# Creating a Human-Annotated Health Record Dataset with Limited Resources

**Krištof Anetta**
**xanetta@fi.muni.cz**
Natural Language Processing Centre
Faculty of Informatics, Masaryk University

December 9, 2023

# Introduction

# Introduction

- Getting health records from hospitals is hard
  - Getting annotated ones is next to impossible
- We have a 42-million-word dataset of oncology health records and we are trying to get as much of it annotated
- A balanced set of ca. 50,000 words was chosen to be annotated by humans

# Visualization of what we want

1 Resekát  [AnatomicalSite_laterality] levá  [AnatomicalSite_name] mamma:  resekát  [AnatomicalSite_name] mléčné žlázy  o rozměrech 80-79-45 [Lab_name] [Lab_value] [Lab_unit] mm.

2 [AnatomicalSite_laterality] Ventrálně  je volně pohyblivá  [AnatomicalSite_name] fascie.

3 Na řezu je  [AnatomicalSite_name] mamma  prostoupena hrubou fibrózou [SignSymptom], v níž se nachází

ostře ohraničené suspektní ložisko [SignSymptom], obdobného vzhledu jako  fibrosa [SignSymptom], s patrnými

prokrvácenými punkčními kanály [SignSymptom], které má přibližné rozměry [Lab_name] 19-20 [Lab_value] mm [Lab_unit].

4 [AnatomicalSite_laterality] Mediální  okraj je cca  5 [Lab_value] mm [Lab_unit], [AnatomicalSite_laterality] ventrální  a  [AnatomicalSite_laterality] kraniální  12 [Lab_value]

mm [Lab_unit], [AnatomicalSite_laterality] kaudálně  navazuje hrubá fibróza [SignSymptom].

# Before annotation

## Preannotation

- Easy concepts can be found with rule-based methods
- It is much faster and more accurate when they are already present and annotators only verify them (but they must verify them)
- We preannotated:
  - Medication names
  - Medical abbreviations

# Medication: SÚKL databases

| | | | | | | |
|---|---|---|---|---|---|---|
| 8949 | | EZETROL | 10MG | TBL NOB | 100 | |
| 9709 | | SOLU-MEDROL | 40MG/ML | INJ PSO LC | 40MG+1M | |
| 9710 | | SOLU-MEDROL | 62,5MG/M | INJ PSO LC | 125MG+2M | |
| 9711 | | SOLU-MEDROL | 62,5MG/M | INJ PSO LC | 500MG+7, | |
| 9712 | | SOLU-MEDROL | 62,5MG/M | INJ PSO LC | 1000MG+1 | |
| 9844 | | TORECAN | 6,5MG | TBL OBD | 50 | |
| 10032 | | PIRACETAM AL | 800MG | TBL FLM | 60 | |
| 10033 | | PIRACETAM AL | 800MG | TBL FLM | 120 | |
| 10045 | | AGNUCASTON | | TBL FLM | 30 | |
| 10046 | | AGNUCASTON | | TBL FLM | 60 | |
| 10047 | | AGNUCASTON | | TBL FLM | 100 | |
| 10052 | | AGNUCASTON | | TBL FLM | 300 | |
| 10055 | | TABACUM | 31CH-200( | GRA | 4G | |
| 10063 | | BROMHEXIN KM | 8MG/ML | POR GTT S | 1X30ML | |
| 10073 | | ECHINACEA ANGUSTIFOLIA | 31CH-200( | GRA | 1X4G | |
| 10087 | | LOBELIA INFLATA | 31CH-200( | GRA | 1X4G | |
| 10111 | | DHC CONTINUS | 120MG | TBL MRL | 56 | |

# Abbreviations: Web resources

Seznam zkratek

Přehled používaných zkratek

**A**

| | |
|---|---|
| A., a. | arterie |
| AA | alergická anamnéza |
| AAA | aneurysma abdominální aorty |
| AAT | antikoantrotomie |
| Ab | abort (potrat) |
| AB | arteria brachiális |
| AB l. dx. | arteria brachialis vpravo |
| AB l. sin. | arteria brachialis vlevo |
| ABD, abd. | abdukce |
| ABF | aortobifemorální |
| ABI | index kotník - paže |
| ABR | acidobazická rovnováha |

| Používané zkratky ve zdravotnické dokumentaci EÚ | | | |
|---|---|---|---|
| **Lékařské názvy** | | | |
| št. | štíná žláza | bilat. | oboustranný |
| ko | kourola | unilat. | jednostranný |
| vyš. | vyšetření | dx. | vpravo, pravý |
| skut. | akutní | sin | vlevo, levý |
| chron. | chronický | L. | lumbální (bedemí) |
| subklin. | subklinický | lob dx., PL | pravý lalok štítné žlázy |
| klin. | klinický | lob sin., LL | levý lalok štítné žlázy |
| dekomp. | de kompenzovaný | | M-mammy |
| opak.recid. | opakované, opakovaný | | P- pubické ochlupení |
| inic. | iniciálně, na začátku | Tanner | G - genitál |
| palp. | palpačně, pohmatem | | A- axilární ochlupení |
| v.s. | pravděpodobně | PMV | psychomotorický vývoj |
| domin. | dominantní | FG sróre | skóre dle Ferrimana a Gallwayové |
| stac. | stacionární, nezměněný | PMR | psychomotorická retardace |
| lab. | laboratorní | | |

| Zkratka | Výklad zkratky |
|---|---|
| Ab | Protilátka (z angl. Antibody) |
| ACE | acetyl cholin esteráza |
| ACE | angiotensin konvertující enzym |
| ACEI | inhibitory acetylcholinesterázy |
| ACTH | Adrenokortikotropní hormon |
| AChR | Acetylcholinové receptory |
| AD | autozomálně dominantní |
| ADD | Porucha pozornosti (attention deficit disorder) |
| ADH | Antidiuretický hormon, vazopresin |

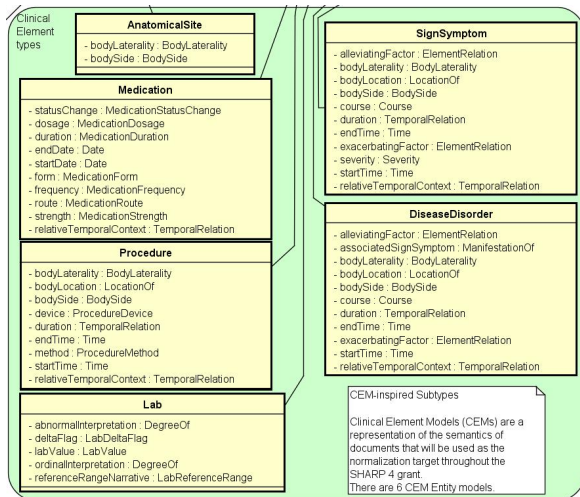`https://nlp.fi.muni.cz/projekty/ehr_analysis/zkratky/`

# Annotation

## Annotators

- 11 students
- Confidence considerations:
    - Since they are not experts, each record is annotated by at least 2 students
    - When they both annotate a string in the same way, confidence is high
    - But even when their views differ, they may both be right - we mark those as 50% confidence

# Annotation schema: Apache cTAKES types



The referential semantics schema used in Apache cTAKES [1]

## Annotation schema

| cTAKES type | Our schema |
|---|---:|
| AnatomicalSite | AnatomicalSite_laterality |
| | AnatomicalSite_name |
| Medication | Medication_dosage |
| | Medication_name |
| | Medication_strength |
| Procedure | Procedure |
| Lab | Lab_name |
| | Lab_unit |
| | Lab_value |
| SignSymptom | SignSymptom |
| DiseaseDisorder | DiseaseDisorder |
| | DateTime |
| | Abbreviation |
| | Negation |

# Annotation in BRAT

# Annotator options in BRAT

# Annotator's manual

**Entities to be annotated**

You can view a sample annotation ⇨here.

- AnatomicalSite
  - names of body parts and locations on the body
  - every AnatomicalSite annotation is either of these two:
    - **AnatomicalSite_name**: the name itself, e.g.

```
našla v pravém prsu bulku
```

    - **AnatomicalSite_laterality**: further specification of location, e.g. v

```
našla v pravém prsu bulku
```

- **DiseaseDisorder**
  - names of diseases and disorders, e.g.

```
léčena xareltem inf mononukleoza v 15 letech
```

- **SignSymptom**
  - medical occurrences which are not names of diseases and disorders but can indicate their presence or absence, e.g.

```
při bolestech svalů, teplotě
```

    If unsure whether it is an official *disorder name* or only a *symptom*, annotate as **SignSymptom**.
- **Procedure**
  - name of a procedure or process (diagnostic or therapeutic) carried out by medical personnel, e.g.

```
benefit adjuvantní chemoterapie minimální
```

# Annotation statistics

| Stage | Annotation count |
|---|---:|
| Initial state of health records | 0 |
| Rule-based preannotations | 4,266 |
| Preannotations handed to annotators | 9,368 |
| New annotations entered by annotators | 22,798 |
| Total number of human-verified or human-entered annotations | 32,166 |
| Total number of tokens with human-verified or human-entered annotation | 45,032 |

# Example of differences

Annotator 1 (record was at the beginning of their dataset)



Annotator 2 (record was in the middle)

# Preliminary NER training

# Preliminary NER training

- Stanford NER
- Cross-evaluation
  - 5 different subsets of data - train on 4, evaluate on 1 as gold standard, calculate average

# NER Results

| True \ Predicted | Abbreviation | AnatomicalSite_laterality | AnatomicalSite_name | DateTime | DiseaseDisorder | Lab_name | Lab_unit | Lab_value | Medication_dosage | Medication_name | Medication_strength | Negation | O | Procedure | SignSymptom |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Abbreviation | 1133 | 1 | 0 | 3 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 2 | 258 | 0 | 0 |
| AnatomicalSite_laterality | 0 | 76 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 105 | 0 | 0 |
| AnatomicalSite_name | 4 | 1 | 210 | 0 | 3 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 197 | 0 | 1 |
| DateTime | 2 | 0 | 0 | 596 | 0 | 0 | 0 | 0 | 2 | 0 | 0 | 0 | 219 | 0 | 0 |
| DiseaseDisorder | 29 | 0 | 9 | 0 | 49 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 132 | 1 | 8 |
| Lab_name | 8 | 0 | 3 | 0 | 0 | 36 | 0 | 0 | 0 | 0 | 0 | 0 | 77 | 0 | 1 |
| Lab_unit | 4 | 0 | 0 | 0 | 0 | 0 | 124 | 4 | 0 | 0 | 0 | 0 | 116 | 0 | 0 |
| Lab_value | 0 | 0 | 0 | 0 | 0 | 3 | 163 | 0 | 0 | 0 | 0 | 2 | 174 | 0 | 0 |
| Medication_dosage | 2 | 0 | 0 | 2 | 0 | 0 | 0 | 1 | 36 | 0 | 0 | 0 | 64 | 0 | 0 |
| Medication_name | 6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 84 | 0 | 0 | 54 | 0 | 0 |
| Medication_strength | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 57 | 0 | 36 | 0 | 0 |
| Negation | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 59 | 154 | 0 | 0 |
| O | 174 | 2 | 31 | 33 | 4 | 8 | 6 | 6 | 2 | 0 | 0 | 21 | 8861 | 6 | 44 |
| Procedure | 5 | 1 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 139 | 58 | 0 |
| SignSymptom | 0 | 0 | 4 | 0 | 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 129 | 0 | 15 |

True

Predicted

| Class | Precision | Recall | F1 score |
|---|---|---|---|
| Abbreviation | 0.885 | 0.707 | 0.781 |
| AnatomicalSite_laterality | 0.871 | 0.460 | 0.595 |
| AnatomicalSite_name | 0.871 | 0.426 | 0.568 |
| DateTime | 0.936 | 0.635 | 0.753 |
| DiseaseDisorder | 0.555 | 0.235 | 0.320 |
| Lab_name | 0.627 | 0.297 | 0.385 |
| Lab_unit | 0.778 | 0.389 | 0.511 |
| Lab_value | 0.832 | 0.394 | 0.524 |
| Medication_dosage | 0.637 | 0.286 | 0.390 |
| Medication_name | 0.959 | 0.554 | 0.701 |
| Medication_strength | 0.822 | 0.605 | 0.673 |
| Negation | 0.616 | 0.354 | 0.416 |
| O (no annotation) | 0.803 | 0.969 | 0.878 |
| Procedure | 0.767 | 0.319 | 0.439 |
| SignSymptom | 0.440 | 0.077 | 0.124 |
| **Weighted average** | **0.819** | **0.813** | **0.789** |

# Future directions

# Future directions

- 50,000 words is not enough for Transformer training
- Bootstrapping more data:
  - Iterative annotation - train smaller models, annotate larger data, evaluate errors, correct
  - Eventually, the whole 40,000,000 corpus could be annotated - lower quality, but sufficient size for LLMs

# Iterative bootstrapping

Thank you for your attention!