

# Typing Massive JSON Datasets

Dario Colazzo  
Université Paris Sud -  
INRIA

Giorgio Ghelli  
Università di Pisa

Carlo Sartiani  
Università della  
Basilicata

# Outline

- Introduction & Motivation
- Data model & Type language
- Typing approach
- Conclusions and future work

# INTRODUCTION & MOTIVATION

# Cloud Computing

- Cloud computing is a very popular computing paradigm
  - Clusters of low-end, unreliable, cheap machines
- Applications
  - Data Storage
    - DropBox, SkyDrive, Google Drive, iCloud
  - Data Analysis
    - Facebook, Yahoo!, Google

# Programming Data Intensive Applications

- Traditional programming languages (e.g., Java, C++, C#, etc.)
  - The programmers must deal with the details of the specific cloud architecture
    - map, reduce, local sort, combiner, etc. in Hadoop/Java
- Languages for the cloud
  - A mix between scripting languages and declarative database programming languages
  - They hide the details of the underlying architecture (e.g., Sawzall, Pig Latin)

# Cloud Languages and Static Analysis

- Support for static analysis and typechecking is usually limited
- Types for input data are optional
- Many controls are deferred at runtime
- There is no validation of the input data against a schema
- Consequences
  - Jobs can raise dynamic type errors

# Example: Pig Latin

- Pig Latin types
  - Base types
    - int, long, double, bytearray, chararray
  - Map types
    - Used for describing records
  - Tuple types
    - Record types without types for the fields (only field labels are specified)
  - Bag types

# Example: Pig Latin

- Input.txt: (yahoo,25) (facebook,15) (twitter,7)

- A program with a wrong schema:

```
data = LOAD 'input.txt' AS (query:INT,count:CHARARRAY);  
data2 = FOREACH data GENERATE TOKENIZE(data.count);  
STORE data2 INTO 'EXAMPLE1';
```

- This program is deemed as type-correct
- A run-time error is raised, as `TOKENIZE` only accepts strings as its input



# Objectives

- To automatically derive succinct and precise schema information from large JSON datasets
  - Succinctness has a direct impact on efficiency and effectiveness of typechecking
  - Avoiding type errors related to imprecise types
  - A concise but precise supertype
- To extend and improve existing type systems of cloud languages

# DATA MODEL AND TYPE LANGUAGE

# Data Model

$o ::= \{l : v, \dots, l : v\}$       Objects

$v ::=$

	$o$	
	$[v, \dots, v]$	Arrays
	$v_s$	Simple values
	$\epsilon$	Empty value

$v_s ::=$     **true** | **false** |  $s$  |  $c$  |  $n$

# Type Language

$T ::=$	$B$	
	$\{l : T, \dots, l : T\}$	Closed record type
	$T \cdot T$	List concatenation
	$T + T$	Union type
	$T \circ T$	Record concatenation
	$T^* \mid T^+ \mid T^? \mid \epsilon$	
$B ::=$	$String \mid Bool \mid Char \mid Number$	Base types

# Type Semantics

$\llbracket \epsilon \rrbracket$	$\triangleq$	$\{\epsilon\}$
$\llbracket \{l_1 : T_1, \dots, l_n : T_n\} \rrbracket$	$\triangleq$	$\{\{m_1 : u_1, \dots, m_n : u_n\} \mid$ $\exists \pi : 1..n \rightarrow 1..n. \forall i \in [1, n] :$ $\pi(i) = h \implies l_i = m_h \wedge u_h \in \llbracket T_i \rrbracket\}$
$\llbracket T_1 \cdot T_2 \rrbracket$	$\triangleq$	$\llbracket T_1 \rrbracket \cdot \llbracket T_2 \rrbracket$
	$\triangleq$	$\{L_1 \cdot L_2 \mid L_1 \in \llbracket T_1 \rrbracket, L_2 \in \llbracket T_2 \rrbracket\}$
$\llbracket T_1 + T_2 \rrbracket$	$\triangleq$	$\llbracket T_1 \rrbracket \cup \llbracket T_2 \rrbracket$
$\llbracket T_1 \circ T_2 \rrbracket$	$\triangleq$	$\llbracket T_1 \rrbracket \circ \llbracket T_2 \rrbracket$
	$\triangleq$	$\{o_i \circ o_j \mid o_i \in \llbracket T_1 \rrbracket, o_j \in \llbracket T_2 \rrbracket\}$
$\llbracket T^* \rrbracket$	$\triangleq$	$\llbracket T \rrbracket^*$
$\llbracket T^+ \rrbracket$	$\triangleq$	$\llbracket T \rrbracket^+$
$\llbracket T^? \rrbracket$	$\triangleq$	$\llbracket T \rrbracket^?$

# Type Equivalence

- We need a type equivalence notion for later use
  - Polynomial time
- We base our type equivalence on a set of syntactical subtyping rules
  - Correct but not complete
  - Polynomial time by using dynamic programming
- When two types are not comparable, we resort to a lexicographical comparison between types

# Subtyping Rules

$\epsilon$	$\forall$	$\epsilon$	
$\epsilon$	$\forall$	$T^*$	
$\{l_1 : T_1, \dots, l_n : T_n\}$	$\forall$	$\{m_1 : U_1, \dots, m_n : U_n\}$	if $\exists \pi : 1..n \rightarrow 1..n. \forall i \in [1, n] :$ $l_i = m_{\pi(i)} \wedge T_i \lesssim U_{\pi(i)}$
$T_1 \cdot T_2$	$\forall$	$U_1 \cdot U_2$	if $T_1 \lesssim U_1$ and $T_2 \lesssim U_2$
$T_1$	$\forall$	$U_2 + U_3$	if $T_1 \lesssim U_2$ or $T_1 \lesssim U_3$ with $T_1 \neq V_1 + V_2$
$T_1 + T_2$	$\forall$	$U$	if $T_1 \lesssim U$ and $T_2 \lesssim U$
$T$	$\forall$	$U^*$	if $T \lesssim U$
$T^*$	$\forall$	$U^*$	if $T \lesssim U^*$
$T_1 \cdot T_2$	$\forall$	$U^*$	if $T_1 \lesssim U^*$ and $T_2 \lesssim U^*$
$T_1 \circ T_2$	$\forall$	$U_1 \circ U_2$	if $\exists \pi : 1..2 \rightarrow 1..2. \forall i \in [1, 2] : T_i \lesssim U_{\pi(i)}$

# TYPING APPROACH



# Typing Algorithm

- Our typing algorithm comprises two stages
- First stage
  - We analyze each JSON object and infer a precise type for it
- Second stage
  - We fuse the collection of types obtained from the first stage so to get a more succinct type
  - Fusion is governed by a set of fusion rules

# First Stage

- A Map/Reduce job
- In the Map phase we infer a type for each JSON object
  - We also record cardinality information (as in WordCount)
- In the Reduce phase equivalent types are grouped together and cardinality is updated
- The output of this job is a collection of pairs  $\langle T_i; \text{card}_i \rangle$ 
  - $T_i$  is a type
  - $\text{card}_i$  is the number of object in the dataset having type  $T_i$
  - $\bigcup_i T_i$  is the type for the objects in the dataset

# Map/Reduce Job

MAP(*JSONObj*  $o$ ; *Optional Type*  $T$ )

```
1  if ( $T == \text{NULL}$ ) or not ISMEMBER( $o, T$ )  
2      return  $\langle \text{INFER}(o); 1 \rangle$   
3  else return  $\langle T; 1 \rangle$ 
```

REDUCE( $\langle \text{Type } T; \text{IntList } list \rangle$ )

```
1  int  $card = 0$   
2  for each  $i \in list$   
3       $card = card + 1$   
4  return  $\langle T; card \rangle$ 
```

# Typing Inference Rules

(TYPETRUEBOOL)

---

$$\vdash \mathbf{true} : Bool$$

(TYPENUMBER)

---

$$\vdash n : Number$$

(TYPECHAR)

---

$$\vdash c : Char$$

(TYPEREC)

$$\begin{array}{l} \forall i = 1, \dots, n : \quad \vdash l_i : String \\ \forall i, j = 1, \dots, n : \quad i \neq j \implies l_i \neq l_j \\ \forall i = 1, \dots, n : \quad \vdash v_i : T_i \end{array}$$

---

$$\vdash \{l_1 : v_1, \dots, l_n : v_n\} : \{l_1 : T_1, \dots, l_n : T_n\}$$

(TYPEFALSEBOOL)

---

$$\vdash \mathbf{false} : Bool$$

(TYPESTRING)

---

$$\vdash s : String$$

(TYPEARRAY)

---

$$\vdash v_i : T_i \quad i = 1, \dots, n$$

---

$$\vdash [v_1, \dots, v_n] : T_1 \cdot \dots \cdot T_n$$

# First Stage Example

- 4 JSON objects

```
{  id : 1,  
   age : 14,  
   admin : false,  
   name : "John Smith",  
   phone : 31324378}
```

```
{  id : 2,  
   name : "Edmond Dantes",  
   email : "ed@mc.com",  
   admin : true}
```

```
{  id : 3,  
   name : "Mattia Pascal",  
   admin : false,  
   age : 37,  
   phone : "+333743227"  
   email : "mp@pir.net" }
```

```
{  id : 4,  
   name : "Amanda Clarke",  
   age : 26,  
   admin : false,  
   phone : 2123142222}
```

# First Stage Example

- Map phase inferred types

$$T_1 = \{ \begin{array}{l} id : Number, \\ age : Number, \\ admin : Bool, \\ name : String, \\ phone : Number \end{array} \}$$
$$T_3 = \{ \begin{array}{l} id : Number, \\ name : String, \\ admin : Bool, \\ age : Number, \\ phone : String, \\ email : String \end{array} \}$$
$$T_2 = \{ \begin{array}{l} id : Number, \\ name : String, \\ email : String, \\ admin : Bool \end{array} \}$$
$$T_4 = \{ \begin{array}{l} id : Number, \\ name : String, \\ age : Number, \\ admin : Bool, \\ phone : Number \end{array} \}$$

- Reduce phase output:  $\langle T_1; \{2\} \rangle$ ,  $\langle T_2; \{1\} \rangle$ ,  $\langle T_3; \{1\} \rangle$

# Second Stage

- Types obtained from the first stage are fused together
  - More succinct types
  - Loss of precision
- Fusion is performed according to fusion rules
- A fusion rule  $\langle T_1 \mid T_2 \rangle \rightarrow T_3$  is a rewriting rule such that
  - $T_1 + T_2 \approx T_3$
  - $|T_1 + T_2| \geq |T_3|$

# Fusion Rules

- Fusion rules should be easy to check and to evaluate
- We have a provisional set of rules
  - Regular expression rules
    - Simplification rules
    - Subtyping rules
    - General rules
  - Record type rules



# Regular Expression Rules

## Simplification rules

- 1)  $T \mid \epsilon \rightarrow T$  if  $T$  is nullable
- 2)  $T \mid \epsilon \rightarrow T?$
- 3)  $T+ \mid \epsilon \rightarrow T^*$
- 4)  $T \cdot \epsilon \mid U \rightarrow T + U$
- 5)  $\epsilon \cdot T \mid U \rightarrow T + U$
- 6)  $T \mid T \rightarrow T$

## Subtyping rules

- 7)  $T \mid U \rightarrow U$  if  $T \sim U$  and  $|T| \geq |U|$
- 8)  $T \mid U \rightarrow U$  if  $T \lesssim U$  and  $|T| \geq |U|$

## General rules

- 9)  $T \mid T \cdot U \rightarrow T \cdot U?$
- 10)  $T \cdot U \mid T \cdot V \rightarrow T \cdot (U + V)$
- 11)  $T \mid U \cdot T \rightarrow U? \cdot T$
- 12)  $U \cdot T \mid V \cdot T \rightarrow (U + V) \cdot T$

# Record Type Rules

- 13)  $\{l_1 : T_1\} \circ U \mid \{l_1 : T_2\} \circ V \rightarrow \{l_1 : T_1 + T_2\} \circ (U + V)$
- 14)  $T \mid T \circ U \rightarrow T \circ U?$
- 15)  $T \circ U \mid T \circ V \rightarrow T \circ (U + V)$

# Type Fusion Algorithm

- An heuristic algorithm that focuses on types with low cardinality
- Types with cardinality greater than a given threshold are ignored to improve the overall precision
- Types are ordered by ascending cardinality into a priority queue
- At each iteration
  - The type with the lowest cardinality is popped
  - The type is compared with the other types in the queue to see if a fusion rule is applicable
- The algorithm halts when no more fusions are possible or a given threshold is satisfied

# Example

- Output of the reduce phase:  $\langle T_1; \{2\} \rangle$ ,  $\langle T_2; \{1\} \rangle$ ,  $\langle T_3; \{1\} \rangle$
- Record types are splitted into concatenations of single-field record types

$$T_1 = \begin{cases} \{ id : Number \} \circ \\ \{ age : Number \} \circ \\ \{ admin : Bool \} \circ \\ \{ name : String \} \circ \\ \{ phone : Number \} \end{cases}$$
$$T_2 = \begin{cases} \{ id : Number \} \circ \\ \{ name : String \} \circ \\ \{ email : String \} \circ \\ \{ admin : Bool \} \end{cases}$$
$$T_3 = \begin{cases} \{ id : Number \} \circ \\ \{ name : String \} \circ \\ \{ admin : Bool \} \circ \\ \{ age : Number \} \circ \\ \{ phone : String \} \circ \\ \{ email : String \} \end{cases}$$

# Example

- $T_2$  is selected and fused with  $T_3$  by using Rule (14)
- Fused type  $V_1$

$$V_1 = \begin{aligned} & \{id : Number\} \circ \\ & \{name : String\} \circ \\ & \{email : String\} \circ \\ & \{admin : Bool\} \circ \\ & (\{age : number\} \circ \{phone : String\})? \end{aligned}$$

# Example

- $T_1$  is selected and fused with  $V_1$  by using Rule (15)
- Fused type  $V_2$

$$V_2 = \{id : Number\} \circ \{name : String\} \circ \{admin : Bool\} \circ ((\{age : Number\} \circ \{phone : String\}) + (\{email : String\} \circ (\{age : Number\} \circ \{phone : Number\})))$$

# Example

- $V_2$  is further simplified by applying fusion rules to its union types
- Rule (15) is applied to factorize  $\{age:Number\}$
- Rule (13) is then <sup>applied</sup> to *phone* fields
- Output type

$$V_3 = \{id : Number\} \circ \{name : String\} \circ \{admin : Bool\} \circ (\{age : Number\} \circ (\{phone : String + Number\} \circ \{email : String\}))?)?$$

# CONCLUSIONS AND FUTURE WORK



# Research Status

- A very preliminary work
- We are working on
  - Defining a richer set of fusion rules
  - Improving the fusion algorithm
  - Designing a Map/Reduce fusion algorithm
  - Adapting our algorithm to the type languages of existing cloud languages

# Conclusions

- An approach for typing massive datasets of JSON objects
- Based on fusion rules
- Adaptable to Map/Reduce