# Passim vs Tracer

| | Tracer | Passim |
|---|---|---|
| Language | Java | Java / Scala |
| Compiler | ant | sbt |
| Dependencies | <ul><li>java8</li><li>medusa (comes together with tracer)</li></ul> | <ul><li>java8</li><li>apache spark</li><li>hadoop</li></ul> |
| Cross-platform? | Yes | No |
| Restrictions | — | Does't work on OS X (due to Hadoop's nativelib) |
| Open-source? | Partly; the repo is restricted to registered users, but they welcome contributions | Yes |
| Source code | eTRAP's GitLab<br><br>http://vcs.etrap.eu/tracer-framework/tracer | GitHub<br><br>https://github.com/dasmiq/passim |
| Manual | https://tracer.gitbook.io/-manual/<br><br>Good detailed manual, but still incomplete.<br><br>Also, after compiling Tracer you can generate a javadoc with a short description of every class used in the program. | The only manual is a README in Passim's GitHub repo, but it's not too bad.<br><br>No minimal working example. |
| Community | eTRAP, active community | Used in several projects like Kitab, but no active community |
| Support | You have to register on http://www.etrap.eu/redmine/login and wait for account approval to raise an issue.<br><br>Approval seems to never happen (still pending a month after account creation) | Very easy to open an issue on GitHub, but they aren't answered. |
| Built-in visualisation | Yes (TraVIZ) | No |

| Built-in alignment | Yes | Yes |
|---|---|---|
| | Outputs aligned passage ids together with their absolute overlap (nr of shared features) to the *.link* file | Produces results in *.json* or *.parquet* format, which contain a document/line from the right-hand, "target" text, along with information about corresponding passages in the left-hand, "witness" or "source" text. |
| Algorithm/scoring | Their own similarity score | Single-link(age) clustering |

# Passim

- D. A. Smith, David R. Cordell, and David A. Ryan Abby Mullen. 2015. Computational methods for uncovering reprinted texts in antebellum newspapers. American Literary History, 27: E1 – E15 (pdf)
- MacLaughlin, Ansel, and David A. Smith. "Content-based Models of Quotation." (pdf)
- Romanello, Matteo. *Detecting Text Reuse in Newspapers Data with Passim*. No. POST_TALK. 2018. (slides)
- https://github.com/dasmiq/passim
- Eclipse Public License

## Requirements & installation

1. Java 8 (not higher!)
2. Scala
3. Sbt ( a build tool for JVM languages)
   a. **NB!** Build Passim with sbt before installing Apache-Spark. Spark comes with Java 11, which results in build conflicts when you are compiling Passim with sbt.
4. Apache-Spark
   a. Comes with Java 11
   b. Depends on Hadoop, which has native libraries not supported in OS X →
      i. Won't run on a Mac

## Key features (from documentation)

- Detects and aligns similar passages in text.
- Can be run either in query mode, to find quoted passages from a reference text, or all-pairs mode, to find all pairs of passages within longer documents with substantial alignments.
- The input to passim is a set of *documents* (whole books, pages of books, whole issues of newspapers, individual newspaper articles, etc.) Minimally, a document consists of an identifier string and a single string of text content.
- Sometimes it is useful to group documents into *series*. Text reuse within the series will be ignored.

- Output: clusters and pairwise alignments between all matching passages
- Uses the [Smith–Waterman alignment algorithm](#) to find all pairs of passages within longer documents (source and derived)
- Parameters:

| Parameter | Default value | Description |
|---|---|---|
| `--n` | 5 | N-gram order for text-reuse detection |
| `--max-series` | 100 | Maximum document frequency of n-grams used. |
| `--min-match` | 5 | Minimum number of matching n-grams between two documents. |
| `--relative-overlap` | 0.5 | Proportion that two different aligned passages from the same document must overlap to be clustered together, as measured on the longer passage. |

## Algorithm

### Step 1: search of candidate document pairs

**Goal:** reduce total number of comparisons to perform

**1.1 shingling:**
- efficient document indexing via n-grams
- document → set of 5-word sequences (5-grams)
- singleton n-grams are filtered out (> 50%)

**1.2 extraction and filtering of candidate pairs**
- suppress repeated n-grams within same series
- suppress n-grams leading to > 5k document pairs
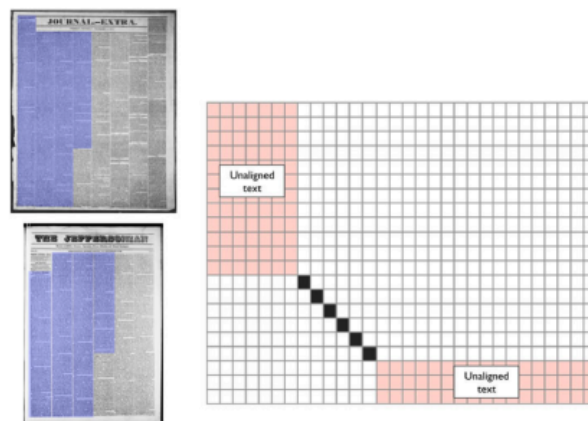- filter out document pairs with < 5 shared n-grams

**5-grams**

devant le verdict de l'expert sur la responsabilité et réclame, en raison de la violence

⇩

1. devant le verdict de l'expert
2. le verdict de l'expert sur
3. verdict de l'expert sur la
4. de l'expert sur la responsabilité
5. ....
6.

### Step 2: local document alignment

**Goal:** given a set of document pairs, obtain a set of aligned passage pairs



*source of illustration: Smith et al. 2015, Fig. 4, p. E9

**Step 3: passage clustering**
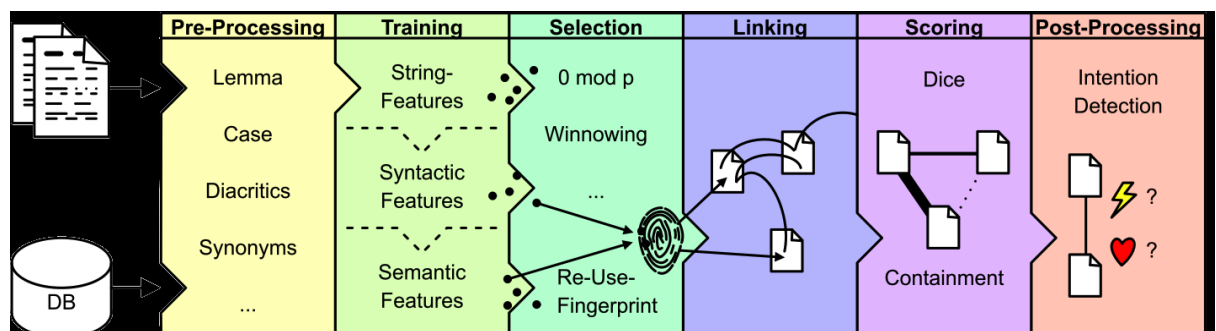
Goal: group similar passages into clusters
- "single-link" clustering
- overlap threshold 50%

**Single-link(age) clustering**
- A method of hierarchical clustering. It is based on grouping clusters in bottom-up fashion (agglomerative clustering), at each step combining two clusters that contain the closest pair of elements not yet belonging to the same cluster as each other (wiki)
- The similarity of two clusters is the similarity of their *most similar* members (Stanford NLP)

## Tracer

- **Büchler, M.**, Burns, P. R., Müller, M., **Franzini, E.**, **Franzini, G.** (2014) 'Towards a Historical Text Re-use Detection', In: Biemann, C. and Mehler, A. (eds.) *Text Mining, Theory and Applications of Natural Language Processing*. Springer International Publishing Switzerland. (pdf)
- Documentation: https://gfranzini.gitbooks.io/tracer/content/
- Academic Free License 3.0 (AFL)


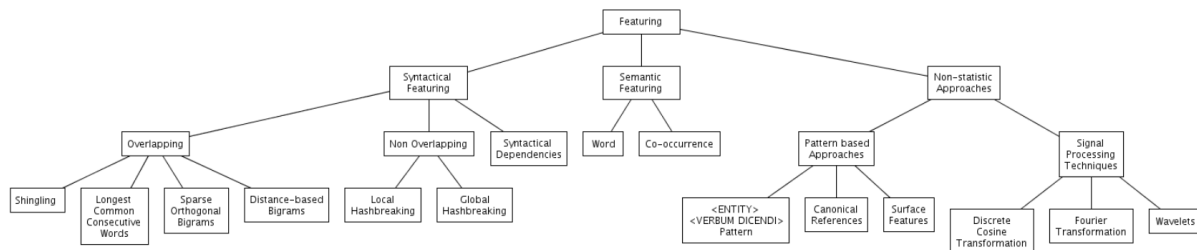
Requirements & installation

1. Any OS type (Windows/OS X/Linux)
2. Java 8 (not higher!)
3. Apache Ant (a build tool)
   a. Only use 'ant' command when compiling Tracer, using it with any flags (like "verbose") will later result in "Unable to access jarfile" error.
4. Troubleshooting: https://gfranzini.gitbooks.io/tracer/content/support/troubleshooting/

Key features (from documentation)

- A framework of ca. 700 algorithms, whose features can be combined to create the optimal model for detecting those words, sentences and ideas that have been reused across texts
- Helps to identify plagiarism in a text, as well as verbatim and near verbatim quotations, paraphrase and even allusions
- reuse results *can* be visualised in a more readable format via TRAViz

- **language independent** and has been successfully tested on Ancient Greek, Arabic, Coptic, English, German, Hebrew, Latin and Tibetan
- Supports intralinking (reused passages within a document) and interlinking (matches between different documents)
- N-gram features, can be characters or words, can be pruned and weighed
  - Overlapping (shingling)
  - Non-overlapping (hash-breaking)
  - Distance-based (good for detecting paraphrases)



## Scoring

The *Scoring* step of TRACER assigns a weight to a reuse pair based on an internal scoring metric. The resemblance score is calculated as follows:

$$\theta_{\Theta}^{R'}(s_i, s_j) = \theta_{\Theta}^{O}(s_i, s_j) \cdot \theta_{\Theta}^{R}(s_i, s_j) = \frac{|\overrightarrow{s_i} \cap \overrightarrow{s_j}|}{|\overrightarrow{s_i} \cup \overrightarrow{s_j}|}$$

The resemblance score *θ* (theta) is the quotient of the reuse overlap (∩) and the reuse union ( ∪ ), where the reuse overlap represents the digital fingerprint of reuse units that appear in two candidate strings/lines, and the reuse union represents the number of all digital fingerprints.

# Experiments

A GitHub repository with preprocessed data and results (both Cath Almaine and Welsh law): https://github.com/ancatmara/text-reuse-test

Could only test Tracer, because Passim doesn't compile on Mac, and trying to run it from a Linux virtual machine under Mac/Windows ended up in memory error. An actual Linux server/a powerful laptop with Linux as the main OS is required to test Passim.

## Tracer

java -Xmx600m -Deu.etrap.medusa.config.ClassConfig=conf/tracer_config.xml -jar tracer.jar > FILENAME.OUT

- Your data must be a single .txt file in tracer/data/corpora/*some-folder*
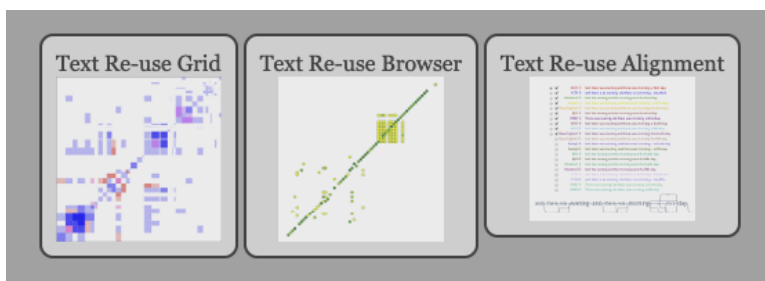- Check and update tracer/conf/tracer_config.xml before running

Data format

1. The first column contains unique sentence IDs.
2. The second column contains the sentence.
3. The third column can contain the date of the file creation in the YYYY-MM-DD format or the word NULL. If you use NULL, make sure it is written in upper-case.
4. The fourth column contains the book or section from which the sentence is taken. This information is crucial for the visualisation shown in the Text Reuse Browser figure in postprocessing. The top drop-down menu you see there will list the information you provide in this fourth column.

```
11000001  Cath Almaine etir Laigniu ocus Uí Néill.  NULL  cath_almaine_1
11000002  teirtíd Decembris rocuired in cath-sa. Cauis in chatha-sa .i. in bóroma romaith Fínnachta do Mo-Ling a
11000003  Nís-tucsat Laigin do Loingsech mac Áengusa ocus ní tucsat do Chongal Chinnmagair, cia rofuilngetar mór
11000004  Ba trom trá la Fergal sin .i. Laigin do nemchomall a n-gellta fris, co rofhuacrad sluaiged dírecra dímc
11000005  Ba fata trá robás ocin tinól-sain; uair is ed at-beired cach fer do Leith Chuinn cosa roiched in fuacra
11000006  Donn Bó immorro mac baintrebthaige eiside d'Fheraib Rois, ocus ní dechaid lá ná aidche a taig a máthar
11000007  Ni raba i n-Éirinn uile bud gríbda nó bud ségainne inás, ocus is uad bud ferr rann espa ocus ríg-scéla
11000008  Is é bud ferr do glés ech ocus do innsma shleg ocus d'fhige fholt, ocus bud ferr fri aichne 'na einech
11000009  Áille macaib Donn Bó báid, binne a laíde luaidit beoil, áine ócaib Inse Fáil, rotócaib táin trilse a th
11000010  Nír léic dno a máthair Donn m-Bó la Fergal co tucad Máel mac Faílbe meic Erannáin meic Chrimthainn, con
11000011  Tochomla dno Fergal for sét.  NULL  cath_almaine_1
```
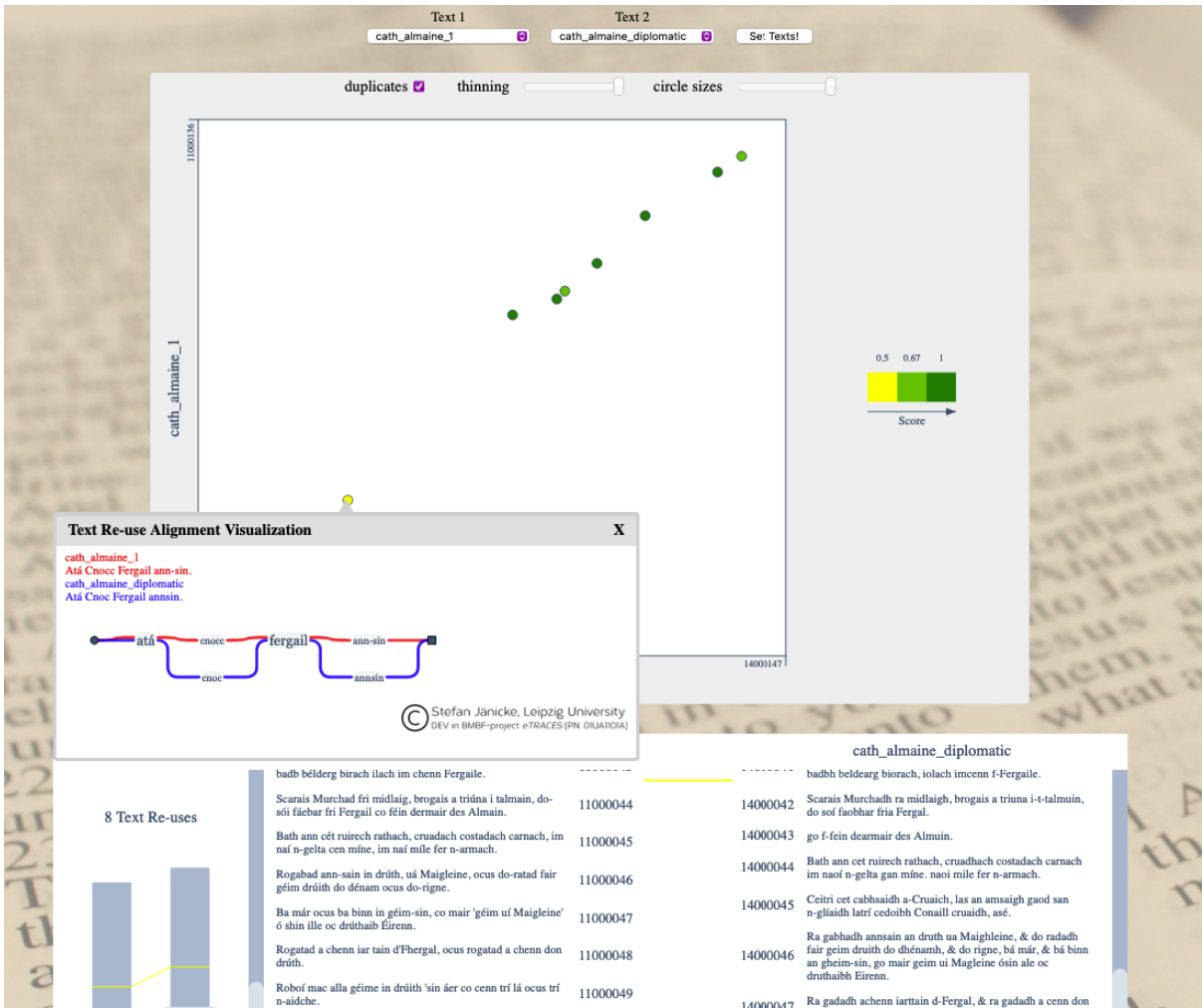
# TraViz

A built-in visualisation package for Tracer
● Text re-use grid and alignment only work for comparing two texts; they seem not to work at all if there were more than 2 texts in the experiment
● Text re-use browser is also binary, but you can select the pair of texts to visualise
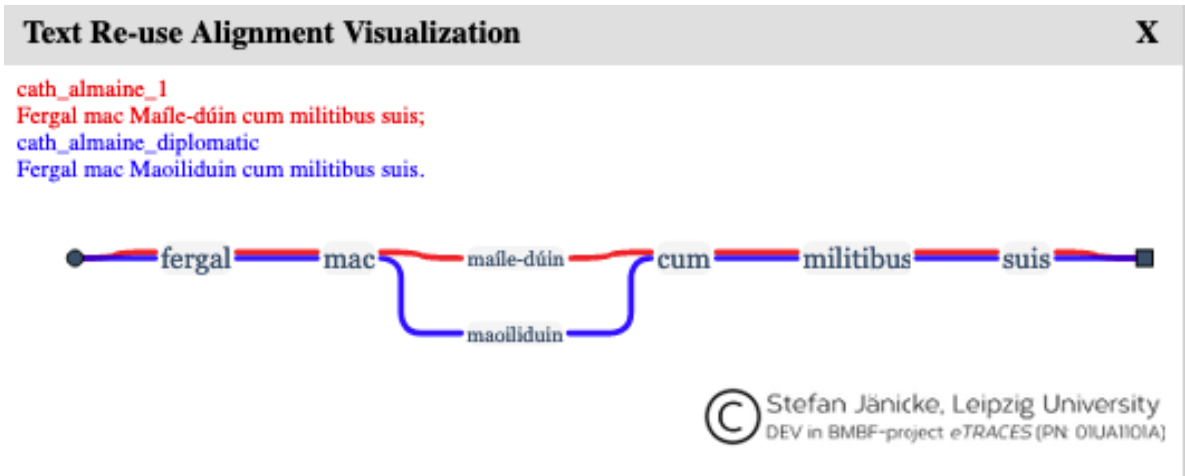


# Cath Almaine Experiment

● 4 different editions, one of them diplomatic
● Minimal unit for comparison — a sentence (semi-automatic splitting)
● Basic TRACER configuration
  ○ No synonyms
  ○ No lemmatisation
  ○ No syntactic parsing
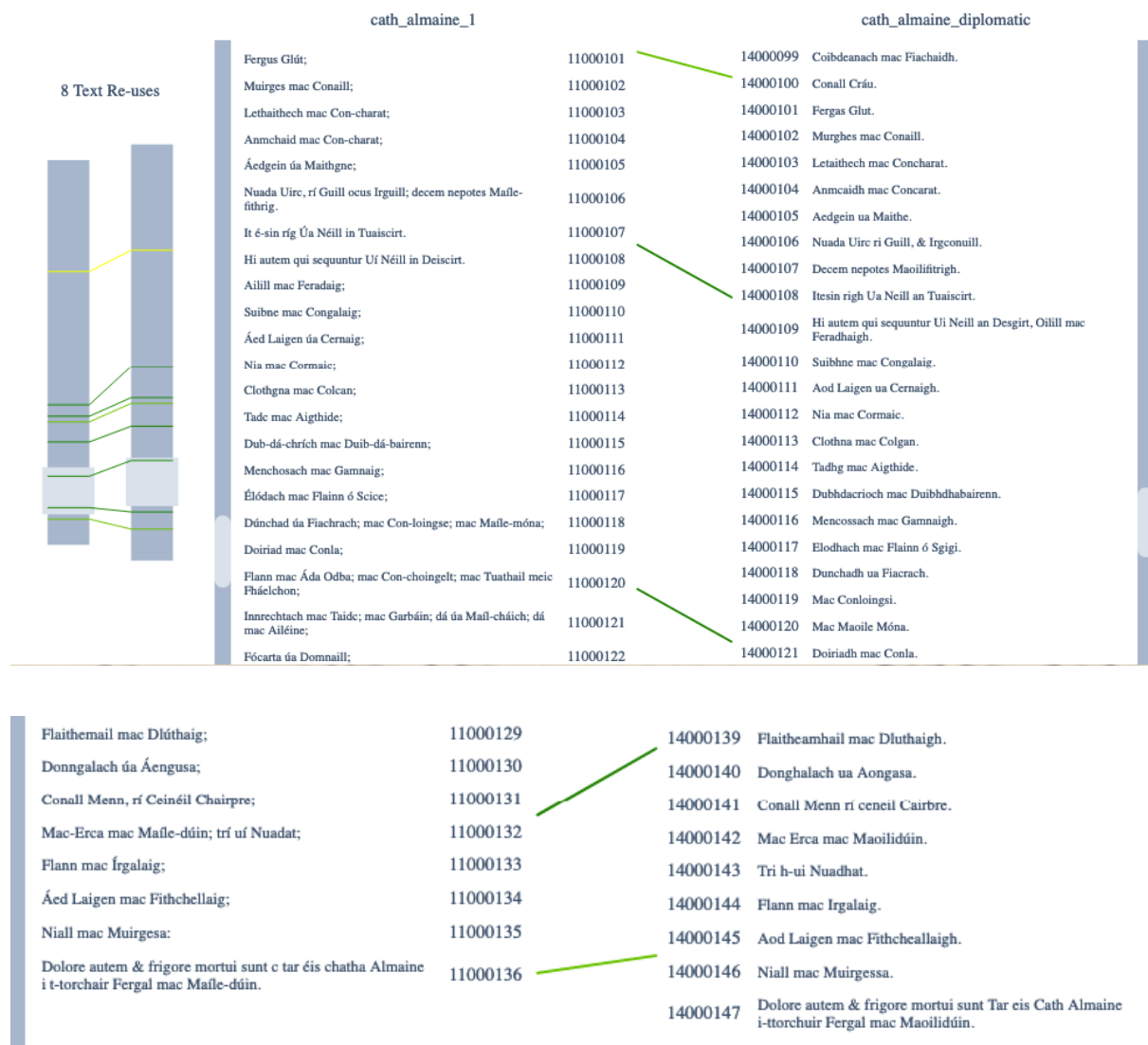  ○ No diacritic removal

# Text reuse browser in Traviz



A graph for each re-use case (appears when you click on a point on the scatterplot)

## Dynamic alignment



### Paths

Names of the elements in red will be different from experiment to experiment, although structured in the same way.

#### To text reuse scores

tracer > data > corpora > your_corpus_name > TRACER_DATA > experiment_parameters > algorithm_chosen > second folder in a list > first folder in a list > threshold_chosen > **.score file**

#### To html with visualisation

tracer > data > corpora > your_corpus_name > TRACER_DATA > experiment_parameters > algorithm_chosen > second folder in a list > first folder in a list > threshold_chosen > -TraVizPosprocessingImpl- > html > **index.html**