

AI powered Data Curation & Publishing Virtual Assistant

*Optimize interoperability & quality of health data to increase data sharing and reuse
across Clinical Registries & Personal Data Intermediaries*

WP 5 Workshop

Annotation Process and NLP Workflow

Markus Kreuzthaler, Sareh Aghaei, Kris Collins, Todor Primov, Stefan Schulz



Funded by
the European Union

- Large parts of EHR content only available as free-text
- Manual markup of concepts and relations of and between text elements (words, word sequences) in clinical narratives
- Purpose:
 - Training / domain-fine-tuning NLP approaches for information extraction for AIDAVA
 - Gold standard for evaluating these approaches
- Annotation vocabulary
 - Concept and relation names/codes used for annotation, consistent with AIDAVA reference ontology
- Pre-annotation
 - Use of existing text analytics pipelines to facilitate the annotation process
- Annotation guideline
 - Set of instructions for training annotators
 - Consistent and comparable annotation results

Task 4.3 Current state: guidelines, tools & training



Call: HORIZON-HLTH-2021-TOOL-06
Topic: HORIZON-HLTH-2021-TOOL-06-03
Funding Scheme: HORIZON Research and Innovation Actions (RIA)

Grant Agreement no: 101057062



AI powered Data Curation & Publishing Virtual Assistant

Deliverable No. 4.3
Update to Annotation guidelines, tools & training

Preparation phase until December 2022

- Access to **some** clinical narratives
- Overview of the annotation process
- **Selection of the annotation tool**
- **First draft** of annotation guidelines



Leading healthcare terminology, worldwide



HL7[®] FHIR[®]

Pilot phase until **May** 2023

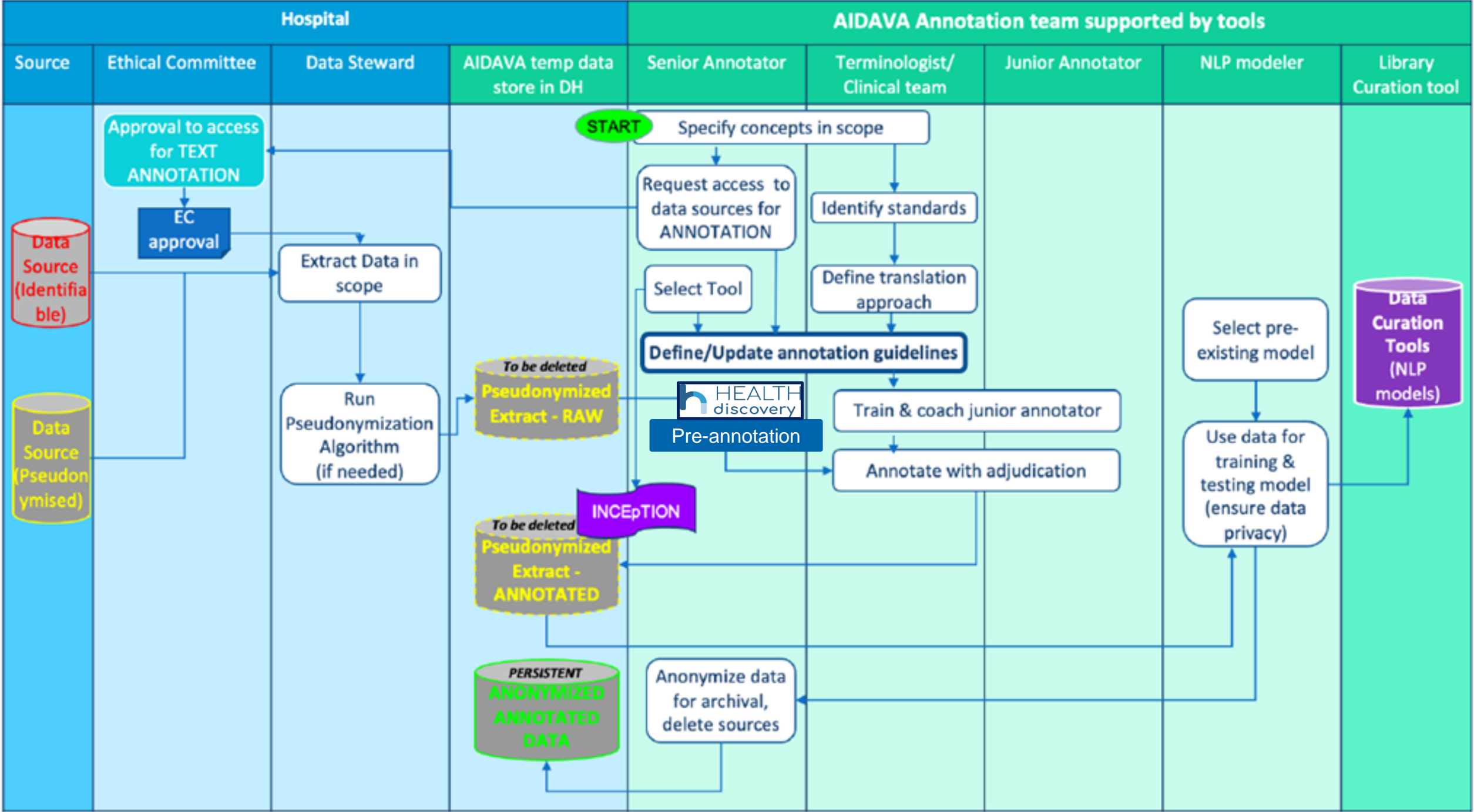
- Revise annotation guidelines
- Access to **use case oriented** clinical narratives
- Feedback from annotators
- First steps for pre-annotation tooling



h HEALTH
discovery



ontotext

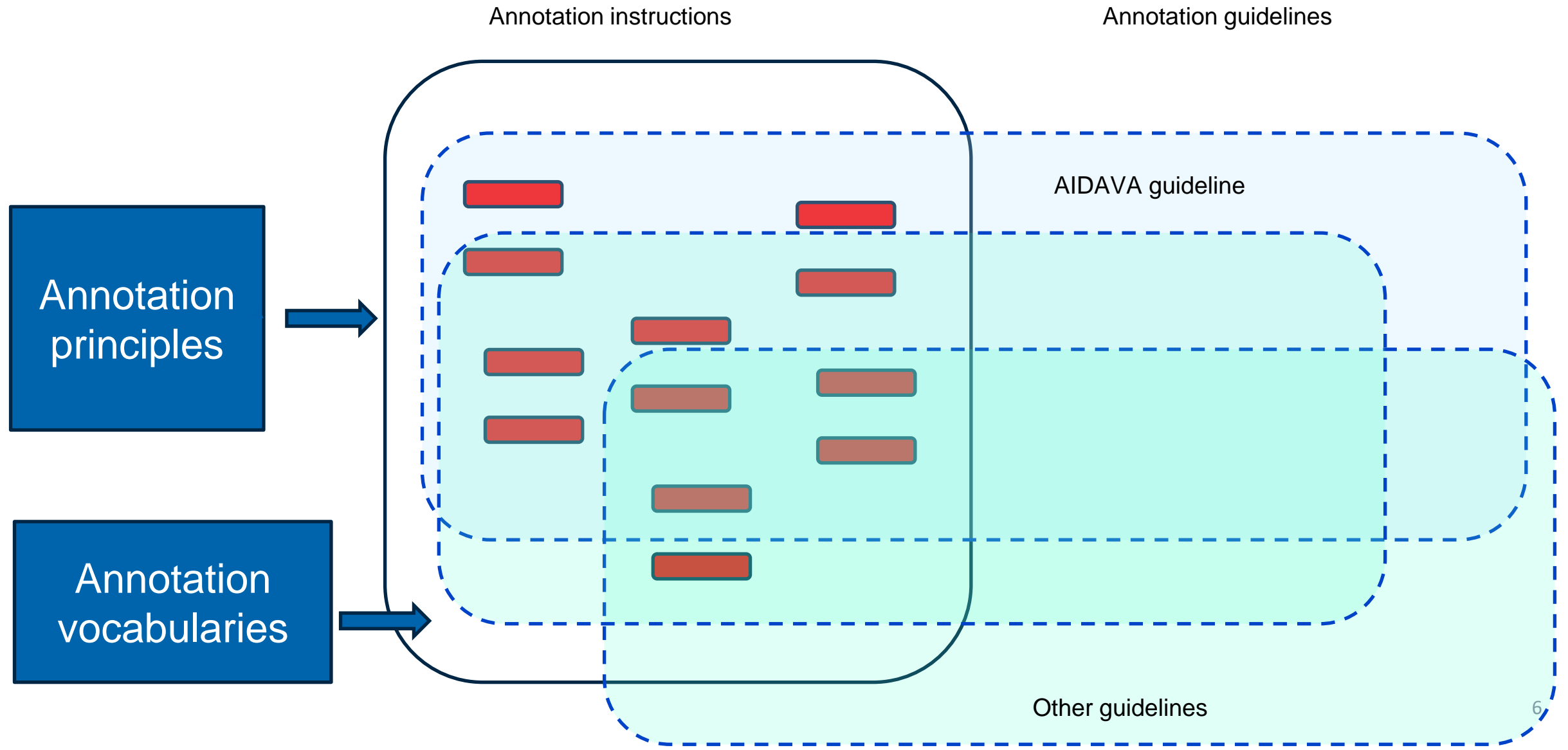


From annotation principles to exemplification

- Semantic annotation
 - Following annotation principles
- Assignment of codes and relations to
 - Documents
 - Database entries
- Codes
 - SNOMED CT subsets
 - LOINC? – mapping at a later stage
- User friendly relation “predicates”
 - Harmonized with SNOMED and FHIR



Annotation: Principles, instructions, guidelines



Guideline – Example 1

- Annotation principle:

Annotate the span that corresponds to the meaning represented by a concept in the annotation vocabulary

-> instruction: give preference to SNOMED CT precoordinated content

↗ anno:laterality ↘		
	85421007 Structure of right half of body (body structure)	
	28576007 Open fracture of femur (disorder)	
Open fracture of	left	femur

↗ anno:site ↘		
397181002 Open fracture (disorder)		734143007 Structure of left thumb (body structure)
Open fracture	of	left thumb

Guideline – Example 2

- Annotation principle:

Annotate , what is in the text, not what you interpret

-> instruction: don't interpret temporal sequence as causality unless explicitly stated

↗ anno:after ↘	
25064002 Headache (finding)	110030002 Concussion injury of brain (disorder)
Headache	following brain concussion

↗ anno:dueTo ↘	
162049009 Left flank pain (finding)	45816000 Pyelonephritis (disorder)
Left flank pain	in pyelonephritis

Guideline – User-friendly predicates



- During annotation: choice of predicates dependent on domain and range constrained to be learned by annotators
- Automatic transformation to standards-conforming expressions after annotation

Example binary predicates for relation annotations with their translation into SNOMED CT and FHIR syntax. “INV” = inverse, “||” = concatenation. The relation paths marked with [a] refer to the concatenation of SNOMED CT relations, those marked with [b] to roughly equivalent FHIR elements

anno:	Domain	Target path	Range
site	‘sct:Clinical finding’	[a] ‘sct:Finding site’ [b] INV(fhir:Condition.code) fhir:Condition.body	‘sct:Body structure’
site	‘sct:Procedure’	[a] ‘sct:Procedure - Direct’ [b] INV(fhir:Procedure.code) fhir:Procedure.body	‘sct:Body structure’
inFamily	‘sct:Clinical finding’	[b] INV(fhir:FamilyMemberHistory.condition) fhir:FamilyMemberHistory.relationship [a] INV(‘sct:Associated finding’) ‘sct:Subject relationship context’	‘sct:Person’
verification status	‘sct:Clinical finding’	[b] INV(fhir:Condition.code) fhir:Condition.verificationStatus [a] INV(‘sct:Associated finding’) ‘sct:Finding context’	Qualifier value’ (cf. Tab. 1)

Task 4.3 Manual Annotations

Exemplification “tickets”



- Growing repository of annotation examples (“tickets”)

Table 10: Example of smoking behaviour

Input	Smoking: none, stopped in 2000, smoked for 5 years before that
<div><div>INCEpTION screenshot</div></div>	
Adjudication Description	<p>Two predicted patterns (stop at a specific date, and time duration) in smoking data appear in the input text. Therefore, instructions X and Y are followed to annotate ‘stopped in 2000’ and ‘smoked for 5 years’ as ‘Date ceased smoking’ and ‘Total time smoked’, respectively. Also, the numbers 2000 and 5 were coded as decimals. Moreover, the time-unit needs to be identified (as instructed in Section 6.2), so ‘years’ correspond to ‘year (qualifier value)’.</p> <p>As shown in Table 5, the predicate ‘value’ has observable entity and decimal as its domain and range, respectively. Thus, we use the predicate ‘value’ between the identified observable entities and their corresponding decimals. The same explanation applies to ‘valueUnit’ between the observable value and the qualifier value (i.e., year).</p>

Task 4.3 Manual Annotations

From exemplification tickets to full document annotations



INCEpTION Projects Dashboard Help Administration admin Log out 29 min

1-62 / 62 lines [doc 6 / 16]

Tumor progression

28 DEKURS DER TUMORERKRANKUNG

29

30 **Surgical procedure** **Histologic test** **Histologic test** **Molecular genetic test**
Operation(en), Histologie(n), Immunhistologie(n), Molekulare(s) Profil(e):

31 **Partial mastectomy** **Excision of axillary lymph node** **Excision** **Procedure status** 2009
kurative TE und axill. Dissektion, Nachresektion (***** 2009, *** - CHIR-KLINIK)

32 **Histologic test** **Infiltrating duct carcinoma of breast** **maximal** 248530000 **Diameter of lump** **unit** **Centimeter** 3
Histo : IDC (max . DM : 3cm),

American Joint Committee on Cancer pT2 (qualifier value) **American Joint Committee on Cancer pN1a (qualifier value)** **February 2012**
p T2 N1a (2/12)

American Joint Committee on Cancer grade G2 **American Joint Committee on Cancer R0** **Clinical stage finding**
, MX, G2 , R0 , Stadium klinisch:

Proliferation marker protein Ki-67 **value** 0.2
Mib-1 : bis 20%

33 **Histologic test** **Solid ductal carcinoma in situ of breast** **American Joint Committee on Cancer pTis**
Histo : u. DCIS, solider Wachstumstyp, p Tis , G2

34 **Immune (qualifier value)** **High** **High**
Histologic test **Oestrogen receptor positive tumour** **Progesterone receptor positive tumour** **HER2-positive carcinoma of breast**
Immunhisto : ER hoch pos , PR hoch pos , Her2neu 1+

Annotation

← →

Layer

Custom MCN

Text

axill. Dissektion

No links or relations connect to this annotation.

Comment

Concept

234262008 | Excision of axillar

Short

Excision of axillary lymph node

Determined by **reference ontology**

- Subsets of SNOMED CT with the maximum coverage for the prioritized items of the use cases

When and how to access the reference ontology?

- Annotators use just the SNOMED CT browser, abstraction to reference ontology content done automatically?
- Annotators have only access to reference ontology?
- Reference ontology imported into annotation tool?

Task 4.3 Manual Annotations

Pre-annotation

Using existing NLP pipeline: **Averbis Health Discovery**

Texteingabe

Text Analyse Ergebnisse

Anatomy

ClinicalSection

ClinicalSectionKeyword

✓ Concept

Date

Diagnosis

DocumentAnnotation

EstrogenReceptor

HER2

Morphology

PatientInformation

ProgesteroneReceptor

✓ TNMGrading

✓ TNMMetastasis

✓ TNMNode

✓ TNMTumor

DEKURS DER TUMORERKRANKUNG

Operation(en),Histologie(n),Immunhistologie(n),Molekulare(s)Profil(e):
kurative TE und axill. Dissektion, Nachresektion (*****2009, *** ***** - CHIR-KLINIK)
Histo: IDC (max. DM: 3cm), p T2 N1a (2/12), MX, G2, R0, Stadium klinisch: Mib-1: bis 20%
Histo: u. DCIS, solider Wachstumstyp, p Tis, G2
Immunhisto: ER hoch pos, PR hoch pos, Her2neu 1+

Bestrahlung(en):
postop. RTX Restbrust re. 60 GY (***09-***10)
pall. RTX LWS (***12-)

Med. TU Therapie(n):
adj. HT mit Arimidex (***09-***12)
(***10-***12) ABCSG-18
PD mit Skelettmetastasen;
supp. Thx. mit Denosumab/Placebo i. R. d. ABCSG-Studie

search

Filter

TNMGrading

G2

begin: 233

end: 235

value: G2

Concept

G2

begin: 233

end: 235

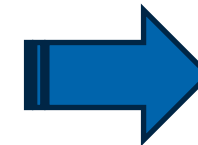
conceptID: 1228850007

dictCanon: G2

matchedTerm: G2

source: SNOMED_CT_GIT

uniqueID: SNOMED_CT_GIT:1228850007



INCEpTION recommender functionality

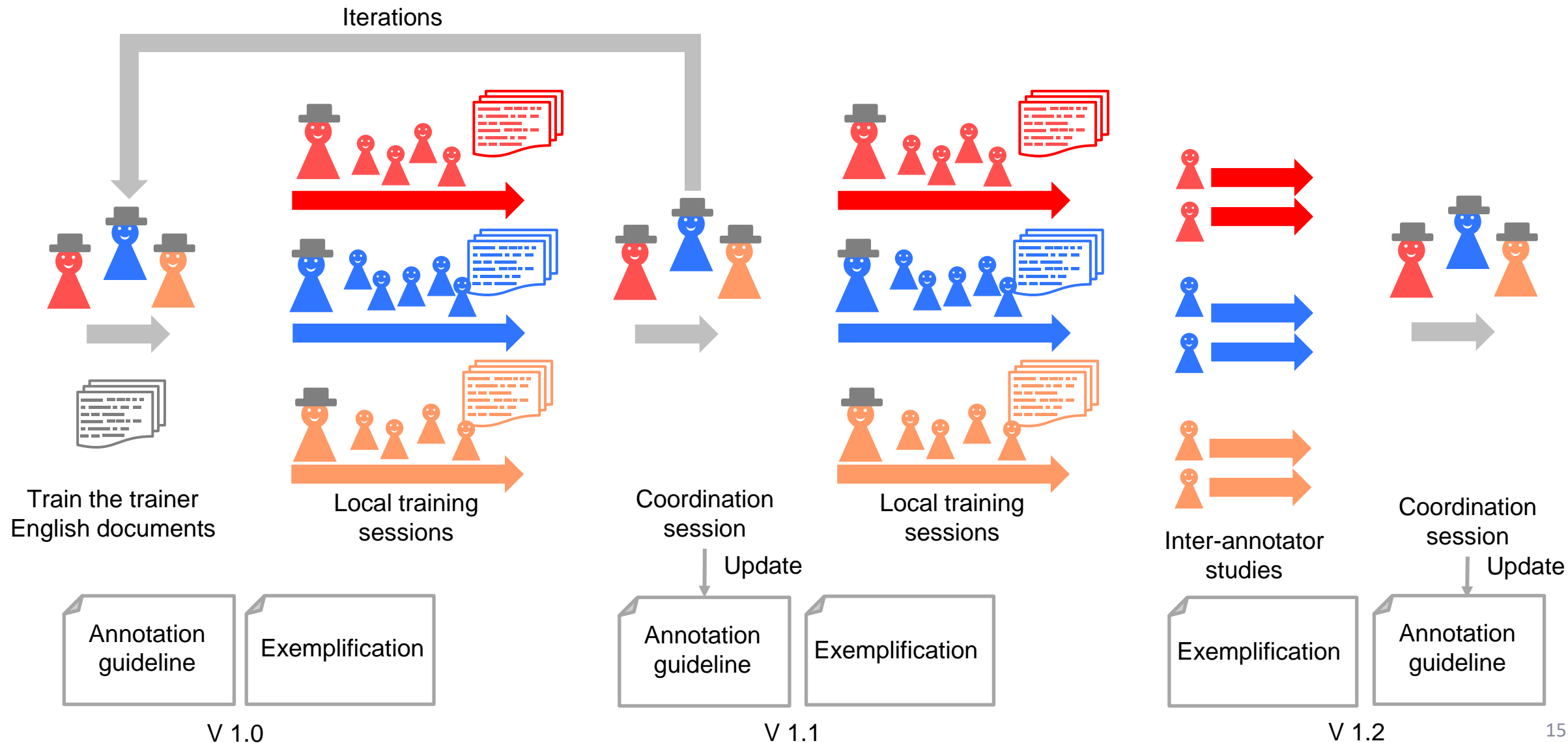
Details

Save Cancel

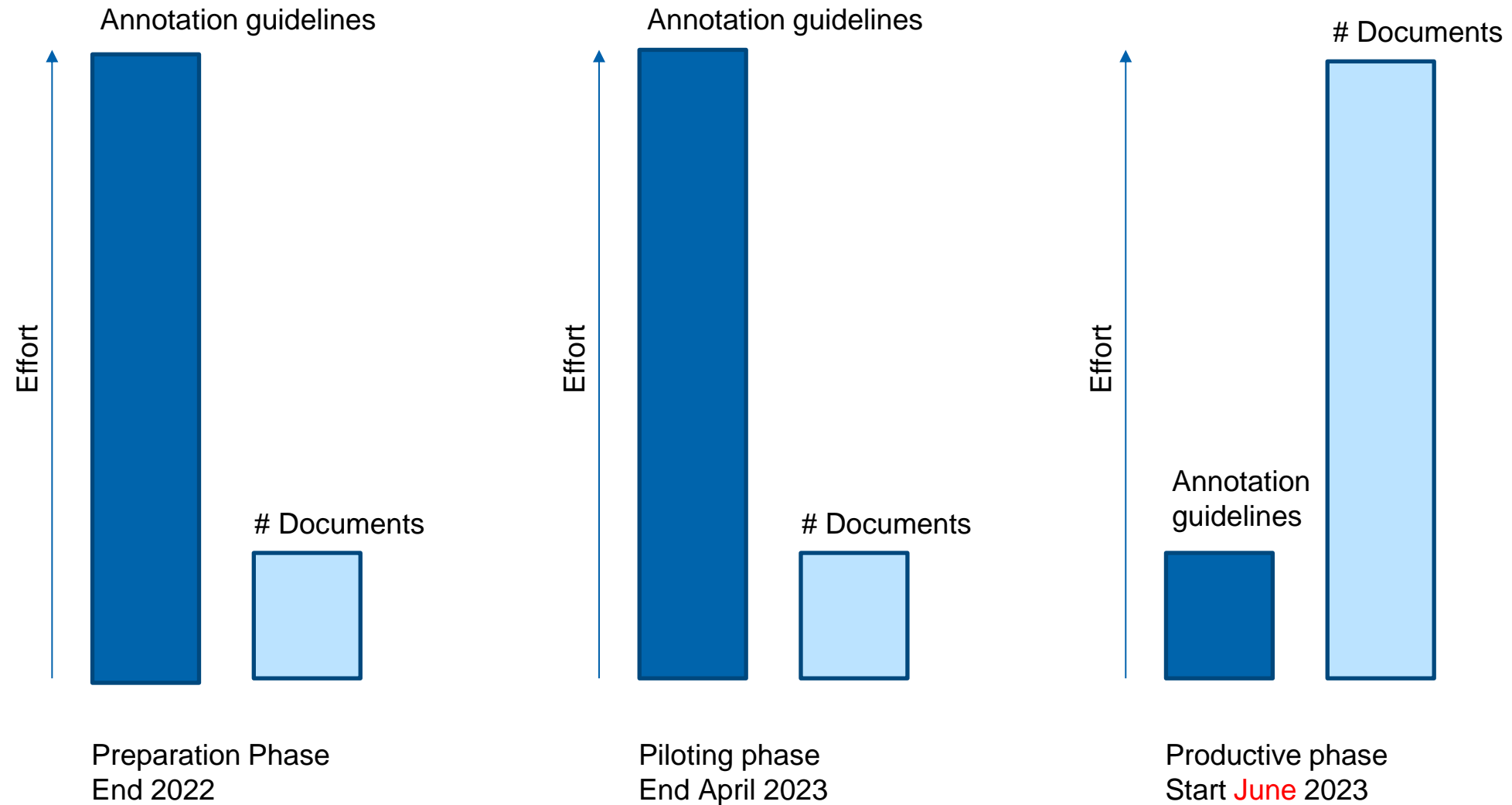
Name	<input type="text" value="[Custom MCN@Concept] String Matcher"/>		<input checked="" type="checkbox"/> auto-generate
	<input checked="" type="checkbox"/> Enabled		
Layer	<input type="text" value="Custom MCN"/>		
Feature	<input type="text" value="Concept"/>		
Tool	<input type="text" value="String Matcher"/>		
Activation strategy	Score threshold	<input type="text" value="0,0"/>	
	<input type="checkbox"/> Always active (no evaluation)		
Max. recommendations	<input type="text" value="3"/>		
States used for training	<input checked="" type="checkbox"/> Annotation not started yet (new)		
	<input checked="" type="checkbox"/> Annotation in progress		
	<input checked="" type="checkbox"/> Annotation finished		
	<input checked="" type="checkbox"/> Document not available for annotation (locked)		
	<input type="checkbox"/> Case insensitive		

Task 4.3 Manual Annotations

Update process in the production phase



Annotation work load



- **Canonical form** of semantic representation of clinical narrative
- **Maximally consistent** inter-annotator agreement (IAA)
- **Annotation graph** as the primary knowledge graph (KG)



Principles of ontology-based annotation of clinical narratives[★]

Stefan Schulz^{1,2,*}, Warren Del-Pinto³, Lifeng Han³, Markus Kreuzthaler¹,
Sareh Aghaei Dinani¹ and Goran Nenadic³

¹*Institute for Medical Informatics, Statistics and Documentation, Medical University of Graz, Austria*

²*Averbis GmbH, Freiburg, Germany*

³*Department of Computer Science, University of Manchester, UK*

Task 4.3 Manual Annotations

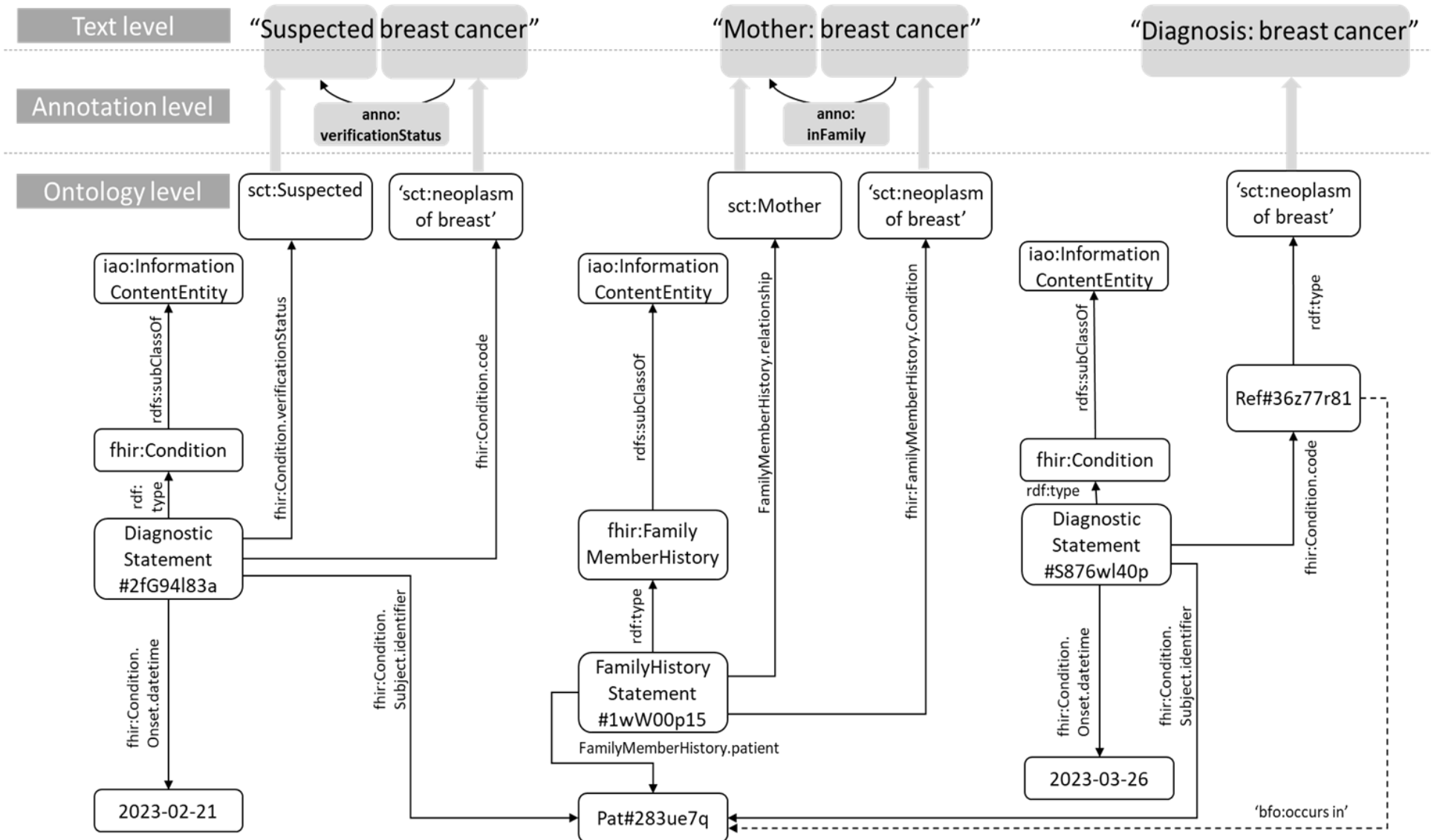
Scientific output - Synergies



- Technical University of Munich (Germany)
 - Research Consortium DIFUTURE
 - Methods Platform GeMTex
- University of Erlangen - Nürnberg (Germany)
- German Research Centre for Artificial Intelligence
- University of Manchester (UK)
- University of Murcia (Spain)
- University of Chiang Mai (Thailand)

- Discussion

From annotation to ontology-based knowledge graph



Example coreference



t_1

126926005
|Neoplasm of breast
(disorder)|



NoB34u73axn4us

owl:sameAs

t_2

254837009
|Malignant neoplasm
of breast (disorder)|

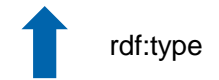


MoBkj88935elp

owl:sameAs

t_3

278054005 |Lobular carcinoma
of breast (disorder)|

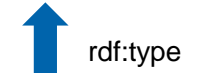


LCBrrp009g65t

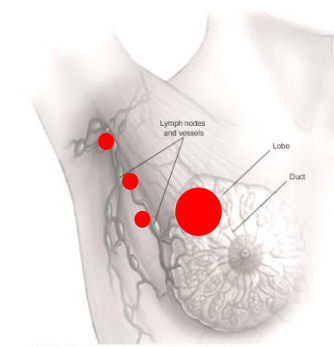
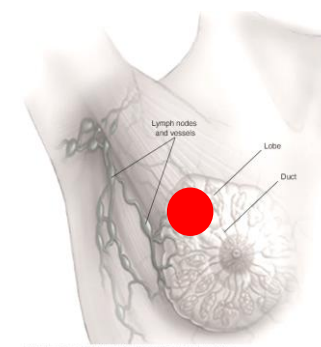
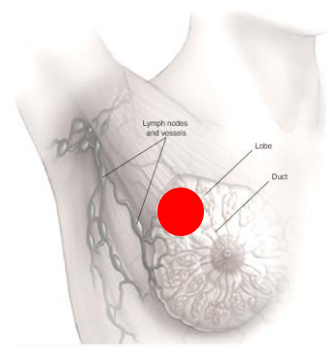
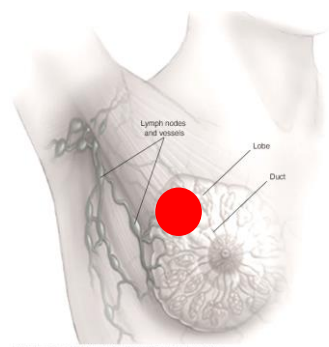
owl:sameAs

t_4

713609000 |Invasive
carcinoma of breast (disorder)|
AND
278054005 |Lobular carcinoma
of breast (disorder)|



ILCB4tz5ppklI



Lobular
carcinoma of breast

Invasive
Lobular carcinoma of breast

Thank you