# Automated subject indexing with Annif and Finto AI
## Putting DIY automated subject indexing into production

Osma Suominen, Mona Lehtinen, Juho Inkinen
DCMI Virtual AI panel discussion
13 October 2021

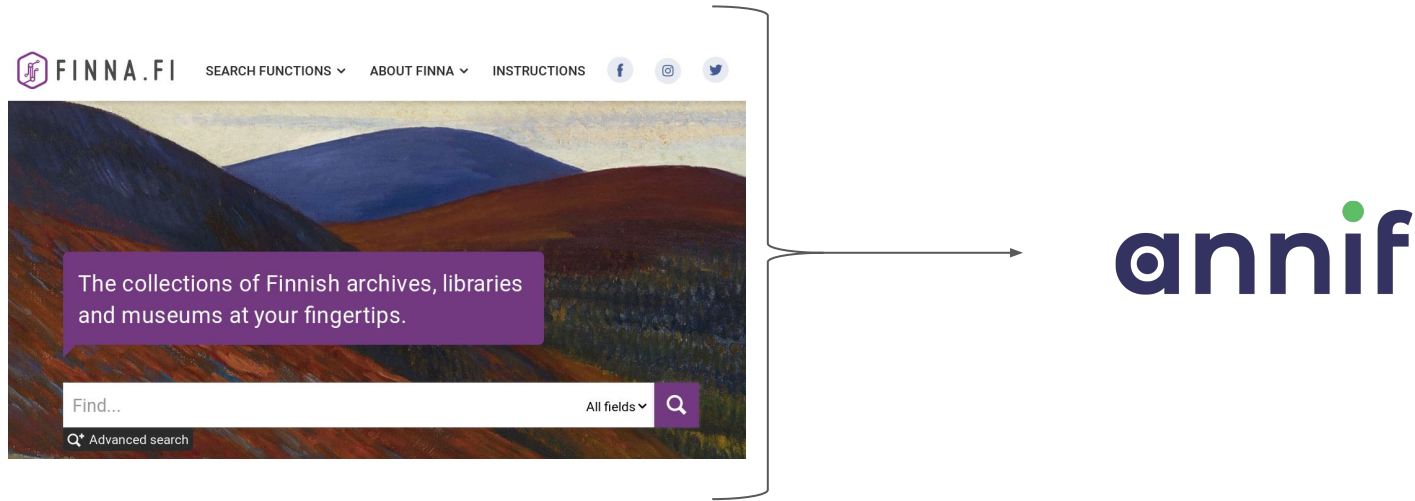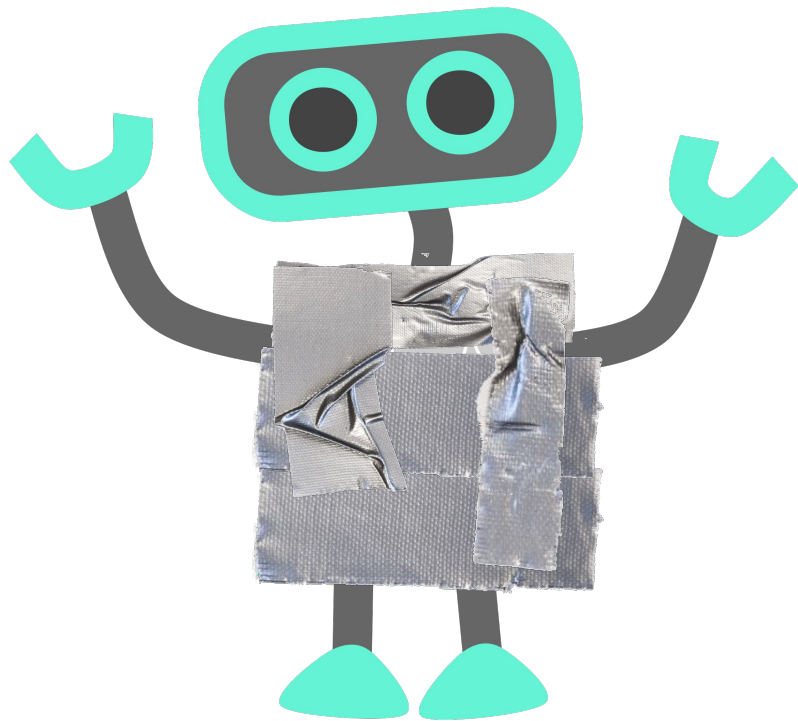NATIONAL LIBRARY
OF FINLAND

# Outline

1. Development of Annif

2. Quality of automated subject indexing

3. Community building

4. Annif deployments

5. Lessons learned

# 1. Development of Annif

# Machine learning using existing metadata

Early prototype (2017) got people excited

# Goals 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 software

annif

# Lexical vs. associative algorithms for subject indexing

**lexical** approaches (e.g.: Maui, MLLM, STWFSA)

match the **terms** in a document
to **terms** in a controlled vocabulary

*"**Renewable resources** are a part of Earth's **natural**
environment and the largest components of its ecosphere."*

yso:p14146
"renewable natural resources"

Lexical approaches need comparatively little training data.

**associative** approaches (e.g.: fastText, Omikuji, SVC)

learn which **subjects** are correlated with which **words**
in documents, based on training data



Associative approaches need a lot more
training data in order to cover each subject.

# 2. Quality of automated subject indexing

# 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.

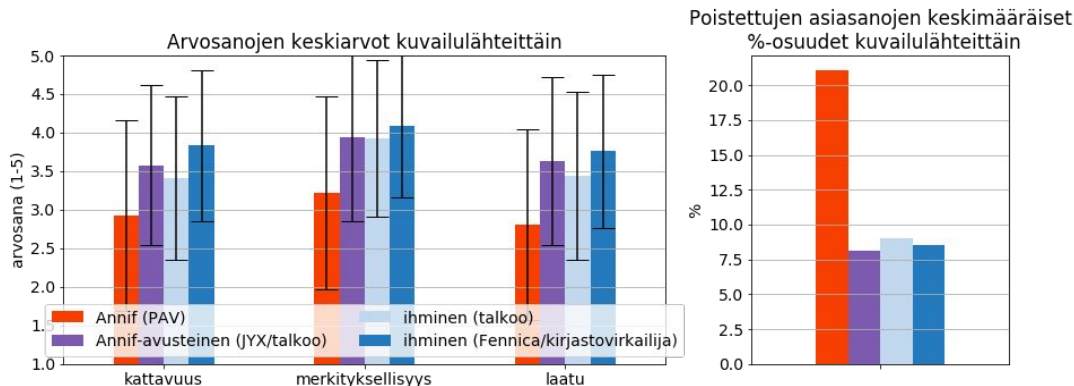The best ensemble algorithm (red bars) was not quite on the level of human indexers in quality scores (left), and significantly more of its suggestions were rejected (right).



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

# 3. Community building

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

**INPUT TEXT**

Why AI ≠ Automated Indexing: What Is and Is Not Possible

Automated indexing is only as good as the training set, or rules that are available for the domain. It's important to learn what type of content a pre-trained algorithm has been trained on. Consider what type of content is readily available to train an algorithm—what's popular and what's available. Scholarly and historical content is not available in consumable formats at the large volume that is required for machine learning. There are exceptions such as science and medicine where large well documented collections are available. This panel will discuss the current state of automated categorization covering domains including research data, art history, and scientific publishing. The goal is to provide practical advice on how to take meaningful steps towards building the infrastructure needed for sustainable automated indexing.

**PROJECT (VOCABULARY AND LANGUAGE)**

YSO NN ensemble English ▾

**MAX # OF SUGGESTIONS**

10  15  20

Get suggestions →    annif

**SUGGESTED SUBJECTS**

- automation
- machine learning
- indexing
- learning contents
- algorithms
- availability
- industrial automation
- artificial intelligence
- content
- oxy-combustion

# 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!



NATIONAL LIBRARY
OF FINLAND

**ZBW** Leibniz-Informationszentrum
Wirtschaft
Leibniz Information Centre
for Economics

# 4. Annif deployments

# JYX repository, University of Jyväskylä
Students upload their Master's and doctoral theses, Annif suggests subjects*

**Keywords**

**Keyword suggestions**
*Choose valid keywords by clicking*

☐ information management systems [YSO]
☐ metadata [YSO]
☐ connections (technical systems) [YSO]
☐ content management [YSO]
☐ multimedia (information technology) [YSO]
☐ digital libraries [YSO]
☐ XML [YSO]
☐ semantic web [YSO]
☐ open source code [YSO]
☐ open data [YSO]
☐ user-centeredness [YSO]
☐ archives (memory organisations) [YSO]
☐ seeking [YSO]
☐ Works [YSO]
☐ cloud services [YSO]
☐ electronic publications [YSO]

Implemented using
DSpace &
GLAMpipe
by Ari Häyrinen

*from YSO =
General Finnish
Ontology

**Your own keywords**
*Comma separated list*

keyword 1, keyword 2

# Finto AI - automated subject indexing tool and API service



ai.finto.fi

# 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

# 5. Lessons learned

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

# Start by experimentation, move slowly towards production



image credit: @kettudolls (IG)

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

**Keywords**

**Keyword suggestions**
*Choose valid keywords by clicking*

- information management systems [YSO]
- metadata [YSO]
- connections (technical systems) [YSO]
- content management [YSO]
- multimedia (information technology) [YSO]
- digital libraries [YSO]
- XML [YSO]
- semantic web [YSO]
- open source code [YSO]
- open data [YSO]
- user-centeredness [YSO]
- archives (memory organisations) [YSO]
- seeking [YSO]
- Works [YSO]
- cloud services [YSO]
- electronic publications [YSO]

**Your own keywords**
*Comma separated list*

keyword 1, keyword 2

What is this?
What should I do here?

Maybe it's better to leave these alone…

# Thank you!

Juho Inkinen

Mona Lehtinen

Osma Suominen