# Automated Subject Indexing and Classification using Annif

Osma Suominen

HELDIG Summit, 23 October 2018

# Extrablad till
# ÅBO UNDERRÄTTELSER

**N:o 2.**

## Finlands oavhängighet.

Landtdagen omfattar regeringens proklamation om Finlands fullständiga oavhängighet och ansluter sig till hufvudprinciperna i regeringens program för tryggandet af landets nya ställning. Beslutet fattades med 100 röster emot 88 vilka tillföllo ett av socialdemokraterna formulerat förslag.

Pris 25 penni.

YSA
Allärs

YSO
KOKO

€£$

# Subject indexing is a hard problem



long tail

## for **humans:**

- **Subjectivity**: when two people index the same document, only ~⅓ of the subjects are the same

- **Many concepts**: tens of thousands of concepts to pick from

- **Vocabulary changes:** new concepts are added, existing ones are renamed and redefined

## for **machines:**

- **Long tail phenomenon:** even with large amounts of training data, most subjects are only used a small number of times

- **Many concepts**: requires complex models that are computationally intensive

- **Difficult to evaluate**: hard to tell "somewhat bad" answers from really wrong ones without human evaluation

- **Vocabulary changes**: models must be retrained

# Approach

**Automating our own processes**

**Creating generic tools for many contexts**



**vs.**

# Enter **Annif**

*Feed your subject indexing robot with bibliographic metadata!*

# Machine learning requires training data

Give feedback

# FINNA.FI

## The material of Finnish archives, libraries and museums with a single search

ⓘ

| Find... | All fields | ▾ | 🔍 |

Q⁺ Advanced Search

## Finna Street

Uncover historical images of where you are!

✛ Show Images

## Search...

Images

Unrestricted collections

📄 Archive collections

📖 Library collections

🏺 Museum collections

Shortcuts:

Other Finna websites

Organisations providing Finna content

## Collection highlights

# Finna API



All Finna metadata is  !

TITLE  General Finnish upper ontology YSO
       YSO - General Finnish ontology

SUBJECT  general concepts

DESCRIPTION  General Finnish Upper Ontology YSO is a trilingual ontology consisting mainly of general concepts. YSO has been founded on the basis of concepts in Finnish cultural sphere. As an indexing tool it is best applicable when indexed material is interdiscliplinary and its themes vary to a great extent.

**Resource counts by type**

| Type | Count |
|------|-------|
| Concept | 29031 |
| • Individual concept | 1890 |
| • Hierarchical concept | ~~101~~ |
| • General concept | 25940 |
| Collection | 241 |

**~30 000 concepts** that can be used for subject indexing

**Term counts by language**

| Language | Preferred terms | Alternate terms | Hidden terms |
|----------|-----------------|-----------------|--------------|
| English | 28566 | 3245 | 11657 |
| Finnish | 29019 | 11491 | 14288 |
| Swedish | 28582 | 13072 | 11079 |

W3C SKOS

# Annif prototype (2017)

# Indexing Wikipedia by topics

Finnish Wikipedia has 410 000 articles (620 MB as raw text)

Automated subject indexing took 7 hours on a laptop, using the Annif prototype

1-3 topics per article (average ~2)

# Indexing Wikipedia by topics

Finnish Wikipedia has 410 000 articles (620 MB as raw text)
Automated subject indexing took 7 hours on a laptop
1-3 topics per article (average ~2)

**Examples:** (random sample)

| Wikipedia article | YSO topics |
|---|---|
| Ahvenuslammi (Urjala) | shores |
| Brasilian Grand Prix 2016 | race drivers, formula racing, karting |
| Guy Topelius | folk poetry researcher, saccharin |
| HMS Laforey | warships |
| Liigacup | football, football players |
| Pää Kii | ensembles (groups), pop music |
| RT-21M Pioneer | missiles |
| Runoja | pop music, recording (music recordings), compositions (music) |
| Sjur Røthe | skiers, skiing, Nordic combined |
| Veikko Lavi | lyricists, comic songs |

# Most common topics in Finnish Wikipedia

# Most common topics in Finnish Wikipedia



Chart categories (top to bottom):
- football players
- ice hockey players
- sports teams
- warships
- football
- pop music
- ensembles (groups)
- recording (music recordin…
- Olympics
- Finnish championships
- heavy rock
- solemnising a marriage
- matches (sports contests)
- Nordic combined
- summer Olympics
- track and field athletes
- formula racing
- goalkeepers
- film directors
- ice hockey

X-axis: 0, 2 000, 4 000, 6 000, 8 000

# People vs. Robots Workshop



20 documents
40 librarians
45 minutes

...

**225 indexing results**
- 11 per document
- 5.5 per person

# Average similarity of subject sets

## 33.39 %

Using Rolling similarity, a.k.a. F1 score, to compare subject sets

# Annif prototype vs. new Annif

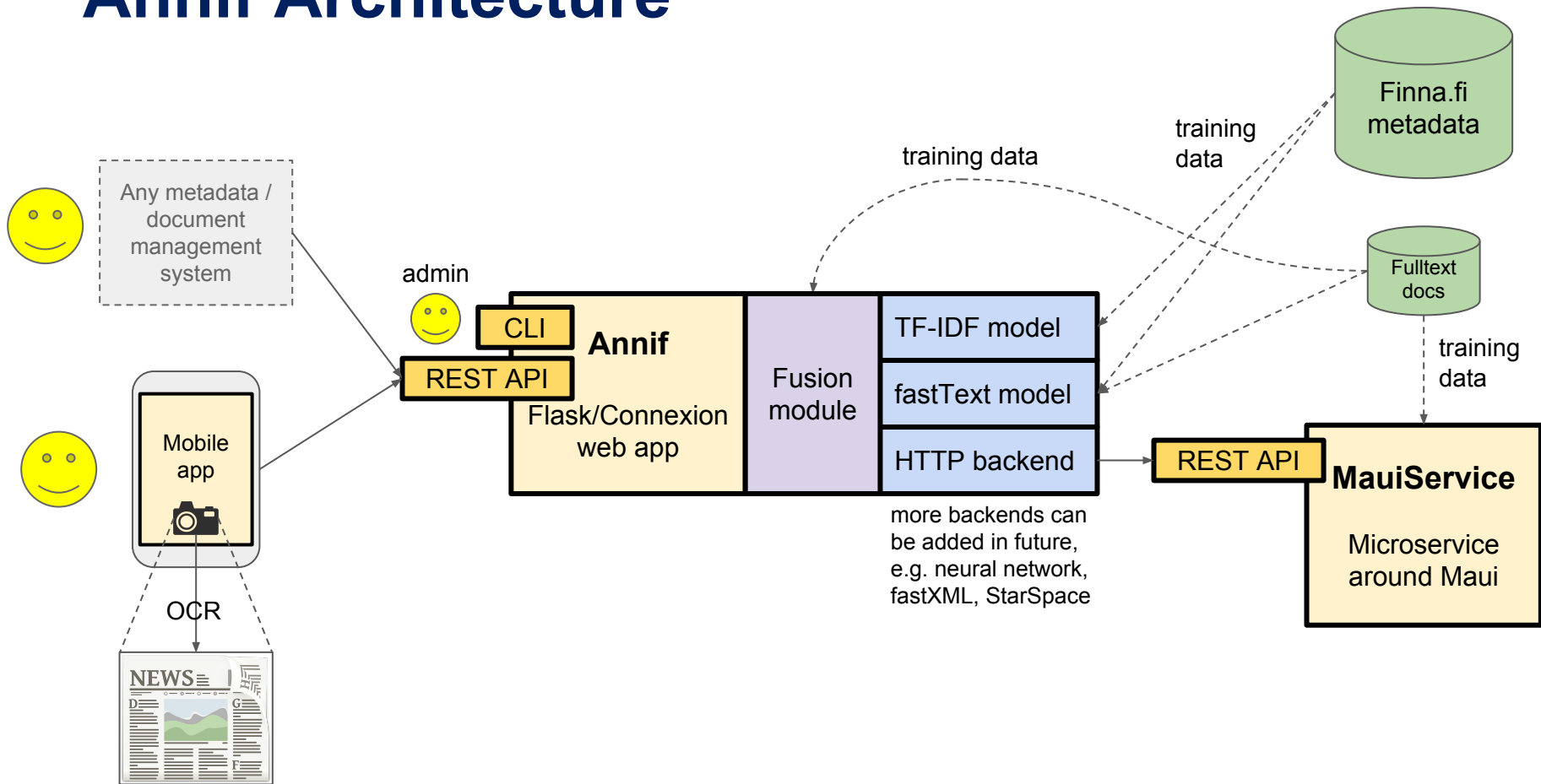| | Prototype (2017) | New Annif (2018→) |
|---:|---|---|
| *architecture* | loose collection of scripts | Flask web application |
| *coding style* | quick and dirty | solid software engineering |
| *backends* | Elasticsearch index | TF-IDF, fastText, Maui ... |
| *language support* | Finnish, Swedish, English | any language supported by NLTK |
| *vocabulary support* | YSO, GACS ... | YSO, YKL, others coming |
| *REST API* | minimal | extended (e.g. list projects) |
| *user interface* | web form for testing | http://dev.annif.org |
| *mobile app* | HTML/CSS/JS based | (native Android app?) |
| *open source license* | CC0 | Apache License 2.0 |

# Annif Architecture

# Backends / Algorithms

- **TF-IDF similarity**

  Baseline bag-of-words similarity measure. Implemented with the Gensim library.

- **fastText** by Facebook Research

  Machine learning algorithm for text classification.

  Uses word embeddings (similar to word2vec) and resembles a neural network architecture.

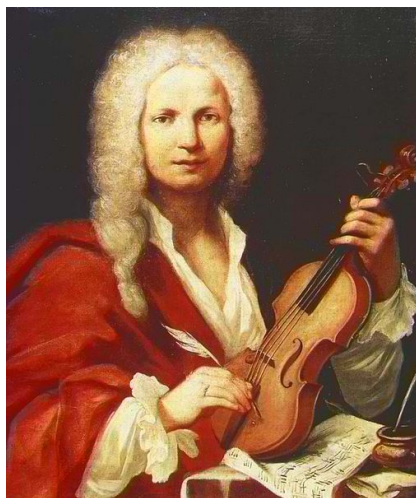  Promises to be good for e.g. library classifications (DDC, UDC, YKL…)

- **HTTP backend** for accessing MauiService REST API

  MauiService is a microservice wrapper around the Maui automated indexing tool.

  Based on traditional Natural Language Processing techniques - finds terms within text.

# Backend configuration

**Backends may be used alone, or in combinations (ensembles)**





**Current challenge:** Which fusion method works best for combining results from multiple backends?

An experiment testing different fusion methods

# Command line interface

**Load a vocabulary to be used by one or more models:**

```
$ annif loadvoc yso-en yso-en.tsv
```

**Train a model:**

```
$ annif train tfidf-en yso-finna-en.tsv.gz
```

**Analyze a document:**

```
$ annif analyze tfidf-en <berries.txt
<http://www.yso.fi/onto/yso/p772>     strawberry            0.39644203283656165
<http://www.yso.fi/onto/yso/p18109>   wild strawberry       0.37539359094384245
<http://www.yso.fi/onto/yso/p25548>   stolons               0.3261554545369906
<http://www.yso.fi/onto/yso/p6749>    berry cultivation     0.2394291077460799
<http://www.yso.fi/onto/yso/p10631>   questionnaire survey  0.22714475653823335
<http://www.yso.fi/onto/yso/p6821>    farms                 0.21725243067995587
<http://www.yso.fi/onto/yso/p3294>    customers             0.216395821347059
<http://www.yso.fi/onto/yso/p1834>    work motivation       0.21612376226244975
<http://www.yso.fi/onto/yso/p8531>    customership          0.21536113638508098
<http://www.yso.fi/onto/yso/p19047>   corporate clients     0.21412270159920782
```

**Evaluate a model using several measures (e.g. recall, precision, F1 score, NDCG):**

```
$ annif eval tfidf-en directory-with-gold-standard-docs/
```
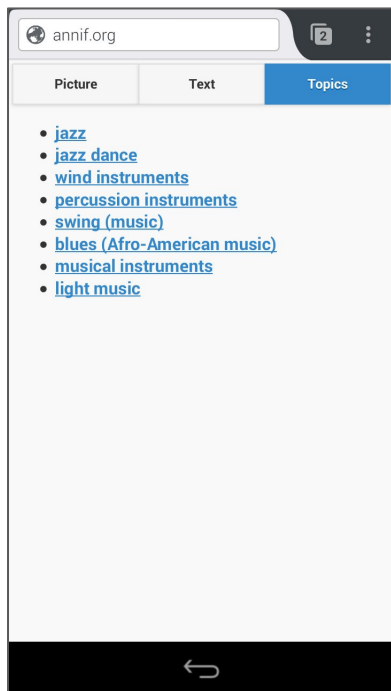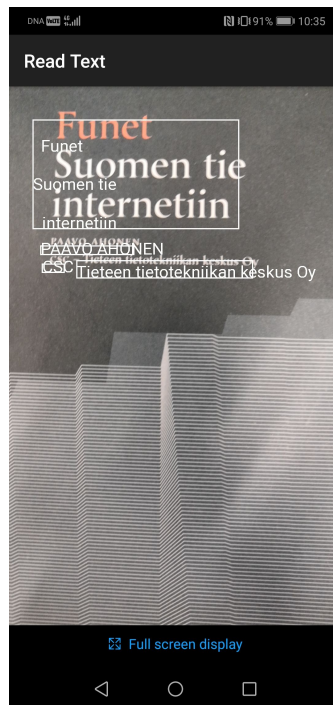
# REST API

## Main operations:

| | |
|---|---|
| GET /projects/ | list available projects |
| GET /projects/<project_id> | show information about a project |
| POST /projects/<project_id>/analyze | analyze text and return subjects |
| POST /projects/<project_id>/explain | analyze text and return subjects, with explanations indicating why they were chosen |
| POST /projects/<project_id>/train | train the model by giving a document and gold standard subjects |

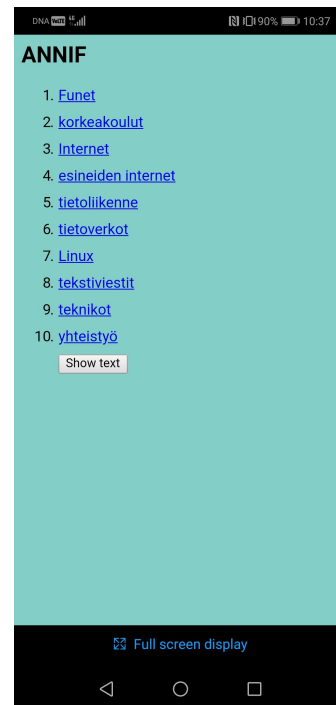Defined using a Swagger / OpenAPI specification

# Mobile apps



Prototype web app,
ocr.space cloud OCR
m.annif.org



Prototype Android app with OCR on the device
(by Okko Vainonen)

# Test corpora

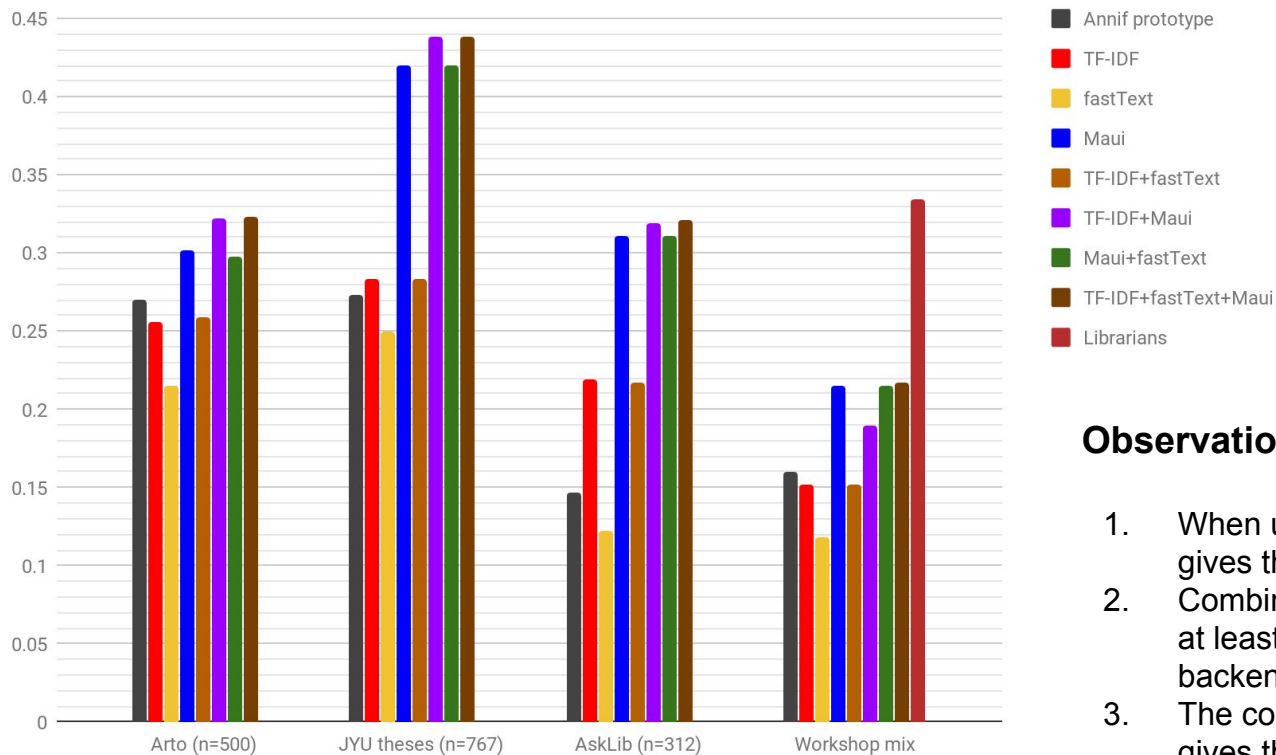Full text documents indexed with YSA/YSO for training and evaluation

- Articles from Arto database (n=6287)
  Both scientific research papers and less formal publications. Many disciplines.

- Master's and Doctoral theses from Jyväskylä University (n=7400)
  Long, in-depth scientific documents. Many disciplines.

- Question/Answer pairs from an Ask a Librarian service (n=3150)
  Short, informal questions and answers about many different topics.

Available on GitHub: https://github.com/NatLibFi/Annif-corpora
(for the first two corpora, only links to PDFs are provided for copyright reasons)

# Evaluation of different backends

F-measure similarity scores against a gold standard



**Observations:**

1. When using just one backend, Maui often gives the best results
2. Combinations (ensembles) usually give at least as good results as single backends
3. The combination of all three backends gives the best results

# Annif on GitHub

Python 3.5+ code base
Apache License 2.0

Fully unit tested (98% coverage)
PEP8 style guide compliant
Usage documentation in the wiki

https://github.com/NatLibFi/Annif

# Apply Annif on your own data!

Choose an indexing vocabulary

Prepare a corpus from your existing metadata

Load the corpus into Annif

Use it to index new documents

# **Lessons learned** (so far)

1. Good quality training data is key for training and evaluation
   Don't expect good results if you don't have the data it takes

2. Gold standard subjects are useful, but human evaluation is necessary
   Subject indexing is inherently subjective; comparing to a single gold standard can be misleading

3. All algorithms have strong and weak points
   Combinations work better than any algorithm by itself

4. Surprising amount of interest also from non-library organizations
   Archives, media organizations, book distributors … automation is better done together!

# Thank you!
## Questions?

osma.suominen@helsinki.fi  -  @OsmaSuominen

Website: http://annif.org
Code: https://github.com/NatLibFi/Annif
Test corpora: https://github.com/NatLibFi/Annif-corpora

These slides: **https://tinyurl.com/annif-heldig**