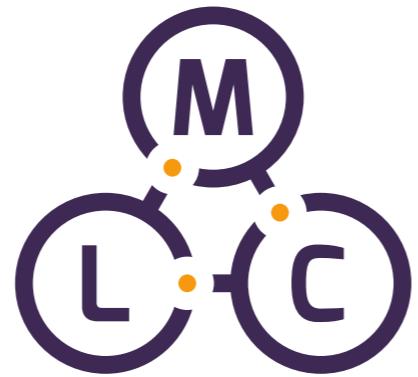


Natural Language Processing

Jiří Materna



Machine
Learning
College

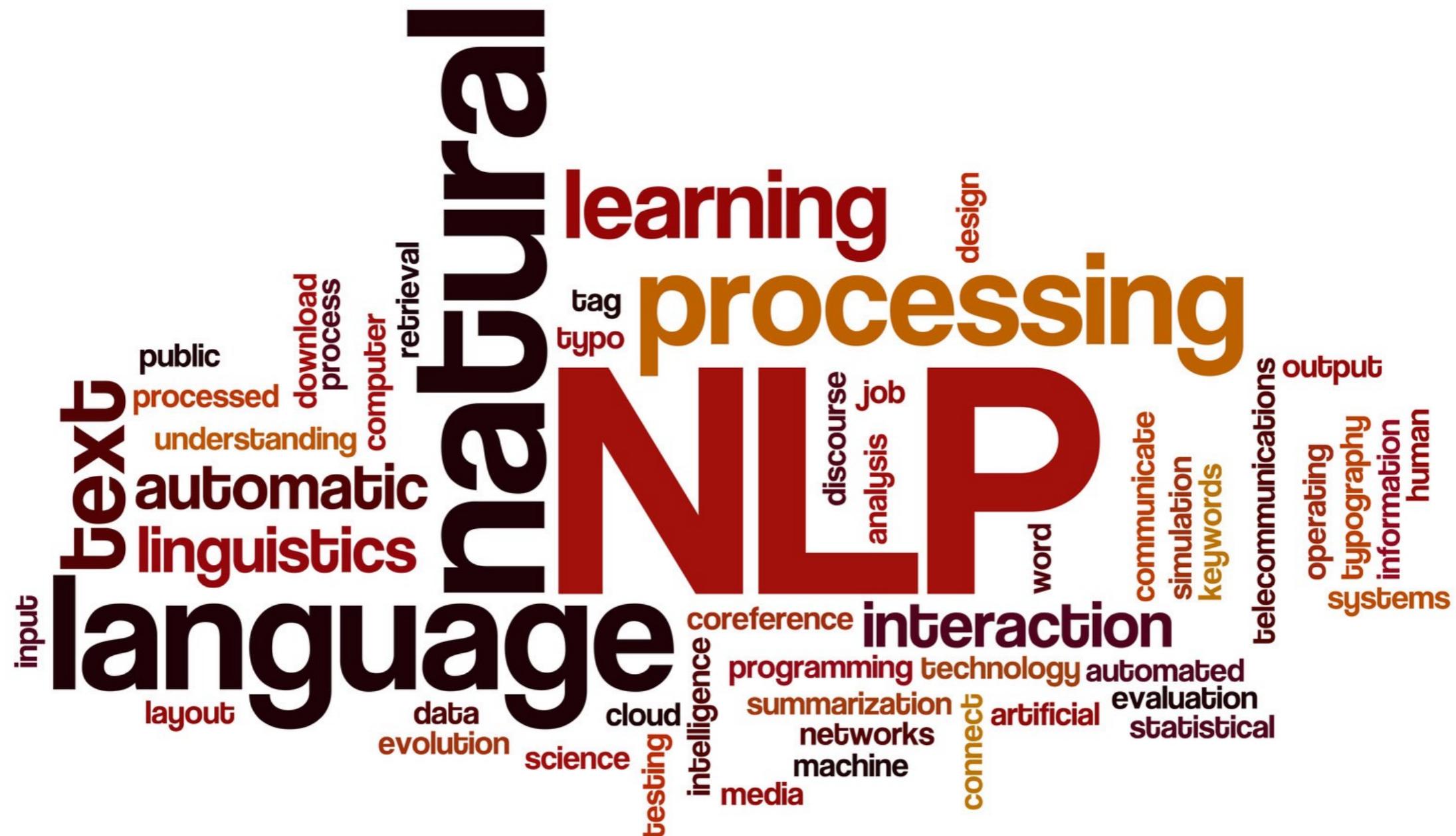
About me

- Ph.D. in Natural Language Processing and Artificial Intelligence at Masaryk University
- 10 years at Seznam.cz (last 8 years as Head Of Research)
- Founder and lecturer at ML College
- Founder and co-organizer of ML Prague
- ML Freelance and consultant

Outline

- Introduction to natural language processing
- Computational linguistics
- Text document vectorization
- Practical document classification task
- Language modeling
- Practical tasks on language modeling
- Word embeddings
- Text generating
- Practical tasks on language modeling
- Transformers

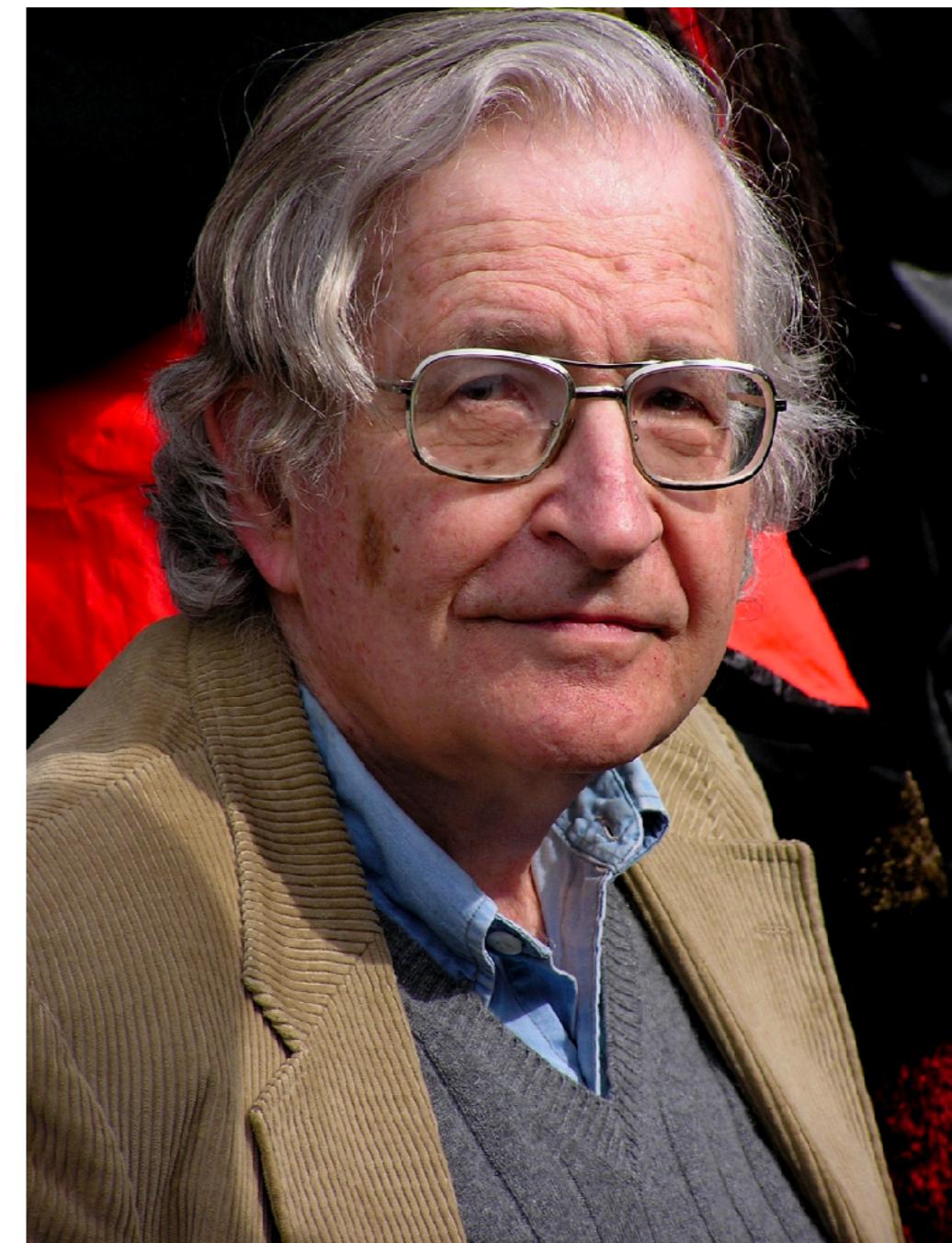
What is Natural Language Processing?



Norvig vs. Chomsky



source: <https://www.commarts.com>



source: <https://citaty.net>

Token & tokenization

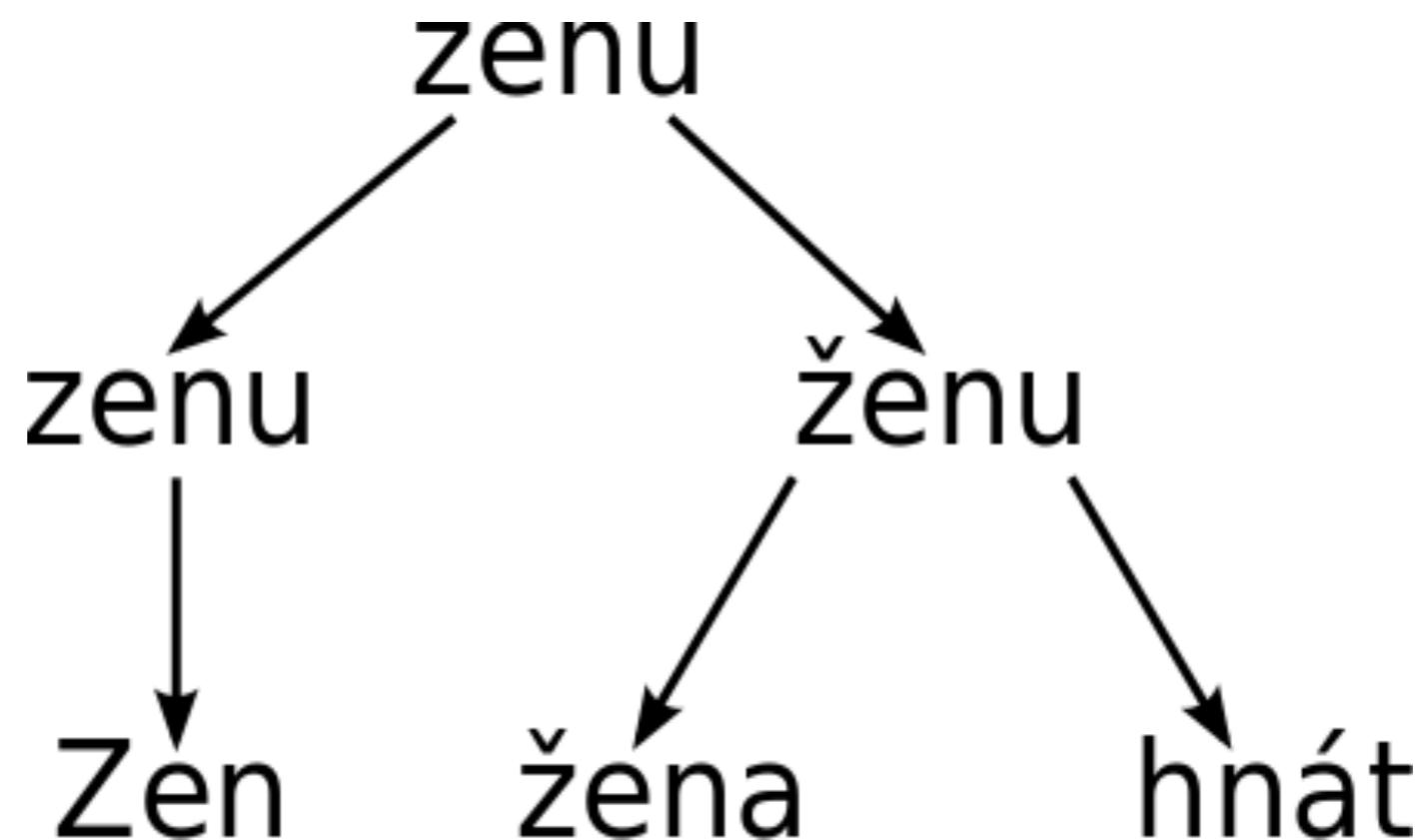
This is a non-trivial English sentence: Ludolph's number is approx. 3.14.

Python library: <http://www.nltk.org/>

Stemming & lemmatization

Original	Stemming	Lemmatization
compensation	compens	compensation
compensations	compens	compensation
mouse	mous	mouse
mice	mice	mouse

Ambiguity in lemmatization



Stemming & lemmatization

English:

<https://tartarus.org/martin/PorterStemmer/>

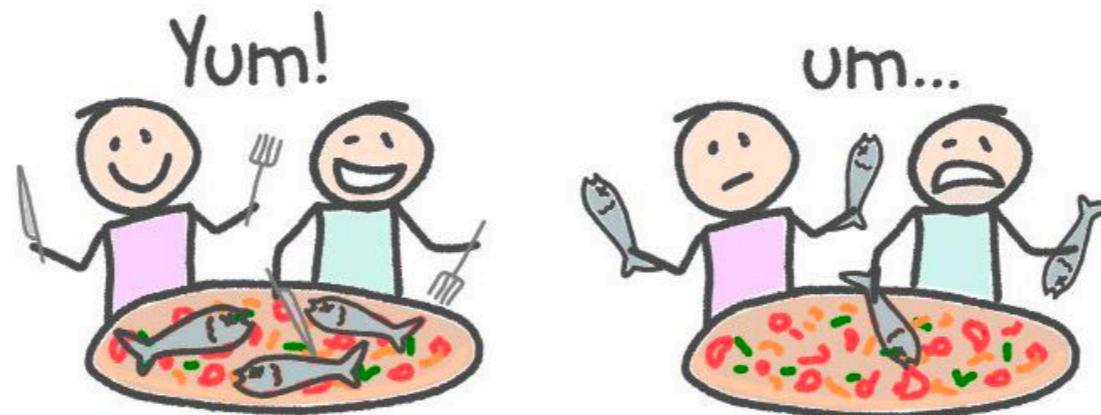
<http://www.cis.uni-muenchen.de/~schmid/tools/TreeTagger/>

Czech:

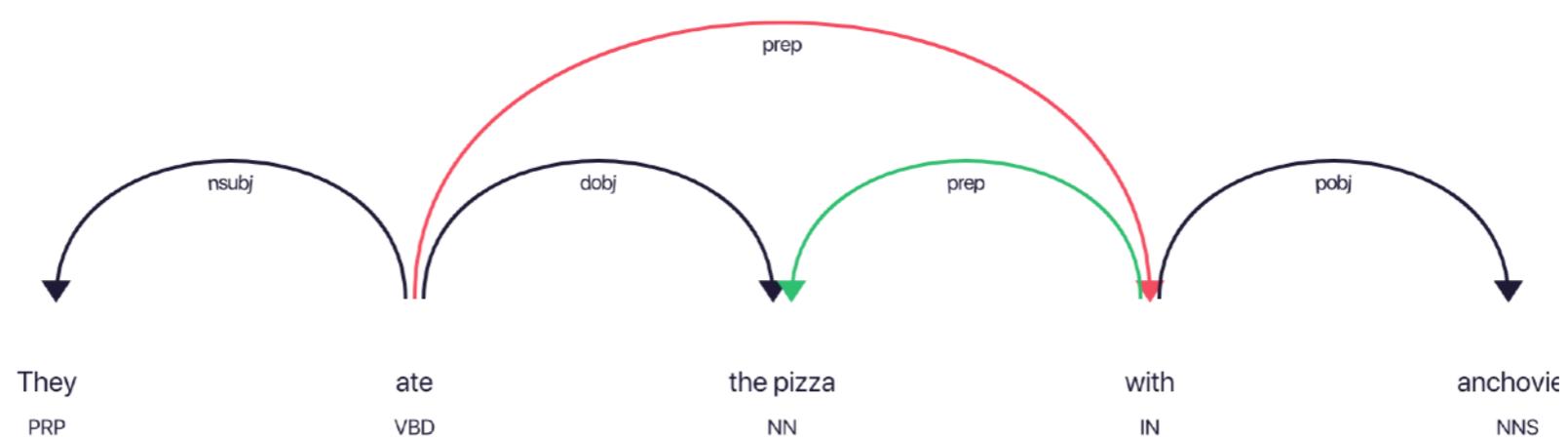
<http://ufal.mff.cuni.cz/morphodita>

Parsing

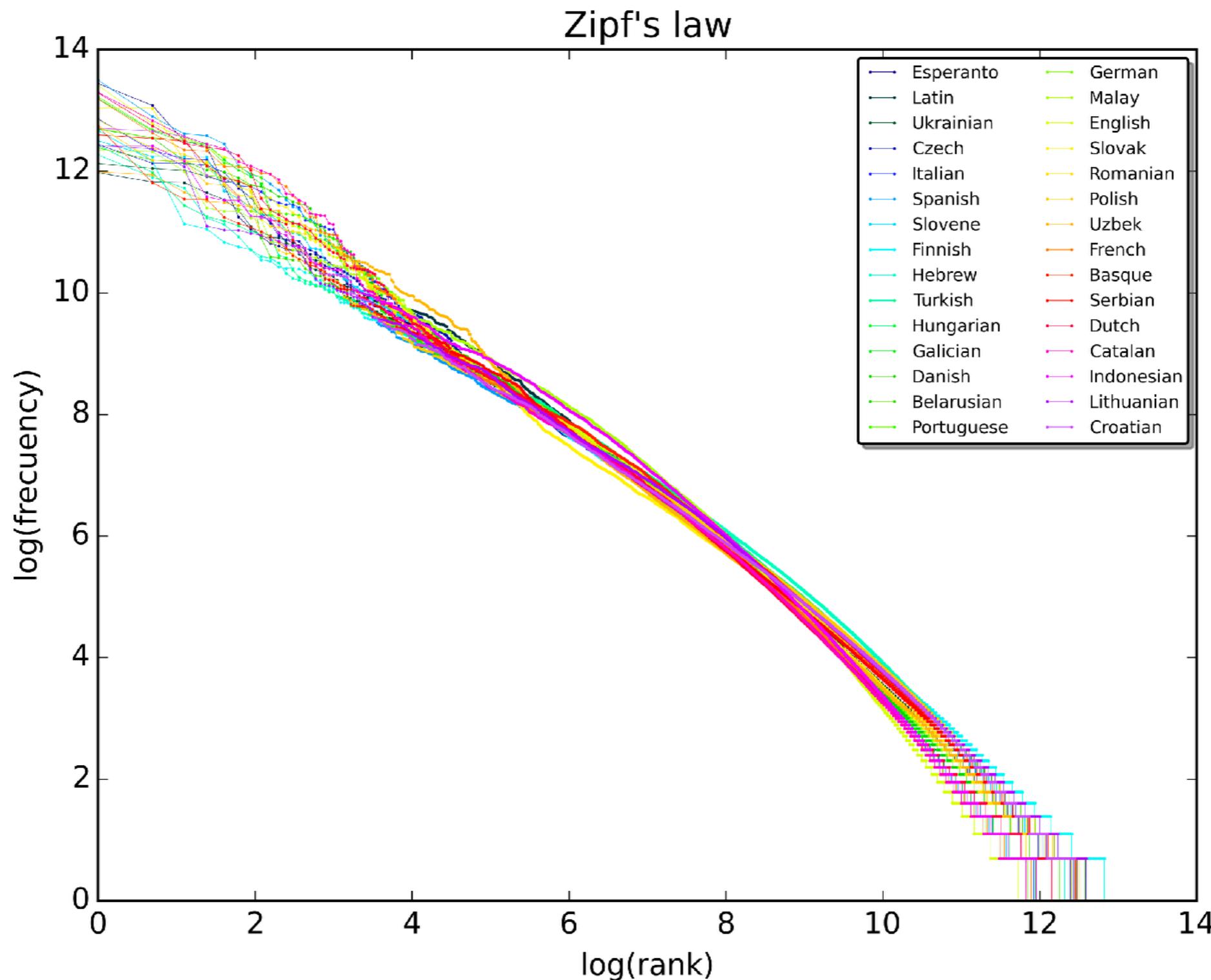
They ate the pizza with anchovies



Creative Commons Attribution-NonCommercial 2.5
James Constable, 2010



Zipf's law & long tail



Publicly available corpora

British National Corpus: <http://www.natcorp.ox.ac.uk/>

Common Crawl: <http://commoncrawl.org/the-data/get-started/>

Wikipedia: <https://dumps.wikimedia.org/>

Feature extraction for NLP

1. *the man walked the dog*
2. *the man took the dog to the park*
3. *the dog went to the park*

[dog, man, park, the, to, took, walked, went]

1. [1, 1, 0, 1, 0, 0, 1, 0]
2. [1, 1, 1, 1, 1, 1, 0, 0]
3. [1, 0, 1, 1, 1, 0, 0, 1]

$$\text{idf}(t, D) = \log \frac{N}{|\{d \in D : t \in d\}|}$$

1. [1, 1, 0, 2, 0, 0, 1, 0]
2. [1, 1, 1, 3, 1, 1, 0, 0]
3. [1, 0, 1, 2, 1, 0, 0, 1]

1. [0, 0.18, 0, 0, 0, 0, 0.48, 0]
2. [0, 0.18, 0.18, 0, 0.18, 0.48, 0, 0]
3. [0, 0, 0.18, 0, 0.18, 0, 0, 0.48]

— . . .

NLP Introduction task

01-text-classification-introduction.ipynb

Language models

- spell checking
- speech recognition
- machine translation
- text generation
- ...

n-gram models

$$P(w_1, w_2, \dots, w_n) = P(w_1)P(w_2|w_1)\dots P(w_n|w_1, \dots, w_{n-1})$$

$$= \prod_i P(w_i|w_1, w_2 \dots w_{i-1})$$

$$\approx \prod_i P(w_i|w_{i-k}, w_{i-k-1} \dots w_{i-1})$$

$$P(w_i|w_{i-k}, w_{i-k-1} \dots w_{i-1}) = \frac{\text{count}(w_{i-k}, w_{i-k-1} \dots w_{i-1}, w_i)}{\text{count}(w_{i-k}, w_{i-k-1} \dots w_{i-1})}$$

Language model smoothing

- Laplace smoothing (plus one)

$$P(w_i|w_{i-1}) = \frac{\text{count}(w_{i-1}, w_i) + 1}{\text{count}(w_{i-1}) + V}$$

- interpolation
- Good-Turing
- Witten-Bell
- ...

Perplexity

$$PP(W) = P(w_1, w_2, \dots, w_N)^{-\frac{1}{N}}$$

$$= \sqrt[N]{\frac{1}{P(w_1, w_2, \dots, w_N)}}$$

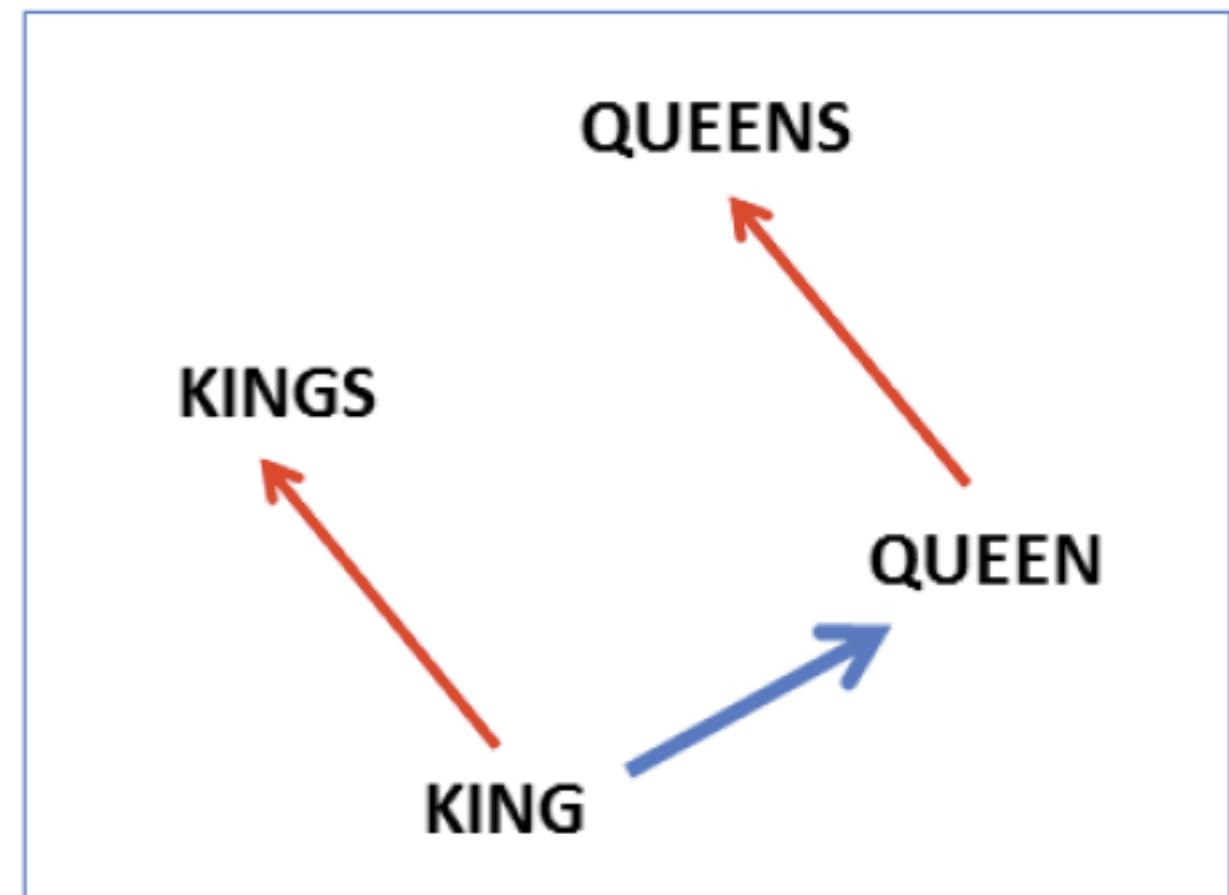
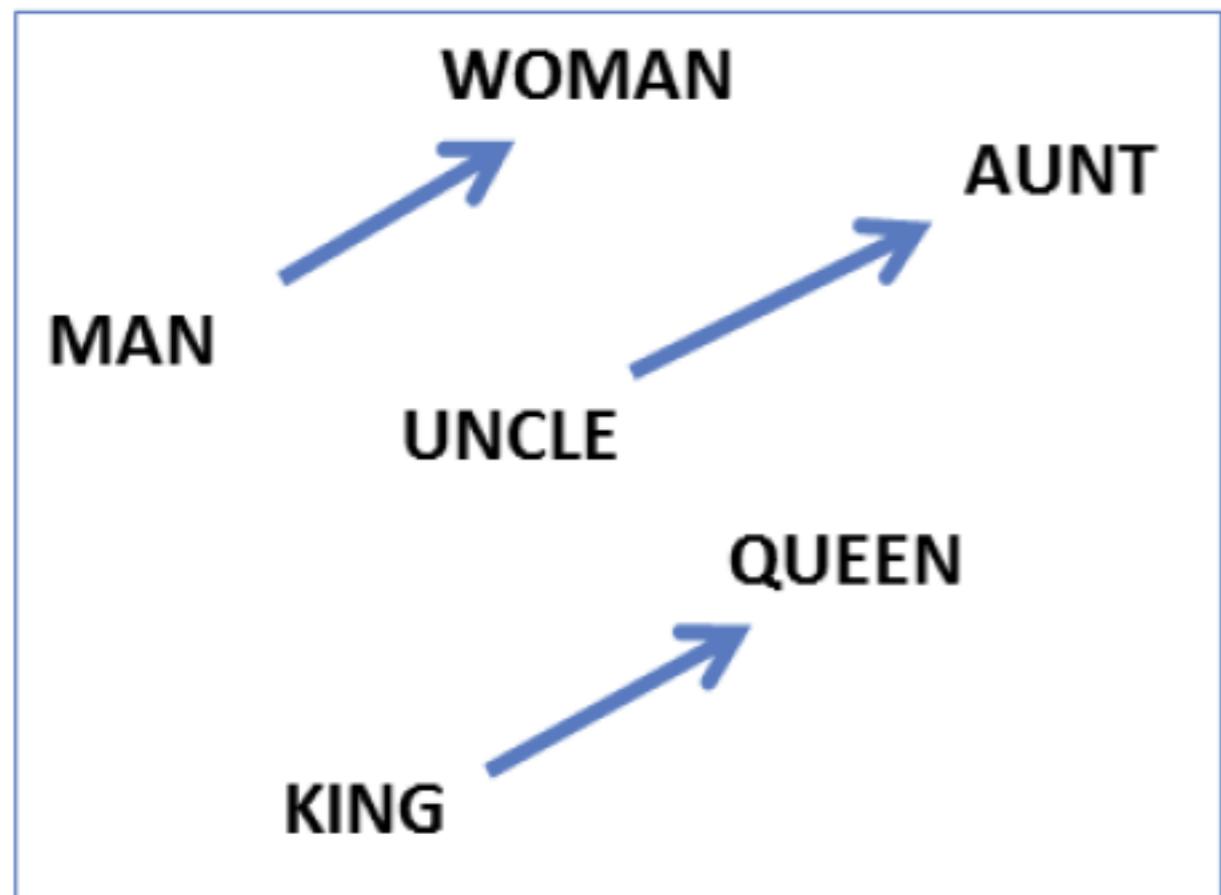
$$= \sqrt[N]{\prod_{i=1}^N \frac{1}{P(w_i | w_1, \dots, w_{i-1})}}$$

$$= 2^{-\frac{1}{N} \sum_{i=1}^N \log_2 P(w_i | w_1, \dots, w_{i-1})}$$

Language detection using language models

02-Language-detection-assignment.ipynb

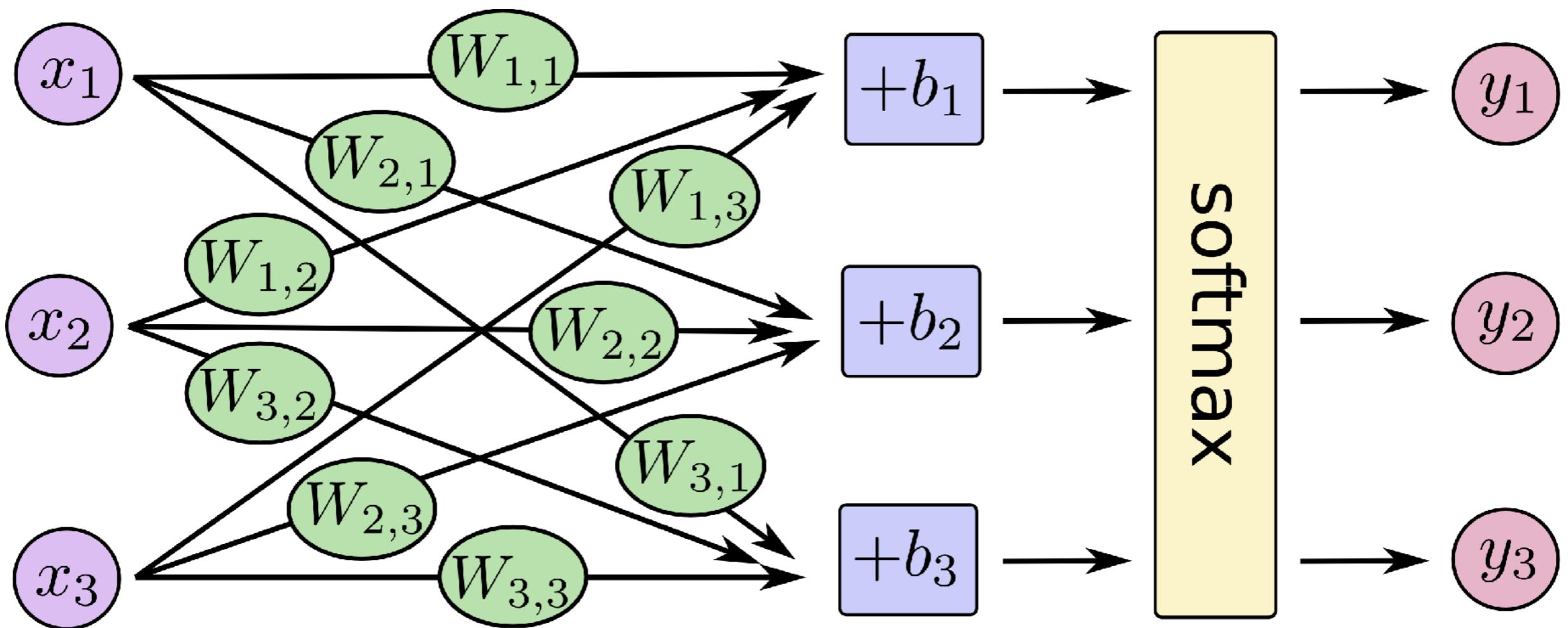
word2vec



king is to **kings** as **queen** to ?.

$$v(\mathbf{kings}) - v(\mathbf{king}) = v(\mathbf{queens}) - v(\mathbf{queen})$$

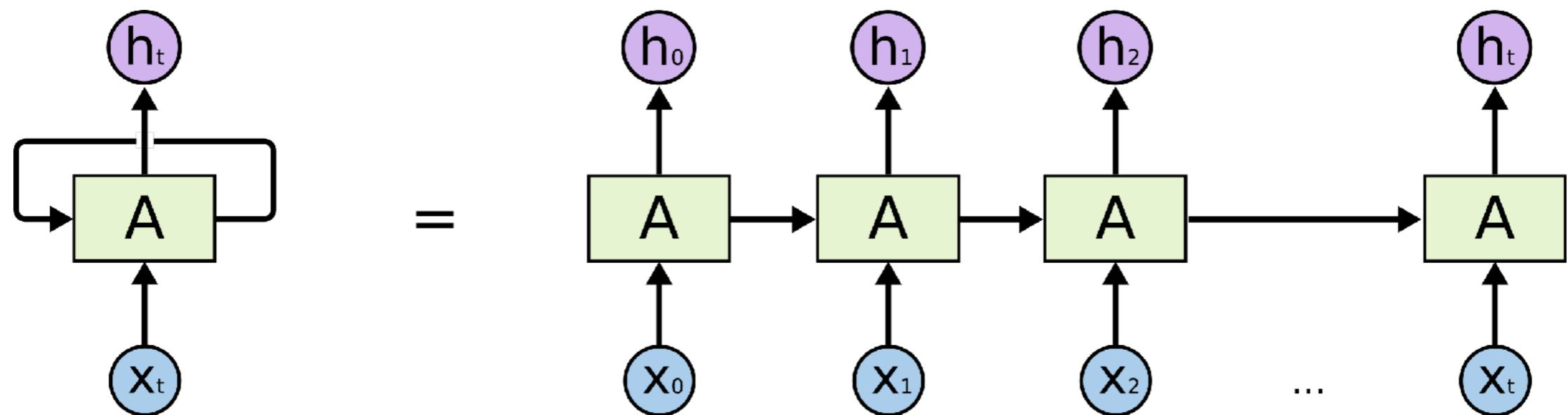
Feed-Forward Neural Network



source: <https://www.tensorflow.org>

Recurrent Neural networks

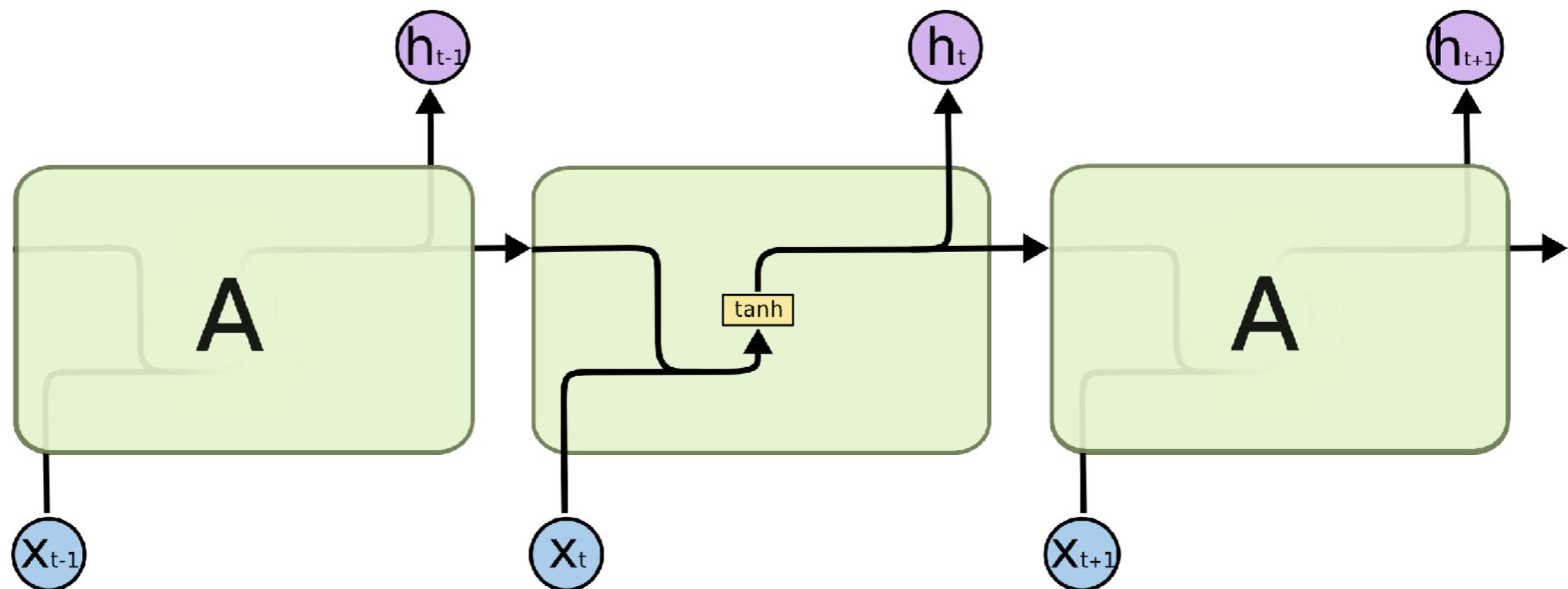
1/2



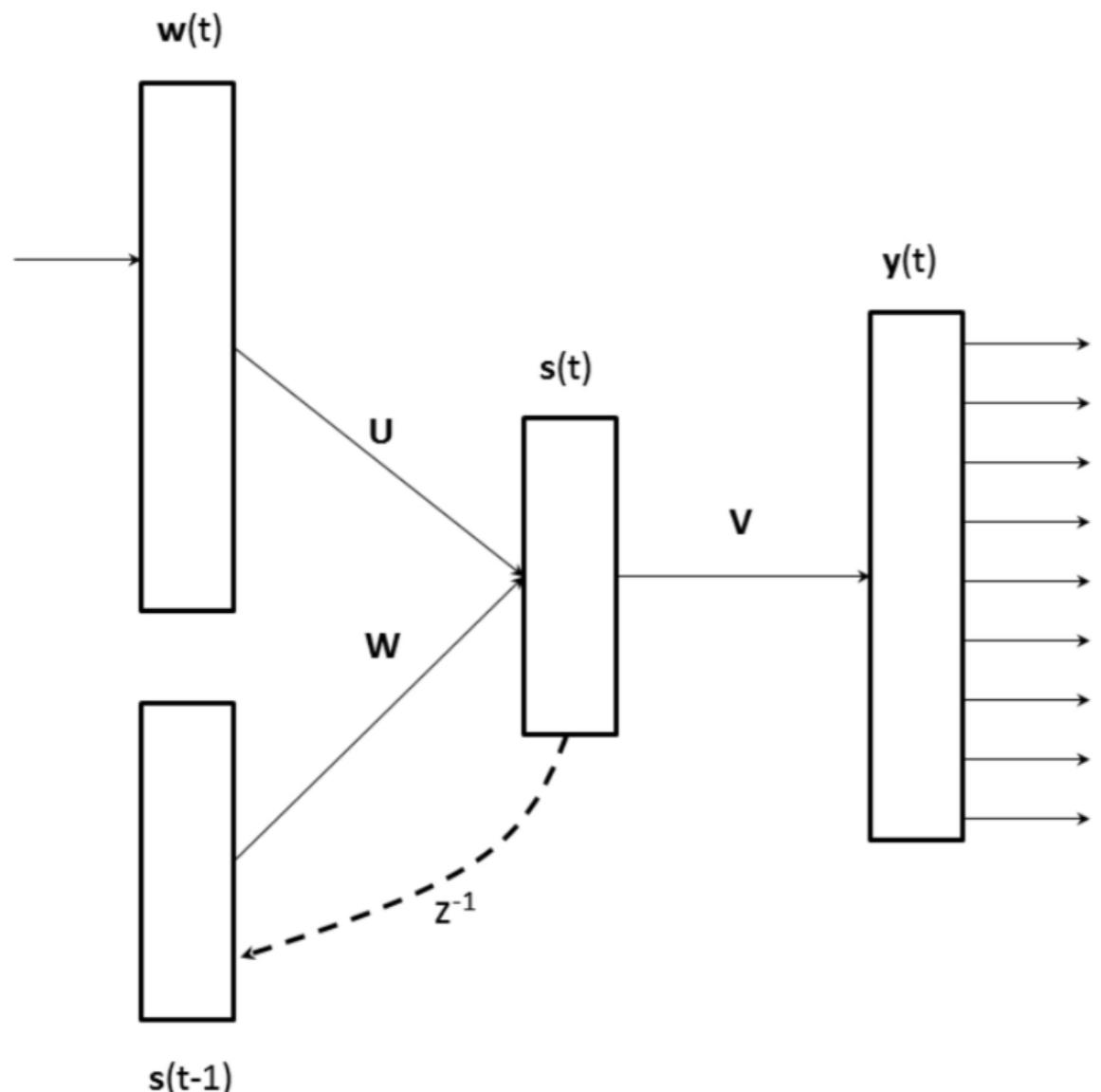
source: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Recurrent Neural Networks

2/2



Recurrent Neural Network Language Modeling Toolkit

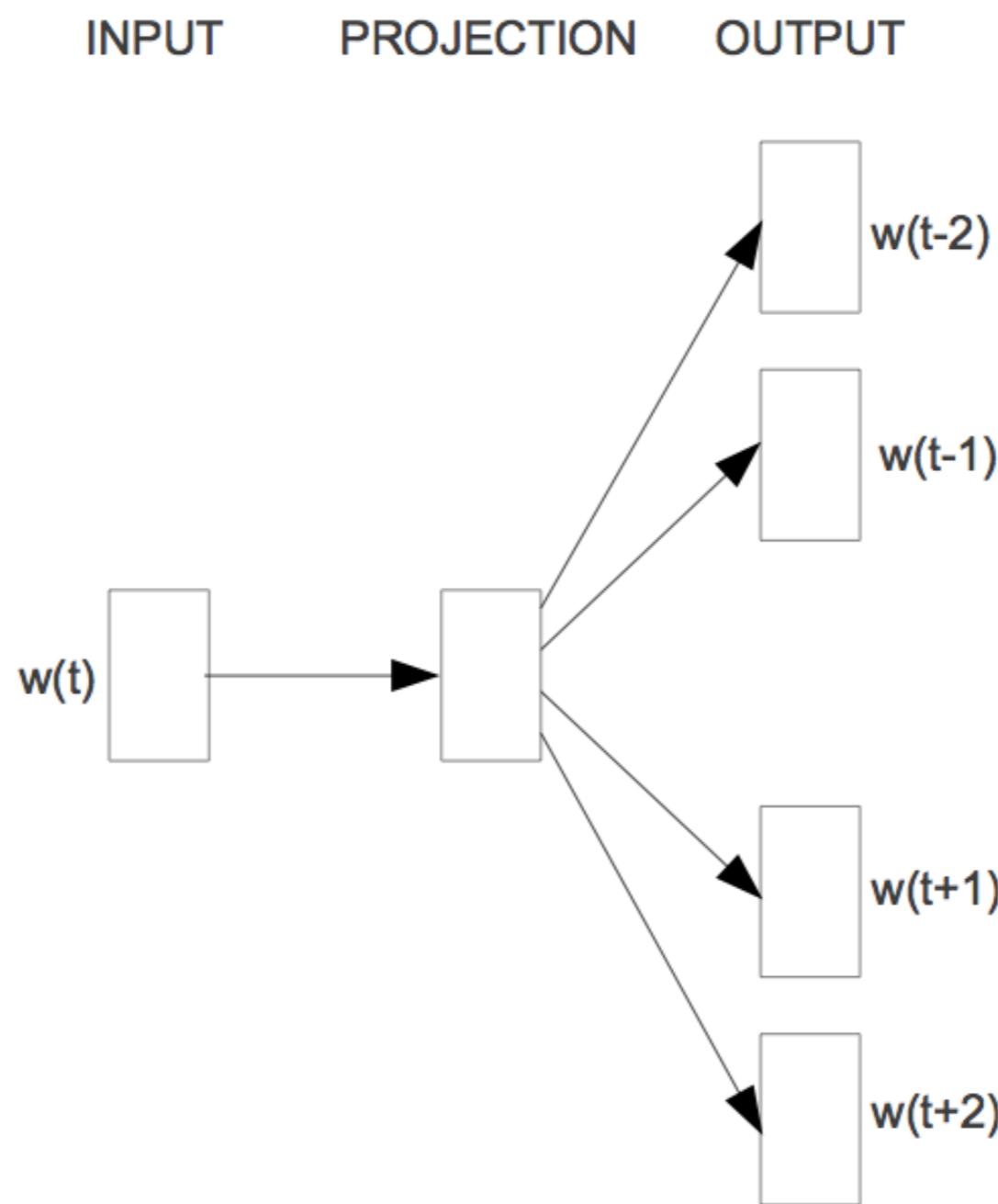


$$\mathbf{s}(t) = f(\mathbf{U}\mathbf{w}(t) + \mathbf{W}\mathbf{s}(t-1))$$

$$\mathbf{y}(t) = g(\mathbf{V}\mathbf{s}(t)),$$

$$f(z) = \frac{1}{1 + e^{-z}}, \quad g(z_m) = \frac{e^{z_m}}{\sum_k e^{z_k}}.$$

The skip-gram model



Experiments with word2vec

04-Word2vec-in-gensim.ipynb

05-Review-classification-w2v-assignment.ipynb

Language models for text generating

Nacházíte se: Úvod > Oddělení > Krásná literatura > Poezie > Česká a slovenská poezie > Elektronická kniha Poezie umělého světa



Poezie umělého světa [E-kniha]

Jiří Materna



Hodnotilo 7 uživatelů, zatím žádné recenze, [napsat vlastní recenzi](#)

Popis: Elektronická kniha, 50 stran, bez zabezpečení DRM,  ePUB,  Mobi,  PDF, česky - [více](#)



Stáhnout



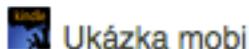
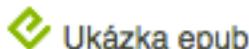
Zdarma

K dispozici pro **okamžité** stáhnutí

Ke stažení

Anotace

Všechny básně v této knize byly automaticky vygenerovány počítačem za pomocí umělých neuronových sítích. Neuronová síť sama o sobě nic neumí a je třeba ji natrénovat pro činnost, kterou má vykonávat.



LISTOPAD

usínám, pláču, umírám, přemýšlím
co cítíš ty?
cítím tvou slabost
a whisky

NOVEMBER

I am falling asleep, crying, dying, thinking
what do you feel?
I feel your weakness
and whisky

SPRAVEDLNOST

na tvou dekadentní duši
ráno i v poledne
bůh má připravenou kuší

JUSTICE

for your decadent soul
in the morning, in the evening
the god has prepared a crossbow

Metaphores

...tělo plné červánků...

...body full of blush of dawn...

...tak vzácný jako listí...

...as rare as leaves of trees...

Language models for text generating

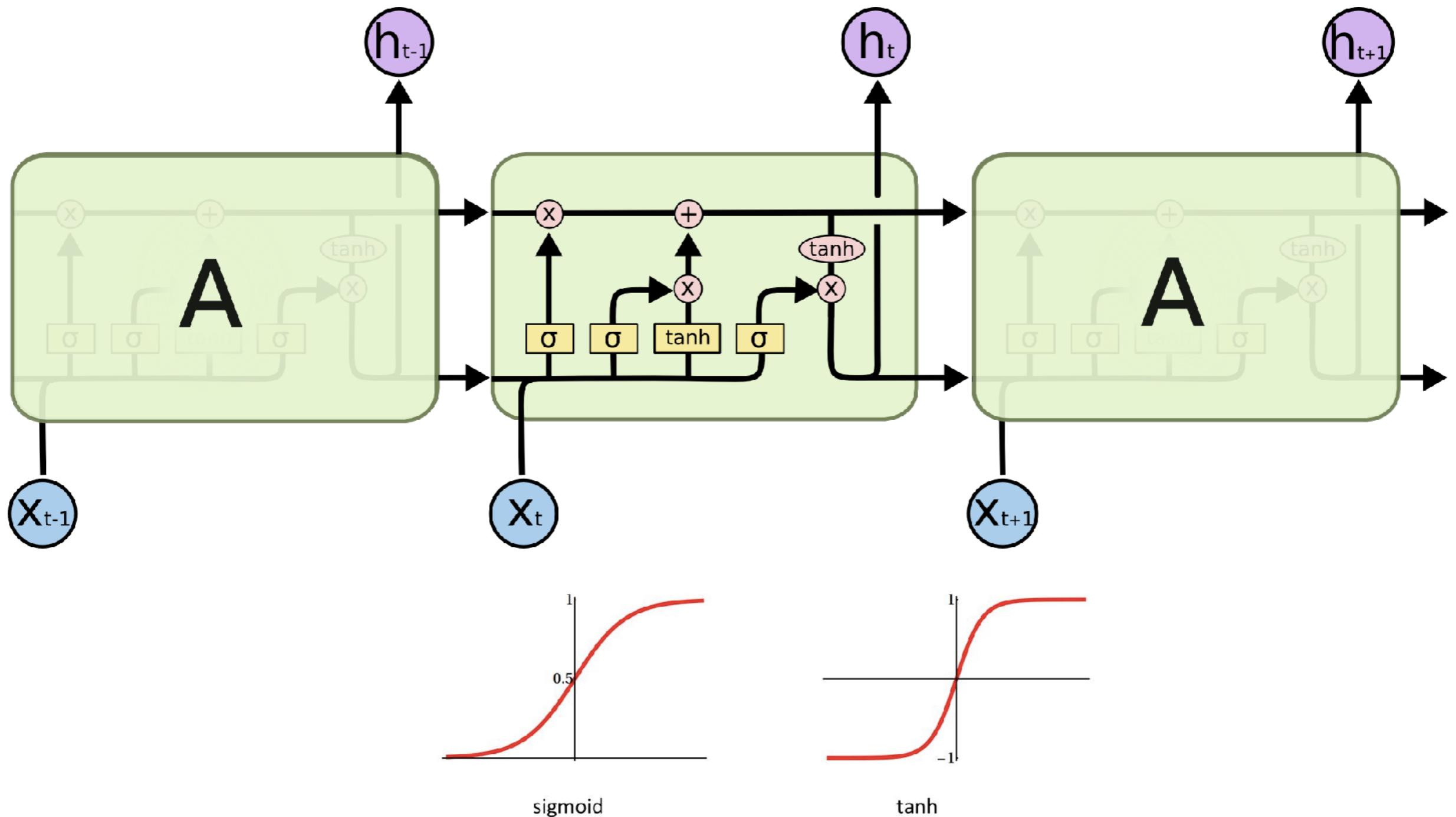
$$P(\text{maso} | \text{máma}, \text{mele}) = 0.5$$

$$P(\text{Emu} | \text{máma}, \text{mele}) = 0.3$$

$$P(\text{tátu} | \text{máma}, \text{mele}) = 0.2$$

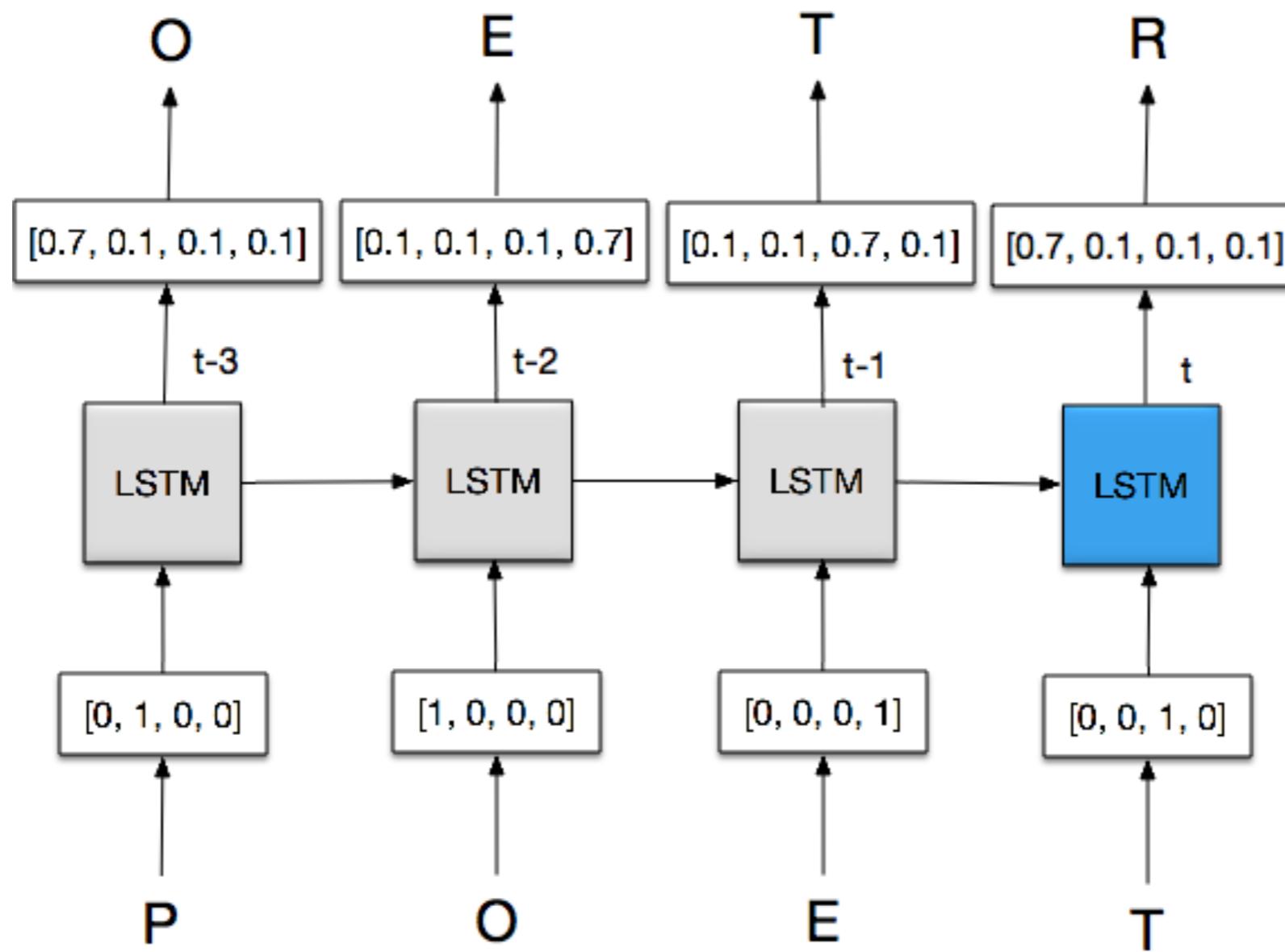
```
t ~ Uniform(0, 1)
s = 0
for v in Vocabulary:
    s += v.prob
    if t < s:
        return v.word
```

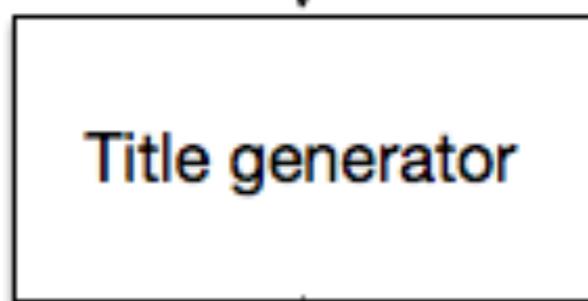
Long Short-Term Memory



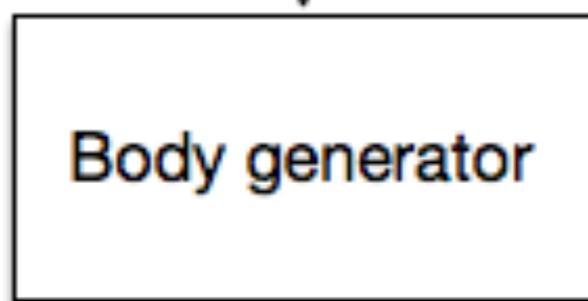
Zdroj: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

LSTM language model





AUTUMN SONG



why don't you kill yourself?
a phone call isn't hope
this planet is still your home
your time is still going on

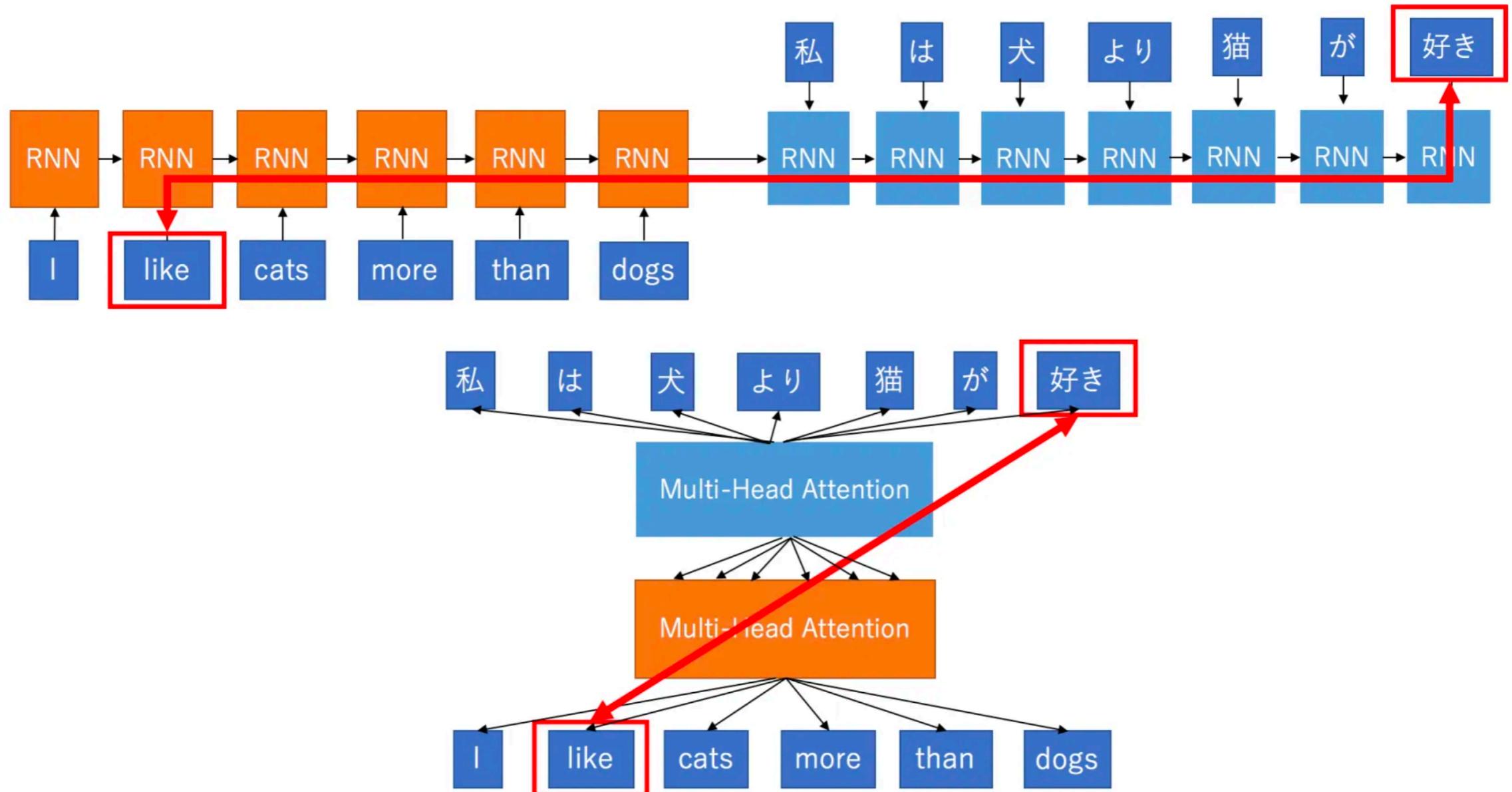
na na na...

I know, I'm revealing a book of dreams
I'll find out I'm nothing more than this

LSTM review generator

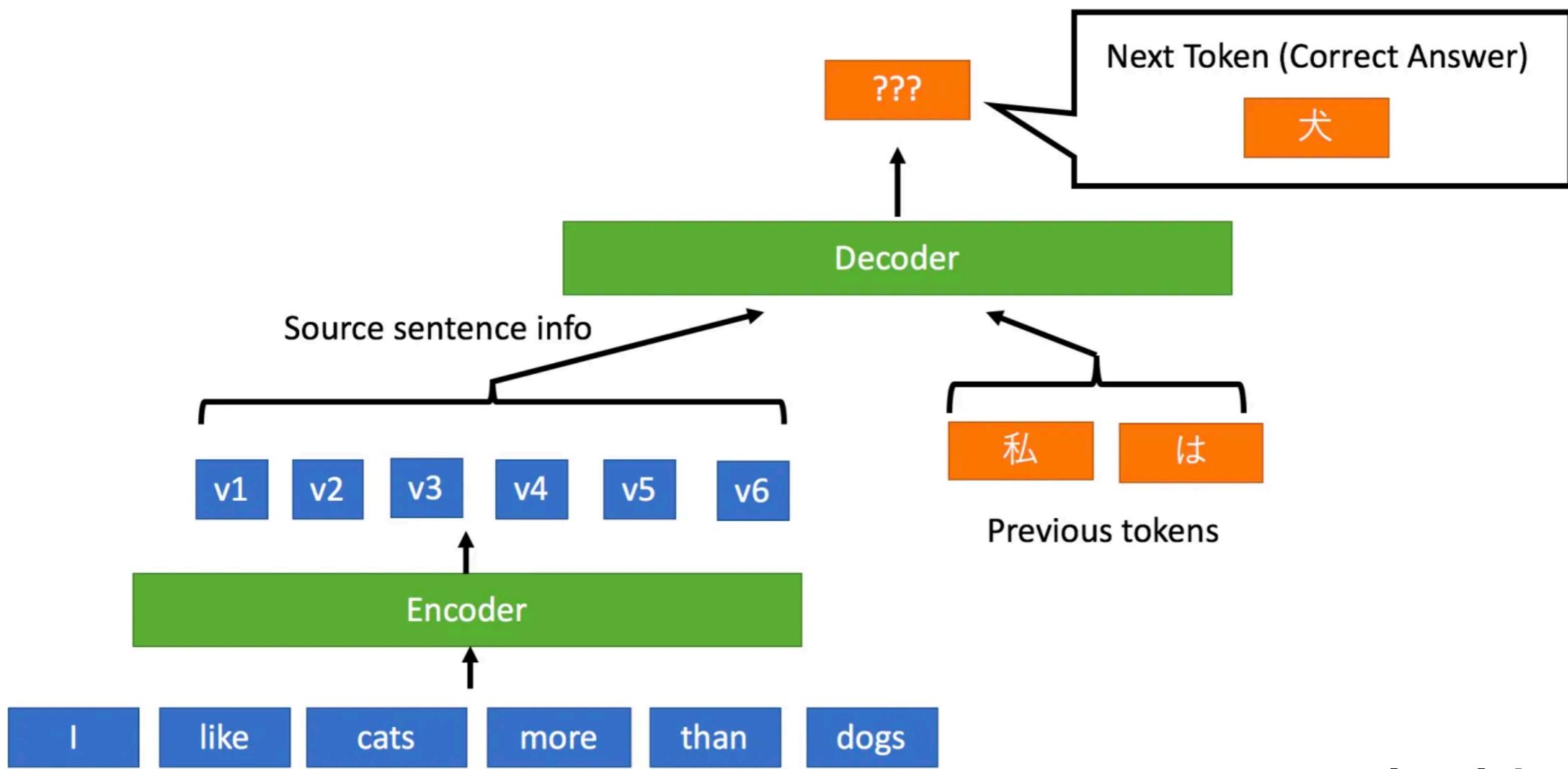
06-Review-generator.ipynb

Transformer



source: www.mlexplained.com

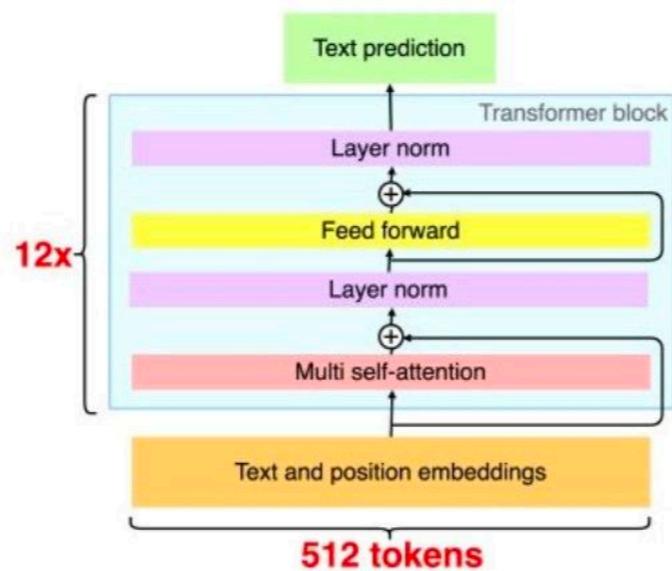
Translation with Transformers



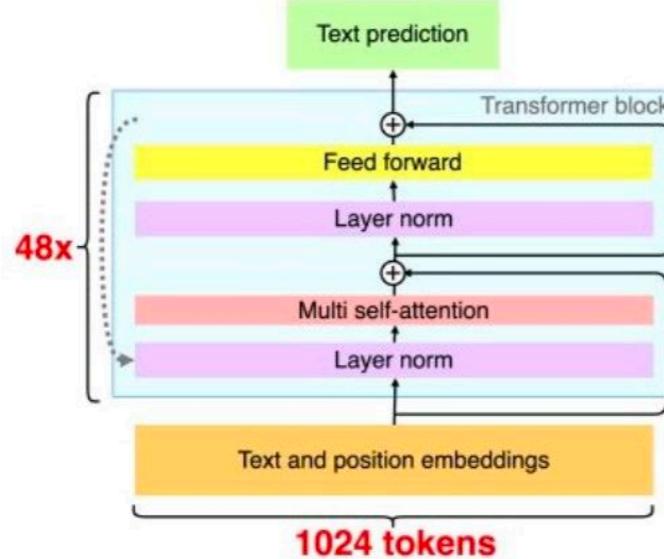
GPT Evolution

| GPT-1 vs GPT-2 vs GPT-3

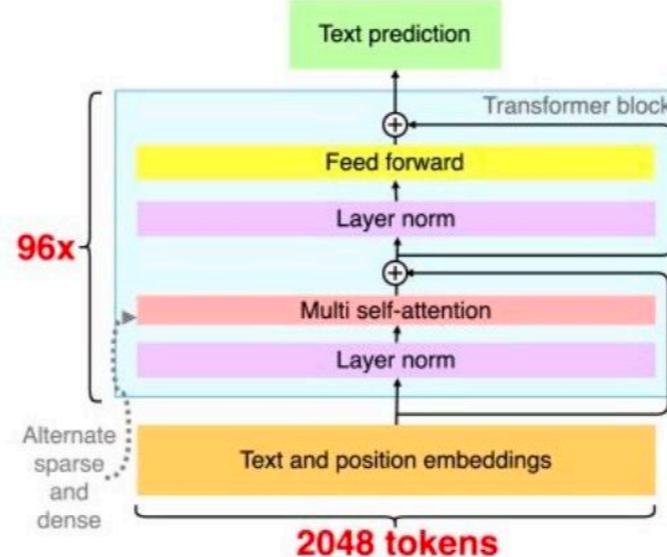
GPT-1



GPT-2



GPT-3



GPT-4

?

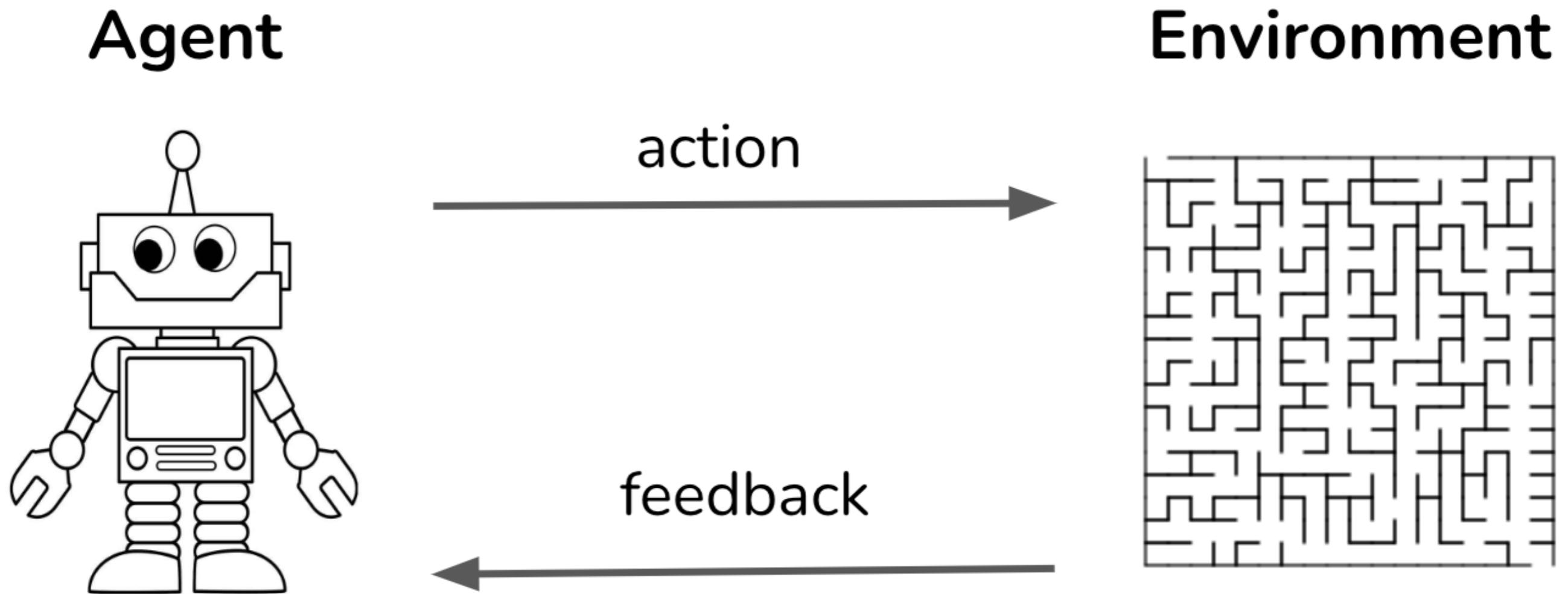
32k tokens

GPT-3 language model

- 499 billions of training tokens
- 179 billions of trainable parameters
- 355 GPU-years of training time
- \$4.6 M estimated training cost

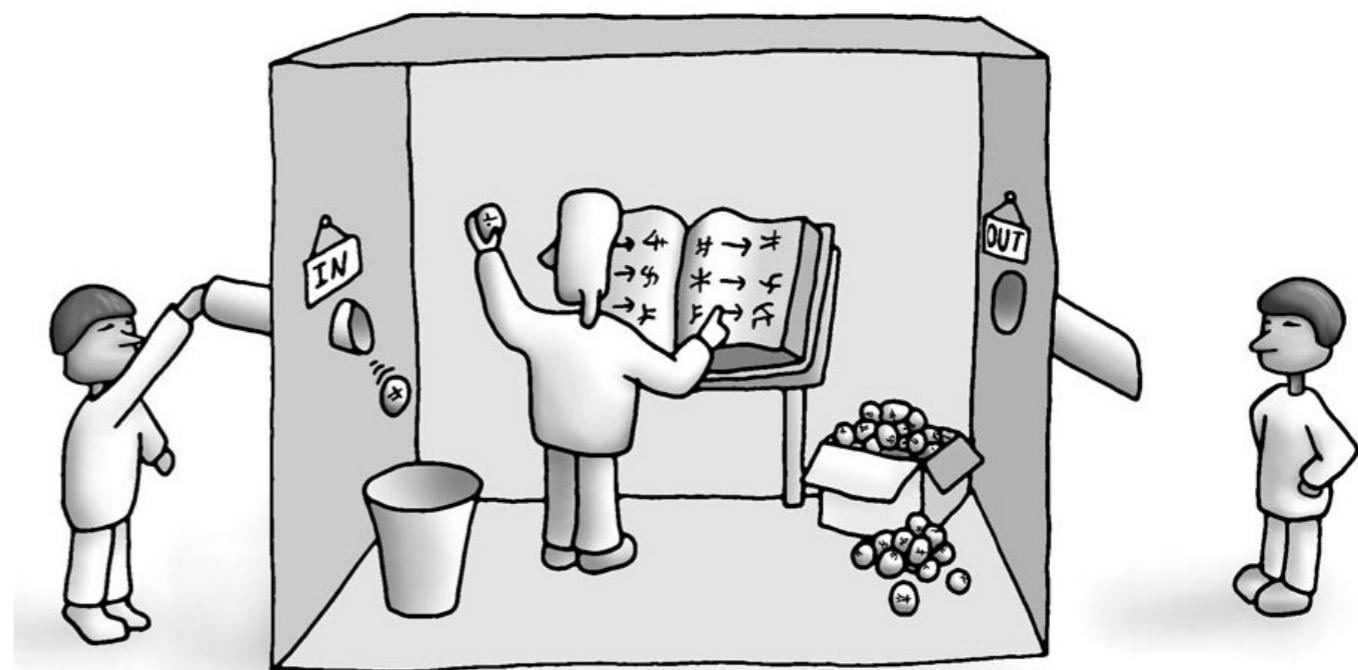
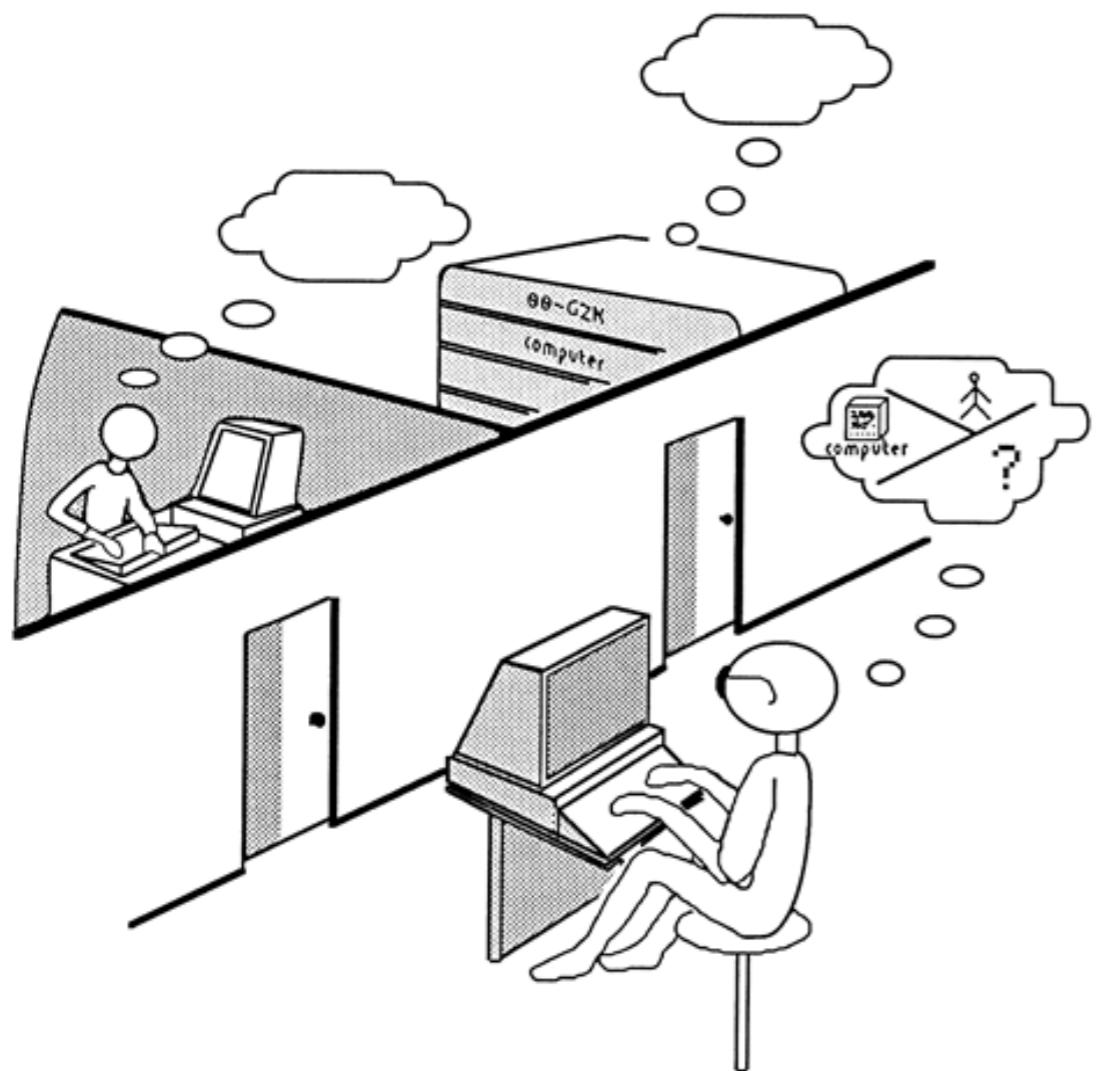
API: <https://platform.openai.com/playground>

From GPT to ChatGPT



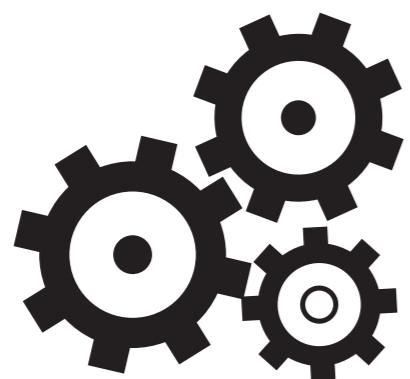
API: <https://chat.openai.com>

Turing Test and Chinese Room Argument



What next?

<https://www.mlcollege.com/en/#courses>



Machine Learning Prague

ML MACHINE LEARNING
MU meetups

Thank you for your attention

e-mail: jiri@mlcollege.com

Web: www.mlcollege.com

Twitter: @JiriMaterna

Facebook: <https://www.facebook.com/maternajiri>

LinkedIn: <https://www.linkedin.com/in/jirimaterna/>