# Anonymisation and De-identification of Medical Documents

Averbis GmbH, David Hübner

# De-identification: Scope

- **De-identification**: Removal of protected health information (**PHI**)
- In the **EU**, we currently lack an actionable definition of PHI, so we resort to the Health Insurance Portability and Accountability Act (**HIPAA**)

## List of HIPAA PHI Identifier

averbis
text analytics

- Names
- Geographic Identifier
- Dates
- Phone numbers
- Fax numbers
- Email addresses

- Social Security Number
- Medical Record Numbers
- Health Insurance Beneficiary Number
- Account numbers
- Certificate Numbers
- Device Identifiers

- Vehicle Identifiers
- Website URLs
- IP Addresses
- Biometric Identifiers
- Facial Images
- Other Identifiers

# De-identification: Scenario

# De-identification: Technical approach

- Use natural language processing (**NLP**) to identify PHI in clinical documents
- NLP models are trained on manually labeled datasets (>175.000 annotations)
- They can identify PHI information based on
  - **Context** ("*Dear Mr …*")
  - **World knowledge** ("*Berlin*")
  - **Layout** (to some degree)
- Key technology: Usage of **Transformer**-based architectures, but fast enough to run on-premise (no LLMs)

# De-identification: Key Results

- **Recall > 98%:** We can find 98-99% of all PHI
- **Precision > 98%:** If we mark a Token as PHI, it is correct to an accuracy of over 98%
- A dedicated **Recall-optimized** pipeline achieves an average of **99.3% Recall**



Token-level micro-average evaluation

# De-identification: From annotations to de-identified documents

Once PHI was identified, we can modify the original documents.

Common strategies:

- **Redaction**: Replace PHI by **XXX** or by a tag, e.g. **<date>**
- **Substitution**: Substitute PHI with **synthetic information**
  - *Advantages*: Enhanced security - It is unclear whether any remaining information was originally PHI
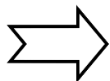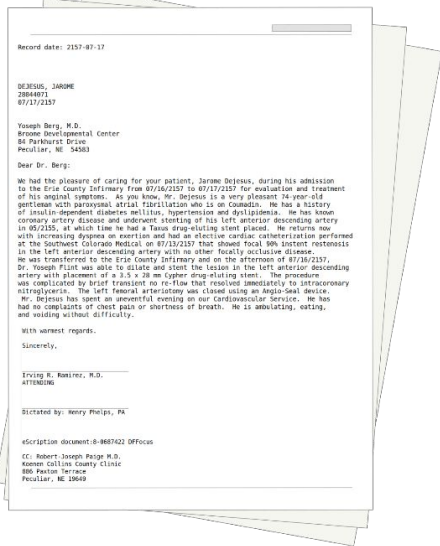  - *Challenge*: Consistent replacements are difficult

# De-identification = Anonymization? Not quite.

- Redaction/Substitution of PHI makes it much more difficult to identify the patients
- However, distinct medical diagnoses and contextual cues still may allow to identify the patient (see **k-anonymity**)
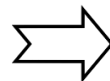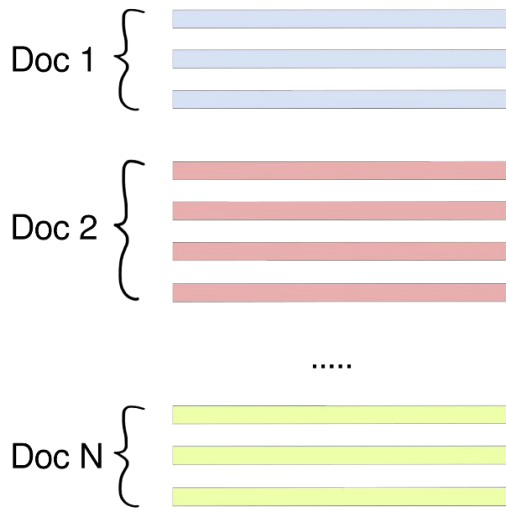- To overcome this, we developed the **Snippet Workflow**

# Snippet Workflow for retrieving anonymized texts

Goal: Workflow to retrieve texts that are **fully anonymized** but still valuable for ML model training
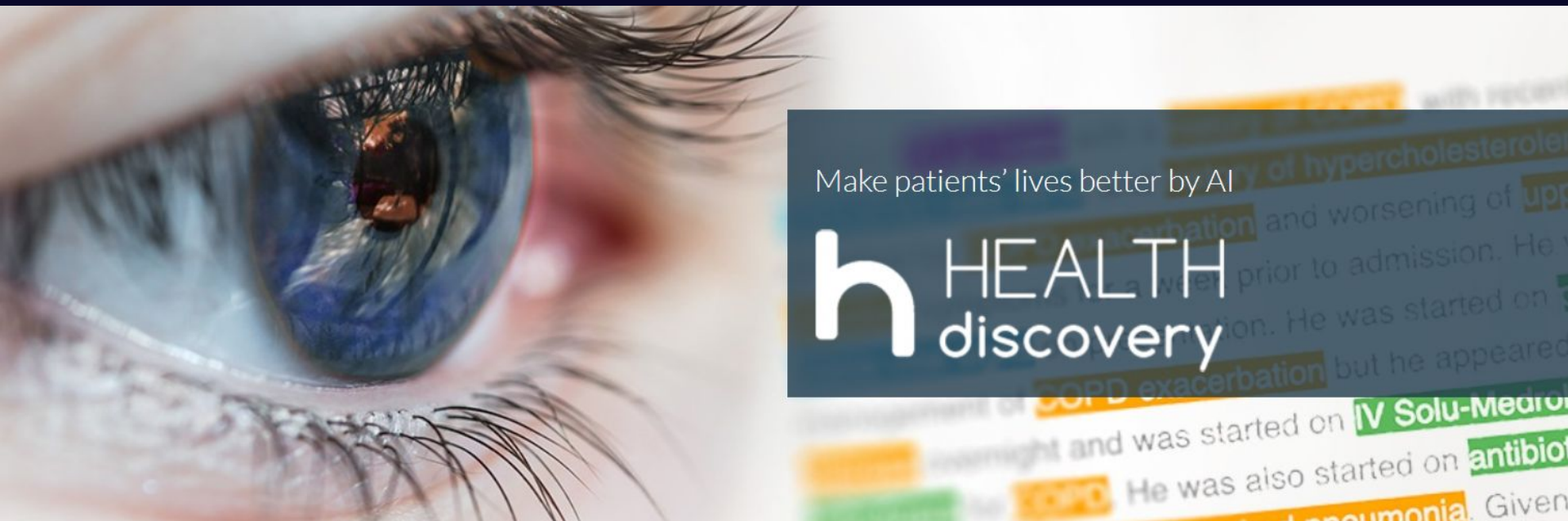


Clinical documents

Relevant de-identified snippets (e.g. sentences) per document

Shuffled snippets yielding fully anonymized text snippets

Doc 1

Doc 2

.....

Doc N

.....

averbis
text analytics

Make patients' lives better by AI

h HEALTH discovery

David Hübner
Lead ML Engineer at Averbis GmbH
Email: david.huebner@averbis.com