

Praktikum - CRF Mapping Tool

Victor Christen
mam08bfa@studserv.uni-leipzig.de

16. September 2013

Inhaltsverzeichnis

1	Einleitung	2
2	CRF Excel Spezifikation	2
3	Realisierung	5
3.1	Datenschema	5
3.2	Allgemeiner Ablauf	6
3.3	Diff-Berechnung	6
4	Funktionsweise	8

1 Einleitung

Die Resultate einer Studie basieren auf der Auswertung der in jeder Phase erfassten Daten bzgl. eines Patienten oder Probanden. Diese Daten werden mittels eines CRF (Check Report Form) Bogen erfasst. Diesbezüglich enthält ein solcher Bogen eine Menge von Datenfeldern (Items), die für die spätere Auswertung relevant sind. Zu Beginn einer jeden Studie, werden deshalb für jede Phase einer Studie die relevanten Items inklusive ihrer Eigenschaften, Regeln und inhaltlichen Strukturierung ermittelt. Die Menge an CRF Bögen für die verschiedenen Studien können in einem klinischen Datenmanagementsystem verwaltet werden, wie z.B. OpenClinica. Im Allgemeinen ist die Entwicklung eines CRF Bogens ein fortlaufender Prozeß, da im Verlauf der Studie fehlende oder unzureichend definierte Items identifiziert werden. Diese müssen in einer aktualisierten Version ergänzt bzw. geändert werden. Sei nun folgendes Szenario gegeben: Es wurden mittels eines CRF Bogens die relevanten Daten für die Phase der Aufnahme eines Patienten erfasst. Es wurde festgestellt, dass die Auswertung der Daten umständlich ist, so dass eine Anpassung des CRF Bogens erfolgen muss. Dennoch sollen die erfassten Daten weiter verwendet werden. Hierfür ist ein Mapping, eine Abbildung der alten Items auf die neuen Items, notwendig, um die alten Daten in das Schema des neuen CRF Bogens zu transformieren.

Das Ziel dieser Arbeit ist die Entwicklung einer Applikation für die Generierung eines Mappings für eine Menge von CRF Versionen. Ein Mapping repräsentiert eine Abbildung von Items. Mithilfe eines Mappings ist eine effiziente Identifikation von gelöschten, hinzugefügten und geänderten Items erkennbar. Die Identifikation der Änderungen ist notwendig, um die Daten entsprechend auf das aktuelle Schema der aktuellen Version zu abbilden. Eine Problematik bzgl. der Änderung von Items ist die Relevanz der Änderung. Die Änderungen von Eigenschaften eines Items, die sich auf das Layout bzgl. des CRF Bogens beschränken, verursachen keine Änderung des medizinischen Konzepts, sodass existierende Daten auf das neue Item gemappt werden können. Die Transformation der Daten eines Items, deren medizinisches Konzept sich geändert hat, ist nicht ohne weiteres möglich.

2 CRF Excel Spezifikation

Der Import eines CRF Bogens in das OpenClinica System erfolgt durch eine Excel Datei, die die enthaltenden Datenfelder mit ihren Eigenschaften, Abschnitten und potentiellen Regeln für die Items in Form von einzelnen Tabellen enthält. Das Schema der Excel Mappe entspricht dem Klassendiagramm in Abbildung 1.

Die Klassen entsprechen den einzelnen Tabellen und die Attribute einer Klasse repräsentieren die Spalten der Tabelle. Die **CRF** Tabelle beinhaltet Informationen bzgl. des Namens des CRF Bogens, der Version, einer Beschreibung sowie Bemerkungen bzgl. der Erstellung.

Die **Sections** Tabelle definiert die layoutspezifische Struktur der Items des CRF's. Items, die das gleiche **section_label** Attribut aufweisen, werden auf der gleichen Seite aufgeführt. Jede Section wird durch das **section_label** Attribut identifiziert. Visuell wird jeder Abschnitt durch eine Überschrift getrennt, die durch den **section_title** spezifiziert wird.

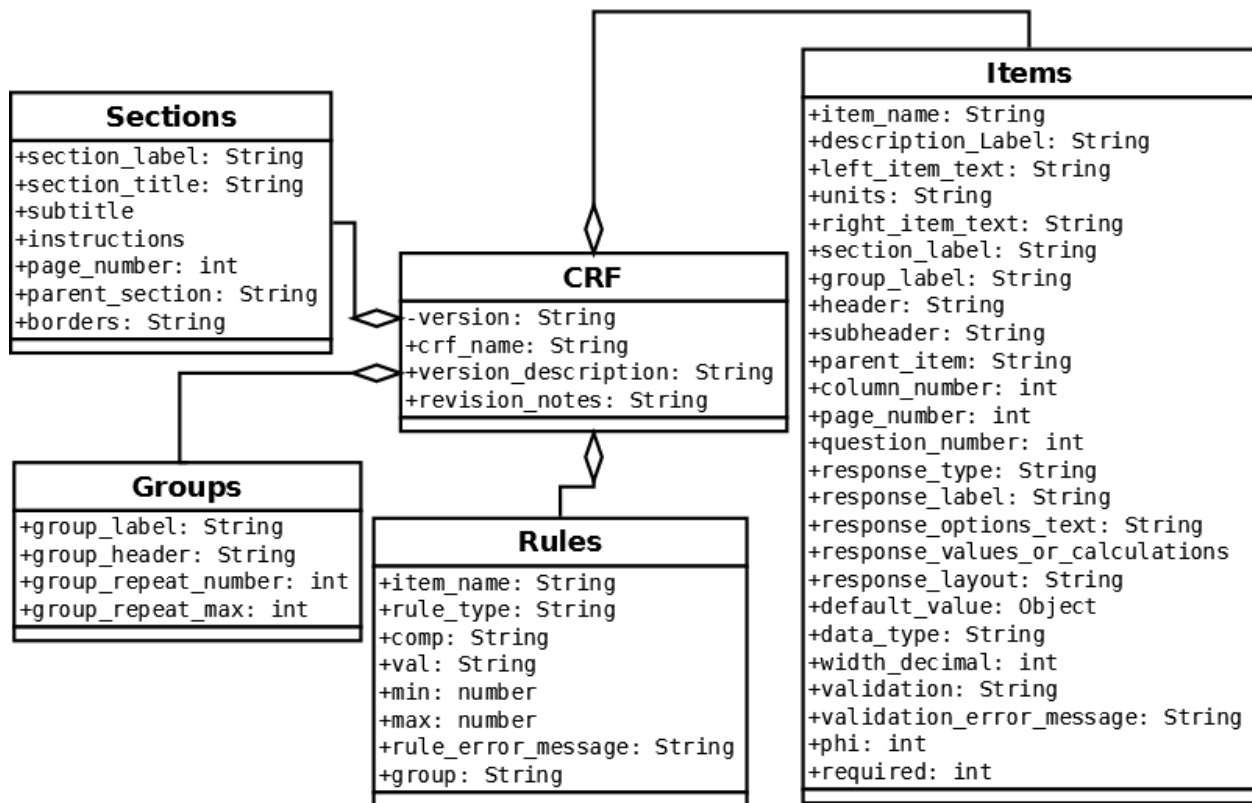


Abbildung 1: Datenschema der Excel Arbeitsmappe für einen CRF Bogen

Die **Groups** Tabelle ermöglicht die Gruppierung von Items. Items, die das gleiche **group_label** besitzen, erscheinen zusammen im CRF. Die Daten, die für die Auswertung notwendig sind, werden durch die Items eines CRF Bogens definiert. Die Spezifikation eines Items beinhaltet die Eigenschaften die in der Tabelle 1 aufgelistet sind.

Spaltenname	Beschreibung
ITEM_NAME	Identifiziert für ein Item
DESCRIPTION_LABEL	Beschreibung
LEFT_ITEM_TEXT	Text, der im CRF auf der linken Seite dargestellt werden soll
UNITS	Definition des Typs der Antwort
RIGHT_ITEM_TEXT	Text, der im CRF auf der rechten Seite dargestellt wird
SECTION_LABEL	Zuordnung zu einer SECTION
GROUP_LABEL	Zuordnung zu einer Gruppe
HEADER	Überschrift für das Item
SUBHEADER	Unterüberschrift für das Item
PARENT_ITEM	vorangegangenes Item referenziert mittels ITEM_NAME
COLUMN_NUMBER	horizontale Anordnung des Items, Reihenfolge gem. Nr.
PAGE_NUMBER	Seite des Items bei Papierform des CRF's
QUESTION_NUMBER	Nr. der Frage, erscheint vor LEFT_ITEM_TEXT
RESPONSE_TYPE	Antworttyp basiert auf den Standard HTML Elementen - text, textarea, single-select, radio, multi-select, checkbox, calculation, group-calculation, file, instant-calculation
RESPONSE_LABEL	Label, das bei einmaliger Definition von anderen Items referenzierbar ist, welche die gleichen Antwortmgl. besitzen referenzierbar durch andere Items
RESPONSE_OPTIONS_TEXT	kommaseparierter Liste, die die mgl. Antwortwerte beinhaltet, die für den Einzutragenden interpretierbar sind
RESPONSE_VALUES_OR_CALCULATIONS	Wenn der RESPONSE_TYPE keine Berechnung ist, beinhaltet dieses Feld die Werte, die in der Datenbank gespeichert werden können. Bei einer Berechnung ist ein Ausdruck anzugeben, der Werte anderer Items des CRF's verwendet, um einen Wert zu ermitteln.
RESPONSE_LAYOUT	Ausrichtung der Antwortoptionen (Blank, Horizontal, Vertical)
DEFAULT_VALUE	Default Wert für das RESPONSE_OPTIONS_TEXT Feld
DATA_TYPE	Datentyp (ST, INT, REAL, DATE, PDATE, FILE)
WIDTH_DECIMAL	max. Anzahl der Zeichen und max. Anzahl der Dezimalstellen
VALIDATION	Validierungsausdruck
VALIDATION_ERROR_MESSAGE	Meldung bei Validierungsfehler
PHI	Signierung, ob Anzeige von geschützten Gesundheitsdaten (blank,0,1)
REQUIRED	Signierung, ob Wert vorhanden sein muss (blank,0,1)

Tabelle 1: Übersicht der Spalten der Items Tabelle

3 Realisierung

Die Applikation ist in Java implementiert und verwendet für das Auslesen der Excel Dateien die Apache POI ¹ Bibliothek.

Der Unterabschnitt 3.1 beschreibt das zugrundeliegende Datenschema der Applikation. Im Unterabschnitt 3.2 wird der prinzipielle Ablauf der Applikation beschrieben. Abschließend wird detailliert auf die Berechnung der Unterschiede bzgl. einer Menge von CRF Versionen eingegangen.

3.1 Datenschema

Um einer potentiellen Veränderung der Excel Datei Spezifikation entgegenzuwirken, wurde ein eigenes Datenschema für eine CRF Version konzipiert. Das Schema ist in Abbildung 2 dargestellt.

Die `CRFVersion` Klasse aggregiert die `Sections`, `Groups` und `Items` Tabelle in Form einer Hashmap mit dem jeweiligen `label` Attribut als Schlüssel und dem dazugehörigen Objekt der Klasse `Section`, `Group` bzw. `Item` als Wert. Die Klassen `Section`, `Group` und `Item` sind prinzipiell gleich aufgebaut. Als Identifier besitzen sie ein `(item|section|group)_label` Attribut und eine Hashmap `propertyMap`, die die Werte aller Spalten einer `Section`, einer `Group`, bzw. eines `Item` beinhaltet. Die Namen der Spalten fungieren als Schlüssel für die `propertyMap` und werden in der jeweiligen `(Item|Section|Group)Constants` Klasse gespeichert. Somit ist durch die lose Koppelung der Spalten zu einem Objekt mittels der `propertyMap` und der Konstanten, eine Veränderung der CRF Excelspezifikation trivial realisierbar.

¹<http://poi.apache.org/>

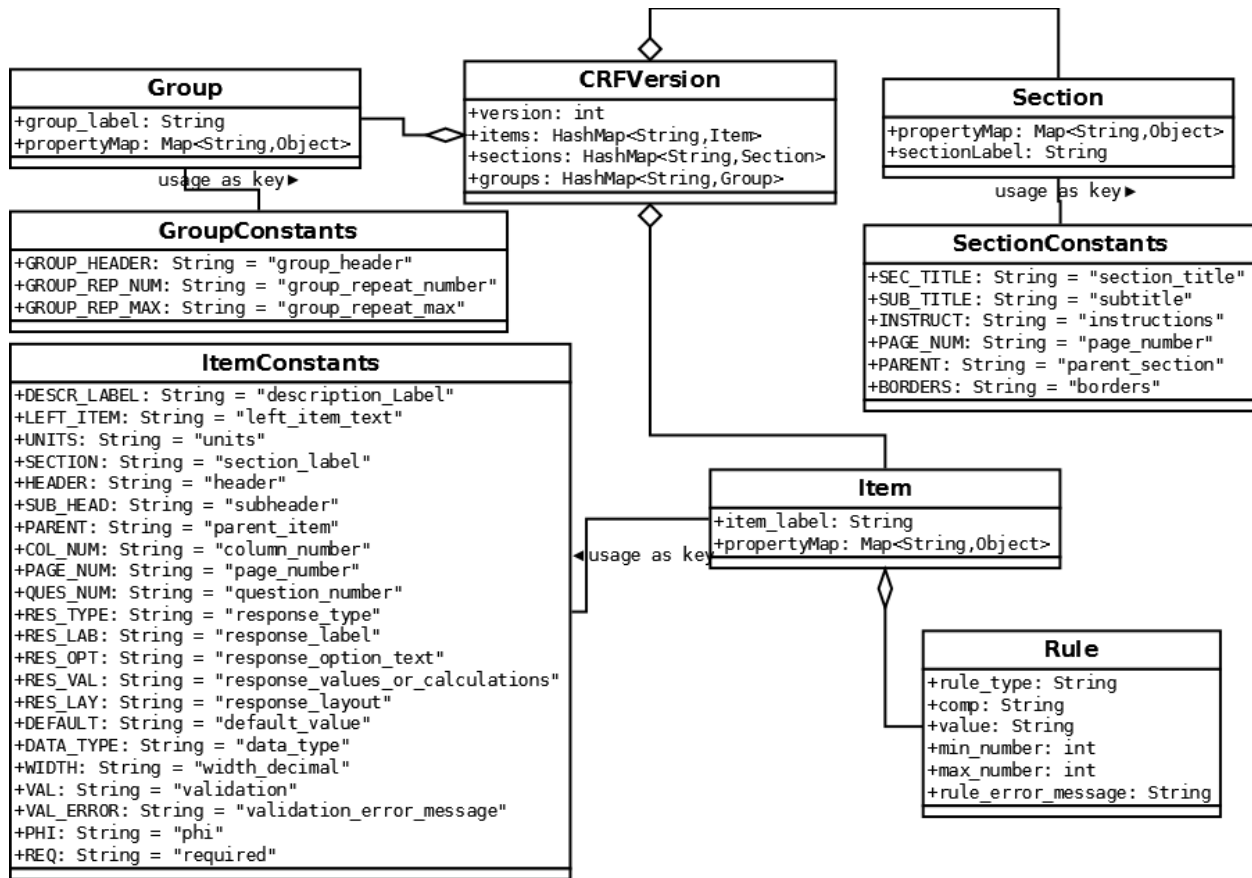


Abbildung 2: Klassendiagramm für eine CRFVersion

3.2 Allgemeiner Ablauf

Der allgemeine Ablauf der Applikationslogik ist in Abbildung 3 dargestellt. Zu Beginn spezifiziert der Anwender eine Menge von CRF Versionen. Die Dateien werden mittels der POI Bibliothek ausgelesen und zu Objekten des eigenen Datenschemas transformiert. Anschließend werden iterativ die Änderungen der CRF Versionen ermittelt und die Änderungen nach ihrer Relevanz klassifiziert. Bei jeder Iteration werden zwei aufeinander folgend Versionen miteinander verglichen und in der nächsten Iteration werden die alte zweite und die darauffolgende Version miteinander verglichen. Des Weiteren werden in jeder Iteration die Änderungen der Eigenschaften der Items bzgl. der Intensität ihres Ausmaßes klassifiziert. Die Änderungen werden in Form eines Baumes visualisiert, indem die aktuelle Version als Bezugspunkt für alle Änderungen der CRF Versionen dient.

3.3 Diff-Berechnung

Die Menge der Änderungsoperationen wird als Diff bezeichnet. Die Ermittlung des Diff's für eine Menge von CRF Versionen basiert auf dem fortlaufenden Vergleich von zwei CRF

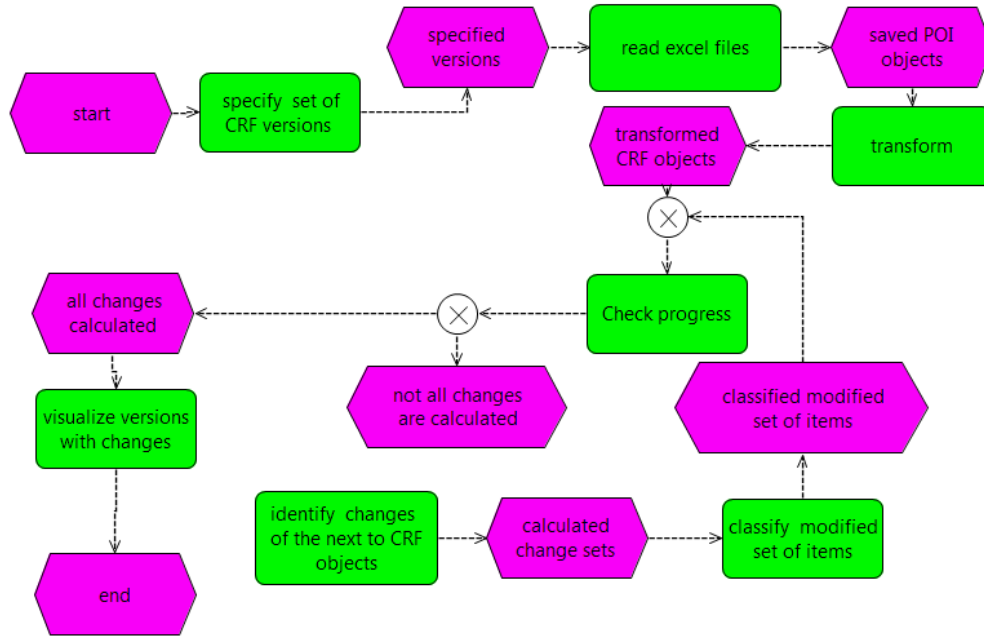


Abbildung 3:

Versionen. Im Folgenden wird das Verfahren für die Ermittlung näher beschrieben. Die Eingabe der Methode sind zwei CRF Versionen. Das Ergebnis ist eine Menge von gleichen, gelöschten, hinzugefügten und geänderten Items. Die geänderten Items werden als Mapping von altem Item zu neuem Item repräsentiert. Des Weiteren wird für jedes Mapping die Menge der geänderten Eigenschaften sowie der alte und der neue Wert gespeichert. Das Vorgehen ist im Algorithmus 1 dargestellt. Jedes Item der alten Version wird mit jedem Item der neuen Version mittels des `item.labels` verglichen. Wenn diese gleich sind, sind die Items gleich und ein Repräsentant wird zu der Menge *eqItemSet*, die die Menge der gleichen Items repräsentiert, hinzugefügt. Die Menge der gelöschten Items *delItemSet* und der hinzugefügten Items *addItemSet*, lässt sich durch die Mengendifferenz von den Mengen der Items der alten Version und der gleichen Items bzw. der Items der neuen Version und der gleichen Items ermitteln. Da geänderte Items sich durch einen anderen Suffix beim `item_label` unterscheiden, befinden sich den Mengen *delItemSet* und *addItemSet* eventuell geänderte Items. Um diese zu ermitteln, wird mithilfe eines Fuzzy- Matchers die Ähnlichkeit für jedes Item der *delItemSet* Menge zwischen jedem Item der *addItemSet* Menge ermittelt. Als geändertes Item wird das Paar angesehen, dessen Ähnlichkeit größer als ein Threshold $\delta(0.6)$ ist und das die höchste Ähnlichkeit aufweist. Anschließend wird für jedes Paar die Menge der geänderten Eigenschaften, die die Spalten repräsentieren, ermittelt.

Algorithm 1: Diff Berechnung für zwei CRF Versionen

```
input : oldVersion, newVersion
output: eqItemSet, addItemSet, delItemSet, modItemMap, propertyMap
1 for oldItem  $\in$  oldVersion do
2   for newItem  $\in$  newVersion do
3     if oldItem.item_label = newItem.item_label then
4       eqItemSet  $\leftarrow$  add(newItem)
5 delItemSet = oldVersion  $\setminus$  eqItemSet
6 addItemSet = newVersion  $\setminus$  eqItemSet
7 for addItem  $\in$  addItemSet do
8   topMapping  $\leftarrow$   $\emptyset$  for delItem  $\in$  delItemSet do
9     similarity = fuzzyMatch(addItem, delItem)
10    if similarity  $\geq$   $\delta$  then
11      topMapping  $\leftarrow$  add(similarity, (delItem, addItem))
12    if topMapping  $\neq$   $\emptyset$  then
13      (delItem, addItem) = bestMapping(topMapping)
14      modItemMap  $\leftarrow$  add(delItem, addItem)
15      addItemSet = addItemSet  $\setminus$  { addItem }
16      delItemSet = delItemSet  $\setminus$  { delItem }
17 for (oldItem, newItem)  $\in$  modItemMap do
18   oldPropertySet = oldItem.properties
19   newPropertySet = newItem.properties
20   modProperties = diff(oldPropertySet, newPropertySet)
21   propertyMap  $\leftarrow$  add(modProperties)
```

4 Funktionsweise

Die Applikation ist in drei Aspekte unterteilt, die Spezifikation der CRF Versionen, die Betrachtung der einzelnen CRF Versionen und die Betrachtung der Unterschiede der CRF Versionen im Bezug auf die aktuelle Version.

Die Spezifikation der CRF Version erfolgt im Menü File/Open. Für die Betrachtung der Unterschiede müssen mindestens zwei Versionen eines CRF Bogens selektiert werden.

Der Anwender ist in der Lage sich die einzelnen CRF Versionen einzeln zu betrachten (siehe Abb. 4), indem er eine Version in der ComboBox selektiert und hinzufügt. Die selektierte Version wird als Baum, deren Blätter die Items repräsentieren, auf der linken Seite der Applikation dargestellt. Bei der Selektion eines Items können die Eigenschaften im unteren Teil der Anwendung betrachtet werden.

Die Änderungen werden auf der rechten Seite des Programms (siehe Abb. 5). Die Items der aktuellen Version werden in einem Baum dargestellt. Des Weiteren werden alle Änderungen

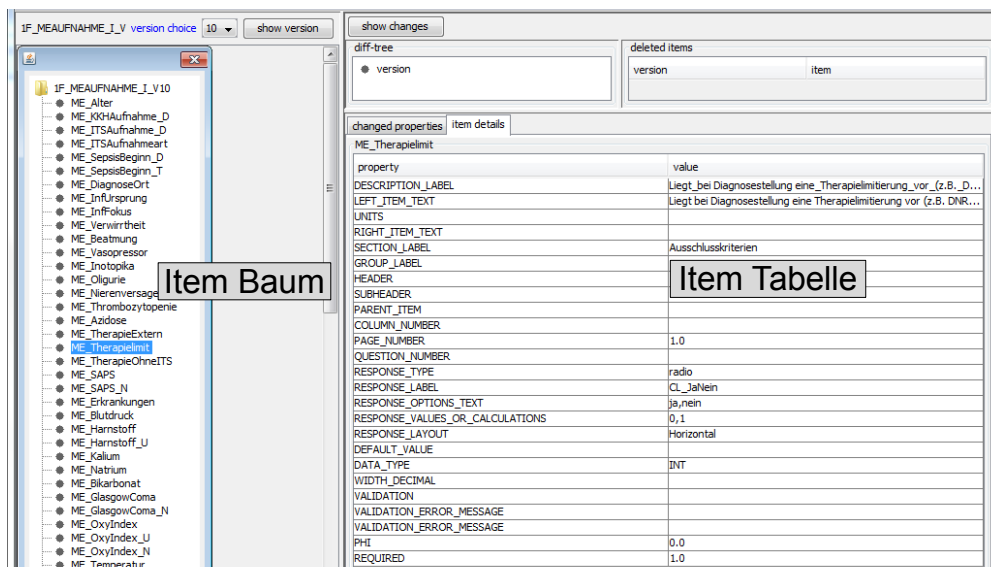


Abbildung 4: Darstellung der Items einer CRF Version in einem Baum inklusive der Ansicht der Eigenschaften eines Items in einer Tabelle

die ein Item betreffen chronologisch in einer Tabelle als Blattknoten des betreffenden Items dargestellt. Die Tabelle enthält die Versionsnummer von der vorherigen Version, den Namen des alten Items und die Art der Änderung. Bei der Selektion eines Eintrages der Tabelle, werden die geänderten Eigenschaften im unteren Teil in einer Tabelle präsentiert. Diese Tabelle enthält den Namen der geänderten Eigenschaft, den alten und neuen Wert und die Relevanz der Änderung. Des Weiteren werden die gelöschten Items in einer separaten Tabelle chronologisch nach ihrem Löschezitpunkt präsentiert, da diese in der aktuellen Version nicht vorhanden sind.

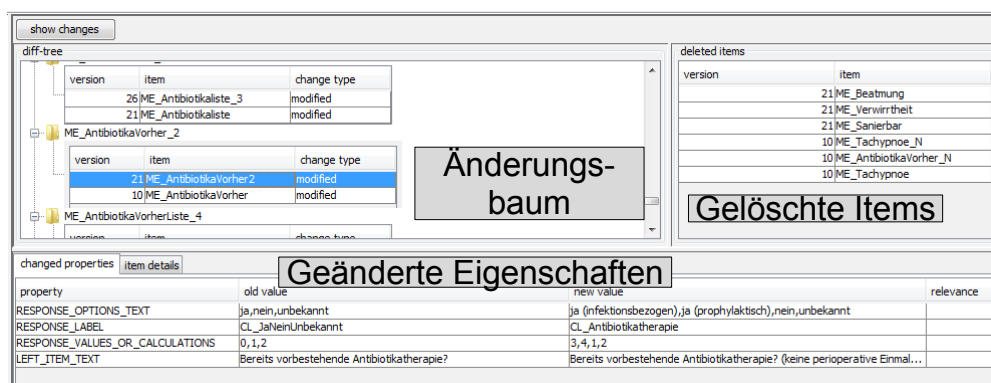


Abbildung 5: Darstellung der Änderung der geänderten Items und der geänderten Eigenschaften