

# Machine Learning and Dewey Decimal Classification

**Freddy Wetjen**  
National Library of Norway

# Outline

## Machine learning and Dewey classification attempts in the National Library of Norway (NLN)

- Why?
- How ?
- Results



# What is *Machine Learning* at NLN?

- NLN has a machine learning lab
- Hands-on experiences with AI technology
- We work with AI and ML on different fields and media types
- AI and ML are tested with all major media types (Film,photo,text,sound..)
- Used for categorization, classification,recognition and discovery
- Build small applications to show the power of machine learning
- Identify strengths and weaknesses of the technology
- Close cooperation with Stanford University Library



AI is not a new technology and certainly not a new way of problem solving.

Machine learning models have improved much in the last five years

The concept of manual knowledge modelling in AI systems is almost gone

Instead, we have introduced the data science concept into machine learning and AI; we let the system build its own knowledge model although carefully selecting the «learning material».

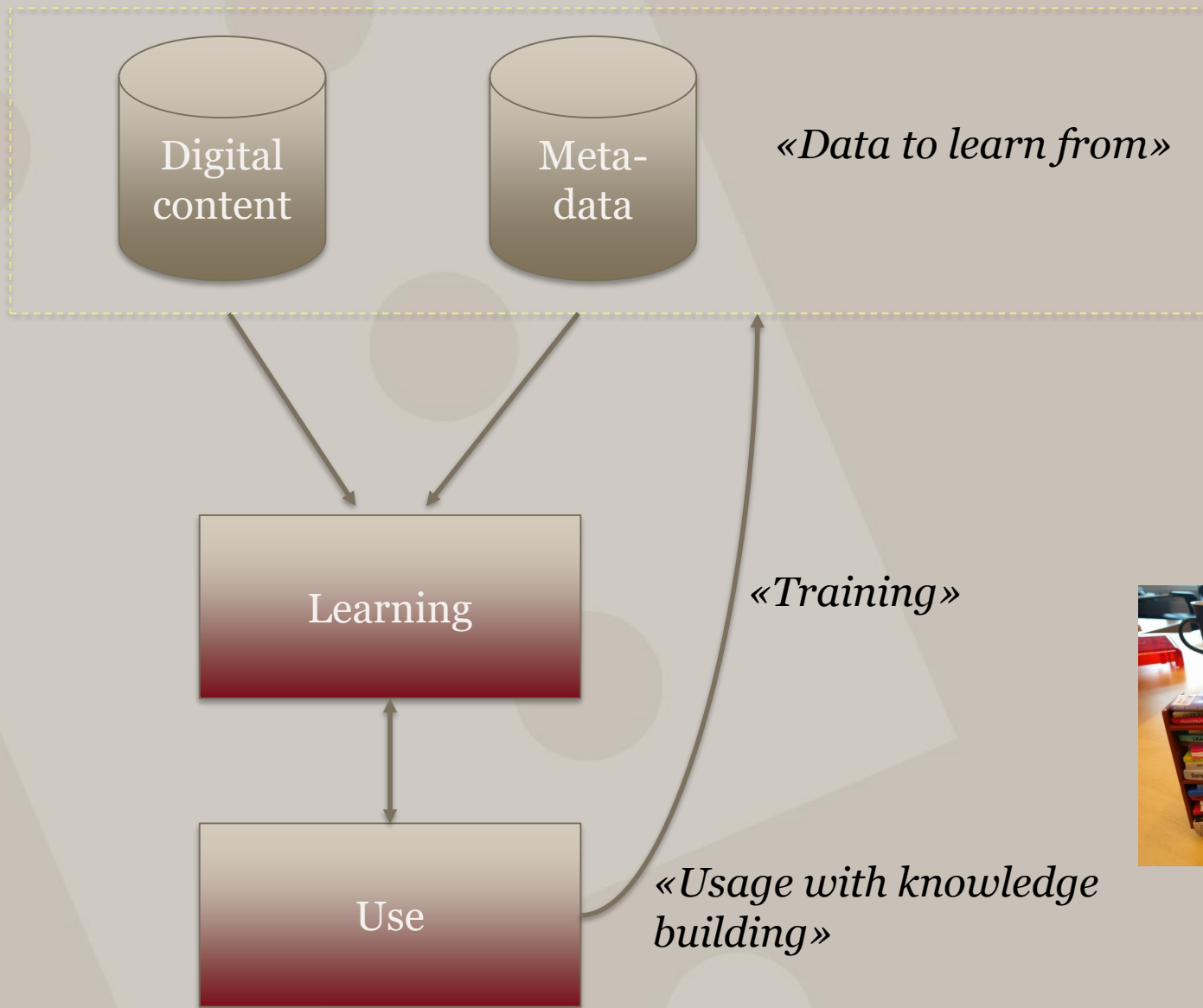
AI methods gets widely available through open frameworks such as Tensorflow, Pytorch, gensim etc.

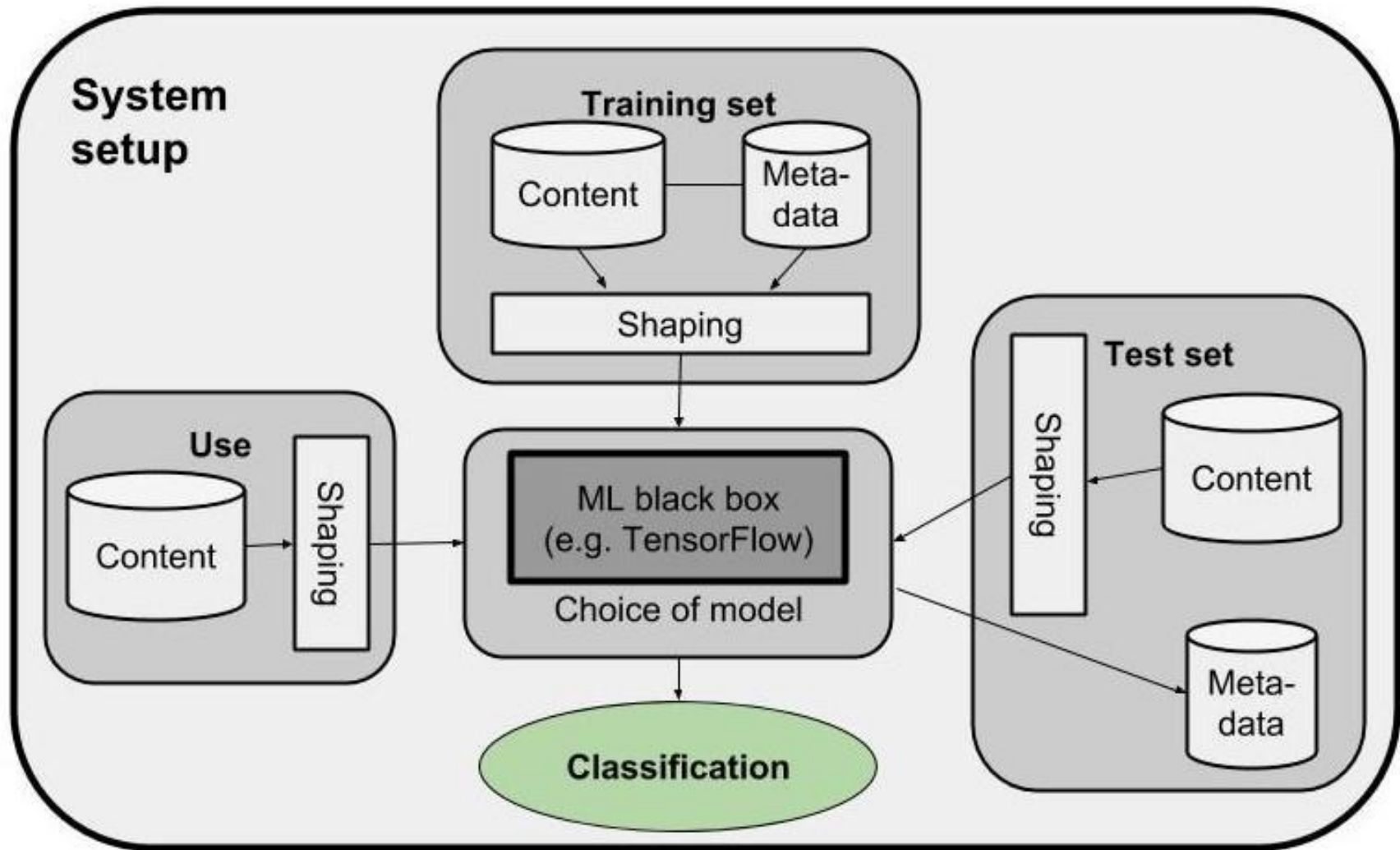
Increasing demand for data science specialists and programmers with knowledge and understanding of ML algorithms

# From *programs* to *rules* to *learning*

- Tradition in programming
  - If-then-else
  - Control and precision
  - Deterministic
- Machine Learning
  - Learning from example data
  - Learning as an automatized task
  - Approximate
  - Non deterministic









# Prerequisites

- Computing power
  - Less power, more time
- Software
  - Mature open-source community
- Training and test data
  - Supervised learning requires high quality labeled data
  - Digital content with metadata (libraries)
- Skills in ML



# Why ML at NLN?



# NLN going digital - ambition

- Mass digitization
  - The complete collection is supposed to be digitized (2006)
  - Most of the published books close to 50 % of all newspaper editions are digitized
- Digital library
  - A complete library at the user's fingertips
  - Search in everything, access to everything
  - UX improvements wanted

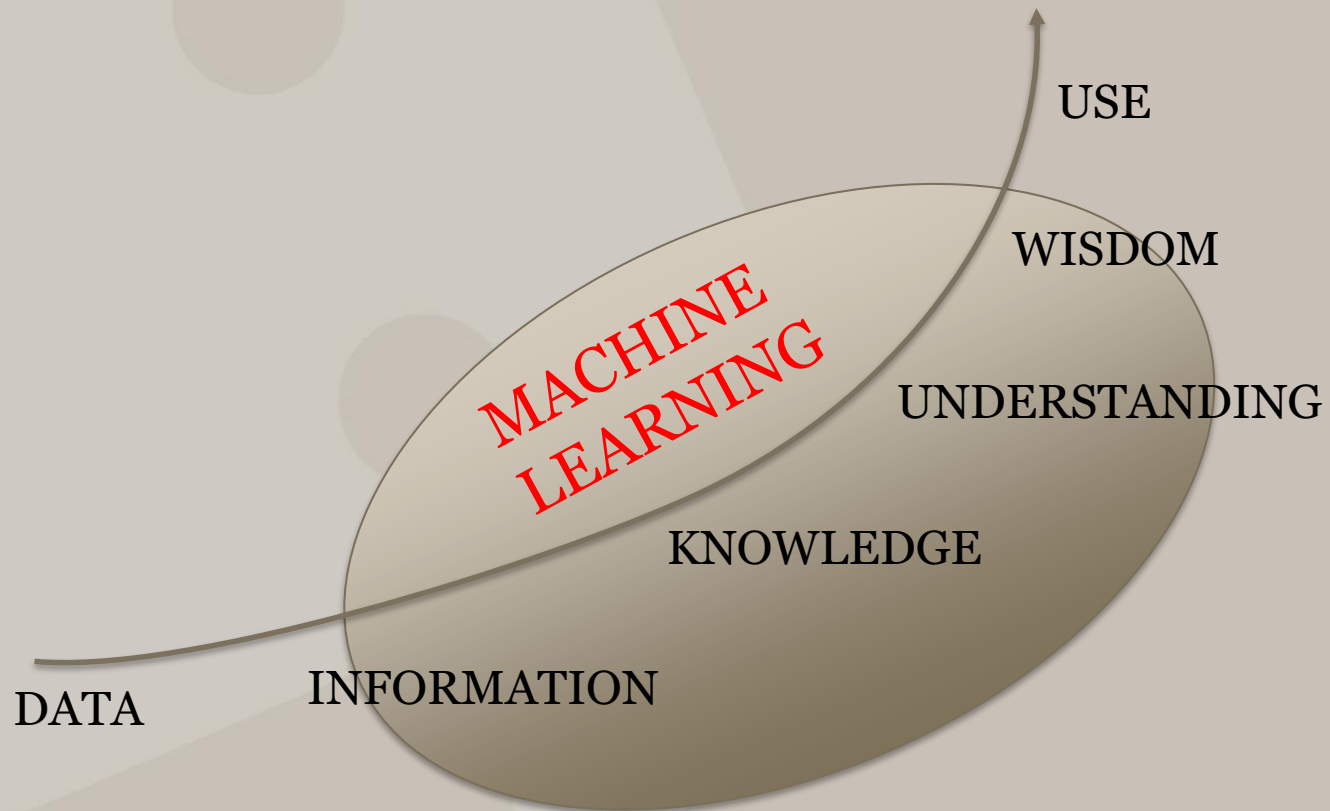


# NLN is the perfect playground

- Massive digital content in all forms
- Good metadata for some data
- User data (user behaviour)
- Good domain understanding, high level of digital skills
- Mature digitalisation technology



# ML helps us being a library

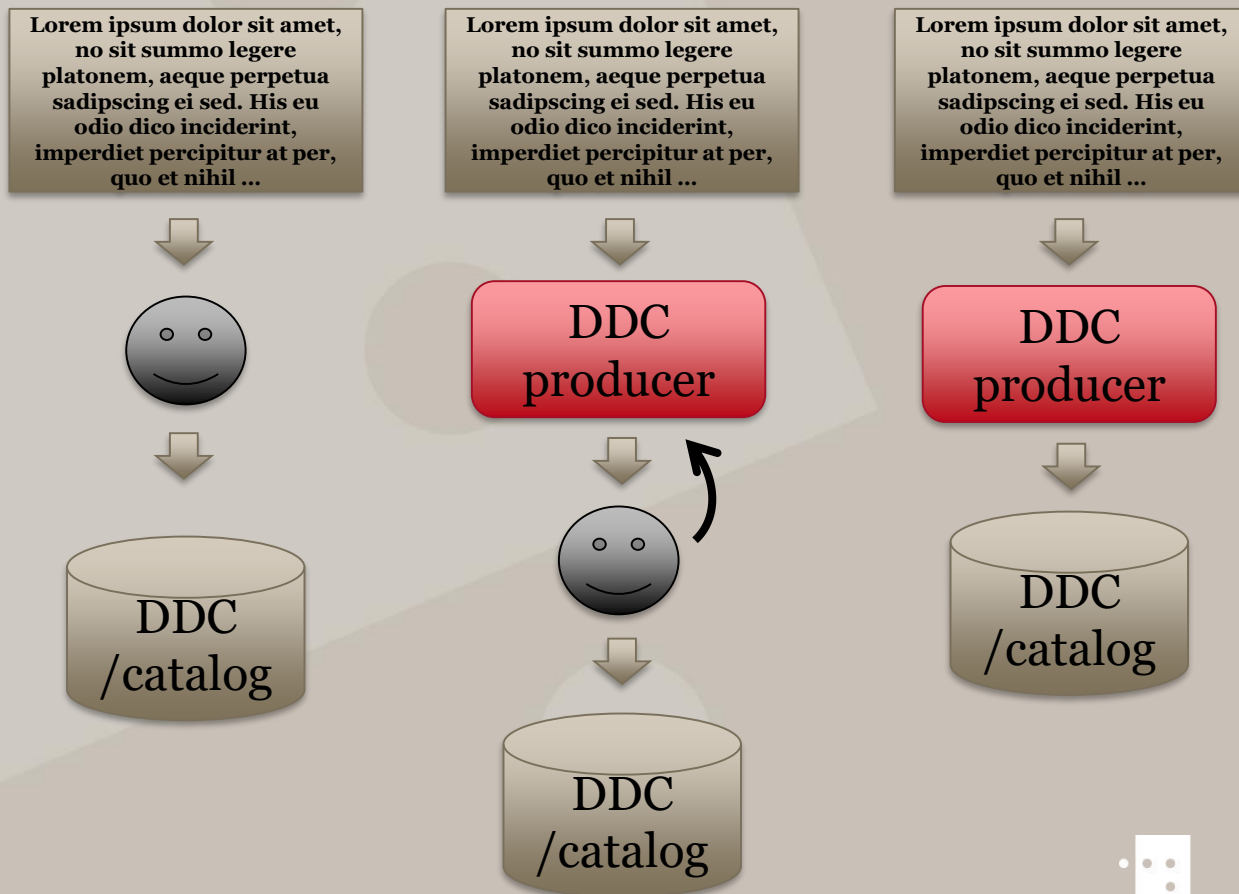


# Various experiments carried out

- Grouping of literature
  - Poetry, Cooking, Sci-Fi, Crime...
- Identifying grey literature
- Speech to text
- Analyzing still images and moving images (video), identifying objects
- Image and video search and identification
- Finding persons, places, organizations and more in text – and relationships between those
- Speaker identification
- Sound fingerprinting



# Ambition: Alternative workflows







# Dewey Decimal Classification experiments with their results



## Using NORART as an example..

- NORART is a hub for access to published Nordic and Norwegian scientific articles
- All articles have dewey classification assigned
- Librarians are classifying all articles
- Time consuming intellectual work
- Carefully selecting publications of particular dewey classification to create train and test sets.
- Working with carefully selected data and testing
- Design of algorithms, parameters, data sets



# Approach

- Define scope for DDC
  - Classes, layers
- Define training set
  - Size
  - Content (articles)
  - Existing metadata
- Define test set
  - Size
  - Content (articles)
  - Existing metadata



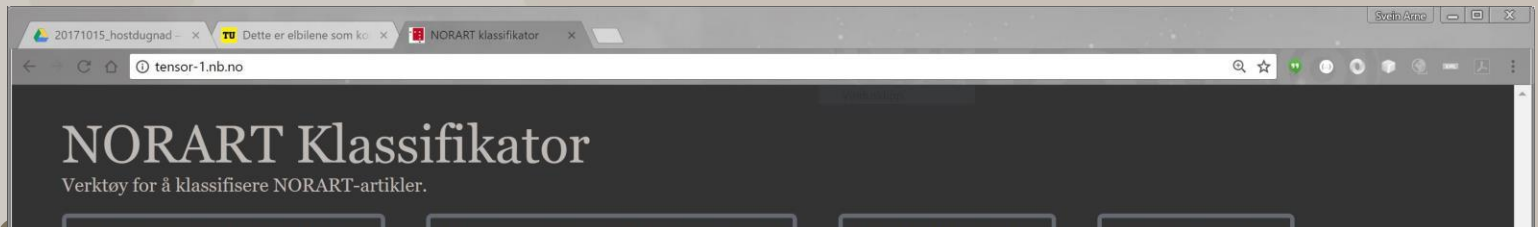
# Constraints

- Limited no of DDC classes
- Only 3, 4, 5 and 6 levels
- More levels, less content per class
  
- Focus example: Automatic DDC identification of NORART scientific articles and content terms



# Example of learning/test definition

<b>L=3</b>	<b>50</b>	<b>100</b>	<b>200</b>	<b>400</b>
Test size	10	20	30	40
Real content only	Yes	Yes	Yes	Yes
Size of artificial content	5/10	10/20	20/40	40/80



tensor-1.nb.no/link.html

# NORART Klassifikator

Verktøy for å klassifisere NORART-artikler.

KLASSIFISERING AV FULLTEKST    KLASSIFISERING AV URL ELLER LINK    LAGE SAMMENDRAG    DOKUMENTASJON

Linken er PDF

Artikkelen lastes ned automatisk. Noen nettsteder tillater ikke dette, så om artikkelen ikke ble lastet den ned, så kan du prøve og kopiere den fra nettstedet og lime den inn manuelt. Det kan gjøres [her](#).

For feilsøk:

**Sjekkliste for tekstklassifikasjon:**

- Er hele linken limt inn?
- Er dette en DOI lenke? - Denne versjonen av klassifisereren støtter ikke DOI. Teksten kan kopieres og klassifiseres [her](#)
- Skal du linke til en PDF, må du linke direkte til PDFen. OBS: Den støtter ikke innebygde PDF-lesere, så om siden benytter en slik en, må du kopiere inn teksten manuelt [her](#).

Her er resultatet fra:  
[https://www.idunn.no/nkt/2015/02/dannelse\\_og\\_folkeopplysning\\_i\\_kulturpolitikk\\_og\\_kulturpolit](https://www.idunn.no/nkt/2015/02/dannelse_og_folkeopplysning_i_kulturpolitikk_og_kulturpolit)

- Dewey-nr: **001** er litt sannsynlig. Klassebetegnelsen er 'Kunnskap'.
- Dewey-nr: **306** er litt sannsynlig. Klassebetegnelsen er 'Kulturelle og sosiale institusjoner'.
- Dewey-nr: **330** er mindre sannsynlig. Klassebetegnelsen er 'Samfunnsøkonomi'.

Riktig Dewey-nr:

transactionExport.xls    Charge Controller .docx    Komplette kvitterin...PDF    20171009\_110334.jpg    Vis alle

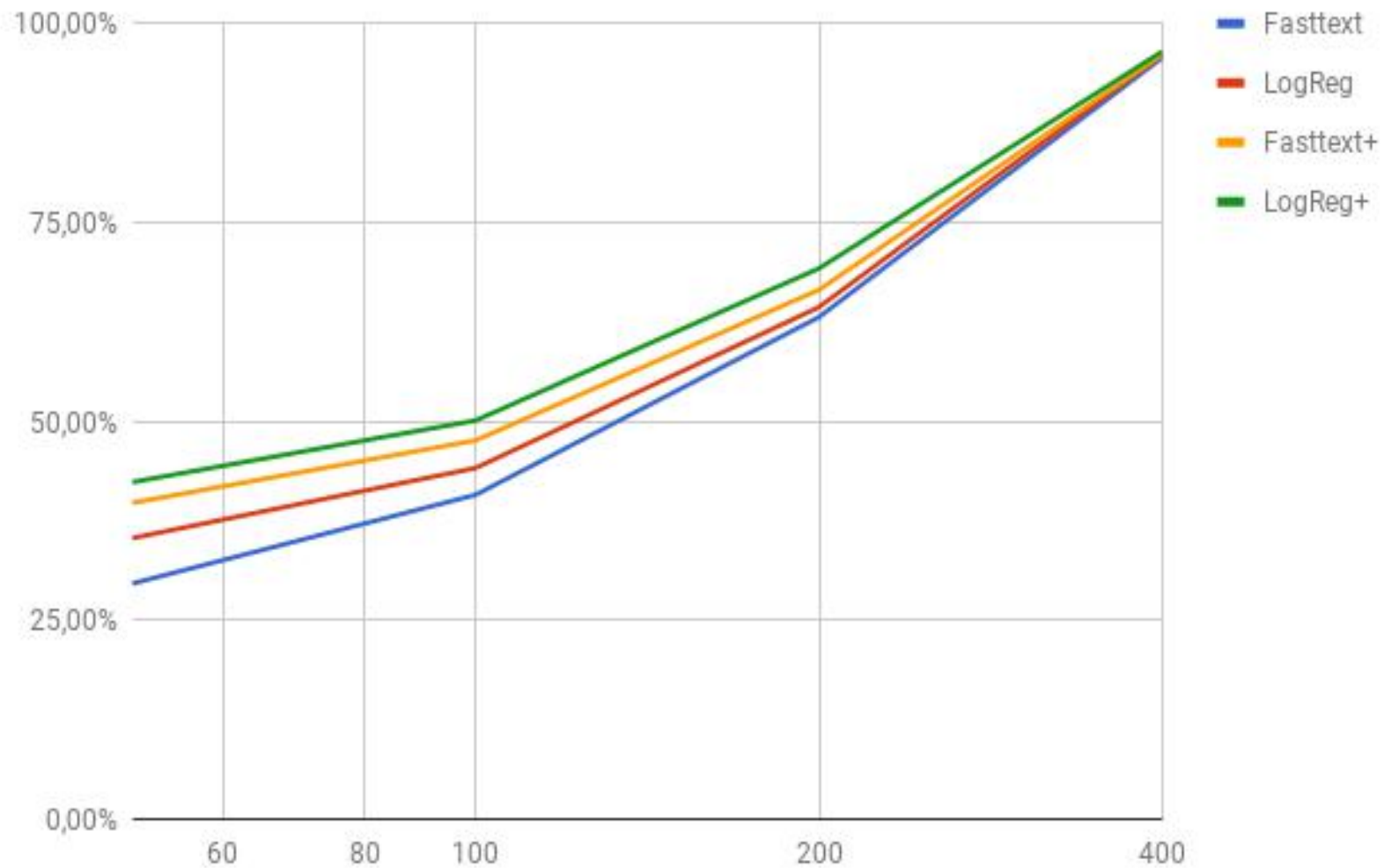
## Btw: Artificial documents

- Used to improve the size of the training set
- «New» articles are produced by interchanging words between articles with the same DDC, or by replacing words/terms with synonyms
- Care taken not to insert bias; Not an easy task to avoid. Using artificial documents has its downside

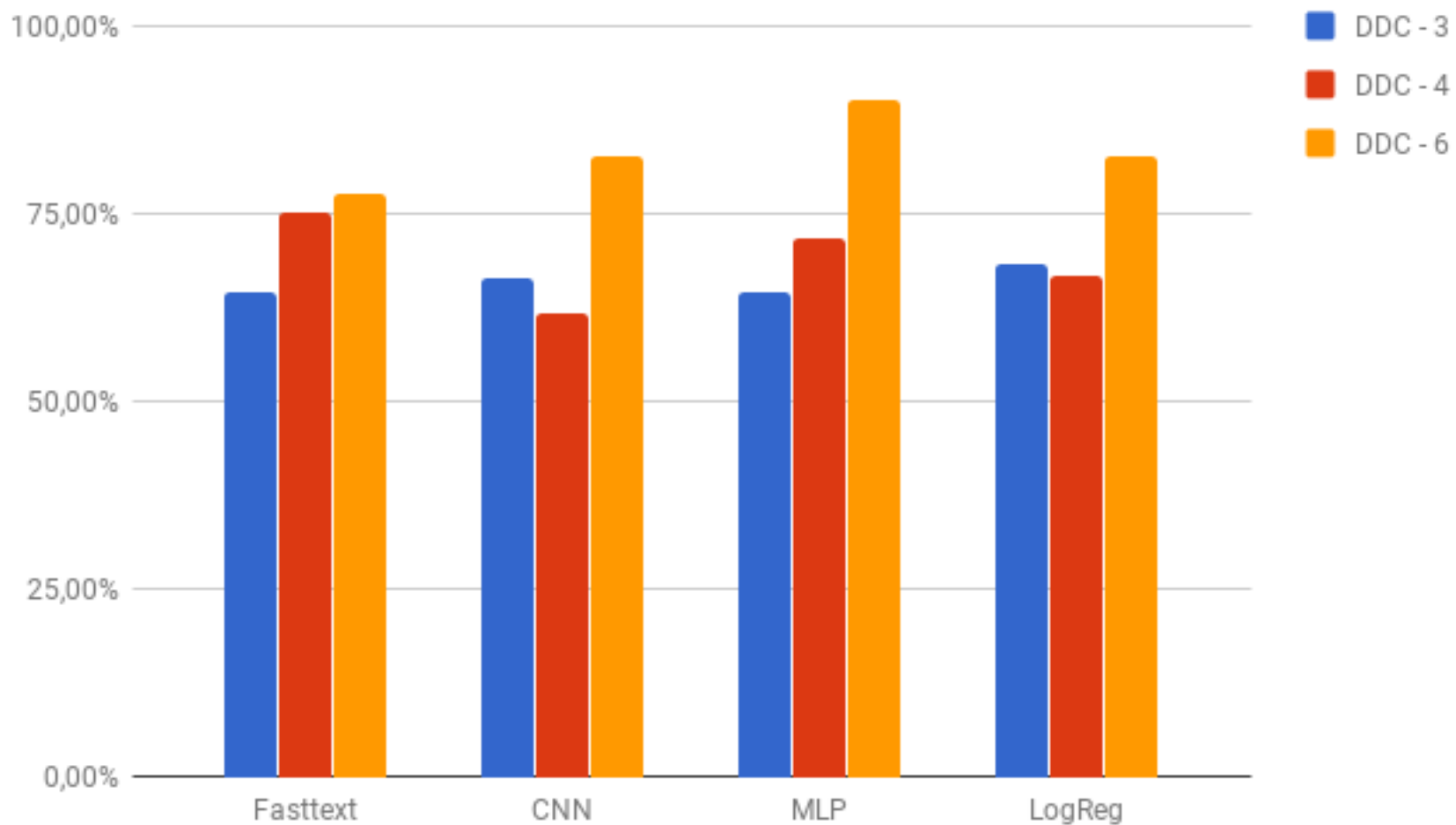




## Precision vs size of training set



## Precision vs number of DDC digits



# Improvements

- Reinforced learning
  - Continuous improvement
  - Corrections from skilled librarians
  - Use of user behaviour
- Change of models



# Conclusions

- Supervised learning on text and metadata from libraries works
- Relatively high precision in prediction of DDC
- Artificial documents helps
- Need for more training data
- Overall, modern ML will play a major role in digital libraries



Thanks for listening  
[freddy.wetjen@nb.no](mailto:freddy.wetjen@nb.no)