

# Automatic extraction of breast cancer information from clinical reports

Claudia Bretschneider (Siemens), Matthias Hammon, Paul Gass (University Hospital Erlangen), Sonja Zillner (Siemens), Daniel Sonntag (DFKI)

## Background

The majority of clinical data is only available in unstructured text documents. Thus, their automated usage in data-based clinical application scenarios, like quality assurance and clinical decision support by treatment suggestions, is hindered because it requires high manual annotation efforts. In this work, we introduce a system for the automated processing of clinical reports of mamma carcinoma patients that allows for the automatic extraction and seamless processing of relevant textual features. Its underlying information extraction pipeline employs a rule-based grammar approach that is integrated with semantic technologies to determine the relevant information from the patient record.

## Our approach

We tested the system on our use case of mammography. We process an anonymised subset of the overall corpus; a de-identification tool was used. This corpus comprises **8,766 clinical texts reporting on 2,096 patients**, where the types of texts range from pathological texts (n=6,884) to operation reports (n=274) and radiology reports (n=1,608) over a time period of 15 years. **For the evaluation, we selected a high quality subset of 92 patients, for which the text records are complete in the corpus**, so that the final development corpus consisted of 6,932 reports and the evaluation corpus of 1,834 texts.

## Results

Using this integrated approach of information extraction, semantic modelling and rule-based decision support, we can extract the textual features with an accuracy of 0.69 for the most complex feature, the HER2 status, and up to 0.90 for the lymph node status. Using this information as input, it is possible to predict therapeutic measures with an accuracy of 0.59. Further error reduction for the IE results is planned to be integrated into the interactive faceted search application based on the IE results.

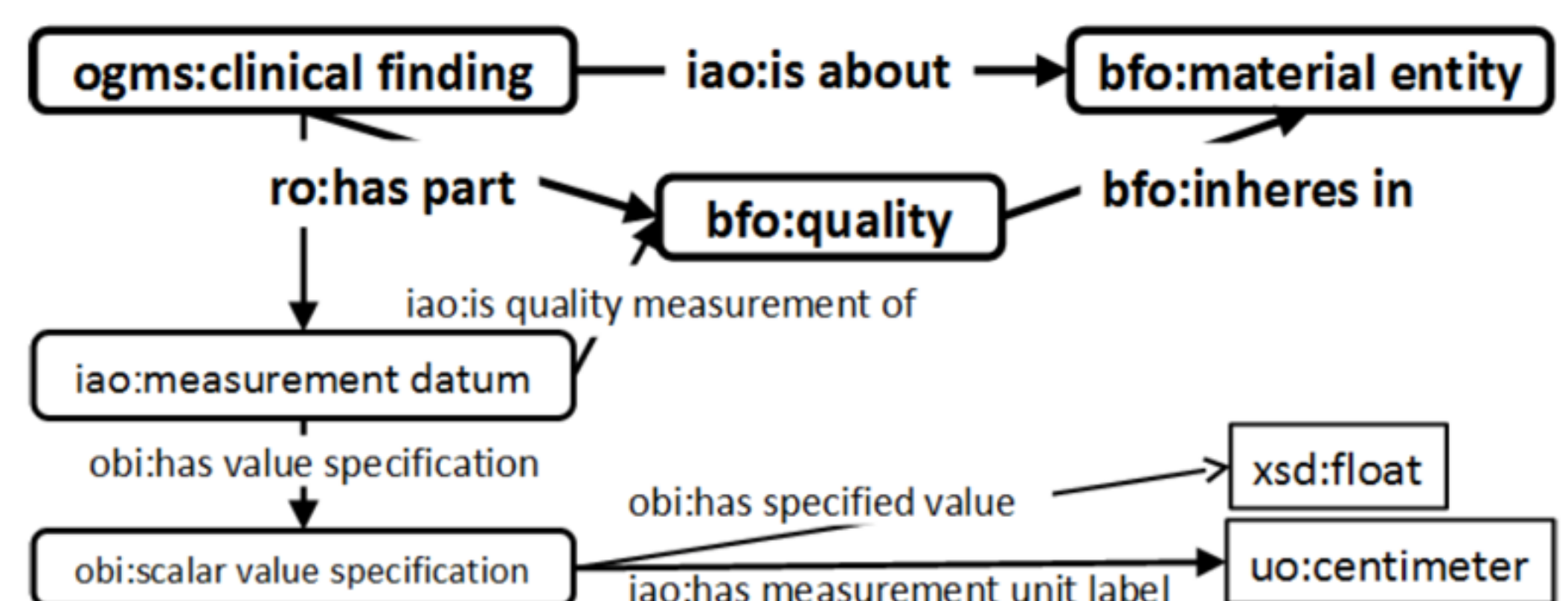


Figure 1. Adapted model for clinical information (MCI)

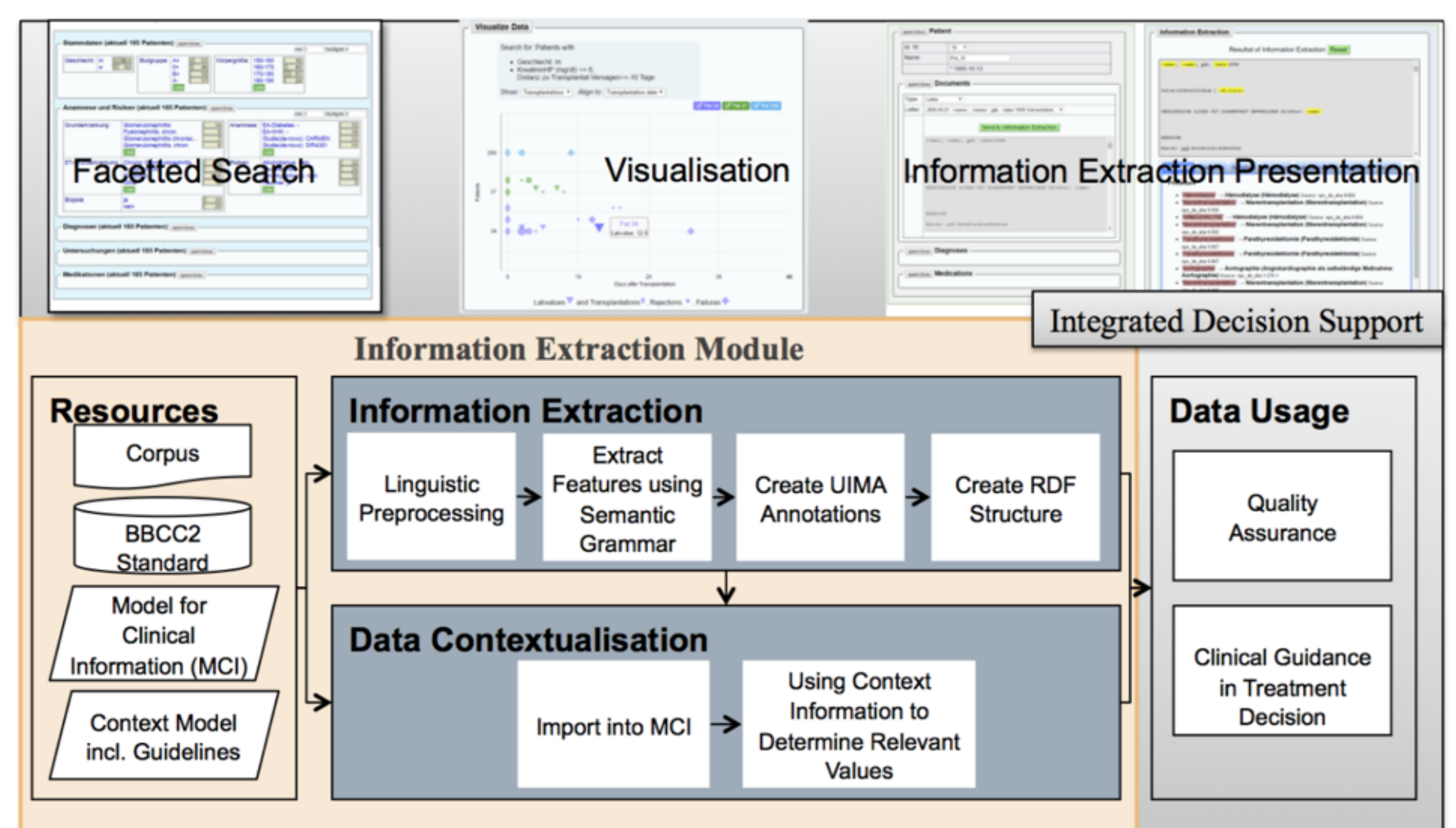


Figure 2. UIMA and SOLR integration for decision support

Table I

FEATURES TO BE EXTRACTED (UNDERLINED> FROM THE CLINICAL REPORTS AND TO AVOID FALSE POSITIVES (CROSSED OUT).

<p><b>(1) Type of operation conducted (OP)</b>            Freitext Therapie: <u>Ablatio</u> re, SNB (blau) ggf Axilla            Axilladissektion besprochen, da <del>XXXXXXXX</del> die primaer empfohlene <u>Ablatio mammae</u> strikt ablehnte.</p>
<p><b>(2) Tumor size (Size)</b>            Klinik: Mammakarzinom links, cT1 cN0, <u>1,7</u> cm.            I.: Maximal <u>1,2</u> cm großes mäßig differenziertes invasiv-duktales Mammakarzinom            Nach kaudal Abstand <del>0,3cm</del>.</p>
<p><b>(3) Grading of tumor (Grading)</b>            Anteile eines <u>mäßig differenzierten</u>, vorwiegend solide wachsenden Mammakarzinoms (linke Mamma, lt. Klinik).            Klinik: Mammakarzinom links, IDC <u>G2</u>.            Vorgeschichte: <del>XXXXXXXXXXXX</del>, eine <del>XXXXXXXXXXXX</del> Patientin (G1/P1)</p>
<p><b>(4) Lymph node status (LK)</b>            pT2 pN1mi (1/13) L0 V0 Pn0</p>
<p><b>(5) Hormone receptor state (Hormone)</b>            Östrogenhormonrezeptoren: &gt; 80 % (IRS 12/12)            Progesteronhormonrezeptoren: <u>negativ</u></p>
<p><b>(6) HER2 state (HER2)</b>            HER 2-Onkogen-Protein-Expression: 0 (<u>negativ</u>).            HER 2-Onkogen-Protein-Expression: <u>Score 1+</u>, somit <u>negativ</u>.            Ratio HER2/CEN17 = &gt; 5</p>
<p><b>(7) Lymphatic spread (Lymph)</b>            pT3 pN0 <u>L0</u> V0 Pn1 G2 R1            Kapsel <u>1+</u>, Kapsel 2+3 L2, Kapsel 4+5 L3, Kapsel 6+7 [...]            Lymphangiosis carcinomatosa</p>

