



UNIVERSITY OF  
EASTERN FINLAND

# **Machine learning, data mining, and the world's largest manuscript traditions**

**Pasi Hyytiäinen, Classical Texts in Digital Media II, Venice, 18.6.2025**



# In this presentation

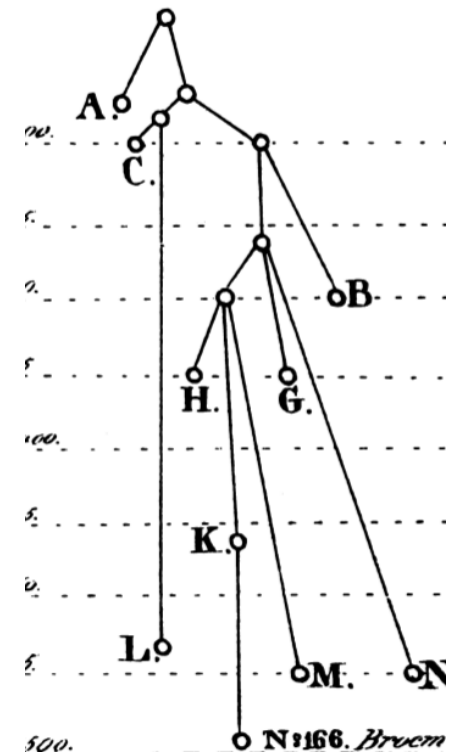
- Stemmatology?
- Challenges
- Computer-assistance...
- The world's largest manuscript traditions
  - Taking all data into account
    - Machine learning
    - Data mining



# Stemmatology

- The classical texts were transmitted through the medium of manuscripts
  - There are different versions of a given literary work
- Stemmatology uses tree (or network) structures to depict the genealogical relationships between manuscripts
- The goal is to uncover the starting point of the manuscript tradition, the initial or the original text

Schema Cognationis  
Codicum manusc.





# Challenges

- Efficiency
  - Traditional Lachmannian stemmatology
    - Extremely time-consuming
      - Transcribing, collating, and dividing texts into variation units



- A the fox jumped over the hedge
- B -
- C the cat jumped over the fence
- D a man saw that the fox jumped over the hedge
- E a man saw that the fox jumped over the fence

**Solution 1****Solution 2**

–	the fox jumped over the hedge	–	the fox jumped	over the hedge
–	–	–	–	–
–	the cat jumped over the fence	–	the cat jumped	over the fence
a man saw that	the fox jumped over the hedge	a man saw that	the fox jumped	over the hedge
a man saw that	the fox jumped over the fence	a man saw that	the fox jumped	over the fence



# Challenges

- Efficiency
  - Traditional Lachmannian stemmatology
    - Extremely time-consuming
      - Transcribing, collating, and dividing texts into variation units
      - Only shared errors (innovations, secondary readings) can provide proof of common ancestry
      - Distinguishing monogenetic variants
      - Constructing stemmata based on shared errors
        - » Stemma codicum



# Challenges

## ■ Efficiency

### – Traditional L

#### • Extremely t

– Transcrik

– Only sha

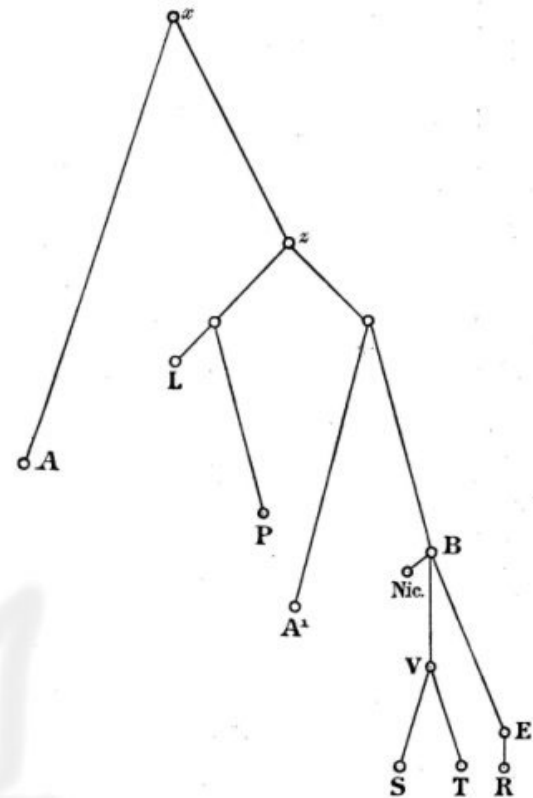
commor

– Distingui

– Construc

» Stem

1000  
1100  
1200  
1300  
1400  
1500



gy

nto variation units

ry readings) can provide proof of

errors



# Challenges

- Contamination
  - “where contamination exists, the science of stemmatics in the strict sense breaks down.” (Paul Maas 1880-1964)
  - The scribes used more than one exemplar
  - A traditional stemma cannot depict a contaminated tradition
    - Contamination mixes the manuscript histories and ancestral lineages





# Computer-assistance...

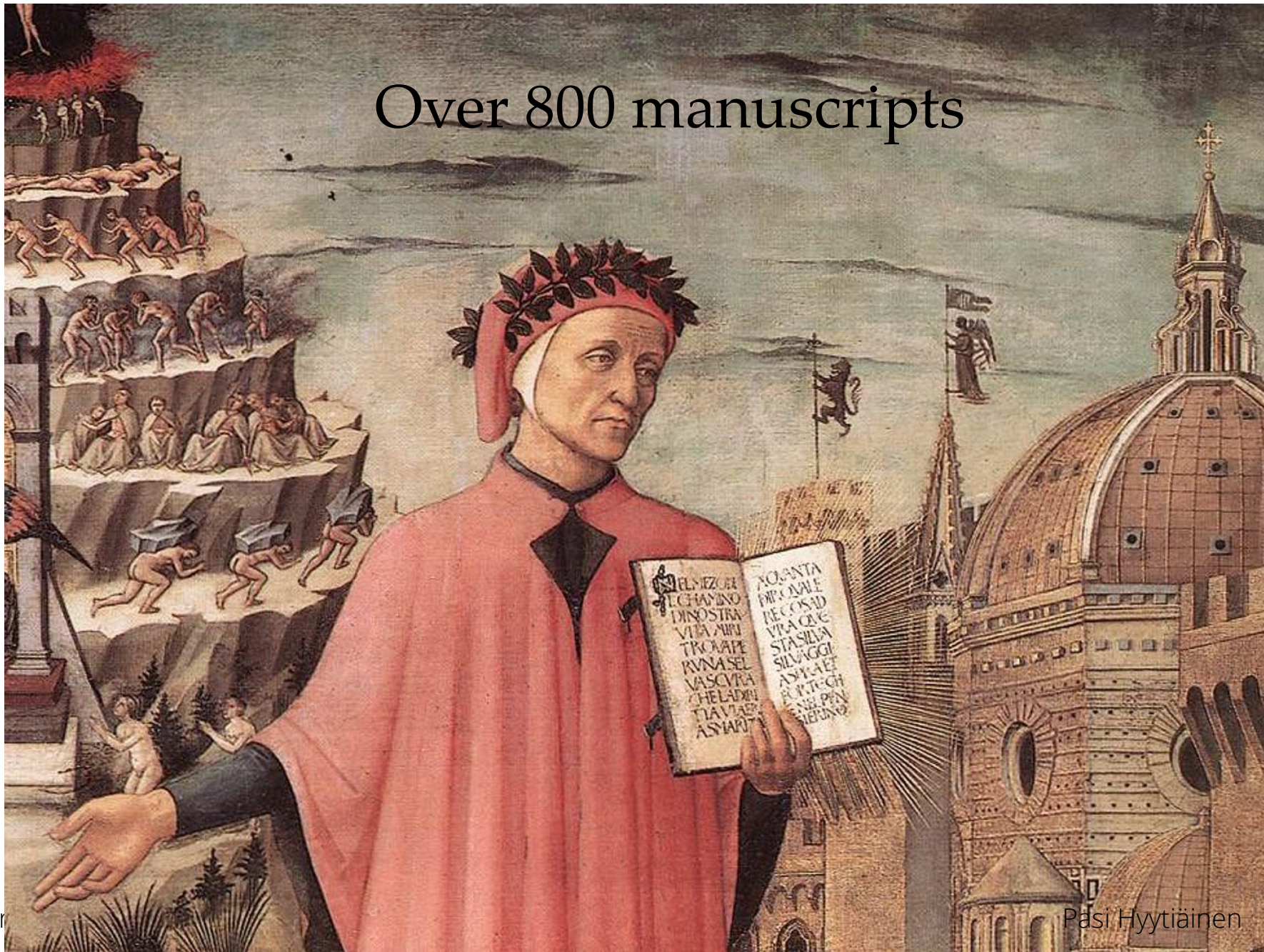
- Digital methods have been used to study manuscript relations from the late 1960s (Griffith 1969)
- Resolving some of the challenges
  - Collations can be automated (Collate)
  - Trees can be automatically generated (phylogenetic methods)
    - Peter Robinson, Chris Howe etc.
- We need new methods in the face of large manuscript traditions



# The world's largest manuscript traditions



# Over 800 manuscripts





Over 2000 manuscripts



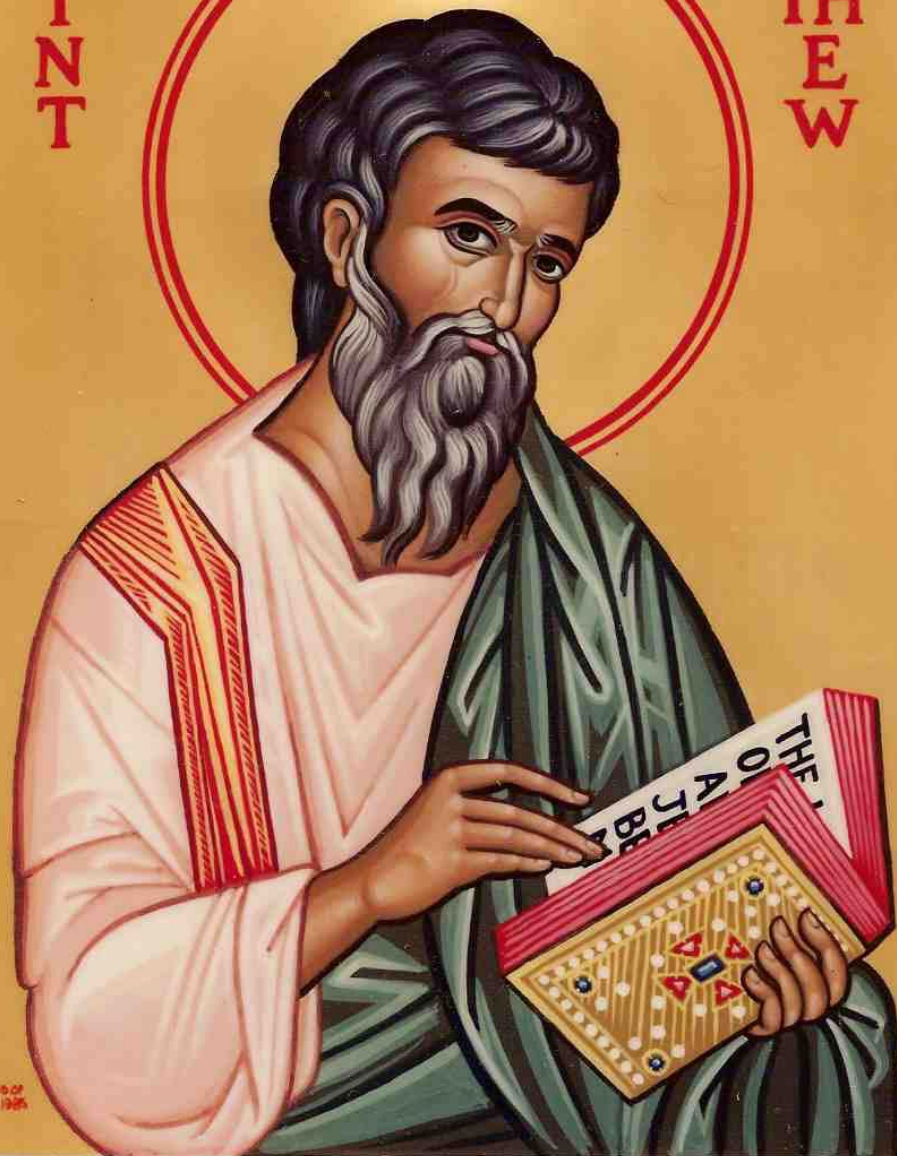
# HOMER



SA  
I  
N  
T

2112 manuscripts

MAT  
THE  
W



BY THE HAND OF  
SARON 1985



M-66

सं १

# Several thousands

2

श्रीदुर्गाजादिषु पंचाशत् गणपतिभ्यो नमः ॥ ३० ॥ अग्निमीळे पुरोहितं यज्ञस्य दे-  
 वमृत्विजो होतारं रत्नधातमम् ॥ अग्निः पूर्वसिद्धिर्षिभिरीशो नूतनैरुत ॥ स देवो एह  
 वक्षति ॥ अग्निनाशयिमश्रुत्वोषमेव दिवेदिव ॥ यज्ञसंवीरुतमे ॥ अग्नेयं यज्ञमध्व-  
 रं विश्वतः परिभूरसि ॥ स इ देवेषु गच्छति ॥ अग्निर्होता कुर्विक्रतुः सत्यश्चित्रश्रव-  
 स्तमः ॥ देवो देवेभिरागमत् ॥ १ ॥ यद्गदाश्रुषेत्नमग्ने भद्रं करिष्यसि ॥ तवेत्तत्सत्यम-  
 गिरः ॥ उपत्वाग्ने दिवेदिवदोषावस्तर्धियावयं ॥ नमो भरंत एमसि ॥ राजंत मधुरा-  
 णां नोपामृतस्य दीदिविं ॥ वदं मानस्येदमे ॥ सनः पिते वसू नवे ग्रेसूपाय नो भवा सुवस्था-  
 नः स्वस्तये ॥ २ ॥ वायवायाहिदशंत मेसोमा अरुक्ताः तिषां पाहि श्रुधी हव ॥ वाय





# The world's largest manuscript traditions

- How do we handle these large traditions?
  - The ultimate goal is to take all data into account
- First step
  - Digital photography
    - All manuscripts need to be digitally photographed



# The world's largest manuscript traditions

- Second step
  - Transcribing the texts into digital form
  - This step must be automated
    - Machine learning (and AI)
      - We generate a model for a given tradition (Latin, Greek, Church Slavonic...)
        - » We teach a program to differentiate text regions, lines, words, and letters from a photo
          - » Transkribus software

## Молитвы оутреннія .

Слава тебѣ Бже нашъ, слава тебѣ .

Црю нѣный, оутѣшителю, дше истины,  
йже вездѣ сый, й вса исполнай, со-  
кровище блгнхъ, й жизни подателю,  
приди й вселса въ ны, й ѡчисти ны  
ѡ всакїа скверны, й спаси Бже дш-  
шы наша .

Стый бже, стый крѣпкій, стый безс-  
мертный, помилди насъ . **Трижды .**

Слава Оцѣ, й Снѣ, й стѡмѣ Дхѣ, й  
нынѣ й приснѣ, й во вѣки вѣкѡвъ,  
аминь .

Престѡа Трѣце помилди насъ: Гди  
ѡчисти грѣхи наша: Вѣко прости  
беззаконїа наша; Стый посѣти  
й исцѣли немощи наша, й мене  
твоегѡ ради .

Гди помилди, **трижды: Слава: й нынѣ:**  
Оче нашъ, йже еси на нѣбсѣхъ, да  
стїтса

## Молитвы оутреннія .

Има твоє: да приидеть црѣвїе твоє:  
да вѡдетъ кола твоа, ѡакш на нѣси,  
й на земли . Хлѣбъ нашъ насѡщный  
даждь намъ днесъ . й ѡстави намъ  
долги наша, ѡкоже й мы ѡставлаемъ  
должникѡмъ нашимъ . Й не введи насъ  
во искушенїе, но избави насъ ѡ лѡкавѡш .

**Таже: Троицны сѡа: Трѡнаръ гласъ ѡ:**

Воставше ѡ сна припадаемъ ти вѣже,  
й аггггскѡю пѣснью копїемъ ти силне:  
стѣ, стѣ, стѣ еси Бже, богородицею  
помилди насъ . **Слава:**

ѡ Одра й сна воздвигль ма еси Гди,  
оумъ мой просѣти й сердце, й оустнѣ  
мой ѡ керзи, во еже пѣти та стѡа  
троице: стѣ, стѣ, стѣ еси Бже, вѣцею  
помилди насъ . Й нынѣ .

Напраснѣ едїа приидеть, й коегѡ  
ждо дѣлїнїа ѡвнѡжатса, но страхѡмъ

## Молитвы оутреннія .

Слава тебѣ бже нашъ, слава тебѣ .

Црю нѣный, оутѣшителю, дше истины,  
йже вездѣ сый, й вса исполнай, со-  
кровище блгнхъ, й жизни подателю,  
приди й вселса въ ны, й ѡчисти ны  
ѡ всакїа скверны, й спаси бже дш-  
шы наша .

Стый бже, стый крѣпкій, стый безс-  
мертный, помилди насъ . **Трижды .**

Слава Оцѣ, й Снѣ, й стѡмѣ Дхѣ й  
нынѣ й приснѣ, й во вѣки вѣкѡвъ,  
аминь .

Престѡа Трѣце помилди насъ: Гди  
ѡчисти грѣхи наша: Вѣко прости  
беззаконїа наша; Стый посѣти  
й исцѣли немощи наша, й мене  
твоегѡ ради .

Гди помилди, **трижды: Слава: й нынѣ:**

Back to overview

# Valamo 0.7

☆ Share Options

## Model description



by pasi.hyytiainen@valamo.fi

- Training Set Size
- Training pages
- Validation pages
- Words
- Lines

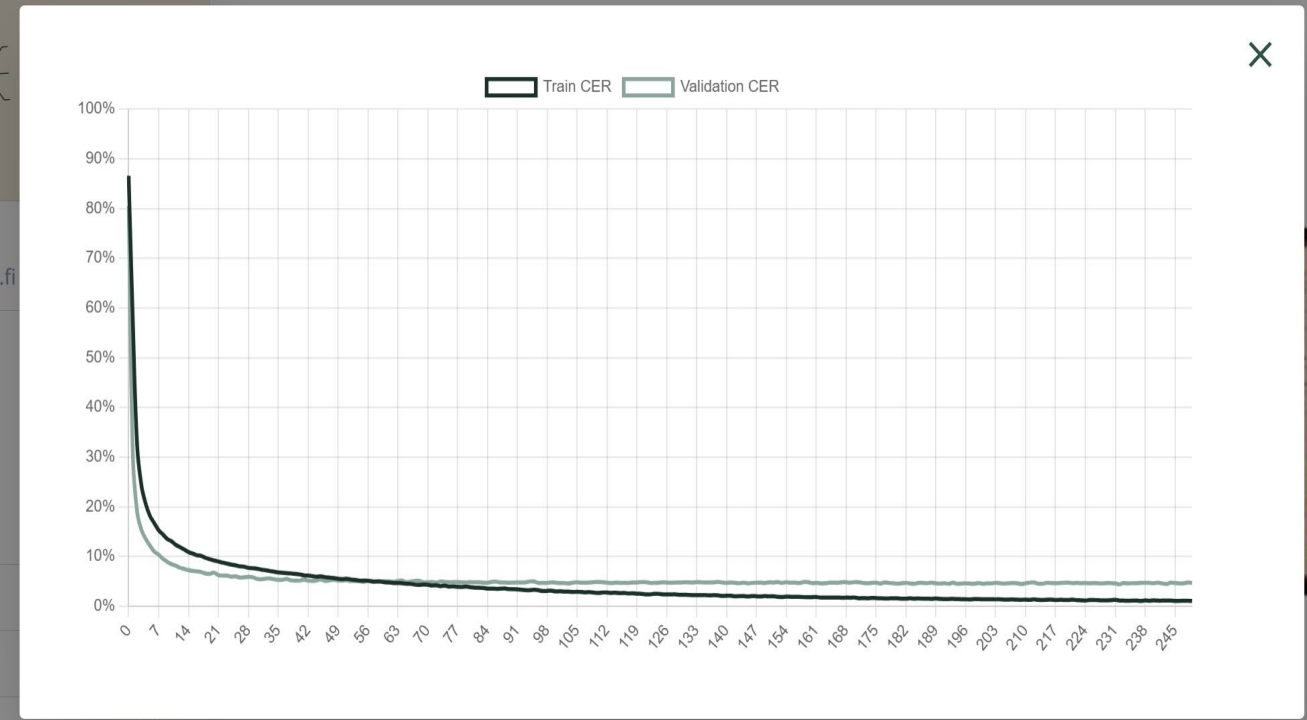
% CER (Accuracy)

Centuries

# Model ID 339993

See training chart

See training parameters



View all

AI Model: Valamo 0.7

Drop files here or browse files



# The world's largest manuscript traditions

- Third step
  - Depicting the relations between the texts using distance values
  - Data mining techniques
    - No need to collate and divide texts into variation units
    - We calculate distances between pairs of texts
    - Introducing a very simple and fast method
      - K-shingling

# K-Shingling

a man saw that the fox jumped over the fence

a man	man saw	saw that	that the	the fox	fox jumped	jumped over	over the	the fence
fox jumped	the fence	a man	over the	that the	saw that	the fox	jumped over	man saw

# K-Shingling

a man saw that the fox jumped over the fence

a man saw that the fox jumped over the hedge

the  
fox

jumped  
over

saw  
that

that  
the

a  
man

fox  
jumped

man  
saw

over  
the

the  
fence

fox  
jumped

the  
hedge

a  
man

over  
the

that  
the

saw  
that

the  
fox

jumped  
over

man  
saw

# K-Shingling

Set1	the fox	jumped over	saw that	that the	a man	fox jumped	man saw	over the	the fence
------	------------	----------------	-------------	-------------	----------	---------------	------------	-------------	--------------

Set2	fox jumped	the hedge	a man	over the	that the	saw that	the fox	jumped over	man saw
------	---------------	--------------	----------	-------------	-------------	-------------	------------	----------------	------------

Word bigram	Set 1	Set 2
fox jumped	1	1
man saw	1	1
jumped over	1	1
that the	1	1
saw that	1	1
the fox	1	1
a man	1	1
over the	1	1
the hedge	0	1
the fence	1	0

Intersection = 8

Union = 10

Sørensen-Dice Coefficient (SDC)

$$\frac{2 \times \text{intersection}}{\text{sum of the number of elements in each set}}$$

$$= 16 / 18 = 0.888 = 88 \%$$



```
Project ▾
  relate C:\Users\K...
  .venv library r...
    Include
  > Lib
  > Scripts
  .gitignore
  pyenvn.cfg
  main.py
  External Libraries
  Scratches and Co...

main.py x
262
263 'L844' 'L950' 'L2211')
264
265 length = 2
266
267 4 usages
268 def shingle(string):
269     shingle = [string[i:i + length] for i in range(len(string) - (length - 1))]
270     return shingle
271
272 1 usage
273 def sorensen_dice(string_a, string_b):
274     intersection = len(set(shingle(string_a)) & set(shingle(string_b)))
275     sum_of_elements = len(set(shingle(string_a))) + len(set(shingle(string_b)))
276     sorensen_dice_coefficient = ((intersection * 2 / sum_of_elements) * 100)
277     return sorensen_dice_coefficient
278
279 1 usage
280 def matrix(sorensen_dice, texts1, texts2):
281     DM = [(sorensen_dice(a,b)) for a in texts1] for b in texts2]
282     matrix = pd.DataFrame(DM, columns=names, index=names)
283     pd.set_option('display.width', None)
284     return matrix
285
286 result = matrix(sorensen_dice, texts1, texts2)
287 print(result)
```



# The world's largest manuscript traditions

- Fourth step
  - Constructing trees or networks using the distance matrix
    - Tree inference algorithms
      - UPGMA
      - Neighbour-joining
    - Network algorithms
      - Neighbour-Net
      - Split decomposition





UNIVERSITY OF  
EASTERN FINLAND

# Thank you!

οἱ σκαπταροῦσθες ἐνεαυόντα  
μετατοθεῖν αὐτὸν ἐν πολλοῖς τεκμηρίοις  
τεσσεράκοντῶν ἡμερῶν  
ὅπταν ὀνομαζοῦσθαι αὐτοῖς καὶ ἐφῶν

**A manuscript scholar,  
theologian, digital humanist,  
textual critic, and stemmatologist.**



ὅτι ἰσθάνησθαι ἐβαπτίσεν ὑδάτι  
ὑμεῖσδε ἐν ἰνιαγίῳ βαπτίσθησθαι  
καὶ ὀμμελλεταί λαμβάνειν  
ὄυ μεταπολλάστας ἡμέρας  
(ἐφῶσθαι ἐν τῆ κοστῆς)