




The TMF and its Metadata and Linked Data Forum

15th International Open Forum on Metadata



Matthias Löbe¹, Johannes Drepper², Sylvie MN Nguoungo³, Jürgen Stausberg³, Philippe Verplancke⁴

¹ Institute for Medical Informatics (IMISE), Universität Leipzig

² TMF - Technology, Methods, and Infrastructure for Networked Medical Research, Berlin

³ IBE, Ludwig-Maximilians-Universität München

⁴ XClinical GmbH, München

Overview

- + Introduction
- + Current projects in the Telematics Platform with regard to metadata
- + The Forum Metadata and Linked Data
- + Implications for the application of ISO 11179-3 edition 3 in a community-oriented approach

Starting point: Trial Item Manager (2008)

- + RDF/OWL-based tool for representing items in clinical trials
- + Planned as an in-house solution only
- + Ontological data model was simple but easily extensible

The screenshot displays the Trial Item Manager interface. On the left, a 'Navigation' pane shows a 'Component tree' with a hierarchical structure of clinical trial components. The 'Creatinine' component is selected and expanded, showing its sub-items: Grade 4 (> 6.0 x N), Grade 1 (N - 1.5 x N), Grade 0 (within normal range WNR), Inputfield_Length1 character, Grade 3 (> 3.0 x N - 6.0 x N), Grade 2 (> 1.5 N - 3.0 x N), Constipation, and Neutrophils.

The main window, titled 'Creatinine', shows the configuration for the 'Trial Item: Creatinine'. It includes sections for 'Characteristics', 'Components', 'Types', and 'Components containing this component'. The 'Characteristics' section shows a 'Label' of 'Creatinine' and a 'Database field' dropdown. The 'Components' section lists several components, including 'Grade 4 (> 6.0 x N)', 'Grade 1 (N - 1.5 x N)', 'Grade 0 (within normal range WNR)', 'Inputfield_Length1 character', 'Grade 3 (> 3.0 x N - 6.0 x N)', and 'Grade 2 (> 1.5 N - 3.0 x N)'. The 'Types' section shows 'Generic Component' and 'Trial Item'. The 'Components containing this component' section lists several trial items, including 'cycle 1 (Please give highest grade.)', 'cycle 6 (TG-3) - Give highest grade.', 'cycle 5 (TG-3) - Give highest grade.', 'cycle 3 (TG-2/3) - Give highest grade.', 'cycle 2 (Please give highest grade.)', and 'cycle 4 (TG-2/3) - Give highest grade.'.

What we intended to replace ...

+ A poor man's trial specification

Vorbekannte Tumorerkrankungen, ggfs mit ICD-10: Suchmöglichkeit ICD-10 C00-D48

Chemotherapie in Anamnese: J/N/weiß nicht

Strahlentherapie in Anamnese: J/N/weiß nicht

[Sonstige Tumorbehandlungen: J/N, bei J Freitext zum Spezifizieren **!zunächst rauslassen!**]

f Relevante Begleiterkrankungen: s. Datei „relevante Begleiterkrankungen“ (der einige Organ \rightarrow Spezifität)

[Tumor-Histologie n. ICD O.3: Freitext mit genauer Bezeichnung / ~~zunächst ersetzen durch~~

Tumor-Grading: G 1-4

\searrow Art des Primärtumors

Rezeptorstatus, ~~HER2-Status~~
HER2-Status

[Patholog. Institut + Einsendenummer)???) **zunächst rauslassen**

Klin. Bzw. Pathologischer TNM-Status: Pulldown für T - N - M

We had built it, why didn't they come?

- + We'll never know for sure, but there is some indication ...
 - + (The software wasn't self-explaining to use and had english labels)
 - + (Data managers were not familiar with the terminology used)
 - + When there was more content, data managers had no idea how to decide which data elements were superior to others
 - + The problem addressed was only a brick in an integrated solution
 - + SOPs for data management require CRFs to conform to an Excel template and a certain powerpoint layout
 - + Changes made later during database setup were not synchronized
 - + Re-use was limited, because in a specific trial, question texts and validation rules are very special
 - + No political support („It was always done this way...“)
 - + Community had too few active contributors
 - + Few community features available

Existing Clinical Metadata Repositories

| Tool | 11179 | Content | Tools | API | Community | Open Source |
|------------------------|-------|---------|-------------------------|-----|-----------|-------------|
| caDSR (USA) | V2 | ✓ | ✓ | ✓ | 🗄️ | 🗄️ |
| UK Cancergrid (GB) | V2 | ? | ✓ | ✓ | 🔍 | ✓ |
| METeOR (Australien) | V2 | ✓ | 🗄️ | ✓ | 🔍 | 🔍 |
| CIHI (Kanada) | V2 | ✓ | Currently not available | | | |
| USHIK (USA) | V2 | ✓ | 🗄️ | 🔍 | 🗄️ | 🔍 |
| MDR (Deutschland) | V3 | ✓★ | ✓★ | ✓★ | ✓★ | ✓★ |

Web-based Repository - Software as a Service - Community Approach - Bottom-up-Harmonization

★Planned

Metadata Repository for Clinical and Epimeidiological Research (MDR)

- + Establishment of a national service for providing harmonized data elements and Case Report Forms
 - + Based on draft edition 3 of ISO 11179
 - + Foundation in Top-Level-Ontology GFO
 - + Import for ODM and ClAML files
 - + Bottom-up community approach
 - + GWT prototype expected in September
- + Reference to metadata:
 - + obvious

Top TMF Projects Utilizing Metadata

- + MDR
- + EHR₄CR
- + KISREK
- + Biobank Registry/ P₂B₂
- + ID-Tools
- + e-Archiving
- + DRT
- + Cloud₄Health

HIS-based Support for Patient Recruitment for Clinical Trials (KISREK)

- + Support for patient recruitment by integrating software tools into the hospital information system routine workflow (trial registry, query engine, screening list, notification service)
- + 5 HIS vendors: Agfa Orbis, Siemens Soarian, Siemens medico, Siemens ISH*med and KAOS
- + Reference to metadata:
 - + Work package 3: detailed report about the suitability of HIS routine data for recruitment
 - + Developed a list of widely-used inclusion/exclusion criteria
 - + Decomposition of free-text into „computable criteria“ showed:
 - + 50% correspond to a single datum in the HIS
 - + 30% correspond to two or three dates
 - + 50% of all inclusion/exclusion criteria are documented in the HIS in principle
 - + But: in many cases incomplete or not in time, decomposition is time-consuming
 - + Most suitable: master data, diagnosis, procedures, lab values, observations

Electronic Health Records for Clinical Research (EHR₄CR)



- + EU project with 33 partners to build a distributed technical platform accessing local data warehouses
- + 4 usage scenarios:
 - + Protocol feasibility: Leverage clinical data to design viable trial protocols and estimate recruitment (cohort estimation)
 - + Patient recruitment: Detect patients eligible for trials to better utilize recruitment potential
 - + Clinical trial execution: Re-use routine clinical data to pre-populate trial CRFs
 - + Pharmacovigilance: Detect adverse events and collect/transmit relevant information
- + Reference to metadata:
 - + Development of a central „Pivot Ontology“ of 100 data elements for eligibility
 - + Semantic mapping from local data to the pivot ontology
 - + Local data elements are immutable

Biobank Registry/ P2B2

- + National registry for biobanks
- + BioMedBridges: EU project providing interoperable services
- + Researchers want to maintain control of their data
 - + De-centralized peer infrastructure
 - + Query tool to request samples
- + Reference to metadata:
 - + Core Data Set
 - + Domain data sets
 - + Basic Biobanking Ontology (BBO)

ID Tools

- + Data security and privacy are big issues in clinical research
- + PID service creates an pseudonym (unique identifier) for a set of patient identification data (similar to HIPAA)
- + PSD service creates an second order pseudonym to be managed by a trusted third party
- + Reference to metadata:
 - + Management of personal data and identifiers
 - + Referent tracking

Long-term archiving (LABIMI/F)

- + Archiving of biomedical research data
 - + Genomic data
 - + Imaging data
- + Need for vocabularies to describe!
 - + Preservation
 - + Provenance
 - + Curation
- + Reference to metadata:
 - + Dublin Core Metadata
 - + LOINC, MeSH, SNOMED CT, UMLS

Integrated Data Repository Toolkit (IDRT)

- + Provides tools and services around the Harvard i2b2 Data Warehouse software
 - + Wizard for semi-automatic installation
 - + ETL import jobs for SQL, CSV and ODM files
 - + Standard terminologies like ICD-10, OPS, LOINC, MedDRA
 - + Data security and privacy via pseudonymization service
- + Reference to metadata:
 - + Metadata editor to provide mappings and alignments for data elements in the i2b2 ontology cell
 - + NCBO BioPortal as ontology source under testing

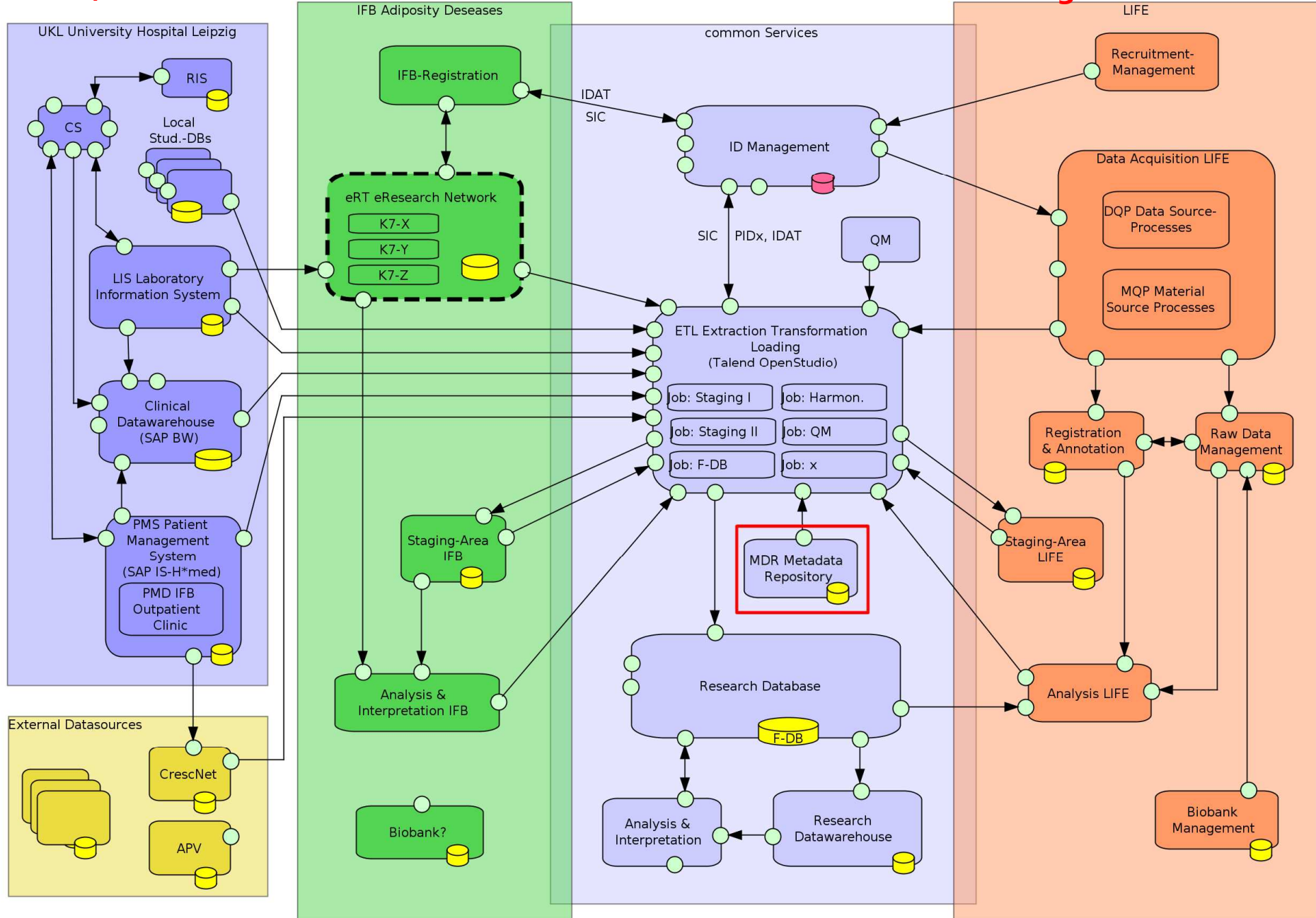
- + Cloud-Computing in Healthcare
 - + Secondary Use of unstructured data (text-analysis)
 - + Data Warehouse technologies in the cloud
 - + Establishment of an infrastructure
- + Use Cases:
 - + Early detection of adverse events
 - + Cost-effectiveness of therapies
- + Reference to metadata:
 - + Mapping named-entities in discharge letters to SNOMED CT (possibly Observable Entities)

HIS/eSource

Clinical Trials

Central Services

Long-term Cohorts



Forum Metadata and Linked Data

- + Founded June 2011 with a focus on medical research
- 1. Concepts and Methods:
 - + Metadata Models (ISO 11179, CDISC ODM, EN 13606, HL7 CDA)
 - + Metadata Artifacts (Std. Values Sets, UCUM, 21090)
 - + Metadata Annotations (med. Terminologies, DC, SKOS)
- 2. Representations und Implementations:
 - + Metadata Element Sets (CDASH , HITSP Data Dict., NINDS CDE)
 - + Metadata Registries (QA, harmonization, consistency, versioning)
 - + Metadata Implementations (Data Integration, Linked Data)

Discussion on 11179



- + Is ISO 11179 the Swiss Army Knife?
 - + Sophisticated data model
 - + More expressive than ODM, Archetypes , CDA
- + Some limitations for our use case:
 - + Missing features for clinical DM: Order of Data Elements or Value Meanings, repeated occurrences, single choice domains, default values, null values, mandatory fields, cross field checks
 - + No classes for modeling document hierarchies or groups of data elements belonging together
 - + No composite data elements

Special Challenges for a Community-based Approach

- + Users must be able to enter arbitrary data, else the MDR won't attract them
 - + What happens to redundant data (**duplicates**)? Which options exist for curating underspecified data elements?
 - + Which **user rights**, roles and views are needed and appropriate ?
 - + How could **modifications** be tracked and visualized? What implications arise from moving or **deleting metadata items** that are interconnected?
 - + How can **harmonization** be supported (reviewed data elements, core data sets)?

Metrics for Excellence

- + Quality of the specification
 - + Level detail (optional attributes)
 - + Consistency (property -> dimensionality -> units)
 - + Update frequency (especially if more than one user is involved)
- + Rating manually by the creator or the community
 - + Adjust to the expertise of the rater
 - + Consensus through community voting
- + Frequency of Use (for instance in other research projects) – the “common”
 - + Adjust to the importance of that project (locally, nationwide, number of subjects)
- + Reference to standards:
 - + Medical terminologies: ICD, OPS, LOINC, SDTM, SNOMED CT
 - + Artifact standards: UCUM measurement units, Null Flavors, ISO 21090 datatypes
 - + Contained in Core Data Sets: NINDS CDE, HL7 Value Sets, UK Biobank, ...
 - + Contained in validated instruments: assessments, scales, scores)

Metrics for Similarity

- + What will “Equivalence” or “Similarity” mean with regard to metadata items?
 - + Trivial approach: items are equal if their parts are equal
 - + Alternatives: items that are conceptually similar
 - + Variants: items with different representational values
 - + Derivations: items derived by some rule
 - + Versions: chronological view on the item’s track record
 - + Mappings: transformations between data elements
- + Most wanted: an ontology for Data Element Concepts
 - + Object Classes and Properties as well
- + Designations have no influence on similarity

Interfaces for Search and Personalization

- + Before one can decide on quality, we need a list of data element candidates
- + Currently, that means a textual search in designations, definitions and other text fields
 - + Solves morphological problems
 - + Problem of synonyms and homonyms persists
 - + Import data elements may have misleading or even no designations
 - + Components of a data element can have very similar names (data element concept, conceptual domain)
- + Facetted search: refinements
 - + Metadata objects, kind of usage
 - + Research projects, institutions
 - + User Profiles, classifications
 - + And combination of these

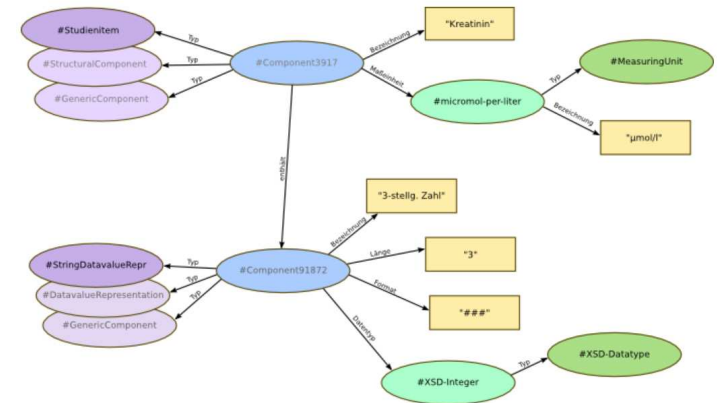
The screenshot shows the Metadata Repository interface. The browser address bar displays "http://mdr.tmf-ev.de". The page title is "Metadata Repository". The navigation menu includes "Itemkorb", "Profil", "Hilfe", "Impressum", and "Log out". The main content area shows a search for "Kreatinin" under the path "Home > Laborwerte". A dropdown menu is open, showing "Alle Kontexte (MDR)" with options: "IFB Sepsis", "IFB Adipositas", "ZKS Leipzig (KIS)", and "...". The search results table is as follows:

| Bezeichnung | Einheit | Code | |
|---------------------|-------------------|------|-------------------------------------|
| Kreatinin | mg/dl | | |
| Kreatinin im Serum | $\mu\text{mol/l}$ | | <input type="checkbox"/> |
| Kreatinin (min/24h) | mg/dl | | <input type="checkbox"/> |
| Kreatinin | $\mu\text{mol/l}$ | | <input checked="" type="checkbox"/> |

Additional interface elements include "Typen:" with checkboxes for "alle", "Datenelemente", "Datenelementkonzepte", "Konzeptuelle Domänen", "Wertedomänen - metrisch", and "Wertedomänen - Codeliste"; "Sprachvarianten:" with checkboxes for "alle Sprachen", "deutsch", and "englisch"; and a "Suchen" button.

Semantic Web Representation

- + Core Data Sets should be part of Linked Data cloud
- + We should provide a RDF serialization
- + SPARQL endpoint for querying
- + Use of Domain Ontologies: OCRE, OBI
- + Use of standard vocabularies:
 - + FOAF, SKOS , SIOC, SWAN
 - + Dublin Core (DC, DCE, DCT)
 - + Data Catalog Vocabulary (DCAT)
 - + Provenance, Geo, People, Org, Relations



Thank you!