# Accessing the Republic.
# Entity extraction from the resolutions of the Dutch States-General

Marijn Koolen,[1,2] Esger Renkema,[1,2] Nienke Groskamp,[1] Frank Smit,[1]
Jirsi Reinders,[1] Ronald Sluijter,[1] Rik Hoekstra,[1,2] Jorid Oddens[1]

[1]Huygens Institute, Amsterdam, Netherlands
[2]DHLab - KNAW Humanities Cluster, Amsterdam Netherlands

DH Benelux 2024 - 6 June 2024 - Leuven, Belgium

# REPUBLIC

- REsolutions PUBLished In a Computational environment
    - Resolutions of the States General of the Dutch Republic (1576-1796)
    - Long serial publication (220 years), over 500,000 pages
    - ~60,000 daily meeting, ~1 million resolutions (propositions and decisions)

# Motivation

- Reasons to tag entities
  - Additional access points: alternative paths to navigate between documents
  - Contextual information: quick assessment of relevance/interest for user
- Reasons not to tag entities
  - If quality is low, it can annoy users, induce a lack of trust
  - Added value may not outweigh required effort
- Which entities?
  - Which entities occur?
  - Which are interesting?
  - Which are taggable?

# Entities in the Resolutions

- Tagging 8 types of entities
  - Person: person name including any attributions (title, job, legal status, ...)
  - Attribution: person attribution (title, job, legal status, ...) if refers to specific entity
  - Organisation: any organisation (incl. region name when it refers to governing body)
  - Committee: members of the States General tasked to investigate a matter
  - (Geo)Location - Political entity: name of a geolocation when it refers to the place
  - Date: absolute and relative dates (of submitted propositions and previous resolutions)
  - Resolution reference: references to specific earlier resolutions
  - Other: any remaining names

# Tagging Project

- Ground Truth
  - 1631 tagged resolutions of printed volumes 1705-1796
    - 370,560 tokens, 23,875 entities
  - 513 tagged paragraphs of handwritten volumes 1597-1702
    - 28,387 tokens, 2,347 entities
- Automated tagging
  - Train and evaluate taggers per entity type

# Nested Entities and Ground Truth

- Entities can be highly complex, with multiple levels of nesting
  - E.g. person name + attribution
  - Attribution can contain an organisation which can contain another organisation which can contain a location
- Examples
  - Henricus Gerhardus de Beveren Esveld, Predikant in de Gereformeerde Gemeente te Schoondyke onder het Classis van Walcheren, be roepen zynde tot Predikant in de Gemeente te Enkhuisen
  - Jan van Reusen — Solliciteur van den heer Thibaut heer van St. Aechtekercke, Burgermeester der Stadt middelburch cum socijs, taeckende de Dijckagie genaempt de Polder benoorden Aerdenburch

# Training NER Taggers

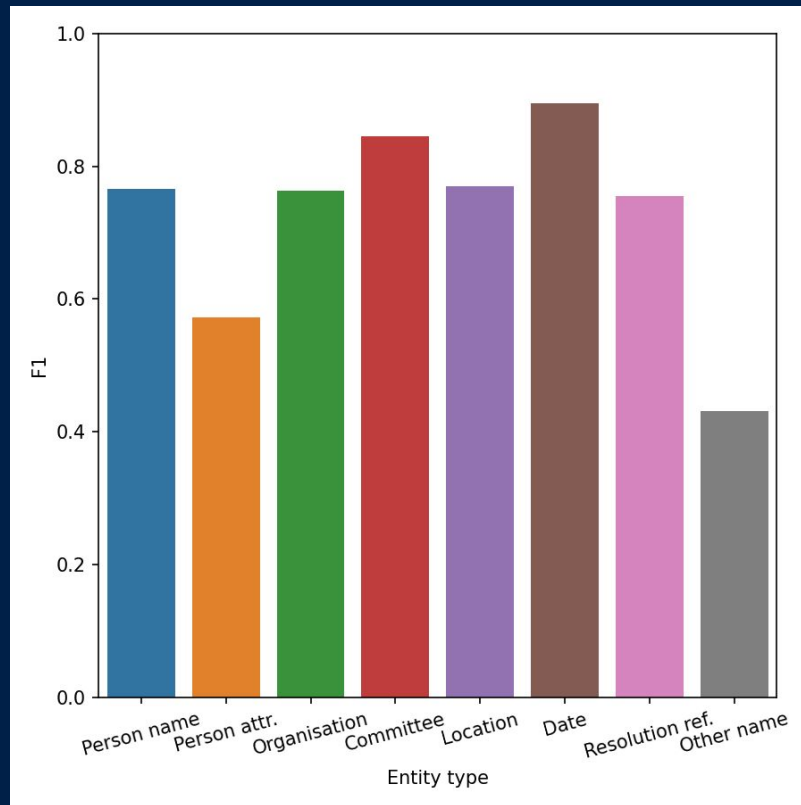- Combine/compare types of embeddings
  - Character-level embeddings, trained on resolution corpus (150 million words)
  - FastText embeddings, trained on resolution corpus
  - GysBERT (Manjavacas & Fonteyn 2022)
- Use Python Flair package (Akbik et al. 2019)
  - https://flairnlp.github.io/flair/
  - Agnostic to type of embeddings: Character-, word-, sentence-level
  - Combine via stacked embeddings!

# Combined Model or Model per Type?

- Single model
  - Advantages:
    - Need to train only one model, no tag conflicts in applying
  - Disadvantages:
    - Model choices may not be optimal for all entity types
    - Difficult to determine hierarchy of nested entities (ground truth is flattened)
- One per type
  - Advantages:
    - pick optimal model per type, allow for partial overlapping entities
    - Ground truth contains all information
  - Disadvantages:
    - Need decide how to deal with tag conflicts (partial overlap)

# Quantitative Evaluation - Best Model Per Type

# Evaluation

| Entity type | Tag repr. | Embeddings | Prec. | Recall | F1 | Support |
|---|---|---|---|---|---|---|
| Person | single type | GysBERT + Char. | 0.81 | 0.69 | 0.75 | 405 |
| Person attr. | single type | GysBERT + FastText + Char. | 0.57 | 0.56 | 0.56 | 573 |
| Organisation | single type | GysBERT + FastText + Char. | 0.82 | 0.71 | 0.76 | 283 |
| Committee | single type | Char. | 1.00 | 0.73 | 0.85 | 41 |
| Location | single type | GysBERT + FastText | 0.79 | 0.76 | 0.77 | 570 |
| Date | single type | Char. | 0.90 | 0.88 | 0.89 | 249 |
| Resolution ref. | all types | GysBERT + FastText | 0.82 | 0.70 | 0.75 | 57 |
| Other names | single type | GysBERT + FastText + Char. | 0.63 | 0.26 | 0.36 | 47 |

All best models use RNN instead of linear layer, CRF for prediction to capture dependencies in outputs

# Tagging All Resolutions

- 1.5 million paragraphs
- 8 million entities
- Mostly persons and attributions

| Entity type | # distinct | # total |
|---|---|---|
| Person | 1,159,672 | 1,929,235 |
| Person attribution | 1,176,039 | 1,743,086 |
| Organisation | 287,022 | 743,860 |
| Committee | 70,518 | 135,198 |
| Location | 617,542 | 2,551,180 |
| Date | 411,477 | 873,202 |
| Resolution reference | 28,990 | 189,865 |
| Other names | ~ | ~ |
| Total | 3,751,260 | 8,165,626 |

# Data Curation of tagged entities

- ◊ Identification of entities to recognise

- ◊ Resolution of entity descriptions to recognised entities

# Data Curation of tagged entities

◇ Identification of entities to recognise

thousands!

very many!

millions!

✦ Resolution of entity descriptions to recognised entities

Where and how
to apply manual input ?

# Sources of variation

+ multiple names / lack of formulaicity
+ political developments / evolution through time
+ inconsistent spelling / abbreviations
+ text recognition errors

# Sources of ambiguity

+ multiplicity of names / descriptions
+ intra-textual references
    „ the matter in question "
    „ aforementioned town "

# Sources of variation

- multiple names / lack of formulaicity
- political developments / evolution through time
- inconsistent spelling / abbreviations
- text recognition errors

} manual intervention

} fuzzy matching, automatic re-writing

# Sources of ambiguity

- multiplicity of names / descriptions
- intra-textual references
  „the matter in question"
  „aforementioned town"

} may be resolved by examining context

# Data curation:

## Traditional model

Sources of variation
- multiple names / lack of formulaicity
- political developments / evolution through time  } manual intervention
- inconsistent spelling / abbreviations
- text recognition errors  } fuzzy matching, automatic re-writing

Sources of ambiguity
- multiplicity of names / descriptions
- intra-textual references  } may be resolved by examining context
  "the matter in question"
  "aforementioned town"

text text text text text → [ NER tagger ]

↓

Entity descriptions → [ Data cleanup ] → Entity candidates → [ Manual curation ] → Named Entities

# Data curation:

## Traditional model

✦ AI ✦ magic ✦

text text text text text $\longrightarrow$ [ NER tagger ]

$\downarrow$

Entity descriptions $\longrightarrow$ algorithms, "regex", snakes? [ Data cleanup ] $\longrightarrow$ Entity candidates $\longrightarrow$ ...endless spreadsheets... [ Manual curation ] $\longrightarrow$ Named Entities
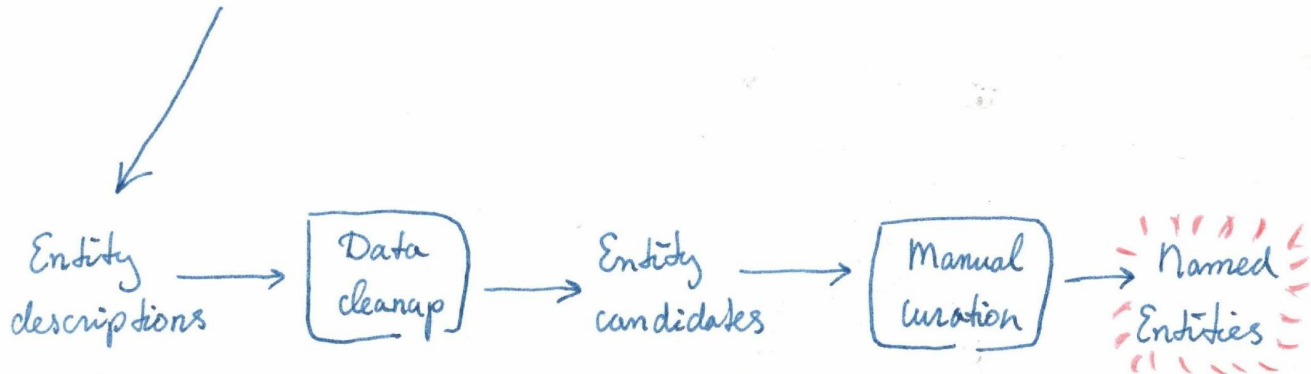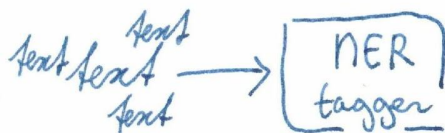
---

## Sources of variation

+ multiple names / lack of formulaicity
+ political developments / evolution through time  } manual intervention
+ inconsistent spelling / abbreviations
+ text recognition errors  } fuzzy matching, automatic re-writing

## Sources of ambiguity

+ multiplicity of names / descriptions
+ intra-textual references  } may be resolved by examining context

   "the matter in question"
   "aforementioned town"

# Data curation:

## Traditional model

text text text text text → [ NER tagger ] ○—→ ground truth prepared by volunteers

*✦ AI magic ✦*

Entity descriptions → [ Data cleanup ] → Entity candidates → [ Manual curation ] → *Named Entities*

algorithms, "regex", snakes?

...endless spreadsheets...

plain old-fashioned computer programming

domain-specific knowledge

---

## Sources of variation

- multiple names / lack of formulaicity
- political developments / evolution through time

  } manual intervention

- inconsistent spelling / abbreviations
- text recognition errors

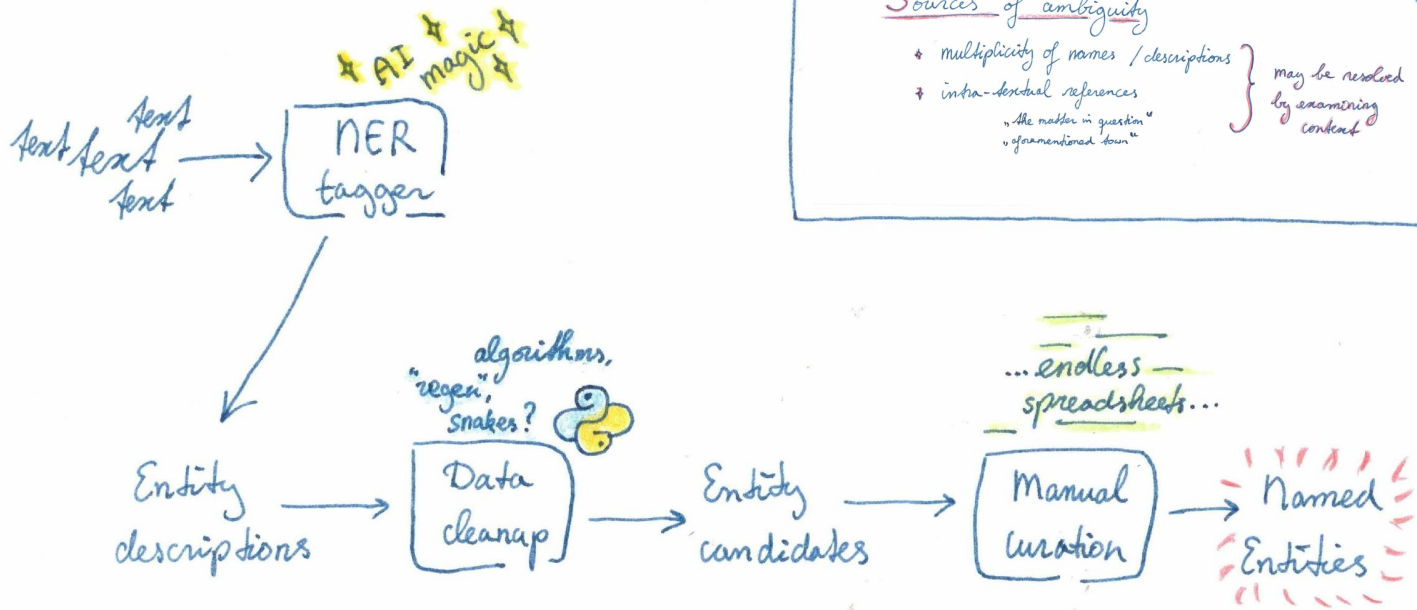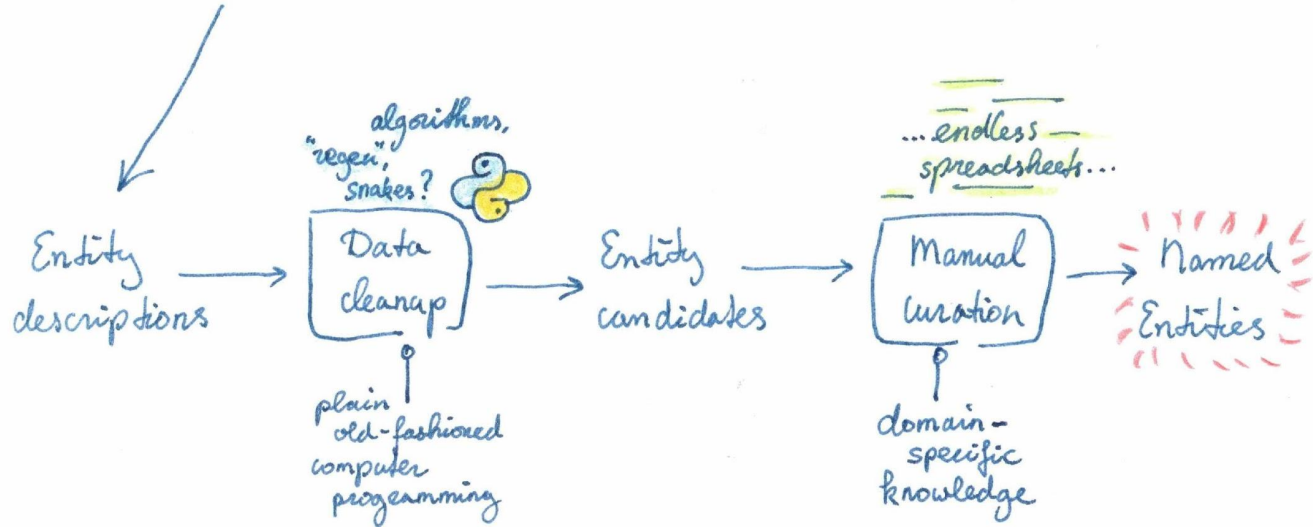  } fuzzy matching, automatic re-writing

## Sources of ambiguity

- multiplicity of names / descriptions
- intra-textual references

  } may be resolved by examining content

  "the matter in question"
  "aforementioned town"

# Data curation:

## Traditional model

text text text text → [NER tagger] ○ → ground truth prepared by volunteers

✦ AI ✦ magic ✦

Entity descriptions → [Data cleanup] → Entity candidates → [Manual curation] → Named Entities

*must be content-agnostic*

*algorithms, "regex", snakes?*

*limited context dependence*

plain old-fashioned computer programming

*only practical on small data sets*

...endless spreadsheets...

*cannot be repeated*

*too much work*

domain-specific knowledge

*decisions not documented*

---

### Sources of variation

✦ multiple names / lack of formulaicity
✦ political developments / evolution through time
} manual intervention

✦ inconsistent spelling / abbreviations
✦ text recognition errors
} fuzzy matching, automatic re-writing

### Sources of ambiguity

✦ multiplicity of names / descriptions
✦ intra-textual references
} may be resolved by examining context

"the matter in question"
"aforementioned town"

# Data curation:

## Traditional model



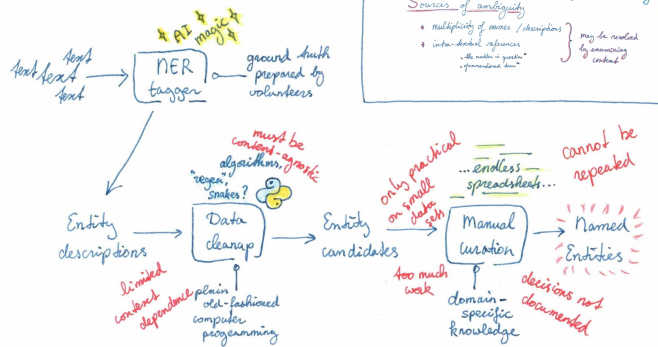text text text → NER tagger — ground truth prepared by volunteers

**Sources of variation**
- multiple names / lack of familiarity
- political developments / evolution through time
- inconsistent spelling / abbreviations } manual intervention
- text recognition errors } fuzzy matching, automatic re-writing

**Sources of ambiguity**
- multiplicity of names / descriptions } may be resolved
- intra-textual references } by examining
  - "the noble – greater" context
  - "aforementioned two"

Entity descriptions

+ AI magic +

"regex", "snakes"?

must be content-agnostic algorithms

Data cleanup

limited content dependence, plain old-fashioned computer programming

→ Entity candidates →

only practical on small data sets

...endless spreadsheets...

Manual curation

too much work

domain-specific knowledge

cannot be repeated

Named Entities

decisions not documented

# Data curation:

## Integrated model

text text text → NER tagger → Entity descriptions →

ground truth

# Data curation:

## Traditional model

AI magic

text text text → NER tagger —○ ground truth prepared by volunteers

↓

Entity descriptions →

**Sources of variation**
- multiple names / lack of formularity
- political developments / evolution through time } manual intervention
- inconsistent spelling / abbreviations
- bad recognition events } fuzzy matching, automatic re-editing

**Sources of ambiguity**
- multiplicity of names / descriptions } may be resolved by examining context
- inter-textual references
  - "the author in question"
  - "aforementioned town"

must be content-agnostic
algorithms,
"regex", snakes?

Data cleanup → Entity candidates →

...endless spreadsheets...

only practical on small data sets

too much work

Manual curation →

cannot be repeated

Named Entities

decisions not documented

limited content dependent
plain old-fashioned computer programming

domain-specific knowledge

# Data curation:

## Integrated model

logical criteria ← Visual examination

↑                    ↑

text text text text text → NER tagger → Entity descriptions → automated sorting → description groups

○

ground truth

# Data curation:

## Traditional model

Sources of variation
- multiple names / lack of formatuity
- political development / evolution through time  } manual intervention
- inconsistent spelling / abbreviations
- text recognition errors  } fuzzy matching, automatic re-writing

Sources of ambiguity
- multiplicity of names / descriptions  } may be resolved by examining context
- inter-textual references
  - "the matter in question"
  - "aforementioned item"

text text text → [ NER tagger ] — ground truth prepared by volunteers

Entity descriptions →

*AI magic*

must be content-agnostic

algorithms, "regex", "snakes"?

[ Data cleanup ] → Entity candidates →

limited context dependent plain old-fashioned computer programming

...endless spreadsheets...

only practical on small data sets

too much work

[ Manual curation ] → Named Entities

domain-specific knowledge

cannot be repeated

decisions not documented

---

# Data curation:

## Integrated model

logical criteria ← Visual examination

text text text text → [ NER tagger ] → Entity descriptions → [ automated sorting ] → description groups

ground truth

| cut-off

Named Entities

# Data curation:

## Traditional model

*AI magic*

text text text → [NER tagger] —○ ground truth prepared by volunteers

Entity descriptions

must be content-agnostic
algorithms
"regex", snakes?
limited content dependent
plain old-fashioned computer programming

→ [Data cleanup] → Entity candidates

only practical on small data sets
too much work
domain-specific knowledge

...endless spreadsheets...

→ [Manual curation] → *Named Entities*

cannot be repeated
decisions not documented

### Sources of variation
- multiple names / lack of formularity
- political development / evolution through time  } manual intervention
- inconsistent spelling / abbreviations
- text recognition errors  } fuzzy matching, automatic re-visiting

### Sources of ambiguity
- multiplicity of names / descriptions
- intra-textual references  } may be resolved by examining content
  "the matter in question"
  "aforementioned item"

# Data curation:

## Integrated model

forms a separate data set!

logical criteria

development cycle

Visual examination

no data entry!

### automated path

text text text text → [NER tagger] → Entity descriptions → [automated sorting] → description groups

ground truth

content-sensitive!
self-documenting!
moderately fast!

cut-off

*Named Entities*

## Example: Admiralities of the Republic

Once sorted in the bin „Admiralty",
    divide again on the following keywords:

Dokkum

Harlingen

Vrieslandt

Rotterdam

Op de Maze

Amst [a-z] *

# Example: Admiralities of the Republic

Once sorted in the bin „Admiralty",
divide again on the following keywords:

Dokkum
Harlingen  } Admiralty of Frisia
Vrieslandt

Rotterdam  } Admiralty of Rotterdam
Op de Mase

Amst[a-z]*  — Admiralty of Amsterdam

⤷ Narrow semantic domain enables
improvements over OCR/HTR

Possible entities
are known
in advance.

# Example : sorting local governments

Borgermren ende Regierders der Stadt Amsterdam

Syndicus ende Raedt ~~der~~ van Geneve

hooch Baillu ende Schepenen van Gent

Regenten van Helmont, Quartiere van Peellandt,
Meyerye van 's Hertogenbosch

# Example : sorting local governments

Borgermren ende Regiersders der Stadt Amsterdam

Syndicus ende Raedt ~~tan~~ van Geneve

hooch Baillu ende Schepenen van Gent

Regenten van Helmont, Quartiere van Peellandt,
Meyerye van 's Hertogenbosch

# Example:   sorting local governments

functions in local governments {

Borgermren ende Regierders der Stadt Amsterdam → Amsterdam

Syndicus ende Raedt ~~der~~ van Geneve → Genève

hooch Baillu ende Schepenen van Gent → Gent

Regenten van Helmont, Quartiere van Peellandt,
Meyerye van 's Hertogenbosch } → Helmond

Attributions ———→ Organisations ←——— Locations

# Conclusions

- Domain-specific NER tagging and training
  - Attributions are difficult
  - Ambiguity in instructions and in entities themselves
- Curation
  - Identification and resolution of entity references by successive grouping
  - (Logical) criteria for grouping form a separate dataset
  - Nesting of entity types and partial overlap can be powerful tools
- Analysis
  - Decomposing complex entities allows for combining multiple dimensions of analysis
- Project results will be published by December 2024!

# Thank You!

We also thank the volunteers for their contributions to this project!

Questions?