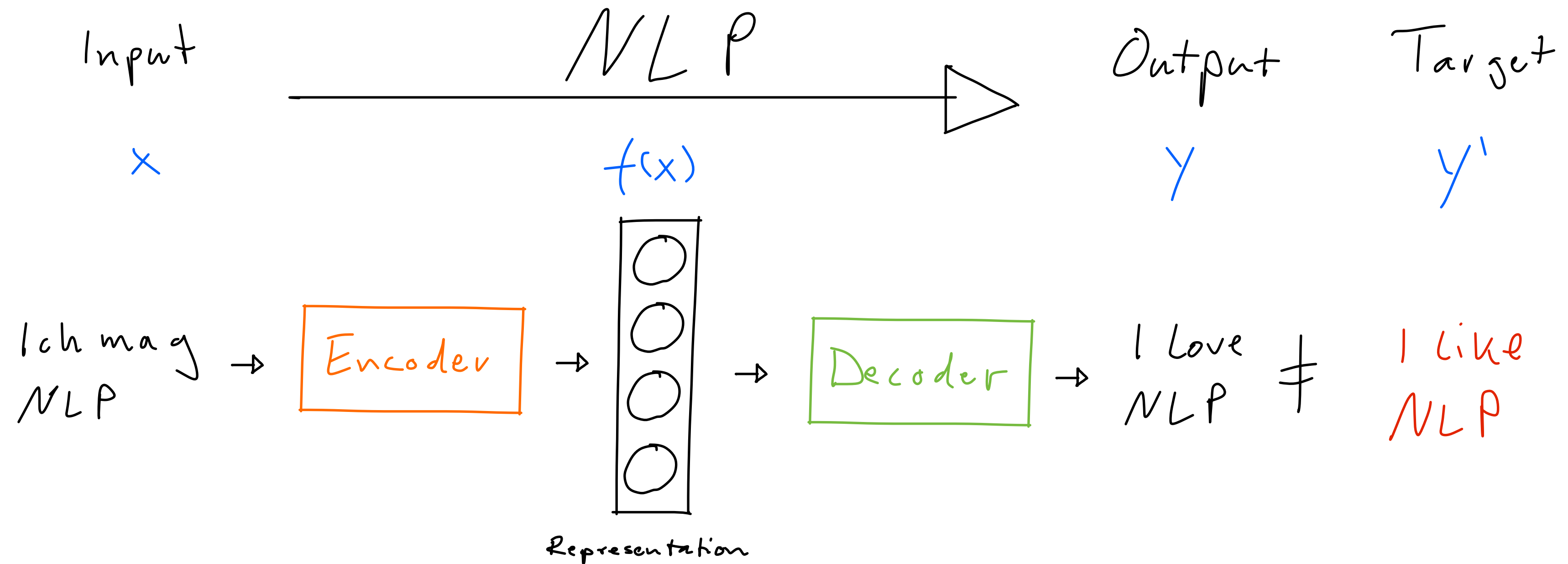# Word Embeddings and Topic Modelling

Tim Rocktäschel & Sebastian Riedel
COMP0087 Natural Language Processing

# NLP in a Nutshell

Input       $NLP$       Output    Target

$x$           $f(x)$          $y$       $y'$

Ich mag NLP → **Encoder** → → **Decoder** → I Love NLP ≠ I like NLP

Representation

| Sentences | Manual | symbolic | softmax | labels | supervised |
|-----------|--------|----------|---------|--------|------------|
| Documents | Learnt | sparse | Sequence decoder | sentences | self supervised |
| Languages | Pretrained | dense | tree decoder | trees | weak |
| Domains | RNN | attention | | graphs | unsup. |
| Databases | CNN | | | | semi-sup. |

# Classify This!

Train $\left\{\begin{array}{l} \text{blah} \quad \text{blah} \quad \text{football} \quad \longrightarrow A \\ \text{blah} \quad \text{blah} \quad \text{stocks} \quad \longrightarrow B \end{array}\right.$

Test $\quad$ blah blah hockey $\longrightarrow$ ?

# Machine Sees this

Train {
blah    blah   symbol 1     → A

blah    blah   symbol 23   → B

Test    blah   blah   symbol 42   → ?

# Word Representations



Encoder

Decoder

hockey

stocks

football

Lookup

+

argmax

Biz

X

f(x)

Y

One-hot vectors can't capture similarities

5

# Fixing one hot vectors

we have **this**



but we want **this**



6

# NLP in a Nutshell

Input       $NLP$          Output    Target

$x$           $f(x)$          $y$    $y' = g(x)$

Ich mag NLP → Encoder → [Representation] → Decoder → I Love NLP $\neq$ I like NLP

Easy to get

Reuse this

Representation

Expensive to annotate

use input as training signal
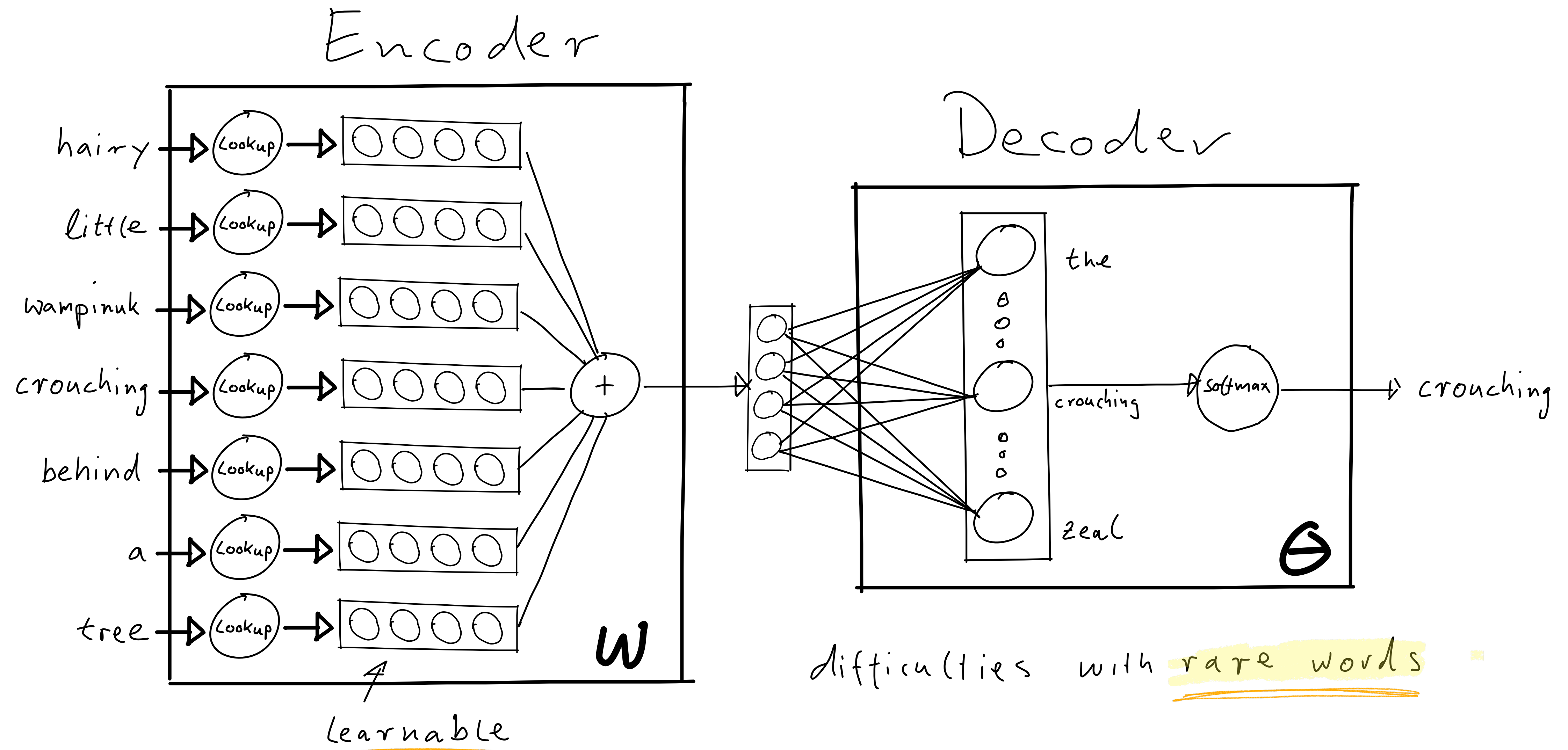
What do we get with $g(x) = x$ ?

# Wampinuk

Marco saw a hairy little **wampinuk** crouching behind a tree

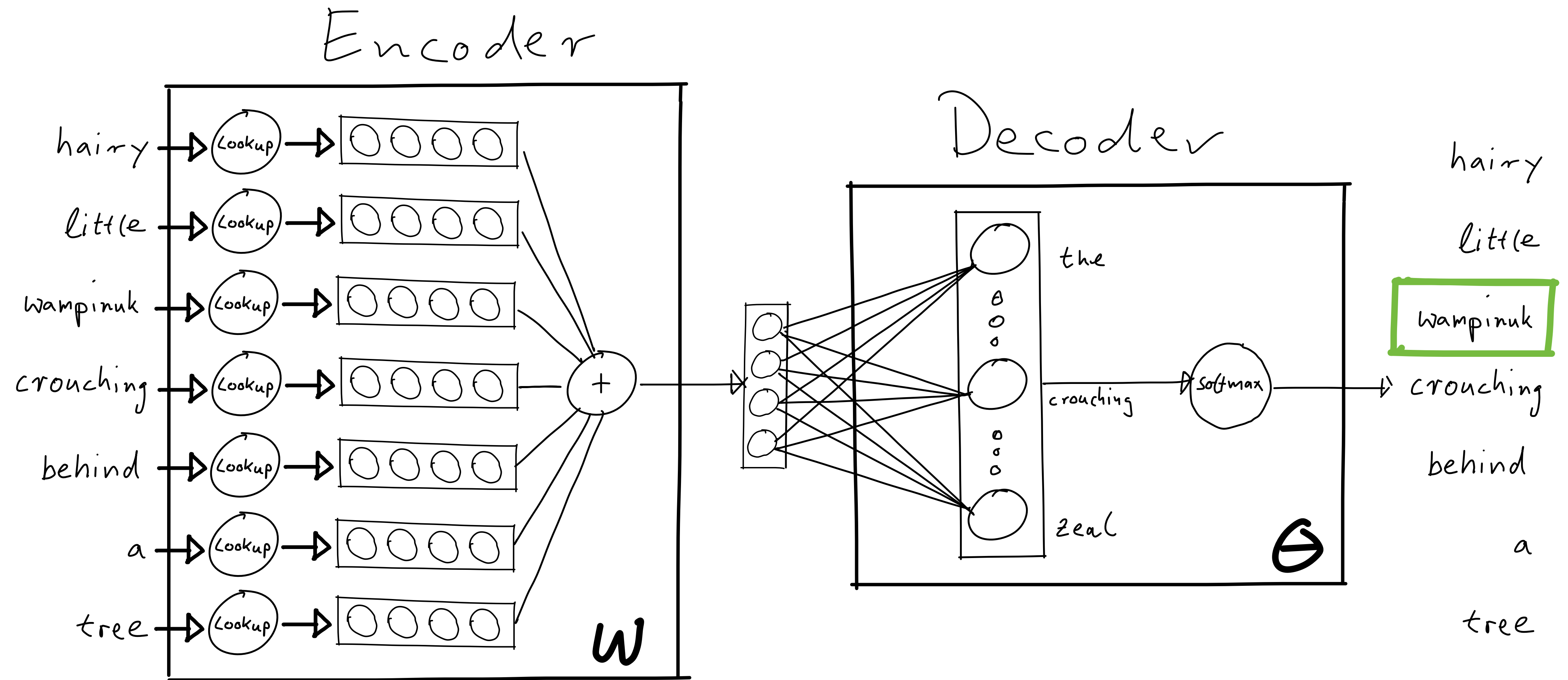You shall know aword by the company it keeps (Firth, J. R. 1957:11)

# Word2Vec: Continuous Bag of Words



Encoder

hairy → Lookup →
little → Lookup →
wampinuk → Lookup →
crouching → Lookup →
behind → Lookup →
a → Lookup →
tree → Lookup →

W

learnable

Decoder

the
crouching
zeal

Softmax → crouching

θ

difficulties with rare words

9

# Word2Vec: Skip Gram

# Skip Gram

# Skip Gram



Encoder

Decoder

crouching → Lookup

W

the

crouching

zeal

Expensive

softmax

hairy

little

wampinuk

behind

a

tree

# Binary Classifier



Encoder

Decoder

hairy

little

wampinuk

the

crouching → Lookup

crouching

softmax

zeal

behind

a

tree

W

θ

# Binary Classifier

# Negative Sampling

Encoder

Decoder

randomly sampled word

Write

crouching → Lookup

Write

W

15

# Word2vec 2D Projections

Skiing
Parasailing
Sailing
KiteBoarding
Swimming
Running
Mountain Biking
Scuba Diving
Motorcycling

Backpacking
Hiking
Camping

Fishing
Seafood

Scientists

Physicists

SoulCycle

Computer Engineering
Engineers

Designers

**Good**
●

Parks
Burning Man

Photographers
Travel Photography

Modern Art

**Where do you
expect "Bad"?**

Travel
Street Art
Museums

Data Mining

Art Galleries

Lean Startups

Speakeasies
Wine Bars
Beer Gardens
Gastropubs
Hipsters

Entrepreneurs

App Developers
Data Scientists
Big Data
Internet of Things

International Development

Open Source

IPAs
Craft Beers

Cabernet Sauvignon
Wine Tasting

Gin

Augmented Reality

Cocktails
Bourbon

Lounges
Beach Volleyball

Soccer

Meditation
Yoga
Spirituality

Poker
NFL
Basketball

Video Games

Foreign Language Learning

Coffee
Lattes
Tea

Farmers Market
Paintball

Mandarin

Cappuccinos
Chocolates

Food Trucks

Soul Food
Milk Shakes
Desserts
Mac & Cheese

Techno Music
Raves

Rap

Hip Hop

16

http://blog.yhat.com/posts/words2map.html

# King - Man + Woman ≈ Queen

# Classify This!

Train {
blah    blah    football    $\rightarrow A$

blah    blah    stocks    $\rightarrow B$
}

Test    blah    blah    hockey    $\rightarrow$ ?
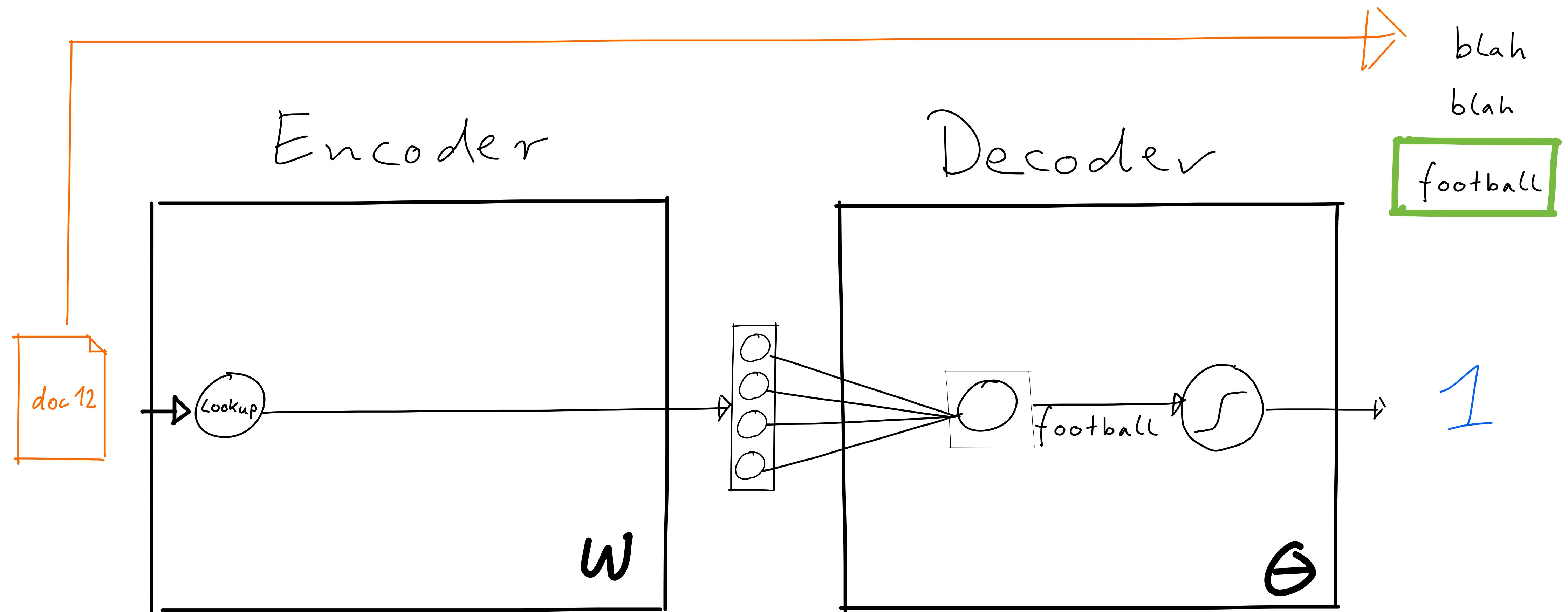
# Cluster This!

blah     blah     football     → ?
.
blah     blah     stocks     → ?
.

Test          blah     blah     hockey     → ?

Why useful ?

# Topic Modelling as Self Supervision



Encoder | Decoder

doc 12

Lookup

W

football

blah blah

football

$\theta$

$1$

$X$  represent the doc such that the "football" classifier says $1$ and the "stocks" classifier $0$  $Y'$

# Matrix View



$\theta_{\bullet, football}$

$2.0$
$-1.0$
$0.0$

$W_{12, \bullet}$

$3.0 \quad 1.5 \quad -1.5$

$W$

$4.5$

$A = W\theta$

$A_{12, football} = W_{12, \bullet} \; \theta_{\bullet, football}$

21

# Matrix View

$\theta_{\bullet, football}$

| | 2.0 | |
| | -1.0 | |
| | 0.0 | |

$\theta$

$W_{12, \bullet}$

| | 3.0 1.5 -1.5 | |
| **W** | | |

| | | |
| | 0.9 | |
| | | $Y = \sigma(W\theta)$ |

$A_{12, football} = W_{12, \bullet} \ \theta_{\bullet, football}$

$Y_{12, football} = \sigma(A_{12, football})$

# Term Document Matrix

$$\theta_{\bullet, football}$$

| | |
|---|---|
| 2.0 | ⊖ |
| -1.0 | |
| 0.0 | |

$W_{12, \bullet}$

| |
|---|
| 3.0  1.5  -1.5 |

**W**

| | | |
|---|---|---|
| | 0 | |
| | 1 | |
| | 0 | |
| | 0 | |
| 0  1  0  1  0  0 | 1 | 0  0  1  0  1  0 |
| | 0 | |
| | 0 | |
| | 1 | |

**Y'**

$$A_{12, football} = W_{12, \bullet} \; \theta_{\bullet, football}$$

$$Y_{12, football} = \sigma\left(A_{12, football}\right)$$

# Binary Matrix Factorization

$$\sigma \left( \begin{array}{c} \boxed{\begin{array}{ccc} 3.0 & 1.5 & -1.5 \end{array}} \\ \mathbf{w} \end{array} \times \boxed{\begin{array}{c} 2.0 \\ -1.0 \\ 0.0 \end{array} \quad \ominus} \right)$$

$$\stackrel{!}{\approx}$$

$$\boxed{\begin{array}{ccccccc|ccccccc} & & & & & & 0 & & & & & & \\ & & & & & & 1 & & & & & & \\ & & & & & & 0 & & & & & & \\ & & & & & & 0 & & & & & & \\ \hline 0 & 1 & 0 & 1 & 0 & 0 & 1 & 0 & 0 & 1 & 0 & 1 & 0 \\ \hline & & & & & & 0 & & & & & & \\ & & & & & & 0 & & & & & & \\ & & & & & & 0 & & & & & & \\ & & & & & & 1 & & & & & & Y' \end{array}}$$

$$A_{12,\,football} = W_{12,\bullet}\,\Theta_{\bullet,\,football}$$

$$Y_{12,\,football} = \sigma\big(A_{12,\,football}\big)$$

$$\stackrel{!}{\approx}\ Y'_{12,\,football}$$

# Matrix Factorization

not
very
interpretable

$W$

3.0  1.5  -1.5    $X$

$\approx$

use counts

2.0
-1.0
0.0

0
1
0
0

0  2  0  1  0  0  3  0  0  2  1  0  4  0

0
0

$Y'$

$A_{12,\,football} = W_{12,\cdot}\,\theta_{football}$

$\approx Y'_{12,\,football}$

# Nonnegative Matrix Factorization



Force normalization for probabilistic Topic Models (e.g. LDA)

$$A_{12, football} = W_{12, \bullet} \, \theta_{\bullet, football}$$

$$\stackrel{!}{\approx} Y'_{12, football}$$

$$\theta \geq 0$$
$$W \geq 0$$

Enforce these constraints during training

# References

- Word Embeddings

  - J&M Chapter 6

  - Goldberg Chapter 5

  - Efficient Estimation of Word Representations in Vector Space, Mikolov et al, ICLR Workshop 2013

  - GloVe: Global Vectors for Word Representation, Pennington et al., EMNLP 2014

- Topic Models

  - An Introduction to Latent Semantic Analysis, TK Landauer

  - Probabilistic Latent Semantic Analysis, T Hofmann

  - Probabilistic Topic Models, Blei

  - Exploring Topic Coherence over Many Models and Many Topics, K Stevens, EMNLP 2012