An exploration of dependent types for data-centric programming

Introduction

Type systems define rules that assign types to expressions of a language, and can be backed by formalizations called **type theories**. The first (Russell, 1908) was devised to be a logical foundation of mathematics, though it had issues. **Dependently-typed** programming languages restore this goal with the "propositions as types" motto (Wadler, 2015), in which types correspond to mathematical propositions and programs to their proofs. As such, they are used to both write programs and formally verify them.

Objectives

- Study dependent types from the ground up;
- Apply a dependentlytyped language to develop a verified data frame implementation.

Untyped λ -calculus

Most type theories have been developed within the λ -calculus, a formal system for computation. Its untyped variant has the following syntax and evaluation:

 λ -term E ::= x (variable) | EE (application)

Data-centric programming

In data-centric programming, the primary objective is to represent, query and transform data. There are many systems built with this purpose, such as SQL.

For this project, we decided to explore the Lean 4 programming language by implementing a framework for data frames akin to the pandas library (McKinney, 2010). We started by adapting a similar library for Idris 2 by Tejiščák (2020) and building upon it.

$\lambda x. E$ (abstraction)

Definition (β -reduction). (($\lambda x.M$)N) $\rightarrow_{\beta} M[x := N]$

The untyped λ -calculus is Turing-complete (Turing, 1936), and as such, can be used to write any program. This also means evaluation of terms may never halt, so a term may not have a *normal form*. This makes it unfit as a proof system, since it can lead to inconsistency.

Simply-typed λ -calculus (STLC)

Adding types to λ -terms forces their evaluation to halt.

type $T ::= \alpha$ (constant) $E ::= \dots$ $| T \rightarrow T$ (function) $| \lambda x : T. E$

For each syntactic variant of a λ -term, we define a type assignment rule through Gentzen's sequent calculus:

 $x: T \in \Gamma$ $\Gamma \vdash M: S \to T$ $\Gamma \vdash N: S$ $\Gamma, x: S \vdash M: T$

 $\Gamma \vdash x:T$ $\Gamma \vdash MN:T$ $\Gamma \vdash \lambda x.M:S \rightarrow T$ This calculus is equivalent to the implicational intuitionistic propositional calculus. Variations of STLC add polymorphism through type quantification, allowing types to depend on other types.

- Data frame representation
 - def exampleDF : DF
 (["name" :- String
 - , "age" :- Int
 - , "country" :- Option String

]) :=											
	<mark>(</mark>	"Artur"	ł	"Bob"	ł	"Claire"	ł	"Blorg"	ł	"Zorg"	ł
	II	32	ł	54	ł	41	ł	101	ł	99	ł
	II	some "BR"	ł	some "US"	ł	some "FR"	ł	none	ł	none	¦)

Queries and transformations

def alienAges := exampleDF

|>.where (Option.isNone <\$> col "country" (by repeat constructor))
|>.select (

- ("name" :- col "name" (by repeat constructor))
- ("age" :- col "age" (by repeat constructor))
- ("century" :- over (col "age" (by repeat constructor))

(FromInt.fromInt 100)))

|>.orderBy [.asc (col "age" (by repeat constructor))]

|>.get "century" (by repeat constructor)

Conclusion

Writing verified code is tricky, as the compiler might not always infer proofs. However, Lean's extensive macro system allows creating simpler interfaces as domainspecific languages, which improves the process. Our framework could be used to define critically important transformations, to later be translated to other languages.

Dependently-typed λ -calculus

Dependent type theories allow types to also depend on terms. This way, functions can be from terms to terms, types to terms and vice-versa, and types to types. One is the Calculus of Constructions (Coquand and Huet, 1988). Below is code in Lean 4 with a proof for $A \lor B \rightarrow B \lor A$:

```
theorem or_comm : A v B → B v A
| Or.inl a => Or.inr a
| Or.inr b => Or.inl b
```

References

- Coquand, Thierry and Gérard Huet (1988). "The calculus of constructions".
- McKinney, Wes (2010). "Data Structures for Statistical Computing in Python".
- Russell, Bertrand (1908). "Mathematical Logic as Based on the Theory of Types".
- Tejiščák, Matúš (2020). idris-data-frame.
- Turing, Alan M. (1936). "On Computable Numbers, with an Application to the Entscheidungsproblem".
- Wadler, Philip (Nov. 2015). "Propositions as types".

Eduardo Sandalo Porto (eduardo.sandalo@usp.br) Supervisor: Prof. Ana Cristina Vieira de Melo Department of Computer Science – University of São Paulo

