# Machine learning for sequences
# Text and time series
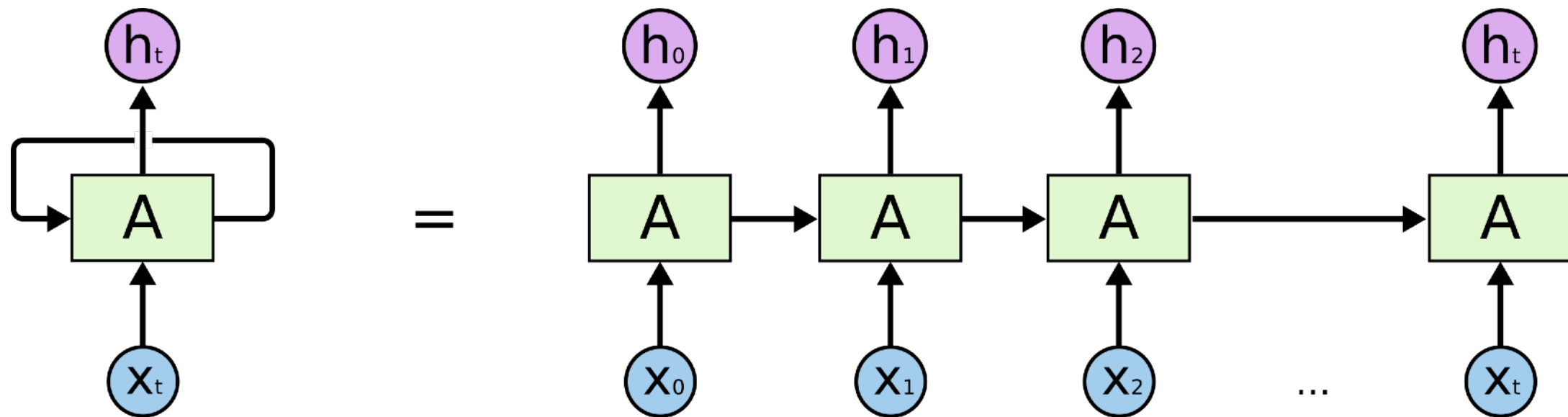
Romain Tavenard (Université de Rennes 2)

# Machine learning for structured data (continued)

- Two options

  1. Cast the data to tabular

     - Representation based on global features (*eg.* bag of words for text or images)

  2. Use structural information in the model

     - images: 2d convolutions

     - sequences: recurrent models, 1d convolutions, temporal kernels
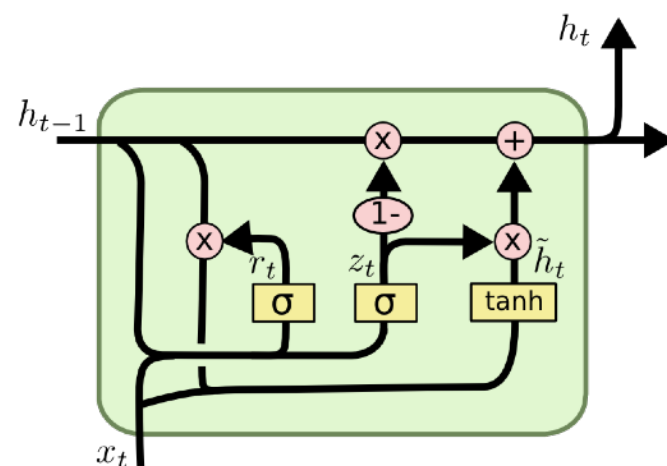
# NN with recurrent units



Source: [Christopher Olah's blog](#)

# NN with recurrent units

- Variants that work well in practice

  - Long Short Term Memory (LSTM)

  - Gated Recurrent Unit (GRU)

- Principle

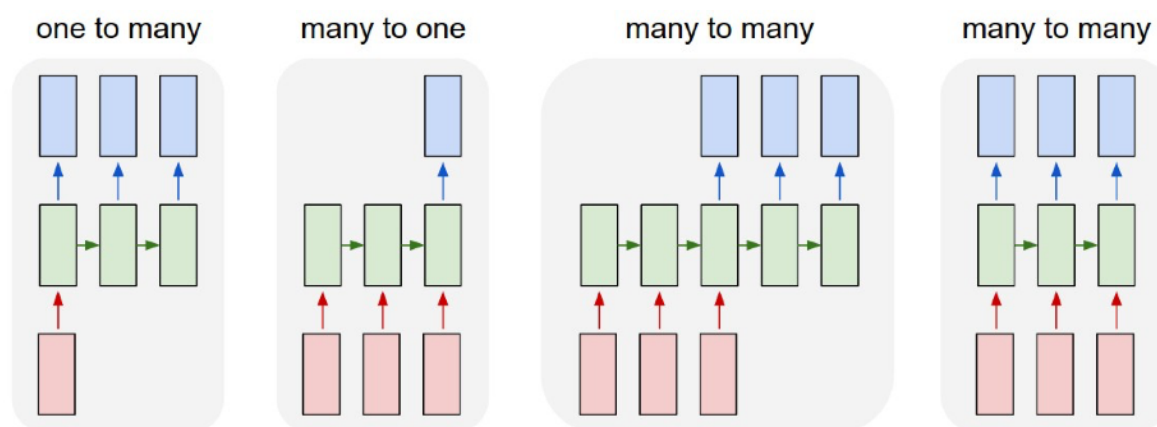  - At each time step, keep only part of the information



$$z_t = \sigma \left( W_z \cdot [h_{t-1}, x_t] \right)$$

$$r_t = \sigma \left( W_r \cdot [h_{t-1}, x_t] \right)$$

$$\tilde{h}_t = \tanh \left( W \cdot [r_t * h_{t-1}, x_t] \right)$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

Illustration: GRU cell, source: <u>Christopher Olah's blog</u>

# NN with recurrent units



one to many    many to one    many to many    many to many

```
PANDARUS:
Alas, I think he shall be come approached and the day
When little srain would be attain'd into being never fed,
And who is but a chain and subjects of his death,
I should not sleep.

Second Senator:
They are away this miseries, produced upon my soul,
Breaking and strongly should be buried, when I perish
The earth and thoughts of many states.

DUKE VINCENTIO:
Well, your wit is in the care of side and that.

Second Lord:
They would be ruled after this chamber, and
my fair nues begun out of the fact, to be conveyed,
Whose noble souls I'll have the heart of the wars.
```

Sample text generated by a RNN
trained on Shakespeare words



Sample LaTeX generated by a RNN
trained on a book of algebraic geometry

Source: Andrej Karpathy's blog, http://karpathy.github.io/2015/05/21/rnn-effectiveness/

# NN with recurrent units

- Efficient at modelling sequential dependencies

- Long-term dependencies: use LSTM or GRU

- Recent alternative (not covered here): Transformer modules

  - Assign importance weights to all items in a sequence

# Dealing with text in practice

- Raw text data is challenging to handle

  - typos

  - what is a term?

  - lots of variants for a term

    - verb conjugation

    - plural form

    - *etc.*

  - synonyms

# Dealing with text in practice

- Raw text data is challenging to handle

  - typos preprocessing

  - what is a term? tokenization

  - lots of variants for a term stemming

    - verb conjugation

    - plural form

    - *etc.*

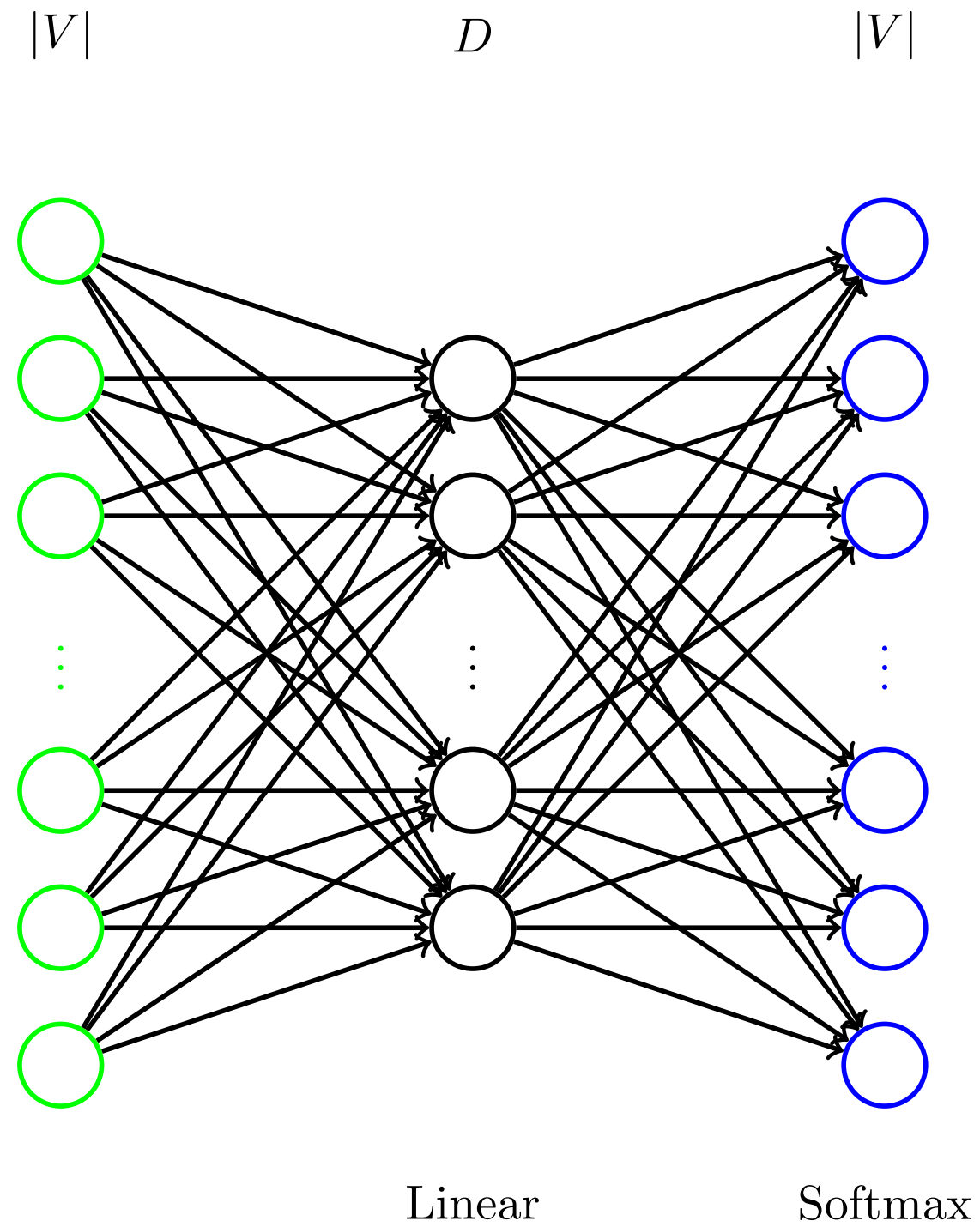  - synonyms word embeddings

# Word embeddings

- Basic idea
  - 1 term = 1 point in multidimensional space
  - Goal: define a space such that similar terms are close

- Reference embedding
  - word2vec

# word2vec

$|V|$  $D$  $|V|$



Linear  Softmax

# word2vec: Continuous Bag of Words (CBOW)

Le **chien mange un os dans sa gamelle**.

- **Input**

  - A bag-of-word representation (binary encoding) of the target word's neighborhood

- Classification task : predict the **target** (middle word)

- Generating a training sample

  1. Sample a word at random in a text
  2. Provide its fixed-length neighborhood

- Why CBOW ?

  - Hidden layer is a Continuous representation of the input Bag of Word

# word2vec: skip-gram

Le **chien** mange un **os** dans sa gamelle.

- **Input**
  - A word
- Classification task: predict **a neighborhood word**

- Generating a training sample
  1. Sample a word at random in a text
  2. Sample a word from its neighbourhood at random

- Why skip-gram ?
  - Associate word pairs (like in bi-gram)
  - Allow skips
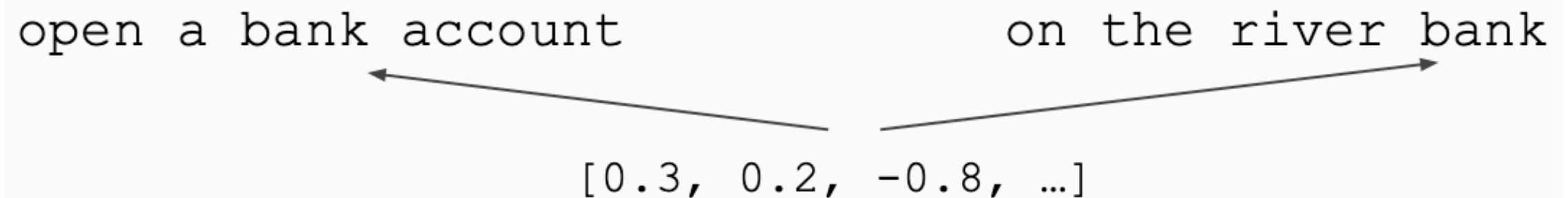
# More about word2vec

- According to authors
  - Skip-gram
    - works well with small amount of the training data
    - represents well even rare words or phrases
  - CBOW
    - several times faster to train than the skip-gram
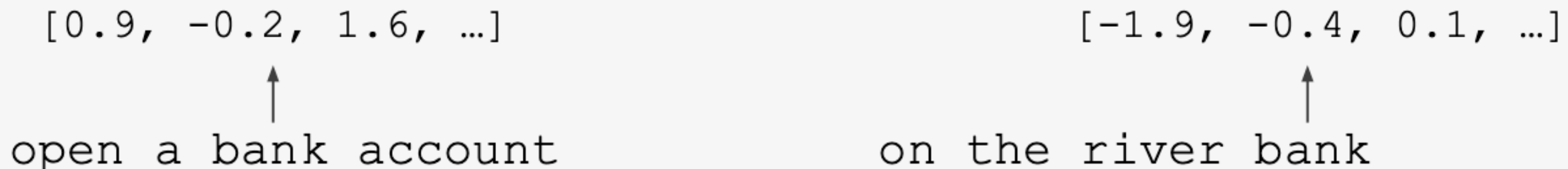    - slightly better accuracy for the frequent words

# Limitations
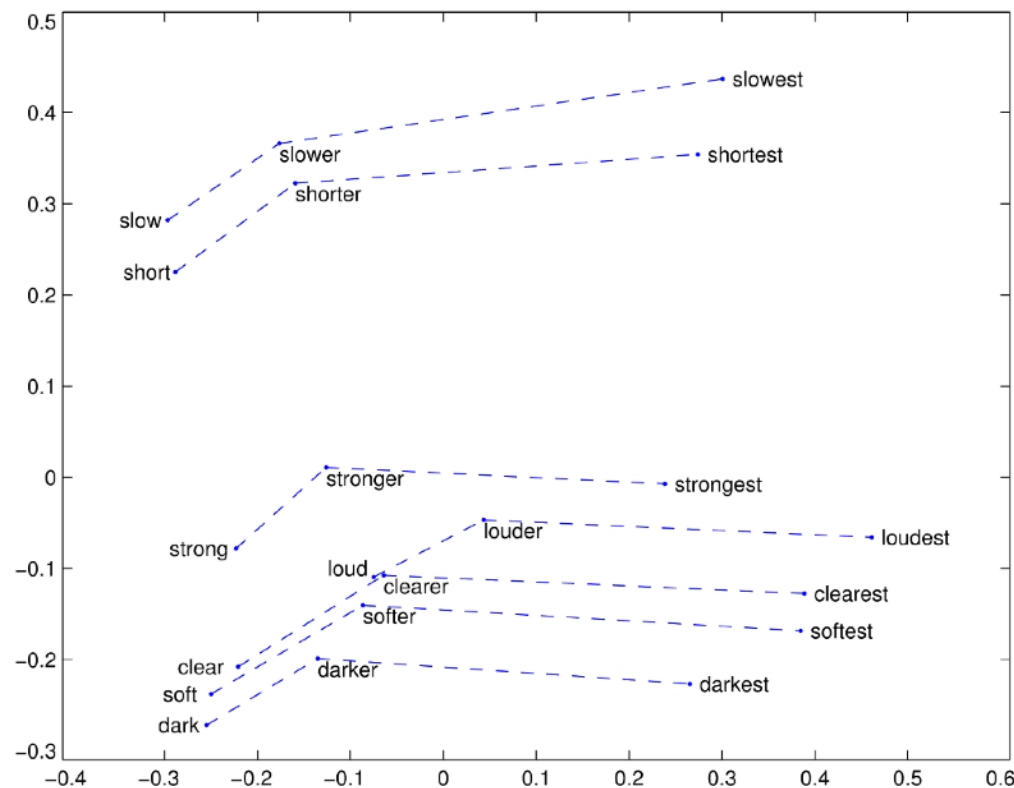
- Word embeddings are applied in a context free manner

```
open a bank account              on the river bank
              ↖                  ↗
                [0.3, 0.2, -0.8, …]
```

- Solution: Train contextual representations on text corpus

```
[0.9, -0.2, 1.6, …]                    [-1.9, -0.4, 0.1, …]
         ↑                                      ↑
open a bank account              on the river bank
```
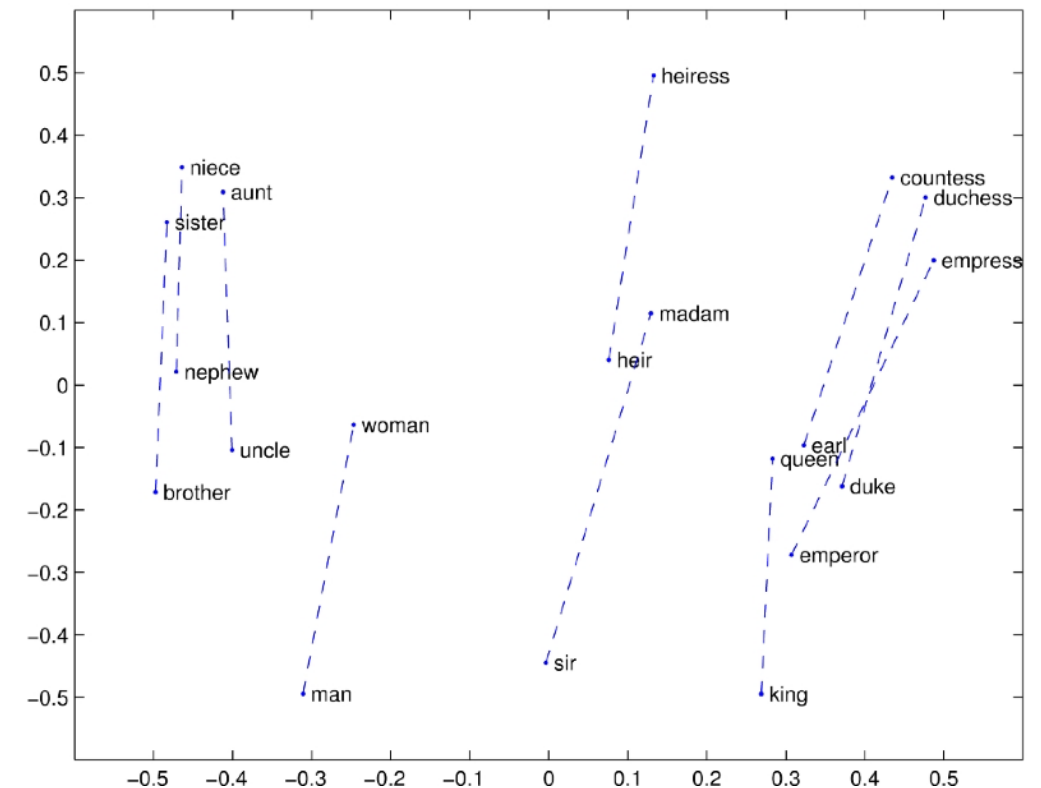
- Example: BERT embeddings (based on Transformers)

# Embedding visualisation



Source: Stanford NLP

- 2d-3d projections (PCA)
  - https://projector.tensorflow.org