One Vector Space To Rule Them All: Seq2Seq

Sessions as documents



Dog catcher 🔍

Dog bit me 🔍

Lost dog 🔍

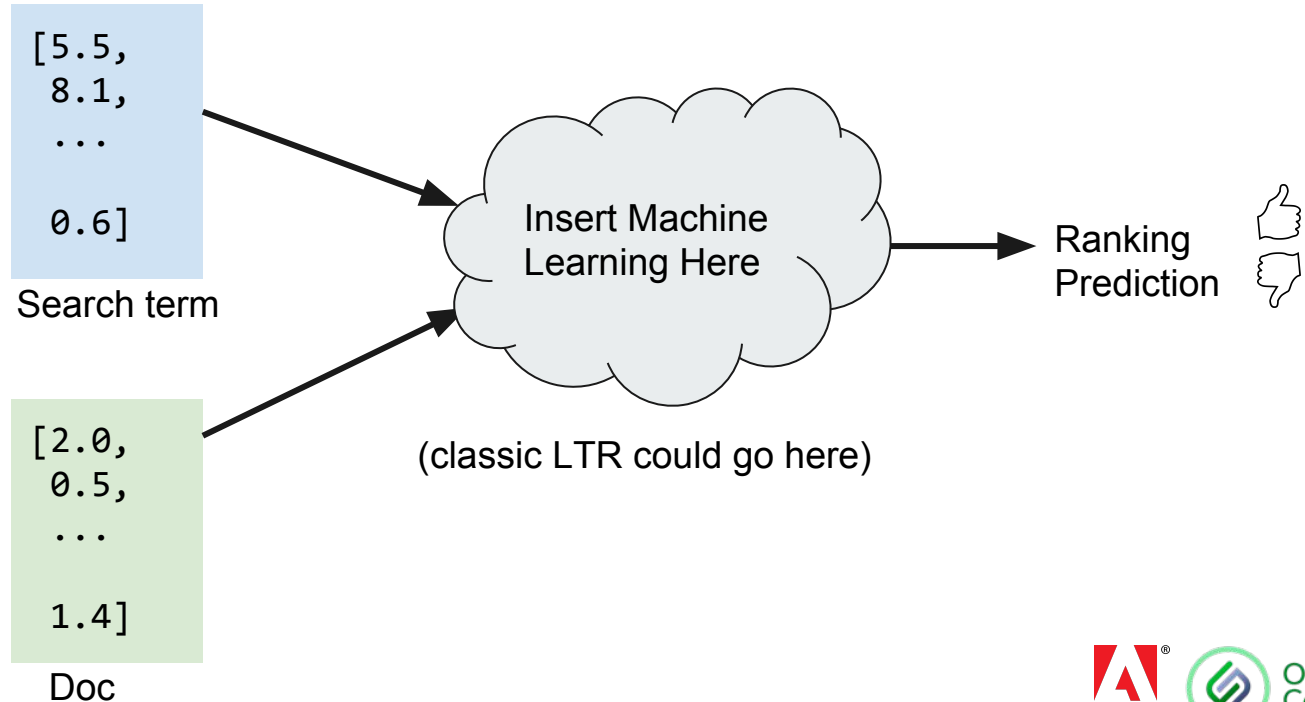
→
word2vec/doc2vec

[5.5,
8.1,
...
0.6]

Embeddings
built for just
query terms



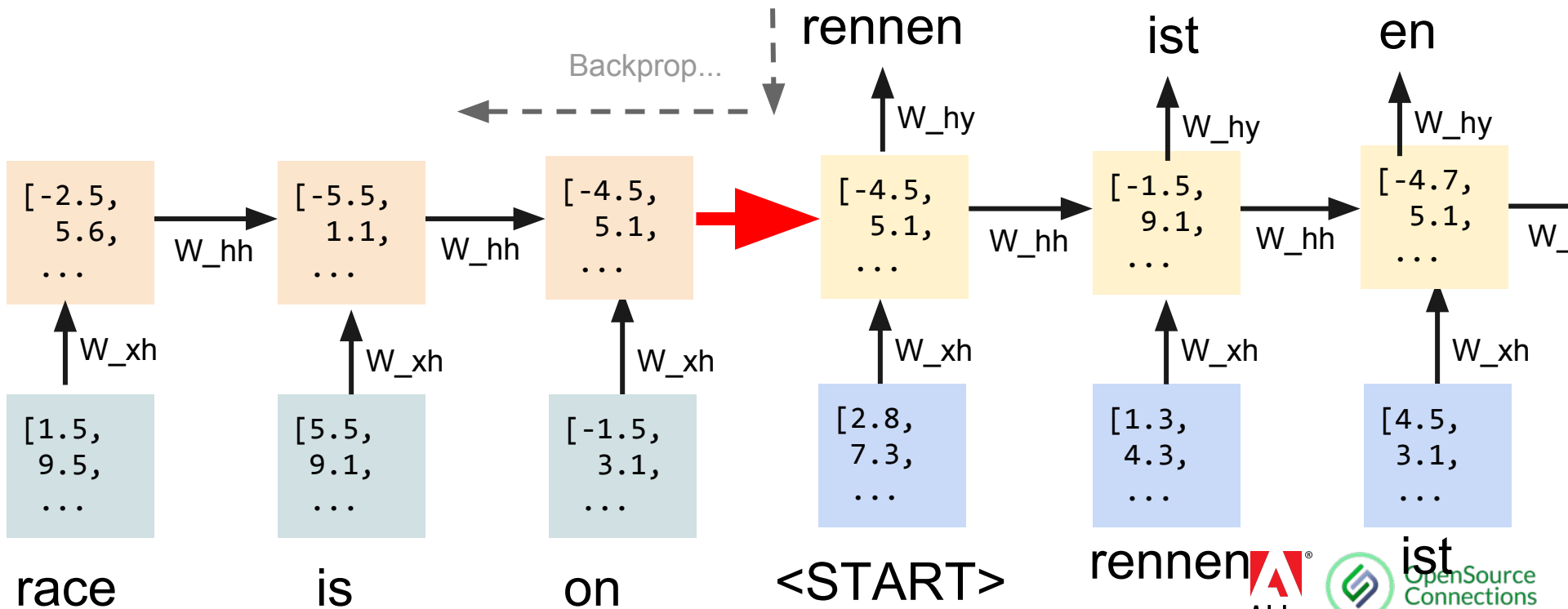
Can we *translate* between searcher & corpus?



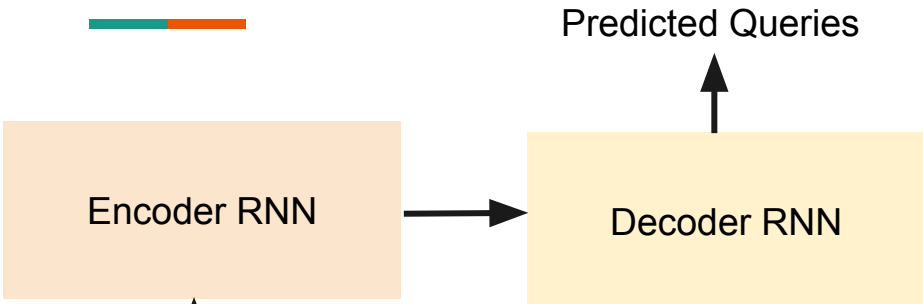
Translation: RNN decoder/encoder

Encoder RNN

Decoder RNN



'Translate' documents to queries?

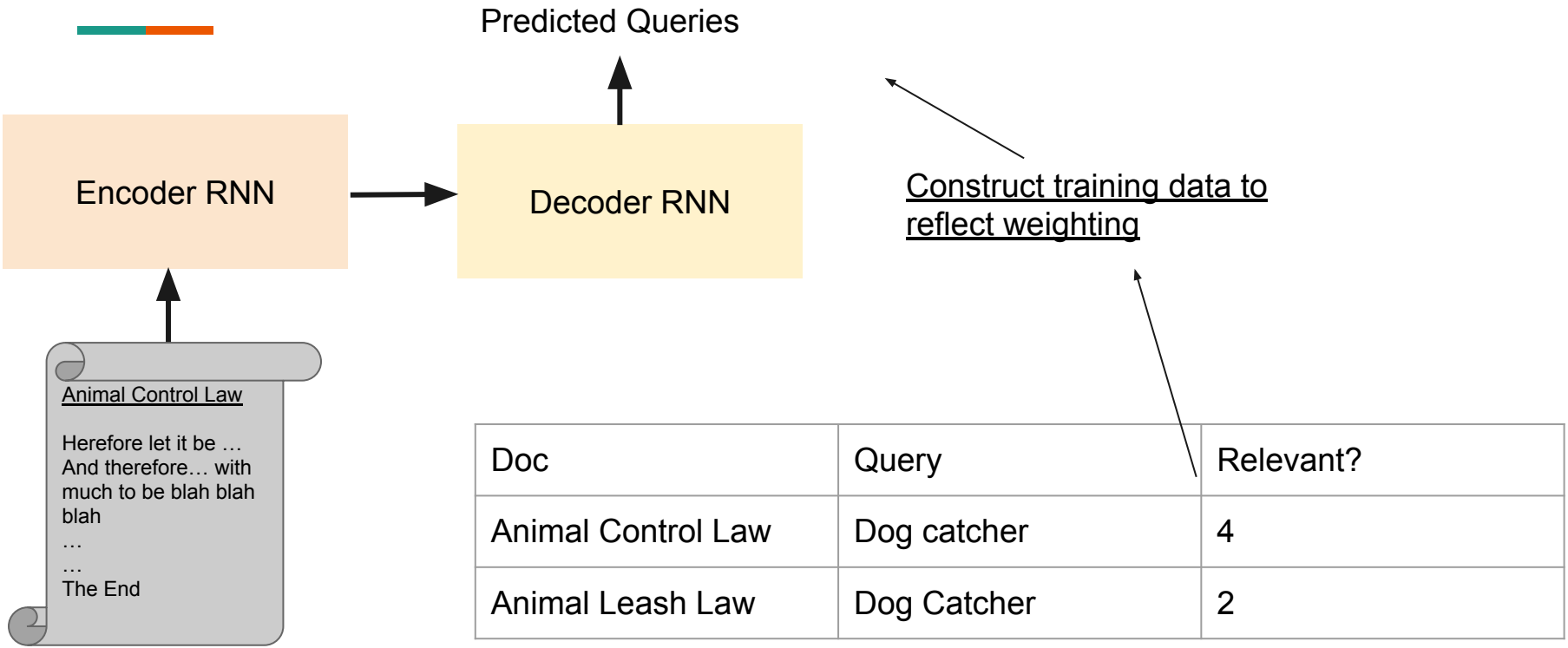


Animal Control Law
Herefore let it be ...
And therefore... with
much to be blah blah
blah
...
...
The End

Using judgments/clicks as training data:

Doc	Query	Relevant?
Animal Control Law	Dog catcher	1
Littering Law	Dog Catcher	0

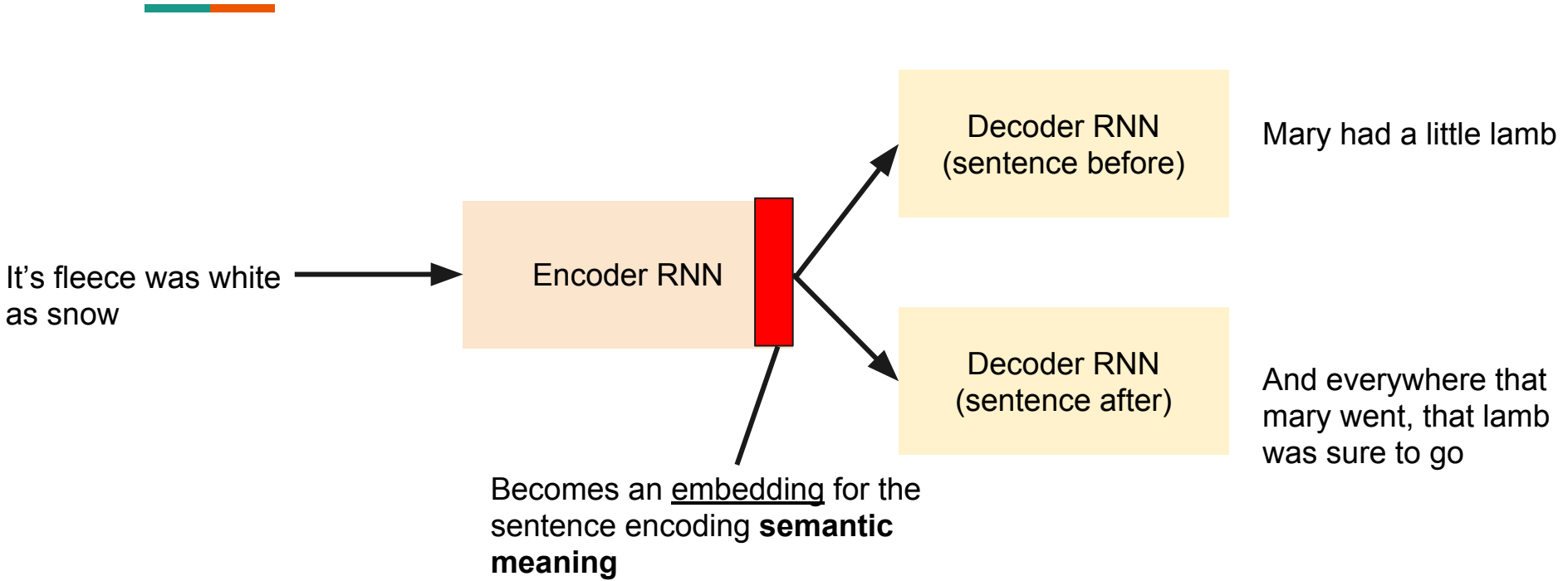
Can we use graded judgments?



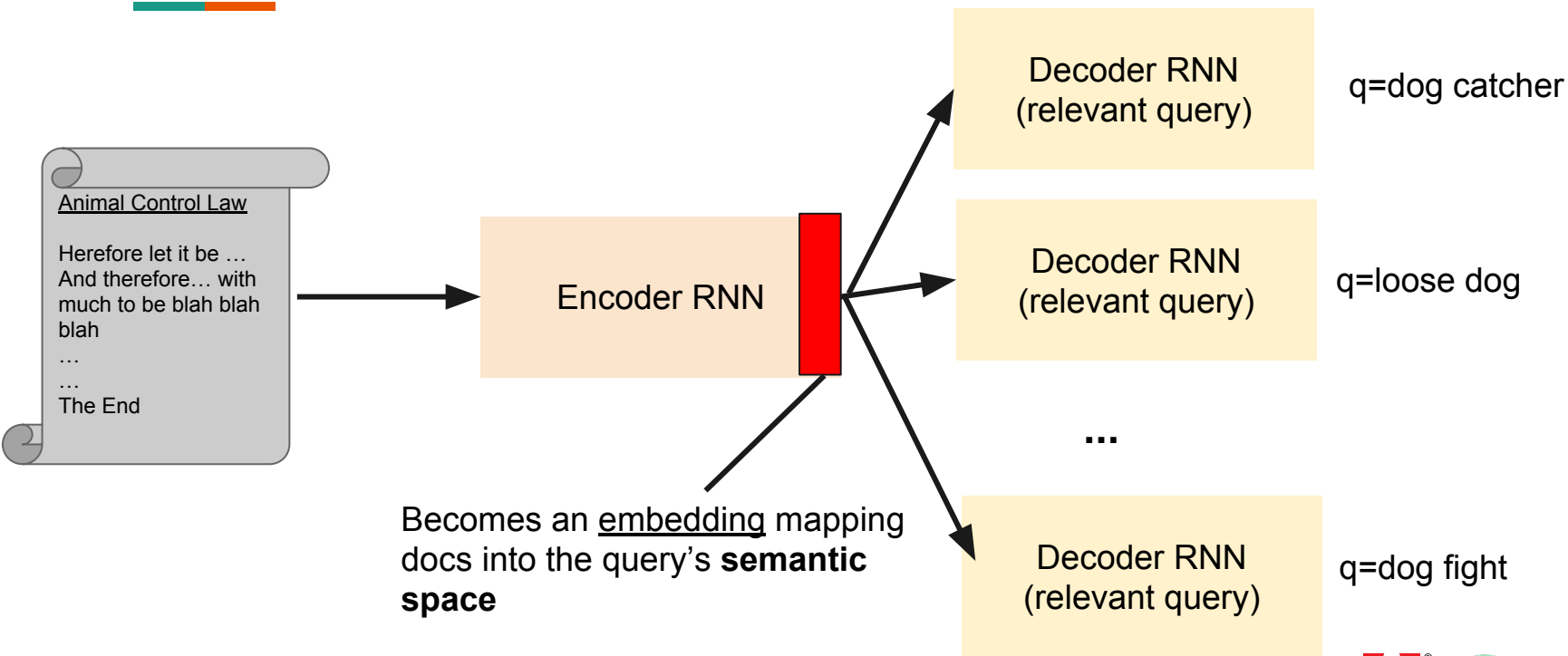
Animal Control Law

Herefore let it be ...
And therefore... with
much to be blah blah
blah
...
...
The End

Skip-Thought Vectors (a kind of 'sentence2vec')



Skip-Thought Vectors for queries? ('doc2queryVec')





Demo



The Frontier

Anything2vec

Deep learning is especially good at learning representations of :

- Images
- User History
- Songs
- Video
- Text

If everything can be 'tokenized' into some kind of token space for retrieval

Everything can be 'vectorized' into some embedding space for retrieval / ranking



Solr/ES community needs to get better first-class vector support

Similarity API and vectors ?

long computeNorm(**FieldInvertState** state)

SimWeight computeWeight(**float** boost, **CollectionStatistics**
collectionStats, **TermStatistics**... *termStats*)

SimScorer simScorer(**SimWeight** weight, **LeafReaderContext** context)



Unsupervising the future

With supervised we can't get much further than the quality of the data itself

Can we do better with unsupervised learning?



It's not magic, it's math!

Join the Search Relevancy Slack Community

<http://o19s.com/slack>

(projects (Elastic LTR!), chat, conferences (Haystack!), book authors, and more...)