# Topic 11: Wikipedia: Surface form extraction (HARD!)

Philipp Heinisch
Ricardo Usbeck
*final presentation*

February 1, 2018

# Table of Contents

# Knowledge extraction from text



  *Blaise Pascal (19 June 1623 to 19 August 1662) was a French mathematician, physicist, inventor, writer and Catholic theologian[...]*
  *Pascal was an important mathematician[...]*

Belongs *Blaise Pascal* and/ or *Pascal* to the URI
`http://dbpedia.org/resource/Blaise_Pascal`?

# Knowledge extraction from text



_Blaise Pascal_ (19 June 1623 to 19 August 1662) was a
French mathematician, physicist, inventor, writer and Catholic
theologian[...]
_Pascal_ was an important mathematician[...]

Belongs _Blaise Pascal_ and/ or _Pascal_ to the URI
http://dbpedia.org/resource/Blaise_Pascal?

**Task:** Knowledge bases contain labels for resources (rdfs:label). [...]
The goal of this task is to detect candidates for labels using Wikipedia. [...]

# Knowledge extraction from text



_Blaise Pascal_ (19 June 1623 to 19 August 1662) was a French mathematician, physicist, inventor, writer and Catholic theologian[...]
_Pascal_ was an important mathematician[...]

Belongs _Blaise Pascal_ and/ or _Pascal_ to the URI
`http://dbpedia.org/resource/Blaise_Pascal`?
**Task:** Knowledge bases contain labels for resources (`rdfs:label`). [...]
The goal of this task is to detect candidates for labels using Wikipedia. [...]

# Solve this task with the KATANA algorithm

- $\exists$ list of candidates ($+$ knowledge about these)
- $\exists$ list of extracted labels with extracted knowledge from text

## KATANA!

Match the labels to the candidates $=$ calculate the score for each candidate-label $\Rightarrow$ the highest score wins!

# Solve this task with the KATANA algorithm

- $\exists$ list of candidates ($+$ knowledge about these)
- $\exists$ list of extracted labels with extracted knowledge from text

## KATANA!

Match the labels to the candidates $=$ calculate the score for each candidate-label $\Rightarrow$ the highest score wins!

How good is the KATANA algorithm?

# Solve this task with the KATANA algorithm

- $\exists$ list of candidates ($+$ knowledge about these)
- $\exists$ list of extracted labels with extracted knowledge from text

### KATANA!

Match the labels to the candidates $=$ calculate the score for each candidate-label $\Rightarrow$ the highest score wins!

How good is the KATANA algorithm?

# KATANA Algorithm/ Formulas

## Given

Knowledge base $KB$ $(s, p, o)$, our extracted triples $ext$ from natural text with labels, find out the matching URI-candidate $c_s$ from $\{c_1, ..., c_n\}$

## Determine the ambiguity of a fact to a given subject $s$

$$\psi(p, o) = 1 - \frac{1}{|\{s | (s, p, o) \in KB\}|}$$

# KATANA Algorithm/ Formulas

## Given

Knowledge base $KB$ $(s, p, o)$, our extracted triples $ext$ from natural text with labels, find out the matching URI-candidate $c_s$ from $\{c_1, ..., c_n\}$

## Determine the ambiguity of a fact to a given subject $s$

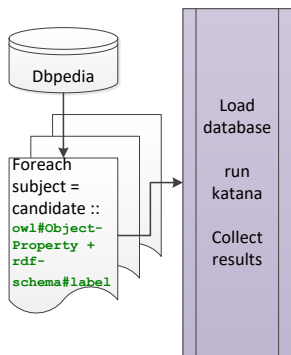$\psi(p, o) = 1 - \frac{1}{|\{s|(s,p,o) \in KB\}|}$

## Determine the score (matching grade) candidate $\rightarrow$ label ($\rightarrow \lambda$)
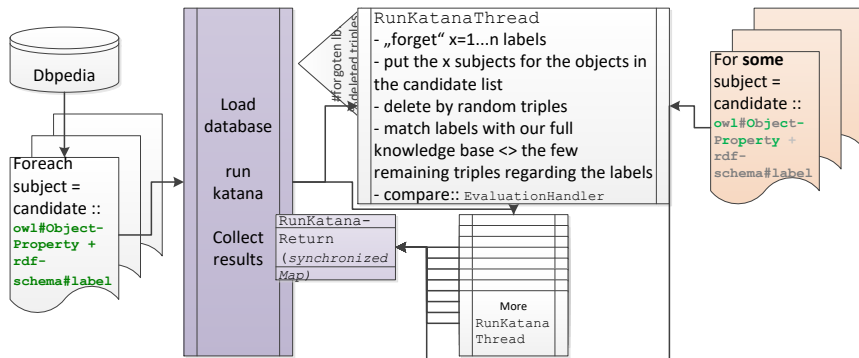
for a certain candidate $c$:

$M(c, s) = \{(p, o)|(\lambda, p, o) \in ext \cap (c, p, o) \in KB\}$

$$score(c, s) = \begin{cases} 0 & M(c, s) = \emptyset \\ 1 - \prod_{(p,o) \in M(c,s)} \psi(p, o) & M(c, s) \neq \emptyset \end{cases}$$

# KATANA Algorithm/ Formulas

## Given

Knowledge base $KB$ $(s, p, o)$, our extracted triples $ext$ from natural text with labels, find out the matching URI-candidate $c_s$ from $\{c_1, ..., c_n\}$

## Determine the ambiguity of a fact to a given subject $s$

$$\psi(p, o) = 1 - \frac{1}{|\{s|(s,p,o) \in KB\}|}$$

## Determine the score (matching grade) candidate $\rightarrow$ label $(\rightarrow \lambda)$

for a certain candidate $c$:
$$M(c, s) = \{(p, o)|(\lambda, p, o) \in ext \cap (c, p, o) \in KB\}$$

$$score(c, s) = \begin{cases} 0 & M(c, s) = \emptyset \\ 1 - \prod_{(p,o) \in M(c,s)} \psi(p, o) & M(c, s) \neq \emptyset \end{cases}$$
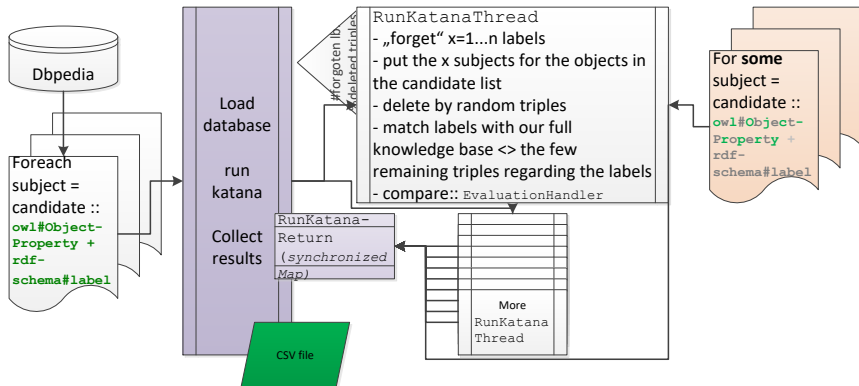
# My application

# My application

# My application follows the command pattern



## Available commands without their parameters

- Environment commands: `help`, `exit`
- Database commands: `load`, `edit`, `print`
- KATANA (evaluation) commands: `katana`

## Demo in the end

... if there is time...

# My application follows the command pattern



## Available commands without their parameters

- Environment commands: `help`, `exit`
- Database commands: `load`, `edit`, `print`
- KATANA (evaluation) commands: `katana`

## Demo in the end

... if there is time...
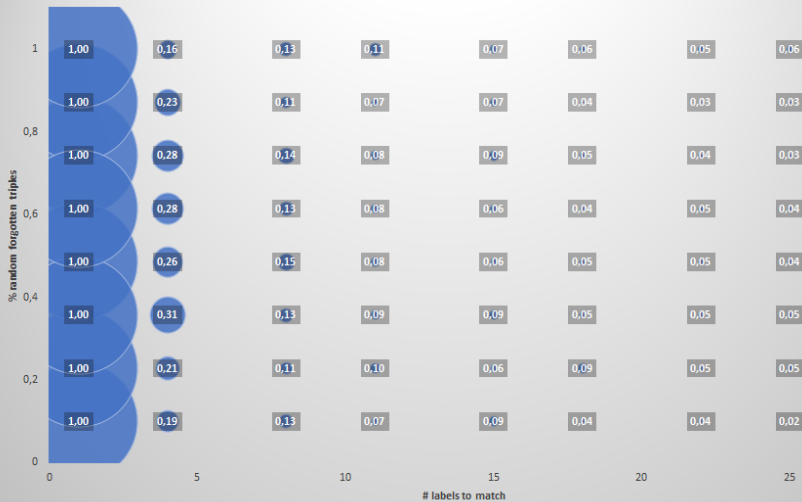
# Evaluation

## Result measurement function (`EvaluationHandler::calculateAccuracy()`)

$L$ is set of guessed labels, $L' \subseteq L$ set of wrong guessed labels, $C$ set of correct labels and $C' \subseteq C$ set of labels that doesn't appear in $L$:

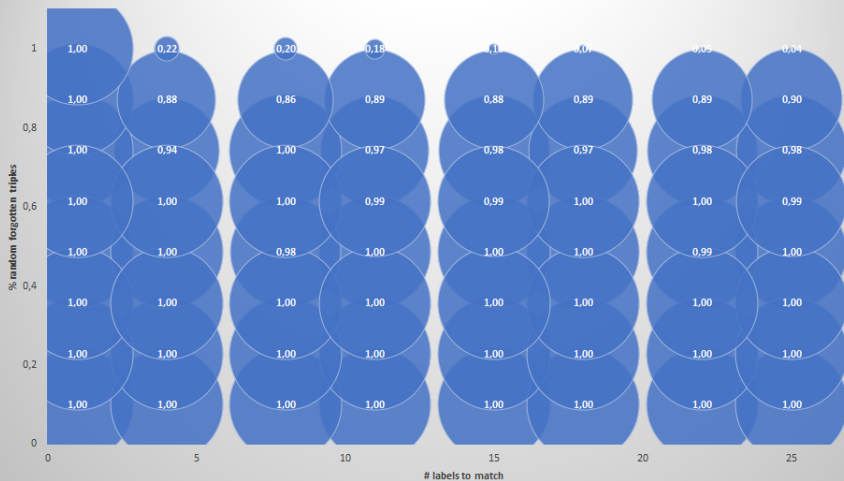Size of bubbles $s = \frac{2 - \frac{|L'|}{|L|} - \frac{|C'|}{|C|}}{2}$, $0 \leq s \leq 1$

# Evaluation



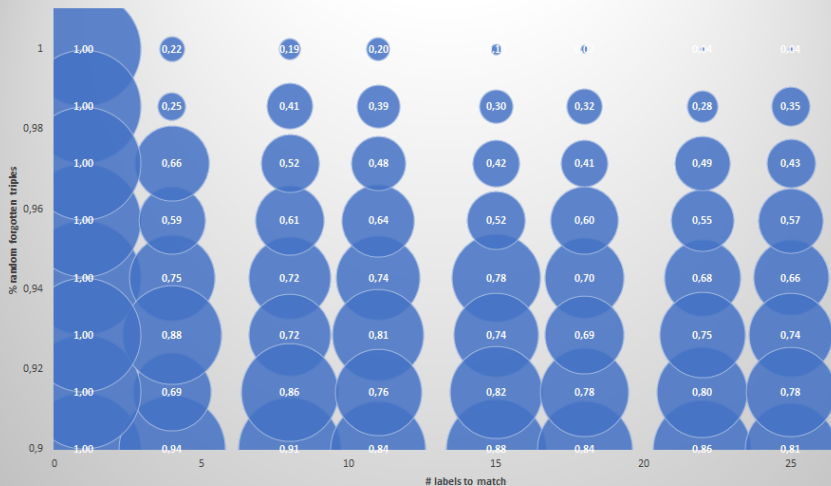Result random label matching (domain Scientist)

# Evaluation



Result KATANA label matching (domain scientist)

# Evaluation



Result KATANA label matching (domain Scientist)

# Discussion of evaluation results

1. KATANA leads to much more better results than the random matching
2. KATANA is very precise ($\geq 80\%$) until 90% data-loss [in a big data set, too]
   - lots of properties are nearly unambiguously, e.g.: birth date, spouse, (wikiPageExternalLink), ...
   - to match 1 label in the real world, you would have $\geq 1$ candidates (selective range)

https://github.com/dice-group/KATANA
branch: Philipp_Heinisch_master