# Counting Types
# for Massive JSON Datasets

Mohamed-Amine Baazizi, Dario Colazzo,
*Giorgio Ghelli*, Carlo Sartiani

# Counting types

- Can types count?
- Should they?

# The problem

▸ Type inference for massive JSON datasets

▸ We infer this type

```
{       title :  Str ;
        text : [ Str ] + Null  ;
        author : { address:T? ; affil:T? ;… } ?
        abstract : Str ?

}
```

▸ How «optional» is the author?

# Let us count

$\{$ title : $Str^{1000}$ ;

text : $([\ Str^{8000}\ ]^{800} + Null^{200})^{1000}$ ;

author : $\{$ add:$T^{300}$ ?; affil:$T^{300}$ ?;… $\}^{800}$ ?

abstract : $Str^{20}$ ?

$\}^{1000}$

# Correlation

- $\{ \text{addr:T}^{300}; \text{ aff:T}^{300}; r{:}T^{800} \}^{800}$

- $\{ \text{addr:T}^{300}; \text{ aff:T}^{300}; r{:}T^{300} \}^{300} + \{ r{:}T^{500} \}^{500}$

- $\{ \text{addr:T}^{300}; r{:}T^{300} \}^{300} + \{ \text{aff:T}^{300}; r{:}T^{500} \}^{500}$

- $\{ \text{addr:T}^{300}; r{:}T^{500} \}^{500} + \{ \text{aff:T}^{300}; r{:}T^{300} \}^{300}$

- $\{ \text{addr:T}^{300}; r{:}T^{300} \}^{300} + \{ \text{aff:T}^{300}; r{:}T^{300} \}^{300} + \{ r{:}T^{200} \}^{200}$

# The type system

- $B ::= Null^i \mid Int^i \mid Str^i \mid Bool^i$

- $R ::= \{ l : T , \ldots, l : T \}^i$

- $A ::= [ T ]^i$

- $S ::= B \mid R \mid A$

- $T ::= S \mid 0 \mid T + T$

- $v : S$

- $v_1, \ldots, v_n :^\flat T$

# The type inference algorithm

$$\frac{v_1 : S_1, \dots, v_n : S_n}{\{l_1 : v_1, \dots, l_n : v_n\} : \{l_1 : S_1, \dots, l_n : S_n\}^1}$$

$$\frac{v_1, \dots, v_n \ : ^{\flat} T}{[v_1, \dots, v_n] : [T]^1}$$

# The type inference algorithm

- Multiset lossless rule:

$$\frac{v_1 : S_1 \quad \dots \quad v_n : S_n}{v_1, \dots, v_n : ^\flat S_1 + \cdots + S_n}$$

- Merging rule:

$$\frac{v_1 : S_1 \quad \dots \quad v_n : S_n}{v_1, \dots, v_n : ^\flat reduce(E)(S_1, \dots, S_n)}$$

# Parametric reduction

$$Reduce(\mathbb{T}_1, \mathbb{T}_2, E) =$$
$$\oplus( \quad \{\!| \; Merge(\mathbb{S}_1, \mathbb{S}_2, E) \mid \mathbb{S}_1 \in \circ\mathbb{T}_1, \mathbb{S}_2 \in \circ\mathbb{T}_2, E(\mathbb{S}_1, \mathbb{S}_2) \; |\!\}^m$$
$$\cup^m \; \{\!| \; \mathbb{S}_1 \mid \mathbb{S}_1 \in \circ\mathbb{T}_1, \; \nexists\mathbb{S}_2 \in \circ\mathbb{T}_2. \, E(\mathbb{S}_1, \mathbb{S}_2) \; |\!\}^m$$
$$\cup^m \; \{\!| \; \mathbb{S}_2 \mid \mathbb{S}_2 \in \circ\mathbb{T}_2, \; \nexists\mathbb{S}_1 \in \circ\mathbb{T}_1. \, E(\mathbb{S}_1, \mathbb{S}_2) \; |\!\}^m \quad )$$

# The type inference algorithm

Map

$$\frac{v_1 : T_1 \quad \ldots \quad v_n : T_n}{v_1, \ldots, v_n :^\flat reduce(E)(T_1, \ldots, T_n)}$$

Combine/Reduce

# The meaning of counting

- Nested counting: how?
  - [ [2,2,2] , [2] , [2,[3,3],2] , [2,2] ] , [ [ [3,3] ] ]

- Cumulative counting
  - [ [ $Int^8$ + [$Int^4$ ]$^2$ ]$^5$ ]$^2$

- More precise subtypes
  - [ [ $Int^8$ + [$Int^2$ ]$^1$ ]$^4$ ]$^1$ + [ [ [$Int^2$ ]$^1$ ]$^1$ ]$^1$
  - [ [ $Int^6$ ]$^3$ + [ $Int^2$ + [$Int^2$ ]$^1$ ]$^1$ ]$^1$ + [ [ [$Int^2$ ]$^1$ ]$^1$ ]$^1$
  - [ [$Int^3$]$^1$ + [$Int^1$]$^1$ + [$Int^2$ + [$Int^2$ ]$^1$ ]$^1$ ] + [$Int^2$]$^1$ ]$^1$ + [[ [$Int^2$ ]$^1$ ]$^1$ ]$^1$

# The meaning of counting

- The formal machinery to interpret counting types:
  - $[\![\ \text{Int}\ ]\!] = \{\ \text{Nat}\ \}$
  - $[\![\ \text{Int}\ ]\!]^3 = \{\ \{i, j, k\}^\flat\ |\ i \in \text{Nat}, j \in \text{Nat}, k \in \text{Nat}\}$

  - $[\![\ T + U\ ]\!] = [\![\ T\ ]\!] \cup [\![\ U\ ]\!]$
  - $[\![\ T + U\ ]\!] = \{\ M1 \cup M2\ |\ (M1, M2)\ \text{in}\ [\![\ T ]\!] \times [\![ U ]\!]\ \}$
  - Hence T not a subtype of T+U

# For example: twitter data

{ contributors: $(\text{Null}^{9,599,980} + [\text{Num}^{20}]^{20})^{9,600,000}$ ?;

retweeted : $\text{Bool}^{9,600,000}$ ?;

retweeted_status $\{\ldots\} : \{\ldots\}^{1,200,000}$ ?;

deleted : $\{\ldots\}^{300,000}$ ?;

$\}^{9,900,000}$

# For example: twitter data

{ con: …$^{7,200,000}$; ret: Bool$^{7,200,000}$;…}$^{7,200,000}$

+{ con: …$^{1,200,000}$; ret: Bool$^{1,200,000}$;…}$^{1,200,000}$

+{ con: …$^{1,040,000}$; ret: Bool$^{1,040,000}$; r_s: {}$^{1,040,000}$;…}$^{1,040,000}$

+{ con: …$^{160,000}$; ret: Bool$^{160,000}$; r_s: {}$^{160,000}$;…}$^{160,000}$

+{ deleted: { }$^{300,000}$;…}$^{300,000}$ :

# To sum up

- An algorithm to summarize JSON data:
  - Easy
  - Well defined semantics
  - Parametric
  - Parallel
  - Yielding quantitative information

- Quantitative information from indexes to types

- What else may a counting type do?