Towards Knowledge-Based Assistance for Scholarly Editing

Jana Kittelmann

MLU Halle-Wittenberg

Christoph Wernhard TU Dresden

AITP 2016 Obergurgl, 6 April 2016

Extended version of the talk slides, 19 April 2016

- 1. Scholarly Editing
- 2. Relevant Knowledge Sources
- 3. KBSET An Experimental Platform
- 4. Coupling Fuzzy and Symbolic Knowledge
- 5. Access Predicates
- 6. Conclusion

1. Scholarly Editing

- 2. Relevant Knowledge Sources
- 3. KBSET An Experimental Platform
- 4. Coupling Fuzzy and Symbolic Knowledge
- 5. Access Predicates
- 6. Conclusion

Scholarly Editing Scholarly Editing as Scientific Discipline

- Some other/related names/concepts:
 - Editionswissenschaft, Editionsphilologie, Editorik
 - Critique génétique
 - Textual criticism
- Emerged in the 1850s from reconstruction of ancient and medieval texts
- Outcome: critical edition
- Concerns
 - tracing and presenting text genesis
 - identifying a "definitive" version
 - presentation
 - bridging temporal and cultural distance to reader
 - "objective editions are not possible"

Scholarly Editing Summary Editions (Regestausgaben) of Correspondences

- Cases with **too much material to transcribe and present in full** Example: 20.000 letters to Goethe – successively published since the 1980s
- "Flat" forms of making accessible
 - involved persons
 - locations
 - dates
 - mentioned works
 - historic events
 - indexes

Scholarly Editing Separation of Descriptive and Procedural Markup: TEI

• Specification of XML elements and attributes for descriptive markup



Scholarly Editing TEI: Example



PDF-Download

S. 1 2 3 · 5

vorheriger Brief nächster Brief

Man pare vient de vacaren par le Innee Dalgovarhi une decorde requisée des Alles Alugganis, Tanenan, Cartetten pare al file. L'Annee Revis trompé Talgor O Cantodor Wind jamai élé de ce competend Main Cartella le pare cerit à non pare, que la al for file perferient des derols important et comm. Norten est pour des most van onles le la la faite de man Martin Color est pour

Faksimile Dipl. Umschrift Lesefassung Metadaten Entitäten XML

Mon père vient de recevoir par le Prince Dolgorouki une Seconde requète des Mssrs. Huggenin, Tournon, Castillon pere et fils. Le Prince s'étoit trompé Sulzer et Lambert n'ont jamais été de ce complot. Mais Castillon le père cert à mon pere, que lui et son fils possedoient des Secrets importans et comme Berlin est pour eux un vrai enfer il supplie qu'on leur procure ici de bonnes places etc. mon pere a envoié leur requete au comte d'Orlow mais je doute qu'on y fase réflexion. Le Prof. Louis

Faksimile Dipl. Umschrift Lesefassung Metadaten Entitäten XML

Scholarly Editing TEI: Remarks

- TEI P5 2.9.2 (2015) <correspDesc>
- TEI P5 (2007) Entity descriptions: <person>, <place>, <date>
- Stand-off markup with W3C XInclude

1. Scholarly Editing

2. Relevant Knowledge Sources

- 3. KBSET An Experimental Platform
- 4. Coupling Fuzzy and Symbolic Knowledge
- 5. Access Predicates
- 6. Conclusion

Relevant Knowledge Sources Wikipedia, Wikidata



Relevant Knowledge Sources Gemeinsame Normdatei ["Common Authority File"] (GND)

- Persons, organizations, works, ...
- 3 M persons, 120 M facts
- Ontology with 60 classes
 - Free (CC0) 0 0 0 DNB, Katalog der Deutsch... × + ☆ 白 🛛 🖈 💁 😐 🟴 ♠ (<) ●</p> https://portal.dnb.de/opac.htm?method=simpleSearch&colMode=true&ouery=idn%3D11879941 C + - Q Search 10 GB RDF LEIPZIG FRANKFURT AM MAIN English Kontakt A-Z Förderer Datenschutz Impressum Hilfe Mein Konto ↓ Katalog KATALOG DER DEUTSCHEN NATIONALBIBLIOTHEK → Einfache Suche Gesamter Bestand Musikarchiv Exilsammlungen Buchmuseum → Erweiterte Suche → Browsen (DDC) → Suchformular zurücksetzen → Suchverlauf idn=11879941X Finden → S Expertensuche ? → Meine Auswahl → Hilfe Datenshop Ergebnis der Suche nach: idn=11879941X Mein Konto Treffer 1 yon 1 Ablieferung von d Aktionen Netzoublikationen 📩 In meine Auswahl übernehmen → Informationsvermittlung A Druckansicht MARC21-XML-Repräsentation dieses Link zu diesem http://d-nb.info/gnd/11879941X Datensatzes Datensatz RDF/XML-Repräsentation dieses Login → Datensatzes Person Sulzer, Johann Georg Dokumentation Linked Data **Akademischer Grad** Professor Korrekturanfrage Über die Deutsche Nationalbibliothek Geschlecht männlich → Nachweis der Ouelle Andere Namen Sulzer, Jean George Sultzer, Johann Georg → Zugehöriger Artikel in Wikipedia Sulzer, Johann George Sulzer, Johann George F Teilen Sulzer, Johannes Georg

Relevant Knowledge Sources GND Example

Weitere Angaben	Professor der Mather Akademie der Wisse	athematik: 1775 Direktor der phil. Klasse der Visse 200				
Beziehungen zu Personen	Sulzer-Keusenhoff, V Keusenhof, Johann A	Link zu diesem Datensatz	http://d-nb.info/gnd/11879941X			
	Graff, Elisabetha Sor	Person	Sulzer, Johann Georg			
Systematik	4.7p Personen zu Ph	Akademischer Grad	Professor			
Тур	Person (piz)	Geschlecht	männlich			
Autor von	 Publikationen Kurzer Begriff Sulzer, Johani Versuch von c Sulzer, Johani der Ausg. Zür von Olaf Breici 	Andere Namen	Sulzer, Jean George Sulzer, Johann George Sulzer, Johanne George Sulzer, Johannes Georg Sulzer, Johannes Georg Sulzer, Johanne G. Sulzer, Sulzer, Sulzer,			
Beteiligt an	5 Publikationen	Quelle	M; B 1996; ADB 37 (1894), S. 144-147 M; B 1996			
	Winterthur : S 2. Allgemeine Th Ressource] Berlin : Direct 3	Zeit	Lebensdaten: 1720-1779			
		Land	Schweiz (XA-CH); Deutschland (XA-DE)			
		Geografischer Bezug	Geburtsort: Winterthur Sterbeort: Berlin Wirkungsort: Magdeburg Wirkungsort: Berlin			
Thema in	1 Publikation 1. Johann Georg Berlin : Akad.	Beruf(e)	Philosoph Pädagoge Theologe Mathematiker Hochschullehrer			

Relevant Knowledge Sources GeoNames

- 2.8 M locations, 10 M names
- Free (CC-BY)
- Table format

moabit all countries + search show on map [advanced search] search show on map [advanced search] variable Country Feature class Latitude Longitude Moabit Germany, Berlin- Moabit,Maabit,	● ○ ○ ③ CeoNames Fulltextsearch × + e ^R ↑ ④ ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● ●							
Name Country Feature class Latitude Longitude 1 Moabit Image: Section of Berlin-Moabit, Maabit, Moabit, Maabit, Maabit, Moabit, Maabit, Moabit, Maabit, Moabit, Maabit, Moabit, Maabit, Maa		moabit	arch show on m	all countries	2 records	found for "moabit"		
1 Moabit Sertin- Berlin- Moabit, Maabit, Maabit, Maabit, Maabit, Maabit, Maabit, Maabit, Moabit, Maabit,		Name	Country	Feature class	Latitude	Longitude		
2 Alt-Moabit Germany, Berlin road N 52° 31' 28'' E 13° 20' 58''	1 P	<mark>Moabit</mark> ම Berlin- Moabit,Maabit,Moabit,Maa6iт,Moa6ит	Germany, Berlin > Berlin, Stadt > Berlin	section of populated place population 70,911	N 52° 31' 34''	E 13° 20' 20''		
	2 ®	<u>Alt-Moabit</u>	<u>Germany</u> , Berlin	road	N 52° 31' 28"	E 13° 20' 58"		

Relevant Knowledge Sources YAGO, DBPedia

- Combined fact bases from Wikipedia, GeoNames, ...
- Developed in computer science
- 5-10 M Objects, 100-3000 M facts
- 700-350.000 classes, based on Wikipedia and WordNet
- Mulit-lingual
- Free licenses
- RDF

- 1. Scholarly Editing
- 2. Relevant Knowledge Sources

3. KBSET – An Experimental Platform

- 4. Coupling Fuzzy and Symbolic Knowledge
- 5. Access Predicates
- 6. Conclusion

KBSET: Introduction Addressed Issues in Scholarly Editing

- Incorporation of automated techniques, e.g.
 - named entity identification
 - statistics-based methods for analysis
- Providing explicit relationship to
 - external knowledge bases
 - formal semantics
- High-quality presentations
 - without expensive transformations and stylesheets
- Loose coupling of object text and markup
 - markup by different authors
 - automatically generated markup

KBSET: Introduction Some AI Aspects Reflected in Scholarly Editing

ΑI

- General background knowledge
- Position of the agent in the environment
- Temporal order
- Incompletely sensed/understood environment
- Coming to decisions about actions to take

SE

- GND, GeoNames
- Position in the text
- Order of word occurrences
- Incompletely understood text
- Coming to decisions about denotations of phrases, about annotations to insert

KBSET: Introduction The KBSET System

- "Knowledge-Based Support for Scholarly Editing and Text Processing"
- Free software: GNU Public License
- With comprehensive **example** (draft) Max Stirner: *Geschichte der Reaction*, Vol. 1, 1852

KBSET: Introduction Guiding Principles

- All phases of editing should be supported
 - 1) Creating the extended object text
 - 2) Generating **intermediate representations** for examination by humans or machines
 - 3) Generating final presentations
- High quality is required for all phases, e.g.
 - good tools for text creation
 - precisely identified persons
 - professional layout
- Consequences:
 - incorporation of special techniques and special systems
 - automated techniques, adjustable by humans

KBSET: Introduction Overview



KBSET: Inputs Processing of Inputs



KBSET: Inputs Embedding into Emacs



KBSET: Inputs System Perspective on Knowledge Bases

- KBSET is implemented in SWI-Prolog
- ... with theorem provers in mind, but currently making substantial use of
 - set abstraction (findall, setof)
 - sorting by term order
 - indexing on first argument
- Preprocessing for efficient access
 - extracting relevant data
 - GND: persons born before 1850 420 k instead of 3 M
 - indexed access predicates

KBSET: Inputs System Perspective on Text Representation

- Sequence of units: word | space | punctuation | command
 - allow to associate information, e.g. about identified entities
 - mapping to/from sequence of characters

KBSET: Entity Identification Entity Identification



KBSET: Entity Identification Identification of Persons

emacs	
Aus den Briefen von Sulzer an Hirzel, 1743-1778 Magdeb. d. 24. Febr. 45 [] Es ist mir lieb, daß Sie in Berlin gewesen, wie gefällt's ihnen dort. Haben Sie keine Gelehrte gesprochen, als Gleim und Spaldi [] Ich bin mit vielen Stüken unter Lan Gedichten, insonderheit aber [] der Od Hrn. Heß gar nicht zufrieden. -U: demo_03.txt Top (6,25) Git-n Gleim, Johann Wilhelm Ludwig (1719-1803) Lehrer, Schriftsteller, Sekretär http://de.wikipedia.org/wiki/Joha http://d-nb.info/gnd/118717758 Not explicitly blocked • Order of candidates	points ,)
No stop word No common noun Not common location name No common location name No common first name Linked to person in context: Sulzer, Johann Georg (1720–1 The preferred name is in the text Linked to more than 50 others Born 1719, matching context year 1760 In the German Wikipedia Linked to 76 others Gleim, Mathias Leberecht Caspar (1725–1783) Oberamtmann http://d-nb.info/gnd/1029929505	779)
-U:**- *kbset-info* Top (1,0) (Fundamental)	F

"Assistance" is Required Here

● ○ ○ X emacs	
Aus den Briefen von Sulzer an Hirzel, 1743-1778	
Magdeb. d. 24. Febr. 45 [] Es ist mir lieb, daß Sie in Berlin gewesen, wie gefällt's ihnen dort. Haben Sie keine Gelehrte gesprochen, als Gleim und Spalding? [] Ich bin mit vielen Stüken unter Langens Gedichten, insonderheit aber [] der ode an Hrn. Heß gar nicht zufrieden. -U: demo_03.txt Top (7,38) Git-master (Text Lange, Joachim (1670-1744) Evangelischer Theologe, Grammatiker, Pietist	: Kbset Fill)
http://de.wikipedia.org/wiki/Joachim_Lange	 By default the wrong
http://d-nb.info/gnd/118569376	candidate is prioritized
No stop word	
No common noun Not commonly used in lowercase No common location name	
No common first name	
The preferred name is in the text Born 1670, matching context year 1760 In the German Wikipedia Linked to 27 others)
Lange, Karl Heinrich (1703–1753) Bibliothekar, Lyriker, Pfarrer, Schriftstelle http://de.wikipedia.org/wiki/Karl_Heinrich_La http://d-nb.info/gnd/11670358X Not explicitly blocked	er Inge
-U:**- *kbset-into* Top (1,0) (Fundamental)	Г

Entry in the Assistance Document

00	X emacs
File Edit Options Buffe	rs Tools Help
def_assistance	<pre>(demo_assistance, [entity(person, [name='Lange', date0fBirth='1711', profession0roccupation='Schriftsteller'], [near_word_in=['0de']])]).</pre>
-U: demo_ass	sistance_03.pl 94% (511,0) Git:master (Prolog)
Lange, Joachim Evange http:// http://	 Prolog syntax, re-loadable Label for grouping and activation of entries
	• Entry: entity(<i>Type</i> , <i>Identifier</i> , [<i>Context</i>])
	• Identifier must uniquely determine the entity
	■ w.r.t. the KB, without technical "ID"
	In the German Wikipedia Linked to 27 others
Lange, Karl He Bibliot http:// http://	inrich (1703-1753) hekar, Lyriker, Pfarrer, Schriftsteller /de.wikipedia.org/wiki/Karl_Heinrich_Lange /d-nb.info/gnd/11670358X Not explicitly blocked info* Top (1 0) (Eundamental)
v. Kosec	

Correction after Adaption by "Assistance"

● ● ● ● X emacs
File Edit Options Buffers Tools Kbset Help
Aus den Briefen von Sulzer an Hirzel, 1743-1778
Magdeb. d. 24. Febr. 45 [] Es ist mir lieb, daß Sie in Berlin gewesen, wie gefällt's ihnen dort. Haben Sie keine Gelehrte gesprochen, als Gleim und Spalding? [] Ich bin mit vielen Stüken unter Mangens Gedichten, insonderheit aber [] der Ode an Hrn. Heß gar nicht zufrieden. -U: demo_03.txt Top (7,38) Git-master (Text Kbset Fill)
Lange, Samuel Gotthold (1711-1781)
http://de.wikipedia.org/wiki/Samuel_Gotthold_Lange http://de.wikipedia.org/wiki/Samuel_Gotthold_Lange http://d-nb.info/gnd/119023784 Net_explicitly blocked
Explicitly specified in context near_word="Ode"
No common noun Not commonly year in lawarana
No common loca The right candidate is now prioritized as No common fir: The preferred "explicitly specified" Born 1711, ma In the German Wikipedia Linked to 5 others
Lange, Joachim (1670–1744) Evangelischer Theologe, Grammatiker, Pietist http://de.wikipedia.org/wiki/Joachim_Lange http://d-nb.info/gnd/118569376
-U:**- *kbset-info* Top (1,0) (Fundamental)

KBSET: Entity Identification Further Possibilities in Assistance Documents

- Supplementing
 - attribute values
 - entities
- Excluding words as entity designators

KBSET: Entity Identification Dates: Parsing and Defaulting

0	00			X emacs			
Fil	le Edit (Options Buffers To	ools Kbset Help	\frown			
A	us dei	n Briefen v	on Sulzer an Hirzel	1743-1778			
	lagdeb] rie ge elehr edich rn. He 1: 745-02	n Briefen w . (d. 24. Fel Es ist mir fâllt's ihn te gesproch Ich bin mit ten, insond GB gar nich demo_03.tx 2-24, Wedne ??45-02-2 1743 comb	on Sulzer an Hirze br. 45 lieb, daß Sie in I en dort. Haben Sie en, als Gleim und g vielen Stüken unte erheit aber [] of t zufrieden. t Top (3,10) sday 4 parsed ined	Berlin gewes e keine Spalding? er Langens ler Ode an Git-master	en, (Text Kbset	Fill)	
-0	<u>.</u> **_	*kbset-inf	o * All (1,0)	(Fundament	al)		

31

Detailed Information on Locations



KBSET: Entity Identification Associated with Occurrences of Words

- In contrast to n-grams (sequences) of words
- Local context is considered
 - preceding and succeeding words
 - already identified entities

KBSET: Entity Identification Comparison with a Popular Entity Recognizer

- Stanford Named Entity Recognizer
 - statistics-based machine learning [Finkel et al., 2005]
 - free, since 2006, here version 3.3.1 (Jan 2014)
 - no identification, just recognizing the entity type!

... in/O Berlin/I-LOC gewesen/0,/O wie/O
gefällt/0's/O ihnen/O dort/0./O Haben/O Sie/O keine/O
Gelehrte/O gesprochen/0,/O als/O Gleim/I-PER und/O
Spalding/I-PER ?/O ...

- KBSET Vanilla configuration
 - GND until year of birth 1850
 - context year 1789
 - word list includes old orthography

KBSET: Entity Identification Comparison with the Stanford Named Entity Recognizer

Recognized occurrences of person designators in Stirner, *Geschichte der Reaction*, Vol. 1, 1852



- ldentification incorrect
- Due to old orthography
- Not recognized by KBSET Assisted hard to identify or not in GND extract

Runtimes: KBSET 25 sec, SNER 20 sec incl. 10 sec classifier loading

KBSET: Document Combination **Document Combination**



KBSET: Document Combination **LATEX / PDF Output**

Auszugsweise Übersetzung von COMTE: Cours de Philosophie Positive, Bd. 5, S. 583–584.

Heinrich (VIII., England, König) (1491-1547) Rom Karl (V., Heiliges Römisches Reich, Kaiser) (1500-1558) Franz (L. Frankreich, König) (1494-1547)

Auszugsweise Übersetzung von ebd., Bd. 5, S. 585-589.

Österreich

Luther, Martin (1483-1546)

Ueber die Stellung des Königthums zur Revolution.

226

desselben bildete, auch das Seinige bei.

Diese Unterwerfung der geistlichen Gewalt unter die weltliche, diese Grundstörung des ganzen bisherigen Lebenssystems gehörte übrigens dem ganzen weltlichen Europa gemeinsam an, und als Heinrich VIII. von Rom sich lossagte, waren Carl V. und Franz I. nicht weniger emancipirt als er.

Die Entstehung der modernen Reaktion.

An dieser ersten revolutionären Umwandlung haben also die katholischen Völker nicht weniger Theil genommen als die protestantischen.

Automatically generated

- margin notes for entities
- indexes
- hyperlinks
 - within the document
 - to Wikipedia, GND, etc.

, in Oesterreich und im Grunde Zeiten Luthers e Herren, nicht rotestantischen

nr calvinischen le Folge, indem iterwerfung un-238] bis dahin, nun die einzige des drohenden on sehen mußte. m katholischen iche und naturllstand, sodann

den Kuckschritt bezweckte.

Es braucht übrigens nicht erst bewiesen zu werden, daß diese wachsende Opposition gegen den Fortschritt der menschlichen Entwickelung, fern davon, nur dem modernen, gallicanischen oder spanischen Katholicismus eigen zu sein, auf eine viel gründlichere

KBSET: Document Combination External Annotations (Stand-off Markup)

Auszugsweise Übersetzung von COMTE: Cours de Philosophie Positive, Bd. 5, S. 583–584.

Heinrich (VIII., England, König) (1491-1547)

Rom

Karl (V., Heiliges Römisches Reich, Kaiser) (1500-1558)

Franz (I., Frankreich, König) (1494-1547)

Diese Unterwerfung der geistlichen Gewalt unter die weltliche, diese Grundstörung des ganzen bisherigen Lebenssystems gehörte übrigens dem ganzen weltlichen Europa gemeinsam an, und als Heinrich VIII. von Rom sich lossagte, waren Carl V. und Franz I. nicht weniger emancipirt als er.

Die Entstehung der modernen Reaktion.

Auszugsweise Übersetzung von ebd., Bd. 5, S. 585-589.

An dieser ersten revolutionären Umwandlung haben also die Lathelischen Väller nicht menioer Theil genommen ols die protes-

X emacs Diese Unterwerfung der geistlichen Gewalt unter die weltliche, diese Grundstörung des ganzen bisherigen Lebenssystems gehörte übrigens dem ganzen weltlichen Europa gemeinsam an. und als Heinrich VIII, von Rom sich lossagte, waren Carl V. und Franz I. nicht weniger emancipirt als er. ₩ 410^ Rome, Charles-Ouint et Francois I \xsecnc{\xugi Die Entstehung der modernen Reaktion} % 411v que les peuples catholiques n'aient tout aussi An dieser ersten revolutionären Umwandlung haben also die katholischen Völker nicht weniger Theil genommen als die -U:--- reaction 01 16.tex 77% (??,0) Git:master (LaTeX Kbset Fill)-----\xabout{source}{txt='Diese Unterverfung der geistlichen Gewalt unter die'} {Auszugsweise Übersetzung ven \volcite{5{[S.~583--584]{comte:cours}.} \xabout{source}{txt='An dieser ersten revolutionären Umwandlung haben also'} {Auszugsweise Übersetzung von \volcite{5}[S.~585--589]{comte:cours}.} -U:--- annot_reaction_01.tex 83% (388,0) Git:master (LaTeX Fill)-----

KBSET: Document Composition Some Future Issues on Document Composition

- **Semantics-based** conditions to specify positions to be modified in the object text, e.g. "in the chapters about ..."
- Relating to concepts of aspect-oriented programming:

Position	Joint point
Set of positions	Pointcut
Specifier of a set of positions	Pointcut designator
Action to be performed at all positions in a set	Advice
Effecting execution of advices	Weaving

KBSET Further Implemented Functionality

- Persons characterized by function: "Bishop of Chartres"
- Consideration of **document structure**
- Keyword extraction

- 1. Scholarly Editing
- 2. Relevant Knowledge Sources
- 3. KBSET An Experimental Platform

4. Coupling Fuzzy and Symbolic Knowledge

- 5. Access Predicates
- 6. Conclusion

Coupling Fuzzy and Symbolic Knowledge Use of Features in the Named Entity Identification of KBSET

Gleim, Johann Wilhelm Ludwig (1719-1803) Lehrer, Schriftsteller, Sekretär http://de.wikipedia.org/wiki/Johann_Wilhelm_Ludwig_Gleim http://d-nb.info/gnd/118717758

Not explicitly blocked

No stop word No common noun Not commonly used in lowercase No common location name No common first name

Linked to person in context: Sulzer [...] The preferred name is in the text Linked to more than 50 others Born 1719, matching context year 1760 In the German Wikipedia Linked to 76 others

not_explicitly_blocked explicitly_specified_in_context(_) followed_by_matching_roman_number(_) preceded_by_matching_first_names(_) explicitly_specified preceded_by_matching_first_names_initials(_) followed_by_matching_extra_names(_) followed_by_matching_extension(_) occupation_mentioned(_) no_stopword no common substantive no_common_downcase no_common_geoname no common firstname already_identified_in_context linked_to_person_in_context(_) referenced_by_preferred_name linked_to_many_others(_) born_in_span_before_year_in_context(_, _) in_wikipedia_de linked_to_others(_)

Coupling Fuzzy and Symbolic Knowledge Simple Plausibility Vector Model Currently Used in KBSET

 $plausibility(denotes(WordOccurrence, Entity)) = \langle V_1, \dots, V_n \rangle \equiv$ $value(feature_1, WordOccurrence, Entity) = V_1$ $\land \dots$ $\land value(feature_n, WordOccurrence, Entity) = V_n$

- Vectors $\langle V_1, \ldots, V_n
 angle$ are compared lexically
- For given *WordOccurrence* entities are arranged in equivalence classes
 - if the first is a singleton, *WordOccurrence* is taken as "identified"
- Feature values can depend on
 - previous entity identifications
 - context of *WordOccurrence*
- Vectors $\langle V_1,\ldots,V_n
 angle$ also serve as justifications

- 1. Scholarly Editing
- 2. Relevant Knowledge Sources
- 3. KBSET An Experimental Platform
- 4. Coupling Fuzzy and Symbolic Knowledge

5. Access Predicates

6. Conclusion

Access Predicates

Knowledge Sources have to be Preprocessed for Applications

• Given are source facts

```
rdf_triple(p1, name, 'Sulzer').
rdf_triple(p1, year_of_birth, 1720).
```

• The knowledge is accessed from the application in "directed" ways

```
name_to_year_of_birth(+N, -Y) :-
    name_to_person(+N, -P),
    person_to_year_of_birth(+P, -Y).
```

• It seems useful to precompute indexed access predicates

name_to_person(+N, -P)
person_to_year_of_birth(+P, -Y)

Access Predicates Tasks to be Automated – with Provers

- Determine required access predicates from given queries
- Rewrite queries in terms of access predicates
- Rewrite to subqueries for different knowledge sources

Access Predicates Definability as Validity and Definientia as Craig Interpolants

• Some second-order entailments can be reduced to first-order entailments:

 $\exists p \, F[p] \models \forall q \, G[q] \quad \text{iff} \quad F[p'] \models G[q'], \qquad p',q' \, \, \text{fresh}$

• **Definability** of p within F can be expressed as follows:

There is a Hx s.th. $F \models \forall x \, px \leftrightarrow Hx$,p not in Hxiff There is a Ha s.th. $F \models pa \leftrightarrow Ha$,p not in Haiff There is a Ha s.th. $\exists p \, F \land pa \models Ha \models \neg(\exists p \, F \land \neg pa)$,p not in Ha

 $a \ {\rm is} \ {\rm fresh}$

- Craig interpolation allows construction of definientia Ha from proofs
- Generalizations:
 - $\hfill \ensuremath{\,\bullet\)}$ complex formulas instead of p
 - \blacksquare specific predicates and constants allowed in Ha
 - specific polarity of predicate occurrences in Ha (Lyndon interpolation)

Access Predicates Very Simple Example

 $\begin{array}{ll} accessor_spec \stackrel{\text{def}}{=} \\ (\forall pn \ b(p) \rightarrow (person_name(p,n) \leftrightarrow person_name_bf(p,n))) & \land \\ (\forall pn \ b(n) \rightarrow (person_name(p,n) \leftrightarrow person_name_fb(p,n))) & \land \\ (\forall pn \ person_name(p,n) \rightarrow b(p) \land b(n)). \end{array}$

 $rewrite_1 \stackrel{\text{\tiny def}}{=}$

 $definiens(person_name(p, n), \\ accessor_spec \land b(n), \\ [person_name_bf, person_name_fb]).$

 $rewrite_1$ expands into a valid implication

 $(\exists b person_name \ accessor_spec \land b(n) \land person_name(p,n)) \rightarrow \\ \neg (\exists b person_name \ accessor_spec \land b(n) \land \neg person_name(p,n)).$

Recall: $\exists p F \land pa \models Ha \models \neg (\exists p F \land \neg pa), p \text{ not in } Ha$

A Craig interpolant for $rewrite_1$ is

 $person_name_fb(p, n).$

Access Predicates Subqueries for Different Knowledge Sources

• Craig interpolation (inductive interpolation [Craig, 1957, Lemma 2]) can compute H_i , each with different restrictions on allowed predicates and constants such that

 $F \models \forall x \, Gx \leftrightarrow H_1 x \lor \ldots \lor H_n x$

Access Predicates An Example with two Dependent Atoms

 $\begin{array}{l} \textit{rewrite}_{2} & \stackrel{\text{def}}{=} \\ \textit{definiens}(\exists p \ \textit{person_name}(p, n) \land \textit{person_yob}(p, y), \\ & accessor_spec \land b(n), \\ & [person_name_bf, person_name_fb, person_yob_bf, person_yob_fb]). \end{array}$

A Craig interpolant of $rewrite_2$ is

 $\exists x \ person_yob_bf(x, y) \land person_name_fb(x, n).$

Access Predicates An Example with a Referential Constraint

 $person_spec \stackrel{\text{def}}{=} (\forall p \ person(p) \to b(p)) \land (\forall pn \ person_name(p, n) \to person(p)).$ $rewrite_3 \stackrel{\text{def}}{=} definiens(\exists p \ person_name(p, n),$

person_spec ∧ *accessor_spec*, [*person*, *person_name_bf*, *person_name_fb*]).

A Craig interpolant of *rewrite*₃ is

 $\exists x \ person_name_bf(x, n) \land person(x).$

Access Predicates Implementation Framework "ToyElim 2"

- Addressed Issues
 - construction of complex formalizations
 - machine evaluation of these
 - reproducible computational tasks as by-product
- Prolog-based system
- Supported core operations for first-order logic
 - proving
 - interpolant computation
 - second-order quantifier elimination
- Macros
 - formula labels
 - to specify e.g. definiens(Q,F,S), is_transitive(P)
- **LATEX** formula pretty printer
- Support for **brief syntax**: *px*
- TPTP and DIMACS import/export
 - interface to first-order provers and SAT solvers

Access Predicates Prover Used for Interpolation

- CM prover (1992,1997,2015): PTTP/SETHEO/PROTEIN/leanCoP-like
 - model elimination / connection method / clausal tableaux
- Extraction of first-order interpolants
 - variant of the Smullyan/Fitting method
 - no change of the core prover needed

Access Predicates

CM Prover: Performance on the CASC-25 (2015) FOF Problems

Without Ed	quality	With Equality	
Prover	Solved from 150	Prover Solved from	n 250
Vampire 2.6	144	Vampire 4.0	241
Vampire 4.0	139	Vampire 2.6	227
iProver	125	E	194
E	122	ET	184
ET	119	CVC4	146
CVC4	111	iProver	97
∪ CM lean, low-5, st	d, hd-2 102	Prover9	82
CM lean	94	ePrincess	78
CM low-5	93	leanCoP	74
iProverModulo	91	∪ CM std, lean, low-5, lem-hd	56
CM std	90	CM std	48
leanCoP	85	CM lean	46
CM hd-2	85	CM low-5	44
ePrincess	35	CM lem-hd	42
Prover9	29	iProverModulo	36
Muscadet	18	Geo-III	22
Geo-III	15	Muscadet	19

CM: 300 sec, 3GB, 3.50GHz, inputs from TPTP 6.3.0, with SWI Prolog CASC: 300 sec, 32GB, 2.40GHz, axiom preloading and analysis allowed

Access Predicates Clausal Tableau for the Very Simple Example



Access Predicates Related Works

- Relativized quantifiers to ensure "evaluability" [Marx, 2007], [Bárány et al., 2013], [Nash et al., 2010]
- Relativized quantifiers to associate binding patterns [Benedikt et al., 2014]
- Comparing interpolants w.r.t. query plan cost estimations [Toman and Weddell, 2011, Hudek et al., 2015]

Access Predicate On the ToDo List

- Rewriting target languages:
 - what are useful properties for evaluation by proving techniques?
 - what properties of interpolants can be ensured with specific calculi?
- Auxiliary access predicates with compound definitions
- **Global "selection conditions"** should be propagated e.g. persons born before 1850
- Some prover control seems useful:
 - preference of "smaller" proofs
 - preference by ordering on predicate names

- 1. Scholarly Editing
- 2. Relevant Knowledge Sources
- 3. KBSET An Experimental Platform
- 4. Coupling Fuzzy and Symbolic Knowledge
- 5. Access Predicates
- 6. Conclusion

Conclusion

Interesting Aspects from the Viewpoint of Scholarly Editing

- Inclusion of automated techniques like named entity identification
- Embedding and automated use of large external KBs like GND
- Combination of KBs with adjustments to achieve precise results

Focusing on the exceptions, where automated techniques fail

- Inclusion of external and generated markup
- High-quality presentations with low cost

Conclusion

Involved Languages/Logics and Methods for them to Develop Further

• Internal access language

- ordered solution sets
- support for justifications
- automated generation of access predicates ⇒ interpolation

• Assistance language to adjust entity identification

- focus on exceptions
 - \Rightarrow non-monotonic reasoning
- Assistance language to control document combination
 - specifying sets of text positions
 - specifying modifications to be performed at these
 - \Rightarrow similarites to aspect-oriented programming
- Semantics-based modularization by forgetting about subvocabularies
 ⇒ second-order quantifier elimination

References

[Benedikt et al., 2014] Benedikt, M., ten Cate, B., and Tsamoura, E. (2014).

Generating low-cost plans from proofs.

In PODS'14 – Proc. 33rd ACM SIGMOD-SIGACT-SIGART Symp. on Principles of Database Systems, pages 200–211. ACM.

[Bárány et al., 2013] Bárány, V., Benedikt, M., and ten Cate, B. (2013).

Rewriting guarded negation queries.

In Mathematical Foundations of Computer Science 2013 (MFCS 2013), volume 8087 of LNCS, pages 98–110. Springer.

[Craig, 1957] Craig, W. (1957).

Three uses of the Herbrand-Gentzen theorem in relating model theory and proof theory.

JSL, 22(3):269-285.

[Finkel et al., 2005] Finkel, J. R., Grenager, T., and Manning, C. (2005).

Incorporating non-local information into information extraction systems by Gibbs sampling.

In Proc. 43nd Ann. Meeting of the Association for Computational Linguistics (ACL 2005), pages 363–370. ACL.

[Hudek et al., 2015] Hudek, A., Toman, D., and Wedell, G. (2015).

On enumerating query plans using analytic tableau.

In TABLEAUX 2015, volume 9323 of LNCS (LNAI), pages 339-354. Springer.

[Marx, 2007] Marx, M. (2007).

Queries determined by views: Pack your views.

In PODS '07, pages 23-30. ACM.

[Nash et al., 2010] Nash, A., Segoufin, L., and Vianu, V. (2010).Views and queries: Determinacy and rewriting.*TODS*, 35(3).

[Toman and Weddell, 2011] Toman, D. and Weddell, G. (2011). Fundamentals of Physical Design and Query Compilation. Synthesis Lectures on Data Management. Morgan and Claypool.