# LARGE-SCALE INFORMATION EXTRACTION FROM NEUROSCIENTIFIC LITERATURE

Marco Antognini

Spring 2015

# MOTIVATION

The ultimate goal of the Blue Brain Project is to reverse engineer the mammalian brain.

➤ Interactions between brain regions?
➤ What about effects of specific cells?
➤ And connections with other organs?
➤ . . . and many many more questions.

The neuroscientific literature already holds many answers but. . .

# MOTIVATION

The ultimate goal of the Blue Brain Project is to reverse engineer the mammalian brain.

➤ Interactions between brain regions?
➤ What about effects of specific cells?
➤ And connections with other organs?
➤ . . . and many many more questions.

The neuroscientific literature already holds many answers but. . .

⟹. . . we need a tool to extract specific information from this colossal amount of text.

# UIMA & RUTA: IN BRIEF

# UIMA

➤ Unstructured Information Management Architecture

➤ Engines produce annotations (metadata)

➤ Engines are combined to form pipelines

➤ Most engines are written in Java

➤ Not limited to neuroscience topics

Assuming we want to find all sentences in this text...

Terminologies which lack semantic connectivity hamper the effective search in biomedical fact databases and document retrieval systems. We here focus on the integration of two such isolated resources, the term lists from the protein fact database UNIPROT and the indexing vocabulary MESH from the bibliographic database MEDLINE.

# UIMA – METADATA

The resulting annotations, presented in JSON format:

```
"DocumentAnnotation" : [
  { "begin" : 0,    "end" : 328,  "language" : "en" }
],
"Sentence" : [
  { "begin" : 0,    "end" : 135,
    "componentId" :
        "de.julielab.types.OpenNLPSentenceDetector" },
  { "begin" : 136,  "end" : 32
    "componentId" :
        "de.julielab.types.OpenNLPSentenceDetector" }
]
```

Bluima regroups UIMA engines, focusing on neuroscientific engines:

➤ Preprocessing (sentence, tokenizer, PoS, . . . )
➤ Linnaeus (species recognition)
➤ Oscar (chemistry)
➤ Brain regions & their relations
➤ Proteins
➤ Measures and units
➤ . . .

# RUTA

➤ Rule-based Text Annotation

➤ Scripting language

➤ Can integrate engine in script

➤ Makes it a bit easier to write UIMA pipeline

Basic RUTA script tagging dogs as ANIMAL:

```
DECLARE Animal;
W{REGEXP("dog") -> MARK(Animal)};
```

Integration of UIMA engines:

```
ENGINE SentenceAnnotator;
ENGINE TokenAnnotator;
ENGINE PosTagAnnotator;

Document{-> EXEC(SentenceAnnotator)};
Document{-> EXEC(TokenAnnotator)};
Document{-> EXEC(PosTagAnnotator)};
```

Integration of UIMA engines:

```
ENGINE SentenceAnnotator;
ENGINE TokenAnnotator;
ENGINE PosTagAnnotator;

Document{-> EXEC(SentenceAnnotator)};
Document{-> EXEC(TokenAnnotator)};
Document{-> EXEC(PosTagAnnotator)};
```

BUT the engines' settings are stored in an unfriendly XML file...

# ISSUES

With both raw UIMA and RUTA scripts, some additional issues:

➤ Every user has to go through the hassle of installing UIMA/RUTA

➤ Not so trivial to run pipelines

➤ Manage external resources (install, update, remove, . . . )

➤ Manage versions of engines & pipelines

# SHERLOK

# SHERLOK – TEXT-MINING SERVICE

In a few words:

➤ RESTful service: 1 server for many users

➤ Based on UIMA $\implies$ existing engines compatible

➤ Based on RUTA $\implies$ allowing powerful scripts

➤ Makes it easy to configure engines and pipelines

➤ Automatic resource management

➤ Versioning of pipelines

# SHERLOK – TEXT-MINING SERVICE

In a few words:

➤ RESTful service: 1 server for many users

➤ Based on UIMA $\implies$ existing engines compatible (almost)

➤ Based on RUTA $\implies$ allowing powerful scripts

➤ Makes it easy to configure engines and pipelines

➤ Automatic resource management

➤ Versioning of pipelines

Sometimes, the Java implementation needs some minor refactoring.

What does it mean?

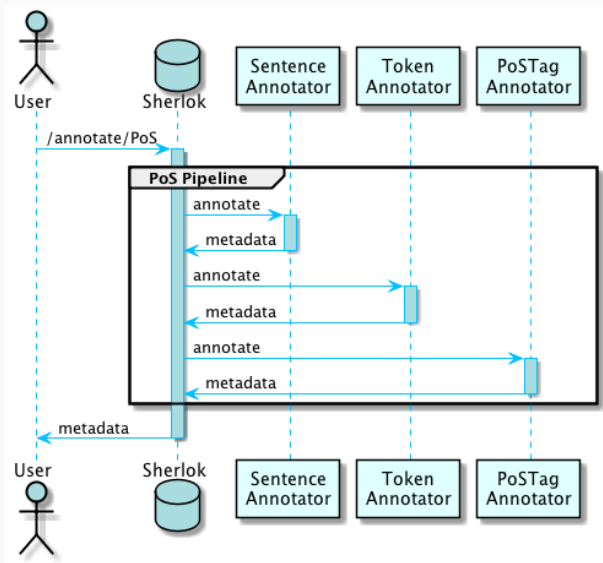Basically, the HTTP protocol is all is needed to communicate with Sherlok.

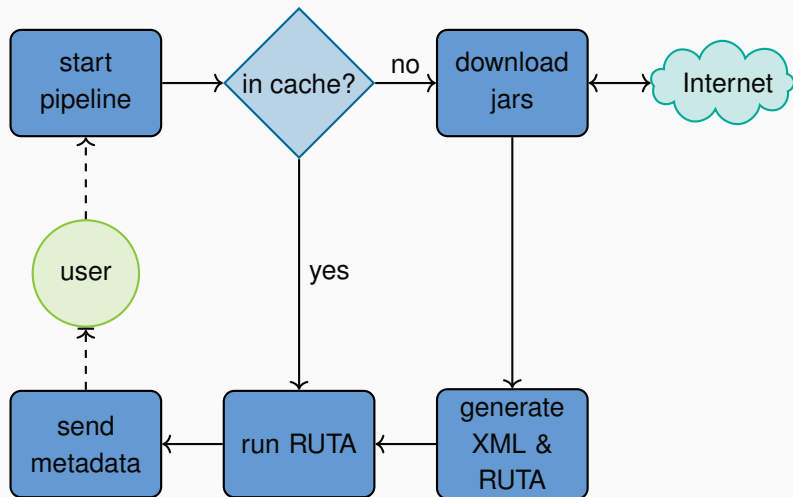$\implies$ Can easily be embedded in any programming language or tool!

# SHERLOK – RESTFUL

What does it mean?

Basically, the HTTP protocol is all is needed to communicate with Sherlok.

$\implies$ Can easily be embedded in any programming language or tool!

---

To annotate some text with the PoS pipeline:
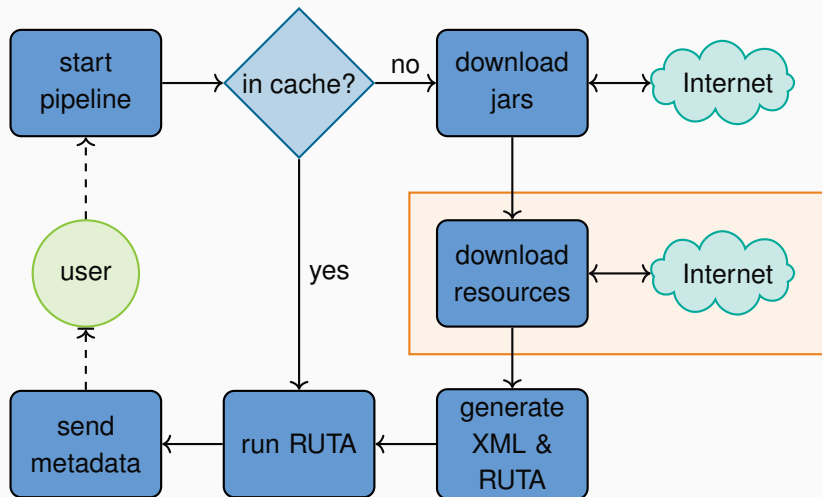
```
GET /annotate/bluima.pos?text=<...>
```

```json
{
  "name": "bluima.sentence", "version": "1.0.1",
  "dependencies": [
    { "value": "ch.epfl.bbp.nlp:bluima_opennlp:1.0.1" }
  ], "config": {
    "bluima": {
      "type": "git", "ref": "master",
      "url":
        "https://github.com/BlueBrain/bluima_resources.git" }
  }, "engines": [ {
      "name": "SentenceAnnotator",
      "class": "ch.epfl.bbp.uima.ae.SentenceAnnotator",
      "parameters": {
        "modelFile":
          "$bluima/opennlp/sentence/SentDetectGenia.bin.gz"
      }
  } ]
}
```

# SHERLOK – PIPELINE CONFIGURATION

```
{
  "name": "countries", "version": "1",
  "description": "Example that annotates countries",
  "config": {
    "countries": {
      "type": "http", "mode": "ruta",
      "url": "https://example.com/countries.txt"
    }
  }, "script": [
    "WORDLIST CountriesList = '$countries';",
    "DECLARE Country;",
    "Document{-> MARKFAST(Country, CountriesList)};"
  ]
}
```

# DEMO!

# Questions?

Check it out

http://sherlok.io