



Aalto University  
School of Science



SUOMALAISEN  
KIRJALLISUUDEN  
SEURA



**HELDIG**  
Helsinki Centre for Digital Humanities

HELSINKI INSTITUTE FOR  
SOCIAL SCIENCES  
AND HUMANITIES

# Aineistojen muuntaminen ja jalostaminen linkitetyksi avoimeksi dataksi: 1,2 miljoona kirjettä ja 100 000 toimijaa

**Suomalaisen Kirjallisuuden Seura**  
**27.5.2025**

Petri Leskinen, Aalto-yliopisto  
Senka Drobac, Helsingin yliopisto

# Sisältö

- Lähde materiaali
- Datamuunno
- Datan rikastaminen
- Tietomalli



Aalto University  
School of Science



SUOMALAISEN  
KIRJALLISUUDEN  
SEURA



**HELDIG**  
Helsinki Centre for Digital Humanities

HELSINKI INSTITUTE FOR  
SOCIAL SCIENCES  
AND HUMANITIES

# Lähdemateriaali

# Lähdemateriaali

- **Lukuisilta eri osapuolilta**
  - Datan määrä vaihtelee, 25 - 360.000 kirjettä
  - Ajallinen rajaus vuoteen 1917
- **Useissa eri formaateissa**
  - Word, Excel, CSV, JSON, RDF
- **Tiedon kattavuus vaihtelee**
  - Aika, kirjeiden määrä, lähetys- ja vastaanottopaikat
  - Henkilöiden biografiset yksityiskohdat



Aalto University  
School of Science



SUOMALAISEN  
KIRJALLISUUDEN  
SEURA

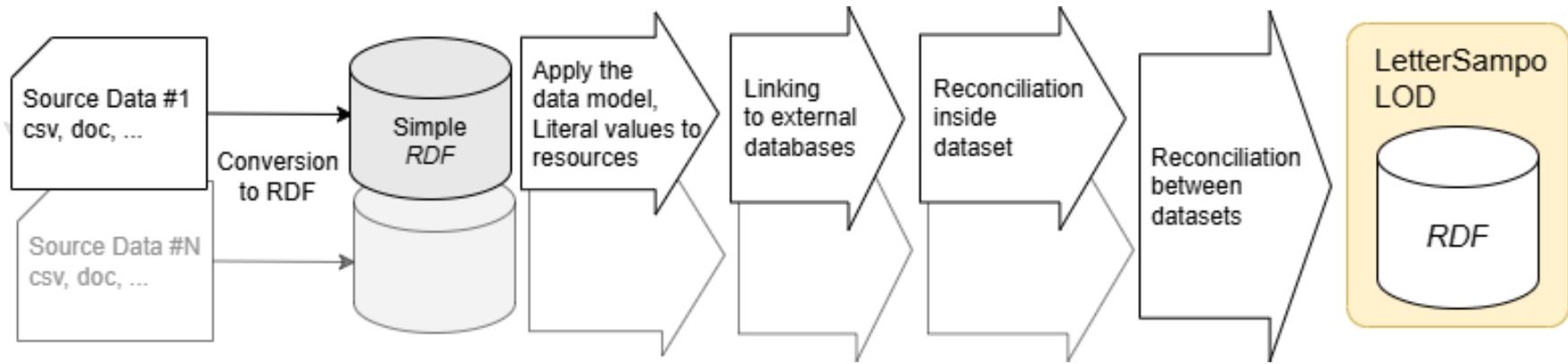


**HELDIG**  
Helsinki Centre for Digital Humanities

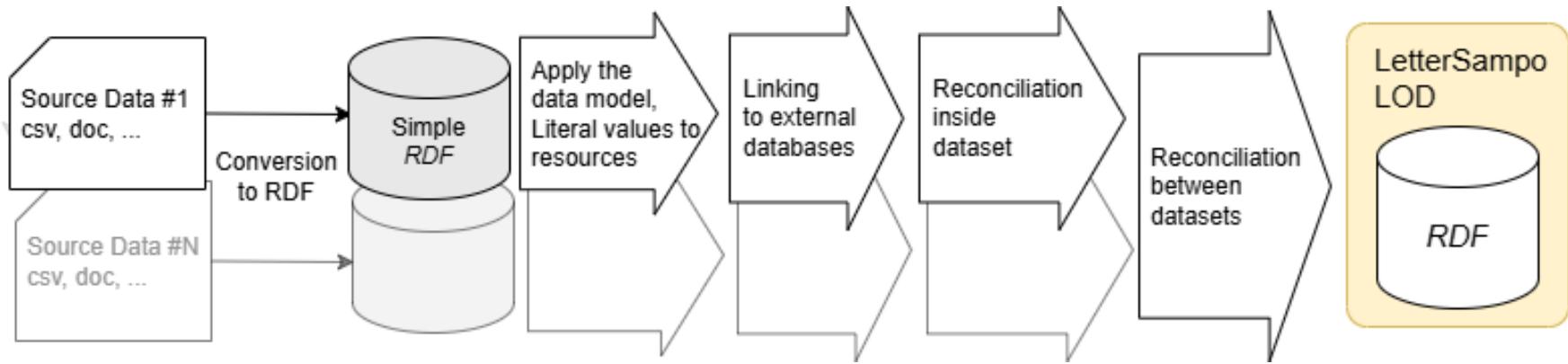
HELSINKI INSTITUTE FOR  
SOCIAL SCIENCES  
AND HUMANITIES

# Datamuunno

# Process pipeline

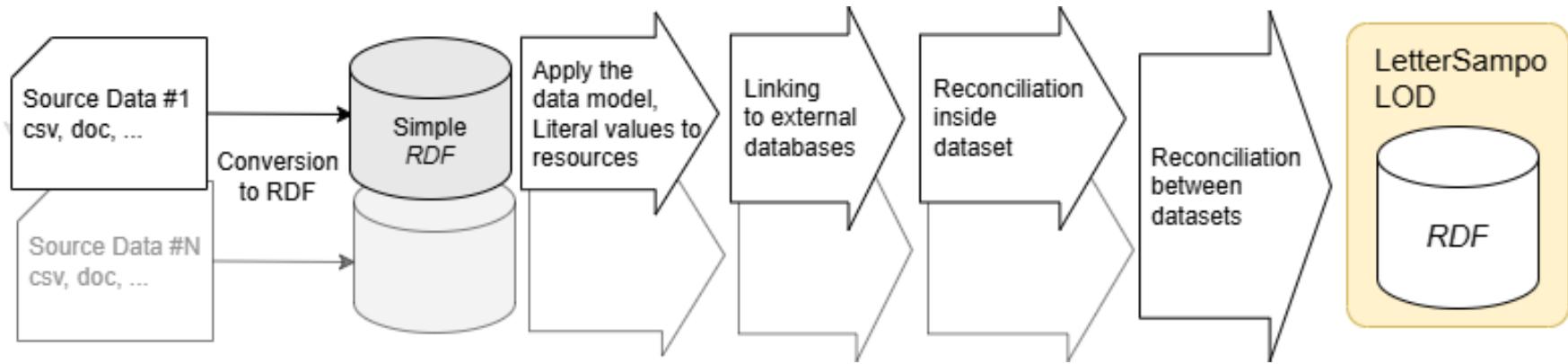


# Process pipeline



- **Data transformation**
- Harmonization
- Deduplication
- Data Linkage

# Process pipeline



- **Data transformation**
  - **From various data formats:**  
**Word, Excel, CSV, JSON, RDF**
- Harmonization
- Deduplication
- Data Linkage

# Word files

- **Different document structure**
  - Categories not in the same order, different titles
  - Data for several people in one document
- **Inconsistent formatting**
  - Mixing tabs/spaces
  - Text in multi lines
  - Other correspondences
  - Numerous small inconsistencies

# Multi lines

Number ↓ of letters

Suomen Eläinsuojeluyhdistys	1929	1
Suomen Evankelinen Seura	1918	1
Suomen ev.lut. Pyhäkouluyhdistys	1922	1
Suomen Historiallinen Seura	1912-1925	13
<b>Suomen kenraalikuvernööri (:n toimisto), /</b>	<b>1900-1916</b>	<b>30</b>
# liitteinä suomennoksia	<del>1900-1916</del>	<del>30</del>
Suomen kielen sanakirjaosakeyhtiö	1916	2
Suomen Kirkkohistoriallinen Seura	1902	1
Suomen Kirkollisviraston leski- ja orpokassa	1900-1919	2
Suomen Kirkon Pappisliitto	1918-1926	6
Suomen Kirkon Seurakuntatyön Keskusliitto (liite)	1918-1934	23
<b>Suomen Kirkon Sisälähetyssseura (liite) /</b>	<b>1917</b>	
<b>## (ks. myös Aarnisalo, Otto ja Suomen Piilia-seurat)</b>	<del>1917</del>	<del>1</del>
Suomen Kotiseutututkimuksen keskusvaliokunta	1917	1

# Multi lines

Detschy, Serafine	s.a.	1
Dillon, E.J.D.	1891, s.a.	6
liitt. Dillon, E.J.D. → Mr Dillon	s.a., 1 kpl	
Dingeldey, Ludvig	s.a.	3
Dmitrieff, Hélène	1906, s.a.	4
Donner, Otto (kts. myös Donner, Minette)	1885-1904, s.a.	11
Donner (o.s. Munck), Wilhelmina Sofia	1899	1
Ch. (Minette) & Donner, O.		
Dühren, C.J. von	1904	1

Donner (o.s. Munck) ; Wilhelmina Sofia /

Ch. (Minette) ; Donner, O.

# Gallen-Kallela Museum / FINNA

Results of a database query  
using Finna API Service

```
▼ events:
  ▼ valmistus:
    ▼ 0:
      type: "valmistus"
      date: "valmistusaika 02.08.1907, valmistusaika 1907"
      methods: []
      methodsExtended: []
      materials: []
      materialsExtended: []
    ▼ places:
      ▼ 0:
        placeName: "Saksa, Warnemünde"
        type: "URI"
        id: "http://www.yso.fi/onto/yso/p105087"
        ▼ ids:
          0: "http://www.yso.fi/onto/yso/p105087"
        ▼ details:
          0: "place_id_type_URI"
      ▼ actors:
        ▼ 0:
          name: "Gallen-Kallela, Akseli"
          role: "vastaanottaja"
        ▼ 1:
          name: "Öhquist, Johannes"
          role: "lähettiläjä"
```

# Albert Edelfelts brev API

```
{  
    "url": "http://edelfelt.sls.fi/api/letters/1/",  
    "id": 1,  
    "title": "Kiala den 16 Juni 67.",  
    "events": [ "http://edelfelt.sls.fi/api/events/1983/", ... ],  
    "web_url": "http://edelfelt.sls.fi/brev/1/kiala-den-16-juni-67/",  
    "date": "1867-06-16",  
    "urn": "URN:NBN:fi:sls-537-1403107521400",  
    "locations": [ "http://edelfelt.sls.fi/api/locations/273/" ]  
}  
...  
{  
    "url": "http://edelfelt.sls.fi/api/events/1995/",  
    "id": 1995,  
    "title": "Alexandra Edelfelt och Carl Albert Edelfelt har varit på resa; de har besökt Åbo och Stockholm.",  
    "letter": "http://edelfelt.sls.fi/api/letters/2/",  
    "persons": [ "http://edelfelt.sls.fi/api/persons/466/",  
                "http://edelfelt.sls.fi/api/persons/786/"  
    ],  
    "mentioned_locations": [ "http://edelfelt.sls.fi/api/locations/71/",  
                            "http://edelfelt.sls.fi/api/locations/78/"  
    ]  
}
```

# Simple RDF

The diagram illustrates the workflow for creating Simple RDF. It starts with 'Source Data' (CSV, DOC) which is converted to 'Simple RDF'. This is then processed through several steps: 'Apply the data model, Literal values to resources', 'Linking to external databases', 'Reconciliation inside dataset', and 'Reconciliation between datasets'. The final output is published to 'LetterSampo LOD' (RDF). Annotations in red highlight specific concepts:

- Actors**: Points to the 'Linking to external databases' and 'Reconciliation inside dataset' steps.
- Timespan**: Points to the 'Reconciliation between datasets' step.
- Letters**: Points to the 'Reconciliation between datasets' step.
- Contributor**: Points to the 'archival\_organization\_as\_recorded' value.
- Fonds**: Points to the 'fonds\_as\_recorded' value.
- Additional archival metadata**: Points to the 'dct:source' value.

```
records:NL_00009996 a :MetadataRecord ;  
    sender "Borg, Einar",  
           "Heikel, O. J.",  
           "Järventaus, H. A.",  
           "Kauppinen, K.",  
           "ym." ;  
    :sender_as_recorded "Borg, Einar ; Heikel, O. J. ; Järventaus, H.A. ; Kauppinen, K. ; ym." ;  
    :recipient "Juho Rudolf Koskimies (Forsman)" ;  
    :recipient_as_recorded "Juho Rudolf Koskimies (Forsman)" ;  
    :sending_date "1916" ;  
    :number_of_letters "5" ;  
    :archival_organization_as_recorded "National Library of Finland" ; Contributor  
    :fonds_as_recorded "Coll. 108 Koskimies, Juho R." ; Fonds  
    :source_filename "COLL108.docx" ;  
    :note "# KK" ;  
    :archival_location "COLL.108.27" ;  
    dct:source <http://ldf.fi/coco/source/nationallibrary> .
```



Aalto University  
School of Science



SUOMALAISEN  
KIRJALLISUUDEN  
SEURA

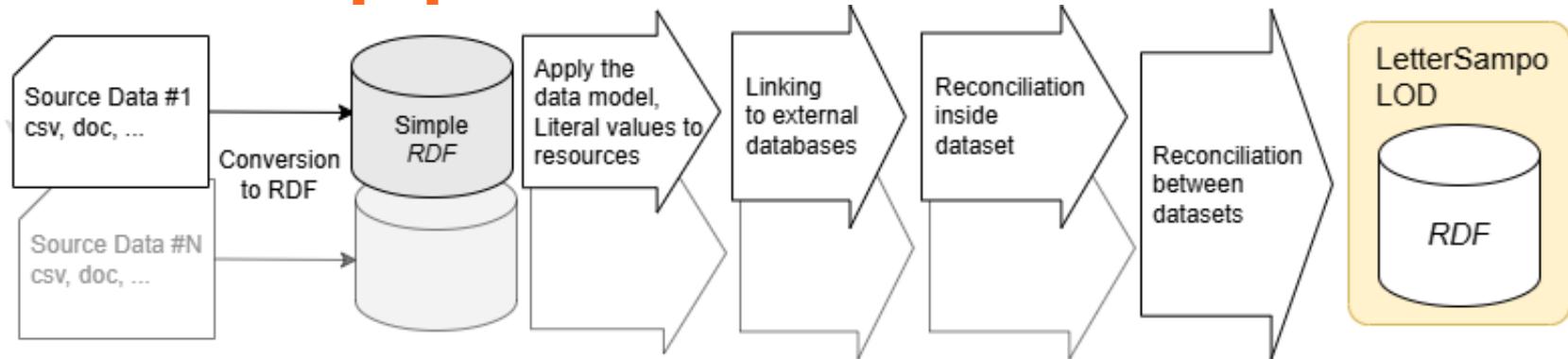


**HELDIG**  
Helsinki Centre for Digital Humanities

HELSINKI INSTITUTE FOR  
SOCIAL SCIENCES  
AND HUMANITIES

# Datan rikastaminen

# Process pipeline

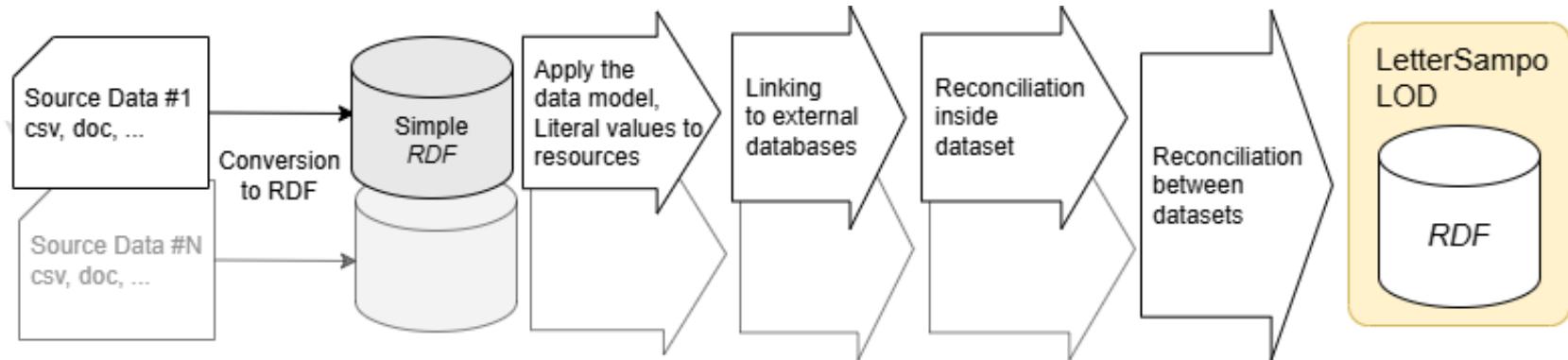


- Data transformation
- **Harmonization**
  - **Into one uniform data schema**
- Deduplication
  - Merge same actor mentioned in different datasets
- Data Linkage
  - To external data sources
  - Enrich with e.g. biographical details and family relations

# Esimerkkejä

- **Henkilö nimet**
  - "Haartman, Axel o. Hedvig f. Stolpe"
    - *Haartman, Axel*
    - *Haartman (Stople), Hedvig*
- **Elin- tai toiminta-aika**
  - Forsström, Petter (1877-1967)
  - Kekoni, Agnes (*fl. 1876-1944*)
  - Förlags Ab Söderström & Co (*fl. 1800-1949*)
- **Arkistokokoelmat**
  - "Hilja Haahden arkisto"
    - *Haahti, Hilja*
  - "SLSA 1319 Släkten Munsterhjelms arkiv"
    - *Munsterhjelm, Hjalmar; Munsterhjelm, Anders Lorentz; Munsterhjelm, Anders Gustaf; Munsterhjelm, Sofie Johanna; ... (linked to Finna)*

# Process pipeline

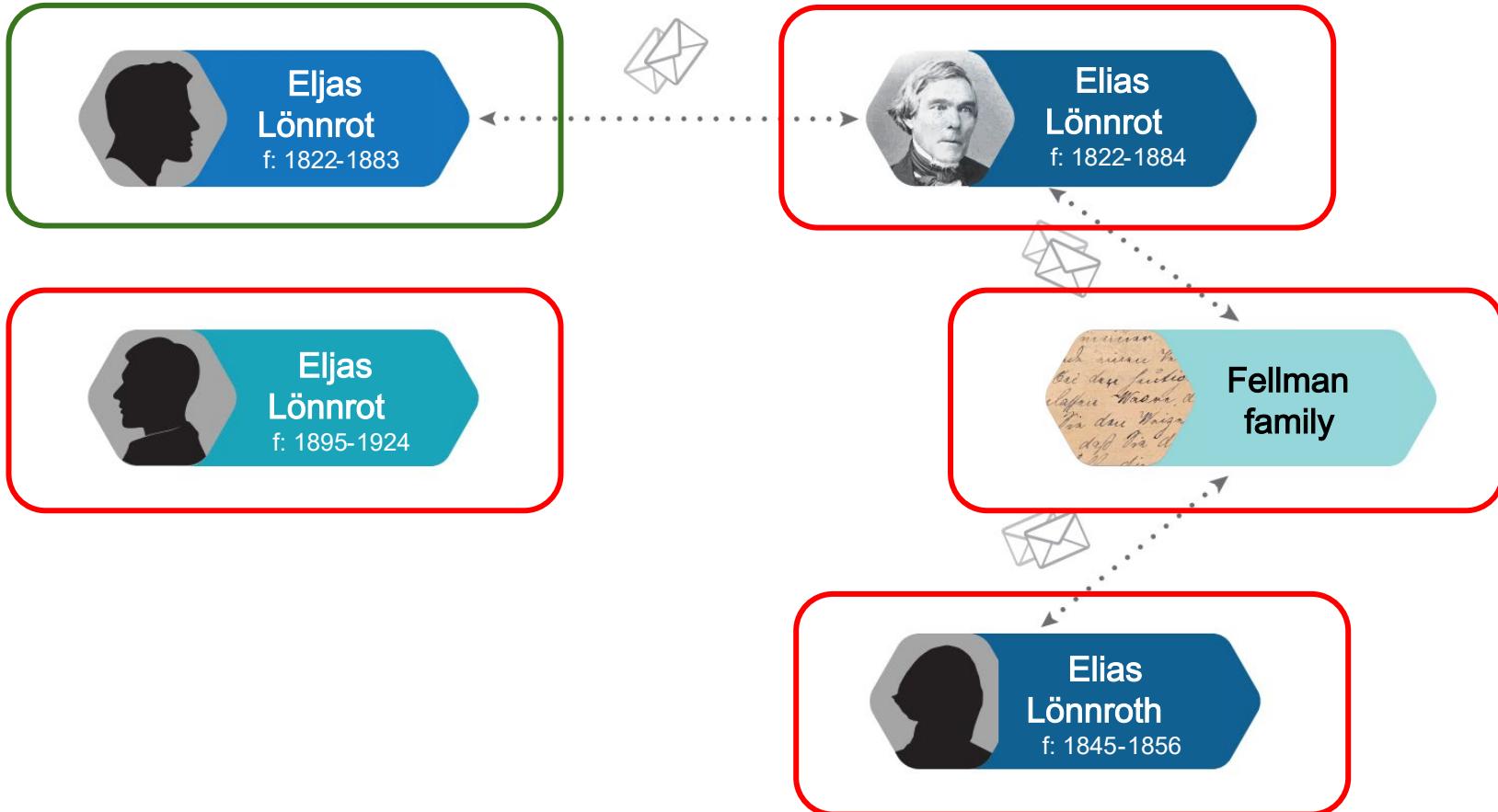


- Data transformation
- Harmonization
- **Deduplication**
  - **Merge same actor mentioned in different datasets**
- Data Linkage
  - To external data sources
  - Enrich with e.g. biographical details and family relations

# Toimija- ja paikkaontologiat

- **Henkilö, Organisaatio, Perhe/Suku, tuntematon**
  - Ulkopuoleiset tietolähteet
    - *SampoSampo*
    - *Biografiliset tiedot, perhe- ja ihmisiin suhteet*
- **Paikat**
  - Ulkopuoleiset tietolähteet
    - *Koordinaatit, hierrickia*

# Actor deduplication



# Actor deduplication

Eva Maria  
Acke



Eva Topelius



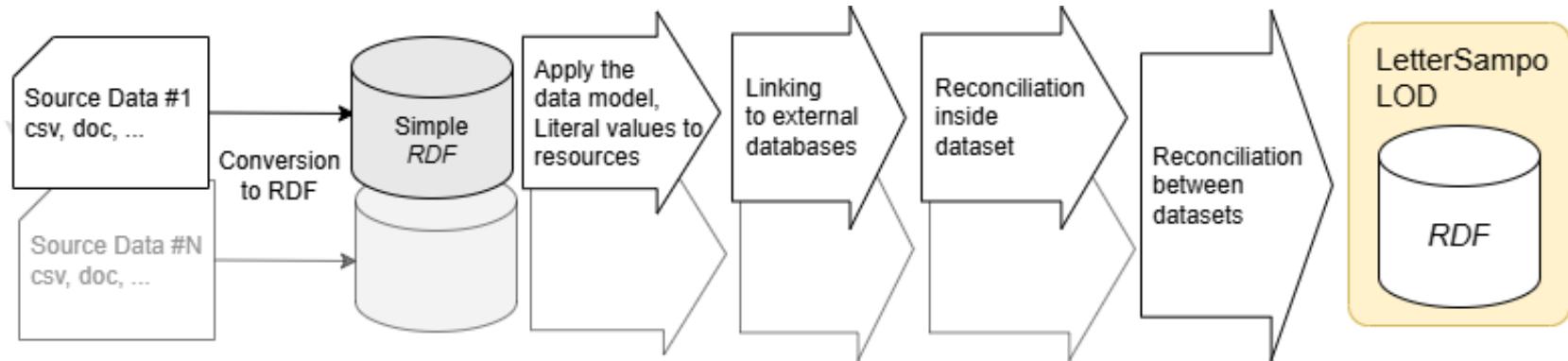
E. M.  
Topelius



**Alternative names:**

Eva Maria Acke, Eva Maria Topelius, Eva Topelius Acke, Eja Topelius, Eva Maria

# Process pipeline



- Data transformation
- Harmonization
- Deduplication
- **Data Linkage**
  - To external data sources
  - Enrich with e.g. biographical details and family relations

# Vocabulary of Finnish Actors

- [SampoSampo](#)

Historical Finnish People, Organizations, and Places

Biographical information about 230000 historical people related to Finland

Name variations, gender, years and places of birth and death ...

- Collected from multiple datasources:

1. AcademySampo (23000 entries)

2. BiographySampo (22500)

3. Wikidata (84000)

4. Kanto (29000)

5. BookSampo (9800)

6. ParliamentSampo (2100)

7. Norssi Highschool Alumni (2600)

8. Union List of Artist Names (1500)

9. Wikitree (12000)

10. Geneanet (144000)

*Parts of CoCo data providing biographical details*

11. Edelfelt Letters (3300)

12. Snellman Letters (3300)

13. Åbo Akademi University Library letters (1570)

# Vocabulary of Finnish Actors

## Eva Topelius (Q4455334)

Finnish painter  
Eva Maria Topelius | E. M. Topelius | E. Topelius | Eva Maria Acke

- In more languages  
Configure

Language	Label	Description	Also known as
English	Eva Topelius	Finnish painter	Eva Maria Topelius E. M. Topelius E. Topelius Eva Maria Acke
Finnish	Eva Topelius	suomalainen taiteilija	E. Topelius E. M. Topelius Eva Maria Topelius Eva Maria Andersson Eva Maria Acke Eva Topelius Acke Eva Topelius-Acke
Swedish	Eva Topelius	finlandssvensk målare	E. Topelius E. M. Topelius Eva Maria Topelius Eva Maria Acke



### Eva Maria Topelius

- Born in 1855 - Nykarleby
- Deceased March 23, 1929 - Vaxholm, Sverige, aged 74 years old

1 file available

## Parents

- Zacharias Topelius, *Verklig statsråd* 1818-1898
- Maria Emilie Lindqvist 1821-1885

## KÄYTETTÄVÄ NIMENMUOTO

## Topelius, Eva, 1855-1929 [edit](#)

### TYYPPI

henkilö

### HENKILÖN TUNNISTE

Asteri ID: 000122880 [edit](#)

### SYNTYMÄAIIKA

1855

### KUOLINAIIKA

1929

### SYNTYMÄPAIKKA

[Ususkaarlepyy \(yso-paikat\)](#)

### KUOLINPAIKKA

Vaxholm, Ruotsi

### HENKILÖN AMMATTI TAI TEHTÄVÄ

[taidemaalari](#) ([Metatietosanasto](#))

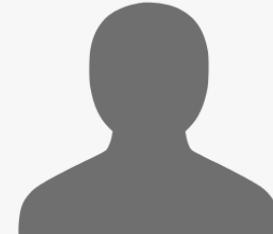
### HENKILÖN TOIMINTA-ALA

[maalaustaide](#) (YSO)

### HENKILÖN LIITTYVÄ HENKILÖ

[Topelius-Acke, Eva, 1855-1929](#)

## Eva Topelius-Andersson



\* [Ususkaarlepyy](#) 4.9.1855 † [Vaxholm](#) 1929

taidemaalari

Lähteet: Kansallisbiografia, Kansallisbiografian sisällysluettelo

# Challenges

- General issues of record linkage, e.g., namesakes, name variations
- Errors in source data
  - Typos
  - Uncertainty in letter dates
  - Erroneous dates
- No (machine readable) data available about some actors

# Record Linkage, Current Implementation

- Linkage to external databases
- Deduplication between datasets based on external matches
- Data Enrichment from external sources
  - Biographical details
  - Occupations
  - Personal relations
  - Images

# Record Linkage, Current Implementation

- Using Python package [SPLink](#):
  - Unsupervised learning
  - Term Frequency learning
  - Alternative Python packages: [Dedupe](#), [Record Linkage](#)
- Similarity based on:
  - Three different name variations considering
    - *Family Names*
    - *Given Names*
    - *Middle Names*
    - *Initials*
  - Quantiles of letter activity period
    - *Compared to lifespans of external actors*
  - Existing linkage
    - *Cases of Snellman and Edelfelt*



Aalto University  
School of Science



SUOMALAISEN  
KIRJALLISUUDEN  
SEURA

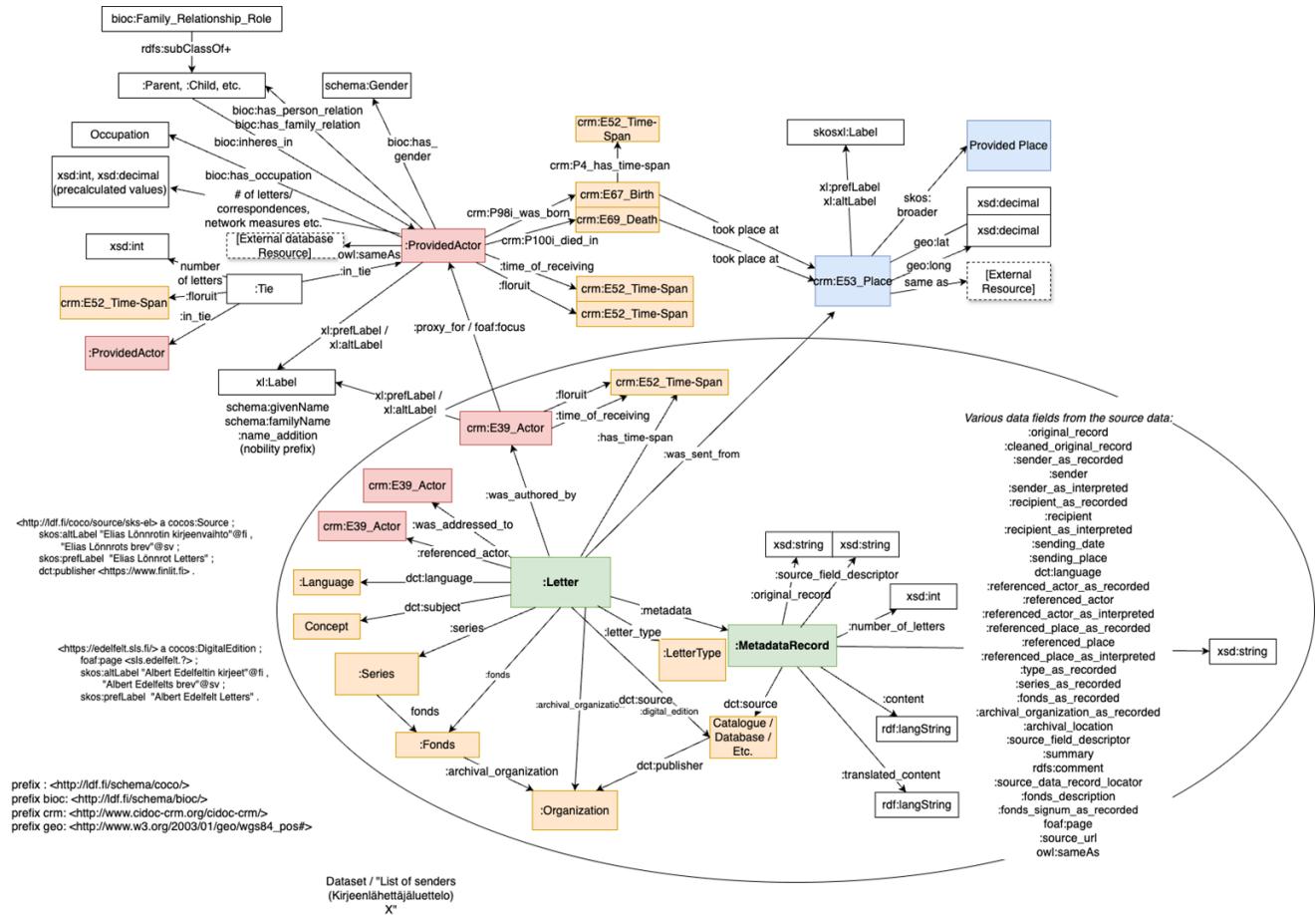


**HELDIG**  
Helsinki Centre for Digital Humanities

HELSINKI INSTITUTE FOR  
SOCIAL SCIENCES  
AND HUMANITIES

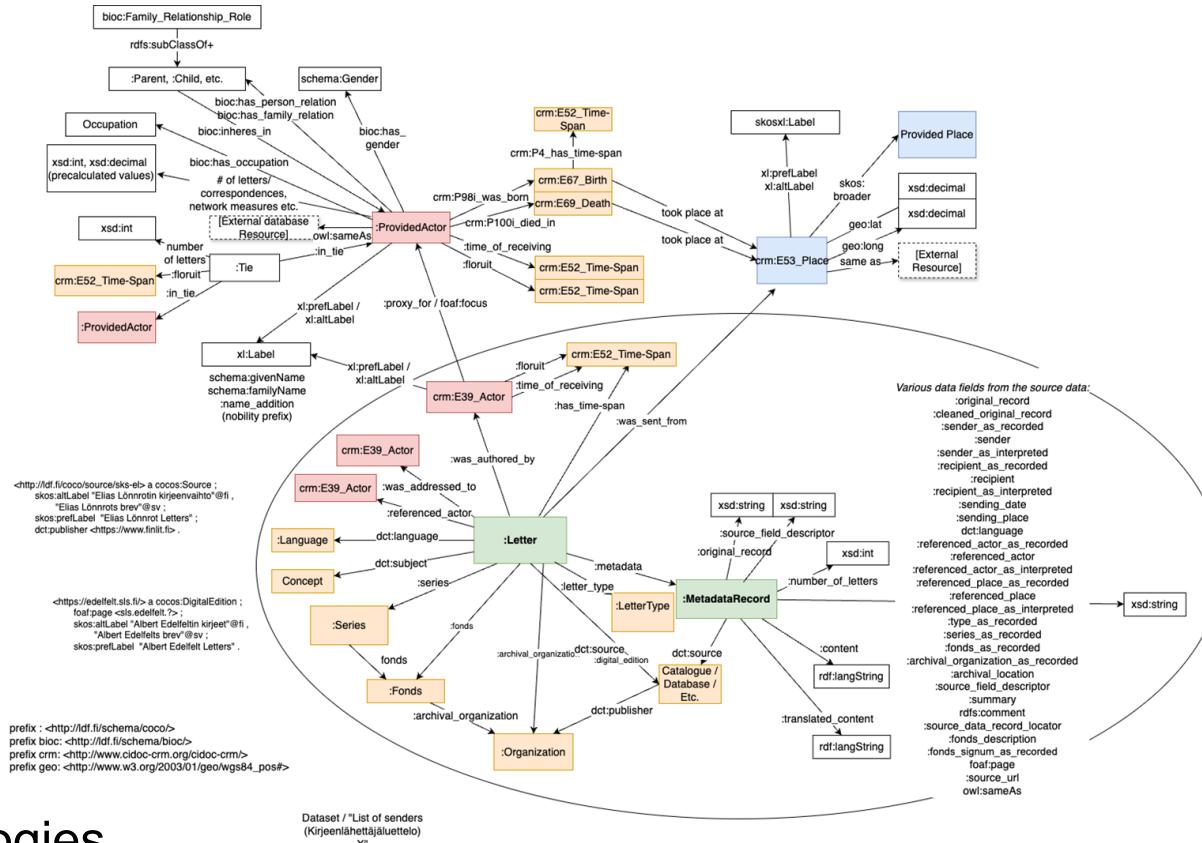
# Tietomalli

# Tietomalli



# Tietomalli

- Metadata Records
- Letters
- Actors
- Fonds
- Places
- Contributors
- Occupations
- Lettertypes
- Times
- ... supporting ontologies



# Tietomalli

```
letters:NL__00009996_0000001 a :Letter ;
    :estimated_year 1916 ;
    :fonds fonds:f2169107577047552292 ;
    :has_time-span <http://ldf.fi/coco/times/time_1916-01-01T00:00:00-1916-12-31T23:59:59-
1916-01-01T00:00-1916-12-31T23:59:59> ;
    :metadata records:NL__00009996 ;
    :original_data_provider sources:nationallibrary ;
    :type :missing_value ;
    :was_addressed_to actors:p2733526637996731453 ;
    :was_authored_by actors:p1641415298093300589,
actors:p1687294493541452016,
actors:p1886281086476616465,
actors:p4082113826364268026,
actors:p5971669413030914771 ;
    dct:source sources:nationallibrary ;
    skos:prefLabel "1916: Borg, Einar; Heikel, O. J.; Järventaus, H. A.; Kauppinen, K.; ym. ->
Juho Rudolf Koskimies (Forsman) [NL__00009996]" .
```



Aalto University  
School of Science



SUOMALAISEN  
KIRJALLISUUDEN  
SEURA



**HELDIG**  
Helsinki Centre for Digital Humanities

HELSINKI INSTITUTE FOR  
SOCIAL SCIENCES  
AND HUMANITIES

# Kiitos!