# TRACER TUTORIAL: TEXT REUSE DETECTION

## Preprocessing

Marco Büchler, Emily Franzini and Greta Franzini

**TRACER**

**TEXT REUSE** DETECTION MACHINE

**eTRAP**

Electronic Text Reuse Acquisition Project

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN

## TABLE OF CONTENTS

1. Download **TRACER** from `http://etrap.eu/tracer/` to your storage folder, e.g. `/roedel/mbuechler`
2. Using the command line, navigate to your storage folder with the `cd` command
3. **Unzip** archive: `gunzip tracer.tar.gz`
4. **Untar** archive: `tar -xvf tracer.tar`
5. Change to the TRACER folder: `cd TRACER`
6. **Open the configuration file** with `vim conf/tracer_config.xml`
7. Configure your input file:

```xml
<property name="SENTENCE_FILE_NAME" value="data/corpora/Bible/KJV.txt" />
```

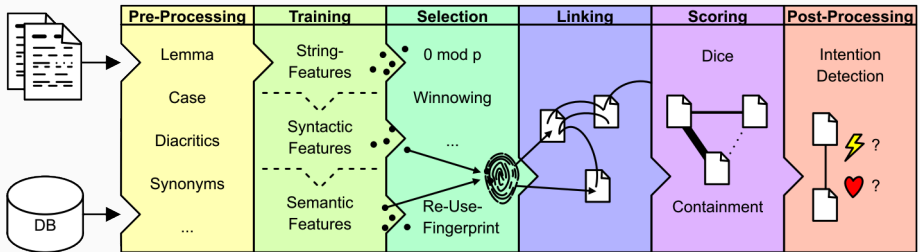Start the tool with the command:

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

**Explanation:**

- `-Xmx600m` (up to 600 MB memory);
- `-Dfile.encoding` sets the encoding of your input file (optionally);
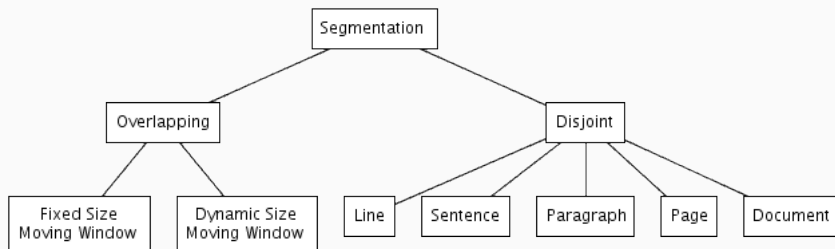- `-Deu.etrap.medusa.config.ClassConfig` (configuration file).
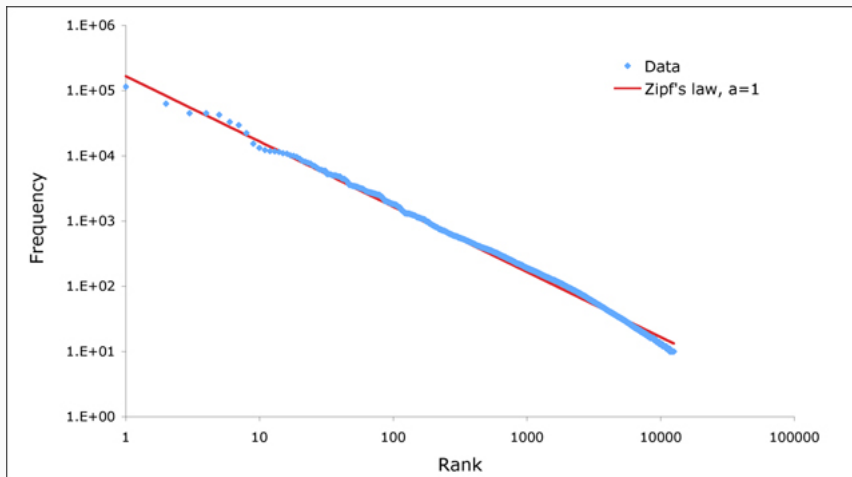
# WHAT IS PREPROCESSING?

**2** *De titulo vide Ag 155 sqq.*     **3** omnis homines \**Char.
gramm. I* 149, 17 *Diom. gramm. I* 305, 29   omnis . . . student *Prisc.
gramm. II* 358, 15   omnis . . . praestare *Non. p.* 371, 11   omnis . . .
animalibus *Char.     gramm. I 140,* 1 *Eugraph. Ter. Eun.* 232
. omneis *Char.* omnes *Eugraph.* qui . . . animalibus *Arus.
gramm. VII* 508, 4   praestare ceteris animalibus *Diom. gramm. I
* 813, 11   **5** pecora . . . finxit *Arus. gramm. VII* 496, 27   quae
. . . finxit *Non. p.* 309, 11 *Victorin. rhet. p.* 160, 36 *Prisc.
gramm. III* 370, 18   ventri oboedientia *Sen. epist.* 8 (60), 4   oboe-
dientes *Sen.*     **6** sed . . . sita est *Serv. Aen.* 2, 452 *georg.* 1, 198
   sed . . . utimur *Lact. inst.* 2, 12, 12     **7** animi . . . utimur *Hier.
ad Gal.* 5, 16 *p.* 410 ad *Eph.* 5, 33 *p.* 537     animi . . . commune
est\* *Hier. adv. Iovin.* 2, 10 *Aug. civ.* 9, 9   animae *Hier. adv. Iovin.*
   **8** utimur] vivere *Hier. ad Gal.* alterum nobis . . . commune est
*Serv. Aen.* 5, 81   **9** videtur] esse videtur X N M T m videtur esse
B K H D F l s n   **10** et . . . efficere *Victorin. rhet. p.* 160, 33   **17** nam
. . . opus est *Don. Ter. Andr.* 334 *Prisc. gramm. III* 226 3 288 , 17

What do you associate with **preprocessing**?

- **Approx. 50%** of all words occur only once
- **Approx. 16%** of all words occur only twice
- **Approx. 8%** of all words occur three times
- ...
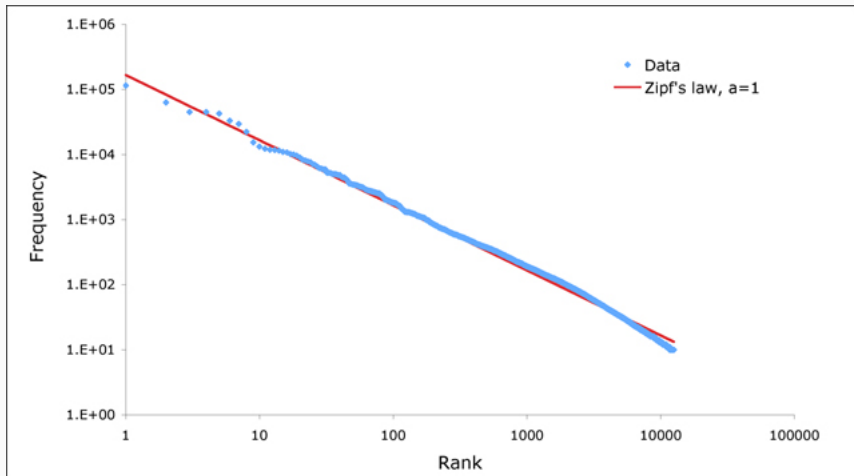- **Approx. 90%** of all words in a corpus occur 10 times or less

$$s(f) = \frac{1}{f * (f + 1)} \qquad s^n(f) = \sum_{f=1}^{n} \frac{1}{f * (f + 1)}$$

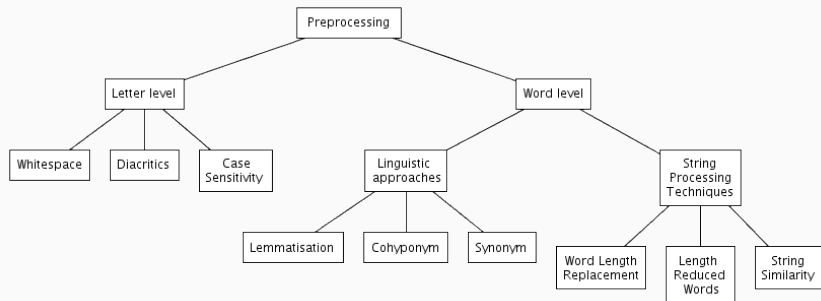- The **top 300-700 most** frequent words cover already about **50%** of all tokens (depending language)

What does lemmatisation mean for this plot?

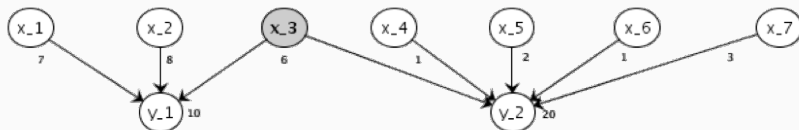# PREPROCESSING TECHNIQUES

E.g. lemmatisation
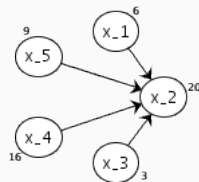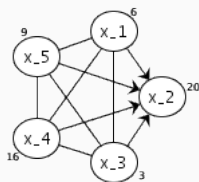
E.g. synonyms, string similarity

# HACKING

**Tasks:**

- Run on your texts ...
    1. ... without preprocessing
    2. ... 1) + lemmatisation
    3. ... 2) + synonym replacement

**Questions:**

- Compare the input file with the `*.prep` file for all preprocessing techniques. Which methods seem to work best for you? Which make no sense for the dataset?
- Compare all `*.meta` files containing some numbers! How many words have changed and through which method?
- (optional and advanced) What is the number of word types for each preprocessing technique (can be derived from the first column of `*.prep.inv`).

```xml
<category name="eu.etrap.tracer.preprocessing.WordLevelPreprocessingImpl">
    <property name="boolLemmatisation" value="false" />
    <property name="boolReplaceSynonyms" value="false" />
    <property name="boolReplaceStringSimilarWords" value="false" />
    <property name="boolRemoveDiachritics" value="false" />
    <property name="boolMakeAllLowerCase" value="false" />
    <property name="boolReplaceWordByWordLength" value="false" />
    <property name="boolReplaceByReducedString" value="false" />
    <property name="intMinWordLengthThreshold" value="5" />
    <property name="intNGramSize" value="5" />
    <property name="weigthByLogLikelihoodRatio" value="false"/>;
</category>
```

**Hint:**

- The configuration file can be found in:
  $TRACER_HOME/conf/tracer_conf.xml
- All values show `false`.

```xml
<category name="eu.etrap.tracer.preprocessing.WordLevelPreprocessingImpl">
    <property name="boolLemmatisation" value="false" />
    <property name="boolReplaceSynonyms" value="false" />
    <property name="boolReplaceStringSimilarWords" value="false" />
    <property name="boolRemoveDiachritics" value="true" />
    <property name="boolMakeAllLowerCase" value="false" />
    <property name="boolReplaceWordByWordLength" value="false" />
    <property name="boolReplaceByReducedString" value="false" />
    <property name="intMinWordLengthThreshold" value="5" />
    <property name="intNGramSize" value="5" />
    <property name="weigthByLogLikelihoodRatio" value="false"/>;
</category>
```

**Hint:**

- `boolRemoveDiachritics` is switched on by value `true`.

```xml
<category name="eu.etrap.tracer.preprocessing.WordLevelPreprocessingImpl">
    <property name="boolLemmatisation" value="true" />
    <property name="boolReplaceSynonyms" value="false" />
    <property name="boolReplaceStringSimilarWords" value="false" />
    <property name="boolRemoveDiachritics" value="true" />
    <property name="boolMakeAllLowerCase" value="false" />
    <property name="boolReplaceWordByWordLength" value="false" />
    <property name="boolReplaceByReducedString" value="false" />
    <property name="intMinWordLengthThreshold" value="5" />
    <property name="intNGramSize" value="5" />
    <property name="weigthByLogLikelihoodRatio" value="false"/>;
</category>
```

**Hint:**

- `boolLemmatisation` is switched on by value `true`.

- **Lemmatisation** can be configured by:
  `<property name="BASEFORM_FILE_NAME"`
  `value="data/corpora/Bible/Bible.lemma" />`

```xml
<category name="eu.etrap.tracer.preprocessing.WordLevelPreprocessingImpl">
    <property name="boolLemmatisation" value="true" />
    <property name="boolReplaceSynonyms" value="true" />
    <property name="boolReplaceStringSimilarWords" value="false" />
    <property name="boolRemoveDiachritics" value="true" />
    <property name="boolMakeAllLowerCase" value="false" />
    <property name="boolReplaceWordByWordLength" value="false" />
    <property name="boolReplaceByReducedString" value="false" />
    <property name="intMinWordLengthThreshold" value="5" />
    <property name="intNGramSize" value="5" />
    <property name="weigthByLogLikelihoodRatio" value="false"/>;
</category>
```

**Hint:**

- `boolReplaceSynonyms` is switched on by value `true`.
- **Synonyms** can be configured by:
  `<property name="SYNONYMS_FILE_NAME" value="data/corpora/Bible/Bible.syns" />`

```xml
<category name="eu.etrap.tracer.preprocessing.WordLevelPreprocessingImpl">
    <property name="boolLemmatisation" value="true" />
    <property name="boolReplaceSynonyms" value="true" />
    <property name="boolReplaceStringSimilarWords" value="true" />
    <property name="boolRemoveDiachritics" value="true" />
    <property name="boolMakeAllLowerCase" value="false" />
    <property name="boolReplaceWordByWordLength" value="false" />
    <property name="boolReplaceByReducedString" value="false" />
    <property name="intMinWordLengthThreshold" value="5" />
    <property name="intNGramSize" value="5" />
    <property name="weigthByLogLikelihoodRatio" value="false"/>;
</category>
```
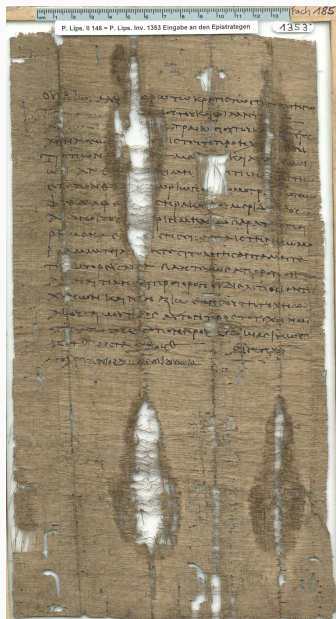
**Hint:**

- `boolReplaceStringSimilarWords` is switched on by value `true`.

- Thresholds:
  `<property name="SYNONYMS_FILE_NAME"`
  `value="data/corpora/Bible/Bible.syns" />`

P. Lips. II 146 = P. Lips. Inv. 1353 Eingabe an den Epistrategen

Οὐιβίῳ Ἀλεξά[ν]δρῳ τῷ κρατίστῳ ἐπιστρατήγῳ
παρὰ Ἀντ[ωνίου Δ]όμγου τοῦ καὶ Φιλαντι[νό]ου
Ἀντωνίο[υ Ῥωμανο]ῦ Τραιανείου τοῦ κα[ὶ Στρα]τείου
Ἀντινοέως. [οὐκ ἂν] εἰς τοῦτο προήχθ[η]γ, ἐπι-
τρόπων [μέγιστ]ε, μέ[τριος] καὶ ἀπράγμων
ὢν ἄνθρ[ωπος,] εἰ μὴ [ὕβρι]ν τὴν μ[εγ]ίστην
ἐπεπόνθ[ειν ὑπὸ] Ὠρίωνο[ς κ]ωμογρα[μ]ματέως
Φ[ι]λαδελφεί[ας τῆ]ς Ἡρακλείδου μερίδο[ς] τοῦ
Ἀρσινοΐτου. [οὗ χά]ριν μην[ύ]ω παρὰ τ[ὰ ἁ]πει-
ρημένα ἑα[υτὸ]ν ἐνσείσαντα εἰς τὴν κωμο-
γραμματείαγ [μ]ήτε σιτολογήσαντα μήτε
πρ[α]κτορεύσαντα παντελῶς ἄπορον ὄν[τ]α.
δι᾿ ἣν αἰτίαν καὶ πρότερον οὐ διέλιπον ἐντυγ-
χάνων καὶ νῦν ἀξιῶ, ἐάν σου τῇ τύχῃ δόξ[ῃ],
ἀκοῦσαί μου π[ρ]ὸς αὐτὸν πρὸς τὸ τυχεῖν με
τῆς ἀπὸ σοῦ [μι]σοπονήρου ἐγδ[ι]κίας, ἵν᾿ ὦ ὑπὸ [σ]οῦ
κατὰ πάντα βεβοηθ(ημένος). διευτύχει
Ἀντώνιος Δόμγος ἐπιδέδωκα.

# CONCLUSION AND REVISION

**Statement:**

- "My lemmatisation tool $<\texttt{XYZ}>$ is able to compute the base forms of 80% of all tokens in a corpus."

**Good or bad?**

**Fact file:**

- Language variants
- Different writing styles
- (Some) dialects
- Diachritics
- OCR errors

**Question:** What's the difference for you?

**Fact file:**

- Language variants
- Different writing styles
- (Some) dialects
- Diachritics
- OCR errors

**Question:** What do you think is the difference for the computer?

- Cleaning and harmonising the data.
- When working with a new corpus -not only the language but also the same language in different epochs or geographical regions- cleaning/harmonising the data can take up to 70% of the overall time.

**Preprocessing mantra:**
Garbage in, garbage out

# Finito!

## Team

Marco Büchler, Greta Franzini and Emily Franzini.

## Visit us

🌐 `http://www.etrap.eu`

✉ `contact@etrap.eu`

TRACER
TEXT REUSE DETECTION MACHINE

eTRAP
Electronic Text Reuse Acquisition Project

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN

SPONSORED BY THE

Federal Ministry
of Education
and Research

The theme this presentation is based on is licensed under a Creative Commons Attribution-ShareAlike 4.0 International License. Changes to the theme are the work of eTRAP.