

High-Precision Semantics Extraction in STEM

Jan Frederik Schaefer
Supervisor: Michael Kohlhase

FAU Erlangen-Nürnberg

CICM 2020 — Doctoral Program
remotely from Erlangen, Germany

My Situation

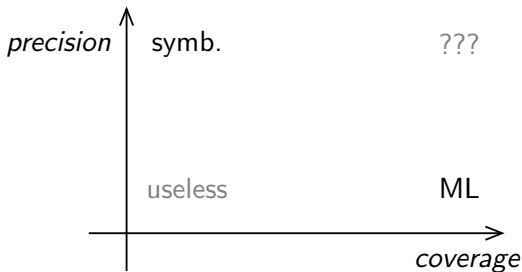
- Finishing up my master studies
 - Did already some research
 - Symbolic natural language semantics
 - Controlled natural languages for mathematics
 - Going to start my PhD soon
 - Supervisor: Michael Kohlhase *kwarc group*
 - Tentative topic: *high-precision semantics extraction in STEM*
 - Topic still very flexible
- **Any feedback appreciated!**

Motivation

- We have large corpora of STEM knowledge *e.g. arxiv*
 - Computers can make it more accessible:
 - Unit conversion
 - Applicable theorem search
 - Screen readers
 - ...
 - Such services require/benefit from semantic information
 - Authors often don't provide much *semantic T_EX macros*
- Semantics extraction

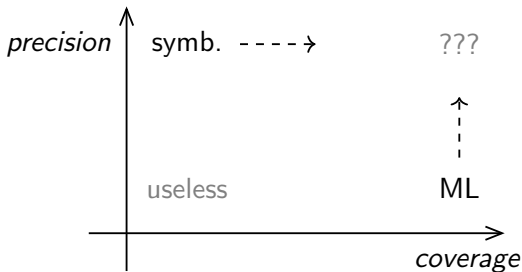
Approaches to Semantics Extraction

- Machine learning–based
 - Training data?
 - Low precision
- Symbolic
 - Low coverage



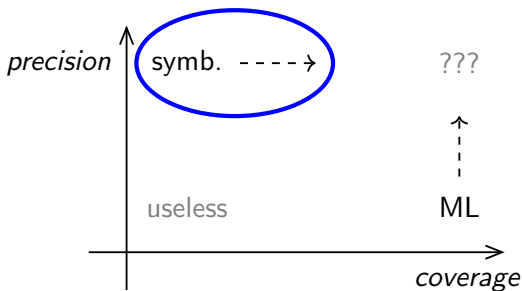
Approaches to Semantics Extraction

- Machine learning-based
 - Training data?
 - Low precision
- Symbolic
 - Low coverage



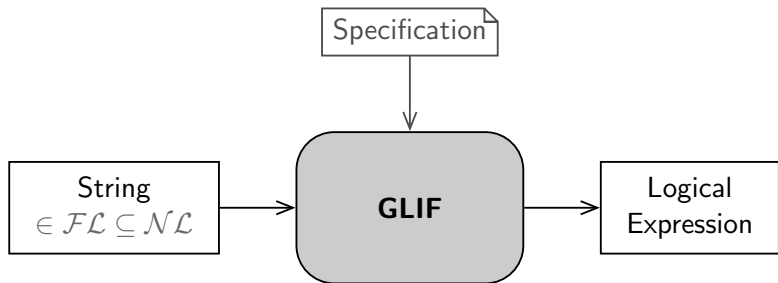
Approaches to Semantics Extraction

- Machine learning-based
 - Training data?
 - Low precision
- Symbolic
 - Low coverage



- Symbolic approaches offer high precision
- Often, high precision is more important than coverage
“each has the mass $\frac{1}{2}m$ ” vs *“each has the mass 1.64ft”*
- We already have a tool: GLIF
- Use statistical methods to increase coverage

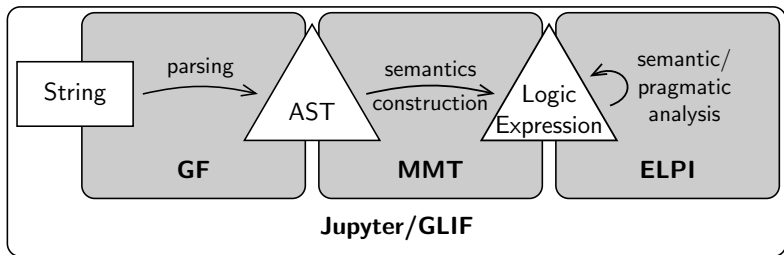
GLIF (Grammatical Logical Inference Framework)



Use cases:

- Designing controlled natural languages
- Prototyping approaches to natural-language semantics

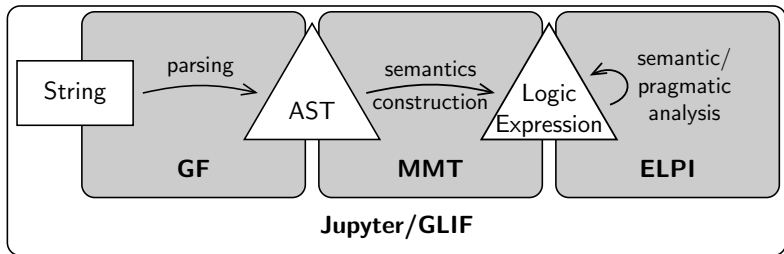
GLIF (Grammatical Logical Inference Framework)



Idea: Combine existing frameworks

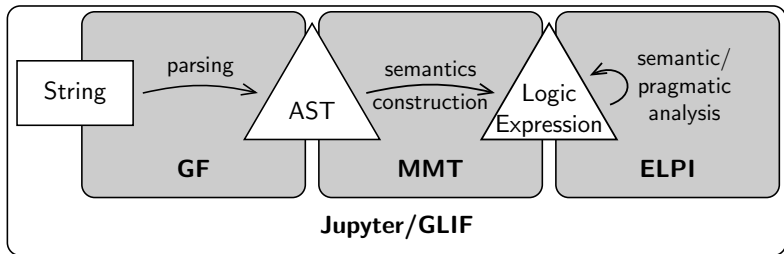
- GF (Grammatical Framework): \mathcal{NL} grammars
- MMT: Logic, knowledge representation
- ELPI (\supseteq λ Prolog): Inference
- Jupyter: Intuitive UI

GLIF (Grammatical Logical Inference Framework)



"... has a mass of 2m" \rightarrow $AST_1 \rightarrow \lambda x.mass(x, quant(2, \mathbf{meters}))$
 \rightarrow $AST_2 \rightarrow \lambda x.mass(x, mul(2, \mathbf{mVar}))$

GLIF (Grammatical Logical Inference Framework)



"... has a mass of 2m" \rightarrow $AST_1 \rightarrow \lambda x.mass(x, \text{quant}(2, \text{meters}))$

"... has a mass of 2m" \rightarrow $AST_2 \rightarrow \lambda x.mass(x, \text{mul}(2, \text{mVar}))$

"Therefore, A is clopen."

Different options:

① Use a dynamic parser

DynGenPar

② Generate lexicon automatically:

clopen_Adj : Adjective = "clopen"
clopen : $\iota \rightarrow o$

③ Replace lexicon entries with tokens:

"Therefore, A is ADJ-1."

Blanking out Unparsable Parts

This may be impossible:

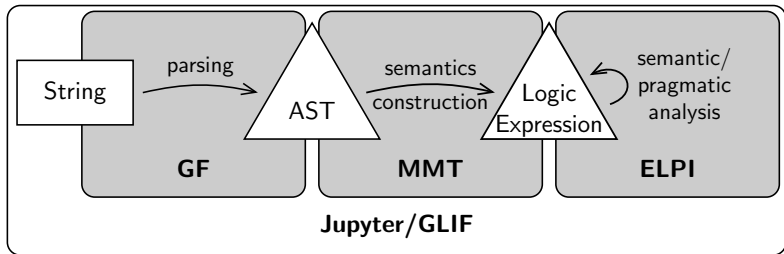
“Let X denote a data set where all entries $x \in X$ are normalized as described above.”

This would still be useful:

“Let X denote a data set where SUB-CLAUSE.”

Semantic Representation?

- Open question
- VIP, Naproche use DRT
- DRT not really supported by GLIF/MMT yet



Late Disambiguation

- ① Syntactic disambiguation:

"2m"

→ *unit?*

- ② Semantic disambiguation:

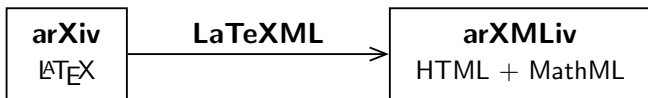
"a mass of 2m"

→ *not a unit*

- ③ Later semantic disambiguation:

*"it has a length of 2m, where m
is the length of a module"*

→ *not a unit*



- $> 10^6$ documents
- aims to preserve any semantic information from L^AT_EX sources
- have some experience with processing **arXMLiv** documents

<https://sigmathling.kwarc.info/>

Work Plan

- ① Prototype GLIF pipeline
 - Target: variable declarations and uses
 - Use generated lexicon
- ② Prototype pipeline for corpus work
 - Load document
 - Enter pre-processed sentences into GLIF pipeline
 - Export results
- ③ Introduce blanking out
- ④ Scaling
 - Larger grammar
 - More semantic phenomena
- ⑤ Build example semantic services
- ⑥ Can we replace more with ML?
Can the results be used as training data?

Discussion

- Is it desirable?
- Could this work?
- Other ideas?
- Anything else?

I think so

I haven't started yet → any feedback is welcome!