



**aLta**  
2023



# Part 2: Clinical NLP in Practice

Vlada Rozova



## Hands-on practice

Link to the dataset:



Link to the Jupyter Notebook:



## Building a named-entity recognition (NER) tool for clinical text data

- Where to start and what to look out for?
- What are the common approaches and what to expect?
- How do they compare on real-world data?

# Let's first take a look at the data

Link to the dataset:



## What is this dataset made up of?

- A .csv file with report metadata
- A folder with .txt files — reports themselves
- A folder with .ann files — annotations

Folder Navigation: <base>		
Name	Size	Modified
annotations	2.5 KB	2023-10-20
reports	48.6 KB	2023-11-02
LICENSE.txt	727 B	2023-01-26
SHA256SUMS.txt	9.1 KB	2023-09-30
annotation.conf		
chifir_metadata.csv		

Now time to head  
over to the Jupyter  
Notebook and  
explore

Link to the  
Jupyter  
Notebook:



## Outline

- Some background on the CHIFIR dataset
- Exploratory data analysis:
  - How many reports? What do they look like?
  - How to annotate data yourself?
    - Where to start? What to look out for?
- Three types of NERs:
  - How are they different? What works better on CHIFIR?
- Other tips & considerations