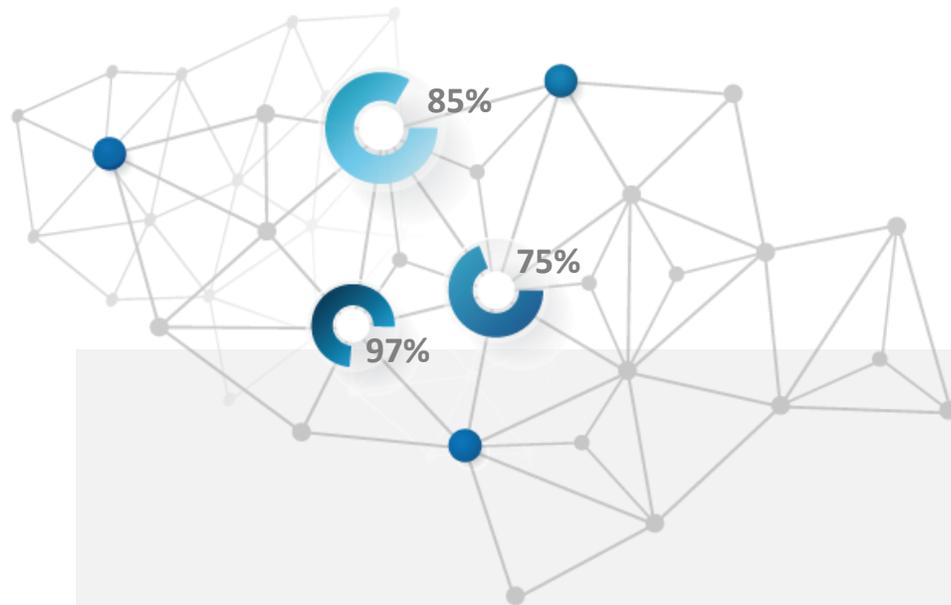


Anonymisierung und De-Identifizierung von medizinischen **Dokumenten** mit dem Averbis DeID-Tool

Stefan Schulz

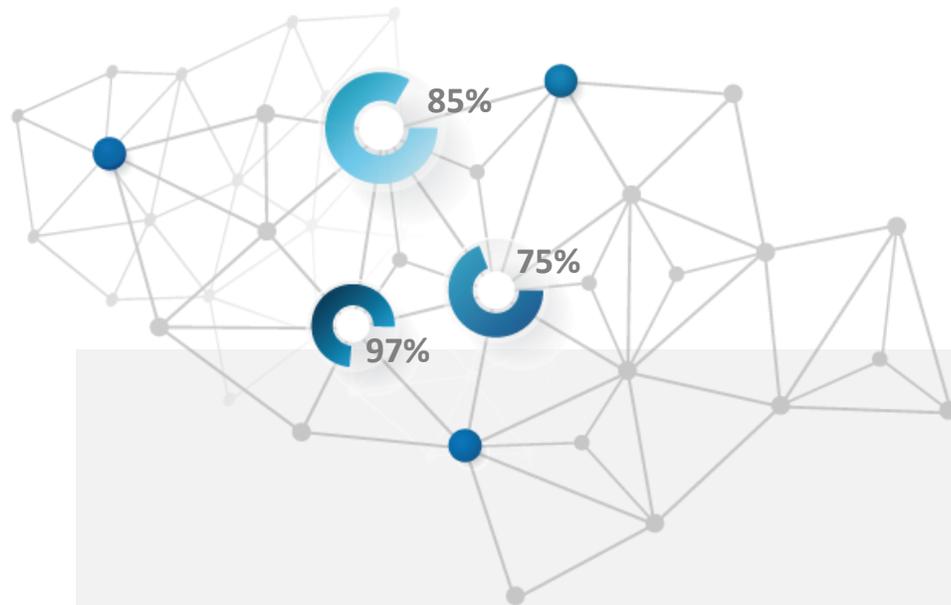
11. Juli 2019



averbis
text analytics

averbis

text analytics



Technologien: Natural Language Processing,
Machine Learning

Industrien: Healthcare, Pharma, IP

Produkte: Information Discovery,
Health Discovery,
Patent Monitor

Firmensitz: Freiburg im Breisgau

Gründungsjahr: 2007

VISION



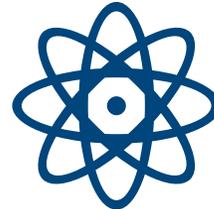
Softwarelösungen, die letztendlich die Unterschiede zwischen strukturierten und unstrukturierten Daten aufheben.



FIND



Turn text into actionable information



UNDERSTAND



Understand and automate cognitive processes



PREDICT



Enable better business critical decisions

KUNDEN

HEALTHCARE



RHÖN-KLINIKUM AG



PATENTS/IP

CENTREDOC



XPAT



PHARMA



Deloitte.



OTHERS



BMW Group



| BertelsmannStiftung

WISSENSCHAFTLICHE PROJEKTE

- EU: DEBUG-IT, MANTRA, SEMCARE, euCASES, ASSESS-CT
- BMBF:
 - MI-I: SMITH, MIRACUM, DIFUTURE
 - TOPOS, XplOit
- ZIM: PH³
- BMWI: Cloud4Health, Klinische Datenintelligenz (KDI)
- THESEUS: RadMining



TECHNOLOGIEN

De-Identifizierung von Dokumenten

TEXT MINING

'A **Lung Adenocarcinoma** biomarker panel comprising an microRNA'



MACHINE LEARNING

Oncology Therapy

- Use of Cabazitaxel in patients with metastatic NSCLC (WO-2015036507-A1)
- Treatment of Melanoma (EP-3021848-A2)



SEMANTIC SEARCH

Lung Cancer



'A **Lung Adenocarcinoma** biomarker panel comprising an microRNA'



TERMINOLOGIES

Lung Cancer

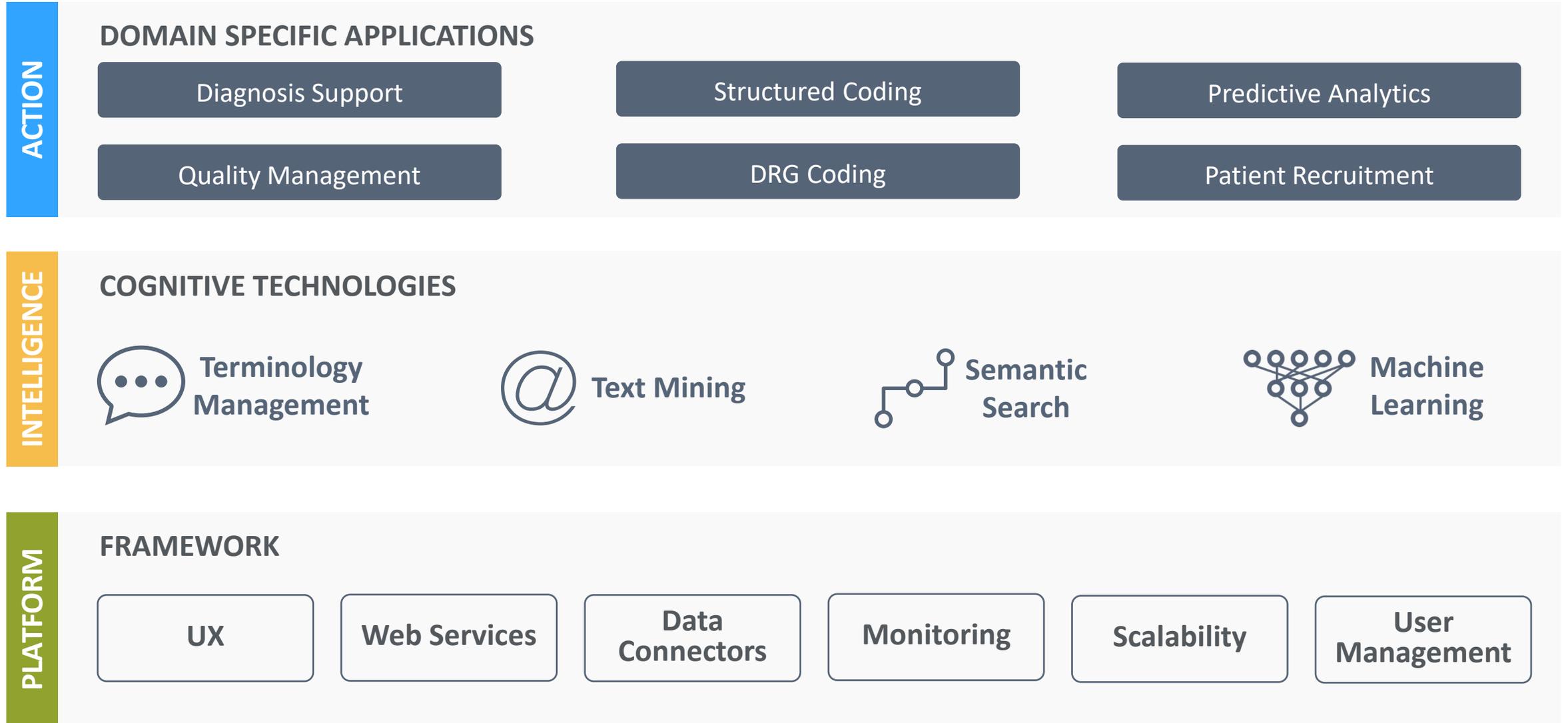
Lung Adenocarcinoma

Pulmonal Mass



h HEALTH
discovery

HEALTH DISCOVERY PLATFORM



TEXT MINING – DISCHARGE PIPELINE

Diagnosis **Drug** **Entity** **Labor** **LaboratoryValue** **Medication** **PreNegation** **Sentence** **Token**

Sehr geehrter Herr Kollege,

wir berichten über o.g. Patientin:

Diagnosen:

1. Koronare Herzkrankheit
2. Hypertrophische obstruktive Kardiomyopathie
3. Mitralklappeninsuffizienz Grad 1-2
4. Trikuspidalklappeninsuffizienz Grad 2-3
5. Arterielle Hypertonie
6. Chronisch venöse Insuffizienz

Anamnese:

Aktuelle Anamnese: Die Aufnahme der Patientin erfolgte über die Notaufnahme bei instabiler Angina pectoris mit Ausstrahlung in den linken Arm. Bei der Aufnahmeuntersuchung zeigte sich im EKG ein intermittierende absolute Arrhythmie bei Vorhofflimmern. Keine Ödeme.

Labor:

20.05.02 11:10 Uhr: Leukozyten 5,9 Tsd/µl; Erythrozyten 3,73 Mio/µl; Hämoglobin 11,8 g/dl; Hämatokrit 34,8 %; MCV 93,4 fl; MCH (HbE) 31,6 pg; MCHC 33,8 g/dl; Thrombozyten 342 Tsd/µl; Quick 100 %; Kalium 4,5 mmol/l; Natrium 137 mmol/l; Harnstoff 22 mg/dl; Bilirubin gesamt 4,9 mg/dl; GOT 43 U/l; GPT 33 U/l

Therapieempfehlung:

ASS 100mg 0-0-1
Concor 5 mg 1-0-0
Norvasc 5 mg 1-0-0
Pantozol 40 0-0-1
Delix 5 mg 0-0-1

Mit kollegialen Grüßen

Diagnoses

- ICD10 Codes
- Context

Lab Values

- Parameters
- Values
- Units
- Normalisation
- Interpretation*

Drugs

- Ingredients
- Brand Names
- Strengths
- Dose Forms
- Dose Schemes

Temporal Aspects

- Length of Stay
- Creation Date

AVERBIS DE-ID (DE-IDENTIFIZIERUNG VON COMPUTERLESBAREN TEXTEN)



Optical
character
recognition

Papier



Bild (pdf, jpg)



Text (txt, doc,)

Klinischer Befund:

72-jährige Patientin in ausreichendem AZ, Körpergröße 1,57 m, Körpergewicht 72,8 kg. Haut und sichtbare Schleimhäute gut durchblutet. Kein Ikterus, keine Zyanose. Mundhöhle o.B.. Keine Struma, keine peripheren Lymphknotenschwellungen. Über Herz und Lungen war der klinische Befund unauffällig. RR bds. 160/80 mm Hg, Pulsfrequenz regelmäßig. Leber und Milz nicht vergrößert. Nierenlager klopfeschmerzfrei. Keine Varizen, keine Oedeme. Fußpulse beiderseits tastbar. MER seitengleich.

OCR-ARTEFAKTE

Englischer Text

Laserdruck, Scan Graustufen, 8 Bit Farbtiefe



Jude the Obscure is the last of Thomas Hardy's novels, begun as a magazine serial and first published in book form in 1895.

Its hero Jude Fawley is a lower-class young man who dreams of becoming a scholar.

The two other main characters are his earthy wife, Arabella, and his intellectual cousin, Sue. Themes include class, scholarship, religion, marriage, and the modernization of thought and society.

Englischer Text

Jude the Obscure is the last of Thomas Hardy's novels, begun as a magazine serial and first published in book form in 1895.

Its hero Jude Fawley is a lower-class young man who dreams of becoming a scholar.

The two other main characters are his earthy wife, Arabella, and his intellectual cousin, Sue. Themes include class, scholarship, religion, marriage, and the modernization of thought and society.

Englischer Text

Kopie auf Thermokopierer mit Artefakten



Jude the Obscure is the last of Thomas Hardy's novels, begun as a magazine serial and first published in book form in 1895.

Its hero Jude Fawley is a lower-class young man who dreams of becoming a scholar.

The two other main characters are his earthy wife, Arabella, and his intellectual cousin, Sue. Themes include class, scholarship, religion, marriage, and the modernization of thought and society.

Englischer Text

V Jude the Obscure is the last of 'l'jhomns | begun as a magazine serial and first published in book [arm in 1895.

IIS]^i6I`O Jude Fawley 'iS' am luwuii-Cl|a\$\$^|hy<jling'|nian ivillio dreams ofbecoming a scholar.

The two other main qhdrabiers | wife, Arabella, and his intellectual cousin, Sue. Themes includepiassqvScholarship, religion, marriage, and the modernization of thought and society.

AVERBIS DE-ID (DE-IDENTIFIZIERUNG VON COMPUTERLESBAREN TEXTEN)



Erkennt und maskiert identifizierende Inhalte von medizinischen Texten

Grundbegriffe:

- Personenidentifizierende Merkmale (IDAT)
- De-Identifizierung: Identifizierende Inhalte werden erkannt und maskiert:
- Anonymisierung: Inhalte sind nicht mehr den Originaltexten zuzuordnen
- Tagging: automatische oder manuelle Markierung einer Textpassage, z.B. <name>Stefan Schulz</name>
- Pseudonymisierung: Über eine Vertrauensstelle ist eine Re-Identifizierung möglich
nicht: Ersatz identifizierender Merkmale wie Namen durch "Fantasiennamen"

HIPAA „Safe Harbor“ identifiers*

- (1) Names
- (2) Geographic subdivisions
- (3) Dates (except years)
- (4) Phone numbers
- (5) Vehicle identifiers
- (6) Fax numbers
- (7) Device identifiers and serial numbers
- (8) Mail addresses
- (9) URLs
- (10) Social security numbers
- (11) IP addresses
- (12) Medical record numbers
- (13) Biometric identifiers
- (14) Health plan beneficiary numbers
- (15) Photos
- (16) Account numbers
- (17) Certificate/license numbers
- (18) Any other identifying number

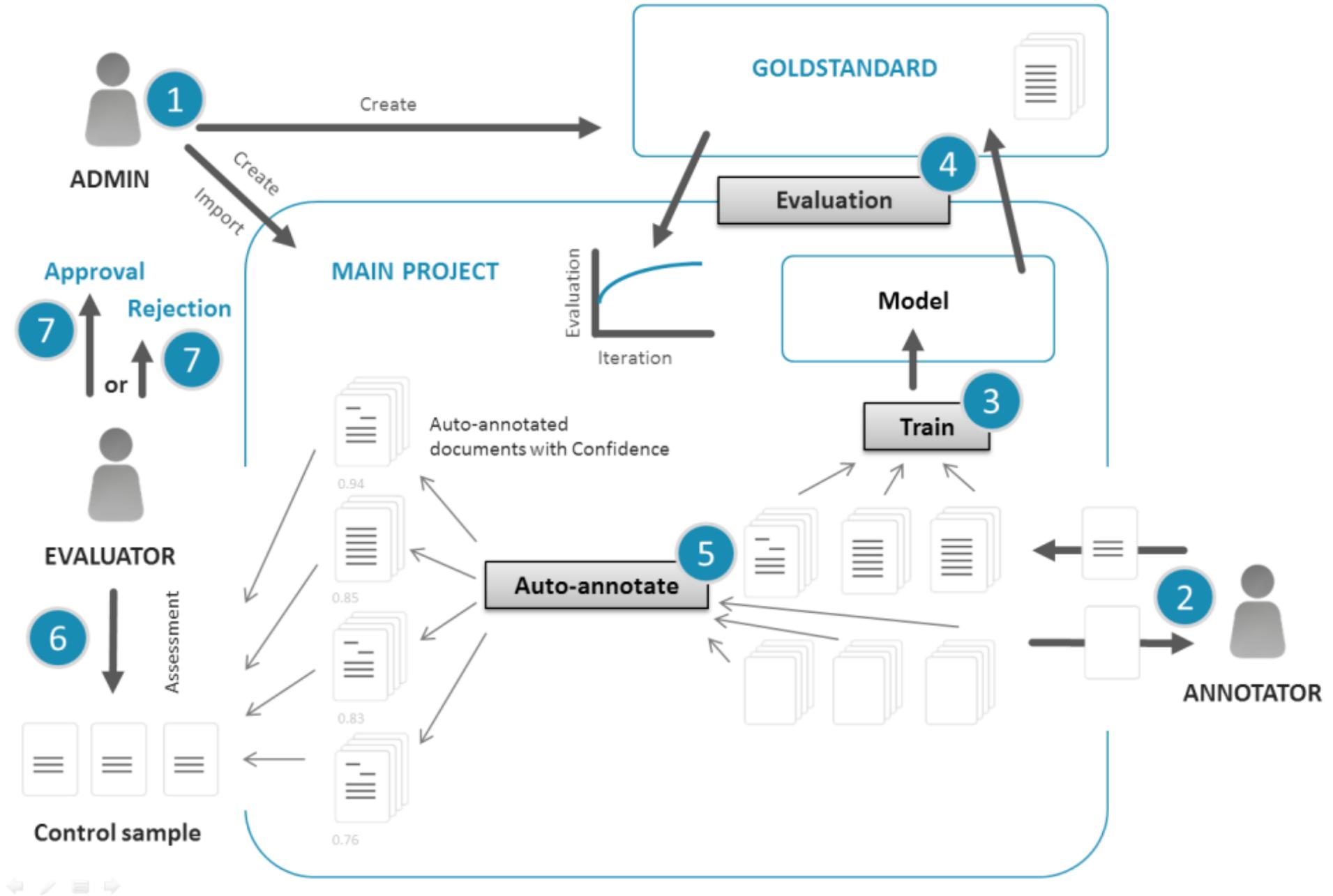
DE-IDENTIFIZIERUNGSPROZESS IN AVERBIS DE-ID

Drei Personengruppen sind beteiligt:

- Administrator: Verwaltung von Benutzerzugängen und Dokumentenkorpora, Import von Dokumenten, Export deidentifizierter Versionen
- Annotator: können nur die Dokumente der ihnen zugewiesenen Projekte bearbeiten. Sie annotieren dabei Dokumente und initiieren iterativ Training und die Evaluation des Modells.
- Evaluator: begutachtet die automatisch annotierten, sowie die manuell nachbearbeiteten Dokumente und entscheidet: Freigabe oder Nachbearbeitung.

Ablauf:

- Dokumente werden importiert
- Dokumente werden vom System vorannotiert (CRF Tagger)
- Annotationen werden von Annotatoren überprüft und modifiziert
- Modell wird trainiert, dadurch wird die automatische Annotation verbessert
- Dokumente werden exportiert (anonymisiert oder mit Tags versehen)



DE-IDENTIFIZIERUNG: TECHNIKEN

Metadata

Name Lists

Patterns

Machine Learning

Document PHI List Error List

Städtisches Krankenhaus Braunschweig
Abteilung Innere 3
Station DaVinci

Ärztlicher Direktor: Frau Prof. Marta Müller

Arztzimmer Tel. 0121/123 Fax. 01217223

Buxtehude, den 17. März 2012
Unser Zeichen: bk/ij

Herrn
Dr. med Martin Müller
Willy-Brandt-Allee 12
01292 Berlin

Betrifft Patienten Luise Kisselbach, geb. 17.07.1921

Sehr geehrter Kollege Müller,

wir berichten Ihnen nachfolgend über o.g. Patienten, der sich am 21. 04.2003 in unserer Behandlung befand.

Export:

Annotations-Tags:

Betrifft Patienten <name>Luise Kisselbach</name>, geb. <date>17.7.1921</date>

Sehr geehrter Kollege <name>Müller</name>

De-identifiziert:

Betrifft Patienten XXXXXXXXXX, geb. XXXXXXXXXX

Sehr geehrter Kollege <name>XXXXXX</name>

AVERBIS DEID



Erkennt und maskiert identifizierende Inhalte von medizinischen Texten

Universitätsklinikum Bad Krozingen Abteilung Neurologie Prof. Dr. Dr. Lutz Schmeisser
Ärztlicher Direktor der Neurologie Telefon 01258/5698245 Fax 01258/5698246

Universitätsklinikum Bad Krozingen
Abteilung: Neurologie
Station 3

Herrn Dr. Schröder Neurologie, Station 3 Universitätsklinikum Bad Krozingen
Nachrichtlich: Herrn Dr. Ludwig Wedel Neurologische Station Rehabilitationskliniken
Münstertal Belchenstraße 7 56895 Münstertal

Bad Krozingen, am 28.02.2016

Sehr geehrter Herr Dr. Wedel,

Wir überweisen Ihnen Herrn Jakob Bleifuss, geb. am 17.12.1990, wohnhaft in Veilchengasse 27a, 56895 Heitersheim.

Jakob Bleifuss stürzte am 15.2.16 beim Skifahren abseits der Piste schwer und blieb bewusstlos mit dem Kopf in einer Schneewehe stecken. Seine Begleiter befreiten ihn aus dem Schnee und reanimierten Herrn Bleifuss bis zum Eintreffen der Rettungskräfte.

Aufgrund der unterbrochenen Sauerstoffzufuhr kam es zu einer schweren Schädigung des Gehirns. Herr Bleifuss ist zum jetzigen Zeitpunkt (Stand 27.02.2016) nicht in der Lage zu sprechen und zu schlucken, weshalb am 17.2.16 eine PEG-Sonde zur künstlichen Ernährung gelegt wurde. Außerdem hat er Schwierigkeiten bei der Feinkoordination der Hände und kann nur mit Rollator ein paar Schritte gehen.

Annotations:		
<input type="checkbox"/>	RE_LI_ML Div.	"Universitätsklini..."
<input type="checkbox"/>	RE_LI_ML Name	"Prof. Dr. Dr. Lut..."
<input type="checkbox"/>	ML Div.	"Ärztlicher Direktor"
<input type="checkbox"/>	RE_ML Con.	"01258/5698245"
<input type="checkbox"/>	RE_ML Con.	"01258/5698246"
<input type="checkbox"/>	RE_LI_ML Div.	"Universitätsklini..."
<input type="checkbox"/>	RE_ML Div.	"Station 3 Herrn"
<input type="checkbox"/>	RE_LI_ML Name	"Dr. Schröder"
<input type="checkbox"/>	ML Div.	"Neurologie"
<input type="checkbox"/>	RE_LI_ML Div.	"Station 3 Univers..."
<input type="checkbox"/>	RE_LI_ML Name	"Dr. Ludwig Wedel"
<input type="checkbox"/>	ML Div.	"Neurologische"
<input type="checkbox"/>	RE_LI_ML Div.	"Rehabilitationskl..."
<input type="checkbox"/>	RE_LI_ML Loc.	"Belchenstraße 7 5..."
<input type="checkbox"/>	LI_ML Date	"28.02.2016"
<input type="checkbox"/>	RE_LI_ML Name	"Dr. Wedel"

Original

Krankenhaus der Samariter Holzhausen

Röntgenabteilung, Vorstand Prim. Univ. Prof. Dr.Dr. Gotthard Vogler

CT Abdomen und kl. Becken

Name: Mustafa Üstün, * 21.06.67

Aufnahmezahl: 1933309807

Abteilung: Chirurgie

Station: A31. OG. Viszeralchirurgie B /

Zi: 119

dikt. Arzt: OA Dr. Huber Karina

WinA. 06/07/2011

Getaggt

<location>Krankenhaus der Samariter</location>

<location>Holzhausen</location>

Röntgenabteilung, Vorstand <name>Prim. Univ. Prof. Dr.Dr. Gotthard Vogler</name>

CT Abdomen und kl. Becken

Name: <name>Mustafa Üstün</name>, *
<date>21.06.67</date>

Aufnahmezahl: <id>1933309807</id>

Abteilung: Chirurgie

Station: <division>A31. OG. Viszeralchirurgie B</division> /

Zi: 119

dikt. Arzt: <name>OA Dr. Huber Karina</name>

WinA. <date>06/01/2011</date>

Getaggt

<location>Krankenhaus der Samariter</location>

<location>Holzhausen</location>

Röntgenabteilung, Vorstand <name>Prim. Univ. Prof. Dr.Dr. Gotthard Vogler</name>

CT Abdomen und kl. Becken

Name: <name>Mustafa Üstün</name>, *
<date>21.06.67</date>

Aufnahmezahl: <id>1933309807</id>

Abteilung: Chirurgie

Station: <division>A31. OG. Viszeralchirurgie B</division> /

Zi: 119

dikt. Arzt: <name>OA Dr. Huber Karina</name>

WinA. <date>06/01/2011</date>

Anonymisiert

XXXXXXXXXX XXXXXXXXXXXX

Röntgenabteilung, Vorstand Prim. Univ. Prof. Dr.Dr. XXXXXXXX
XXXXXXXXXXXXXXXXXX

CT Abdomen und kl. Becken

Name: XXXXXXXX XXXXXXXXXXXX, * X.X.X

Aufnahmezahl: XXXXXXXXXXXXX

Abteilung: Chirurgie

Station: XXXXXXXXXXXX XXXXXXXXXXXXX

Zi: 119

dikt. Arzt: OA Dr. XXXXXXXXXXXX XXXXXXXXXXXX

WinA. XX/XX/XXXX

Getaggt

<location>Krankenhaus der Samariter</location>

<location>Holzhausen</location>

Röntgenabteilung, Vorstand <name>Prim. Univ. Prof. Dr.Dr. Gotthard Vogler</name>

CT Abdomen und kl. Becken

Name: <name>Mustafa Üstün</name>, *
<date>21.06.67</date>

Aufnahmezahl: <id>1933309807</id>

Abteilung: Chirurgie

Station: <division>A31. OG. Viszeralchirurgie B</division> /

Zi: 119

dikt. Arzt: <name>OA Dr. Huber Karina</name>

WinA. <date>06/01/2011</date>

Anonym und Verfremdet (separater Prozess)

Kantonsspital Friedrichshafen

Röntgenabteilung, Vorstand Prim. Univ. Prof. Dr.Dr. **Gerhard Voigtländer**

CT Abdomen und kl. Becken

Name: **Manuel Überreuter**, * **1.07.69**

Aufnahmezahl: **9983209971**

Abteilung: Chirurgie

Station: **Station Sauerbruch**

Zi: 119

dikt. Arzt: OA Dr. **Heilmann Kristina**

WinA. **16/07/2013**

GÜTE DER AUTOMATISCHEN VERFAHREN

- Text-Mining ist kein 100%-Verfahren
- Abhängig von den Spezifika des Dokumentenkopus
- Güte der automatischen Verfahren i.A. mit manueller De-Identifizierung vergleichbar
- Studie Radiologie Erlangen :
 - Training verbessert kontinuierlich die De-Identifizierungsperformance
 - Ab etwa 100 korrigierten Texten wurde die Detektion als zuverlässig angesehen



AUSBLICK AVERBIS DE-ID FÜR NAKO ENDPOINTVALIDIERUNG

Derzeit Neuimplementierung als Komponente von Averbis Health Discovery

- 1. Release Ende des Jahres

Nutzerwünsche können mit eingearbeitet werden

- Features
- Batch-Processing
- Post-Processing

KONTAKT

Stefan Schulz, Averbis GmbH: stefan.schulz@averbis.com