



# **Annif and Finto AI: DIY automated subject indexing from prototype to production**

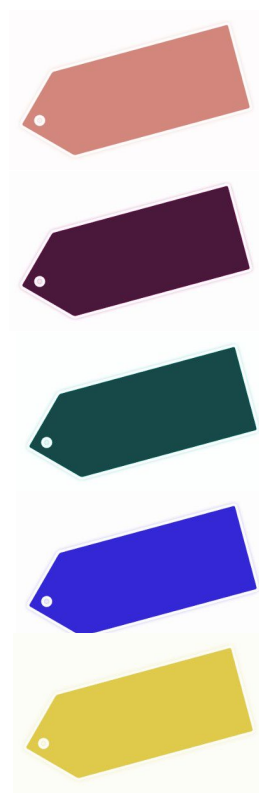
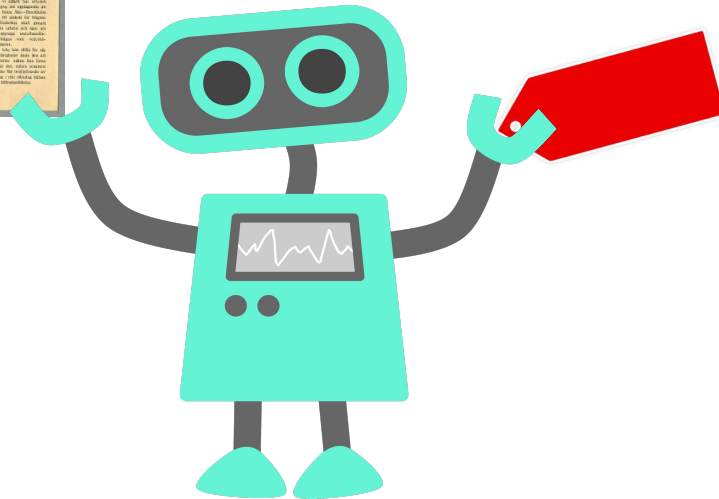
Osma Suominen, Mona Lehtinen, Juho Inkinen  
SWIB20, 23 November 2020



# Outline

1. Development of Annif
2. Quality of automated subject indexing
3. Community building
4. Annif deployments
5. Lessons learned

# 1. Development of Annif



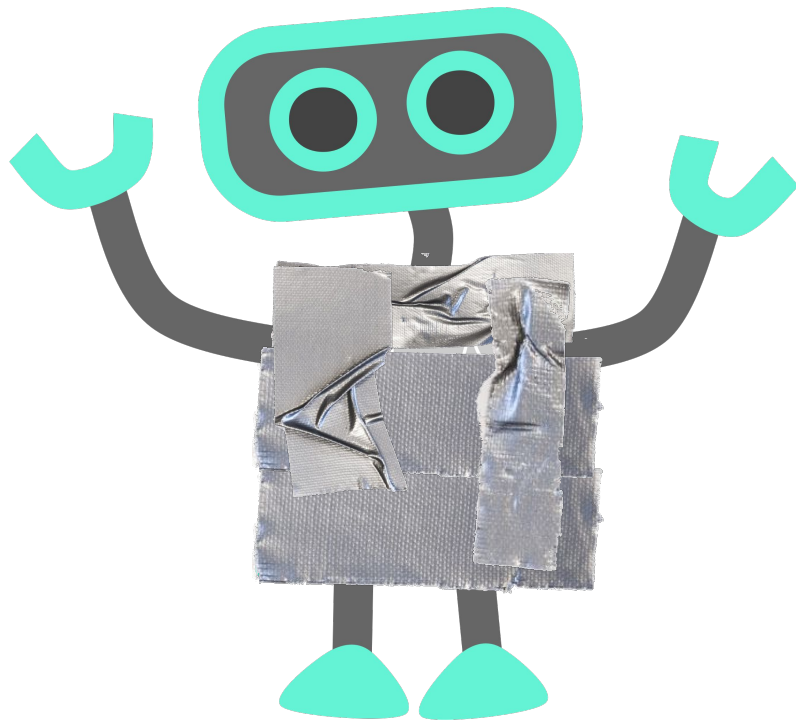
**YSO**, General Finnish Ontology  
with 30,000+ subjects

# Machine learning using existing metadata



annif

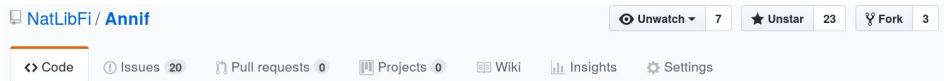
Early prototype (2017) got people excited



# Starting points for Annif implementation (2018 → )

1. multilingual
2. independent of indexing vocabulary
3. support different subject indexing algorithms
4. CLI, Web user interface and REST API
5. community-oriented open source





Annif is a multi-algorithm automated classification and subject indexing tool for libraries, archives and museums. This repository is used for developing a production version of the system, based on ideas from the initial prototype. <http://annif.org> Edit

[subject-indexing](#) [python](#) [machine-learning](#) [code4lib](#) [classification](#) [rest-api](#) [flask-application](#) [connexion](#) [Manage topics](#)

766 commits 7 branches 48 releases 5 contributors [View license](#)

Branch: master New pull request

[Create new file](#) [Upload files](#) [Find file](#) [Clone or download](#)

osma add Zenodo DOI badge	Latest commit d832514 2 days ago
<a href="#">annif</a>	refactor: split off JSON input to document corpus conversion in rest ... 2 days ago
<a href="#">tests</a>	CLI unit test for trying to learn when backend doesn't support it 2 days ago
<a href="#">.codeclimate.yml</a>	more comprehensive Code Climate configuration a year ago
<a href="#">.codecov.yml</a>	Codecov should ignore setup.py 10 months ago
<a href="#">.coveragerc</a>	Generate Codecov reports 2 years ago
<a href="#">.gitignore</a>	Add virtualenv (default? de-facto?) folder to gitignore 15 days ago
<a href="#">.lgtm.yml</a>	Add LGTM configuration excluding fasttext 5 months ago
<a href="#">.scrutinizer.yml</a>	Try to fix pipenv/pip compatibility issue pypa/pipenv#2924 within Scr... 5 months ago
<a href="#">.travis.yml</a>	install deb packages using apt addon (even though they're unnecessary... a month ago

# Annif on GitHub

Python 3.6+ code base

Apache License 2.0

Fully unit tested (99% coverage)

PEP8 style guide compliant

<https://github.com/NatLibFi/Annif>



**annif 0.49.0**

```
pip install annif
```

Automated subject indexing and classification tool

Python package on PyPI



**natlibfi / annif**

Docker images on Quay.io



## 2. Quality of automated subject indexing

# Document collections for training and evaluation

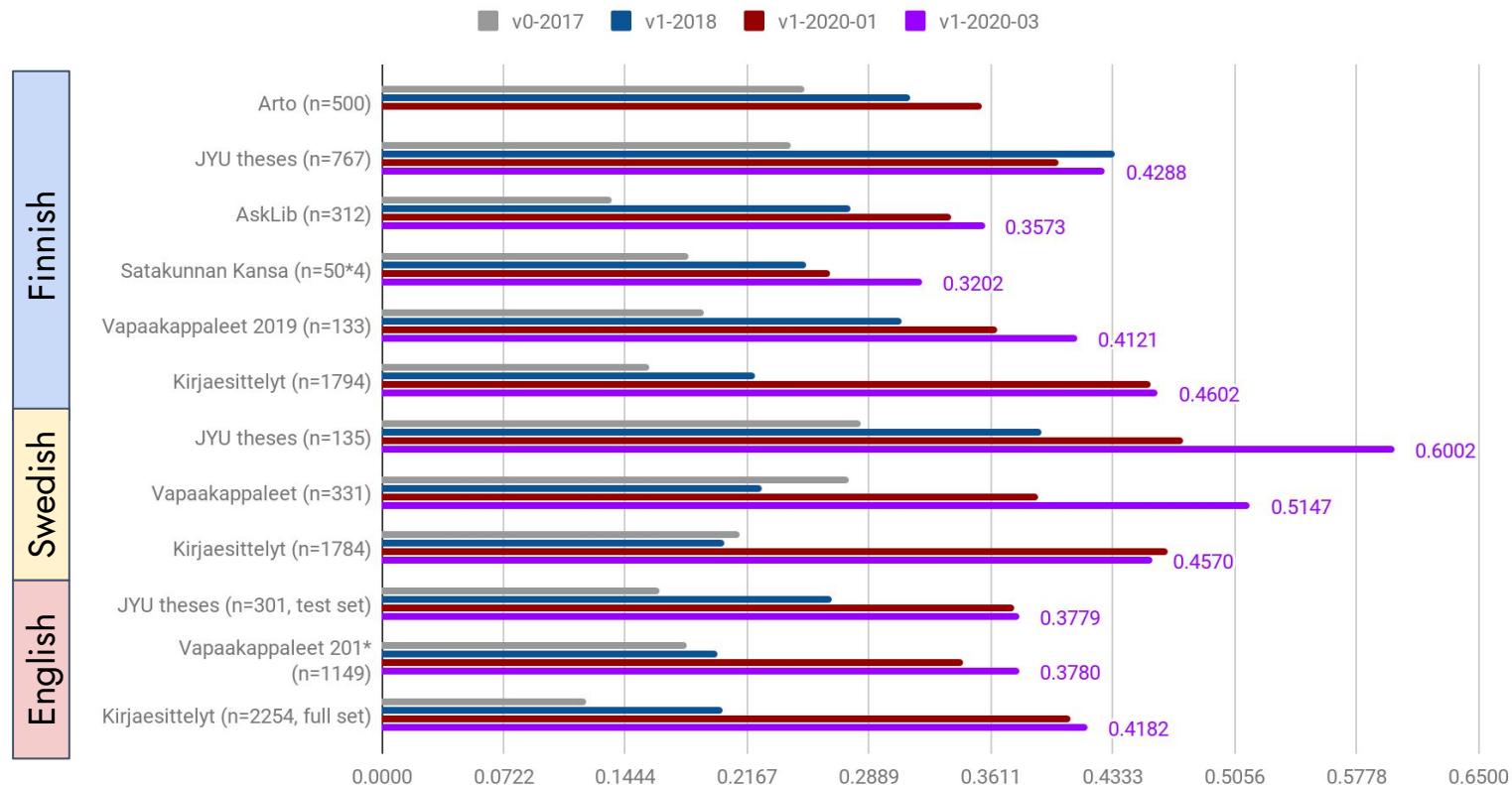
1. Metadata records from Finna.fi discovery system
2. Ask a Librarian question-answer pairs
3. Master's and Doctoral theses from University of Jyväskylä
4. Book descriptions from publishers (via Kirjavälitys Oy)
5. E-books from our electronic deposit system
6. ...

Converted to Annif corpus format & split into train/validate/test subsets

The ones we could republish are in the [Annif-corpora](#) repository GitHub

# Comparison to “gold standard”

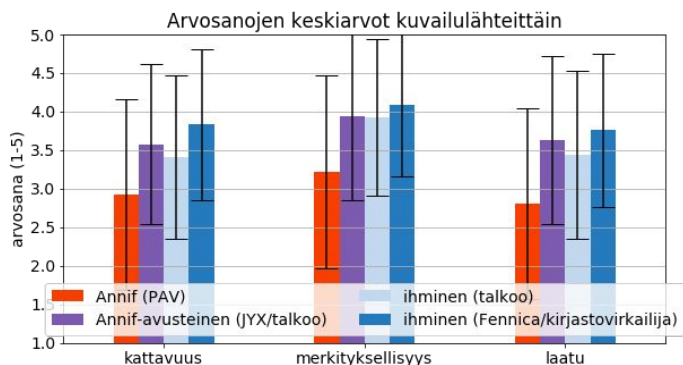
F1@5 scores for different test corpora and Annif API/model versions



# Assessment by evaluators

At a workshop in 2019, **48 evaluators** evaluated subjects for **50 documents**. Subjects were given by either human indexers or four different algorithms.

The best ensemble algorithm (red bars) was not quite on the level of human indexers in quality scores (left), and significantly more of its suggestions were rejected (right).



Poistettujen asiasanojen keskimääräiset %-osuudet kuvailulähteittäin

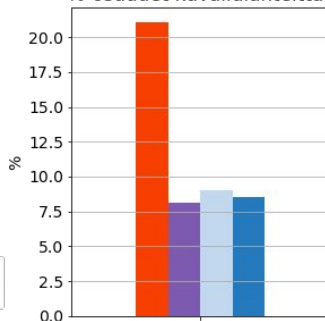


Photo: Mikko Lappalainen.

Lehtinen M., Inkinen J. & Suominen O. (2019). Aaveita koneessa: Automaattisen sisällönkuvailun arviointia Kirjastoverkkopäivillä 2019. [Tietolinja, 2019\(2\)](http://urn.fi/URN:NBN:fi-fe2019120445612). <http://urn.fi/URN:NBN:fi-fe2019120445612>

# Annif-Leiki Comparison at Finnish Broadcasting Company Yle

- Annif vs Leiki (commercial service) tagging compared by 28 human evaluators at Yle
- About 100 Finnish and Swedish articles and their tags
  - business, science, culture, sport

**Finnish:** Annif **slightly better** than Leiki

**Swedish:** Annif **substantially** better than Leiki

Arvioi asiasanojen osuvuutta

Voit katsoa asiasanan kuvauksen käymällä osoitteessa "[https://meta.api.yle.fi/v1/concepts.json?app\\_id=HIEKKALAATIKKO&app\\_key=HIEKKALAATIKKO&yle\\_id=YLE\\_ID](https://meta.api.yle.fi/v1/concepts.json?app_id=HIEKKALAATIKKO&app_key=HIEKKALAATIKKO&yle_id=YLE_ID)" ottamalla YLE\_ID:n asiasanan suluista

	essential keskeinen	ok ok	non relevant epärelevantti	wrong väärä
professorit (18-211223)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
aurinkokunnat (18-8786)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Markku Wilenius (18-295234)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
fossiiliset polttoaineet (18-2396)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
ilmasto (18-215534)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
tulevaisuus (18-211385)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Evaluation scale

Tags

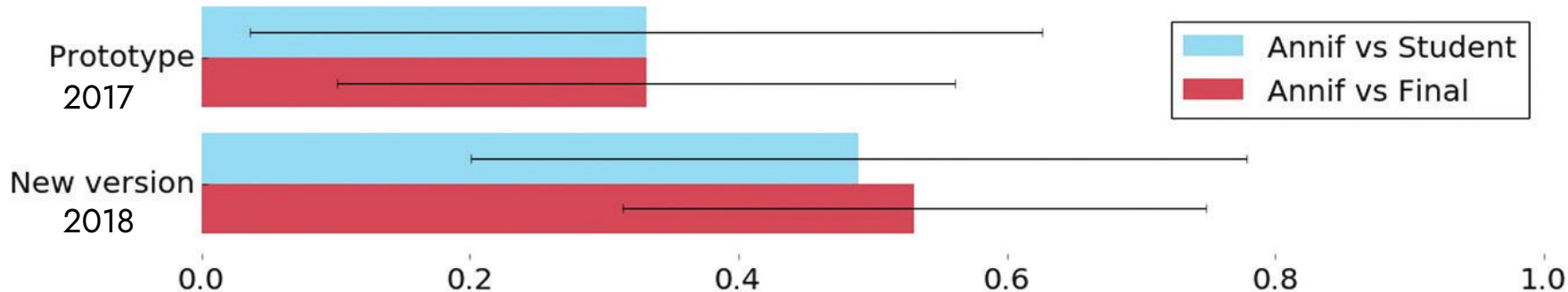
# Evaluating in the context of an indexing workflow

## JYX repository, University of Jyväskylä:

F1 similarity between Annif suggestions and the subjects

a) chosen by the student (blue)

b) confirmed by the JYX librarian (red)



Suominen, O., 2019. Annif: DIY automated subject indexing using multiple algorithms. *LIBER Quarterly*, 29(1), pp.1–25. DOI: <http://doi.org/10.18352/lq.10285>

### 3. Community building

# Web site with form for testing at [annif.org](https://annif.org)

## INPUT TEXT

SWIB focuses on Linked Open Data (LOD) in libraries and related organizations. It is well established as an event where IT staff, developers, librarians, and researchers from over the world meet and mingle and learn from each other. The topics of talks and workshops at SWIB revolve around opening data, linking data and creating tools and software for LOD production scenarios. These areas of focus are supplemented by presentations of research projects in applied sciences, industry applications, and LOD activities in other areas.

As usual, SWIB20 will be organized by the ZBW - Leibniz Information Centre for Economics and the North Rhine-Westphalian Library Service Centre (hbz). The conference language is English.

Would you like to share your experiences working on an interesting service, research topic or project – not just what you did, but also how you did it?

For this SWIB rendition we adjusted the formats to the online environment:

Presentations (15 minutes plus 5 q&a)

Practical workshops or tutorials (maximum 120 min)

www.annif.org

## PROJECT (VOCABULARY AND LANGUAGE)

YSO NN Ensemble English ▼

## MAX # OF SUGGESTIONS

10 15 20

Get suggestions →

**annif**

## SUGGESTED SUBJECTS

- data acquisition
- libraries
- linked open data
- information retrieval
- metadata
- data mining
- Finland
- information and communications technology
- information technology
- social media



# Wiki documentation on GitHub

- [issues](#)
- [pull requests](#)

Welcome to the Annif wiki!

- [Getting started](#)
- [System requirements](#)
- [Optional features and dependencies](#)
- [Usage with Docker](#)
- [Architecture](#)
- [Commands](#)
- [Web user interface](#)
- [Corpus formats](#)
  - [Document corpus formats](#)
  - [Subject vocabulary formats](#)
- [Project configuration](#)
- [Analyzers](#)
- [Achieving good results](#)
- [Reusing preprocessed training data](#)
- [Running as a WSGI service](#)
- [Backends/Algorithms supported by Annif](#)
  - Regular backends for automated subject indexing and classification
    - [Backend: TF-IDF](#)
    - [Backend: fastText](#)
    - [Backend: Omikujj](#)
    - [Backend: Maui](#)
    - [Backend: vw\\_multi](#)
  - Fusion/Ensemble backends that combine results from other backends
    - [Backend: Ensemble](#)

Welcome to the [Annif](#) users' mailing list / web forum! This list can be used for

- general discussion about Annif, its features and usage scenarios
- asking for help with installing or running Annif
- future directions for Annif
- announcements for new versions and other Annif-related news

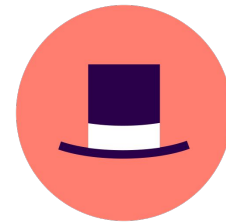
# annif-users forum on Google Groups



- |  |                                |   |              |   |
|--|--------------------------------|---|--------------|---|
|  | Remi Malessa 2                 | <b>Loading OCLC FAST Vocabulary</b> – Actually, my server has around 32MB of memory and it s  | 14.55        | ☆ |
|  | osma.s..... , sale...@g... 4   | <b>Annif presentation and workshop at SWIB20 online conference</b> – Thanks Osma for sharing  | 10. marrask. |   |
|  | mona.l...@..., ... Annif... 13 | <b>Annif tutorial at the 2020 DCMI Virtual Conference</b> – Hi All, There is also a second online h   | 7. lokak.    | ☆ |
|  | sara.v...@kb.nl                | <b>Paper on Annif for categorizing laws</b> – Dear all, Some of you may have read my blogpost o   | 8. syysk.    | ☆ |
|  | haig...@... , osma.s..... 2    | <b>Running Annif as a WSGI service - connexion module error</b> – Hi Thomas! It looks like the P  | 7. syysk.    | ☆ |
|  | juho.k..... , osma.s..... 4    | <b>Results with Annif in employment service context</b> – Hi Juho! juho.k...@gmail.com kirjoitti 2  | 31. elok.    | ☆ |
|  | stephane5... , osma.s..... 2   | <b>Matter of vocabulary</b> – Hi Stephane, Since you have a SKOS vocabulary, you can just load it   | 26. elok.    | ☆ |
|  | osma.s...@helsinki.fi          | <b>ANN: Annif 0.49 released</b> – Annif 0.49 has been released! <a href="https://github.com/NatLibFi/Ann">https://github.com/NatLibFi/Ann</a> | 30. heinäk.  | ☆ |

# Hands-on Annif tutorial

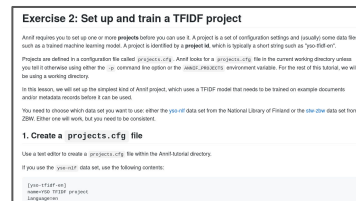
for those who want to use Annif on their own



**SWIB19**  
Semantic Web in Libraries

**DCMI Virtual, 2020**  
September 14th-25th, 2020

**SWIB20**  
Semantic Web in Libraries



Videos and exercises freely available on YouTube & GitHub!



## 4. Annif deployments

# JYX repository, University of Jyväskylä

Students upload their Master's and doctoral theses, Annif suggests subjects\*

## Keywords

<p><b>Keyword suggestions</b> Choose valid keywords by clicking</p>	<ul style="list-style-type: none"><li><input type="checkbox"/> information management systems [YSO]</li><li><input type="checkbox"/> metadata [YSO]</li><li><input type="checkbox"/> connections (technical systems) [YSO]</li><li><input type="checkbox"/> content management [YSO]</li><li><input type="checkbox"/> multimedia (information technology) [YSO]</li><li><input type="checkbox"/> digital libraries [YSO]</li><li><input type="checkbox"/> XML [YSO]</li><li><input type="checkbox"/> semantic web [YSO]</li><li><input type="checkbox"/> open source code [YSO]</li><li><input type="checkbox"/> open data [YSO]</li><li><input type="checkbox"/> user-centeredness [YSO]</li><li><input type="checkbox"/> archives (memory organisations) [YSO]</li><li><input type="checkbox"/> seeking [YSO]</li><li><input type="checkbox"/> Works [YSO]</li><li><input type="checkbox"/> cloud services [YSO]</li><li><input type="checkbox"/> electronic publications [YSO]</li></ul>
<p><b>Your own keywords</b> Comma separated list</p>	<input type="text" value="keyword 1, keyword 2"/>

Implemented using  
DSpace &  
[GLAMpipe](#)  
by Ari Häyrinen

\*from YSO =  
General Finnish  
Ontology

# Osuva repository, University of Vaasa

## Trepo repository, University of Tampere

### Theseus repository, Finnish universities of applied sciences

Same idea as JYX: students upload their theses,  
Annif suggests subjects

Pilot started with Osuva in March  
2020, others followed later.

DSpace extension implemented  
by Anis Moubarik.

Asiasanat:

**Annif-ehdotukset**

<input type="checkbox"/> working abroad	<input type="checkbox"/> families (groups)
<input type="checkbox"/> career development	<input type="checkbox"/> managers and executives
<input type="checkbox"/> career	<input type="checkbox"/> human resources
<input type="checkbox"/> adaptation (change)	<input type="checkbox"/> work
<input type="checkbox"/> expatriates	<input type="checkbox"/> returnees (immigrants)

Lisää

Lisää

Syötä asiasanat, jokainen asiasana omaan kenttäänsä. Paina siis jokaisen asiasanan jälkeen Lisää-nappia. Kirjoita tarvittava määrä asiasanan alkua, jolloin ennakoiva tekstinsyöttö ehdottaa asiasanoja. Muista myös valita yllä olevasta laatikosta Annif-ehdotukset, jotka perustuvat edellisessä vaiheessa syöttämäsi kokotekstin sisältöön.

# Finto AI - automated subject indexing tool and API service



[About](#) [Feedback](#)

[suomeksi](#) [på svenska](#)

Finto AI suggests subjects for a given text. It's based on Annif, a tool for automated subject indexing. [Read more...](#)

## API service

Finto AI is also an API service that can be integrated to other systems.

[Lisätietoja](#) | [OpenAPI-kuvaus](#)

### Enter text to be indexed

In computer science, artificial intelligence (AI), sometimes called machine intelligence, is intelligence demonstrated by machines, in contrast to the natural intelligence displayed by humans and animals. Leading AI textbooks define the field as the study of "intelligent agents": any device that perceives its environment and takes actions that maximize its chance of successfully achieving its goals.[1] Colloquially, the term "artificial intelligence" is often used to describe machines (or computers) that simulate or associate with the human mind, such as expert systems and the new generation of AI.

As machines become increasingly capable, tasks that were once thought to require intelligence are often removed from the definition of intelligence. The Turing test, named after Alan Turing's 1950 paper in Tesler's Theorem says "AI is whatever it does that computers do better than people do".[2] Computer recognition is frequently excluded from things considered to be AI,[3] having become a routine technology.[6] Modern machine capabilities generally classified as AI include successfully understanding human speech,[7] competing at the highest level in strategic game systems (such as chess and Go),[8] autonomously operating cars, intelligent routing in content delivery networks, and military simulations.

[ai.finto.fi](http://ai.finto.fi)

**Launched in  
May 2020**

### Subject indexing

Vocabulary and text language

YSO English

Maximum # of suggestions

10

15

20

Get subject suggestions

### Suggestions

- [artificial intelligence](#)
- [machine learning](#)
- [intelligence \(mental properties\)](#)
- [information technology](#)
- [computational science](#)
- [computer science](#)
- [computers](#)
- [computer-assisted teaching](#)
- [learning](#)
- [automation](#)

# Subject indexing for electronic deposits

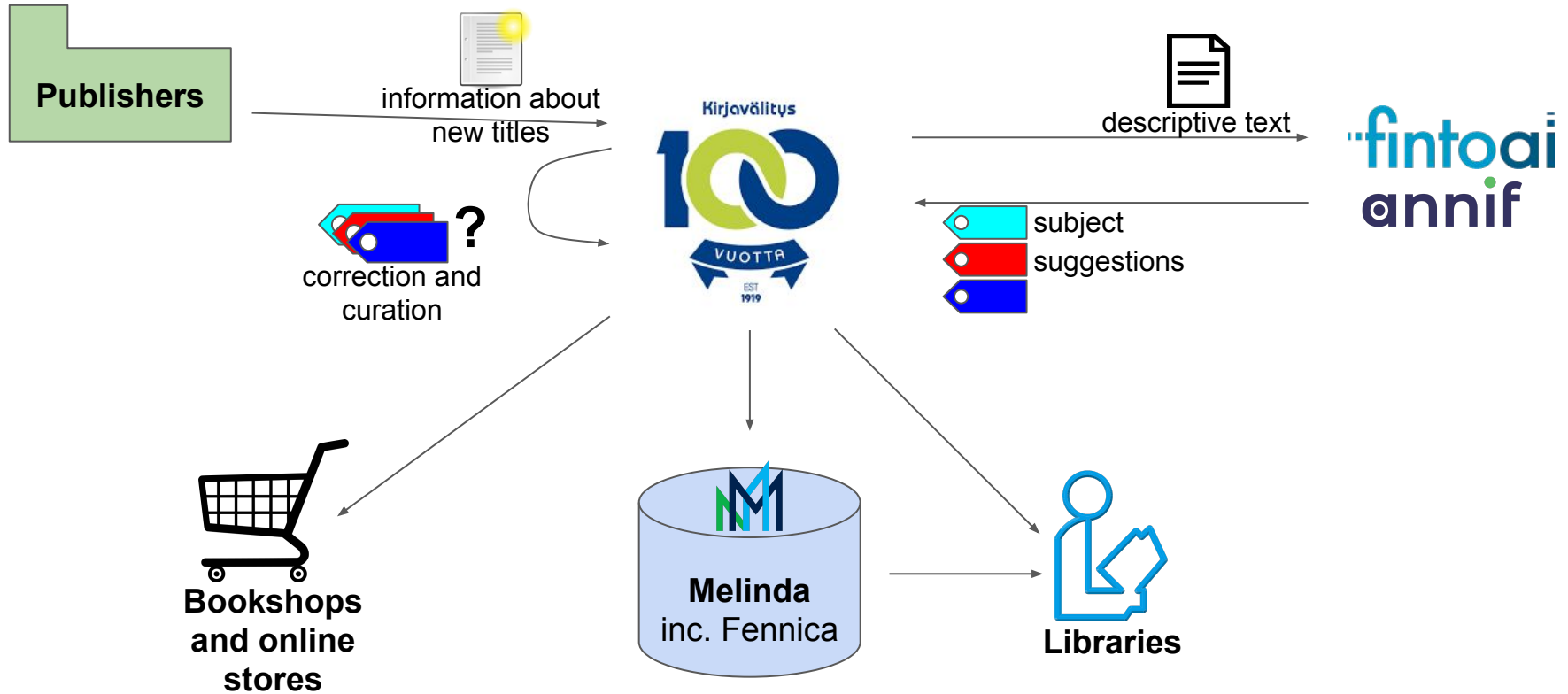
In November 2020, the National Library of Finland started using **Finto AI** to suggest subjects when processing electronic deposits submitted through the individual submission form.

Implementation: Erik Lindgren,  
Mikko Merioksa, Satu Niininen

The screenshot shows the submission form for digital publications. At the top left is the logo of the National Library of Finland (Kansalliset Kirjasto) with the text '1640 KANSALLISET KIRJASTO'. To the right is the title 'Digitaalisten julkaisujen arkistointi' and language selection buttons for 'fi', 'sv', and 'en'. Below the header, there are two main sections: 'Luovuttajan tiedot' (Submitter information) and 'Julkaisun tiedot' (Publication information). The 'Luovuttajan tiedot' section includes fields for 'Yhteysthenkilö \*', 'Sähköpostiosoite \*', 'Puhelinnumero \*', and 'Organisaatio', each with a red asterisk indicating it is mandatory. A red note says 'Tämä on pakollinen kenttä.' (This is a mandatory field). There is also a checkbox for 'Muista nämä tiedot.' (Remember these details). The 'Julkaisun tiedot' section includes 'Julkaisujen lukumäärä' (Number of publications) with radio buttons for 'Luovutan yhden julkaisun' (I submit one publication) and 'Luovutan useita julkaisuja' (I submit several publications). It also has 'Julkaisun tyyppi' (Publication type) with radio buttons for 'kirja' (book), 'muotti' (format), 'äänite' (audio), and 'muu' (other). Finally, there is a 'Perustiedot' (Basic information) section with a field for 'ISBN (viivoilla)' (ISBN (with dashes)).



# Kirjavälitys Oy - logistics company serving bookstores and libraries

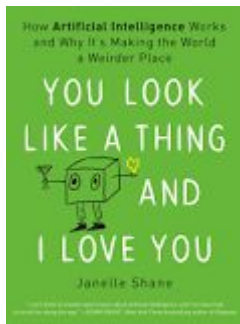


# 5. Lessons learned

# Subject indexing is hard.

Humans have different perspectives and make understandable mistakes.

Algorithms make very silly mistakes.



Case in point:

*Image recognition algorithms will frequently identify **giraffes** in pictures where there are none.*

(Janelle Shane: You Look Like a Thing and I Love You)

Algorithms may be used **alone**, or in combinations, **ensembles**  
**Ensembles are nearly always better** than individual algorithms



# Lessons from evaluation

- The different evaluation approaches are complementary. (see Golub et al., 2016)  
Not a good idea to look at just a single measure.
- Continuous and elusive process: it never stops...

Golub, K., Soergel, D., Buchanan, G., Tudhope, D., Hiom, D., and Lykke, M. 2016. A framework for evaluating automatic indexing or classification in the context of retrieval. *Journal of the Association for Information Science and Technology*, 67(1): 3-16.

Start by experimentation, move slowly towards production



image credit: @kettutatinukkeilee

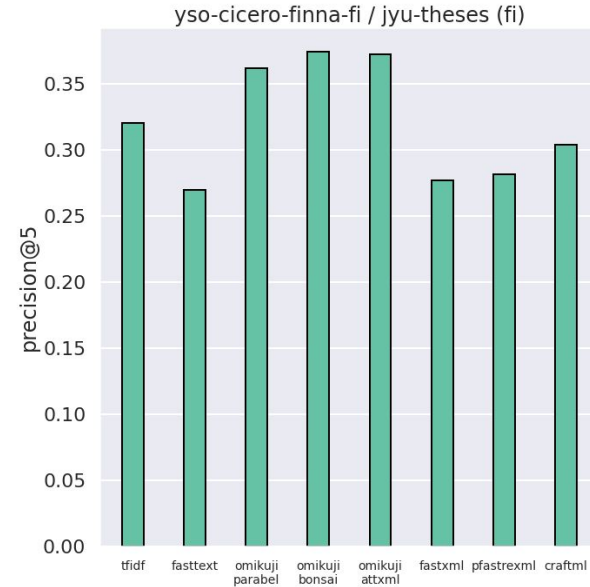
# With an API service such as Finto AI, implementing semi-automated indexing becomes easy; explaining it to users can be more challenging

## Keywords

<p><b>Keyword suggestions</b></p> <p><i>Choose valid keywords by clicking</i></p>	<ul style="list-style-type: none"><li><input type="checkbox"/> information management systems [YSO]</li><li><input type="checkbox"/> metadata [YSO]</li><li><input type="checkbox"/> connections (technical systems) [YSO]</li><li><input type="checkbox"/> content management [YSO]</li><li><input type="checkbox"/> multimedia (information technology) [YSO]</li><li><input type="checkbox"/> digital libraries [YSO]</li><li><input type="checkbox"/> XML [YSO]</li><li><input type="checkbox"/> semantic web [YSO]</li><li><input type="checkbox"/> open source code [YSO]</li><li><input type="checkbox"/> open data [YSO]</li><li><input type="checkbox"/> user-centeredness [YSO]</li><li><input type="checkbox"/> archives (memory organisations) [YSO]</li><li><input type="checkbox"/> seeking [YSO]</li><li><input type="checkbox"/> Works [YSO]</li><li><input type="checkbox"/> cloud services [YSO]</li><li><input type="checkbox"/> electronic publications [YSO]</li></ul>
<p><b>Your own keywords</b></p> <p><i>Comma separated list</i></p>	<input type="text" value="keyword 1, keyword 2"/>



# Collaboration is valuable! (1)

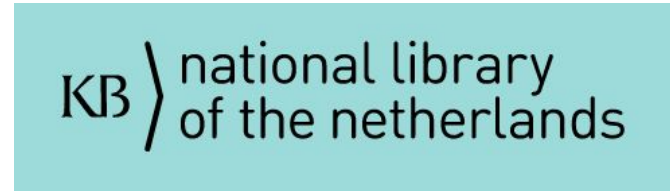


CSC has tested many state of the art text classification algorithms for us. They discovered Omikuji, which is by far the best individual algorithm in Annif currently.

[High-Performance Digitisation](#) project 2018-2020, funded by INEA



# Collaboration is valuable! (2)



[1,2]



[1] Martijn Kleppe, Sara Veldhoen, Meta van der Waal-Gentenaar, Brigitte den Oudsten, & Dorien Haagsma. (2019). Exploration possibilities Automated Generation of Metadata. DOI: <http://doi.org/10.5281/zenodo.3375192>

[2] Romein, C.A., Gruijter, M.D., & Veldhoen, S. (2020). The Datafication of Early Modern Ordinances. DH Benelux Journal, issue 2, 2020. <https://journal.dhbenelux.org/journal/issues/002/article-23-romein/article-23-romein.html>

[3] Lehtonen, T., Piukkula, J. Automaattinen asiasanoitus Radio- ja televisio-ohjelmatietokanta Ritvassa. Informaatiotutkimus 39 (1), 2020. DOI: <https://doi.org/10.23978/inf.88107>

# Thank you!



Juho Inkinen



Mona Lehtinen



Osma Suominen

[annif.org](https://annif.org)

These slides: <https://tinyurl.com/annif-swib20>