

Minimal module extraction from *DL-Lite* ontologies

Roman Kontchakov

School of Computer Science and Inf. Systems, Birkbeck College, London

<http://www.dcs.bbk.ac.uk/~roman>

joint work with

Luca Pulina, Frank Wolter and Michael Zakharyashev

Large-scale ontologies

- Life-sciences, healthcare, and other knowledge intensive areas depend on having a **common language** for gathering and sharing knowledge
- Such a common language is provided by **reference terminologies**
- Examples:
 - SNOMED CT (Systematized Nomenclature of Medicine Clinical Terms)
 - NCI (National Cancer Institute Ontology)
 - FMA (Foundational Model of Anatomy)
 - GALEN
 - ...
- Typical size: at least **50,000** terms and axioms
- Trend towards axiomatising reference terminologies in **(‘lightweight’) description logics**

Description logic \mathcal{ALCQT}

Vocabulary:

- individuals a_0, a_1, \dots
(e.g., john, mary)
- concept names A_0, A_1, \dots
(e.g., Person, Female)
- role names R_0, R_1, \dots
(e.g., hasChild, loves)
- roles

$$R ::= R_i \mid R_i^-$$

- concepts

$$C ::= A_i \mid \neg C \mid C_1 \sqcap C_2 \mid \exists R.C \mid \forall R.C \mid \geq qR.C$$

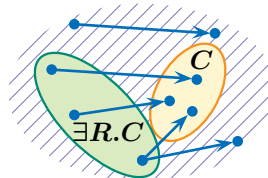
$\mathcal{I} = (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$ an **interpretation**

$$a_i^{\mathcal{I}} \in \Delta^{\mathcal{I}}$$

$$A_i^{\mathcal{I}} \subseteq \Delta^{\mathcal{I}}$$

$$R_i^{\mathcal{I}} \subseteq \Delta^{\mathcal{I}} \times \Delta^{\mathcal{I}}$$

$$(R_i^-)^{\mathcal{I}} = \{(y, x) \mid (x, y) \in R_i^{\mathcal{I}}\}$$



$$\forall R.C \quad \equiv \quad \neg(\exists R.\neg C)$$

'there are at least q distinct R -successors that are in C '

Description logic *ALCQI* (cont.)

knowledge base \mathcal{K} = TBox \mathcal{T} + ABox \mathcal{A}

- \mathcal{T} is a set of **terminological axioms** of the form $C \sqsubseteq D$
- \mathcal{A} is a set of **assertional axioms** of the form $C(a)$ and $R(a, b)$

Reasoning: – satisfiability \mathcal{K}

is there a *model* \mathcal{I} for \mathcal{K} ($\mathcal{I} \models C \sqsubseteq D$ iff $C^{\mathcal{I}} \subseteq D^{\mathcal{I}}$)

– subsumption $\mathcal{K} \models C \sqsubseteq D$

$\mathcal{I} \models C \sqsubseteq D$, for each \mathcal{I} with $\mathcal{I} \models \mathcal{K}$

– instance checking $\mathcal{K} \models C(a)$

$a^{\mathcal{I}} \in C^{\mathcal{I}}$, for each \mathcal{I} with $\mathcal{I} \models \mathcal{K}$

– query answering $\mathcal{K} \models q(\vec{a})$, $q(\vec{a})$ a positive existential formula

$\mathcal{I} \models q(\vec{a})$ (as a first-order structure), for each \mathcal{I} with $\mathcal{I} \models \mathcal{K}$

OWL 1.0 DL is based on *SHOIQ(D)*,

OWL 2.0 on *SROIQ(D)*

ALCQI + role inclusions + nominals + transitive roles + concrete domains

SHOIQ(D) + role chains + disjoint roles + self (diagonal)

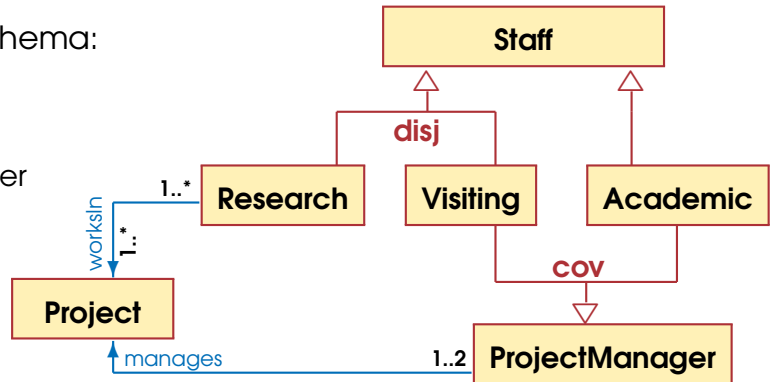
Developing and maintaining ontologies

- **versions:**
comparing **logical consequences** over some common vocabulary Σ
not the syntactic form of the axioms (as in `diff`)
- **refinement:**
adding new axioms but **preserving** the relationships
between terms of a certain part Σ of the vocabulary
- **reuse:**
importing an ontology and using its vocabulary Σ as originally defined
(relationships between terms of Σ should not change)
- **module extraction:**
computing a subset \mathcal{M} (ideally as small as possible) of an ontology \mathcal{T} that
'says' the same about Σ as \mathcal{T}

new types of reasoning problems

DL-Life: Description Logic for Databases

A fragment of a conceptual schema:



Translating into DL:

$\exists \text{manages.T} \sqsubseteq \text{ProjectManager}$

$\exists \text{manages}^-.T \sqsubseteq \text{Project}$

$\text{Project} \sqsubseteq \exists \text{manages}^-.T$

$\geq 3 \text{manages}^-.T \sqsubseteq \perp$

$\text{Research} \sqcap \text{Visiting} \sqsubseteq \perp$

$\text{Academic} \sqsubseteq \text{ProjectManager}$

$\text{ProjectManager} \sqsubseteq \text{Academic} \sqcup \text{Visiting}$

...

DL-Life_{bool}

$R ::= P \mid P^-$

$B ::= \perp \mid A_i \mid \exists R \mid \geq q R$

$C ::= B \mid \neg C \mid C_1 \sqcap C_2 \mid C_1 \sqcup C_2$

TBox axioms: $C_1 \sqsubseteq C_2$ **ABox assertions:** $C(a), R(b, c)$

Essentially positive existential queries: $\exists \vec{y} \varphi(\vec{x}, \vec{y})$, built from $C(t), R(t, t'), \wedge, \vee$

Σ -entailment and Σ -inseparability

Let \mathcal{T}_1 and \mathcal{T}_2 be TBoxes and Σ a **signature** (concept and role names)

When do \mathcal{T}_1 and \mathcal{T}_2 'say' the same about Σ ?

- \mathcal{T}_1 **Σ -concept entails** \mathcal{T}_2 if, for all Σ -concept inclusions $C \sqsubseteq D$,

$$\mathcal{T}_1 \preceq_{\Sigma}^c \mathcal{T}_2$$

$$\mathcal{T}_1 \models C \sqsubseteq D \text{ implies } \mathcal{T}_2 \models C \sqsubseteq D$$

- \mathcal{T}_1 **Σ -query entails** \mathcal{T}_2 if, for all Σ -queries $q(\vec{x})$ and ABoxes \mathcal{A} ,

$$\mathcal{T}_1 \preceq_{\Sigma}^q \mathcal{T}_2$$

$$(\mathcal{T}_1, \mathcal{A}) \models q(\vec{a}) \text{ implies } (\mathcal{T}_2, \mathcal{A}) \models q(\vec{a}), \text{ for all } \vec{a}$$

- ...

- \mathcal{T}_1 **Σ -model entails** \mathcal{T}_2 if, for all Σ -interpretations \mathcal{I} ,

$$\mathcal{T}_1 \preceq_{\Sigma}^m \mathcal{T}_2$$

$$\exists \mathcal{I}_1 \supseteq \mathcal{I} \mathcal{I}_1 \models \mathcal{T}_1 \text{ implies } \exists \mathcal{I}_2 \supseteq \mathcal{I} \mathcal{I}_2 \models \mathcal{T}_2$$

- \mathcal{T}_1 and \mathcal{T}_2 are S_{Σ} (**concept/query/model**) **inseparable** if

$$\mathcal{T}_1 \equiv_{\Sigma}^S \mathcal{T}_2$$

$$\mathcal{T}_1 \preceq_{\Sigma}^S \mathcal{T}_2 \text{ and } \mathcal{T}_2 \preceq_{\Sigma}^S \mathcal{T}_1$$

Σ -inseparability: Examples

Example 1. $\Sigma = \{\text{Lecturer, Course}\}$

$$\mathcal{T}_1 = \emptyset, \quad \mathcal{T}_2 = \{\text{Lecturer} \sqsubseteq \exists \text{teaches}, \exists \text{teaches}^- \sqsubseteq \text{Course}\}$$

- Is $\mathcal{T}_1 \equiv_{\Sigma}^c \mathcal{T}_2$?
- Is $\mathcal{T}_1 \equiv_{\Sigma}^q \mathcal{T}_2$?

Take $\mathcal{A} = \{\text{Lecturer}(a)\}$, $q = \exists y \text{Course}(y)$. Then $(\mathcal{T}_1, \mathcal{A}) \not\models q$ but $(\mathcal{T}_2, \mathcal{A}) \models q$

Example 2. $\Sigma = \{\text{Lecturer}\}$

$$\mathcal{T}_1 = \emptyset, \quad \mathcal{T}_2 = \{\text{Lecturer} \sqsubseteq \exists \text{teaches}, \text{Lecturer} \sqcap \exists \text{teaches}^- \sqsubseteq \perp\}$$

- Is $\mathcal{T}_1 \equiv_{\Sigma}^c \mathcal{T}_2$?
- Is $\mathcal{T}_1 \equiv_{\Sigma}^q \mathcal{T}_2$?

Take $\mathcal{A} = \{\text{Lecturer}(a)\}$, $q = \exists y \neg \text{Lecturer}(y)$. Then $(\mathcal{T}_1, \mathcal{A}) \not\models q$ and $(\mathcal{T}_2, \mathcal{A}) \models q$

Σ -inseparability: Examples (cont.)

Example 3. Let \mathcal{T}_1 contain the axioms

Research $\sqsubseteq \exists \text{worksIn}$,	$\exists \text{worksIn}^- \sqsubseteq \text{Project}$,
Project $\sqsubseteq \exists \text{manages}^-$,	$\exists \text{manages} \sqsubseteq \text{Academic} \sqcup \text{Visiting}$,
$\exists \text{teaches} \sqsubseteq \text{Academic} \sqcup \text{Research}$,	Academic $\sqsubseteq \exists \text{teaches} \sqcap \leq 1 \text{teaches}$,
Research $\sqcap \text{Visiting} \sqsubseteq \perp$,	$\exists \text{writes} \sqsubseteq \text{Academic} \sqcup \text{Research}$,

$\mathcal{T}_2 = \mathcal{T}_1 \cup \{\text{Visiting} \sqsubseteq \geq 2 \text{writes}\}$ and $\Sigma = \{\text{teaches}\}$

- $\mathcal{T}_1 \equiv_{\Sigma}^c \mathcal{T}_2$
 $\mathcal{T}_2 \models \text{Visiting} \sqsubseteq \text{Academic}$, but nothing new in the signature Σ

- $\mathcal{T}_1 \not\equiv_{\Sigma}^q \mathcal{T}_2$:
 $\mathcal{A} = \{\text{teaches}(a, b), \text{teaches}(a, c)\}$
 $q = \exists x ((\exists \text{teaches})(x) \wedge (\leq 1 \text{teaches})(x))$
'is there anybody who teaches precisely one module?'



$(\mathcal{T}_1, \mathcal{A}) \not\models q$

$(\mathcal{T}_2, \mathcal{A}) \models q$

Σ -entailment: semantic criteria

Let Q be a set of numerical parameters and Σ a signature

ΣQ -concepts B : $A_i \in \Sigma$ and $(\geq q R)$ with $q \in Q$ and $R \in \Sigma$

ΣQ -type \mathbf{t} is a set of ΣQ -concepts containing B or $\neg B$ (but not both), for all B

For a TBox \mathcal{T} ,

a ΣQ -type \mathbf{t} is **\mathcal{T} -realisable** if \mathbf{t} is satisfied in a model of \mathcal{T}

(i.e., there is a \mathcal{I} of \mathcal{T} and a point w in it such that $w \in B^{\mathcal{I}}$ iff $B \in \mathbf{t}$)

a set Ξ of ΣQ -types is **precisely \mathcal{T} -realisable** if

there is a model of \mathcal{T} realising precisely the types from Ξ

Theorem. Let Q denote the set of parameters occurring in $\mathcal{T}_1 \cup \mathcal{T}_2$

\mathcal{T}_1 **Σ -concept entails** \mathcal{T}_2 iff every \mathcal{T}_1 -realisable ΣQ -type is **\mathcal{T}_2 -realisable**

\mathcal{T}_1 **Σ -query entails** \mathcal{T}_2 iff every precisely \mathcal{T}_1 -realisable set Ξ of ΣQ -types is **precisely \mathcal{T}_2 -realisable**

Σ -inseparability: complexity

Theorem.

- Deciding Σ -concept and Σ -query inseparability is Π_2^P -complete
- Deciding Σ -model inseparability is NEXPTIME-complete
- Can be simpler for various fragments of $DL\text{-Lite}_{bool}$
E.g. deciding Σ -concept and Σ -query inseparability for $DL\text{-Lite}_{horn}$ is
coNP-complete

NB. Π_2^P -completeness means that the problem can be encoded as satisfiability of $\forall\exists$ quantified Boolean formulas

Various QBF solvers can be used to check Σ -concept and Σ -query inseparability

NB. Inseparability is much harder for \mathcal{ALC} and other non-'Lite' DLs
(2EXPTIME-complete for \mathcal{ALC} , undecidable for \mathcal{ALCQIO})

Encoding Σ -concept entailment in QBF

Let \mathcal{T} be a TBox, Q a set of numerical parameters and \mathbf{t} a $\mathbf{sig}(\mathcal{T})Q$ -type

$$\begin{array}{l} \text{'}\mathbf{t}_0 \text{ is } \mathcal{T}\text{-realisable with } \mathbf{t}_1, \dots, \mathbf{t}_n \text{ being witnesses'} \\ \text{propositional formula} \end{array} = \Phi_{\mathcal{T}}(b_0, b_1, \dots, b_n)$$

b_j is the vector of all propositional variables B^* of the type \mathbf{t}_j

Then the condition

'every \mathcal{T}_1 -realisable ΣQ -type \mathbf{t} is \mathcal{T}_2 -realisable'

is described by the following QBF

$$\forall b_0^{\Sigma Q} \left[\exists b_0^{\mathcal{T}_1 \setminus \Sigma Q} \exists b_1^{\mathcal{T}_1} \dots \exists b_{n_1}^{\mathcal{T}_1} \Phi_{\mathcal{T}_1}(b_0^{\Sigma Q} \cdot b_0^{\mathcal{T}_1 \setminus \Sigma Q}, b_1^{\mathcal{T}_1}, \dots, b_{n_1}^{\mathcal{T}_1}) \rightarrow \right. \\ \left. \exists b_0^{\mathcal{T}_2 \setminus \Sigma Q} \exists b_1^{\mathcal{T}_2} \dots \exists b_{n_2}^{\mathcal{T}_2} \Phi_{\mathcal{T}_2}(b_0^{\Sigma Q} \cdot b_0^{\mathcal{T}_2 \setminus \Sigma Q}, b_1^{\mathcal{T}_2}, \dots, b_{n_2}^{\mathcal{T}_2}) \right]$$

($b_0^{\Sigma Q}$ is the ΣQ -part of b_0 and $b_0^{\mathcal{T}_i \setminus \Sigma Q}$ contains the rest of the variables)

Experiments

TBox instances (standard Department Ontology + ICNARC)

series	description	no. of instances	axioms		basic concepts		
			\mathcal{T}_1	\mathcal{T}_2	\mathcal{T}_1	\mathcal{T}_2	Σ
NN	\mathcal{T}_1 does not Σ -concept entail \mathcal{T}_2	840	59–308	74–396	47–250	49–300	5–103
YN	\mathcal{T}_1 Σ -concept but not Σ -query entails \mathcal{T}_2	504	56–302	77–382	44–246	58–298	6–89
YY	\mathcal{T}_1 Σ -query entails \mathcal{T}_2	624	43–178	43–222	40–158	40–188	5–64

QBF solvers

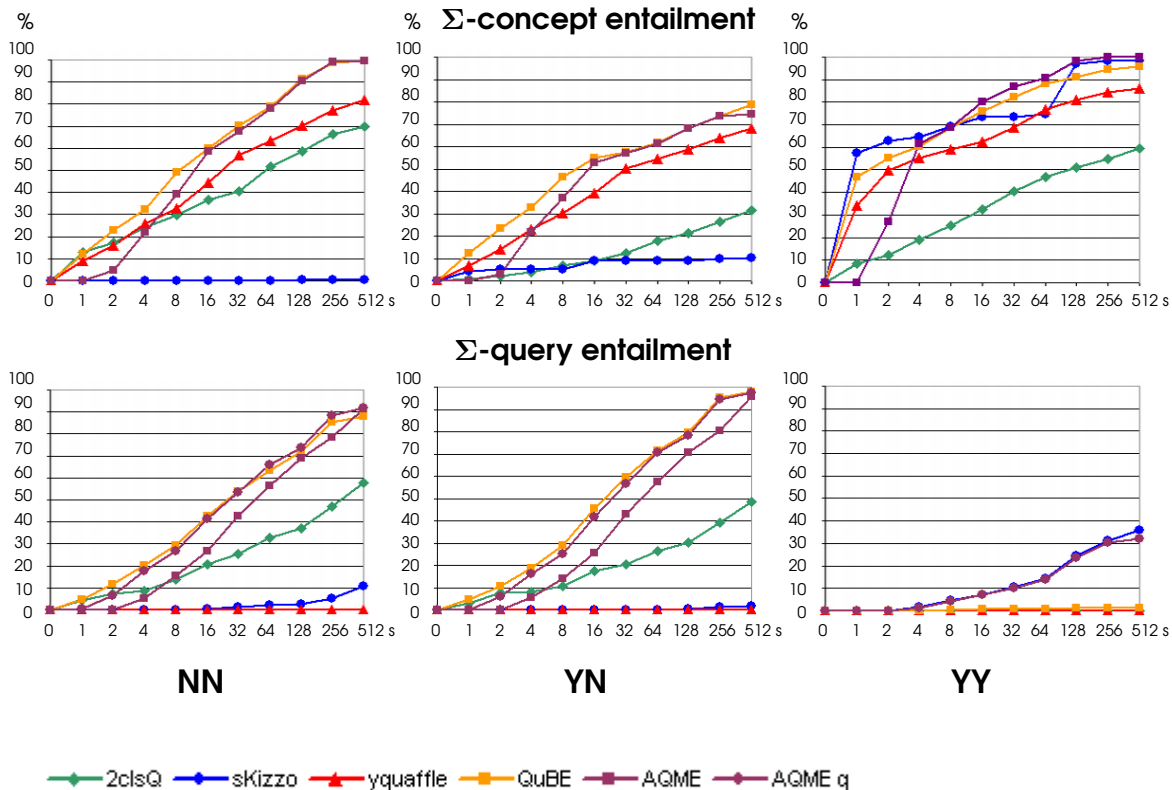
- sKizzo 0.8.2
- 2clsQ
- yQuaffle
- QuBE 6.4
- **AQME**

series	Σ -concept entailment QBF		Σ -query entailment QBF	
	variables	clauses	variables	clauses
NN	1,469–48,631	2,391–74,621	1,715–60,499	5,763–1,217,151
YN	1,460–46,873	2,352–71,177	1,755–59,397	7,006–1,122,361
YY	1,006–16,033	1,420–23,363	1,202–20,513	2,963–204,889

number of clauses is **linear**

quadratic
(in the number of roles)

Experimental results: percentage of solved instances



What is a module?

Let S be an inseparability relation, \mathcal{T} a TBox and Σ a signature.

$\mathcal{M} \subseteq \mathcal{T}$ is

- an S_Σ -module of \mathcal{T} if $\mathcal{M} \equiv_\Sigma^S \mathcal{T}$
- a self-contained S_Σ -module of \mathcal{T} if $\mathcal{M} \equiv_{\Sigma \cup \text{sig}(\mathcal{M})}^S \mathcal{T}$
- a depleting S_Σ -module of \mathcal{T} if $\emptyset \equiv_{\Sigma \cup \text{sig}(\mathcal{M})}^S \mathcal{T} \setminus \mathcal{M}$

\mathcal{M} is a minimal module of \mathcal{T} if it can't be made smaller

Facts:

- depleting \equiv_Σ^q -module \Rightarrow self-contained \equiv_Σ^q -module \Rightarrow \equiv_Σ^q -module
- self-contained \equiv_Σ^c -module \Rightarrow \equiv_Σ^c -module
- There is precisely **one** minimal depleting \equiv_Σ^q -module
- There may be (exponentially) many minimal modules of other types

Modules for $\Sigma = \{\text{Publisher}\}$

- (1) $\text{Publisher} \sqsubseteq \exists \text{pubHasDistrib}$
- (2) $\exists \text{pubHasDistrib}^- \sqsubseteq \text{Distributor}$
- (3) $\text{Publisher} \sqsubseteq \neg \text{Distributor}$
- (4) $\exists \text{pubHasDistrib} \sqsubseteq \text{Publisher}$
- (5) $\text{Publisher} \sqsubseteq \leq 1 \text{ pubHasDistrib}$
- (6) $\text{Role} \sqsubseteq \neg \text{Distributor}$
- (7) $\text{User} \sqsubseteq \neg \text{Distributor}$
- (8) $\text{Publisher} \sqsubseteq \exists \text{pubAdmedBy}$
- (9) $\exists \text{pubAdmedBy}^- \sqsubseteq \text{AdmUser} \sqcup \text{BookUser}$
- (10) $\text{AdmUser} \sqsubseteq \text{User}$

- (11) $\text{BookUser} \sqsubseteq \text{User}$
- (12) $\text{User} \sqsubseteq \exists \text{hasRole}$
- (13) $\exists \text{hasRole}^- \sqsubseteq \text{Role}$
- (14) $\text{Role} \sqsubseteq \neg \text{Publisher}$
- (15) $\text{User} \sqsubseteq \neg \text{Publisher}$
- (16) $\text{Role} \sqsubseteq \neg \text{User}$
- (17) $\text{User} \sqsubseteq \exists \text{userAdmedBy}$
- (18) $\exists \text{userAdmedBy}^- \sqsubseteq \text{AdmUser}$
- (19) $\exists \text{userAdmedBy} \sqsubseteq \text{User}$
- (20) $\exists \text{pubAdmedBy} \sqsubseteq \text{Publisher}$

the minimal S_{Σ}^c -module is \emptyset

minimal S_{Σ}^q -modules of \mathcal{T} : \mathcal{M}_D , \mathcal{M}_R and \mathcal{M}_U

the minimal depleting S_{Σ}^q -module is \mathcal{T}

Module extraction algorithms

- minimal S_Σ -module

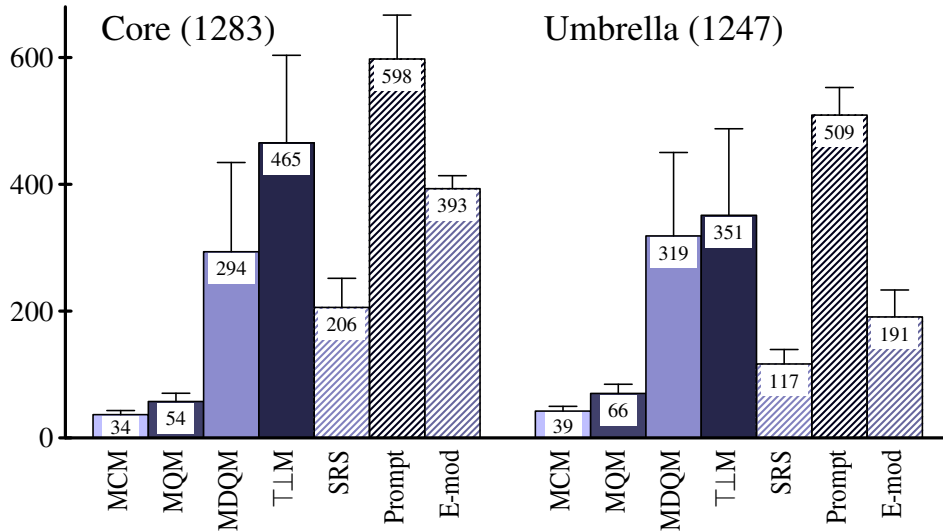
```
input  $\mathcal{T}, \Sigma$ 
 $\mathcal{M} := \mathcal{T}$ 
for each  $\alpha \in \mathcal{M}$  do
  if  $\mathcal{M} \setminus \{\alpha\} \equiv_\Sigma^S \mathcal{M}$  then  $\mathcal{M} := \mathcal{M} \setminus \{\alpha\}$ 
end for
output  $\mathcal{M}$ 
```

NB: depends on the order of axioms in \mathcal{T}

- minimal depleting S_Σ -module

```
input  $\mathcal{T}, \Sigma$ 
 $\mathcal{T}' := \mathcal{T}; \Gamma := \Sigma; \mathcal{W} := \emptyset$ 
while  $\mathcal{T}' \setminus \mathcal{W} \neq \emptyset$  do
  choose  $\alpha \in \mathcal{T}' \setminus \mathcal{W}$ 
   $\mathcal{W} := \mathcal{W} \cup \{\alpha\}$ 
  if  $\mathcal{W} \not\equiv_\Gamma^S \emptyset$  then
     $\mathcal{T}' := \mathcal{T}' \setminus \{\alpha\}; \mathcal{W} := \emptyset; \Gamma := \Gamma \cup \text{sig}(\alpha)$ 
  endif
end while
output  $\mathcal{T} \setminus \mathcal{T}'$ 
```

Practical minimal module extraction



Module sizes and standard deviation for $|\Sigma| = 10$