# Big Data Architectures
# For Machine Learning and Data Mining
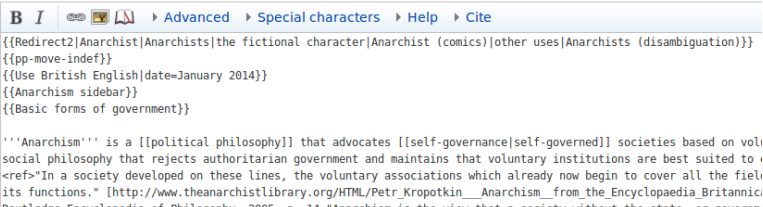
April 26, 2017

Presentation of the data sets

# Presentation of the data sets from ClueWeb

- Location: HDFS
- How to get there? **hadoop fs -ls /corpora/clueweb/09** or **hadoop fs -ls /corpora/clueweb/12**
- Example: In /**corpora**/**clueweb**/**09**/**ClueWeb09**/**_Arabic_1**/ **ar0000**/**00.warc.gz** are websites as html
- Typical sizes of one file: 133.421.986 Bytes, for all in Clueweb09: 3.9 TByte, for all in Clueweb12: 4.5 TByte

```
<LINK type="text/css" href="theme/CSS/theme.css"  rel=stylesheet>
<LINK type="text/css" href="theme/JS/theme.css"  rel=stylesheet>
<link rel="sitemap" href="sitemap/site_map.php">
<script type="text/javascript" src="theme/JS/JSCookTree.js"></SCRIPT>
<script type="text/javascript" src="theme/JS/theme.js"></SCRIPT>
```

# Presentation of the data sets from Wikipedia

- Location: HDFS
- How to get there? **hadoop fs -ls /corpora/enwiki/20160501/ enwiki-20160501-pages-articles.xml.bz2**
- Example: In **Wikipedia Anarchism** → **Edit** are texts with links
- Typical sizes: 12.896.598.324 Bytes

# Presentation of the data sets from the academic graph

- Location: HDFS
- How to get there? **hadoop fs -ls /corpora/corpus-microsoft-academic-graph**
- Example (shortened): graph in text
- Typical sizes: 26.6 GByte

| Paper ID | Orig. paper title | Paper publish year | Paper publish date | Paper Docum. Object Identif. | Orig. venue name | Normal. venue name | Confer. ID mapped to venue name | Paper rank |
|---|---|---|---|---|---|---|---|---|
| | | | | | | | | |