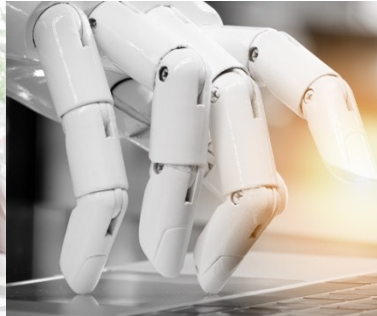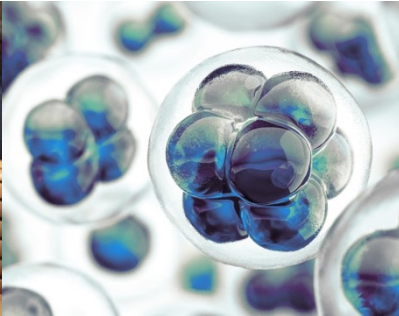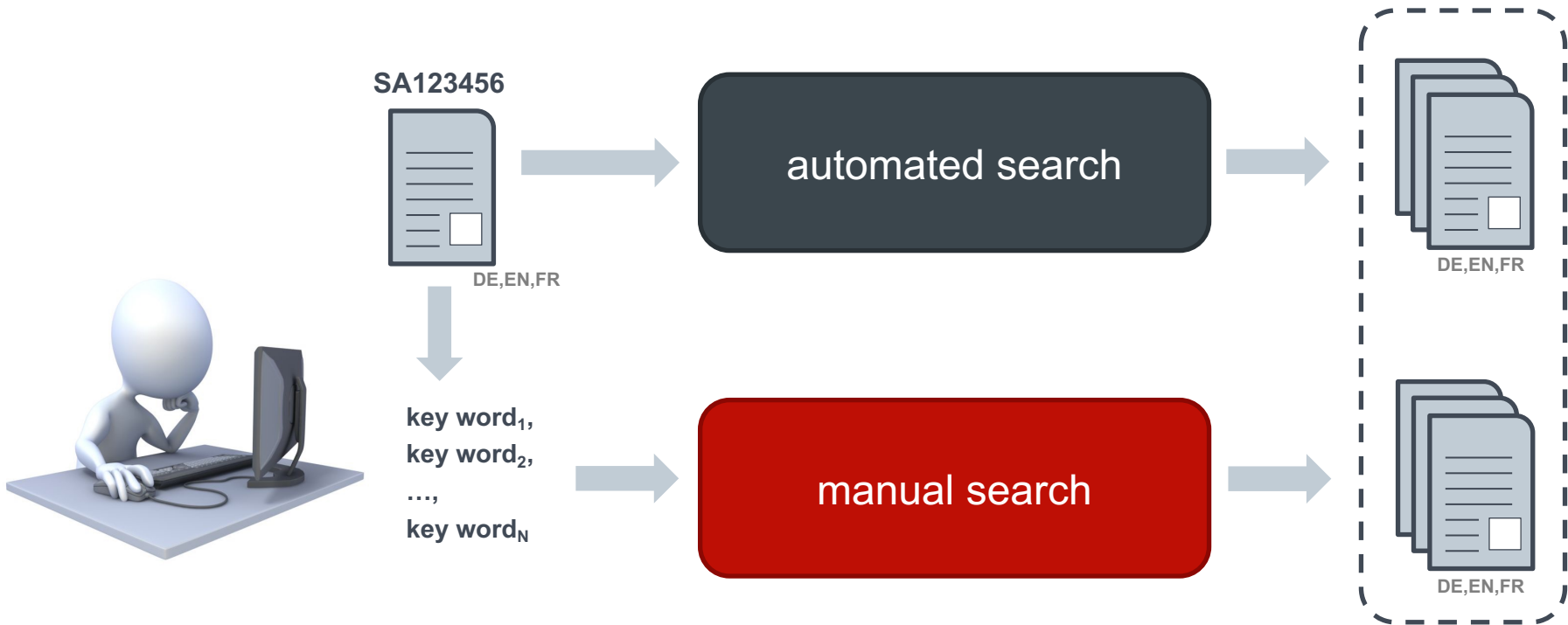# Query Terms Suggestion

# Patent Search

**SA123456**

automated search
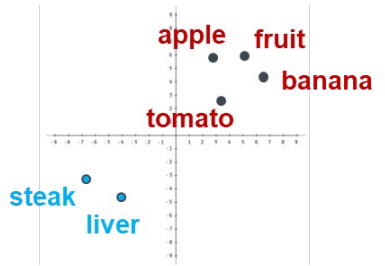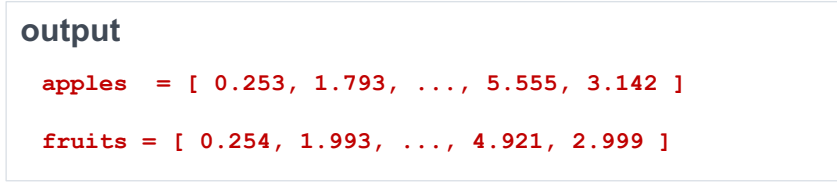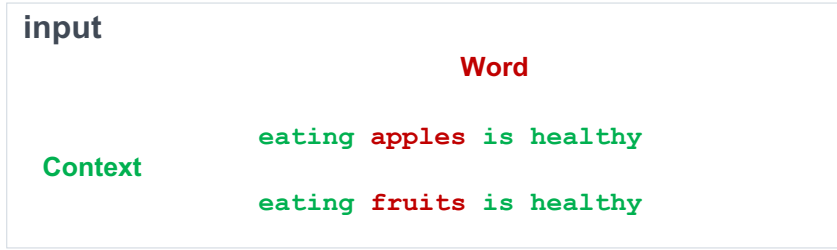
DE,EN,FR

DE,EN,FR

**key word₁,**
**key word₂,**
**…,**
**key wordₙ**

manual search

DE,EN,FR

# Query Terms Suggestion

**Keyboard** → **Semantic Similarity Search** →

**keypad**
**touch screen**
**mouse**
**touch pad**
**joystick**
**joypad**

# Word2Vec

**input**

Context

**Word**

eating **apples** is healthy

eating **fruits** is healthy

**output**

apples  = [ 0.253, 1.793, ..., 5.555, 3.142 ]

fruits = [ 0.254, 1.993, ..., 4.921, 2.999 ]

apple  fruit

banana

tomato

steak  liver

Input layer

Hidden layer

Output layer

apples

$x_k$  $\mathbf{W}_{V \times N}$  $h_i$  $\mathbf{W}'_{N \times V}$

$\mathbf{W}'_{N \times V}$

$\mathbf{W}'_{N \times V}$

$y_{1,j}$  eating

$y_{2,j}$  is

$y_{C,j}$  healthy

$V$-dim

$N$-dim

$C \times V$-dim

[
  0.253,
  1.793,
  ...,
  5.555,
  3,142
]

4

# PreProcessing - text quality control

- A range of filters were developed to detect problematic words

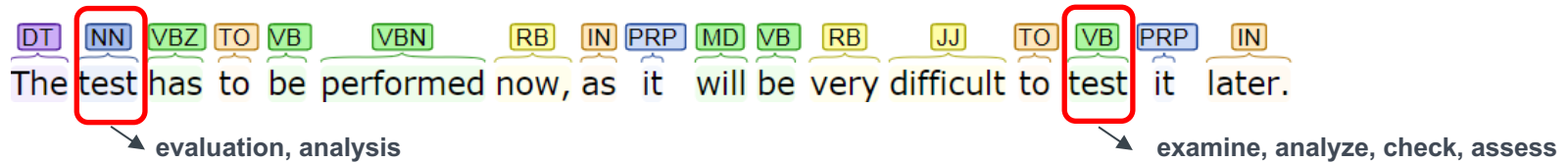| **Examples:** | OCR errors | e.g. *l1elp, rnyself, 1ever, deSense* | **Rule-based** |
| | Space deletion | e.g. *dieArbeit, thenumber* | **CRF, Random Forest** |
| | Space insertion | e.g. *Austral_ia, bio_logy* | **Rule-based** |
| | DNA/Protein sequences | e.g. *AGGATTTCTAAAC, MVFPMWTLKR* | **Rule-based** |

## What to do with these cases ?

- Erroneous words are tagged with a PROB flag

- Number of erroneous words per sentence are counted

- **≥ 15% ?** Too much noise! The sentence is ignored for training of the models

- **< 15% ?** Sentence is used for training, but the resulting model is cleaned up from PROB flagged words

# PreProcessing - Part-of-speech & N-gram detection

**POS tagging**

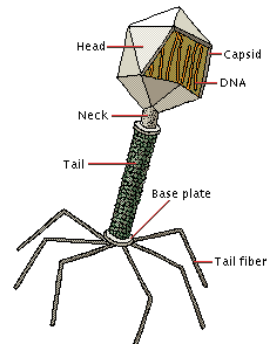- Assignment of POS tags for clear distinguishment of word forms



evaluation, analysis

examine, analyze, check, assess

**Bigrams**

- Example: "windshield wiper" or "car wash"

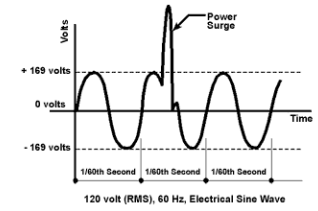- Detected by applying a *chi square* algorithm: How often occur words alone *vs.* together?

**FR – n-grams**

- **3-grams** bridged by *de, d', à, du, en, des, au, aux, a*

- **4-grams** bridged by *de la, d' un, à la, d' une*

# How to deal with word ambiguity?

# CPC to the rescue!



**A43B** – CHARACTERISTICS OF FOOTWEAR …



**G06F** - ELECTRIC DIGITAL DATA PROCESSING …



**H02H** – EMERGENCY PROTECTIVE CIRCUIT ARRANGEMENTS …



**C12N** - MICROORGANISMS OR ENZYMES…

**F16B -** DEVICES FOR FASTENING …

**A23L -** FOODS, FOODSTUFFS …

# Deciding for the CPC level depth

**section**
**G**

Physics

**class**
**G06**

Computing, Calculating, Counting

**sub-class**
**G06F**

Electrical Digital Data Processing

- Trade-off between training corpus size affecting quality and word disambiguation

- Will be input parameter at request time, or result in multiclass output

- What is the optimal level in the CPC tree (1-digit, 3-digit, 4-digit) ?

- Do we have to host 8, 124, or 663 vector sets for similarity search ?

# Deciding for the CPC level depth

**Query:** *" theme "*
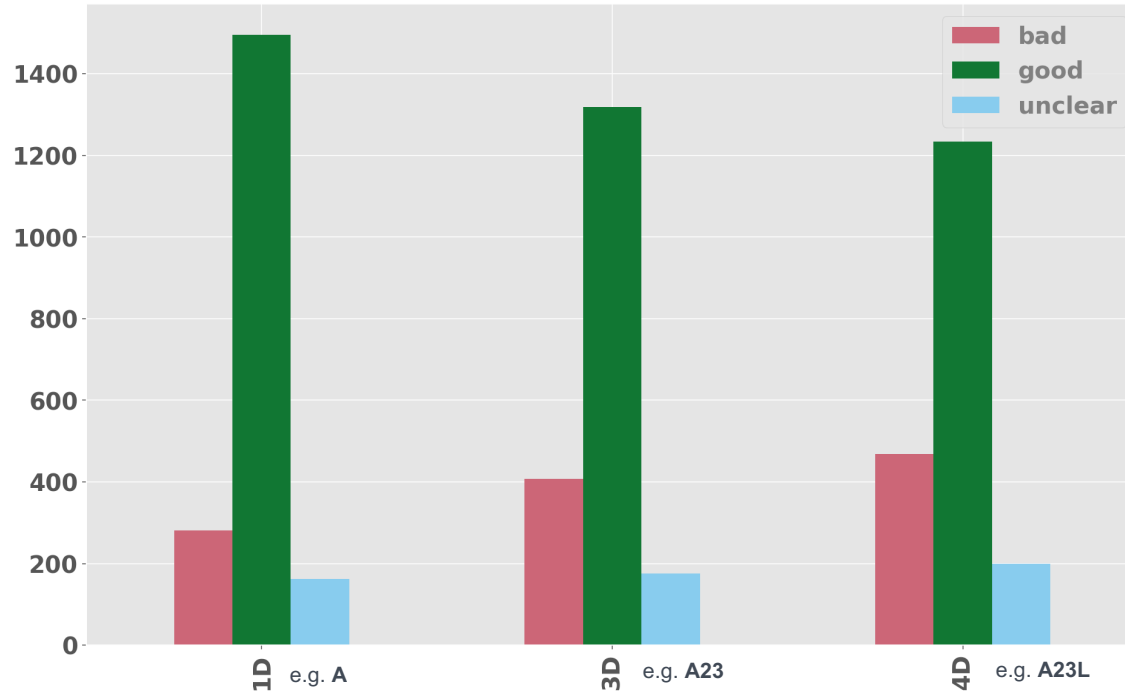
against models

**G**

**G06**

**G06F**



pre-selected query term

common terms of all three models

proposed term
(not found in top 15 of all models)

# CPC level depth results



- Users prefer suggestions from section model
- Clear trend observable of decreasing *good* and increasing *bad* suggestions
- Least models to host and least complex input parameters

**Section models (=8)**
**X**
**Office languages (=3)**
**=**
**24 independent models**

# Some Statistics EN

- Data ingested from 2015-01-01  (7+ years)

- 2,440,964 documents total
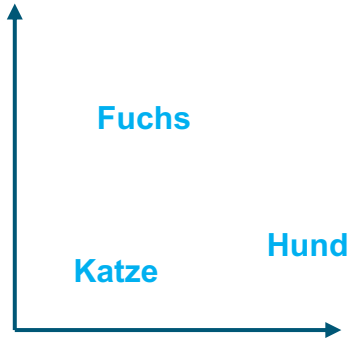
- 204 Gb of ingested data

| processed sentences | | vector counts | |
|---|---|---|---|
| A | 175,933,077 | A | 1,674,493 |
| B | 110,707,118 | B | 869,773 |
| C | 142,638,316 | C | 1,589,179 |
| D | 6,491,405 | D | 153,821 |
| E | 16,958,187 | E | 222,295 |
| F | 42,404,148 | F | 360,287 |
| G | 277,512,083 | G | 1,556,453 |
| H | 233,961,244 | H | 1,044,503 |
| Total | **1,006,605,578** | | |

# Finding similar terms in other languages
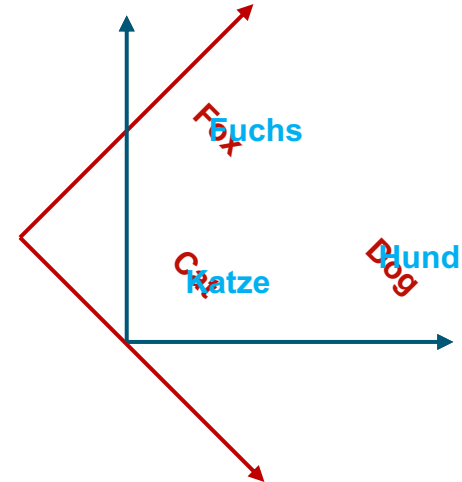
**DE word2vec vector space**

**EN word2vec vector space**

**Aligned word2vec vector spaces**



**SVD**

## SVD

Given known correspondences, i.e. pre-calculated word translations

1. Translate vector space
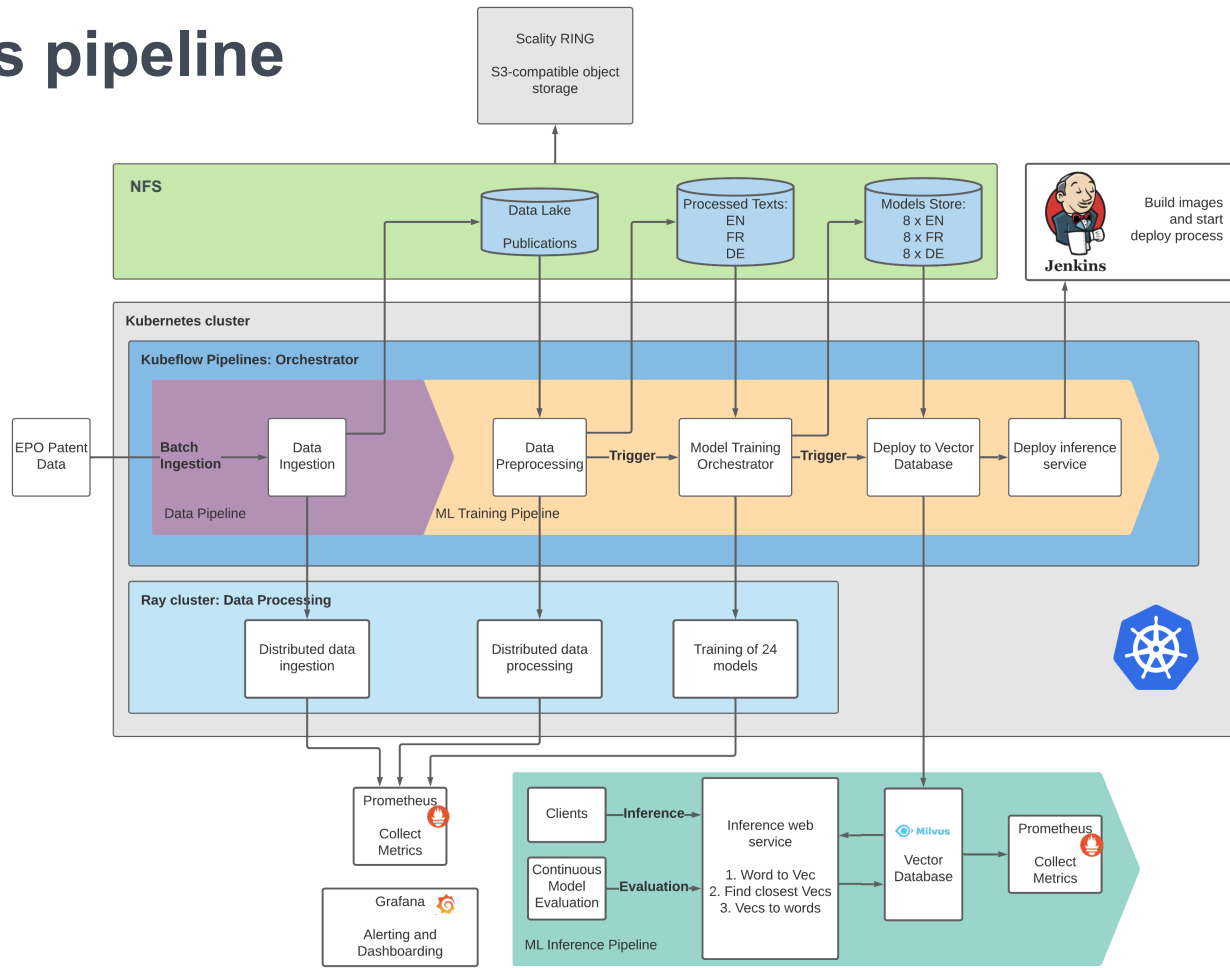2. Rotate vector space
    → allows cross-lingual queries

# Single term query in different language

```
Word searched was --|>>   chair:NN   <<|--

nearest neighbours are:

        chaise:NN                 - (distance: 1358.786376953125, id: 4592)
        fauteuil:NN               - (distance: 1438.371337890625, id: 3112)
        siège:NN                  - (distance: 1637.8919677734375, id: 690)
        lit:NN                    - (distance: 1663.5936279296875, id: 1503)
        fauteuil_roulant:NN           - (distance: 1745.59912109375, id: 4469)
        dossier:NN                - (distance: 1807.46533203125, id: 1736)
        repose-pieds:NN               - (distance: 1813.4537353515625, id: 9678)
```

# MLOps pipeline

# Tool Integration

# Tool Integration

# Thank you very much for your attention!

**Acknowledgements**

Daniel Schneider

Hennadii Stas

Abdelkader Kouhli