# corpus.tools

Jan Michelfeit

Lexical 文 Computing

jan.michelfeit@sketchengine.co.uk

6th Sketch Engine Workshop
Herstmonceux, August 10, 2015

# Our strategy

- academia a very fertile environment

# Our strategy

- academia a very fertile environment
- but projects often end up abandoned

# Our strategy

- academia a very fertile environment
- but projects often end up abandoned
- adopt useful tools and continue developing and supporting them

# Our strategy

- academia a very fertile environment
- but projects often end up abandoned
- adopt useful tools and continue developing and supporting them
- business-friendly open source licences

- "boilerplate" content removal tool
- boilerplate = headers, footers, navigation etc. in web pages
- leaves only large contiguous chunks of text (whole paragraphs)

# Chared



- character encoding detection tool
- webpage metadata often not reliable
- a small amount of noise can pollute keyword reports
- languge-aware detection more successful than language-agnostic

# SpiderLing



- $=$ web-SPIDER for LINGuistics
- utilizes both JusText and Chared
- prioritizes sites rich in text
- goal: maximize number of words per megabyte downloaded

# Onion



- $=$ ONe Instance ONly
- duplicate (and near-duplicate) removal tool
- based on word n-grams
- typically runs on paragraphs
- similarity threshold, structure and *n* configurable

# Unitok



- $=$ UNIversal TOKenizer
- based on regular expressions
- URLs, e-mail addresses
- abbreviations, clictics (do – n't, ca – n't)
- configuration files for a variety of languages
- versioning promotes replicability
- includes uninorm, UNIversal NORMalizer

# NoSketch Engine



- like Sketch Engine, but without Word Sketches
- manatee – backend, storage, indexing and query evaluation
- bonito – user interface for querying
- corpus architect – user corpora – not included
- google group for support