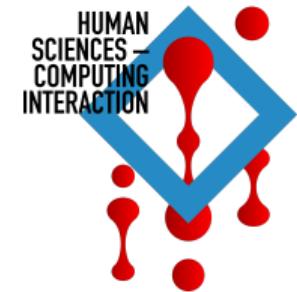




Text similarity and alignment in the study of Finnic oral folk poetry

Maciej Janicki

July 6, 2022





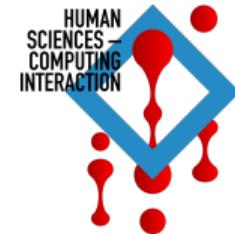
THE FILTER PROJECT

Formulaic intertextuality, thematic networks and poetic variation
across regional cultures of Finnic oral poetry

A cooperation of:

- ① **Finnish Literature Society**

Kati Kallio, Jukka Saarinen



- ② **University of Helsinki, HSCI group**

Eetu Mäkelä, Maciej Janicki

- ③ **Estonian Literary Museum**

Mari Sarv, Liina Saarlo

| || ||| || |||| Estonian Literary Museum



THE CORPUS

Text collections:

- ① Suomen Kansan Vanhat Runot (SKVR)
- ② Eesti Regilaulude Andmebaas (ERAB)
- ③ SKS's corpus of unpublished poems (JR)
- ④ Printed works: Kalevala, Kanteletar, Kalevipoeg

Characteristics:

- written in Finnic tetrameter
- mostly 19th-early 20th century
- languages: Finnish, Estonian, Karelian, Ingrian, Votic, Seto, ...

Lappalai e on kytt silm 
Piti viikoista vihuva,
Kauvon aijasta katsetta.
Huol'itti tulista nuolda
Tul'izilla j ndevill ,
Pe sk n pienill  sulilla,
Varposen vivuttsimilla.
SKVR I1 17

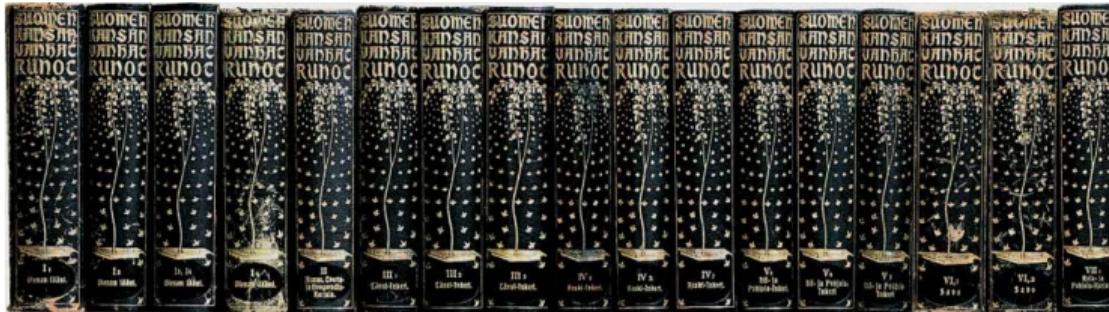
M istke, m istke, mehed noored
Teadke, teadke naesed targad
Arvake k la  eksed
Mis see seal m ella kasvab
Tamme seal m ella kasvab
E S VII 313 (67)



SIZE

Currently (SKVR + ERAB + JR + Kalevala + Kanteletar + Kalevipoeg):

- 275,181 documents (poems),
- 4,369,489 verses (tokens),
- 3,477,074 unique verses.



(– 15% of our corpus)



SIZE

Currently (SKVR + ERAB + JR + Kalevala + Kanteletar + Kalevipoeg):

- 275,181 documents (poems),
- 4,369,489 verses (tokens),
- 3,477,074 unique verses.





GOAL: ALIGNMENT

Yht' ei kuttsun Lemmingäistä.

*Rujot (ne) reillä reissaabi,
Rammat rattšahin ajeli,
20_Sogiat venozin soudi.*

Lemmingäin on poiga_pilo
Pillojah on piilemässä,#10
Pahojah pagenemassa.#11
"Hoib om moamo, kandajańi,
25_Armas maijon andańi,
Ihalan imettäjäńi,
Ettsib om miul pelvoi paid[al],
Ennembä neidona kuvottu,
Kassabapeän#12 on kalkuteltu,
30_Kannabas paloni_paid[a]."

<http://runoregi.rahtiapp.fi/poemdif?nro1=skvr07108030&nro2=skvr07108080>

15 Kuttšu veri-sogeat,
Ruiot re'ellä rembuttelı,
Rammat rattšahin ajeli,
Sogeat venosin souti,
Yht' ei kuttsu Lemmingästä.
20_Lemmingäne on piilopoiga
Pilloja on piilemässä,
Pahoja pagenemassa.
"Hoi on moammoni, kantajani,
Armas maion antajani,
25_Ihala imettäjäńi,
Tuos miull' sot'isobani,
Kannas paloini-paita!"

- 1 align similar texts line-by-line
- 2 compute similarity measures summarizing the alignment
- 3 automatically find pairs of similar texts in the corpus
- 4 get large-scale overviews of the identified similarities



METHODS

- ① Verse similarity: bigram-based embedding + cosine similarity
- ② Poem similarity: weighted edit distance and alignment



VERSE SIMILARITY

Bag-of-bigrams representation:

| | ei | ip | ib | pä | bä | ä_ | lö | öy | öu | so | yä | uv | op | |
|------------------------|----|----|----|----|----|----|----|----|----|----|----|----|----|-----|
| eipä_löyä_väinämöistä | 1 | 1 | 0 | 1 | 0 | 2 | 1 | 1 | 0 | 0 | 1 | 0 | 0 | ... |
| eibä_löuvä_väinämöistä | 1 | 0 | 1 | 0 | 1 | 2 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | ... |
| eipä_sopi_väinämöistä | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | ... |

- using d most frequent bigrams (e.g. $d = 300$)
- insensitive to word order!



VERSE SIMILARITY

Cosine similarity:

| | ei | ip | ib | pä | bä | ä_ | lö | öy | öu | so | yä | uv | op | |
|---|------------------------|----|----|----|----|----|----|----|----|----|----|----|----|---|
| x | eipä_löyä_väinämöistä | 1 | 1 | 0 | 1 | 0 | 2 | 1 | 1 | 0 | 0 | 1 | 0 | 0 |
| y | eibä_löuvä_väinämöistä | 1 | 0 | 1 | 0 | 1 | 2 | 1 | 0 | 1 | 0 | 0 | 1 | 0 |

$$\begin{aligned}\cos \angle(\mathbf{x}, \mathbf{y}) &= \frac{\mathbf{x} \cdot \mathbf{y}}{|\mathbf{x}| \cdot |\mathbf{y}|} = \frac{\sum_{i=1}^d x_i y_i}{\sqrt{\sum_{i=1}^d x_i^2} \sqrt{\sum_{i=1}^d y_i^2}} \\ &= \frac{1 \cdot 1 + 1 \cdot 0 + 0 \cdot 1 + 1 \cdot 0 + 0 \cdot 1 + 2 \cdot 2 + \dots}{\sqrt{1^2 + 1^2 + 0^2 + 1^2 + 0^2 + 2^2 + \dots} \sqrt{1^2 + 0^2 + 1^2 + 0^2 + 2^2 + \dots}} \\ &\approx 0.81\end{aligned}$$



VERSE SIMILARITY

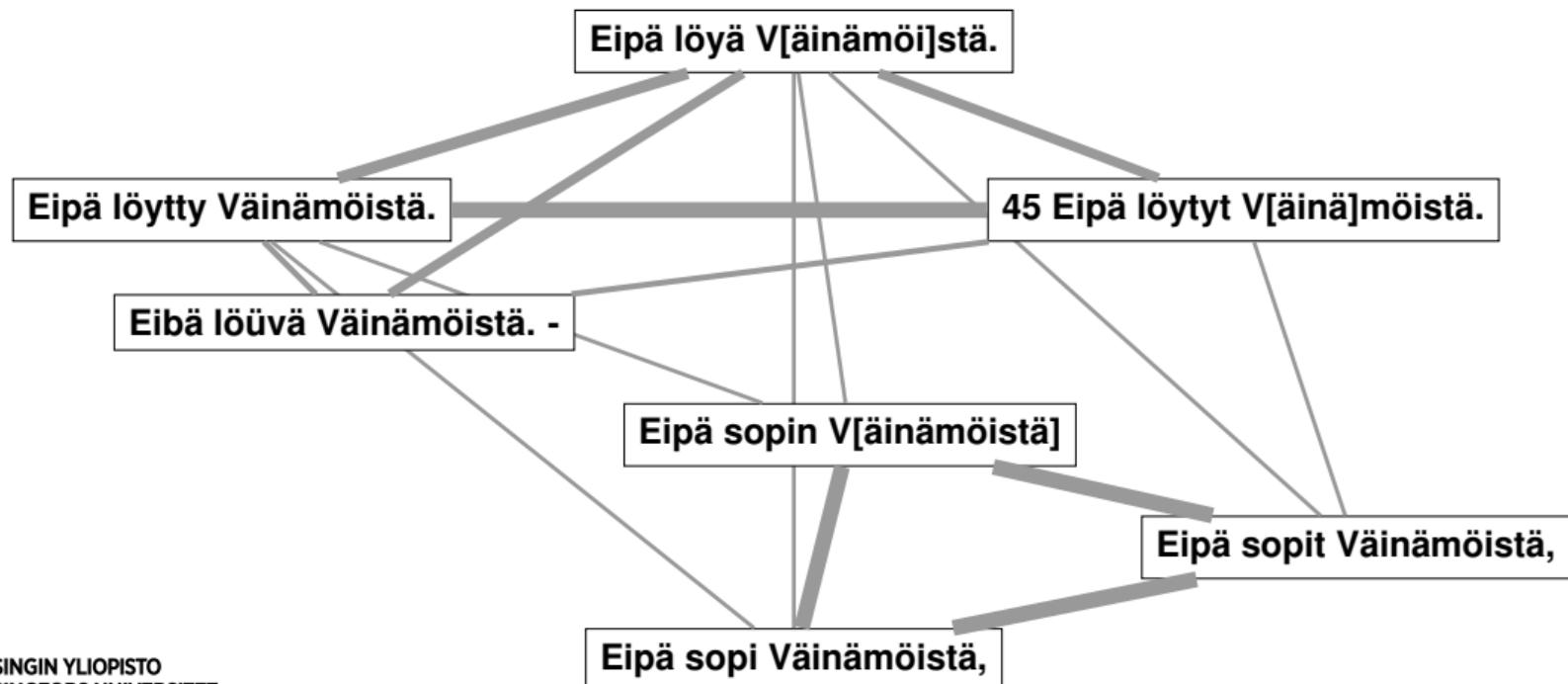
Query: Eipä löyä V[äinämöi]stä

Result:

| | |
|-------------------------------|----------|
| Eipä löytty Väinämöistä. | 0.863636 |
| 45 Eipä löytyt V[äinä]möistä. | 0.826869 |
| Eibä löuvä Väinämöistä. - | 0.826869 |
| Eipä löütöt Väinämöistä. | 0.800197 |
| Eipä sopi Väinämöistä, | 0.76277 |
| Eipä sopin V[äinämöistä] | 0.755742 |

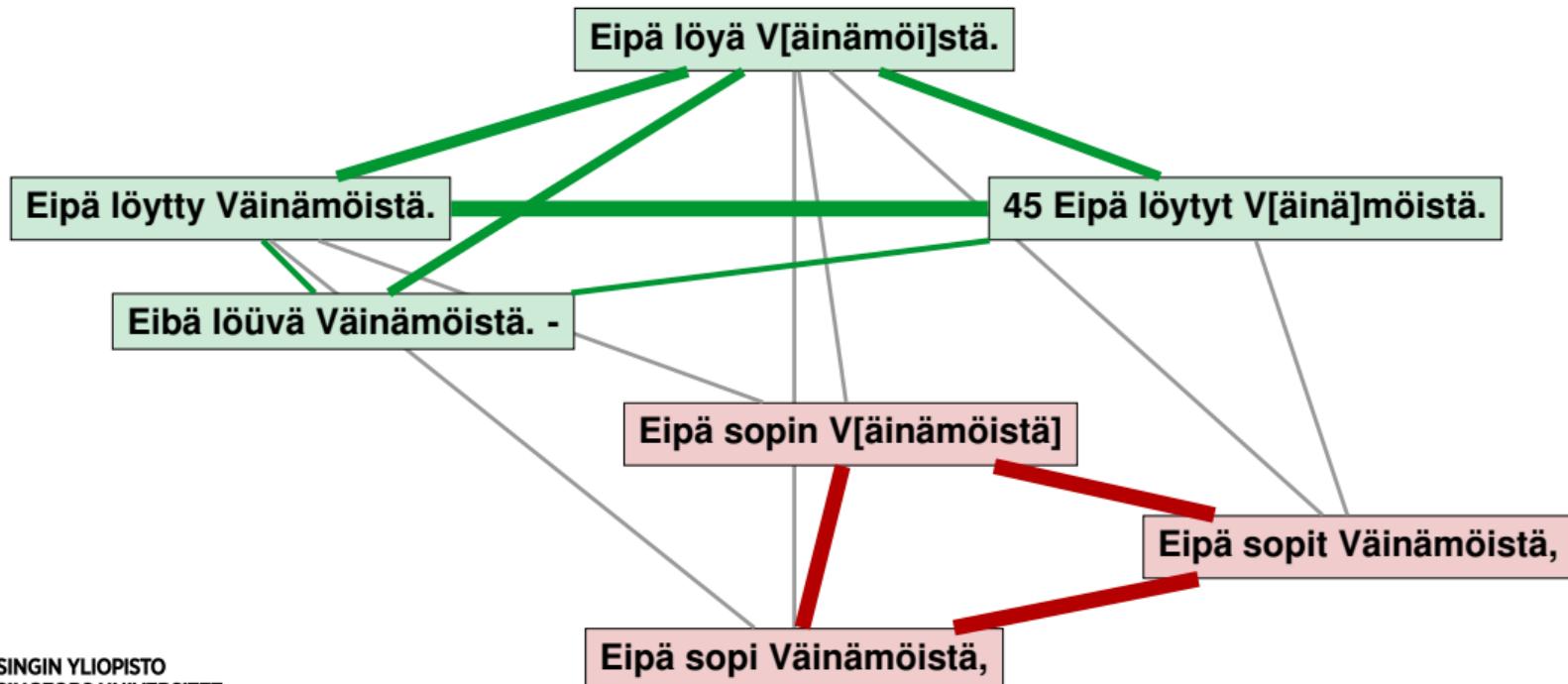


CLUSTERING





CLUSTERING





CLUSTERING¹

Eipä löyä V[äinämöi]stää.

[\[back to poem\]](#) [\[CSV\]](#) [\[map\]](#)

Cluster

SKVR.11.1.³² Eipä löyä V[äinämöi]stää.

skvr01100910
Vienna — Jyskyjärvi
1835 Lönnrot, Elias

1. Kertovat runot
1. Epäilka
1. Kansioon sotaa
1. Kansioon synty
1. Muodostusyruru
1. Sammon taonta
1. Viisimisen ammunta

SKVR.11.2.³² Eihä löyvä Väinämöistä. -

skvr01100920
Vienna — Jyskyjärvi
1872 Borenius, A. A.

1. Kertovat runot
1. Epäilka
1. Lähetä sotaa
1. Kansioon synty
1. Kilpailenta
1. Sammon taonta
1. Viisimisen ammunta

SKVR.11.16.³³ Eipä löytty Väinämöistä.

skvr01100160
Vienna — Jyskyjärvi
1889 Meriläinen, Heikki

1. Kertovat runot
1. Epäilka
1. Muodostusyruru
1. Viisimisen ammunta
7. Erilaisia pieniä runuja
1. Sanomiskuja
1. Janakkos on jokseen määrä, eikä miehen miehessä

SKVR.11.38.³⁴ Eipä löytöt Väinämöistä.

skvr01100580
Vienna — Vuokkiniemi
1871 Borenius, A. A.

1. Kertovat runot
1. Epäilka
1. Muodostusyruru
1. Sammon ryöttö
1. Sammon taonta
1. Viisimisen ammunta
2. Levyt
- 1.2. Väinöslautut
- 1.2. Loulut kodittomudesta ja kulkeneesta
1. Vienselle mäelle joitumatt

SKVR.11.38.a.³⁴ 45 Eipä löytty V[äinämöi]stää.

- 51 Talvet lylyn livulla;
52 Eipä löyä V[äinämöi]stää.
53 Min seppol[an pajahan]:
5 similar i sie se[ppo] Ilm[ollinen],
55 55 Takoja iän [ikuinen],

¹ Maciej Janicki, Kati Kallio, and Mari Sarv. “Exploring Finnic written oral folk poetry through string similarity”. In: *Digital Scholarship in the Humanities* (in press).



METHODS

- ① Verse similarity: bigram-based embedding + cosine similarity
- ② Poem similarity: weighted edit distance and alignment



EDIT DISTANCE

Given two words, e.g. **plotting** and **poetry**, find the minimum number of *edit operations* needed to transform one word into the other.

An *edit operation* is a:

- insertion
- deletion
- substitution

of one character.

The optimal transformation is:

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| p | I | o | t | t | i | n | g |
| p | o | e | t | r | y | | |

 with edit distance 5.



EDIT DISTANCE

| | y_1 | y_2 | y_3 | y_4 | y_5 | y_6 |
|-------|-------|-------|-------|-------|-------|-------|
| | p | o | e | t | r | y |
| x_1 | p | | | | | |
| x_2 | i | | | | | |
| x_3 | o | | | | | |
| x_4 | t | | | | | |
| x_5 | t | | | | | |
| x_6 | i | | | | | |
| x_7 | n | | | | | |
| x_8 | g | | | | | |

$$d_{i,j} = \min \left\{ \begin{array}{l} d_{i-1,j} + 1 \\ d_{i,j-1} + 1 \\ d_{i-1,j-1} + \delta(x_i, y_j) \end{array} \right\}$$

where:

$$\delta(x_i, y_j) = \begin{cases} 0 & \text{if } x_i = y_j \\ 1 & \text{otherwise} \end{cases}$$



EDIT DISTANCE

| | y_1 | y_2 | y_3 | y_4 | y_5 | y_6 |
|-------|-------|-------|-------|-------|-------|-------|
| | p | o | e | t | r | y |
| x_1 | p | 1 | | | | |
| x_2 | i | 2 | | | | |
| x_3 | o | 3 | | | | |
| x_4 | t | 4 | | | | |
| x_5 | t | 5 | | | | |
| x_6 | i | 6 | | | | |
| x_7 | n | 7 | | | | |
| x_8 | g | 8 | | | | |

$$d_{i,j} = \min \left\{ \begin{array}{l} d_{i-1,j} + 1 \\ d_{i,j-1} + 1 \\ d_{i-1,j-1} + \delta(x_i, y_j) \end{array} \right\}$$

where:

$$\delta(x_i, y_j) = \begin{cases} 0 & \text{if } x_i = y_j \\ 1 & \text{otherwise} \end{cases}$$



EDIT DISTANCE

| | y_1 | y_2 | y_3 | y_4 | y_5 | y_6 |
|-------|-------|-------|-------|-------|-------|-------|
| | p | o | e | t | r | y |
| x_1 | p | 1 | 0 | | | |
| x_2 | i | 2 | | | | |
| x_3 | o | 3 | | | | |
| x_4 | t | 4 | | | | |
| x_5 | t | 5 | | | | |
| x_6 | i | 6 | | | | |
| x_7 | n | 7 | | | | |
| x_8 | g | 8 | | | | |

$$d_{i,j} = \min \left\{ \begin{array}{l} d_{i-1,j} + 1 \\ d_{i,j-1} + 1 \\ d_{i-1,j-1} + \delta(x_i, y_j) \end{array} \right\}$$

where:

$$\delta(x_i, y_j) = \begin{cases} 0 & \text{if } x_i = y_j \\ 1 & \text{otherwise} \end{cases}$$



EDIT DISTANCE

| | | y_1 | y_2 | y_3 | y_4 | y_5 | y_6 |
|-------|----------|----------|----------|----------|----------|----------|----------|
| | | p | o | e | t | r | y |
| x_1 | p | 0 | 1 | 2 | 3 | 4 | 5 |
| x_2 | i | 1 | 0 | 1 | 2 | 3 | 4 |
| x_3 | o | 2 | | | | | |
| x_4 | t | 3 | | | | | |
| x_5 | t | 4 | | | | | |
| x_6 | i | 5 | | | | | |
| x_7 | n | 6 | | | | | |
| x_8 | g | 7 | | | | | |
| | | 8 | | | | | |

$$d_{i,j} = \min \left\{ \begin{array}{l} d_{i-1,j} + 1 \\ d_{i,j-1} + 1 \\ d_{i-1,j-1} + \delta(x_i, y_j) \end{array} \right\}$$

where:

$$\delta(x_i, y_j) = \begin{cases} 0 & \text{if } x_i = y_j \\ 1 & \text{otherwise} \end{cases}$$



EDIT DISTANCE

| | | y_1 | y_2 | y_3 | y_4 | y_5 | y_6 |
|-------|----------|----------|----------|----------|----------|----------|----------|
| | | p | o | e | t | r | y |
| | | 0 | 1 | 2 | 3 | 4 | 5 |
| x_1 | p | 1 | 0 | 1 | 2 | 3 | 4 |
| x_2 | i | 2 | 1 | 1 | 2 | 3 | 4 |
| x_3 | o | 3 | | | | | |
| x_4 | t | 4 | | | | | |
| x_5 | t | 5 | | | | | |
| x_6 | i | 6 | | | | | |
| x_7 | n | 7 | | | | | |
| x_8 | g | 8 | | | | | |

$$d_{i,j} = \min \left\{ \begin{array}{l} d_{i-1,j} + 1 \\ d_{i,j-1} + 1 \\ d_{i-1,j-1} + \delta(x_i, y_j) \end{array} \right\}$$

where:

$$\delta(x_i, y_j) = \begin{cases} 0 & \text{if } x_i = y_j \\ 1 & \text{otherwise} \end{cases}$$



EDIT DISTANCE

| | y_1 | y_2 | y_3 | y_4 | y_5 | y_6 |
|-------|----------|----------|----------|----------|----------|----------|
| | p | o | e | t | r | y |
| x_1 | p | 0 | 1 | 2 | 3 | 4 |
| x_2 | i | 1 | 0 | 1 | 2 | 3 |
| x_3 | o | 2 | 1 | 1 | 2 | 3 |
| x_4 | t | 3 | 2 | 1 | 2 | 3 |
| x_5 | t | 4 | | | | |
| x_6 | i | 5 | | | | |
| x_7 | n | 6 | | | | |
| x_8 | g | 7 | | | | |
| | | 8 | | | | |

$$d_{i,j} = \min \left\{ \begin{array}{l} d_{i-1,j} + 1 \\ d_{i,j-1} + 1 \\ d_{i-1,j-1} + \delta(x_i, y_j) \end{array} \right\}$$

where:

$$\delta(x_i, y_j) = \begin{cases} 0 & \text{if } x_i = y_j \\ 1 & \text{otherwise} \end{cases}$$



EDIT DISTANCE

| | y_1 | y_2 | y_3 | y_4 | y_5 | y_6 |
|-------|----------|----------|----------|----------|----------|----------|
| | p | o | e | t | r | y |
| x_1 | p | 0 | 1 | 2 | 3 | 4 |
| x_2 | i | 1 | 0 | 1 | 2 | 3 |
| x_3 | o | 2 | 1 | 1 | 2 | 3 |
| x_4 | t | 3 | 2 | 1 | 2 | 3 |
| x_5 | t | 4 | 3 | 2 | 2 | 2 |
| x_6 | i | 5 | | | | |
| x_7 | n | 6 | | | | |
| x_8 | g | 7 | | | | |
| | | 8 | | | | |

$$d_{i,j} = \min \left\{ \begin{array}{l} d_{i-1,j} + 1 \\ d_{i,j-1} + 1 \\ d_{i-1,j-1} + \delta(x_i, y_j) \end{array} \right\}$$

where:

$$\delta(x_i, y_j) = \begin{cases} 0 & \text{if } x_i = y_j \\ 1 & \text{otherwise} \end{cases}$$



EDIT DISTANCE

| | | y_1 | y_2 | y_3 | y_4 | y_5 | y_6 |
|-------|----------|----------|----------|----------|----------|----------|----------|
| | | p | o | e | t | r | y |
| | | 0 | 1 | 2 | 3 | 4 | 5 |
| x_1 | p | 1 | 0 | 1 | 2 | 3 | 4 |
| x_2 | i | 2 | 1 | 1 | 2 | 3 | 4 |
| x_3 | o | 3 | 2 | 1 | 2 | 3 | 4 |
| x_4 | t | 4 | 3 | 2 | 2 | 3 | 4 |
| x_5 | t | 5 | 4 | 3 | 3 | 2 | 3 |
| x_6 | i | 6 | | | | | |
| x_7 | n | 7 | | | | | |
| x_8 | g | 8 | | | | | |

$$d_{i,j} = \min \left\{ \begin{array}{l} d_{i-1,j} + 1 \\ d_{i,j-1} + 1 \\ d_{i-1,j-1} + \delta(x_i, y_j) \end{array} \right\}$$

where:

$$\delta(x_i, y_j) = \begin{cases} 0 & \text{if } x_i = y_j \\ 1 & \text{otherwise} \end{cases}$$



EDIT DISTANCE

| | | y_1 | y_2 | y_3 | y_4 | y_5 | y_6 |
|-------|----------|----------|----------|----------|----------|----------|----------|
| | | p | o | e | t | r | y |
| | | 0 | 1 | 2 | 3 | 4 | 5 |
| x_1 | p | 1 | 0 | 1 | 2 | 3 | 4 |
| x_2 | i | 2 | 1 | 1 | 2 | 3 | 4 |
| x_3 | o | 3 | 2 | 1 | 2 | 3 | 4 |
| x_4 | t | 4 | 3 | 2 | 2 | 3 | 4 |
| x_5 | t | 5 | 4 | 3 | 3 | 2 | 3 |
| x_6 | i | 6 | 5 | 4 | 4 | 3 | 3 |
| x_7 | n | 7 | | | | | |
| x_8 | g | 8 | | | | | |

$$d_{i,j} = \min \left\{ \begin{array}{l} d_{i-1,j} + 1 \\ d_{i,j-1} + 1 \\ d_{i-1,j-1} + \delta(x_i, y_j) \end{array} \right\}$$

where:

$$\delta(x_i, y_j) = \begin{cases} 0 & \text{if } x_i = y_j \\ 1 & \text{otherwise} \end{cases}$$



EDIT DISTANCE

| | | y_1 | y_2 | y_3 | y_4 | y_5 | y_6 |
|-------|----------|----------|----------|----------|----------|----------|----------|
| | | p | o | e | t | r | y |
| | | 0 | 1 | 2 | 3 | 4 | 5 |
| x_1 | p | 1 | 0 | 1 | 2 | 3 | 4 |
| x_2 | i | 2 | 1 | 1 | 2 | 3 | 4 |
| x_3 | o | 3 | 2 | 1 | 2 | 3 | 4 |
| x_4 | t | 4 | 3 | 2 | 2 | 2 | 3 |
| x_5 | t | 5 | 4 | 3 | 3 | 2 | 3 |
| x_6 | i | 6 | 5 | 4 | 4 | 3 | 3 |
| x_7 | n | 7 | 6 | 5 | 5 | 4 | 4 |
| x_8 | g | 8 | | | | | |

$$d_{i,j} = \min \left\{ \begin{array}{l} d_{i-1,j} + 1 \\ d_{i,j-1} + 1 \\ d_{i-1,j-1} + \delta(x_i, y_j) \end{array} \right\}$$

where:

$$\delta(x_i, y_j) = \begin{cases} 0 & \text{if } x_i = y_j \\ 1 & \text{otherwise} \end{cases}$$



EDIT DISTANCE

| | | y_1 | y_2 | y_3 | y_4 | y_5 | y_6 |
|-------|----------|----------|----------|----------|----------|----------|----------|
| | | p | o | e | t | r | y |
| | | 0 | 1 | 2 | 3 | 4 | 5 |
| x_1 | p | 1 | 0 | 1 | 2 | 3 | 4 |
| x_2 | i | 2 | 1 | 1 | 2 | 3 | 4 |
| x_3 | o | 3 | 2 | 1 | 2 | 3 | 4 |
| x_4 | t | 4 | 3 | 2 | 2 | 2 | 3 |
| x_5 | t | 5 | 4 | 3 | 3 | 2 | 3 |
| x_6 | i | 6 | 5 | 4 | 4 | 3 | 3 |
| x_7 | n | 7 | 6 | 5 | 5 | 4 | 4 |
| x_8 | g | 8 | 7 | 6 | 6 | 5 | 5 |

$$d_{i,j} = \min \left\{ \begin{array}{l} d_{i-1,j} + 1 \\ d_{i,j-1} + 1 \\ d_{i-1,j-1} + \delta(x_i, y_j) \end{array} \right\}$$

where:

$$\delta(x_i, y_j) = \begin{cases} 0 & \text{if } x_i = y_j \\ 1 & \text{otherwise} \end{cases}$$



EDIT DISTANCE

| | y_1 | y_2 | y_3 | y_4 | y_5 | y_6 |
|-------|----------|----------|----------|----------|----------|----------|
| | p | o | e | t | r | y |
| | 0 | ← | ← | ← | ← | ← |
| x_1 | p | ↑ | ↖ | ← | ← | ← |
| x_2 | l | ↑ | ↑ | ↖ | ← | ← |
| x_3 | o | ↑ | ↑ | ↖ | ↖ | ← |
| x_4 | t | ↑ | ↑ | ↑ | ↖ | ← |
| x_5 | t | ↑ | ↑ | ↑ | ↖ | ← |
| x_6 | i | ↑ | ↑ | ↑ | ↑ | ↖ |
| x_7 | n | ↑ | ↑ | ↑ | ↑ | ↖ |
| x_8 | g | ↑ | ↑ | ↑ | ↑ | ↑ |

$$d_{i,j} = \min \left\{ \begin{array}{l} d_{i-1,j} + 1 \\ d_{i,j-1} + 1 \\ d_{i-1,j-1} + \delta(x_i, y_j) \end{array} \right\}$$

where:

$$\delta(x_i, y_j) = \begin{cases} 0 & \text{if } x_i = y_j \\ 1 & \text{otherwise} \end{cases}$$

p l o t t i n g
p o e t r y



EDIT DISTANCE

| | y_1 | y_2 | y_3 | y_4 | y_5 | y_6 |
|-------|----------|----------|----------|----------|----------|----------|
| | p | o | e | t | r | y |
| 0 | ← | ← | ← | ← | ← | ← |
| x_1 | p | ↑ ↗ | ← | ← | ← | ← |
| x_2 | l | ↑ ↗ | ← | ← | ← | ← |
| x_3 | o | ↑ ↗ | ↖ | ← | ← | ← |
| x_4 | t | ↑ ↗ | ↖ | ↖ | ← | ← |
| x_5 | t | ↑ ↗ | ↑ ↗ | ↖ | ← | ← |
| x_6 | i | ↑ ↗ | ↑ ↗ | ↑ ↗ | ↖ | ← |
| x_7 | n | ↑ ↗ | ↑ ↗ | ↑ ↗ | ↑ ↗ | ↖ |
| x_8 | g | ↑ ↗ | ↑ ↗ | ↑ ↗ | ↑ ↗ | ↑ ↗ |

$$d_{i,j} = \min \left\{ \begin{array}{l} d_{i-1,j} + 1 \\ d_{i,j-1} + 1 \\ d_{i-1,j-1} + \delta(x_i, y_j) \end{array} \right\}$$

where:

$$\delta(x_i, y_j) = \begin{cases} 0 & \text{if } x_i = y_j \\ 1 & \text{otherwise} \end{cases}$$

p l o t t i n g
p o e t r y



MAXIMUM-WEIGHT ALIGNMENT

Similar as edit distance, but:

- insertions and deletions have weight 0,
- $\delta(x_i, y_j)$ is the **similarity** between x_i and y_j (between 0 and 1),
- we are looking for an alignment with **maximum** total weight.

For example, let:

$$\delta(x_i, y_j) = \begin{cases} 1 & \text{if } x_i = y_j \\ 0.8 & \text{for } (i, y) \text{ and } (r, l) \\ 0.5 & \text{for } (t, r), (t, l) \text{ and } (e, i) \\ 0 & \text{otherwise} \end{cases}$$

The optimal transformation is:  with total weight 4.3



MAXIMUM-WEIGHT ALIGNMENT

The optimal transformation is: p I o t t i n g with total weight 4.3
 p o e t r y

$$s_{\text{raw}} = 4.3 \quad - \text{raw sim.}$$

$$s_l = \frac{s_{\text{raw}}}{|\mathbf{x}|} = \frac{4.3}{8} \approx 0.54 \quad - \text{left-normalized sim.}$$

$$s_r = \frac{s_{\text{raw}}}{|\mathbf{y}|} = \frac{4.3}{6} \approx 0.72 \quad - \text{right-normalized sim.}$$

$$s = \frac{2}{\frac{1}{s_l} + \frac{1}{s_r}} = \frac{2 \cdot s_{\text{raw}}}{|\mathbf{x}| + |\mathbf{y}|} \approx 0.61 \quad - \text{avg. normalized sim.}$$



| | | | | |
|----------------------------|-----|---|-----|-----|
| yht_ei_kuttsun_lemmingästä | .88 | 1 | .71 | .93 |
| ruiot_reellä_rembutteli | | | .84 | .52 |
| rammat_rattsahin_ajeli | | | .57 | .86 |
| sogeat_venosin_souti | | | | |
| yht_ei_kuttsu_lemmingästä | | | | |
| lemmingäne_on_pilopoiga | | | | |
| piiloja_on_piilemässä | | | | |
| pahoja_pagenemassa | | | | |

$$d_{i,j} = \max \left\{ \begin{array}{l} d_{i-1,j} \\ d_{i,j-1} \\ d_{i-1,j-1} + \delta(x_i, y_j) \end{array} \right\}$$



| | | | | | |
|---------------------------|---|-----|---|------------------------------|--|
| kuttsu_verisogeat | 0 | 0 | 0 | yht_ei_kuttsun_lemmingäistä | |
| ruiot_reellä_rembutteli | 0 | | | ruijot_ne_reillä_reissuaabi | |
| rammat_rattsahin_ajeli | 0 | | | rammat_rattsahin_ajeli | |
| sogeat_venosin_souti | 0 | | | sogiat_venozin_soudi | |
| yht_ei_kuttsu_lemmingästä | 0 | .88 | | lemmingäin_on_pilopoga_piilo | |
| lemmingäne_on_pilopoga | 0 | | | pillojah_on_piilemässä | |
| piiloja_on_piilemässä | 0 | | | pahojah_pagenemassa | |
| pahoja_pagenemassa | 0 | | | | |

.71
.84 .52
.57 .86
.93

$$d_{i,j} = \max \left\{ \begin{array}{l} d_{i-1,j} \\ d_{i,j-1} \\ d_{i-1,j-1} + \delta(x_i, y_j) \end{array} \right\}$$



| | | | | |
|---------------------------|---|-----|-----------------------------|---------|
| kuttsu_verisogeat | 0 | 0 | yht_ei_kuttsun_lemmingäistä | 0 |
| ruiot_reellä_rembutteli | 0 | 0 | rujot_ne_reillä_reissuaabi | 0 |
| rammat_rattsahin_ajeli | 0 | 0 | rammat_rattsahin_ajeli | 0 |
| sogeat_venosin_souti | 0 | 0 | sogiat_venozin_soudi | 0 |
| yht_ei_kuttsu_lemmingästä | 0 | .88 | 1 | .71 |
| lemmingäne_on_pilopoiga | 0 | | | .84 .52 |
| piiloja_on_piilemässä | 0 | | | .57 .86 |
| pahoja_pagenemassa | 0 | | | .93 |

$$d_{i,j} = \max \left\{ \begin{array}{l} d_{i-1,j} \\ d_{i,j-1} \\ d_{i-1,j-1} + \delta(x_i, y_j) \end{array} \right\}$$



$$d_{i,j} = \max \left\{ \begin{array}{c} d_{i-1,j} \\ d_{i,j-1} \\ d_{i-1,j-1} + \delta(x_i, y_j) \end{array} \right\}$$



| | | | | | |
|---------------------------|---|---|---|-----------------------------|-------------|
| kuttsu_verisogeat | 0 | 0 | 0 | yht_ei_kuttsun_lemmingäistä | .88 |
| ruiot_reellä_rembutteli | 0 | 0 | 0 | rujot_ne_reillä_reissuaabi | |
| rammat_rattsahin_ajeli | 0 | 0 | 0 | rammat_rattsahin_ajeli | |
| sogeat_venosin_souti | 0 | 0 | 0 | sogiat_venozin_soudi | .71 |
| yht_ei_kuttsu_lemmingästä | 0 | 0 | 0 | pillojah_on_pillemässä | |
| lemmingäne_on_pilipoiga | 0 | 0 | 0 | lemmingäin_on_poiga_pillo | |
| piiloja_on_piilemässä | 0 | 0 | 0 | piilojah_on_pillemässä | .52 |
| pahoja_pagenemassa | 0 | 0 | 0 | pahojah_pagenemassa | .57 .86 .93 |

$$d_{i,j} = \max \left\{ \begin{array}{l} d_{i-1,j} \\ d_{i,j-1} \\ d_{i-1,j-1} + \delta(x_i, y_j) \end{array} \right\}$$



| | | | | | | |
|---------------------------|---|-----|---|----------------------------|-----|--|
| kuttsu_verisogeat | 0 | 0 | 0 | yht_ei_kuttsun_lemmingästä | | |
| ruiot_reellä_rembutteli | 0 | 0 | 0 | rujot_ne_reillä_reissuaabi | | |
| rammat_rattsahin_ajeli | 0 | 1 | 0 | rammat_rattsahin_ajeli | | |
| sogeat_venosin_souti | 0 | 0 | 1 | sogiat_venozin_soudi | | |
| yht_ei_kuttsu_lemmingästä | 0 | .88 | | | | |
| lemmingäne_on_pilopoiga | 0 | | | | .52 | |
| piiloja_on_piilemässä | 0 | | | | .86 | |
| pahoja_pagenemassa | 0 | | | | .93 | |

$$d_{i,j} = \max \left\{ \begin{array}{l} d_{i-1,j} \\ d_{i,j-1} \\ d_{i-1,j-1} + \delta(x_i, y_j) \end{array} \right\}$$



| | | | | | | | | |
|---------------------------|---|-----|-----|---|------|------|------|------|
| kuttsu_verisogeat | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| ruiot_reellä_rembutteli | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| rammat_rattsahin_ajeli | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| sogeat_venosin_souti | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| yht_ei_kuttsu_lemmingästä | 0 | .88 | .88 | 1 | 1.71 | 1.71 | 1.71 | 1.71 |
| lemmingäne_on_piiopoiga | 0 | | | | | | .52 | |
| piiloja_on_piilemässä | 0 | | | | | | .86 | |
| pahoja_pagenemassa | 0 | | | | | | .93 | |

$$d_{i,j} = \max \left\{ \begin{array}{l} d_{i-1,j} \\ d_{i,j-1} \\ d_{i-1,j-1} + \delta(x_i, y_j) \end{array} \right\}$$



| | | | | | | | | |
|---------------------------|---|-----|-----|---|------|------|------|------|
| | | | | | | | | |
| kuttsu_verisogeat | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| ruiot_reellä_rembutteli | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| rammat_rattsahin_ajeli | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| sogeat_venosin_souti | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| yht_ei_kuttsu_lemmingästä | 0 | .88 | .88 | 1 | 1.71 | 1.71 | 1.71 | 1.71 |
| lemmingäne_on_piiłopoiga | 0 | .88 | .88 | 1 | 1.71 | 2.55 | 2.55 | 2.55 |
| piiloja_on_piilemässä | 0 | | | | | .86 | | |
| pahoja_pagenemassa | 0 | | | | | | .93 | |

$$d_{i,j} = \max \left\{ \begin{array}{l} d_{i-1,j} \\ d_{i,j-1} \\ d_{i-1,j-1} + \delta(x_i, y_j) \end{array} \right\}$$



| | | | | | | | | |
|---------------------------|---|-----|-----|---|------|------|------|------|
| | | | | | | | | |
| kuttsu_verisogeat | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| ruiot_reellä_rembutteli | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| rammat_rattsahin_ajeli | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| sogeat_venosin_souti | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 |
| yht_ei_kuttsu_lemmingästä | 0 | .88 | .88 | 1 | 1.71 | 1.71 | 1.71 | 1.71 |
| lemmingäne_on_piiopoiga | 0 | .88 | .88 | 1 | 1.71 | 2.55 | 2.55 | 2.55 |
| piiloja_on_piilemässä | 0 | .88 | .88 | 1 | 1.71 | 2.55 | 3.41 | 3.41 |
| pahoja_pagenemassa | 0 | | | | | | | .93 |

$$d_{i,j} = \max \left\{ \begin{array}{l} d_{i-1,j} \\ d_{i,j-1} \\ d_{i-1,j-1} + \delta(x_i, y_j) \end{array} \right\}$$



| | | | | | | | | |
|---------------------------|---|-----|-----|---|------|------|------|------|
| | | | | | | | | |
| kuttsu_verisogeat | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| ruiot_reellä_rembutteli | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| rammat_rattsahin_ajeli | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| sogeat_venosin_souti | 0 | 0 | 0 | 1 | 1.71 | 1.71 | 1.71 | 1.71 |
| yht_ei_kuttsu_lemmingästä | 0 | .88 | .88 | 1 | 1.71 | 1.71 | 1.71 | 1.71 |
| lemmingäne_on_piiopoiga | 0 | .88 | .88 | 1 | 1.71 | 2.55 | 2.55 | 2.55 |
| piiloja_on_piilemässä | 0 | .88 | .88 | 1 | 1.71 | 2.55 | 3.41 | 3.41 |
| pahoja_pagenemassa | 0 | .88 | .88 | 1 | 1.71 | 2.55 | 3.41 | 4.34 |

$$d_{i,j} = \max \left\{ \begin{array}{l} d_{i-1,j} \\ d_{i,j-1} \\ d_{i-1,j-1} + \delta(x_i, y_j) \end{array} \right\}$$



| | | | | | | | | |
|---------------------------|---|-----|-----|---|------|------|------|------|
| | | | | | | | | |
| kuttsu_verisogeat | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| ruiot_reellä_rembutteli | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| rammat_rattsahin_ajeli | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| sogeat_venosin_souti | 0 | 0 | 0 | 1 | 1.71 | 1.71 | 1.71 | 1.71 |
| yht_ei_kuttsu_lemmingästä | 0 | .88 | .88 | 1 | 1.71 | 1.71 | 1.71 | 1.71 |
| lemmingäne_on_piiopoiga | 0 | .88 | .88 | 1 | 1.71 | 2.55 | 2.55 | 2.55 |
| piiloja_on_piilemässä | 0 | .88 | .88 | 1 | 1.71 | 2.55 | 3.41 | 3.41 |
| pahoja_pagenemassa | 0 | .88 | .88 | 1 | 1.71 | 2.55 | 3.41 | 4.34 |

$$d_{i,j} = \max \left\{ \begin{array}{l} d_{i-1,j} \\ d_{i,j-1} \\ d_{i-1,j-1} + \delta(x_i, y_j) \end{array} \right\}$$
$$s_{\text{raw}} = 4.34$$
$$s_l = \frac{4.34}{8} \approx 0.54$$
$$s_r = \frac{4.34}{7} = 0.62$$
$$s = \frac{2 \cdot 4.34}{8 + 7} \approx 0.58$$



kuttsu_verisogeat
 ruiot_reellä_rembutteli
 rammat_rattsa hin_ajeli
 sogeat_venosin_souti
 yht_ei_kuttsu_lemmingästä
 lemmingäne_on_pilopoiga
 piloja_on_piilemässä
 pahoja_pagenemassa

→ → → → → 0
 → → → / → → ↑ yht_ei_kuttsun_lemmingästä
 → → → ↑ → → ↑ ruiot_ne_reillä_reissa abi
 → → → / → → ↑ rammat_rattsa hin_ajeli
 → → → / ↑ → → ↑ sogiat_venozin_soudi
 → → / → ↑ ↑ → → ↑ lemmingäin_on_poiga_piloo
 → / ↑ → ↑ ↑ → → ↑ piloja h_on_pillemässä
 / ↑ ↑ → ↑ ↑ → → ↑ pahoja_h_pagenemassa

$$d_{i,j} = \max \left\{ \begin{array}{l} d_{i-1,j} \\ d_{i,j-1} \\ d_{i-1,j-1} + \delta(x_i, y_j) \end{array} \right\}$$

$$\begin{aligned}
 s_{\text{raw}} &= 4.34 \\
 s_l &= \frac{4.34}{8} \approx 0.54 \\
 s_r &= \frac{4.34}{7} = 0.62 \\
 s &= \frac{2 \cdot 4.34}{8 + 7} \approx 0.58
 \end{aligned}$$



kuttsu_verisogeat
ruiot_reellä_rembutteli
rammat_rattsahtin_ajeli
sogeat_venosin_souti
yht_ei_kuttsu_lemmingästä
lemmingäne_on_piiłopoiga
piiłoja_on_piiłemässä
pahoja_pagenemassa

$$d_{i,j} = \max \left\{ \begin{array}{c} d_{i-1,j} \\ d_{i,j-1} \\ d_{i-1,j-1} + \delta(x_i, y_j) \end{array} \right\}$$

$$S_{\text{raw}} = 4.34$$

$$s_l = \frac{4.34}{8} \approx 0.54$$

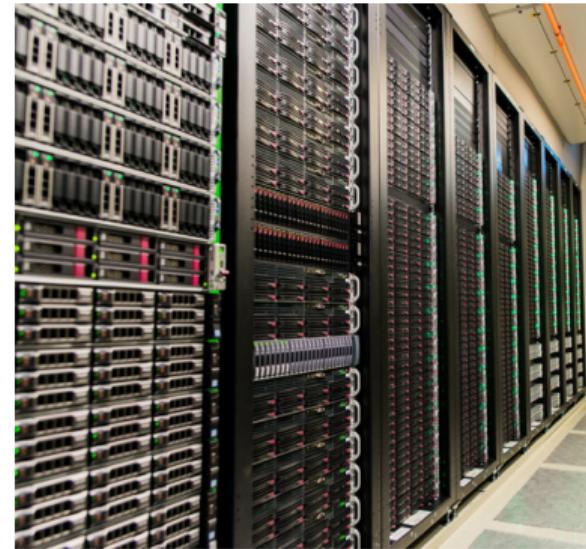
$$S_r = \frac{4.34}{7} = 0.62$$

$$s = \frac{2 \cdot 4.34}{8+7} \approx 0.58$$



MAXIMUM-WEIGHT ALIGNMENT

- compute on all pairs from among 275,181 poems
- optimization: use vector and matrix operations
 - compute entire rows at once
 - compute the alignment between one poem and all others at once
- run on a GPU
(a \$13,000 NVidia Tesla V100 with 32 GB memory)
- done within 64h!
- result: over 16 million pairs of similar poems
- criteria: $s_{\text{raw}} > 2$ and $\min\{s_l, s_r\} > 0.1$





RESULT: RUNOREGI

SKVR VII1 803.

skvr07108030
Laatokan Karjala (Raja-Karjala) — Suistamo
1897 Borenius, A. A.

Metadata

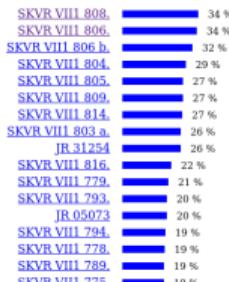
COL Bor.
ID 803.
INF Muuanto. Iivana Shemeikka ("Jehkin iivana"), 57 v.
LOC Suistamo.
OSA VIII
SGN III A, s. 304.
TMP -7/6 97.

Themes

¹ Kertovat runot
¹ Epitika
¹ Lemminkäisen virsi

Similar poems

[compare all]



| | |
|--|--------------------|
| 1884 Laatokan Karjala (Raja-Karjala) — Suistamo | Basilier, Hj |
| 1884 Laatokan Karjala (Raja-Karjala) — Suistamo | Basilier, Hj |
| 1884 Laatokan Karjala (Raja-Karjala) — Suistamo | Krohn, Kaarle |
| 1897 Laatokan Karjala (Raja-Karjala) — Suistamo | Borenius, A. A. |
| 1900 Laatokan Karjala (Raja-Karjala) — Suistamo | Härkönen, Iivo |
| 1915 Laatokan Karjala (Raja-Karjala) — Suistamo | Hannikainen, Lauri |
| 1897 Laatokan Karjala (Raja-Karjala) — Korpiseläki Potschtaroff, Wasille | |
| 1888 Laatokan Karjala (Raja-Karjala) — Suistamo | Relander, O. |
| 1935 Laatokan Karjala (Raja-Karjala) — Suistamo | Kähni, Martta |
| 1900 Laatokan Karjala (Raja-Karjala) — Korpiseläki Härkönen, Iivo | |
| 1884 Laatokan Karjala (Raja-Karjala) — Suistamo | Krohn, Kaarle |
| 1884 Laatokan Karjala (Raja-Karjala) — Suistamo | Basilier, Hj |
| 1935 Laatokan Karjala (Raja-Karjala) — Suistamo | Haavio, Martti |
| 1884 Laatokan Karjala (Raja-Karjala) — Suistamo | Krohn, Kaarle |
| 1900 Laatokan Karjala (Raja-Karjala) — Suistamo | Härkönen, Iivo |
| 1884 Laatokan Karjala (Raja-Karjala) — Suistamo | Basilier, Hj |
| 1909 Laatokan Karjala (Raja-Karjala) — Suistamo | Niemi, A. R. |

Text

[show verse-level themes] [show shared verses matrix]

Lemmingäzes lauletab.

- ² Savupa soarelli palaabi,
- ³ Niemenbä kylgyvöt kyööbi.
- ⁴ Toivoimba#1 palimozen tuleksi#2,#3,
- ⁵ Suur oli palimozen tuleksi#3;
- ⁶ 5 Tolvolb on#4 sodisavuksi,
- ⁷ Pień oli#5 sodisavuksi,
- ⁸ Osmatta on#6 olutta keitti,
- ⁹ *Kallervööba#7 kal'oivettä*
- ¹⁰ Yheksässä#8 ozranjyvässä,
- ¹¹ 10 Kaheksas#9 kagranjyvässä,
- ¹² Tulijillaba vierahilla.
- ¹³ *Laittobi viestit viizijillä,
- ¹⁴ Kutsutpa kuuzilla jageli*,
- ¹⁵ Kuttsuba rujot, kuttus rammatt,
- ¹⁶ 15 Kuttsubon perisoglat,
- ¹⁷ Kuttsuba ...
- ¹⁸ Yht' ei kuttsun Lemmingäistä.
- ¹⁹ *Rujot (ne) reillä reissuaabli,
- ²⁰ Rammat rattshän ajell,
- ²¹ 20 Sogiat venozin soudi,*
- ²² Lemmingän on poiga pillo
- ²³ Pillojah on pillemässä,#10
- ²⁴ Pahojah pagenemassa. #11
- ²⁵ *Hoib om moamo, kandajafäli,
- ²⁶ 25 Armas majton andajafäli,
- ²⁷ Ihlan imettäjäni,
- ²⁸ Ettsib om miul pelvi paidjal,
- ²⁹ Enneembä neidomä kuvottu,
- ³⁰ Kassabapeain#12 on kalkuteltu,
- ³¹ 30 Kamabas paloni paidja]."
- ³² "Minne (sie) lähtet, pojuguvoi?"
- ³³ "Lähtembä Väinöläni pidothe,
- ³⁴ Jumaliston juomingihe,
- ³⁵ Suuremuba synnin syömingihe."
- ³⁶ 35 "(N)Eläs (sie) lähte, pojuguvoi#13,
- ³⁷ Stielä sul on 3 surmoa."



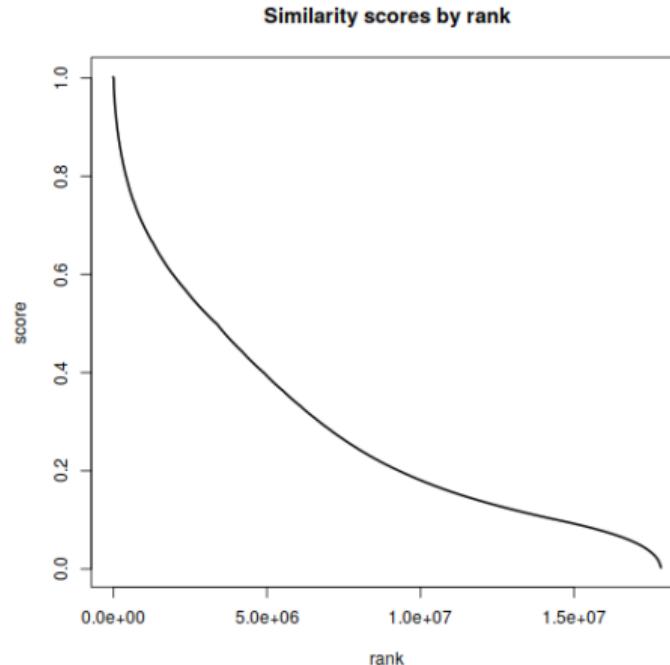
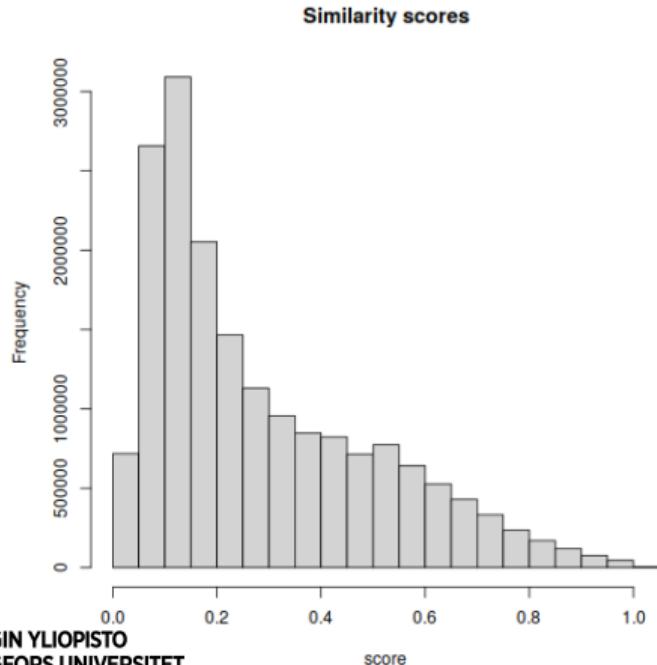
QUANTITATIVE ANALYSIS

Research questions:

- ① Types of relationships between similar texts
same singer / parish / poem type? temporal and geographical distance?
- ② Oral tradition ↔ printed works
- ③ Building type indices
does the similarity computation help improve and evaluate them?



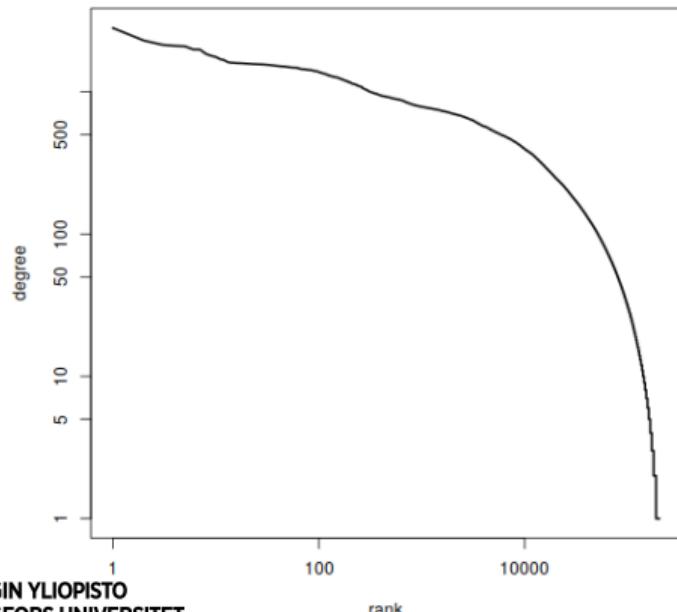
OVERVIEW



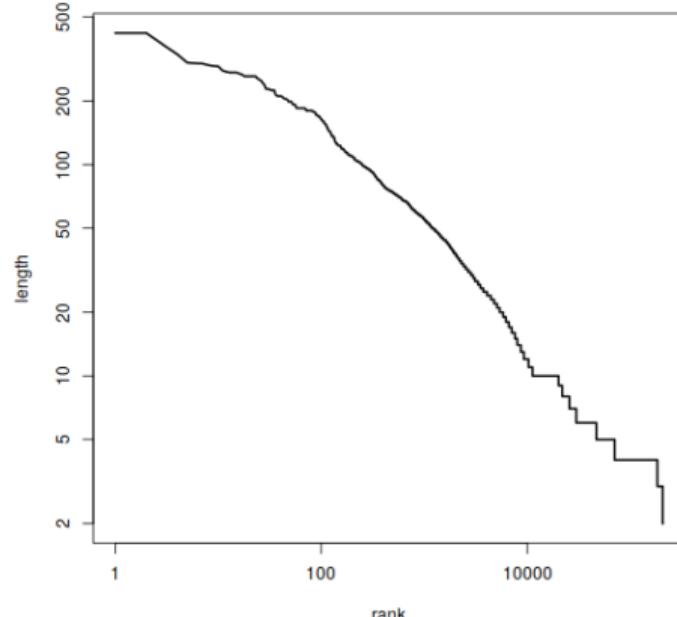


OVERVIEW

Node degree by rank



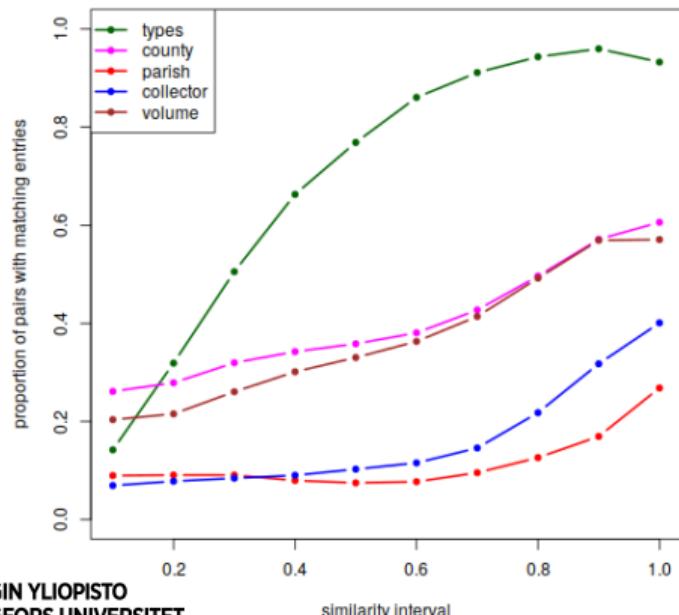
Poem length in nearly-identical poem pairs



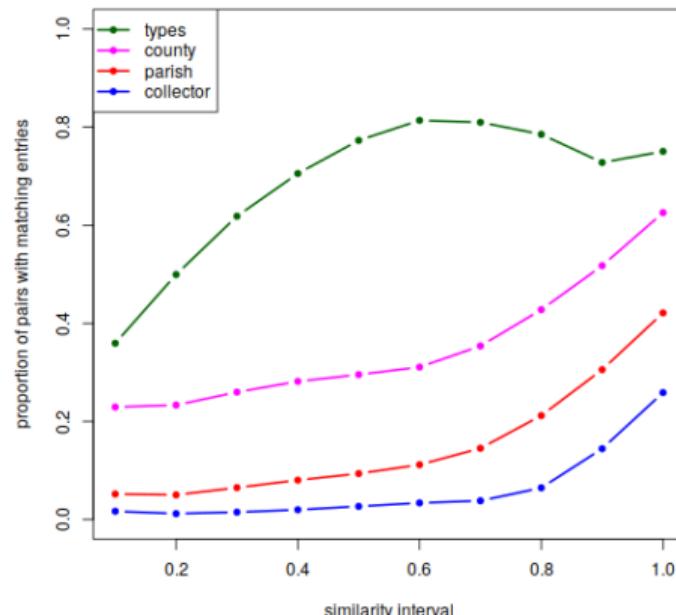


OVERVIEW

SKVR: Matching metadata entries for pairs of similar poems



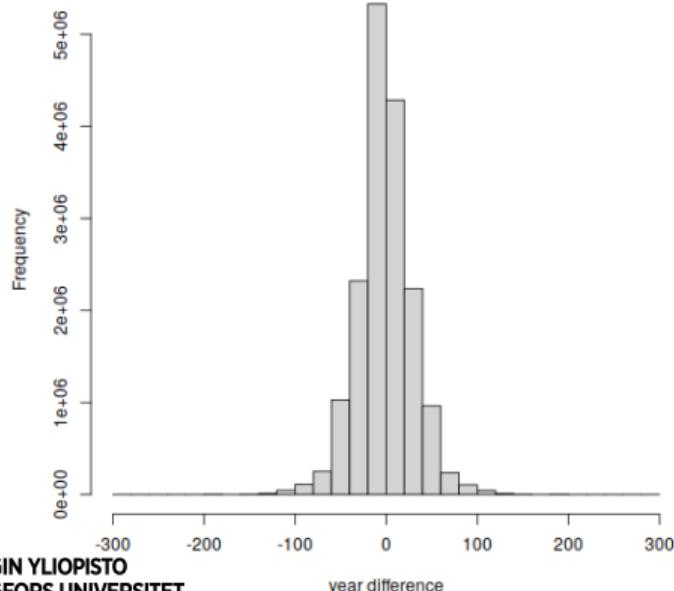
ERAB: Matching metadata entries for pairs of similar poems



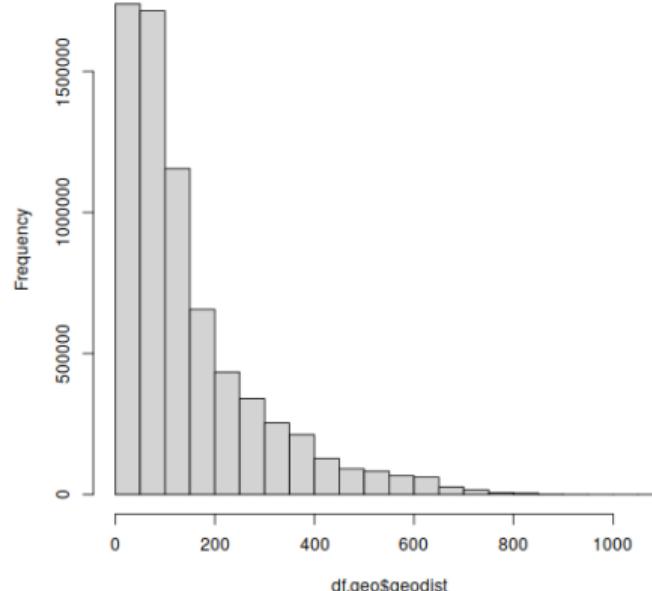


OVERVIEW

Histogram of year difference for similar poems



Histogram of distance (km) of similar poems





OUTLIER ANALYSIS

Similar poems with the largest temporal distance:

- charms (e.g. 1 2 3 4 5)
- proverbs (e.g. 1 2 3)
- epics: ‘The spider song’ (1)

Similar poems with the largest geographic distance:

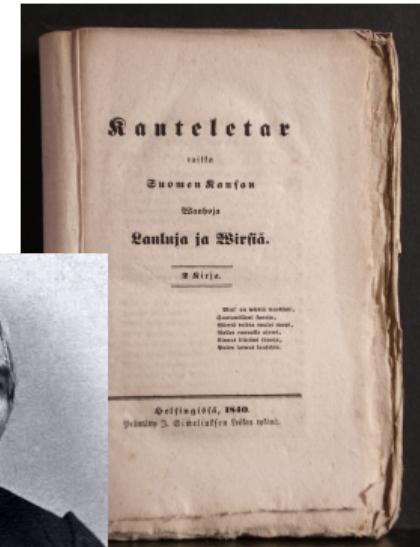
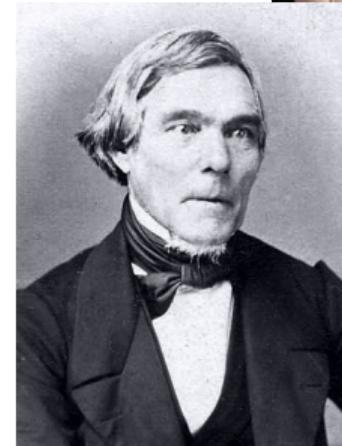
- charms (e.g. 1 2)
- lullabies (e.g. 1 2)
- children’s songs (e.g. 1 2)



ORAL-LITERARY RELATIONSHIPS

Example: Kanteletar

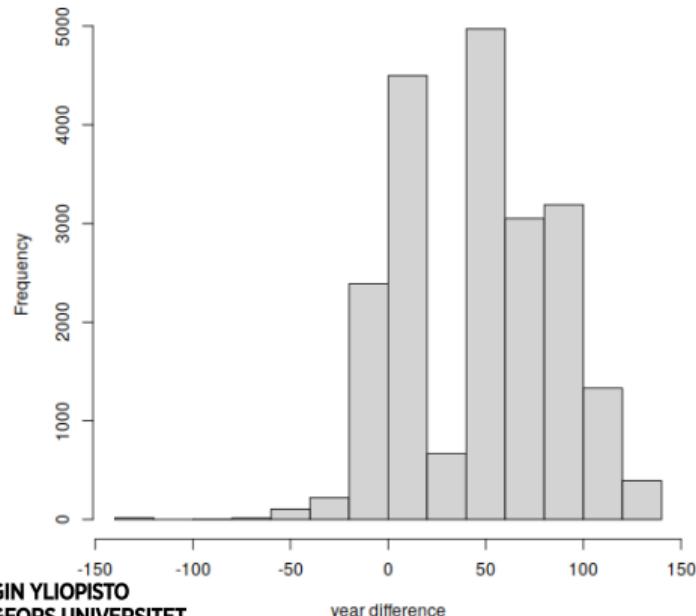
- a large collection of folk poems compiled by Elias Lönnrot
- 676 poems
- published in 1840
- close to oral tradition
- very well-known and influential



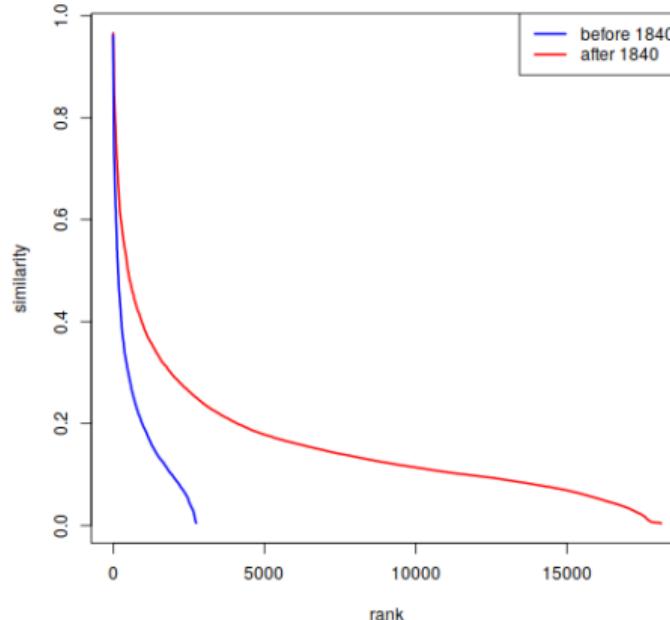


ORAL-LITERARY RELATIONSHIPS

Year differences between "Kanteletar" and similar poems

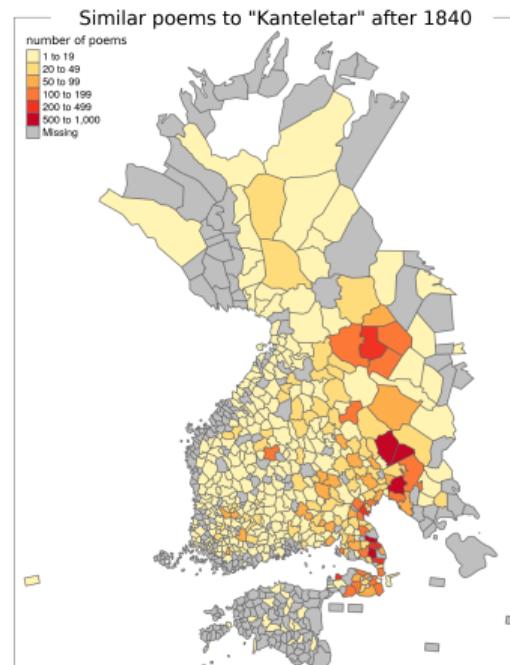
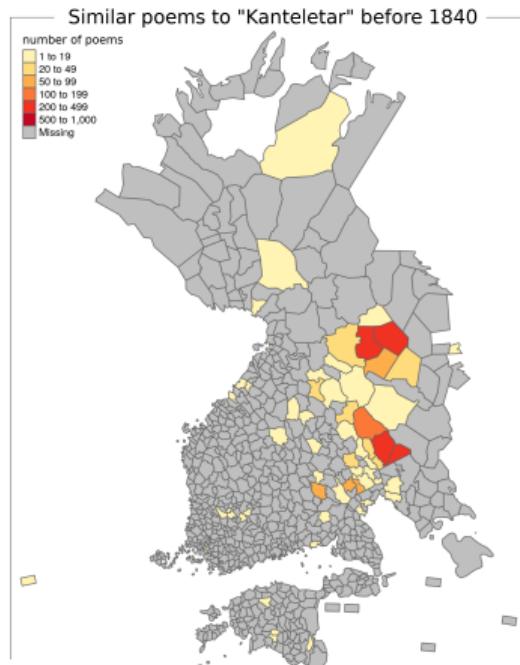


Distribution of similarity related to Kanteletar





ORAL-LITERARY RELATIONSHIPS





ORAL-LITERARY RELATIONSHIPS

Poems similar to “Kanteletar” by collector:

| before 1840 | | | 1840 and after | | |
|----------------------|------|-----------|---------------------|------|-----------|
| collector | n | examples | collector | n | examples |
| Lönnrot, Elias | 1798 | 1 2 3 4 5 | Europaeus, D. E. D. | 1008 | 1 2 3 4 5 |
| Gottlund, K. A. | 207 | 1 2 3 4 5 | Alava, Vihtori | 207 | 1 2 3 4 5 |
| Cajan, J. Fr. | 185 | 1 2 3 4 5 | Ahlqvist A. E. | 185 | 1 2 3 4 5 |
| Keckman, K. N. | 128 | 1 2 3 4 5 | Lönnrot, Elias | 128 | 1 2 3 4 5 |
| Ganander, Christfrid | 89 | 1 2 3 4 5 | Krohn, Kaarle | 89 | 1 2 3 4 5 |
| ... | | | ... | | |



ORAL-LITERARY RELATIONSHIPS

Poems similar to “Kanteletar” by singer:

| before 1840 | | |
|---------------------|----|-----------|
| collector | n | examples |
| Arhippa Perttuńe | 78 | 1 2 3 4 5 |
| Juhana Kainulainen | 34 | 1 2 3 4 5 |
| Martiska Karjalaińi | 13 | 1 2 3 4 5 |
| Matleena Kuivalatar | 11 | 1 2 3 4 5 |
| Mateli Kuivalatar | 9 | 1 2 3 4 5 |
| ... | | |

| 1840 and after | | |
|---------------------|-----|-----------|
| collector | n | examples |
| Larin Paraske | 401 | 1 2 3 4 5 |
| Liisa Itko | 76 | 1 2 3 4 5 |
| Arhippaińi Miihkali | 46 | 1 2 3 4 5 |
| Vyötär Tapanainen | 45 | 1 2 3 4 5 |
| Mari Laakko | 35 | 1 2 3 4 5 |
| ... | | |



ENHANCING THE TYPE INDICES

Proposing type annotations for unannotated poems:

| SKVR | | | ERAB | | |
|-------------------------|-----|-----------|-----------------------|-------|-----------|
| type | n | examples | type | n | examples |
| The maid to be ransomed | 19 | 1 2 3 4 5 | The Martinmas song | 741 | 1 2 3 4 5 |
| Words of fire | 16 | 1 2 3 4 5 | Marry me! | 630 | 1 2 3 4 5 |
| Jacob de la Gardie | 15 | 1 2 3 4 5 | Harvesting song | 369 | 1 2 3 4 5 |
| Words of dry weather | 12 | 1 2 3 4 5 | Wirble-warble! | 300 | 1 2 3 4 5 |
| Golden maiden | 12 | 1 2 3 4 5 | Song of communal work | 314 | 1 2 3 4 5 |
| ... | | | ... | | |
| total | 711 | | total | 32806 | |



ENHANCING THE TYPE INDICES

(Nearly) identical poems with different type annotations?

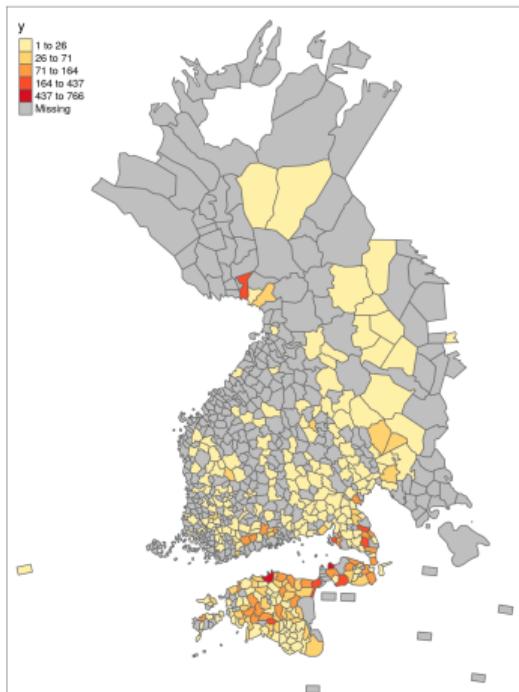
| Type 1 | Type 2 | n | examples |
|---------------------------|----------------------------------|----|-----------|
| Shoo, pig, to hunt foxes | I go to the forest to hunt foxes | 70 | 1 2 3 4 5 |
| The Lord will wake you up | Sleep, sleep my grassbird | 57 | 1 2 3 4 5 |
| Growing of hair | Wolves and sheep | 36 | 1 2 3 4 5 |
| Sing as long as you live | Sing as long as you live! | 34 | 1 2 3 4 5 |
| The sad one's home | Maiden's town | 33 | 1 2 3 4 5 |
| ... | | | |



SIMILARITIES ACROSS COLLECTIONS

(table excludes ERAB collections from Ingria)

| SKVR type | ERAB type | n | examples |
|------------------------------|---------------------------|-----|-----------|
| Shoo, goat, to the cattle | Shoo, goat, to the cattle | 780 | 1 2 3 4 5 |
| The maid to be ransomed | The maid to be ransomed | 412 | 1 2 3 4 5 |
| I served the farmer my years | The taskmaster's farmhand | 94 | 1 2 3 4 5 |
| The thief takes the jewelry | A horse stolen | 29 | 1 2 3 4 5 |
| Singer's reward | Singer's reward | 42 | 1 2 3 4 5 |
| ... | | | |





CONCLUSIONS

- *Relational* view on the collections
- Large-scale patterns have non-trivial justifications
 - different stability for different kinds of texts
 - collection density varies depending on time and place
 - influence of particular collectors / singers / works
- Desirable: a methodology **between** close and distant reading
- Proposed workflow:
 - ① select poem pairs according to some criteria based on metadata,
 - ② browse through aligned views of representative examples.
- Generalize? (characteristics: size of the collections, amount of similarity)