

# Knowledge Representation and Processing

Florian Rabe (for a course given with Michael Kohlhase)

Computer Science, University Erlangen-Nürnberg, Germany

Summer 2020

# Administrative Information

# Format

## Zoom

- ▶ lectures and exercises via zoom
- ▶ participants muted by default for simplicity
- ▶ interaction strongly encouraged      We don't want to lecture —  
we want to have a conversation during which you learn
- ▶ let's try out zoom
  - ▶ use reactions to say yes no, ask for break etc.
  - ▶ feel free to annotate my slides
  - ▶ talk in the chat

## Recordings

- ▶ maybe prerecorded video lectures or recorded zoom meeting
- ▶ to be decided along the way

# Background

## Instructors

- ▶ Prof. Dr. Michael Kohlhase  
Professor of Knowledge Representation and Processing
- ▶ PD Dr. Florian Rabe  
same research group

## Course

- ▶ This course is given for the first time
- ▶ Always a little bit of an experiment      cutting edge vs. unpolished
- ▶ Could become signature course of our research group      same name!

# Prerequisites

## Required

- ▶ basic knowledge about formal languages, context-free grammars  
but we'll do a quick revision here

## Helpful

- ▶ Algorithms and Data Structures mostly as a contrast to this lecture
- ▶ Basic logic we'll revise it slightly differently here
- ▶ all other courses as examples of how knowledge pervades all of CS

## General

- ▶ Curiosity this course is a bit unusual
- ▶ Interest in big picture  
this course touches on lots of things from all over CS

# Examination and Grading

## Suggestion

- ▶ grade determined by single exam
- ▶ written or oral depends on number of students
- ▶ some acknowledgment for practical exercises

to be finalized next week

## Exam-relevant

- ▶ anything mentioned in notes
- ▶ anything discussed in lectures

neither is a superset of the other!

# Materials and Exam-Relevance

## Textbook

- ▶ does not exist
- ▶ normal for research-near specialization courses

## Notes

- ▶ textbook-style but not as comprehensive
- ▶ developed along the way

## Slides

- ▶ not comprehensive
- ▶ used as visual aid, conversation starters

# Communication

## Open for questions

- ▶ open door policy in our offices if the lockdown ever ends
- ▶ always room for questions during lectures
- ▶ for personal questions, contact me during/after lecture or by email
- ▶ forum at <https://fsi.cs.fau.de/forum/154-Wissensrepraesentation-und-Verarbeitung>

## Materials

- ▶ official notes and slides as pdf:  
<https://kwarc.info/teaching/WuV/>  
will be updated from time to time
- ▶ watch me prepare the materials: <https://github.com/florian-rabe/Teaching/tree/master/WuV>  
pull requests and issues welcome



# Exercises

## Learning Goals

- ▶ Get acquainted with state of the art of practice
- ▶ Try out real tools

## Homeworks

- ▶ one major project as running example
- ▶ homeworks building on each other

build one large knowledge-based system  
details on later slides

# Overview and Essential Concepts

# Representation and Processing

Common pairs of concepts:

| Representation  | Processing |
|-----------------|------------|
| Static          | Dynamic    |
| Situation       | Change     |
| Be              | Become     |
| Data Structures | Algorithms |
| Set             | Function   |
| State           | Transition |
| Space           | Time       |

## Data and Knowledge

$2 \times 2$  key concepts

| Syntax    | Data      |
|-----------|-----------|
| Semantics | Knowledge |

- ▶ Data: any object that can be stored in a computer  
Example:  $((49.5739143, 11.0264941), "2020 - 04 - 21 T16 : 15 : 00 CEST")$
- ▶ Syntax: a system of rules that describes which data is **well-formed**  
Example: "a pair of (a pair of two IEEE double precision floating point numbers) and a string encoding of a time stamp"
- ▶ Semantics: system of rules that determines the meaning of well-formed data
- ▶ Knowledge: combination of some data with its syntax and semantics

# Knowledge is Elusive

## Representation of key concepts

- ▶ Data: using primitive objects  
implemented as bits, bytes, strings, records, arrays, ...
- ▶ Syntax: (context-free) grammars, (context-sensitive) type systems  
implemeted as inductive data structures
- ▶ Semantics: functions for evaluation, interpretation, of well-formed data  
implemented as recursive algorithms on the syntax
- ▶ Knowledge: elusive  
emerges from applying and interacting with the semantics

## Semantics as Translation

- ▶ Knowledge can be captured by a higher layer of syntax
- ▶ Then semantics is translation into syntax

| Data syntax        | Semantics function        | Knowledge syntax    |
|--------------------|---------------------------|---------------------|
| SPARQL query       | evaluation                | result set          |
| SQL query          | evaluation                | result table        |
| program            | compiler                  | binary code         |
| program expression | interpreter               | result value        |
| logical formula    | interpretation in a model | mathematical object |
| HTML document      | rendering                 | graphics context    |

# Heterogeneity of Data and Knowledge

- ▶ Capturing knowledge is difficult
- ▶ Many different approaches to semantics
  - ▶ fundamental formal and methodological differences
  - ▶ often captured in different fields, conferences, courses, languages, tools
- ▶ Data formats equally heterogeneous
  - ▶ ontologies
  - ▶ programs
  - ▶ logical proofs
  - ▶ databases
  - ▶ documents

# Challenges of Heterogeneity

## Challenges

- ▶ collaboration across communities
- ▶ translation across languages
- ▶ conversion between data formats
- ▶ interoperability across tools

## Sources of problems

- ▶ interoperability across formats/tools major source of
  - ▶ complexity
  - ▶ bugs
- ▶ friction in project team due to differing preferences, expertise
- ▶ difficult choice between languages/tools with competing advantages
  - ▶ reverting choices difficult, costly
  - ▶ maintaining legacy choices increases complexity



## Aspects of Knowledge

- ▶ Tetrapod model of knowledge      **active research by our group**
- ▶ classifies approaches to knowledge into five aspects

| Aspect         | KRLs (examples)                                    |
|----------------|--|
| ontologization | ontology languages (OWL), description logics (ALC) |
| concretization | relational databases (SQL, JSON)                   |
| computation    | programming languages (C)                          |
| deduction      | logics (HOL)                                       |
| narration      | document languages (HTML, LaTeX)                   |

## Relations between the Aspects

Ontology is distinguished: capture the knowledge that the other four aspects share



## Complementary Advantages of the Aspects

| Aspect          | objects                   | advantage                  | characteristic<br>joint advantage<br>of the other as-<br>pects | application                              |
|-----------------|---------------------------|----------------------------|--|--|
| ded.<br>comp.   | formal proofs<br>programs | correctness<br>efficiency  | ease of use<br>well-<br>definedness                            | verification<br>execution                |
| concr.<br>narr. | concrete objects<br>texts | tangibility<br>flexibility | abstraction<br>formal seman-<br>tics                           | storage/retrieval<br>human understanding |

| Aspect pair                | characteristic advantage                          |
|----------------------------|---|
| ded./comp.<br>narr./concr. | rich meta-theory<br>simple languages              |
| ded./narr.<br>comp./concr. | theorems and proofs<br>normalization              |
| ded./concr.<br>comp./narr. | decidable well-definedness<br>Turing completeness |

# Structure of the Course

## Aspect-independent parts

- ▶ general methods that are shared among the aspects
- ▶ to be discussed as they come up

## Aspects-specific parts

- ▶ one part (about 2 weeks) for each aspect
- ▶ high-level overview of state of the art
- ▶ focus on comparison/evaluation of the aspect-specific results

## Structure of the Exercises

### One major project

- ▶ representative for a project that a CS graduate might be put in charge of
- ▶ challenging heterogeneous data and knowledge
- ▶ requires integrating/combining different languages, tools

unique opportunity in this course because knowledge is everywhere

### Concrete project

- ▶ develop a univis-style system for a university
- ▶ lots of heterogeneous knowledge
  - ▶ course and program descriptions
  - ▶ legal texts
  - ▶ websites
  - ▶ grade tables
  - ▶ transcript generation code
- ▶ build a completely functional system applying the lessons of the course

# Ontological Knowledge

## Components of an Ontology

8 main declarations

- ▶ **individual** — concrete objects that exist in the real world, e.g., "Florian Rabe" or "WuV"
- ▶ **concept** — abstract groups of individuals, e.g., "instructor" or "course"
- ▶ **relation** — binary relations between two individuals, e.g., "teaches"
- ▶ **properties** — binary relations between an individuals and a concrete value (a number, a date, etc.), e.g., "has-credits"
- ▶ **concept assertions** — the statement that a particular individual is an instance of a particular concept
- ▶ **relation assertions** — the statement that a particular relation holds about two individuals
- ▶ **property assertions** — the statement that a particular individual has a particular value for a particular property
- ▶ **axioms** — statements about relations between concepts, e.g., "instructor"  $\sqsubseteq$  "person"

# Divisions of an Ontology

## Abstract vs. concrete

- ▶ TBox: concepts, relations, properties, axioms  
everything that does not use individuals
- ▶ ABox: individuals and assertions

## Named vs. unnamed

- ▶ Signature: individuals, concepts, relations, properties  
together called entities or resources
- ▶ Theory: assertions, axioms



# Comparison of Terminology

| Here       | OWL             | Description logics | ER model    | UML                | semantics via logics |
|------------|-----------------|--------------------|-------------|--------------------|----------------------|
| individual | instance        | individual         | entity      | object, instance   | constant             |
| concept    | class           | concept            | entity-type | class              | unary predicate      |
| relation   | object property | role               | role        | association        | binary predicate     |
| property   | data property   | (not common)       | attribute   | field of base type | binary predicate     |

|                         |                |             |
|-------------------------|----------------|-------------|
| domain                  | individual     | concept     |
| type theory, logic      | constant, term | type        |
| set theory              | element        | set         |
| database                | row            | table       |
| philosophy <sup>1</sup> | object         | property    |
| grammar                 | proper noun    | common noun |

<sup>1</sup>as in <https://plato.stanford.edu/entries/object/>

## Ontologies as Sets of Triples

| Assertion          | Triple         |               |              |
|--------------------|----------------|---------------|--------------|
|                    | Subject        | Predicate     | Object       |
| concept assertion  | "Florian Rabe" | is-a          | "instructor" |
| relation assertion | "Florian Rabe" | "teaches"     | "WuV"        |
| property assertion | "WuV"          | "has credits" | 7.5          |

Efficient representation of ontologies using RDF and RDFS  
standardized special entities.

## Special Entities

RDF and RDFS define special entities for use in ontologies:

- ▶ "rdfs:Resource": concept of which all individuals are an instance and thus of which every concept is a subconcept
- ▶ "rdf:type": relates an entity to its type:
  - ▶ an individual to its concept (corresponding to is-a above)
  - ▶ other entities to their special type (see below)
- ▶ "rdfs:Class": special class for the type of classes
- ▶ "rdf:Property": special class for the type of properties
- ▶ "rdfs:subClassOf": a special relation that relates a subconcept to a superconcept
- ▶ "rdfs:domain": a special relation that relates a relation to the concepts of its subjects
- ▶ "rdfs:range": a special relation that relates a relation/property to the concept/type of its objects

Goal/effect: capture as many parts as possible as RDF triples.

## Declarations as Triples using Special Entities

| Assertion                   | Triple     |                   |                 |
|-----------------------------|------------|-------------------|-----------------|
|                             | Subject    | Predicate         | Object          |
| individual                  | individual | "rdf:type"        | "rdfs:Resource" |
| concept                     | concept    | "rdf:type"        | "rdf:Class"     |
| relation                    | relation   | "rdf:type"        | "rdf:Property"  |
| property                    | property   | "rdf:type"        | "rdf:Property"  |
| concept assertion           | individual | "rdf:type"        | concept         |
| relation assertion          | individual | relation          | individual      |
| property assertion          | individual | property          | value           |
| for special forms of axioms |            |                   |                 |
| $c \sqsubseteq d$           | $c$        | "rdfs:subClassOf" | $d$             |
| $\text{dom } r \equiv c$    | $r$        | "rdfs:domain"     | $c$             |
| $\text{rng } r \equiv c$    | $r$        | "rdfs:range"      | $c$             |

## An Example Ontology Language

see syntax of BOL in the lecture notes

# Semantics as Translation

## Example: Syntax of Arithmetic Language

Syntax: represented as formal grammar

### Numbers

|                  |          |
|------------------|----------|
| $N ::= 0 \mid 1$ | literals |
| $\mid N + N$     | sum      |
| $\mid N * N$     | product  |

### Formulas

|                    |                  |
|--------------------|------------------|
| $F ::= N \doteq N$ | equality         |
| $\mid N \leq N$    | ordering by size |

Implementation as inductive data type

## Example: Semantics of Arithmetic Language

Semantics: represented as translation into known language

Problem: Need to choose a known language first

Here: unary numbers represented as strings

Built-in data (strings and booleans):

|                     |            |
|---------------------|------------|
| $S ::= \varepsilon$ | empty      |
| (Unicode)           | characters |
| $B ::= \text{true}$ | truth      |
| false               | falsity    |

Built-in operations to work on the data:

- ▶ concatenation of strings  $S ::= \text{conc}(S, S)$
- ▶ replacing all occurrences of  $c$  in  $S_1$  with  $S_2$   
 $S ::= \text{replace}(S_1, c, S_2)$
- ▶ equality test:  $B ::= S_1 == S_2$
- ▶ prefix test:  $B ::= \text{startsWith}(S_1, S_2)$



## Example: Semantics of Arithmetic Language

Represented as function from syntax to semantics

- ▶ mutually recursive, inductive functions for each non-terminal symbol
- ▶ compositional: recursive call on immediate subterms of argument

For numbers  $n$ : semantics  $\llbracket n \rrbracket$  is a string

- ▶  $\llbracket 0 \rrbracket = \varepsilon$
- ▶  $\llbracket 1 \rrbracket = \text{"|"}$
- ▶  $\llbracket m + n \rrbracket = \text{conc}(\llbracket m \rrbracket, \llbracket n \rrbracket)$
- ▶  $\llbracket m * n \rrbracket = \text{replace}(\llbracket m \rrbracket, \text{"|"}, \llbracket n \rrbracket)$

For formulas  $f$ : semantics  $\llbracket f \rrbracket$  is a boolean

- ▶  $\llbracket m \dot{=} n \rrbracket = \llbracket m \rrbracket == \llbracket n \rrbracket$
- ▶  $\llbracket m \leq n \rrbracket = \text{startsWith}(\llbracket n \rrbracket, \llbracket m \rrbracket)$

## Semantics of BOL

| Aspect         | kind of semantic language | semantic language |
|----------------|---------------------------|-------------------|
| deduction      | logic                     | SFOL              |
| concretization | database language         | SQL               |
| computation    | programming language      | Scala             |
| narration      | natural language          | English           |

see details of each translation in the lecture notes

## General Definition

A semantics by translation consists of

- ▶ syntax: a formal language  $I$
- ▶ semantic language: a formal language  $L$   
different or same aspect as  $I$
- ▶ semantic prefix: a theory  $P$  in  $L$   
formalizes fundamentals that are needed to represent  $I$ -objects
- ▶ interpretation: translates every  $I$ -theory  $T$  to an  $L$ -theory  $P, \llbracket T \rrbracket$

## Common Principles

Properties shared by all semantics of BOL

not part of formal definition, but best practices

- ▶ *I*-declaration translated to *L*-declaration for the same name
- ▶ ontologies translated declaration-wise
- ▶ one inductive function for every kind of complex *I*-expression
  - ▶ individuals, concepts, relations, properties, formulas
  - ▶ maps *I*-expressions to *L*-expressions
- ▶ atomic cases (base cases): *I*-identifier translated to *L*-identifier of the same name or something very similar
- ▶ complex cases (step cases): compositional

## Compositionality

Case for operator  $*$  in interpretation function compositional iff interpretation of  $*(e_1, \dots, e_n)$  only depends on the interpretation of the  $e_i$

$$\llbracket *(e_1, \dots, e_n) \rrbracket = \llbracket * \rrbracket (\llbracket e_1 \rrbracket, \dots, \llbracket e_n \rrbracket)$$

for some function  $\llbracket * \rrbracket$

Example:  $;$ -operator of BOL in translation to FOL

- ▶ translation:  $\llbracket R_1; R_2 \rrbracket = \exists m : \iota. \llbracket R_1 \rrbracket(x, m) \wedge \llbracket R_2 \rrbracket(m, y)$
- ▶ special case of the above via
  - ▶  $* = ;$ ;
  - ▶  $n = 2$
  - ▶  $\llbracket ; \rrbracket = (p_1, p_2) \mapsto \exists m : \iota. p_1(x, m) \wedge p_2(m, y)$
- ▶ Indeed, we have  $\llbracket R_1; R_2 \rrbracket = \llbracket ; \rrbracket (\llbracket R_1 \rrbracket, \llbracket R_2 \rrbracket)$

## Compositionality (2)

Translation compositional iff

- ▶ one translation function for each non-terminal all written  $\llbracket - \rrbracket$
- ▶ each defined by one induction on syntax  
i.e., one case for production  
mutually recursive
- ▶ all cases compositional

Substitution theorem: a compositional translation satisfies

$$\llbracket E(e_1, \dots, e_n) \rrbracket = \llbracket E \rrbracket(\llbracket e_1 \rrbracket, \dots, \llbracket e_n \rrbracket)$$

for

- ▶ every expression  $E(N_1, \dots, N_n)$  with non-terminals  $N_i$
- ▶ some function  $\llbracket E \rrbracket$  that only depends on  $E$

## Compositionality (3)

$$\llbracket E(e_1, \dots, e_n) \rrbracket = \llbracket E \rrbracket(\llbracket e_1 \rrbracket, \dots, \llbracket e_n \rrbracket)$$

for every expression  $E(N_1, \dots, N_n)$  with non-terminals  $N_i$

Now think of

- ▶ variable  $x_i$  of type  $N_i$  instead of non-terminal  $N_i$
- ▶  $E(x_1, \dots, x_n)$  as expression with free variables  $x_i$  of type  $N_i$
- ▶ expressions  $e$  derived from  $N$  as expressions of type  $N$
- ▶  $E(e_1, \dots, e_n)$  as result of substituting  $e_i$  for  $x_i$
- ▶  $\llbracket E \rrbracket(x_1, \dots, x_n)$  as (semantic) expression with free variables  $x_i$

Then both sides of equations act on  $E(x_1, \dots, x_n)$ :

- ▶ left side yields  $\llbracket E(e_1, \dots, e_n) \rrbracket$  by
  - ▶ first substitution  $e_i$  for  $x_i$
  - ▶ then semantics  $\llbracket - \rrbracket$  of the whole
- ▶ right side yields  $\llbracket E \rrbracket(\llbracket e_1 \rrbracket, \dots, \llbracket e_n \rrbracket)$  by
  - ▶ first semantics  $\llbracket - \rrbracket$  of all parts
  - ▶ then substitution  $\llbracket e_i \rrbracket$  for  $x_i$

semantics commutes with substitution

# Non-Compositionality

## Examples

- ▶ deduction: cut elimination, translation from natural deduction to Hilbert calculus
- ▶ computation: optimizing compiler, e.g., loop unrolling
- ▶ concretization: query optimization, e.g., turning a WHERE of a join into a join of WHEREs,
- ▶ narration: ambiguous words are translated based on context

## Typical sources

- ▶ subcases in a case of translation function
  - ▶ based on inspecting the arguments, e.g., subinduction
  - ▶ based on context
- ▶ custom-built semantic prefix



# Type Systems

## Breakout Question

Is this an improvement over BOL?

|              |
|--------------|
| Declarations |
|--------------|

|         |                   |                           |                         |
|---------|-------------------|---------------------------|-------------------------|
| $D ::=$ | <b>individual</b> | $ID : C$                  | typed atomic individual |
|         | <b>concept</b>    | $ID$                      | atomic concept          |
|         | <b>relation</b>   | $ID \subseteq C \times C$ | typed atomic relation   |
|         | <b>property</b>   | $ID \subseteq C \times T$ | typed atomic property   |

rest as before

## Actually, when is a language an improvement?

Criteria: **orthogonal, often mutually exclusive**

- ▶ syntax design trade-off
  - ▶ expressivity: easy to express knowledge  
e.g., big grammar, extra production for every user need
  - ▶ simplicity: easy to implement/interpret  
e.g., few, carefully chosen productions
- ▶ semantics: specify, implement, document
- ▶ intended users
  - ▶ skill level
  - ▶ prior experience with related languages
  - ▶ amount of training needed
- ▶ long-term plans: re-answer the above question but now
  - ▶ maintainability: syntax was changed, everything to be redone
  - ▶ scalability: expressed knowledge content has reached huge sizes

# Church vs. Curry Typing

|  | intrinsic   | extrinsic   |
|--|---|---|
| $\lambda$ -calculus by<br>type is<br>typing is a<br>objects have<br>types interpreted as | Church<br>carried by object<br>function objects $\rightarrow$ types<br>unique type<br>disjoint sets | Curry<br>given by environment<br>relation objects $\times$ types<br>any number of types<br>unary predicates |
| type given by<br>example   | part of declaration<br><b>individual</b> "WuV" : "course"   | additional axiom<br><b>individual</b> "Wuv",<br>"WuV" is-a "course"   |
| examples   | SFOL, SQL<br>most logics, functional PLs<br><br>many type theories                                  | OWL, Scala, English<br>ontology, OO,<br>natural languages<br>set theories                                   |

# Type Checking

|  | intrinsic  | extrinsic   |
|--|--|---|
| type is<br>typing is a<br>objects have   | carried by object<br>function objects $\rightarrow$ types<br>unique type   | given by environment<br>relation objects $\times$ types<br>any number of types  |
| type given by<br>example   | part of declaration<br><b>individual</b> "WuV" : "course"  | additional axiom<br><b>individual</b> "Wuv",<br>"WuV" is-a "course"   |
| type inference for $x$<br>type checking<br>subtyping $A <: B$<br>typing decidable<br>typing errors | uniquely infer $A$ from $x$<br>inferred=expected<br>cast from $A$ to $B$<br>yes unless too expressive<br>static (compile-time) | find minimal $A$ with $x : A$<br>prove $x : A$<br>$x : A$ implies $x : B$<br>no unless restricted<br>dynamic (run-time) |
| advantages   | easy<br>unique type inference  | flexible<br>allows subtyping  |

# Curry Typing in BOL

| language     | objects           | types                          | typing relation          |
|--------------|-------------------|--------------------------------|--------------------------|
| Syntax       | individuals       | concepts                       | $i$ is-a $c$             |
| Semantics in |                   |                                |                          |
| FOL          | type $\iota$      | predicates $c \subseteq \iota$ | $c(i)$ true              |
| SQL          | table Individuals | tables containing ids          | id of $i$ in table $c$   |
| Scala        | String            | hash sets of strings           | $c.\text{contains}(i)$   |
| English      | proper nouns      | common nouns                   | " $i$ is a $c$ " is true |

# Subtyping

Subtyping works best with Curry Typing

- ▶ explicit subtyping as in  $\mathbb{N} <: \mathbb{Z}$
- ▶ comprehension/refinement as in  $\{x : \mathbb{N} \mid x \neq 0\}$
- ▶ operations like union and intersection on types
- ▶ inheritance between classes, in which case subclass = subtype
- ▶ anonymous record types as in  $\{x : \mathbb{N}, y : \mathbb{Z}\} <: \{x : \mathbb{N}\}$

## A General Definition of a Type System

A **type system** consists of

- ▶ a collection, whose elements are called **objects**,
- ▶ a collection, whose elements are called **intrinsic types**,
- ▶ a function assigning to every object  $x$  its **intrinsic type**  $I$ , in which case we write  $x : I$ ,
- ▶ for some intrinsic types  $I$ 
  - ▶ an intrinsic type  $E_I$
  - ▶ a relation  $\in_I$  between objects with intrinsic types  $I$  and  $E_I$ , called the **extrinsic typing** relation for  $I$ .



# Examples

| System      | intrinsic types           | $E_I$             | $\in_I$                         |
|-------------|---------------------------|-------------------|---------------------------------|
| pure Church | one per type              | none              | none                            |
| pure Curry  | objects $O$ , types $T$   | $E_O = T$         | $\in_O = :$                     |
| FOL         | one per type              | none              | none                            |
| Scala       | <i>Any</i> , <i>Class</i> | $E_{Any} = Class$ | $\in_{Any} = \text{isInstance}$ |
| BOL         | <i>Ind</i> , <i>Conc</i>  | $E_{Ind} = Conc$  | $\in_{Ind} = \text{is-a}$       |
| set theory  | <i>Set</i> , <i>Prop</i>  | $E_{Set} = Set$   | $\in_{Set} = \in$               |

## Breakout Question

What do the following have in common?

- ▶ Java class
- ▶ SQL schema for a table
- ▶ logical theory (e.g., Monoid)

## Breakout Question

What do the following have in common?

- ▶ Java class
- ▶ SQL schema for a table
- ▶ logical theory (e.g., Monoid)

all are (essentially) abstract data types

## Abstract Data Types: Motivation

Recall subject-centered representation of assertion triples:

```
individual "FlorianRabe"  
  is-a "instructor" "male"  
  "teach" "WuV" "KRMT"  
  "age" 40  
  "office" "11.137"
```

Can we use types to force certain assertions to occur together?

- ▶ Every instructor should teach a list of courses.
- ▶ Every instructor should have an office.

## Abstract Data Types: Motivation

Inspires **subject-centered types**, e.g.,

```
concept instructor
  teach course*
  age: int
  office: string
```

```
individual "FlorianRabe": "instructor"
  is-a "male"
  teach "WuV" "KRMT"
  age 40
  office "11.137"
```

Incidental benefits:

- ▶ no need to declare relations/properties separately
- ▶ reuse relation/property names  
distinguish via qualified names: `instructor .age`

## Abstract Data Types: Motivation

Natural next step: inheritance

```
concept person  
  age: int
```

```
concept male <: person
```

```
concept instructor <: person  
  teach course*  
  office: string
```

```
individual "FlorianRabe": "instructor"  $\sqcap$  "male"  
  "teach" "WuV" "KRMT"  
  "age" 40  
  "office" "11.137"
```

our language quickly gets a very different flavor

## Abstract Data Types: Examples

Prevalence of abstract data types:

| aspect         | language | abstract data type                    |
|----------------|----------|---------------------------------------|
| ontologization | UML      | class                                 |
| concretization | SQL      | table schema                          |
| computation    | Scala    | class, interface                      |
| deduction      | various  | theory, specification, module, locale |
| narration      | various  | emergent feature                      |

same idea, but may look very different across languages

# Abstract vs. Concrete Types

**Concrete** type: values are

- ▶ given by their internal form,
- ▶ defined along with the type, typically built from already-known pieces.

examples: products, inductive data types

**Abstract** type: values are

- ▶ given by their externally visible properties,
- ▶ defined in any environment that understands the type definition.

main example: abstract data types



# Abstract Data Types: Examples

| aspect         | type           | values                            |
|----------------|----------------|-----------------------------------|
| computation    | abstract class | instances of implementing classes |
| concretization | table schema   | table rows                        |
| deduction      | theory         | models                            |

Values depend on the environment in which the type is used:

- ▶ class defined in one specification language (e.g., UML),  
implementations in programming languages Java, Scala, etc.  
available values may depend on run-time state
- ▶ theory defined in logic,  
models defined in set theories, type theories, programming  
languages  
available values may depend on philosophical position

## Abstract Data Types: Definition

Given some type system, an **abstract data type** (ADT) is

- ▶ a **flat** type

$$\{c_1 : T_1 [= t_1], \dots, c_n : T_n [= t_n]\}$$

where

- ▶  $c_i$  are distinct names
  - ▶  $T_i$  are types
  - ▶  $t_i$  are optional definitions; if given,  $t_i : T_i$  required
- ▶ or a **mixin** type

$$A_1 * \dots * A_n$$

for ADTs  $A_i$ .

Languages may or may not make ADTs additional types of the type system

## Abstract Data Types: Class Definitions

A class definition in OO:

```
abstract class  $a$  extends  $a_1$  with ... with  $a_m$  {  
   $c_1 : T_1$   
   $\vdots$   
   $c_n : T_n$   
}
```

Corresponding ADT definition:

$$a = a_1 * \dots * a_m * \{c_1 : T_1, \dots, c_n : T_n\}$$

The usual terminology:

- ▶  $a$  **inherits** from  $a_i$
- ▶  $a_i$  are **super- $X$**  or **parent- $X$**  of  $a$  where  $X$  is whatever the language calls its ADTs (e.g.,  $X$ =class)

# Abstract Data Types: Flattening

The **flattening**  $A^b$  of an ADT  $A$  is

- ▶ if  $A$  is flat:  $A^b = A$
- ▶  $(A_1 * \dots * A_n)^b$  is union of all  $A_i^b$   
where duplicate field names are handled as follows
  - ▶ same name, same type, same or omitted definition: merge  
details may be much more difficult
  - ▶ otherwise: ill-formed

# Abstract Data Types: Subtleties

We gloss over several major issues:

- ▶ How exactly do we merge duplicate field names? Does it always work? **implement abstract methods, override, overload**
- ▶ Is recursion allowed, i.e., can I define an ADT  $a = A$  where  $a$  occurs in  $A$ ?  
**common in OO-languages: use  $a$  in the types of its fields**
- ▶ What about ADTs with type arguments?  
**e.g., generics in Java, square-brackets in Scala**
- ▶ Is mutual recursion between fields in a flat type allowed?  
**common in OO-languages**
- ▶ Is  $*$  commutative? What about dependencies between fields?  
**no unique answers**  
**incarnations of ADTs subtly different across languages**

## Context-Sensitive Syntax

## Definition

A **language system** consists of

- ▶ context-free syntax
- ▶ distinguished non-terminal symbol  $\mathcal{V}$   
words called **vocabularies**
- ▶ some distinguished non-terminal symbols  $\mathcal{E}$   
words called  **$\mathcal{E}$ -expressions**
- ▶ unary predicate  $\text{wft}(\Theta)$  on vocabularies  $\Theta$   
well-formed vocabulary  $\Theta$
- ▶ unary predicates  $\text{wff}_{\Theta}^{\mathcal{E}}(E)$  well-formed  $\mathcal{E}$ -expressions  $E$

# Typical Structure

## Vocabularies

- ▶ lists of declarations

## Declarations

- ▶ named
- ▶ at least one for each expression kind
- ▶ may contain other expressions e.g., type, definition
- ▶ may contain nested declarations e.g., fields in an ADT

## Expressions

- ▶ inductive data type
- ▶ relative to vocabulary names occur as base cases
- ▶ formulas as special case



## Vocabularies and Expressions

| Aspect         | vocabulary $\Theta$ | expression kinds $\mathcal{E}$                   |
|----------------|---------------------|--|
| Ontologization | ontology            | individual, concept, relation, property, formula |
| Concretization | database schema     | cell, row, table, formula                        |
| Computation    | program             | term, type, object, class, ...                   |
| Logic          | signature, theory   | term, type, formula, ...                         |
| Narration      | dictionary          | phrases, sentences, texts                        |

## Examples

See notes made during the lecture for examples

# Concrete Knowledge and Typed Ontologies

# Motivation

## Main ideas

- ▶ Ontology abstractly describes concepts and relations
- ▶ Tool maintains concrete data set
- ▶ Focus on efficiently
  - ▶ identifying (i.e., assign names)
  - ▶ representing
  - ▶ processing
  - ▶ querying

large sets of concrete data

## Recall: TBox-ABox distinction

- ▶ TBox: general parts, abstract, fixed  
main challenge: correct modeling of domain
- ▶ ABox: concrete individuals and assertions about them, growing  
main challenge: aggregate them all

# Concrete Data

## Concrete is

- ▶ Base values: integers, strings, booleans, etc.
- ▶ Collections: sets, multisets, lists (always finite)
- ▶ Aggregations: tuples, records (always finite)
- ▶ User-defined concrete data: enumerations, inductive types
- ▶ Advanced objects: finite maps, graphs, etc.

## Concrete is not

- ▶ Free symbols to be interpreted by a model  
exception: foreign function interfaces
- ▶ Variables (free or bound)  $\lambda$ -abstraction, quantification
- ▶ Symbolic expressions  
formulas, algorithms  
Exceptions:
  - ▶ expressions of inductive type
  - ▶ application of built-in functions
  - ▶ queries that return concrete data

## Breakout question

What is the difference between

- ▶ an OWL ontology
- ▶ an SQL database

## Two Approaches

### Based on **untyped** (Curry-typed) ontology languages

- ▶ Representation based on **knowledge graph**
- ▶ Ontology written in BOL-like language
- ▶ Data maintained as **set of triples** tool = triple store
- ▶ Typical language/tool design
  - ▶ ontology and query language **separate** e.g., OWL, SPARQL
  - ▶ triple store and query engine integrated e.g., Virtuoso tool

### Based on **typed** (Church-typed) ontology languages

- ▶ Representation based on **abstract data types**
- ▶ Ontology written as database schema
- ▶ Data maintained as **tables** tool = (relational) database
- ▶ Typical language/tool design
  - ▶ ontology and query language **integrated** e.g., SQL
  - ▶ table store and query engine integrated e.g., SQLite tool

# Evolution of Approaches

## Our usage is non-standard

- ▶ Common
  - ▶ ontologies = untyped approach, OWL, triples, SPARQL
  - ▶ databases = typed approach, tables, SQL
- ▶ Our understanding: two approaches evolved from same idea
  - ▶ triple store = untyped database
  - ▶ SQL schema = typed ontology

## Evolution

- ▶ Typed-untyped distinction minor technical difference
- ▶ Optimization of respective advantages causes speciation
- ▶ Today segregation into different
  - ▶ jargons
  - ▶ languages, tools
  - ▶ communities, conferences
  - ▶ courses



## Curry-typed concrete data

### Central data structure = knowledge graph

- ▶ nodes = individuals  $i$ 
  - ▶ identifier
  - ▶ sets of concepts of  $i$
  - ▶ key-value sets of properties of  $i$
- ▶ edges = relation assertions
  - ▶ from subject to object
  - ▶ labeled with name of relation

### Processing strengths

- ▶ store: as triple set
- ▶ edit: Protege-style or graph-based
- ▶ visualize: as graph different colors for concepts, relations
- ▶ query: match, traverse graph structure
- ▶ untyped data simplifies integration, migration

## Church-typed concrete data

### Central data structure = relational database

- ▶ tables = abstract data type
- ▶ rows = objects of that type
- ▶ columns = fields of ADT
- ▶ cells = values of fields

### Processing strengths

- ▶ store: as CSV text files, or similar
- ▶ edit: SQL commands or table editors
- ▶ visualize: as table view
- ▶ query: relational algebra
- ▶ typed data simplifies selecting, sorting, aggregating

# Identifiers

## Curry-Typed Knowledge Graph

- ▶ concept, relation, property names given in TBox
- ▶ individual names attached to nodes

## Church-Typed Database

- ▶ table, column names given in schema
- ▶ row identified by distinguished key
- options:
  - ▶ preexistent characteristic column or column set
  - ▶ added upon insertion
    - ▶ UUID string
    - ▶ incremental integers
- ▶ column/row identifiers formed by qualifying with table name

# Axioms

## Curry-Typed Knowledge Graph

- ▶ traditionally very expressive axioms
- ▶ yields inferred assertions
- ▶ triple store must do consequence closure to return correct query results
- ▶ not all axioms supported by every triple store

## Church-Typed Database

- ▶ typically no axioms
- ▶ instead consistency constraints, triggers
- ▶ allows limited support for axioms without calling it that way
- ▶ stronger need for users to program the consequence closure manually

## Breakout question

When using typed concrete data,  
how to fully realize abstract data types

- ▶ nesting: ADTs occurring as field types
- ▶ inheritance between ADTs
- ▶ mixins

## ADTs in Typed Concrete Data

### Nesting: field $a : A$ in ADT $B$

- ▶ field types must be base types,  $a : A$  not allowed
- ▶ allow  $ID$  as additional base type
- ▶ use field  $a : ID$  in table  $B$
- ▶ store value of  $a$  in table  $A$

### Inheritance: $B$ inherits from $A$

- ▶ add field  $parent_A$  to table  $B$
- ▶ store values of inherited fields of  $B$  in table  $A$

general principle: all objects of type  $A$  stored in same table

### Mixin: $A * B$

- ▶ essentially join of tables  $A$  and  $B$  on common fields
- ▶ some subtleties depending on ADT flattening

## Open/Closed World

- ▶ Question: is the data complete?
  - ▶ closed world: yes
  - ▶ open world: not necessarily
- ▶ Dimensions of openness
  - ▶ existence of individual objects
  - ▶ assertions about them
- ▶ Sources of openness
  - ▶ more exists but has not yet been added
  - ▶ more could be created later
- ▶ Orthogonal to typed/untyped distinction, but in practice
  - ▶ knowledge graphs use open world
  - ▶ databases use closed world

Open world is natural state, closing adds knowledge

## Closing the World

### Derivable consequences

- ▶ induction: prove universal property by proving for each object
- ▶ negation by failure: atomic property false if not provable
- ▶ term-generation constraint: only nameable objects exist

### Enabled operations

- ▶ universal set: all objects
- ▶ complement of concept/type
- ▶ defaults: assume default value for property if not otherwise asserted

### Monotonicity problem

- ▶ monotone operation: bigger world = more results
- ▶ examples: union, intersection,  $\exists R.C$ , join, IN conditions
- ▶ counter-examples: complement,  $\forall R.C$ , NOT IN conditions

technically, non-monotone operations in open world dubious



# Primitive Types and Encoding Data

# Primitive Types and Encoding Data

## Motivation

# Data Interoperability

## Situation

- ▶ languages systems focus on different aspects  
frequent need to exchange data
- ▶ generally, lots of aspect/language-specific objects  
proofs, programs, tables, sentences
- ▶ but same/similar **primitive** data types used across systems  
should be easy to exchange

## Problem

- ▶ crossing system barriers usually require interchange language  
serialize as string and reparsing
- ▶ interchange languages typically untyped XML, JSON, YAML, ...

## Solution

- ▶ standardize primitive data types
- ▶ standardize encoding in interchange languages

## Primitive vs. Declared

### Primitive Types

- ▶ built into the language
- ▶ assumed to exist a priori fundamentals of nature
- ▶ fixed semantics (usually interpreted by identity function)

### Triple Structure: 3 kinds of named objects

- ▶ the type eg: 'int'
- ▶ values of the type eg: 0, 1, -1, ...
- ▶ operations on type eg: addition, multiplication, ...

|                 | primitive            | declared       |
|-----------------|----------------------|----------------|
| introduced by   | language designer    | user           |
| introduced in   | grammar              | vocabulary $V$ |
| visible in      | all vocabularies     | $V$ only       |
| semantics given | explicitly           | implicitly     |
| ... by          | translation function | axioms         |

## Examples

### Typical primitive types

- ▶ natural numbers ( $= \mathbb{N}$ )
- ▶ arbitrary precision integers ( $= \mathbb{Z}$ )
- ▶ fixed precision integers (32 bit, 64 bit, ...)
- ▶ floating point (float, double, ...)
- ▶ Booleans
- ▶ characters (ASCII, Unicode)
- ▶ strings

### Observation:

- ▶ essentially the same in every language  
including whatever language used for semantics
- ▶ semantics by translation straight-forward  
watch out for different endiannesses for ints,  
charsets and encodings for characters/strings

## Quasi-Primitive = Declared in standard library

### Standard library

- ▶ present in every language      assumed empty vocabulary by default
- ▶ one fixed vocabulary
  - ▶ implicitly included into every other vocabulary
  - ▶ implicitly fixed by any translation between vocabularies
- ▶ objects technically declared
- ▶ but practically part of primitive objects

### Examples

- ▶ sufficiently expressive languages
  - ▶ push many primitive objects to standard library      never all
  - ▶ simplifies language, especially when defining operations  
strings in C, BigInteger in Java, inductive type for  $\mathbb{N}$
- ▶ inexpressive languages
  - ▶ many primitives      SQL, spreadsheet software
  - ▶ few (quasi)-primitives      few operations available in OWL

## Treatment in this Course

### BOL syntax and semantics so far

- ▶ primitive objects omitted in syntax
- ▶ assumed reasonable collection available
- ▶ assumed same (quasi-)primitive objects in semantic languages  
irrelevant if interpreting primitive objects as primitive or quasi-primitive

largely justified by practical languages

### But what exactly is the standard?

- ▶ will present possible solution
- ▶ uses special ontology language just for specifying primitive objects
  - ▶ name
  - ▶ type
  - ▶ semantics  
typically narrative; alternatively deductive, computational
- ▶ current research, not standard practice

# Encoding Primitive Types

## Problem

- ▶ quickly encounter primitive types not supported by common languages
- ▶ need to encode them using existing types  
typically as strings, ints, or products/lists thereof

## Examples

- ▶ date, time, color, location on earth
- ▶ graph, function
- ▶ picture, audio, video
- ▶ physical quantities (1*m*, 1*in*, etc.)
- ▶ gene, person

Breakout questions: What primitive types do we need for univis?



## Failures of Encodings

### Y2K bug

- ▶ date encoded as tuple of integers, using 2 digits for year
- ▶ needed fixing in year 2000
- ▶ estimated \$300 billion spent to change software
- ▶ possible repeat: in 2038, number of seconds since 1970-01-01 (used by Unix to encode time as integer) overflows 32-bit integers

### Genes in Excel

- ▶ 2016 study found errors in 20% of spreadsheets accompanying genomics journal papers
- ▶ gene names encoded as strings but auto-converted to other types by Excel
  - ▶ "SEPT2" (Septin 2) converted to September 02
  - ▶ REKIN identifiers, e.g., "2310009E13", converted to float  $2.31E + 1$

## Failures of Encodings (2)

### Mars Climate Orbiter

- ▶ two components exchanged physical quantity
- ▶ specification required encoding as number using unit Newton seconds
- ▶ one component used wrong encoding (with pound seconds as unit)
- ▶ led to false trajectory and loss of \$300 million device

### Shellshock

- ▶ bash allowed gaining root access from 1998 to 2014
- ▶ function definitions were encoded as source code
- ▶ not decoded at all; instead, code simply run (as root)
- ▶ allowed appending ";" ... to function definitions

SQL injection similar: complex data encoded as string, no decoding

# Research Goal for Aspect-Independent Data in Tetrapod

## Standardization of Common Data Types

- ▶ Ontology language optimized for declaring types, values, operations  
semantics must exist but can be extra-linguistic
- ▶ Vocabulary declaring such objects  
should be standardized, modular, extensible

## Standardization of Codecs

- ▶ Fixed small set of primitive objects  
should be (quasi-)primitive in every language  
not too expressive, possibly untyped
- ▶ Standard codecs for translating common types to interchange languages

## Codec for type $A$ and int. lang. $L$

- ▶ coding function  $A$ -values  $\rightarrow L$ -objects
  - ▶ partial decoding function  $L$ -objects  $\rightarrow A$ -values
  - ▶ inverse to each other
- in some sense

# Overview

Next steps

1. Data types
2. Data interchange languages
3. Codecs

# Primitive Types and Encoding Data

## Data Types

## Breakout Question

What types do we need?

## Atomic Data Types: basic

### typical in IT systems

- ▶ fixed precision integers (32 bit, 64 bit, ...)
- ▶ IEEE float, double
- ▶ Booleans
- ▶ Unicode characters
- ▶ strings

could be list of characters but usually bad idea

### typical in math

- ▶ natural numbers ( $= \mathbb{N}$ )
- ▶ arbitrary precision integers ( $= \mathbb{Z}$ )
- ▶ rational, real, complex numbers
- ▶ graphs, trees

clear: language must be modular, extensible

## Atomic Data Types: advanced

### general purpose

- ▶ date, time, color, location on earth
- ▶ picture, audio, video

### domain-specific

- ▶ physical quantities (*1m*, *1in*, etc.)
- ▶ gene, person
- ▶ semester, course id, ...

clear: language must be modular, extensible



## Complex Data Types

- ▶ relatively easy if all primitive types atomic      `int`, `string`, etc.
- ▶ but need to allow for complex types

Two kinds

- ▶ type operators: take **only type arguments**, return types
  - ▶ type operator  $\times$
  - ▶ takes two types  $A, B$
  - ▶ returns type  $A \times B$
- ▶ dependent types: take **also data arguments**, return types
  - ▶ dependent type operator *vector*
  - ▶ takes natural number  $n$ , type  $A$
  - ▶ returns type  $A^n$  of  $n$ -tuples over  $A$

dependent types much more complicated, less uniformly used  
harder to standardize

# Collection Data Types

## Homogeneous Collection Types

- ▶ sets
- ▶ multisets (= bags)
- ▶ lists      all unary type operators, e.g. *list A* is type of lists over *A*
- ▶ fixed-length lists (= Cartesian power, vector *n*-tuple)  
dependent type operator

## Heterogeneous Collection Types

- ▶ lists
- ▶ fixed-length lists (= Cartesian power, *n*-tuple)
- ▶ sets
- ▶ multisets (= bags)  
all atomic types, e.g., *list* is type of lists over any objects

## Aggregation Data Types

### Products

- ▶ Cartesian product of some types  $A \times B$   
values are pairs  $(x, y)$     numbered projections  $_1, _2$  — order relevant
- ▶ labeled Cartesian product (= record)  $\{a : A, b : B\}$   
values are records  $\{a = x, b = y\}$   
named projections  $a, b$  — order irrelevant

### Disjoint Unions

- ▶ disjoint union of some types  $A \uplus B$   
values are  $inj_1(x), inj_2(y)$  numbered injections  $_1, _2$  — order relevant
- ▶ labeled disjoint union  $a(A) | b(B)$   
values are constructor applications  $a(x), b(y)$   
named injections  $a, b$  — order irrelevant

labeled disjoint unions uncommon  
but recursive labeled disjoint union = inductive data type

## Subtyping

- ▶ relatively easy if all data types disjoint
- ▶ better with subtyping      open problem how to do it nicely

### Subtyping Atomic Types

- ▶  $\mathbb{N} <: \mathbb{Z}$
- ▶ ASCII <: Unicode

### Subtyping Complex Types

- ▶ covariance subtyping (= vertical subtyping) same for disjoint unions

$$A <: A' \Rightarrow \text{list } A <: \text{list } A'$$

$$A_i <: A'_i \Rightarrow \{\dots, a_i : A_i, \dots\} <: \{\dots, a_i : A'_i, \dots\}$$

- ▶ structural subtyping (= horizontal subtyping)

$$\{a : A, b : B\} :> \{a : A, b : B, c : C\}$$

$$a(A)|b(B) <: a(A)|b(B)|c(C)$$

## A Basic Language for Typed Data

Let BDL be given by

### Types

|         |       |              |         |     |          |     |        |                   |
|---------|-------|--------------|---------|-----|----------|-----|--------|-------------------|
| $T ::=$ | $int$ | $ $          | $float$ | $ $ | $string$ | $ $ | $bool$ | base types        |
|         | $ $   | $list$       | $T$     |     |          |     |        | homogeneous lists |
|         | $ $   | $(ID : T)^*$ |         |     |          |     |        | record types      |
|         | $ $   | $\dots$      |         |     |          |     |        | additional types  |

### Data

|         |                             |                               |                                   |
|---------|-----------------------------|-------------------------------|-----------------------------------|
| $D ::=$ | $(64 \text{ bit integers})$ |                               |                                   |
|         | $ $                         | $(IEEE \text{ double})$       |                                   |
|         | $ $                         | $"(Unicode \text{ strings})"$ |                                   |
|         | $ $                         | $true \mid false$             |                                   |
|         | $ $                         | $D^*$                         | lists                             |
|         | $ $                         | $(ID = D)^*$                  | records                           |
|         | $ $                         | $\dots$                       | constructors for additional types |

## BDL Extended with Named ADTs

|  |                  |
|--|------------------|
| $V ::= D^*$                                      | Vocabularies     |
| $D ::= \mathbf{adt} \ t \ \{\mathbf{ID} : T^*\}$ | ADT definitions  |
| $\mathbf{datum} \ d : T = D$                     | data definitions |

## Types

|               |                          |
|---------------|--------------------------|
| $T ::= \dots$ | as before                |
| $t$           | reference to a named ADT |

## Data

|                            |                            |
|----------------------------|----------------------------|
| $D ::= \dots$              | as before                  |
| $d$                        | reference to a named datum |
| $t\{(\mathbf{ID} = D)^*\}$ | ADT elements               |

# Primitive Types and Encoding Data

## Data Representation Languages

## Overview

### General Properties

- ▶ general purpose or domain-specific
- ▶ typed or untyped
  - typical: Church-typed but no type operators, quasi untyped
- ▶ text or binary serialization
- ▶ libraries for many programming languages
  - ▶ data structures
  - ▶ serialization (data structure to string)
  - ▶ parsing (string to data structure, partial)

### Candidates

- ▶ XML: standard on the web, notoriously verbose
- ▶ JSON: minimal, more human-friendly text syntax
  - newer than XML, very dominant in JS ecosystem
- ▶ YAML: line/indentation-based



## Breakout Question

What is the difference between XML, JSON, and YAML?

# Typical Data Representation Languages

XML, JSON, YAML essentially the same

except for concrete syntax

## Atomic Types

- ▶ integer, float, boolean, string
- ▶ need to read fine-print on precision

## (Not Very) Complex Types

- ▶ heterogeneous lists
- ▶ records

a single type for all lists

a single type for all records

## Example: JSON

JSON:

```
{
  "individual": "FlorianRabe",
  "age": 40,
  "concepts": ["instructor", "male"],
  "teach": [
    {"name": "WuV", credits: 7.5},
    {"name": "KRMT", credits: 5}
  ]
}
```

Weirdnesses:

- ▶ atomic/list/record called basic/array/object
- ▶ record field names are arbitrary strings, must be quoted
- ▶ records use : instead of =

## Example: YAML

inline syntax: same as JSON but without quoted field names

alternative: indentation-sensitive syntax

```
individual: "Florian Rabe"
```

```
age: 40
```

```
concepts:
```

```
  - instructor
```

```
  - male
```

```
teach:
```

```
  - name: WuV
```

```
    credits: 7.5
```

```
  - name: KRMT
```

```
    credits: 5
```

Weirdnesses:

- ▶ atomic/list/record called scalar/collection/structure
- ▶ records use : instead of =

## Example: XML

Weird structure but very similar

- ▶ elements both record (= attributes) and list (= children)
- ▶ elements carry name of type (= tag)

```
<person individual="Florian Rabe" age="40">  
  <concepts>  
    <concept>instructor </concept>  
    <concept>male</concept>  
  </concepts>  
  <teach>  
    <course name="WuV" credits="7.5" />  
    <course name="KRMT" credits="5" />  
  </teach>  
</person>
```

- ▶ Good: `<person>`, `<course>`, `<concept>` give type of object  
easier to decode
- ▶ Bad: values of record fields must be string  
concepts cannot be given in attribute  
integers, Booleans, whitespace-separated lists coded as strings

# Structure Sharing

## Problem

- ▶ Large objects are often redundant specially when machine-produced
- ▶ Same string, URL, mathematical object occurs in multiple places
- ▶ Handled in memory via pointers
- ▶ Size of serialization can explode

## Solution 1: in language

- ▶ Add definitions to language common part of most languages anyway
- ▶ Users should introduce name whenever object used twice
- ▶ Problem: only works if
  - ▶ duplication anticipated
  - ▶ users introduced definition
  - ▶ duplication within same context

structure-sharing most powerful if across contexts

## Structure Sharing (2)

### Solution 2: in tool

- ▶ Use factory methods instead of constructors
- ▶ Keep huge hash set of all objects
- ▶ Reuse existing object if already in hash set
- ▶ Advantages
  - ▶ allows optimization
  - ▶ transparent to users
- ▶ Problem: only works if
  - ▶ for immutable data structures
  - ▶ if no occurrence-specific metadata e.g., source references

### In data representation language

- ▶ Allow any subobject to carry identifier
- ▶ Allow identifier references as subobjects
  - allows preserving structure-sharing in serialization

supported by XML, YAML

# Primitive Types and Encoding Data

## Codecs



## General Definition

Throughout this section, we fix a data representation language  $L$ .

$L$ -words called codes

Given a data type  $T$ , a codec for  $T$  consists

- ▶ coding function:  $c : T \rightarrow L$
- ▶ partial decoding function:  $d : L \rightarrow^? T$
- ▶ such that

$$d(c(x)) = x$$

## Codec Operators

Given a data type operator  $T$  taking  $n$  type arguments,  
a codec operator  $C$  for  $T$

- ▶ takes  $n$  codecs  $C_i$  for  $T_i$
- ▶ returns a codec  $C(C_1, \dots, C_n)$  for  $T(T_1, \dots, T_n)$

## Exercise 4

We fix strings as the data representation language  $L$ .

Then,

1. Jointly specify
  - ▶ additional BDL types and constructors for univis-specific data
  - ▶ codecs and codec operators for all types resp. type operators
2. Individually, in any programming language, implement
  - ▶ data structures for BDL
  - ▶ string codecs (operators) for all BDL base types (operators)
3. Use your codecs to exchange example data with your fellow students, who used different implementations and different programming languages.

## Codecs for Base Types

We define codecs for the base types using strings as the data representation language  $L$ .

Easy cases:

- ▶ `StandardFloat`: as specified in IEEE floating point standard
- ▶ `StandardString`: as themselves
- ▶ `StandardBool`: as *true* or *false*
- ▶ `StandardInt` (64-bit): decimal digit-sequences as usual

## Breakout Question

How to encode unlimited precision integers?

## Codecs for Unlimited Precision Integers

Encode  $z \in \mathbb{Z}$

- ▶  $L$  is strings: decimal digit sequence as usual
- ▶  $L$  is JSON:
  - ▶ IntAsInt: decimal digit sequence as usual  
JSON does not specify precision  
but target systems may get in trouble
  - ▶ IntAsString: string containing decimal digit sequence  
safe but awkward
  - ▶ IntAsDecList: list of decimal digits  
safe but awkward
  - ▶ IntAsList1: as list of digits for base  $2^{64}$   
OK, but we can do better
  - ▶ IntAsList2: as list of
    - ▶ integer for the number of digits, sign indicate sign of  $z$
    - ▶ list of digits of  $|z|$  for base  $2^{64}$

Question: Why is this smart?

## Codecs for Unlimited Precision Integers

Encode  $z \in \mathbb{Z}$

- ▶  $L$  is strings: decimal digit sequence as usual
- ▶  $L$  is JSON:
  - ▶ IntAsInt: decimal digit sequence as usual  
JSON does not specify precision  
but target systems may get in trouble
  - ▶ IntAsString: string containing decimal digit sequence  
safe but awkward
  - ▶ IntAsDecList: list of decimal digits  
safe but awkward
  - ▶ IntAsList1: as list of digits for base  $2^{64}$   
OK, but we can do better
  - ▶ IntAsList2: as list of
    - ▶ integer for the number of digits, sign indicate sign of  $z$
    - ▶ list of digits of  $|z|$  for base  $2^{64}$

Question: Why is this smart?

Can use lexicographic ordering for size comparison

## Codecs for Lists

Encode list  $x$  of elements of type  $T$

- ▶  $L$  is strings: e.g., comma-separated list of  $T$ -encoded elements of  $x$   
(un)escape comma in  $T$ -encoded elements
- ▶  $L$  is JSON:
  - ▶ ListAsString: like for strings above
  - ▶ ListFromArray: JSON array of  $T$ -encoded elements of  $x$



## Additional Types

Example: semester

Extend BDL:

Types

$T ::= \text{Sem}$                       semester

Data

$D ::= \text{sem}(\text{int}, \text{bool})$     i.e., year + summer?

Define standard codec:

$\text{sem}(y, \text{true}) \rightsquigarrow \text{"SSY"}$

$\text{sem}(y, \text{false}) \rightsquigarrow \text{"WSY"}$

where  $Y$  is encoding of  $y$

## Additional Types (2)

Example: timestamps

Extend BDL:

Types

$T ::= \text{timestamp}$

Data

$D ::= (\text{productions for dates, times, etc.})$

Standard codec: encode as string as defined in ISO 8601

# Primitive Types and Encoding Data

## Data Interchange

## Design

### 1. Specify type system, e.g., BDL

- ▶ types
- ▶ constructors
- ▶ operations

can be done in appropriate type theory

### 2. Pick data representation language $L$

### 3. Specify codecs for type system and $L$

- ▶ at least one codec per base type
- ▶ at least one codec operator per type operator

on paper

### 4. Every system implements

- ▶ type system (as they like) typically aspect-specific constraints
- ▶ codecs as specified
- ▶ function mapping types to codecs

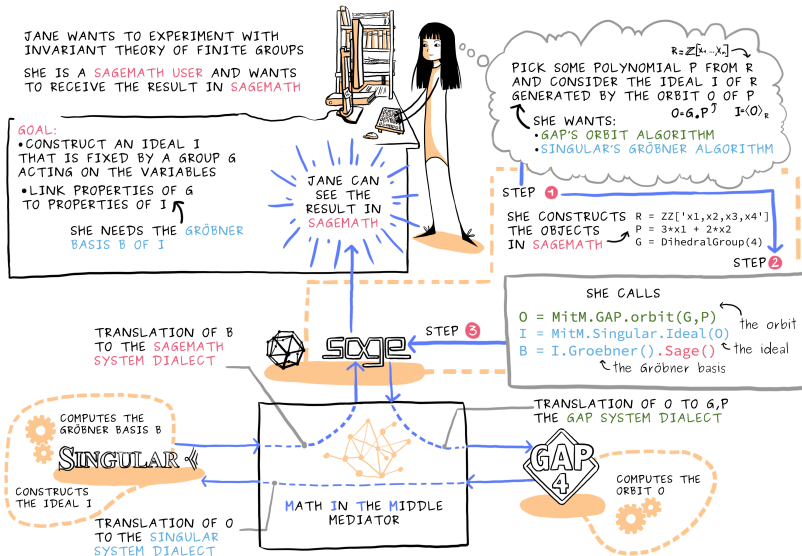
### 5. Systems can exchange data by encoding-decoding

type-safe because codecs chosen by type

## Example

Implementation in Scala part of course resources

# Example Application: OpenDreamKit research project



## Integrating BOL and BDL

### OWL-near option

- ▶ use BDL to define the primitive types of BOL
- ▶ use those as types of BOL properties
- ▶ Curry-typing throughout easy: just merge the grammars

### SQL-near option

- ▶ use BDL to define the primitive types of BOL
- ▶ also add ADTs
- ▶ Church typing more prominent  
open question: ADTs in addition to or instead of BOL concepts

We assume the latter for now without spelling out the details.

## BDL-Mediated Interoperability

### Idea

- ▶ define data types in BDL or similar typed ontology language
- ▶ use ADTs
- ▶ generate corresponding
  - ▶ class definitions for programming languages PL
  - ▶ table definitions in SQL
- ▶ use codecs to convert automatically when interchanging data between PL and SQL

one class per ADT

one table per ADT

### Open research problem

no shiny solution yet that can be presented in lectures



## Codecs in ADT Definitions

SQL table schema = list of fields where field is

- ▶ name
- ▶ type only types of database supported

BDL semantic table schema = list of fields where field is

- ▶ name
- ▶ type  $T$  of **type system** independent of database
- ▶ codec for  $T$  using primitive objects of database as codes  
see research paper [https://kwarc.info/people/frabe/Research/WKR\\_](https://kwarc.info/people/frabe/Research/WKR_)

Codec could be chosen automatically, but we want to allow multiple users a choice of codecs for the same type.

## Example

Ontology based on BDL-ADTs with additional codec information:

```
schema Instructor
```

|          |             |       |                           |
|----------|-------------|-------|---------------------------|
| name:    | string      | codec | StandardString            |
| age:     | int         | codec | StandardInt               |
| courses: | list Course | codec | CommaSepList CourseAsName |

```
schema Course
```

|           |          |       |                  |
|-----------|----------|-------|------------------|
| name:     | string   | codec | StandardString   |
| credits:  | float    | codec | StandardFloat    |
| semester: | Semester | codec | SemesterAsString |

Generated SQL tables:

```
CREATE TABLE Instructor
```

```
(name string , age int , courses string)
```

```
CREATE TABLE Course
```

```
(name string , credits float , semester string)
```

## Open Problem: Non-Compositionality

### Sometimes optimal translation is non-compositional

- ▶ example translates *list*-type in ADT to comma-separated string in DB
- ▶ better break up *list B* fields in type *A* into separate table with columns for *A* and *B*

### Similar problems

- ▶ a pair type in an ADT could be translated to two separate columns
- ▶ an option type in an ADT could be translated to a normal column using SQL's NULL value

## Open Problem: Querying

- ▶ General setup  $L = \text{SQL}$ 
  - ▶ write SQL-style queries using BDL
  - ▶ automatically encode values when writing to database from PL
  - ▶ automatically decode query results when reading from DB
- ▶ But queries using semantic operations cannot always be translated to DB
  - ▶ operation  $IsSummer : Semester \rightarrow bool$  in BDL
  - ▶ query `SELECT * FROM course WHERE  $IsSummer(semester)$`
  - ▶ how to map  $IsSummer$  to SQL?
- ▶ Ontology operations need commuting operations on codes
  - ▶ given  $f : T_1 \rightarrow T_2$  in BDL and respective codecs  $C_{1,2}$ ,
  - ▶  $L$ -function  $f'$  commutes with  $f$  iff

$$\begin{array}{ccc} T_1 & \xrightarrow{f} & T_2 \\ C_1.\text{encode} \downarrow & & \uparrow C_2.\text{decode} \\ L & \xrightarrow{f'} & L \end{array}$$

## Exercise 5, part 1

We build on the implementation of BDL and codecs from Exercise 4 and on the database schemas from Exercise 3.

1. Extend the implementation to BDL+ADT (see Slide 101).
2. Extend
  - ▶ codecs and codec operators with identifiers  $l ::= (\text{strings})$
  - ▶ ADT fields with codec expressions  $c ::= l \mid l(c_1 \dots, c_n)$and write a function that maps  $c$  to the corresponding codec.

## Exercise 5, part 2

3. Write a function that takes a vocabulary (= a list of ADT definitions with codec expressions) and generates an SQL schema for it. Use the type returned by the codec as the database type.
4. Write a function that takes an element  $d$  of an ADT and generates the SQL (or CSV) representation of  $d$  with all field values encoded by the corresponding codec.
5. Write a function that takes an ADT name and a SQL or CSV object and applies decoding to build the corresponding ADT element.
6. Test this by
  - ▶ writing some of your univis table schemas as ADTs and some example values as ADT elements,
  - ▶ exchanging these with a database and/or via CSV with fellow students' implementations.

# Querying

# Querying

## Overview



## General Ideas

- ▶ Recall
  - ▶ syntax = context-free grammar
  - ▶ semantics = translation to another language
- ▶ Example: BOL translated to SQL, SFOL, Scala, English
- ▶ Querying = use semantics to answer questions about syntax

Note:

- ▶ Not the standard definition of querying
- ▶ Design of a new Tetrapod-level notion of querying
  - ongoing research
- ▶ Subsumes concepts of different names from the various aspects

# Propositions

syntax with propositions =  
designated non-terminals for propositions

Examples:

| aspect               | basic propositions                       |
|----------------------|--|
| ontology language    | assertions, concept equality/subsumption |
| programming language | equality for some types                  |
| database language    | equality for base types                  |
| logic                | equality for all types                   |
| natural language     | sentences                                |

Aspects vary critically in how propositions can be formed

- ▶ any program in computation
- ▶ quantifiers in deductions
- ▶  $\exists$  in databases

undecidable

# Propositions as Queries

Propositions allow defining queries

|                | Query                           | Result                |
|----------------|---------------------------------|-----------------------|
| deduction      | proposition                     | yes/no                |
| concretization | proposition with free variables | true ground instances |
| computation    | term                            | value                 |
| narration      | question                        | answer                |

# Semantics of Propositions

syntax with propositions =  
designated non-terminals for propositions

needed to ask queries

semantics with theorems =  
designates some propositions as theorems or contradictions

needed to answer queries

Note:

- ▶ A propositions may be neither theorem nor contradiction.
- ▶ We say that language has negation if:  
 $F$  theorem iff  $\neg F$  contradiction and vice versa.

We write  $\vdash F$  if  $F$  is theorem.

# Querying

## Deductive Queries

## Definition

We assume

- ▶ a semantics  $\llbracket - \rrbracket$  from  $I$  to  $L$
- ▶  $I$  has propositions
- ▶ there is an operation  $\text{True}$  that maps translations of  $I$ -propositions to  $L$ -propositions
- ▶  $L$  has semantics with propositions

We define

- ▶ a deductive query is an  $I$ -proposition  $p$
- ▶ the result is
  - ▶ yes if  $\text{True}[\llbracket p \rrbracket]$  is a theorem of  $L$
  - ▶ no if  $\text{True}[\llbracket p \rrbracket]$  is a contradiction in  $L$

## Breakout question

What can go wrong?

## Problem: Inconsistency

In general, (in)consistency of semantics

- ▶ Some propositions may be both a theorem and a contradiction.
- ▶ In that case, queries do not have a result.

In practice, however:

- ▶ If this holds for some propositions, it typically holds for all of them.
- ▶ In that, we call  $L$  inconsistent.
- ▶ We usually assume  $L$  to be consistent.



## Problem: Incompleteness

In general, (in)completeness of semantics

- ▶ We cannot in general assume that every proposition in  $L$  is either a theorem or a contradiction.
- ▶ In fact, most propositions are neither.
- ▶ So, queries do not necessarily have a result.
- ▶ We speak of incompleteness.

Note: not the same as the usual (in)completeness of logic

In practice, however:

- ▶ It may be that  $L$  is complete for all propositions in the image of  $\text{True}[\![ - ]\!]$ .
- ▶ This is the case if  $L$  is simple enough  
typical for ontology languages

## Problem: Undecidability

In general, (un)decidability of semantics:

- ▶ We cannot in general assume that it is decidable whether a proposition in  $L$  is a theorem or a contradiction.
- ▶ In fact, it usually isn't.
- ▶ So, we cannot necessarily compute the result of a query.
- ▶ However: If we have completeness, decidability is likely.

run provers for  $F$  and  $\neg F$  in parallel

In practice, however:

- ▶ It may be that  $L$  is decidable for all propositions in the image of  $\text{True}[\![ - ]\!]$ .
- ▶ This is the case if  $I$  is simple enough

typical for ontology languages

## Problem: Inefficiency

In general, (in)efficiency of semantics:

- ▶ Answering deductive queries is very slow.
- ▶ Even if we are complete and decidable.

In practice, however:

- ▶ Decision procedures for the image of  $\text{True}[\![ - ]\!]$  may be quite efficient.
- ▶ Dedicated implementations for specific fragments.
- ▶ This is the case if  $I$  is simple enough  
typical for ontology languages

# Querying

## Contexts and Free Variables

## Concepts

Recall the analogy between grammars and typing:

| grammars                             | typing                         |
|--------------------------------------|--------------------------------|
| non-terminal                         | type                           |
| production                           | constructor                    |
| non-terminal on left of production   | return type of constructor     |
| non-terminals on right of production | arguments types of constructor |
| terminals on right of production     | notation of constructor        |
| words derived from non-terminal $N$  | expressions of type $N$        |

We will now add contexts and substitutions.

## Contexts

Given a context-free language  $I$ , we define:

- ▶ A **context**  $\Gamma$  is of the form  $x_1 : N_1, \dots, x_n : N_n$  where the
  - ▶  $x_i$  are names
  - ▶  $N_i$  are non-terminals

We write this as  $\vdash_I \Gamma$ .

- ▶ A **substitution** for  $\Gamma$  is of the form  $x_1 := w_1, \dots, x_n := w_n$  where the
  - ▶  $x_i$  are as in  $\Gamma$
  - ▶  $w_i$  derived from the corresponding  $N_i$

We write this as  $\vdash_I \gamma : \Gamma$ .

- ▶ An **expression in context**  $\Gamma$  of type  $N$  is a word  $w$  derived from  $N$  using additionally the productions  $N_i ::= x_i$ .

We write this as  $\Gamma \vdash_I w : N$ .

- ▶ Given  $\Gamma \vdash w : N$  and  $\vdash \gamma : \Gamma$  as above, the **substitution of**  $\gamma$  in  $w$  is obtained by replacing every  $x_i$  in  $w$  with  $w_i$ . We write this as  $w[\gamma]$ .

## Contexts under Compositional Translation

Consider a compositional semantics  $\llbracket - \rrbracket$  from  $I$  to  $L$  between context-free languages.

- ▶ Every  $\vdash_I w : N$  is translated to some  $\vdash_L \llbracket w \rrbracket : N'$  for some  $N'$ .
- ▶ Compositionality ensures that  $N'$  is the same for all  $w$  derived from  $N$ .
- ▶ We write  $\llbracket N \rrbracket$  for that  $N'$ .
- ▶ Then we have

$$\vdash_I w : N \quad \text{implies} \quad \vdash_L \llbracket w \rrbracket : \llbracket N \rrbracket$$

Now we translate contexts, substitutions, and variables as well:

$$\llbracket x_1 : N_1, \dots, x_n : N_n \rrbracket := x_1 : \llbracket N_1 \rrbracket, \dots, x_n : \llbracket N_n \rrbracket$$

$$\llbracket x_1 := w_1, \dots, x_n := w_n \rrbracket := x_1 := \llbracket w_1 \rrbracket, \dots, x_n := \llbracket w_n \rrbracket$$

$$\llbracket x \rrbracket := x$$

Then we have

$$\Gamma \vdash_I w : N \quad \text{implies} \quad \llbracket \Gamma \rrbracket \vdash_L \llbracket w \rrbracket : \llbracket N \rrbracket$$

## Substitution under Compositional Translation

From previous slide:

$$\llbracket x_1 : N_1, \dots, x_n : N_n \rrbracket := x_1 : \llbracket N_1 \rrbracket, \dots, x_n : \llbracket N_n \rrbracket$$

$$\llbracket x_1 := w_1, \dots, x_n := w_n \rrbracket := x_1 := \llbracket w_1 \rrbracket, \dots, x_n := \llbracket w_n \rrbracket$$

$$\llbracket x \rrbracket := x$$

$$\Gamma \vdash_I w : N \quad \text{implies} \quad \llbracket \Gamma \rrbracket \vdash_L \llbracket w \rrbracket : \llbracket N \rrbracket$$

We can now restate the substitution theorem as follows:

$$\llbracket E[\gamma] \rrbracket = \llbracket E \rrbracket \llbracket \llbracket \gamma \rrbracket \rrbracket$$



# Querying Concretized Queries

## Definition

We assume

- ▶ as for deductive queries
- ▶ semantics must be compositional

We define

- ▶ a concretized query is an  $l$ -proposition  $p$  in context  $\Gamma$
- ▶ a **single** result is a
  - ▶ a substitution  $\vdash_I \gamma : \Gamma$
  - ▶ such that  $\vdash_L \text{True}[\![p[\gamma]]\!]$
- ▶ the **result set** is the set of all results

## Example

1. BOL ontology:

*concept male, concept person, axiom male  $\sqsubseteq$  person,*  
*individual FlorianRabe, assertion FlorianRabe isa male*

2. Query  $x : \textit{individual} \vdash_{BOL} x \text{ is-a } \textit{person}$
3. Translation to SFOL:  $x : \iota \vdash_{SFOL} \textit{person}(x)$
4. SFOL calculus yields theorem  $\vdash_{SFOL} \textit{person}(\textit{FlorianRabe})$
5. Query result  $\llbracket \gamma \rrbracket = x := \textit{FlorianRabe}$
6. Back-translating the result to BOL:  $\gamma = x := \textit{FlorianRabe}$   
back translation is deceptively simple:  
translates SFOL-constant to BOL-individual of same name

## Breakout question

What can go wrong?

## Problem: Open World

In general, semantics uses open world:

- ▶ open world: result contains **all known** results  
same query might yield more results later
- ▶ closed world: result set contains **all** results

always relative to concrete database for  $L$

In practice, however,

- ▶ system explicitly assumes closed world    typical for databases
- ▶ users aware of open world and able to process results correctly

## Problem: Infinity of Results

In general, there may be infinitely many results:

- ▶ e.g., query for all  $x$  such that  $\vdash x$ ,

In practice, however,

- ▶ systems pull results from finite database e.g., SQL, SPARQL
- ▶ systems enumerate results, require user to explicitly ask for more e.g., Prolog

## Problem: Back-Translation of Results

In general,  $\llbracket - \rrbracket$  may be non-trivial to invert

- ▶ easy to obtain  $\llbracket p \rrbracket$  in context  $\llbracket \Gamma \rrbracket$       just apply semantics
- ▶ possible to find substitutions

$$\vdash_L \delta : \llbracket \Gamma \rrbracket \quad \text{where} \quad \llbracket \Gamma \rrbracket \vdash_L \text{True}[\llbracket p \rrbracket][\delta]$$

easiest case: just look them up in database

- ▶ but how to translate  $\delta$  to  $l$ -substitutions  $\gamma$  with

$$\vdash_l \gamma : \Gamma \quad \text{where} \quad \llbracket \Gamma \rrbracket \vdash_L \text{True}[\llbracket p[\gamma] \rrbracket]$$

substitution theorem: pick such that  $\llbracket \gamma \rrbracket = \delta$

the more  $\llbracket - \rrbracket$  does, the harder to invert

In practice, however:

- ▶ often only interested in concrete substitutions
- ▶ translation of concrete data usually identity

But: practice restricted to what works even if more is needed

# Querying

## Computational Queries



## Definition

We assume

- ▶ the same as for deductive queries
- ▶ semantics has equality/equivalence  $\doteq$

We define

- ▶ a computational query is an  $l$ -expression  $e$
- ▶ the result is an  $l$ -expression  $e'$  so that  $\vdash_L \llbracket e \rrbracket \doteq \llbracket e' \rrbracket$

intuition:  $e'$  is the result of evaluating  $e$

If semantics is compositional,  $e$  may contain free variables

evaluate to themselves

## Problem: Back-Translation of Results

In general,  $\llbracket - \rrbracket$  may be non-trivial to invert

- ▶ easy to obtain  $E := \llbracket e \rrbracket$
- ▶ possible to find  $E'$  with  $\vdash_L E' \doteq E$  by working in the semantics
- ▶ non-obvious how to obtain  $e'$  such that  $\llbracket e' \rrbracket = E'$

In practice, however:

- ▶ evaluation meant to simplify, i.e., only useful if  $E'$  very simple
- ▶ simple  $E'$  usually in the image of  $\llbracket - \rrbracket$
- ▶ typical case:  $E'$  is concrete data and  $e' = E'$       called a value

## Problem: Non-Termination

In general, computation of  $E'$  from  $E$  might not terminate

- ▶ while-loops
- ▶ recursion
- ▶  $(\lambda x.x\ x)(\lambda x.x\ x)$  with  $\beta$ -rule
- ▶ simplification rule  $x \cdot y \rightsquigarrow y \cdot x$

similar: distributivity, associativity

In practice, however:

- ▶ image of  $\llbracket - \rrbracket$  part of terminating fragment

But: if  $I$  is Turing-complete or undecidable, general termination not possible

## Problem: Lack of Confluence

In general, there may be multiple  $E'$  that are simpler than  $E$

- ▶ there may be multiple rules that apply to  $E$
- ▶ e.g.,  $f(g(x))$ 
  - ▶ call-by-value: first simplify  $g(x) \rightsquigarrow y$ , then  $f(y) \rightsquigarrow z$
  - ▶ call-by-name: first plug  $g(x)$  into definition of  $f$ , then simplify
- ▶ Normal vs. canonical form
  - ▶ normal:  $\vdash_L E \doteq E'$
  - ▶ canonical: normal and  $\vdash_L E_1 \doteq E_2$  iff  $E'_1 = E'_2$ 
    - equivalent expressions have identical evaluation
    - allows deciding equality

In practice, however:

- ▶ image of  $\llbracket - \rrbracket$  part of confluent fragment
- ▶ typical: evaluation to a value is canonical form
  - works for BDL-types but not for, e.g., function types

# Querying Narrative Queries

## Definition

We assume

- ▶ semantics into natural language

We define

- ▶ a narrative query is an  $L$ -question about some  $I$ -expressions
- ▶ the result is the answer to the question

## Problem: Unimplementable

very expressive = very difficult to implement

- ▶ Natural language understanding
  - ▶ no implementable syntax of natural language  
needs restriction to controlled natural language
  - ▶ specifying semantics hard even when controlled
- ▶ Knowledge base for question answering needed
  - ▶ very large  
must include all common sense
  - ▶ might be inconsistent  
common sense often is
  - ▶ finding answers still very hard

In practice, however:

- ▶ accept unreliability  
attach probability measures to answers
- ▶ implement special cases  
e.g., lookup in databases like Wikidata
- ▶ search knowledge base for related statements  
Google, Watson

Querying

Syntactic Querying



# Search

- ▶ “search” not systematically separated from “querying”
  - ▶ often interchangeable
  - ▶ querying tends to imply formal languages for queries with well-specified semantics e.g., SQL
  - ▶ search tends to imply less targeted process e.g., Google
- we will not distinguish between the two

# Syntactic vs. Semantic Querying

## Semantic querying

- ▶ Query results specified by vocabulary  $V$  but (usually) not contained in it
- ▶ Query answered using semantics of language
- ▶ Challenge: apply semantics to find results
  - ▶ deductive query  $\vdash f : \text{prop}$  requires theorem prover
  - ▶ computation query  $\vdash e : E$  requires evaluator
  - ▶ concrete query  $\Gamma \vdash f : \text{prop}$  requires enumerating all substitutions, running theorem prover/evaluator on all of them

what we've looked at so far

## Syntactic querying

- ▶ Query is an expression  $e$
- ▶ Result is set of occurrences of  $e$  in  $V$
- ▶ Independent of semantics
- ▶ Much easier to realize

## Challenges for Syntactic Search

Easier to realize → scale until new challenges arise

- ▶ large vocabularies
  - ▶ narrative: all text documents in a domain  
e.g., all websites, all math papers
  - ▶ deductive: large repositories of formalization in proof assistants  
10<sup>5</sup> theorems
  - ▶ computational: package managers for all programming languages
  - ▶ concrete: all databases in a domain  
TBs routine
- ▶ incremental indexing: reindex only new/changed parts
- ▶ incremental search to handle large result sets  
pagination
- ▶ sophisticated techniques for
  - ▶ indexing: to allow for fast retrieval
  - ▶ similarity: to select likely results
  - ▶ quality: to rank selected results
- ▶ integration of some semantic parts

## Overview

- ▶ Deduction
  - ▶ semantic: theorem proving called search
  - ▶ syntactic: text search
- ▶ Concretization
  - ▶ semantic: complex query languages (nestable queries)  
SQL, SPARQL
  - ▶ syntactic: search by identifier (linked data)
- ▶ Computation
  - ▶ semantic: interpreters called execution
  - ▶ mixed: IDEs search for occurrences, dependencies
  - ▶ syntactic: search in IDE, package manager
- ▶ Narration:
  - ▶ semantic: very difficult
  - ▶ syntactic: bag of words search

## Abstract Definition: Document

### **Document** =

- ▶ file or similar resource that contains vocabularies
- ▶ often with comments, metadata
- ▶ different names per aspect
  - ▶ deduction: formalization, theory, article
  - ▶ computation: source files
  - ▶ concretization: database, ontology ABox
  - ▶ narrative: document, web site

### **Library** =

- ▶ collection of documents
- ▶ usually structured into folders, files or similar
- ▶ often grouped by user access e.g., git repository
- ▶ vocabularies interrelated within and across libraries

## Abstract Definition: Document Fragment

**Fragment** = subdivision of documents into nested semantic units

Examples

- ▶ deductive: theory, section, theorem, definition, proof step, etc.
- ▶ computational: class, function, command, etc.
- ▶ concrete: table, row, cell
- ▶ narrative: section, paragraph, etc.

Assign unique **fragment URI**, e.g., LIB/DOC?FRAG where

- ▶ LIB: base URI of library e.g., repository URL
- ▶ DOC: path to document within library e.g., folder structure, file name
- ▶ FRAG: id of fragment within document e.g., class name/method name

## Abstract Definition: Index(er)

**Indexer** consists of

- ▶ data structure  $O$  for indexable objects  
aspect-specific index design  
e.g., words, syntax trees
- ▶ function that maps library to index  
the indexing

**Index entry** consists of

- ▶ object that occurred in the library
- ▶ URI of the containing fragment
- ▶ information on where in the fragment it was found

**Index** = set of index entries

## Abstract Definition: Query and Result

Given

- ▶ indexer  $I$  with data structure  $O$
- ▶ set of libraries
- ▶ union of their indexes computed once, queried often

**Query** = object  $\Gamma \vdash^I q : O$

**Result** consists of

- ▶ index entry with object  $o$
- ▶ substitution for  $\Gamma$  such that  $q$  matches  $o$   
definition of “match” index-specific, e.g.,  $q[\gamma] = o$

**Result set** = set of all results in the index



## Bag of Words Search

### Definition:

- ▶ Index data structure = sequences of words (n-grams) up to a certain length
- ▶ Query = bag of words bag = multiset
- ▶ Match: (most) words in query occur in same n-gram or n-grams near each other

### Example implementations

- ▶ internet search engines for websites
- ▶ Elasticsearch: open source engine for custom vocabularies

### Mostly used for narrative documents

- ▶ can treat concrete values as words e.g., numbers
- ▶ could treat other expressions as words works badly

# Symbolic Search

## Definition:

- ▶ Index data structure = syntax tree (of any grammar) of expressions  $o$  with free/bound variables
- ▶ Query = expression  $q$  with free (meta-)variables
- ▶ Match:  $q[\gamma] =_{\alpha} o$ , i.e., up to variable renaming

## Example implementation

- ▶ MathWebSearch  
see separate slides on MathWebSearch in the repository

## Mostly used for formal documents

- ▶ deductive
- ▶ computational

# Knowledge Graph Search

## Definition:

- ▶ Index data structure = assertion forming node/edge in a knowledge graph
- ▶ Index = big knowledge graph  $G$
- ▶ Query = knowledge graph  $g$  with free variables
- ▶ Match:  $g[\gamma]$  is part of  $G$

## Example implementations

- ▶ SPARQL engines without consequence closure  
i.e., the most common case in practice
- ▶ graph databases

Mainly used for ABoxes of untyped ontologies

## Value Search

Definition:

- ▶ Index data structure = BDL values  $v$
- ▶ Query = BDL expression  $q$  with free variables
- ▶ Match:  $q[\gamma] = v$

Example implementations

- ▶ no systematic implementation yet
- ▶ special cases part of most database systems

Could be used for values occurring in any document

- ▶ all aspects
- ▶ may need to decode/encode before putting in index

## Cross-Aspect Occurrences

### Observation

- ▶ libraries are written in one primary aspect
- ▶ indexer focuses on one aspect and kind of object
- ▶ but documents may contain indexable objects of any index

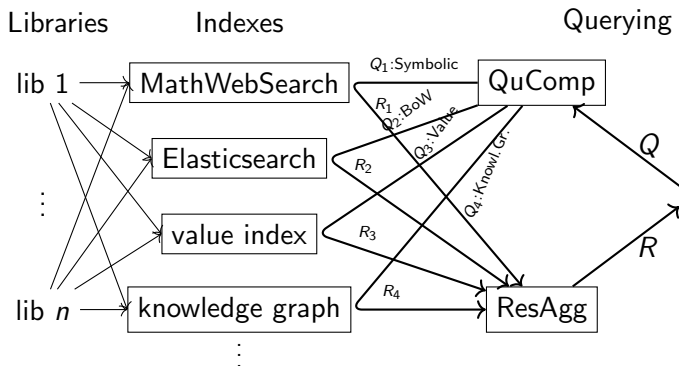
## Cross-Aspect Occurrences: Examples

- ▶ Any library can contain
  - ▶ metadata on fragments
    - ▶ relation assertions induce knowledge graph structure between fragments
    - ▶ property assertions contain values narrative, symbolic objects, or values
  - ▶ cross-references to fragments of any other library
  - ▶ narrative comments
- ▶ Narrative text may contain symbolic expressions  
STEM documents
- ▶ Database table may have columns containing
  - ▶ text
  - ▶ encoded BDL values
  - ▶ symbolic expression (often as strings)
- ▶ Symbolic fragments may contain database references  
e.g., when using database for persistent memoization

## A New Indexing Design

recent paper [https://kwarc.info/people/frabe/Research/BKR\\_mdql\\_20.pdf](https://kwarc.info/people/frabe/Research/BKR_mdql_20.pdf)

with K. Bercic



## A New Indexing Design (2)

Tricky question: What is the query language that allows combining queries for each index?

Easy:

- ▶ query = conjunction of atomic queries
- ▶ each atom queries one index
- ▶ QuComp splits into atoms
- ▶ ResAgg take intersection of results

Better: allow variables to be shared across atoms

open research question



## A New Indexing Design: Example

Consider

- ▶ table of graphs with human-recognizable names and arc-transitivity property indexed into
  - ▶ value index for the graph sparse6 codec
  - ▶ Boolean computed property for the arc-transitivity in knowledge graph
  - ▶ text index for name
- ▶ papers from arXiv in narrative index indexed into
  - ▶ narrative index for text
  - ▶ MathWebSearch for formulas
  - ▶ knowledge graph for metadata

Query goal: find arc-transitive graphs mentioned by name in articles with h-index greater than 50

## Integrating Semantic Querying

### Word search

- ▶ find multi-meaning words for only one meaning  
“normal” in math
- ▶ special treatment of certain queries e.g., “weather” in Google

### Symbolic search

- ▶ match query  $e \doteq e'$  against occurrence  $e' \doteq e$
- ▶ similarly: associativity, commutativity, etc.
- ▶ slippery slope to deductive queries

### Value search

- ▶ match query 1.5 against interval  $1.4 \pm 0.2$
- ▶ match query  $5 \cdot x$  against 25
- ▶ slippery slope to computational queries

frontiers of research — in our group: for STEM documents

# Semantics

# Semantics

## Overview

# Motivation

Recall:

| Syntax    | Data      |
|-----------|-----------|
| Semantics | Knowledge |

Representing

- ▶ syntax = formal language
  - ▶ grammar context-free part
  - ▶ type system context-sensitive well-formedness
- ▶ data = words in the syntax
  - ▶ set of vocabularies
  - ▶ set of typed expressions for each vocabulary
- ▶ semantics = ???
- ▶ knowledge = emergent property of having well-formed words with semantics

So far: semantics by translation

# Relative Semantics

## Semantics by Translation

- ▶ Two syntaxes
  - ▶ object-language  $I$  e.g., BOL
  - ▶ meta-language  $L$  e.g., SFOL, Scala, SQL, English
- ▶ Semantics of  $L$  assumed fixed captures what we already know
- ▶ Semantics of  $I$  by translation into  $L$   
semantics of  $I$  relative to existing semantics of  $L$

Problem: just kicking the can?

# Discussion of Relative Semantics

## Advantages

- ▶ a few meta-languages yield semantics for many languages
- ▶ easy to develop new languages
- ▶ good connection between syntax and semantics via compositionality, substitution theorem

## Disadvantages

- ▶ does not solve the problem once and for all
- ▶ impractical without implementation of semantics of meta-language
- ▶ meta-languages typically much more expressive than needed for object-languages
- ▶ translations can be difficult, error-prone

Also needed: absolute semantics

## Absolute vs. Relative Semantics

Absolute = self-contained, no use of meta-language  $L$

### Get off the ground

- ▶ semantics for a few important meta-languages  
e.g., FOL, assembly language, set theory
- ▶ relative semantics for all other languages, e.g.,
  - ▶ model theory: logic  $\rightarrow$  set theory
  - ▶ compilation: Scala  $\rightarrow$  JVM  $\rightarrow$  assembly

### Redundant semantics

- ▶ common to give
  - ▶ relative and absolute semantics for same syntax
  - ▶ multiple relative semantics      translations to different aspects
  - ▶ sometimes even maybe multiple absolute ones
- ▶ Allows understanding syntax from multiple perspectives
- ▶ Allows cross-checking      show equivalence of two semantics



## No Perfect Model for Absolute Semantics

- ▶ Machine-actionable requires reduction to finite set of rules  
whatever a rule is
- ▶ Does not work for most domains
  - ▶ practical argument: any practically interesting system has too many rules  
cf. physics, e.g., three-body problem already chaotic
  - ▶ theoretical argument: no language can fully model itself  
cf. Gödel's incompleteness theorems
- ▶ Imperfect representation of intended semantics required  
focus on some aspect

Big question: what aspects to focus on?

## Querying as a Guide

### Idea

- ▶ Very difficult to choose aspects for absolute semantics
- ▶ Turn problem around
  - ▶ ask what the practical purpose of the semantics could be
  - ▶ then choose aspects that allow realizing that purpose

Meta-remark: That's why we did relative semantics and querying first in this course even though absolute semantics conceptually belongs at the beginning.

### Querying as the Purpose

- ▶ Before: identified different kinds of querying  
focussing on different aspects of knowledge
- ▶ Now: each induces a kind of absolute semantics

# Semantics

## Absolute Semantics

# Deductive Semantics

## Definition

- ▶ A system that determines which propositions are theorems  
called a calculus
- ▶ Languages called logics
- ▶ Implementations called theorem provers

## More precisely

- ▶ Judgment:  $\vdash F$  for “ $F$  is theorem”
- ▶ Set of rules for deriving judgments

## Examples

- ▶ Natural deduction for first-order logic
- ▶ Axiomatic set theory for (most of) mathematics

# Redundant Deductive Semantics

## Multiple deductive semantics

- ▶ Proof theory: absolute
- ▶ Model theory: relative via translation to set theory  $L$   
write  $\models F$  for  $\vdash_L \text{True}[[F]]$
- ▶ Logic translation: relative via translation into standard logics, e.g., SFOL

## Equivalence Theorems

- ▶ Soundness:  $\vdash F$  implies  $\models F$
- ▶ Completeness:  $\models F$  implies  $\vdash F$  accordingly for other translations

# Computational Semantics

## Definition

- ▶ A system that evaluates expressions to values
- ▶ Languages typically called programming languages
- ▶ Implementations called interpreters, evaluators

## More precisely

- ▶ Judgment:  $\vdash E \rightsquigarrow V$  for “ $E$  evaluates to  $V$ ”
- ▶ Set of rules for deriving  $E \rightsquigarrow E_1 \rightsquigarrow E_2 \rightsquigarrow \dots$
- ▶ Often more complex judgments using context containing heap, stack, IO channels, local variables

## Examples

- ▶ Any interpreted language Python, bash, ...
- ▶ Machine language interpretation rules built into microchips

# Redundant Computation Semantics

## Multiple computational semantics

- ▶ Specification: absolute as rules on paper
- ▶ Interpreter: absolute as implementation
- ▶ Compiler: relative via translation to assembly  $L$

write  $\models E \rightsquigarrow V$  for  $\vdash_L \llbracket E \rrbracket \rightsquigarrow \llbracket V \rrbracket$

- ▶ Cross-compilation: relative via translation into other languages

Church-Turing thesis: always possible

## Equivalence Theorems

- ▶ Correctness of compiler:  $\vdash E \rightsquigarrow V$  iff  $\models E \rightsquigarrow V$

accordingly for other translations

## Concrete Semantics

### Definition

- ▶ A system that finds known ground instances of propositions
- ▶ Languages often called query languages  
inspired our, more general use of the word
- ▶ Implementations focusing on caching finite sets of ground instances called triple stores, databases

### More precisely

- ▶ Judgment:  $\vdash F[\gamma]$  for “ $\gamma$  is known ground instance of  $F$ ”
- ▶ Set of rules/sets/tables for finding all  $\gamma$

### Examples

- ▶ SQL for Church-typed ontologies with ADTs
- ▶ SPARQL for Curry-typed ontologies
- ▶ Prolog for first-order logic



## Yes/No vs. Wh-Questions

Deductive/concrete semantics may be a bit of a misnomer

- ▶ Queries about  $\vdash F$  are yes/no questions
  - ▶ specialty of deductive semantics
  - ▶ but maybe only because everything else is ever harder to do deductively
- ▶ Queries about ground instances of  $\Gamma \vdash F$  are Wh questions
  - ▶ specialty of concrete databases
  - ▶ for the special case of retrieving finite results sets from a fixed concrete store
  - ▶ only situation where Wh questions are easy

But Yes/no and Wh questions exist in all aspects.

# Redundant Concrete Semantics

## Multiple concrete semantics

- ▶ Specification: absolute as rules on paper
- ▶ Database: absolute by custom database
- ▶ Database: relative via translation to assembly  $L$

## Equivalence Theorems

- ▶ typically: choose one, no redundancy, no equivalence theorems
- ▶ infinite results: easy on paper, hard in database
- ▶ open world: are all known ground instances in database?

# Narrative Semantics

## Definition

- ▶ Describes how to answer (some) questions
- ▶ Implementations tend to be AI-complete, hypothetical
- ▶ In practice, information retrieval = find related documents

## More precisely?

- ▶ Not much theory, wide open research problem
- ▶ Some natural language document with interspersed definitions, formulas
- ▶ Maybe judgment:  $\vdash Q?A$  for “A is answer to Q”

## Examples

- ▶ “W3C Recommendation OWL 2” and Google
- ▶ “ISO/IEC 14882: 1998 Programming Language C++” and Stroustrup’s book
- ▶ Mathematics textbooks and mathematicians

# Semantics

## An Abstract Definition

# Languages

## A formal system / consists of

- ▶ a set of vocabularies  $\text{Voc}^I$
- ▶ for every  $V \in \text{Voc}^I$ , a set  $\text{Exp}^I(V)$  of expressions
- ▶ a typing relation  $\vdash_V^I e : E$  between  $e, E \in \text{Exp}^I(V)$   
define:  $\text{Exp}_V^I(E) = \{e \in \text{Exp}_V^I \mid \vdash_V^I e : E\}$

convention: leave out superscript  $I$ , subscript  $V$  if clear

## A formal system with propositions

- ▶ additionally has a distinguished expression  $\text{prop}$
- ▶ define  $F$  is proposition if  $\vdash_V F : \text{prop}$

## A formal system with equality

- ▶ additionally has a distinguished proposition  $e_1 \doteq_E e_2$  whenever  $\vdash e_i : E$

in the sequel: fix  $I$  as above

## Deductive Semantics

A deductive semantics for  $I$  consists of

- ▶ for every  $V$ , a subset  $\text{Thm}_V^I \subseteq \text{Exp}_V^I(\text{prop})$  of theorems  
write  $\vdash_V^I F$  for  $F \in \text{Thm}_V^I$

# Curry-Howard

Define deductive semantics as a special case of typing

- ▶ propositions as types
- ▶ proofs as expressions
- ▶ add typing rules such that  $\vdash P : F$  captures the statement “ $P$  is proof of  $F$ ”
- ▶ define:  $\vdash F$  iff there is  $P$  such that  $\vdash P : F$

# Computational Semantics

A computational semantics for  $\lambda$  consists of

- ▶ for every  $V$ , a function  $\text{Eval}'_V : \text{Exp}'_V \rightarrow \text{Exp}'_V$
- ▶ the image of  $\text{Eval}$  called the values

write  $\vdash'_V e \rightsquigarrow e'$  for  $e' = \text{Eval}'_V(e)$

If we also have typing, we say

- ▶ subject reduction: if  $\vdash e : E$ , then  $\vdash \text{Eval}(e) : E$

If we also have equality and deductive semantics, we say

- ▶ normal forms:
  - ▶  $\text{Eval}'_V$  idempotent, i.e.,  $\text{Eval}'_V(x) = x$  if  $x$  value
  - ▶  $\vdash'_V e \doteq_E \text{Eval}'_V(e)$
- ▶ canonical forms:  $\vdash'_V e_1 \doteq_E e_2$  iff  $\text{Eval}'_V(e_1) = \text{Eval}'_V(e_2)$



# Interdefinability

Given a computational semantics, define a deductive one:

- ▶ distinguished expression  $\vdash \text{true} : \text{prop}$ ,
- ▶  $\vdash F$  iff  $\text{Eval}(F) = \text{true}$   
implies decidability, so usually only possible for some  $F$

Given a deductive semantics, define computational one:

- ▶  $\text{Eval}(e)$  is some  $e'$  such that  $\vdash e \doteq e'$   
trivially normal, but usually not canonical

Both kinds of semantics add different value. We usually want both.

## Why Abstract?

### Our definitions are abstract

- ▶  $\text{Exp}$ ,  $\text{Thm}$ ,  $\text{Eval}$  just assumed as sets/functions
- ▶ No requirement how they are constructed
  - ▶ inductive structure of expressions optional
  - ▶ both absolute and relative semantics are special cases

### A more concrete definition might demand

- ▶  $\text{Exp}_V$  defined by grammar
- ▶ type system defined by
  - ▶ calculus for  $\vdash_V e : E$
  - ▶ alternatively: trivial type system where all non-terminals  $N$  are expressions too and  $\vdash E : N$  iff  $E$  derived from  $N$
- ▶  $\text{Thm}_V$  defined by calculus for  $\vdash_V F$
- ▶  $\text{Eval}_V$  defined by calculus for  $\vdash_V e \rightsquigarrow e'$

## Syntax with Contexts

If we want to talk about contexts, too, we need to expand all of the above.

### Syntax with contexts

- ▶ contexts: for every  $V$ , a set  $\text{Cont}'_V$  write  $\vdash_V \Gamma$
- ▶ substitutions: for  $\Gamma, \Delta \in \text{Cont}_V$ , a set  $\text{Subs}_V(\Gamma, \Delta)$   
write  $\vdash_V \gamma : \Gamma \rightarrow \Delta$

### Expressions in context

- ▶ expressions: sets  $\text{Exp}_V(\Gamma)$
- ▶ substitution application: functions  $\text{Exp}(\gamma) : \text{Exp}(\Gamma) \rightarrow \text{Exp}(\Delta)$  for  $\gamma \in \text{Subs}(\Gamma, \Delta)$  write  $\text{Exp}(\gamma)(e)$  as  $e[\gamma]$

### Typing in context

- ▶ expressions: sets  $\text{Exp}_V(\Gamma, E)$ , written as  $\Gamma \vdash_V e : E$
- ▶ substitution preserves types: if  $\Gamma \vdash e : E$  and  $\vdash \gamma : \Gamma \rightarrow \Delta$ , then  $\Delta \vdash e[\gamma] : E[\gamma]$

## Contexts: General Definition

We can leave contexts abstract or spell out a concrete definition:

- ▶ contexts  $\Gamma$  are of the form

$$x_1 : E_1, \dots, x_n : E_n$$

where  $E_i \in \text{Exp}(x_1 : E_1, \dots, x_{i-1} : E_{i-1})$

- ▶ for  $\Gamma$  as above, substitutions  $\Gamma \rightarrow \Delta$  are of the form:

$$x_1 = e_1, \dots, x_n = e_n$$

where  $\Delta \vdash e_i : E_i[x_1 = e_1, \dots, x_{i-1} = e_{i-1}]$

This works uniformly for any formal system. But most formal systems are a bit more restrictive, e.g., by requiring that all  $E_i$  are types.

# Semantics with Contexts

## Deductive semantics

- ▶ define: theorem sets  $\text{Thm}_V(\Gamma)$       write  $F \in \text{Thm}_V(\Gamma)$  as  $\Gamma \vdash_V F$
- ▶ such that theorems are preserved by substitution:  
if  $\Gamma \vdash_V F$  and  $\vdash \gamma : \Gamma \rightarrow \Delta$ , then  $\Delta \vdash_V F[\gamma]$

## Computational semantics

- ▶ define: evaluation functions  $\text{Eval}_V(\Gamma) : \text{Exp}_V(\Gamma) \rightarrow \text{Exp}_V(\Gamma)$   
write  $e' = \text{Eval}_V(\Gamma)(e)$  as  $\Gamma \vdash_V e \rightsquigarrow e'$
- ▶ extend to substitutions:  
 $\text{Eval}_V(\Delta)(\dots, x = e, \dots) = \dots, x = \text{Eval}_V(\Delta)(e), \dots$
- ▶ require that evaluation is preserved by substitution  $\vdash \gamma : \Gamma \rightarrow \Delta$   
 $\text{Eval}_V(\Delta)(e[\gamma]) = \text{Eval}_V(\Delta)(e)[\text{Eval}_V(\Delta)(\gamma)]$   
substitution theorem for  $\text{Eval}$  as a translation from  $I$  to itself

# Concrete Semantics

## General definitions for substitutions

- ▶ write  $\cdot$  for empty context/substitution
- ▶ ground expression is expression in empty context  
also called closed; then opposite is open
- ▶ ground substitution:  $\vdash \gamma : \Gamma \rightarrow \emptyset$  no free variables after substitution
- ▶ if we have computational semantics:  
value substitution is ground substitution where all expressions are values
- ▶ if we have deductive semantics:  
true instance of  $\Gamma \vdash F : \text{prop}$  is  $\gamma$  such that  $\vdash F[\gamma]$

## A concrete semantics for $/$ consists of

- ▶ for every  $\Gamma \vdash_V^I F : \text{prop}$ , a set  $\text{Inst}_V^I(\Gamma, F)$  of ground substitutions  
write  $\Gamma \vdash_V \gamma : F$  for  $\gamma \in \text{Inst}_V(\Gamma, F)$

## Interdefinability

Given concrete semantics, define a deductive one

- ▶ for ground  $F$ ,  $\text{Inst}(\cdot, F)$  is either  $\{\cdot\}$  or  $\{\}$
- ▶  $\vdash F$  iff  $\text{Inst}(\cdot, F) = \{\cdot\}$   
but concrete semantics usually cannot find all substitutions for all  $F$

Given concrete semantics, define a computational one

- ▶  $\vdash e \rightsquigarrow e'$  iff  $(x = e') \in \text{Inst}(x : E, e \doteq_E x)$   
but concrete semantics usually cannot find that substitution for all  $e$

Given deductive semantics, define a concrete one

- ▶  $\text{Inst}(\Gamma, F) = \{\vdash \gamma : \Gamma \rightarrow \cdot \mid \vdash F[\gamma]\}$   
but deductive semantics usually does not allow computing that set

Given computational semantics, define a concrete one

- ▶  $\text{Inst}(\Gamma, F) = \{\text{Eval}(\cdot, \gamma) \mid \vdash \gamma : \Gamma \rightarrow \cdot, \vdash F[\gamma] \rightsquigarrow \text{true}\}$
- ▶ allows restricting results to value substitutions  
composition of previous inter-definitions, inherits both problems

# Translations

A translation  $T$  from formal system  $I$  to formal system  $L$  consists of

- ▶ function  $\text{Voc}^T : \text{Voc}^I \rightarrow \text{Voc}^L$
- ▶ family of functions  $\text{Exp}_V^T : \text{Exp}_V^I \rightarrow \text{Exp}_{\text{Voc}^T(V)}^L$

## Desirable properties

- ▶ Should satisfy type preservation:

$$\vdash_V^I e : E \quad \text{implies} \quad \vdash_{\text{Voc}^T(V)}^L \text{Exp}_V^T(e) : \text{Exp}_V^T(E)$$

intuition: what we have, is preserved

- ▶ Might satisfy conservativity:

$$\vdash_{\text{Voc}^T(V)}^L e' : \text{Exp}_V^T(E) \quad \text{implies} \quad \vdash_V^I e : E \text{ for some } e$$

intuition: nothing new is added



# Translation of Contexts

## Translations extend to contexts and substitutions

- ▶  $\text{Cont}^T(\dots, x : E, \dots) = \dots, x : \text{Exp}^T(E), \dots$
- ▶  $\text{Subs}^T(\dots, x = e, \dots) = \dots, x = \text{Exp}^T(e), \dots$
- ▶  $\text{Exp}^T(x) = x$  for all variables

## Desirable properties for arbitrary contexts

- ▶ Type preservation:

$$\Gamma \vdash_V^I e : E \quad \text{implies} \quad \text{Cont}_V^T(\Gamma) \vdash_{\text{Voc}^T(V)}^L \text{Exp}_V^T(e) : \text{Exp}_V^T(E)$$

- ▶ Conservativity:

$$\text{Cont}_V^T(\Gamma) \vdash_{\text{Voc}^T(V)}^L e' : \text{Exp}_V^T(E) \quad \text{implies} \quad \Gamma \vdash_V^I e : E \text{ for some } e$$

## Compositionality

Define: a translation is compositional iff we can show the substitution theorem for it

Given

$$\Gamma \vdash_V^I e : E \quad \vdash_V^I \gamma : \Gamma \rightarrow \Delta$$

we have that

$$\text{Cont}_V^T(\Delta) \vdash_{\text{Voc}^T(V)}^L \text{Exp}_V^T(e[\gamma]) \doteq_{\text{Exp}_V^T(E[\gamma])} \text{Exp}_V^T(e)[\text{Subs}_V^T(\gamma)]$$

Simplify: write  $T(-)$  for  $\text{Voc}^T(-)$ ,  $\text{Exp}_V^T(-)$ ,  $\text{Cont}_V^T(-)$ ,  $\text{Subs}_V^T(-)$

$$T(\Delta) \vdash_{T(V)}^L T(e[\gamma]) \doteq_{T(E[\gamma])} T(e)[T(\gamma)]$$

## Relative Semantics

Given

- ▶ formal systems  $I$  and  $L$
- ▶ semantics for  $L$
- ▶ translation  $T$  from  $I$  to  $L$

define semantics for  $I$

- ▶ deductive: define

$$\vdash_V^I F \quad \text{iff} \quad \vdash_{\text{Voc}^T(V)}^L \text{Exp}_V^T(F)$$

- ▶ computational: define

$$\vdash_V^I e \rightsquigarrow e' \quad \text{iff} \quad \vdash_{\text{Voc}^T(V)}^L \text{Exp}_V^T(e) \rightsquigarrow \text{Exp}_V^T(e')$$

both work accordingly with a context  $\Gamma$

- ▶ concrete: define

$$\Gamma \vdash_V^I \gamma : F \quad \text{iff} \quad \text{Cont}_V^T(\Gamma) \vdash_{\text{Voc}^T(V)}^L \text{Subs}_V^T(\gamma) : \text{Exp}_V^T(F)$$

## Equivalence of Semantics

Two semantics  $\vdash^1$  and  $\vdash^2$  for  $I$  are equivalent if they are equal in the abstract sense.

- ▶ deductive:  $\vdash^1 F$  iff  $\vdash^2 F$
- ▶ computational:  $\vdash^1 e \rightsquigarrow e'$  iff  $\vdash^2 e \rightsquigarrow e'$
- ▶ concrete:  $\Gamma \vdash^1 \gamma : F$  iff  $\Gamma \vdash^2 \gamma : F$

Example for deductive semantics:

- ▶  $\vdash^1$  absolute semantics by calculus  
e.g., natural deduction for SFOL
- ▶  $\vdash^2$  relative semantics by translation e.g.,  $L$  is set theory,  $T$  is model theory of SFOL
- ▶ Assume proofs-as-expressions
- ▶ Then:
  - ▶ type preservation = soundness
  - ▶ conservativity = completeness

## Exercise 6: Relative Deductive Semantics for BOL

- ▶ Implement a translation from BOL to untyped FOL
  - you can drop properties, types, and values
  - so that only one type of individuals is needed
- ▶ Use TPTP syntax for FOL      see <http://www.tptp.org/>
- ▶ Translate an example ontology
  - pick any ontology with a non-trivial consequence closure
- ▶ Use a theorem prover for first-order logic to implement a relative deductive semantics for BOL
  - Vampire and E are standard choices
  - see also <http://www.tptp.org/cgi-bin/SystemOnTPTP>
- ▶ Test by example whether your semantics yields the correct consequence closure

## Example: Relative Computational Semantics for BOL

Scala, SQL semantics evaluates

- ▶ concept  $c$  to
  - ▶ SQL: table of individuals      result of running query  $\llbracket c \rrbracket$
  - ▶ Scala: hashset of individuals      result of running program  $\llbracket c \rrbracket$
- ▶ propositions to booleans      accordingly

Technically, results not in image of  $\llbracket - \rrbracket$

Fix: add productions for all values

$F ::= \text{true} \mid \text{false}$       truth values  
 $C ::= \{I, \dots, I\}$       finite concepts

## Equivalence with respect to Semantics

So far: equivalence of **two semantics** wrt **all queries**

Related concept: equivalence of **two queries** wrt **one semantics**

- $F, G$  deductively equivalence:

$$\vdash F \quad \text{iff} \quad \vdash G$$

may be internalized by syntax as proposition  $F \leftrightarrow G$

- $F, G$  concretely equivalent:

$$\vdash F[\gamma] \quad \text{iff} \quad \vdash G[\gamma]$$

for all ground substitutions  $\gamma$       weaker than  $\Gamma \vdash F \leftrightarrow G$

- closed  $e, e'$  computationally equivalent:

$$\vdash e \rightsquigarrow v \quad \text{iff} \quad \vdash e' \rightsquigarrow v$$

may be internalized by syntax as proposition  $e \doteq e'$

## Equivalence with respect to Semantics (2)

Interesting variants of computational semantics

- ▶ open  $e, e'$  extensionally equivalent:

$$\vdash e[\gamma] \rightsquigarrow v \quad \text{iff} \quad \vdash e'[\gamma] \rightsquigarrow v$$

for all ground substitutions  $\gamma$

equal inputs produce equal outputs

weaker than  $\Gamma \vdash e \doteq e'$  — intensional equivalence

- ▶ machines  $M, M'$  observationally equivalent:  
produce equal sequences of outputs for the same sequence of  
inputs                      e.g., automata, objects in OO-programming  
choice of semantics defines legal optimizations in compiler



# Semantics

## Absolute Semantics for BOL

## Judgments

Typing:

$$\Gamma \vdash_V^{BOL} e : E$$

Deduction:

$$\Gamma \vdash_V^{BOL} F$$

Propositions prop:

- ▶  $C \sqsubseteq D, C \equiv D$
- ▶ all three kinds of assertions

Notation:

- ▶ We drop the superscript  $BOL$  everywhere.
- ▶ We drop the subscript  $v$  unless we need to use  $V$ .
- ▶ We drop the context  $\Gamma$  unless we need to use/change  $\Gamma$ .

# Typing

Trivial intrinsic typing (Church)  $\vdash e :^{int} E$

- ▶  $E$  is a non-terminal
- ▶  $e$  an expression derived from  $E$

Refined by extrinsic typing (Curry)  $\vdash e :^{ext} E$

- ▶  $e$  is an individual, i.e.,  $\vdash e :^{int} I$
- ▶  $E$  is a concept, i.e.,  $\vdash E :^{int} C$   
where  $I$  and  $C$  are the non-terminals from the grammar
- ▶  $e$  has concept  $E$ , i.e.,  $\vdash e \text{ is-a } E$

## Propositions as Types

Say also  $\vdash p : f$  for proofs  $p$  of proposition  $f$   
in particular:  $x : f$  in contexts to make local assumptions

Notation:

$\Gamma, f$  instead of  $\Gamma, p : f$

sufficient if we only state the rules, not build proofs

## Lookup Rules

The main rules that need to access the vocabulary:

$$\frac{f \text{ in } V}{\vdash_V f}$$

for assertions or axioms  $f$

Assumptions in the context are looked up accordingly:

$$\frac{x : f \text{ in } \Gamma}{\Gamma \vdash f}$$

## Rules for Subsumption and Equality

Subsumption is an order with respect to equality:

$$\overline{\vdash c \sqsubseteq c}$$

$$\frac{\vdash c \sqsubseteq d \quad \vdash d \sqsubseteq e}{\vdash c \sqsubseteq e}$$

$$\frac{\vdash c \sqsubseteq d \quad \vdash d \sqsubseteq c}{\vdash c \equiv d}$$

Equal concepts can be substituted for each other:

$$\frac{\vdash c \equiv d \quad x : C \vdash f : \text{prop} \quad \vdash f[x := c]}{\vdash f[x := d]}$$

This completely defines equality.

## Rules relating Instancehood and Subsumption

$$\frac{\vdash i \text{ is-a } c \quad \vdash c \sqsubseteq d}{\vdash i \text{ is-a } d}$$

Read:

- ▶ if
  - ▶  $i \text{ is-a } c$
  - ▶  $c \sqsubseteq d$
- ▶ then  $i \text{ is-a } d$

$$\frac{x : I, x \text{ is-a } c \vdash x \text{ is-a } d}{\vdash c \sqsubseteq d}$$

Read:

- ▶ if
  - ▶ assuming an individual  $x$  and  $x \text{ is-a } c$ , then  $x \text{ is-a } d$
- ▶ then  $c \sqsubseteq d$

## Induction

Consider from before

$$\frac{x : I, x \text{ is-a } c \vdash x \text{ is-a } d}{\vdash c \sqsubseteq d}$$

Question: Do we allow proving the hypothesis by checking for each individual  $x$ ? induction



## Induction

Consider from before

$$\frac{x : I, x \text{ is-a } c \vdash x \text{ is-a } d}{\vdash c \sqsubseteq d}$$

Question: Do we allow proving the hypothesis by checking for each individual  $x$ ? induction

- Open world: no

## Induction

Consider from before

$$\frac{x : I, x \text{ is-a } c \vdash x \text{ is-a } d}{\vdash c \sqsubseteq d}$$

Question: Do we allow proving the hypothesis by checking for each individual  $x$ ? induction

- ▶ Open world: no
- ▶ Closed world: yes

$$\frac{\Gamma, x : I \vdash f : \text{prop} \quad \Gamma[x := i] \vdash f[x := i] \text{ for every ind. } i}{\Gamma, x : I \vdash f}$$

effectively applicable if only finitely many individuals

# Rules for Union and Intersection of Concepts

Union as the least upper bound:

$$\frac{\overline{\vdash c \sqsubseteq c \sqcup d} \quad \overline{\vdash d \sqsubseteq c \sqcup d}}{\vdash c \sqsubseteq h \quad \vdash d \sqsubseteq h} \quad \frac{}{\vdash c \sqcup d \sqsubseteq h}$$

Dually, intersection as the greatest lower bound:

$$\frac{\overline{\vdash c \sqcap d \sqsubseteq c} \quad \overline{\vdash c \sqcap d \sqsubseteq d}}{\vdash h \sqsubseteq c \quad \vdash h \sqsubseteq d} \quad \frac{}{\vdash h \sqsubseteq c \sqcap d}$$

## Rules for Existential and Universal

Easy rules:

- Existential

$$\frac{\vdash irj \quad \vdash j \text{ is-a } c}{\vdash i \text{ is-a } \exists r.c}$$

- Universal

$$\frac{\vdash i \text{ is-a } \forall r.c \quad \vdash irj}{\vdash j \text{ is-a } c}$$

Other directions are trickier:

- Existential

$$\frac{\vdash i \text{ is-a } \exists r.c \quad j : I, irj, j \text{ is-a } c \vdash f}{\vdash f}$$

- Universal

$$\frac{j : I, irj \vdash j \text{ is-a } c}{\vdash i \text{ is-a } \forall r.c}$$

## Selected Rules for Relations

Inverse:

$$\frac{\vdash irj}{\vdash jr^{-1}i}$$

Composition:

$$\frac{\vdash irj \quad \vdash js k}{\vdash i(r;s)k}$$

Transitive closure:

$$\frac{}{\vdash ir^*i} \quad \frac{\vdash irj \quad \vdash jr^*k}{\vdash ir^*k}$$

Identity at concept  $c$ :

$$\frac{\vdash iis\text{-}ac}{\vdash i\Delta_c i}$$

# Semantics

## Equivalence of BOL Semantics

## Overview

Now 5 semantics for BOL

- ▶ absolute deductive via calculus
- ▶ relative deductive via SFOL
- ▶ relative computational via Scala
- ▶ relative concrete via SQL
- ▶ relative narrative via English

Moreover, these are interdefinable.

e.g., Scala translation also induces deductive semantics

Can compare equivalence

- ▶ for every pair of semantics
- ▶ for every kind of equivalence (deductive, concrete, computational)

Question: Which of them hold?

## Questions

For example, consider:

- ▶ Are the absolute semantics and the Scala semantics deductively equivalent?
- ▶ Assuming BOL and SQL have the base types and values: Are the absolute semantics and the SQL semantics concretely equivalent?



## Deductive Semantics

Are these two BOL semantics deductively equivalent

- ▶ absolute deductive semantics
- ▶ relative deductive semantics via translation  $\llbracket - \rrbracket$  to SFOL

Soundness:  $\vdash_V^{BOL} f$  implies  $\vdash_{\llbracket V \rrbracket}^{SFOL} \llbracket f \rrbracket$

- ▶ induction on derivations of  $\vdash_V^{BOL} f$
- ▶ one case per rule                      induction rule from above not sound
- ▶ several pages of work but straightforward and relatively easy

## Deductive Semantics

Are these two BOL semantics deductively equivalent

- ▶ absolute deductive semantics
- ▶ relative deductive semantics via translation  $\llbracket - \rrbracket$  to SFOL

Completeness:  $\vdash_V^{BOL} f$  implied by  $\vdash_{\llbracket V \rrbracket}^{SFOL} \llbracket f \rrbracket$

works if we add missing rules

- ▶ induction on SFOL derivations does not work
  - ▶ SFOL more expressive than BOL
  - ▶  $\llbracket - \rrbracket$  not surjective
- ▶ instead show that  $\llbracket - \rrbracket$  preserves consistency of vocabularies
 

no universal recipe how to do that
- ▶ then a typical proof uses  $V$  extended with  $\neg f$ 
  - ▶ if  $V$  inconsistent,  $\vdash_V f$  for all  $f$ , done
  - ▶ if  $V$  consistent and  $V + \neg f$  inconsistent, then  $\vdash_V f$ , done
  - ▶ if  $V + \neg f$  consistent, so is  $\llbracket V + \neg f \rrbracket$ , which contradicts  $\vdash_{\llbracket V \rrbracket}^{SFOL} \llbracket f \rrbracket$

## Computational Semantics

Are these two BOL semantics deductively equivalent

- ▶ absolute deductive semantics
- ▶ relative deductive semantics via translation  $\llbracket - \rrbracket$  to Scala

Soundness:  $\vdash_V^{BOL} f$  implies  $\vdash_{\llbracket V \rrbracket}^{Scala} \llbracket f \rrbracket \rightsquigarrow true$

- ▶ Problem: Absolute semantics performs consequence closure, e.g.,
  - ▶ transitivity of  $\sqsubseteq$
  - ▶ relationship between  $\sqsubseteq$  and is-a
- ▶ Scala semantics does so only if we explicitly implemented it
 

we didn't

same problem for SQL semantics

## Computational Semantics

Are these two BOL semantics deductively equivalent

- ▶ absolute deductive semantics
- ▶ relative deductive semantics via translation  $\llbracket - \rrbracket$  to Scala

Completeness:  $\vdash_V^{BOL} f$  implied by  $\vdash_{\llbracket V \rrbracket}^{Scala} \llbracket f \rrbracket \rightsquigarrow true$

- ▶ absolute semantics leaves closed world optional
- ▶ Scala uses closed worlds
  - e.g., used to compute  $c \sqsubseteq d$  by checking all individuals
- ▶ complete only if we add induction rule