

The 30th Symposium on Implementation and Application of Functional Languages

University of Massachusetts Lowell, 5th - 7th Sept. 2018

IFL 2018 Draft Pre-Proceedings

Contents

1	<i>Pure Functional Epidemics</i> Jonathan Thaler, Thorsten Altenkirch and Peer-Olaf Siebers	4
2	<i>Delta Debugging Type Errors with a Blackbox Compiler</i> Joanna Sharrad, Olaf Chitil and Meng Wang	16
3	<i>HiPERJiT: A Profile-Driven Just-in-Time Compiler for Erlang</i> Konstantinos Kallas and Konstantinos Sagonas	27
4	<i>Spine-local Type Inference</i> Christopher Jenkins and Aaron Stump	39
5	<i>Verifiably Lazy: Verified Compilation of Call-by-Need</i> George Stelle and Darko Stefanovic	51
6	<i>esverify: Verifying Dynamically-Typed Higher-Order Functional Programs by SMT Solving</i> Christopher Schuster, Sohum Banerjea and Cormac Flanagan	62
7	<i>MIL, a Monadic Intermediate Language for Implementing Functional Languages</i> Mark Jones, Justin Bailey and Theodore Cooper	74
8	<i>Task Oriented Programming and the Internet of Things</i> Mart Lubbers, Pieter Koopman and Rinus Plasmeijer	86
9	<i>On Optimizing Bignum Toom-Cook Multiplication</i> Shamil Dzhatdoyev, Marco Morazan, Nicholas Nelson and Josephine Des Rosiers	98
10	<i>Towards Escaping the Compilation Phase with Implicit Parameters</i> Mark Lemay	105
11	<i>Verifying Embedded Attribute Grammars</i> Florent Balestrieri, Alberto Pardo and Marcos Viera	112
12	<i>Extended Memory Reuse – an optimisation for reducing memory allocations and deallocations in reference counted code</i> Hans-Nikolai Viemann, Artjoms Sinkarovs and Sven-Bodo Scholz	123

13	<i>Continuation Passing Style and Primitive Recursive Functions</i>	
	Ian Mackie	132
14	<i>Extracting Verified Constraints from Coq-embedded Array DSLs</i>	
	Artjoms Sinkarovs, Sven-Bodo Scholz and Hans-Nikolai Vießmann	134
15	<i>A DSL embedded in Rust</i>	
	Kyle Headley	138

Pure Functional Epidemics

An Agent-Based Approach

1 Jonathan Thaler
2 Thorsten Altenkirch
3 Peer-Olaf Siebers

4 jonathan.thaler@nottingham.ac.uk
5 thorsten.altenkirch@nottingham.ac.uk
6 peer-olaf.siebers@nottingham.ac.uk
7 University of Nottingham
8 Nottingham, United Kingdom

9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58

59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116

ABSTRACT

Agent-Based Simulation (ABS) is a methodology in which a system is simulated in a bottom-up approach by modelling the micro interactions of its constituting parts, called agents, out of which the global system behaviour emerges.

So far mainly object-oriented techniques and languages have been used in ABS. Using the SIR model of epidemiology, which simulates the spreading of an infectious disease through a population, we demonstrate how to use pure Functional Reactive Programming to implement ABS. With our approach we can guarantee the reproducibility of the simulation at compile time and rule out specific classes of run-time bugs, something that is not possible with traditional object-oriented languages. Also, we found that the representation in a purely functional format is conceptually quite elegant and opens the way to formally reason about ABS.

KEYWORDS

Functional Reactive Programming, Monadic Stream Functions, Agent-Based Simulation

ACM Reference Format:

Jonathan Thaler, Thorsten Altenkirch, and Peer-Olaf Siebers. 2019. Pure Functional Epidemics: An Agent-Based Approach. In *Proceedings of International Symposium on Implementation and Application of Functional Languages (IFL'18)*. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

1 INTRODUCTION

The traditional approach to Agent-Based Simulation (ABS) has so far always been object-oriented techniques, due to the influence of the seminal work of Epstein et al [17] in which the authors claim "[..] object-oriented programming to be a particularly natural development environment for Sugarscape specifically and artificial societies generally [...]" (p. 179). This work established the metaphor in the ABS community, that *agents map naturally to objects* [33], which still holds up today.

In this paper we challenge this metaphor and explore ways of approaching ABS in a pure (lack of implicit side-effects) functional way using Haskell. By doing this we expect to leverage the benefits of pure functional programming [23]: higher expressivity through

declarative code, being polymorph and explicit about side-effects through monads, more robust and less susceptible to bugs due to explicit data flow and lack of implicit side-effects.

As use case we introduce the SIR model of epidemiology with which one can simulate epidemics, that is the spreading of an infectious disease through a population, in a realistic way.

Over the course of three steps, we derive all necessary concepts required for a full agent-based implementation. We start with a Functional Reactive Programming (FRP) [49] solution using Yampa [22] to introduce most of the general concepts and then make the transition to Monadic Stream Functions (MSF) [37] which allow us to add more advanced concepts of ABS to pure functional programming.

The aim of this paper is to show how ABS can be implemented in *pure* Haskell and what the benefits and drawbacks are. By doing this we give the reader a good understanding of what ABS is, what the challenges are when implementing it and how we solve these in our approach.

The contributions of this paper are:

- We present an approach to ABS using *declarative* analysis with FRP in which we systematically introduce the concepts of ABS to *pure* functional programming in a step-by-step approach. Also this work presents a new field of application to FRP as to the best of our knowledge the application of FRP to ABS (on a technical level) has not been addressed before. The result of using FRP allows expressing continuous time-semantics in a very clear, compositional and declarative way, abstracting away the low-level details of time-stepping and progress of time within an agent.
- Our approach can guarantee reproducibility already at compile time, which means that repeated runs of the simulation with the same initial conditions will always result in the same dynamics, something highly desirable in simulation in general. This can only be achieved through purity, which guarantees the absence of implicit side-effects, which allows to rule out non-deterministic influences at compile time through the strong static type system, something not possible with traditional object-oriented approaches. Further, through purity and the strong static type system, we can rule out important classes of run-time bugs e.g. related to dynamic typing, and the lack of implicit data-dependencies

117 which are common in traditional imperative object-oriented
 118 approaches.

- 119 • Using pure functional programming, we can enforce the
 120 correct semantics of agent execution through types where
 121 we demonstrate that this allows us to have both, sequen-
 122 tial monadic behaviour, and agents acting *conceptually* at
 123 the same time in lock-step, something not possible using
 124 traditional object-oriented approaches.

125 In Section 2 we define Agent-Based Simulation, introduce Func-
 126 tional Reactive Programming, Arrowized programming and Monadic
 127 Stream Functions, because our approach builds heavily on these
 128 concepts. In Section 3 we introduce the SIR model of epidemiology
 129 as an example model to explain the concepts of ABS. The heart of
 130 the paper is Section 4 in which we derive the concepts of a pure
 131 functional approach to ABS in three steps, using the SIR model.
 132 Section 5 discusses related work. Finally, we draw conclusions and
 133 discuss issues in Section 6 and point to further research in Section
 134 7.

2 BACKGROUND

2.1 Agent-Based Simulation

135 Agent-Based Simulation is a methodology to model and simulate
 136 a system where the global behaviour may be unknown but the
 137 behaviour and interactions of the parts making up the system is
 138 known. Those parts, called agents, are modelled and simulated, out
 139 of which then the aggregate global behaviour of the whole system
 140 emerges.

141 So, the central aspect of ABS is the concept of an agent which
 142 can be understood as a metaphor for a pro-active unit, situated in
 143 an environment, able to spawn new agents and interacting with
 144 other agents in some neighbourhood by exchange of messages.

145 We informally assume the following about our agents [28, 42, 50]:

- 146 • They are uniquely addressable entities with some internal
 147 state over which they have full, exclusive control.
- 148 • They are pro-active, which means they can initiate actions
 149 on their own e.g. change their internal state, send messages,
 150 create new agents, terminate themselves.
- 151 • They are situated in an environment and can interact with
 152 it.
- 153 • They can interact with other agents situated in the same
 154 environment by means of messaging.

155 Epstein [16] identifies ABS to be especially applicable for analysing
 156 "spatially distributed systems of heterogeneous autonomous actors
 157 with bounded information and computing capacity". They exhibit
 158 the following properties:

- 159 • Linearity & Non-Linearity - actions of agents can lead to
 160 non-linear behaviour of the system.
- 161 • Time - agents act over time, which is also the source of their
 162 pro-activity.
- 163 • States - agents encapsulate some state, which can be accessed
 164 and changed during the simulation.
- 165 • Feedback-Loops - because agents act continuously and their
 166 actions influence each other and themselves in subsequent
 167 time-steps, feedback-loops are the common in ABS.

- 168 • Heterogeneity - agents can have properties (age, height,
 169 sex,...) where the actual values can vary arbitrarily between
 170 agents.
- 171 • Interactions - agents can be modelled after interactions with
 172 an environment or other agents.
- 173 • Spatiality & Networks - agents can be situated within e.g. a
 174 spatial (discrete 2D, continuous 3D,...) or complex network
 175 environment.

2.2 Functional Reactive Programming

176 Functional Reactive Programming is a way to implement systems
 177 with continuous and discrete time-semantics in pure functional lan-
 178 guages. There are many different approaches and implementations
 179 but in our approach we use *Arrowized FRP* [24, 25] as implemented
 180 in the library Yampa [11, 22, 31].

181 The central concept in Arrowized FRP is the Signal Function
 182 (SF), which can be understood as a *process over time* which maps an
 183 input- to an output-signal. A signal can be understood as a value
 184 which varies over time. Thus, signal functions have an awareness
 185 of the passing of time by having access to Δt which are positive
 186 time-steps, the system is sampled with.

$$\begin{aligned} \text{Signal } \alpha &\approx \text{Time} \rightarrow \alpha \\ \text{SF } \alpha \beta &\approx \text{Signal } \alpha \rightarrow \text{Signal } \beta \end{aligned}$$

187 Yampa provides a number of combinators for expressing time-
 188 semantics, events and state-changes of the system. They allow to
 189 change system behaviour in case of events, run signal functions and
 190 generate stochastic events and random-number streams. We shortly
 191 discuss the relevant combinators and concepts we use throughout
 192 the paper. For a more in-depth discussion we refer to [11, 22, 31].

193 *Event.* An event in FRP is an occurrence at a specific point in
 194 time, which has no duration e.g. the recovery of an infected agent.
 195 Yampa represents events through the *Event* type, which is program-
 196 matically equivalent to the *Maybe* type.

197 *Dynamic behaviour.* To change the behaviour of a signal function
 198 at an occurrence of an event during run-time, (amongst others) the
 199 combinator *switch* :: $SF a (b, Event c) \rightarrow (c \rightarrow SF a b) \rightarrow SF a b$ is
 200 provided. It takes a signal function, which is run until it generates an
 201 event. When this event occurs, the function in the second argument
 202 is evaluated, which receives the data of the event and has to return
 203 the new signal function, which will then replace the previous one.
 204 Note that the semantics of *switch* are that the signal function, into
 205 which is switched, is also executed at the time of switching.

206 *Randomness.* In ABS, often there is the need to generate stochas-
 207 tic events, which occur based on e.g. an exponential distribution.
 208 Yampa provides the combinator *occasionally* :: $RandomGen g \Rightarrow g$
 209 $\rightarrow \text{Time} \rightarrow b \rightarrow SF a (Event b)$ for this. It takes a random-number
 210 generator, a rate and a value the stochastic event will carry. It gen-
 211 erates events on average with the given rate. Note that at most
 212 one event will be generated and no 'backlog' is kept. This means
 213 that when this function is not sampled with a sufficiently high
 214 frequency, depending on the rate, it will lose events.

233 Yampa also provides the combinator `noise :: (RandomGen g, Random b) => g -> SF a b`, which generates a stream of noise by returning
 234 a random number in the default range for the type `b`.
 235

236 *Running signal functions.* To *purely* run a signal function Yampa
 237 provides the function `embed :: SF a b -> (a, [(DTime, Maybe a)]) -> [b]`, which allows to run an SF for a given number of steps where
 238 in each step one provides the Δt and an input `a`. The function then
 239 returns the output of the signal function for each step. Note that the
 240 input is optional, indicated by `Maybe`. In the first step at $t = 0$, the
 241 initial `a` is applied and whenever the input is `Nothing` in subsequent
 242 steps, the last `a` which was not `Nothing` is re-used.
 243

2.3 Arrowized programming

244 Yampa's signal functions are arrows, requiring us to program with
 245 arrows. Arrows are a generalisation of monads, which in addition
 246 to the already familiar parameterisation over the output type, allow
 247 parameterisation over their input type as well [24, 25].
 248

249 In general, arrows can be understood to be computations that
 250 represent processes, which have an input of a specific type, process
 251 it and output a new type. This is the reason why Yampa is using
 252 arrows to represent their signal functions: the concept of processes,
 253 which signal functions are, maps naturally to arrows.
 254

255 There exists a number of arrow combinators, which allow arrowized
 256 programming in a point-free style but due to lack of space
 257 we will not discuss them here. Instead we make use of Paterson's
 258 do-notation for arrows [34], which makes code more readable as it
 259 allows us to program with points.
 260

261 To show how arrowized programming works, we implement a
 262 simple signal function, which calculates the acceleration of a falling
 263 mass on its vertical axis as an example [38].
 264

```
fallingMass :: Double -> Double -> SF () Double
fallingMass p0 v0 = proc _ -> do
  v <- arr (+v0) <<< integral -< (-9.8)
  p <- arr (+p0) <<< integral -< v
  returnA -< p
```

265 To create an arrow, the `proc` keyword is used, which binds a variable after which the `do` of Paterson's do-notation [34] follows. Using
 266 the signal function `integral :: SF v v` of Yampa, which integrates the
 267 input value over time using the rectangle rule, we calculate the
 268 current velocity and the position based on the initial position `p0`
 269 and velocity `v0`. The `<<<` is one of the arrow combinators, which
 270 composes two arrow computations and `arr` simply lifts a pure function
 271 into an arrow. To pass an input to an arrow, `-<` is used and `<-`
 272 to bind the result of an arrow computation to a variable. Finally to
 273 return a value from an arrow, `returnA` is used.
 274

2.4 Monadic Stream Functions

275 Monadic Stream Functions (MSF) are a generalisation of Yampa's
 276 signal functions with additional combinators to control and stack
 277 side effects. An MSF is a polymorphic type and an evaluation function,
 278 which applies an MSF to an input and returns an output and a
 279 continuation, both in a monadic context [36, 37]:
 280

```
newtype MSF m a b = MSF {unMSF :: MSF m a b -> a -> m (b, MSF m a b)}
```

281 MSFs are also arrows, which means we can apply arrowized
 282 programming with Paterson's do-notation as well. MSFs are imple-
 283 mented in Dunai, which is available on Hackage. Dunai allows us
 284 to apply monadic transformations to every sample by means of
 285

286 combinators like `arrM :: Monad m => (a -> m b) -> MSF m a b` and
 287 `arrM_ :: Monad m => m b -> MSF m a b`. A part of the library Dunai
 288 is BearRiver, a wrapper, which re-implements Yampa on top of
 289 Dunai, which enables one to run arbitrary monadic computations
 290 in a signal function. BearRiver simply adds a monadic parameter `m`
 291 to each SF, which indicates the monadic context this signal function
 292 runs in.
 293

294 To show how arrowized programming with MSFs works, we
 295 extend the falling mass example from above to incorporate monads.
 296 In this example we assume that in each step we want to accelerate
 297 our velocity `v` not by the gravity constant anymore but by a random
 298 number in the range of 0 to 9.81. Further we want to count the
 299 number of steps it takes us to hit the floor, that is when position `p`
 300 is less than 0. Also when hitting the floor we want to print a debug
 301 message to the console with the velocity by which the mass has hit
 302 the floor and how many steps it took.
 303

304 We define a corresponding monad stack with `IO` as the innermost
 305 Monad, followed by a `RandT` transformer for drawing random
 306 numbers and finally a `StateT` transformer to count the number of
 307 steps we compute. We can access the monadic functions using `arrM`
 308 in case we need to pass an argument and `_arrM` in case no argument
 309 to the monadic function is needed:
 310

```
type FallingMassStack g = StateT Int (RandT g IO)
type FallingMassMSF g   = SF (FallingMassStack g) () Double
fallingMassMSF :: RandomGen g => Double -> Double -> FallingMassMSF g
fallingMassMSF v0 p0 = proc _ -> do
  -- drawing random number for our gravity range
  r <- arrM_ (lift $ lift $ getRandomR (0, 9.81)) -< ()
  v <- arr (+v0) <<< integral -< (-r)
  p <- arr (+p0) <<< integral -< v
  -- count steps
  arrM_ (lift (modify (+1))) -< ()
  if p > 0
    then returnA -< p
    -- we have hit the floor
  else do
    -- get number of steps
    s <- arrM_ (lift get) -< ()
    -- write to console
    arrM (liftIO . putStrLn) -< "hit floor with v " ++ show v ++
      " after " ++ show s ++ " steps"
  returnA -< p
```

311 To run the `fallingMassMSF` function until it hits the floor we
 312 proceed as follows:
 313

```
runMSF :: RandomGen g => g -> Int -> FallingMassMSF g -> IO ()
runMSF g s msf = do
  let msfReaderT = unMSF msf ()
    msfStateT = runReaderT msfReaderT
    msfRand = runStateT msfStateT s
    msfIO = runRandT msfRand g
  (((p, msf'), s'), g') <- msfIO
  when (p > 0) (runMSF g' s' msf')
```

314 Dunai does not know about time in MSFs, which is exactly what
 315 BearRiver builds on top of MSFs. It does so by adding a `ReaderT`
 316 `Double`, which carries the Δt . This is the reason why we need one
 317 extra lift for accessing `StateT` and `RandT`. Thus `unMSF` returns a
 318 computation in the `ReaderT Double` Monad, which we need to peel
 319 away using `runReaderT`. This then results in a `StateT Int` computa-
 320 tion, which we evaluate by using `runStateT` and the current number
 321 of steps as state. This then results in another monadic computation
 322 of `RandT` Monad, which we evaluate using `runRandT`. This finally
 323 returns an `IO` computation, which we simply evaluate to arrive at
 324 the final result.
 325



Figure 1: States and transitions in the SIR compartment model.

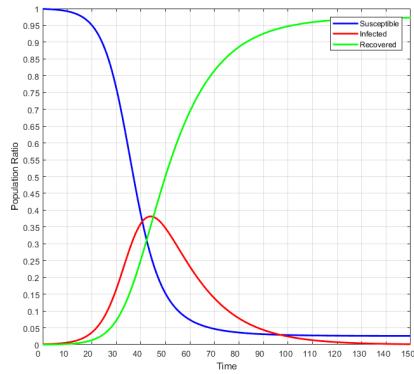


Figure 2: Dynamics of the SIR compartment model using the System Dynamics approach. Population Size $N = 1,000$, contact rate $\beta = \frac{1}{5}$, infection probability $\gamma = 0.05$, illness duration $\delta = 15$ with initially 1 infected agent. Simulation run for 150 time-steps.

3 THE SIR MODEL

To explain the concepts of ABS and of our pure functional approach to it, we introduce the SIR model as a motivating example and use-case for our implementation. It is a very well studied and understood compartment model from epidemiology [27], which allows to simulate the dynamics of an infectious disease like influenza, tuberculosis, chicken pox, rubella and measles spreading through a population [15].

In this model, people in a population of size N can be in either one of three states *Susceptible*, *Infected* or *Recovered* at a particular time, where it is assumed that initially there is at least one infected person in the population. People interact with each other *on average* with a given rate of β per time-unit and become infected with a given probability γ when interacting with an infected person. When infected, a person recovers *on average* after δ time-units and is then immune to further infections. An interaction between infected persons does not lead to re-infection, thus these interactions are ignored in this model. This definition gives rise to three compartments with the transitions seen in Figure 1.

This model was also formalized using System Dynamics (SD) [39]. In SD one models a system through differential equations, allowing to conveniently express continuous systems, which change over time, solving them by numerically integrating over time, which gives then rise to the dynamics. We won't go into detail here and provide the dynamics of such a solution for reference purposes, shown in Figure 2.

An Agent-Based approach

The approach of mapping the SIR model to an ABS is to discretize the population and model each person in the population as an individual agent. The transitions between the states are happening due to discrete events caused both by interactions amongst the agents and time-outs. The major advantage of ABS is that it allows to incorporate spatiality as shown in Section 4.3 and simulate heterogeneity of population e.g. different sex, age. This is not possible with other simulation methods e.g. SD or Discrete Event Simulation [51].

According to the model, every agent makes *on average* contact with β random other agents per time unit. In ABS we can only contact discrete agents thus we model this by generating a random event on average every $\frac{1}{\beta}$ time units. We need to sample from an exponential distribution because the rate is proportional to the size of the population [5]. Note that an agent does not know the other agents' state when making contact with it, thus we need a mechanism in which agents reveal their state in which they are in *at the moment of making contact*. This mechanism is an implementation detail, which we will derive in our implementation steps. For now we only assume that agents can make contact with each other somehow.

4 DERIVING A PURE FUNCTIONAL APPROACH

In [45] two fundamental problems of implementing an ABS from a programming-language agnostic point of view is discussed. The first problem is how agents can be pro-active and the second how interactions and communication between agents can happen. For agents to be pro-active, they must be able to perceive the passing of time, which means there must be a concept of an agent-process, which executes over time. Interactions between agents can be reduced to the problem of how an agent can expose information about its internal state which can be perceived by other agents.

Both problems are strongly related to the semantics of a model and the authors show that it is of fundamental importance to match the update-strategy with the semantics of the model - the order in which agents are updated and actions of agents are visible can make a big difference and need to match the model semantics. The authors identify four different update-strategies, of which the *parallel* update-strategy matches the semantics of the agent-based SIR model due to the underlying roots in the System Dynamics approach. In the parallel update-strategy, the agents act *conceptually* all at the same time in lock-step. This implies that they observe the same environment state during a time-step and actions of an agent are only visible in the next time-step - they are isolated from each other, see Figure 3.

Also, the authors [45] have shown the influence of different deterministic and non-deterministic elements in ABS on the dynamics and how the influence of non-determinism can completely break them down or result in different dynamics despite same initial conditions. This means that we want to rule out any potential source of non-determinism, which we achieve by keeping our implementation pure. This rules out the use of the IO Monad and thus any potential source of non-determinism under all circumstances because we would lose all compile time guarantees about reproducibility. Still we will make use of the Random Monad, which indeed allows

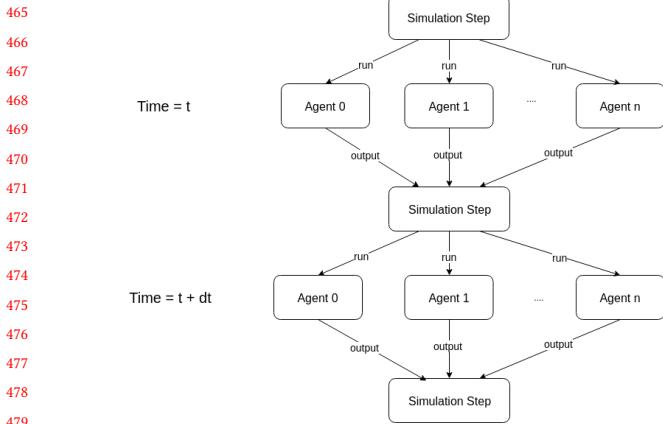


Figure 3: Parallel, lock-step execution of the agents.

side-effects but the crucial point here is that we restrict side-effects only to this type in a controlled way without allowing general unrestricted effects like in traditional object-oriented approaches in the field.

In the following, we derive a pure functional approach for an ABS of the SIR model in which we pose solutions to the previously mentioned problems. We start out with a first approach in Yampa and show its limitations. Then we generalise it to a more powerful approach, which utilises Monadic Stream Functions, a generalisation of FRP. Finally we add a structured environment, making the example more interesting and showing the real strength of ABS over other simulation methodologies like System Dynamics and Discrete Event Simulation¹.

4.1 Functional Reactive Programming

As described in the Section 2.2, Arrowized FRP [24] is a way to implement systems with continuous and discrete time-semantics where the central concept is the signal function, which can be understood as a process over time, mapping an input- to an output-signal. Technically speaking, a signal function is a continuation which allows to capture state using closures and hides away the Δt , which means that it is never exposed explicitly to the programmer, meaning it cannot be messed with.

As already pointed out, agents need to perceive time, which means that the concept of processes over time is an ideal match for our agents and our system as a whole, thus we will implement them and the whole system as signal functions.

4.1.1 Implementation. We start by defining the SIR states as ADT and our agents as signal functions (SF) which receive the SIR states of all agents from the previous step as input and outputs the current SIR state of the agent. This definition, and the fact that Yampa is not monadic, guarantees already at compile, that the agents are isolated from each other, enforcing the *parallel* lock-step semantics of the model.

```
518 data SIRState = Susceptible | Infected | Recovered
519
```

¹The code of all steps can be accessed freely through the following URL: <https://github.com/thalerjonathan/phd/tree/master/public/purefunctionalepidemics/code>

```

523 type SIRAgent = SF [SIRState] SIRState
524
525 sirAgent :: RandomGen g => g -> SIRState -> SIRAgent
526 sirAgent g Susceptible = susceptibleAgent g
527 sirAgent g Infected = infectedAgent g
528 sirAgent _ Recovered = recoveredAgent
529
530
531
```

Depending on the initial state we return the corresponding behaviour. Note that we are passing a random-number generator instead of running in the Random Monad because signal functions as implemented in Yampa are not capable of being monadic.

We see that the recovered agent ignores the random-number generator because a recovered agent does nothing, stays immune forever and can not get infected again in this model. Thus a recovered agent is a consuming state from which there is no escape, it simply acts as a sink which returns constantly *Recovered*:

```

532 recoveredAgent :: SIRAgent
533 recoveredAgent = arr (const Recovered)
534
535
536
537
538
```

Next, we implement the behaviour of a susceptible agent. It makes contact *on average* with β other random agents. For every *infected* agent it gets into contact with, it becomes infected with a probability of γ . If an infection happens, it makes the transition to the *Infected* state. To make contact, it gets fed the states of all agents in the system from the previous time-step, so it can draw random contacts - this is one, very naive way of implementing the interactions between agents.

Thus a susceptible agent behaves as susceptible until it becomes infected. Upon infection an *Event* is returned which results in switching into the *infectedAgent* SF, which causes the agent to behave as an infected agent from that moment on. When an infection event occurs we change the behaviour of an agent using the Yampa combinator *switch*, which is quite elegant and expressive as it makes the change of behaviour at the occurrence of an event explicit. Note that to make contact *on average*, we use Yampas *occasionally* function which requires us to carefully select the right Δt for sampling the system as will be shown in results.

Note the use of $iPre :: a \rightarrow SF a a$, which delays the input signal by one sample, taking an initial value for the output at time zero. The reason for it is, that we need to delay the transition from susceptible to infected by one step due to the semantics of the *switch* combinator: whenever the switching event occurs, the signal function into which is switched will be run at the time of the event occurrence. This means that a susceptible agent could make a transition to recovered within one time-step, which we want to prevent, because the semantics should be that only one state-transition can happen per time-step.

```

573 susceptibleAgent :: RandomGen g => g -> SIRAgent
574 susceptibleAgent g
575   = switch
576     -- delay switching by 1 step to prevent against transition
577     -- from Susceptible to Recovered within one time-step
578     (susceptible g >> iPre (Susceptible, NoEvent))
579     (const (infectedAgent g))
580   where
581     susceptible :: RandomGen g
582       => g -> SF [SIRState] (SIRState, Event ())
583     susceptible g = proc as -> do
584       makeContact <- occasionally g (1 / contactRate) () -< ()
585       if isEvent makeContact
586         then (do
587           -- draw random element from the list
588           a <- drawRandomElemSF g -< as
589           case a of
590             Infected -> do
591               ...
592             Susceptible -> do
593               ...
594             Recovered -> do
595               ...
596             
```

```

581      -- returns True with given probability
582      i <- randomBoolSF g infectivity -< ()
583      if i
584          then returnA -< (Infected, Event ())
585          else returnA -< (Susceptible, NoEvent)
586      else returnA -< (Susceptible, NoEvent)

```

To deal with randomness in an FRP way, we implemented additional signal functions built on the *noiseR* function provided by Yampa. This is an example for the stream character and statefulness of a signal function as it allows to keep track of the changed random-number generator internally through the use of continuations and closures. Here we provide the implementation of *randomBoolSF*. *drawRandomElemSF* works similar but takes a list as input and returns a randomly chosen element from it:

```

594 randomBoolSF :: RandomGen g => g -> Double -> SF () Bool
595 randomBoolSF g p = proc _ -> do
596     r <- noiseR ((0, 1) :: (Double, Double)) g -< ()
597     returnA -< (r <= p)

```

An infected agent recovers *on average* after δ time units. This is implemented by drawing the duration from an exponential distribution [5] with $\lambda = \frac{1}{\delta}$ and making the transition to the *Recovered* state after this duration. Thus the infected agent behaves as infected until it recovers, on average after the illness duration, after which it behaves as a recovered agent by switching into *recoveredAgent*. As in the case of the susceptible agent, we use the *occasionally* function to generate the event when the agent recovers. Note that the infected agent ignores the states of the other agents as its behaviour is completely independent of them.

```

598 infectedAgent :: RandomGen g => g -> SIRAgent
599 infectedAgent g
600     = switch
601         -- delay switching by 1 step
602         (infected >> iPre (Infected, NoEvent))
603         (const recoveredAgent)
604     where
605         infected :: SF [SIRState] (SIRState, Event ())
606         infected = proc _ -> do
607             recEvt <- occasionally g illnessDuration () -< ()
608             let a = event Infected (const Recovered) recEvt
609             returnA -< (a, recEvt)

```

For running the simulation we use Yampas function *embed*:

```

610 runSimulation :: RandomGen g => g -> Time -> DTime
611                 -> [SIRState] -> [[SIRState]]
612 runSimulation g t dt as
613     = embed (stepSimulation sfs as) ((), dts)
614     where
615         steps    = floor (t / dt)
616         dts     = replicate steps (dt, Nothing)
617         n       = length as
618         (rngs, _) = rngSplits g n []
619         sfs     = zipWith sirAgent rngs as

```

What we need to implement next is a closed feedback-loop - the heart of every agent-based simulation. Fortunately, [11, 31] discusses implementing this in Yampa. The function *stepSimulation* is an implementation of such a closed feedback-loop. It takes the current signal functions and states of all agents, runs them all in parallel and returns this step's new agent states. Note the use of *notYet*, which is required because in Yampa switching occurs immediately at $t = 0$. If we don't delay the switching at $t = 0$ until the next step, we would enter an infinite switching loop - *notYet* simply delays the first switching until the next time-step.

```

620 stepSimulation :: [SIRAgent] -> [SIRState] -> SF () [SIRState]
621 stepSimulation sfs as =
622     dpSwitch

```

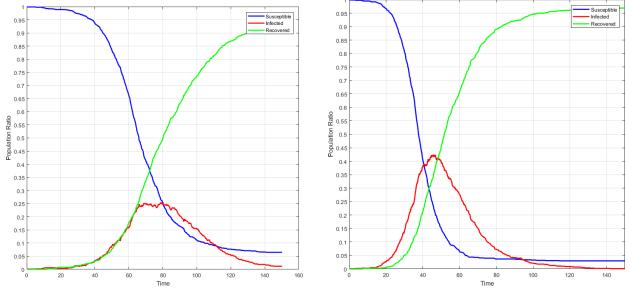
(a) $\Delta t = 0.1$ (b) $\Delta t = 0.01$

Figure 4: FRP simulation of agent-based SIR showing the influence of different Δt . Population size of 1,000 with contact rate $\beta = \frac{1}{5}$, infection probability $\gamma = 0.05$, illness duration $\delta = 15$ with initially 1 infected agent. Simulation run for 150 time-steps with respective Δt .

```

623     -- feeding the agent states to each SF
624     (\_ sfs' -> (map (\sf -> (as, sf)) sfs'))
625     -- the signal functions
626     sfs
627     -- switching event, ignored at t = 0
628     (switchingEvt >> notYet)
629     -- recursively switch back into stepSimulation
630     stepSimulation
631     where
632         switchingEvt :: SF ((), [SIRState]) (Event [SIRState])
633         switchingEvt = arr (\(_ , newAs) -> Event newAs)

```

Yampa provides the *dpSwitch* combinator for running signal functions in parallel, which has the following type-signature:

```

634 dpSwitch :: Functor col
635     -- routing function
636     => (forall sf. a -> col sf -> col (b, sf))
637     -- SF collection
638     -> col (SF b c)
639     -- SF generating switching event
640     -> SF (a, col c) (Event d)
641     -- continuation to invoke upon event
642     -> (col (SF b c) -> d -> SF a (col c))
643     -> SF a (col c)

```

Its first argument is the pairing-function, which pairs up the input to the signal functions - it has to preserve the structure of the signal function collection. The second argument is the collection of signal functions to run. The third argument is a signal function generating the switching event. The last argument is a function, which generates the continuation after the switching event has occurred. *dpSwitch* returns a new signal function, which runs all the signal functions in parallel and switches into the continuation when the switching event occurs. The *d* in *dpSwitch* stands for decoupled which guarantees that it delays the switching until the next time-step: the function into which we switch is only applied in the next step, which prevents an infinite loop if we switch into a recursive continuation.

Conceptually, *dpSwitch* allows us to recursively switch back into the *stepSimulation* with the continuations and new states of all the agents after they were run in parallel.

4.1.2 Results. The dynamics generated by this step can be seen in Figure 4.

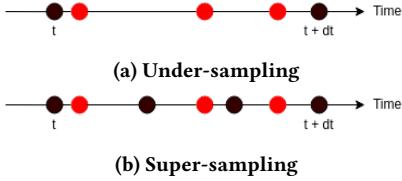


Figure 5: A visual explanation of under-sampling and super-sampling. The black dots represent the time-steps of the simulation. The red dots represent virtual events which occur at specific points in continuous time. In the case of under-sampling, 3 events occur in between the two time steps but *occasionally* only captures the first one. By increasing the sampling frequency either through a smaller Δt or super-sampling all 3 events can be captured.

By following the FRP approach we assume a continuous flow of time, which means that we need to select a *correct* Δt , otherwise we would end up with wrong dynamics. The selection of a correct Δt depends in our case on *occasionally* in the *susceptible* behaviour, which randomly generates an event on average with *contact rate* following the exponential distribution. To arrive at the correct dynamics, this requires us to sample *occasionally*, and thus the whole system, with small enough Δt which matches the frequency of events generated by *contact rate*. If we choose a too large Δt , we loose events, which will result in wrong dynamics as can be seen in Figure 4a. This issue is known as under-sampling and is described in Figure 5.

For tackling this issue we have two options. The first one is to use a smaller Δt as can be seen 4b, which results in the whole system being sampled more often, thus reducing performance. The other option is to implement super-sampling and apply it to *occasionally*, which would allow us to run the whole simulation with $\Delta t = 1.0$ and only sample the *occasionally* function with a much higher frequency.

4.1.3 Discussion. We can conclude that our first step already introduced most of the fundamental concepts of ABS:

- Time - the simulation occurs over virtual time which is modelled explicitly, divided into *fixed* Δt , where at each step all agents are executed.
- Agents - we implement each agent as an individual, with the behaviour depending on its state.
- Feedback - the output state of the agent in the current time-step t is the input state for the next time-step $t + \Delta t$.
- Environment - as environment we implicitly assume a fully-connected network (complete graph) where every agent 'knows' every other agent, including itself and thus can make contact with all of them.
- Stochasticity - it is an inherently stochastic simulation, which is indicated by the random-number generator and the usage of *occasionally*, *randomBoolSF* and *drawRandomElemSF*.
- Deterministic - repeated runs with the same initial random-number generator result in same dynamics. This may not come as a surprise but in Haskell we can guarantee that

property statically already at compile time because our simulation runs *not* in the IO Monad. This guarantees that no external, uncontrollable sources of non-determinism can interfere with the simulation.

- Parallel, lock-step semantics - the simulation implements a *parallel* update-strategy where in each step the agents are run isolated in parallel and don't see the actions of the others until the next step.

Using FRP in the instance of Yampa results in a clear, expressive and robust implementation. State is implicitly encoded, depending on which signal function is active. By using explicit time-semantics with *occasionally* we can achieve extremely fine grained stochastics by sampling the system with small Δt : we are treating it as a truly continuous time-driven agent-based system.

A very severe problem, hard to find with testing but detectable with in-depth validation analysis, is the fact that in the *susceptible* agent the same random-number generator is used in *occasionally*, *drawRandomElemSF* and *randomBoolSF*. This means that all three stochastic functions, which should be independent from each other, are inherently correlated. This is something one wants to prevent under all circumstances in a simulation, as it can invalidate the dynamics on a very subtle level, and indeed we have tested the influence of the correlation in this example and it has an impact. We left this severe bug in for explanatory reasons, as it shows an example where functional programming actually encourages very subtle bugs if one is not careful. A possible but not very elegant solution would be to simply split the initial random-number generator in *sirAgent* three times (using one of the splitted generators for the next split) and pass three random-number generators to *susceptible*. A much more elegant solution would be to use the Random Monad which is not possible because Yampa is not monadic.

So far we have an acceptable implementation of an agent-based SIR approach. What we are lacking at the moment is a general treatment of an environment and an elegant solution to the random number correlation. In the next step we make the transition to Monadic Stream Functions as introduced in Dunai [37], which allows FRP within a monadic context and gives us a way for an elegant solution to the random number correlation.

4.2 Generalising to Monadic Stream Functions

A part of the library Dunai is BearRiver, a wrapper which re-implements Yampa on top of Dunai, which should allow us to easily replace Yampa with MSFs. This will enable us to run arbitrary monadic computations in a signal function, solving our problem of correlated random numbers through the use of the Random Monad.

4.2.1 Identity Monad. We start by making the transition to BearRiver by simply replacing Yampas signal function by BearRivers', which is the same but takes an additional type parameter m , indicating the monadic context. If we replace this type-parameter with the Identity Monad, we should be able to keep the code exactly the same, because BearRiver re-implements all necessary functions we are using from Yampa. We simply re-define the agent signal function, introducing the monad stack our SIR implementation runs in:

```
type SIRMonad = Identity
type SIRAgent = SF SIRMonad [SIRState] SIRState
```

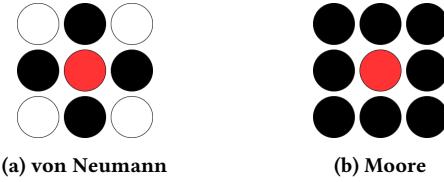


Figure 6: Common neighbourhoods in discrete 2D environments of Agent-Based Simulation.

4.2.2 *Random Monad*. Using the Identity Monad does not gain us anything but it is a first step towards a more general solution. Our next step is to replace the Identity Monad by the Random Monad, which will allow us to run the whole simulation within the Random Monad with the full features of FRP, finally solving the problem of correlated random numbers in an elegant way. We start by re-defining the SIRMonad and SIRAgent:

```
type SIRMonad g = Rand g
type SIRAgent g = SF (SIRMonad g) [SIRState] SIRState
```

The question is now how to access this Random Monad functionality within the MSF context. For the function *occasionally*, there exists a monadic pendant *occasionallyM* which requires a MonadRandom type-class. Because we are now running within a MonadRandom instance we simply replace *occasionally* with *occasionallyM*.

```
occasionallyM :: MonadRandom m => Time -> b -> SF m a (Event b)
-- can be used through the use of arrM and lift
randomBoolM :: RandomGen g => Double -> Rand g Bool
-- this can be used directly as a SF with the arrow notation
drawRandomElemSF :: MonadRandom m => SF m [a] a
```

4.2.3 *Discussion*. Running in the Random Monad solved the problem of correlated random numbers and elegantly guarantees us that we won't have correlated stochastics as discussed in the previous section. In the next step we introduce the concept of an explicit discrete 2D environment.

4.3 Adding an environment

So far we have implicitly assumed a fully connected network amongst agents, where each agent can see and 'knows' every other agent. This is a valid environment and in accordance with the System Dynamics inspired implementation of the SIR model but does not show the real advantage of ABS to situate agents within arbitrary environments. Often, agents are situated within a discrete 2D environment [17] which is simply a finite $N \times M$ grid with either a Moore or von Neumann neighbourhood (Figure 6). Agents are either static or can move freely around with cells allowing either single or multiple occupants.

We can directly map the SIR model to a discrete 2D environment by placing the agents on a corresponding 2D grid with an unrestricted neighbourhood. The behaviour of the agents is the same but they select their interactions directly from the shared read-only environment, which will be passed to the agents as input. This allows agents to read the states of all their neighbours which tells them if a neighbour is infected or not. To show the benefit over the System Dynamics approach and for purposes of a more interesting approach, we restrict the neighbourhood to Moore (Figure 6b).

We also implemented this spatial approach in Java using the well known ABS library RePast [32], to have a comparison with a state of the art approach and came to the same results as shown in Figure 7. This supports, that our pure functional approach can produce such results as well and compares positively to the state of the art in the ABS field.

4.3.1 *Implementation*. We start by defining the discrete 2D environment for which we use an indexed two dimensional array. Each cell stores the agent state of the last time-step, thus we use the *SIRState* as type for our array data. Also, we re-define the agent signal function to take the structured environment *SIREnv* as input instead of the list of all agents as in our previous approach. As output we keep the *SIRState*, which is the state the agent is currently in. Also we run in the Random Monad as introduced before to avoid the random number correlation.

```
type Disc2dCoord = (Int, Int)
type SIREnv      = Array Disc2dCoord SIRState
type SIRAgent g = SF (Rand g) SIREnv SIRState
```

Note that the environment is not returned as output because the agents do not directly manipulate the environment but only read from it. Again, this enforces the semantics of the *parallel* update-strategy through the types where the agents can only see the previous state of the environment and see the actions of other agents reflected in the environment only in the next step.

Note that we could have chosen to use a StateT transformer with the *SIREnv* as state, instead of passing it as input, with the agents then able to arbitrarily read/write, but this would have violated the semantics of our model because actions of agents would have become visible within the same time-step.

The implementation of the susceptible, infected and recovered agents are almost the same with only the neighbour querying now slightly different.

Stepping the simulation needs a new approach because in each step we need to collect the agent outputs and update the environment for the next next step. For this we implemented a separate MSF, which receives the coordinates for every agent to be able to update the state in the environment after the agent was run. Note that we need use *mapM* to run the agents because we are running now in the context of the Random Monad. This has the consequence that the agents are in fact run sequentially one after the other but because they cannot see the other agents actions nor observe changes in the shared read-only environment, it is *conceptually* a *parallel* update-strategy where agents run in lock-step, isolated from each other at conceptually the same time.

```
simulationStep :: RandomGen g => [(SIRAgent g, Disc2dCoord)]
                  -> SIREnv -> SF (Rand g) () SIREnv
simulationStep sfsCoords env = MSF (\_ -> do
  let (sfs, coords) = unzip sfsCoords
  -- run agents sequentially but with shared, read-only environment
  ret <- mapM ('unMSF` env) sfs
  -- construct new environment from all agent outputs for next step
  let (as, sfs') = unzip ret
  env' = foldr (\(a, coord) envAcc -> updateCell coord a envAcc)
               env (zip as coords)
  sfsCoords' = zip sfs' coords
  cont      = simulationStep sfsCoords' env'
  return (env', cont))
updateCell :: Disc2dCoord -> SIRState -> SIREnv -> SIREnv
```

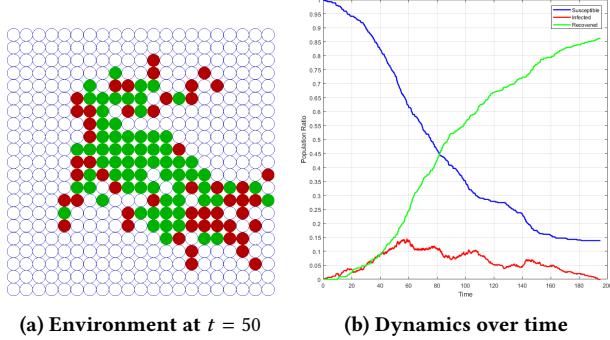


Figure 7: Simulating the agent-based SIR model on a 21x21 2D grid with Moore neighbourhood (Figure 6b), a single infected agent at the center and same SIR parameters as in Figure 2. Simulation run until $t = 200$ with fixed $\Delta t = 0.01$. Last infected agent recovers around $t = 194$. The susceptible agents are rendered as blue hollow circles for better contrast.

4.3.2 Results. We implemented rendering of the environments using the gloss library which allows us to cycle arbitrarily through the steps and inspect the spreading of the disease over time visually as seen in Figure 7.

Note that the dynamics of the spatial SIR simulation which are seen in Figure 7b look quite different from the reference dynamics of Figure 2. This is due to a much more restricted neighbourhood which results in far fewer infected agents at a time and a lower number of recovered agents at the end of the epidemic, meaning that fewer agents got infected overall.

4.3.3 Discussion. By introducing a structured environment with a Moore neighbourhood, we showed the ABS ability to place the heterogeneous agents in a generic environment, which is the fundamental advantage of an agent-based approach over other simulation methodologies and allows us to simulate much more realistic scenarios.

Note, that an environment is not restricted to be a discrete 2D grid and can be anything from a continuous N-dimensional space to a complex network - one only needs to change the type of the environment and agent input and provide corresponding neighbourhood querying functions.

4.4 Additional Steps

ABS involves a few more advanced concepts, which we don't fully explore in this paper due to lack of space. Instead we give a short overview and discuss them without presenting code or going into technical details.

4.4.1 Synchronous Agent Interactions. Synchronous agent interactions are necessary when an arbitrary number of interactions between two agents need to happen instantaneously within the same time-step. The use-case for this are price negotiations between multiple agents where each pair of agents needs to come to an agreement in the same time-step [17]. In object-oriented programming, the concept of synchronous communication between

agents is implemented directly with method calls. We have implemented synchronous interactions in an additional step. We solved it pure functionally by running the signal functions of the transacting agent pair as often as their protocol requires but with $\Delta t = 0$, which indicates the instantaneous character of these interactions.

4.4.2 Event-Driven Approach. Our approach is inherently time-driven where the system is sampled with fixed Δt . The other fundamental way to implement an ABS in general, is to follow an event-driven approach [30], which is based on the theory of Discrete Event Simulation [51]. In such an approach the system is not sampled in fixed Δt but advanced as events occur, where the system stays constant in between. Depending on the model, in an event-driven approach it may be more natural to express the requirements of the model.

In an additional step we have implemented a rudimentary event-driven approach, which allows the scheduling of events. Using the flexibility of MSFs we added a State transformer to the monad stack, which allows queuing of events into a priority queue. The simulation is advanced by processing the next event at the top of the queue, which means running the MSF of the agent which receives the event. The simulation terminates if there are either no more events in the queue or after a given number of events, or if the simulation time has advanced to some limit. Having made the transition to MSFs, implementing this feature was quite straight forward, which shows the power and strength of the generalised approach to FRP using MSFs.

4.4.3 Conflicts in Environment. The semantics of the agent-based SIR model allowed a straight-forward implementation of the parallel update-strategy. This is not easily possible when there could be conflicts in the environment e.g. moving agents where only a single one can occupy a cell. Most models in ABS [17] solve this by implementing a *sequential* update-strategy [45], where agents are run after another but can already observe the changes by agents run before them in the same time-step. To prevent the introduction of artefacts due to a specific ordering, these models shuffle the agents before running them in each step to average the probability for a specific agent to be run at a fixed position.

It is possible to implement a *sequential* update-strategy using the State Monad but functional programming might offer other conflict resolving mechanisms as well because of immutable data and its different nature of side-effects. One approach could be to still run the agents isolated from each other without a State Monad but in case of conflicts, to randomly select a winner and re-run other conflicting agents signal functions until there is no more conflict. As long as the underlying monadic context is robust to re-runs, e.g. the Random Monad, this is no problem. We argue that such an approach is conceptually and semantically cleaner and easier implemented in functional programming than in traditional object-oriented approaches.

5 RELATED WORK

The amount of research on using pure functional programming with Haskell in the field of ABS has been moderate so far. Most of the papers are related to the field of Multi Agent Systems and look

1045 into how agents can be specified using the belief-desire-intention
 1046 paradigm [13, 26, 44].

1047 The author of [4] investigated in his master thesis Haskells parallel
 1048 and concurrency features to implement (amongst others) *HLogo*,
 1049 a Haskell clone of the ABS simulation package NetLogo, where
 1050 agents run within the IO Monad and make use of Software Trans-
 1051 actional Memory for a limited form of agent-interactions.

1052 A library for DES and SD in Haskell called *Aivika* 3 is described
 1053 in the technical report [43]. It is not pure, as it uses the IO Monad
 1054 under the hood and comes only with very basic features for event-
 1055 driven ABS, which allows to specify simple state-based agents with
 1056 timed transitions.

1057 Using functional programming for DES was discussed in [26]
 1058 where the authors explicitly mention the paradigm of FRP to be
 1059 very suitable to DES.

1060 A domain-specific language for developing functional reactive
 1061 agent-based simulations was presented in [48]. This language called
 1062 FRABJOUS is human readable and easily understandable by domain-
 1063 experts. It is not directly implemented in FRP/Haskell but is com-
 1064 piled to Yampa code which they claim is also readable. This supports
 1065 that FRP is a suitable approach to implement ABS in Haskell. Unfor-
 1066 tunately, the authors do not discuss their mapping of ABS to FRP
 1067 on a technical level, which would be of most interest to functional
 1068 programmers.

1069 Object-oriented programming and simulation have a long history
 1070 together as the former one emerged out of Simula 67 [12] which
 1071 was created for simulation purposes. Simula 67 already supported
 1072 Discrete Event Simulation and was highly influential for today's
 1073 object-oriented languages. Although the language was important
 1074 and influential, in our research we look into different approaches,
 1075 orthogonal to the existing object-oriented concepts.

1076 Lustre is a formally defined, declarative and synchronous dataflow
 1077 programming language for programming reactive systems [19].
 1078 While it has solved some issues related to implementing ABS in
 1079 Haskell it still lacks a few important features necessary for ABS.
 1080 We don't see any way of implementing an environment in Lustre as
 1081 we do in our approach in Section 4.3. Also the language seems not
 1082 to come with stochastic functions, which are but the very building
 1083 blocks of ABS. Finally, Lustre does only support static networks,
 1084 which is clearly a drawback in ABS in general where agents can be
 1085 created and terminated dynamically during simulation.

1086 There exists some research [14, 41, 47] of using the functional
 1087 programming language Erlang [3] to implement ABS. The language
 1088 is inspired by the actor model [1] and was created in 1986 by Joe
 1089 Armstrong for Eriksson for developing distributed high reliability
 1090 software in telecommunications. The actor model can be seen as
 1091 quite influential to the development of the concept of agents in ABS,
 1092 which borrowed it from Multi Agent Systems [50]. It emphasises
 1093 message-passing concurrency with share-nothing semantics, which
 1094 maps nicely to functional programming concepts. The mentioned
 1095 papers investigate how the actor model can be used to close the con-
 1096 ceptual gap between agent-specifications, which focus on message-
 1097 passing and their implementation. Further they also showed that
 1098 using this kind of concurrency allows to overcome some problems
 1099 of low level concurrent programming as well. Despite the natural
 1100 mapping of ABS concepts to such an actor language, it leads to

1103 simulations, which despite same initial starting conditions, might
 1104 result in different dynamics each time due to concurrency.

6 CONCLUSIONS

1105 Our FRP based approach is different from traditional approaches in
 1106 the ABS community. First it builds on the already quite powerful
 1107 FRP paradigm. Second, due to our continuous time approach, it
 1108 forces one to think properly of time-semantics of the model and
 1109 how small Δt should be. Third it requires one to think about agent
 1110 interactions in a new way instead of being just method-calls.

1111 Because no part of the simulation runs in the IO Monad and we
 1112 do not use *unsafePerformIO* we can rule out a serious class of bugs
 1113 caused by implicit data-dependencies and side-effects, which can
 1114 occur in traditional imperative implementations.

1115 Also we can statically guarantee the reproducibility of the simu-
 1116 lation, which means that repeated runs with the same initial con-
 1117 ditions are guaranteed to result in the same dynamics. Although
 1118 we allow side-effects within agents, we restrict them to only the
 1119 Random Monad in a controlled, deterministic way and never use
 1120 the IO Monad, which guarantees the absence of non-deterministic
 1121 side effects within the agents and other parts of the simulation.

1122 Determinism is also ensured by fixing the Δt and not making
 1123 it dependent on the performance of e.g. a rendering-loop or other
 1124 system-dependent sources of non-determinism as described by
 1125 [38]. Also by using FRP we gain all the benefits from it and can use
 1126 research on testing, debugging and exploring FRP systems [35, 38].

1127 Also we showed how to implement the *parallel* update-strategy
 1128 [45] in a way that the correct semantics are enforced and guaranteed
 1129 already at compile time through the types. This is not possible in
 1130 traditional imperative implementations and poses another unique
 1131 benefit over the use of functional programming in ABS.

Issues

1132 Currently, the performance of the system does not come close to
 1133 imperative implementations. We compared the performance of
 1134 our pure functional approach as presented in Section 4.3 to an
 1135 implementation in Java using the ABS library RePast [32]. We ran
 1136 the simulation until $t = 100$ on a 51×51 (2,601 agents) with $\Delta t =$
 1137 0.1 (unknown in RePast) and averaged 8 runs. The performance
 1138 results make the lack of speed of our approach quite clear: the pure
 1139 functional approach needs around 72.5 seconds whereas the Java
 1140 RePast version just 10.8 seconds on our machine to arrive at $t = 100$.
 1141 It must be mentioned, that RePast does implement an event-driven
 1142 approach to ABS, which can be much more performant [30] than a
 1143 time-driven one as ours, so the comparison is not completely valid.
 1144 Still, we have already started investigating speeding up performance
 1145 through the use of Software Transactional Memory [20, 21], which
 1146 is quite straight forward when using MSFs. It shows very good
 1147 results but we have to leave the investigation and optimization of
 1148 the performance aspect of our approach for further research as it is
 1149 beyond the scope of this paper.

1150 Despite the strengths and benefits we get by leveraging on FRP,
 1151 there are errors that are not raised at compile time, e.g. we can
 1152 still have infinite loops and run-time errors. This was for exam-
 1153 ple investigated in [40] where the authors use dependent types to
 1154 avoid some run-time errors in FRP. We suggest that one could go
 1155

1156
 1157
 1158
 1159
 1160

1161 further and develop a domain specific type system for FRP that
 1162 makes the FRP based ABS more predictable and that would support
 1163 further mathematical analysis of its properties. Furthermore,
 1164 moving to dependent types would pose a unique benefit over the
 1165 traditional object-oriented approach and should allow us to express
 1166 and guarantee even more properties at compile time. We leave this
 1167 for further research.

1168 In our pure functional approach, agent identity is not as clear
 1169 as in traditional object-oriented programming, where there is a
 1170 quite clear concept of object-identity through the encapsulation of
 1171 data and methods. Signal functions don't offer this strong identity
 1172 and one needs to build additional identity mechanisms on top e.g.
 1173 sending messages to specific agents.

1174 We can conclude that the main difficulty of a pure functional
 1175 approach evolves around the communication and interaction be-
 1176 tween agents, which is a direct consequence of the issue with agent
 1177 identity. Agent interaction is straight-forward in object-oriented
 1178 programming, where it is achieved using method-calls mutating the
 1179 internal state of the agent, but that comes at the cost of a new class
 1180 of bugs due to implicit data flow. In pure functional programming
 1181 these data flows are explicit but our current approach of feeding
 1182 back the states of all agents as inputs is not very general. We have
 1183 added further mechanisms of agent interaction which we had to
 1184 omit due to lack of space.

1186 7 FURTHER RESEARCH

1187 We see this paper as an intermediary and necessary step towards
 1188 dependent types for which we first needed to understand the po-
 1189 tential and limitations of a non-dependently typed pure functional
 1190 approach in Haskell. Dependent types are extremely promising
 1191 in functional programming as they allow us to express stronger
 1192 guarantees about the correctness of programs and go as far as al-
 1193 lowing to formulate programs and types as constructive proofs,
 1194 which must be total by definition [2, 29, 46].

1195 So far no research using dependent types in agent-based sim-
 1196 ulation exists at all. In our next paper we want to explore this for
 1197 the first time and ask more specifically how we can add dependent
 1198 types to our pure functional approach, which conceptual implica-
 1199 tions this has for ABS and what we gain from doing so. We plan
 1200 on using Idris [6] as the language of choice as it is very close to
 1201 Haskell with focus on real-world application and running programs
 1202 as opposed to other languages with dependent types e.g. Agda and
 1203 Coq.

1204 We hypothesize that dependent types could help ruling out even
 1205 more classes of bugs at compile time and encode invariants and
 1206 model specifications on the type level, which implies that we don't
 1207 need to test them using e.g. property-testing with QuickCheck. This
 1208 would allow the ABS community to reason about a model directly
 1209 in code. We think that a promising approach is to follow the work
 1210 of [7–10, 18] in which the authors utilize GADTs to implement
 1211 an indexed monad, which allows to implement correct-by-
 1212 construction software.

- 1214 • In the SIR implementation one could make wrong state-
 1215 transitions e.g. when an infected agent should recover, nothing
 1216 prevents one from making the transition back to suscep-
 1217 tible.

1219 Using dependent types it should be possible to encode invari-
 1220 ants and state-machines on the type level, which can prevent
 1221 such invalid transitions already at compile time. This would
 1222 be a huge benefit for ABS because many agent-based models
 1223 define their agents in terms of state-machines.

- 1224 • An infected agent recovers after a given time - the transi-
 1225 tion of infected to recovered is a timed transition. Nothing
 1226 prevents us from *never* doing the transition at all.
 1227 With dependent types we should be able to encode the pass-
 1228 ing of time in the types and guarantee on a type level that an
 1229 infected agent has to recover after a finite number of time
 1230 steps.
- 1231 • In more sophisticated models agents interact in more com-
 1232 plex ways with each other e.g. through message exchange
 1233 using agent IDs to identify target agents. The existence of an
 1234 agent is not guaranteed and depends on the simulation time
 1235 because agents can be created or terminated at any point
 1236 during simulation.
 1237 Dependent types could be used to implement agent IDs as
 1238 a proof that an agent with the given id exists *at the current*
 1239 *time-step*. This also implies that such a proof cannot be used
 1240 in the future, which is prevented by the type system as it is
 1241 not safe to assume that the agent will still exist in the next
 1242 step.
- 1243 • In our implementation, we terminate the SIR model always
 1244 after a fixed number of time-steps. We can informally reason
 1245 that restricting the simulation to a fixed number of time-
 1246 steps is not necessary because the SIR model *has to* reach a
 1247 steady state after a finite number of steps. This means that
 1248 at that point the dynamics won't change any more, thus one
 1249 can safely terminate the simulation. Informally speaking,
 1250 the reason for that is that eventually the system will run
 1251 out of infected agents, which are the drivers of the dynamic.
 1252 We know that all infected agents will recover after a finite
 1253 number of time-steps *and* that there is only a finite source
 1254 for infected agents which is monotonously decreasing.
 1255 Using dependent types it might be possible to encode this
 1256 in the types, resulting in a total simulation, creating a corre-
 1257 spondence between the equilibrium of a simulation and the
 1258 totality of its implementation. Of course this is only possible
 1259 for models in which we know about their equilibria *a priori*
 1260 or in which we can reason somehow that an equilibrium
 1261 exists.

1262 ACKNOWLEDGMENTS

1263 The authors would like to thank I. Perez, H. Nilsson, J. Greensmith,
 1264 M. Baerenz, H. Vollbrecht, S. Venkatesan, J. Hey and the anonymous
 1265 referees from IFL 2018 for constructive feedback, comments and
 1266 valuable discussions.

1269 REFERENCES

- 1271 [1] Gul Agha. 1986. *Actors: A Model of Concurrent Computation in Distributed Systems*.
 MIT Press, Cambridge, MA, USA.
- 1272 [2] Thorsten Altenkirch, Nils Anders Danielsson, Andres Loeh, and Nicolas Oury.
 1273 2010. Pi Sigma: Dependent Types Without the Sugar. In *Proceedings of the 10th Inter-
 1274 national Conference on Functional and Logic Programming (FLOPS'10)*. Springer-
 1275 Verlag, Berlin, Heidelberg, 40–55. https://doi.org/10.1007/978-3-642-12251-4_5

- [3] Joe Armstrong. 2010. Erlang. *Commun. ACM* 53, 9 (Sept. 2010), 68–75. <https://doi.org/10.1145/1810891.1810910>
- [4] Nikolaos Bezirgiannis. 2013. *Improving Performance of Simulation Software Using Haskell's Concurrency & Parallelism*. Ph.D. Dissertation. Utrecht University - Dept. of Information and Computing Sciences.
- [5] Andrei Borshchev and Alexei Filippov. 2004. From System Dynamics and Discrete Event to Practical Agent Based Modeling: Reasons, Techniques, Tools. Oxford.
- [6] Edwin Brady. 2013. Idris, a general-purpose dependently typed programming language: Design and implementation. *Journal of Functional Programming* 23, 05 (2013), 552–593. <https://doi.org/10.1017/S095679681300018X>
- [7] Edwin Brady. 2013. Programming and Reasoning with Algebraic Effects and Dependent Types. In *Proceedings of the 18th ACM SIGPLAN International Conference on Functional Programming (ICFP '13)*. ACM, New York, NY, USA, 133–144. <https://doi.org/10.1145/2500365.2500581>
- [8] Edwin Brady. 2016. *State Machines All The Way Down - An Architecture for Dependently Typed Applications*. Technical Report. <https://www.idris-lang.org/drafts/sms.pdf>
- [9] Edwin Brady and Kevin Hammond. 2010. Correct-by-Construction Concurrency: Using Dependent Types to Verify Implementations of Effectful Resource Usage Protocols. *Fundam. Inf.* 102, 2 (April 2010), 145–176. <http://dl.acm.org/citation.cfm?id=1883634.1883636>
- [10] Edwin C. Brady. 2011. Idris systems programming meets full dependent types. In *Proc. 5th ACM workshop on Programming languages meets program verification, PLPV '11*. ACM, 43–54.
- [11] Antony Courtney, Henrik Nilsson, and John Peterson. 2003. The Yampa Arcade. In *Proceedings of the 2003 ACM SIGPLAN Workshop on Haskell (Haskell '03)*. ACM, New York, NY, USA, 7–18. <https://doi.org/10.1145/871895.871897>
- [12] Ole-johan Dahl. 2002. The birth of object orientation: the simula languages. In *Software Pioneers: Contributions to Software Engineering, Programming, Software Engineering and Operating Systems Series*. Springer, 79–90.
- [13] Tanja De Jong. 2014. *Suitability of Haskell for Multi-Agent Systems*. Technical Report. University of Twente.
- [14] Antonello Di Stefano and Corrado Santoro. 2005. Using the Erlang Language for Multi-Agent Systems Implementation. In *Proceedings of the IEEE/WIC/ACM International Conference on Intelligent Agent Technology (IAT '05)*. IEEE Computer Society, Washington, DC, USA, 679–685. <https://doi.org/10.1109/IAT.2005.141>
- [15] Richard H. Enns. 2010. *It's a Nonlinear World* (1st ed.). Springer Publishing Company, Incorporated.
- [16] Joshua M. Epstein. 2012. *Generative Social Science: Studies in Agent-Based Computational Modeling*. Princeton University Press. Google-Books-ID: 6jPiuMbKKJ4C.
- [17] Joshua M. Epstein and Robert Axtell. 1996. *Growing Artificial Societies: Social Science from the Bottom Up*. The Brookings Institution, Washington, DC, USA.
- [18] Simon Fowler and Edwin Brady. 2014. Dependent Types for Safe and Secure Web Programming. In *Proceedings of the 25th Symposium on Implementation and Application of Functional Languages (IFL '13)*. ACM, New York, NY, USA, 49:49–49:60. <https://doi.org/10.1145/2620678.2620683>
- [19] N. Halbwachs, P. Caspi, P. Raymond, and D. Pilaud. 1991. The synchronous data flow programming language LUSTRE. *Proc. IEEE* 79, 9 (Sept. 1991), 1305–1320. <https://doi.org/10.1109/5.97300>
- [20] Tim Harris, Simon Marlow, Simon Peyton-Jones, and Maurice Herlihy. 2005. Composable Memory Transactions. In *Proceedings of the Tenth ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP '05)*. ACM, New York, NY, USA, 48–60. <https://doi.org/10.1145/1065944.1065952>
- [21] Tim Harris and Simon Peyton Jones. 2006. Transactional memory with data invariants. <https://www.microsoft.com/en-us/research/publication/transactional-memory-data-invariants/>
- [22] Paul Hudak, Antony Courtney, Henrik Nilsson, and John Peterson. 2003. Arrows, Robots, and Functional Reactive Programming. In *Advanced Functional Programming*, Johan Jeuring and Simon L. Peyton Jones (Eds.). Number 2638 in Lecture Notes in Computer Science. Springer Berlin Heidelberg, 159–187. https://doi.org/10.1007/978-3-540-44833-4_6
- [23] Paul Hudak, John Hughes, Simon Peyton Jones, and Philip Wadler. 2007. A History of Haskell: Being Lazy with Class. In *Proceedings of the Third ACM SIGPLAN Conference on History of Programming Languages (HOPL III)*. ACM, New York, NY, USA, 12–1–12–55. <https://doi.org/10.1145/1238844.1238856>
- [24] John Hughes. 2000. Generalising Monads to Arrows. *Sci. Comput. Program.* 37, 1–3 (May 2000), 67–111. [https://doi.org/10.1016/S0167-6423\(99\)00023-4](https://doi.org/10.1016/S0167-6423(99)00023-4)
- [25] John Hughes. 2005. Programming with Arrows. In *Proceedings of the 5th International Conference on Advanced Functional Programming (AeFP '04)*. Springer-Verlag, Berlin, Heidelberg, 73–129. https://doi.org/10.1007/11546382_2
- [26] Peter Jankovic and Ondrej Such. 2007. *Functional Programming and Discrete Simulation*. Technical Report.
- [27] W. O. Kermack and A. G. McKendrick. 1927. A Contribution to the Mathematical Theory of Epidemics. *Proceedings of the Royal Society of London A: Mathematical, Physical and Engineering Sciences* 115, 772 (Aug. 1927), 700–721. <https://doi.org/10.1098/rspa.1927.0118>
- [28] C. M. Macal. 2016. Everything you need to know about agent-based modelling and simulation. *Journal of Simulation* 10, 2 (May 2016), 144–156. <https://doi.org/10.1057/jos.2016.7>
- [29] James McKinna. 2006. Why Dependent Types Matter. In *Conference Record of the 33rd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '06)*. ACM, New York, NY, USA, 1–1. <https://doi.org/10.1145/1111037.1111038>
- [30] Ruth Meyer. 2014. Event-Driven Multi-agent Simulation. In *Multi-Agent-Based Simulation XV (Lecture Notes in Computer Science)*. Springer, Cham, 3–16. https://doi.org/10.1007/978-3-319-14627-0_1
- [31] Henrik Nilsson, Antony Courtney, and John Peterson. 2002. Functional Reactive Programming, Continued. In *Proceedings of the 2002 ACM SIGPLAN Workshop on Haskell (Haskell '02)*. ACM, New York, NY, USA, 51–64. <https://doi.org/10.1145/581690.581695>
- [32] Michael J. North, Nicholson T. Collier, Jonathan Ozik, Eric R. Tatara, Charles M. Macal, Mark Bragen, and Pam Sydelko. 2013. Complex adaptive systems modeling with Repast Simphony. *Complex Adaptive Systems Modeling* 1, 1 (March 2013), 3. <https://doi.org/10.1186/2194-3206-1-3>
- [33] Michael J. North and Charles M. Macal. 2007. *Managing Business Complexity: Discovering Strategic Solutions with Agent-Based Modeling and Simulation*. Oxford University Press, USA. Google-Books-ID: gRATDAAAQBAJ.
- [34] Ross Paterson. 2001. A New Notation for Arrows. In *Proceedings of the Sixth ACM SIGPLAN International Conference on Functional Programming (ICFP '01)*. ACM, New York, NY, USA, 229–240. <https://doi.org/10.1145/507635.507664>
- [35] Ivan Perez. 2017. Back to the Future: Time Travel in FRP. In *Proceedings of the 10th ACM SIGPLAN International Symposium on Haskell (Haskell 2017)*. ACM, New York, NY, USA, 105–116. <https://doi.org/10.1145/3122955.3122957>
- [36] Ivan Perez. 2017. *Extensible and Robust Functional Reactive Programming*. Doctoral Thesis. University of Nottingham, Nottingham.
- [37] Ivan Perez, Manuel Baerenz, and Henrik Nilsson. 2016. Functional Reactive Programming, Refactored. In *Proceedings of the 9th International Symposium on Haskell (Haskell 2016)*. ACM, New York, NY, USA, 33–44. <https://doi.org/10.1145/2976002.2976010>
- [38] Ivan Perez and Henrik Nilsson. 2017. Testing and Debugging Functional Reactive Programming. *Proc. ACM Program. Lang.* 1, ICFP (Aug. 2017), 2:1–2:27. <https://doi.org/10.1145/3110246>
- [39] Donald E. Porter. 1962. *Industrial Dynamics*. Jay Forrester. M.I.T. Press, Cambridge, Mass.; Wiley, New York, 1961. xv + 464 pp. Illus. \$18. *Science* 135, 3502 (Feb. 1962), 426–427. <https://doi.org/10.1126/science.135.3502.426-a>
- [40] Neil Sculthorpe and Henrik Nilsson. 2009. Safe Functional Reactive Programming Through Dependent Types. In *Proceedings of the 14th ACM SIGPLAN International Conference on Functional Programming (ICFP '09)*. ACM, New York, NY, USA, 23–34. <https://doi.org/10.1145/1596550.1596558>
- [41] Gene I. Sher. 2013. *Agent-Based Modeling Using Erlang Eliminating The Conceptual Gap Between The Programming Language & ABM*.
- [42] Peer-Olaf Siebers and Uwe Aickelin. 2008. Introduction to Multi-Agent Simulation. *arXiv:0803.3905 [cs]* (March 2008). <http://arxiv.org/abs/0803.3905> arXiv: 0803.3905.
- [43] David Sorokin. 2015. *Aivika 3: Creating a Simulation Library based on Functional Programming*.
- [44] Martin Sulzmann and Edmund Lam. 2007. *Specifying and Controlling Agents in Haskell*. Technical Report.
- [45] Jonathan Thaler and Peer-Olaf Siebers. 2017. The Art Of Iterating: Update-Strategies in Agent-Based Simulation. Dublin.
- [46] Simon Thompson. 1991. *Type Theory and Functional Programming*. Addison Wesley Longman Publishing Co., Inc., Redwood City, CA, USA.
- [47] Carlos Varela, Carlos Abalde, Laura Castro, and Jose GulÁjas. 2004. On Modelling Agent Systems with Erlang. In *Proceedings of the 2004 ACM SIGPLAN Workshop on Erlang (ERLANG '04)*. ACM, New York, NY, USA, 65–70. <https://doi.org/10.1145/1022471.1022481>
- [48] Ivan Vendrov, Christopher Dutschyn, and Nathaniel D. Osgood. 2014. Frabjous: A Declarative Domain-Specific Language for Agent-Based Modeling. In *Social Computing, Behavioral-Cultural Modeling and Prediction*, William G. Kennedy, Nitin Agarwal, and Shanchieh Jay Yang (Eds.). Number 8393 in Lecture Notes in Computer Science. Springer International Publishing, 385–392. https://doi.org/10.1007/978-3-319-05579-4_47
- [49] Zhanyong Wan and Paul Hudak. 2000. Functional Reactive Programming from First Principles. In *Proceedings of the ACM SIGPLAN 2000 Conference on Programming Language Design and Implementation (PLDI '00)*. ACM, New York, NY, USA, 242–252. <https://doi.org/10.1145/349299.349331>
- [50] Michael Wooldridge. 2009. *An Introduction to MultiAgent Systems* (2nd ed.). Wiley Publishing.
- [51] Bernard P. Zeigler, Herbert Praehofer, and Tag Gon Kim. 2000. *Theory of Modeling and Simulation: Integrating Discrete Event and Continuous Complex Dynamic Systems*. Academic Press. Google-Books-ID: REzmYOQmHuQC.

Delta Debugging Type Errors with a Blackbox Compiler

Joanna Sharrad
University of Kent
Canterbury, UK
jks31@kent.ac.uk

Olaf Chitil
University of Kent
Canterbury, UK
oc@kent.ac.uk

Meng Wang
University of Bristol
Bristol, UK
meng.wang@bristol.ac.uk

ABSTRACT

Debugging type errors is a necessary process that programmers, both novices and experts alike, face when using statically typed functional programming languages. All compilers often report the location of a type error inaccurately. This problem has been a subject of research for over thirty years. We present a new method for locating type errors: We apply the Isolating Delta Debugging algorithm coupled with a blackbox compiler. We evaluate our implementation for Haskell by comparing it with the output of the Glasgow Haskell Compiler; overall we obtain positive results in favour of our method of type error debugging.

CCS CONCEPTS

- Software and its engineering → Software testing and debugging;
- Theory of computation → Program analysis;

KEYWORDS

Type Error, Error diagnosis, Blackbox, Delta Debugging, Haskell

ACM Reference Format:

Joanna Sharrad, Olaf Chitil, and Meng Wang. 2018. Delta Debugging Type Errors with a Blackbox Compiler. In *Proceedings of International Symposium on Implementation and Application of Functional Languages (IFL'18)*. ACM, New York, NY, USA, Article 1, 11 pages. <https://doi.org/0000001.0000001>

1 INTRODUCTION

Compilers for Haskell, OCaml and many other statically typed functional programming languages produce type error messages that can be lengthy, confusing and misleading, causing the programmer hours of frustration during debugging. One role of such a type error message is to tell the programmer the location of a type error within the ill-typed program. Although there has been over thirty years of research [8, 22] on how to improve the way we locate type conflicts and present them to the programmer, type error messages can be misleading. We can trace the cause of inaccurate type error location to an advanced feature of functional languages: type inference.

A typical Haskell or OCaml program contains only little type information: definitions of data types, some type signatures for top-level functions and possibly a few more type annotations. Type inference works by generating constraints for the type of every

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

IFL'18, August 2019, Lowell, MA, USA

© 2019 Association for Computing Machinery.
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00
<https://doi.org/0000001.0000001>

expression in the program and solving these constraints. An ill-typed program is just a program with type constraints that have no solution. Because the type checker cannot know which program parts and thus constraints are correct, that is, agree with the programmer's intentions, it may start solving incorrect constraints and therefore assume wrong types early on. Eventually, the type checker may note a type conflict when considering a constraint that is actually correct.

1.1 Variations of an Ill-Typed Programs

Consider the following Haskell program from Stuckey et al. [17]:

```
1 insert x [] = x
2 insert x (y:ys) | x > y     = y : insert x ys
3                           | otherwise = x : y : ys
```

The program defines a function that shall insert an element into an ordered list, but the program is ill-typed. Stuckey et al. state that the first line is incorrect and should instead look like below:

```
1 insert x [] = [x]
```

The Glasgow Haskell Compiler (GHC) version 8.2.2 wrongly gives the location of the type error as (part of) line two.

```
2 insert x (y:ys) | x > y     = y : insert x ys
```

Let us see how GHC comes up with this wrong location. GHC derives type constraints and immediately solves them as far as possible. It roughly traverses our example program line by line, starting with line 1. The type constraints for line 1 are solvable and yield the information that `insert` is of type $\alpha \rightarrow [\beta] \rightarrow \alpha$. Subsequently in line 2 the expression $x > y$ yields the type constraint that x and y must have the same type, so together with the constraints for the function arguments x and $(y:ys)$, GHC concludes that `insert` must be of type $\alpha \rightarrow [\alpha] \rightarrow \alpha$. Finally, the occurrence of `insert x ys` as subexpression of $y : insert x ys$ means that the result type of `insert` must be the same list type as the type of its second argument. So `insert x ys` has both type $[\alpha]$ and type α , a contradiction reported as type error.

Our program contains no type annotations or signature, meaning we have to infer all types. Surely adding a type signature will ensure that GHC returns the desired type error location? Indeed for

```
1 insert :: Ord a => a -> [a] -> [a]
2 insert x [] = x
3 insert x (y:ys) | x > y     = y : insert x ys
4                           | otherwise = x : y : ys
```

GHC identifies the type error location correctly:

```
2 insert x [] = x
```

However, a recent study showed that type signatures are often wrong [23]. Wrong type signatures are the cause of 30% of all type errors! GHC trusts that a given type signature is correct and hence for

```
1 insert :: Ord a => a -> [a] -> a
2 insert x [] = x
3 insert x (y:ys) | x > y     = y : insert x ys
4                   | otherwise = x : y : ys
```

GHC wrongly locates the cause in line 2 again:

```
2 insert x (y:ys) | x > y     = y : insert x ys
```

In summary we see that the order in which type constraints are solved determines the reported type error location. There is no fixed order to always obtain the right type error location and requiring type annotations in the program does not help.

As a consequence researchers developed type error slicing [7, 16], which determines a minimal unsatisfiable type constraint set and reports all program parts associated with these constraints as type error slice. However, practical experience showed that these type error slices are often quite big [7] and thus they do not provide the programmer with sufficient information for correcting the type error. Our aim is to determine a smaller type error location, a single line in the program.

1.2 Our Method

Our method is based on the way programmers systematically debug errors without additional tools. The programmer removes part of the program, or adds previously removed parts back in. They check for each such variant of the program whether the error still exists or has gone. By doing this systematically, the programmer can determine a small part of the program as the cause of the error.

This general method was termed *Delta Debugging* by Zeller [24]. Specifically, we apply the *Isolating Delta Debugging algorithm*, which determines two variants of the original program that capture a minimal difference between a correct and erroneous variant of the program. Eventually our method produces the following result:

Result of our type error location method

```
1 insert x [] = x
2 insert x (y:ys) | x > y     = y : insert x ys
3                   | otherwise = x : y : ys
```

This program listing with different highlighting shows that the type error location is in line 1 and that line 1 and 2 together cause the type error; that is, even without line 3 this program is ill-typed.

The Isolating Delta Debugging algorithm has two prerequisites; An input that can repeatedly be modified and a means of inquiring if these modifications were successful. We fulfil the first prerequisite by employing the raw source code of the programmer's ill-typed program. We then work directly on the program text rather than the abstract syntax tree. We make modifications that generate new

variants of the program ready for testing to see if they remain ill-typed. To examine if they are indeed ill-typed or not, we employ the compiler as a black box. We do not use any location information included in any type error message of the compiler. This black box satisfies the second prerequisite of the Isolating Delta Debugging algorithm.

Once implemented in our tool Gramarye, we can apply our method to any ill-typed program, no matter how many type errors it contains, to locate one type error. Once our approach has the correct location, the programmer can fix it and reuse the tool to find further type errors.

Our tool Gramarye works on Haskell programs and uses the Glasgow Haskell Compiler as a black box. We evaluated Gramarye against the Glasgow Haskell Compiler using thirty programs containing single type errors and eight hundred and seventy programs generated to include two type errors.

Our paper makes the following contributions:

- We describe how to apply the Isolating Delta Debugging algorithm to type errors (Section 2).
- We use the compiler as a true black box; it can easily be replaced by a different compiler (Section 3.2).
- We implement the method in a tool called Gramarye that directly manipulates Haskell source code (Section 3.3).
- We evaluate our method against the Glasgow Haskell Compiler (Section 4).

Our evaluation shows an improvement in reporting type errors for many programs and demonstrates that our approach has promise in the field of type error debugging.

2 AN ILLUSTRATION OF OUR METHOD

Figure 1 gives an overview of the Gramarye framework. It indicates the steps taken to locate type errors in an ill-typed program.

We start with a single ill-typed Haskell program. This program must contain a type error; otherwise we reject it. Here we work with the original ill-typed program of the Introduction.

From this program, we obtain two programs that the *Isolating Delta Debugging* algorithm will work with. One is the ill-typed program, from which the algorithm removes lines that are irrelevant for the type error. The algorithm aims to minimise this program. The other program is the empty program, which is definitely well-typed and smaller than the ill-typed program. The algorithm moves lines from the ill-typed program to the well-typed program; the algorithm aims to maximise the well-typed program. The well-typed program shall always be smaller than the ill-typed program. So we start with:

Step 1: well-typed program

1
2
3

Step 1: ill-typed program

```
1 insert x [] = x
2 insert x (y:ys) | x > y     = y : insert x ys
3                   | otherwise = x : y : ys
```

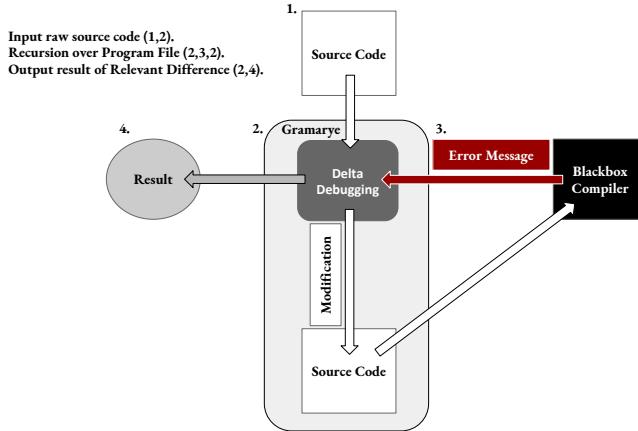


Figure 1: The Gramarye Framework

Now we move a line, we pick line 3, from the ill-typed program to the well-typed program, obtaining two new program variants:

Step 1: modified well-typed program

```

1
2
3 | otherwise = x : y : ys

```

Step 1: modified ill-typed program

```

1 insert x [] = x
2 insert x (y:ys) | x > y     = y : insert x ys
3

```

We then send these two programs to the black box compiler for type checking:

- Step 1: modified well-typed program: unresolved.
- Step 1: modified ill-typed program: ill-typed.

The modified well-typed program is not a syntactically valid Haskell program; the compiler yields a parse error. So note that our black box compiler yields one of three possible results:

- (1) unresolved; compiler yields a non-type error
- (2) ill-typed; compiler yields a type error
- (3) well-typed; compilation successful

We cannot use an unresolved program for locating a type error, but each of the other two possible results are useful. Our modified ill-typed program is smaller than our original ill-typed program. We now know that the modified variant is ill-typed too, so we can replace our ill-typed program for the next step:

Step 2: well-typed program

```

1
2
3

```

Step 2: ill-typed program

```

1 insert x [] = x
2 insert x (y:ys) | x > y     = y : insert x ys
3

```

The algorithm now repeats: Again we move a single line from the ill-typed program to the well-typed program. Let us pick line 2:

Step 2: modified well-typed program

```

1
2 insert x (y:ys) | x > y     = y : insert x ys
3

```

Step 2: modified ill-typed program

```

1 insert x [] = x
2
3

```

Again we type check these two programs:

- Step 2: modified well-typed program: well-typed.
- Step 2: modified ill-typed program: well-typed.

Because both variants are well-typed and bigger than the previous well-typed, we can use either of them as new well-typed program. We pick the modified well-typed program and thus obtain;

Step 3: well-typed program

```

1
2 insert x (y:ys) | x > y     = y : insert x ys
3

```

Step 3: ill-typed program

```

1 insert x [] = x
2 insert x (y:ys) | x > y     = y : insert x ys
3

```

The well-typed and ill-typed programs differ by only a single line, and hence our algorithm terminates.

The final result is that the difference between well-typed and ill-typed program, here line 1, is the location of the type error. Because the ill-typed program contains only lines 1 and 2 of the original program, we also know that only lines 1 and 2 are needed to make the program ill-typed. Thus we obtain the output shown in the Introduction. If we want to add a compiler type error message for further explanations, we can pick the one we received for the minimised ill-typed program. The error message may be clearer than for the original, larger program.

The isolating delta debugging method is non-deterministic. Often different choices lead to the same final result, but not always. Zeller argues that this non-determinism does not matter and that one result provides insightful debugging information to the programmer [25]. Hence his algorithm is deterministic, making arbitrary choices. Our implementation follows his algorithm and for our example makes the choices described here.

The algorithm is based on an ordering of programs, where a program is just a sequence of strings. A program P_1 is less or equal a program P_2 if they have the same number of lines and for every line, the line content is either the same for both programs, or the line is empty in P_1 . All programs that we consider are between the well-typed and ill-typed programs that we start with. The final well-typed and ill-typed programs have minimal distance, that is, they either differ by just one line or programs between them yield a non-type error at compilation (and thus are not syntactically valid programs).

In this example, in each step, we moved only a single line from the ill-typed to the well-typed program. For programs with hundreds of lines, this simple approach would be expensive in time due to, too many programs needing consideration. Hence we use the full Isolating Delta Debugging algorithm which starts with moving either the first or second half of the program from the ill-typed to well-typed program. If both modified programs are unresolved, then we increase the granularity of modifications from moving half the program to moving a quarter of the program. In general, every time both modified programs are unresolved, we half the size of our modifications. This increase of granularity can continue until only a single line is modified.

Zeller analysed the complexity of the isolating delta debugging algorithm. In our case complexity is the number of calls to the black box compiler in relation to the number of lines of the original program. In the worst case, if most calls yield unresolved, the number of calls is quadratic. In the best case, when no call yields unresolved, the number of calls is logarithmic [25].

3 IMPLEMENTATION

As illustrated in figure 1, our Gramarye tool has four components;

- Delta Debugging.
- Blackbox Compiler.
- Source Code Modification.
- Result Processing.

We shall next describe each of the components in greater detail.

3.1 Delta Debugging

The first component of our tool is *Delta Debugging*, a method formalised by Zeller [6, 24–26], that can be described as a systematic replication of the scientific approach of *Hypothesis-Test-Result* [25]. When programmers debug, they first use the error message to narrow the cause (hypothesis), then modify the source code and recompile (test), and lastly use the outcome of the recompilation (result) to see whether the modifications were successful. Zeller presents two delta debugging methods. He refers to them as Simplifying and Isolating [25].

3.1.1 Simplifying Delta Debugging.

The algorithm determines a smaller variant of the given failing (ill-typed) program. The result is minimal in that removing any further part makes the program pass (well-typed). The algorithm works by removing parts of a failing program until it no longer fails. The last failing variant of the program is the result of the algorithm. The minimality allows us to surmise that the parts of the program left must be the cause of our error. The Simplifying Delta Debugging algorithm has the same disadvantage as program slicing algorithms for type errors [7, 16]: it often reports a rather program. The second Delta Debugging algorithm, *Isolating*, aims to reduce the size of the reported program slice further.

3.1.2 Isolating Delta Debugging.

Isolating Delta Debugging incorporates the Simplifying algorithm to generate a minimal set of source code that contains an error. As well as employing the use of the simplifying algorithm, the isolating algorithm produces its own minimal set of source code; one that does not hold an error. The isolating minimal set is created by taking the working program and adding sections until the program reports an error. The aspect of having two minimal sets, one that contains the error and one that does not, is our reason for choosing the latter algorithm over the former. Focusing on the output of both minimal sets, we should receive a smaller result.

The Isolating Delta Debugging algorithm is composed of two parts; granularity and 'program replacement'. Granularity has the task of supplying which number of lines of the source code we are moving between our ill-typed and well-typed programs. Initially, granularity is set at two and applied in combination with the length of the program. The initial setting of granularity means it resembles a binary chop algorithm and when applied divides our program in half. After the initial application granularity is increased, decreased or remains static as the Delta Debugging algorithm iterates. We present a brief demonstration of how granularity works, independently from the algorithm, using a generic four-line program that contains a type error on line 4 below;

We first convert our lines of code into a list format;

[1, 2, 3, 4]

Our granularity currently equals 2. The initial divide splits our list in half;

[1, 2] [3, 4]

Isolating Delta Debugging checks the leading half first. These are the line numbers we shall modify in our program.

[1, 2]

The program is type checked with a blackbox compiler. The result is that there is no type error so, we check the second half;

[3, 4]

Again, type checking returns a well-typed result. We split the granularity and set it to 1, dividing our initial list into lists containing just one line number;

[1] [2] [3] [4]

Using the blackbox compiler we type check each list starting from the leading list. The lists containing 1, 2 and 3 are all well-typed. We then type-check the last list containing line 4 and it returns a type error.

Line [4] contains a [type error](#).

The increasing and decreasing of the granularity depends on the result category. We return a category when checking the success of the program modifications against our blackbox compiler; which we explain in more detail in section 3.2. Zeller does not use a blackbox compiler and as such assumes the use of a ‘testing function’ [24] to place the results into the following categories;

- (1) Test is undetermined (UNRESOLVED, ?)
- (2) Test has an error (FAIL, ×)
- (3) Test succeeds (PASS, ✓)

Our tool, on the other hand, categorises them slightly differently. Restricting the categories further due to the nature of only wanting to discover the position of type errors;

- (1) Test returns *unresolved* (non-type error, ?)
- (2) Test returns a *type error* (ill-typed, ×)
- (3) Test succeeds (well-typed; compilation was successful, ✓)

Program replacement also uses these categories to determine the path the algorithm takes after each iteration. The program files that the Isolating Delta Debugging algorithm uses are fluid. In our case, they start with the ‘Ill’ and ‘Well-Typed’ programs. As we iterate over the algorithm, the modified ill-typed and well-typed programs replace our initial programs depending on the result of the modifications.

Keeping our terminology from the example program (section 2) the changes seen in algorithm 1 are completed depending on the result of type checking with the blackbox compiler.

3.2 A Blackbox Compiler

We use a compiler as a blackbox, an entity of which we only know the input and the output. Anything that happens within the blackbox remains a mystery to us. Compilers naturally lend themselves to this usage, taking an input (source code), and returning an output; a successfully compiled program or error. The compiler we chose to use as a blackbox is the Glasgow Haskell Compiler (GHC),

ALGORITHM 1: Granularity and ‘Program Replacement’ section of DD

```
testModProgWell = test(modProgWell)
testModProgIll = test(modProgIll)
if testModProgIll == IllTyped && granularity == 2 then
| progIll = modProgIll
else if testModProgIll == WellTyped then
| progWell = modProgIll
| granularity = 2
else if testModProgWell == IllTyped then
| progIll = modProgWell
| granularity = 2.
else if testModProgWell == WellTyped then
| progWell = modProgWell
| granularity = max (granularity - 1) 2
else
| if n >= lineDiff progWell progIll then
| | terminate algorithm
| else
| | granularity = min (2*n) (lineDiff progWell progIll)
| | repeat algorithm
```

which is widely used by the Haskell community. As we can exploit GHC to gather type checking information without the need to alter the compiler itself, we can keep our tool separate. Not modifying the compiler has many benefits; changes made by the compiler developers will not affect the way our method works, users of our tool can avoid downloading a specialist compiler, and do not have the hassle of patching an existing one. Avoiding modification of the compiler also means that though we decided to employ Haskell in our initial investigation, our method is not restricted to this language, giving scope to expand to other functional languages such as OCaml.

We employ our blackbox compiler by using it as a type checker. During each iteration of the Isolating Delta Debugging algorithm, we determine the status of our modified ill-typed and well-typed programs as described in section 2. When using the blackbox compiler, our tool receives the same output a programmer would when they are using GHC. Though the result of compiling with GHC gives a message that includes many factors, we are only interested in if our programs are well-typed, using this information to attach the categories we discuss in section 3.1. Using our example from section 2, type checking both programs would give us the following messages and applied categories;

Ill-Typed Initial Error Message

Occurs check: cannot construct the infinite [type](#): a ~ [a]
....
Category: Ill-Typed

Well-Typed Initial Error Message

0
Category: Well-Typed

The Isolating Delta Debugging algorithm receives the attached categories and uses them to determine which path to apply as presented in algorithm 1. Depending on the route taken, we modify the source code of our programs in different ways, and again send them to the blackbox compiler for further type checking before reiterating over the whole method again. Where our programs source code is modified is automated by the Isolating Delta Debugging algorithm, but the idea of directly changing the raw code is solely inspired by how programmers manually debug.

3.3 Source Code Manipulation

When programmers manually debug they edit their source code directly, looking at where the error is suggested to occur and making changes in the surrounding area. We are also directly manipulating the source code, modifying our programs using the line numbers determined by the Isolating Delta Debugging algorithm. One significant bonus to the strategy of directly changing the source code is that it keeps our approach very simple. As we do not work on the Abstract Syntax Tree (AST) we do not need to parse our source code with each modification, allowing us to avoid making changes to an existing compiler or creating our own parser. Not editing the AST also means we can stay true to the programmer's original program, keeping personal preferences in layout intact by using empty lines as placeholders.

Our overall concern is the inaccurate reporting of the line number a type error occurred on, and as such, our tool works on a line-by-line based approach. As observed in section 2, we do this by adding and removing lines of source code. On completion of the algorithm we are left with two final programs; One program has all ill-typed source code removed, and the other only contains well-typed code. As we have directly modified the source code to achieve these two final programs we can use them to find the line number of where our type error appears by calculating the difference between the two.

3.4 Processing the Results

The idea is if one program is well-typed and the other ill-typed, then the source of the type error lies within the difference of the two; the relevant difference [25]. After the Isolating Delta Debugging algorithm has completed, two programs are left. The two final programs are used to create two lists generated by adding the line number of each empty line in the program. If a line number does not make an appearance in both of these generated lists we report it as a relevant difference. Working on whole lines of code in which we can report line numbers, also means we can easily evaluate how successfully we are in locating type errors.

4 EVALUATION

In the illustration of our method we have shown how we can successfully locate the correct line number of a type error. However, though positive for the example program we have used throughout, a more thorough evaluation was needed to be undertaken to show the strength of our method in type error locating.

We chose to evaluate our method against a benchmark of programs specially engineered to contain type errors. The programs

collated by Chen and Erwig [3] were used to assess their Counter-Factual approach to type error debugging. In all, there are one hundred and twenty-one programs in the CE benchmark, but not all had what Chen and Erwig called the 'oracle', the knowledge of where the type error lay. Though a program could have many correct solutions to remove a type error, we needed to know the correct location of where the type error occurred to evaluate accurately; so we removed all programs that did not specify the exact cause. To make our evaluation more compact, programs that were ill-typed in similar ways were also removed, reducing our set of test programs to thirty. However, as we also wanted to see whether our method could report multiple type errors, we took these thirty test programs and generated a further eight hundred and seventy programs to use in evaluation.

Our evaluation answers the following questions;

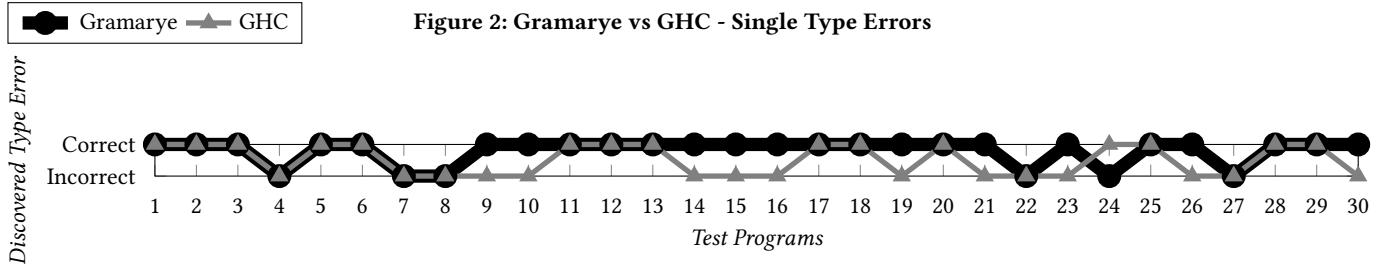
- (1) When applying our method to Haskell source code that contains a single type error; Do we show improvement in locating the errors compared to the Glasgow Haskell Compiler? (Section 4.1)
- (2) If we add multiple(two) type errors in our Haskell source code; Do we show improvement in locating these errors compared to the Glasgow Haskell Compiler? (Section 4.2)
- (3) Does our method return a smaller set of type error locations compared to the Glasgow Haskell Compiler? Specifically, a single precise line number of where the type error occurred. (Section 4.3)

Answering these questions involved creating a series of tests. To evaluate these tests we chose to compare our approach against GHC 8.2.2. We are using GHC as a blackbox compiler within our own tool, but as we use it solely as a type checker we do not have any knowledge of the line numbers it reports, and thus it has no interference with our evaluation. GHC and our tool take the CE benchmarks, and type checks each one; this results in a set of suggested line numbers where the cause of the type error could occur. To judge the success of locating the type error in the tests we have chosen to use the same criteria as Wand [22]. Wand states that even if we get multiple locations returned, the method is classed as a success if the exact location of the type error is within these. As, both our tool and GHC can report multiple line numbers for one type error; we use Wands criteria to allow us to take into consideration all line numbers returned, and not just the first.

4.1 Singular Type Error Evaluation

- (1) *When applying our method to Haskell source code that contains a single type error; Do we show improvement in locating the errors compared to the Glasgow Haskell Compiler?*

The first set of test programs contain one single type error; if the line number reported matched the 'oracle' response, then our result was accurate. In figure 2, we can see an overview of the outcome. The graph shows all thirty ill-typed programs and whether Gramarye and GHC correctly discovered the position of the type error. The results of our approach were positive. Out of the thirty ill-typed programs we accurately located 24 (80%) of the type errors, compared to 15(50%) from GHC.



In some cases, multiple line numbers were returned but still contained the correct errored line. We found the primary cause of multiple line numbers was due to statements that relied on each other or line breaks. Examples of this are If-Else or Let-In statements or lines that wrap around; the latter of which we present below;

Listing 1: Layout over two rows

```

1 doRow (y:ys) r = (if y < r && y > (r-dy) then '*' 
2           else ' ') : doRow r ys

```

In this example, our tool correctly identifies the line number even though we are returned two to choose from, but this was not the same for GHC, who suggests the first line in the above program as causing the issue. These issues caused by the programmer's layout decisions are one direction for future work.

In all our method using the initial evaluation criteria, has a 31% success rate over locating type errors in Haskell source code than GHC, but when programming we can often end up with multiple errors in our programs. A second evaluation of programs containing more than one type error would be advantageous.

4.2 Multiple Type Errors Evaluation

(2) *If we add multiple(two) type errors in our Haskell source code; Do we show improvement in locating these errors compared to the Glasgow Haskell Compiler?*

In the context of our evaluation, testing multiple type errors is represented by having more than one self-contained error within an ill-typed program. Self-contained type errors within a program mean we have two separate functions that do not interact with each other, with both functions contain a single type error. In listing 2, the first function has an error on line 2 and the second function on line 6, but neither type error affects the other;

Listing 2: Multiple Type Error Example

```

1 addList ls s = if s `elem` ls then ls else s : ls
2 v5 = addList "a" ["b"]
3
4 sumLists = sum2 . map sum2
5 sum2 [] = []
6 sum2 (x:xs) = x + sum2 xs

```

Listing 2, is just one of the programs we generated that contains multiple type errors. We created these by merging the CE benchmark programs. Each set of programs includes the original

source code with the addition of another CE program attached to the bottom. In all, we generated eight hundred and seventy new ill-typed programs to test. The success criteria for reporting an accurate discovery of the position of a type error in an ill-typed program that contains multiple errors is similar to what we used for singular errors. The only difference being, that though we have two errors per program we only need one error to be reported to deem a success.

Table 1, shows one set of results from a merged file. The first column lists the program number that we are using as the base and the second column indexes the number of the program we merged to the end of the source code. Under the Gramarye and GHC columns, we use ticks and crosses to denote if either correctly reports a type errors location, under this, we total the amount of correct matches as a percentage, the higher of which shows a greater success.

With this particular combination of CE benchmark programs, we can see that Gramarye finds 50% more type error positions than GHC. However, this is not always the case. Table 2, provides the total results for all of our combination of programs. Column one lists the base program, and the last two columns show the percentage of how accurate our tool and GHC were at locating type errors.

In total, we can see that Gramarye finds 3% fewer type errors in our multiple programs than GHC, this is not surprising. The Isolating Delta Debugging algorithm restricts Gramarye to always locating just one type error, the first it has come across. Once it has found this error, the algorithm assumes the job is complete and does not check any further. Currently, the programmer has to repeatedly use the tool after each implemented fix, working on each type error separately. We feel the removal of this limitation, would close the gap between Gramarye and GHC considerably, however, being restricted to working on one error at a time could also prove to be beneficial. Our evaluation of allowing a return of only a precise line shows this is the case.

4.3 Precise Type Error Evaluation

(3) *Does our method return a smaller set of type error locations? Specifically, a single precise line number of where the type error occurred.*

Though our criteria for success allowed us to check multiple returned line numbers for the correct type error position, reporting large amount of locations to the programmer is not ideal. As we aimed to return just a singular line number as the cause of the

Table 1: Testing a program with two type errors.

Original Program	Merged Program	Gramarye	GHC
15	1	✓	✓
15	2	✓	✓
15	3	✓	✗
15	4	✓	✗
15	5	✓	✓
15	6	✓	✓
15	7	✓	✗
15	8	✗	✗
15	9	✓	✗
15	10	✓	✗
15	11	✓	✓
15	12	✓	✓
15	13	✓	✓
15	14	✓	✗
15	16	✓	✗
15	17	✓	✓
15	18	✓	✓
15	19	✓	✗
15	20	✓	✓
15	21	✓	✗
15	22	✗	✗
15	23	✓	✗
15	24	✓	✓
15	25	✓	✗
15	26	✓	✗
15	27	✗	✗
15	28	✓	✓
15	29	✓	✓
15	30	✓	✗
Total		89.66%	44.83%

type error, an additional evaluation criteria allowed us to pinpoint how specific our tool was compared to GHC. All of the programs we tested had a single type error on a distinct line; our new rule specified that if either Gramarye or GHC returned a single accurate location, then they were classed as having a "precise success".

Table 3 shows all the program files that had a single type error; a tick denotes if either Gramarye or GHC accurately reports a single line number as being the cause of the type error. A report of multiple lines means a cross is displayed, even if a report of a correctly located type error was within them.

Our method had a positive outcome when locating a single line as the cause of the fault. Gramarye reported accurately 16 times (53%), with, GHC doing slightly worse at 12 times(40%).

When evaluating programs that included multiple self-contained type errors, we had a slightly different criteria, judging "precise success" under the following rules;

- A single line number containing the location of error one.

```

1 addList ls s = if s `elem` ls then ls else s : ls
2 v5 = addList "a" ["b"]
3
4 sumLists = sum2 . map sum2
5 sum2 [] = []
6 sum2 (x:xs) = x + sum2 xs

```

Table 2: Overall testing of programs with two type errors.

Program	Gramarye	GHC
1	72.41%	96.55%
2	65.52%	96.55%
3	68.97%	96.55%
4	75.86%	51.72%
5	68.97%	96.55%
6	72.41%	93.19%
7	65.52%	48.28%
8	62.07%	48.28%
9	75.86%	55.17%
10	58.62%	93.10%
11	65.52%	93.10%
12	68.97%	96.55%
13	68.97%	96.55%
14	75.86%	00.00%
15	89.66%	44.83%
16	65.52%	51.72%
17	68.97%	96.55%
18	65.52%	93.10%
19	82.76%	51.72%
20	68.97%	96.55%
21	82.76%	44.83%
22	31.03%	48.28%
23	72.41%	51.72%
24	55.17%	89.66%
25	68.97%	96.55%
26	72.41%	55.17%
27	65.52%	51.72%
28	65.52%	89.66%
29	65.52%	96.55%
30	65.52%	51.72%
Total		68.39% 71.03%

```

5 sum2 [] = []
6 sum2 (x:xs) = x + sum2 xs

```

- A single line number containing the location of error two.

```

1 addList ls s = if s `elem` ls then ls else s : ls
2 v5 = addList "a" ["b"]
3
4 sumLists = sum2 . map sum2
5 sum2 [] = []
6 sum2 (x:xs) = x + sum2 xs

```

- Two line numbers containing the location of both error one and two.

```

1 addList ls s = if s `elem` ls then ls else s : ls
2 v5 = addList "a" ["b"]
3
4 sumLists = sum2 . map sum2
5 sum2 [] = []
6 sum2 (x:xs) = x + sum2 xs

```

Table 3: "precise success" on single type errors.

Program	Gramarye	GHC
1	✓	✓
2	✗	✓
3	✓	✗
4	✗	✗
5	✗	✓
6	✗	✓
7	✗	✗
8	✗	✗
9	✓	✗
10	✓	✗
11	✗	✓
12	✓	✓
13	✓	✓
14	✗	✗
15	✗	✗
16	✓	✗
17	✓	✗
18	✓	✓
19	✗	✗
20	✓	✓
21	✓	✗
22	✗	✗
23	✓	✗
24	✗	✗
25	✓	✓
26	✓	✗
27	✗	✗
28	✓	✓
29	✗	✓
30	✓	✗
Total	53.33%	40.00%

All other results, even those that include the correct location, are recorded as failing the "precise success" criteria of discovering type errors. Table 4 represents the test programs that contained two type errors. The name of the original program along with the percentage of type error locations deemed to be a "precise success" are shown.

Analysing Table 4 we can see that our method is again successful in reporting the correct type error location using just one line number with 43% accuracy compared to GHC at 15%. GHC tends to report as many line numbers it feels are associated with the type error, very much like slicing. However, our evaluation shows that it may be more useful and accurate for the programmer to receive only one location at a time.

Overall, our evaluation has proven positive towards our method of type error debugging. From the testing, our strength lies in the reporting of singular type errors, be that one per program or the reporting of one instance of type error amongst many. Our results compared to GHC when testing more than one type error in a program suggests an algorithm that improves upon locating multiple types errors at a time could be beneficial. However, we

Table 4: "precise success" on programs with two type errors.

Program	Gramarye	GHC
1	48.28%	37.93%
2	44.83%	34.48%
3	48.28%	6.90%
4	51.72%	00.00%
5	44.83%	41.38%
6	37.93%	34.48%
7	44.83%	00.00%
8	41.38%	00.00%
9	51.72%	00.00%
10	48.28%	00.00%
11	48.28%	37.93%
12	48.28%	48.28%
13	48.28%	37.93%
14	34.48%	00.00%
15	13.79%	00.00%
16	44.83%	00.00%
17	51.72%	00.00%
18	41.38%	31.03%
19	10.34%	00.00%
20	44.83%	34.48%
21	72.41%	00.00%
22	20.69%	00.00%
23	41.38%	00.00%
24	37.93%	00.00%
25	48.28%	34.48%
26	48.28%	00.00%
27	44.83%	3.45%
28	41.38%	34.48%
29	37.93%	34.48%
30	48.28%	00.00%
Total	42.99%	15.06%

believe that several locations for one error is an unnecessary burden on the programmer, and a preference of accurate location over broad suggestion is preferential.

5 RELATED WORK

Type error debugging has taken many forms over the past thirty years; we will not be able to cover them all of them. Some core categories within type error debugging include: Slicing [7, 14, 18], Interaction [4, 5, 15, 16, 21], Type Inference Modification [1, 11], and working with Constraints [13, 27]. However, these solutions are complicated to implement. Some expect reliance on the compiler developers to accept the changes, for the programmer to patch their version or to use a particular compiler. Others do not provide an implementation to use and in the cases where there is an implementation, it is not maintained to work with the latest version of the programming language [9]. We, however, counter these by providing our approach within a tool, used separately from the compiler that employs the Delta Debugging algorithm to locate the type errors.

Delta Debugging, the name for two algorithms, one that simplifies and another that isolates, sparked our interest due to it's

closeness to debugging techniques that programmers use[6, 24–26]. There is only one other demonstration of the application of the Simplifying Delta Debugging algorithm to types, namely an application in the Liquid Haskell type checker [19]. Their approach differs from ours in that we concentrate on the Isolating Delta Debugging technique, combining the algorithm with direct modification of the programs source code, and using the Glasgow Haskell Compilers type checker as a blackbox.

Prior works that mention using the idea of a black box include; using the compiler's type inferencer as a black box to construct a type tree to use to debug the program [20], and having an SMT solver as a blackbox to return the satisfiable set of constraints to show type errored expressions[12]. SEMINAL, a tool which uses the type checker as a black box is the closest to our approach [9, 10]. Unlike our method though, SEMINAL along with previous solutions of using a blackbox compiler, makes modifications to an existing compiler. Though SEMINAL is also passing information, a patch is required for it to work with the OCaml compiler. Another difference between our tool, Gramarye, and SEMINAL is that SEMINAL modifies the Abstract Syntax Tree (AST), unlike our strategy of working directly on the source code itself.

Other approaches that talk about altering source code are a constraint-free tool inspired by SEMINAL, but though the author refers to source code modification the implementation also works with the AST [14]. Another tool TypeHope also discusses changing the source code of a program to stay true to how a programmer debugs. However, again, the implementation edits the AST [2]. At this point, as far as the authors know, modifying source code directly is a new approach in the type error debugging field.

6 CONCLUSION AND FUTURE WORK

Our method combines the Isolating Delta Debugging algorithm, a black box compiler and direct source code modification to locate type errors. Our tool Gramarye implements the method for Haskell using the Glasgow Haskell compiler as a black box. From our evaluation, we have gathered positive results that support our method for type error debugging. For single type errors our tool gives a 31% improvement over GHC. However, for two separate type errors in a single program GHC was 3% more successful. When applied to our aim of returning only a single line number for type errors, our method proved positive with 53% for locating singular type errors, and 43% when applied to a program that contained two type errors. A significant practical advantage of our method is that our tool Gramarye has only a small GHC-specific component and thus can easily be modified for other programming languages and compilers.

In the future we will be looking at where Gramarye did well and what its points of failure were. We will then use the outcome of the investigation to improve our algorithm for type error debugging. We will study closer the non-determinism of our method: can we sometimes determine whether one choice is better than another? After we have improved our method to determine the correct line number, we intend to increase the granularity of the tool further to eventually modify programs by single characters instead of lines, thus identifying subexpressions that cause type errors. On the theoretical side, there is clearly a close link between our method and methods described in the literature that perform type error

slicing based on minimal unsolvable constraint sets. We want to formalise that link.

Additional improvements to the tool outside of the algorithm would also be useful. An improved GUI, though not necessary for seeing if our approach is beneficial, does open up the options of not only combining with other methodologies that rely on interaction but also testing with real-life participants. We also would like to conduct empirical research of our solution in combination with evaluating against collected student programs to cement our strategy.

REFERENCES

- [1] Karen L Bernstein and Eugene W Stark. 1995. *Debugging type errors (full version)*. Technical Report. State University of New York at Stony Brook, Stony Brook, NY 11794-4400 USA. <https://pdfs.semanticscholar.org/814c/164c88ba7dd22e7e501ccda1951586a3117b.pdf>
- [2] Bernd Braßel. 2004. Typehope: There is hope for your type errors. In *Int. Workshop on Implementation of Functional Languages*. https://www.informatik.uni-kiel.de/~mh/wlp2004/final_papers/paper13.ps
- [3] Sheng Chen and Martin Erwig. 2014. Counter-factual typing for debugging type errors. In *The 41st Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL '14, San Diego, CA, USA, January 20-21, 2014*. 583–594. <https://doi.org/10.1145/2535838.2535863>
- [4] Sheng Chen and Martin Erwig. 2014. Guided Type Debugging. In *Functional and Logic Programming - 12th International Symposium, FLOPS 2014, Kanazawa, Japan, June 4-6, 2014. Proceedings*. 35–51. https://doi.org/10.1007/978-3-319-07151-0_3
- [5] Olaf Chitil. 2001. Compositional Explanation of Types and Algorithmic Debugging of Type Errors. In *Proceedings of the Sixth ACM SIGPLAN International Conference on Functional Programming (ICFP '01), Firenze (Florence), Italy, September 3-5, 2001*. 193–204. <https://doi.org/10.1145/507635.507659>
- [6] Holger Cleve and Andreas Zeller. 2005. Locating causes of program failures. In *27th International Conference on Software Engineering (ICSE 2005), 15-21 May 2005, St. Louis, Missouri, USA*. 342–351. <https://doi.org/10.1145/1062455.1062522>
- [7] Christian Haack and Joe B. Wells. 2004. Type error slicing in implicitly typed higher-order languages. *Sci. Comput. Program.* 50, 1-3 (2004), 189–224. <https://doi.org/10.1016/j.scico.2004.01.004>
- [8] Gregory F. Johnson and Janet A. Walz. 1986. A Maximum-Flow Approach to Anomaly Isolation in Unification-Based Incremental Type Inference. In *Conference Record of the Thirteenth Annual ACM Symposium on Principles of Programming Languages, St. Petersburg Beach, Florida, USA, January 1986*. 44–57. <https://doi.org/10.1145/512644.512649>
- [9] Benjamin S. Lerner, Matthew Flower, Dan Grossman, and Craig Chambers. 2007. Searching for type-error messages. In *Proceedings of the ACM SIGPLAN 2007 Conference on Programming Language Design and Implementation, San Diego, California, USA, June 10-13, 2007*. 425–434. <https://doi.org/10.1145/1250734.1250783>
- [10] Benjamin S. Lerner, Dan Grossman, and Craig Chambers. 2006. Seminal: searching for ML type-error messages. In *Proceedings of the ACM Workshop on ML, 2006, Portland, Oregon, USA, September 16, 2006*. 63–73. <https://doi.org/10.1145/1159876.1159887>
- [11] Bruce J McAdam. 1999. On the unification of substitutions in type inference. *Lecture notes in computer science* 1595 (1999), 137–152. https://link.springer.com/chapter/10.1007/3-540-48515-5_9
- [12] Zvonimir Pavlinovic. 2014. General Type Error Diagnostics Using MaxSMT. (2014). <https://pdfs.semanticscholar.org/1c14/7bc9f51cc950596dbc3e7cc5121202d160da.pdf>
- [13] Vincent Rahli, Joe B. Wells, John Pirie, and Fairouz Kamareddine. 2015. Skalpel: A Type Error Slicer for Standard ML. *Electr. Notes Theor. Comput. Sci.* 312 (2015), 197–213. <https://doi.org/10.1016/j.entcs.2015.04.012>
- [14] Thomas Schilling. 2011. Constraint-Free Type Error Slicing. In *Trends in Functional Programming, 12th International Symposium, TFP 2011, Madrid, Spain, May 16-18, 2011, Revised Selected Papers*. 1–16. https://doi.org/10.1007/978-3-642-32037-8_1
- [15] Eric L. Seidel, Ranjit Jhala, and Westley Weimer. 2016. Dynamic witnesses for static type errors (or, ill-typed programs usually go wrong). In *Proceedings of the 21st ACM SIGPLAN International Conference on Functional Programming, ICFP 2016, Nara, Japan, September 18-22, 2016*. 228–242. <https://doi.org/10.1145/2951913.2951915>
- [16] Peter J. Stuckey, Martin Sulzmann, and Jeremy Wazny. 2003. Interactive type debugging in Haskell. In *Proceedings of the ACM SIGPLAN Workshop on Haskell, Haskell 2003, Uppsala, Sweden, August 28, 2003*. 72–83. <https://doi.org/10.1145/871895.871903>

- [17] Peter J. Stuckey, Martin Sulzmann, and Jeremy Wazny. 2004. Improving type error diagnosis. In *Proceedings of the ACM SIGPLAN Workshop on Haskell, Haskell 2004, Snowbird, UT, USA, September 22-22, 2004*. 80–91. <https://doi.org/10.1145/1017472.1017486>
- [18] Frank Tip and T. B. Dinesh. 2001. A slicing-based approach for locating type errors. *ACM Trans. Softw. Eng. Methodol.* 10, 1 (2001), 5–55. <https://doi.org/10.1145/366378.366379>
- [19] A Tondwalkar. 2016. *Finding and Fixing Bugs in Liquid Haskell*. Master's thesis, University of Virginia. <https://pdfs.semanticscholar.org/79b4/22959847253c40aff25c228205372d9ebc60.pdf>
- [20] Kanae Tsushima and Kenichi Asai. 2012. An Embedded Type Debugger. In *Implementation and Application of Functional Languages - 24th International Symposium, IFL 2012, Oxford, UK, August 30 - September 1, 2012, Revised Selected Papers*. 190–206. https://doi.org/10.1007/978-3-642-41582-1_12
- [21] Kanae Tsushima and Olaf Chitil. 2014. Enumerating Counter-Factual Type Error Messages with an Existing Type Checker. In *16th Workshop on Programming and Programming Languages, PPL2014*. <http://kar.kent.ac.uk/49007/>
- [22] Mitchell Wand. 1986. Finding the Source of Type Errors. In *Conference Record of the Thirteenth Annual ACM Symposium on Principles of Programming Languages, St. Petersburg Beach, Florida, USA, January 1986*. 38–43. <https://doi.org/10.1145/512644.512648>
- [23] Baijun Wu and Sheng Chen. 2017. How Type Errors Were Fixed and What Students Did?. In *ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications*.
- [24] Andreas Zeller. 1999. Yesterday, My Program Worked. Today, It Does Not. Why?. In *Software Engineering - ESEC/FSE'99, 7th European Software Engineering Conference, Held Jointly with the 7th ACM SIGSOFT Symposium on the Foundations of Software Engineering, Toulouse, France, September 1999, Proceedings*. 253–267. https://doi.org/10.1007/3-540-48166-4_16
- [25] Andreas Zeller. 2009. *Why Programs Fail - A Guide to Systematic Debugging, 2nd Edition*. Academic Press. <http://store.elsevier.com/product.jsp?isbn=9780123745156&pageName=search>
- [26] Andreas Zeller and Ralf Hildebrandt. 2002. Simplifying and Isolating Failure-Inducing Input. *IEEE Trans. Software Eng.* 28, 2 (2002), 183–200. <https://doi.org/10.1109/32.988498>
- [27] Danfeng Zhang, Andrew C Myers, Dimitrios Vytiniotis, and Simon Peyton-Jones. 2015. *Diagnosing Haskell type errors*. Technical Report. Technical Report <http://hdl.handle.net/1813/39907>, Cornell University. <https://pdfs.semanticscholar.org/d32f/81a5c1706e225e2255b72c1e4b41f799e8f1.pdf>

Received May 2018

HiPERJiT: A Profile-Driven Just-in-Time Compiler for Erlang

Konstantinos Kallas
University of Pennsylvania
Philadelphia, USA
kallas@seas.upenn.edu

Konstantinos Sagonas
Uppsala University
Uppsala, Sweden
kostis@it.uu.se

ABSTRACT

We introduce HiPERJiT, a profile-driven Just-in-Time compiler for the BEAM ecosystem based on HiPE, the High Performance Erlang compiler. HiPERJiT uses runtime profiling to decide which modules to compile to native code and which of their functions to inline and type-specialize. HiPERJiT is integrated with the runtime system of Erlang/OTP and preserves aspects of Erlang’s compilation which are crucial for its applications: most notably, tail-call optimization and hot code loading at the module level. We present HiPERJiT’s architecture, describe the optimizations that it performs, and compare its performance with BEAM, HiPE and Pyrlang. HiPERJiT offers performance which is about two times faster than BEAM and almost as fast as HiPE, despite the profiling and compilation overhead that it has to pay compared to an ahead-of-time native code compiler. But there also exist programs for which HiPERJiT’s profile-driven optimizations allow it to surpass HiPE’s performance.

1 INTRODUCTION

Erlang is a concurrent functional programming language with features that support the development of scalable concurrent and distributed applications, and systems with requirements for high availability and responsiveness. Its main implementation, the Erlang/OTP system, comes with a byte code compiler for its virtual machine, called BEAM, which produces portable and reasonably efficient code. For applications with requirements for better performance, an ahead-of-time native code compiler called HiPE (High Performance Erlang) can be selected. In fact, byte code and native code can happily coexist in the Erlang/OTP runtime system.

Despite this flexibility, the selection of the modules of an application to compile to native code is currently manual. Perhaps it would be better if the system itself could decide on which parts to compile to native code in a just-in-time fashion. Moreover, it would be best if this process was guided by profiling information gathered during runtime, and was intelligent enough to allow for the continuous run-time optimization of the code of an application.

This paper makes a first big step in that direction. We have developed HiPERJiT, a profile-driven Just-in-Time compiler for the BEAM ecosystem based on HiPE. HiPERJiT employs the tracing support of the Erlang runtime system to profile the code of bytecode-compiled modules during execution and choose whether to compile these modules to native code, currently employing all optimizations that HiPE also performs by default. For the chosen modules, it additionally decides which of their functions to inline and/or specialize based on runtime type information. We envision that the list of additional optimizations to perform based on profiling information will be extended in the future, especially if HiPERJiT grows to become a JIT compiler which performs lifelong feedback-directed optimization of programs. Currently, JIT compilation is triggered

only once for each loaded instance of a module, and the profiling of their functions is stopped at that point.

The main part of this paper presents the architecture of HiPERJiT and the rationale behind some of our design decisions (Section 3), the profile-driven optimizations that HiPERJiT performs (Section 4), and the performance it achieves (Section 5). Before all that, in the next section, we review the current landscape of Erlang compilers. Related work is scattered throughout.

2 ERLANG AND ITS COMPILERS

In this section, we present a brief account of the various compilers for Erlang. Our main aim is to set the landscape for our work and present a comprehensive high-level overview of the various compilers that have been developed for the language, rather than a detailed technical one. For the latter, we refer the reader to the sites of these compilers and to publications about them.

2.1 JAM and BEAM

Technically, JAM and BEAM are virtual (abstract) machines for Erlang, not compilers. However, they both have lent their name to compilers that generate byte code for these machines. JAM is the older one, but since 1998, when the Erlang/OTP system became open source, the system has been exclusively based on the BEAM.

The BEAM compiler, is a fast, module-at-a-time compiler that produces relatively compact byte code, which is then loaded into the Erlang Run-Time System and expanded to indirectly threaded code for the VM interpreter (called “emulator” in the Erlang community). In the process, the BEAM loader specializes and/or merges byte code instructions. By now, the BEAM comes with a well-engineered VM interpreter offering good performance for many Erlang applications. As a result, some other languages (e.g., Elixir) have chosen to use BEAM as their target. In recent years, the term “BEAM ecosystem” has been used to describe these languages, and also to signify that the BEAM is not important only for Erlang.

2.2 HiPE and ErLLVM

The High-Performance Erlang (HiPE) compiler [13, 23] is an ahead-of-time native code compiler for Erlang. It has backends for SPARC, x86 [21], x86_64 [20], PowerPC, PowerPC 64 and ARM, has been part of the Erlang/OTP distribution since 2002, and is mature by now. HiPE aims to improve the performance of Erlang applications by allowing users to compile their time-critical modules to native code, and offer flexible integration between interpreted and native code. Since 2014, the ErLLVM compiler [24], which generates native code for Erlang using the LLVM compiler infrastructure [15] (versions 3.5 or higher), has also been integrated into HiPE as one of its backends, selectable by the `to_llvm` compiler option. In

general, ErLLVM offers similar performance to the native HiPE backends [24].

Figure 1 shows how the HiPE compiler fits into Erlang/OTP. The compilation process typically starts from BEAM byte code. HiPE uses three intermediate representations, namely Symbolic BEAM, Icode, and RTL, and then generates target-specific assembly either directly or outsources its generation to ErLLVM.

Icode is a register-based intermediate language for Erlang. It supports an infinite number of registers which are used to store arguments and temporaries. All values in Icode are proper Erlang terms. The call stack is implicit and preserves registers. Finally, as Icode is a high-level intermediate language, bookkeeping operations (such as heap overflow checks and context switching) are implicit.

RTL is a generic three-address register transfer language. Call stack management and the saving and restoring of registers before and after calls are made explicit in RTL. Heap overflow tests are also made explicit and are propagated backwards to limit performance overhead. Registers in RTL are separated in tagged and untagged, where the untagged registers hold raw integers, floating-point numbers, and addresses.

RTL code is then translated to machine-specific assembly code (or to LLVM IR), virtual registers are assigned to real registers, symbolic references to atoms and functions are replaced by their real values in the running system, and finally, the native code is loaded into the Erlang Run-Time System (ERTS) by the HiPE loader.

ERTS allows seamless interaction between interpreted and native code. However, there are differences in the way interpreted and native code execution is managed, therefore a transition from native to interpreted code (or vice versa), triggers a *mode switch*. Mode switches occur in function calls, returns, and exception throws between natively-compiled and interpreted functions. HiPE was designed with the goal of no runtime overhead as long as the execution mode stays the same, so mode switches are handled by instrumenting calls with special mode switch instructions and by adding extra call frames that cause a mode switch after each function return. Frequent mode switching can negatively affect execution performance. Therefore, it is recommended that the most frequently executed functions of an application are all executed in the same mode [23].

2.3 BEAMJIT and Pyrlang

Unsurprisingly, JIT compilation has been investigated in the context of Erlang several times in the past, both distant and more recent. For example, both Jerico [12], the first native code compiler for Erlang (based on JAM) circa 1996, and early versions of HiPE contained some support for JIT compilation. However, this support never became robust to the point that it could be used in production.

More recently, two attempts to develop tracing JITs for Erlang have been made. The first of them, BEAMJIT [6], is a tracing just-in-time compiling runtime for Erlang. BEAMJIT uses a tracing strategy for deciding which code sequences to compile to native code and the LLVM toolkit for optimization and native code emission. It extends the base BEAM implementation with support for profiling, tracing, native code compilation, and support for switching between these three (profiling, tracing, and native) execution modes. Performance-wise, back in 2014, BEAMJIT reportedly managed to reduce the execution time of some small Erlang benchmarks

by 25–40% compared to BEAM, but there were also many other benchmarks where it performed worse than BEAM [6]. Moreover, the same paper reported that “HiPE provides such a large performance improvement compared to plain BEAM that a comparison to BEAMJIT would be uninteresting” [6, Sect. 5]. At the time of this writing (May 2018), BEAMJIT is not yet a complete implementation of Erlang; for example, it does not yet support floats. Although work is ongoing in extending BEAMJIT, its overall performance, which has improved, does not yet surpass that of HiPE.¹ Since BEAMJIT is not available, not even as a binary, we cannot compare against it.

The second attempt, Pyrlang [11], is an alternative virtual machine for the BEAM byte code which uses RPython’s meta-tracing JIT compiler [5] as a backend in order to improve the sequential performance of Erlang programs. Meta-tracing JIT compilers are tracing JIT compilers that trace and optimize an interpreter instead of the program itself. The Pyrlang paper [11] reports that Pyrlang achieves average performance which is 38.3% faster than BEAM and 25.2% slower than HiPE, on a suite of sequential benchmarks. Currently, Pyrlang is a research prototype and not yet a complete implementation of Erlang, even for the sequential part of the language. On the other hand, unlike BEAMJIT, Pyrlang is available, and we will directly compare against it in Section 5.

2.4 Challenges of Compiling Erlang

The Erlang programming language comes with some special characteristics which make its efficient compilation quite challenging. On the top of the list is *hot code loading*: the requirement to be able to replace modules, on an individual basis, while the system is running and without imposing a long stop to its operation. The second characteristic, which is closely related to hot code loading, is that the language makes a semantic distinction between module-local calls, which have the form `f(...)`, and so called *remote calls* which are module-qualified, i.e., have the form `m:f(...)`, and need to look up the most recently loaded version of module `m`, even from a call within `m` itself.

These two characteristics, combined with the fact that Erlang is primarily a concurrent language in which a large number of processes may be executing code from different modules at the same time, effectively impose that compilation happens in a module-at-a-time fashion, without opportunities for cross-module optimizations. The alternative, i.e., performing cross-module optimizations, implies that the runtime system must be able to quickly undo optimizations in the code of some module that rely on information from other modules, whenever those other modules change. Since such undoings can cascade arbitrarily deep, this alternative is not very attractive engineering-wise.

The runtime system of Erlang/OTP supports hot code loading in a particular, arguably quite ad hoc, way. It allows for *up to two* versions of each module (m_{old} and $m_{current}$) to be simultaneously loaded, and redirects all remote calls to $m_{current}$, the most recent of the two. Whenever m_{new} , a new version of module m , is about to be loaded, all processes that still execute code of m_{old} are abruptly terminated, $m_{current}$ becomes the new m_{old} , and m_{new} becomes $m_{current}$. This quite elaborate mechanism is implemented by the

¹Lukas Larsson, private communication, May 2018.

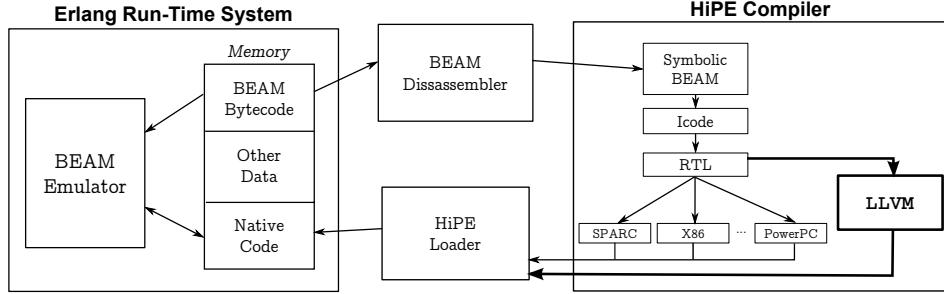


Figure 1: Architecture of some Erlang/OTP components: ERTS, HiPE and ErLLVM (from [24]).

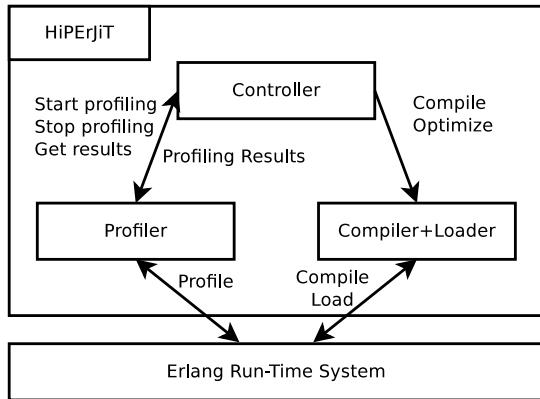


Figure 2: High-level architecture of HiPERJiT.

code loader with support from ERTS, which has control over all Erlang processes.

Unlike BEAMJIT and Pyrlang, we decided, at least for the time being, to leave the Erlang Run-Time System unchanged as far as hot code loading is concerned. This also means that, like HiPE, the unit of JIT compilation of HiPERJiT is the entire module.

3 HIPERJIT

The functionality of HiPERJiT can be briefly described as follows. It profiles executed modules, maintaining runtime data such as execution time and call graphs. It then decides which modules to compile and optimize based on the collected data. Finally, it compiles and loads the JIT-compiled modules in the runtime system. Each of these tasks is handled by a separate component.

Controller The central controlling unit which decides which modules should be profiled and which should be optimized based on runtime data.

Profiler An intermediate layer between the controller and the low-level Erlang profilers. It gathers profiling information, cleans up and organizes it, and sends it back to the controller for further processing.

Compiler+Loader An intermediate layer between the controller and HiPE. It compiles the modules chosen by the controller and then loads them into the runtime system.

The architecture of the HiPERJiT compiler can be seen in Fig. 2.

3.1 Controller

The controller, as stated above, is the fundamental component of HiPERJiT. It chooses the modules to profile, and uses the profiling data to decide which modules to compile and optimize. It is essential that no user input is needed to drive decision making. Our design extends a lot of concepts from the work on Jalapeño JVM [1].

Traditionally, many JIT compilers use a lightweight call frequency method to drive compilation [3]. This method maintains a counter for each function and increments it every time the function is called. When the counter surpasses a specified threshold, JIT compilation is triggered. This approach, while having very low overhead, does not give precise information about the program execution.

Instead, HiPERJiT makes its decision based on a simple cost-benefit analysis. Compilation for a module is triggered when its predicted future execution time when compiled, would be less than its future execution time when interpreted combined with the compilation time:

$$\text{FutureExecTime}_c < \text{FutureExecTime}_i + \text{CompTime}$$

Of course, it is not possible to predict the future execution time of a module or the time needed to compile it, so some estimations need to be made. First of all, we need to estimate future execution time. The results of a study about the lifetime of UNIX processes [8] show that the total execution time of UNIX processes follows a Pareto (heavy-tailed) distribution. The mean remaining waiting time of this distribution is analogous to the amount of time that has passed already. Thus, assuming that the analogy between a module and a UNIX process holds, we consider future execution time of a module to be equal to its execution time until now $\text{FutureExecTime} = \text{ExecTime}$. In addition, we consider that compiled code has a constant relative speedup to the interpreted code, thus $\text{ExecTime}_c * \text{Speedup}_c = \text{ExecTime}_i$. Finally, we consider that the compilation time for a module depends linearly on its size, so $\text{CompTime} = C * \text{Size}$. In short, the condition to check is:

$$\frac{\text{ExecTime}_i}{\text{Speedup}_c} < \text{ExecTime}_i + C * \text{Size}$$

If this condition holds, the module is “worth” compiling.

We conducted two experiments in order to find suitable values for the condition parameters Speedup_c and C . In the first one, we executed a set of Erlang applications before and after compiling their modules to native code. The average speedup we measured was 2 and we used it as an estimate for the Speedup_c parameter. In

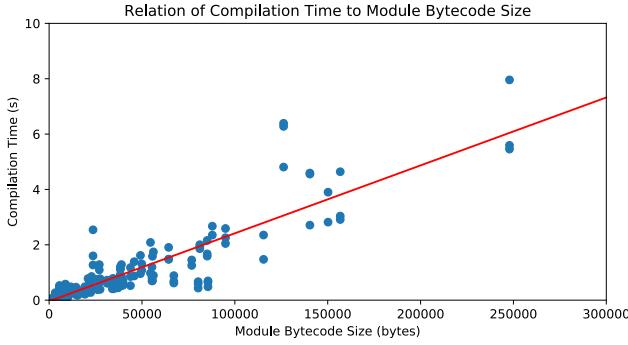


Figure 3: Estimating compilation time as a function of module bytecode size.

the second experiment, we compiled a set of modules of varying sizes, measured their compilation time, and fitted a line to the set of given points using the Least Squares method (cf. Fig. 3).

The line has a slope of $2.5e^{-5}$ so that is also the estimated compilation cost per byte of byte code. It is worth mentioning that, in reality, the compilation time does not depend only on module size, but on many other factors (e.g., branching, exported functions, etc.). However, we consider this estimate to be adequate for our goal.

Finally, HiPERJiT uses a feedback-directed scheme to improve the precision of the compilation decision condition when possible. HiPERJiT measures the compilation time of each module it compiles to native code and stores it in a persistent key-value storage, where the key is a pair of the module name and the MD5 hash value of its byte code. If the same version of that module is considered for compilation at a subsequent time, HiPERJiT will use the stored compilation time measurement in place of *CompTime*.

3.2 Profiler

The profiler is responsible for efficiently profiling executed code. Its architecture, which can be seen in Fig. 4, consists of:

- The profiler core, which receives the profiling commands from the controller and transfers them to ERTS. It also receives the profiling results from the individual profilers and transfers them back to the controller.
- The router, which receives all profiling messages from ERTS and routes them to the correct individual profiler.
- The individual profilers, which handle profiling messages, aggregate them, and transfer the results to the profiler core. Each individual profiler handles a different subset of the execution data, namely function calls, execution time, argument types, and process lifetime.

We designed the profiler in a way that facilitates the addition and removal of individual profilers.

3.2.1 Profiling Execution Time. In order to profile the execution time that is being spent in each function, the profiler receives a time-stamped message from ERTS for each function call, function return, and process scheduling action. Messages have the form $\langle Action, Timestamp, Process, Function \rangle$, where *Action* can be any of the following values: *call*, *return_to*, *sched_in*, *sched_out*.

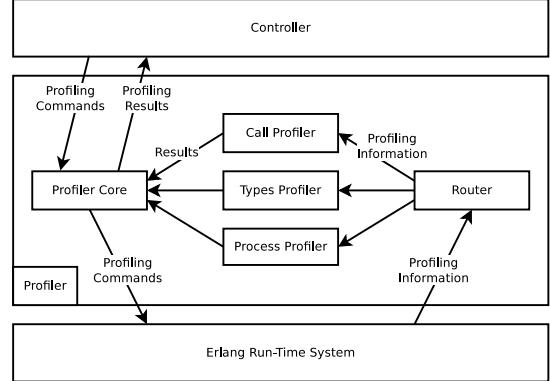


Figure 4: Profiler architecture and its components.

For each sequence of trace messages that refer to the same process *P*, i.e., a sequence of form $[(\text{sched_in}, T_1, P, F_1), (A_2, T_2, P, F_2), \dots, (A_{n-1}, T_{n-1}, P, F_{n-1}), (\text{sched_out}, T_n, P, F_n)]$, we can compute the time spent at each function by finding the difference of the timestamps in each pair of adjacent messages, so $[(T_2 - T_1, F_1), (T_3 - T_2, F_2), \dots, (T_n - T_{n-1}, F_{n-1})]$. The sum of all the differences for each function gives the total execution time spent in each function. An advantage of this method is that it allows us to also track the number of calls between each pair of functions, which is later used for profile-driven call inlining (Section 4.2).

The execution time profiler is able to handle the stream of messages sent during the execution of sequential programs. However, in concurrent programs, the rate at which messages are sent to the profiler can increase uncontrollably, thus flooding its mailbox, leading to high memory consumption.

In order to tackle this problem, we implemented a probing mechanism that checks whether the number of unhandled messages in the mailbox of the execution time profiler exceeds a threshold (currently 1 million messages), in essence checking whether the arrival rate of messages is higher than their consumption rate. If the mailbox exceeds this size, no more trace messages are sent from ERTS until the profiler handles the remaining messages in the mailbox.

3.2.2 Type Tracing. The profiler is also responsible for type tracing. As we will see, HiPERJiT contains a compiler pass that uses type information for the arguments of functions in order to create type-specialized, and hopefully better performing, versions of functions. This type information is extracted from type specifications in Erlang source code (if they exist) and from runtime type information which is collected and aggregated as described below.

For modules that the profiler has selected for profiling, the arguments of function calls to this module are also profiled. The ERTS low-level profiling functionality returns the complete argument values, so in principle their types can be determined. However, functions could be called with complicated data structures as arguments, and determining their types precisely requires a complete traversal. As this could lead to a significant performance overhead, the profiler approximates their types by a limited-depth term traversal. In other words, *depth-k type abstraction* is used. Every

type T is represented as a tree with leaves that are either singleton types or type constructors with zero arguments, and internal nodes that are type constructors with one or more arguments. The depth of each node is defined as its distance from the root. A depth- k type T_k is a tree where every node with depth $\geq k$ is pruned and over-approximated with the top type (`any()`). For example, in Erlang's type language [18], which supports unions of singleton types, the Erlang term `{foo,{bar,[{a,17},{a,42}]}}` has type `{'foo',{bar,list({'a},17|42)}}`, where '`A`' denotes the singleton atom type `A`. Its depth-1 and depth-2 type abstractions are `{'foo',{any(),any()}}` and `{'foo',{bar,list(any())}}`.

The following two properties should hold for depth- k type abstractions:

- (1) $\forall k. k \geq 0 \Rightarrow T \sqsubseteq T_k$ where \sqsubseteq is the subtype relation.
- (2) $\forall t. t \neq T \Rightarrow \exists k. \forall i. i \geq k \Rightarrow t \not\sqsubseteq T_i$

The first property guarantees correctness while the second allows us to improve the approximation precision by choosing a greater k , thus allowing us to trade performance for precision or vice versa.

Another problem that we faced is that Erlang collections can contain elements of different types. Traversing them in order to find their complete type could also create undesired overhead. As a result, we decided to optimistically estimate the collection element types. Although Erlang collections can contain elements of many different types, this is rarely the case in practice. In most programs, collections usually contain elements of the same type. This, combined with the fact that HiPERJiT uses the runtime type information to specialize some functions for specific type arguments, gives us the opportunity to be more optimistic while deducing the types of values. Therefore, we decided to approximate the types of Erlang collections (currently lists and maps) by considering only a small subset of their elements.

What is left is a way to generalize and aggregate the type information acquired through the profiling of many function calls. Types in Erlang are internally represented using a subtyping system [18], thus forming a type lattice. Because of that, there exists a supremum operation that can be used to aggregate type information that has been acquired through different traces.

3.2.3 Profiling Processes. Significant effort has been made to ensure that the overhead of HiPERJiT does not increase with the number of concurrent processes. Initially, performance was mediocre when executing concurrent applications with a large number (> 100) of spawned processes because of profiling overhead. While profiling a process, every function call triggers an action that sends a message to the profiler. The problem arises when many processes execute concurrently, where a lot of execution time is needed by the profiler process to handle all the messages being sent to it. In addition, the memory consumption of the profiler skyrocketed as messages arrived with a higher rate than they were consumed.

To avoid such problems, HiPERJiT employs *probabilistic profiling*. The idea is based on the following observation. Massively concurrent applications usually consist of a lot of sibling processes that have identical functionality. Because of that, it is reasonable to sample a part of the hundreds (or even thousands) of processes and only profile them to get information about the system as a whole.

This sample should not be taken at random, but it should rather follow the principle described above. We maintain a process tree by

monitoring the lifetime of all processes. The nodes of the process tree are of the following type:

```
-type ptnode() :: {pid(), mfa(), [ptnode()]}.
```

The `pid()` is the process identifier, the `mfa()` is the initial function of the process, and the `[ptnode()]` is a list of its children processes, which is empty if it is a leaf. The essence of the process tree is that sibling subtrees which share the same structure and execute the same initial function usually have the same functionality. Thus we could only profile a subset of those subtrees and get relatively accurate results. However, finding subtrees that share the same structure can become computationally demanding. Instead, we decided to just find leaf processes that were spawned using the same MFA and group them. So, instead of tracing all processes of each group, we just profile a randomly sampled subset. The sampling strategy that we used is very simple but still gives satisfactory results. When the processes of a group are more than a specified threshold (currently 50) we sample only a percentage (currently 10%) of them. The profiling results for these subsets are then scaled accordingly (i.e., multiplied by 10) to better represent the complete results.

3.3 Compiler+Loader

This is the component that is responsible for the compilation and loading of modules. It receives the collected profiling data (i.e., type information and function call data) from the controller, formats them, and calls an extended version of the HiPE compiler to compile the module and load it. The HiPE compiler is extended with two additional optimizations that are driven by the collected profiling data; these optimizations are described in Section 4.

There are several challenges that HiPERJiT has to deal with, when loading an Erlang module, as also mentioned in Section 2.4. First of all, Erlang supports hot code loading, so a user can reload the code of a module while the system is running. If this happens, HiPERJiT automatically triggers the process of re-profiling this module.

Currently ERTS allows only up to two versions of the same module to be simultaneously loaded. This introduces a rather obscure but still possible race condition between HiPERJiT and the user. If the user tries to load a new version of a module after HiPERJiT has started compiling the current one, but before it has completed loading it, HiPERJiT could load JIT-compiled code for an older version of the module after the user has loaded the new one. Furthermore, if HiPERJiT tries to load a JIT-compiled version of a module m when there are processes executing m_{old} code, this could lead to processes being killed. Thus, HiPERJiT loads a module m as follows. First, it acquires a module lock, not allowing any other process to reload this specific module during that period. Then, it calls HiPE to compile the target module to native code. After compilation is complete, HiPERJiT repeatedly checks whether any process executes m_{old} code. When no process executes m_{old} code, the JIT-compiled version is loaded. Finally, HiPERJiT releases the lock so that new versions of module m can be loaded again. This way, HiPERJiT avoids that loading JIT-compiled code leads to processes being killed.

Of course, the above scheme does not offer any progress guarantees, as it is possible for a process to execute m_{old} code for an indefinite amount of time. In that case, the user would be unable to

reload the module (e.g., after fixing a bug). In order to avoid this, HiPERJiT stops trying to load a JIT-compiled module after a user-defined timeout (the default value is 10 seconds), thus releasing the module lock.

Let us end this section with a brief note about decompilation. Most JIT compilers support decompilation, which is usually triggered when a piece of JIT-compiled code is not executed frequently enough anymore. The main benefit of doing this is that native code takes up more space than byte code, so decompilation often reduces the memory footprint of the system. However, since code size is not a big concern in today's machines, and since decompilation can increase the number of mode switches (which are known to cause performance overhead) HiPERJiT currently does not support decompilation.

4 PROFILE-DRIVEN OPTIMIZATIONS

In this section, we describe two optimizations that HiPERJiT performs, based on the collected profiling information, in addition to calling HiPE: type specialization and inlining.

4.1 Type Specialization

In dynamic languages, there is typically very little type information available during compilation. Types of values are determined at runtime and because of that, dynamic language compilers generate code that handles all possible value types by adding type tests to ensure that operations are performed on terms of correct type. Also all values are tagged, which means that their runtime representation contains information about their type. The combination of these two features leads to dynamic language compilers generating less efficient code than those for statically typed languages.

This problem has been tackled with static type analysis [16, 23], runtime type feedback [10], or a combination of both [14]. Runtime type feedback is essentially a mechanism that gathers type information from calls during runtime, and uses this information to create specialized versions of the functions for the most commonly used types. On the other hand, static type analysis tries to deduce type information from the code in a conservative way, in order to simplify it (e.g., eliminate some unnecessary type tests). In HiPERJiT, we employ a combination of these two methods.

4.1.1 Optimistic Type Compilation. After profiling the argument types of function calls as described in Section 3.2.2, the collected type information is used. However, since this information is approximate or may not hold for all subsequent calls, the optimistic type compilation pass adds type tests that check whether the arguments have the appropriate types. Its goal is to create specialized (optimistic) function versions, whose arguments are of known types.

Its implementation is straightforward. For each function f where the collected type information is non-trivial:

- The function is duplicated into an optimized $f\$opt$ and a “standard” $f\$std$ version of the function.
- A header function that contains all the type tests is created. This function performs all necessary type tests for each argument, to ensure that they satisfy any assumptions upon which type specialization may be based. If all tests pass, the header calls $f\$opt$, otherwise it calls $f\$std$.

- The header function is inlined in every local function call of the specified function f . This ensures that the type tests happen on the caller’s side, thus improving the benefit from the type analysis pass that happens later on.

4.1.2 Type Analysis. Optimistic type compilation on its own does not offer any performance benefit. It simply duplicates the code and forces execution of possibly redundant type tests. Its benefits arise from its combination with type analysis and propagation.

Type analysis is an optimization pass performed on Icode that infers type information for each program point and then simplifies the code based on this information. It removes checks and type tests that are unnecessary based on the inferred type information. It also removes some boxing and unboxing operations from floating-point computations. A brief description of the type analysis algorithm [18] is as follows:

- (1) Construct the call graph for all the functions in a module and sort it topologically based on the dependencies between its strongly connected components (SCCs).
- (2) Analyze the SCCs in a bottom-up fashion using a constraint-based type inference to find the most general success typings [19] under the current constraints.
- (3) Analyze the SCCs in a top-down order using a data-flow analysis to propagate type information from the call sites to module-local functions.
- (4) Add new constraints for the types, based on the propagated information from the previous step.
- (5) If a fix-point has not been reached, go back to step 2.

Initial constraints are mostly generated using the type information of Erlang Built-In Functions (BIFs) and functions from the standard library which are known to the analysis. Guards and pattern matches are also used to generate type constraints.

After type analysis, code optimizations are performed. First, all redundant type tests or other checks are completely removed. Then some instructions, such as floating-point arithmetic, are simplified based on the available type information. Finally, control flow is simplified and dead code is removed.

It is important to note that type analysis and propagation is restricted to the module boundary. No assumptions are made about the arguments of the exported functions, as those functions can be called from modules which are not yet present in the system or available for analysis. Thus, modules that export only few functions benefit more from this analysis as more constraints are generated for their local functions and used for type specializations.

4.1.3 An Example. We present a simple example that illustrates type specialization. Listing 1 contains a function that computes the power of two values, where the base of the exponentiation is a number (an integer or a float) and the exponent is an integer.

```
-spec power(number(), integer(), number()) -> number().
power(_V1, 0, V3) -> V3;
power(V1, V2, V3) -> power(V1, V2-1, V1*V3).
```

Listing 1: The source code of a simple power function.

Without the optimistic type compilation, HiPE generates the Icode shown in Listing 2.

```

power/3(v1, v2, v3) ->
12:
  v5 := v3
  v4 := v2
  goto 1
1:
  _ := redtest()      (primop)
  if is_{integer,0}(v4) then goto 3 else goto 10
3:
  return(v5)
10:
  v8 := '-'(v4, 1)   (primop)
  v9 := '*'(v1, v5)  (primop)
  v5 := v9
  v4 := v8
  goto 1

```

Listing 2: Generated Icode for the power function.

If we consider that this function is mostly called with a floating point number as the first argument, then HiPE with optimistic type compilation generates the Icode in Listing 3. Note that `power$std` has the same Icode as `power` without optimistic type compilation.

```

power/3(v1, v2, v3) ->
16:
  _ := redtest()      (primop)
  if is_float(v1) then goto 3 else goto 14
3:
  if is_integer(v2) then goto 5 else goto 14
5:
  if is_float(v3) then goto 7 else goto 14
7:
  power$opt/3(v1, v2, v3)
14:
  power$std/3(v1, v2, v3)

power$opt/3(v1, v2, v3) ->
24:
  v5 := v3
  v4 := v2
  goto 1
1:
  _ := redtest() (primop)
  if is_{integer,0}(v4) then goto 3 else goto 20
3:
  return(v5)
20:
  v8 := '-'(v4, 1)  (primop)
  _ := gc_test<3>() (primop) -> goto 28, #fail 28
28:
  fv10 := unsafe_untag_float(v5) (primop)
  fv11 := unsafe_untag_float(v1) (primop)
  _ := fclearerror()           (primop)
  fv12 := fp_mul(fv11, fv10)   (primop)
  _ := fcheckerror()           (primop)
  v5 := unsafe_tag_float(fv12) (primop)
  v4 := v8
  goto 1

```

Listing 3: Generated Icode for the power function with optimistic type compilation.

As it can be seen, optimistic type compilation has used type propagation and found that both `v1` and `v3` are always floats. The code was then modified so that they are untagged unsafely in every iteration and multiplied using a floating point instruction, whereas without optimistic type compilation they are multiplied using the standard general multiplication function, which checks the types of both arguments and untags (and unboxes) them before performing the multiplication.

4.2 Inlining

Inlining is the process of replacing a function call with the body of the called function. This improves performance in two ways. First, it mitigates the function call overhead. Second, and most important, it enables more optimizations to be performed later, as most optimizations usually do not cross function boundaries. That is the reason why inlining is usually performed in early phases of compilation so that later phases become more effective.

However, aggressive inlining has several potential drawbacks, both in compilation time, as well as in code size increase (which in turn can also lead to higher execution times because of caching effects). However, code size is less of a concern in today's machines, except of course in application domains where memory is not abundant (e.g., in IoT or embedded systems).

In order to get a good performance benefit from inlining, the compiler must achieve a fine balance between inlining every function call and not inlining anything at all. Therefore, the most important issue when inlining is choosing which function calls to inline. There has been a lot of work on how to make this decision with compile-time [22, 25, 26, 28] as well as run-time information in the context of JIT compilers [4, 7, 9, 27]. HiPERJiT makes its inlining decisions based on run-time information, but also keeps its algorithm fairly lightweight so as to not to impose a big overhead.

4.2.1 Inlining Decision. Our mechanism borrows two ideas from previous work, the use of call frequency [4] and call graphs [27] to guide the inlining decision. Recall that HiPERJiT compiles whole modules at a time, thus inlining decisions are made based on information from all the functions in a module. Due to hot code loading, inlining across modules is not performed.

Finding the most profitable, performance-wise, function calls to inline and also the most efficient order in which to inline them is a heavy computational task and thus we decided to use heuristics to greedily decide which function calls to inline and when. The call frequency data that are used is the number of calls between each pair of function, as also described in Section 3.2.

Before describing our mechanism in detail we briefly describe the main idea and rationale behind its design.

Decisions are based on the assumption that call sites which are visited the most are the best candidates for inlining. Because of that our mechanism greedily chooses to inline the function pair (F_1, F_2) with most calls from F_1 to F_2 .

Greedy inlining might introduce performance overheads as some calls might be inlined multiple times. For example, consider the case where our mechanism decides to inline the following pairs in that order $[(F_1, F_2), (F_2, F_3), (F_1, F_3)]$. In that case function F_3 would be inlined two times in the function F_2 , once in the original F_2 and once in the inlined copy of F_2 in F_1 . This could have been prevented if F_3 were first inlined in F_2 and then F_2 were inlined in F_1 . However, the analysis needed to achieve this is not cost effective. Thus, the inlining decision is made greedily, despite occasionally inlining a function more than once.

Of course, inlining has to be restrained so that it does not happen for every function call in the program. We achieve this by restricting the maximum code size of each module. Small modules are allowed to grow up to double their size, while larger ones are only allowed to grow by 10%.

4.2.2 Implementation. Our inlining implementation is fairly standard except for some details specific to Icode. The steps below describe the inlining process. In case the inlined call is a tail-call only the first two steps need to be performed.

- (1) Variables and labels in the body of the callee are updated to fresh ones so that there is no name clash with the variables and labels in the body of the caller.
- (2) The arguments of the function call are moved to the parameters of the callee.
- (3) Tail-calls in the body of the callee are transformed to a normal call and a return
- (4) Return instructions in the body of the callee are transformed into a move and a goto instruction that moves the return value to the destination of the call and jumps to the end of the body of the callee.

It is worth mentioning that Erlang uses cooperative scheduling which is implemented by making sure that processes are executed for a number of reductions and then yield. When the reductions of a process are depleted, the process is suspended and another process is scheduled instead. Reductions are implemented in Icode by the `redtest()` primitive operation (primop) in the beginning of each function body; cf Listing 2. When inlining a function call, the call to `redtest()` in the body of the callee is also removed. This is safe to do, because inlining is bounded anyway.

5 EVALUATION

In this section, we evaluate the performance of HiPERJiT against BEAM, which serves as a baseline for our comparison, and against two systems that also aim to surpass the performance of BEAM. The first of them is the HiPE compiler.² The second is the Pyrlang³ meta-tracing JIT compiler for Erlang, which is a research prototype and not a complete implementation; we report its performance on the subset of benchmarks that it can handle. We were not able to obtain a version of BEAMJIT to include in our comparison, as this system is still (May 2018) not publicly available. However, as mentioned in Section 2.3, BEAMJIT does not achieve performance that is superior to HiPE anyway.

We conducted all experiments on a laptop with an Intel Core i7-4710HQ @ 2.50GHz CPU and 16 GB of RAM running Ubuntu 16.04.

5.1 Profiling Overhead

Our first set of measurements concerns the profiling overhead that HiPERJiT imposes. To obtain a rough worst-case estimate, we used a modified version of HiPERJiT that profiles programs as they run but does not compile any module to native code. Note that this scenario is very pessimistic for HiPERJiT, because the common case is that JIT compilation will be triggered for some modules, HiPERJiT will stop profiling them at that point, and these modules will then most likely execute faster, as we will soon see. But even if native code compilation were not to trigger for any module, it would be very easy for HiPERJiT to stop profiling after a certain time period has passed or some other event (e.g., the number of calls that have been profiled has exceeded a threshold) has occurred.

²For both BEAM and HiPE, we used the ‘master’ branch of Erlang/OTP 21.0.

³We used the latest version of Pyrlang (<https://bitbucket.org/hrc706/pyrlang/overview>) built using PyPy 5.0.1.

In any case, our measurements showed that the overhead caused by probabilistic profiling is around 10%. More specifically, the overhead we measured ranged from 5% (in most cases) up to 40% for some heavily concurrent programs.

We also separately measured the profiling overhead on all concurrent benchmarks with and without probabilistic profiling, to measure the benefit of probabilistic over standard profiling. The average overhead that standard profiling imposes on concurrent benchmarks is 19%, while the average slowdown of probabilistic profiling on such benchmarks is 13%.

5.2 Evaluation on Small Benchmark Programs

The benchmarks we used come from the ErLLVM benchmark suite⁴, which has been previously used for the evaluation of HiPE [13], ErLLVM [24], BEAMJIT [6], and Pyrlang [11]. The benchmarks can be split in two sets: (1) a set of small, relatively simple but quite representative Erlang programs, and (2) the set of Erlang programs from the Computer Language Benchmarks Game (CLBG)⁵ as they were when the ErLLVM benchmark suite was created.

Comparing the performance of a Just-in-Time with an ahead-of-time compiler on small benchmarks is tricky, as the JIT compiler also pays the overhead of compilation during the program’s execution. Moreover, a JIT compiler needs some time to warm up. For this reason, we report three numbers for HiPERJiT: the speedup achieved when considering all overheads, the speedup achieved when disregarding the first run, and the speedup achieved in the last 50% of many runs, when a steady state has been reached. We also use two configurations of HiPE: one with the maximum level of optimization (o3), and one where static inlining (using the `{inline_size, 25}` compiler directive) has been performed besides o3.

The speedup of HiPERJiT, HiPE, and Pyrlang compared to BEAM for the small benchmarks is shown in Figs. 5 and 6. (We have split them into two figures based on scale of the y-axis, the speedup over BEAM.) Note that all speedups we report are averages of several different executions. The overall average speedup for each configuration is summarized in Table 1.

Table 1: Speedup over BEAM for the small benchmarks.

Configuration	Speedup
HiPE	2.08
HiPE + Static inlining	2.17
Pyrlang	1.05
HiPERJiT	1.85
HiPERJiT w/o 1st run	2.08
HiPERJiT last 50% runs	2.33

Overall, the performance of HiPERJiT is almost two times better than BEAM and Pyrlang, but slightly worse than HiPE. However, there also exist five benchmarks (barnes, nrev, fib, smith and tak) where HiPERJiT, in its steady state (the last 50% of runs), surpasses HiPE’s performance. This strongly indicates that the profile-driven

⁴<https://github.com/cstavr/erllvm-bench>

⁵<http://benchmarksgame.alioth.debian.org/>

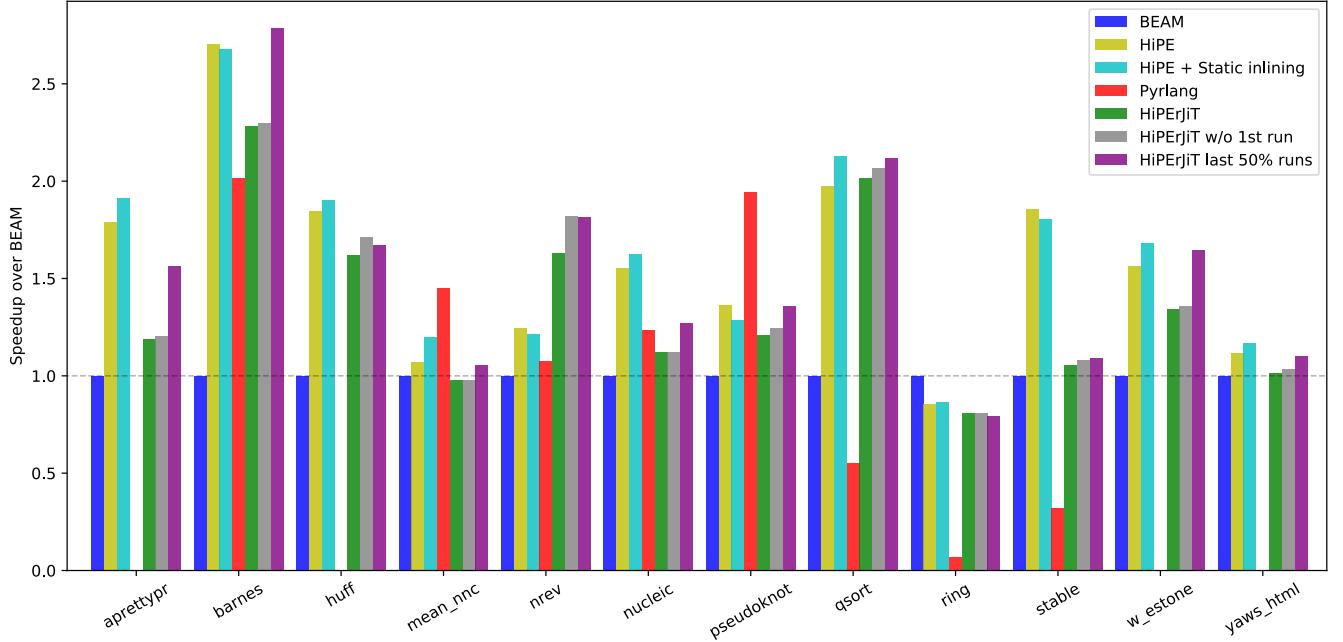


Figure 5: Speedup over BEAM on small benchmarks.

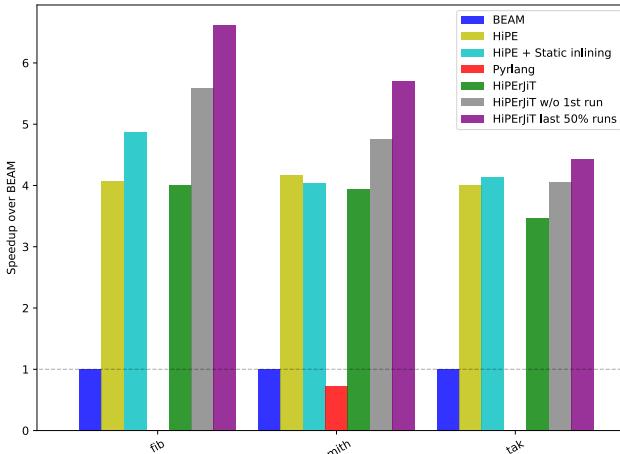


Figure 6: Speedup over BEAM on small benchmarks.

optimizations that HiPERJiT performs lead to more efficient code, compared to that of an ahead-of-time native code compiler performing static inlining.

Out of those five benchmarks, the most interesting one is *smith*. It is an implementation of the Smith-Waterman DNA sequence matching algorithm. The reason why HiPERJiT offers better speedup over HiPE is that profile-driven inlining manages to inline functions `alpha_beta_penalty/2` and `max/2` in `match_entry/5`, which is the most time-critical function of the program, thus allowing further optimizations to improve performance. Static inlining on the other

Table 2: Speedup over BEAM for the CLBG programs.

Configuration	Speedup
HiPE	2.05
HiPE + Static inlining	1.93
HiPERJiT	1.78
HiPERJiT w/o 1st run	2.14
HiPERJiT last 50% runs	2.26

hand does not inline `max/2` in `match_entry/5` even if one chooses very large values for the inliner's thresholds.

On the *ring* benchmark (Fig. 5), which is heavily concurrent, HiPERJiT performs worse than both HiPE and BEAM. To be more precise, all systems perform worse than BEAM on this benchmark. The main reason is that the majority of the execution time is spent on message passing (around 1.5 million messages are sent per second), which is part of BEAM's runtime system. Finally, there are a lot of process spawns and exits (around 20 thousand per second), which leads to considerable profiling overhead, as HiPERJiT maintains and constantly updates the process tree.

Another interesting benchmark is *stable* (also Fig. 5), where HiPERJiT performs slightly better than BEAM but visibly worse than HiPE. This is mostly due to the fact that there are a lot of process spawns and exits in this benchmark (around 240 thousand per second), which leads to significant profiling overhead.

The speedup of HiPERJiT and HiPE compared to BEAM for the CLBG benchmarks is shown in Figs. 7 and 8. The overall average speedup for each configuration is shown in Table 2.

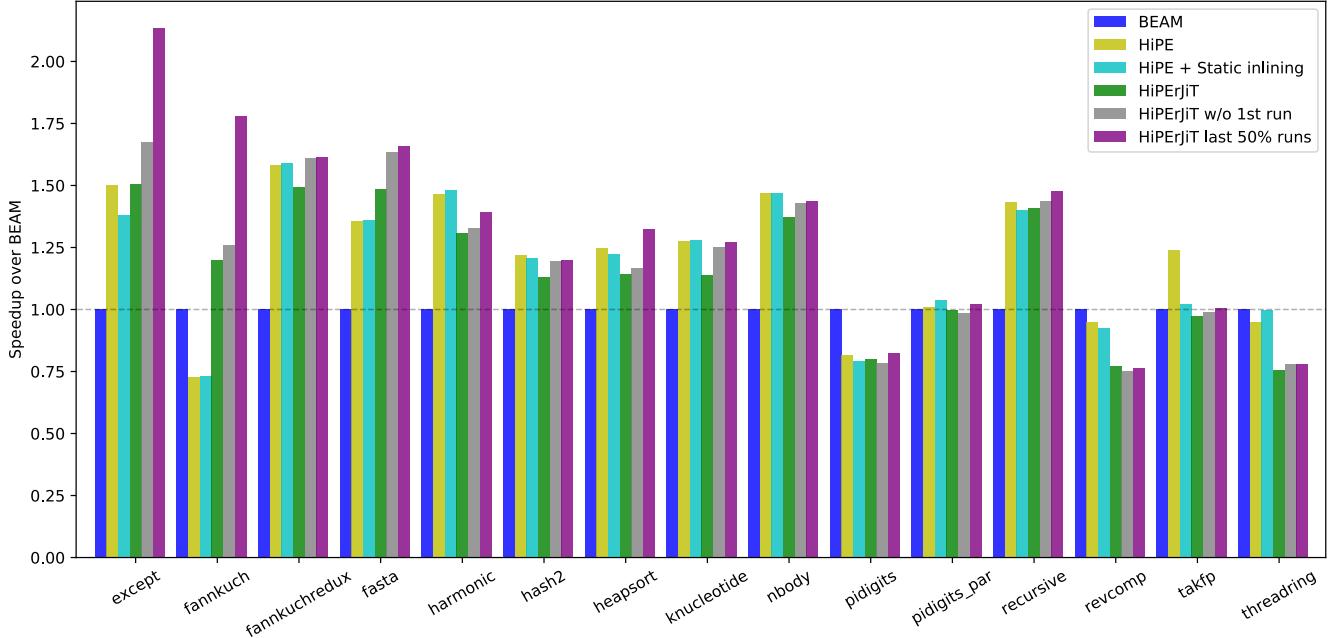


Figure 7: Speedup over BEAM on the CLBG programs.

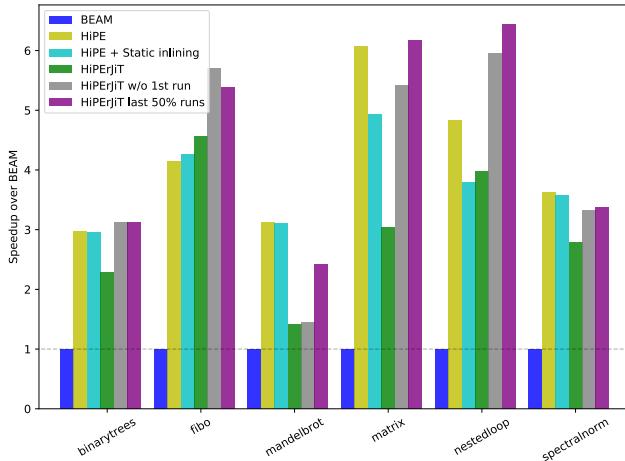


Figure 8: Speedup over BEAM on the CLBG programs.

As with the small benchmarks, the performance of HiPERJiT mostly lies between BEAM and HiPE. When excluding the first run, HiPERJiT outperforms both BEAM and HiPE in several benchmarks (binarytrees, except, fannkuch, fannkuchredux, fasta, fibo, nestedloop, and recursive). Finally, if we only consider the steady state of the JiT Compiler (last 50% of the runs), HiPERJiT outperforms both HiPE and BEAM in one more benchmark (matrix). However, HiPERJiT's performance on some benchmarks (revcomp, takfp, and threadring) is worse than both BEAM and HiPE because the profiling overhead is higher than the benefit of compilation.

5.3 Evaluation on a Bigger Program

Besides small benchmarks, we also evaluate the performance of HiPERJiT on a program of considerable size and complexity, as results in small or medium-sized benchmarks may not always provide a complete picture for the expected performance of a compiler. The Erlang program we chose, the Dialyzer [17] static analysis tool, is both big (about 30,000 LOC) and complex and highly concurrent [2]. It has also been heavily engineered over the years and comes with hard-coded knowledge of the set of 26 modules it needs to compile to native code upon its start to get maximum performance for most use cases. Using an appropriate option, the user can disable this native code compilation phase, which takes more than a minute on the laptop we use, and in fact this is what we do to get measurements for BEAM and HiPERJiT.

We use Dialyzer in two ways. The first builds a Persistent Lookup Table (PLT) containing cached type information for all modules under erts, compiler, crypto, hipe, kernel, stdlib and syntax_tools. The second analyzes these applications for type errors and other discrepancies. The speedups of HiPERJiT and HiPE compared to BEAM for the two use ways of using Dialyzer are shown in Table 3.

The results show that HiPERJiT achieves performance which is better than BEAM's but worse than HiPE's. There are various reasons for this. First of all, Dialyzer is a complex application where functions from many different modules of Erlang/OTP (50–100) are called with arguments of significant size (e.g., the source code of the applications that are analyzed). This leads to considerable tracing and bookkeeping overhead. Second, some of the called modules contain very large, often compiler-generated, functions and their compilation takes considerable time (the total compilation time is about 70 seconds, which is a significant portion of the total time).

Table 3: Speedups over BEAM for two Dialyzer use cases.

Configuration	Dialyzer Speedup	
	Building PLT	Analyzing
HiPE	1.73	1.70
HiPE + Static inlining	1.75	1.72
HiPERJiT	1.51	1.30
HiPERJiT w/o 1st run	1.58	1.35
HiPERJiT last 50% runs	1.62	1.37

Finally, HiPERJiT does not compile all modules from the start, which means that a percentage of the time is spent running interpreted code and performing mode switches which are more expensive than same-mode calls. In contrast, HiPE has hard-coded knowledge of “the best” set of modules to compile to native code before the analysis starts.

6 CONCLUDING REMARKS AND FUTURE WORK

We have presented HiPERJiT, a profile-driven Just-in-Time compiler for the BEAM ecosystem based on the HiPE native code compiler. It offers performance which is better than BEAM’s and comparable to HiPE’s, on most benchmarks. Aside from performance, we have been careful to preserve features such as hot code loading that are considered important for Erlang’s application domain, and have made design decisions that try to maximize the chances that HiPERJiT remains easily maintainable and in-sync with components of the Erlang/OTP implementation. In particular, besides employing the HiPE native code compiler for most of its optimizations, HiPERJiT uses the same concurrency support that the Erlang Run-Time System provides, and relies upon the tracing infrastructure that it offers. Thus, it can straightforwardly profit from any improvements that may occur in these components.

Despite the fact that the current implementation of HiPERJiT is quite robust and performs reasonably well, profile-driven JIT compilers are primarily engineering artifacts and can never be considered completely “done”. Besides making the current optimizations more effective and implementing additional ones, one item which is quite high on our “to do” list is investigating techniques that reduce the profiling overhead, especially in heavily concurrent applications. For the current state of HiPERJiT, the profiling overhead is bigger than we would like it to be but, on the other hand, it’s not really a major concern because once HiPERJiT decides to compile a module to native code the profiling of its functions stops and the overhead drops to zero. But it will become an issue if HiPERJiT becomes a JIT compiler that performs lifelong feedback-directed optimization of programs, which is a direction that we want to pursue.

REFERENCES

- [1] Matthew Arnold, Stephen Fink, David Grove, Michael Hind, and Peter F Sweeney. 2000. Adaptive Optimization in the Jalapeño JVM. In *Proceedings of the 15th ACM SIGPLAN Conference on Object-oriented Programming, Systems, Languages, and Applications (OOPSLA ’00)*. ACM, New York, NY, USA, 47–65. <https://doi.org/10.1145/353171.353175>
- [2] Stavros Aronis and Konstantinos Sagonas. 2013. On Using Erlang for Parallelization – Experience from Parallelizing Dialyzer. In *Trends in Functional Programming, 13th International Symposium, TFP 2012, Revised Selected Papers (LNCS)*, Vol. 7829. Springer, 295–310. https://doi.org/10.1007/978-3-642-40447-4_19
- [3] John Aycock. 2003. A Brief History of Just-in-time. *ACM Comput. Surv.* 35, 2 (2003), 97–113. <https://doi.org/10.1145/857076.857077>
- [4] Andrew Ayers, Richard Schoeller, and Robert Gottlieb. 1997. Aggressive Inlining. In *Proceedings of the ACM SIGPLAN 1997 Conference on Programming Language Design and Implementation (PLDI ’97)*. ACM, New York, NY, USA, 134–145. <https://doi.org/10.1145/258915.258928>
- [5] Carl Friedrich Bolz, Antonio Cuni, Maciej Fijalkowski, and Armin Rigo. 2009. Tracing the Meta-level: PyPy’s Tracing JIT Compiler. In *Proceedings of the 4th Workshop on the Implementation, Compilation, Optimization of Object-Oriented Languages and Programming Systems (ICOOLPS ’09)*. ACM, New York, NY, USA, 18–25. <https://doi.org/10.1145/1565824.1565827>
- [6] Frej Dreibhammar and Lars Rasmussen. 2014. BEAMJIT: A Just-in-time Compiling Runtime for Erlang. In *Proceedings of the Thirteenth ACM SIGPLAN Workshop on Erlang (Erlang ’14)*. ACM, New York, NY, USA, 61–72. <https://doi.org/10.1145/2633448.263350>
- [7] Brian Grant, Matthai Philipose, Markus Mock, Craig Chambers, and Susan J Eggers. 1999. An Evaluation of Staged Run-time Optimizations in DyC. In *Proceedings of the ACM SIGPLAN 1999 Conference on Programming Language Design and Implementation (PLDI ’99)*. ACM, New York, NY, USA, 293–304. <https://doi.org/10.1145/301618.301683>
- [8] Mor Harchol-Balter and Allen B Downey. 1997. Exploiting Process Lifetime Distributions for Dynamic Load Balancing. *ACM Trans. Comput. Syst.* 15, 3 (1997), 253–285. <https://doi.org/10.1145/263326.263344>
- [9] Kim Hazelwood and David Grove. 2003. Adaptive Online Context-sensitive Inlining. In *Proceedings of the International Symposium on Code Generation and Optimization: Feedback-directed and Runtime Optimization (CGO ’03)*. IEEE Computer Society, Washington, DC, USA, 253–264. <http://dl.acm.org/citation.cfm?id=776261.776289>
- [10] Urs Hözle and David Ungar. 1994. Optimizing Dynamically-dispatched Calls with Run-time Type Feedback. In *Proceedings of the ACM SIGPLAN 1994 Conference on Programming Language Design and Implementation (PLDI ’94)*. ACM, New York, NY, USA, 326–336. <https://doi.org/10.1145/178243.178278>
- [11] Ruochen Huang, Hidehiko Masuhara, and Tomoyuki Aotani. 2016. Improving Sequential Performance of Erlang Based on a Meta-tracing Just-In-Time Compiler. In *Post-Proceeding of the 17th Symposium on Trends in Functional Programming*. https://ftp2016.org/papers/TFP_{-}2016_{-}paper_{-}16.pdf
- [12] Erik Johansson and Christer Jonsson. 1996. *Native Code Compilation for Erlang*. Technical Report. Computing Science Department, Uppsala University.
- [13] Erik Johansson, Mikael Pettersson, and Konstantinos Sagonas. 2000. A High Performance Erlang System. In *Proceedings of the 2nd ACM SIGPLAN International Conference on Principles and Practice of Declarative Programming (PPDP ’00)*. ACM, New York, NY, USA, 32–43. <https://doi.org/10.1145/351268.351273>
- [14] Madhukar N Kedlaya, Jared Roesch, Behnam Robatmili, Mehrdad Reshad, and Ben Hardekopf. 2013. Improved Type Specialization for Dynamic Scripting Languages. In *Proceedings of the 9th Symposium on Dynamic Languages (DLS ’13)*. ACM, New York, NY, USA, 37–48. <https://doi.org/10.1145/2508168.2508177>
- [15] Chris Lattner and Vikram Adve. 2004. LLVM: A Compilation Framework for Lifelong Program Analysis & Transformation. In *Proceedings of the International Symposium on Code Generation and Optimization: Feedback-directed and Runtime Optimization (CGO ’04)*. IEEE Computer Society, Washington, DC, USA, 75–. <http://dl.acm.org/citation.cfm?id=977395.977673>
- [16] Tobias Lindahl and Konstantinos Sagonas. 2003. Unboxed Compilation of Floating Point Arithmetic in a Dynamically Typed Language Environment. In *Proceedings of the 14th International Conference on Implementation of Functional Languages (IFL ’02)*. Springer Berlin Heidelberg, Berlin, Heidelberg, 134–149. <http://dl.acm.org/citation.cfm?id=1756972.1756981>
- [17] Tobias Lindahl and Konstantinos Sagonas. 2004. Detecting Software Defects in Telecom Applications Through Lightweight Static Analysis: A War Story. In *Programming Languages and Systems: Second Asian Symposium, APLAS 2004, Taipei, Taiwan, November 4–6, 2004. Proceedings*, Wei-Ngan Chin (Ed.). Springer-Verlag, Berlin, Heidelberg, 91–106. https://doi.org/10.1007/978-3-540-30477-7_7
- [18] Tobias Lindahl and Konstantinos Sagonas. 2005. TypEr: A Type Annotator of Erlang Code. In *Proceedings of the 2005 ACM SIGPLAN Workshop on Erlang (ERLANG ’05)*. ACM, New York, NY, USA, 17–25. <https://doi.org/10.1145/1088361.1088366>
- [19] Tobias Lindahl and Konstantinos Sagonas. 2006. Practical Type Inference Based on Success Typings. In *Proceedings of the 8th ACM SIGPLAN International Conference on Principles and Practice of Declarative Programming (PPDP ’06)*. ACM, New York, NY, USA, 167–178. <https://doi.org/10.1145/1140335.1140356>
- [20] Daniel Luna, Mikael Pettersson, and Konstantinos Sagonas. 2004. HiPE on AMD64. In *Proceedings of the 2004 ACM SIGPLAN Workshop on Erlang*. ACM, New York, NY, USA, 38–47. <https://doi.org/10.1145/1022471.1022478>
- [21] Mikael Pettersson, Konstantinos Sagonas, and Erik Johansson. 2002. The HiPE/x86 Erlang Compiler: System Description and Performance Evaluation.

- In *Functional and Logic Programming, 6th International Symposium, FLOPS 2002, Proceedings (LNCS)*, Vol. 2441. Springer, Berlin, Heidelberg, 228–244. https://doi.org/10.1007/3-540-45788-7_14
- [22] Simon Peyton Jones and Simon Marlow. 2002. Secrets of the Glasgow Haskell Compiler Inliner. *J. Funct. Program.* 12, 5 (2002), 393–434. <https://doi.org/10.1017/S0956796802004331>
 - [23] Konstantinos Sagonas, Mikael Pettersson, Richard Carlsson, Per Gustafsson, and Tobias Lindahl. 2003. All You Wanted to Know About the HiPE Compiler (but Might Have Been Afraid to Ask). In *Proceedings of the 2003 ACM SIGPLAN Workshop on Erlang (Erlang '03)*. ACM, New York, NY, USA, 36–42. <https://doi.org/10.1145/940880.940886>
 - [24] Konstantinos Sagonas, Chris Stavrakakis, and Yiannis Tsioris. 2012. ErLLVM: an LLVM Backend for Erlang. In *Proceedings of the Eleventh ACM SIGPLAN Workshop on Erlang Workshop*. ACM, New York, NY, USA, 21–32. <https://doi.org/10.1145/2364489.2364494>
 - [25] Andre Santos. 1995. *Compilation by transformation for non-strict functional languages*. Ph.D. Dissertation. University of Glasgow, Scotland. <https://www.microsoft.com/en-us/research/publication/compilation-transformation-non-strict-functional-languages/>
 - [26] Manuel Serrano. 1997. Inline expansion: When and how?. In *Programming Languages: Implementations, Logics, and Programs: 9th International Symposium, PLILP '97 Including a Special Track on Declarative Programming Languages in Education Southampton, UK, September 3–5, 1997 Proceedings*. Springer Berlin Heidelberg, Berlin, Heidelberg, 143–157. <https://doi.org/10.1007/BFb0033842>
 - [27] Toshio Suganuma, Toshiaki Yasue, and Toshio Nakatani. 2002. An Empirical Study of Method In-lining for a Java Just-in-Time Compiler. In *Proceedings of the 2nd Java Virtual Machine Research and Technology Symposium*. USENIX Association, Berkeley, CA, USA, 91–104. <http://dl.acm.org/citation.cfm?id=648042.744889>
 - [28] Oscar Waddell and R Kent Dybvig. 1997. Fast and effective procedure inlining. In *Static Analysis*, Pascal Van Hentenryck (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 35–52.

Spine-local Type Inference

Christopher Jenkins

University of Iowa

christopher-jenkins@uiowa.edu

Aaron Stump

University of Iowa

aaron-stump@uiowa.edu

ABSTRACT

We present *spine-local* type inference, a method of partially inferring omitted type annotations for System F based on local type inference. *Local type inference* relies on bidirectional inference rules to propagate type information into and out of adjacent nodes of the AST and it restricts type-argument inference to occur only within a single node. Spine-local inference relaxes the restriction on type-argument inference, allowing it to occur only within an *application spine*, and improves upon it by using *contextual type-argument inference*. As our goal is to explore the design space of local type inference, we show that, relative to other variants, spine-local type inference better supports desirable features such as first-class curried applications and partial type applications, and it has the ability to infer types for some terms not otherwise possible. Our approach enjoys usual properties of a bidirectional system of having a specification for our inference algorithm and predictable requirements for typing annotations, and in particular maintains some of the advantages of local type inference such as a relatively simple implementation and a tendency to produce good-quality error messages when type inference fails.

CCS CONCEPTS

- Software and its engineering → Language features;

KEYWORDS

bidirectional typechecking, polymorphism, type errors

ACM Reference Format:

Christopher Jenkins and Aaron Stump. 2019. Spine-local Type Inference. In *Proceedings of International Symposium on Implementation and Application of Functional Languages (IFL'18)*. ACM, New York, NY, USA, 12 pages.

1 INTRODUCTION

Local type inference[18] is a simple yet effective partial technique for inferring types for programs. In contrast to complete methods of type inference such as the Damas-Milner system[3] which can type programs without any annotations by restricting the language of types, *partial* methods require the programmer to provide some type annotations and, in exchange, are suitable for use in programming languages with rich type features such as impredicativity and subtyping[15, 18], dependent types[25], and higher-rank types[16], where complete type inference may be undecidable.

Local type inference is also contrasted with *global* inference methods (usually based on unification) which are able to infer more missing annotations by solving typing constraints generated from the entire program. Though more powerful, global inference methods can also be more difficult for programmers to use when type inference fails, as they can generate type errors whose root cause

is distant from the location the error is reported[11]. Local type inference address this issue by only propagating typing information between adjacent nodes of the abstract syntax tree (AST), allowing programmers to reason *locally* about type errors. It achieves this by using two main techniques: *bidirectional type inference rules* and *local type-argument inference*.

The first of these techniques, bidirectional type inference, is not unique to local type inference ([5, 16, 23, 26] are just a few examples), and uses two main judgment forms, often called *synthesis* and *checking* mode. When a term t synthesizes type T , we view this typing information as coming up and out of t and as available for use in typing nearby terms; when t checks against type T (called in this paper the *contextual type*), this information is being pushed down and in to t and is provided by nearby terms.

The second of these techniques, local type-argument inference, finds the missing types arguments in polymorphic function applications by using only the type information available at an application node of the AST. For a simple example, consider the expression $\text{id } z$ where id has type $\forall X. X \rightarrow X$ and z has type \mathbb{N} . Here we can perform *synthetic* type-argument inference by synthesizing the type of z and comparing this to the type of id to infer we should instantiate type variable X with \mathbb{N} .

Local type inference has a number of desirable properties. Using just these two techniques it can in practice infer a good deal of type annotations, and often those that are required to remain are predictable and coincide with programmers' expectations that they serve as useful and machine-checked documentation[8, 18]. Without further instrumentation, local type inference already tends to report type errors close to where further annotations are required; more recently, it has been used as the basis for developing autonomous type-driven debugging and error explanations[19]. The type inference algorithms of [15, 18] admit a specification for their behavior, helping programmers understand why certain types were inferred without requiring they know every detail of the type-checker's implementation. Add to this its relative simplicity and robustness when extended to richer type systems and it seems unsurprising that it has been a popular choice for type inference in programming languages.

Unfortunately, local type inference can fail even when it seems like there should be enough typing information available locally. Consider trying to check that the expression $\text{pair } (\lambda x. x) z$ has type $((\mathbb{N} \rightarrow \mathbb{N}) \times \mathbb{N})$, assuming pair has type $\forall X. \forall Y. X \rightarrow Y \rightarrow (X \times Y)$. The inference systems presented in [15, 18] will fail here because the argument $\lambda x. x$ does not synthesize a type. The techniques proposed in the literature of local type inference for dealing with cases similar to this include classifying and avoiding such "hard-to-synthesize" terms[8] and utilizing the partial type information provided by polymorphic functions[15]; the former was dismissed as unsatisfactory by the same authors that introduced it and the latter is of no help in this situation, since the type of

pair tells us nothing about the expected type of $\lambda x. x$. What we need in this case is *contextual* type-argument inference, utilizing the information available from the contextual type of the result of the application to know argument $\lambda x. x$ is expected to have type $\mathbb{N} \rightarrow \mathbb{N}$.

Additionally, languages using local type inference usually use fully-uncurried applications in order to maximize the notion of “locality” for type-argument inference, improving its effectiveness. The programmer can still use curried applications if desired, but “they are second-class in this respect.”[18]. It is also usual for type arguments to be given in an “all or nothing” fashion in such languages, meaning that even if only one cannot be inferred, all must be provided. We believe that currying and partial type applications are useful idioms for functional programming and wish to support them as first-class language features.

1.1 Contributions

In this paper, we explore the design space of local type inference in the setting of System F[6, 7] by developing *spine-local* type inference, an approach that both expands the locality of type-argument inference to an *application spine* and augments its effectiveness by using the *contextual* type of the spine. In doing so, we

- show that, compared to other variants of local type inference, we better support first-class currying and partial type applications, and can infer types for some “hard-to-synthesize” terms that these systems would not type;
- provide a specification for contextual type-argument inference with respect to which we show our algorithm is sound and complete
- give a weak completeness theorem for our type system with respect to fully unannotated System F terms, indicating the conditions under which the programmer can expect type inference succeeds and where additional annotations are required when it fails.

Spine-local type inference is implemented in Cedille[20], a functional programming language with higher-order and impredicative polymorphism, dependent types, and dependent intersections. Though the setting for this paper is much simpler, we are optimistic that spine-local type inference will serve as a good foundation for type inference in Cedille that makes using its rich type features more convenient for programmers.

The rest of this paper is organized as follows: in Section 2 we cover the syntax and some useful terminology for our setting; in Section 3 we present the type inference rules constituting a specification for contextual type-argument inference, consider its annotation requirements, and illustrate its use, limitations, and the type errors it presents to users; in Section 4 we show the prototype-matching algorithm implementing contextual type-argument inference; and in Section 5 we discuss how this work compares to other approaches to type inference.

2 INTERNAL AND EXTERNAL LANGUAGE

Type inference can be viewed as a relation between an *internal* language of terms, where all needed typing information is present, and an *external* language, in which programmers work directly and

where some of this information can be omitted for their convenience. Under this view, type inference for the external language not only associates a term with some type but also with some *elaborated* term in the internal language in which all missing type information has been restored. In this section, we present the syntax for our internal and external languages as well as introduce some terminology that will be used throughout the rest of this paper.

2.1 Syntax

We take as our internal language explicitly typed System F (see [7]); we review its syntax below:

Types	$S, T, U, V ::= X, Y, Z \mid S \rightarrow T \mid \forall X. T$
Contexts	$\Gamma ::= \cdot \mid \Gamma, X \mid \Gamma, x:T$
Internal Terms	$e, p ::= x \mid \lambda x:T. e \mid \Lambda X. e \mid e e' \mid e[T]$

Types consist of type variables, arrow types, and type quantification, and typing contexts consist of the empty context, type variables (also called the context’s *declared type variables*), and term variables associated with their types. The internal language of terms consists of variables, λ -abstractions with annotations on bound variables, Λ -abstractions for polymorphic terms, and term and type applications. Our notational convention in this paper is that term meta-variable e indicates an elaborated term for which all type arguments are known, and p indicates a *partially* elaborated term where some elaborated type arguments are type meta-variables (discussed in Section 3).

The external language contains the same terms as the internal language as well as bare λ -abstractions – that is, λ -abstractions missing an annotation on their bound variable:

External Terms	$t, t' ::= x \mid \lambda x:T. t \mid \lambda x. t \mid \Lambda X. t \mid t t' \mid t[T]$
-----------------------	---

Types and contexts are the same as for the internal language.

2.2 Terminology

In both the internal and external languages, we say that the *applicand* of a term or type application is the term in the function position. A *head* is either a variable or abstraction (term or type), and an *application spine*[2] (or just *spine*) is a view of an application as consisting of a head (called the *spine head*) followed by a sequence of (term and type) arguments. The *maximal application* of a sub-expression is the spine in which it occurs as an applicand, or just the sub-expression itself if it does not. For example, spine $x[S] y z$ is the maximal application of itself and its applicand sub-expressions x , $x[S]$, and $x[S] y$, with x as head of the spine. Predicate $App(t)$ indicates term t is some term or type application (in either language) and we define it formally as $(\exists t_1, t_2. t = t_1 t_2) \vee (\exists t', S. t = t'[S])$.

Turning to definitions for types and contexts, function $DTV(\Gamma)$ calculates the set of *declared type variables* of context Γ and is defined recursively by the following set of equations:

$$\begin{aligned} DTV(\cdot) &= \emptyset \\ DTV(\Gamma, X) &= DTV(\Gamma) \cup \{X\} \\ DTV(\Gamma, x:T) &= DTV(\Gamma) \end{aligned}$$

Predicate $WF(\Gamma, T)$ indicates that type T is *well-formed* under Γ – that is, all free type variables of T occur as declared type variables in Γ (formally $FV(T) \subseteq DTV(\Gamma)$).

3 TYPE INFERENCE SPECIFICATION

The typing rules for our internal language are standard for explicitly typed System F and are omitted (see Ch. 23 of [17] for a thorough discussion of these rules). We write $\Gamma \vdash e : T$ to indicate that under context Γ internal term e has type T . For type inference in the external language, Figure 1 shows judgment \vdash_δ which consists mostly of standard (except for *AppSyn* and *AppChk*) bidirectional inference rules with elaboration to the internal language, and Figure 2 shows the specification for contextual type-argument inference. Judgment \vdash^P in Figure 2b handles traversing the spine and judgment \vdash^I in Figure 2c types its term applications and performs type-argument inference (both synthetic and contextual). Figure 2a gives a “shim” judgment \vdash^I which bridges the bidirectional rules with the specification for rhetorical purposes (discussed below). While these rules are not fully syntax-directed, they are subject-directed, meaning that for each judgment the shape of the term we are typing (i.e. the *subject* of typing) uniquely determines which rule applies.

Bidirectional Rules. We now consider more closely each judgment form and its rules starting with \vdash_δ , the point of entry for type inference. The two modes for type inference, checking and synthesizing, are indicated resp. by \vdash_\Downarrow (suggesting pushing a type down and into a term) and \vdash_\Uparrow (suggesting pulling a type up and out of a term). Following the notational convention of Peyton Jones et al.[16] we abbreviate two inference rules that differ only in their direction to one by writing \vdash_δ , where δ is a parameter ranging over $\{\Downarrow, \Uparrow\}$. We read judgment $\Gamma \vdash_\Uparrow t : T \rightsquigarrow e$ as: “under context Γ , term t synthesizes type T and elaborates to e ,” and a similar reading for checking mode applies for \vdash_\Downarrow . When the direction does not matter, we will simply say that we can *infer* t has type T .

Rule *Var* is standard. Rule *Abs* says we can infer missing type annotation T on a λ -abstraction when we have a contextual arrow type $T \rightarrow S$. Rules *AAbs* and *TAbs* say that Λ - and annotated λ -abstractions can have their types either checked or synthesized. *TApp* says that a type application $t[S]$ has its type inferred in either mode when the applicand t synthesizes a quantified type. The reason for this asymmetry between the modes of the conclusion and the premise is that even when in checking mode, it is not clear how to work backwards from type $[S/X]T$ to $\forall X. T$.

AppSyn and *AppChk* are invoked on maximal applications and are the first non-standard rules. To understand how these rules work, we must 1) explain the “shim” judgment \vdash^I serving as the interface for spine-local type-argument inference and 2) define meta-language function *MV*. Read $\Gamma; T_? \vdash^I t t' : T \rightsquigarrow (p, \sigma)$ as “under context Γ and with (optional) contextual type $T_?$, we partially infer application $t t'$ has type T with elaboration p and solution σ ,” where σ is a substitution mapping a subset of the meta-variables (i.e. the originally omitted type arguments) in p to contextually-inferred type arguments.

In rule *AppSyn*, $?$ is provided to \vdash^I indicating no contextual type is available. We constrain σ to be the identity substitution (notation σ_{id}) and insist the elaborated term has no unsolved meta-variables, matching our intuition that all type arguments must be inferred synthetically. In rule *AppChk*, we provide the contextual type to \vdash^I and check (implicitly) that it equals σT and (explicitly) that all remaining meta-variables in p are solved by σ , then elaborate σp (the replacement of each meta-variable in p with its mapping in

σ). Shared by both is the second premise of the (anonymous) rule introducing \vdash^I that σ solves precisely the meta-variables of the partially inferred type T for application $t t'$.

Meta-variables. What do we mean by the “meta-variables” of partial elaborations and types? When t is a term application with some type arguments omitted in its spine, then (under context Γ) its partial elaboration p from type-argument inference fills in each missing type argument with either a type (well-formed under Γ) or with a *meta-variable* (a type variable not declared in Γ) depending on whether it was inferred through synthesizing the types of the term arguments of t . For example, if $t = \text{pair } (\lambda x. x) z$ and we wanted to check that it has type $T = \langle (\mathbb{N} \rightarrow \mathbb{N}) \times \mathbb{N} \rangle$ under a typing context Γ associating *pair* with type $\forall X. \forall Y. X \rightarrow Y \rightarrow \langle X \times Y \rangle$ and z with type \mathbb{N} , then we could derive

$$\Gamma; T \vdash^I t : \langle X \times \mathbb{N} \rangle \rightsquigarrow (\text{pair}[X][\mathbb{N}] (\lambda x:\mathbb{N}. x) z, [\mathbb{N} \rightarrow \mathbb{N}/X])$$

(assuming some base type \mathbb{N} , some family of base types $\langle S \times T \rangle$ for all types S and T , and assuming X is not declared in Γ .) Looking at the partial elaboration of t , we would see that type argument X was inferred from its contextual type $\langle (\mathbb{N} \rightarrow \mathbb{N}) \times \mathbb{N} \rangle$ and that Y was inferred from the synthesized types of the arguments to *pair* (in this case z).

Meta-variables never occur in a judgment formed by \vdash_δ , only in the judgments of Figure 2. In particular, these rules enforce that meta-variables in a partial elaboration p can occur *only* as type arguments in its spine, not within its head or term arguments. This restriction guarantees *spine-local* type-argument inference and helps to narrow the programmer’s focus when debugging type errors. Furthermore, meta-variables correspond to omitted type arguments *injectively*, significantly simplifying the kind of reasoning needed for debugging type errors. We make this precise by defining meta-language function *MV*($\Gamma, _$) which yields the set of meta-variables occurring in its second argument with respect to the context Γ . *MV* is overloaded to take both types and elaborated terms for its second argument: for types we define $\text{MV}(\Gamma, T) = \text{FV}(T) - \text{DTV}(\Gamma)$, the set of free variables in T less the declared type variables of Γ ; for terms, $\text{MV}(\Gamma, p)$ is defined recursively by the following equations:

$$\text{MV}(\Gamma, p) = \emptyset \quad \text{when } \neg \text{App}(p)$$

$$\text{MV}(\Gamma, p[X]) = \text{MV}(\Gamma, p) \cup \{X\} \quad \text{when } X \notin \text{DTV}(\Gamma)$$

$$\text{MV}(\Gamma, p[S]) = \text{MV}(\Gamma, p) \quad \text{when } \text{WF}(\Gamma, S)$$

$$\text{MV}(\Gamma, p e) = \text{MV}(\Gamma, p)$$

Using our running example where the subject t is $\text{pair } (\lambda x. x) z$ we can now show how the meta-variable checks are used in rules *AppSyn* and *AppChk*. We have for our partially elaborated term that $\text{MV}(\Gamma, \text{pair}[X][\mathbb{N}] (\lambda x:\mathbb{N}. x) z) = \{X\}$ and also for our type that $\text{MV}(\Gamma, \langle X \times \mathbb{N} \rangle) = \{X\}$. If we have a derivation of the judgment above formed by \vdash^I we can then derive with rule *AppChk*

$$\Gamma \vdash_\Downarrow t : \langle (\mathbb{N} \rightarrow \mathbb{N}) \times \mathbb{N} \rangle \rightsquigarrow \text{pair}[\mathbb{N} \rightarrow \mathbb{N}][\mathbb{N}] (\lambda x:\mathbb{N}. x) z$$

because substitution $[\mathbb{N} \rightarrow \mathbb{N}/X]$ solves the remaining meta-variable X in the elaborated term and type, and when utilized on the partially inferred type $\langle X \times \mathbb{N} \rangle$ yields the contextual type for the term. However, we would not be able to derive with rule *AppSyn*

$$\Gamma \vdash_\Uparrow t : \langle (\mathbb{N} \rightarrow \mathbb{N}) \times \mathbb{N} \rangle \rightsquigarrow \text{pair}[\mathbb{N} \rightarrow \mathbb{N}][\mathbb{N}] (\lambda x:\mathbb{N}. x) z$$

$$\begin{array}{c}
\boxed{\Gamma \vdash_{\delta} t : T \rightsquigarrow e} \\
\frac{}{\Gamma \vdash_{\delta} x : \Gamma(x) \rightsquigarrow x} Var \\
\frac{\Gamma, x : T \vdash_{\delta} t : S \rightsquigarrow e}{\Gamma \vdash_{\delta} \lambda x : T. t : T \rightarrow S \rightsquigarrow \lambda x : T. e} AAbs \\
\frac{\Gamma, X \vdash_{\delta} t : T \rightsquigarrow e}{\Gamma \vdash_{\delta} \Lambda X. t : \forall X. T \rightsquigarrow \Lambda X. e} TAbs \\
\frac{\Gamma, x : T \vdash_{\Downarrow} t : S \rightsquigarrow e}{\Gamma \vdash_{\Downarrow} \lambda x. t : T \rightarrow S \rightsquigarrow \lambda x : T. e} Abs \\
\frac{\Gamma, x : T \vdash_{\delta} t : S \rightsquigarrow e}{\Gamma \vdash_{\delta} \lambda x : T. t : T \rightarrow S \rightsquigarrow \lambda x : T. e} AAbs \\
\frac{\Gamma, X \vdash_{\delta} t : T \rightsquigarrow e}{\Gamma \vdash_{\delta} \Lambda X. t : \forall X. T \rightsquigarrow \Lambda X. e} TAbs \\
\frac{\Gamma \vdash_{\Downarrow} t[S] : [S/X]T \rightsquigarrow e[S]}{\Gamma \vdash_{\delta} t[S] : [S/X]T \rightsquigarrow e[S]} TApp \\
\frac{\Gamma; ? \vdash^I t t' : T \rightsquigarrow (e, \sigma_{id}) \quad MV(\Gamma, e) = \emptyset}{\Gamma \vdash_{\uparrow\uparrow} t t' : T \rightsquigarrow e} AppSyn \\
\frac{\Gamma; ? \vdash^I t t' : T \rightsquigarrow (p, \sigma) \quad MV(\Gamma, p) = dom(\sigma)}{\Gamma \vdash_{\Downarrow} t t' : \sigma T \rightsquigarrow \sigma p} AppChk
\end{array}$$

Figure 1: Bidirectional inference rules with elaboration

$$\begin{array}{c}
(a) \text{ Shim (specification)} \\
\frac{\Gamma \vdash^P t t' : T \rightsquigarrow (p, \sigma) \quad MV(\Gamma, T) = dom(\sigma) \quad T_? \in \{?, \sigma T\}}{\Gamma; T_? \vdash^I t t' : T \rightsquigarrow (p, \sigma)} \\
(b) \boxed{\Gamma \vdash^P t : T \rightsquigarrow (p, \sigma)} \\
\frac{\neg App(t) \quad \Gamma \vdash_{\uparrow\uparrow} t : T \rightsquigarrow e}{\Gamma \vdash^P t : T \rightsquigarrow (e, \sigma_{id})} PHead \quad \frac{\Gamma \vdash^P t : \forall X. T \rightsquigarrow (p, \sigma)}{\Gamma \vdash^P t[S] : [S/X]T \rightsquigarrow (p[S], \sigma)} PTApp \quad \frac{\Gamma \vdash^P t : T \rightsquigarrow (p, \sigma) \quad \Gamma \vdash^{\cdot} (p : T, \sigma) \cdot t' : T' \rightsquigarrow (p', \sigma')}{\Gamma \vdash^P t t' : T' \rightsquigarrow (p', \sigma')} PApp \\
(c) \boxed{\Gamma \vdash^{\cdot} (p : T, \sigma) \cdot t' : T' \rightsquigarrow (p', \sigma')} \\
\frac{\sigma'' \in \{\sigma, [S/X] \circ \sigma\} \quad WF(\Gamma, S) \quad \Gamma \vdash^{\cdot} (p[X] : T, \sigma'') \cdot t' : T' \rightsquigarrow (p', \sigma')}{\Gamma \vdash^{\cdot} (p : \forall X. T, \sigma) \cdot t' : T' \rightsquigarrow (p', \sigma')} PForall \quad \frac{MV(\Gamma, \sigma S) = \emptyset \quad \Gamma \vdash_{\Downarrow} t' : \sigma S \rightsquigarrow e'}{\Gamma \vdash^{\cdot} (p : S \rightarrow T, \sigma) \cdot t' : T \rightsquigarrow (([U/Y] p) e', \sigma')} PChk \\
\frac{MV(\Gamma, \sigma S) = \overline{Y} \neq \emptyset \quad \Gamma \vdash_{\uparrow\uparrow} t' : \overline{[U/Y]} \sigma S \rightsquigarrow e'}{\Gamma \vdash^{\cdot} (p : S \rightarrow T, \sigma) \cdot t' : \overline{[U/Y]} T \rightsquigarrow (([U/Y] p) e', \sigma')} PSyn
\end{array}$$

Figure 2: Specification for contextual type-argument inference

since we do not have σ_{id} as our solution and we have meta-variable X remaining in our partial elaboration and type. Together, the checks in $AppSyn$ and $AppChk$ ensure that meta-variables are never passed up and out of a maximal application during type inference.

Specification Rules. Judgment \vdash^I serves as an interface to spine-local type-argument inference. In Figure 2a it is defined in terms of the specification for contextual type-argument inference given by judgments \vdash^P and \vdash^{\cdot} ; we call it a “shim” judgment because in Figure 4a we give for it an alternative definition using the algorithmic rules in which the condition $MV(\Gamma, T) = dom(\sigma)$ is not needed. Its purpose, then, is to cleanly delineate what we consider specification and implementation for our inference system.

Though the details of maintaining spine-locality and performing synthetic type-argument inference permeate the inference rules for \vdash^P and \vdash^{\cdot} , these rules form a specification in that they fully abstract away the details of contextual type-argument inference, describing how solutions are used but omitting how they are generated. Spine-locality in particular contributes to our specification’s perceived complexity – what would be one or two rules in a fully-uncurried language with all-or-nothing type argument applications is broken down in our system into multiple inference rules to support currying and partial type applications.

Judgment \vdash^P contains three rules and serves to dig through a spine until it reaches its head, then work back up the spine typing its term and type applications. The reading for it is the same as for \vdash^I , less the optional contextual type. Rule $PHead$ types the spine head t by deferring to $\vdash_{\uparrow\uparrow}$; our partial solution is σ_{id} since no meta-variables are present in a judgment formed by $\vdash_{\uparrow\uparrow}$. $PTApp$ is similar to $TApp$ except it additionally propagates solution σ . Rule $PApp$ is used for term applications: first it partially synthesizes a type T for the applicand and then it uses judgment \vdash^{\cdot} to ensure that a function of type T can be applied to the argument t' .

Judgment \vdash^{\cdot} performs synthetic and contextual type-argument inference and ensures that term applications with omitted type arguments are well-typed. We read $\Gamma \vdash^{\cdot} (p : T, \sigma) \cdot t' : T' \rightsquigarrow (p', \sigma')$ as “under context Γ , elaborated applicand p of partial type T together with solution σ can be applied to term t' ; the application has type T' and elaborates p' with solution σ' .”

Contextual type-argument inference happens in rule $PForall$, which says that when the applicand has type $\forall X. T$ we can choose to guess any well-formed S for our contextual type argument by picking $\sigma'' = [S/X] \circ \sigma$ (indicating σ'' contains all the mappings present in σ and an additional mapping S for X), or choose to attempt to synthesize it later from an argument by picking $\sigma'' = \sigma$. *The details of which S to guess, or whether we should guess at all, are*

not present in this specifical rule. In both cases, we elaborate the applicand to $p[X]$ of type T and check that it can be applied to t' – we do this even when we guess S for X to maintain the invariant that all partial elaborations p and solutions σ we generate satisfy $\text{dom}(\sigma) \subseteq MV(\Gamma, p)$, needed when checking in the (specifical) rule for \vdash^I that these guessed solutions are ultimately justified by the contextual type (if any) of our maximal application.

We illustrate the use of $PForall$ with an example: if the subject of (i.e. input to) judgment \vdash^I is

$$(\text{pair} : \forall X. \forall Y. X \rightarrow Y \rightarrow \langle X \times Y \rangle, \sigma_{id}) \cdot (\lambda x. x)$$

then after two uses of rule $PForall$ where we guess $\mathbb{N} \rightarrow \mathbb{N}$ for X and decline to guess for Y the subject would be:

$$(\text{pair}[X][Y] : X \rightarrow Y \rightarrow \langle X \times Y \rangle, [\mathbb{N} \rightarrow \mathbb{N}/X]) \cdot (\lambda x. x)$$

After working through omitted type arguments, \vdash^I requires that we eventually reveal some arrow type $S \rightarrow T$ to type a term application. When it does we have two cases, handled resp. by $PChk$ and $PSyn$: either the domain type S of applicand p together with solution σ provide enough information to fully know the expected type for argument t' (i.e. $MV(\Gamma, \sigma p) = \emptyset$), or else they do not and we have some non-empty set of unsolved meta-variables \bar{Y} in S corresponding to type arguments we must synthesize. Having full knowledge, in $PChk$ we check t' has type σS ; otherwise, in $PSyn$ we try to solve meta-variables \bar{Y} by synthesizing a type for t' and checking it is instantiation $[U/Y]$ (vectorized notation for the simultaneous substitution of types \bar{U} for \bar{Y}) of σS . Once done, we conclude with result type $[\bar{U}/\bar{Y}] T$ and elaboration $([\bar{U}/\bar{Y}] p) e$ for the application, as the meta-variables \bar{Y} of p corresponding to omitted type arguments have now been fully solved by type-argument synthesis. Together, $PChk$ and $PSyn$ prevent meta-variables from being passed down to term argument t' , as we require that it either check against or synthesize a well-formed type.

We illustrate the use of rule $PSyn$ with an example: suppose that under context Γ the subject of judgment \vdash^I is

$$(\text{pair}[X][Y] (\lambda x : \mathbb{N}. x) : Y \rightarrow \langle X \times Y \rangle, [\mathbb{N} \rightarrow \mathbb{N}/X]) \cdot z$$

and furthermore that $\Gamma \vdash_I z : \mathbb{N}$. Then we have instantiation $[\mathbb{N}/Y]$ from synthetic type-argument inference and use it to produce for the application the result type $[\mathbb{N}/Y] \langle X \times Y \rangle = \langle X \times \mathbb{N} \rangle$ and the elaboration $\text{pair}[X][\mathbb{N}] (\lambda x : \mathbb{N}. x) z$. Note that synthesized type arguments are used *eagerly*, meaning that the typing information synthesized from earlier arguments can in some cases be used to infer the types of later arguments in *checking* mode (see Section 3.2). This is reminiscent of *greedy* type-argument inference for type systems with subtyping[1, 4], which is known to cause unintuitive type inference failures due to sub-optimal type arguments (i.e. less general wrt to the subtyping relation) being inferred. As System F lacks subtyping, this problem does not affect our type inference system and we can happily utilize synthesized type arguments eagerly (see Section 5).

3.1 Soundness, Weak Completeness, and Annotation Requirements

The inference rules in Figure 2 for our external language are *sound* with respect to the typing rules for our internal language (i.e. explicitly typed System F), meaning that elaborations of typeable external terms are typeable at the same type¹:

THEOREM 1. (*Soundness of \vdash_S*):
If $\Gamma \vdash_S t : T \rightsquigarrow e$ then $\Gamma \vdash e : T$.

Our inference rules also enjoy a trivial form of completeness that serves as a sanity-check with respect to the internal language: since any fully annotated internal language term e is in the external language, we expect that e should be typeable using the typing rules for external terms:

THEOREM 2. (*Trivial Completeness of \vdash_S*):
If $\Gamma \vdash e : T$ then $\Gamma \vdash_S e : T \rightsquigarrow e$

A more interesting form of completeness comes from asking *which* external terms can be typed – after all, this is precisely what a programmer needs to know when trying to debug type inference failures! Since our external language contains terms without any annotations and our type language is impredicative System F, we know from [24] that type inference is in general undecidable. Therefore, to state a completeness theorem for type inference we must first place some restrictions on the set of external terms that can be the subject of typing.

We start by defining what it means for t to be a *partial erasure* of internal term e . The grammar given in Section 2 for the external language does not fully express where we hope our inference rules will restore missing type information. Specifically, the rules in Figures 1 and 2 will try to infer annotations on bare λ -abstractions and only try to infer missing type arguments that occur in the applicand of a term application. For example, given (well-typed) internal term $x[S_1][S_2] y[T]$ and external term $x y$, our inference rules will try to infer the missing type arguments S_1 and S_2 but *will not* try to infer the missing T .

A more artificial restriction on partial erasures is that the sequence of type arguments occurring between two terms in an application can only be erased in a right-to-left fashion. For example, given internal term $x[S_1][S_2] y[T_1][T_2] z$, the external term $x y[T_1] z$ is a valid erasure (S_1 and S_2 are erased between x and y , and between y and z rightmost T_2 is erased), but term $x[S_2] y[T_2] z$ is not. This restriction helps preserve soundness of the external type inference rules by ensuring that every explicit type argument preserved in an erasure of an internal term e instantiates the same type variable it did in e ; it is artificial because we could instead have introduced notation for “explicitly erased” type arguments in the external language, such as $x[_] y$, to indicate the first type argument has been erased, or allow type arguments to be given *by name* such as $x[_X_2=S_2] y$, but chose not to do so to simplify the presentation of our inference rules and language.

The above restrictions for partial erasures are made precise by the functions $\lfloor _ \rfloor$ and $\lfloor _ \rfloor_a$ which map an internal term e to sets of

¹A complete list of proofs for non-trivial theorems can be found in the proof appendix at http://homepage.cs.uiowa.edu/~cwjnkns/Papers/JSL18_Spine-local/proof-appendix.pdf

partial erasures $\lfloor e \rfloor$. They are defined mutually recursively below:

$$\lfloor \lambda x:T. e \rfloor = \{ \lambda x:T. t \mid t \in \lfloor e \rfloor \} \cup \{ \lambda x. t \mid t \in \lfloor e \rfloor \}$$

$$\lfloor \Lambda X. e \rfloor = \{ \Lambda X. t \mid t \in \lfloor e \rfloor \}$$

$$\lfloor e e' \rfloor = \{ t t' \mid t \in \lfloor e \rfloor_a \wedge t' \in \lfloor e' \rfloor \}$$

$$\lfloor e[S] \rfloor = \{ t[S] \mid t \in \lfloor e \rfloor \}$$

$$\lfloor e[S] \rfloor_a = \{ t \mid t \in \lfloor e \rfloor_a \} \cup \{ t[S] \mid t \in \lfloor e \rfloor \}$$

$$\lfloor e \rfloor_a = \lfloor e \rfloor \text{ otherwise}$$

We are now ready to state a weak completeness theorem for typing terms in the external language which over-approximates the annotations required for type inference to succeed (we write $\forall \bar{X}. T$ to mean some number of type quantifications over type T)

THEOREM 3. (*Weak completeness of \vdash_{\uparrow} :*)

Let e be a term of the internal language and t be an external language term such that $t \in \lfloor e \rfloor$. If $\Gamma \vdash e : T$ then $\Gamma \vdash_{\uparrow} t : T \rightsquigarrow e$ when the following conditions hold for each sub-expression e' of e , corresponding sub-expression t' of t , and corresponding sub-derivation $\Gamma' \vdash e' : T'$ of $\Gamma \vdash e : T$:

- (1) If $e' = \lambda x : S. e''$ for some S and e'' , then $t' = \lambda x : S. t''$ for some t''
- (2) If e' is a maximal term application in e and if $\Gamma' \vdash^P t' : T'' \rightsquigarrow (p, \sigma_{id})$ for some T'' and p , then $MV(\Gamma, p) = \emptyset$.
- (3) If e' is a term application and $t' = t_1 t_2$ for some t_1 and t_2 , and if $\Gamma' \vdash^P t_1 : T'' \rightsquigarrow (p, \sigma_{id})$ for some T'' and p , then $T'' = \forall \bar{X}. S_1 \rightarrow S_2$ for some S_1 and S_2 .
- (4) If e' is a type application and $t' = t''[S]$ for some t'' and S , and $\Gamma' \vdash^P t'' : T'' \rightsquigarrow (p, \sigma_{id})$ for some T'' and p , then $T'' = \forall X. S'$ for some S' .

Theorem 3 only considers synthetic type-argument inference, and in practice condition (1) is too conservative thanks to contextual type-argument inference. Though a little heavyweight, our weak completeness theorem can be translated into a reasonable guide for where type annotations are required when type synthesis fails. Conditions (3) and (4) suggest that when the applicand of a term or type application already partially synthesizes some type, the programmer should give enough type arguments to at least reveal it has the appropriate shape (resp. a type arrow or quantification). (2) indicates that type variables that do not occur somewhere corresponding to a term argument of an application should be instantiated explicitly, as there is no way for synthetic type-argument inference to do so. For example, in the expression $f z$ if f has type $\forall X. \forall Y. Y \rightarrow X$ there is no way to instantiate X from synthesizing argument z . Finally, condition (1) we suggest as the programmer's last resort: if the above advice does not help it is because some λ -abstractions need annotations.

Note that in conditions (2), (3), and (4) we are not circularly assuming type synthesis for sub-expressions of partial erasure t succeeds in order to show that it succeeds for t , only that if a certain sub-expression can be typed then we can make some assumptions about the shape of its type or elaboration. Conditions (3) and (4) in particular are a direct consequence of a design choice we made for our algorithm to maintain injectivity of meta-variables to

omitted type arguments. As an alternative, we could instead refine meta-variables when we know something about the shape of their instantiation. For example, if we encountered a term application whose applicand has a meta-variable type X , we know it must have some arrow type and could refine X to $X_1 \rightarrow X_2$, where X_1 and X_2 are fresh meta-variables. However, doing so means type errors may now require non-trivial reasoning from users to determine why some meta-variables were introduced in the first place.

Still, we find it somewhat inelegant that our characterization of annotation requirements for type inference is not fully independent of the inference system itself. For programmers using these guidelines, this implies that there must be some way to interactively query the type-checker for different sub-expressions of a program during debugging. Fortunately, many programming languages offer just such a feature in the form of a REPL, meaning that in practice this is not too onerous a requirement to make.

Theorem 3 only states when an external term will synthesize its type, but what about when a term can be *checked* against a type? It is clear from the typing rules in Figure 1 that some terms that fail to synthesize a type may still be successfully checked against a type. Besides typing bare λ -abstractions (which can only have their type checked), checking mode can also reduce the annotation burden implied by condition (2) of Theorem 3: consider again the example $f z$ where f has type $\forall X. \forall Y. Y \rightarrow X$. If instead of attempting type synthesis we were to check that it has some type T then we would not need to provide an explicit type argument to instantiate X .

From these observations and our next result, we have that checking mode of our type inference system can infer the types of strictly more terms than can synthesizing mode – whenever a term synthesizes a type, it can be checked against the same type.

THEOREM 4. (*Checking extends synthesizing*):

$$\text{If } \Gamma \vdash_{\uparrow} t : T \rightsquigarrow e \text{ then } \Gamma \vdash_{\downarrow} t : T \rightsquigarrow e$$

3.2 Examples

Successful Type Inference. We conclude this section with some example programs for which the type inference system in Figures 1 and 2 will and will not be able to type. We start with the motivating example from the introduction of checking that the expression $\text{pair } (\lambda x. x) z$ has type $\langle (\mathbb{N} \rightarrow \mathbb{N}) \times \mathbb{N} \rangle$, which is not possible in other variants of local type inference. For convenience, we assume the existence of a base type \mathbb{N} and a family of base types $\langle S \times T \rangle$ for all types S and T . These assumptions are admissible as we could define these types using lambda encodings. A full derivation for typing this program is given in Figure 3, including the following abbreviations:

$$\begin{aligned} I_X X Y &= X \rightarrow Y \rightarrow \langle X \times Y \rangle \\ \Gamma &= \text{pair} : \forall X. \forall Y. I_X X Y, z : \mathbb{N} \\ \sigma &= [\mathbb{N} \rightarrow \mathbb{N}/X] \\ p &= \text{pair}[X][\mathbb{N}] (\lambda x : \mathbb{N}. x) z \end{aligned}$$

To type this application $\text{pair } (\lambda x. x) z$ we first dig through the spine, reach the head pair , and synthesize type $\forall X. \forall Y. I_X X Y$. No meta-variables are generated by judgment \vdash_{\uparrow} and thus there can be no meta-variable solutions, so we generate solution σ_{id} .

Next we type the first application, $\text{pair } (\lambda x. x)$, shown in sub-derivation \mathcal{D}_1 . In the first invocation of rule $PForall$ we guess

$$\begin{array}{c}
 \frac{\Gamma \vdash_{\uparrow} \text{pair} : \forall X. \forall Y. I_X X Y \rightsquigarrow \text{pair}}{\Gamma \vdash^P \text{pair} : \forall X. \forall Y. I_X X Y \rightsquigarrow (\text{pair}, \sigma_{id})} \text{Var} \\
 \frac{\Gamma \vdash^P \text{pair} : \forall X. \forall Y. I_X X Y \rightsquigarrow (\text{pair}, \sigma_{id}) \quad \mathcal{D}_1}{\Gamma \vdash^P \text{pair} (\lambda x. x) : Y \rightarrow \langle X \times \mathbb{N} \rangle \rightsquigarrow (\text{pair}[X][\mathbb{N}] (\lambda x: \mathbb{N}. x), \sigma)} \text{PHead} \\
 \frac{\Gamma \vdash^P \text{pair} (\lambda x. x) : Y \rightarrow \langle X \times \mathbb{N} \rangle \rightsquigarrow (\text{pair}[X][\mathbb{N}] (\lambda x: \mathbb{N}. x), \sigma) \quad \mathcal{D}_2}{\Gamma \vdash^P \text{pair} (\lambda x. x) z : \langle X \times \mathbb{N} \rangle \rightsquigarrow (p, \sigma)} \text{PApp} \\
 \frac{\Gamma \vdash^I \text{pair} (\lambda x. x) z : \langle X \times \mathbb{N} \rangle \rightsquigarrow (p, \sigma) \quad MV(\Gamma, X \times \mathbb{N}) = \text{dom}(\sigma)}{\Gamma \vdash_{\downarrow} \text{pair} (\lambda x. x) z : \langle (\mathbb{N} \rightarrow \mathbb{N}) \times \mathbb{N} \rangle \rightsquigarrow \text{pair}[\mathbb{N} \rightarrow \mathbb{N}][\mathbb{N}] (\lambda x: \mathbb{N}. x) z} \text{MV}(\Gamma, X \times \mathbb{N}) = \text{dom}(\sigma) \\
 \frac{\Gamma \vdash^I \text{pair} (\lambda x. x) z : \langle X \times \mathbb{N} \rangle \rightsquigarrow (p, \sigma) \quad MV(\Gamma, p) = \text{dom}(\sigma)}{\Gamma \vdash_{\downarrow} \text{pair} (\lambda x. x) z : \langle (\mathbb{N} \rightarrow \mathbb{N}) \times \mathbb{N} \rangle \rightsquigarrow \text{pair}[\mathbb{N} \rightarrow \mathbb{N}][\mathbb{N}] (\lambda x: \mathbb{N}. x) z} \text{AppChk}
 \end{array}$$

$$\mathcal{D}_1 = \Gamma \vdash^* (\text{pair}: \forall X. \forall Y. I_X X Y, \sigma_{id}) \cdot (\lambda x. x) : Y \rightarrow \langle X \times Y \rangle \rightsquigarrow (\text{pair}[X][Y] (\lambda x: \mathbb{N}. x), \sigma) \quad \text{PChk}$$

$$\frac{\Gamma \vdash^* (\text{pair}[X][Y] : I_X X Y, \sigma) \cdot (\lambda x. x) : Y \rightarrow \langle X \times Y \rangle \rightsquigarrow (\text{pair}[X][Y] (\lambda x: \mathbb{N}. x), \sigma) \quad MV(\Gamma, \sigma X) = \emptyset}{\Gamma \vdash^* (\text{pair}[X]: \forall Y. I_X X Y, \sigma) \cdot (\lambda x. x) : Y \rightarrow \langle X \times Y \rangle \rightsquigarrow (\text{pair}[X][Y] (\lambda x: \mathbb{N}. x), \sigma)} \text{PForall}$$

$$\frac{\Gamma \vdash^* (\text{pair}: \forall X. \forall Y. I_X X Y, \sigma_{id}) \cdot (\lambda x. x) : Y \rightarrow \langle X \times Y \rangle \rightsquigarrow (\text{pair}[X][Y] (\lambda x: \mathbb{N}. x), \sigma)}{\mathcal{D}_2} \text{PForall}$$

$$\mathcal{D}_2 = \Gamma \vdash^* (\text{pair}[X][Y] (\lambda x: \mathbb{N}. x) : Y \rightarrow \langle X \times Y \rangle, \sigma) \cdot z : \langle X \times \mathbb{N} \rangle \rightsquigarrow (\text{pair}[X][\mathbb{N}] (\lambda x: \mathbb{N}. x) z, \sigma) \quad \text{PSyn}$$

where	$I_X X Y = X \rightarrow Y \rightarrow \langle X \times Y \rangle$	$\Gamma = \text{pair}: \forall X. \forall Y. I_X X Y, z: \mathbb{N}$
	$\sigma = [\mathbb{N} \rightarrow \mathbb{N}/X]$	$p = \text{pair}[X][\mathbb{N}] (\lambda x: \mathbb{N}. x) z$

Figure 3: Example typing derivation with the specification rules

solution σ for X , and in the second invocation we decline to guess an instantiation for Y (in this example we could have also guessed \mathbb{N} for Y as this information is also available from the contextual type, but choose not to in order to demonstrate the use of all three rules of \vdash^*). Then using rule $PChk$ we check argument $\lambda x. x$ against $\sigma X = \mathbb{N} \rightarrow \mathbb{N}$. This is the point at which the local type inference systems of [15, 18] will fail: as a bare λ -abstraction this argument will not synthesize a type, and the expected type X as provided by the applicand pair alone does not tell us what the missing type annotation should be. However, by using the information provided by the contextual type of the entire application we know it must have type $\mathbb{N} \rightarrow \mathbb{N}$. The resulting partial type of the application is $Y \rightarrow \langle X \times Y \rangle$, and we propagate solution σ to the rest of the derivation. Note that we elaborate the argument $\lambda x. x$ of this application to $\lambda x: \mathbb{N}. x$ – we never pass down meta-variables to term arguments, keeping type-argument inference local to the spine.

In sub-derivation \mathcal{D}_2 we type $(\text{pair} (\lambda x. x)) z$ (parentheses added) where our applicand has partial type $Y \rightarrow \langle X \times Y \rangle$. We find that we have unsolved meta-variable Y as the expected type for z , so we use rule $PSyn$ and synthesize the type \mathbb{N} for z . Using solution $[\mathbb{N}/Y]$, we produce $\langle X \times \mathbb{N} \rangle$ for the resulting type of the application and elaborate the application to $\text{pair}[X][\mathbb{N}] (\lambda x: \mathbb{N}. x) z$, wherein type argument Y is replaced by \mathbb{N} in the original elaborated applicand $\text{pair}[X][Y] (\lambda x: \mathbb{N}. x)$.

Finally, in rule $AppChk$ we confirm that the meta-variables remaining in our partial type synthesis of the application is precisely those for which we knew the solutions contextually. For this example, the only remaining meta-variable in both the partially synthesized type and elaboration is X , which is also the only mapping in σ , so type inference succeeds. We use σ to replace all occurrences

of X with $\mathbb{N} \rightarrow \mathbb{N}$ in the type and elaboration and conclude that term $\text{pair} (\lambda x. x) z$ can be checked against type $\langle (\mathbb{N} \rightarrow \mathbb{N}) \times \mathbb{N} \rangle$.

The next example shows how eager use of *synthetic* type-argument inference can type some terms not possible in other variants of local type inference. Consider checking that the expression $\text{rapp } x \lambda y. y$ has type \mathbb{N} , where rapp has type $\forall X. \forall Y. X \rightarrow (Y \rightarrow Y) \rightarrow Y$ and x has type \mathbb{N} . From the contextual type we know that Y should be instantiated to \mathbb{N} , and when we reach application $\text{rapp } x$, we learn that X should be instantiated to \mathbb{N} from the synthesized type of x . Together, this gives us enough information to know that argument $\lambda y. y$ should have type $\mathbb{N} \rightarrow \mathbb{N}$. Such eager instantiation is neither novel, nor is it necessarily desirable when extended to richer type languages or more powerful methods of inference (see Section 5), but in our setting it is a useful technique that we can use to infer types for expressions like the one above.

Type Inference Failures. To see where type inference can fail, we again use $\text{pair} (\lambda x. x) z$ but now ask that it *synthesize* its type. Rule $AppSyn$ insists that we make no guesses for meta-variables (as there is no contextual type for the application that they could have come from), so we would need to synthesize a type for argument $\lambda x. x$ – but our rules do not permit this! In this case the user can expect an error message like the following:

```
expected type: ?X
error: We are not in checking mode, so bound
variable x must be annotated
```

where $?X$ indicates an unsolved meta-variable corresponding to type variable X in the type of pair . The situation above corresponds to condition (1) of Theorem 3: in general, if there is not enough information from the type of an applicand and the contextual type of the application spine in which it occurs to fully know the expected

types of arguments that are λ -abstractions, then such arguments require explicit type annotations.

We next look at an example corresponding to condition (2) of Theorem 3, namely that the type variables of a polymorphic function that do not correspond to term arguments in an application should be instantiated explicitly. Here we will assume a family of base types $S + T$ for every type S and T , a variable `right` of type $\forall X. \forall Y. Y \rightarrow (X + Y)$, and a variable `z` of type \mathbb{N} . In trying to synthesize a type for the application `right z` the user can expect an error message like:

```
synthesized type: (?X + N)
error: This maximal application has unsolved
      meta-variables
```

indicating that type variable X requires an explicit type argument be provided. Fortunately for the programmer, and unlike the local type inference systems of [15, 18], our system supports partial explicit type application, meaning that X can be instantiated without also explicitly (and redundantly) instantiating Y^2 . On the other hand, local type inference systems for System F_≤[15, 18] can succeed in typing `right z` *without* additional type arguments, as they will instantiate X to the minimal type (with respect to their subtyping relation) `Bot`. Partial type application, then, is more useful for our setting of System F where picking an instantiation in this situation would be rather arbitrary.

A more subtle point of failure for our algorithm corresponds to conditions (3) and (4) of Theorem 3. Even when the head and all arguments of an application spine can synthesize their types, the programmer may still be required to provide some additional type arguments. Consider the expression `bot z`, where `bot : \forall X. X` and `z : \mathbb{N}. Even with some contextual type for this expression, type inference still fails because the rules in Figure 2c require that the type of the applicand of a term application reveals some arrow, which $\forall X. X$ does not. The programmer would be met with the following error message:`

```
applicand type: ?X
error: The type of an applicand in a term
      application must reveal an arrow
```

prompting the user to provide an explicit type argument for X . To make expression `bot z` typeable, the programmer could write `bot[\mathbb{N} \rightarrow \mathbb{N}] z`, or even `bot[\forall Y. Y \rightarrow Y] z – our inference rules are able to solve meta-variables introduced by explicit (and even synthetically-inferred) type arguments, as long as there is enough information to reveal a quantifier or arrow in the type of a term or type applicand.`

For our last type error example, we consider the situation where the programmer has written an ill-typed program. Local type inference enjoys the property that type errors can be understood *locally*, without any “spooky action” from a distant part of the program. In particular, with local type inference we would like to avoid error messages like the following:

```
synthesized type: B → B
expected type: ?X := N → N
error: type mismatch
```

²As discussed in Section 3.1, our system is artificially sensitive to the order of quantifiers, and the corresponding `left : \forall X. \forall Y. X \rightarrow (X + Y)` does not enjoy the same benefit, but could with named type arguments.

From this error message alone the programmer has no indication of why the expected type is $\mathbb{N} \rightarrow \mathbb{N}$! In our type inference system we have expanded the distance information travels by allowing it to flow from the contextual type of an application to its arguments – so can programmers expect now to see such error messages?. As an example, the error message above could have been generated when checking the expression pair $(\lambda x : \mathbb{B}. x) z$ against type $\langle (\mathbb{N} \rightarrow \mathbb{N}) \times \mathbb{N} \rangle$, specifically when inferring the type of the first argument. Fortunately, our notion of locality is still quite small and we can easily demystify the reason our system expected a different type:

```
synthesized type: B → B
expected type: ?X := N → N
contextual match: <?X × ?Y> := <(\mathbb{N} → \mathbb{N}) × \mathbb{N}>
```

where `contextual match` tells the programmer to compare the partially synthesized and contextual return types of the application spine to determine why X was instantiated to $\mathbb{N} \rightarrow \mathbb{N}$. A similar field, `synthetic match`, could tell the programmer that the type of an earlier argument informs the expected type of current one.

4 ALGORITHMIC INFERENCE RULES

The type inference system presented in Section 3 do not constitute an algorithm. Though the rules forming judgment \vdash indicate *where* and *how* we use contextually-inferred type arguments, they do not specify *what* their instantiations are or even *whether* this information is available to use, and it is not obvious how to work backwards from the second premise in Figure 2a to develop an algorithm.

Figure 4 shows the algorithmic rules implementing contextual type-argument inference. The full algorithm for spine-local type inference, then, consists of the rules in Figure 1 with the shim judgment \vdash^I as defined in Figure 4a. At the heart of the implementation is our prototype matching algorithm; to understand the details of how we implement contextual type-argument inference, we must first discuss this algorithm and the two new syntactic categories it introduces, prototypes and decorated types.

4.1 Prototype Matching

Figure 4d lists the rules for the prototype matching algorithm. We read the judgment $\bar{X} \Vdash^{\perp\!\!\perp} T := P \Rightarrow (\sigma, W)$ as: “solving for meta-variables \bar{X} , we match type T to prototype P and generate solution σ and decorated type W ,” and we maintain the invariant that $\text{dom}(\sigma) \subseteq \bar{X}$. Meta-variables can only occur in T , thus these are *matching* (not unification) rules. The grammar for prototypes and decorated types is given below:

Prototypes $P ::= ? \mid T \mid ? \rightarrow P$	Decorated Types $W ::= T \mid S \rightarrow W \mid \forall X = X. W \mid \forall X = S. W$
	$\mid (X, ? \rightarrow P)$

Prototypes carry the contextual type of the maximal application of a spine. In the base case they are either the uninformative $?$ (as in `AppSyn`), indicating no contextual type, or they are informative of type T (as in `AppChk`). In this way, prototypes generalize the syntactic category $T_?$ we introduced earlier for optional contextual types. We use the last prototype former $? \rightarrow$ as we work our way down an application spine to track the expected arity of its head. For example, if we wished to check that the expression `id suc x`

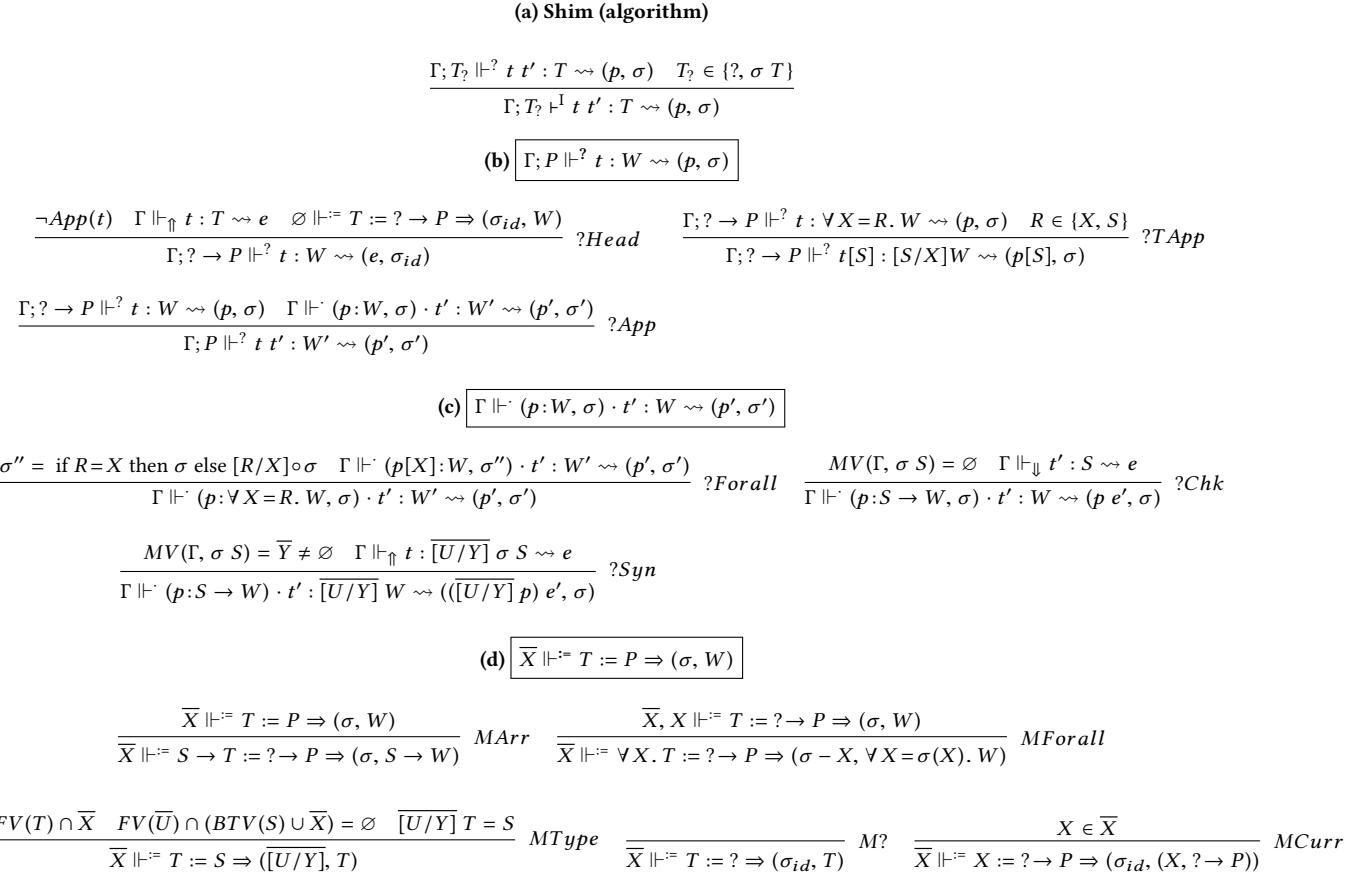


Figure 4: Algorithm for contextual type argument inference

has type \mathbb{N} , then when we reached the head id using the rules in Figure 4b we would generate for it prototype $? \rightarrow ? \rightarrow \mathbb{N}$

Decorated types consist of types (also called *plain-decorated* types), an arrow with a regular type as the domain (as prototypes only inform us of the result type of a maximal application, not of the types of arguments), quantified types whose bound variable X may be decorated with the type to which we expect to instantiate it, and “stuck” decorations. On quantifiers, decoration $X = X$ indicates that P did not inform us of an instantiation for X – we sometimes abbreviate the two cases as $\forall X = R. W$, where $R \in \{X, S\}$ and $S \neq X$.

To explain the role of stuck decorations, consider again $id \ succ$. Assuming id has type $\forall X. X \rightarrow X$, matching this with prototype $? \rightarrow ? \rightarrow \mathbb{N}$ generates decorated type $\forall X = X. X \rightarrow (X, ? \rightarrow \mathbb{N})$, meaning that we only know that X will be instantiated to some type that matches $? \rightarrow \mathbb{N}$. Stuck decorations occur when the expected arity of a spine head (as tracked by a given prototype) is greater than the arity of the head’s synthesized type and are the mechanism by which we propagate a contextual type to a head that is “over-applied” – a not-uncommon occurrence in languages with curried applications!

Turning to the prototype matching algorithm in Figure 4d, rule *MArr* says that we match an arrow type and prototype when we can

match their codomains. Rule *MType* says that when the prototype is some type S we must find an instantiation $\overline{[U/Y]}$ (where $\overline{Y} \subseteq \overline{X}$) such that $\overline{[U/Y]} T = S$, and rule *M?* says that any type matches with $?$ with no solutions generated (thus we call $?$ “uninformative”). In rule *MForall* we match a quantified type with a prototype by adding bound variable X to our meta-variables and matching the body T to the same prototype; the substitution in the conclusion, $\sigma - X$, is the solution generated from this match less its mapping for X , which is placed in the decoration $X = \sigma(X)$. For example, matching $\forall X. \forall Y. X \rightarrow Y \rightarrow X$ with prototype $? \rightarrow ? \rightarrow \mathbb{N}$ generates decorated type $\forall X = \mathbb{N}. \forall Y = Y. X \rightarrow Y \rightarrow X$. Finally, rule *MCurr* applies when there is incomplete information (in the form of $? \rightarrow P$) on how to instantiate a meta-variable; we generate a stuck decoration with identity solution σ_{id} .

We conclude by showing that our prototype matching rules are *functional*; when \overline{X} , T , and P are considered as inputs then there is at most one output pair (σ, W) :

THEOREM 5. (*Function-ness of $\Vdash^=$*):

Given \overline{X} , T , and P , if $\overline{X} \Vdash^= T := P \Rightarrow (\sigma, W)$ and $\overline{X} \Vdash^= T := P \Rightarrow (\sigma', W')$, then $\sigma = \sigma'$ and $W = W'$

It should also be clear that these rules are *syntax-directed* (i.e. that the rules are non-overlapping and that the inputs to premises and the outputs to the conclusions are uniquely determined), meaning they can be straightforwardly translated to an algorithm.

4.2 Decorated Type Inference

We now discuss the rules in Figures 4b and 4c which implement contextual type-argument inference (as specified by Figures 2b and 2c) by using the prototype matching algorithm. We begin by giving a reading for judgments $\Vdash^? - \text{read } \Gamma; P \Vdash^? t : W \rightsquigarrow (p, \sigma)$ as: “under context Γ and with prototype P , t synthesizes decorated type W and elaborates p with solution σ ,” where σ again represents the contextually-inferred type arguments.

In rule *AppSyn* we required that the solution generated by \vdash^I in its premise is σ_{id} ; in *AppChk* we (implicitly) required that the contextual type is equal to σT ; and now with the algorithmic definition for \vdash^I we appear to be requiring in both that the decorated type generated by $\Vdash^?$ is a plain-decorated type T . With the algorithmic rules, these are not requirements but *guarantees* that the specification makes of the algorithm:

LEMMA 1. *Let $\text{arr}_P(P)$ be the number of prototype arrows prefixing P and $\text{arr}_W(W)$ be the number of decorated arrows prefixing W . If $\Gamma; P \Vdash^? t : W \rightsquigarrow (p, \sigma)$ then $\text{arr}_W(W) \leq \text{arr}_P(P)$*

THEOREM 6. (*Soundness of $\Vdash^?$ wrt $\Vdash^{:=}$*):

If $\Gamma; P \Vdash^? t : W \rightsquigarrow (p, \sigma)$ then $MV(\Gamma, p) \Vdash^{:=} T := P \Rightarrow (\sigma, W)$

Assuming prototype inference succeeds, when we specialize P in Theorem 6 to $?T$ we have immediately by rule *M?* that $\sigma = \sigma_{id}$; when we specialize it to some contextual type T' for an application, then by the premise of *MType* we have $\sigma T = T'$. Theorem 1 and 6 together tell us that we generate plain-decorated types in both cases, as in particular we cannot have leading (decorated) arrows or stuck decorations with prototypes $?T$.

Next we discuss the rules forming judgment $\Vdash^?$ in Figure 4b, constituting the algorithmic version of the rules in Figure 2b. In rule *?Head*, after synthesizing a type T for the application head we match this type against expected prototype $? \rightarrow P$ (we are guaranteed the prototype has this shape since a maximal term application begins all derivations of $\Vdash^?$). No meta-variables occur in T initially – as we perform prototype matching these will be generated by rule *MForall* from quantified type variables in T and their solutions will be left as quantifier decorations in the resulting decorated type W . We are justified in requiring that matching T to $? \rightarrow P$ generates empty solution σ_{id} since we have in general that the meta-variables solved by our prototype matching judgment are a subset of the meta-variables it was asked to solve:

LEMMA 2. *If $\bar{X} \Vdash^{:=} T := P \Rightarrow (\sigma, W)$ then $\text{dom}(\sigma) \subseteq \bar{X}$*

In *?TApp*, we can infer the type of a type application $t[S]$ when t synthesizes a decorated type $\forall X = R. W$ and R is either an uninformative decoration X or is precisely S (that is, the programmer provided explicitly the type argument the algorithm contextually inferred). We synthesize $[S/X]W$ for the type application, where we extend type substitution to decorated types by the following

recursive partial function:

$$\begin{aligned} \sigma S \rightarrow W &= (\sigma S) \rightarrow (\sigma W) \\ \sigma \forall X = R. W &= \forall X = R. \sigma W \\ \sigma (X, ? \rightarrow P) &= W && \text{if } \emptyset \Vdash^{:=} \sigma(X) := ? \rightarrow P \Rightarrow (\sigma_{id}, W) \end{aligned}$$

This definition is straightforward except for the case dealing with stuck decorations. Here, σ (representing type arguments given explicitly or inferred synthetically) may provide information on how to instantiate X and this must match our current (though incomplete) information $? \rightarrow P$ about this instantiation. For example, if we have decorated type $W = X \rightarrow (X, ? \rightarrow \mathbb{N})$, then $[\mathbb{N} \rightarrow \mathbb{N}/X]W$ would require we match $\mathbb{N} \rightarrow \mathbb{N}$ with $? \rightarrow \mathbb{N}$ and matching would generate (plain) decorated type $(\mathbb{N} \rightarrow \mathbb{N}) \rightarrow \mathbb{N} \rightarrow \mathbb{N}$

The definition of substitution on decorated types is partial since prototype matching may fail (consider if we used substitution $[\mathbb{N}/X]$ in the above example instead). When a decorated type substitution σW appears in the conclusion of our algorithmic rules, such as in *?TApp* or *?Syn*, we are implicitly assuming an additional premise that the result is defined.

The last rule for judgment $\Vdash^?$ is *?App*, and like *PApp* it benefits from a reading for judgment $\Vdash^?$ occurring in its premise. We read $\Gamma \Vdash^? (p : W, \sigma) \cdot t' : W' \rightsquigarrow (p', W')$ as: “under Γ , elaborated applicand p of decorated type W together with solution σ can be applied to t' ; the application has decorated type W' and elaborates p' with solution σ' .” Thus, *?App* says that to synthesize a decorated type for a term application $t' t$ we synthesize the decorated type of the applicand t and ensure that the resulting elaboration p , along with its decorated type and solution, can be applied to t' .

We now turn to the rules for the last judgment $\Vdash^?$ of our algorithm. Rule *?Forall* clarifies the non-deterministic guessing done by the specifical rule *PForall*: the contextually-inferred type arguments we build during contextual type-argument inference are just the accumulation of quantified type decorations. The solution σ'' we provide to the second premise of *?Forall* contains mapping $[R/X]$ if R is an informative decoration, and as we did in rule *PForall* we provide elaborated term $p[X]$ to track the contextually-inferred type arguments separately from those synthetically inferred.

Rule *?Chk* works similarly to *PChk*: when the only meta-variables in the domain S of our decorated type are solved by σ , we can check that argument t' has type σS . In rule *?Syn* we have some meta-variables \bar{Y} in S not solved by σ – we synthesize a type for the argument, ensure that it is some instantiation $[\bar{U}/\bar{Y}]$ of σS , and use this instantiation on the meta-variables in p as well as the decorated codomain type W , potentially unlocking some stuck decoration to reveal more arrows or decorated type quantifications.

We conclude this section by noting that the specifical and algorithmic type inference system are equivalent, in the sense that they type precisely the same set of terms:

THEOREM 7. (*Soundness of \Vdash_δ wrt \vdash_δ*):

If $\Gamma \Vdash_\delta t : T \rightsquigarrow e$ then $\Gamma \vdash_\delta t : T \rightsquigarrow e$

THEOREM 8. (*Completeness of \Vdash_δ wrt \vdash_δ*):

If $\Gamma \vdash_\delta t : T \rightsquigarrow e$ then $\Gamma \Vdash_\delta t : T \rightsquigarrow e$

(where \Vdash_δ indicates \vdash^I is defined as in Figure 4a)

Taken together, Theorems 7 and 8 justify our claim that the rules of Figure 2 constitute a specification for contextual type-argument inference – it is not necessary that the programmer know the notably more complex details of prototype matching or type decoration to understand how some type arguments are inferred contextually. Indeed, the judgment \vdash provides more flexibility in reasoning about type inference than does \Vdash , as in rule *PForall* we may freely decline to guess a contextual type argument even when this would be justified and instead try to learn it synthetically. In contrast, algorithmic rule *?Forall* requires that we use any informative quantifier decoration. We use this flexibility when giving guidelines for the required annotations in Section 3.1 for typing external terms, as the required conditions for typeability in Theorem 3 would be further complicated if we could not restrict ourselves to using only synthetic type-argument inference.

5 DISCUSSION & RELATED WORK

5.1 Local Type Inference and System F_{\leq}

Local Type Inference. Our work is most influenced by the seminal paper by Pierce and Turner[18] on local type inference that describes its broad approach, including the two techniques of bidirectional typing rules and local type-argument inference and the design-space restriction that polymorphic function applications be fully-uncurried to maximize the benefit of these techniques. In their system, either all term arguments to polymorphic functions must be synthesized or else all type arguments must be given – no compromise is available when only a few type arguments suffice to type an application, be they provided explicitly or inferred contextually. Our primary motivation in this work was addressing these issues – improving support for first-class currying and partial type application, and using the contextual type of an application for type-argument inference – while maintaining some of the desirable properties of local type inference and staying in the spirit of their approach.

Colored Local Type Inference. Odersky, Zenger, and Zenger[15] extend the type system of Pierce and Turner by allowing *partial* type information to be propagated downwards when inferring types for term arguments. Their insight was to internalize the two modes of bidirectional type inference to the structure of types themselves, allowing different parts of a type to be synthetic or contextual. In contrast, we use an “all or nothing” approach to type propagation: when we encountered a term argument for which we have incomplete information, we require that it fully synthesize its type. On the other hand, their system uses only the typing information provided by the application head, whereas we combine this with the contextual type of an application, allowing us to type some expressions their system cannot. The upshot of the difference in these systems is that spine-local type inference uses *more* contextual information and colored local type inference uses contextual information *more cleverly*.

The syntax for prototypes in our algorithm was directly inspired by the prototypes used in the algorithmic inference rules for [15]. Our use of prototypes complements theirs; ours propagates the partial type information provided by contextual type of an application

spine to its head, whereas theirs propagates the partial type information provided by an application head to its arguments. In future work, we hope to combine these two notions of prototype to propagate *partially* the type information coming from the application’s contextual type *and* head to its arguments.

Subtyping. Local type inference is usually studied in the setting of System F_{\leq} which combines impredicative parametric polymorphism and subtyping. The reason for this is two-fold: first, a partial type inference technique is needed as complete type inference for F_{\leq} is undecidable[21]; second, global type inference systems fail to infer principal types in F_{\leq} [9, 12, 14], whereas local type inference is able to promise that it infers the “locally best”[18] type arguments (i.e. the type arguments minimizing the result type of the application, relative to the subtyping relation). The setting for our algorithm is System F, so the reader may ask whether our developments can be extended gracefully to handle subtyping. We believe the answer is yes, though with some modification on how synthetic type arguments are used.

In rule *PSyn* in Figure 2c, meta-variables \bar{Y} are instantiated to types \bar{U} immediately. In the presence of subtyping this would make our rules *greedy*[1, 4] and we would not be able to guarantee synthetic type-argument inference produced locally best types, possibly causing type inference to fail later in the application spine. To illustrate this, consider the expression $\text{rapp } x \text{ neg}$, assuming $\text{rapp} : \forall X. \forall Y. X \rightarrow (X \rightarrow Y) \rightarrow Y, x : \mathbb{N}, \text{neg} : \mathbb{Z} \rightarrow \mathbb{Z}$, and some subtyping relation \leq where $\mathbb{N} \leq \mathbb{Z}$. Greed causes us to instantiate X with \mathbb{N} , but in order to type the expression we would need to instantiate it to \mathbb{Z} instead!

To correct this, we could instead collect these constraints and solve them only when the function is fully applied to its arguments (i.e., when we reach a stuck decoration). This mirrors the requirement in [18] that constraints are solved at fully uncurried applications, maintaining currying but losing a syntactically-obvious location for synthetic type-argument inference.

We would also need to justify our use of contextual type-argument inference for checking the types of term arguments. Happily, this does not appear to be an intractable problem like greed: unlike in synthesis mode, checking mode for applications in [18] does not require that the synthesized type arguments minimize the result type of the application, so there is greater freedom in choosing the instantiations for contextually-inferred type arguments. Hosoya and Pierce note in [8] that the optimal instantiations for these type arguments are ones that “maximize the expected type corresponding to the [argument],” as the type that the programmer meant for the argument (if type correct) will be a subtype of this. Though the informal approach they proposed (and later dismissed) for inferring the types of hard-to-synthesize terms differs from ours in the use of a “slightly ad-hoc” analysis of arguments, it anticipated contextual type-argument inference and suggests the way forward for using it in the presence of subtyping.

5.2 Bidirectional Type Inference and System F

Predicative Polymorphism. The popularity of bidirectional type inference extends well beyond local inference methods. Dunfield and Krishnaswami[5] introduced a relatively simple and elegant type inference system for predicative System F using a dedicated application judgment that instantiates type arguments at term applications. Their application judgment was the direct inspiration for our own, though there are some significant differences between the two. First, our rules distinguish between checking the argument of an application with a fully known expected type and synthesizing its argument when incomplete information is available to keep meta-variables *spine-local*, whereas in their approach meta-variables and typing constraints are passed downwards to check term arguments. Our system also contains the additional judgment form \vdash^P that theirs does not, again to contain meta-variables within an application spine.

Approaches to type inference for System F (impredicative and predicative alike) often make use of some form of subsumption rule to decrease the required type annotations in terms. A popular basis for such rules is the “more polymorphic than” subtyping relation introduced by Odersky and Läufer in [13] which stratifies polymorphic and monomorphic types and is able to perform deeply nested monomorphic type instantiation. This line of work also includes [5] above as well as work by Peyton Jones et. al. [16], both of which are able to infer arbitrary-rank types in the setting of predicative System F. In contrast our type inference algorithm supports more powerful impredicative polymorphism at the cost of significant increase in required type annotations.

Impredicative Polymorphism. The “more-polymorphic-than” subtyping relation for impredicative System F is undecidable[21], so type inference systems wishing to use a subsumption rule in this setting must make some compromises. With **ML^F** [10] Le Botlan and Rémy develop a type language with bounded type quantification and an inference system using type *instantiation* (a covariant restriction of subtyping). *Boxy type inference*[23] by Vytiniotis et al. uses an idea similar to [15] of propagating partial type information (though with a very different implementation) to allow inference for polymorphic types only in checking mode; its later development in **FPH**[22] both simplifies the specification for type inference and extends boxy types to synthesis mode to allow “boxy monotypes” to be inferred for polymorphic functions. These inference systems add additional constructs (resp. bounded quantifications and boxy types) to System F types in their specification, whereas we reserve our new constructs (decorated types and prototypes) for the algorithmic rules only. Our use of first-order matching when typing applications of and arguments to polymorphic functions can be viewed as a crude form of subtyping via shallow type instantiation – it is significantly easier for programmers to understand but at the same time significantly less powerful than the subtyping used in the type inference systems above.

ACKNOWLEDGMENTS

We thank Larry Diehl, Anthony Cantor, and Ernesto Copello for their feedback on earlier versions of this paper which helped us clarify some points of terminology and improve the readability of the more technical sections of the paper. We gratefully acknowledge

NSF support under award 1524519, and DoD support under award FA9550-16-1-0082 (MURI program).

REFERENCES

- [1] Luca Cardelli. 1997. An implementation of F<. <http://citeseervx.ist.psu.edu/viewdoc/summary?doi=10.1.1.48.6158>
- [2] Iliaio Cervesato and Frank Pfenning. 2003. A Linear Spine Calculus. *J. Log. Comput.* 13 (2003), 639–688.
- [3] Luis Damas and Robin Milner. 1982. Principal Type-schemes for Functional Programs. In *Proceedings of the 9th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '82)*. ACM, New York, NY, USA, 207–212. <https://doi.org/10.1145/582153.582176>
- [4] Joshua Dunfield. 2009. Greedy Bidirectional Polymorphism. In *Proceedings of the 2009 ACM SIGPLAN Workshop on ML (ML '09)*. ACM, New York, NY, USA, 15–26. <https://doi.org/10.1145/1596627.1596631>
- [5] Joshua Dunfield and Neelakantan R. Krishnaswami. 2013. Complete and Easy Bidirectional Typechecking for Higher-rank Polymorphism. *SIGPLAN Not.* 48, 9 (Sept. 2013), 429–442. <https://doi.org/10.1145/2544174.2500582>
- [6] Jean-Yves Girard. 1986. The system F of variable types, fifteen years later. *Theoretical Computer Science* 45 (1986), 159 – 192. [https://doi.org/10.1016/0304-3975\(86\)90044-7](https://doi.org/10.1016/0304-3975(86)90044-7)
- [7] Jean-Yves Girard, Paul Taylor, and Yves Lafont. 1989. *Proofs and Types*. Cambridge University Press, New York, NY, USA.
- [8] Haruo Hosoya and Benjamin Pierce. 1999. How Good is Local Type Inference? (07 1999).
- [9] A. J. Kennedy. 1996. *Type Inference and Equational Theories*. Technical Report LIX, ECOLE POLYTECHNIQUE, 91128 PALAISEAU CEDEX.
- [10] Didier Le Botlan and Didier Rémy. 2003. MLF: Raising ML to the Power of System F. *SIGPLAN Not.* 38, 9 (Aug. 2003), 27–38. <https://doi.org/10.1145/944746.944709>
- [11] Bruce McAdam. 2002. Trends in Functional Programming. Intellect Books, Exeter, UK, UK, Chapter How to Repair Type Errors Automatically, 87–98. <http://dl.acm.org/citation.cfm?id=644403.644412>
- [12] Martin Odersky. 2002. Inferred Type Instantiation for GJ. Note sent to the types mailing list.
- [13] Martin Odersky and Konstantin Läufer. 1996. Putting Type Annotations to Work. In *Proceedings of the 23rd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '96)*. ACM, New York, NY, USA, 54–67. <https://doi.org/10.1145/237721.237729>
- [14] Martin Odersky, Martin Sulzmann, and Martin Wehr. 1999. Type Inference with Constrained Types. *Theor. Pract. Object Syst.* 5, 1 (Jan. 1999), 35–55. [https://doi.org/10.1002/\(SICI\)1096-9942\(199901/03\)5:1<35::AID-TAPO4>3.0.CO;2-4](https://doi.org/10.1002/(SICI)1096-9942(199901/03)5:1<35::AID-TAPO4>3.0.CO;2-4)
- [15] Martin Odersky, Christoph Zenger, and Matthias Zenger. 2001. Colored Local Type Inference. *SIGPLAN Not.* 36, 3 (Jan. 2001), 41–53. <https://doi.org/10.1145/373243.360207>
- [16] Simon Peyton Jones, Dimitrios Vytiniotis, Stephanie Weirich, and Mark Shields. 2005. Practical type inference for arbitrary-rank types. 17 (January 2005), 1–82. <https://www.microsoft.com/en-us/research/publication/practical-type-inference-for-arbitrary-rank-types/>
- [17] Benjamin C. Pierce. 2002. *Types and Programming Languages* (1st ed.). The MIT Press.
- [18] Benjamin C. Pierce and David N. Turner. 2000. Local Type Inference. *ACM Trans. Program. Lang. Syst.* 22, 1 (Jan. 2000), 1–44. <https://doi.org/10.1145/345099.345100>
- [19] Hubert Plociniczak. 2016. *Decrypting Local Type Inference*. Ph.D. Dissertation. École polytechnique fédérale de Lausanne. <http://dx.doi.org/10.5075/epfl-thesis-6741>
- [20] Aaron Stump. 2017. The calculus of dependent lambda eliminations. *J. Funct. Program.* 27 (2017), e14. <https://doi.org/10.1017/S0956796817000053>
- [21] J. Tiuryn and P. Urzyczyn. 1996. The subtyping problem for second-order types is undecidable. In *Proceedings 11th Annual IEEE Symposium on Logic in Computer Science*. 74–85. <https://doi.org/10.1109/LICS.1996.561306>
- [22] Dimitrios Vytiniotis, Simon Peyton Jones, and Stephanie Weirich. 2008. FPH: First-class polymorphism for Haskell. In Proceedings of the 13th ACM SIGPLAN International Conference on Functional Programming. <https://www.microsoft.com/en-us/research/publication/fph-first-class-polymorphism-for-haskell/>
- [23] Dimitrios Vytiniotis, Stephanie Weirich, and Simon Peyton Jones. 2006. Boxy Types: Inference for Higher-rank Types and Impredicativity. *SIGPLAN Not.* 41, 9 (Sept. 2006), 251–262. <https://doi.org/10.1145/1160074.1159838>
- [24] J. B. Wells. 1998. Typability and Type Checking in System F Are Equivalent and Undecidable. *ANNALS OF PURE AND APPLIED LOGIC* 98 (1998), 111–156.
- [25] Hongwei Xi and Frank Pfenning. 1999. Dependent Types in Practical Programming. In *Proceedings of the 26th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '99)*. ACM, New York, NY, USA, 214–227. <https://doi.org/10.1145/292540.292560>
- [26] Ningning Xie and Bruno C. d. S. Oliveira. 2018. Let Arguments Go First. In *Programming Languages and Systems*, Amal Ahmed (Ed.). Springer International Publishing, Cham, 272–299.

Verifiably Lazy

Verified Compilation of Call-by-Need

George Stelle

CCS-7

Los Alamos National Laboratory

Los Alamos, New Mexico, United States

Computer Science

University of New Mexico

Albuquerque, New Mexico, United States

stelleg@lanl.gov

Darko Stefanovic

Computer Science

University of New Mexico

Albuquerque, New Mexico, United States

darko@unm.edu

ABSTRACT

Call-by-need semantics are widely used in the programming language Haskell. Unfortunately, unlike their call-by-value counterparts, there are no verified compilers for call-by-need. In this paper we present the first verified compiler for call-by-need semantics. We use recent work on a simple call-by-need abstract machine as a way of reducing the formalization burden. We discuss some of the difficulties in verifying call-by-need, and show how we overcome them.

ACM Reference Format:

George Stelle and Darko Stefanovic. 2019. Verifiably Lazy: Verified Compilation of Call-by-Need. In *Proceedings of International Symposium on Implementation and Application of Functional Languages (IFL'18)*. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Non-strict languages, such as Haskell, rely heavily on call-by-need semantics to ensure efficient execution. Without the memoization of results provided by call-by-need, Haskell would be prohibitively inefficient, often exponentially slower than its call-by-value counterparts. Haskell is one of the purest languages in wide use for software development, making reasoning about correctness in Haskell easier than most languages. It is for this reason that we would like to have a compiler that gives formal guarantees about preservation of call-by-need semantics. We wish to ensure that any reasoning we do about our non-strict functional programs is preserved through compilation, and that their execution will not be prohibitively inefficient.

Unfortunately, one of the challenges for formalization of non-strict compilers is that the semantics of call-by-need abstract machines tend to be complex, incorporating complex optimizations into the semantics, requiring preprocessing of terms, and closures of variable sizes [6, 13]. Recently we developed a particularly simple

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

IFL'18, September 2018, Lowell, MA, USA

© 2019 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM... \$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

abstract machine for call-by-need, the $\mathcal{C}\mathcal{E}$ machine [15]. In addition to being exceedingly simple to implement and reason about, the machine shows performance often comparable to state-of-the-art.

Verified compilers provide powerful guarantees about the code they generate and its relation to the corresponding source code [4, 8, 10]. In particular, for higher order functional languages, they can ensure that the non-trivial task of compiling lambda calculus and its derivatives to machine code is implemented correctly, preserving source semantics. The amortized return on investment for verified compilers is high: any reasoning about any program which is compiled with a verified compiler is provably preserved.

Existing verified compilers have focused on call-by-value semantics [4, 8, 10]. This semantics has the property of being historically easier to implement than call-by-need, and therefore likely easier to reason about formally. In this paper, we build on recent work developing a simple method for implementing call-by-need semantics, which has enabled us to implement, and reason formally about the correctness of, call-by-need. We use the Coq proof assistant [2] to implement and prove the correctness of our compiler. We start with a source language of λ calculus with de Bruijn indices:

$$\begin{aligned} t ::= & t \, t \mid x \mid \lambda \, t \\ & x \in \mathbb{N} \end{aligned}$$

Our source semantics is the big-step operational semantics of the $\mathcal{C}\mathcal{E}$ machine, which uses shared environments to share results between instances of a bound variable. To strengthen the result, and relate it to a better-known semantics, we also show that the call-by-name $\mathcal{C}\mathcal{E}$ machine implements Curien's call-by-name calculus of closures.

It may surprise the reader to see that we do not start with a better known call-by-need semantics; we address this concern in Section 8. We hope that the proof of compiler correctness, along with the proof that our call-by-name version of the semantics implements Curien's call-by-name semantics, convinces the reader that we have indeed implemented a call-by-need semantics, despite not using a better definition of call-by-need.

For our target, we define a simple instruction machine, described in Section 5. This simple target allows us to describe the compiler and proofs concisely for the paper, while still allowing flexibility in eventually verifying a compiler down to machine code for some set of real hardware, e.g., x86, ARM, or Power.

Our main results is a proof that whenever the source semantics evaluates to a value, the compiled code evaluates to the same value. While there are stronger definitions of what qualifies as a verified compiler, we argue that this is sufficient in Section 8. This main result, along with the proof that the call-by-name version of our semantics implements Curien’s calculus of closures, are the primary contributions of this paper. We are unaware of any existing verified non-strict compilers, much less a verified compiler of call-by-need.

The paper is structured as follows. In Section 2 we give the necessary background. In Section 3 we describe the source syntax and semantics (the big-step \mathcal{CE} semantics) in detail. We also use this section to define a call-by-name version of the semantics, and show that it implements Curien’s calculus of closures [5]. In Section 4 we describe the small-step \mathcal{CE} semantics and its relation to the big-step semantics. In Section 5 we describe the instruction machine syntax and semantics. In Section 6 we describe the compilation from machine terms to assembly language. In Section 7 we describe how the evaluation of compiled programs is related to the small-step \mathcal{CE} semantics. We compose this proof with the proof that the small-step semantics implement the big-step semantics to show that the instruction machine implements the big-step semantics. In Section 8 we discuss threats to validity, future work, and related work. We conclude in Section 9. The Coq source code with all the definitions and proofs described in this document is available at https://github.com/stelleg/cem_coq.

2 BACKGROUND

Programming languages can roughly be broken into two camps: those with *strict* and those with *non-strict* semantics. A strict language is one in which arguments at a call-site are always evaluated, while a non-strict language only evaluates arguments when they are needed. One can further break non-strict into two categories: call-by-name and call-by-need. Call-by-name is essentially evaluation by substitution: an argument term or closure is substituted for every instance of a corresponding variable. This has the downside that it can result in exponential slowdown due to repeated work: every variable dereference must re-evaluate the corresponding argument. Call-by-need is an evaluation strategy devised to address this shortcoming. By sharing the result of argument evaluation between instances of a variable, one avoids duplicated work. Unsurprisingly, call-by-need is the default semantics implemented by compilers for non-strict languages like Haskell [13].

Also perhaps unsurprisingly, call-by-need implementations tend to be more complicated than their strict counterparts. For example, even attempts at simple call-by-need abstract machines such as the Three Instruction Machine [6] require lambda lifting and shared directions, both of which make formal reasoning more difficult. Our \mathcal{CE} machine avoids these complications by using shared environments to share evaluation results between instances of a variable. We showed in previous work that in addition to being simpler to implement and reason about, performance of this approach can often compete with the state of the art [15].

With recent improvements in higher order logics, machine verification of algorithms has become a valuable tool in software development. Instead of relying heavily on tests to reason about the correctness of programs, verification can prove that algorithms implement

their specification for *all* inputs. Having the ability to write down both the specification and the proof in a machine-checked logic removes the vast majority of bugs found in hand-written proofs, ensuring far higher confidence in correctness than other standard methods. Other approaches, such as fuzz testing, have confirmed that verified programs remove effectively all bugs [16].

This approach applies particularly well to compilers. Often, the specification for a compiler is relatively easy to write down: source level semantics for some languages are exceedingly straightforward to specify. In addition, writing tests for compilers that cover all cases is even more hopeless than most domains, due to the size and complexity of the domain and codomain. The amortized return on investment is high: all reasoning about programs compiled with a verified compiler is provably preserved.

Due to the complexities discussed above involved in implementing lazy languages, existing work has focused on compiling strict languages [4, 8, 10]. In this work we use the simple \mathcal{CE} machine as a base for a verified compiler of a lazy language, using the Coq proof assistant.

As with many areas of research, the devil is in the details. What exactly does it mean to claim a compiler is verified? Essentially, a verified compiler of a functional language is one that preserves computation of values. That is, we have a one directional implication: *if the source semantics computes a value, then the compiled code computes an equivalent value* [4]. The important thing to note is that the implication is only in one direction. If the source semantics never terminates, this class of correctness theorem says nothing about the behavior of the compiled code. This has consequences for Turing-complete source languages. If we are unsure if a source program terminates, and wish to run it to check experimentally if it does, if we run the compiled code and it returns a value, we cannot be certain that it corresponds to a value computed in the source semantics.

While in theory one could solve this by proving the implication the other direction, that is, *if the compiled code computes a value then the source semantics computes an equivalent value*, in practice this is prohibitively difficult. Essentially, the induction rules for the abstract machine make constructing such a proof monumentally tricky.

One approach for getting around this issue is to try and capture the divergent behavior by defining a diverging semantics explicitly [12]. Then we can safely say that *if the source semantics diverges according to our diverging semantics, then the compiled code also diverges*.

For this paper, we choose to take the approach of [4] and define verification as the first implication above, focusing on the case in which the source semantics evaluates to a value. This is still a strong result: any source program that has meaning compiles to an executable with equivalent meaning. In addition, if we ever choose to extend the language with a type system that ensures termination, or some notion of progress, then we can use this proof in combination with our verification proof to prove the other direction.

3 $\mathcal{C}\mathcal{E}$ BIG-STEP SEMANTICS

In this section we define our big-step source semantics. A big-step semantics has the advantage of powerful, easy-to-use induction properties. This eases reasoning about many program properties. We shall also define a small-step semantics and prove that it implements the big-step semantics, but by showing that our implementation preserves the big-step semantics, we prove preservation of any inductive reasoning on the structure of evaluation tree.

As discussed in Section 1, our source syntax is lambda calculus with de Bruijn indices. De Bruijn indices count the number of intermediate lambdas between them and their binding lambda.

$$\begin{aligned} t ::= & t \mid x \mid \lambda t \\ & x \in \mathbb{N} \end{aligned}$$

Throughout the rest of this paper, we will use Coq code directly to formalize some of the definition and theorems in the paper. This has the advantage of avoiding ambiguities and correctness issues involved in translating to LATEX. For example, the Coq code defining lambda calculus with de Bruijn indices is:

```
Inductive tm : Type :=  
| var : nat → tm  
| lam : tm → tm  
| app : tm → tm → tm.
```

The essence of the $\mathcal{C}\mathcal{E}$ semantics is that we implement a shared environment, and use its structure to share results of computations. This allows for a simple abstract machine which can operate on lambda calculus directly, which is uncommon among call-by-need abstract machines [6, 7, 9, 13]. This simplifies formalization, as we do not need to prove that these transformations are semantics-preserving. Another advantage to the $\mathcal{C}\mathcal{E}$ machine is that it has constant sized closures, obviating the need to reason about re-allocating the results of computation and adding indirections due to closure size changes from thunk to value [13]. We operate on closures, which combine terms with pointers into the shared environment, which is implemented as a heap. Every heap location contains a cell, which consists of a closure and a pointer to the next environment location, which we will refer to as the environment continuation. Variable dereferences index into this shared environment structure, and if/when a dereferenced location evaluates to a value, the original closure (potentially a thunk or closure not evaluated to WHNF) will be replaced with said value. The binding of a new variable extends the shared environment structure with a new cell. This occurs during application, which evaluates the left hand side to an abstraction, then extend the environment with the argument term closed under the environment pointer of the application. The app ensures that two variables bound to the same argument closure will point to the same location in the shared environment. Because they point to the same location by construction of the shared environment, we can update that location with the value computed at the first variable dereference, and then each subsequent dereference will point to this value. The variable rule applies the update by indexing into the shared environment structure and replacing the closure at that location with the resulting value. The big-step semantics is described in the standard way by logical relation in Figure 1. A few notes about definitions used in the semantics: the update $\Psi \times v$ function replaces the closure at location x in heap Ψ with value v . It is

```
Fixpoint clu (v env:nat) (h:heap) : option (nat * cell) :=  
match lookup env h with  
| None ⇒ None  
| Some (cl c a) ⇒ match v with  
| S n ⇒ clu n a h  
| O ⇒ Some (env, cl c a)  
end  
end.  
  
Inductive step : conf → conf → Type :=  
| Id : ∀ M x y z Φ Ψ v e,  
  clu y e Φ = Some (x, {M, z}) →  
  ⟨Φ⟩M ↴ ⟨Ψ⟩v →  
  ⟨Φ⟩close (var y) e ⇄ ⟨update Ψ x v⟩v  
| Abs : ∀ N Φ e,  
  ⟨Φ⟩close (lam N) e ⇄ ⟨Φ⟩close (lam N) e  
| App : ∀ M N B' Φ Ψ Y f e ne ae,  
  isfresh (domain Ψ) f →  
  f > 0 →  
  ⟨Φ⟩close M e ⇄ ⟨Ψ⟩close (lam B) ne →  
  ⟨Ψ, f {close N e, ne}⟩close B f ⇄ ⟨Y⟩close (lam B') ae →  
  ⟨Φ⟩close (app M N) e ⇄ ⟨Y⟩close (lam B') ae  
where " c1 ↴ c2 " := (step c1 c2).
```

Figure 1: Big Step $\mathcal{C}\mathcal{E}$ Semantics

worth noting that while the closures in the heap cells are mutable, the shared environment structure is never mutated. This property is crucial when reasoning about variable dereferences. The curly brackets (or the $c1$ constructor, depending on context) construct a heap cell containing a closure and a environment continuation. The app constructor represents an application, while the var constructor is a variable constructor and lam is the constructor for an abstraction. The clu function looks up a variable index in the shared environment structure by following environment continuation pointers. The $\langle\Phi\rangle c$ notation constructs a configuration from a heap Φ and a closure c . Note that we require that fresh heap locations are greater than zero. This is required for reasoning about compilation to the instruction machine, which we will return to in Section 5. The isfresh relation is defined as $f \notin \text{dom}(\Psi)$. Of course, for a real implementation, this is far too strong a constraint, as it doesn't allow any sort of heap re-use. We return to this issue in Section 8, and discuss how this could be relaxed to either allow reasoning about garbage collection or direct heap-reuse.

The fact that our natural semantics is defined on lambda calculus with de Bruijn indices differs from most existing definitions of call-by-need, such as Ariola's call-by-need [1] or Launchbury's lazy semantics [9]. These semantics are defined on lambda calculus with named variables. While it should be possible to relate our semantics to these¹, the comparison is certainly made more difficult by this disparity. A more fruitful relation to semantics operating on lambda

¹Both of these well known existing semantics have known problems that arise during formalization, as discussed in Section 8

calculus with named variables would likely be relating Curien's calculus of closures to call-by-name semantics implemented with substitution. We return to this discussion in Section 8.

As mentioned in Section 1, these big-step semantics do not explicitly include a notion of nontermination. Instead, nontermination would be implied by the negation of the existence of an evaluation relation. This prevents reasoning directly about nontermination in an inductive way, but for the purpose of our primary theorem this is acceptable.

One interesting property of defining an inductive evaluation relation in a language such as Coq is that we can do computation on the evaluation tree. In other words, the evaluation relation given above defines a data type, one that we can do computation on in standard ways. For example, we could potentially compute properties such as size and depth, which would be related to operational properties of compiled code. We hope in future work to explore this approach further.

Finally, given a term t , we define the initial configuration as $\langle\rangle\text{close } t \ 0$. As discussed, the choice of the null pointer for the environment pointer is not arbitrary, but chosen across our semantics uniformly to represent failed environment lookup.

3.1 Call-By-Name

In this section we define a call-by-name variant of our big-step semantics and prove that it is an implementation of Curien's call-by-name calculus of closures [5].

See Figure 3 for the definition of our call-by-name semantics. Note that the only change from our call-by-need semantics is that we do not update the heap location with the result of the dereferenced computation. This is the essence of the difference between call-by-name and call-by-need.

A well known existing call-by-name semantics is Curien's calculus of closures [5]. Refer to Figure 2 for a formalization of this semantics. This semantics defines closures as a term, environment pair, where an environment is a list of closures. Abstractions are in weak head normal form, variables index into the environment using the `nth_error` standard library function, and applications evaluate the left hand side to a value, then extend the environment of the value with the closure of the argument.

We define a heterogeneous equivalence relation between our shared environment and Curien's environment. Effectively, this relation is the proposition that the shared environment structure is a linked list implementation of the environment list in Curien's semantics. This is defined inductively, and we require that every closure reachable in the environment is also equivalent. We say two closures are equivalent if their terms are identical and their environments are equivalent.

Given these definitions, we can prove that our call-by-name semantics implement Curien's call by name semantics:

THEOREM 3.1. *If a closure c in Curien's call-by-name semantics is equivalent to a configuration c' , and c steps to v , then there exists a v' that our call-by-name semantics steps to from c' that is equivalent to v .*

```
Inductive closure := | close : tm → list closure → closure.
Definition env := list closure.
```

```
Inductive step : closure → closure → Type :=
| Abs : ∀ b e,
  step (close (lam b) e) (close (lam b) e)
| Var : ∀ x e v c,
  nth_error e x = Some c →
  step c v →
  step (close (var x) e) v
| App : ∀ m n b e v mve,
  step (close m e) (close (lam b) mve) →
  step (close b (close n e; mve)) v →
  step (close (app m n) e) v.
```

Figure 2: Curien's Calculus of Closures

```
Inductive step : conf → conf → Type :=
| Id : ∀ M x y z Φ Ψ v e,
  clu y e Φ = Some (x, {M, z}) →
  ⟨Φ⟩M ↳ ⟨Ψ⟩v →
  ⟨Φ⟩close (var y) e ↳ ⟨Ψ⟩v
| Abs : ∀ N Φ e,
  ⟨Φ⟩close (lam N) e ↳ ⟨Φ⟩close (lam N) e
| App : ∀ N M B Φ Ψ f e ne ae,
  isfresh (domain Ψ) f →
  f > 0 →
  ⟨Φ⟩close M e ↳ ⟨Ψ⟩close (lam B) ne →
  ⟨Ψ, f {close N e, ne}⟩close B f ↳ ⟨Υ⟩close (lam B') ae →
  ⟨Φ⟩close (app M N) e ↳ ⟨Υ⟩close (lam B') ae
where " c1 ↳ c2 " := (step c1 c2).
```

Figure 3: Call-by-Name \mathcal{CE} Semantics

PROOF OUTLINE. The proof proceeds by induction on Curien's step relation. The abstraction rule is a trivial base case. The variable lookup rule uses a helper lemma that proves by induction on the variable that if the two environments are equivalent and the `nth_error` function looks up a closure, then the `clu` function will look up an equivalent closure. The application rule uses a helper lemma proves that a fresh allocation will keep any equivalent environments equivalent, and that the new environment defined by the fresh allocation will be equivalent to the extended environment of Curien's semantics.

By proving that Curien's semantics is implemented by the call-by-name variant of our semantics, we provide further evidence that our call-by-need is a meaningful semantics. While eventually we would like to prove that the call-by-need semantics implements an optimization of the call-by-name, we leave that for future work.

One important note is that nowhere do we require that a term being evaluated is closed under its environment. Indeed, it's possible that a term with free variables can be evaluated by both semantics to a value as long as a free variable is never dereferenced. This theme

```

Inductive step : transition state :=
| Upd :  $\forall \Phi b e l s, st \Phi (inr l::s) (close (lam b) e) \rightarrow _s$ 
| st  $\Phi (inr l::s) (close (lam b) e) \rightarrow _s$ 
| st  $(update \Phi l (close (lam b) e)) s (close (lam b) e) \rightarrow _s$ 
| Var :  $\forall \Phi s v l c e e', clu v e \Phi = Some (l, cl c e') \rightarrow$ 
| clu v e  $\Phi = Some (l, cl c e') \rightarrow$ 
| st  $\Phi s (close (var v) e) \rightarrow _s$ 
| st  $\Phi (inr l::s) c \rightarrow$ 
| Abs :  $\forall \Phi b e f c s, isfresh (domain \Phi) f \rightarrow$ 
| f > 0  $\rightarrow$ 
| st  $\Phi (inl c::s) (close (lam b) e) \rightarrow _s$ 
| st  $((f, cl c e):: \Phi) s (close b f) \rightarrow$ 
| App :  $\forall \Phi e s n m, st \Phi s (close (app m n) e) \rightarrow _s$ 
| st  $\Phi (inl (close n e):: s) (close m e) \rightarrow _s$ 
where "  $c1 \rightarrow _s c2 := (step c1 c2)$ ".

```

Figure 4: Small Step \mathcal{CE} Semantics

will recur through the rest of the paper, so it is worth keeping in mind.

4 \mathcal{CE} SMALL-STEP SEMANTICS

In this section we discuss the small-step semantics of the \mathcal{CE} machine, and show that it implements the big-step semantics of Section 3. This is a fairly straightforward transformation implemented by adding a stack. The source language is the same, and we simply add a stack to our configuration (and call it a state). The stack elements are either argument closures (inl) or update markers (inr). Update markers are pushed onto the stack when a variable dereferences that location in the heap. When they are popped by an abstraction, the closure at that location is replaced by said abstraction, so that later dereferences by the same variable in the same scope dereference the value, and do not repeat the computation. Argument closures are pushed onto the stack by applications, with the same environment pointer duplicated in the current closure and the argument closure. Argument closures are popped off the stack by abstractions, which allocate a fresh memory location, write the argument closure to it, write the environment continuation as the current environment pointer, then enter the body of the abstraction with the fresh environment pointer. This is the mechanism used for extending the shared environment structure. The semantics is defined formally in Figure 4. The update function replaces a closure in a heap cell with the specified closure.

Note that the presentation given here differs slightly from our previous presentation [15], which inlined the lookup into the machine steps. This is to simplify formalization and relation to the big-step semantics, but does not change the semantics of the machine. As a tradeoff, it does make the relation to the instruction machine in the later sections slightly more involved, but it is generally a superficial change.

4.1 Relation to Big Step

Here we prove that the small-step semantics implements the big-step semantics of Section 3. This requires first a notion of reflexive transitive closure, which we define in the standard way. We also make use of the fact that the reflexive transitive closure can be defined equivalently to extend from the left or right.

LEMMA 4.0.1. *If the big-step semantics evaluates from one configuration to another, then the reflexive transitive closure of the small-step semantics evaluates from the same starting configuration with any stack to the same value configuration with that same stack.*

PROOF OUTLINE. The proof proceeds by induction on the big-step relation. We define our induction hypothesis so that it holds for all stacks, which gives us the desired case of the empty stack as a simple specialization. The rule for abstractions is the trivial base case. Var rule applies as the first step, and the induction hypothesis applies to the stack with the update marker on it. To ensure that the Upd rule applies we use the fact that the big-step semantics only evaluates to abstraction configurations, and the fact that the reflexive transitive closure can be rewritten with steps on the right. For the Application rule, we take advantage of the fact that we can append two evaluations together, as well as extend a reflexive transitive closure from the left or the right. As with the Var rule we use the fact that the induction rule is defined for all stacks to ensure we evaluate the left hand side to a value with the argument on the top of the stack. Finally, we extend the environment with the argument closure, and evaluate the result to a value by the second induction hypothesis.

Adding a stack in this fashion is a standard approach to converting between big step and small-step semantics. Still, we appreciate that this approach applies here in a straightforward way.

5 INSTRUCTION MACHINE

Here we describe in full the instruction machine syntax and semantics. We choose a simple stack machine with a Harvard architecture (with separate instruction and heap memory). We use natural numbers for pointers, though it shouldn't be too difficult to replace these with standard-sized machine words, e.g., 64 bits, making the stack and malloc operations partial. Our stack is represented as a list of pointers, though again it should be a relatively straightforward exercise to represent the stack in main memory. With the fixed machine word size, we would need to make push operations partial to represent stacks this way. We define our machine to have only four registers: an instruction pointer, an environment pointer, and two scratch registers. Our instruction set is minimal, consisting only of a conditional jump instruction, pop and push instructions, a mov instruction, and a new instruction for allocating new memory. Note that for our program memory, we have pointers to basic blocks, but for simplicity of proofs we choose to not increment the instruction pointer within a basic block. Instead, the instruction pointer is constant within a basic block, only changing between basic blocks. In fact, we represent the program as a list of basic blocks, with pointers indexing into the list. This has the advantage of letting us easily reason about sublists and their relation to terms. As with other design decisions, this also should be fairly unproblematic for formalization

```

Inductive Reg := 
| IP
| EP
| R1
| R2.

Inductive WO := 
| WR : Reg → WO
| WM : Reg → nat → WO.

Infix "%" := WM (at level 30).

Inductive RO := 
| RW : WO → RO
| RC : nat → RO.

Inductive Instr : Type := 
| push: RO → Instr
| pop : WO → Instr
| new : nat → WO → Instr
| mov : RO → WO → Instr.

Inductive BasicBlock : Type := 
| instr: Instr → BasicBlock → BasicBlock
| jump: option (R0*Ptr) → RO → BasicBlock.

```

Definition Program := list BasicBlock.

Figure 5: Instruction Machine Syntax

to a more realistic hardware design. The full syntax of the machine is given in Figure 5.

We separate read (RO) and write (WO) operands. Write operands can be registers (WR) or memory (WM). Read operands can be any write operand (RW) or a constant (RC). For reading, we have a read relation, which takes a read operand and a state and is inhabited when the third argument can be read from that read operand in that state. Similarly, a write relation is inhabited when writing the second argument into the first in a state defined by the third argument results in the state defined by the fourth argument.

The machine semantics should be fairly unsurprising. A State, constructed by `st`, consists of a register file, program memory, a stack, and a heap. The `push` instruction takes a read operand and pushes it onto the stack. The `pop` instruction pops the top of the stack into a write operand. The `mov` instruction moves a machine word from a read operand to a write operand. The `jump` instruction is parameterized by an optional pair, which, if present, reads the first element of the pair from a read operand, checks if it is zero, and if so sets the IP to the second element of the pair, which is a constant pointer. If the condition is not zero, then it sets the IP to the instruction pointer contained in the second jump argument. If we pass nothing as the first argument, then it becomes an unconditional set of the IP to the value read from the second argument. Note that the second argument is a read operand, so it can either be a constant or read from a register or memory. This means it can be

effectively either a direct or indirect jump, both of which are used in the compilation of lambda terms. The new instruction allocates a contiguous block of new memory and writes the resulting pointer to the fresh memory into a write operand. We take the approach of not choosing a particular allocation strategy. Instead, we follow existing approaches and parameterize our proof on the existence of such functionality [4]. For simplicity, we assume that the allocation function returns completely fresh memory, though it should be possible to modify this assumption to be less restrictive, i.e., let it re-use heap locations that are no longer live. The complete semantics of the machine is given in Figure 6. Note that we separate instruction steps and basic block steps. Recall that a basic block is a sequence of instructions that ends with a jump. The `step_bb` relation will execute the instructions in the basic block in order, then set the IP in accordance with the jump semantics. The `step` relation dereferences a basic block at the current IP, and if executing the basic block results in a new state, then the machine executes to that state.

6 COMPILER

In this section we describe the compiler, which compiles lambda terms with de Bruijn indices to programs. The compiler proceeds by recursion on lambda terms, keeping a current index into the program to ensure correct linking without a separate pass. For variables, when we get to zero we push the current environment pointer and a null instruction pointer to denote the update marker to the location of the closure being entered. Then we `mov` the closure at that location into R1 and EP, and `jump` to R1, recalling that the `jump` sets the IP. For nonzero variables, we replicate traversing the environment pointer i times before loading the closure. For applications, we calculate the program location of the argument basic block, and push that and the current environment pointer onto the stack, effectively pushing an argument closure on top of the stack. We then `jump` to the left hand side of the application, as is standard for push-enter evaluation. For abstractions, we use a conditional jump depending on whether the top of the stack is a null pointer (and therefore an update marker) or a valid instruction pointer (and therefore an argument). If it is an update marker, we update the heap location defined by the update marker with the current value instruction pointer and the current environment pointer. We must point to the first of the three abstractions basic blocks, as this value could later update another heap location as well. In the case that the top of the stack was a valid instruction pointer, we allocate a new chunk of 3 word of memory, and `mov` the argument closure into it, with the current environment pointer as the environment continuation. We then set our current environment pointer to this fresh location. This is the process by which we extend our shared environment structure in the instruction machine. Finally, we perform an unconditional jump to the next basic block, which is the first basic block of the compiled body of the lambda. As this is an unconditional jump to the next basic block, for real machine code this jump can be omitted.

Being able to define the full compiler this simply is crucial to this verification project. Other, more sophisticated implementations of call-by-need, such as the STG machine, are much harder to implement and reason about. It is worth noting that despite this simplicity, initial tests suggest that performance is not as horrible as one might suspect, and is often competitive with state of the art [15].

```

Inductive step_bb : BasicBlock → State → State → Type :=
| step_push: forall rf p s h ro bb v sn,
  read ro (st rf p s h) v →
  step_bb bb (st rf p (v:: s) h) sn →
  step_bb (instr(push ro) bb) (st rf p s h) sn
| step_pop: forall rf p s h wo bb w s' sn,
  write wo w (st rf p s h) s' →
  step_bb bb s' sn →
  step_bb (instr(pop wo) bb) (st rf p (w:: s) h) sn
| step_new: forall rf p s h wo bb w s' n sn,
  (forall i, i < n → not (In (i+w) (domain h))) →
  w > 0 →
  write wo w (st rf p s (zeroes n w ++ h)) s' →
  step_bb bb s' sn →
  step_bb (instr(new n wo) bb) (st rf p s h) sn
| step_mov: forall s bb ro wo s' v sn,
  read ro s v → write wo v s' →
  step_bb bb s' sn →
  step_bb (instr(mov ro wo) bb) s sn
| step_jump0: forall ro k j s s',
  read ro s 0 →
  write (WR IP) k s s' →
  step_bb (jump (Some (ro, k)) j) s s'
| step_jumpS: forall ro k j s s' l k',
  l > 0 →
  read ro s l →
  read j s k →
  write (WR IP) k s s' →
  step_bb (jump (Some (ro, k')) j) s s'
| step_jump: forall ro s s' l,
  read ro s l →
  write (WR IP) l s s' →
  step_bb (jump None ro) s s'
.
.

Inductive step : transition State :=
| enter: forall rf p s h k bb sn,
  read IP (st rf p s h) k →
  nth_error p k = Some bb →
  step_bb bb (st rf p s h) sn →
  step (st rf p s h) sn.

```

Figure 6: Instruction Machine Semantics

As with the relation discussed in Section 3, note that even when compiling, we do not require that a term is closed to compile it. Indeed, we will happily generate code that if entered, will attempt to dereference the null pointer, leaving the machine stuck. Because we are only concerned with proving that we implement the source semantics in the case that it evaluates to a value, this is not a problem. If we wanted to strengthen our proof further, we would try to show that if the source semantics gets stuck trying to dereference a free variable, the implementation would get stuck in the same way, both failing to dereference a null pointer.

```

Infix ";" := instr(at level 30, right associativity).

Fixpoint var_inst (i : nat) : BasicBlock := match i with
| 0 ⇒ push EP ;
  push (RC 0) ;
  mov (EP%0) R1 ;
  mov (EP%1) EP ;
  jump None R1
| S i ⇒ mov (EP%2) EP ;
  var_inst i
end.

Fixpoint compile(t : tm) (k : nat) : Program := match t with
| var v ⇒ [var_inst v]
| app m n ⇒ let ms := compile m (1+k) in
  let nk := 1+k+length ms in
  push EP ;
  push (RC nk) ;
  jump None (RC (1+k)) :: ms ++
  compile n nk
| lam b ⇒ pop R1 ;
  jump (Some (RW (WR R1), (1+k))) (RC (2+k)) :: (*Update*)
  pop R1 ;
  mov (RC k) (R1%0) ;
  mov EP (R1%1) ;
  jump None (RC k) :: (*Take*)
  new 3 R2 ;
  mov R1 (R2%0) ;
  pop (R2%1) ;
  mov EP (R2%2) ;
  mov R2 EP ;
  jump None (3+k) :: compile b (3+k)
end.

```

Figure 7: Compiler Definition

7 COMPILER CORRECTNESS

In this section we define a relation between the state of the small-step semantics and the state of the instruction machine semantics, and show that the instruction machine implements the small-step semantics under that relation.

In general, we implement closures as instruction pointer, environment pointer pairs. For the instruction pointers, we relate them to terms via the compile function defined in Section 6. Essentially, we require that the instruction pointer points to a list of basic blocks that the related term compiles to. For the current closure, we relate the instruction pointer register in the instruction machine to the current term in the small-step source semantics. The environment pointers of each machine are more similar. Given a relation between the heaps of the two machines, we define the relation between two environment pointers as existing in the relation of the heaps, or both

being the null pointers. While it should be possible to avoid this special case, during the proof it became apparent that not having the special case made the proof significantly harder. This forces us to add the constraint to all machines that pointers are non-null, which for real hardware shouldn't be an issue.

We use null pointers in two crucial ways. One is to explicitly define the root of the shared environment structure in both the source semantics and the machine semantics. The other use is for instruction pointers. To differentiate between update markers and pointers to basic blocks, we use a null pointer to refer to an update marker, and a non-null pointer as an instruction pointer for an argument closure. Note that in fact, while the null pointers in heaps required us to only allocate non-null fresh locations in the heaps of our semantics, using null pointers to denote update markers requires no change to our program generation, due to the fact that an argument term of an application cannot occur at position 0 in the program.

The relation between the heaps of the small-step source semantics and the instruction machine is the trickiest part of the state relation. Note that for each location in the source semantics heap, we have a cell with a closure and environment continuation pointer. Naturally, the instruction machine represents these as three pointers: two for the closure (the instruction pointer and environment pointer) and one for the environment continuation. The easiest approach turned out to be to use the structure of the heap constructs to define a one-to-three mapping between this single cell and the three machine words. The structure used for each of the heaps is a list of pointer, value bindings. We use the ordering of these bindings in the list to define a one binding to three binding mapping between the source heap and the machine heap. We define an analogue to the standard In list relation that defines when an element is in a list for this heap relation, proceeding recursively on the inductive relation structure. This allows us to define a notion of which pairs of each type of closure are in the heap, along with their respective locations. Due to the ordering in which they are allocated in the heap during evaluation, each pair of memory allocations corresponds to an equivalent cell. We use this property as a heap equivalence property that is preserved through evaluation: every binding pair in the heap relation property described above defines equivalent closures and environment continuations. For the relation between our stacks, we define a similar notion. For update markers, we require that every update marker points to related environments (they are two pointers that exist in the heap relation). For argument closures, we require that the closures are equivalent (the instruction pointer and environment pointer are equivalent to their respective counterparts in the small-step semantics).

In summary, we require that the current closure in the small-step semantics is equivalent to the closure represented by the instruction pointer, environment pointer pair, and that the stacks and the heaps are equivalent. See Figure 8 for some of the formal definitions of these relations.

Now that we have our relation between heaps, we can state our primary lemma.

LEMMA 7.0.1. *Given that an instruction machine state i is related to a small-step semantics state s , and that small-step semantics state steps to a new state s' , the instruction machine will step in zero or more steps to a related state i' .*

```

Inductive heap_rel : cesm.heap → im.Heap → Type :=
| heap_nil : heap_rel [] []
| heap_cons : ∀ l l' ne ch ih ip ep ine e t,
  l ≠ domain ch →
  l' ≠ domain ih → S l' ≠ domain ih → S (S l') ≠ domain ih →
  l > 0 → l' > 0 →
  heap_rel ch ih →
  heap_rel
  ((l, cl (close t e) ne):: ch)
  ((l', ip)::( S l', ep)::( S (S l'), ine):: ih).
Definition prog_eq (p : Ptr) (pr : Program) (t : tm) :=
let subpr := compile t p in
subpr = firstn (length subpr) (skipn p pr).

Inductive env_eq (ch : cesm.heap) (ih : im.Heap) (r : heap_rel ch ih)
: nat → Ptr → Type :=
| e0 : env_eq ch ih r 0 0
| es : ∀ l e ne t il ip ep nep,
  in_heap_rel ch ih r l e ne t il ip ep nep →
  env_eq ch ih r l il.

Inductive clos_eq (ch : cem.heap)
  (ih : im.Heap)
  (r : heap_rel ch ih)
  (p : Program):
closure → Ptr → Ptr → Type :=
| c_eq : ∀ t e ip ep,
  prog_eq ip p t →
  env_eq ch ih r e ep →
  clos_eq ch ih r p (close t e) ip ep.

Inductive state_rel (cs : cesm.state) (is : im.State) : Type :=
| str : ∀ r,
  heap_eq(st_hp cs) (st_h is) r
  (st_p is) →
  clos_eq(st_hp cs) (st_h is) r
  (st_p is) (st_cl cs)
  (rff (st_rf is) IP)
  (rff (st_rf is) EP) →
  stack_eq(st_hp cs) (st_h is) r
  (st_p is) (st_st cs) (st_s is) →
  state_rel cs is.

```

Figure 8: Select Relation Definitions Between Small Step $\mathcal{C}\mathcal{E}$ and the Instruction Machine

PROOF OUTLINE. Our proof proceeds by case analysis on the step rules for the small-step semantics. We'll focus on the second half of the proof, that i' is related to s' . The proofs that i evaluates to i' follow fairly directly from the compiler definition given in Section 6. For the Var rule, because we need to proceed by induction, we have to define a separate lemma and proceed by induction on a basic block while forgetting the program, as the induction hypothesis is invalid in the presence of the program. We then use the lemma

to show that evaluation of a compiled variable implements the evaluation of the variable in the small-step semantics. In particular, we use the null environment as a base case for our induction, as we know the only way lookup could fail is if both environment pointers are null, but that cannot be the case due to the fact that we know that the small-step semantics must have successfully looked up its environment pointer in the heap. Therefore the only option is for both environment pointers to exist in the heap relation, which when combined with the heap equivalence relation in the outer proof gives us the necessary property that the environment continuations are equivalent. Finally, because the last locations reached must have been in the heap relation, we know they are equivalent environment pointers, and therefore the stack relation is preserved when we push the update marker onto the heap. For the App rule, we use the definition of our compiler to prove that the argument term and argument instruction pointer are equivalent and that the left hand side term and instruction pointer are also equivalent. They share an environment pointer which is equivalent by the fact that the application closures are related. This proves that the stack relation is preserved as well as the current closure, while the heap is unchanged. For the Lam rule, we allocate a fresh variable and because of our stack relation we can be sure that the closures that we allocate are equivalent, as well as the environment continuations, as they are taken from the previous current continuation. Because of how we define it, the new allocations are equivalent under our heap relation, and preserve heap equivalence. Finally, the Upd rule trivially preserves the stack and current closure relations, and for proving that the relation is preserved for the heap, we proceed by induction on the heap relation. In addition, we must prove a supporting lemma that all environment relations are preserved by the update.

To paraphrase Pascal: given more time, I would have written a shorter proof.

We now have a proof that the small-step semantics implements the big-step semantics, and a proof that the instruction machine implements the small-step semantics. We can now combine these to get our correct compiler theorem.

THEOREM 7.1. *If a term t placed into the initial configuration for the big-step semantics evaluates to a value configuration v , then the instruction machine starting in the initial state with $\text{compile } \emptyset t$ as its program will evaluate to a related state v' .*

PROOF OUTLINE. We first require that the relation defined between the small-step semantics state and the instruction machine state holds for the initial configurations. This follows fairly directly from the definition of the initial conditions and the compile function. Second, we have by definition of reflexive transitive closure that the previous lemma implies that if the reflexive transitive closure of the small-step relation evaluates in zero or more steps from a state c to a state v , then a related state of the instruction machine c' will evaluate to a state v' which is related to v . We use these two facts, along with the proof that the small-step implements the big-step for any stack, specialized on the empty stack, to prove our theorem.

See Figure 9 for a formal statement of Theorem 7.1.

It is worth recalling exactly what the relation implies about the two value states. Namely, in addition to the value closures being

```
Theorem compile_correct (t : db.tm) v : cem.step (cem.I t) v
sigT( v', refl_trans_clos im.step(im.I (compile t)) v' *)
state_rel (cesm.st (cem.conf_h v) nil (cem.conf_c v)) v').
```

Figure 9: Formal Definition of Theorem 7.1

equivalent, their heaps and environments are equivalent, so that every reachable closure in the environment is equivalent between the two.

8 DISCUSSION

Here we reflect on what we have accomplished, including threats to validity, future work, related work, and general discussion of the results.

One thing we'd like to communicate is the difficulty we had in writing comprehensible proofs. The reader is discouraged from attempting to understand the proofs in any way by reading the Coq tactic source code. While we attempted to keep our definitions and lemmas as clean and comprehensible as possible, we found it extremely difficult to do the same with tactics. Partially this may be a failure on our part to become more familiar with the tactic language of Coq, but we suspect that the imperative nature of tactic proofs prevents composability of tactic meta-programs.

Another lesson was the importance of good induction principles. For example, in Section 8.1, we will discuss the issue of only proving the implication of correctness in one direction. This is effectively a product of the power of the inductive properties of high level semantics, which makes them so much easier to reason about. Indeed, this lesson resonates with the purpose of the paper, which is that we'd like to reason about high level semantics, because they are so much easier to reason about due to their pleasant inductive properties, and have that reasoning preserved through compilation.

8.1 Threats to Validity

There are a few potential threats to validity that we address in this section. The first is the one mentioned in Section 2, that we only show that our compiler is correct in the case of termination of the source semantics. In other words, if the source semantics doesn't terminate, we can say nothing about how the compiled code behaves. This means that we could in theory get termination when the source semantics never does.

One argument in defense of our verification is that *we generally only care about preservation of semantics for preserving reasoning about our programs*. In other words, if we have a program that we can't reason about, and therefore may not terminate, we care less about having a proof that semantics are preserved. Of course, this is a claim about most uses of program analysis. There are possible analyses that could say things along the lines of *if* the source program terminates, then we can conclude x . We claim these cases are rare, and therefore the provided proof of correctness can still be applied to most use cases.

Another potential threat to validity is the use of a high level instruction machine language. While we claim that its high level and simplicity should make it possible to show that a set of real ISAs implement this instruction machine, we haven't formally verified this

step. We certainly agree that this would make for valuable future work, and hope that the reader agrees that nothing in the design of our high level instruction machine would prevent such work.

As a dual to the issue of a high level instruction machine language, some readers may take issue with calling lambda calculus with de Bruijn indices a "source language". Indeed, we do not advocate writing programs in such a language. Still, the conversion between lambda calculus with named variables and lambda calculus with de Bruijn indices is a well understood topic, it is not particularly interesting in the context of a verified compiler. Indeed, as noted below, semantics using named variables and substitution can be hard to get right [3, 11], so we stand by our decision to use a semantics based on lambda calculus with de Bruijn indices. One potential approach for future work would be to prove that a call-by-name semantics using substitution is equivalent to Curien's calculus of closures, which when combined with a proof that our call-by-need implements our call-by-name, would prove the compiler implements the semantics of a standard lambda calculus with named variables.

A third threat to the validity of this work is the question of whether we have really proved that we have implemented *call-by-need*. The question naturally arises of what exactly it means to prove an implementation of call-by-need is correct. There are certainly well-established semantics [1, 9], so one option would be to directly prove that the \mathcal{CE} semantics implements one of those existing semantics. Unfortunately, recent work has shown that both of these have small issues that arise when formalized that require fixes. Indeed, we did stray down this path a good ways and discovered one of these issues which has been previously described in the literature [11]. This raises the question of whether or not semantics that aren't obviously correct are a good base for what it means to be a call-by-need semantics. Instead, we have chosen to relate our call-by-name semantics formally to a semantics that is obviously correct, Curien's calculus of closures. Along with the tiny modification required for memoization of results, we hope that we have convinced the reader that it is *extremely likely* that the memoization of results is correct. Of course, further evidence such as examples of correct evaluation would go further to convince the reader, and for that we encourage readers to play with a toy implementation at https://github.com/stelleg/cem_pearl. Finally, a more convincing result would be a proof that the call-by-need semantics implement the call-by-name semantics.

Yet another threat to validity is our approach (or lack of approach) to heap-reuse. For simplicity, we have assumed that our fresh locations are fresh with respect to all existing bindings in the heap. Of course, this is unsatisfactory when compared to real implementations. It would be preferable to have our freshness constraint relaxed to only be fresh with respect to live bindings on the heap. We believe that this modification should be possible, at the cost of increased complexity in the proofs.

8.2 Future Work

In addition to some of the future work discussed as ways of addressing issues in Section 8.1, there are some additional features that we think make for exciting areas of future work.

One such area is reasoning about preservation of operational properties such as time and space requirements. This would enable reasoning about time and space properties at the source level and ensuring that these are preserved through compilation. In addition, there is the possibility of verified optimizations, where one can prove that some optimizations are both *correct*, in that they provably preserve semantics, and *true* optimizations, in that they only improve performance with respect to some performance model. By defining a baseline compiler and proving that it preserved operational properties such as time and space usage, one would have a good platform for which to apply this class of optimizations, resulting in a full compiler that verifiably preserves bounds on time and space consumption. As with correctness, reasoning about operational properties is often likely to be easier in the context of the easy-to-reason-about high level semantics, and having that reasoning provably preserved would be extremely valuable.

Those familiar with Coq may have noticed that in general we operate in the Type universe instead of Prop. This is largely due to the desire to eventually do computation on the proof structures defined here, towards the goal of proving operational properties. For example, by operating on the big-step relation data structure, one could compute time and space requirements for a program evaluation.

Another exciting area of future work is powerful proofs of type preservation through compilation. While there has been existing work on type-preserving compilers, fully verified compilers like this one provide such a strong property that type-safety should fall out directly.

8.3 Related Work

Chlipala implements a compiler from a STLC to a simple instruction machine in [4]. In many ways it is more sophisticated than our work: it converts to CPS, performs closure conversion, and proves a similar compiler correctness theorem to the one we've proved here. The primary difference is that we've defined a call-by-need compiler, which forces us to reason about updating thunks in the heap, a challenge not shared by call-by-value implementations.

Breitner formalizes Launchbury's natural semantics and proves an optimization is sound with respect to the semantics [3, 9]. By relating his formalization with ours, these projects could be combined to prove a more sophisticated lazy compiler correct: one with non-trivial optimizations applied.

CakeML [8] is a verified compiler for a large subset of the Standard ML language formalized in HOL4 [14]. Like Chlipala's work, this is a call-by-value language, though they prove correctness down to an x86 machine model, and are working with a much larger real-world source language. They also make divergence arguments along the lines of [12], strengthening their correctness theorem in the presence of nontermination. It's also worth noting that like [10], they are also formalizing a frontend to the compiler.

9 CONCLUSION

We have presented the first verified compiler of a non-strict lambda calculus. In addition to proving that our call-by-need semantics is preserved through compilation, we have proved that Curien's calculus of closures is implemented by our call-by-name semantics. We

argue that this provides compelling evidence that our compiler is a true verified compiler of call-by-need.

We hope that this work can serve as a foundation for future work on real-world verified compilers. While it is clearly a toy compiler, we have reason to believe that performance is acceptable and can be further improved [15]. In combination with efforts to formalize semantics of real-world languages like Haskell, we hope that this work can help us move towards fully verified non-strict programs.

REFERENCES

- [1] Z. M. Ariola, J. Maraist, M. Odersky, M. Felleisen, and P. Wadler. A call-by-need lambda calculus. In *Proceedings of the 22nd ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, pages 233–246. ACM, 1995.
- [2] B. Barras, S. Boutin, C. Cornes, J. Courant, J.-C. Filliatre, E. Gimenez, H. Herbelin, G. Huet, C. Munoz, C. Murthy, et al. *The Coq proof assistant reference manual: Version 6.1*. PhD thesis, Inria, 1997.
- [3] J. Breitner. *Lazy Evaluation: From natural semantics to a machine-checked compiler transformation*. PhD thesis, Karlsruher Institut für Technologie, 2017.
- [4] A. Chlipala. A certified type-preserving compiler from lambda calculus to assembly language. In *PLDI '07: Proceedings of the 28th ACM SIGPLAN Conference on Programming Language Design and Implementation*, volume 42, pages 54–65. ACM, 2007.
- [5] P.-L. Curien. An abstract framework for environment machines. *Theoretical Computer Science*, 82(2):389–402, 1991.
- [6] J. Fairbairn and S. Wray. TIM: A simple, lazy abstract machine to execute supercombinators. In *Functional Programming Languages and Computer Architecture*, pages 34–45. Springer, 1987.
- [7] T. Johnsson. Efficient compilation of lazy evaluation. In *SIGPLAN Notices*, 1984.
- [8] R. Kumar, M. O. Myreen, M. Norrish, and S. Owens. Cakeml: A verified implementation of ml. In *Proceedings of the 41st ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, POPL '14, pages 179–191, New York, NY, USA, 2014. ACM. ISBN 978-1-4503-2544-8. doi: 10.1145/2535838.2535841. URL <http://doi.acm.org/10.1145/2535838.2535841>.
- [9] J. Launchbury. A natural semantics for lazy evaluation. In *Proceedings of the 20th ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, pages 144–154. ACM, 1993.
- [10] X. Leroy. The compcert c verified compiler. *Documentation and users manual. INRIA Paris-Rocquencourt*, 2012.
- [11] K. Nakata and M. Hasegawa. Small-step and big-step semantics for call-by-need. *Journal of Functional Programming*, 19(6):699–722, 2009.
- [12] S. Owens, M. O. Myreen, R. Kumar, and Y. K. Tan. Functional big-step semantics. In P. Thiemann, editor, *Programming Languages and Systems*, pages 589–615, Berlin, Heidelberg, 2016. Springer Berlin Heidelberg. ISBN 978-3-662-49498-1.
- [13] S. L. Peyton Jones. Implementing lazy functional languages on stock hardware: The spineless tagless G-machine. *Journal of functional programming*, 2(2):127–202, 1992.
- [14] K. Slind and M. Norrish. A brief overview of hol4. In *International Conference on Theorem Proving in Higher Order Logics*, pages 28–32. Springer, 2008.
- [15] G. Stelle, D. Stefanovic, S. L. Olivier, and S. Forrest. Cactus environment machine: Shared environment call-by-need. In *Proceedings of the 17th ACM symposium on Trends in Functional Programming*, page to appear. ACM, 2017.
- [16] X. Yang, Y. Chen, E. Eide, and J. Regehr. Finding and understanding bugs in c compilers. In *Proceedings of the 32Nd ACM SIGPLAN Conference on Programming Language Design and Implementation*, PLDI '11, pages 283–294, New York, NY, USA, 2011. ACM. ISBN 978-1-4503-0663-8. doi: 10.1145/1993498.1993532. URL <http://doi.acm.org/10.1145/1993498.1993532>.

ESVERIFY: Verifying Dynamically-Typed Higher-Order Functional Programs by SMT Solving

Christopher Schuster

University of California, Santa Cruz
cschuste@ucsc.edu

Sohum Banerjea

University of California, Santa Cruz
sobanerj@ucsc.edu

Cormac Flanagan

University of California, Santa Cruz
cormac@ucsc.edu

ABSTRACT

Program verifiers statically check correctness properties of programs with annotations such as assertions and pre- and postconditions. Recent advances in SMT solving make it applicable to a wide range of domains, including program verification. In this paper, we describe `ESVERIFY`, a program verifier for JavaScript based on SMT solving, supporting functional correctness properties comparable to languages with refinement and dependent function types. `ESVERIFY` supports both higher-order functions and dynamically-typed idioms, enabling verification of programs that static type systems usually do not support. To verify these programs, we represent functions as universal quantifiers in the SMT logic and function calls as instantiations of these quantifiers. To ensure that the verification process is decidable and predictable, we describe a bounded quantifier instantiation algorithm that prevents matching loops and avoids ad-hoc instantiation heuristics. We also present a formalism and soundness proof of this verification system in the Lean theorem prover and a prototype implementation.

CCS CONCEPTS

- Theory of computation → Pre- and post-conditions; Logic and verification;
- Software and its engineering → Language types; Functional languages;

KEYWORDS

program verification, functional programming, SMT solving

ACM Reference Format:

Christopher Schuster, Sohum Banerjea, and Cormac Flanagan. 2019. `ESVERIFY`: Verifying Dynamically-Typed Higher-Order Functional Programs by SMT Solving. In *Proceedings of International Symposium on Implementation and Application of Functional Languages (IFL'18)*. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

The goal of program verification is to statically check programs for properties such as robustness, security and functional correctness across all possible inputs. For example, program verifiers might statically verify that the result of a sorting routine is both sorted and is a permutation of the input.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

IFL'18, August 2019, Lowell, MA, USA

© 2019 Association for Computing Machinery.
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

In this paper, we present `ESVERIFY`, a program verification system for JavaScript. JavaScript, a dynamically-typed functional scripting language, was chosen as target because it is a relatively simple language compared to other dynamic programming languages, and because its broad user base suggests many beneficial use cases for static analysis.

JavaScript programs often include idioms and patterns that do not adhere to standard typing rules. For instance, the latest edition of the JavaScript/ECMAScript standard [ECMA-262 2017] introduces *promises*: a promise can be composed with any object, as long as that object has a "then" method. Since `ESVERIFY` does not rely on static types, we can easily accommodate these idioms.

JavaScript programs also often use higher-order functions. In order to support verification of these functions, we introduce syntax to constrain function values in terms of their pre- and postconditions. Similarly, other JavaScript values such as numbers, strings, arrays and classes can be used in assertions including invariants on array contents and class instances.

We briefly discuss the prototype implementation, `ESVERIFY`, describe its design and give an overview of the verification process. If a JavaScript program uses features unsupported by `ESVERIFY`, it will be rejected early; otherwise, verification conditions are generated based on the source code and annotations. Each verification condition is transformed according to a quantifier instantiation algorithm and then checked by an SMT solver. For refuted verification conditions, a dynamic test is synthesized based on the counterexample from the SMT solver and executed to improve error reporting with concrete witnesses. The implementation, including source code¹ and a live demo² are available online.

Additionally, we formally define a JavaScript-inspired, statically verified but dynamically typed language, called λ^S . Functions in λ^S are annotated with pre- and postconditions, rendered as logical propositions. These propositions can include operators, refer to variables in scope, denote function results with uninterpreted function calls, and constrain the pre- and postconditions of function values. The verification rules for λ^S involve checking verification conditions for validity. This checking is performed by an SMT solver augmented with decidable theories for linear integer arithmetic, equality, data types and uninterpreted functions. The key difficulty is that verification conditions can include quantifiers, as function definitions in the source program correspond to universally quantified formulas in verification conditions. Unfortunately, SMT solvers may not always perform the right instantiations, and therefore quantifiers imperil the decidability of the verification process [Ge and de Moura 2009; Reynolds et al. 2013]. We ensure that the verification process remains decidable and predictable, by

¹Implementation Source Code: <https://github.com/levjj/esverify>

²Online live demo of `ESVERIFY`: <https://esverify.org/try>

proposing a bounded quantifier instantiation algorithm such that function calls in the source program act as hints (“triggers”) that instantiate these quantifiers. The algorithm only performs a bounded number of trigger-based instantiations and thereby avoids brittle instantiation heuristics and matching loops. Using this decision procedure for verification conditions, we show that verification of λ^S is sound, i.e. verifiable λ^S programs do not get stuck. The proof is formalized in the Lean Theorem Prover and available online³.

To evaluate the expressiveness of this approach, we also include a brief comparison with static refinement types [Freeman and Pfenning 1991]. While a formalization and proof is beyond the scope of this paper, refined base type and dependent function types can be translated to assertions such that the resulting program is verifiable if the original program is well-typed. This suggests that `ESVERIFY` is at least as expressive as a language with refinement types.

To summarize, the main contributions of this paper are

- (1) an approach for verifying dynamically-typed, higher-order JavaScript programs,
- (2) a bounded quantifier instantiation algorithm that enables trigger-based instantiations without heuristics or matching loops,
- (3) a prototype implementation called `ESVERIFY`, and
- (4) a formalization of the verification rules and a proof of soundness in the Lean theorem prover.

The structure of the rest of the paper is as follows: Section 2 illustrates common use cases and relevant features of `ESVERIFY`, Section 3 outlines the verification process and the design of the implementation, Section 4 formally defines quantifier instantiation, as well as the syntax, semantics, verification rules and a soundness theorem for a core language λ^S , Section 5 compares the program verification approach to refinement type systems, Section 6 discusses related work, and finally Section 7 concludes the paper.

2 VERIFYING JAVASCRIPT PROGRAMS

Our program verifier, `ESVERIFY`, targets a subset of ECMAScript-/JavaScript. By supporting a dynamically-typed scripting language, `ESVERIFY` is unlike existing verifiers for statically-typed programming languages. We do not aim to support complex and advanced JavaScript features such as prototypical inheritance and metaprogramming, leaving these extensions for future work. We do aim to support both functional as well as object-oriented programming paradigms, as commonly present in JavaScript. This paper, however, focuses on pure functional JavaScript programs with higher-order functions.

2.1 Annotating JavaScript with Assertions

`ESVERIFY` extends JavaScript with syntax for annotating functions with pre- and postconditions, as well as loop invariants and static assertions. We write these with as *pseudo function calls* with standard Javascript syntax. While some program verification systems specify these in comments [Flanagan et al. 2002], this approach enables a better integration with existing tooling support such as refactoring tools and syntax highlighters.

³Formal Definitions and Proofs in Lean: <https://github.com/levjj/esverify-theory/>

```

1  function max(a, b) {
2    requires(typeof(a) === 'number');
3    requires(typeof(b) === 'number');
4    ensures(res => res >= a);
5    ensures(res => res >= b); // does not hold
6    if (a >= b) {
7      return a;
8    } else {
9      return a; // bug
10   }
11 }
```

Listing 1: A JavaScript function `max` annotated with pre- and postconditions.

```

1  function sumTo (n) {
2    requires(Number.isInteger(n) && n >= 0);
3    ensures(res => res === (n + 1) * n / 2);
4    let i = 0;
5    let s = 0;
6    while (i < n) {
7      invariant(Number.isInteger(i) && i <= n);
8      invariant(Number.isInteger(s));
9      invariant(s === (i + 1) * i / 2);
10     i++;
11     s = s + i;
12   }
13   return s;
14 }
```

Listing 2: A JavaScript function that shows $\sum_{i=0}^n i = \frac{(n+1)\cdot n}{2}$. Loop invariants are not inferred and need to be specified explicitly for all mutable variables in scope.

The assertion language is embedded in JavaScript but does not support all of JavaScript’s semantics. In particular, it is restricted to pure expressions and cannot define new functions.

2.2 `max`: A Simple Example

Listing 1 shows an example of an annotated JavaScript program. The calls to `requires` and `ensures` in lines 2–5 are only used for verification purposes and excluded from evaluation. Instead of introducing custom type annotations, the values of function arguments are constrained using the standard JavaScript `typeof` operator.

Due to a bug in line 9, the `max` function does not return the maximum of the arguments if `b` is greater than `a`, violating the postcondition in line 5.

2.3 Explicit Loop Invariants

For programs without loops or recursion, any correctness property can be checked precisely. However, the potential behavior of programs with loops or recursion cannot be determined statically and is “overapproximated” by `ESVERIFY`. Therefore, correct programs may be rejected if the program lacks a sufficiently strong loop invariant or pre- or postcondition, but verified programs are guaranteed to not violate an assertion regardless of the number

```

1  function inc (x) {
2    requires(Number.isInteger(x));
3    ensures(y => Number.isInteger(y) && y > x);
4    // implicit: ensures(y => y === x + 1);
5    return x + 1;
6  }
7  function twice (f, n) {
8    requires(spec(f, (x) => Number.isInteger(x),
9                  (x,y) => Number.isInteger(y) &&
10                 y > x));
11   requires(Number.isInteger(n));
12   ensures(res => res >= n + 2);
13   return f(f(n));
14 }
15 const n = 3;
16 const m = twice(inc, n); // 'inc' satisfies spec
17 assert(m > 4);          // statically verified

```

Listing 3: The higher-order function `twice` restricts its function argument `f` with a maximum precondition and a minimum postcondition. The function `inc` has its body as implicit postcondition and therefore satisfies this `spec`.

of iterations or function calls. Listing 2 shows a JavaScript function that computes the sum of the first `n` natural numbers with a while loop. The loop requires annotated invariants for mutable variables including their types and bounds⁴. ESVERIFY internally uses standard SMT theorems for integer arithmetic to establish $\sum_{i=0}^n i = \frac{(n+1)\cdot n}{2}$ at each iteration and thereby verify the whole function.

There is extensive prior work on automatically inferring loop invariants [Furia and Meyer 2010]. However, this topic is orthogonal to the program verification approach presented in this paper.

2.4 Higher-order Functions

In order to support function values as arguments and results, ESVERIFY introduces a `spec` construct in pre-, postconditions and assertions. Listing 3 illustrates this syntax in lines 8–10 of the higher-order `twice` function. The argument `f` needs to be a function that satisfies the given constraints, and therefore ESVERIFY must compare the pre- and postconditions of `inc` with this `spec` at the call-site `twice(inc, n)` in line 16. The explicitly stated postcondition of `inc` does not satisfy this `spec`, but ESVERIFY implicitly strengthens the postcondition of `inc` by inlining its function body `n + 1`. It is important to note that recursive functions are only inlined by one level. This means that recursive functions, similarly to loop invariants, need to be explicitly annotated with adequate pre- and postconditions.

2.5 Arrays and Objects

In addition to floating point numbers and integers, ESVERIFY also supports other standard JavaScript values such as boolean values, strings, functions, arrays and objects. However, ESVERIFY restricts how objects and arrays can be used. In particular, mutation of arrays and objects is not currently supported. Additionally, objects

⁴Here, `Number.isInteger(i)` ensures that `i` is an actual integer, while `typeof(i) === 'number'` is also true for floating point numbers.

```

1  function f (a) {
2    requires(a instanceof Array);
3    requires(a.every(e => e > 3));
4    requires(a.length >= 2);
5    assert(a[0] > 2); // holds
6    assert(a[1] > 4); // does not hold
7    assert(a[2] > 0); // does not hold
8  }

```

Listing 4: ESVERIFY includes basic support for immutable arrays. The elements of an array can be described with `every`.

can either be immutable *dictionaries* that map string keys to values or instances of a user-defined class with a fixed set of fields and without inheritance.

The elements of an array can be described with a quantified proposition, corresponding to the standard array method `every`. This is illustrated in Listing 4.

Despite these restrictions, it is possible to express complex recursive data structures. For example, Listing 5 shows a user-defined linked list class that is parameterized by a predicate. Here, the `each` field is actually a function that returns `true` for each element in the linked list. Mapping over the elements of the list with a function `f` requires that `f` can be invoked with elements that satisfy `this.each` and that return values of `f` satisfy the new predicate `newEach`. This demonstrates how generic data structures can be used to verify correctness in a similar way to parameterized types. It is important to note that function calls in an assertion context are uninterpreted, so the call `newEach(y)` in line 19 only refers to the function return value but does not actually invoke the function.

2.6 Dynamic Programming Idioms: Promises

JavaScript programs often include functions that have polymorphic calling conventions. A common example is the jQuery library which provides a function “\$” whose behavior varies greatly depending on the arguments: given a function argument, the function is scheduled for deferred execution, while other argument types find and return portions of the current webpage.

Even the JavaScript standard itself uses dynamic programming idioms to provide a more convenient programming interface. For example, the latest edition of the ECMAScript standard [ECMA-262 2017] includes *Promises* [Liskov and Shriram 1988] and specifies a polymorphic `Promise.resolve()` function. It performs different actions depending on whether it is called with a promise, an arbitrary non-promise object with a method called “`then`”, or any other arbitrary non-promise object. ESVERIFY can accurately express these kinds of specifications in pre- and postconditions as shown in Listing 6, while standard type systems need to resort to code changes (for example, sum types and injections).

2.7 Complex Programs: MergeSort

We also demonstrate non-trivial programs such as MergeSort and verify their functional correctness⁵. The implementation is purely functional and uses a linked list data type that is defined as a class.

⁵The source code of a MergeSort algorithm in ESVERIFY is available at <https://esverify.org/try#msort>.

```

1  class List {
2    constructor (head, tail, each) {
3      this.head = head; this.tail = tail; this.each = each;
4    }
5    invariant () {
6      // this.each is a predicate that is true for this element and the rest of the list
7      return spec(this.each, x => true, (x, y) => pure() && typeof(y) === 'boolean') &&
8        (true && this.each)(this.head) && // same as 'this.each(this.head)' but without binding 'this'
9        (this.tail === null || (this.tail instanceof List && this.each === this.tail.each));
10     }
11   }
12   function map (f, lst, newEach) {
13     // newEach needs to be a predicate
14     // (a pure function without precondition that returns a boolean)
15     requires(spec(newEach, x => true, (x, y) => pure() && typeof(y) === 'boolean'));
16     // the current predicate 'this.each' must satisfy the precondition of 'f'
17     // and the return value of 'f' needs to satisfy the new predicate 'newEach'
18     requires(lst === null || spec(f, x => (true && lst.each)(x), (x, y) => pure() && newEach(y)));
19     requires(lst === null || lst instanceof List);
20     ensures(res => res === null || (res instanceof List && res.each === newEach));
21     ensures(pure()); // necessary as recursive calls could otherwise invalidate the class invariant
22     return lst === null ? null : new List(f(lst.head), map(f, lst.tail, newEach), newEach);
23   }

```

Listing 5: Custom linked list class with a field `each` which is a function that is true for all elements. Mapping over the list results in a new list whose elements satisfy a new predicate analogous to a map function in a parametrized type system.

```

1  class Promise {
2    constructor (value) { this.value = value; }
3  }
4  function resolve (fulfill) {
5    // "fulfill" is promise, thenable or
6    // a value without a "then" property
7    requires(fulfill instanceof Promise ||
8      spec(fulfill.then, () => true,
9           () => true) ||
10        !('then' in fulfill));
11    ensures(res => res instanceof Promise);
12    if (fulfill instanceof Promise) {
13      return fulfill;
14    } else if ('then' in fulfill) {
15      return new Promise(fulfill.then());
16    } else {
17      return new Promise(fulfill);
18    }
19  }

```

Listing 6: The standard `Promise.resolve()` function in JavaScript has complex polymorphic behavior. This simplified mock definition illustrates how ESVERIFY enables such dynamic programming idioms.

Interestingly, about 48 out of a total 99 lines are verification annotations, including invariants, pre- and postconditions and the predicate function `isSorted`. `isSorted` is primarily used in specifications, but the implementations of `merge` and `sort` also include calls to it. These calls are used as *triggers*, hints to the underlying SMT solver that do not contribute to the result. In other verified languages such as Dafny [Leino 2013], `isSorted` would correspond

to a “ghost function”, but ESVERIFY does not currently differentiate between verification-only and regular implementation functions.

2.8 JavaScript as Theorem Prover

A simple induction proof over natural numbers can be written as a while loop as previously shown in Listing 2. This idea can be generalized by using the `spec` construct to reify propositions.

In particular, the postcondition of a function need not only describe its return value; it can also state a proposition such that a value that satisfies the function specification acts as proof of this proposition – analogous to the Curry-Howard isomorphism. Such a “function” can then be supplied as argument to higher-order functions to build up longer proofs. For an example, Listing 7 includes a proof written in JavaScript showing that any locally increasing integer-ranged function is globally increasing. This example was previously used to illustrate refinement reflection in LiquidHaskell [Vazou et al. 2018].

3 IMPLEMENTATION

The sources of the ESVERIFY prototype implementation⁶ are available online. Because the implementation itself is written in TypeScript, a dialect of JavaScript, it can be used in a browser. In fact, we also implemented a simple browser-based editor with ESVERIFY checking. Its source code⁷ as well as a live demo⁸ are available online. Alternative user interfaces such as integrations with Vim and Emacs also exist.

⁶Implementation Source Code: <https://github.com/levjj/esverify/>

⁷Browser-based ESVERIFY Editor Source: <https://github.com/levjj/esverify-web/>

⁸Online live demo of ESVERIFY: <https://esverify.org/try>

```

1  function proof_f_mono (f, proof_f_inc, n, m) {
2    // f is a function from non-negative int to int
3    requires(spec(f,
4      (x) => Number.isInteger(x) && x >= 0,
5      (x, y) => Number.isInteger(y) && pure())));
6    // proof_f_inc states that f is increasing
7    requires(spec(proof_f_inc,
8      x => Number.isInteger(x) && x >= 0,
9      x => f(x) <= f(x + 1) && pure()));
10   requires(Number.isInteger(n) && n >= 0);
11   requires(Number.isInteger(m) && m >= 0);
12   requires(n < m);
13   // show that f is increasing for arbitrary n,m
14   ensures(f(n) <= f(m));
15   ensures(pure()); // no side effects
16   proof_f_inc(n); // instantiate proof for n
17   if (n + 1 < m) {
18     // invoke induction hypothesis (I.H.)
19     proof_f_mono(f, proof_f_inc, n + 1, m);
20   }
21 }
22 function fib (n) {
23   requires(Number.isInteger(n) && n >= 0);
24   ensures(res => Number.isInteger(res));
25   ensures(pure());
26   if (n <= 1) {
27     return 1;
28   } else {
29     return fib(n - 1) + fib(n - 2);
30   }
31 }
32 // A proof that fib is increasing
33 function proof_fib_inc (n) {
34   requires(Number.isInteger(n) && n >= 0);
35   ensures(fib(n) <= fib(n + 1));
36   ensures(pure());
37   fib(n); // unfolds fib at n
38   fib(n + 1);
39   if (n > 0) {
40     fib(n - 1);
41     proof_fib_inc(n - 1); // I.H.
42   }
43   if (n > 1) {
44     fib(n - 2);
45     proof_fib_inc(n - 2); // I.H.
46   }
47 }
48 function proof_fib_mono (n, m) {
49   requires(Number.isInteger(n) && n >= 0);
50   requires(Number.isInteger(m) && m >= 0);
51   requires(n < m);
52   ensures(fib(n) <= fib(m));
53   ensures(pure());
54   proof_f_mono(fib, proof_fib_inc, n, m);
55 }

```

Listing 7: A proof about monotonous integer functions in JavaScript and an instantiation for `fib`. This example was previously used to illustrate refinement reflection in the statically-typed LiquidHaskell system [Vazou et al. 2018].

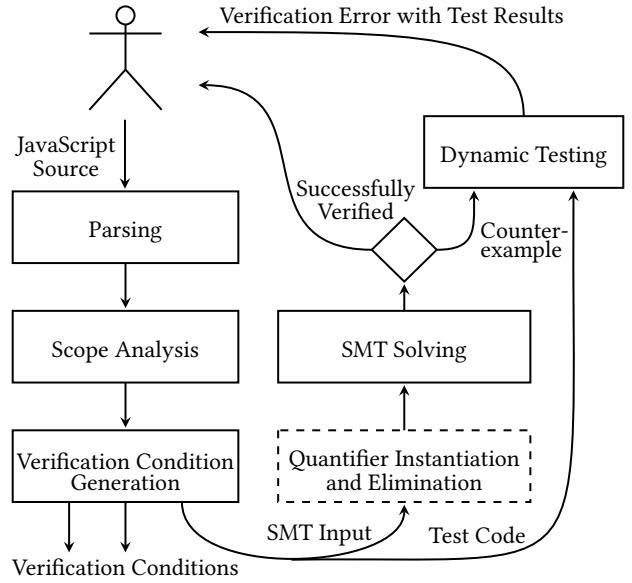


Figure 1: The basic verification workflow: ESVERIFY generates, statically checks and dynamically tests verification conditions based on the provided JavaScript source code.

The basic verification process and overall design of ESVERIFY is depicted in Figure 1.

The first step of the process involves **parsing** the source code and restricting the input language to a subset of JavaScript that is supported by ESVERIFY. Some of these restrictions may be lifted in future versions of ESVERIFY, such as support for regular expressions or a variable number of function arguments. However, some JavaScript features would require immense complexity for accurate verification due to their dynamic character and interactions with the rest of the program, such as metaprogramming with `eval` or `new Function()`, or have been deprecated in newer versions of strict mode JavaScript such as `argumentscallee`, `this` outside of functions or the `with` statement.

The parser also differentiates between the syntax of expressions and assertions. For example, `spec` can only be used in assertions while function definitions can only appear in the actual program implementation.

During the second step, **scope analysis** determines variable scopes and rejects programs with scoping errors or references to unsupported global objects. In addition to user-provided definitions, it includes a whitelist of supported globals such as `Array`, `Math` and `console`. The analysis also takes mutability into account. For example, mutable variables cannot be referenced in class invariants, and the `old(x)` syntax in a postcondition requires `x` to be mutable.

The main **verification** step then traverses the entire source program. At each statement and expression, the current verification context is used to generate verification conditions and augment the context for subsequent statements and expressions. Specifically, the verification context includes a logical proposition that acts as

precondition, a synthesized unit test with *holes*, and a set of variables with unknown values. Generated verification conditions combine this context with an assertion, such as a function postcondition. Section 4.4 describes the verification rules in more detail.

The proposition of a verification condition is then transformed with a **quantifier instantiation** procedure. Quantifiers are instantiated based on triggers and remaining quantifiers erased⁹. The resulting quantifier-free proposition can be checked precisely by SMT solving, ensuring that the verification process remains predictable. However, this approach to quantifier instantiation requires the programmer to provide explicit triggers in terms of function calls. Given this tradeoff, the step is optional and can be skipped.

The next step of the verification process involves checking the verification condition with an **SMT solver** such as z3 [de Moura and Bjørner 2008] or CVC4 [Barrett and Berezin 2004; Barrett et al. 2011]. If the solver cannot refute the proposition, the verification succeeded. Otherwise, the returned model includes an assignment of free variables that acts as a counterexample.

Given a counterexample and a synthesized unit test with holes, the verification condition can be **dynamically tested**. This serves two purposes. On the one hand, the test might not be able reproduce an error or assertion violation. This indicates that the static analysis did not accurately model the actual program behavior due to a loop invariant or assertion not being sufficiently strong. In this case, the programmer can use the variable assignment in the counterexample to better understand the shortcomings of the analysis and improve those annotations. On the other hand, the test might lead to an error or assertion violation. In this case, the programmer is presented with both a verification error message as well as a concrete witness that assists in the debugging process similar to existing test generators [Tillmann and de Halleux 2008].

4 FORMALISM

In order to reason about `ESVERIFY`, this section introduces a formal development of λ^S , a JavaScript-inspired, statically verified but dynamically typed language, and shows that its verification is sound.

The verification rules of λ^S use verification conditions whose validity is checked according to a custom decision procedure, so this section first defines a language for logical propositions and axiomatizes their validity. Then, the decision procedure including quantifier instantiation is given and proven to be sound. Finally, the syntax and semantics of λ^S are defined and its verification rules are shown to be sound.

The definitions, axioms and theorems in this section are also formalized in the Lean theorem prover and available online¹⁰.

4.1 Logical Foundation

Figure 2 formally defines the syntax of propositions, terms, values and environments. Propositions can use terms, connectives \neg , \wedge and \vee , symbols pre_1 , pre_2 , pre , $post$ and universal quantifiers. Here, terms are either values, variables, unary or binary operations or uninterpreted function calls. Finally, values include boolean and integer constants as well as *closures* which are opaque values that will be explained in section 4.3.

⁹The formal definition of the quantifier instantiation algorithm is given in Section 4.2.
¹⁰Formal Definitions and Proofs in Lean: <https://github.com/levjj/esverify-theory/>

$\phi \in \text{Propositions}$	$::= \tau \mid \neg\phi \mid \phi \wedge \phi \mid \phi \vee \phi \mid pre_1(\otimes, \tau) \mid$ $pre_2(\oplus, \tau, \tau) \mid pre(\tau, \tau) \mid post(\tau, \tau) \mid \forall x. \phi$
$\tau \in \text{Terms}$	$::= v \mid x \mid \otimes \tau \mid \tau \oplus \tau \mid \tau(\tau)$
$\otimes \in \text{UnaryOperators}$	$::= \neg \mid isInt \mid isBool \mid isFunc$
$\oplus \in \text{BinaryOperators}$	$::= + \mid - \mid \times \mid / \mid \wedge \mid \vee \mid = \mid <$
$v \in \text{Values}$	$::= \text{true} \mid \text{false} \mid n \mid \langle f(x) \text{ req } R \text{ ens } S \{e\}, \sigma \rangle$
$\sigma \in \text{Environments}$	$::= \emptyset \mid \sigma[x \mapsto v]$
$n \in \mathbb{N}$	$f, x, y, z \in \text{Variables}$

Figure 2: Syntax of logical propositions used in the verifier.

The validity judgement about propositions $\vdash \phi$ is axiomatized as follows.

The term “true” is trivially valid. $\vdash \phi$

Axiom 1. $\vdash \text{true}$.

Validity of conjunctions and disjunctions follows standard inference rules.

Axiom 2. Iff both $\vdash \phi_1$ and $\vdash \phi_2$ then $\vdash \phi_1 \wedge \phi_2$.

Axiom 3. If $\vdash \phi_1$ then $\vdash \phi_1 \vee \phi_2$.

Axiom 4. If $\vdash \phi_2$ then $\vdash \phi_1 \vee \phi_2$.

Axiom 5. If $\vdash \phi_1 \vee \phi_2$ then $\vdash \phi_1$ or $\vdash \phi_2$.

For negation, we assume that valid propositions do not include contradictions and we also assume the law of excluded middle.

Axiom 6. $\vdash \phi \wedge \neg\phi$ is not true.

Axiom 7. $\vdash \phi \vee \neg\phi$.

For convenience, we include a notation for implication based on disjunction and negation.

Notation 1 (Implication). $(\phi_1 \Rightarrow \phi_2) \stackrel{\text{def}}{=} (\neg\phi_1 \vee \phi_2)$.

Instead of defining evaluation of terms, we include a few axioms that suffice for the remaining formal development. For example, we assume that a term is valid if it evaluates to the value “true” and thereby satisfies the following equality.

Axiom 8. Iff $\vdash \tau$ then $\vdash \tau = \text{true}$.

The meaning of unary and binary operators is specified in terms of a partial function δ , e.g. $\delta(+, 2, 3) = 5$.

Axiom 9. Iff $\delta(\otimes, v_x) = v$ then $\vdash v = \otimes v_x$.

Axiom 10. Iff $\delta(\oplus, v_x, v_y) = v$ then $\vdash v = v_x \oplus v_y$.

The propositions $pre_1(\otimes, v_x)$ and $pre_2(\otimes, v_x, v_y)$ can be used to reason about the domain of operators.

Axiom 11. If $\vdash pre_1(\otimes, v_x)$ then $(\otimes, v_x) \in \text{dom}(\delta)$.

Axiom 12. If $\vdash pre_2(\oplus, v_x, v_y)$ then $(\oplus, v_x, v_y) \in \text{dom}(\delta)$.

The constructs $pre(f, x)$ and $post(f, x)$ in propositions denote the pre- and postcondition of a function f when applied to a given argument x but are left uninterpreted for now.

Definition 1 (Substitution). $\phi[x \mapsto v]$ denotes the proposition in which free occurrences of x in ϕ are replaced by v .

A universally quantified proposition is true for all values and can be instantiated with any term.

Axiom 13. If $\vdash \phi[x \mapsto v]$ for all values v , then $\vdash \forall x. \phi$.

Axiom 14. If $\vdash \forall x. \phi$ then $\vdash \phi[x \mapsto \tau]$ for all terms τ .

A valid proposition is not necessarily closed. In fact, a free variable in a valid proposition is assumed to be implicitly universally quantified.

Axiom 15. If x is free in ϕ and $\vdash \phi$ then $\vdash \forall x. \phi$.

An environment σ mapping variables to values can be used as a substitution for terms and propositions.

Definition 2 (Lookup). $\sigma(x)$ looks up x in the environment σ .

Definition 3 (Substitution with Environment). $\sigma(\tau)$ and $\sigma(\phi)$ substitute free variables in τ and ϕ with values according to σ .

An environment σ is a model for a proposition if the substituted proposition is valid. Note that this definition of models is unconventional but it facilitates our subsequent formal development.

Notation 2 (Model). $\sigma \models \phi \stackrel{\text{def}}{=} \vdash \sigma(\phi)$

$\boxed{\sigma \models \phi}$

It is important to note that the validity judgement may not be decidable for all propositions due to the use of quantifiers, so in addition to this (undecidable) validity judgement, we also introduce a notion of satisfiability by an SMT solver.

Definition 4 (Satisfiability). $\text{Sat}(\phi)$ denotes that the SMT solver found a model that satisfies ϕ .

$\boxed{\text{Sat}(\phi)}$

Theorem 1. If ϕ is quantifier-free, then $\text{Sat}(\phi)$ terminates and $\text{Sat}(\phi)$ iff $\sigma \models \phi$ for some model σ .

PROOF. SMT solving is not decidable for arbitrary propositions but the QF-UFLIA fragment of quantifier-free formulas with equality, linear integer arithmetic and uninterpreted function is known to be decidable [Christ et al. 2012; Nelson and Oppen 1979]. \square

4.2 Quantifier Instantiation Algorithm and Decision Procedure

As described in the previous sections, verification of λ^S involves checking the validity of verification conditions that include quantifiers. Quantifier instantiation in SMT solvers is an active research topic [Ge and de Moura 2009; Reynolds et al. 2013], and often requires heuristics or explicit matching triggers. However, heuristics can cause unpredictable results and trigger-based instantiation might lead to matching loops. This section describes a bounded quantifier instantiation algorithm that avoids matching loops and brittle heuristics, thus providing us with a predictable decision procedure for verification conditions.

The syntax of verification conditions (VCs) is shown in Figure 3. VCs used in verification rules are similar to propositions introduced in the previous section. However, universal quantifiers have explicit matching patterns to indicate that instantiation requires a trigger. Accordingly, the construct $\text{call}(x)$ acts as an instantiation

$$\begin{aligned} P \in \text{VerificationConditions} ::= \\ \tau \mid \neg P \mid P \wedge P \mid P \vee P \mid \text{pre}_1(\otimes, \tau) \mid \text{pre}_2(\oplus, \tau, \tau) \mid \\ \text{pre}(\tau, \tau) \mid \text{post}(\tau, \tau) \mid \text{call}(\tau) \mid \forall x. \{\text{call}(x)\} \Rightarrow P \mid \exists x. P \end{aligned}$$

Figure 3: Syntax of verification conditions (VCs).

trigger that does not otherwise affect validity of propositions, i.e. $\text{call}(x)$ can always be assumed to be true. Intuitively, $\text{call}(x)$ represents a function call or an asserted function specification while $\forall x. \{\text{call}(x)\} \Rightarrow P$ corresponds to a function definition or an assumed function specification.

The complete decision procedure for VCs including quantifier instantiation is shown in Figure 4.

To make the definition more concise, we first define contexts $P^+[\circ]$ and $P^-[\circ]$ for a VC with positive and negative polarity regarding negation. Using this definition, the sets of call triggers in positive and negative positions can be defined as triggers for which there exists a context of the right polarity.

The procedure lift^+ matches universal quantifiers in positive and existential quantifiers in negative positions. In both cases, an equivalent VC without the quantifier can be obtained by renaming the quantified variable to a fresh and thereby implicitly universally quantified variable. It is important to note that the matching pattern $\text{call}(y)$ now becomes available to instantiate other quantifiers. The lifting is repeated until no more such quantifiers can be found.

The procedure instantiateOnce^- performs one round of trigger-based instantiation, within which each universal quantifier with negative polarity is instantiated with all triggers in negative position. All such instantiations are conjoined with the original quantifier.

Both lifting and instantiations are repeated for n rounds by the recursive instantiate^- procedure. During the final step, erase^- removes all remaining triggers and quantifiers in negative positions.

The overall decision procedure $\langle P \rangle$ performs n rounds of instantiations where n is the maximum level of quantifier nesting. The original VC P is valid if the resulting proposition cannot be refuted by SMT solving.

VCs P can syntactically include both existential and universal quantifiers in both positive and negative positions. However, we can show that VCs generated by the quantifier have existential quantifiers only in negative positions.

Theorem 2 (Decision Procedure Termination). If P does not contain existential quantifiers in negative positions, the decision procedure $\langle P \rangle$ terminates.

PROOF. The lift^+ function eliminates a quantifier during each recursive call and therefore terminates when there are no more matching quantifiers in the formula. instantiateOnce^- and erase^- are both non-recursive and trivially terminate. Since the maximum level of nesting is finite, $\langle P \rangle$ performs only a finite number of instantiations. With existential quantifiers only in negative positions, the erased and lifted result is quantifier-free, so according to Lemma 1, the final SMT solving step terminates and so does the whole decision procedure. \square

$$\begin{aligned}
P^+[\circ] &::= \circ \mid \neg P^-[\circ] \mid P^+[\circ] \wedge P \mid P \wedge P^+[\circ] \mid P^+[\circ] \vee P \mid P \vee P^+[\circ] & calls^+(P) &\stackrel{\text{def}}{=} \{ call(\tau) \mid P = P^+[call(\tau)] \} \\
P^-[\circ] &::= \neg P^+[\circ] \mid P^-[\circ] \wedge P \mid P \wedge P^-[\circ] \mid P^-[\circ] \vee P \mid P \vee P^-[\circ] & calls^-(P) &\stackrel{\text{def}}{=} \{ call(\tau) \mid P = P^-[call(\tau)] \} \\
lift^+(P) &\stackrel{\text{def}}{=} \begin{array}{ll} \text{match } P \text{ with} \\ P^+[\forall x.\{call(x)\} \Rightarrow P'] \rightarrow lift^+(P^+[call(y) \Rightarrow P'[x \mapsto y]]) \\ P^-[\exists x.P'] \rightarrow lift^+(P^-[P'[x \mapsto y]]) \\ \text{otherwise} \rightarrow P \end{array} & (y \text{ fresh}) \\
instantiateOnce^-(P) &\stackrel{\text{def}}{=} P \left[P^- \left[(\forall x.\{call(x)\} \Rightarrow P') \right] \mapsto P^- \left[(\forall x.\{call(x)\} \Rightarrow P') \wedge \bigwedge_{call(\tau) \in calls^-(P)} P'[x \mapsto \tau] \right] \right] \\
erase^-(P) &\stackrel{\text{def}}{=} P^- \left[(\forall x.\{call(x)\} \Rightarrow P'') \right] \mapsto P^-[\text{true}], P^+ \left[call(\tau) \right] \mapsto P^+[\text{true}], P^- \left[call(\tau) \right] \mapsto P^-[\text{true}] \\
instantiate^-(P, n) &\stackrel{\text{def}}{=} \text{if } (n > 0) \text{ then } erase^-(lift^+(P)) \text{ else } instantiate^-(instantiateOnce^-(lift^+(P)), n - 1) \\
\langle P \rangle &\stackrel{\text{def}}{=} \text{let } n = \text{maximum level of quantifier nesting of } P \text{ in } \neg Sat(\neg instantiate^-(P, n)) \quad \boxed{\langle P \rangle}
\end{aligned}$$

Figure 4: The decision procedure lifts, instantiates and finally eliminates quantifiers. The number of iterations is bounded by the maximum level of quantifier nesting.

$$\begin{aligned}
e \in \text{Expressions} &::= \text{let } x = \text{true in } e \mid \text{let } x = \text{false in } e \mid \text{let } x = n \text{ in } e \mid \\
&\quad \text{let } f(x) \text{ req } R \text{ ens } S = e \text{ in } e \mid \text{let } y = \otimes x \text{ in } e \mid \\
&\quad \text{let } z = x \oplus y \text{ in } e \mid \text{let } y = f(x) \text{ in } e \mid \text{if } (x) e \text{ else } e \mid \text{return } x \\
R, S \in \text{Specs} &::= \tau \mid \neg R \mid R \wedge R \mid R \vee R \mid \text{spec } \tau(x) \text{ req } R \text{ ens } S \\
\kappa \in \text{Stacks} &::= (\sigma, e) \mid \kappa \cdot (\sigma, \text{let } y = f(x) \text{ in } e)
\end{aligned}$$

Figure 5: Syntax of λ^S programs. Function definitions have pre- and postconditions written as simple logical propositions with the `spec` syntax for higher-order functions.

Definition 5 (Proposition Translation). $prop(P)$ denotes a proposition such that triggers and matching patterns in P are removed and existential quantifiers $\exists x. P$ translated to $\neg \forall x. \neg prop(P)$.

Theorem 3 (Quantifier Instantiation Soundness). If P has no existential quantifiers in negative positions, then $\langle P \rangle$ implies $\vdash prop(P)$.

PROOF. By Axiom 15, $list^+$ preserves equisatisfiability. Note that any conjuncts inserted by $instantiateOnce^-$ could also be obtained via classical (not trigger-based) instantiation. Furthermore, since $erase^-$ only removes quantifiers in negative positions and (inconsequential) triggers, the resulting propositions are implied by the original non-erased VC. Finally, with existential quantifiers only in negative positions, the erased and lifted result is quantifier-free. Therefore, Axiom 1 can be used to show that a valid VC according to the decision procedure is also valid without trigger-based instantiation¹¹. \square

4.3 Syntax and Operational Semantics of λ^S

Figure 5 defines the syntax of λ^S . Programs are assumed to be in A-normal form [Flanagan et al. 1993] and the dynamic semantics uses

¹¹A complete proof is available at: <https://github.com/levjj/esverify-theory/>

environments and stack configurations. This formalism avoids substitution in expressions and assertions in order to simplify subsequent proofs.

Here, a function definition $\text{let } f(x) \text{ req } R \text{ ens } S = e_1$ in e_2 is annotated with a precondition R and a postcondition S . These specifications can include terms τ such as constants, program variables, and uninterpreted function application $\tau(\tau)$ as well as logical connectives, and a special syntax “`spec` $\tau(x)$ req R ens S ” that describes the pre- and postcondition of a function value.

As an example, the JavaScript function `inc(n)` in Listing 3 could be expressed in λ^S as follows:

$\text{let } inc(x) \text{ req } \text{isInt}(x) \text{ ens } (\text{isInt}(inc(x)) \wedge inc(x) > x) = x + 1 \text{ in } ..$

Since λ^S is a pure functional language, the postcondition can refer to the function result with an uninterpreted invocation such as `inc(x)` instead of introducing additional variables (as in `ESVERIFY`).

Values in λ^S are either constants, such as integers and boolean literals, or closures consisting of a function definition and an environment σ (see Figure 2).

The operational semantics of λ^S are specified by a small-step evaluation relation over stack configurations κ , as shown in Figure 6. Most noteworthy, the callee function name is added to the environment at each call to enable recursion, and function pre- and postconditions are not checked or enforced during evaluation.

The evaluation of a stack configuration terminates either by getting stuck or by reaching a successful finishing configuration.

Definition 6 (Evaluation Finished). If $\text{isFinished}(\kappa)$ then there exists σ and x such that $\kappa = (\sigma, \text{return } x)$ and $x \in \sigma$.

4.4 Program Verification

The actual verification of λ^S expressions is defined in terms of the verification judgement $P \vdash e : Q$ shown in Figure 7. For simplicity, we adopt the convention that the free variables of P , denoted $FV(P)$, must be exactly the set of variables in scope at e .

$(\sigma, \text{let } y = v \text{ in } e) \hookrightarrow (\sigma[x \mapsto v], e)$	where $v \in \{\text{true}, \text{false}, n\}$	[E-VAL]
$(\sigma, \text{let } f(x) \text{ req } R \text{ ens } S = e_1 \text{ in } e_2) \hookrightarrow (\sigma[f \mapsto \langle f(x) \text{ req } R \text{ ens } S \{e_1\}, \sigma\rangle], e_2)$		[E-CLOSURE]
$(\sigma, \text{let } y = \otimes x \text{ in } e) \hookrightarrow (\sigma[y \mapsto v], e)$	where $v = \delta(\otimes, \sigma(x))$	[E-UNOP]
$(\sigma, \text{let } z = x \oplus y \text{ in } e) \hookrightarrow (\sigma[z \mapsto v], e)$	where $v = \delta(\oplus, \sigma(x), \sigma(y))$	[E-BINOP]
$(\sigma, \text{let } z = f(y) \text{ in } e) \hookrightarrow (\sigma_f[g \mapsto \sigma(f), x \mapsto \sigma(y)], e_f) \cdot (\sigma, \text{let } z = f(y) \text{ in } e)$	where $\sigma(f) = \langle g(x) \text{ req } R \text{ ens } S \{e_f\}, \sigma_f \rangle$	[E-CALL]
$(\sigma, \text{return } z) \cdot (\sigma_2, \text{let } y = f(x) \text{ in } e_2) \hookrightarrow (\sigma_2[y \mapsto \sigma(z)], e_2)$		[E-RETURN]
$(\sigma, \text{if } (x) e_1 \text{ else } e_2) \hookrightarrow (\sigma, e_1)$	if $\sigma(x) = \text{true}$	[E-IF-TRUE]
$(\sigma, \text{if } (x) e_1 \text{ else } e_2) \hookrightarrow (\sigma, e_2)$	if $\sigma(x) = \text{false}$	[E-IF-FALSE]
$\kappa \cdot (\sigma, \text{let } y = f(x) \text{ in } e) \hookrightarrow \kappa' \cdot (\sigma, \text{let } y = f(x) \text{ in } e)$	if $\kappa \hookrightarrow \kappa'$	[E-CONTEXT]

Figure 6: Operational semantics

Given a known precondition P and a expression e , a verification rule checks potential VCs and generates a postcondition Q . This postcondition contains a hole \bullet for the evaluation result of e . Since λ^S is purely functional, P still holds after evaluating e , so we call Q the *marginal postcondition*. The *strongest postcondition* is actually $P \wedge Q[\bullet]$.

As an example, a unary operation such as $\text{let } y = \otimes x \text{ in } e$ is verified with the rule vc-unop. It requires x to be a variable in scope, i.e. a variable that is free in the precondition P . To avoid name clashes, the result y should not be free. Additionally, the VC $\langle P \Rightarrow \text{pre}(\otimes, x) \rangle$ needs to be valid for all free variables (such as x). This check ensures that the value of x is in the domain of the operator \otimes . The rules vc-binop, vc-if, etc. follow analogously.

For function applications $f(x)$, an additional *call*(x) trigger is assumed in order to instantiate quantified formulas that correspond to the function definition or specification of the callee.

The most complex rule concerns the verification of function definitions, such as $\text{let } f(x) \text{ req } R \text{ ens } S = e_1 \text{ in } e_2$. Here, the precondition R , the specification of f and the marginal postcondition $Q_1[f(x)]$ together have to imply the annotated postcondition S . Any recursive calls of f appearing in its function body will instantiate its own (non-recursive) specification, while subsequent calls of f in e_2 will use a postcondition that is strengthened by the generated marginal postcondition. This corresponds to expanding or inlining the function definition by one level at each non-recursive callsite.

The special syntax $\text{spec } \tau(x) \text{ req } R \text{ ens } S$, as used in verification rules, user-provided pre- and postconditions, is a notation that desugars to a universal quantifier.

Notation 3 (Function Specifications). $\text{spec } \tau(x) \text{ req } R \text{ ens } S \equiv \text{isFunc}(\tau) \wedge \forall x. \{ \text{call}(x) \} \Rightarrow ((R \Rightarrow \text{pre}(\tau, x)) \wedge (\text{post}(\tau, x) \Rightarrow S))$

That is, if a function call instantiates this quantifier, the precondition R of the *spec* satisfies the precondition of f and the postcondition S of the *spec* is implied by the postcondition of f . For a concrete function call, this means that R needs to be asserted by the calling context and S can be assumed at the callsite.

4.5 Soundness

Given the decision procedure and the verification rules described in the previous sections, it is possible to show that verified programs can be evaluated to completion without getting stuck.

First, it is important to note that quantifiers in VCs only appear in certain positions.

Lemma 1. If P is a proposition with existential quantifiers only in positive positions, then each VC P' used in the derivation tree of $\vdash e : Q$ has existential quantifiers only in negative positions.

PROOF. All VCs in the verification rules shown in Figures 7 are implications $\langle P'' \Rightarrow Q'' \rangle$. In each of these implications, there are no existential quantifiers in Q'' , as user-supplied postconditions S have no existential quantifiers. Additionally, all propositions P'' on the left-hand side have existential quantifiers only in positive positions, since newly introduced existential quantifiers in marginal postconditions are always in positive positions. \square

From Lemma 1 and Theorem 2, it follows that verification always terminates, ensuring a predictable verification process.

As mentioned in section 4.1, the axiomatization of logical proposition does not include evaluation and treats terms $\tau(\tau)$ as uninterpreted function calls. However, for a proof of verification soundness, if it can be shown that a given closure applied to an argument value evaluates to a value, it will be necessary to establish facts about this function application.

Axiom 16. If $(\sigma[f \mapsto \langle f(x) \text{ req } R \text{ ens } S \{e\}, \sigma\rangle, x \mapsto v_x], e) \hookrightarrow^* (\sigma', y)$ and $\sigma'(y) = v$ then $\vdash \langle f(x) \text{ req } R \text{ ens } S \{e\}, \sigma\rangle(v_x) = v$.

Similarly, axioms about *pre*(f, x) and *post*(f, x) can be added for concrete values of f and x .

Axiom 17. Iff $\sigma[f \mapsto \langle f(x) \text{ req } R \text{ ens } S \{e\}, \sigma\rangle, x \mapsto v_x] \models R$ then $\vdash \text{pre}(\langle f(x) \text{ req } R \text{ ens } S \{e\}, \sigma\rangle, v_x)$.

Axiom 18. If $\vdash \sigma : Q_1$ and $Q_1 \wedge \text{spec } f(x) \text{ req } R \text{ ens } S \wedge R \vdash e : Q_2[\bullet]$ and $\sigma[f \mapsto \langle f(x) \text{ req } R \text{ ens } S \{e\}, \sigma\rangle, x \mapsto v_x] \models Q_2[f(x)] \wedge S$ then $\vdash \text{post}(\langle f(x) \text{ req } R \text{ ens } S \{e\}, \sigma\rangle, v_x)$.

Axiom 19. If $\vdash \sigma : Q_1$ and $Q_1 \wedge \text{spec } f(x) \text{ req } R \text{ ens } S \wedge R \vdash e : Q_2[\bullet]$ and $\vdash \text{post}(\langle f(x) \text{ req } R \text{ ens } S \{e\}, \sigma\rangle, v_x)$ then $\sigma[f \mapsto \langle f(x) \text{ req } R \text{ ens } S \{e\}, \sigma\rangle, x \mapsto v_x] \models Q_2[f(x)] \wedge S$

We can now show that verifiable expressions according to the rules in Figure 7 can be evaluated without getting stuck, i.e. all reachable configurations finish normally or can be evaluated further.

Theorem 4 (Verification Safety). If $\text{true} \vdash e : Q$ and $(\emptyset, e) \hookrightarrow^* \kappa$ then $\text{isFinished}(\kappa)$ or $\kappa \hookrightarrow \kappa'$ for some κ' .

$$\begin{array}{lcl}
Q[\bullet] \in \text{PropositionContexts} & ::= & P \mid \eta[\bullet] \mid \neg Q[\bullet] \mid Q[\bullet] \wedge Q[\bullet] \mid Q[\bullet] \vee Q[\bullet] \mid \text{pre}_1(\otimes, \eta[\bullet]) \mid \text{pre}_2(\oplus, \eta[\bullet], \eta[\bullet]) \mid \\
& & \text{pre}(\eta[\bullet], \eta[\bullet]) \mid \text{post}(\eta[\bullet], \eta[\bullet]) \mid \text{call}(\eta[\bullet]) \mid \forall x. \{\text{call}(x)\} \Rightarrow Q[\bullet] \mid \exists x. Q[\bullet] \\
\eta[\bullet] \in \text{TermContexts} & ::= & \bullet \mid \tau \mid \otimes \eta[\bullet] \mid \eta[\bullet] \oplus \eta[\bullet] \mid \eta[\bullet](\eta[\bullet])
\end{array}
\quad \boxed{P \vdash e : Q}$$

$$\frac{x \notin FV(P) \quad v \in \{\text{true}, \text{false}, n\} \quad P \wedge x = v \vdash e : Q}{P \vdash \text{let } x = v \text{ in } e : \exists x. x = v \wedge Q} \text{VC-VAL} \quad \frac{x \in FV(P) \quad y \notin FV(P) \quad \langle P \Rightarrow \text{pre}(\otimes, x) \rangle \quad P \wedge y = \otimes x \vdash e : Q}{P \vdash \text{let } y = \otimes x \text{ in } e : \exists y. y = \otimes x \wedge Q} \text{VC-UNOP}$$

$$\frac{x \in FV(P) \quad y \in FV(P) \quad z \notin FV(P) \quad \langle P \Rightarrow \text{pre}(\oplus, x, y) \rangle \quad P \wedge z = x \oplus y \vdash e : Q}{P \vdash \text{let } z = x \oplus y \text{ in } e : \exists z. z = x \oplus y \wedge Q} \text{VC-BINOP}$$

$$\frac{f \notin FV(P) \quad x \notin FV(P) \quad f \neq x \quad x \in FV(R) \quad FV(R) \subseteq FV(P) \cup \{f, x\} \quad FV(S) \subseteq FV(P) \cup \{f, x\} \quad P \wedge \text{spec } f(x) \text{ req } R \text{ ens } S \wedge R \vdash e_1 : Q_1 \quad \langle P \wedge \text{spec } f(x) \text{ req } R \text{ ens } S \wedge R \wedge Q_1[f(x)] \Rightarrow S \rangle \quad P \wedge \text{spec } f(x) \text{ req } R \text{ ens } (Q_1[f(x)] \wedge S) \vdash e_2 : Q_2}{P \vdash \text{let } f(x) \text{ req } R \text{ ens } S = e_1 \text{ in } e_2 : \exists f. \text{spec } f(x) \text{ req } R \text{ ens } (Q_1[f(x)] \wedge S) \wedge Q_2} \text{VC-FUNC}$$

$$\frac{f \in FV(P) \quad x \in FV(P) \quad y \notin FV(P) \quad \langle P \wedge \text{call}(x) \Rightarrow \text{isFunc}(f) \wedge \text{pre}(f, x) \rangle \quad P \wedge \text{call}(x) \wedge \text{post}(x) \wedge y = f(x) \vdash e : Q}{P \vdash \text{let } y = f(x) \text{ in } e : \exists y. \text{call}(x) \wedge \text{post}(f, x) \wedge y = f(x) \wedge Q} \text{VC-APP}$$

$$\frac{x \in FV(P) \quad \langle P \Rightarrow \text{isBool}(f) \rangle \quad P \wedge x \vdash e_1 : Q_1 \quad P \wedge \neg x \vdash e_2 : Q_2}{P \vdash \text{if } (x) e_1 \text{ else } e_2 : (x \Rightarrow Q_1) \wedge (\neg x \Rightarrow Q_2)} \text{VC-ITE} \quad \frac{x \in FV(P)}{P \vdash \text{return } x : x = \bullet} \text{VC-RETURN}$$

Figure 7: The judgement $P \vdash e : Q$ verifies the expression e given a known context P .

PROOF. Due to the complex quantifier instantiation of the decision procedure, we first show verification safety for a similar but undecidable verification judgement without quantifier instantiation. With theorem 3, soundness of this verification judgement also implies soundness with trigger-based quantifier instantiation. The verification safety proof for this alternate judgement uses a standard progress/preservation proof strategy, where the notion of verifiability is extended to stack configurations. Here, a given runtime stack is considered verifiable if, at each stack frame, the expression is verifiable with σ as precondition¹². \square

4.6 Extensions

The core language λ^S includes higher-order functions but does not address other language features supported by EVERIFY, such as imperative programs and complex recursive data types.

Extending λ^S for imperative programs would entail syntax, semantics and verification rules for allocating, mutating and referencing values stored in the heap. Most noteworthy, loops and recursion invalidate previous facts about heap contents and therefore require precise invariants. This issue can be addressed with segmentation logic and dynamic frames [Smans et al. 2009].

Additionally, λ^S can be extended to support “classes” as shown in Figure 8. These classes are immutable and more akin to recursive data types as they do not support inheritance. Each class definition consists of an ordered sequence of fields and an invariant S that is specified in terms of a free variable $this$. The class invariant can be used to express complex recursive data structures such as the parameterized linked list shown in Section 2.5.

$$\begin{array}{lll}
c \in \text{ClassNames} & fd \in \text{FieldNames} & this \in \text{Variables} \\
v \in \text{Values} & ::= \dots \mid C(\bar{v}) \\
e \in \text{Expressions} & ::= \dots \mid \text{let } y = \text{new } C(\bar{x}) \text{ in } e \mid \text{let } y = x.fd \text{ in } e \\
\tau \in \text{Terms} & ::= \dots \mid \tau.fd \mid C(\bar{\tau}) \\
P \in \text{VCs} & ::= \dots \mid fd \text{ in } \tau \mid \tau \text{ instanceof } C \mid \\
& \quad \text{access}(\tau) \mid \forall x. \{\text{access}(x)\} \Rightarrow P \\
D \in \text{ClassDefs} & ::= \text{class } C(\bar{fd}) \text{ inv } S
\end{array}$$

$$\frac{x \in FV(P) \quad y \notin FV(P) \quad \langle P \Rightarrow fd \text{ in } x \rangle \quad P \wedge y = x.fd \vdash e : Q}{P \vdash \text{let } y = x.fd \text{ in } e : \exists y. y = x.fd \wedge Q} \text{VC-CLASS-INV}$$

$$\frac{x \in FV(P) \quad y \notin FV(P) \quad \text{class } C(\bar{fd}) \text{ inv } S \in \bar{D} \quad \langle P \wedge this = C(\bar{x}) \wedge this \text{ instanceof } C \Rightarrow S \rangle \quad P \wedge y = C(\bar{x}) \wedge y \text{ instanceof } C \vdash e : Q}{P \vdash \text{let } y = \text{new } C(\bar{x}) \text{ in } e : \exists y. y = C(\bar{x}) \wedge y \text{ instanceof } C \wedge Q} \text{VC-CLASS-DEF}$$

Figure 8: Extending the verification rules of λ^S with simple immutable classes with class invariants.

Similarly to function calls, a trigger $\text{access}(x)$ is inserted into verification conditions at each field access to instantiate the class invariant. However, unlike function definitions, class definitions \bar{D} are global. Therefore, we augment verification conditions such that for each class $C(\bar{fd})$ inv $S \in \bar{D}$ the following quantifier is

¹²A complete proof is available at: <https://github.com/levjj/esverify-theory/>

$t \in \text{TypedExpressions}$	$::= \dots \mid \text{let } f(x : T) : T = t \text{ in } t$
$T \in \text{Types}$	$::= \{x : B \mid R\} \mid x : T \rightarrow T$
$B \in \text{BaseTypes}$	$::= \text{Bool} \mid \text{Int}$
$\Gamma \in \text{TypeEnvironments}$	$::= \emptyset \mid \Gamma, x : T$
	$\boxed{\Gamma \vdash t : T}$
$x, f \notin \text{dom}(\Gamma)$	$FV(T_x) \subseteq \text{dom}(\Gamma)$
$\Gamma, x : T_x, f : (x : T_x \rightarrow T) \vdash t_1 : T_1$	$\Gamma, x : T_x \vdash T_1 <: T$
$\Gamma, x : T_x, f : (x : T_x \rightarrow T) \vdash t_2 : T_2$	
	$\frac{}{\Gamma \vdash \text{let } f(x : T) : T = t_1 \text{ in } t_2 : T_2} \text{-FN}$
$B_1 = B_2$	$\langle \llbracket \Gamma \rrbracket \wedge R \Rightarrow S \rangle$
	$x \notin \text{dom}(\Gamma)$
	$\frac{}{\Gamma \vdash \{x : B_1 \mid R\} <: \{x : B_2 \mid S\}} \text{-ST-REF}$
$\Gamma \vdash T'_x <: T_x$	$\Gamma, x : T'_x \vdash T <: T'$
	$\frac{}{\Gamma \vdash (x : T_x \rightarrow T) <: (x : T'_x \rightarrow T')} \text{-ST-FUN}$
$\llbracket \emptyset \rrbracket \stackrel{\text{def}}{=} \text{true}$	$\boxed{\llbracket \Gamma \rrbracket}$
$\llbracket \Gamma, x : T \rrbracket \stackrel{\text{def}}{=} \llbracket \Gamma \rrbracket \wedge \llbracket x : T \rrbracket$	
$\llbracket \tau : \{x : \text{Bool} \mid R\} \rrbracket \stackrel{\text{def}}{=} \text{isBool}(\tau) \wedge R[x \mapsto \tau]$	$\boxed{\llbracket \tau : T \rrbracket}$
$\llbracket \tau : \{x : \text{Int} \mid R\} \rrbracket \stackrel{\text{def}}{=} \text{isInt}(\tau) \wedge R[x \mapsto \tau]$	
$\llbracket \tau : (x : T_x \rightarrow T) \rrbracket \stackrel{\text{def}}{=} \text{spec } \tau(x) \text{ req } \llbracket x : T_x \rrbracket \text{ ens } \llbracket \tau(x) : T \rrbracket$	

Figure 9: Selected typing and subtyping rules of a statically typed language λ^T . Functions are annotated with types where refinements R are analogous to specifications in λ^S .

assumed:

$$\forall x. \{access(x)\} \Rightarrow (x \text{ instanceof } C \Rightarrow (\overline{x \text{ has } fd} \wedge S[\text{this} \mapsto x]))$$

This quantifier is instantiated by an $access(x)$ trigger and the instantiated formula includes both the class invariant as well as a description of its fields.

5 COMPARISON WITH REFINEMENT TYPES

Despite being dynamically typed, the verification rules shown in Figure 7 resemble static typing rules. In this section, we provide a brief comparison of this program verification approach with static type checking.

An comprehensive formalization of refinement and dependent type systems and a formal proof that examines their expressiveness is beyond the scope of this paper. However, by describing a translation of types to assertions and investigating concrete examples, we enable a comparison of ESVERIFY with systems such as LiquidHaskell [Vazou et al. 2014] and conjecture that it is at least as expressive.

First, we assume a language λ^T similar to λ^S but with type annotations instead of pre- and postconditions. Figure 9 shows an excerpt of such a language. Here, a type is either a dependent function type or a refined base type where refinements R are consistent with specifications in λ^S .

Given such a language, the typing rule for function definitions (t-FN) checks the function body t_1 and compares its type T_1 with

the annotated return type T . This return type might refer to the function argument in order to support dependent types. However, other free variables in refinements can break hygiene, so t-FN restricts free variables in user-provided types T_x and T accordingly.

Figure 9 also shows the subtyping relation. Most relevantly, subtyping of refined base types requires checking an implication between the refinements. Therefore, the typing rules also involve a translation of the type environment Γ to a logical formula, where refinements become propositions and function types translate to the $spec$ syntax.

Intuitively, the implication used for refinements also extends to translated function types, so if $\Gamma \vdash T <: T'$ then for all terms τ , $\llbracket \Gamma \rrbracket \wedge \llbracket \tau : T \rrbracket$ implies $\llbracket \tau : T' \rrbracket$.

As an example, the following λ^T expression is well-typed as the return type is a subtype of the argument type:

$$\text{let } f(g : (x : \{x : \text{Int} \mid x > 3\} \rightarrow \{y : \text{Int} \mid y > x\})) : \\ (x : \{x : \text{Int} \mid x > 4\} \rightarrow \{y : \text{Int} \mid y \geq x\}) = g \text{ in } \dots$$

Translated into ESVERIFY, we obtain a program with **spec** in pre- and postcondition:

```
function f (g) {
  requires(spec(g, x => x > 3, (x, y) => y > x));
  ensures(r =>
    spec(r, x => x > 4, (x, y) => y >= x));
  return g;
}
```

This program is verifiable with the quantifier instantiation algorithm described in section 4.2. The second **spec** is translated to a universal quantifier in positive position that will be lifted, introducing a free global variable x . This also exposes a $call(x)$ trigger in negative position that now instantiates the quantifier in the antecedent. The resulting proposition can now be checked without further instantiations by comparing the argument and return propositions of both functions for all possible values of x .

Extended to expressions, this suggests that the translation of λ^T to λ^S programs preserves verifiability, such that all well-typed λ^T programs translate to verifiable λ^S programs.

Conjecture 1 (Translated well-typed expressions are verifiable). If $\llbracket t \rrbracket$ is the translation of a λ^T expression t to λ^S , then $\Gamma \vdash t : T$ implies $\llbracket \Gamma \rrbracket \vdash \llbracket t \rrbracket : Q$ for some Q .

A formal proof of this conjecture needs to take quantifier instantiation and the quantifier nesting bound into account. This is beyond the scope of this paper, but might be addressed in future work.

Coincidentally, such a translation would also enable seamless interleaving of statically-typed λ^T expressions in dynamically-typed λ^S programs in a sound way. This might be a step towards a full spectrum type system.

6 RELATED WORK

There have been decades of prior work on software verification. In particular, static verification of general purpose programming languages based on pre- and postconditions has previously been explored in verifiers such as ESC/Java [Flanagan et al. 2002; Leino

2001], JaVerT [Fragoso Santos et al. 2018], Dafny [Leino 2010, 2013, 2017] and LiquidHaskell [Vazou et al. 2017, 2014, 2018].

ESC/Java [Flanagan et al. 2002] proposed the idea of using undecidable but SMT-solvable logic to provide more powerful static checking than traditional type systems. Their proposed extended static checking gave up on soundness to do so, citing the utility of tools that on balance contributed to finding bugs.

JaVert [Fragoso Santos et al. 2018] is a more recent program verifier for JavaScript. It supports object-oriented programs but, in contrast to ESVERIFY, does not support higher-order functions.

Dafny [Leino 2010] pushed this further by seeking to provide a full verification language, with support for both functional and imperative programming. Dafny provides programmers with advanced constructs for verified programming, such as ghost functions and parameters, termination checking, quantifiers, and reasoning about the heap. However, Dafny also requires function calls in an assertion contexts to satisfy the precondition instead of treating these as uninterpreted calls. Therefore, Dafny does not support higher-order proofs such as those shown in Section 2.8. Additionally, quantifier instantiation in Dafny is often implicit and based on heuristics, which allows for brittle and unpredictable verification.

In trying to find a compromise, with predictable checking but also a larger scope than traditional type systems, LiquidHaskell is most closely related to ESVERIFY. In fact, the refinement type system discussed in Section 5 loosely resembles its formalization by Vazou et al. [Vazou et al. 2014].

More recently [Vazou et al. 2018], LiquidHaskell introduced *refinement reflection*, which enables external proofs in a similar way as the `spec` construct in ESVERIFY, and *proof by logical evaluation* which is a close cousin to our quantifier instantiation algorithm in Section 4.2 (though it is not based on triggering matching patterns). In contrast to LiquidHaskell, ESVERIFY is not based on static type checking and thus also supports dynamically-typed programming idioms such as dynamic type checks instead of injections.

Finally, the decision procedure described in Section 4.2 involves trigger-based quantifier instantiation which has been studied by extensive prior work [Dross et al. 2016; Ge and de Moura 2009; Leino and Pit-Claudel 2016; Reynolds et al. 2013]. The instantiation in ESVERIFY is specifically bounded in order to prevent matching loops, but further research could provide this kind of instantiation as a built-in feature of off-the-shelf SMT solvers.

7 CONCLUSION

This paper introduced ESVERIFY, a program verifier for dynamically-typed JavaScript programs. ESVERIFY supports both dynamic programming idioms as well as higher-order functional programs, and thus has an expressiveness comparable to and potentially greater than common refinement type systems.

Internally, the verifier relies on a bounded quantifier instantiation algorithm, thus avoiding brittle heuristics, and SMT solving, yielding concrete counterexamples for verification errors.

We showed that this approach to program verification is sound by formalizing the quantifier instantiation algorithm and the verification rules in the Lean Theorem prover.

While ESVERIFY enables verification of non-trivial programs such as MergeSort, it lacks termination checking and support for object-oriented programming. However, it would be possible to combine it with an external termination checker for total correctness [Sereni and Jones 2005], and to extend it with reasoning about the heap, such as dynamic frames [Smans et al. 2009].

Finally, while the approach presented in this paper is purely static, future work might use runtime checks similar to hybrid and gradual type checking [Ahmed et al. 2011; Knowles and Flanagan 2010; Siek and Taha 2006] to enable sound execution of programs that are not fully verified.

REFERENCES

- Amal Ahmed, Robert Bruce Findler, Jeremy G. Siek, and Philip Wadler. 2011. Blame for All. In *POPL '11*.
- Clark Barrett and Sergey Berezin. 2004. CVC Lite: A New Implementation of the Cooperating Validity Checker. In *CAV'04*.
- Clark Barrett, Christopher L. Conway, Morgan Deters, Liana Hadarean, Dejan Jovanović, Tim King, Andrew Reynolds, and Cesare Tinelli. 2011. CVC4. In *CAV'11*.
- Jürgen Christ, Jochen Hoenicke, and Alexander Nutz. 2012. SMTInterpol: An Interpolating SMT Solver. In *SPIN'12*.
- Leonardo de Moura and Nikolaj Bjørner. 2008. Z3: An Efficient SMT Solver. In *TACAS'08*.
- Claire Dross, Sylvain Conchon, Johannes Kanig, and Andrei Paskevich. 2016. Adding decision procedures to SMT solvers using axioms with triggers. *Journal of Automated Reasoning* (2016).
- ECMA-262. 2017. *ECMAScript 2017 Language Specification* (6 / 2017 ed.).
- Cormac Flanagan, K. Rustan M. Leino, Mark Lillibridge, Greg Nelson, James B. Saxe, and Raymie Stata. 2002. Extended Static Checking for Java. In *PLDI'02*.
- Cormac Flanagan, Amr Sabry, Bruce F. Duba, and Matthias Felleisen. 1993. The Essence of Compiling with Continuations. In *PLDI'93*.
- José Fragoso Santos, Petar Maksimović, Daiva Naudžiūnienė, Thomas Wood, and Philippa Gardner. 2018. JaVerT: JavaScript verification toolchain. *POPL'18* (2018).
- Tim Freeman and Frank Pfenning. 1991. Refinement Types for ML. In *PLDI '91*.
- Carlo Alberto Furia and Bertrand Meyer. 2010. Inferring Loop Invariants Using Post-conditions. In *Fields of Logic and Computation*.
- Yeting Ge and Leonardo de Moura. 2009. Complete Instantiation for Quantified Formulas in Satisfiability Modulo Theories. In *CAV'09*.
- Kenneth Knowles and Cormac Flanagan. 2010. Hybrid Type Checking. *TOPLAS* (2010).
- K. Rustan M. Leino. 2001. Extended Static Checking: A Ten-Year Perspective. In *Informatics - 10 Years Back. 10 Years Ahead*.
- K. Rustan M. Leino. 2010. Dafny: An Automatic Program Verifier for Functional Correctness. In *Logic for Programming, Artificial Intelligence, and Reasoning (LPAR'10)*.
- K. Rustan M. Leino. 2013. Developing Verified Programs with Dafny. In *International Conference on Software Engineering, (ICSE'13)*.
- K. Rustan M. Leino. 2017. Accessible Software Verification with Dafny. *IEEE Software* (2017).
- K. Rustan M. Leino and Clément Pit-Claudel. 2016. Trigger selection strategies to stabilize program verifiers. In *CAV'16*.
- B. Liskov and L. Shriram. 1988. Promises: Linguistic Support for Efficient Asynchronous Procedure Calls in Distributed Systems. In *PLDI'88*.
- Greg Nelson and Derek C. Oppen. 1979. Simplification by Cooperating Decision Procedures. *TOPLAS* (1979).
- Andrew Reynolds, Cesare Tinelli, Amit Goel, Sava Krstić, Morgan Deters, and Clark Barrett. 2013. Quantifier Instantiation Techniques for Finite Model Finding in SMT. In *CADE'13*.
- Damien Sereni and Neil D. Jones. 2005. Termination Analysis of Higher-order Functional Programs. In *APLAS'05*.
- Jeremy G. Siek and Walid Taha. 2006. Gradual Typing for Functional Languages. In *Workshop on Scheme and Functional Programming*.
- Jan Smans, Bart Jacobs, and Frank Piessens. 2009. Implicit Dynamic Frames: Combining Dynamic Frames and Separation Logic. In *ECCOP 2009*.
- Nikolai Tillmann and Jonathan de Halleux. 2008. Pex—White Box Test Generation for .NET. In *TAP'08*.
- Niki Vazou, Leonidas Lampropoulos, and Jeff Polakow. 2017. A tale of two provers: verifying monoidal string matching in liquid Haskell and Coq. *International Symposium on Haskell*.
- Niki Vazou, Eric L. Seidel, Ranjit Jhala, Dimitrios Vytiniotis, and Simon Peyton-Jones. 2014. Refinement Types for Haskell. In *ICFP'14*.
- Niki Vazou, Anish Tondwalkar, Vikraman Choudhury, Ryan Scott, Ryan Newton, Philip Wadler, and Ranjit Jhala. 2018. Refinement Reflection: Complete Verification with SMT. In *POPL'18*.

MIL, a Monadic Intermediate Language for Implementing Functional Languages

Mark P. Jones
Portland State University
Portland, Oregon, USA
mpj@pdx.edu

Justin Bailey
Portland, Oregon, USA
jgbailey@gmail.com

Theodore R. Cooper
Portland State University
Portland, Oregon, USA
ted.r.cooper@gmail.com

ABSTRACT

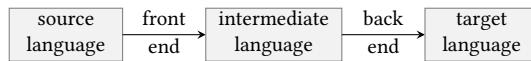
This paper describes MIL, a “monadic intermediate language” that is designed for use in optimizing compilers for strict, strongly typed functional languages. By using a notation that exposes the construction and use of closures and algebraic datatype values, for example, the MIL optimizer is able to detect and eliminate many unnecessary uses of these structures prior to code generation. One feature that distinguishes MIL from other intermediate languages in this area is the use of a typed, parameterized notion for basic blocks. This both enables new optimization techniques, such as the ability to create specialized versions of basic blocks, and leads to a new approach for implementing changes in data representation.

ACM Reference Format:

Mark P. Jones, Justin Bailey, and Theodore R. Cooper. 2019. MIL, a Monadic Intermediate Language for Implementing Functional Languages. In *Proceedings of International Symposium on Implementation and Application of Functional Languages (IFL'18)*. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

1 INTRODUCTION

Modern compilers rely on intermediate languages to simplify the task of translating programs in a rich, high-level source language to equivalent, efficient implementations in a low-level target. For example, a compiler may use two distinct compilation steps, with a *front end* that translates source programs to an intermediate language (IL) and a *back end* that translates from IL to the target:



This splits the compilation process into smaller, conceptually simpler components, each of which is easier to implement and maintain. Moreover, with a well-defined IL, it is possible for the front and back end components to be developed independently and to be reused in other compilers with the same source or target.

The biggest challenge in realizing these benefits is in identifying a suitable intermediate language. A good IL, for example, should retain some high-level elements, allowing a relatively simple translation from source to IL, and avoiding a premature commitment to

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

IFL'18, August 2019, Lowell, MA, USA

© 2019 Association for Computing Machinery.
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00
<https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

low-level details such as data representation and register allocation. At the same time, it should also have some low-level characteristics, exposing implementation details that are hidden in the original source language as opportunities for optimization, but ultimately also enabling a relatively simple translation to the target.

In this paper, we describe MIL, a “monadic intermediate language”, inspired by the monadic metalanguage of Moggi [12] and the monad comprehensions of Wadler [16], that has been specifically designed for use in implementations of functional programming languages. MIL bridges the tension between the high-level and low-level perspectives, for example, by including constructs for constructing, defining, and entering closures. These objects, and associated operations, are important in the implementation of higher-order functions, but the details of their use are typically not explicit in higher-level languages. At the same time, MIL does not fix a specific machine-level representation for closures, and does not specify register-level protocols for entering closures or for accessing their fields, delegating these decisions instead to individual back ends. This combination of design choices enables an optimizer for MIL, for example, to detect and eliminate many unnecessary uses of closures before the resulting code is passed to the back end.

MIL was originally developed as a platform for experimenting with the implementation and optimization of functional languages [2]. Since then, it has evolved in significant ways, both in the language itself—which now supports a type system with multi-value returns, for example—and its implementation—including an aggressive whole-program optimizer and novel representation transformation and specialization components. The implementation can be used as a standalone tool, but it can also be integrated into larger systems. We currently use MIL, for example, as one of several intermediate languages in a compiler for the Habit programming language [15] with the following overall structure:



In the initial stage of this compiler, Habit programs are desugared to LC, a simple functional language that includes lambda and case constructs (LC is short for “LambdaCase” and is named for those two features) but lacks the richer syntax and semantics of Habit. LC programs are then translated into MIL, processed by the tools described in this paper, and used to generate code for LLVM [9, 10], which is a popular, lower-level IL. There is, for example, no built-in notion of closures in LLVM, so it would be difficult to implement analogs of MIL closure optimizations at this stage of the compiler. However, the LLVM tools do implement many important, lower-level optimizations and can ultimately be used to generate assembly code (“asm” in the diagram) for an executable program.

1.1 Outline and Contributions

We begin the remaining sections of this paper with a detailed overview of the syntax and features of MIL (Section 2). To evaluate its suitability as an intermediate language, we have built a practical implementation that uses MIL as the target language for a simple front end (Section 3); as a framework for optimization (Section 5), and as a source language for an LLVM back end (Section 6). Along the way, we describe new global program transformations techniques for implementing changes in data representation (Section 4).

MIL has much in common with numerous other ILs for functional languages, including approaches based on CPS [1, 8], monads [12, 16], and other functional representations [5, 11]; these references necessarily only sample a very small portion of the literature. Our work makes some new contributions by exploring, for example, a parameterized form of basic block, and the new optimizations for “derived blocks” that this enables (Section 5.5), as well as new techniques for representation transformations (Section 4). An additional contribution of our work is the combination of many small design and engineering decisions that, together, result in a coherent and effective IL design, scaling to a practical system that can be used either as a standalone tool or as part of a larger compiler.

2 AN OVERVIEW OF MIL

We begin with an informal introduction to MIL, highlighting fundamental concepts, introducing terminology, and illustrating the concrete syntax for programs. We follow Haskell conventions [14] for basic syntax—from commenting and use of layout, to lexical details, such as the syntax of identifiers. That said, while we expect some familiarity with languages like Haskell, we will not assume deep knowledge of any specific syntax or features.

2.1 Types in MIL

MIL is a strongly typed language and all values that are manipulated in MIL must have a valid type. MIL provides several builtin types, including `Word` (machine words) and `Bool` (single-bit booleans). MIL also supports first class functions with types of the form $[d_1, \dots, d_n] \rightarrow [r_1, \dots, r_m]$ that map a tuple of input values with types d_i to a tuple of results with types r_j . (The prefixes `d` and `r` were chosen as mnemonics for domain and range, respectively.) Curried functions can also be described by using multiple \rightarrow types. For example, the type `[Word] → [[Word] → [Flag]]` could be used for a curried function that compares `Word` values.

MIL programs can also include definitions for parameterized algebraic datatypes, as illustrated by the following familiar examples:

```
data List a = Nil | Cons a (List a)
data Pair a b = MkPair a b
```

Each definition introduces a new type name (`List` and `Pair` in these examples) and one or more *data constructors* (`Nil`, `Cons`, and `MkPair`). Each data constructor has an *arity*, which is the number of fields it stores. For example, `Nil` has arity 0 and `Cons` has arity 2.

MIL also supports `bitdata` types [4] that specify a bit-level representation for all values of the type. In other respects, however, values of these types can be used in the same way as values of algebraic datatypes. With the following definitions, for example: Every value of type `Bool` is either `False` or `True`, with either option

being represented by a single bit; and a `PCI` address can be described by a 16 bit value with subfields of width 8, 5, and 3:

```
bitdata Bool = False [B0] | True [B1]
bitdata PCI = PCI [bus::Bit 8 | dev::Bit 5 | fun::Bit 3]
```

`Bitdata` types allow programmers to match the bit-level representation of data structures that are used in low-level systems applications, which is a key domain for Habit. Ultimately, however, `bitdata` values are represented by sequences of `Word` values (Section 4.2). For completeness, we mention that MIL also supports several other builtin types for systems programming, including the `Bit n` types seen here that represent bit vectors of width n .

2.2 MIL programs

The core of a MIL program is a sequence of block, closure, and top-level definitions, as suggested by the following grammar:

```
prog ::= def1 ... defn           -- program
def ::= b[v1, ..., vn] = c        -- block
      | k{v1, ..., vn} [u1, ..., um] = c    -- closure
      | [v1, ..., vn] ← t                      -- top-level
      | ...                                         -- other
```

The last line here uses “ \dots ” as a placeholder for other forms of definition, such as those introducing new types (see the previous section). The full details, however, are not central to this paper.

2.3 Blocks

Blocks in MIL are like basic blocks in traditional optimizing compilers except that each block also has a list of zero or more parameters. Intuitively, these can be thought of as naming the “registers” that must be loaded with appropriate values before the block is executed. The syntax of block definitions is reminiscent of Haskell’s monadic notation (without the `do` keyword), as in the following schematic example for a block `b` with arguments u_1, \dots, u_n :

```
b[u1, ..., un] = v1 ← t1 -- run t1, capture result in v1
...
vn ← tn -- run tn, capture result in vn
e
```

The body of this block is executed by running the statements $v_i \leftarrow t_i$ in sequence, capturing the result of each computation t_i in a variable v_i so that it can be referenced in later steps. The expression `e` on the last line may be either a tail expression (such as a `return` or an unconditional jump to another block) or a `case` or `if` construct that performs a conditional jump.

MIL relies on type checking to ensure basic consistency properties (for example, to detect calls with the wrong number or types of arguments). The code fragment above assumes that each t_i returns exactly one result but it is possible to use computations that return zero or more results at the same time, as determined by the type of t_i . In these cases, the binding for t_i must specify a corresponding list of variables, enclosed in brackets, as in $[w_1, \dots, w_m] \leftarrow t_i$.

2.4 Tail Expressions

MIL allows several different forms of “tail” expressions (each of which may be used in place of the expressions t_1, \dots, t_n in the preceding discussion), as summarized in the following grammar.

```

t ::= return [a1,...,an] -- monadic return
| b[a1,...,an] -- basic block call
| p((a1,...,an)) -- primitive call
| C(a1,...,an) -- allocate data value
| C i a -- select component
| k{a1,...,an} -- allocate closure
| f @ [a1,...,an] -- enter closure

```

Every tail expression represents a computation that (potentially) returns zero or more results:

- `return [a1,...,an]` returns the values a₁,...,a_n immediately with no further action. The most frequently occurring form produces a single result, a, and may be written without brackets as `return a`.
- `b[a1,...,an]` is a call to a block b with the given arguments. When a block call occurs at the end of a code sequence, it can potentially be implemented as a tail call (i.e., a jump). This is why we refer to expressions in the grammar for t as *tail expressions* or sometimes *generalized tail calls*.
- `p((a1,...,an))` is a call to a primitive p with the given arguments. Primitives typically correspond to basic machine instructions (such as add, and, or shr), but can also be used to link to external functions. Primitive calls are written with double parentheses around the argument list so that they are visibly distinguished from block and constructor calls.
- `C(a1,...,an)` allocates a data value with constructor C and components a₁,...,a_n, where n is the arity of C.
- `C i a` is a selector that returns the ith component of the value in a, assuming that it was constructed with C. For example, having executed `v ← C(a1,...,an)`, a subsequent `C i v` call (with $0 \leq i < n$) will return a_{i+1}. (Note that i must be a constant and that the first component has index 0.)
- `k{a1,...,an}` allocates a closure (i.e., the representation of a function value) with a code pointer specified by k and arguments a₁,...,a_n. This is only valid in a program that includes a corresponding definition for k. Closures are an important feature of MIL and are the topic of further discussion throughout this paper.
- `f @ [a1,...,an]` is a call to an unknown function, f, with arguments a₁,...,a_n. For this to be valid, f must hold a closure for a function that expects exactly n arguments. The process of executing an expression of this form is often referred to as *entering* a closure. Mirroring the syntax for `return`, we write `f @ a` in the case where there is exactly one argument. As another special case, an unknown function f can be applied to an empty list of arguments, as in `f @ []`; this can be used, for example, to invoke a (monadic) thunk.

2.5 Atoms

In the preceding examples, the symbols a, a₁,...,a_n, and f represent arbitrary *atoms* that are either numeric constants (integer literals) or variable names. The latter correspond either to *temporaries* (i.e., local variables holding parameter values or intermediate results) or to variables defined at the top level (see Section 2.9):

```

a ::= i -- a constant (i.e., an integer literal)
| v -- a variable name

```

Atoms are the only form of expression that can be used as arguments in tail expressions. If a program requires the evaluation of a more complex expression, then it must be computed in steps using temporaries to capture intermediate results. For example, the following block calculates the sum of the squares of its inputs using the mul and add primitives for basic arithmetic:

```

sumSqr[x,y] = u ← mul((x,x)) -- compute x*x in u
v ← mul((y,y)) -- compute y*y in v
add((u,v)) -- return the sum

```

2.6 Block Types

Blocks are not first class values in MIL, so they cannot be stored in variables, passed as inputs to a block, or returned as results. Nevertheless, it is still useful to have a notion of types for describing the inputs and results of each block, and we use the notation `[d1,...,dn] >>= [r1,...,rm]` for this where d_i and r_j are the types of the block's inputs and results, respectively. For example, the type of `sumSqr` can be declared explicitly as follows:¹

```
sumSqr :: [Word,Word] >>= [Word]
```

Blocks may have polymorphic types, as in the following examples, the second of which also illustrates a block with multiple results:

```

idBlock :: forall (a::type). [a] >>= [a]
idBlock[x] = return x

```

```

dupBlock :: forall (a::type). [a] >>= [a, a]
dupBlock[x] = return [x, x]

```

The explicit `forall` quantifiers here give a precise way to document polymorphic types but they are not required; those details can be calculated automatically instead using kind and type inference.

It is important to distinguish block types, formed with `>>=`, from the first-class function types using `→` that were introduced in Section 2.1. The notations are similar because both represent a kind of function, but there are also fundamental differences. In particular, blocks are executed by jumping to a known address and not by entering a closure, as is done with a first-class function.

2.7 Code Sequences

The syntax for the code sequences on the right of each block definition is captured in the following grammar:

```

c ::= [v1,...,vn] ← t; c -- monadic bind
| assert a C; c -- data assertion
| t -- tail call
| if v then bc1 else bc2 -- conditional
| case v of alts -- case construct
alts ::= {alt1;...;altn} -- alternatives
alt ::= C → bc -- match against C
| _ → bc -- default branch
bc ::= b[a1,...,an] -- block call

```

As illustrated previously, any given code sequence begins with zero or more *monadic bind* steps (or *statements*) of the form `vs ← t` (where vs is either a single temporary or a bracketed list). Each such step describes a computation that begins by executing the tail

¹The same block type is also used for the mul and add primitives.

t , binding any results to the temporaries in vs , and then continuing with the entries in vs now in scope. The grammar also includes an `assert` construct; its purpose will be illustrated later.

Where space permits, we write the steps in a code sequence using a vertical layout, relying on indentation to reflect its structure. However, we will also use a more compact notation, matching the preceding grammar with multiple statements on a single line:

```
v1 ← t1; v2 ← t2; ...; vn ← tn; e
```

However they are written, every code sequence is ultimately terminated by an expression e that is either an unconditional jump (i.e., a tail expression, t); or a conditional jump (i.e., an `if` or `case` construct). Note that, like a traditional basic block, there are no labels in the middle of code sequences. As a result, we can only enter a code sequence at the beginning and, once entered, we cannot leave until we reach the end expression e .

2.8 Conditional Jumps

An `if` construct ends a code sequence with a conditional jump; it requires an atom, v , of type `Flag` as the test expression and block calls for each of the two branches. Beyond these details, `if` constructs work in the usual way. For example, the following two blocks can be used to calculate the maximum of two values u and v :

```
b1[u,v] = t ← primGt((u,v)) -- primitive for u>v
          if t then b2[u] else b2[v]
b2[x]   = return x
```

Other conditional jumps can be described with a `case v of alts` construct, using the value of v —the “discriminant”—to choose between the alternatives in `alts`. Specifically, a `case` construct is executed by scanning the list of alternatives and executing the block call `bc` for the first one whose constructor, C , matches v . The following code sequence corresponds to an `if-then-else` construct for a variable v of type `Bool`: the code will execute `b1[x,y]` if v was built using `True()`, or `b2[z]` if v was built using `False()`:

```
case v of { True → b1[x,y]; False → b2[z] }
```

We can also write the previous expression using a default branch (signaled by an underscore) instead of the `case` for `False`:

```
case v of { True → b1[x,y]; _ → b2[z] }
```

Once again, we often write examples like this with a vertical layout, eliding the punctuation of semicolons and braces:

```
case v of
  True → b1[x,y]
  _     → b2[z]
```

Each alternative in a `case` only tests the outermost constructor, C , of the discriminant, v . If the constructor matches but subsequent computation requires access to components of v , then they must be referenced explicitly using `C i v` selectors, where i is the appropriate field number. The following blocks, for example, can be used to compute the length of a list:

```
length[list] = loop[0,list]
loop[n,list] = case list of
  Nil → done[n]
  Cons → step[n,list]
done[n]      = return n
```

```
step[n,list] = assert list Cons
              tail ← Cons 1 list
              m ← add((n,1))
              loop[m,tail]
```

The `length` block passes the given `list` to `loop`, with an initial (accumulating) parameter \emptyset . The `loop` block uses a `case` to examine the `list`. A `Nil` value represents an empty list, in which case we branch to `done` and immediately return the current value of n . Otherwise, `list` must be a `Cons`, and we jump to `step`, which uses a `Cons 1 list` selector to find the tail of the list. After calculating an incremented count in m , `step` jumps back to examine the rest of the `list`. The `assert` here indicates that `list` is known to be a `Cons` node; this information ensures that the `Cons` selector in the next line is valid, and can also be leveraged during optimization.

2.9 Top-level Definitions

MIL programs can use definitions of the form $[v_1, \dots, v_n] \leftarrow t$ to introduce top-level bindings for variables called v_1, \dots, v_n . Variables defined in this way are initialized to the values that are produced by executing the tail expression t . In the common case where a single variable is defined, we omit the brackets on the left of the \leftarrow symbol and write $v_1 \leftarrow t$. Note that variables introduced in a definition like this are not the same as *global variables* in an imperative programming language because their values cannot be changed by a subsequent assignment operation.

Top-level definitions like this are typically used to provide names for constants or data structures (including closures, as shown below). For example, the following top-level definitions, each with a data constructor on the right hand side of the \leftarrow symbol, will construct a static list data structure with two elements:

```
list1 ← Cons(1, list2)
list2 ← Cons(2, nil)
nil   ← Nil()
```

It is possible to specify types for variables introduced in a top-level definition using declarations of the form $v_1, \dots, v_n :: t$. For example, we might declare types for the `list` variables defined previously by writing:

```
list1, list2, nil :: List Word
```

However, there is no requirement to provide types for top-level variables because they can also be inferred in the usual way.

2.10 Closures and Closure Definitions

First-class functions in MIL are represented by *closure* values that are constructed using tail expressions of the form $k\{a_1, \dots, a_n\}$. Here, a_1, \dots, a_n are atoms that will be stored in the closure and accessed when the closure is entered (i.e., applied to an argument) and the code identified by k is executed. For each different k that is used in a given program, there must be a corresponding *closure definition* of the form $k\{v_1, \dots, v_n\} [u_1, \dots, u_m] = t$. Here, the v_i are variables representing the values stored in the closure; the u_i are variables representing the arguments that are required when the closure is entered; and t is a tail expression that may involve any of these variables. Conceptually, each such closure definition corresponds to the code that must be executed when a closure is entered, loading stored values and arguments as necessary in to

registers before branching to the code as described by t . Following our usual pattern, in the case where there is just one argument, u_1 , the brackets may be omitted and the definition can be written $k\{v_1, \dots, v_n\} u_1 = t$. The ability to pass multiple (or zero) arguments u_i to a closure, however, is important because it allows us to work with values whose representation may be spread across multiple components; we will return to this in Section 4.2.

As an example, given the definition $k\{n\} x = \text{add}((n, x))$, we can use a closure with code pointer k and a stored free variable n to represent the function $(\lambda x \rightarrow n + x)$ that will return the value $n+x$ whenever it is called with an argument x .

The type of any closure can be written in the form:

$$\{t_1, \dots, t_n\} [d_1, \dots, d_m] \rightarrow [r_1, \dots, r_p]$$

where the t_i are the types of the stored components, the d_j are the types of the inputs (the domain), and the r_k are the types of the results (the range). This is actually a special case of an *allocator type* $\{t_1, \dots, t_n\} t$, which corresponds to a value of type t whose representation stores values of the types listed in the braces. For example, the type of the closure k defined above is $\{\text{Word}\} [\text{Word}] \rightarrow [\text{Word}]$, while the type of the `Cons` constructor for lists can be written $\{a\} \text{List } a \rightarrow \text{List } a$.

2.11 Example: Implementing map in MIL

In this section, we illustrate how the components of MIL described previously can be used together to provide an implementation for the familiar `map` function. In a standard functional language, we might define this function using the following two equations:

$$\begin{aligned} \text{map } f \text{ Nil} &= \text{Nil} \\ \text{map } f \text{ (Cons } y \text{ ys)} &= \text{Cons } (f \text{ } y) \text{ (map } f \text{ ys)} \end{aligned}$$

Using only naive techniques, this function can be implemented by the following MIL code (presented here in two columns) with one top-level, two closure, and three block definitions:

$$\begin{array}{ll} \begin{array}{l} \text{map } f \leftarrow k_0\{\} \\ k_0\{ \} \text{ } f \text{ } = \text{k}_1\{f\} \\ k_1\{f\} \text{ } xs \text{ } = \text{b}_0[f, xs] \\ \text{b}_0[f, xs] \text{ } = \text{case } xs \text{ of} \\ \quad \text{Nil} \rightarrow b_1[] \\ \quad \text{Cons} \rightarrow b_2[f, xs] \\ b_1[] \text{ } = \text{Nil}() \end{array} & \begin{array}{l} b_2[f, xs] \text{ } = \text{assert } xs \text{ Cons} \\ y \leftarrow \text{Cons } 0 \text{ } xs \\ ys \leftarrow \text{Cons } 1 \text{ } xs \\ z \leftarrow f @ y \\ m \leftarrow \text{map } @ f \\ zs \leftarrow m @ ys \\ \text{Cons}(z, zs) \end{array} \end{array}$$

This implementation starts with a top-level definition for `map`, binding it to a freshly constructed closure $k_0\{\}$. When the latter is entered with an argument f , it captures that argument in a new closure structure $k_1\{f\}$. No further work is possible until this second closure is entered with an argument xs , which results in a branch to the block b_0 , and a test to determine whether the list value is either a `Nil` or `Cons` node. In the first case, the `map` operation is completed by returning `Nil` in block b_1 . Otherwise, we execute the code in b_2 , extracting the head, y , and tail, ys , from the argument list. This is followed by three closure entries: the first calculates the value of $f \text{ } y$, while the second and third calculate `map f ys`, with one closure entry for each argument. The results are then combined using `Cons(z, zs)` to produce the final result for the `map` call.

The MIL definition of `map` reveals concrete implementation details, such as the construction of closures, that are not immediately visible in the original code. For an intermediate language, this is

exactly what we need to facilitate optimization, and we will revisit this example in Section 5.4, showing how the code in b_2 can be rewritten to avoid unnecessary construction of closures.

2.12 Notes on Formalization of MIL

Although we do not have space here for many details, we have developed both a formal type system and an abstract machine for MIL. The former has served as a guide in the implementation of the MIL typechecker, which is useful in practice for detecting errors in MIL source programs (and, occasionally, for detecting bugs in our implementation). Perhaps the most interesting detail here is the appearance of a type constructor, M —representing the underlying monad in which MIL computations take place—in rules like the following (for type checking bindings in code sequences):

$$\frac{A \vdash t : M[r_1, \dots, r_n] \quad A, v_1 : r_1, \dots, v_n : r_n \vdash c : Ma}{A \vdash ([v_1, \dots, v_n] \leftarrow t; c) : Ma}$$

An interesting direction for future work is to allow the fixed M to be replaced with a type variable, m , and to perform a monadic effects analysis by collecting constraints on m .

The design of an abstract machine for MIL has also had practical impact, guiding the implementation of a bytecode interpreter that is useful for testing. In the future, the abstract machine may also provide a formal semantics for verifying the program rewrites used in the MIL optimizer (Section 5).

3 COMPIILING LC TO MIL

In this section, we discuss the work involved in compiling programs written in LC—a simple, high-level functional language—into MIL. This serves two important practical goals: First, by using MIL as the target language, we demonstrate that it has sufficient features to serve as an intermediate language for a functional language with higher-order functions, pattern matching, and monadic operations. Second, in support of testing, it is easier to write programs in LC and compile them to MIL than to write MIL code directly.²

3.1 Translating LC Types to MIL

The task of translating LC types to corresponding MIL types is almost trivial: the two languages have the same set of primitive types and use the same mechanisms for defining new data and bitdata types. The only complication is in dealing with function types in LC, which map a single input to a single result, instead of the tuples of inputs and results that are used in MIL. To bridge this gap, we define the LC function type, written using a conventional infix \rightarrow symbol, as an algebraic datatype:

$$\text{data } d \rightarrow r = \text{Func } ([d] \rightarrow [r])$$

With this definition, we now have three different notions of function types for MIL: \rightarrow and \Rightarrow describe first-class function values while $\gg=$ is for block types (Section 2.6). The following definitions show how all of these function types can be used together in a MIL implementation of the LC identity function, $(\lambda x \rightarrow x)$:

²Our implementation actually accepts combinations of MIL and LC source files as input; this allows users to mix higher-level LC code that is translated automatically to MIL with handwritten MIL code. The latter is useful in practice for implementing libraries of low-level primitives that cannot be expressed directly in LC.

```

k   :: {} [a] → [a] -- a closure definition for
k{} x = return x    -- the identity function

b :: [] >> [a → a] -- create a closure value and
b[] = c ← k{}       -- package it as a → function
Func(c)

id :: a → a          -- set top-level name id to the
id ← b[]             -- value produced by b[]

```

A legitimate concern here is that every use of an LC function (\rightarrow) requires extra steps to wrap or unwrap the Func constructor around a MIL function (\Rightarrow). Fortunately, we will see that it is possible to eliminate these overheads (Section 4.1).

3.2 Translating LC Code to MIL

There is a large body of existing work on compilation of functional languages, much of which can be easily adapted to the task of translating LC source programs in to MIL. As a simple example, to compile a lambda expression like $\lambda x \rightarrow e$ with free variables fvs , we just need to: pick a new closure name, k ; add a definition $k\{fvs\} x = e'$ to the MIL program (where e' is compiled code for e); and then use $k\{fvs\}$ in place of the original lambda expression. Beyond these general comments, we highlight the following details from our implementation:

- To simplify code generation, we use a lambda lifting transformation [6] to move locally defined recursive definitions to the top-level (possibly with added parameters).
- Our code generator is based on compilation schemes for “compiling with continuations” [8]. As a rough outline (in pseudo-Haskell notation), it can be described as a function:

```
compile :: Expr → (Tail → CM Code) → CM Code
```

The first argument is the LC expression, exp , that we want to compile. The second argument is a continuation that takes a MIL tail expression, t , corresponding to exp , and embeds it in a MIL code sequence for the rest of the program. For instance, if exp is the lambda expression $\lambda x \rightarrow e$ in the example at the start of this section, then t will be the closure allocator $k\{fvs\}$. Note that CM in the type above represents a “compilation monad”, with operations for generating new temporaries, and for adding new block, closure, or top-level definitions to the MIL program as compilation proceeds.

- Continuation based techniques require special care with conditionals like $if\ c\ then\ t\ else\ f$. In particular, we need to ensure that the $(Tail \rightarrow CM\ Code)$ continuation is not applied separately to each of the tail expressions for t and f , which could lead to duplicated code. Fortunately, it is easy to avoid this in MIL by placing the continuation code in a new block, serving as a “join point” [11], and then having the code in each branch end with a jump to that block.

4 REPRESENTATION TRANSFORMATIONS

Although MIL is more broadly applicable, our current implementation assumes a whole-program compilation model. One benefit of this is that we can consider transformations that require changes across many parts of input programs. In the following subsections,

we describe three specific transformations of this kind, all implemented in our toolset, that change the representation of data values in MIL programs. Each of these typically requires modifications in both code and type definitions. For example, if a program uses values of type T that can be more efficiently represented as values of type T' , then implementing a change of representation will require updates, not only to code that manipulates values of type T , but also to any type definitions that mention T . Our tools allow these transformations to be applied in any order or combination, running the MIL type checker after each pass as a sanity check (although the original types are not preserved, each transformation is expected to preserve typeability). In addition, because these transformations can create new opportunities for optimization, it is also generally useful to (re)run the MIL optimizer (Section 5) after each pass.

4.1 Eliminating Single Constructor Types

Our first application for representation transformations deals with ‘single constructor’, non-recursive algebraic datatypes with definitions of the form: $data\ T\ a_1\ …\ a_n = C\ t'$. Types like this are often used to create wrappers that avoid type confusion errors. For example, by defining $data\ Time = Millis\ Word$, we ensure that values of type $Word$ are not used accidentally where $Time$ values are actually required. And, as discussed in Section 3.1, we also use a type of this form to map between the \rightarrow and \Rightarrow function types when compiling LC to MIL. These types are useful because they can enforce correct usage in source programs. But for compilation purposes, the extra constructors—like $Millis$ and $Func$ —add unnecessary runtime overhead, and can block opportunities for optimization. To avoid this, we can rewrite the MIL program to replace every tail of the form $C\ 0\ a$ or $C(a)$ with $return\ a$ (effectively treating selection and construction as the identity function) and every case v of $C \rightarrow bc$ that ‘matches’ on C with a direct block call bc . In addition, the original definition of T can be discarded, but every remaining type of the form $T\ t_1\ …\ t_n$ must be replaced with the corresponding instance $[t_1/a_1, …, t_n/a_n]t'$ of t' .

4.2 Representation Vectors

Our second application was initially prompted by the use of `bitdata` types in MIL (see Section 2.1) but is also useful with some algebraic datatypes. In the early stages of compilation, we manipulate values of `bitdata` types like `Bool` or `PCI` with the same pattern matching and constructor mechanisms as algebraic datatypes. At some point, however, the compiler must transition to the bit-level representation that is specified for each of these types. This means, for example, that values of type `Bool` should be represented by values of type `Flag`. Similarly, `PCI` values might be represented using `Word` values, with the associated selectors for `bus`, `dev`, and `fun` replaced by appropriate bit-twiddling logic using `shift` and `mask` operations.

To specify representation changes like this, we define a function that maps every MIL type t to a suitable *representation vector*: a list of zero or more types that, together, provide a representation for t . We use a vector, rather than a single type, to accommodate types that require multi-word representations. A `Bit 64` value, for example, cannot be stored in one 32 bit word, but can be supported by using a representation vector `[Word, Word]` with two `Word` values. Similarly, a zero length vector, `[]`, can be used for `Bit 0` and,

indeed, for any other singleton type, such as the standard unit type, $()$. This reflects the fact that, if a type only has one possible value, then it does not require a runtime representation.

The resulting representation vectors can be used to guide a program transformation that replaces each variable v of a type t with a list of zero or more variables v_1, \dots, v_n , where n is the length of the representation vector for t . This was the primary technical motivation for using tuples of values throughout MIL because it allows us to rewrite tails like $f @ a$ and top-level definitions like $v \leftarrow t$, for example, as $f @ [a_1, a_2]$ or $[v_1, v_2] \leftarrow t$ when a and v are each represented by two words. Of course, additional rewrites are needed for selectors and primitive calls that use values whose representation is changed. A convenient way to manage this is to replace the operations in question with a calls to new blocks that will be inlined and optimized with the surrounding code.

Support for the representation transformation described here is a distinguishing feature of our toolset: to our knowledge, no other current system implements a transformation of this kind.

4.3 Specializing Polymorphic Definitions

Our third application, included to satisfy a requirement of the LLVM backend (Section 6), is a transformation that eliminates polymorphic definitions and generates specialized code for each monomorphic instance that is needed in a given program. One problem is that this transformation cannot be used in some programs that rely on polymorphic recursion [13]. Another concern is that specialization has the potential to cause an explosion in code size. In practice, however, specialization has proved to be an effective tool in domain-specific [3] and general-purpose functional language compilers [17], and even in implementations of overloading [7].

Representation transformation is relevant here when we consider the task of generating code for specific monomorphic instances of polymorphic functions. As part of this process, our implementation of specialization also eliminates all uses of parameterized datatypes. A program that uses a value of type $\text{Pair } \text{PCI PCI}$, for example, might be transformed in to code that uses a value of a new type $\text{data Pair}_1 = \text{MkPair}_1 \text{ PCI PCI}$, that is generated by the specializer. This provides an interesting opportunity for using type-specific representations. For example, the Pair_1 type described here should hold two 16-bit PCI values, so it could easily be represented using a single, 32-bit Word with no need for heap allocation. We have already started to explore the possibility of inferring bit-level representations for a wide-range of types like this, and expect to include support for this in a future release of our toolset.

5 OPTIMIZATION

A common goal in the design of an intermediate language is to support optimization: program transformations that preserve program behavior but improve performance, reduce memory usage, etc. For MIL programs, we have identified a collection of rules that describe how some sections of code can be rewritten to obtain better performance. A small but representative set of these rules is presented in Figure 1. The full set has been used to build an optimizer for MIL that works by repeatedly applying rewrites to input programs. Individual rewrites typically have limited impact, but using many

rewrites in sequence can yield significant improvements by reducing code size, eliminating unnecessary operations, and replacing expensive operations with more efficient equivalents.

The table in Figure 1 has three columns that provide, for each rule: a short name; a rewrite; and, in several cases, a set of side conditions that must be satisfied in order to use the rule. The rewrites are written in the form $e \Rightarrow e'$ indicating that an expression matching e should be replaced by the corresponding e' . To avoid ambiguity, we use two variants of this notation with \Rightarrow_c for rewrites on code sequences and \Rightarrow_t for rewrites on tails. In the rest of this section, we will walk through each of the rewrites in Figure 1 in more detail (Sections 5.1–5.5), describe additional optimizations that are supported by our implementation (Section 5.6), and reflect on the overall effectiveness of our optimizer for MIL (Section 5.7).

5.1 Using the Monad Laws

The first group of rules are by standard laws for monads (as described, for example, by Wadler [16]), but are also recognizable as traditional compiler optimizations. Rule (1), for example, implements a form of *copy* or *constant propagation*, depending on whether the atom a is a variable or a constant. The notation $[a/x]c$ here represents the substitution of a for all free occurrences of x in the code sequence c ; concretely, instead of creating an extra variable, x , to hold the value of a , we can just use that value directly. Rule (2) corresponds to the right monad law, which can also be seen as eliminating an unnecessary return at the end of a code sequence and potentially as creating a new opportunity for a tail call. Finally, Rules (3) and (4) are based on the associativity law for monads, but can also be understood as descriptions of function inlining rules; we use the terms *prefix* and *suffix* to distinguish between cases where the block being inlined appears at the beginning or the end of the code sequence. To better understand the relationship with the associativity law, note that a naive attempt to inline the block b in the code sequence $v \leftarrow b[x]; c$, using the definition of b in the figure, would produce $v \leftarrow (v_0 \leftarrow t_0; t_1); c$. With the previous grammar for MIL code sequences, this is not actually a valid expression. Using associativity, however, it can be flattened/rewritten as the code sequence on the right hand side of the rule. As is always the case, unrestricted use of inlining can lead to an explosion in code size with little or no benefit in performance. To avoid such problems, our implementation uses a typical set of heuristics—limiting inlining to small blocks or to blocks that are only used once, for example—to strike a good balance between the benefits and potential risks of an aggressive inlining strategy.

5.2 Eliminating Unnecessary Code

The second group of rewrites in Figure 1 provide ways to simplify MIL programs by trimming unnecessary code. Rule (5), for example, uses the results of a simple analysis to detect tail expressions t that do not return (e.g., because they call a primitive that terminates the program, or enter an infinite loop). In these situations, any code that follows t can be deleted without a change in semantics.

Rules (6) and (7) allow us to detect, and then, respectively, to eliminate code that has no effect. Rule (6), sets the result variable name for a statement to underscore to flag situations where the original variable is not used in the following code. A subsequent

Name	Rewrite		Conditions
Using the monad laws , where block b is defined by $b[x] = v_0 \leftarrow t_0; t_1$			
1) Left monad law (constant propagation)	$x \leftarrow \text{return } a; c \Rightarrow_c [a/x]c$		-
2) Right monad law (tail call introduction)	$x \leftarrow t; \text{return } x \Rightarrow_c t$		-
3) Prefix inlining	$v \leftarrow b[x]; c \Rightarrow_c v_0 \leftarrow t_0; v \leftarrow t_1; c$		-
4) Suffix inlining	$v \leftarrow t; b[x] \Rightarrow_c v \leftarrow t; v_0 \leftarrow t_0; t_1$		-
Eliminating unnecessary code			
5) Unreachable code elimination	$v \leftarrow t; c \Rightarrow_c t$		t does not return
6) Wildcard introduction	$v \leftarrow t; c \Rightarrow_c _ \leftarrow t; c$		v is not free in c
7) Dead tail elimination	$_ \leftarrow t; c \Rightarrow_c c$		t is pure
8) Common subexpression elimination	$t \Rightarrow_t \text{return } v$		{v=t}
Using algebraic identities (focusing here on bitwise and and writing M, N, and P for arbitrary integer constants)			
9) Identity laws	$\text{and}((v, 0)) \Rightarrow_t \text{return } 0$		-
10) Idempotence	$\text{and}((v, v)) \Rightarrow_t \text{return } v$		-
11) Constant folding	$\text{and}((M, N)) \Rightarrow_t \text{return } (M \& N)$		-
12) Commutativity	$\text{and}((M, v)) \Rightarrow_t \text{and}((v, M))$		-
13) Associative folding $(u \& M) \& N = u \& (M \& N)$	$\text{and}((v, N)) \Rightarrow_t \text{and}((u, M \& N))$		{v=and((u, M))}
14) Distributive folding (1) $(u M) \& N = (u \& N) (M \& N)$	$\text{and}((v, N)) \Rightarrow_c v' \leftarrow \text{and}((u, N)) \text{ or } (v', M \& N)$		{v=or((u, M))}
15) Distributive folding (2) $((u M) \& N) P = (u \& N) ((M \& N) P)$	$\text{or}((w, P)) \Rightarrow_c v' \leftarrow \text{and}((u, N)) \text{ or } (v', (M \& N) P)$		{v=or((u, M)), w=and((v, N))}
Known structures , where closure k is defined by $k\{x\} y = t$			
16) Known constructor	$\text{case } v \text{ of } \dots; C \rightarrow b'[\dots]; \dots \Rightarrow_c b'[\dots]$		{c=C(...)}
17) Known closure	$f @ y \Rightarrow_t t$		{f=k{x}}
Derived blocks , where block b is defined by $b[x] = v_0 \leftarrow t_0; t_1$			
18) Known structure	$b[v] \Rightarrow_t b'[y]$ where $b'[y] = x \leftarrow C(y); v_0 \leftarrow t_0; t_1$		{v=C(y)}
19) Trailing enter	$f \leftarrow b[x]; f @ a \Rightarrow_c b'[x, a]$ where $b'[x, a] = v_0 \leftarrow t_0; f \leftarrow t_1; f @ a$		-
20) Trailing case	$v \leftarrow b[x]; \text{case } v \text{ of } \dots \Rightarrow_c b'[x, \dots]$ where $b'[x, \dots] = v_0 \leftarrow t_0; v \leftarrow t_1; \text{case } v \text{ of } \dots$		-

Figure 1: A representative set of rewrite rules for MIL optimization

use of Rule (7) will then eliminate the statement altogether if the associated tail expression t has no externally visible effects (for example, if t is a call to a pure primitive function like `add` or `mul`, or if t is a closure or data allocator). Our presentation of this process using two separate rules reflects the fact that our optimizer actually applies these two rules in separate passes over the abstract syntax.

Rule (8) uses a local dataflow analysis to detect situations where the result of a previous computation can be reused. To implement this, our optimizer calculates a set of “facts”, each of which is a statement of the form $v=t$, as it traverses the statements in each code sequence. Starting with an empty set, the optimizer will add (or “generate”) a new fact $v=t$ for every statement $v \leftarrow t$ that it encounters with a pure t. At the same time, for each statement $v \leftarrow t$, it will also remove (or “kill”) any facts that mention v because the variable that they reference will no longer be in scope. The rightmost column in Figure 1 documents the facts that are required to apply each rewrite. In this case, given $v=t$, we can

avoid recalculating t and just return the value v that it produced previously. (In practice, the `return` introduced here will often be eliminated later using Rules (1) or (2).)

5.3 Using Algebraic Identities

The next group (Rules (9)–(15)) take advantage of (mostly) well-known algebraic identities to simplify programs using the builtin primitives for arithmetic, logic, and comparison operations. We only show a small subset of the (more than 100) rewrites of this kind that are used in our implementation, all of which involve the bitwise and primitive. In combination with additional rewrites involving bitwise or and shift operations, these rules are very effective in simplifying the bit-twiddling code that is generated when manipulating or constructing Habit-style bitdata types [4, 15].

Rules (9), (10), and (11), for example, each eliminate a use of `and` in a familiar special case. Rule (12) is not useful as an optimization by itself, but instead is a first step in rewriting expressions in to

$b_2[f, xs] = \text{assert } xs \text{ Cons}$ $y \leftarrow \text{Cons } \emptyset \text{ xs}$ $ys \leftarrow \text{Cons } 1 \text{ xs}$ $z \leftarrow f @ y$ $m \leftarrow \text{map } @ f$ $zs \leftarrow m @ ys$ $\text{Cons}(z, zs)$	$b_2[f, xs] = \text{assert } xs \text{ Cons}$ $y \leftarrow \text{Cons } \emptyset \text{ xs}$ $ys \leftarrow \text{Cons } 1 \text{ xs}$ $z \leftarrow f @ y$ $m \leftarrow k_1\{f\}$ $zs \leftarrow m @ ys$ $\text{Cons}(z, zs)$	$b_2[f, xs] = \text{assert } xs \text{ Cons}$ $y \leftarrow \text{Cons } \emptyset \text{ xs}$ $ys \leftarrow \text{Cons } 1 \text{ xs}$ $z \leftarrow f @ y$ $m \leftarrow k_1\{f\}$ $zs \leftarrow b_0[f, ys]$ $\text{Cons}(z, zs)$	$b_2[f, xs] = \text{assert } xs \text{ Cons}$ $y \leftarrow \text{Cons } \emptyset \text{ xs}$ $ys \leftarrow \text{Cons } 1 \text{ xs}$ $z \leftarrow f @ y$ $zs \leftarrow b_0[f, ys]$ $\text{Cons}(z, zs)$
(a)	(b)	(c)	(d)

Figure 2: An example illustrating the optimization of known closures (Rule (17))

a canonical form that enables subsequent optimizations. In this case, the rewrite ensures that, for an and with one constant and one unknown argument, the constant is always the second argument. This explains, for example, why we do not need a variant of Rule (9) for tails of the form $\text{and}((\emptyset, v))$, and also why we do not need four distinct variations of the pattern in Rule (13) where an unknown u is combined via bitwise ands with two constants M and N . In this case, we rely on facts produced by the dataflow analysis to determine that the value of v in $\text{and}((v, N))$ was calculated using a prior $\text{and}((u, M))$ call. The rewrite here replaces one and with another, so it may not result in an immediate program optimization. However, there are two ways in which this *might* open up opportunities for subsequent rewrites: One possibility is that $M \& N$ may be zero, in which case we will be able to eliminate the and using Rule (9). Another possibility is that, by rewriting the call to and in terms of u , we may eliminate all references to v and can then eliminate the statement defining v as dead code using Rules (6) and (7).

In a similar way, Rules (14) and (15) take advantage of standard distributivity laws to implement more complex rewrites. By using these rules in combination, we can rewrite any expression involving bitwise ands and ors of an unknown u with an arbitrary sequence of constants into an expression of the form $(u \& M) | N$ for some constants M and N , with exactly one use of each primitive. These rewrites are potentially dangerous because they increase the size of the MIL code, replacing one primitive call on the left of the rewrite with two on the right. In practice, however, these rules often turn the definitions of the variables v and w that they reference in to dead code that can be eliminated by subsequent rewrites.

5.4 Known Structures

MIL provides case and @ constructs that work with arbitrary data values and closures, respectively. But the general operations are not needed in situations where we are working with known structures. Rule (16), for example, eliminates a conditional jump if the constructor, C , that was used to build v is already known: we can just make a direct jump using the alternative for C (or the default branch if there is no such alternative), and delete all other parts of the original case construct.

In a similar way, Rule (17) can be used to avoid a general closure entry operation when the specific closure is known. To see how this works in practice, consider the example in Figure 2, with the original code for b_2 in our implementation of map (Section 2.11) in Column (a). From the other definitions in this program, we know that

m is a reference to the closure $k_0\{\}$, and that $k_0\{\} f = k_1\{f\}$. By allowing our dataflow analysis to derive the fact $m = k_0\{\}$ from its top-level definition, we can use Rule (17) to rewrite the tail defining m , as shown in Column (b). After this transformation, the local dataflow analysis will reach the definition of zs with a list of facts that includes $m = k_1\{f\}$, and so we can apply Rule (17) again to obtain the code in Column (c). This removes the only reference to m , and allows subsequent uses of Rules (6) and (7) will eliminate its definition, producing the code definition in Column (d). This example shows that the original implementation of map would have allocated a fresh closure, $k_1\{f\}$, for every element of the input list. The transformations we have applied here, however, eliminate this overhead and substitute a more efficient, direct recursive call.

5.5 Derived Blocks

Optimizing compilers often use collections of basic blocks, connected together in control flow graphs, as a representation for programs. Most of the rules that we have described so far are traditionally considered local optimizations because they only consider the code within a single block. While much can be accomplished using local optimizations, it is also useful to take a more global view of the program, and to use optimizations that span multiple blocks. In other words, in addition to the *content* of individual blocks, we would also like to account for the *context* in which they appear.

One interesting way that we have been able to handle this in MIL is by using the code of existing blocks to generate new versions—which we refer to as *derived blocks*—that are specialized for use in a particular context. The final group of rewrites in Figure 1 corresponds to different strategies for generating derived blocks that we have found to be effective in the optimization of MIL programs. Of course, adding new blocks to a program increases program size and does not immediately provide an optimization. But, in practice, the addition of new derived blocks often opens new opportunities for optimization—by bringing the construction and matching of a data value into the same block, for example—and the original source block often becomes dead code that will be removed from the program once any specialized versions have been generated.

Rule (18) illustrates one way of using derived blocks to take advantage of information produced by our local dataflow analysis and to obtain results that typically require a global analysis. The techniques that we describe here can be applied very broadly, but, for this paper, we restrict ourselves to a special case: a call to a block b that has just one parameter and a very simple definition.

In this situation, if the argument, v , to b is known to have been constructed using a tail $C(y)$, then we can replace the call $b[v]$ with a call of the form $b'[y]$. Here, b' is a new block that begins with a statement that recomputes $v \leftarrow C(y)$ and then proceeds in the same way as the original block b . In theory, this could result in an ‘optimized’ program that actually allocates twice as many $C(y)$ objects as necessary; that would obviously not be a good outcome. In practice, however, this transformation often enables subsequent optimizations both in the place where the original $b[v]$ call appeared (for example, the statement that initialized v may now be dead code) and in the new block, the latter resulting from a new fact, $v=C(y)$, that can be propagated through the code for b' . Note that Rule (18) also extends naturally to calls with multiple parameters and to cases where one or more of those parameters is a known closure; in that case the process of generating a new derived block has much the same effect as specializing a higher-order function to a known function argument.

Rules (19) and (20) deal with situations where a block is called and its result is immediately used as a closure or data value, respectively. The code on the left side of these rewrites essentially *forces* the allocation of a closure or data object in b , just so that value can be returned and then, most likely, discarded after one use. The right sides deal with this by introducing a tail call in the caller and then turning the use of whatever value is produced in to a ‘trailing’ action in the new block. As in other examples, this does not produce an immediate optimization. However, these rules generally lead to useful improvements in practice, enabling new optimizations by bringing the construction and use of v in to the same context.

5.6 Additional Optimizations

Beyond the rewrites kind described in previous sections, our optimizer implements several other program transformations that help to improve code quality. One of the most important of these in practice—because it also performs a strongly-connected components analysis on the program to prioritize the order in which rewrites are applied—is a “tree-shaking” analysis. This automatically removes definitions from a program if they are not reachable from the program’s entry points. In addition to sections of library code that are not used in a given application, this also helps to clean up after other optimizations by eliminating single-use blocks whose definitions have been inlined or blocks that have been replaced with new derived versions. The optimizer also attempts to recognize and merge duplicated definitions, and to rewrite block definitions (and all corresponding uses) to eliminate unused parameters. Unused stored fields in closure definitions can also be eliminated in this way, but we cannot remove arguments in closure definitions: even if they are not used in the code for a particular closure, they must be retained for compatibility with other closures of the same type.

One other detail in our implementation is that we run the MIL type checker after every use of the optimizer. This has two practical benefits: (1) All of our optimizations are required to preserve typing, so running the type checker provides a quick sanity check and may help to detect errors in the optimizer. (2) The type checker will automatically reconstruct type information for each part of the program, so we do not need to deal with those details in the implementations of individual rewrites.

5.7 Reflections on Optimization

The design of an optimizer compiler inevitably requires some judicious compromises. After all, the problem that it is trying to solve—to generate truly optimal versions of any input program—is uncomputable, and so, at some point, it must rely instead on heuristics and incomplete strategies. Even if an optimizer delivers good results on a large set of programs, there is still a possibility that it will perform poorly on others. Subtle interactions between different optimization techniques may prevent the use of key transformations in some situations and instead lead to expanded code size or degraded performance. With those caveats in mind, we have, so far, been very satisfied with the performance of our optimizer for MIL.

As a concrete example, the diagrams in Figure 3 outline the structure of a MIL implementation of the Habit “prioset” example [15, Section 5]. Part (a) here is for the original MIL implementation, generated directly from 56 lines of LC code; it is too small to be readable, but does convey that the original program—910 lines of MIL code—is quite complex. Part (b) shows the result obtained after 1,217 separate rewrite steps in the MIL optimizer, resulting in 140 lines of MIL code. This version of the program still has non-trivial control flow, but its overall structure is much simpler and we can start to see details such as loop structures and a distinction between the blue nodes (representing blocks) and the red nodes (representing closure definitions). By comparison, there are red and blue nodes scattered throughout the diagram in Part (a), which suggests that our optimizations have been effective in eliminating many uses of closures in the original program. (The remaining red nodes in Part (b) are only there because the program exports the definitions of two functions, `insertPriority` and `removePriority` as first-class values; in a program that only uses fully-applied calls to these functions, even those nodes would be eliminated.)

The results that we see with this example are representative of our experience across a range of test programs, and suggest that the MIL optimizer can work well in practice. That said, we plan to do more extensive studies, including performance benchmarking, and to use those results to further tune and refine our implementation. (As a new language, we do not currently have a pre-existing set of benchmarks, but we do hope to grow such a library as our implementation matures and gains users.)

One particularly effective aspect of our implementation that is already clear is the decision to perform optimization at multiple levels throughout the compiler pipeline. An initial use of the MIL optimizer takes care of relatively high-level rewrites, such as inlining of higher-order functions to eliminate the costs of constructing closures. It would be harder to apply this kind of optimization at later stages once the operations for function application and closure construction have been decomposed in to sequences of lower level instructions. A subsequent use of the optimizer, after representation transformations have been applied, enables the compiler to find newly exposed opportunities for optimization, and to further simplify the generated MIL code before it is translated in to LLVM, as described in the next section. Finally, the use of LLVM itself allows us to exploit the considerable effort that has been invested in that platform to perform lower-level optimizations, and—although our focus to date has been on IA32-based systems—also provides a path for targeting other architectures.

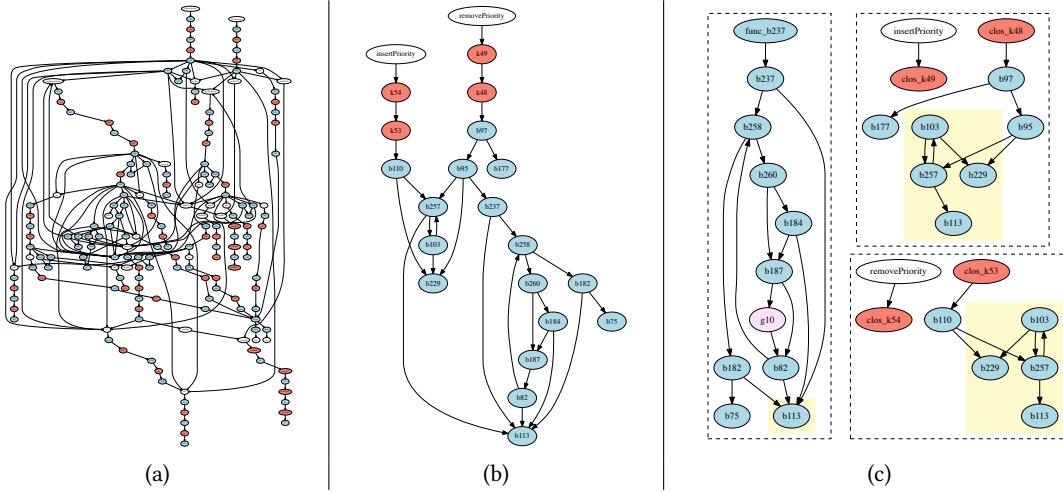


Figure 3: Control flow graph examples.

6 COMPILING MIL TO LLVM

In this section, we explain how MIL programs can be translated into corresponding LLVM programs, which can then be subjected to further optimization and used to generate executable binaries. Beyond the specific practical role that it serves in our Habit compiler, this also provides another test for MIL’s suitability as an intermediate language: it is important, not only that we are able to generate executable programs from our IL, but also that we are able to do so without introducing overhead or undoing any of the improvements that were made as a result of optimizations on the IL.

6.1 Translating MIL Types to LLVM

Before generating LLVM code, we use representation transformations to provide Word-based implementations for bitdata types (Section 4.2) and to eliminate polymorphism and parameterized datatypes (Section 4.3). In the resulting programs, MIL types like Word and Flag are easily mapped to LLVM types such as i32 and i1. The only types that require special attention are for functions (which we cover in this section) and algebraic datatypes (which are handled in a similar manner).

Every MIL value of type $[d_1, \dots, d_m] \rightarrow [r_1, \dots, r_n]$ will be a closure that can be represented by a block of memory that includes a code pointer and provides space, as needed, for stored fields:

entry	...
-------	-----

We can describe structures of this form using three LLVM types with mutually recursive definitions:

```
%clo = type { %fun }
%fun = type {r1, ..., rn} (%ptr, d1, ..., dm)
%ptr = type %clo*
```

Here, %clo is a structure type that describes the layout of the closure. Its only component is the code pointer of type %fun: no additional components are listed because the number and type of fields is a property of individual closure definitions, not the associated function type. In the generated code, functions of this type will be represented by pointers of type %ptr. Given such a pointer, the

implementation can read the code pointer from the start of the closure and invoke the function, passing in the closure pointer and the argument values corresponding to the domain types d_i. If that function needs access to stored fields, then it can cast the %ptr value to a more specific type that reflects the full layout for that specific type of closure. Finally, the function can return a new structure containing values for each of the range types r_j. (If there is only one result, then it can be returned directly, without a structure; if there are no results at all, then we can use a void function.)

The techniques described here are standard, but there are still many details to account for. Among other things, this reinforces the importance of performing closure optimizations in MIL, rather than generating LLVM code directly and then hoping, unrealistically, that the LLVM tools will be able to detect the same opportunities for improvement. Instead, we divide the responsibilities for optimization between MIL and LLVM, with each part making important contributions to the overall quality of generated code.

6.2 Translating MIL Code to LLVM

The process of translating MIL statements to LLVM instructions is relatively straightforward. As examples: an and primitive in MIL maps directly to the (identically named) and instruction in LLVM; a closure allocation in MIL is implemented by a call to a runtime library function to allocate space for the closure, followed by a sequence of store instructions to initialize its fields; and so on. As such, we will not discuss the fine details of this translation here.

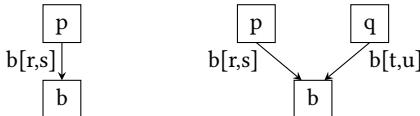
There are, however, some key, higher-level structural mismatches between MIL and LLVM—specifically, parameterization and sharing of blocks—that do need to be addressed. As we have seen, MIL programs are collections of (parameterized) basic blocks that are connected together either by regular block calls or tail calls (in the middle, or at the end, respectively, of a code sequence). By contrast, LLVM programs consist of a collection of (parameterized) functions, each of which has a *control flow graph* (CFG) comprising a single (parameterized) entry point, and a body that is made up from a

collection of (unparameterized) basic blocks. How then should we approach the translation of an arbitrary MIL programs into LLVM?

One approach would be to generate an separate LLVM function for each MIL block (and closure definition). Although LLVM does provide support for tail calls, those features are difficult to use and it might be difficult to ensure that the compiled loops, encoded as tail recursive blocks in MIL, run in constant space.

Our strategy instead is to compile mutually tail recursive blocks directly into loops, accepting that there will be some (small) duplication of blocks in the process. In Figure 3(b), for example, there are several blocks that are reachable from either of the two distinct entry points at the top of the diagram. Our code generator uses some simple heuristics to generate a set of LLVM CFG structures from MIL programs with the following properties: (1) There must be a distinct CFG for every closure definition and for every block that is a program entry point or the target of a non-tail call; (2) If one block is included in a CFG, then all other blocks in the same strongly connected component (SCC) should also be included in the same CFG (this ensures that tail calls can be compiled to jumps); (3) If a single block has multiple entry points from outside its SCC, then it is a candidate entry point for a new CFG (this attempts to reduce duplication of blocks). The result of applying our algorithm to this particular example is shown in Figure 3(c) and is typical of the behavior we see in general: there is some duplication of blocks (the portions highlighted with a yellow background) but the amount of duplicated code is small and has not been a concern in practice.

Our second challenge is in dealing with the mapping from parameterized blocks in MIL to unparameterized blocks in LLVM. It turns out that the number of predecessors is key in determining how to generate code for the body of each block. To understand this, consider the following two diagrams:



In both diagrams, we assume a block b defined by $b[x, y] = c$ for some code sequence c . For the diagram on the left, there is exactly one predecessor, p , which ends with a call to $b[r, s]$. In this situation, there is actually no need for the parameters to b because we already know what values they will take at the only point where b is called. All that it needed is to apply a substitution, replacing the formal parameters, x and y , with the actual parameters r and s , respectively, so that we use the code sequence $[r/x, s/y]c$ for the body of b . (Of course, we also need to account for this substitution on any edges from b to its successors.) For the diagram on the right, there are two predecessors, each of which ends by calling b with (potentially distinct) parameters. In this case, the phi functions that are part of LLVM's SSA representation provide exactly the functionality that we need to ‘merge’ the incoming parameters and we can generate code of the following form for b :

```

x = phi [r,p], [t,q]
y = phi [s,p], [u,q]
... LLVM code for c goes here ...

```

Generating code in SSA form is sometimes considered to be a tricky or complicated step in the construction of an optimizing compiler.

It is fortunate that the translation is relatively straightforward and that we are able to take advantage of phi functions—a key characteristic of the SSA form—quite so directly.

7 CONCLUSIONS

In this paper, we have described the MIL language and toolset, demonstrating (1) that it has the fundamental characteristics needed to qualify as an effective intermediate language for the compilation of functional languages; and (2) that it can enable new techniques for choosing efficient data representations. The design and implementation of any new intermediate language requires considerable engineering effort. As we continue to develop the MIL system ourselves—for example, to explore its potential for compilation of lazy languages—we hope that it will also serve as useful infrastructure for other functional language implementors.

ACKNOWLEDGMENTS

The authors thank the anonymous reviewers for their helpful feedback. This work was supported in part by funding from the National Science Foundation, Award No. CNS-1422979.

REFERENCES

- [1] Andrew W. Appel. 1992. *Compiling with Continuations*. Cambridge University Press, New York, NY, USA.
- [2] Justin Bailey. 2012. *Using Dataflow Optimization Techniques with a Monadic Intermediate Language*. Master’s thesis. Department of Computer Science, Portland State University, Portland, OR.
- [3] Adam Chlipala. 2015. An Optimizing Compiler for a Purely Functional Web-application Language. In *Proceedings of the 20th ACM SIGPLAN International Conference on Functional Programming (ICFP 2015)*. ACM, New York, NY, USA.
- [4] Iavor S. Diatchki, Mark P. Jones, and Rebekah Leslie. 2005. High-level views on low-level representations. In *Proceedings of the 10th ACM SIGPLAN International Conference on Functional Programming, ICFP 2005, Tallinn, Estonia, September 26-28, 2005*. ACM, 168–179.
- [5] Matthew Fluet and Stephen Weeks. 2001. Confication Using Dominators. In *Proceedings of the Sixth ACM SIGPLAN International Conference on Functional Programming (ICFP ’01)*. ACM, New York, NY, USA, 2–13.
- [6] Thomas Johnson. 1985. Lambda Lifting: Transforming Programs to Recursive Equations. In *Proceedings of the IFIP conference on Functional Programming Languages and Computer Architecture (Lecture Notes in Computer Science, 201)*. Springer-Verlag, 190–203.
- [7] Mark P. Jones. 1994. Dictionary-free Overloading by Partial Evaluation. In *ACM SIGPLAN Workshop on Partial Evaluation and Semantics-Based Program Manipulation (PEPM ’94)*.
- [8] Andrew Kennedy. 2007. Compiling with Continuations, Continued. In *Proceedings of the 12th ACM SIGPLAN International Conference on Functional Programming (ICFP ’07)*. ACM, New York, NY, USA, 177–190.
- [9] Chris Lattner. 2002. *LLVM: An Infrastructure for Multi-Stage Optimization*. Master’s thesis. Computer Science Dept., University of Illinois at Urbana-Champaign, Urbana, IL.
- [10] LLVM 2018. The LLVM Compiler Infrastructure. <http://llvm.org>.
- [11] Luke Maurer, Zena Ariola, Paul Downen, and Simon Peyton Jones. 2017. Compiling without continuations. In *ACM Conference on Programming Languages Design and Implementation (PLDI’17)*. ACM, 482–494.
- [12] E. Moggi. 1989. Computational Lambda-calculus and Monads. In *Proceedings of the Fourth Annual Symposium on Logic in Computer Science*. IEEE Press, Piscataway, NJ, USA, 14–23.
- [13] Alan Mycroft. 1984. Polymorphic Type Schemes and Recursive Definitions. In *Proceedings of the 6th Colloquium on International Symposium on Programming*. Springer-Verlag, London, UK, UK, 217–228.
- [14] Simon Peyton Jones (Ed.). 2003. *Haskell 98 Language and Libraries – The Revised Report*. Cambridge University Press.
- [15] The Hasp Project. 2010. The Habit Programming Language: The Revised Preliminary Report. <http://github.com/habit-lang/language-report>.
- [16] Philip Wadler. 1990. Comprehending Monads. In *Proceedings of the 1990 ACM Conference on LISP and Functional Programming (LFP ’90)*. 61–78.
- [17] Stephen Weeks. 2006. Whole-program Compilation in MLton. In *Proceedings of the 2006 Workshop on ML (ML ’06)*. ACM, New York, NY, USA.

Task Oriented Programming and the Internet of Things

Mart Lubbers

Institute for Computing and
Information Sciences
Radboud University
Nijmegen, The Netherlands
mart@cs.ru.nl

Pieter Koopman

Institute for Computing and
Information Sciences
Radboud University
Nijmegen, The Netherlands
pieter@cs.ru.nl

Rinus Plasmeijer

Institute for Computing and
Information Sciences
Radboud University
Nijmegen, The Netherlands
rinus@cs.ru.nl

ABSTRACT

iTasks is an embedded Domain Specific Language (eDSL) implemented in the functional language Clean offering Task Oriented Programming (TOP). It can be used to express workflows that end-users and systems have to execute in collaboration to accomplish a certain goal. One can compile an *iTasks* program to a web-server that coordinates the work described. In this paper we present a novel extension of *iTasks* to control Internet of Things (IoT) devices. IoT devices can dynamically receive and execute IoT tasks, interpretable on the devices' byte-code interpreter, sent from an *iTasks* system. These tasks express workflow within the limitations of IoT devices and are executed as regular *iTasks* tasks. Shared Data Sources (SDSs) are the de facto way of in *iTasks* to provide access to resources in the outside world. Examples of SDSs are files, time, shared memory or a random stream. SDSs in *iTasks* are used to communicate with the IoT devices and their tasks in a transparent way. This hides specific knowledge of the device and the underlying communication technique. The advantage of this approach is that all IoT devices can be programmed on a high abstraction level from a single source. The IoT devices only have to be programmed once – even when repurposing – and the interaction is communication method agnostic and transparent. Hence, it avoids the integration problems that are common in developing IoT applications.

CCS CONCEPTS

- Computer systems organization → Distributed architectures;
- Software and its engineering → Client-server architectures; Functional languages; Domain specific languages;

KEYWORDS

Internet of Things, Functional Programming, Distributed Applications, Task Oriented Programming, Clean

ACM Reference Format:

Mart Lubbers, Pieter Koopman, and Rinus Plasmeijer. 2010. Task Oriented Programming and the Internet of Things. In *Proceedings of International Symposium on Implementation and Application of Functional Languages (IFL '18)*. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

IFL '18, August 2019, Lowell, MA, USA

© 2010 Association for Computing Machinery.
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00
<https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

1 INTRODUCTION

1.1 Internet of Things

The term Internet of Things (IoT) stands for a whole network of (smart) devices that interact with each other and – most of all – interact with the world. Gartner estimated that there will be around 21 billion IoT devices online in 2020¹. IoT devices are already entering households today in the form of smart electricity meters, thermostats, weather stations and so on. Other industries such as healthcare, are also adopting IoT technology. For example, wearable IoT devices are used to monitor patients at home [Islam et al. 2015]. IoT technology is emerging rapidly and is transforming the way people interact with technology and with each other.

Although IoT devices are driven by software, there are typically severe limitations on the processing power on IoT devices. Because they must be cheap, very tiny computers are used with limited memory, e.g. Microcontroller Units (MCUs). The programs are usually stored in flash memory that can only handle a very limited number of write cycles, e.g. 10,000 for the MCU in an Arduino Uno. IoT devices communicate with the internet to share information such as sensor data and to act on demand with actuators [Da Xu et al. 2014]. They track information accurately while using as little power, bandwidth, and memory as possible.

IoT not only encompasses the communication between the devices, it encompasses the system of systems as a whole. The architecture of IoT systems is often divided into layers. Several layered structures have been proposed ranging from three to five layers. For the purpose of this paper we use the three layered architecture.

The first layer is the perception layer and contains the IoT devices with their sensors and actuators. E.g. in a smart electricity meter, this layer would contain the sensors detecting the current drawn and the MCU reading these sensors. There are myriads of different devices available on the market to use in this layer. They are programmed using a variety of programming languages, high or low level, compiled or interpreted.

The second layer is the networking layer and connects the first layer with the third layer. In a smart electricity meter, this is the connection between the meter and the server. Existing general purpose networking techniques are often used to transmit the data. Recently, specialized IoT networks – so called Low Power Low Throughput Networks (LTNs) – are employed as well.

Finally, the third layer – the application layer – provides the end-point for the user to interact with the IoT devices and their data. In a smart electricity meter, this layer is the server aggregating the data and the frontends such as apps.

¹Gartner (November 2015)

All of this seems to be well thought through, but there are major problems with developing entire IoT systems. There is an impedance mismatch between layers which leads to isolated logic for the components and it leads to integration problems. Every layer has to be programmed separately in different languages and on different platforms resulting in varying update roll-out costs. For example, rolling out updates in the sensing layer is expensive since reprogramming MCUs in the field often requires physical access, while updating an app in the application layer is as simple as deploying an updated app.

Reprogramming devices on the fly can be beneficial in systems in which a device fails unexpectedly or is repurposed a lot in a more dynamic system. However – in compiled programs – deploying a different program on a device requires a complete reprogramming.

Interpretation can mitigate this limitation but comes with downsides as well. Sending serialized general purpose programs causes overhead and they have to be stored in memory.

1.2 TOP – Task Oriented Programming

The Task Oriented Programming (TOP) paradigm and the corresponding *iTasks* implementation offer a high abstraction level for defining real world workflow tasks [Plasmeijer et al. 2007]. These tasks are described through an embedded Domain Specific Language (eDSL) hosted in the purely functional programming language Clean [Brus et al. 1987; Plasmeijer et al. 2011]. Tasks are the basic building blocks of the language; they resemble actual work that needs to be done. The language contains combinators – arising from workflow modelling – to combine tasks in a sequential, parallel or conditional way. The *iTasks* system generates a multi-user web application to coordinate the tasks that the end users and the computer systems have to do in collaboration.

The *iTasks* eDSL is highly type-driven and built on generic functions that are created on the fly for the given types. These generic functions provide the basic TOP functionality which means the programmer has to do little to no implementation work on details such as the user interface. If needed, functions can be specialized to offer fine-grained control over the TOP functionality.

The *iTasks* system has shown to be useful in many fields of operation such as incident management [Lijnse 2013].

1.3 Integrating IoT devices with TOP

With TOP, one can describe arbitrary complex collaboration between distributed acting end users and systems without the need to worry about the technical details. In this paper we show how to program all layers of IoT from one single source and thus incorporate IoT devices seamlessly in TOP/*iTasks*. Adding IoT devices to the current *iTasks* system is difficult as it was not designed to cope with devices which are that tiny and that restricted in their communication bandwidth.

A natural way to incorporate new hardware is to use the distributed *iTasks* extension as explained in [Oortgiese et al. 2017]. They lifted *iTasks* from a single server model to a distributed server architecture. For example, Android apps can be created that run an entire *iTasks* core and are able to receive tasks from a different server and execute them. While Android often runs on small ARM

devices, they are a lot more powerful than the average IoT MCU. Although their system is suitable for dynamically sending tasks over platforms with different types of processors, their solution cannot be ported to MCUs because IoT devices are simply not powerful enough to run or store an *iTasks* core. Moreover, sending serialized *iTasks* tasks over an LTN requires too much bandwidth usage.

1.4 Research Contribution

In this paper we present a novel way of controlling IoT devices in TOP/*iTasks* using restricted tasks for IoT devices and special interfaces to Shared Data Sources (SDSs). Our research contributions are:

(1) The restricted tasks run on IoT devices expressed in an extended version of the type safe class-based shallow embedded mTask-eDSL by Koopman and Plasmeijer [2016]. (2) A novel backend for the eDSL generates specialized bytecode programs. Devices are initialized once with a runtime system to dynamically receive and execute this bytecode. This results in an ecosystem where devices can be repurposed without needing physical access. (3) To port a device, only the device specific code needs to be ported. The bulk can be compiled from a shared Runtime System (RTS) written in C. (4) Functionality is added to *iTasks* to integrate the IoT devices through familiar *iTasks* concepts. IoT tasks can be executed as if they were regular *iTasks* tasks and communication with these tasks is achieved via SDSs. Device- or communication-method specific information is hidden for the programmer and user.

1.5 Structure of this Paper

In Section 2, the basic concepts of TOP as offered by the *iTasks* system are introduced together with an IoT application that we use as a running example. Section 3 explains the eDSL techniques and the actual eDSL used to express IoT tasks. The glue needed for the interaction between *iTasks* and IoT devices is discussed in Section 4. Moreover, the example from Section 2 is finished to illustrate the process of building IoT applications. Section 5 shows the bytecode compilation backend for the eDSL to dynamically generate code that can be executed on IoT devices. Section 6 covers the details of the run-time system for the IoT devices. In Section 7, the server-side implementation and integration is discussed, Section 8 describes related work and Section 9 concludes with the conclusion and discussions.

2 BRIEF OVERVIEW OF TOP

Here we present a brief overview of the main concepts of TOP. The details are specific to the TOP implementation *iTasks*. More detailed information can be found in [Plasmeijer et al. 2012].

2.1 Tasks

A Task `a` is a statefull event processor that returns a value of type `:: TaskValue a = NoValue | Value a Bool` in which the `Bool` represents the stability. A `TaskValue` is special since it may change over time because its event handling function is re-evaluated on every event. Once a value is `Stable`, the does not change again. A `TaskValue` can be observed by other tasks and it can affect which other tasks can be started. There are basic tasks and combinators to compose tasks in familiar workflow patterns.

Basic tasks come in two flavours: interactive and non-interactive. Non-interactive basic tasks consist of processing, task value manipulation, external connections, server interaction and communication with the host system.

Interactive tasks provide interaction with a user via a type driven generated web interface. A type used in *iTasks* must have instances for a collection of generic functions that is captured in the class *iTask*. One of these generic functions is the generic web form generation that allows the user to edit a value of that type using the web browser. Basic types have specialization instances for these generic functions and they can be derived for any first-order user-defined type. When desired, derived interfaces can be fine-tuned or specialized instances can be defined.

Some interactive tasks for entering, viewing and updating values of arbitrary types are shown in Listing 1. These functions spawn a web form for the user to work with. The first argument is the title, the second argument contains the display options on the data and the third argument of the view and update variant are the initial value.

```
enterInformation :: String [EnterOption m] → Task m | iTask m
viewInformation :: String [ViewOption m] m → Task m | iTask m
updateInformation :: String [UpdateOption m m] m → Task m | iTask m
```

Listing 1: Interactive tasks

2.2 Task Combinators

Tasks can be combined with monad-*esque* task combinators to express sequential, parallel and conditional workflows. With task combinators one can define how tasks depend on each other, and how the information is passed between them. The resulting combination delivers a new task. Some of the – for this paper – relevant combinators are explained below.

The parallel combinators in Listing 2 combine two or more tasks in such a way that they are offered to the user at the same time, possibly combining the result.

```
(||) infixr 3 :: (Task a) (Task b) → Task b | iTask a & iTask b
(||) infixr 3 :: (Task a) (Task b) → Task a | iTask a & iTask b
(||) infixr 3 :: (Task a) (Task a) → Task a | iTask a & iTask b
(&&) infixr 4 :: (Task a) (Task b) → Task (a, b) | iTask a & iTask b
allTasks :: [Task a] → Task [a] | iTask a
```

Listing 2: Parallel combinators

Listing 3 shows the sequential task combinators. The value of the left-hand side is fed to the right-hand side if one of the *TaskCont* predicates hold. These predicates can be based on the stability of the value, the actual value, an action or an exception. Actions are presented to the user as buttons in the generated web form. Exceptions are thrown by tasks and can also be caught using the *try* construction. The bind combinator ($\gg=$) is implemented as a step that continues only when the left-hand side is stable or the user presses the continue button.

```
(>>) infixl 1 :: (Task a) (a → Task b) → Task b | iTask a & iTask b
(>>) infixl 1 :: (Task a) (Task b) → Task b | iTask a & iTask b
(>>) infixl 1 :: (Task a) [TaskCont a (Task b)] → Task b
| iTask a & iTask b
```

```
:: TaskCont a b
=    OnValue      ((TaskValue a) → Maybe b)
|    OnAction String ((TaskValue a) → Maybe b)
```

```
| ∃e: OnException (e → b) & iTask e
try :: (Task a) (e → Task a) → Task a | iTask a & iTask e & toString e
```

Listing 3: Sequential combinators

2.3 Shared Data Sources

The TOP way of sharing information between tasks are SDSs. SDSs are an abstraction over resources in the broadest sense. SDSs are solely defined by their stateful read and write functions. For example, an SDS can be a file on disk, the system time, a place in memory, lenses on other SDSs, or an external database. There is a publish-subscribe system attached to SDSs which means that a task reading an SDS is automatically notified when the value has changed. This results in low resource usage because tasks do not need to poll. Lenses on SDSs can be used to map functions, combine multiple SDSs, or apply data or notification filters. An SDS is typed by three type variables; the read, the write type and the parametric lens. The parametric lens is ignored for and it is fixed to () in all access functions. However, later it is used and explained in Section 7.4.

There are four atomic functions in the task domain to interact with SDSs (Listing 4). The *get* function retrieves the value, the *set* value sets the value and the *upd* changes the value. Finally the *watch* function is a task that constantly returns the value of the SDS when it is changed by someone.

```
:: SDS p r w
:: Shared a ::= SDS () a a

get :: (SDS () r w) → Task r | iTask r
set :: w      (SDS () r w) → Task w | iTask w
upd :: (r → w) (SDS () r w) → Task w | iTask r & iTask w
watch :: (SDS () r w) → Task r | iTask r
```

Listing 4: SDS functions

Moreover, the interactive *Information (see Section 2.1) tasks are available for SDSs as well. In this way, one can interact with shared data using the generated web forms. Listing 5 shows the type signatures for these functions. The *view* on the SDS data in the web page is automatically updated when the SDS is updated. These functions support the same lenses as their counterparts.

```
viewSharedInformation :: String [...] (SDS r w) → Task r | iTask r
updateSharedInformation :: String [...] (SDS a a) → Task a | iTask a
```

Listing 5: SDS interaction

There are some extra functions available in the realm of SDSs that need some introduction (Listing 6). The $\gg*$ operator combines two SDSs. The *mapRead* function embeds a transformation function atomically in the given SDS. Derived from the *watch* function is the *whileUnchanged* function that – given an SDS and a task – executes the task every time the SDS value changes.

A way to creating a SDS is by using the *withShared* function. This function creates a memory mapped SDS that is only available within scope.

```
(>>*) infixl 6 :: (SDS rx wx) (SDS ry wy) → SDS (rx, ry) (wx, wy)
mapRead       :: (SDS r w) (r → r') → (SDS r' w)
whileUnchanged :: (SDS r w) (r → Task b) → Task b | iTask b
withShared     :: a ((SDS a a) → Task b) → Task b | iTask b
```

Listing 6: Extra SDS functions

2.4 Thermostat Example

As an illustrative running example we introduce a thermostat application written in *iTasks*.

The *iTasks* system provides the interaction with the data sources to the user. All of the communication is going through these data sources which are represented by SDSs. One SDS is created for the current temperature and another one for the target temperature. Moreover, virtual SDSs that represent the on/off state of the heater and the cooler. The current temperature and the corresponding target are combined to virtual Shared Bool. The user can set the target temperatures through the generated web interface. If the temperature drops below the lower limit, the heater should turn on. If it rises over the upper limit, the fan should turn on.

Now suppose we are also able to execute tasks on an IoT device and suppose we have all the actual device interaction captured in the *iotThermostat* task (Section 4.3). Then we could program a thermostat as follows.

```

1 thermostat :: Task Int
2 thermostat =
3   withShared 0      @currentTemp →
4   withShared (18, 22) @targetTemp →
5   let cooler = mapRead (uncurry (>)) (currentTemp >&< mapRead snd targetTemp)
6     heater = mapRead (uncurry (<)) (currentTemp >&< mapRead fst targetTemp)
7   in viewSharedInformation "Current" [viewAs viewAsCelcius] currentTemp
8   -|| updateSharedInformation "Lower limit Upper limit" [] targetTemp
9   -|| viewSharedInformation "Cooler" [] cooler
10  -|| viewSharedInformation "Heater" [] heater
11  -|| iotThermostat currentTemp cooler heater)
12 where
13   viewAsCelcius s = toString s ++ "°C"

```

Listing 7: Thermostat, the TOP part

Lines 3-4 instantiates SDSs in memory representing the current temperature and the target. The initial value for the current temperature (type Shared Int) is irrelevant and set to 0. The initial target temperature lies between 18 and 22 degrees and has type Shared (Int, Int). All communication between the tasks goes via these SDSs.

Lines 7-11 are the tasks that generates the user interface. The interface depends on the type of the SDS. Hence, it is used to change the target temperatures, view the current temperature and view the status of the cooler and the heater. This code – modulo some extra display options – results in the interface in Figure 1. A display option (line 13) is applied on the temperature SDS to pretty print it with the correct unit.



Figure 1: Thermostat user interface

Line 11 contains the IoT logic which is defined in Section 4.3. It is integrated in a task that uses the SDSs to operate the thermostat. The first SDS is the current temperature, the second and third SDS for which lenses on the temperature and the limit are used (lines 5-6). A lot of things happen here, the task it returns compiles to bytecode and executes this code for IoT tasks that run on IoT device.

The devices measures the temperature and control the cooler and the heater and communicate via the SDSs. The cooler and heater status are represented as a lens on the temperature and the limit SDS (lines 5-6).

3 AN eDSL FOR IoT TASKS

Regular *iTasks* tasks are not suitable to run on small devices because of their resource usage. However, a subset of the TOP tasks are natural to the IoT domain and we still want to express them in a type safe and extendable way. EDSLs offer a solution for creating new languages in a host language while benefiting from properties of the host language such as the type system.

3.1 Class-Based Shallow Embedding

There are several basic embedding techniques, such as shallow and deep embedding [Gibbons 2015]. Class-based shallow embedding – or tagless embedding – has the advantages of both shallow and deep embedding [Carette et al. 2009; Svennbergsson and Axelsson 2012]. Here, the language constructs are defined as type classes and a backend is a type with an instance of some of the classes. This means that adding backends is easy and a backend only needs to implement the classes it needs. Moreover, type safety is guaranteed because the types can contain phantom types and constraints can be enforced by the type signatures of the class functions. Lastly, extensions can be added easily. Existing backends do not need to be updated when an extension is added in a new class. Naturally, if the extension is added in an existing class, the backends implementing the class need to be updated.

3.2 IoT eDSL

The mTask eDSL is a class-based shallowly eDSL hosted by Clean created by Koopman and Plasmeijer [2016]. Their backend generates C-code for complete TOP-like programs that can run on an Arduino. The language itself is imperative of nature and programs written in it are suited to run on an MCU. However, this backend generates a self contained system and is not suitable for our purpose because there is no connection whatsoever with the regular *iTasks* system.

In this paper, the mTask eDSL is extended with a new backend for bytecode generation and language extensions that allow run-time assignment of tiny IoT tasks to IoT devices as well as integration of these IoT tasks with regular *iTasks* tasks. To avoid confusion we will call the extended mTask eDSL with the novel backend the IoT eDSL.

The IoT eDSL is a collection of classes implementable by types with two type variables. The type implementing the classes is called the backend. The first type variable (*t*) represents the type of the construction and the second type variable (*r*) the role of the construction. Type constraints are used to make sure the expressions are well typed and to disallow expressions like: lit True +. lit 1. Roles can be Expr, Stmt and Upd to denote expressions, statements and updatables. The roles form a hierarchy that is expressed in class constraints, for example, an updatable can be used as an expression but not the other way around. This hierarchy is shown in Listing 8.

:: Expr = Expr

```
:: Stmt = Stmt
:: Upd = Upd

class isExpr a :: a
instance isExpr Upd, Expr
class isStmt a :: a
instance isStmt Upd, Expr, Stmt
```

Listing 8: Role hierarchy

The constructions in the IoT language can be grouped by roles and extra categories for device access and SDS operations.

3.2.1 Expressions. There are two classes of expressions, namely boolean expressions and arithmetic expressions. The class of arithmetic language constructs also contains the function `lit` that lifts a host-language value into the eDSL domain. All operators are suffixed with a full stop to avoid name clashes with Clean's built-in operators. All standard arithmetic functions are included in the eDSL, but some are omitted for brevity. The boolean and arithmetic expression classes are shown in Listing 9.

```
class arith b where
    lit      :: t          → b t Expr
    (+.) infixl 6 :: (b t r) (b t q) → b t Expr | + t & isExpr r & isExpr q
    ...
class boolExpr b where
    Not     :: (b Bool r) → b Bool Expr | isExpr r
    ...
(=.) infix 4 :: (b a r) (b a q) → b Bool Expr | == a & isExpr r & isExpr q
```

Listing 9: Basic classes for expressions

3.2.2 Statements. There is no loop control needed because IoT tasks are scheduled for execution continuously by the RTS and we only need conditional and sequential statements. Both the sequence operator (`(.)`) and the conditional (IF, `?`) statements are shown in Listing 10. The `?` is a variant of the standard conditional operation where the else clause is empty.

```
class IF b where
    IF      :: (b Bool p) (b t q) (b s r) → v () Stmt | isExpr p
    (?) infix 1 :: (b Bool r) (b t q)           → v () t | isExpr p
class seq b where
    (.) infixr 0 :: (b t r) (b u q)           → b u Stmt
```

Listing 10: Control flow operators

3.2.3 Assignables. The IoT eDSL offers an imperative language and therefore an assignment construction is very natural. Only constructs with the `Upd` role can be assigned to. Examples of constructs with the `Upd` role are variables and General Purpose Input/Output (GPIO) pins.

```
:: DigitalPin = D0 | D1 | D2 | D3 | D4 | D5 | ...
class dIO b :: DigitalPin → b Bool Upd
class assign b where
    (=.) infixr 2 :: (b t Upd) (b t r) → b t Stmt | isExpr r
```

Listing 11: Input/output classes

3.2.4 Variables. Imperative languages have variables that can be assigned to and read from. The classes associated with variables are listed in Listing 12. Variables — and other future decorations — may not be declared inside an expression and therefore the expression must be boxed in the `Main` type. The type signature is complex;

to illustrate the usage, an implementation example for a variable written to an analog pin is given.

```
:: In a b = In infix 0 a b
:: Main a   = {main :: a}
class var b where
    var :: ((b t Upd) → In t (Main (b c r))) → (Main (b c r)) | ...

writeAnalog :: Main (b Bool Stmt)
writeAnalog = var λx=True In {main = dIO D3 =. x}
```

Listing 12: SDSs in the IoT language

3.2.5 Analogy. The `Main` type is an analogy to the `main` function in a C program or the `loop` function in an Arduino program. For example the following IoT task — albeit ugly, but this is the way microcontrollers are programmed — that calculates the factorial of 6 (Listing 13) roughly translates to pseudo C code in Listing 14.

```
var λx=1 In
var λy=6 In
{main =
  (y== 0) ?
    (x= y *. x:.
     y= y -. lit 1)
  }
}
```

Listing 13: IoT task

```
int x = 1;
int y = 6;
void loop(void) {
    if (y == 0){
        return;
    } else {
        x = y * x;
        y = y - 1;
    }
}
```

Listing 14: Pseudo C

4 THE GLUE BETWEEN iTASKS AND IoT TASKS

With a language to express IoT tasks we still need glue to actually put the bytecode on the device at run-time for execution as well as a way to integrate them with *iTasks*. These glue functions and tasks can be divided into three categories, namely *devices*, *tasks* and *SDSs*.

4.1 Glue Functions and Tasks

We introduce the `withDevice` task that is needed to interact with an IoT device (Listing 15). This task — given a specification — connects to a specified device and set it up for usage with the *iTasks* system. The task only requires a specification that implements the `Duplex` class. The `Duplex` class' only member is the *iTasks* task that, given the specification synchronizes the communication channels. Implementations of this class have been made for TCP and TTY/Serial devices. When the device is connected, the further interaction is communication agnostic. The `Device` is an abstract type representing an IoT device and passed to functions interacting with devices. If the programmer wants to use another type of device they only need to change the value they pass to `withDevice`.

```
withDevice :: a (Device → Task b) → Task b | Duplex, iTask a & iTask b
class Duplex a where synFun :: a Device → Task ()
```

Listing 15: Device glue

IoT tasks that are expressed in the IoT eDSL have to be compiled to bytecode first. Next, they are sent to the device for execution. All of this is captured in the `liftIOTTTask` function (Listing 16). This function compiles the IoT task, send it to the device and handle the communication. When the IoT task terminates, the lifted task

becomes stable. In *iTasks*, there is no hard limit in the number of tasks that are assigned to the same system. Hence, multiple IoT tasks that are sent to the same device, are executed after each other. The scheduling of the IoT tasks is done by the RTS on the device, hence the lack of looping functionality in the IoT tasks. An IoT task is always accompanied by a scheduling strategy that is either a oneshot execution or a repeated execution (see Section 4.2).

```
:: Interval = OneShot | OnInterval Int
liftIOTTTask :: (Device, Interval) (Main (ByteCode a Stmt)) → Task ()
```

Listing 16: Task glue

All interaction of the *iTasks* system with the running IoT tasks happens via SDSs. To accommodate this, a class has been added that looks similar to the var class that allows *iTasks* SDS to be used in IoT tasks. SDSs in the *iTasks* world are automatically published to all readers when it is written. Applying this strategy to IoT SDSs could cause a large communication overhead. To mitigate this overhead, a lifted SDS that is written on the server is passed on to the device immediately, but a device writing a SDS needs to publish this explicitly using the added pub class.

```
class sds where
  sds :: ((v t Upd) → In (SDS t t) (Main (v c s))) → (Main (v c s))
    | IOTType t
class pub v where
  pub :: (v t Upd) → v t Stmt
```

Listing 17: SDS glue

4.2 Scheduling

Tasks sent to an IoT device are accompanied by a scheduling strategy. With this strategy they behave like TOP tasks in the sense that they are continuously executed and their values can be observed, albeit through SDSs. The following scheduling strategies are available:

- The OneShot strategy can be used to execute a task only once. In IoT applications, often the status of a peripheral or system has to be queried only once on the request of the user. For example, a pin might be read every hour but on request of the user it is queried an extra time.

However, the user might want to know the temperature at that exact moment, and then they can just send a OneShot task probing the temperature. If the temperature sensor is connected to GPIO analog pin 7, such a task looks like

```
someShare :: Shared Pin
someShare = ...
```

```
liftIOTTTask (dev, OneShot) $ sds λpin=someShare
In {main = pin =. aIO A7 :. pub pin}
```

- Tasks accompanied with the OnInterval Int strategy are executed every given number of milliseconds. This strategy most closely resembles tasks as in the *iTasks* system and fits the use case of periodic measurements. The task shown previously can be used to measure the temperature constantly. Moreover, the strategy can be (ab)used to simulate recursion because variables and SDS are kept between executions on the device. The execution can be terminated by executing

the return instruction. The following Listing shows this with the factorial function.

```
liftIOTTTask (dev, OnInterval 500) (fac 6)
where
  fac n = var λresult=1 In var λy=n In {main =
    IF (y ==. 0) return (result =. y *. result :. y =. y -. lit 1)}
```

4.3 Thermostat Example (2)

Now that we have provided all the tooling, we can finish the example from Section 2.4. It is possible to control the task with only one device but to make the example a bit more interesting, we divide the work over two devices: sensorDevice and coolerDevice.

The first device – connected through TCP – measures the temperature and operates the heater. The temperature is read from analog GPIO pin A0 and written in the currentTemp SDS. The heater is connected to digital GPIO pin D1 and is set according to the value in the heater SDS.

The second device – connected via a serial connection – operates the cooling fan. The cooler is connected to digital GPIO pin D5 and the state of the cooler is read from the cooler SDS. The sensor is executed every 500 milliseconds and the cooler and heater IoT tasks every 1000 milliseconds.

```
1  iotThermostat :: (Shared Int) (Shared Bool) (Shared Bool) → Task ()
2  iotThermostat currentTemp cooler heater
3    = readTempHeat -||| operateCooler
4  where
5    sensorDevice :: TCPSettings
6    sensorDevice = {host="192.168.0.12", port=8888}
7
8    readTempHeat :: Task ()
9    readTempHeat = withDevice sensorDevice λsensor→
10      liftIOTTTask (sensor, OnInterval 500)
11        (sds λx=currentTemp In {main=x =. analogRead A0 :. pub x})
12      -|||
13        liftIOTTTask (sensor, OnInterval 1000)
14          (sds λf=heater In {main=digitalWrite D1 f})
15
16    coolerDevice :: TTYSettings
17    coolerDevice = {devicePath="/dev/ttyUSB0", baudrate=B9600, ...}
18
19    operateCooler :: Task ()
20    operateCooler = withDevice coolerDevice λcoolerOper→
21      liftIOTTTask (coolerOper, OnInterval 1000)
22        (sds λf=cooler In {main=
23          digitalWrite LED1 f :. digitalWrite D5 f})
```

Listing 18: Thermostat, the IoT part

Line 3 is the parallel combination of the two tasks representing the work that needs to be done on each device.

Lines 8-14 show the work that is done on the TCP device. First the device is instantiated with the withDevice task. Then, two IoT tasks are sent to the device that run in parallel. The first IoT task reads and publishes the current temperature and the second operates the heater.

Lines 19-23 describe the work that is done on the cooler device. The IoT task operates the cooler through GPIO pin D5 and shows the status on LED1.

5 COMPILING IoT TASKS

An IoT task is compiled using a multi step process before it is sent to the device for execution.

First, the class functions from the IoT eDSL are implemented for the ByteCode type. This type is a boxed Reader Writer State Transformer (RWST)[Jones 1995] that transforms a compiler state while writing bytecode instructions when evaluated.

```
:: ByteCode t r = BC (RWS () [BC] BCState ())
```

Listing 19: The IoT task eDSL backend type

Secondly, an IoT task is not only defined by its bytecode but also by its SDSs and variables. The RWST generates the appropriate bytecode but it also keeps track of the used SDSs and variables in the state. Furthermore, the state contains fresh identifier streams and the state is kept between compilations to not have common identifiers between tasks.

Finally, all the aforementioned data must be converted to messages that the device can understand. To keep the communication overhead small and the execution fast, the bytecode is assembled. The assembly consists of converting the instructions to bytes and resolving the labels. A device receives two types of messages for an IoT task. First, it receives the specification for all the SDSs and variables. Then, it receives the task containing the bytecode.

5.1 Instruction Set and Representation

The instruction set is defined by the BC type and contains basic instructions for a stack machine such as stack operations, labels, jumping and arithmetics (Listing 20). IoT specific instructions are included to allow interaction with peripherals, variables and SDSs. There is no typing on the instruction level, all types are boxed in the BCValue type. This box can contain any type for which the IOTTType class is defined. The context restriction contains all functions needed to interact with the values (e.g. serialization and de-serialization) to send them to the device and the iTask class to interact with them via a web interface. Instances for the serialization classes are given for the basic types, IoT specific types – e.g. GPIO pins and LEDs – and for the box type itself.

```
:: BC
= BCLab Int      | BCJmp Int      | BCJmpT Int | BCJmpF Int
| BCPop          | BCPush BCValue
| BCAdd          | BCMult         | ...
| BCSdsStore Int | BCSdsFetch Int | BCSdsPub Int
| BCAnalogRead AnalogPin | ...
:: BCValue = $\exists$ : BCValue e & IOTTType e
class toByteCode a :: a  $\rightarrow$  String
class fromByteCode a :: String  $\rightarrow$  (Either String (Maybe a), String)
class IOTTType a | toByteCode, fromByteCode, iTask a
```

Listing 20: IoT task SDSs and variables

Generating the stack machine bytecode is straightforward for the basic imperative language constructs. Listing 21 shows partial implementations for the arithmetic and conditional classes are given for illustration.

```
//From MonadWriter
tell :: w  $\rightarrow$  RWST r w s m ()

instance arith ByteCode where
  lit x           = BC $ tell [BCPush (BCValue x)]
  (+.) (BC x) (BC y) = BC $ x  $\gg$  tell [BCAdd]
```

...

```
instance IF ByteCode where
  (?) b t = ...
  IF (BC b) (BC t) (BC e) = BC $
    freshLabel  $\gg=\lambda$  else  $\rightarrow$  freshLabel  $\gg=\lambda$  endif  $\rightarrow$ 
    b  $\gg$  tell [BCJmpF else]  $\gg$ 
    t  $\gg$  tell [BCJmp endif, BCLab else]  $\gg$ 
    e  $\gg$  tell [BCLab endif]

//Fetch a label from the state
freshLabel :: RWS () [BC] BCState Int
```

Listing 21: Instance for the arithmetic and conditional class

5.2 Shared Data Sources

SDSs and variables used in an IoT task are related concepts. They are represented differently in the compiler and in the *iTasks* system but on the device they are the same thing. They are both stored in a list of BCShares that is stored in the compiler's state (Listing 22). In the compilation process, the method of getting the initial value is different and the sds are synchronized with the referenced *iTasks* SDSs on execution. For a var, the initial value is available but for an sds this initial value must be queried using the get *iTasks* task. A BCShare consists of a device unique identifier and either the initial value or the *iTasks* reference.

```
:: BCShare = {sdsi :: Int, sdsval :: Either BCValue (Shared BCValue)}
```

Listing 22: IoT task SDSs and variables

An sds definition is always of the form sds λ x=someShare In {main = ...}. The compiler adds a BCShare to the list and the lambda variable – named x in this case – is the RWST writing the BCSdsFetch instruction with the identifier embedded. The BCShare is initialized with the default value in the var case. The sds case requires some more work because IoT SDSs in a BCShare record are not typed anymore by the Clean type system in the compiler state but boxed in the BCValue type. Therefore, a lens on the linked *iTasks* SDS is created to map the original type to the box and the other way around. This is a potentially unsafe cast, but the BCValue SDS is not accessible from the outside. It is used to process publications coming from the device. The device cannot change the type and therefore this is safe.

```
mapReadWriteError :: (r  $\rightarrow$  MaybeError String r`  

, w` r  $\rightarrow$  MaybeError String (Maybe w)) (SDS r w)  $\rightarrow$  SDS r` w`  
  

lens :: (SDS t t)  $\rightarrow$  SDS BCValue BCValue | IOTTType t  

lens s = mapReadWriteError  

  ( $\lambda$  t  $\rightarrow$  Ok (BCValue t)  

,  $\lambda$  (BCValue v) t  $\rightarrow$  case fromByteCode (toByteCode v) of  

  (Right (Just t), "") = Ok (Just t)  

  _ = Error "Mismatch in BCValue type"  

) v
```

Listing 23: BCValue lens

5.3 Assignables

Assignables – such as the dIO construct – result in an RWST writing a fetch instruction. However, if it occurs in the left hand side of an assignment, this needs to be rewritten to the store counterpart. To achieve this, the censor function from the Writer monad is used to rewrite the instruction accordingly. This technique is applied for

all assignables and the technique is used for the pub function as well to transform the BCSdsFetch instruction to a BCSdsPublish instruction.

```
instance dIO ByteCode where
    dIO (BC p) = BC $ tell [BCDigitalRead p]

instance assign ByteCode where
    (=>) (BC v) (BC e) = BC (e >>| censor makeStore v)

makeStore [BCSdsFetch i] = [BCSdsStore i]
makeStore [BCDigitalRead i] = [BCDigitalWrite i]
makeStore [...] = [...]
```

Listing 24: Bytecode instance for assignment

5.4 Compilation Example

Using the bytecode backend, a factorial program such as in Listing 25 is compiled to the bytecode given in Listing 26. The bytecode is numbered with the program memory offset. The labels are already resolved to actual addresses. The compiler state returning includes two BCShare values, the second – identifier 2 initialized with the var construct – has the initial value 1 :: Int. The first – identifier 1 initialized with the sds construct – has no initial value because it is a referenced *iTasks* SDS. The initial value is retrieved from the *iTasks* system upon sending the IoT task to the device.

```
fac resultShare n = sds λresult=resultShare In var λy=n In {main =
  IF (y==, 0) return (result=, y *, result:, y=, y -, lit 1)}
```

Listing 25: The factorial IoT task

1. BCSdsFetch 1
2. BCPush (1 :: Int)
3. BCLeq
4. BCJmpF 17
5. BCSdsPublish 2
6. BCReturn
7. BCJmp 38
8. BCSdsFetch 2
9. BCSdsFetch 1
10. BCMul
11. BCSdsStore 2
12. BCSdsFetch 1
13. BCPush (1 :: Int)
14. BCSub
15. BCSdsStore 1
16. End of program

Listing 26: Bytecode for the factorial IoT task

6 RUNTIME SYSTEM

The RTS is a single program that has to be written once on a concrete device for it to work with the new extension. It is responsible for the task scheduling, task execution, communication handling and memory management.

6.1 Interface

The RTS is written in C and only uses standard C functions such that the same code is used for all devices. All device specific functions are hidden in a single header file called the interface. It contains functions for accessing device specific peripherals, communication functions, setup and tear down and it defines the specification. The interface header file is created in a modular way using conditional

macros. This means that – similar to the eDSL – some parts of the interface are optional. Every device has to implement the communication functions but peripherals are optional. The specification of the device is generated from the implemented functions to tell the server upon startup what its capabilities are. Therefore, porting the RTS to a new device only requires implementing the interface. The device specific interface is very simple. For example, all communication happens only via three communication functions.

```
bool   input_available(void);
uint8_t read_byte(void);
void   write_byte(uint8_t b);
```

Listing 27: Device interface communication functions

Implementations are made for *POSIX*, *mbed*², *ChibiOS*³ and *Arduino*⁴ compatible platforms using either TCP or Serial connections.

6.2 Storage

The RTS has to store both statically and dynamically allocated data. The static data contains the interpreter state, the interpreter stack and the communication buffers. The dynamic data consists of the tasks descriptions and bytecode and the SDSs description and value. Tasks consists of the bytecode, an identifier and the scheduling information. An SDS contains the identifier and a typed value.

SDSs and tasks are dynamically added and removed. Some devices have very little memory and therefore space needs to be used optimally. While almost all MCUs support heaps nowadays, the functions for allocating and freeing memory on the heap are not very space optimal and often leave holes if allocations are not freed in a last in first out order. To mitigate this problem, the RTS manages its own – compile time configurable sized – memory in the global data segment. Tasks are stored from the top down and SDSs are stored from the bottom up.

When a task or SDS is removed, all the remaining objects in the memory space are reordered in such a way that there are no holes left. In practice this means that if the first received task is removed, all tasks received later are moved. Obviously, this is quite time intensive, but it cannot be permitted to leave holes in the memory since the memory space is so limited. With this aggressive memory management technique, even an Arduino with only 2K RAM can hold several tasks and SDSs at the same time.

6.3 Runtime System Scheduling

The RTS starts initializing the allocated memory followed by running a device specific setup function. After the setup, the RTS waits for a connection. When the connection has been established it executes the main loop continuously. Only in case of a shutdown request, the memory is cleared and the device specific setup function is executed again so that the device is in the initial state again and waits for a connection.

The main loop (1) checks if there is input on the communication channel. If input is available, it reads and parses this to a message and process the message. All messages – such as new tasks, SDSs or a specification request – are processed immediately. (2) The RTS executes tasks that are ready. This means that all one shot tasks

²<https://mbed.com>

³<https://chibios.org>

⁴<https://arduino.cc>

and all interval tasks for which the interval has passed are executed. The execution happens in a round robin fashion. If a one shot task is executed, it is removed from the memory. If an interval task is executed, its last run time is set to the current time. (3) When one entire loop is finished, the program waits for a — compile time determined — time after which it continues. During this waiting the device idles.

7 ITASKS INTEGRATION

The *iTasks* part of the system — running on the server — is responsible for housekeeping the IoT devices. Only the glue functions are available for the programmer. However, internally, a lot of work needs to be done when a device is connected, a task is sent or an SDS is written.

7.1 Communication with Devices

Every IoT device has an *iTasks* memory based SDS assigned to it that contains their communication channels. These channels form the communication-agnostic interface for all communication to-and-fro the device. Listing 28 shows the type signature. The channels SDS is a tuple containing an incoming channel, an outgoing channel and a stop flag. The synchronization function — started on device connection — synchronizes the communication channels with the device. If messages appear in the first list, the synchronization function relays them to the device. Moreover, if the device sends a message, the synchronization function places it in second list. When the stop flag is set, the synchronization function terminates because the connection with the device is closed. If a new communication method is to be added, a programmer only has to implement the synchronization function for the whole system to work.

```
:: Channels ::= Shared ([MSGRecv], [MSGSend], Bool)
:: MSGRecv = MTTaskAck Int Int | MTSDSACK Int | MTPub Int BCValue | ...
:: MSGSend = MTTask Interval String | MTSDs Int BCValue | ...
```

Listing 28: Device type

7.2 Device Storage

All devices are stored in one global SDS containing a list of Device records. Each record contains everything a device encompasses. Storing all devices in a single SDSs has the advantage that one can inspect all devices from anywhere in the *iTasks* program. This facilitates debugging and error handling. The system knows which devices are connected and which IoT SDSs are available on a device. Moreover, this approach also allows the creation of systems that reconnect devices after a restart because a persistent SDS can be used to store the devices.

Every record stores a reference to the communication channels, the compiler state, the information to setup the synchronization function, a list of IoT tasks and a list of IoT SDSs. If the device is connected, it also stores the task identifier for the synchronization function, and the hardware specification. If the device is erroneously disconnected it can set a descriptive error.

The programmer can only add devices to the system through the *withDevice* function. This function is a wrapper around several asynchronous functions that interact with the device records in the global SDS. This function (1) adds the device to the global device list (2) connects the device (3) executes the task requiring the device

(4) waits for a stable value for the given task (5) requests a shutdown for the device (6) removes the device from the global device list.

Connecting a device is also a multi step process in itself. It (1) clears the channels. (2) starts a message processing task. (3) sends a device specification request and stabilizes when this request has been honored.

The message processing task is a compound task consisting of the device specific synchronization function and a device agnostic message processing function. The message processing function acts upon new messages in the incoming channels. For example, when an MTPub is received, the according SDS in the device record is updated or when an MTTaskAck is received the task in the device record is updated with the appropriate task identifier.

7.3 Synchronizing Shared Data Sources

The server stores a proxy value for every IoT SDS in the form of an *iTasks* SDS. This proxy *iTasks* SDS gives the latest cached value from the device on read on the server and sends an SDS write request to the device when written from the server. Watchers are notified when the device published a new value. This cached value is stored in the device record in the IOTShare type. This IOTShare type contains the identifier, the current value and possibly the reference to the *iTasks* SDS.

```
:: IOTShare =
{ identifier: Int
, value    : BCValue
, iTaskRef : Maybe (Shared BCValue)
}
```

Listing 29: Device SDS type

If the device publishes a new value for an IoT SDS, the processing task updates the proxy value stored in the device record. Moreover, the processing task writes the new value to the reference *iTasks* SDS lens. This lens automatically translates the BCValue box to the correct value and write the actual referenced *iTasks* SDS.

The other way around, when a task updates the referenced *iTasks* SDS, a watcher task — started in the *liftIOTTTask* function (Section 7.4) is notified. The watcher task can then update the proxied value in the device record accordingly. The synchronizing of this proxied value with the actual device is explained in Section 7.5.

7.4 Executing Tasks

The programmer can only execute tasks on a device through the *liftIOTTTask* task. In the same fashion as the *withDevice* function, it wraps several device record modifying tasks.

The function (1) compiles the IoT task to messages (2) places the messages in the channels SDS (3) adds the task to the device record (4) sets the task identifier when the task acknowledgement is received (5) waits for the task to stabilize while watching all the reference *iTasks* SDSs watching a reference *iTasks* task is done by using the *whileUnchanged* function on the Shared BCValue lens. Moreover, it removes the task from the device by sending a MTTaskDel message when the *iTasks* task is terminated. Termination of *iTasks* tasks happens for example if the task is combined with a step (\gg^*) combinator and one of the predicates of the right hand side triggers continuation.

7.5 Utility Lenses on the Device Shared Data Source

All the device information is stored in a single *iTasks* SDS for reasons mentioned in Section 7.2. Unfortunately, there are also two issues with this approach.

First, it is not convenient to edit the global SDS containing the devices if you are only interested in a part of it. For example, updating a single IoT SDS value for a single device requires a complicated update function.

Secondly, a task watching the global SDS might only be interested in a very small section of it but is notified on all changes. For example, watcher tasks are launched in the `liftTask` function. These tasks watch only a single reference *iTasks* SDS to make sure the value is proxied in the device record. However, if another task writes in the device SDS in a different place, the watchers are notified either way. To overcome these issues, parametric lenses were used.

To solve the first type of problems, parametric lenses were introduced [Domoszlai et al. 2014]. The type variable was left untouched in Section 2. This extra type variable represents the parameter that is available during reading and writing. This *p* must be of type `()` in the standard atomic access functions which can be achieved using the `sdsFocus` function. This function fixes the parameter and casts it to `()` so that the SDS is usable for the standard tasks. The SDS can specify or refine the data read or written according to the given value of *p*. Moreover, it can place a filter on the notifications using this parameter. For example, it can be used to create a lens on a tuple — e.g. only giving write access to the first element. If the second element is written, a watcher on the first element is not notified.

```
sdsFocus :: p (SDS p r w) → SDS () r w | iTask p
```

Listing 30: Focussing a parametric lens

The following lens is used.

```
:: IOTParam = Global | Local Device | Share Device Int
deviceStore :: SDS IOTParam [Device] [Device]
```

Listing 31: Global device SDS

Every constructor denotes a different type of view on the root SDS. First, the `Global` constructor is only interested in the entire list of devices. Secondly, the `Local` lens only looks at a single device. Finally, the `Share` view only focusses on a single SDS on a single device.

For these SDS lenses, functions are available for accessing the parts of the global SDS.

```
deviceStoreP :: Shared [Device]
deviceShare :: Device → Shared Device
shareShare :: Device IOTShare → Shared BCValue
```

Listing 32: Utility SDSs

Focussing the `deviceStore` to `Global` gives access to the global SDS. Global watchers are only be notified if the structure of the list changes, i.e. if a device is added or removed.

Accessing a single device is done with the `deviceShare` function. The SDS requires a `Device` record to know on which device to focus. This record is in scope when the `withDevice` function is used. Watchers of a local SDS are notified when something in the

device changes. Writers can change something in the device such as the specification.

Share SDSs can be accessed through the `shareShare` function. This function focusses the global SDS on a single SDS on a single device. It is only notified when that specific proxy SDS value changes.

This brings us to the last unsolved part of the extension, namely synchronizing the proxy values with the device. Albeit not designed for it, parametric lenses can also be used to solve this problem. The parameter is known in the stateful write function in so that other SDS — e.g. device's channels — can also be written.

When a task writes to this SDS, the global SDS knows this through the parameter and propagates the value to the device. When the server or the device changes the SDS, this view is notified. The SDS requires a `Device` and a `IOTShare` record. The `IOTShare` record and therefore the lens is only used internally, for example by the processing function watching the IoT SDSs to synchronize them with their *iTasks* references.

7.6 Thermostat Example (3)

The thermostat example is kept simple for illustration purposes. However, it does not show the full potential of the extension.

For example, it is possible to create redundant systems where IoT tasks on one device are automatically reassigned to another if the device breaks. We show this by extending the example with a redundant cooler. The only function we need to change is the `operateCooler` function, note that the cooler *iTasks* SDS is in scope here. The `withDevice` function throws an exception if a device terminates the connection for unknown reasons and the connection cannot be re-established. This exception can be caught and the work that is done on the device can be moved to another device.

```
operateCooler :: Task ()
operateCooler = try
  (withDevice coolerDevice runCoolerTask)
  (λexc → viewInformation "Exception" [] (toString exc)
   >>| withDevice coolerDevice2 runCoolerTask)
where
  runCoolerTask dev = liftIOTTask (dev, OnInterval 1000
    (sds λf=cooler In {main= digitalWrite LED1 f :. digitalWrite D5 f}))
  coolerDevice2 :: TTYSettings
  coolerDevice2 = {devicePath="/dev/ttyACM0", baudrate=B19200, ...}
```

Listing 33: Redundant cooler

Moreover, IoT tasks can be sent dynamically at runtime to the device. To illustrate this, we show a task with which the user can send tasks to a device. The user selects the task from a list of predefined tasks and provides the execution strategy via the web interface as well. If this task is executed in parallel with the sensor and cooling tasks with the device, the user can use the device as well for miscellaneous tasks. For example to blink a light, or to open window blinds on demand.

```
interact :: MTaskDevice → Task ()
interact device = enterChoice "Choose a task" [ChooseFromList fst] taskList
  -&&- enterInformation "Execution Strategy" []
  >^* [OnAction (Action "Send") $ withValue λ( _, task), strat) →
        Just (task >=> λiottask liftTask (device, strat) iottask)]
  @! ()
```

```

where
taskList :: [(String, Task (Main (ByteCode () Stmt)))]
taskList =
  [ ("faculty", enterInformation "Faculty of what?" [])
    >>> λn → var λx=1 In var λy=n In {main =
      IF (x ==. 0) return (x =. y * . x :. y -. lit 1)})
    , ("count", return $
      sds λx=0 In {main = x =. x +. lit 1 :. pub x})
    , ("blink", enterInformation "Led on which pin?" [])
      >>> λ1 → var λx=True In {main = x =. Not x :. dIO 1 =. x)
    , (... , ...)
  ]

```

Listing 34: Dynamic IoT task sending

8 RELATED WORK

Related research has been conducted on the subject arising from academia and the industry. For example, MCUs such as the Arduino can be remotely controlled very directly using the Firmata-protocol⁵. This protocol is designed to allow control of the peripherals – such as sensors and actuators – directly through commands sent via a communication channel such as a serial port. This allows very fine grained control but with the cost of excessive communication overhead since no code is executed on the device itself, only the peripherals are queried. A Haskell implementation of the protocol is also available⁶. The hardware requirements for running a Firmata client are very low because all the logic is on the server. However, the bandwidth requirements are high and therefore it is not suitable for IoT applications that communicate through LTN networks. Similarly, Grebe and Gill [2016] created *Haskino*, a monadic interface over *hArduino* that allows remote code execution on Arduinos. Their initial tethered solution is based on Firmata but they also propose an untethered approach that is similar to compilation by storing the program in EEPROM. However, there is no communication between the device and the server that programmed it and the solution is very specific to the Arduino ecosystem.

Clean has a history of interpretation, for example, there is a lot of research happening on the intermediate language SAPL. SAPL is a purely functional intermediate language that can be efficiently interpreted. It has interpreters written in C++ [Jansen et al. 2007] and a compiler to JavaScript [Domoszlai et al. 2011]. Compiler backends exist for Clean and Haskell which compile the respective code to SAPL [Domoszlai and Plasmeijer 2012]. The SAPL language is a functional language and therefore requires big stacks and heaps to operate and is therefore not directly suitable for devices with little RAM such as the Arduino which only boasts 2K of RAM. It might be possible to compile the SAPL code into efficient machine language or C but then the system would lose its dynamic properties since the MCU then would have to be reprogrammed every time a new task is sent to the device.

EDSLs have often been used to generate C code for MCU environments. This work uses parts of the existing mTask-eDSL which generates C code to run a TOP-like system on MCUs [Koopman et al. 2018; Koopman and Plasmeijer 2016]. Again, this requires a reprogramming cycle every time the task-specification is changed and there is no interaction with the server. The nature of the embedding technique allows additional backends to be written without

⁵<https://github.com/firmata/protocol>

⁶<https://leventerkok.github.io/hArduino>

touching existing ones. Hence, the eDSL is used for this solution but with a novel backend.

Another eDSL designed to generate low-level programs is called Ivory and uses Haskell as a host language [Elliott et al. 2015]. The language uses the Haskell type-system to make unsafe languages type safe. For example, Ivory has been used in the automotive industry to program parts of an autopilot [Hickey et al. 2014; Pike et al. 2014]. Ivory’s syntax is deeply embedded but the type system is shallowly embedded. This requires several Haskell extensions that offer dependent type constructions. The process of compiling an Ivory program happens in two stages. The embedded code is transformed into an AST that is sent to a chosen backend. In the novel extension, the eDSL is transformed directly into functions and there is no intermediate AST. Moreover, Ivory generates static programs and thus it is necessary to reprogram the devices when they need to be repurposed.

Not all IoT devices run solely compiled code, e.g. the ESP8266 powered NodeMCU is able to run interpreted Lua code. Moreover, there is a variation of Python called micropython that is suitable for running on MCUs. However, the overhead of the interpreter for such rich languages often results into limitations on the program size. It would not be possible to repurpose a device with IoT because implementing this extensibility in the interpreted language leaves no room for the actual programs. Also, some devices only have 2K of ram, which is not enough for this.

9 CONCLUSION AND FUTURE WORK

9.1 Conclusion

This paper introduces a novel system for adding IoT control to the TOP implementation *iTasks*. A new backend for an existing eDSL has been created that compiles the program into bytecode that is interpreted by a client. Clients are written for several MCUs and architectures which are connected through various means of communication such as serial port, WiFi and wired network communication. The bytecode on the devices is interpreted using a stack machine and provides the programmer with interfaces to the peripherals. Communication between IoT tasks and *iTasks* tasks happens via SDSs.

The dynamic nature of the system allows the MCU to be programmed once and used many times for different applications or within the same application. When a device is assigned to an IoT task but the device suddenly becomes unusable, the *iTasks* system can reassign the IoT task to another device that is also suitable for running the task without needing to recompile the code. Adding peripherals is not a time consuming task and does not even require recompilation of clients not having the peripheral.

The new functionality extends the reach of *iTasks* by adding IoT functionality and allowing devices to run IoT tasks. With this extension, a programmer can create an entire IoT application from one source that reaches all layers of the IoT architecture. Components can be updated individually without causing integration problems. Devices can be repurposed just by sending new tasks to it. Most importantly, it gives an insight in the possibilities of adding IoT to TOP systems.

9.2 Extensions

The current system is functional but there are a lot of ideas for future work.

An additional simulation view to the IoT eDSL can be added that works similar to the existing C-backend simulation. The first option is to simulate a device in total, this means that the simulator is implemented the Duplex class. Secondly, it can simulate the behaviour on a higher level by implementing the eDSL classes. At the time of writing work is done on an *iTasks* simulator device of this kind.

True multitasking can be added to the client software. IoT tasks get slices of execution time and have their own stack, this allows IoT tasks to run truly concurrent. Multitasking allows tasks to be truly interruptible by other tasks. Furthermore, this allows for more fine-grained timing control of tasks. However, it influences memory requirements.

Many research topics can be explored in the field of resource management and analysis, both statically at compile time and dynamically at runtime for both peripheral requirements and memory requirements.

9.3 Further Improvements

A lot of improvements can be done in the interpreter as it is implemented very straightforwardly without performance taken into mind.

The current implementation offers a small subset of the *iTasks* combinators in IoT. The subset can be extended to allow for more fine-grained control flow between IoT tasks. Furthermore, more logic can be moved to the device instead of it residing on the server reducing the communication overhead.

At the moment, all data is sent as plain text over the wire and the device cannot know whether it talks to a legitimate server or an attacker. Due to the nature of the system, namely sending code that is executed, security needs to be investigated. Only bytecode for very specific IoT tasks can be sent to the device at the moment which mitigates the risk somewhat. As the language will become more expressive, the security risk increases.

Finally, the robustness of the system can be improved. Tasks residing on a device that disconnects should be kept on the server to allow a swift reconnect and restoration of the tasks. Moreover, an extra specialization of the shutdown can be added that drops the connection but keeps the tasks in memory. This can be extended by allowing devices to send their tasks back to the server. In this way devices can even connect to different servers. Also, tasks can be stored in EEPROM or on external memory to be able to access them even after a reboot or to save memory. For example, EEPROM can be written about ten to a hundred times more often than flash memory.

REFERENCES

- Tom Brus, Marko van Eekelen, Maarten Van Leer, and Marinus Plasmeijer. 1987. Clean – a language for functional graph rewriting. In *Conference on Functional Programming Languages and Computer Architecture*. Springer, 364–384.
- Jacques Carette, Oleg Kiselyov, and Chung-chieh Shan. 2009. Finally tagless, partially evaluated: Tagless staged interpreters for simpler typed languages. *Journal of Functional Programming* 19, 5 (2009), 509–543.
- Li Da Xu, Wu He, and Shancang Li. 2014. Internet of things in industries: a survey. *Industrial Informatics, IEEE Transactions on* 10, 4 (2014), 2233–2243.
- László Domoszlai, Eddy Bruel, and Jan Martin Jansen. 2011. Implementing a non-strict purely functional language in JavaScript. *Acta Universitatis Sapientiae* 3 (2011), 76–98.
- László Domoszlai, Bas Lijnse, and Rinus Plasmeijer. 2014. Parametric lenses: change notification for bidirectional lenses. In *Proceedings of the 26nd 2014 International Symposium on Implementation and Application of Functional Languages*. ACM, 9.
- László Domoszlai and Rinus Plasmeijer. 2012. Compiling Haskell to JavaScript through Clean's core. *Annales Universitatis Scientiarum Budapestinensis de Rolando Eötvös Nominae. Sectio Computatorica* 36 (2012), 117–142.
- Trevor Elliott, Lee Pike, Simon Winwood, Pat Hickey, James Bielman, Jamey Sharp, Eric Seidel, and John Launchbury. 2015. Guilt free ivory. In *ACM SIGPLAN Notices*, Vol. 50. ACM, 189–200.
- Jeremy Gibbons. 2015. Functional programming for domain-specific languages. In *Central European Functional Programming School*. Springer, 1–28.
- Mark Grebe and Andy Gill. 2016. Haskino: A remote monad for programming the arduino. In *International Symposium on Practical Aspects of Declarative Languages*. Springer, 153–168.
- Patrick Hickey, Lee Pike, Trevor Elliott, James Bielman, and John Launchbury. 2014. Building embedded systems with embedded DSLs. In *ACM SIGPLAN Notices*. ACM Press, 3–9. <https://doi.org/10.1145/2628136.2628146>
- Riazul Islam, Daehan Kwak, Humaun Kabir, Mahmud Hossain, and Kyung-Sup Kwak. 2015. The Internet of Things for Health Care: A Comprehensive Survey. *IEEE Access* 3 (2015), 678–708. <https://doi.org/10.1109/ACCESS.2015.2437951>
- Jan Martin Jansen, Pieter Koopman, and Rinus Plasmeijer. 2007. Efficient Interpretation by Transforming Data Types and Patterns to Functions. *Trends in Functional Programming* 7 (2007), 73.
- Mark Jones. 1995. Functional programming with overloading and higher-order polymorphism. In *International School on Advanced Functional Programming*. Springer, 97–136.
- Pieter Koopman, Mart Lubbers, and Rinus Plasmeijer. 2018. A Task-Based DSL for Microcomputers. In *Proceedings of the Real World Domain Specific Languages Workshop 2018*. ACM Press, 1–11. <https://doi.org/10.1145/3183895.3183902>
- Pieter Koopman and Rinus Plasmeijer. 2016. A Shallow Embedded Type Safe Extendable DSL for the Arduino. In *Trends in Functional Programming*. Lecture Notes in Computer Science, Vol. 9547. Springer International Publishing, Cham. DOI: 10.1007/978-3-319-39110-6.
- Bas Lijnse. 2013. *TOP to the rescue: task-oriented programming for incident response applications*. s.n.; UB Nijmegen, S.L.; Nijmegen. OCLC: 833851220.
- Arjan Oortgiese, John van Groningen, Achter Peter, and Plasmeijer Rinus. 2017. A Distributed Dynamic Architecture for Task Oriented Programming. In *Proceedings of the 29nd 2017 International Symposium on Implementation and Application of Functional Languages (IFL '17)*. ACM, New York, NY, USA, To appear.
- Lee Pike, Patrick Hickey, James Bielman, Trevor Elliott, Thomas DuBuisson, and John Launchbury. 2014. Programming languages for high-assurance autonomous vehicles: extended abstract. In *Proceedings of the ACM SIGPLAN 2014 Workshop on Programming Languages meets Program Verification*. ACM Press, 1–2. <https://doi.org/10.1145/2541568.2541570>
- Rinus Plasmeijer, Peter Achten, and Pieter Koopman. 2007. iTasks: executable specifications of interactive work flow systems for the web. *ACM SIGPLAN Notices* 42, 9 (2007), 141–152.
- Rinus Plasmeijer, Bas Lijnse, Steffen Michels, Peter Achten, and Pieter Koopman. 2012. Task-Oriented Programming in a Pure Functional Language. In *Proceedings of the 2012 ACM SIGPLAN International Conference on Principles and Practice of Declarative Programming, PPDP '12*. ACM, Leuven, Belgium, 195–206.
- Rinus Plasmeijer, Marko van Eekelen, and John van Groningen. 2011. Clean language report version 2.2 (2011). <https://clean.cs.ru.nl/download/doc/CleanLangRep.2.2.pdf>
- Josef Svenningsson and Emil Axelsson. 2012. Combining deep and shallow embedding for EDSL. In *International Symposium on Trends in Functional Programming*. Springer, 21–36.

On Optimizing Bignum Toom-Cook Multiplication

Shamil Dzhatdoyev, Nicholas Nelson, Marco Morazán, Josephine Des Rosiers

ACM Reference format:

Shamil Dzhatdoyev, Nicholas Nelson, Marco Morazán, Josephine Des Rosiers. 2019. On Optimizing Bignum Toom-Cook Multiplication. In *Proceedings of International Symposium on Implementation and Application of Functional Languages, Lowell, MA, USA, August 2019 (IFL'18)*, 7 pages.
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

As part of the abstractions they provide, many functional languages offer support for integers of arbitrary size and precision, called *bignums*. Current hardware limitations force programmers to choose between exact integer types such as int or long long which have limited size, or inexact numeric types such as float or double which also have limited size, but also have limited precision when considering larger numbers. Bignums, however, liberate programmers from these limitations. That is, programmers do not need to concern themselves with whether or not an integer fits in a word of memory. Consider for example, algorithms which find arbitrarily large prime numbers. Given that $2^{77,232,917} - 1$, the largest known Mersenne prime as of January 3, 2018 [8] has 23,249,425 digits, it is impossible to represent it precisely using a single hardware word. Therefore, in order to manipulate integers of this magnitude, bignums are used to represent integers as compound data structures occupying multiple hardware words of memory.

Bignums are fundamental in fields such as cryptography (e.g., [13]), network security services (e.g., [11]), and computation algebra systems that involve astronomically large integers (e.g., [12]). The main advantages offered by bignums over standard hardware-limited integers are their size, their precision, and the ease of implementing computational algorithms which require very large integers. In cryptography, for example, bignums are used to represent very large prime numbers utilized to shield hashing algorithms from attacks that can exploit the existence of small factors in a hashing key. The efficient implementation of the bignum multiplication operator is of key importance to performance, owing to the fact that many applications rely heavily on multiplication of integers of arbitrary size.

Of particular interest to implementors of a bignum library for a functional language is how fast multiplication can be done. The traditional base case algorithm (a.k.a the grade school method) is $O(n^2)$, where each digit (i.e., a digit in base B) of the multiplier is multiplied by each digit of the multiplicand. Algorithms by Karatsuba [4] and later by Knuth [5] reduce the complexity of multiplying two

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

IFL'18, August 2019, Lowell, MA, USA

© 2019 Association for Computing Machinery.
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

large bignums to $O(n^{1.585})$. These algorithms achieve this result by splitting each bignum of length $2n$ into two pieces of length n to reduce the product to three multiplications of length n and some adding operations. Later, the Toom-Cook algorithm abstracted this splitting to split the bignums into an arbitrary number of pieces. In particular, the Toom-Cook 3-way algorithm, known also as Toom-3 multiplication, splits each bignum of length $3n$ into three pieces of length n to reduce the product to five multiplications of length n , with some addition operations, and multiplications and divisions by scalars. This reduces the complexity of multiplying two large bignums to $O(n^{1.465})$.

A natural question arises when implementing bignum multiplication: determining the optimal cut-off value to decide when to split two bignums to use the Toom-Cook algorithm and when to use another multiplication algorithm such as the base-case algorithm. This article explores this question in the context of a library being developed for the Green functional programming language, and in particular for the Toom-3 multiplication algorithm. Section 2 reviews related work and the two multiplication algorithms. Section 3 briefly describes the implementation of the Green bignum library. Section 4 presents empirical data for the dimension being studied. Finally, Section 5 presents conclusions and directions for future work.

2 RELATED WORK

To manipulate integers of arbitrary size, *bignums* are used to represent integers as compound data structures. A bignum contains a sequence of *bigits*, integers in $[0..B)$. In practical bignum implementations, the base of the representation, B , is always much larger than 10 and is limited by the hardware. For efficiency reasons, B does not exceed the largest integer representable using one word of memory for the given implementation. A bignum is usually represented as compound data that contains at least three pieces of information:

- (1) The address of a vector of bigits,
- (2) The length of the vector of bigits,
- (3) The sign of the bignum.

A vector is typically used to store the bigits, rather than a list, in order to have efficient implementations of all of the basic arithmetic operators. For example, addition and subtraction operators always start with the least significant bigits of the bignums given as arguments, while division always starts with the most significant bigits of the integers it receives as arguments. Using lists, which do not have random access, to represent bignums would add overhead to one or more of the basic arithmetic operators.

2.1 Multiplication Algorithms

Any implementation of bignums must include a multiplication operator. Bignum multiplication algorithms can be divided into two categories: base-case multiplication algorithms and divide and conquer algorithms. The divide and conquer algorithms are more

efficient for "long" bignums. The best known base-case multiplication algorithm is the classical algorithm described by Knuth. This algorithm is implemented by virtually all systems and libraries that support bignums (e.g., [2]). It is similar to the grade-school method. Each bit of the multiplier is multiplied by each bit of the multiplicand to produce a result that is multiplied by B^p , where p is the position of the bit in the multiplier and B , as before, is the base used to represent the bignums. In practice, vector-based shifting is often employed rather than a multiplication by B^p , as it incurs much less overhead. Instead of allocating many bignums for intermediate products, the results are all directly added to the vector of bits holding the result.

The best known divide and conquer algorithms are the Karatsuba [4] algorithm and the Toom-Cook algorithm [2, 5]. These divide and conquer algorithms split bignums with a length greater than some threshold, k , into smaller pieces until the splitting stops and another multiplication algorithm is used. This second algorithm could either be another divide and conquer algorithm, or a base-case algorithm. The Toom-3 algorithm in particular splits two bignums of length $3n$ into three pieces of length n , using these pieces to solve a polynomial multiplication problem.

Two bignums of length $3n$, u and v can be split into three pieces of length n : u_0, u_1 , and u_2 and v_0, v_1 , and v_2 respectively such that $u = u_0 + u_1B^n + u_2B^{2n}$ and $v = v_0 + v_1B^n + v_2B^{2n}$. From this, two second-degree polynomials arise:

$$\begin{aligned} U(x) &= u_0 + u_1x + u_2x^2 \\ V(x) &= v_0 + v_1x + v_2x^2 \end{aligned}$$

where $u = U(B^n)$ and $v = V(B^n)$.

By multiplying $U(x)$ and $V(x)$ we get a fourth-degree polynomial

$$W(x) = w_0 + w_1x + w_2x^2 + w_3x^3 + w_4x^4$$

where $W(B^n)$ then, is in fact uv . Rather than multiplying $U(B^n)$ with $V(B^n)$ directly, which would constitute 9 multiplications of length n , the product is determined by finding the coefficients of W through steps of evaluation, recursive multiplication, and interpolation.

Since U and V are second-degree polynomials, their product W is a fourth-degree polynomial which can be uniquely identified by five points. These points are carefully picked to reduce to time needed for evaluation of $U(x)$ and $V(x)$, as well as to simplify the interpolation step later. One set of points, determined by Bodrato [1], and also used in GMP [2] are 0, 1, -1, -2, and ∞ . U and V are evaluated at these points, according to the following set of

evaluations

$$\begin{aligned} r &= u_0 + u_2 \\ U(0) &= u_0 \\ U(1) &= r + u_1 \\ U(-1) &= r - u_1 \\ U(-2) &= (U(-1) + u_2) * 2 - u_0 \\ U(\infty) &= u_2 \\ s &= v_0 + v_2 \\ V(0) &= v_0 \\ V(1) &= s + v_1 \\ V(-1) &= s - v_1 \\ V(-2) &= (V(-1) + v_2) * 2 - v_0 \\ v(\infty) &= v_2 \end{aligned}$$

and these values are multiplied recursively to produce the values $W(0), W(1), W(-1), W(-2)$ and $W(\infty)$. The process of evaluation requires 10 additions of length n as well as two multiplications by scalar values, whereas the recursive multiplication requires 5 multiplications of length n . Once these values are known, the polynomial $W(x)$ can be solved by solving the following system of equations:

$$\begin{aligned} W(0) &= w_0 \\ W(1) &= w_0 + w_1 + w_2 + w_3 + w_4 \\ W(-1) &= w_0 - w_1 + w_2 - w_3 + w_4 \\ W(-2) &= w_0 - 2w_1 + 4w_2 - 8w_3 + 16w_4 \\ W(\infty) &= w_4 \end{aligned}$$

This is a system of five equations with three unknown values: w_1, w_2 , and w_3 , and can be solved by some elementary linear algebra. However, more efficient methods of solving for these unknowns are typically used, such as the step-wise interpolation method offered by Bodrato, which calculates the coefficients in two phases:

$$\begin{aligned} w_0 &\leftarrow W(0) \\ w_4 &\leftarrow W(\infty) \\ w_3 &\leftarrow (W(-2) - W(1))/3 \\ w_1 &\leftarrow (W(1) - W(-1))/2 \\ w_2 &\leftarrow W(-1) - w_0 \\ w_3 &\leftarrow (w_2 - w_3)/2 + 2W(\infty) \\ w_2 &\leftarrow w_2 + w_1 - w_4 \\ w_1 &\leftarrow w_1 - w_3 \end{aligned}$$

This interpolation step requires 8 additions or subtractions of length $2n$, as well as 2 multiplications by scalar values and 2 divisions by scalar values. Once all five of the coefficients are determined, the final result is recomposed by adding them into the resultant bignum in the correct position.

To make the algorithm more concrete, consider the multiplication of two bignums represented using base 10, 123654 and 789210. Each bignum has 6 digits and can be visualized as follows:

$$\begin{aligned} u &= 12\ 36\ 54 \\ v &= 78\ 92\ 10 \end{aligned}$$

Using a threshold of 3 bigits to control the splitting means that the numbers must be split as their lengths exceed the threshold. This leads to the following bignums of length 2:

$$\begin{aligned} u_2 &= 12, \quad u_1 = 36, \quad \text{and } u_0 = 54 \\ v_2 &= 78, \quad v_1 = 92, \quad \text{and } v_0 = 10 \end{aligned}$$

Both polynomials are evaluated according to the equations above, giving the values 54, 102, 30, 30, and 12 for $U(0), U(1), U(-1), U(-2)$, and $U(\infty)$ respectively as well as 10, 180, -4, 138, and 78 for $V(0), V(1), V(-1), V(-2)$, and $V(\infty)$ respectively. As none of these values exceed 3 bigits, no more splitting occurs and the values are simply multiplied using a base-case algorithm. These 5 multiplications yield the values 540, 18360, -120, 4140, and 936 for $W(0), W(1), W(-1), W(-2)$, and $W(\infty)$. Interpolation follows the steps outlined above:

$$\begin{aligned} w_0 &\leftarrow 540 \\ w_4 &\leftarrow 936 \\ w_3 &\leftarrow (4140 - 18360)/3 = -4740 \\ w_1 &\leftarrow (18360 - (-120))/2 = 9240 \\ w_2 &\leftarrow -4 - 540 = -660 \\ w_3 &\leftarrow (-660 - (-4740))/2 + 2(936) = 2040 + 1872 = 3912 \\ w_2 &\leftarrow -660 + 9240 - 936 = 7644 \\ w_1 &\leftarrow 9240 - 3912 = 5328 \end{aligned}$$

and the final result is determined by recomposing the value from the coefficients

$$540 + 5328(10^2) + 7644(10^4) + 3912(10^6) + 936(10^8) = 97588973340$$

The key to success is finding a good value for k to control the splitting. If k is too small, then the overhead incurred by splitting is too costly. If k is too large, then too little splitting occurs and the runtime is dominated by the base case algorithm. The GNU GMP documentation states that the cut-off can be set as low as $k = 81$ [2] and the library sets the default value for k in the range 81..99 depending on the CPU being used, for modern CPU architectures. The authors have found no empirical evidence to validate these values as good choices. It is clear, however, that as the base case algorithm is made more efficient, the value of k grows. In addition, some versions of the Toom-Cook algorithm default to the Karatsuba algorithm when the length of the bignums is too short. In this case, as the Karatsuba algorithm is made more efficient, k will also grow. Conversely, as the implementation of the Toom-Cook algorithm is made more efficient that value of k decreases.

3 BIGNUMS IN GREEN

Before examining empirical data, the Green programming language and the Green virtual machine (GVM) are briefly described in this section. Green is an experimental functional, eager, and impure programming language. Although mutation is supported, its use is, in general, discouraged favoring mutation-free programming whenever possible. The Green compiler (implemented in Racket) transforms a program written in source syntax into bytecode for the GVM. Currently, Green is being used to study novel closure representations. In particular, the goal is to explore the efficiency of memoized closures and the elimination of data-structure closures in favor of bytecode closures created through the use of a controlled

form of β -reduction [10]. Programs are λ -lifted [3, 9] to make explicit the free variables of functions and, therefore, the values by which functions/closures are specialized by at runtime in the spirit of β -reduction. These values can, of course, be bignums.

The GVM is, in essence, a bytecode interpreter (written in C++) that manages a heap from which memory is dynamically allocated at runtime including memory for bignums. It has a set of registers that are used for parameter passing. These registers must be set correctly to transfer control to a function being called, as the GVM has no mechanism for returning from a function call. That is, the GVM assumes that there are no delayed operations in the code produced by the compiler as is expected for programs written in CPS. Seamless support for bignums is provided as is done in many functional languages (e.g., [7]). Arithmetic operators automatically produce bignums or fixed-sized integers when necessary without burdening the programmer with the task of guarding against overflow or underflow.

Internally, a bignum is represented as a structure that has three fields, the heap address of a vector of bigits, the length of the vector, and the sign of the bignum. Bigits in the bignum are allocated from least significant to more significant. Macros are used to overload the arithmetic operators to accept bignums, fixed-sized integers, or a combination of inputs. Toom-3 multiplication is implemented only allocating 19 intermediate vectors of bigits. The evaluation of $U(x)$ and $V(x)$ at the five points requires 6 vectors each. Three vectors are allocated for $W(1), W(-1)$, and $W(-2)$ in the recursive multiplication stage. Finally, four vectors are allocated in the interpolation stage for w_1, w_2 , and w_3 , as well as the value of $2W(\infty)$, a value needed to compute w_3 . The values of w_0 and w_4 are computed first and stored directly in the output vector in the first and last third respectively. In addition, rather than defaulting to the base-case algorithm when the inputs to Toom-3 are shorter than the cut-off threshold, Karatsuba multiplication is used with a cut-off of 35.

The size of the vector of bigits for the result is automatically adjusted to eliminate internal fragmentation. For example, when multiplying two bignums of sizes s_1 and s_2 , the result may be of size $s_1 + s_2$ or $s_1 + s_2 - 1$. At the beginning of a multiplication operation, a vector of size $s_1 + s_2$ is allocated for the bigits of the result. If the result only has $s_1 + s_2 - 1$ bigits, then the vector is shortened by one before returning the result. To make this concrete, consider the product of two digits (i.e., base 10). A vector of length two is allocated for the result. If the digits 4 and 3 are multiplied, the result of 12 needs a vector of size 2. However, if 3 and 3 are multiplied, the result 9 needs only a vector of length 1. In this case, the vector is shortened by one and the unused memory is given back to the heap for future allocation.

4 EMPIRICAL RESULTS

This section presents empirical results for determining the ideal cut-off value for Toom-Cook 3-way splitting.

Three benchmarks are used to explore the issue of choosing the splitting threshold. The benchmarks are:

MM This benchmark multiplies two 25x25 square matrices, A , and B , containing random bignums of the same length. The matrices are represented using two-dimensional vectors. The well-known formula, $r_{ij} = a_i b_j$ is used to compute the

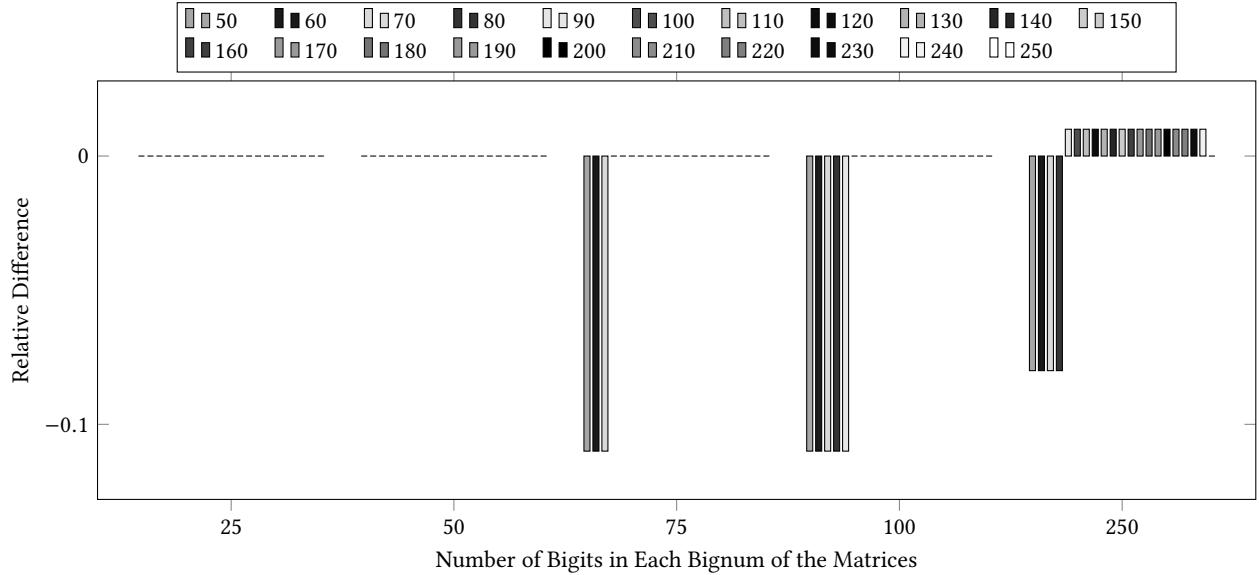


Figure 1: Karatsuba and Toom-3 Running-Time Relative Difference for MM Using Small Bignum Lengths.

result. The entry in the i^{th} row and j^{th} column of the result is the dot product of the i^{th} row of A and the j^{th} column of B [14].

Tetra This benchmark computes the n^{th} tetrahedral number using the binomial formula $T_n = \binom{n+2}{3}$ [6].

RSA This benchmark encrypts and decrypts a message, represented as a number, using the RSA encryption algorithm [13]. The message is a random bignum having 50 bits.

4.1 Toom-Cook Splitting

In this section, the value of k , the length limit to split bignums for multiplication is determined. The goal, concretely, is to determine the default value of k for the library.

The benchmarks are used to conduct experiments with the following inputs:

MM All the elements of the two 25x25 matrices contain random bignums of the same length. The following lengths in bits are used: 25, 50, 75, 100, 250, 500, 1000, 1500, and 2000. These are denoted **MM** bignum-length. For example, for bignums of length 25 the benchmark is denoted **MM 25**.

Tetra This benchmark is used to compute the $10^{10000^{th}}$ and the $10^{100000^{th}}$ tetrahedral numbers.

RSA This benchmark is used with four different public key lengths in bits: 109, 272, 431, and 614.

An experiment consists of running a benchmark with the same input and the same k -value 10 times. For each experiment, the average execution time is recorded. For example, **MM 100** and $k = 50$ is executed 10 times and the average running time is recorded.

To determine the default value of k , the relative difference between Karatsuba multiplication and Toom-3 is examined. For each use of a benchmark, 21 difference values for k are tested: all multiples of 10 in [50..250].

Figures 1 and 2 display the results obtained from 189 experiments using the **MM** benchmark. The x-axis represents the length of all the bignums in the two matrices, whereas the y-axis represents the relative run-time difference over the Karatsuba algorithm. Figure 1 displays the results of using relatively small bignum lengths (i.e., lengths ≤ 250). For the matrices with bignums that are too short to be split, we observe that the running time is, as expected, dominated by the Karatsuba algorithm, and therefore the performance of the Toom-3 algorithm is virtually the same as the Karatsuba algorithm. However, in cases where the bignums being processed are only split once by the algorithm, we see a significant slow-down compared to the Karatsuba algorithm. This can be seen in **MM 75** where k is in [50 – 70], **MM 100** where k is in [50 – 90], and **MM 250** where k is in [50 – 80]. The relative difference ranges from a low of -12% , indicating that Karatsuba multiplication is faster, to a high of 3% . Even though splitting does occur for some of these experiments, in none of these experiments is Toom-3 significantly faster than Karatsuba, suggesting that the overhead involved in the Toom-Cook algorithm is too large to overcome for the multiplication of smaller bignums.

Figure 2 displays the results from using relatively long bignum lengths (i.e., length ≥ 500). These are lengths for which a bignum is split more than once. The speed-up ranges from a low of 5% to a high of 16% . As expected, better performance is observed as the bignums get longer. This is due longer bignums allowing for more splitting and, thus, exploiting the Toom-Cook optimization more. It is important to note, that for the cut-off values in the range [70 – 250], we see the same pattern as with short bignums: for each bignum length, the speed-up observed for every value of k is within 2% of the other values of k . This suggests that for the values of k in this range, the observed performance is very similar and the length of the bignums being multiplied is not a distinguishing factor to select the best cut-off value.

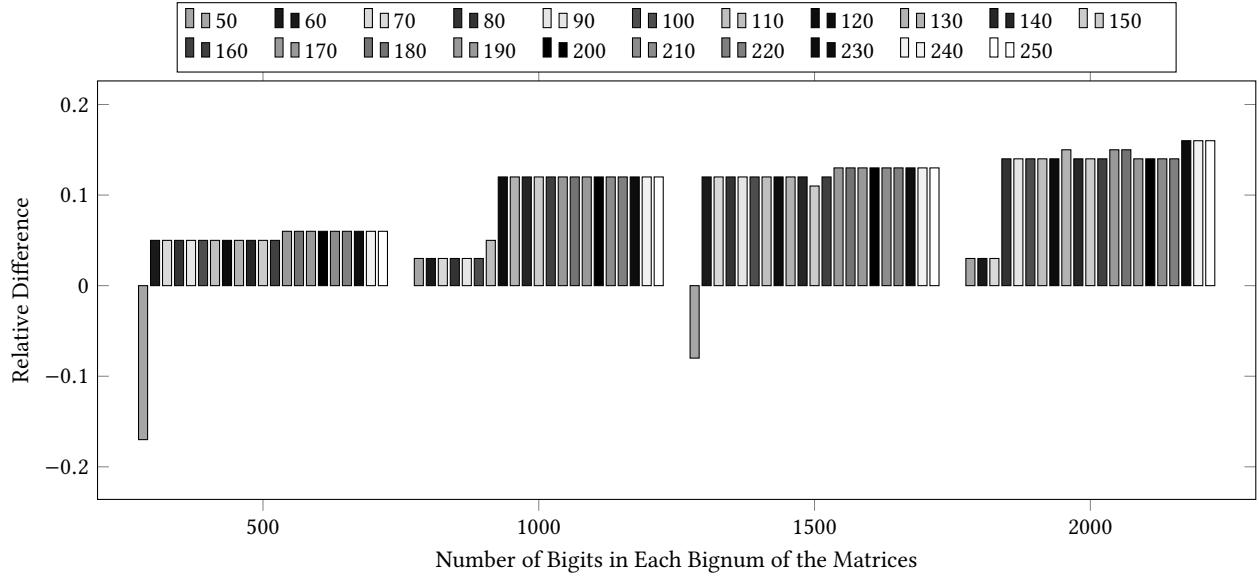


Figure 2: Karatsuba and Toom-3 Running-Time Relative Difference for MM Using Large Bignum Lengths.

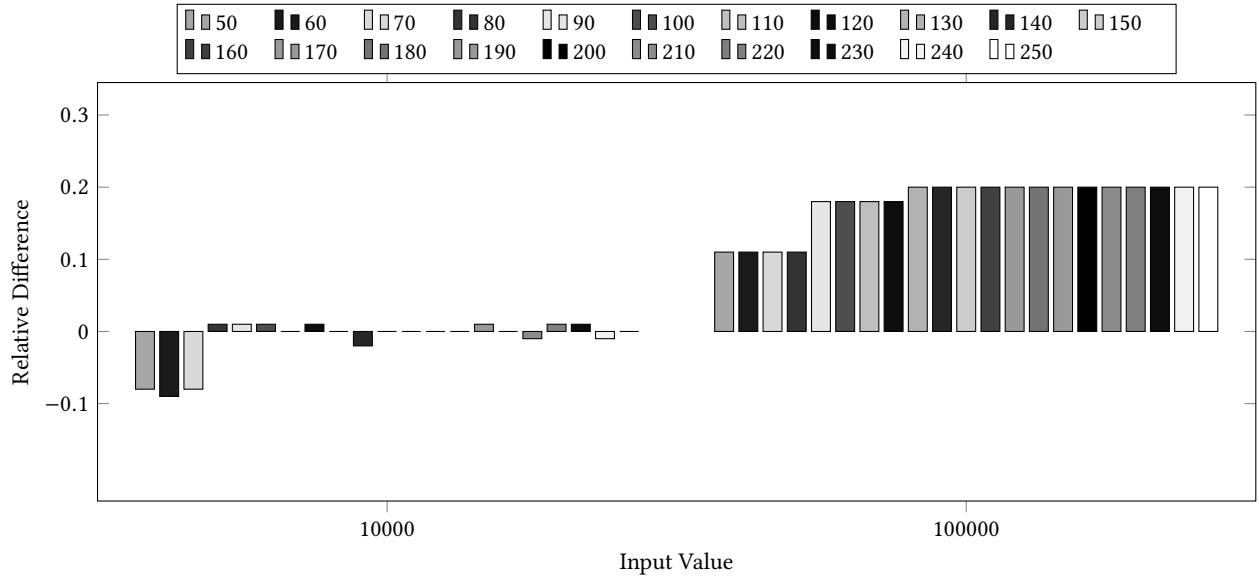


Figure 3: Karatsuba and Toom-3 Running-Time Relative Difference for Tetra

Figure 3 presents the data collected from 42 experiments using the Tetra benchmark. The x-axis here represents the inputs used (i.e., 10^{10000} and 10^{100000}). For the smaller input, 10^{10000} , there is relatively little difference between the Toom-Cook algorithm and the Karatsuba algorithm. However, for the larger input, with the exception of cut-off values of 50 and 60, there is a significant speed-up in performance when using the Toom-Cook algorithm. These speed-ups range from a minimum of 2% to a maximum speed-up of 20%. The k values for which we observe the best speed-ups are those in the range [140 – 250]. We observe a similar pattern to

the one observed using matrix multiplication: when processing shorter bignums, there is no performance benefit to using Toom-3. In addition, as with matrix multiplication, longer bignums being processed yields better performance for all k , and the length of the bignums is not a distinguishing factor to select the best value of k .

Figure 4 displays the results observed from 84 experiments using the RSA benchmark. The x-axis here represents the size of the public key which is varied to have 109, 272, 431, and 614 bits. For public key sizes of 109, 272, and 431 bits, there are negligible differences between the Karatsuba algorithm and Toom-3. However,

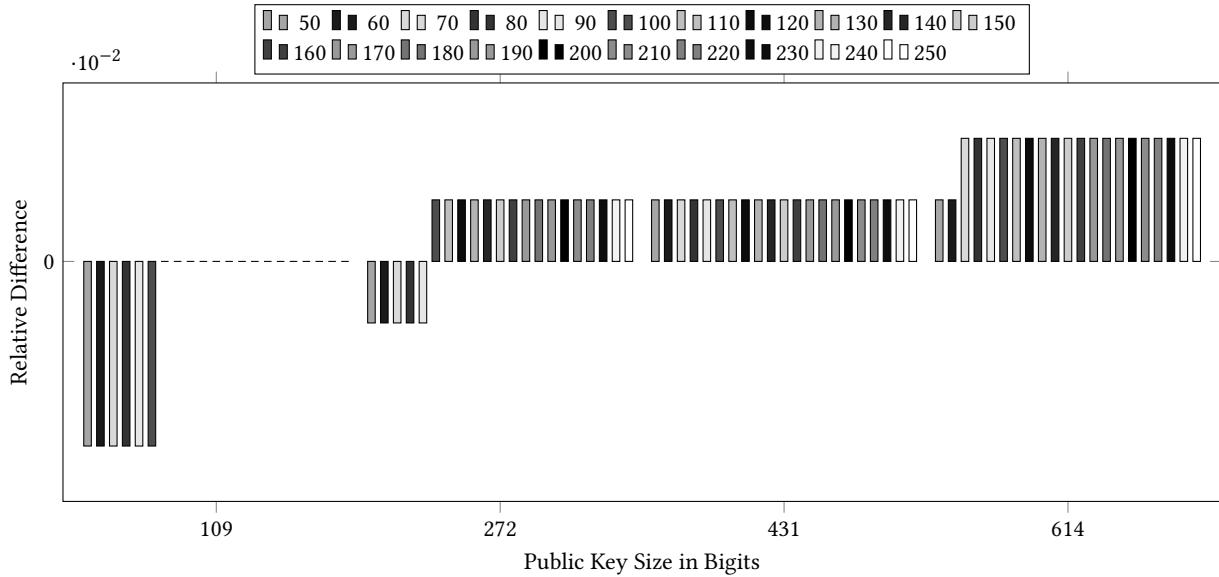


Figure 4: Karatsuba and Toom-3 Running-Time Relative Difference for RSA

for a public key size of 614 bits, there is a slight speed-up for k values between 70 and 250. All of the k values in this range exhibit a less than 1% difference, with the minimum speed-up being 1% and the maximum being 2%.

Based on the data presented so far, we conclude that all the values of k in the range [140 – 250] that were tested perform relatively well when compared to one another. The values in this range display a relative difference that is at most 1% away from the performance of any other value in the range, and from the best performing speed-up. If every bit of performance does not have to be extracted, any value in [140 – 250] may be chosen as the default k value. However, if every bit of performance must be extracted, we must analyze the consistency with which each particular value of k is among the best speed-ups. The value for k which is most consistently among the best speed-ups is 140, which performs the best for 60% of the benchmarks. Even in benchmarks where 140 does not perform the best, it is at most 2% slower than the best performer.

5 CONCLUDING REMARKS

This article explores the implementation and optimization of multiplication for a bignum library for the Green programming language. In particular, it explores the cut-off value for the number of digits in a bignum to decide when to use either the base-case algorithm or Karatsuba algorithm, and when to use the Toom-3 multiplication algorithm. The empirical evidence shows that the default value for this cut-off ought to be 140 as it exhibits the best more consistent performance for the benchmarks studied.

In addition to the conclusions offered by the empirical evidence, this article also contributes a methodology that other implementors of bignum libraries may choose to follow to fine-tune the performance of bignum multiplication. The performance of any bignum library is a delicate balance between architecture, host language, and the efficiency of an implementation. For bignum multiplication

in particular, this means being aware of the limitation of the hardware and host language, as well as the efficiency of the algorithms being implemented. In order to determine the values for the default parameters of a bignum library, an empirical study similar to the one described in this article should be performed.

Future work includes the study of memory-allocation specializations of different divide and conquer multiplication algorithms, including an empirical study on the impact of different implementation choices. Concretely, our goal is to reduce the memory footprint of bignum multiplication.

REFERENCES

- [1] Marco Bodrato. Towards optimal Toom-Cook multiplication for univariate and multivariate polynomials in characteristic 2 and 0. In Claude Carlet and Berk Sunar, editors, *WAIFI'07 proceedings*, volume 4547 of *LNCS*, pages 116–133. Springer, June 2007. <http://bodrato.it/papers/#WAIFI2007>.
- [2] Torbjörn Granlund. *GNU MP: The GNU Multiple Precision Arithmetic Library*, December 2016.
- [3] Thomas Johnsson. Lambda Lifting: Transforming Programs to Recursive Equations. In *Proc. Of a Conference on Functional Programming Languages and Computer Architecture*, pages 190–203, New York, NY, USA, 1985. Springer-Verlag New York, Inc.
- [4] A. Karatsuba and Y. Ofman. Multiplication of Multidigit Numbers on Automata. *Soviet Physics Doklady*, 7:595–596, Jan 1963.
- [5] Donald E. Knuth. *The Art of Computer Programming, Volume 2 (2nd Ed.): Seminumerical Algorithms*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 1981.
- [6] Thomas Koshy. *Catalan Numbers with Applications*. Oxford University Press, 2009.
- [7] Robert Bruce Findler Matthew Flatt and PLT. *The Racket Guide*. PLT Group, v7.0.0.6 edition, July 2018.
- [8] Mersenne.org. 50th Known Mersenne Prime Found! <https://www.mersenne.org/>, July 2018.
- [9] Marco T. Morazán. Bytecode and Memoized Closure Performance. In Jay McCarthy, editor, *Trends in Functional Programming*, pages 58–75, Berlin, Heidelberg, 2014. Springer Berlin Heidelberg.
- [10] Marco T. Morazán and Ulrik P. Schultz. Optimal Lambda Lifting in Quadratic Time. In Olaf Chitil, Zoltán Horváth, and Viktória Zsók, editors, *Implementation and Application of Functional Languages*, pages 37–56, Berlin, Heidelberg, 2008. Springer Berlin Heidelberg.

- [11] Mozilla.org. NSS 3.1 Release Notes. https://www-archive.mozilla.org/projects/security/pki/nss/release_notes_31.html, October 2000.
- [12] Winfried Neun and Herbert Melenk. Very large Gröbner basis calculations. In Richard E. Zippel, editor, *Computer Algebra and Parallelism*, pages 89–99, Berlin, Heidelberg, 1992. Springer Berlin Heidelberg.
- [13] Wade Trappe and Lawrence C. Washington. *Introduction to Cryptography with Coding Theory*. Prentice-Hall, Inc, Upper Saddle River, NJ, USA, second edition, 2005.
- [14] S. Venit and Wayne Bishop. *Elementary Linear Algebra*. Prindle, Weber, and Schmidt, 1985.

Towards Escaping the Compilation Phase with Implicit Parameters

Mark Lemay

Boston University

Boston, Massachusetts

lemay@bu.edu

ABSTRACT

Dependently typed languages have the potential to revolutionize the way software is written. However, there has been very little adoption of these tools and languages in industry. A solution may come from designing dependently typed languages to be used like widely adopted automated tests. This paper reviews some features that make automated tests popular, and determines a design criteria that should be considered by dependently typed languages that target the same users. Then this paper presents two work in progress systems that attempt to fulfill that criteria using implicit parameters.

KEYWORDS

dependent types, usability, implicit parameters

ACM Reference Format:

Mark Lemay. 2019. Towards Escaping the Compilation Phase with Implicit Parameters. In *Proceedings of International Symposium on Implementation and Application of Functional Languages (IFL'18)*. ACM, New York, NY, USA, Article 4, 7 pages. https://doi.org/10.475/123_4

1 INTRODUCTION

Dependently typed languages have the potential to revolutionize the practice of software development by allowing user defined code to be proven correct over precise user defined specifications. Dependent type systems have the some advantages over other formal methods tools. Because dependent types systems make use of the proofs as programs principle, programmers can use tools and paradigms they are already familiar with to form and analyze proofs. Despite this promise, there has been very little mainstream adoption of dependently typed programming languages or tools in the software industry.

To understand why, we should consider the current industry best practices for ensuring the correctness of code against complicated specifications. Best practices state that, as much as possible, software should be tested automatically in a separate test phase of development. This test phase is independent of compilation, so that incorrectly specified tests never prevent other unrelated software improvements. This test phase is also independent from evaluation, with test code usually not included in a release. This has the additional advantage that programs can be experimented

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

IFL'18, August 2019, Lowell, MA, USA

© 2019 Copyright held by the owner/author(s).

ACM ISBN 123-4567-24-567/08/06.

https://doi.org/10.475/123_4

with and debugged against the failing test. Testing frameworks encourage detailed communication about what went wrong and programmers are encouraged to write tests to clearly communicate why a property is important. Further, testing tools are relatively easy for programmer to use since they are implemented in the programmer's software ecosystem of choice. It is not uncommon for a project's total automated tests to run for hours. When testing takes this long, tests are run on a fixed schedule over night or on the weekend. Sometimes tests will be executed in parallel across multiple machines.

Automated tests are deterministic and decidable by design. Since they are very simple from a logical perspective, they could easily be added to the compilation phase of development. However, doing so would make automated tests much harder to use, so it makes sense that this easy compilation extension has never been implemented or adopted.

Most dependently typed languages treat complicated typed specifications much more like types are treated in mainstream programming languages like C or Java, where type checking is a necessary precondition for compilation and execution. Unlike C, satisfying a dependently typed proof constraint can be highly complicated and time consuming. However, Dependently typed systems tout their ability to verify arbitrarily complicated constraints that would normally only be partially checked by automated testing. This stands in contrast to how the industry views testing, automated testing might never have been so widely adopted if it blocked execution. Even worse is if type checking fails by some quirk of the underlying logic, for instance a property holds over a function classically, but not intuitionistically.

There are several workarounds to the high cost of proof in existing dependently typed languages. ATS¹ provides multiple versions of its standard libraries that require different levels of type specificity. Agda and Idris provide holes that allow for partial implementation of terms or types, but holes are not intended to be left in code, and will not be filled automatically (though users may manually request proof searches for a hole). Virtually all dependently typed languages languages allow assuming falsehood as a last resort to working around the type system.

A dependently typed language designed based on the mainstream software testing practices would have the following properties:

- In all popular test systems programs are always runnable when there are test errors. In a practical dependently typed system a program should always be runnable, even when type checking fails.

¹<http://www.ats-lang.org/>

- Since there are no restrictions on what code may be used in tests. This practical dependently typed system should put no restriction on the expressible constraints.
- When there is a test failure, as much information is communicated about the failure. When type checking fails, a practical dependently typed system should clearly state what proofs would be needed for the program to be well typed.
- Just as large test suites can be run passively. We should use the ample downtime in the software development cycle to search for a proof that would guarantee safety.

Designing a dependently typed language to take full advantage of the testing phase of development could have many advantages. While many current dependent type systems are very concerned with what predicates can be decided at compile time, the software industry has made clear that they would accept much less precise information that could be generated over a much longer period of time, if the information is clear and they can decide when and how to deal with failure (or lack of proof).

How should the boundary between passive constraint solutions and dependent types be drawn? Implicit parameters could be used to provide a framework to interface between the compile phase and the test/proof search phase. There are many different frameworks for implicit parameters, however the proposal to use implicits as a way to delegate proof work from compile time appears original.

To this end, we will describe 2 systems that try to fulfill these design goals. The first is a language, called the indexed implicit calculus of constructions, is based on Miquel [4]’s Implicit Calculus of Constructions and it is used to explore theoretical properties related to passive dependently typed constraint solving. The second is based on a work in progress experimental language that approaches the design goals directly. A prototype of the system described in section 4 is available online².

2 RELATED WORK

Implicits have a long history solving different problems in different programming languages, in Haskell and Scala they are used to remove inferable code that would otherwise clutter the syntax. They have been studied theoretically for different purposes in dependent types. In addition, there have been several research programs to make advanced types more pragmatic.

2.1 Typeclass Style Implicits

Perhaps the most popular form of implicits is that of Haskell’s typeclasses (Wadler and Blott [15]). Typeclasses offer a principled way to use ad hoc polymorphism: by binding implicit data to a type, this data can then be made concrete at compile time using a proof search that is guaranteed to terminate.

This mechanism is now standard in pure functional programming languages. It has been added to Agda in Devriese and Piessens [3]. In Devriese and Piessens [3] the authors make the good point that having a full Haskell style proof search is confusing when in the presence of a separate, more powerful, underlying dependently typed logic. Proof search in a dependently typed language should be over dependently typed proofs. The authors go to some lengths to avoid this in their system.

²<https://github.com/marklemay/Escaping-the-Compilation-Phase>

Typeclasses enforce that they are resolved uniquely in any context so that no unpredictable behavior can arise. This is almost the opposite use case for our system since we expect type constraints to be unsatisfied, so that they should never affect run time behavior. When the constraints are satisfied, we do not care about the details of the proof.

2.2 Scala Style Implicits

In contrast to Haskell, Scala has explored true implicit parameters that are more expressive (and more complex and dangerous than Typeclasses) in Odersky et al. [6], Oliveira et al. [7]. This unpredictability is mitigated by tooling that explains how constraints were resolved. Like typeclasses, they need to be satisfied at compile time.

2.3 Implicit Calculus of Constructions

The first system in this paper are based on the Implicit Calculus of Constructions in Miquel [4]. Miquel created a framework to extend pure type systems with an implicit Curry style intersection binder (\vee), the binder only associated variables that could be used in typing annotations, so term level computation was always well defined independent of type checking. Because of the extensive power of the ICC, type checking is undecidable. ICC was further refined in Barras and Bernardo [1] where implicit parameters were instantiated explicitly, making type checking decidable.

2.4 Refinement Typing

Refinement typing takes an existing type system and refines values of a type relative to some decidable language of constraints that are often solved with a dedicated SMT solver. This technique was pioneered in Dependent ML in Xi and Pfenning [17] (a variant of ML with support for dependent types) to resolve array bounds constraints at compile time, freeing the programmer from writing a manual proof. This approach was extended to non dependently typed languages OCaml in Rondon et al. [10] and Haskell in Vazou et al. [14] under the name liquid types. There has been some work to apply refinement types to higher order constraints in Vazou et al. [13]. The downside to this approach is that the programmer may not easily recognize the boundary of the decidable fragment of the logic. For instance Presburger arithmetic is implemented in ATS (the successor of DML), but users are sometimes confused that complicated addition constraints can be solved automatically while simple multiplication constraints require a manual proof.

Refinement systems have the advantage that they have a well defined operational semantics even when the refinement constraint checking fails.

2.5 Gradual Typing

Gradual typing (Siek and Taha [12]) is based on the philosophy that typing should be gradually made more precise as software is developed. When types are too coarse, gradual typing suggests that a run time assertion is inserted to verify the condition, and that these run time checks can be removed as the program’s type becomes more accurate. If the assertion fails, useful blame information is reported about the source of the error (Wadler and Findler [16]). While gradual typing is solving similar problems, most gradual

tying research avoids dependent types (with the notable exception of Ou et al. [8]). This makes sense when you consider that many dependent types correspond to predicates that are not amenable to run time checks (for instance if $f : \mathbb{N} \rightarrow \mathbb{N}$ where refined to $f' : (\mathbb{N} \rightarrow \mathbb{N}) |\exists n : \mathbb{N}. f(n) = 0$, which is undecidable). Furthermore, there is no research I am aware of that advocates for work outside of the compile time and run time phases.

3 INDEXED IMPLICIT CALCULUS OF CONSTRUCTIONS

The first language is based on Miquel's ICC (Miquel [4]) without η expansions, and only the 2 universes of the original calculus of constructions. The novel change is to assume that implicits are resolved from an indexed set D of predefined judgments. D is meant to represent the results of some long running constraint satisfying search.

When describing IIIC indexed as a specific set of judgments, D , we write IIIC_D .

3.1 Syntax and Abbreviations

Variables will be presented with symbols x, y, z, X, Y , and Z

Sorts are defined as

$$s, s_1, s_2 ::= \bullet \mid \square$$

Terms are defined as

$$a, b, c, f, A, B, C, F ::= x \mid s \mid \prod x : A.B \mid \forall x : A.B \mid \lambda x.B \mid FA$$

We will use capitalized metacharacters for syntax that is intended to be a type ($A : \bullet$), and lower case metacharacters that are intended as a term ($a : A$). We will use the common abbreviation of $A \rightarrow B$ for $\Pi x : A.B$ when x does not appear in B . Likewise we will use $A \Rightarrow B$ as an abbreviation for $\Pi_i x : A.B$ when x does not appear in B . This syntax will closely match the syntax in section 4. The standard rules for contexts are assumed and not stated. Note that lambdas do not type their bound variables.

3.2 Typing Judgments

Given a set of judgments D , and set of Sorts $S = \{\bullet, \square\}$ the typing derivations for IIIC_D are

$$\begin{array}{c} \frac{}{\Gamma \vdash \bullet : \square} \\ \frac{\Gamma \vdash A : s}{\Gamma, x : A \vdash x : A} \\ \frac{\Gamma \vdash A : B \quad \Gamma \vdash C : s}{\Gamma, x : C \vdash A : B} \\ \frac{\Gamma \vdash A : s_1 \quad \Gamma, x : A \vdash B : s_2}{\Gamma \vdash (\Pi x : A.B) : s_2} \\ \frac{\Gamma \vdash A : s_1 \quad \Gamma, x : A \vdash B : s_2}{\Gamma \vdash (\forall x : A.B) : s_2} \\ \frac{\Gamma \vdash F : (\Pi x : A.B) \quad \Gamma \vdash a : A}{\Gamma \vdash Fa : B[x := a]} \end{array}$$

$$\begin{array}{c} \frac{\Gamma \vdash F : (\forall x : A.B) \quad (\Gamma \vdash a : A) \in D}{\Gamma \vdash F : B[x := a]} \\ \frac{\Gamma \vdash (\Pi x : A.B) : s \quad \Gamma, x : A \vdash b : B}{\Gamma \vdash (\lambda x.b) : (\Pi x : A.B)} \\ \frac{\Gamma \vdash (\forall x : A.B) : s \quad \Gamma, x : A \vdash b : B \quad x \notin \text{FV}(b)}{\Gamma \vdash b : (\forall x : A.B)} \\ \frac{\Gamma \vdash A : B \quad \Gamma \vdash B' : s \quad \Gamma \vdash B =_{\beta} B'}{\Gamma \vdash A : B'} \end{array}$$

3.3 Relation to ICC

When D is the set of all derivations in ICC without η expansions, IIIC_D is exactly ICC without η expansions. A number of good properties follow immediately from this. For any D ,

- A sequence of reductions in IIIC_D will also be a sequence of reductions in ICC
- A derivation in IIIC_D will correspond to a derivation in ICC
- A term typed in IIIC_D will always beta reduce to a value. IIIC_D is normalizing.
- A term typed in IIIC_D will never “get stuck”
- IIIC_D is consistent, false cannot be defined in an empty context

Unfortunately the type preservation property (subject reduction) does not translate as smoothly

3.4 Example D such that preservation does not hold

We can see one pathological case where preservation does not hold where D is the smallest set such that

$(X : \bullet, x : X \vdash x : X) \in D$ and if $(\Gamma \vdash a : A) \in D$ and $\Gamma \subseteq \Gamma'$ then $(\Gamma' \vdash a : A) \in D$

Because

$$\vdash \lambda X. \lambda x. \lambda f. f$$

$$: \Pi X : \bullet. \Pi x : X. \forall Y : \bullet. \Pi f : X \Rightarrow Y. Y$$

The so this term admits the following typing judgment

$$\vdash (\lambda X. \lambda x. \lambda f. f) (\forall X : \bullet. X \rightarrow X)$$

$$: \Pi x : (\forall X : \bullet. X \rightarrow X). \forall Y : \bullet. \Pi f : (\forall X : \bullet. X \rightarrow X) \Rightarrow Y. Y$$

But if the term is reduced, then the entire expression reduces to

$$\vdash (\lambda X. \lambda x. \lambda f. f) (\forall X : \bullet. X \rightarrow X)$$

$$\rightsquigarrow_{\beta}$$

$$\vdash \lambda x. \lambda f. f$$

which is no longer typeable to

$$: \Pi x : (\forall X : \bullet. X \rightarrow X). \forall Y : \bullet. \Pi f : (\forall X : \bullet. X \rightarrow X) \Rightarrow Y. Y$$

The judgment set needs to satisfy some properties so cases like these are avoided.

3.5 Conditions on D for preservation

The following conditions are apparently sufficient for preservation based on modifying the lemmas in Miquel [4]:

- D is closed under weakening: if $\Gamma \subseteq \Gamma'$ and $(\Gamma \vdash a : A) \in D$ then $(\Gamma' \vdash a : A) \in D$
- D is closed under context substitution: if $(\Gamma, y : C; \Gamma' \vdash a : A) \in D$ and $\Gamma \vdash c : C$ then $(\Gamma; (\Gamma' [y := c]) \vdash a[y := c] : A[y := c]) \in D$
- D is closed under beta reduction, if $(\Gamma \vdash a : A) \in D, \Gamma \rightsquigarrow_{\beta} \Gamma',$ and $A \rightsquigarrow_{\beta} A'$ then $(\Gamma' \vdash a : A') \in D$

It is easy to extend the any set of derivation so the first condition holds. The last condition is not too worrisome because every well typed term normalizes.

The second condition is more problematic, it is not clear how to extend derivations across every replacement in general. It does hint that “closed” solutions that do not rely on the context are particularly valuable.

3.6 Example D satisfying “obvious” equalities

Assuming a standard encoding of equality and it’s reflective proof, say that for all types $\Gamma \vdash A$ and all terms $\Gamma \vdash a : A,$

$$(\Gamma \vdash \text{refl} : a =_A a) \in D$$

This satisfies the conditions above.

It may be possible to expand this example with rewrite rules, allowing for a more convenient use of equality. This seems similar to the results in Blanqui [2] where type checking on the calculus of constructions is extended with rewrite rules.

3.7 Design Criteria

IICC satisfies only some of the design criteria

- Implicit binders give a clear operational semantic even if type constraints are not yet satisfied
- IICC allows for very expressive constraints
- Type constraints cannot be clearly communicated on type check failure. Furthermore since it is a Curry style type system even the limited type checking we are interested in is undecidable.
- The D allows for constraints to be accumulated for use in type checking

4 EXPERIMENTAL SYSTEM

We have also been working on an experimental system that is intended to be more pragmatic. It has the interesting feature that type inference is internalized with implicit constraints. Like before the system is indexed by a set of judgments D that is meant to represent the output of some offline proof search procedure. The key feature of this language is that explicit application occurs independently of implicit binders, and application are validated using implicit equality checks. Where IICC was very open about what could be resolved as an implicit, this system only allows equality constraints to be satisfied, and then only by proofs that are themselves equal to the computational content of `refl`. Another difference from IICC

is that implicit resolution is considered a reduction step, and implicitly bound variables may appear inside terms. Because of this there needs to be both a safe and unsafe evaluation scheme.

This language is implementation driven, with desirable properties (Church-Rosser, preservation, termination) experimentally verified by extensive generative testing. Developing a language like this allows for relatively rapid prototyping.

4.1 Syntax and Abbreviations

Sorts are defined as

$$s, s_1, s_1 ::= \bullet \mid \square$$

Terms are defined as

$$a, b, c, f, A, B, C, F ::= x \mid s \mid \prod x : A.B \mid \Pi_i x : A.B \mid \lambda x : A.B \mid FA$$

$$\mid \lambda_i x : A.b$$

Note that in this system every variable is bound with a type. We will maintain the conventions from section 3. When binders become “irreverent”, they will be marked with a slash. For instance, $\Pi_i x : A.B$ will be written $\Pi_i x : A.B$ if x does not appear in $b : B.$ This is only meta information meant for the programmer, and not part of the actual type system.

Additionally there are term symbols that can be defined as part of the syntax

$$\text{cast}_{A,B}(e, a) \mid A =_C B \mid \text{refl}_A(a)$$

4.2 Safe Reductions

Since typing constraints are used to reduce terms, safe reduction needs to keep track of the typing context.

$$\frac{\Gamma \vdash A : s \quad \Gamma \vdash a : A \quad \Gamma, x : A \vdash b : B}{\Gamma \vdash (\lambda_i x : A.b) \{a\} \rightsquigarrow_{\text{safe}} b[x := a]}$$

The implicit binders λ_i and Π_i as fulfilling the role of λ and Π in the regular lambda calculus.

$$\frac{\Gamma \vdash A : s \quad \Gamma, x : A \vdash b : B \quad \Gamma \vdash c : C}{\Gamma \vdash (\lambda_i x : A.b) c \rightsquigarrow_{\text{safe}} (\lambda_i x : A.b c)}$$

There is a higher priority applications that will propagate past implicits until a λ is found.

$$\frac{\Gamma \vdash A : s \quad \Gamma \vdash c : C \quad \Gamma \vdash C : s \quad \Gamma, x : A \vdash b : B}{\Gamma \vdash (\lambda x : A.b) c \rightsquigarrow_{\text{safe}} (\lambda_i e : A = C.b[x := \text{cast}_{C,A}(e, c)])}$$

When the application finds a λ to bind with an implicit equality constraint is created, the cast blocks application until e is resolved to a value.

$$\frac{\Gamma \vdash A : s \quad \Gamma, x : A \vdash b : B \quad (\Gamma \vdash a : A) \in D}{\Gamma \vdash (\lambda_i x : A.b) \rightsquigarrow_{\text{safe}} (\lambda_i x : A.b[x := a])}$$

Some equality implicits can be satisfied automatically. when this happens the binder is marked for user reporting reasons. They act exactly as implicits normally would. This seems to allow for subject reduction to hold while allowing the programmer to know when a constraint has already been satisfied.

$$\frac{\Gamma; \Gamma' \vdash A \rightsquigarrow_{\text{safe}} A' \quad \Gamma \vdash B[A] : C}{\Gamma \vdash B[A] \rightsquigarrow_{\text{safe}} B[A']}$$

Finally, reductions in sub-terms induces reductions in outer terms. These rules appear to have type preservation for suitable D .

4.3 Optimistic Reductions

When in a “run-time” environment we ignore equality constraints by adding the additional equality rule

$$\frac{}{\Gamma \vdash \text{cast}_{C,A}(e, c) \rightsquigarrow c}$$

This is a reasonable choice since the only syntactic way to satisfy equality is with a proof of refl . Which has the content of identity. adding this rule removes mean that a simple notion of type preservation no longer holds.

4.4 Abbreviations

The system is rich enough to encode some internal notion of equality. However this could also be done with an additional inductive definition.

$$A =_C B \quad \text{:} \equiv \Pi_i P : C \Rightarrow \bullet.P A \Rightarrow P B$$

$$\text{cast}_{A,B}(e : (B =_\bullet A), a : A) \quad \text{:} \equiv e \{\Pi_i x : \bullet.x \Rightarrow B\} \{\lambda_i x : B.x\} \{a\}$$

$$\text{refl}_A(a : A) \quad \text{:} \equiv \lambda_i P : A \Rightarrow \bullet.\lambda_i x : P a.x$$

Note that every binder is implicit, so it will not generate a new tower of equality constraints.

4.5 Typing Judgments

Here is a summery of the typing judgments relative to a set of derivations D that are limited to judgments of equality.

$$\begin{array}{c} \frac{}{\Gamma \vdash \bullet : \square} \\ \frac{\Gamma \vdash A : s}{\Gamma, x : A \vdash x : A} \\ \frac{\Gamma \vdash A : B \quad \Gamma \vdash C : s}{\Gamma, x : C \vdash A : B} \\ \frac{\Gamma \vdash A : s_1 \quad \Gamma, x : A \vdash B : s_2}{\Gamma \vdash (\Pi x : A.B) : s_2} \\ \frac{\Gamma \vdash A : s_1 \quad \Gamma, x : A \vdash B : s_2}{\Gamma \vdash (\lambda_i x : A.B) : s_2} \\ \frac{\Gamma \vdash (\Pi x : A.B) : s \quad \Gamma, x : A \vdash b : B}{\Gamma \vdash (\lambda x : A.b) : (\Pi x : A.B)} \\ \frac{\Gamma \vdash (\Pi x : A.B) : s \quad \Gamma, x : A \vdash b : B}{\Gamma \vdash (\lambda_i x : A.b) : (\Pi_i x : A.B)} \\ \frac{\Gamma \vdash F : (\Pi_i x : A.B) \quad \Gamma \vdash c : C \quad \Gamma, x : A, y : B \vdash y c : E}{\Gamma \vdash F c : (\Pi_i x : A.E)} \\ \frac{\Gamma \vdash F : (\Pi x : A.B) \quad \Gamma \vdash A : \bullet \quad \Gamma \vdash c : C}{\Gamma \vdash F c : (\Pi_i e : A =_\bullet C. B[x := \text{cast}_{C,A}(e, c)])} \end{array}$$

$$\frac{\Gamma \vdash F : (\Pi_i x : A.B) \quad \Gamma \vdash a : A}{\Gamma \vdash F \{a\} : (B[x := a])}$$

$$\frac{\Gamma \vdash F : (\Pi_i x : A.B) \quad (\Gamma \vdash a : A) \in D}{\Gamma \vdash F : (\Pi_i x : A.(B[x := a]))}$$

Like the reduction rules above, we mark when an implicit binder becomes irrelevant. The reason the binder is not simply removed is so that explicit applications are still possible.

$$\frac{\Gamma \vdash A \rightsquigarrow_{\text{safe}} A' \quad \Gamma \vdash a : A}{\Gamma \vdash a : A'}$$

For suitable D type-checking appears decidable.

4.6 Example

Here is a simple term that can be encoded within the system

$$\lambda_i X : \bullet.\lambda_i f : X \rightarrow X.(\lambda_i Y : \bullet.\lambda g : Y \rightarrow Y.\lambda y : Y.g(g y)) f$$

It has type assuming D is empty.

$$\Pi_i X : \bullet.\lambda_i \Pi_i : X \rightarrow X. \Pi_i Y : \bullet.\Pi_i e : Y \rightarrow Y =_\bullet X \rightarrow X.$$

$$\Pi y : Y. \Pi_i e' : Y \rightarrow Y =_\bullet Y \rightarrow Y. \Pi_i e'' : Y =_\bullet Y. Y$$

But assuming the obvious equalities are in D

$$\Pi_i X : \bullet.\lambda_i \Pi_i : X \rightarrow X. \Pi_i Y : \bullet.\Pi_i e : Y \rightarrow Y =_\bullet X \rightarrow X.$$

$$\Pi y : Y. \Pi_i e' : Y \rightarrow Y =_\bullet Y \rightarrow Y. \Pi_i e'' : Y =_\bullet Y. Y$$

This shows how an improved version of this system could gracefully handle type inference. Even when type inference is undecidable.

Also

$$\lambda_i X : \bullet.\lambda_i f : X \rightarrow X.(\lambda_i Y : \bullet.\lambda g : Y \rightarrow Y.\lambda y : Y.g(g y)) f$$

\rightsquigarrow

$$\lambda_i X : \bullet.\lambda_i f : X \rightarrow X. \lambda_i Y : \bullet.\lambda_i e : Y \rightarrow Y =_\bullet X \rightarrow X.. \lambda y : Y.f(f y)$$

Where the safe evaluation protects the calls to f with casts.

4.7 Example: empty context

In the untyped lambda calculus it is possible to encode the a fixed point

$$(\lambda x.x x)(\lambda x.x x)$$

It is possible to encode a version of that term in this system

$$\lambda_i b : (\Pi z : \bullet.z). \lambda_i X : \bullet.$$

$$((\lambda x : X \rightarrow X. (x x) \{b...\}) (\lambda x : X \rightarrow X. (x x) \{b...\}) \{b...\})$$

Where the ellipses merely restate the equality type needed for conversion. This term type-checks to

$$\Pi_i b : (\Pi z : \bullet.z) . \Pi_i X : \bullet.X$$

This should communicate to the programmer that to safely run this function there must be a term of type $\Pi x : X$. Since that term will never be satisfiable, the program can never be run and type checking is “safe”. A programmer may run the program unsafely to observe the non termination first hand.

4.8 Design Criteria

The experimental system satisfies only some of the design criteria

- Because D is limited to equality derivations, there is a clear operational semantics when typing fails
- The system allows for very expressive constraints, but only equalities are resolved automatically
- The system is designed to clearly propagate unsatisfied constraints
- The D allows for equality constraints to be accumulated for use at compile time

While right now D is only allowed over equality constraints we hope to generalize to any type that is term/proof irreverent. There would need to be a way to specify the computational behavior of these terms for use in untyped reduction.

5 FUTURE WORK

There is a large amount of work still to do.

- There are several improvements that could be made to the experimental system
 - The experimental system still does not have a clean way to back solve equality constraints. For instance, roughly $\Gamma; \text{cons} : (\forall A : \bullet. \forall n : \mathbb{N}. A \rightarrow \text{Vect } An \rightarrow \text{Vect } An + 1), ns : (\text{Vect } \mathbb{N} 8) \vdash \text{cons} 3\ ns : \forall A : \bullet. \forall n : \mathbb{N}. (A = \mathbb{N}) \Rightarrow (\text{Vect } An = \text{Vect } \mathbb{N} 3) \Rightarrow \text{Vect } An + 1$ when clearly A and n should be specialized to the only consistent values for them to inhabit.
 - Generalize this style of implicit resolution beyond equality constraints.
 - Adding additional features to test how viable this scheme is as a practical programming language
 - Earlier in development subtyping based on coercing implicits was considered. The feature was more complicated then needed for this demonstration. But it may be reincorporating this subtyping back into the system as the the experimental system become more mature.
 - When the system becomes more stable it would be good to prove the properties that have only been checked experimentally to this point.
- There are many improvements that can be made to IICC
 - Find exactly what a necessary and sufficient condition on D is for IICC for subject reduction to be preserved.
 - It may be possible to loosen the restrictions on D if the terms that are known in advance. For instance, if $d \in D$, d would not need to be closed under every contextual replacement, but only those reachable by the terms and other derivations of D . This should be explored.

- Define a framework such that different indexed sets of derivations can be combined to maintain desirable properties.
- Formalize the proofs for IICC in a theorem proving system.
- Exploring what other advantages a passive implicit system might have
 - It could be used to avoid awkward and unpredictable heuristics that are limited to type checking or elaboration.
 - It may be possible to resolve universe hierarchies as implicit constraints.
 - Passive search could be a good match for type directed program synthesisPolikarpova et al. [9], since there is no need that the synthesis be decidable at compile time, and some approximation of the synthesis can be exported to run time (possibly proposed in Mirman [5]).
 - Implicits could use the constraint for generative testing.
 - The proof search over implicit constraints could still be useful even if it just uncovers some simpler related lemmas.
- Empirically study how software practitioners assure the correctness of their programs and what trade-offs they make.
- Compare with to the Cochis frameworkSchrijvers et al. [11]

6 CONCLUSION

Dependently typed languages will not become mainstream until a number of usability issues are addressed. This paper presents a first step in one direction of typing usability, by trying to adapt to existing software practices.

7 ACKNOWLEDGMENTS

Thanks goes to Hongwei Xi for the helpful feedback, and Stephanie Savir for numerous spelling and grammatical corrections.

REFERENCES

- [1] Bruno Barras and Bruno Bernardo. 2008. The implicit calculus of constructions as a programming language with dependent types. In *International Conference on Foundations of Software Science and Computational Structures*. Springer, 365–379.
- [2] Frédéric Blanqui. 2005. Definitions by rewriting in the Calculus of Constructions. *Mathematical Structures in Computer Science* 15, 1 (2005), 37–92.
- [3] Dominique Devriese and Frank Piessens. 2011. On the bright side of type classes: instance arguments in Agda. *ACM SIGPLAN Notices* 46, 9 (2011), 143–155.
- [4] Alexandre Miquel. 2001. The implicit calculus of constructions extending pure type systems with an intersection type binder and subtyping. In *International Conference on Typed Lambda Calculi and Applications*. Springer, 344–359.
- [5] Matthew Mirman. 2014. *Logic Programming and Type Inference with the Calculus of Constructions*. Ph.D. Dissertation. Carnegie Mellon University Pittsburgh, PA.
- [6] Martin Odersky, Olivier Blanvillain, Fengyun Liu, Aggelos Biboudis, Heather Miller, and Sandro Stucki. 2017. Simplicity: foundations and applications of implicit function types. *Proceedings of the ACM on Programming Languages* 2, POPL (2017), 42.
- [7] Bruno Cds Oliveira, Tom Schrijvers, Wontae Choi, Wonchan Lee, and Kwangkeun Yi. 2012. The implicit calculus: a new foundation for generic programming. In *ACM SIGPLAN Notices*, Vol. 47. ACM, 35–44.
- [8] Ximming Ou, Gang Tan, Yitzhak Mandelbaum, and David Walker. 2004. Dynamic typing with dependent types. In *Exploring new frontiers of theoretical informatics*. Springer, 437–450.
- [9] Nadia Polikarpova, Ivan Kuraj, and Armando Solar-Lezama. 2016. Program synthesis from polymorphic refinement types. *ACM SIGPLAN Notices* 51, 6 (2016), 522–538.
- [10] Patrick M Rondon, Ming Kawaguci, and Ranjit Jhala. 2008. Liquid types. In *ACM SIGPLAN Notices*, Vol. 43. ACM, 159–169.
- [11] Tom Schrijvers, Bruno C d S Oliveira, and Philip Wadler. 2017. Cochis: Deterministic and coherent implicits. (2017).

- [12] Jeremy G Siek and Walid Taha. 2006. Gradual typing for functional languages. In *Scheme and Functional Programming Workshop*, Vol. 6. 81–92.
- [13] Niki Vazou, Alexander Bakst, and Ranjit Jhala. 2015. Bounded refinement types. In *ACM SIGPLAN Notices*, Vol. 50. ACM, 48–61.
- [14] Niki Vazou, Eric L Seidel, Ranjit Jhala, Dimitrios Vytiniotis, and Simon Peyton-Jones. 2014. Refinement types for Haskell. In *ACM SIGPLAN Notices*, Vol. 49. ACM, 269–282.
- [15] Philip Wadler and Stephen Blott. 1989. How to make ad-hoc polymorphism less ad hoc. In *Proceedings of the 16th ACM SIGPLAN-SIGACT symposium on Principles of programming languages*. ACM, 60–76.
- [16] Philip Wadler and Robert Bruce Findler. 2009. Well-typed programs can't be blamed. In *European Symposium on Programming*. Springer, 1–16.
- [17] Hongwei Xi and Frank Pfenning. 1999. Dependent types in practical programming. In *Proceedings of the 26th ACM SIGPLAN-SIGACT symposium on Principles of programming languages*. ACM, 214–227.

Verifying Embedded Attribute Grammars

Florent Balestrieri
fbalestrieri@orange.fr

Alberto Pardo
Instituto de Computación
Universidad de la República
Uruguay
pardo@fing.edu.uy

Marcos Viera
Instituto de Computación
Universidad de la República
Uruguay
mviera@fing.edu.uy

ABSTRACT

In this paper, we present an embedding of attribute grammars in Haskell, that is both modular and type-safe, while providing the user with domain specific error messages.

Our approach involves to delay part of the safety checks to runtime. When a grammar is correct, we are able to extract a function that can be run without expecting any runtime error related to the EDSL.

KEYWORDS

Attribute Grammars, EDSL, Dynamics

ACM Reference Format:

Florent Balestrieri, Alberto Pardo, and Marcos Viera. 2019. Verifying Embedded Attribute Grammars. In *Proceedings of International Symposium on Implementation and Application of Functional Languages (IFL'18)*. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/nmnnnnnn.nmnnnnnn>

1 INTRODUCTION

This paper is in the continuity of the work of Oege de Moor, Kevin Backhouse and S. Doaitse Swierstra who proposed that the combination of attribute grammars [10] and functional programming provide a powerful toolset for the rapid prototyping of domain specific languages. Dedicated attribute grammar systems may lack the flexibility of a general purpose programming language, whereas functional programming languages may lack the structuring mechanisms provided by attribute grammars.

In their paper [4], O. de Moor and others give a compositional semantics of the structuring mechanisms of attribute grammars. The semantics is given as a Haskell program and can be easily extended with new structuring operators, enriching the attribute grammar idiom. The implementation of the attribute grammar operators is elegant and allows us to define attribute grammars modularly by combining smaller attribute grammar fragments together and hence enables code reuse. The main drawback of their approach is that it fails to capture many important well-definedness properties of attribute grammars in the type system, hence many of the operations on attribute grammars are partial.

Previous attempts to use Haskell's type system to make structural checks on attribute grammars relied on complex types making use of many type system extensions of GHC [15], [2, § 8.4,p. 121]. The deal breaker however is that the error messages related to the attribute grammars reveal implementation details that shouldn't be shown to the users of the library: types that are private to the library are leaked out when printing error messages, and are extremely

difficult to understand for someone who is not familiar with the implementation.

In this paper we present an implementation of first-class attribute grammars in Haskell that does not compromise on safety, nor on the simplicity of the public interface nor on the clarity of error messages. It relies on a combination of static and dynamic type-checking. The dynamic type-checking is all done in one single global verification of the attribute grammar. Upon success, it returns a function that can be executed as many times as necessary, without any further dynamic checks. Upon failure, it returns a precise error message related to the attribute grammar being constructed, and by using one of GHC extensions, may even point to the relevant position in the source code.

In Section 2 we introduce the public interface of the library through a simple example. In Section 3 we explain the internals of the library, while in Section 4 we show how some common patterns of attribute grammars can be encoded using its public interface. In Section 5 we discuss some related work and in Section 6 we finally conclude.

2 PUBLIC INTERFACE OF THE LIBRARY

An attribute grammar is a context free grammar to which we add attribution rules, also called semantic rules, that define the value of attributes for each non-terminal of the grammar.

Using attribute grammars with the library involves the following steps.

- Define the elements of a context free grammar: non-terminals, productions, non-terminal children of a production, terminal attributes.
- Define attribute rules and aspects.
- Combine aspects to form a complete attribute grammar.
- Associate grammar elements to a concrete datatype and obtain, after validity checks, a function that evaluates attributes on a value of such datatype.

The whole process is modular. Only the last step requires us to fix the datatype and the corresponding grammar on which the attributes will be computed as well as the concrete rules to compute them.

In the following subsections we will explain the interface of the library through a running example. We will use a very simple expression grammar, with integers, variables and addition, with the following concrete grammar:

```
<expr> ::= <int> | <var> | <expr> "+" <expr>
```

The values of the variables are provided in a given context. For example, the evaluation of the expression $2 + x$ in the context $x = 7$ results in 9.

```

expr   = non_terminal "Expr"
val    = production expr "Val" [] (constT ival nilT)
var    = production expr "Var" [] (constT vname nilT)
add    = production expr "Add" [leftAdd, rightAdd] nilT
leftAdd = child add "leftAdd" expr
rightAdd = child add "rightAdd" expr
ival   = terminal "ival" pInt
vname  = terminal "vname" pString

```

Figure 1: Grammar

2.1 Context free grammar

There is a direct correspondence between context free grammars and mutually recursive algebraic datatypes. When we will want to evaluate an attribute grammar, we will attach the grammar to a concrete datatype. Unlike Haskell datatypes, our grammars are modular and extensible. We can add new datatypes (non-terminals) and new constructors (productions).

We define a Grammar as a set of productions.

```

type Grammar = Set Production
grammar :: [Production] → Grammar

```

The types *NonTerminal*, *Production*, *Child*, *Terminals* are abstract, and represent the different elements of a grammar. The following functions are their smart constructors, which we use in Figure 1 to define the grammar for the expression language:

- *non_terminal* :: *Name* → *NonTerminal*
Defines a non-terminal given its name.
- *production* :: *NonTerminal* → *Name* → [*Child*]
→ *Terminals* → *Production*

Defines a production given the non-terminal to which it belongs, the name of the production and the lists of children (non-terminals) and terminals that compose it.

- *child* :: *Production* → *Name* → *NonTerminal* → *Child*
Defines a child out of the production to which it belongs, its name and non-terminal.
- *terminal* :: *Typeable a* ⇒ *Name* → *Proxy a*
→ *Terminal a*

Defines a terminal given its name, and its type. A type proxy is used to specify the type parameter *a* without giving a value of type *a*. For instance, *pInt* has type *Proxy Int*.

Notice that we provide names¹ to all the ingredients of a grammar; i.e. non-terminals, productions, children and terminals. Also notice that the grammar definition is modular, since all its ingredients are defined separately. Although, typically a production and its children are mutually recursive, like the production *add* and its children *leftAdd* and *rightAdd* of Figure 1².

The list of terminals (*Terminals*) is abstract and is constructed using the functions *nilT* and *constT*.

¹The type *Name* is a synonym of *String*.

²The library provides alternative constructors to define a production at the same time as its children

```

eval = attr S "eval" tEval
env  = attr I "env"  tEnv

```

Figure 2: Attributes

```

nilT  :: Terminals
constT :: Typeable a ⇒ Terminal a → Terminals
        → Terminals

```

Thus, in the example the production *add* has two children (*leftAdd* and *rightAdd*) and does not have non-terminals, while *val* has a terminal *ival*, representing the integer value, and *var* has a terminal *vname*, for the name of the variable.

Accessors. *prod_nt* and *prod_children* respectively return the non-terminals and children of a production, *child_prod* returns the production of a child.

2.2 Attributes

Attributes provide semantics to a grammar. In our case, the semantics we want to define is the evaluation of an expression. We define two attributes:

- **env** distributes the environment from the root of the tree to the leaves;
- **eval** computes the result of the evaluation from the leaves to the root.

The former is said to be an *inherited attribute* because it is computed in a top-down traversal where the information flows from the root to the leaves, while the latter is a *synthesized attribute* because it is computed in a bottom-up traversal, with the information flowing from the leaves to the root.

The type (*Attr k a*) of attributes of kind *k* (e.g. inherited) and type *a* is abstract. We declare new attributes with:

```

attr :: Typeable a ⇒ Kind k → Name → proxy a
      → Attr k a

```

Kinds are given by three abstract types *I* (inherited), *S* (synthesized) and *T* (terminal), and are specified using the homonymous constructors of the type *Kind k*.

```

data Kind k where
  I :: Kind I
  S :: Kind S
  T :: Kind T

```

In our approach terminals are treated as a special kind of attributes, called *terminal attributes*.

```
type Terminal a = Attr T a
```

Thus, the function *terminal*, that defines terminals, applies *attr* to the Kind *T*:

```
terminal = attr T
```

We declare the attributes of our example in Figure 2; *tEval* and *tEnv* have types *Proxy Int* and *Proxy (String :-> Int)*³, respectively.

³We use the type operator (*:>*) as a synonym of *Map*.

2.3 Aspects

Aspects are collections of attribution rules. They are constructed using the functions *inh* and *syn* which define the attribution rules for a single inherited and synthesized attribute, respectively.

```
inh :: (Typeable a) ⇒ Attr I a → [(Child,      AR a)]
      → Aspect
syn :: (Typeable a) ⇒ Attr S a → [(Production, AR a)]
      → Aspect
```

Inherited attributes are computed top-down, thus for a given attribute (*Attr I a*), we declare the attribution rules (*AR a*) to compute the values for each child (*Child*). In the case of synthesized attributes (*Attr S a*), since the information flows bottom-up, the rules are defined per production.

Attribution Rules. *AR* is the applicative functor [13] of *attribution rules*. The following *AR* primitives are used to access attributes:

- $(!)$:: (*Typeable a*) ⇒ *Child* → *Attr S a* → *AR a*
projects a synthesized attribute from a child.
- par :: (*Typeable a*) ⇒ *Attr I a* → *AR a*
projects an inherited attribute from the parent.
- ter :: (*Typeable a*) ⇒ *Attr T a* → *AR a*
projects a terminal attribute.

For example, the following code declares the set of rules to define the way the *eval* attribute is computed:

```
evalA = syn eval [ add |- (+) <⇒ (leftAdd ! eval)
                  & (rightAdd ! eval)
                  , val |- ter ival
                  , var |- slookup <⇒ ter vname
                  & par env ]
```

The operator $(|-)$ is a synonym of the pair constructor $(,)$ that provides a more readable syntax for association lists. For the production *add* we sum the values of the synthesized *eval* attributes from its children. For *val* is the terminal *ival* representing its value. In the case of *var* the value associated to the variable has to be looked-up into the inherited environment.

The *env* attribute is declared using *inh*.

```
envA = inh env [ leftAdd |- par env
                 , rightAdd |- par env ]
```

For both children *leftAdd* and *rightAdd* of the *add* production, the value is just the inherited *env* of the parent. This is a common pattern when defining inherited attributes, called *copy rule*. We capture this pattern with the *copyPs* combinator, which defines an aspect given an attribute and the list of productions where it has to be applied. Thus we can define *envA* in the following way.

```
envA = copyPs env [ add ]
```

In the previous definitions we used the applicative combinators $\&$ and $\&&$, Figure 3 includes an alternative, less verbose, formulation that uses a Template Haskell-based syntactic sugar for applicatives provided in the library.

```
evalA
= syn eval [ add |- [[(leftAdd ! eval) + (rightAdd ! eval)]]
            , val |- ter ival
            , var |- [[fromJust $ Map.lookup (ter vname)
                      (par env)]]]
          ]
envA = copyPs env [ add ]
asp  = evalA # envA
```

Figure 3: Aspects

Merging Aspects. Aspects form a *Monoid*: *mempty* defines an aspect with no rule and *mappend* (for which we defined an infix equivalent operator $(\#)$) combines the rules of two aspects. The aspects of our example (*asp* in Figure 3) are composed by *evalA* and *envA*.

Merging aspects requires their domains to be disjoint, we provide primitives to delete some attribution rules from an aspect.

```
delete_I :: Typeable a
           ⇒ Attr I a → Child      → Aspect → Aspect
delete_S :: Typeable a
           ⇒ Attr S a → Production → Aspect → Aspect
```

2.4 Contexts

Attribute rules and aspects generate a set of constraints that must be satisfied by a grammar in order to be able to be evaluated. The context of an aspect contains information about which attributes it needs and which it provides. Each use of $(!)$, *par* and *ter* generates a *Require* constraint in order to make sure that the corresponding attribute will be defined when we run the attribute grammar. The aspect constructors *syn* and *inh* generate *Ensure* constraints for synthesized and inherited attributes, respectively.

We can inspect the constraints associated with an aspect using the function *print_context* in the interactive Haskell evaluator.

```
print_context :: Aspect → IO ()
```

For example, the context of *envA* denotes that the inherited attribute *env* is defined for *leftAdd* and *rightAdd*, and that the definition of the same inherited attribute is required for the non-terminal *Expr*.

```
$ print_context envA
Require Inherited   : Expr.env
Require Synthesized : -
Require Terminal    : -
Ensure Inherited   : leftAdd.env, rightAdd.env
Ensure Synthesized : -
```

Similarly, the context of *evalA* denotes that the synthesized attribute *eval* is defined for the productions *Add*, *Val* and *Var*, given that: both *eval* and *env* are defined for the non-terminal *Expr* and the terminal attributes *ival* and *vname* are defined for the productions *Val* and *Var*, respectively.

```
$ print_context evalA
Require Inherited   : Expr.env
Require Synthesized : Expr.eval
Require Terminal    : Val.ival, Var.vname
```

```
Ensure Inherited : -
Ensure Synthesized : Add.eval, Val.eval, Var.eval
```

If we merge *evalA* with the aspect *envA*, then we have the same requirements, adding the definition of the inherited attribute *env* for the children *leftAdd* and *rightAdd*.

```
$ print_context asp
Require Inherited : Expr.env
Require Synthesized : Expr.eval
Require Terminal : Val.ival, Var.vname
Ensure Inherited : leftAdd.env, rightAdd.env
Ensure Synthesized : Add.eval, Val.eval, Var.eval
```

Not all aspects are valid, for instance merging two aspects with overlapping domains yields an invalid aspect. In such cases the function *print_context* shows the errors of the aspect. For example, if we merge *evalA* with itself we get an error for each definition of *eval*.

```
$ print_context (evalA # evalA)
AG:Error: # at <location>
merging conflict: Add.eval
merging conflict: Val.eval
merging conflict: Var.eval
```

where the location (*<location>*) is the position in the source code where the aspects containing the duplicated attributes are merged.

Modular attribute grammars allow us to pick and choose any set of productions to form a grammar. An attribute grammar is complete if all the inherited and synthesized attributes *required* by the aspect are *ensured* by the same aspect, and all the terminal aspects are provided by the grammar. The function *print_missing* performs this checks, displaying the unmet requirements.

```
print_missing :: Grammar → Aspect → IO ()
```

For instance, if we check if the aspect *evalA* is complete for our grammar, we get that there is no rule to compute the way the inherited attribute *env*, used to compute *eval* at the production *Var*, is distributed.

```
$ print_missing (grammar [add,val,var]) evalA
Missing Inherited : leftAdd.env, rightAdd.env
Missing Synthesized : -
Missing Terminal : -
```

2.5 Running the Attribute Grammar

An attribute grammar denotes a function on a tree that takes the inherited attributes of the root and returns the synthesized attributes of the root.

```
type Check a = Either Error a
```

```
run :: Typeable t
  ⇒ GramDesc t → InhDesc i → SynDesc s → Aspect
  → Check (t → i → s)
```

In order to run the attribute grammar in a type-safe way, the function *run* returns, in case it is possible, a function of type $t \rightarrow i \rightarrow s$ given:

- A description *GramDesc t* of how the grammar can be extracted from a concrete tree *t*
- A mapping *InhDesc i* between the inherited attributes and a value of type *i* representing them

- A mapping *SynDesc s* between the synthesized attributes and a value of type *s* representing them
- The *Aspect* containing the rules

If any of the mappings is not correct or the aspect is not complete, then the function is not constructed and the errors are returned. It is important to notice that a given grammar (and aspect) can be associated to different concrete tree types, given the mappings are correct.

In Figure 4 we show how the semantic function of the example is extracted for a type of expressions *Expr*. To map the concrete tree type with the grammar, in *exprDesc* we have to define for each production how to obtain its children and terminal attributes from the constructors of the type. The mapping of the inherited attributes *InhDesc* is a monoid; the function *embed* creates a mapping of an attribute from a value of type *i*, given the mapping function. For instance, in Figure 4 we embed a list $[(String, Int)]$ as the attribute *env* using the function *fromList*.

```
embed :: Typeable a ⇒ Attr I a → (i → a) → InhDesc i
```

The mapping of the synthesized attributes *SynDesc* is an attribute functor. The function *project* projects the value of an attribute to create a value of type *s*. In the example we project the *env* attribute to a *Maybe Int*.

```
project :: Typeable a ⇒ Attr S a → SynDesc a
```

Thus we can apply the semantics of our example to a given expression and environment.

```
$ (fromRight runExpr) (Add (Add (Val 2) (Val 3))
                           (Var "x"))
  [("x",20)]
Just 25
```

However, if we try:

```
$ (fromLeft $ run exprDesc (embed env id)
                  (project eval)
                  evalA)
ERRORS
- missing inherited = leftAdd.env,rightAdd.env
$ (fromLeft $ run exprDesc emptyInDesc
                  (project eval)
                  asp)
ERRORS
- missing InhDesc = Expr.env
```

2.6 Extending The Grammar

An important design goal of our library is modular extensibility, both at the attribute and grammar level. Suppose we want to extend the expression language with a *let* construct.

```
<expr> ::= ...
| "let" <var> "=" <expr> "in" <expr>
```

In Figure 5 we show how the new production can be defined, declaring that it belongs to the non-terminal *expr*, it has two children (*exprLet* and *bodyLet*), and one terminal attribute (*vlet*) of type *String*.

Figure 6 shows the definition of the attributes for the extension.

```

data Expr = Val Int | Add Expr Expr | Var String
deriving (Show, Typeable)

exprDesc = gramDesc $ expr |= [ add |= [ leftAdd |= λcase {Add x _ → Just x; _ → Nothing}
, rightAdd |= λcase {Add _ x → Just x; _ → Nothing}
] & []
, val |= [] & [ ival |= λcase {Val x → Just x; _ → Nothing} ]
, var |= [] & [ vname |= λcase {Var x → Just x; _ → Nothing} ]
]

runExpr :: Check (Expr → [(String, Int)] → Maybe Int)
runExpr = run exprDesc (embed env Map.fromList) (Just ⇨ project eval) asp

```

Figure 4: Running the Attribute Grammar with a concrete tree *Expr*

```

elet    = production expr "Let" [exprLet, bodyLet]
        (vlet `consT` nilT)
exprLet = child elet "exprLet" expr
bodyLet = child elet "bodyLet" expr
vlet    = terminal "vlet" pString

```

Figure 5: Grammar Extension

```

evalA2 = syn eval [elet |- bodyLet ! eval] # evalA
envA2 = inh env [exprLet |- par env
, bodyLet |- [[Map.insert ⟨ter vlet⟩
⟨exprLet ! eval⟩
⟨par env⟩]]]
] # envA
asp2 = evalA2 # envA2

```

Figure 6: Attributes of the extension

3 DESIGN OF THE LIBRARY

In this section we introduce the design of the library by showing fragments of its implementation. The complete implementation of the library is available at [3].

3.1 Grammars

An important goal on the design of the library is that the grammars have to be modular. A *Grammar* can be given by a set of productions; this fully specifies a grammar and the representation is unique (up to set equality).

type Grammar = Set *Production*

Productions are identified by their name and non-terminal to which they belong. A *Production* includes a list of children and a set of terminal attributes.

type ProdName = (NonTerminal, Name)

```

data Production = Production {prod_name :: ProdName
, prod_children :: [Child]
, prod_terminals :: Set (Attribute T)}

```

deriving (Eq, Ord)

Non-terminals are identified by their name.

```

type NonTerminalName = Name
newtype NonTerminal = NT NonTerminalName
deriving (Eq, Ord)

A Child has a name, a production to which it belongs, and a non-terminal.

type ChildName = (Production, Name)
data Child = Child {child_name :: ChildName
, child_nt :: NonTerminal}
deriving (Eq, Ord)

```

Notice that the flexibility obtained by this representation comes with the cost of not statically assuring the correctness grammar. For instance, a production could contain a child that is defined to belong to another production. We will see later in this section that the checks made before generating the semantic function prevent incorrect grammars.

3.2 Pure Rules

Our encoding of attributes, rules and their evaluation, which we will present in this subsection, is based on the embedding for modular attribute grammars defined in [4]. We deviate from the original proposal in that we use *Strings* and *Dynamic* values, instead of closed universes, in order to obtain a greater extensibility.

The type *Attr k a*, to represent attributes of kind *k* and type *a*, is a phantom type [7], since the parameter *a* does not appear on the right-hand side of the definition.

```

data Attr k a = Attr {attr_name :: Name
, attr_kind :: Kind k}

```

Attributes are collected in finite maps. We hide the types of the attributes to use them as keys of the maps.

```

data Attribute k where
Attribute :: Typeable a ⇒ Attr k a → Attribute k

```

In order to be able to store the attribute computations in a Map (:->) from the attributes to their values, we hide their types using *Dynamic*.

```

type AttrMap k     = Attribute k :-> Dynamic
type ChildrenAttrs k = Child :-> AttrMap k

```

It is important to emphasize that the use of *Dynamic* and *Attribute k* are not exposed by the library. When using the library primitives the

user will always use concrete types and will refer to the attributes with type $\text{Attr } k \ a$.

A rule is a function that computes the attributes of a production; it is a mapping from *input attributes* to *output attributes*. These kinds of rules are called *pure rules*, in contrast to the monadic rules that we will define in the following subsections.

```
type PureRule = InAttrs → OutAttrs
```

The input attributes of a production are the inherited attributes of the parent, the synthesized attributes of the children and the terminal attributes of the production.

```
type InAttrs = (AttrMap I, ChildrenAttrs S, AttrMap T)
```

The output attributes are the synthesized attributes of the parent and the inherited attributes of the children.

```
type OutAttrs = (AttrMap S, ChildrenAttrs I)
```

As we said in the previous section, an attribute grammar denotes a function on a tree that takes the inherited attributes of the root and returns the synthesized attributes of the root. In our encoding this is represented by the following type:

```
type SemTree = AttrMap I → AttrMap S
```

The semantics of a production can be modeled as a function that takes the semantic functions of the children (subtrees) and the values of the terminal attributes, and returns the semantic function (*SemTree*) of the tree with this production as root.

```
type SemProd = Child :> SemTree → AttrMap T → SemTree
```

The function *sem_prod* ties the knot of attribute computation; it takes a rule and produces the semantics of a production.

```
sem_prod :: PureRule → SemProd
```

```
sem_prod rule semch terminals inh = syn
```

where

```
(syn, inh_children) = rule (inh, syn_children, terminals)
syn_children = applyMap semch extended_inh
extended_inh = Map.union inh_children
            (constantMap emptyAttrs
             (Map.keysSet semch))
```

Given the rule (*rule*), the semantics of the children (*semch*), the terminal attributes (*terminals*) and the inherited attributes of the parent (*inh*), :

- Compute the synthesized attributes (*syn*) and the inherited attributes of the children (*inh_children*) by applying the rule to the inherited attributes of the parent, the synthesized attributes of the children (*syn_children*) and the terminal attributes.
- Compute the synthesized attributes of the children (*syn_children*) by applying the semantics of the children to their inherited attributes (*inh_children*). Since the rule may not produce attribute for all the children of the production, we need to extend the domain of *inh_children* (with empty attributions) to cover all of them.

Notice that the definition is cyclic, since we use *syn_children* to produce *inh_children* and we use *inh_children* to produce *syn_children*.

A *pure aspect*, that groups together rules that define related attributes across multiple productions, is defined as a mapping from productions to pure rules.

```
type PureAspect = Production :> PureRule
```

Thus, the semantics of a pure aspect is defined by applying *sem_prod* to its rules.

```
sem_asp = Map.map sem_prod
```

3.3 Contexts

In order to check that the rules are compatible with a grammar, we must collect information about the rules, like which attributes are used and of which type. While building attribution rules, we build a system of inference rules that we call *Context*. The context is formed of a set of premises and a set of conclusions. The meaning of the context is that the conjunction of premises entails the conjunction of conclusions.

```
data Context = Ctx {ensure_I :: Set Ensure_I
                    , ensure_S :: Set Ensure_S
                    , require_I :: Set Require_I
                    , require_S :: Set Require_S
                    , require_T :: Set Require_T}
```

The premises are composed by the inherited (*require_I*), synthesized (*require_S*) and terminal (*require_T*) attributes that are used by the rules. The conclusions are the inherited (*ensure_I*) and synthesized (*ensure_S*) attributes that are defined by the rules. The premises are captured by the *Require* types and the conclusions by the *Ensure* types, which are all constraints that keep track of the used/defined attribute (hiding its type *a*) and the place of the grammar where it occurs.

```
data Constraint k t where
```

```
Constraint :: Typeable a ⇒ Attr k a → t → Constraint k t
```

```
type Ensure_I = Constraint I Child
```

```
type Ensure_S = Constraint S Production
```

```
type Require_I = Constraint I NonTerminal
```

```
type Require_S = Constraint S NonTerminal
```

```
type Require_T = Constraint T Production
```

Contexts form a monoid, where *mempty* is a *Context* with all its *Constraint* sets empty and *mappend* performs the union of such sets. We define a smart constructor *cstr* to construct singleton constraint sets.

```
cstr :: Typeable a ⇒ Attr k a → t → Set (Constraint k t)
```

```
cstr a x = Set.singleton (Constraint a x)
```

3.4 Checking the Contexts

The check if a *Context* is complete with respect to a *Grammar* is to check whether all *require* constraints are met by matching *ensure* constraints, and all terminals are defined in the grammar.

We define a type *Missing* to represent the unmet ensure constraints.

```
newtype Missing = Missing (Set Ensure_I
                           , Set Ensure_S
                           , Set Ensure_T)
```

For instance, given a set of children *chn* and a set of ensured inherited attributes *ensure*, the function *missing_I* returns the set of missing ensured inherited attributes to meet a given required

inherited attribute r . That is, an ensure constraint for each child with the same non-terminal of the required attribute that does not define it.

```
missing_I :: Set Child → Set Ensure_I → Require_I → Set Ensure_I
missing_I chn ensure r@(Constraint a n) =
  Set.difference match_child match_ensure
  where
    match_chn = Set.map (λc → Constraint a c) $
      Set.filter ((≡ n) o child_nt) chn
    match_ensure = Set.filter ((≡ n) o child_nt o constr_obj)
      ensure
```

The missing ensure constraints are computed as the difference between the following sets:

- *match_chn*, with an ensure constraint for the attribute a for each child whose non-terminal matches with the non-terminal (n) of the require constraint.
- *match_ensure*, the subset of the ensure constraints for children with the same non-terminal n of the require constraint r .

Thus, in the function *missing*, that returns all the unmet ensure constraints, the set of not ensured inherited attributes is computed by applying the function *missing_I* for all the children of the grammar and the ensured inherited attributes of the context to every required inherited attribute of the context. Since every ensure constraint of the context implies a require constraint, the function *require_all_I* returns the set of required inherited attributes (*Set Require_I*) of the context, including all its *Require_I* and *Ensure_I* constraints with the last ones are transformed to *Require_I*.

```
missing :: Grammar → Context → Missing
missing g c = Missing
(foldMap (missing_I (gram_children g) (ensure_I c))
  (require_all_I c)
, foldMap (missing_S g (ensure_S c)) (require_all_S c)
, missing_T g (require_T c))
```

The set of not ensured synthesized attributes is computed by the function *missing_S*, which is similar to *missing_I*. This case of terminal attributes is different, since they are not associated with non-terminals but with productions.

```
missing_T :: Grammar → Set Require_T → Set Ensure_T
missing_T g required =
  Set.difference (Set.filter is_prod required) (ensure_T g)
  where
    is_prod r = constr_obj r `Set.member` g
```

We also remove the constraints that mention productions that are not in the grammar.

3.5 Aspects

To construct a *Context* in parallel with the rule, instead of directly defining pure rules, rules are defined in an applicative *AR*, that comes with primitives to project terminal attributes or attributes from either the parent or a child of the production. Those primitives generate *require* and *ensure* constraints, or throw errors if necessary.

In order to compute rules, we must first check that they are valid. Rules are defined in the context of a production, and may fail if some constraints are not met; e.g. using a child that is not a valid child of the current production. And lastly, we collect constraints. We define the *Aspect Monad*: a *Reader* monad that reads from a *Production*, an *Error* monad that reports errors and a *Writer* monad that constructs a *Context*.

```
newtype A a = A (ReaderT Production
  (ExceptT Error (Writer Context)))
a)
deriving (Functor, Applicative, Monad
  , MonadReader Production
  , MonadError Error
  , MonadWriter Context)
```

Given a production, the monadic computation can be *run* to get the computed value (or an *Error*⁴) and the context constructed so far out of the monad.

```
runA :: A a → Production → (Check a, Context)
runA (A a) = runWriter ∘ runExceptT ∘ runReaderT a
```

Within this monad we can, for example, define a function to add a constraint to the context requiring that a given terminal attribute has to be provided at the current production.

```
require_terminal :: Attr T a → A ()
require_terminal a = do p ← ask
  tell $ mempty { require_T = cstr a p }
```

We can also generate errors if a given child is not valid in the current production.

```
assert_child :: HasCallStack ⇒ Child → A ()
assert_child c = do
  p ← ask
  unless (c ∈ prod_children p)
    throwErrorA (Error_Rule_Invalid_Child c p))
```

If the child c does not belong to the children of the production p , the function *throwErrorA* throws the given error, including the most recent call-site off the call stack. The *CallStack* is used to point to the location in the source code where the error occurs. Then we can define a function to add a constraint requiring an inherited attribute in a given child.

```
require_child :: HasCallStack ⇒ Child → Attr S a → A ()
require_child c a = do
  assert_child c
  tell $ mempty { require_S = cstr a (child_nt c) }
```

Similarly, we can add an ensure constraint.

```
ensure_child :: Child → Attr I a → A ()
ensure_child c a = do
  assert_child c
  tell $ mempty { ensure_I = cstr a c }
```

In order to gather information from the use of input attributes, we define the rules in a specific monad in which the input attributes are accessed through primitives. The *Rule Monad* is a *Reader* monad reading from the input attributes.

⁴The datatype *Error* includes a constructor for each possible error to report.

```
newtype R a = R { runR :: Reader InAttrs a }
deriving (Functor, Applicative, Monad, MonadReader InAttrs)
```

By combining an *Aspect Monad* and a *Rule Monad* we define the *Attribution Rule* applicative functor.

```
newtype AR a = AR { runAR :: A (R a) }
```

```
instance Functor AR where
```

```
fmap f x = pure f <>> x
```

```
instance Applicative AR where
```

```
pure x = AR (pure (pure x))
AR f <>> AR x = AR ((<>>) <>> f <>> x)
```

The *AR* primitives we introduced in the previous section are implemented using the *Aspect Monad* to collect the constraints and the *Rule Monad* to access to the input attributes. For instance, the function *ter* to project a terminal attribute imposes a *require_terminal* of the attribute and looks it up at the terminal attributes mapping.

```
ter :: HasCallStack => Attr T a -> AR a
ter a = AR $ do
  require_terminal a
  return $ fromJust `o` lookupAttr a <>> asks (λ(_,_,_)> t) -> t)
```

When searching the attribute we apply *fromJust* to the result, assuming that it will be found. Also in the function *lookupAttr* we do an unsafe conversion. These are not problems, because the AG will be checked before we can evaluate it.

```
lookupAttr :: Attr k a -> AttrMap k -> Maybe a
lookupAttr a@Attr {} m =
  fromJust `o` fromDynamic <>> Map.lookup (Attribute a) m
```

Another example of *AR* primitive is the operator (!), to project a synthesized attribute from a child. In this case we require that the child must belong to the production and that it has to include the attribute.

```
(!) :: HasCallStack => Child -> Attr S a -> AR a
c ! a = AR $ do
  require_child c a
  return $ fromJust `o` lookupAttr a `o` fromJust `o` Map.lookup c <>>
    asks (λ(_,_ch,_)> ch)
```

A *Rule* is an attribution rule *AR* that computes output attributes.

```
type Rule = AR OutAttrs
```

Empty rules are rules that produce no output; i.e. an empty mapping of synthesized attributes and an empty mapping of inherited attributes of the children. An empty rule has no effect on the context.

```
emptyRule :: Rule
emptyRule = pure (empty, empty)
```

The rules are merged by the union of their attribute mappings.

Merging rules whose domain overlap is an error.

```
mergeRule :: HasCallStack => Rule -> Rule -> Rule
```

```
mergeRule left_rule right_rule
```

```
| ¬ (Set.null duplicate_inh) = AR $ throwErrorA $
  Error_Rule_Merge_Duplicate_I duplicate_inh
| ¬ (Set.null duplicate_syn) = AR $ throwErrorA $
  Error_Rule_Merge_Duplicate_S duplicate_syn
| otherwise = liftA2 mergeOutAttrs left_rule right_rule
```

```
where duplicate_inh = ...
      duplicate_syn = ...
```

```
mergeOutAttrs (s1, ic1) (s2, ic2) =
  (Map.union s1 s2, Map.unionWith Map.union ic1 ic2)
```

An *Aspect* maps productions to rules.

```
newtype Aspect = Aspect (Production :> Rule)
```

Aspects form a monoid, where the binary operator combines the rules of the two aspects by merging them. This is the (#) operator of the library that we introduced in the previous section.

```
instance Monoid Aspect where
```

```
mempty = Aspect empty
mappend (Aspect a1) (Aspect a2) = withFrozenCallStack $
  Aspect $ Map.unionWith mergeRule a1 a2
```

```
(#) = mappend
```

The primitives to define an aspect for a single attribute, given an *AR* with the attribute computation, add to the context the respective ensure constraints. For example, the function *syn1*, to define a synthesized attribute *attr* for a production *prod* with the computation *rule*, returns an *Aspect* that contains a singleton mapping from the production *prod* to an *AR* that adds an *Ensure_S* of *attr* to the context produced by *rule*. Note that we set the production of the rule using *local*. The rule results in a singleton mapping from the attribute *attr* to the result of its computation, which is stored as a *Dynamic* object.

```
syn1 :: Typeable a => Attr S a -> Production -> AR a -> Aspect
syn1 atr prod rule =
  = Aspect $ Map.singleton prod $ AR $ do
    rule' ← local (const prod) (ensure_parent atr >> runAR rule)
    return $ do r ← rule'
              return (Map.singleton (Attribute attr) (toDyn r)
                      , empty)
```

Thus, the primitive *syn* of Section 2.3 merges the aspects resulting from applying *syn1 a* to all the elements of the list.

```
syn :: Typeable a => Attr S a -> [(Production, AR a)] -> Aspect
syn a = foldl (λrs (p, r) -> rs # syn1 a p r) emptyAspect
```

3.6 Checking the Rules

When the rule set is deemed complete, we can check its context with respect to a grammar. The function *checkAspect* takes a grammar and an aspect, returning an *Error* monad that computes the pure aspect and the produced context if no errors are found.

```
checkAspect :: Grammar -> Aspect -> Check (PureAspect, Context)
checkAspect gram rule = do
  check_grammar gram
  let (check_aspect, ctx) = runAspect rule
  pure_asp ← check_aspect
  check_missing (missing gram ctx)
  return (pure_asp, ctx)
```

First, we use *check_grammar* to check that the children have unique names for each production. Then we obtain the pure rule (if it has no errors) and context with *runAspect*. With the function *missing*, from Section 3.4, we check that the context is complete with respect to

the grammar. If some constraints are unmet, *check_missing* throws an error.

The function *runAspect* is implemented in the following way:

```
runAspect :: Aspect → (Check PureAspect, Context)
runAspect (Aspect asp) = (asp_err, ctx)
  where
    asp_ar = traverse runAR asp
    asp_a = liftM (Map.map (runReader ∘ runR)) asp_ar
    (asp_ag, ctx) = runA asp_a ⊥
    errors_ag = fst $ runA (collect_errors $ runAR <\$ asp) ⊥
    asp_err = do errors ← errors_ag
                 if null errors then asp_ag
                 else throwErrorCheck $ Errors errors
```

To compute the pair *(asp_err, ctx)*, we compute:

- *asp_ar*, an Aspect Monad of type *A* (*Production*:→*R OutAttrs*) that results from applying *runAR* to every element of the aspect and evaluate the actions of the Aspect Monads from left to right.
- *asp_a*, an Aspect Monad of type *A PureAspect* that results from running the Rule Monad of the mapping.
- *(asp_ag, ctx)*, a pair with the pure aspect and the context computed by running *asp_a*. Notice that we pass \perp as the production in the Rule Monad. This is possible because the production is not used for rules, since when we build rules we always override the production with a call to *local* (see e.g. the code of *syn1* in Section 3.5).
- *errors_ag*, a list of (possible) errors collected from each production.
- *asp_err*, an error monad that returns the pure aspect *asp_ag* if the list of errors is empty, or throws the errors otherwise.

3.7 Running the Attribute Grammar

In order to run the AG in a type-safe way, we must check that a concrete type is compatible with the context free grammar. We do this verification at runtime, but once and for all. If everything is ok we obtain a safe semantic function. Recall the type of the function *run* of Section 2.5.

```
run :: Typeable t ⇒ GramDesc t → InhDesc i → SynDesc s
      → Aspect → Check (t → i → s)
```

The type *GramDesc t* associates a grammar to a family of concrete types, where *t* is associated to the start symbol of the grammar; i.e. the root of the tree will have type *t*.

```
newtype GramDesc t = GramDesc { runGramDesc :: 
  Check (NonTerminal, Grammar, Child :> TypeRep, GramMap) }
```

A *GramDesc* is an *Error* monad that computes a quadruple containing: the root non-terminal of the grammar, the grammar, a mapping from the children of the grammar to the representation of their corresponding concrete type⁵, and a *GramMap* mapping each non-terminal with the information of its associated type.

```
type GramMap = NonTerminal :> (TypeRep, Dynamic → Match)
```

⁵This is used to check that type of the non-terminal associated with each child corresponds to the type of the child.

For each non-terminal we have to provide the representation of its associated type, and a function that *destructs* a *Dynamic* that wraps a value of this type into a *Match*.

```
data Match = Match { match_prod :: Production
                     , match_child :: Child :> Dynamic
                     , match_terminals :: AttrMap T }
```

The *Match* information includes the associated production, a map from the children of the production to their corresponding values (into *Dynamic*), and an *AttrMap* containing the non-terminals.

The library provides a function *gramDesc* and a couple of combinators (|= and &) to construct a value of type *GramDesc* as we showed in Section 2.5.

In order to hide the *AttrMap* mappings from the user, we have to specify an association between the inherited and synthesized attributes of the root and two given types *i* and *s*, respectively. With *SynDesc* and *InhDesc* we provide an interface to define this association and build conversion functions (*i* → *AttrMap I*) and (*AttrMap S* → *s*). Both are *Writer* monads that accumulate the mentioned attributes, in order to be able to check that they are correct with respect to the rule.

In the case of the synthesized attributes, *SynDesc* accumulates a set of synthesized attributes and returns a conversion function (*AttrMap S* → *a*).

```
newtype SynDesc a = SynDesc { runSynDesc :: 
  Writer (Set (Attribute S)) (AttrMap S → a) }
```

When we *project* a synthesized attribute, the conversion function is to look it up into the *AttrMap*. Notice that again we use *fromJust* because we know that its existence will be checked.

```
project :: Attr S a → SynDesc a
project a = SynDesc $ do
  tell (Set.singleton (Attribute a))
  return $ fromJust ∘ lookupAttr a
```

In the case of the inherited attributes, *InhDesc* accumulates a list of inherited attributes and returns a conversion function (*a* → *AttrMap I*).

```
newtype InhDesc t = InhDesc { runInhDesc :: 
  Writer ([Attribute I]) (t → AttrMap I) }
```

To *embed* an inherited attribute is to apply a function that *extracts* the value of the attribute from *i* and insert it as a *Dynamic* into the *AttrMap*.

```
embed :: Typeable a ⇒ Attr I a → (i → a) → InhDesc i
embed a p = InhDesc $ do
  tell [Attribute a]
  return $ Map.singleton (Attribute a) ∘ toDyn ∘ p
```

To get the conversion functions we have to run the monads and return their result.

```
proj_S :: SynDesc a → AttrMap S → a
proj_S = fst ∘ runWriter ∘ runSynDesc
proj_I :: InhDesc i → i → AttrMap I
proj_I = fst ∘ runWriter ∘ runInhDesc
```

The function *check* returns a computation that checks the whole attribute grammar.

```

check :: GramDesc t → InhDesc i → SynDesc s → Aspect
  → Check (NonTerminal, PureAspect, GramMap)
check g i s r = do
  (root, grammar, gmap) ← check_gramDesc
  (pure_asp, ctx)      ← checkAspect grammar r
  check_inh_unique i
  check_inh_required i root (require_I ctx)
  check_syn_ensured s grammar root (ensure_S ctx)
  return (root, pure_asp, gmap)

```

We perform the following checks:

- *check_gramDesc*, the concrete types associated to every non-terminal correspond to the types of the children of these non-terminals.
- *checkAspect*, the aspect is complete with respect to the grammar.
- *check_inh_unique*, there are no duplicated inherited attributes specified for the root.
- *check_inh_required*, all the required inherited attributes have been specified for the root.
- *check_syn_ensured*, all the synthesized attributes accessed from the root are ensured by the rules.

The function *run* performs the checks and returns the semantic function.

```

run :: Typeable t ⇒ GramDesc t → InhDesc i → SynDesc s
  → Aspect → Check (t → i → s)
run g i s a = do
  (root, pure_asp, gmap) ← check g i s a
  let coalg = Map.map snd gmap
  let sem = sem_coalg (sem_asp pure_asp) coalg root ∘ toDynG g
  return $ λx → proj_S s ∘ sem x ∘ proj_I i
where
  toDynG :: Typeable t ⇒ GramDesc t → t → Dynamic
  toDynG _ = toDyn

```

Once the checks are passed, we proceed to extract the coalgebra of the *Dynamic* tree from the GramMap of the description of the grammar. The coalgebra has the following type.

```
type Coalg = NonTerminal :> (Dynamic → Match)
```

Then we construct the semantic function using the pure aspect and the coalgebra. To apply the semantic function we convert the tree to *Dynamic*. We also apply the conversion functions to and from *AttrMap*.

The function *sem_coalg* iterates the AG-algebra on the tree-coalgebra.

```

sem_coalg :: Production :> SemProd → Coalg
  → NonTerminal → Dynamic → SemTree
sem_coalg alg coalg = sem
where
  sem nt dyn = (alg ! prod) children terms
  where
    Match prod cmap terms = (coalg ! nt) dyn
    children = Map.mapWithKey (λk d → sem (child_nt k) d)
    cmap

```

4 GENERIC RULES

A nice advantage of having first class attribute grammars hosted in a functional language like Haskell, is that we can use it to capture common patterns. Some examples of common patterns in attribute grammars are the copy, collect and chain rules.

The copy rule simply passes down an inherited attribute from the parent to a children. Using the public primitives of the library we can implement a function *copy* that passes down an inherited attribute to a given child.

```

copy :: Typeable a ⇒ Attr I a → Child → Aspect
copy a c = inh a c (par a)

```

We can also apply *copy* to a list of children for which the attribute is to be copied.

```

copyN :: Attr I a → [Child] → Aspect
copyN = many1 copy

```

Or we can copy the inherited attribute of the parent to all the children that have the same non-terminal.

```

copyP :: Attr I a → Production → Aspect
copyP a p = copyN a cs
  where cs = [c | c ← prod_children p, child_nt c ≡ prod_nt p]

```

Or apply the copy rule for a list of productions.

```

copyPs :: Attr I a → [Production] → Aspect
copyPs a = foldr (λp r → copyP a p # r) emptyAspect

```

The collect rule applies a function to collect the values of the synthesized attributes of the children. For instance, we can define a function *collect* that applies a function to the attributes of a list of children to compute a synthesized attribute.

```

collect :: Attr S a → ([a] → a) → Production → Children → Aspect
collect a reduce p cs = syn a p $ reduce <⇒> traverse (!a) cs

```

The chain rule takes a pair of an inherited and a synthesized attribute and thread them through the children of a production: the parent attribute is given to the first child, the attribute of the first child is given to the second and so on, the attribute of the last child is given back to the parent. So, this rule defines the inherited attribute for all the children and the synthesized attribute for the parent.

```

chain :: Typeable a ⇒
  Attr I a → Attr S a → Production → Children → Aspect
chain i s p cs = (inhs i $ zip cs $ par i : ((!s) <⇒> init cs))
  # syn s p (last cs ! s)

```

5 RELATED WORK

Starting with Johnson in [9], there have been several works on encoding attribute grammars in pure functional languages with lazy evaluation.

As we said in the introduction, this work builds on the compositional semantics given by [4] to the structuring mechanisms of attribute grammars. However, our embedding is both more extensible and safer.

Other approaches for attribute grammars in Haskell are based on the use of extensible records [5] or type-level programming techniques [2, 15] which use many type system extensions of GHC.

In this way it is possible to perform correctness checks at compile time, but with the cost of poor quality error messages.

Some other embeddings use a *Zipper* [8] to model attribute grammars [1, 6, 11, 12, 14], but all of them work over grammars represented by concrete datatypes, not being able to solve the expression problem[16].

6 CONCLUSIONS

We have introduced an embedding of attribute grammars as a domain specific language in Haskell. Embedding attribute grammars in a host language has many advantages over the alternative of using a specific tool. In particular the host language and its libraries are readily available to the grammar writer for use when expressing the computation rules of attributes. Another benefit of the embedding is that it becomes possible to run an attribute grammar within a program, and hence compute the semantics of a syntax tree. This extends the host language with the attribute grammar paradigm for writing recursive tree computations.

Our embedding of modular attribute grammars does not compromise on safety, on the simplicity of the public interface and the clarity of error messages. The approach is to have two phases of typechecking: (1) the static phase uses the host typesystem to capture as many invariants as possible, (2) the runtime phase implements a *global* verification of the domain specific program (the attribute grammar) before running it. This global verification captures the remaining constraints that were not caught at compile time, and allows us to precisely identify the source of the problem. In addition, after the runtime global checks, the embedded program may be run without dynamic checks, thus speeding up execution.

We believe that our approach to domain specific errors may be successfully applied to define other domain specific language embeddings, being this a possible direction of future research.

It would be also interesting to measure the impact in performance of moving the checks to runtime, and compare our library with other embeddings.

To continue improving the library is another line of future research. For example, to investigate which checks can be moved back to compile time without compromising the quality of error messages.

REFERENCES

- [1] E. Badouel, R. Tchougong, C. Nkuimi-Jugnia, and B. Fotsing. 2013. Attribute grammars as tree transducers over cyclic representations of infinite trees and their descriptive composition. *Theoretical Computer Science* 480, 0 (2013), 1 – 25. <https://doi.org/10.1016/j.tcs.2013.02.007>
- [2] Florent Balestrieri. 2014. *The Productivity of Polymorphic Stream Equations and the Composition of Circular Traversal*. Ph.D. Dissertation. University of Nottingham.
- [3] Florent Balestrieri, Alberto Pardo, and Marcos Viera. 2018. Repository of the SafeAG library. <https://github.com/balez/Safe-AG>. (2018). Accessed: 2018-06-08.
- [4] Oege de Moor, Kevin Backhouse, and S. Doaitse Swierstra. 2000. First Class Attribute Grammars. *Informatica: An International Journal of Computing and Informatics* 24, 2 (June 2000), 329–341. Special Issue: Attribute grammars and Their Applications.
- [5] Oege de Moor, L. Peyton Jones, Simon, and Van Wyk, Eric. 2000. Aspect-Oriented Compilers. In *Proceedings of the 1st Int. Symposium on Generative and Component-Based Software Engineering*. Springer-Verlag, London, UK, 121–133.
- [6] João Paulo Fernandes, Pedro Martins, Alberto Pardo, João Saraiva, and Marcos Viera. 2016. Memoized Zipper-Based Attribute Grammars. In *Programming Languages - 20th Brazilian Symposium, SBLP 2016, Maringá, Brazil, September 22-23, 2016, Proceedings (Lecture Notes in Computer Science)*, Fernando Castor and Yu David Liu (Eds.), Vol. 9889. Springer, 46–61. https://doi.org/10.1007/978-3-319-45279-1_4
- [7] Ralf Hinze. 2003. Fun with phantom types. In *The Fun of Programming*, Jeremy Gibbons and Oege de Moor (Eds.). Palgrave Macmillan, 245–262.
- [8] Gérard Huet. 1997. The Zipper. *Journal of Functional Programming* 7, 5 (1997), 549–554.
- [9] Thomas Johnsson. 1987. Attribute grammars as a functional programming paradigm. In *Proceedings of the Functional Programming Languages and Computer Architecture*. Springer-Verlag, London, UK, 154–173.
- [10] Donald Knuth. 1968. Semantics of Context-free Languages. *Mathematical Systems Theory* 2, 2 (June 1968). *Correction: Mathematical Systems Theory* 5 (1), March 1971.
- [11] Pedro Martins, João Paulo Fernandes, and João Saraiva. 2013. Zipper-Based Attribute Grammars and Their Extensions. In *Programming Languages - 17th Brazilian Symposium, SBLP 2013, Brasília, Brazil, October 3 - 4, 2013, Proceedings (Lecture Notes in Computer Science)*, André Rauber Du Bois and Phil Trinder (Eds.), Vol. 8129. Springer, 135–149. https://doi.org/10.1007/978-3-642-40922-6_10
- [12] Pedro Martins, João Paulo Fernandes, João Saraiva, Eric Van Wyk, and Anthony Sloane. 2016. Embedding attribute grammars and their extensions using functional zippers. *Sci. Comput. Program.* 132 (2016), 2–28. <https://doi.org/10.1016/j.scico.2016.03.005>
- [13] Conor McBride and Ross Paterson. 2008. Applicative programming with effects. *Journal of Functional Programming* 18, 1 (2008), 1–13.
- [14] Tarmo Uustalu and Varmo Vene. 2005. Comonadic functional attribute evaluation (*Trends in Functional Programming*). Intellect Books (10), 145–162.
- [15] Marcos Viera, S. Doaitse Swierstra, and Wouter Swierstra. 2009. Attribute Grammars Fly First-Class: How to do aspect oriented programming in Haskell. In *Proceedings of the 14th Int. Conf. on Functional Programming*. ACM, New York, USA, 245–256.
- [16] Phil Wadler. 1998. The Expression Problem. (1998). <http://www.daimi.au.dk/~madst/tool/papers/expression.txt> E-mail available online.

Extended Memory Reuse – an optimisation for reducing memory allocations and deallocations in reference counted code

DRAFT

Hans-Nikolai Vießmann
Heriot-Watt University
Edinburgh, UK
hv15@hw.ac.uk

Artjoms Šinkarovs
Heriot-Watt University
Edinburgh, UK
a.sinkarovs@hw.ac.uk

Sven-Bodo Scholz
Heriot-Watt University
Edinburgh, UK
S.Scholz@hw.ac.uk

ABSTRACT

In this paper we present an optimisation for improving code generation for reference counted memory management. While preserving the key benefits of reference counting, i.e., the opportunities for *in-place* updates as well as the deallocation of memory without global garbage collection, the proposed optimisation massively reduces the total number of memory allocations and deallocations made. The key idea is to extend the lifetime of variables to enable reuse of the allocated memory without deallocation and subsequent reallocation. This avoidance of frequent memory deallocations and allocations turns out particularly useful in the context of innermost loops as it effectively leads to a generated code that performs pointer swaps between buffers as it would be implemented manually in codes with explicitly managed memory.

We have implemented the proposed optimisation in the context of the SaC compiler tool chain. The paper provides an algorithmic description of our optimisation and an evaluation of its effectiveness over a collection of benchmarks including a subset of the Livermore Loop benchmarks and the NAS Parallel Benchmarks. We can show that for a large number of the benchmarks with allocations within loops our optimisation reduces the amount of allocations significantly. We also observe that it does neither have a negative impact on the overall memory footprint nor on the overall runtime. Instead, we see some runtime improvements. While sequential and multi-threaded executions typically benefit only mildly, on GPU devices we observe speedups of up to a factor of 4.

KEYWORDS

memory reuse, compiler generated code, compiler optimisation

ACM Reference Format:

Hans-Nikolai Vießmann, Artjoms Šinkarovs, and Sven-Bodo Scholz. 2019. Extended Memory Reuse – an optimisation for reducing memory allocations and deallocations in reference counted code: DRAFT. In *Proceedings of International Symposium on Implementation and Application of Functional Languages (IFL'18)*. ACM, New York, NY, USA, Article 4, 9 pages. https://doi.org/10.475/123_4

1 INTRODUCTION

Reference-counting memory is one possible solution to having implicit memory management within a programming language [4].

IFL'18, August 2019, Lowell, MA, USA

2019. This is the author's version of the work. It is posted here for your personal use. Not for redistribution. The definitive Version of Record was published in *Proceedings of International Symposium on Implementation and Application of Functional Languages (IFL'18)*, https://doi.org/10.475/123_4.

A counter is associated to each data object in the language and, incremented or decremented based upon the number of times the data object is *referenced* within a function scope. When the counter is zero, we know that there are no further reference of the data object and so its associated memory can be freed.

Several languages make use of reference-counting for memory management, some use these internally such as SISAL [2], PYTHON [5], SWIFT [12], and Single-Assignment C (SAC) [7]. While others provide facilities to use reference-counting for user-defined objects like in RUST [13].

There are two advantages to using reference-counting, the first is that we minimise the memory footprint of the application by freeing all data objects immediately the moment their count is zero¹. Secondly, reference-counting offers an elegant solution to *in-place* updating of data objects, known as the aggregate update problem [11]. At the point of modifying a data object, we know that this is the last reference of the object and so its memory can be reused for the new data object rather than allocating new memory.

In some instances though, the *eager*-freeing strategy of reference-counting can cause impede opportunities for *memory reuse*. For instance, we may free a data object just before a new data object with the same memory requirement is to be allocated. Ideally we would want to reuse the existing allocated memory for our new data object and avoid the overheads of freeing and allocating. If we consider loops, these overheads can increase proportionately with the number of iterations.

To some extent these overheads can be kept low by means of heap management techniques that cater for such cases of reallocation. Freed memory typically is not immediately released from the process and the allocation process can favour finding recently freed memory for new allocations. Some reference counting specific techniques for heap management can be found in [6].

While such techniques may be considered sufficient in the context of classical shared memory architectures, on systems that employ accelerators such as GPUs eager-freeing can have more significant performance implications. Typically such systems have separate memories for the host and the accelerator and all memory management is performed by the host [15]. Transfers of data between these memories as well as operations that require involvement of the system cpu and the accelerator such as it occurs during management of accelerator memory induce a major overhead. Consequently, advanced optimisations aim to avoid data transfers, particularly within innermost loops if they contain kernel invocations [10]. In those cases where freeing and allocation operations

¹We assume here that the order of execution is fixed and cannot be changed.

on the accelerator remain within the loops substantial performance degradations can be observed.

In this paper we propose a solution to this eager-freeing problem, in the form of a compiler optimisation that we call *extended memory reuse (EMR)*. This optimisation is built on top of and extends the reference-counting related optimisations used by the SAC compiler. In particular we introduce an inference system that finds opportunities for memory reuse, and that also can handle loops.

We provide our optimisation as a series of code-transformations given as algorithmic specifications. We have implemented the *EMR* optimisation within the SAC compiler sac2c. We use this to evaluate the performance of the optimisation against a large collection of benchmarks, including some from the Livermore Loops Suite [14], the Rodinia Benchmark Suite [3], the NAS Parallel Benchmark Suite [1], and others.

The SAC compiler has support for generating code for both SMA systems and GPGPU CUDA-compatible devices. We evaluate the performance of these code-transformation by running these on a collection of benchmarks, on both a SMA system and a GPGPU. Our results show that in some cases we significantly reduce the number of memory allocation calls *without* affecting the runtime. Additionally we show that for some benchmarks, we have significant speedup for our GPGPU generated code.

2 BACKGROUND

In order to better understand reference-counting, we will look at an example written in the Single-Assignment C (SAC) programming language. SAC is a functional array programming language, which treats arrays as first-class citizens and provides a data-parallel loop construct, the *with-loop*, which is similar to a *map*-operation but additionally supports index ranges.

The *with-loop* is the main mechanism in SAC to *generate* arrays, through the **genarray** operation. It additionally supports array modification (**modarray**) and reduction operations (**fold**). For the rest of the paper we will focus only on **genarray** *with-loops*. The following example demonstrates how we can allocate a 2-dimensional array:

```
a = with {
  (. <= iv <= .) : 42;
} : genarray([2,10], 0);
```

With the *with-loop* we specify an index range which is within the shape of the array (which is [2,10]) and set each element of the array to 42. We use the short-notation (.) which is equivalent to specifying [0,0] <= iv <= [1,9]. If we chose an index range smaller than the array shape, then all these indices would be set to the default value 0.

SAC uses reference-counting internally, meaning that all arrays are annotated with a counter. At compile time we introduce into the code primitives which either increment (**incrc**) or decrement (**decrc**) this counter at runtime. In some instances it is not possible to statically decide on certain actions as it is not clear what the reference count is for a particular array. In this case we generate code which at runtime inspects (**getrc**) the reference count and decides what action to take. We use **alloc** to allocate memory and specify the initial reference count for the array; in this paper we assume that it is always 1, though it can be any positive number value.

```
1 int foo (int[11] a) {
2   b = with {
3     ([1] <= iv <= [5]) : a[iv];
4   } : genarray ([11], 0);
5   s = bar (b);
6   c = with {} : genarray ([11], 2);
7   print (c);
8   for (i = s; i < 10; i++) {
9     d = with {} : genarray ([11], i);
10    print (d);
11  }
12  return s;
13 }
```

Listing 1: SAC code example

A significant advantage of using reference-counting is that array modifications can be done in-place, thereby reusing memory [4]. Further memory reuse optimisation are possible, by for instance introducing parallel updates for both local and non-local reads using static inference [8] and the polyhedral model [9], respectively.

In listing 1 we show an example of a SAC function *foo*. It has a single argument *a* (1-dimensional array) and returns a scalar value. Within the function we generate a new 1-dimensional array *b* using a *with-loop* with a shape of [11]. We make use of a function *bar*, which is not shown here, but takes an array and returns a scalar. We pass *b* to the function which assigns to *s*. We generate a new array *c* with a shape of [11] and print it directly afterwards. Notice that we don't specify an index range, which means we iterate over the shape and assign each element the default value. At the *for-loop* we assign *s* to the iteration counter *i*. The loop iterates till the counter reaches 10, on each iteration a new array *d* is generated using the *i* as the default value for the *with-loop*. We print this array on each iteration, we leave the loop, and finally return our scalar value *s*.

At some point within the SAC compiler, we introduce memory primitives into the code. Using our example from listing 1, we now introduce the memory primitives and show this in listing 2.

The addition of the memory primitives doesn't significantly change the code. The main change is that we introduce a conditional that checks the reference count of *a*. This is done because optimisations within the SAC compiler have identified *a* as a *reuse candidate* for *b*. In general an array is considered a reuse candidate iff it is referenced inside the body of a *with-loop*. The optimisations though were not able to statically reuse the memory as it isn't clear if *a* is used beyond the application of *foo*. Thus at runtime, if *a* is referenced later, we need to allocate *b*, otherwise we can reuse the memory for *b*.

Additionally all of the *with-loops* have been transformed into a series of *for-loops*.

After the conditional, we decrement the reference count of *a*, which may or may not cause *a* to be freed. When we call *bar*, the array *b* is passed, and its reference count is decremented inside the function. When we allocated *b*, we implicitly set its reference count to 1; as no further incrementation of the counter occurred prior to being passed to *bar*, *b* will be freed inside *bar*. We allocate *c* and initialise it, and print it. Like with *b*, it is freed. We now reach our *for-loop* where we allocate and initialise our array *d* before printing

```

1 int foo (int[11] a) {
2     if (getrc (a) == 1) {
3         b = a;
4         for (iv = 0; iv < 1; iv++) b[iv] = 1;
5         for (iv = 6; iv < 11; iv++) b[iv] = 1;
6     } else {
7         b = alloc (...);
8         for (iv = 0; iv < 1; iv++) b[iv] = 1;
9         for (iv = 1; iv < 6; iv++) b[iv] = a[iv];
10        for (iv = 6; iv < 11; iv++) b[iv] = 1;
11    }
12    decrc (a);
13    s = bar (b); /* RC of b decremented inside */
14    c = alloc (...);
15    for (iv = 0; iv < 11; iv++) c[iv] = 2;
16    print (c); /* RC of c decremented inside */
17    for (i = s; i < 10; i++) {
18        d = alloc (...);
19        for (iv = 0; iv < 11; iv++) d[iv] = i;
20        print (d); /* RC of d decremented inside */
21    }
22    return s;
23 }
```

Listing 2: SAC code example with memory primitives

it. Similar to b, once d is passed to `print`, it is freed; this happens on each iteration.

3 MOTIVATION FOR EXTENDED MEMORY REUSE

Our example demonstrate how reference counting works in SAC and also showcases the two challenges that we intend to solve with our *EMR* optimisation.

The first problem is that we allocate two arrays b and c which have the same shape, and free both immediately. We would prefer it if we could reuse the memory of b for c, thereby reducing the number of allocations.

The array b is an example of what we call an *extended reuse candidate (ERC)*. It differs from a normal reuse candidate in that it is *not* referenced inside the body of a *with-loop*. Its only other criteria is that it *must* have the same shape as the *with-loop*.

The second problem is that the *for-loop* repeatedly allocates and frees the array d. In order to solve this we want to *lift* out the initial allocation of d and reuse that memory within the loop; there would be no freeing within the loop. In this instance, this transforms into a buffer-swapping operations:

```

d = alloc (...); / set RC to 2 */
for (i = s; i < 10; i++) {
    D = d;
    for (iv = 0; iv < 11; iv++) D[iv] = i;
    print (D); /* RC of D decremented */
    d = D;
}
decr (d);
```

where D is some new array which is not allocated, but instead used as an alias to the array d. At the end of the *for-loop* we swap

the buffers. Finally, when we exit the *for-loop*, we decrement the reference-count of d, which ultimately frees it.

4 EXTENDED MEMORY REUSE OPTIMISATION

We have designed the *EMR* code-transformations to work over an AST similar to that of SAC [16, 17] with some existing reference-counting related infrastructure [18]. By no means does this imply that we are restricted to that setting, the optimisation could be made to work with any referenced-counted language. An important constraint to note is that at the point where our code-transformations are applied, reference-counting as such has *not yet* been introduced into the AST. Rather, our code-transformation adds values to existing tables which are later used by other code-transformations to introduce reference-counting to the code [8]. The tables contain the reuse candidates for *with-loops*.

Some general assumptions that we make regarding the structure of the AST is that *expressions* exist within the body of a *function definition*, *conditionals*, and *with-loops*. For the sake of our algorithmic specifications, we assume that these expressions are accessible as ordered elements within a list. For simplicity's sake we inspect *conditionals* implicitly, their expressions being exposed to us directly. Additionally, constructs like loops are expressed as *loop-functions*, which are tail-end recursive. Finally, we need to access and potentially modify properties of a given expression, such as the arguments of a *function application*.

We have split the *EMR* code-transformations into three distinct parts. The first inferences *ERCs*, the second filters these, and lastly we lift allocations out of loops. The high-level overview of what the code-transformations are and what they achieve is given:

- (1) extended memory reuse inference
 - (a) analyse function bodies and collect all result arrays from function applications and *with-loops*. These we call our *extended reuse candidates (ERCs)*.
 - (b) annotate all² *with-loops* and all recursive *loop-function* with the current list of collected *ERCs*.
- (2) extended reuse candidate filtering
 - (a) filter out all invalid *ERCs* (we will go into more detail later about what *invalid* means).
- (3) extended loop memory propagation
 - (a) for loops, find *with-loops* with no reuse candidates (normal or extended) and lift the memory allocation out of the loop. Subsequent recursive calls will pass on a reference to some free variable within the *loop-function* body (previously collected as an *ERC*).

4.1 Extended Reuse Candidate Inference

This purpose of this code-transformation is to infer *ERCs*. As previously stated, these candidates are not referenced inside of a *with-loop*, but are any array found within the scope of a given function definition that is of the same type and shape.

Our inference code-transformation (shown in algorithm 1) expects a function definition (F_d) and two tables of mappings. The first table, called FC (function candidates), is only used when we

²except fold/reduce operations

are dealing with a *loop*-function. In this case, we store *all ERCs* up until the recursive call. The second table, called EC (expression candidates), stores all *with-loop ERCs* that we find. Both of these tables are stored globally and will be used by later code-transformations.

We create a local list *C* which we use to store all found candidates. We iterate through all expressions in F_d , looking for array assignments being made by either *with-loops* or function applications. We also implicitly inspect the expressions of if- and else-branches of conditionals.

When we come across a *with-loop*, we search the body for free variables. We need to filter out all candidates from *C* which are referenced within the body of the *with-loop*. Additionally we need to filter out candidates that are of a different shape than given by the *with-loop*. Once this is done, we map our new list \mathcal{ERC} to the *with-loop* expression and store this within our table EC. Finally we append the return values of the *with-loop* to *C*.

When we find a function application, we do one of two things. If F_d is a *loop*-function and the application is the recursive call, we make a copy of our currently collected candidates *C* and map this to F_d , storing this in the FC table. Otherwise we just add the return values of the application to *C*.

We continue updating *C* as we iterate over all the expressions in F_d . We now have stored all mappings of candidates in our two tables which we will use in the next transformations.

4.2 Extended Reuse Candidate Filtering

When inferring the *ERCs*, we over-approximate and collect candidates that are likely invalid. We filter these out using the code-transformation shown in algorithm 2, which moves bottom-up through the *assignments* of the function definition F_d . It looks for and removes *ERCs* that are (1) referenced after the given *with-loop*, and (2) that are already used as a normal reuse candidate. These two criteria exist in order to eliminate the chance of introducing erroneous allocations. In the first instance, if an *ERC* is referenced later, we inadvertently will create a copy of it, meaning we introduce two allocations. The first being the initial allocation and the latter at the point where we intended to reuse memory. In the second case, if the *ERC* has previously been identified as a normal reuse candidate, it is fairly clear that it will be used for memory reuse, so the *ERC* is redundant.

In use two tables of mappings EC and CRC which respectively hold our *ERC* values, and the current reuse candidates for each *with-loop*.

Finally, after filtering our invalid *ERCs*, we truncate the list \mathcal{ERC} down to one *ERC* value. We do this to avoid invalidating our *ERCs* when introducing explicit memory operations into the code. The following example, where we define the function baz shows this:

```
int[11], int[11] baz (int[11] a, int[11] b) {
    print(a);
    print(b);

    k = with {} : genarray ([11], 9); /* ERC: a, b */
    j = with {} : genarray ([11], 6); /* ERC: a, b, k */

    return (k, j);
}
```

Algorithm 1: Extended Reuse Candidate Inference

Input: A function definition F_d
A table FC of (*function definition* \mapsto *list*)
mappings
A table EC of (*expression* \mapsto *list*) mappings

```

1 begin
2     Let C be a list of function-level candidates which
     includes the arguments of  $F_d$ 
3     foreach (var = expr) in  $F_d$  do
4         switch expr do
5             case expr is a with-loop do
6                 Let  $B_f$  be list of FreeVariables(expr)
7                 Let  $W_s$  be the shape of expr
8                 Let  $\mathcal{ERC}$  be Filter( $x \Rightarrow x \notin B_f$ 
9                     and Shape( $x$ ) ==  $W_s$ , C)
10                /*  $\mathcal{ERC}$  contains all candidates
11                   that can be used for expr */
12                Add expr  $\mapsto$   $\mathcal{ERC}$  to EC
13                C = C  $\cup$  var
14            end
15        case expr is a function application do
16            if  $F_d$  is a loop-function
17                and expr is recursive call then
18                    Let  $C_c$  be a copy of C
19                    Add  $F_d \mapsto C_c$  to FC
20                else
21                    Append return value(s) of expr to C
22            end
23        otherwise do           /* do nothing */
24    end
25 end
26
```

Output: The table FC of mappings
The table EC of mappings

The function accepts and returns two arrays of the same shape. Inside we print our two parameters *a* and *b*, and then generate two new arrays *k* and *j* and return these. Our inference has identified the former two arrays as *ERCs* for the latter two arrays. For the last *with-loop* we also identify *k* as an *ERC*, which is clearly an over-approximation. When it comes time to introduce the explicit memory operation, the *ERCs* will be seen as being referenced by the other *with-loop*, which immediately invalids them. This results in only the first *with-loops* reusing the memory of *b*:

```
k = with {} : genarray ([11], 9); /* only reuses: b */
j = with {} : genarray ([11], 6);
```

Really what we want is that both *with-loop* are transformed to reuse the memory of *a* and *b*:

```
k = with {} : genarray ([11], 9); /* reuses: b */
j = with {} : genarray ([11], 6); /* reuses: a */
```

Algorithm 2: Extended Reuse Candidate Filtering

Input: A function definition F_d
 A table EC of ($expression \mapsto list$) mappings
 A table CRC of ($expression \mapsto list$) mappings

```

1 begin
  Let  $\mathcal{VR}$  be a (initially empty) list variable references
  foreach ( $var = expr$ ) in  $F_d$  do /* bottom-up */
    switch  $expr$  do
      case  $expr$  is function application do
         $\mathcal{VR} = arguments\ of\ expr \cup \mathcal{VR}$ 
      end
      case  $expr$  is a with-loop do
        Let  $\mathcal{ERC}$  be  $Filter(x \Rightarrow x \notin CRC[expr]$ 
          and  $x \notin \mathcal{VR}$ ,  $EC[expr]$ )
         $\mathcal{VR} = all\ function\ application\ arguments \in$ 
           $expr\ body \cup \mathcal{VR}$ 
         $\mathcal{VR} = CRC[expr] \cup \mathcal{ERC}_0 \cup \mathcal{VR}$ 
         $EC[expr] = \mathcal{ERC}_0$ 
      end
      otherwise do /* do nothing */
    end
  end
end
```

Output: The table EC of mappings

Because we are bottom-up, we remove k as it is referenced by the return statement. We select a as our main ERC and discard b. The next *with-loop* up we discard a as it is referenced by the previous *with-loop* for reuse, leaving use with b.

4.3 Extended Loop Memory Propagation

The last stage of code-transformations for the *EMR* optimisation is lift array allocations out of loops. As mentioned in section 1, by lifting out allocations we can reduce the number of memory operation done within the loop, thereby reducing performance overheads.

We have split out our code-transformation into two parts, which (1) finds *with-loops* in *loop-functions* which don't have a reuse candidate or ERC , meaning that we allocate in the loop, and creates a new argument for the *loop-function* which it sets as the reuse candidate for the *with-loop* and updates the recursive call; (2) searches for the initial application of the *loop-function* and assigns it either an ERC from the current scope or generates a new variable.

Part 1. We iterate top-down the assignments of the function definition F_d which is a tail-end recursive *loop-function*, as shown in algorithm 3. When we come across a *with-loop*, which check if the it has any reuse candidates or $ERCs$. If it has *none*, we add create a new array variable with the same shape as the return value of the *with-loop* and set this as both an argument to F_d and a reuse candidate of the *with-loop* but updating the table CRC. On the other hand, we collect the reuse candidates and $ERCs$ into $C\mathcal{R}$ which we will use later to find arrays with which to update the recursive call of the *loop-function*. When we find the recursive call, we search

Algorithm 3: Extended Loop Memory Propagation (part 1)

Input: A function definition F_d
 A table FC of mappings ($function\ definition \mapsto list$ of candidates)
 A table EC of mappings ($expression \mapsto list$ of candidates)
 A table TRC of mappings ($expression \mapsto list$ of reuse candidates)

```

1 begin
  Let  $\mathcal{TR}$  be a (initially empty) list of temporary reuse
  candidates
  Let  $\mathcal{CR}$  be a (initially empty) list of collected reuse
  candidates
  if  $F_d$  is a loop-function then
    foreach ( $var = expr$ ) in  $F_d$  do
      switch  $expr$  do
        case  $expr$  is a with-loop do
          if  $EC[expr] = \emptyset$  and  $CRC[expr] = \emptyset$ 
            then
               $V_{rc} = new\ array\ with\ Shape(var)$ 
               $\mathcal{TR} = V_{rc} \cup \mathcal{TR}$ 
               $CRC[expr] = V_{rc}$ 
          else
             $\mathcal{CR} = EC[expr] \cup CRC[expr] \cup C\mathcal{R}$ 
          end
        end
        case  $expr$  is a function application do
          if  $expr$  is the recursive call then
             $\mathcal{ERC}_f = FC[F_d]$ 
             $N_a =$ 
               $FindMatching(\mathcal{TR}, \mathcal{ERC}_f \cup C\mathcal{R})$ 
            Append  $N_a$  to arguments of  $expr$ 
          end
        end
        otherwise do /* do nothing */
      end
    end
    Append  $\mathcal{TR}$  to arguments of  $F_d$ 
  end
end
```

Output: the function definition F_d
 the table TRC of mappings

through the list of *lifted* temporary arrays to find match in both the $ERCs$ of F_d and the our collected reuse candidates in $C\mathcal{R}$. Assuming we find suitable arrays, we update the arguments of the recursive call.

Part 2. In algorithm 4 we go through all function definitions F_d again searching for the initial function application of the *loop-function* we previously updated. When we find it we call we either (1) take a ERC in the current scope, or (2) we create a new array

Algorithm 4: Extended Loop Memory Propagation (part 2)

Input: A function definition F_d
A table FD of mappings
 $(\text{expressions} \leftrightarrow \text{function definitions of loops})$

```

1 begin
2   foreach (var = expr) in  $F_d$  do
3     switch expr do
4       case expr is a function application do
5         if expr ∈ FD and expr is not the recursive
          call then
6           if number of expr arguments < number
              of  $F_d$  arguments then
7              $\mathcal{NA}$  =
               CreateArgs( $\Delta$  number of arguments)
8             Append  $\mathcal{NA}$  to arguments of expr
9           end
10        end
11      end
12    otherwise do /* do nothing */
13  end
14 end
15 end

```

and allocate it. We update the arguments of the *loop*-function application with the new arrays.

5 RESULTS AND ANALYSIS

We now present the results of using the *EMR* optimisation on several benchmarks. We use several benchmarks from different benchmarks suites and compile them for both SMA machines and GPUs (using CUDA) using the sac2c compiler (based on version 1.3.2) with GCC 4.8.5 and CUDA 9.0.

The benchmarks include basic ones such as matrix multiplication and matrix relaxation. We also use a kernel from a real-world application by BGS called the Global Geomagnetic Model which was translated to SAC in [19].

Additionally we include benchmarks from the Livermore Loops suite [14], the Rodinia Benchmark suite [3], and the NAS Parallel Benchmark suite³ [1].

For the Livermore Loop benchmarks, we iterate 10^6 times. We don't use the same problem size for all of them as otherwise runtimes become excessive. A list of benchmarks and their respective problem sizes is given in table 1. Some benchmarks have several implementations, which we call variants throughout this section. For the Livermore Loop benchmarks we have variants such as the C implementation which is nearly identical to that of the C formulation of the loops. We also have naive implementations which may make use of *with-loops* but does so in a trivial manner. We also

³SAC implementations of the Livermore Loops, Rodinia Benchmarks, and the NAS Parallel Benchmarks (as well as miscellaneous benchmarks) are available at <https://github.com/SacBase>.

Suite	Benchmark	Problem Size
Livermore Loops	Loop01	100001 vector
	Loop02	100001 vector
	Loop03	100001 vector
	Loop04	10001 vector
	Loop05	100001 vector
	Loop07	100001 vector
	Loop08-split	1001 vector
	Back-propagation	65536 input units
Rodinia Benchmarks	Hotspot	1024×1024 grid
	Kmeans	34 features, 5 clusters
	Leukocytes	n/a
	LUD	2048×2048 matrix
	Needleman-Wunsch	24×24 grid
	Particle Filter	n/a
	Pathfinder	100×512000 grid
NAS Parallel	SRAD	2048×2048 matrix
	Conjugate gradient	14000, 75000, and 150000
Miscellaneous	Embarrassingly parallel	2^{24}
	Multigrid	128^3
	Matrix Multiplication	1000×1000 matrices
	Matrix Relaxation	10000×10000 matrix
	Gauss-Jordan	500×500 matrix
	FFT	n/a
	BGS Kernel	1000×1000 matrix

Table 1: Benchmarks and their Problem Sizes

have APL-like implementations and finally an implementation that uses SAC's set-comprehension notation.

As we are interested in the effectiveness of the *EMR* optimisation, we count the number of allocations during execution. Our collected values are shown in table 2, with a reduction in allocations given in green and the opposite in red. Yellow is used to indicate that the data set is not complete due to problems with the benchmark⁴. Furthermore, we note that the sequential numbers and the CUDA number won't always be that same as in the CUDA case we only measure the number of allocations on the GPU device, *not* the host machine.

We also measure the wall-clock runtime of the benchmarks to see if a change in the number of allocations affects runtime. We run each benchmark five times and take the average of the runtimes as our measurement. These results, given as speedups, are shown in fig. 1. In most cases there isn't a significant change in runtime. For the Livermore Loop 8 we can see that with the *EMR* optimisation applied we gain up-to a 4x speedup when compiling for CUDA devices, we also get a speedup of 1.5x on the host.

A deeper analysis of the results is given in the respective subsection for each benchmark or suite of benchmarks.

5.1 Livermore Loops

For the majority of the Livermore Loop benchmarks we observe no change in the number of allocations or in runtimes, this also holds for the various variants. For most variants such as TC, APL, and Naive, the existing reuse candidate inference present in the compiler is able to identify a candidate from the arguments of the *loop*-function which in of itself leads to a lifting of allocations out of the loop. This is why our count of number of allocations is so low, given that we iterate 1000000 times. The exception to this is

⁴More often than not, the compiler was not able to generate code for the specified target.

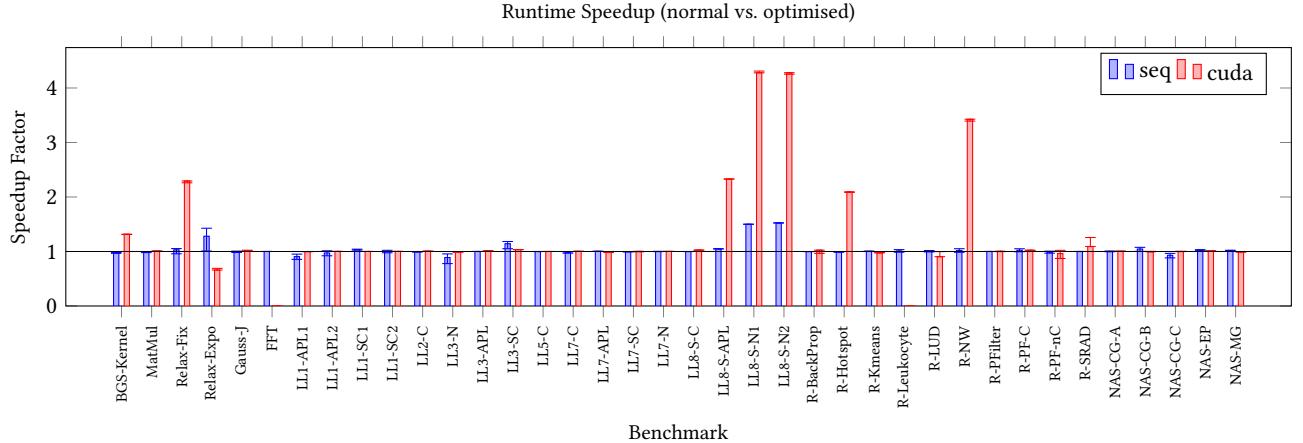


Figure 1: Speedup Results

Benchmarks		Sequential		CUDA	
Name	Variant	Unopt.	Opt.	Unopt.	Opt.
BGS Kernel		2570925	9521	2574279	9532
MatMul		4	4	5	4
Relax	Fix	101	2	102	3
	Expo	101	2	201	201
FFT		16	16		
Gauss-Jordan		1003	504	1501	1002
LL01	APL1	9	9	2	2
	APL2	9	9	2	2
	SC1	9	9	2	2
	SC2	9	9	2	2
LL02	C	9	9	1	1
LL03	N	7	7	1	1
	APL	7	7	1	1
	SC	7	7	1	1
LL05	C	8	8	1	1
LL07	C	1000007	1000007	1	1
	APL	8	8	2	2
	SC	8	8	2	2
	Naive	8	8	2	2
LL08-Split	C	3000007	3000007	1	1
	APL	3000007	10	3000003	6
	N1	3000007	10	3000003	6
	N2	3000007	10	3000003	6
Rodinia Back Propagation		9	9	805	805
Rodinia Hotspot		5001	2	5002	3
Rodinia Kmeans		2427453447	2427453447	929	929
Rodinia Leukocytes		725	725		
Rodinia LUD		2	2	2	2
Rodinia Needleman-Wunsch		4198403	4194310	4097	4
Rodinia Particle Filter		19	19	36	61
Rodinia Pathfinder	COND1	101	3	102	4
	no-COND1	101	101	229	229
Rodinia SRAD		1003	1003	2002	2002
NAS CG	A	28443	28068	28863	28503
	B	152183	150308	154283	152483
	C	302183	300308	304283	302483
NAS EP		100	100	91	91
NAS MG		794	794	588	588

Table 2: Allocation Counts of Targets seq and cuda

with they C variant, which typically uses no *with-loops* and such offers no opportunities for memory reuse.

With Loop 8 Split, whereby we split the 3-dimension input array into two 2-dimensional arrays, we observe a significant reduction in the number of allocations and speedups. The most significant speedup being a 4.1x with variants Naive-1 and Naive-2 when using CUDA. Additionally for the same variants we have a 1.5x speedup

running sequential. The lifting of allocations out of the loop is the primary causes of these speedups, especially on the GPU device.

5.2 Rodinia Benchmarks

For several benchmarks we see no change in number of allocations or runtime. The reasons for these include there not being any opportunities for memory reuse, such as with the LUD benchmark. Furthermore, the existing reuse candidate inference in the compiler is often sufficient, as is the case with the Back-propagation benchmark. Both reasons are why the SRAD and Kmeans benchmark have no observable changes. Unfortunately, we were unable to generated CUDA code for the Leukocytes benchmark. We don't observe any change in the number of allocation here though.

Hotspot. Here we see a significant decrease in the number of allocations for both targets, with a 2x factor speedup for the CUDA case. The reduction is mainly due to the lifting of allocations out of a loop construct which iterates 5000 times.

Needleman-Wunsch. Here we observe a reduction in the number of allocations for both targets with an almost 4x speedup in the CUDA case. The reduction is primarily due to the lifting of allocations out of two loop constructs that both iteration 2049 times each. The significant speedup achieved in the CUDA case is also due to this reduction, partial also because the size of the allocation at each iteration.

Pathfinder. Here we use two variant of the benchmark, the first (which we call *cond1*) makes use of a *with-loop* containing an if-else construct to perform the main path finding operation. The second (*non-cond1*) implements this same if-else construct as *with-loop* partition, meaning that we no longer branch within the *with-loop*. These variants exist because in the CUDA case, the latter variant performs better than the former case [ref?].

We observe that the *cond1* variant has a significant reduction in the number of allocations with no significant change in runtime. The reduction is primarily caused by the lifting of allocations out of the loop construct which iterates over the rows of the grid. In the *no-cond1* variant we see no change whatsoever to the number

of allocations or runtime. The number of allocations with CUDA is much large then with the sequential case, and likely this is caused by how the partitions of the *with*-loop are being turned into CUDA kernels. The underlying structure of the *with*-loop though is similar to that of the *condi* variant, and so we would expect that memory reuse would be possible here as well. The cause of why we do not observe this unclear.

Particle-Filter. With this benchmark we observe an *increase* in the number of allocations for the CUDA case with no significant change in runtime. When inspecting the transformations that are done by the *EMR* optimisation, we observe that were in the sequential case we had inferred some reuse candidates we no longer did so. This opens up the opportunity for extended candidates to be used instead. In this instance though our chosen candidate needs to be additionally re-allocated because of when CUDA kernels are being launched.

5.3 NAS Parallel Benchmarks

For both the EP and MG benchmarks, the *EMR* optimisation has no effect. There is no change in the number of allocations or the runtime. For EP there are no loops and minimal number of opportunities for memory reuse besides.

For CG we observe a reduction in the number of allocations for all variants and targets. We do not find that the runtime of these changes significantly. The lifting of allocation is the primary reason for the reduction in allocations, but there are several instances were the ERCI optimisation is able to increase the number of suitable reuse candidates.

5.4 Miscellaneous Benchmarks

Matrix Multiplication. For matrix multiplication benchmark we only see a change in the number of allocations when compiling for CUDA — this is achieved by the EMRI optimisation, which identifies one opportunity for memory reuse. Though we have one less allocation within a CUDA kernel, the effect on runtime is negligible. For both targets, there is no measured change in runtime.

Matrix Relaxation. We use two variants of matrix relaxation, one that uses fixed-point iteration and the other that uses exponential iteration. For both variants we observe a large decrease in the number of allocations when compiling for both targets, except for the exponential variant compiled for CUDA. The reduction is primarily caused by the application of the ELMP optimisation which has lifted out allocations from the loop construct. We notice that the runtime for the sequential case seems to worsen with the *EMR* optimisation applied, though there is large amount of error. With the CUDA case we observe no significant change in runtime. For the exponential case using CUDA we do not observe a reduction in allocations, as occurs otherwise. This happens because on each iteration we need to check if we have satisfied the loop conditional, which means we must access the latest GPU computed value.

Gauss-Jordon Elimination. For gauss-jordon elimination benchmark we see a large reduction in the number of allocations but no significant change in runtime. The reduction is primarily due to the lifting of allocations out of the loop, which has a fixed iteration

equal to the number of rows within the matrix (500×500). We can see that we reduce the number of allocations by 500.

Fast Fourier Transformation. We were not able to generate code for the CUDA case, as such are results are limited to the sequential case. Here we observe that there is no change in the number of allocation, or in the runtime of the benchmark. On the one hand the existing reuse candidate inference within sac2c has already identified suitable candidates, but also that there is no loop construct with *with*-loops present.

BGS Kernel. For the BGS kernel we can see a significant decrease in the number of allocations and in the CUDA case a $1.5\times$ factor speedup due to that reduction. Within the benchmark there are several loops with *with*-loops and a large number of *with*-loop constructs in general. The lifting of allocations from the loops is the primary reason for the reduction in allocations. One loop in particular iterates up to 30000 times and is responsible for many allocations.

6 RELATED WORK

To be done.

7 CONCLUSION AND FUTURE WORK

To be done.

REFERENCES

- [1] David H Bailey, Eric Barszcz, John T Barton, David S Browning, Robert L Carter, Leonardo Dagum, Rod A Fatoohi, Paul O Frederickson, Thomas A Lasinski, Rob S Schreiber, et al. 1991. The NAS parallel benchmarks. *The International Journal of Supercomputing Applications* 5, 3 (1991), 63–73.
- [2] David C. Cann and Paraskevas Epiridou. 1995. Advanced array optimizations for high performance functional languages. *IEEE Transactions on Parallel and Distributed Systems* 6, 3 (March 1995), 229–239. <https://doi.org/10.1109/71.372771>
- [3] Shuai Che, M. Boyer, Jiayuan Meng, D. Tarjan, J. W. Sheaffer, Sang-Ha Lee, and K. Skadron. 2009. Rodinia: A benchmark suite for heterogeneous computing. In *IEEE International Symposium on Workload Characterization (ISWC '09)*, 44–54. <https://doi.org/10.1109/ISWC.2009.5306797>
- [4] George E. Collins. 1960. A Method for Overlapping and Erasure of Lists. *Commun. ACM* 3, 12 (Dec. 1960), 655–657. <https://doi.org/10.1145/367487.367501>
- [5] Python Software Foundation. 2018. Python 3.7 Language Documentation. <https://docs.python.org/3/c-api/index.html> Online, accessed 14 August 2018.
- [6] Clemens Grelck and Sven-Bodo Scholz. 2008. Efficient Heap Management for Declarative Data Parallel Programming on Multicores. In *3rd Workshop on Declarative Aspects of Multicore Programming (DAMP 2008), San Francisco, CA, USA*, Manuel Hermenegildo, Leaf Peterson, and Neal Glew (Eds.). ACM Press, 17–31.
- [7] C. Grelck and K. Trojahn. 2004. Implicit Memory Management for SAC. *Proceedings of the 16th International Workshop on Implementation and Application of Functional Languages (IFL '04)* (September 2004), 335–348. University of Kiel, Germany.
- [8] Clemens Grelck and Kai Trojahn. 2004. Implicit Memory Management for SaC. In *Implementation and Application of Functional Languages, 16th International Workshop, IFL '04*, Clemens Grelck and Frank Huch (Eds.). University of Kiel, Institute of Computer Science and Applied Mathematics, 335–348. Technical Report 0408.
- [9] Jing Guo, Robert Bernecker, Jeyarajan Thiagarajan, and Sven-Bodo Scholz. 2014. Polyhedral Methods for Improving Parallel Update-in-Place. In *Proceedings of the 4th International Workshop on Polyhedral Compilation Techniques*, Sanjay Rajopadhye and Sven Verdoolaege (Eds.). Vienna, Austria.
- [10] Jing Guo, Jeyarajan Thiagarajan, and Sven-Bodo Scholz. 2011. Breaking the Gpu Programming Barrier with the Auto-parallelising Sac Compiler. In *6th Workshop on Declarative Aspects of Multicore Programming (DAMP'11), Austin, USA*. ACM Press, 15–24. <https://doi.org/10.1145/1926354.1926359>
- [11] Paul Hudak and Adrienne Bloss. 1985. The Aggregate Update Problem in Functional Programming Systems. In *Proceedings of the 12th ACM SIGACT-SIGPLAN Symposium on Principles of Programming Languages (POPL '85)*. ACM, New York, NY, USA, 300–314. <https://doi.org/10.1145/318593.318660>

- [12] Apple Inc. 2018. Swift Language Documentation. <https://docs.swift.org/swift-book/> Online, accessed 14 August 2018.
- [13] Steve Klabnik and Carol Nichols. 2018. *The Rust Programming Language* (2 ed.). No Starch Press.
- [14] Frank H. McMahon. 1986. *The Livermore Fortran Kernels: A computer test of the numerical performance range*. Technical Report UCRL-53745. Lawrence Livermore National Lab., CA, USA.
- [15] NVIDIA Corporation 2018. *CUDA C Programming Guide* (9.2.148 ed.). NVIDIA Corporation. <https://docs.nvidia.com/cuda/cuda-c-programming-guide/> Online, accessed 12 Aug. 2018.
- [16] SaC Development Team 2016. *SaC EBNF Grammar*. SaC Development Team. Available online at <http://www.sac-home.org/doku.php?id=docs:syntax>.
- [17] Sven-Bodo Scholz. 2003. Single Assignment C: efficient support for high-level array operations in a functional setting. *Journal of Functional Programming* 13, 6 (2003), 1005–1059. <https://doi.org/10.1017/S0956796802004458>
- [18] Kai Trojahnner. 2005. *Implicit Memory Management for a Functional Array Processing Language*. Master's thesis, Universität zu Lübeck.
- [19] Hans-Nikolai Vießmann, Sven-Bodo Scholz, Artjoms Šinkarovs, Brian Bainbridge, Brian Hamilton, and Simon Flower. 2015. Making Fortran Legacy Code More Functional. *27th Symposium on Implementation and Application of Functional Languages (IFL '15)* (2015). <https://doi.org/10.1145/2897336.2897348>

Continuation Passing Style and Primitive Recursive Functions

(Draft)

Ian Mackie

ABSTRACT

A tail recursive function is one in which the last operation performed is the recursive call. Such functions can be implemented as a loop, so without the overhead of stack frames needed to handle the function call/return operations. All functions can be written in tail recursive form by using a *continuation passing style* transformation, but this transformation adds an extra parameter (which is the stack frame internalised). This extra parameter however is a higher-order function, and may therefore off-set any advantages. The good news is that sometimes a simple representation of this extra parameter—for example a number or a list—can be found. This paper addresses which functions can be converted to tail recursive form where a simple data structure can be used to represent the extra parameter.

KEYWORDS

CPS, dataflow, primitive recursive functions

ACM Reference Format:

Ian Mackie. 2019. Continuation Passing Style and Primitive Recursive Functions: (Draft). In *Proceedings of International Symposium on Implementation and Application of Functional Languages (IFL'18)*. ACM, New York, NY, USA, 2 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnn>

1 INTRODUCTION

One of the many advantages of functional programming is the ability to write programs that reflect, or directly correspond to, the mathematical definition of the problem. This is because at the heart of most functional programming languages one finds recursive functions as the most dominant programming technique.

From an implementation perspective, recursive functions have an overhead when compared to iteration. Some languages provide this alternative to the programmer and some compilers will try to optimise by converting a recursive function to iteration. However neither is perfect:

- If the programmer is offered iteration, then the algorithm may not be the same as the mathematical definition, taking away one of the main advantages of functional programming in the first place.
- All recursive functions can be converted to iteration, but at a cost.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

IFL'18, August 2019, Lowell, MA, USA

© 2019 Association for Computing Machinery.
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00
<https://doi.org/10.1145/nnnnnnnn.nnnnnnn>

Although every recursive function can be transformed into a tail-recursive one (and hence into a loop), the continuation may not be representable as a simple value such as a number or a list. This paper addresses when the continuation can be represented by a simple value. We show this by passing through some other work on representing functions in a linear calculus, and use results known in that setting. This yields the main result of the paper that if a function is primitive recursive, then the continuation can be represented by a simple value.

Overview. The rest of this draft paper is organised summarised in the following section where we give a brief overview of the idea. In the full paper we will develop the definitions and transformations in detail to show the result formally.

2 LINEAR FUNCTIONS AND CONTINUATIONS

2.1 Linear λ -calculus

Assume a set of variables v_0, v_1, \dots , and let x be an arbitrary variable, then the set of linear λ -terms Λ_L is the least set satisfying

- $x \in \Lambda_L$ (variable)
- $t \in \Lambda_L, x \in \text{FV}(t) \Rightarrow \lambda x.t \in \Lambda_L$ (abstraction)
- $t, u \in \Lambda_L, \text{FV}(t) \cap \text{FV}(u) = \emptyset \Rightarrow tu \in \Lambda_L$ (application)

The typing rules for this restricted calculus is given by the following three rules:

$$\frac{}{x : A \vdash x : A} \quad \frac{\Gamma, x : A \vdash t : B}{\Gamma \vdash \lambda x.t : A \multimap B} \quad \frac{\Gamma \vdash t : A \multimap B \quad \Delta \vdash t : A}{\Gamma, \Delta \vdash tu : B}$$

We adding pairs, and generalise the abstraction rule to allow pair patterns:

Pattern	Variable Constraint	Bound Variables (BV)
x	—	$\{x\}$
$\langle p, q \rangle$	$\text{BV}(p) \cap \text{BV}(q) = \emptyset$	$\text{BV}(p) \cup \text{BV}(q)$

This gives the following typing rules:

$$\frac{}{p : A \vdash p : A} \quad \frac{\Gamma, p : A, q : B \vdash t : C}{\Gamma, \langle p, q \rangle : A \otimes B \vdash t : C} \text{ (Pattern Pair)} \\ \frac{\Gamma, p : A \vdash t : B \quad \Gamma \vdash t : A \quad \Delta \vdash u : B}{\Gamma \vdash \lambda p.t : A \multimap B \quad \Gamma, \Delta \vdash \langle t, u \rangle : A \otimes B} \text{ (Pair)}$$

Examples or terms of this linear calculus include functions to reorganise data (commutativity and associativity):

$$\lambda \langle x, y \rangle. \langle y, x \rangle : A \otimes B \multimap B \otimes A \\ \lambda \langle \langle x, y \rangle, z \rangle. \langle x, \langle y, z \rangle \rangle : (A \otimes B) \otimes C \multimap A \otimes (B \otimes C)$$

We can now increase the computational power by adding numbers and an iterator. The extra rules that we need for typing are:

$$\frac{}{\Gamma \vdash 0 : \text{nat}} \text{ (Zero)} \quad \frac{\Gamma \vdash t : \text{nat}}{\Gamma \vdash S t : \text{nat}} \text{ (Succ)}$$

$$\frac{\Gamma \vdash t : \text{nat} \quad \Delta \vdash u : A \quad \Theta \vdash v : A \multimap A}{\Gamma, \Delta, \Theta \vdash \text{iter } t \ u \ v : A} \quad (\text{Iter})$$

And the extra reduction rules are the following:

$$\begin{array}{lll} (\lambda p.t)v & \rightarrow [p \ll v].t & \text{FV}(v) = \emptyset \\ \text{iter } (S \ t) \ u \ v & \rightarrow \text{iter } t \ (vu) \ v & \text{FV}(v) = \emptyset \\ \text{iter } 0 \ u \ v & \rightarrow u & \text{FV}(v) = \emptyset \\ \hline [x \ll v].t & \rightarrow t[v/x] & \\ [\langle p, q \rangle \ll \langle t, u \rangle].t & \rightarrow [p \ll t].[q \ll u].t & \end{array}$$

where we have introduced and defined a matching construct that unpacks pairs to create usual substitutions. We can now write some more interesting functions, including the following examples:

- Addition, multiplication and exponentiation can be defined as:

$$\begin{array}{lll} \text{add} & = \lambda mn.\text{iter } m \ n \ (\lambda x.S \ x) \\ \text{mult} & = \lambda mn.\text{iter } m \ 0 \ (\text{add } n) \\ \text{exp} & = \lambda mn.\text{iter } n \ (S \ 0) \ (\text{mult } m) \end{array}$$

- When we need to copy or erase:

$$\begin{array}{ll} C & = \lambda x.\text{iter } x \ \langle 0, 0 \rangle \ (\lambda \langle a, b \rangle.(Sa, Sb)) \\ \text{fst} & = \lambda \langle n, m \rangle.\text{iter } m \ n \ (\lambda x.x) \end{array}$$

Where $C : \text{nat} \multimap \text{nat} \otimes \text{nat}$, and $\text{fst} : \text{nat} \otimes \text{nat} \multimap \text{nat}$.

- Note that the functions iterated for copy and erase are closed

Using the linear λ -calculus with an iterator we have great computational power:

- Closed reduction \iff Gödel's System T
- Closed construction \implies Primitive Recursive Functions

Closed reduction means that we only fire an iterator reduction when the iterated function is closed. This means that we preserve syntactical linearity in the calculus (copying free variables would violate this). Closed construction is a more restrictive constraint that means that we are only allowed to create programs with iterators in them when the iterated function is closed to begin with. Using this calculus, we have an interesting way therefore to distinguish between Gödel's System T and Primitive Recursive Functions.

The single implication above showing that closed construction gives primitive recursive functions is due to the fact that there are many primitive recursive functions that are not closed by construction, for example mult can be written two ways:

$$\begin{array}{ll} \text{mult} & = \lambda mn.\text{iter } m \ 0 \ (\text{add } n) \\ \text{mult} & = \lambda mn.\text{iter } m \ \langle n, 0 \rangle \ (\lambda \langle x, y \rangle.(x, \text{add } x \ y)) \end{array}$$

In the second way, the function is now closed, but the function is not linear: x is used twice. Fortunately, we know how to copy numbers with a closed linear function as shown above, so we can create a closed by construction multiplication if we transform using pairs (an extra parameter) and copy numbers.

This leads to a natural question:

- Can we always close functions in this way?

Of course, we have hinted to this already, as only primitive recursive functions can be captured this way. However, it leads to a

interesting sequence of observations. Consider again the function Mult :

$$\text{mult} = \lambda mn.\text{iter } m \ 0 \ (\text{add } n)$$

which we can write as a usual recursive function as follows:

$$\text{mult } 0 \ n = 0$$

$$\text{mult } (S \ m) \ n = \text{add } (\text{mult } m \ n) \ n$$

This can be transformed using continuation passing style to give the following:

$$\text{mult } 0 \ n \ k = k \ 0$$

$$\text{mult } (S \ m) \ n \ k = \text{mult } m \ n \ (\lambda r \rightarrow k(\text{add } r \ n))$$

which corresponds in our calculus to the following function:

$$\text{mult} = \lambda mn.\text{iter } m \ \langle n, \lambda x.x \rangle \ (\lambda \langle x, k \rangle.(x, \lambda r.k(\text{add } r \ x)))$$

Thus the closing and coping transformation above is nothing more than CPS transformation in another setting. Next, if we examine the shape of the continuation:

$$\lambda r.r + n, \lambda r.(r + n) + n = \lambda r.r + 2n, \dots$$

it is well established that this we can represent the continuation as a number: a much simpler data type than the lambda calculus. The alternative version written using an accumulating parameter demonstrates this:

$$\text{mult } 0 \ n \ a = a$$

$$\text{mult } (S \ m) \ n \ a = \text{mult } m \ n \ (\text{add } n \ a)$$

which again can be written in our language as"

$$\text{mult} = \lambda mn.\text{iter } m \ \langle n, 0 \rangle \ (\lambda \langle x, y \rangle.(x, \text{add } x \ y))$$

(Note: we didn't show the copy or projections to keep the examples simple.)

In the full paper we show how the ideas used in this example lead to the main result of this paper:

THEOREM 2.1. *If a function is primitive recursive, then the continuation can be represented by a simple data type, such as a number.*

3 CONCLUSION

In this draft paper we have outlined a result that a continuation passing style transformation can offer a simple representation of the continuation when the function is primitive recursive.

Extracting Verified Constraints from Coq-embedded Array DSLs

— Extended Abstract —

Artjoms Šinkarovs
School of Mathematical and
Computer Sciences
Heriot-Watt University
Edinburgh, Scotland, UK
a.sinkarovs@hw.ac.uk

Sven-Bodo Scholz
School of Mathematical and
Computer Sciences
Heriot-Watt University
Edinburgh, Scotland, UK
s.scholz@hw.ac.uk

Hans-Nikolai Vießmann
School of Mathematical and
Computer Sciences
Heriot-Watt University
Edinburgh, Scotland, UK
hv15@hw.ac.uk

ABSTRACT

Dependently-typed languages make it possible to equip programs with rich and useful properties and verify their consistency by means of type checking. If we want to run these programs efficiently, typically we pre-optimize our programs by erasing computationally irrelevant proof objects.

This paper suggests to pass on such proof objects to the code generator instead of simply eliding them from the programs as they often enable more aggressive program optimisations. To demonstrate the potential of this approach, we look at code generation for an array DSL embedded in Coq. We use the SaC toolchain for generating the eventually executed code and we demonstrate how the Coq-level proofs can be leveraged for helping the SaC compiler to generate more runtime efficient code.

ACM Reference Format:

Artjoms Šinkarovs, Sven-Bodo Scholz, and Hans-Nikolai Vießmann. 2019. Extracting Verified Constraints from Coq-embedded Array DSLs – Extended Abstract – . In *Proceedings of International Symposium on Implementation and Application of Functional Languages (IFL'18)*. ACM, New York, NY, USA, 4 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnn>

1 INTRODUCTION

Array languages come with a number of static properties, commonly known as array algebras. The algebraic laws give rise to the number of program optimisations, and are widely used in the optimising array compilers like APL, SaC, SISAL, and others. For example:

$$\text{shape}(\text{flatten } a) = \text{prod}(\text{shape } a)$$

says that the shape of the flattened array is a product of the shape components of that array. If we use this property as a universal rewrite rule, we may significantly simplify the program, as the right-hand side avoids evaluating `flatten a`.

A number of such properties are typically captured by the type systems of array languages. They are particularly good at tracking relations between the shapes of function arguments and the shape

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

IFL'18, August 2019, Lowell, MA, USA

© 2019 Association for Computing Machinery.
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00
<https://doi.org/10.1145/nnnnnnnn.nnnnnnn>

of the result. This is important as usually, the more precise shape information we have, the better code we can generate.

However, for most of the existing array languages, it will be possible to define a function where argument-result shape relations would be inexpressible in their type systems. Specifically, if a language allows for shape-polymorphic expressions on arrays. For example, it is usually possible to capture the fact that element-wise addition of two arrays preserves the shape of its arguments. On the other hand capturing the fact that transposing an array reverses its shape vector is rarely possible.

In the dependently typed languages, it is surely possible. Moreover, we can prove that a certain array expression has a property of interest. The main difficulty with these languages is that they struggle to generate efficient code.

To solve this problem, in Coq there exists a concept of program extraction. It is possible to translate well-typed terms into a functional languages like Haskell or OCaml. For example:

`Definition dbl x := x + x.`

is extracted in the following OCaml code:

```
type nat =
| O
| S of nat

let rec add n m = match n with
| O -> m
| S p -> S (add p m)

let dbl x = add x x
```

One of the important optimisations that happens during the extraction is erasure of computationally irrelevant proof objects. Consider the following example:

```
Require Import List.
Inductive szlst (n: nat) : Type :=
SzLst : forall l, (@length nat l = n) → szlst n.
```

We define an inductive family `szlst` for sized list of length `n`. It has a single constructor that takes an argument `l` which is a list and a proof that the length of `l` is `n`. Consider now a trivial `prep_one` function that prepends `1` to the given sized list of natural numbers:

```
Lemma h : forall {t l n el},
@length t l = n → @length t (cons el l) = (S n).
intros. simpl. simpl in H. rewrite H. auto.
Qed.
```

```
Definition prep_one {n} (szl: szlst n) : szlst (S n) :=
match szl with
```

```
| SzLst _ 1 pf ⇒ SzLst (S n) (cons 1 1) (h pf)
end.
```

We use a helper lemma *h* to construct a proof that the resulting list has one more element in it. After extraction we get the following code:

```
type szlst =
  nat list

let prep_one _ szl =
  Cons ((S O), szl)
```

Note that all the proof objects that contain information about the length of the argument lists are gone. Computationally, this is exactly what we want. At the same time, in case we want to perform further optimisations on the extracted code, such an information would have been very useful. For example, we could derive the following rewrite rule:

```
szl_length (prep_one a) = S (sz_length a)
```

where *szl_length* returns the length of the sized list.

Now we bring the same idea in the context of the functional language for multi-dimensional arrays. We define a shallowly-embedded DSL in Coq and we propose extraction technique that translates terms of our DSL into functional array language SaC and we turn proven shape relations into rewrite rules that can be used for further program transformations. Since many years, SaC exposed mechanisms to handle user-defined constraint that make it possible for a given function to define shape relations between the arguments and the function result. Since the function bodies are generated from the certified Coq terms, the validity of these constraints has been proven already, therefore they can be actually trusted.

2 ARRAY DSL IN COQ

We define a DSL for multi-dimensional arrays that supports shape-invariant array expressions that are often found in APL or SaC. As we have shown before [3], the following DSL is as expressive as SaC:

$e ::= c$	(constants)
x	(variables)
$\lambda x.e$	(abstractions)
$e e$	(applications)
<i>if e then e else e</i>	(conditionals)
<i>letrec x = e in e</i>	(recursive let)
$e + e, \dots$	(built-in binary)
<i>imap e e</i>	(index map)
$e.e$	(selections)
$ e $	(shape)
$c ::= d$	
$[d, \dots]$	(constant array)
$d ::= 0, 1, \dots$	(natural numbers)

which is a minimalistic functional language appended with three array-specific constructs: the array constructor *imap*, the selection operator *(.)* and the shape operator $|_|$.

As it can be seen, embedding such a language into Coq is straightforward, as all the non-array constructs can be used directly as they defined by Gallina. All we need to do is to describe a type for multi-dimensional arrays and define three array operators as built-in functions.

We start with defining array types. We will model arrays with elements of type t as functions from some index-space to t . In the context of this paper we assume that all the index-spaces are hyper-rectangular, but we will keep the formalism extendable to other index-space types.

The type of an index into d -dimensional array is a vector of d elements which we modeled as a function of type $\text{Fin.t } d \rightarrow \text{nat}$. Therefore the type of indices into d -dimensional arrays is defined as:

```
Inductive index (d:nat) : Type :=
  IIdx : (\text{Fin.t } d \rightarrow \text{nat}) \rightarrow index d.
```

Rectangular d -dimensional index-spaces can be always described with a tuple of d natural numbers. This tuple gives an upper bound, defining an index-space with the valid indices being from $\langle 0, \dots, 0 \rangle$ up to the upper bound. The upper bound is modeled as a function of type $\text{Fin.t } d \rightarrow \text{nat}$. We define an inductive family for index-spaces with a single constructor *Rect* that describes rectangular index-spaces:

```
Inductive ispace (d:nat) : Type :=
  Rect : (\text{Fin.t } d \rightarrow \text{nat}) \rightarrow ispace d.
```

We define a type that captures that the index is within the rectangular index-space. This happens when all the elements of the index-tuple are less than the corresponding elements of the shape tuple:

```
Inductive rect_bound (d:nat)
  : (index d) \rightarrow (ispace d) \rightarrow Type :=
  RectBound : \forall idx_f sh_f,
    (\forall i, idx_f i < sh_f i)
    \rightarrow rect_bound d (IIdx d idx_f) (Rect d sh_f).
```

Finally, the array type for the given dimensionality d , element type t and index-space *is* is given as a function from index and a proof that the index is within *is* into t :

```
Inductive array (d:nat) (t:Type) (is:ispace d) : Type :=
  RectArray : \forall f:(\forall (idx: index d),
    (rect_bound d idx is) \rightarrow t),
    array d t is.
```

Embedding. With the given definitions of the array type, here are the missing operations from our array language:

```
Definition shape {d t f_shp}
  (a: array d t (Rect d f_shp))
  : ispace d
  := Rect d f_shp.
```

```
Definition imap {d t}
  (isp: ispace d)
  (f: \forall idx: index d, rect_bound d idx isp \rightarrow t)
  : array d t is\@p
  := RectArray d t isp f.
```

```

Definition asel {d t sh_a}
  (a: array d t (Rect d sh_a))
  (idx: index d)
  (pf: rect_bound d idx (Rect d sh_a))
: t
:= match a with
| RectArray _ _ _ f => f idx pf
end.

```

As it can be seen, all the operations are wrapper around either inductive constructor or eliminator.

3 SAC CONSTRAINTS

The SaC compiler relies on the validity of various constraints and properties when it comes to achieving the highest levels of code optimisation. Examples are range constraints for guaranteeing safe array accesses or for enabling advanced optimisations such as *with-loop-folding*. Further examples are relations between function compositions such as associativity or distributivity which enable the compiler to perform more advanced partial applications.

Currently, the SaC compiler injects the constraints and properties needed as *symbiotic expressions* [1, 2]. Whenever the type system or the built-in analyses and symbolic evaluation rules are capable of proving these constraints the corresponding optimisations succeed. Otherwise, the optimisations are not applied or some runtime check is inserted which guards two code variants.

While this approach works reasonably well when the compiler can specialise the code quite aggressively, it typically is rather ineffective when the compiled can not be specialised sufficiently due to the lack of static knowledge. Here the Coq-DSL approach proposed in this paper can be most effective. As an added benefit, the ability to prove constraints and properties without specialisation enables high levels of optimisation for separately compiled modules as well. The current setup that relies on heavy specialisation re-optimises all modules that are being imported and, thus, effectively performs whole-world optimisations for every executable.

As an example, let us consider array transposition:

```

Definition trans {d t sh}
  (a : array d t (Rect d sh))
  : array d t (Rect d (rev sh)) :=
RectArray d t (Rect d (rev sh))
(fun idx pf =>
  match idx as i return (rect_bound d i (Rect d (rev sh)) -> t)
  with
  | Idx _ idx_f =>
    fun pf => asel a (Idx d (rev idx_f)) (rev_thm pf)
  end pf).

```

We use a helper function rev that reverses a tuple; and we use the rev_thm theorem that converts the proof to satisfy selection function:

```
Definition rev {d} (f: Fin.t d -> nat): Fin.t d -> nat.
```

```
Theorem rev_thm :
forall {d} {sh idx_f},
rect_bound d (Idx d idx_f) (Rect d (rev sh)) ->
```

```
rect_bound d (Idx d (rev idx_f)) (Rect d sh).
```

From the type of the trans function, we know that the shape of a transpose is a reverse of the shape of the result. However, after the extraction this information is gone.

Due to design of the type system of SaC, the type signature of the non-specialised version of transpose does not maintain any information about the way shapes of arguments and results are related. All that is exposed to the calling site of the functions rev and trans is:

```

int[.] rev (int[.]);
int[*] trans (int[*] a);

```

If we consider a call site of the form

```

... b = trans (a);
c = with /* ... */ : genarray (shape (b), 0);

```

even if we do have local knowledge of the shape of the array a, no knowledge of the shape b can be inferred without specialising and inlining the definitions of trans and rev.

Therefore, evaluation of the expression in the last assignment, in general case has to assume the most generic case, due to the lack of shape. If we would derive the following properties from Coq types:

```

// shape (rev (a)) = shape (a)
// shape (trans (a)) = rev (shape (a))

```

we could rewrite our last assignment as:

```

c = with /* ... */
      : genarray (rev (shape (a)), 0);

```

which offers much better chances to further optimisations.

Finally, not only can we extract rewrite rules from rich function types, but we could also use theorems about objects with simple types. Consider the following example with a proven property about the fibonacci function with type `nat -> nat`:

```

Fixpoint fib n :=
match n with
| 0 => 1
| S n' => match n' with
| 0 => 1
| S n'' => fib n' + fib n''
end
end.

```

```
Lemma fib_prop: forall i, fib i >= i.
```

Proof.

intros.

induction i.

simpl. auto.

simpl.

destruct i. auto.

assert (l: forall a b c, a >= b -> c > 0 -> a + c >= S b).

intros.

omega.

apply l.

apply IHi.

assert (forall i, fib i > 0).

induction i0.

auto.

```

simpl.
destruct i0. auto.
assert (l': forall a b, a > 0 → a + b > 0).
  intros.
  omega.
apply l'.
apply IHi0.
apply H.
Qed.

```

Proving such an expression automatically as a part of compiler optimisations is quite challenging. However, turning the `fib_prop` lemma into a rewrite rule is straight-forward, given that we are extracting `fib` already. Such a knowledge may be useful when dealing with fibonacci heaps.

4 TOWARDS EXTRACTION

Extracting programs to SaC follows a very similar procedure as existing extractions into OCaml or Haskell. Additionally we will use the knowledge that not all the Coq terms are valid within our DSL – only if they are follow the structure of our DSL definition.

In SaC there is no distinction between indices, shapes and normal arrays. Indices and shapes are just 1-dimensional arrays. Therefore our translation turns functions of Coq type `Fin.t n → nat` into SaC type `int[.]`. This is straight-forward to do using an *imap* array constructor:

$$\llbracket \text{fun } f \{n\} (x:\text{Fin.t } n) \Rightarrow e \rrbracket_E = \text{imap } [n] \lambda iv. \llbracket e \rrbracket_{x \mapsto iv, E}$$

Applications of such functions are extracted as selections at the given index:

$$\llbracket f \ e \rrbracket_E = \llbracket f \rrbracket_E \cdot [\llbracket e \rrbracket_E]$$

The values of `Fin.t n` are turned into regular integers, similarly to the way `nat` type is treated. Finally, for meaningful operations with indices, basic operations on `Fin.t n` values have to be translated into corresponding operations on `ints`.

When a function gets a proof object as an argument, by the construction of our DSL we know that it can be only used in two positions: either when constructing an array or when making selections. Therefore, after we verify this, a proof-object argument can

be substituted with a dummy argument throughout the body of the function, as any other elimination of proof objects would require functions outside of the proposed DSL. Other than that, translation of functions, lets, applications and conditionals is straight-forward.

Selections remove the proof object and turn the index into 1-dimensional array:

$$\llbracket \text{asel } a \ idx \ pf \rrbracket_E = \llbracket a \rrbracket_E \cdot \llbracket idx \rrbracket_E$$

The *imap* removes the proof object and turns the shape into 1-dimensional array:

$$\llbracket \text{imap } s (\text{fun } idx \ pf \Rightarrow e) \rrbracket_E = \text{imap } \llbracket s \rrbracket_E \lambda iv. \llbracket e \rrbracket_{idx \mapsto iv}$$

The *shape* operation translates one-to-one into the built-in operation:

$$\llbracket \text{shape } e \rrbracket_E = |\llbracket e \rrbracket_E|$$

As for deriving the rewrite rules, for translated functions involving arguments of array types, we observe shape and dimension equalities and generate one rewrite rule for each identified equality. As for the lemmas outside of function types, we require them to contain a single equality without pre-conditions. For example:

$$\text{Lemma } f_prop : \text{forall } x, \text{shape}(f x) = \text{shape } x.$$

In the full paper we elaborate on the way rewrite rules are defined in SaC. We will present formal description of the extraction and the implementation. Finally we will measure runtime improvements of a few classical array examples.

REFERENCES

- [1] Robert Bernecker, Stephan Herhut, and Sven-Bodo Scholz. 2010. Symbiotic Expressions. In *Implementation and Application of Functional Languages, 21st International Symposium, IFL 2009, South Orange, NJ, USA (Lecture Notes in Computer Science)*, Marco T. Morazan and Sven-Bodo Scholz (Eds.). Springer, 107–126. https://doi.org/10.1007/978-3-642-16478-1_7
- [2] Stephan Herhut, Sven-Bodo Scholz, Robert Bernecker, Clemens Grelck, and Kai Trojahnner. 2008. From Contracts Towards Dependent Types: Proofs by Partial Evaluation. In *19th International Symposium on Implementation and Application of Functional Languages (IFL '07), Freiburg, Germany, Revised Selected Papers (Lecture Notes in Computer Science)*, Olaf Chitil, Zoltan Horváth, and Viktória Zsók (Eds.), Vol. 5083. Springer, 254–273. https://doi.org/10.1007/978-3-540-85373-2_15
- [3] Artjoms Åäinkaroffs and Sven-Bodo Scholz. 2017. A Lambda Calculus for Transfinite Arrays: Unifying Arrays and Streams. *CoRR* abs/1710.03832 (2017). arXiv:1710.03832 <http://arxiv.org/abs/1710.03832>

A DSL embedded in Rust

Draft Paper

Kyle Headley

University of Colorado Boulder

kyle.headley@colorado.edu

ABSTRACT

In this paper we present a DSL parsed by the Rust parser (and macros), type checked by the Rust type checker, and run with compiled Rust code. We focus on type checking, where we use two "languages" provided by Rust: functional-style meta-programming, and a form of declarative logic mirroring type rules. Both are implemented with trait resolution, part of Rust's generics system. Rust's macros are hygienic and can rewrite tokens into expressions, allowing a wide range of syntax available to DSLs. The compiled code is standard, but included to demonstrate the full range of language definition available in Rust.

ACM Reference Format:

Kyle Headley. 2019. A DSL embedded in Rust : Draft Paper . In *Proceedings of International Symposium on Implementation and Application of Functional Languages (IFL'18)*. ACM, New York, NY, USA, 6 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

The Rust language reached version 1.0 in mid-2015, bringing together high-performance, thread-safety, and a minimal runtime system. We ignore those features in this paper, concentrating instead on the macros and trait-based generics, the meta language of Rust. Both of these are expressive enough to be used as their own general-purpose programming languages. This paper is an exploration of these languages in the context of creating a lambda calculus DSL. We briefly introduce macros in Section 2, then move on to an intro of the meta language, which we will call "TraitLang" for the remainder of this paper. We describe our usage of TraitLang in Section 3.

Rust has been enhanced with procedural macros, but we will be using the standard ones. These are a name and a list of rewrite rules; the first rule whose pattern matches the syntax the macro was called on is used to expand it into new syntax. The new syntax can include macro invocations (but not definitions), allowing recursive calls. We use this to mirror BNF grammars, with a different macro for each component of the grammar. Macros are commonly used to transform code snippets, but they have a mode that deals with arbitrary tokens.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

IFL'18, August 2019, Lowell, MA, USA

© 2019 Association for Computing Machinery.
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

TraitLang makes use of Rust generics, but is shaped by traits, a feature similar to typeclasses in other functional languages. Traits normally provide an interface for using data of the type implementing them, but are still useful as a classifier for type-level values. To avoid confusion when discussing these type-level values, we refer to one as a "struct", the keyword used when defining a type in Rust. There are no type-level types, so the word "trait" is sufficient for our purposes. We introduce the basic constructions in Section 3.1

TraitLang is a lazy, untyped language with some features similar to both logic and functional languages. It is declarative and order of declaration doesn't matter, as all items are fully recursive. It is interpreted by the Rust trait resolution algorithms, which are expected to be enhanced in the near future. In this paper we use the original semantics from version 1.0, though the Rust team almost never introduces breaking changes (until the next major version).

TraitLang is pure, since Rust's type-level items do not have access to the object language at all. It is not even possible to get output from a TraitLang program. In this paper, we rely on the type-checker to verify that our programs are well-formed, but demonstrating correct output is beyond our scope, requiring object-level display methods. Because TraitLang is lazy, the well-formed check also requires that type aliases be used. We assume a "fn main() { let x:TypeAlias1; ... }" with each alias used at least once.

Like types in a common language, traits classify structs, but unlike types, a struct can "implement" an unlimited number of traits. Each of these traits may contain associated types specific to its implementation by a struct. This implementation therefore acts as a mapping from one struct to another, one of the ways to define a function. However, to allow first-class functions, we prefer a slightly deeper embedding, which is described in Section 3.2.

Haskell. Type-level programming is available in many other languages. Haskell in particular provides many of the features we make use of in this paper. Haskell users may recognize the semantics of our type-level language, though the syntax is very different.

The next part of the paper walks through our implementation of the lambda calculus with addition as a DSL. We describe parsing in Section 4, type checking in Section 5, and running the program in Section 6. Finally, we look at some related work and conclude with a link to a working version of the code.

2 RUST MACROS

Each Rust macro is a list of rewriting rules, from a pattern matcher to a template for expansion. The matcher includes literals, pattern variables, and repeaters. Pattern variables are prefixed with a \$ and include a "fragment specifier". We will mostly be using token trees (tt), which macro invocations are initially parsed into. Token trees are either a single token, or a parenthesized ((),{},[]) sequence of tokens. We also make use of expr, which signals the Rust parser

```
macro_rules! expr {
    ($0) => (Num(Zero));
    ($a:tt $p:tt) => (App(expr![${$a}],expr![${$p}])); 
    ({^$e:expr}) => ($e);
    ((${$e:tt}+) ) => (expr![${$e}+]);
    ($n:tt + ${$ns:tt}+)
    => (Plus(expr![${$n}],expr![${$ns}+])));
}
```

Figure 1: Selected macro rules (out of order)

to fully parse the match as a Rust expression. Repeaters `$(...)`,⁺ match multiple instances of their inner pattern, with optional separators. The macro expander checks matchers in order, and expands the first one that matches the input. Some selected rules are in Figure 1, explained below.

Later we will use the `expr!` macro to parse syntax into an AST. For now, we use selected portions shown in Figure 1 to introduce Rust macros. The first line contains only a literal in the matcher, to transform a number into its AST representation. (Our AST distinguishes raw natural numbers from syntax.) There are no pattern variables, so the `0` must be matched exactly. The second line shows two pattern variables specified as token trees. This pattern matches any two tokens, and the expander recursively invokes the `expr!` macro on each, placing them within an `App` node in our AST. The third line is used to insert pre-created expressions into our AST. The pattern variable is parsed and used directly. The other tokens are literals and must be matched exactly. Using one of the forms of parenthesis to surround the match allows it to be treated as a token tree before reaching this rule.

The final two lines of Figure 1 show off the macro repeaters. The first is a minimal repeater surrounded by parentheses, for parsing parenthesized expressions. The contents of the parentheses are copied into a recursive invocation. The final line is a more complex version of the same principal, used to put everything after the first `+` into the second section of the `Plus` AST node (after a recursive call). If the matched pattern contained multiple `+`'s, they would evaluate left to right. Our DSL doesn't need to deal with order of operations, but one that did would need a more complex matcher.

3 TRAITLANG

This section describes the use of TraitLang as a functional language. The syntax and programming style are very different from traditional languages, so we take some care in walking through a series of progressively more complex examples. TraitLang is interesting on its own, so we go a bit beyond what is needed to implement our DSL, describing a method for using first-class functions. Type checking mostly makes use of a logical style, but does use supporting functions. The membership function for contexts (described in Section 5.1) is rather verbose, since TraitLang is not well-suited for functions with multiple branches.

3.1 A Hidden Language

When we ignore Rust's main language and focus on the trait language, we are left with four items: declaration of a trait, declaration of a struct, implementation of a trait for a struct, and declaring a type alias, which functions like a let-binding. The basic syntax of

```
trait Nat {}

struct Zero;
impl Nat for Zero {}

struct Succ<N>(N);
impl<N:Nat> Nat for Succ<N> {}

type One = Succ<Zero>;
```

Figure 2: Declaration of Natural numbers in Rust's type-level language

these items in shown in Figure 2, which gives the standard definition of natural numbers. Here we define `Nat` as a trait, which works well at first, but is not sophisticated enough for a formal definition. Structs may implement multiple traits, allowing a later "crate" (Rust package) to implement e.g. trait `Real` for the same zero. We return to this issue later.

Figure 2 continues by declaring a struct called `zero` and implementing `Nat` for it. This is the simplest form of the declarations. More complex is `Succ`, which requires a parameter when the struct is used. In this case `N` may be any other struct, including a recursive `Succ` (though infinite sequences cannot be defined). The second to last line is read "For all `N` such that `N` implements `Nat`, implement `Nat` for `Succ<N>`". Using this definition, the compiler will not give an error when using e.g. `Succ<Red>` (assuming a struct `Red` has been declared), but it would not implement `Nat`. We could have given a trait bound when declaring `Succ`, that is, "struct `Succ<N:Nat>(N)`". Doing so would cause a compiler error on use of `Succ<Red>`. We can use a struct by creating an alias like in the last line. The struct must be concrete, with no type variables.

There are a few syntactic peculiarities in Figure 2. Trait definitions end in curly braces, which are usually filled with object-level function definitions. We will add associated types here later. Formal type parameters, which can appear in any of the four syntactic items, are placed between angle braces and separated by commas. Each one may be required to implement any number of traits, placed after a colon and separated by a `+`. Struct definitions must include each type parameter in parens, which is required for the object-level language, but we will not use it anywhere else. In the second to last line of Figure 2, the formal type parameters are after the `impl`, and their use is after the `Succ`. Usage does not include trait bounds.

3.2 A Functional Language

The full power of a functional language requires having functions. Figure 3 presents the simplest form available, using a trait as a mapping. Like defining `Nat` as a trait above, this form is simpler but limited, and we mainly use it for DSL meta-functions. We describe the syntax and semantics of this form first before moving on to one that allows first-class functions. In the figure, we define addition and subtraction by one.

Figure 3 introduces bounds for trait declarations, associated types, and how to access them. The first line declares a trait `AddOne` that requires any struct it's implemented for to also implement `Nat`. It includes a single associated type named `Result` that must also implement `Nat`. The implementation on the next line shows

```

trait AddOne : Nat { type Result:Nat; }
impl<N:Nat> AddOne for N { type Result = Succ<N>; }

trait SubOne : Nat { type Result:Nat; }
impl<N:Nat> SubOne for Succ<N> { type Result = N; }

type Two = <One as AddOne>::Result;

```

Figure 3: Using traits as mappings

```

trait Func2<A,B> { type Result; }

struct Add;
impl<N:Nat> Func2<Zero,N> for Add { type Result = N; }
impl<N1,N2> Func2<Succ<N1>,N2> for Add where
    N1:Nat, N2:Nat,
    Add : Func2<N1,N2>
{ type Result = Succ<<Add as Func2<N1,N2>>::Result>; }

type Three = <Add as Func2<One,Two>>::Result;

```

Figure 4: Structs that can be used as functions

off the power of variables, implementing `AddOne` for every `Nat`, and providing an associated type dependent upon it. `SubOne` is similar, but note that it is not implemented for every `N`. Every associated type must be defined in order to implement a trait, but traits need not be implemented for every struct. This can be useful to ensure that suitable values are provided to computations. If appropriate, we could implement `SubOne` for every `Nat` by including the line “`impl SubOne for Zero { type Result = Zero; }`”

The last line in Figure 3 uses the unfortunate syntax for accessing an associated type. Both of the traits here have the same associated type name, so we must disambiguate by naming the struct, the trait implemented on the struct, and the associated type of that trait. The syntax is slightly better in Figure 4 where we use structs as functions.

Rust generics use type variables, but there are no trait variables. If we were to continue to use traits as we did above, we would run into problems with generics and first-class functions in our type-level language. Below, we use traits to represent higher-level concepts. For example, there is a trait to mean that a struct is a function, rather than using a trait as a function. Figure 4 declares a trait representing a function of two variables. It then declares a struct `Add` and implements the inductive algorithm for adding two numbers.

Figure 4 introduces trait parameters which are similar to struct parameters. It also introduces the “where” clause, which can be used to add arbitrary requirements to any item. Here, the inductive case for defining `Add` requires `Add` be defined on a smaller structure. Since it is guaranteed by the where clause, we can look up its associated type to use in the definition of the result. Since `Add` is a struct, it can be passed as a parameter to functions just like `One` and `Two` were in the last line. Since the function parameters are declared on the trait, it can be called by any code with a where clause recognizing it as a function. There will be an example later.

```

trait Typed { type Type; }

struct Natural;
impl Typed for Zero { type Type = Natural; }
impl<N:Typed<Type=Natural>> Typed for Succ<N>
{ type Type = Natural; }

trait Func1<A:Typed> { type Result:Typed; }

struct Next;
impl<N:Type=Natural> Func1<N> for Next
{ type Result = Succ<N>; }

type Four = <Next as Func1<Three>>::Result;

```

Figure 5: The typing trait, allowing us to emulate a standard type system

```

struct Apply;
impl<A,B,R> Func2<A,B> for Apply where
    B: Typed, R: Typed,
    A: Func1<B,Result=R>
{ type Result = R; }

type Five = <Apply as Func2<Next,Four>>::Result;

```

Figure 6: A function that takes another as an argument

3.3 A Constraint Language

In this section we describe the final piece of syntax that will allow us to emulate a standard type system in our language. So far, we have been using traits as if they were types. But structs can implement multiple traits, so our functions can be applied to multiple “types”. In order to have one type per struct, we need a mapping. Figure 5 demonstrates using a trait to declare types.

Figure 5 repeats functionality defined above, but in our final form with types. The fourth line can be read: “For all `N` of type `Natural`, the type of `Succ<N>` is `Natural`”. Note that we now have multiple levels of constraint, since we do not need to constrain the associated type. `Func1` requires its argument and result to be `Typed`, but doesn’t require a specific type. Applying it to `Three` in the last line works the same as our prior example, but now the compiler is checking the associated type (required for arguments of `Next`) as well as the trait of `Three`.

We now know all the features we need to use TraitLang as a general-purpose language. Our language is untyped, but uses traits both to add and remove capabilities. Functions were added from mappings in traits, and the ability for `succ` to take any parameter was removed by a constraint. Structs can be used as any value in the language, even when that value is acting as a different feature, like a type or a function. For example, Figure 6 shows a use of first-class functions. Using structs as types allows type-based operations, as in dependent types. We also see a syntactic optimization in the third line. We can constrain the associated type as well as type parameters. This in effect “binds” `R` to the result of `A` applied to `B`, allowing us to use it rather than the longer form used in Figure 4.

```

e :=           expressions
  (e)          parentheses
  n            number
  v            variable
  lam (v:t) e abstraction
  lam (v1:t1)(v2:t2)... e
                multiple abstraction
  e1 + e2      addition
  e1 e2        application
  e1 e2 e3 ... multiple application
t :=           types
  (t)          parentheses
  Number       base type
  t1 -> t2 -> ... arrow type

```

Figure 7: The grammar for our DSL

4 PARSING OUR DSL

The DSL we're implementing is the simply-typed lambda calculus with numbers and addition. Our grammar is standard and shown in Figure 7. Using macro rules for parsing means that we can follow our grammar very closely. We create one macro for each syntax class, and one rule for each syntax form. Our only deviations are in representing numbers and variables, and adding an injection point for easier composition, as described in Section 2. The full parser is shown in Figure 8.

Representing numbers and variables is a pain point of this method. Since we're working in TraitLang, we don't have access to any runtime functionality, only logic and induction. Integers and arithmetic are not available, so we use inductively-defined natural numbers (nats). The parser needs to map number literals to nats, so we need a rule for each number. Variables are available as additional structs, which would still add lines to the code. Also, we need to abstract over variables in our type checking later, but Rust does not give us an easy way to check both equality and inequality. To overcome this, we use nats as variables as well, with AST nodes that distinguish them from numbers.

Many of the rules in Figure 8 were shown previously or are similar to those. We describe some additional complexity here. The multi variable lambda rule has a nested repeater. The inner matches all the var and type tokens, and the outer matches the parenthesized groups. Before it is the first variable and type, which are used to create the lambda AST node. The repeater represents additional variables, which are used to create a nested lambda node with a recursive call. The lambda nesting pattern is convenient in this way, but the application nesting pattern is not. It is the reason we created the injection rule above. The nesting of applications must be as deep initially as the number of parameters, which we don't know. So we create the first AST node and pass it unchanged into the recursive call. The type rules are simple because arrow nests the lame way that lambda does.

5 TYPE CHECKING OUR DSL

This section describes the code for our type checker, divided into two parts. The first deals with functions for context lookup, and is done in the functional style introduced in Section 3. The next part handles static checks of our AST, and are written in a more logical

```

macro_rules! expr {
    ($($e:tt)+) => (expr![ $($e)+ ]);
    ($^$e:expr) => ($e);
    (0) => (Num(Zero));
    (1) => (Num(Succ(Zero)));
    ...
    (x) => (Var(Zero));
    (y) => (Var(Succ(Zero)));
    ...
    (lam ($x:ident : $($t:tt)+) $($ts:tt)+) ($e:tt)+)
        => (Lam(
            expr![ $x ],
            typ![ $($t)+ ],
            expr![ lam $($ts)+ ]+$e+ ]
        ));
    (lam ($x:ident : $($t:tt)+) $($e:tt)+) => (Lam(
        expr![ $x ],
        typ![ $($t)+ ],
        expr![ $($e)+ ]
    )));
    ($n:tt + $($ns:tt)+) => (Plus(expr![ $n ],expr![ $($ns)+ ]));
    ($a:tt $p:tt) => (App(expr![ $a ],expr![ $p ]));
    ($a:tt $p:tt $($ps:tt)+)
        => (expr![ ^App(expr![ $a ],expr![ $p ]) } $($ps),+]);
}
macro_rules! typ {
    ($($ts:tt)+) => (typ![ $ts ]);
    (N) => (Number);
    ($t:tt -> $($ts:tt)+)
        => (Arrow(typ![ $t ],typ![ $($ts)+ ]));
}

```

Figure 8: The parser for our DSL

```

trait NatEq<N> { type Eq; }
impl NatEq<Zero> for Zero { type Eq=True; }
impl<N> NatEq<Succ<N>> for Zero { type Eq=False; }
impl<N> NatEq<Zero> for Succ<N> { type Eq=False; }
impl<N1,N2,E> NatEq<Succ<N1>> for Succ<N2> where
    N2: NatEq<N1,Eq=E>
{ type Eq=E; }

```

Figure 9: Equality function for natural numbers

style. This is valuable because, like for parsing, our code can mirror the rules from the notation of the theory.

5.1 Supporting functions

Context lookup appears in type checking rules, but often as a function over the data structure for simplicity. We mirror that here, but still need a full description of the algorithm. Context lookup involves comparing variables to find our target. As seen above, we have implemented our variables as natural numbers, so we need an equality function for them. This is shown in Figure 9. It follows the method from Figure 3, since we don't make use of first-class functions. There are three base cases for equality with zero, followed by an inductive case. This code is rather elegant, but the membership function that makes use of it is not.

The context membership function (called `contains` in code) requires a data structure and two branch points, one for checking

```

struct EmptyCtx;
struct TypeCtx<Id, Typ, Next>(Id, Typ, Next);

trait Contains<Id> { type Result; }
impl<N> Contains<N> for EmptyCtx
{ type Result=None; }

impl<Check, First, Typ, Next, Eq, R>
Contains<Check> for TypeCtx<First, Typ, Next> where
    Check: NatEq<First, Eq=Eq>,
    Next: Contains2<Eq, Typ, Check, Result=R>,
{ type Result=R; }

trait Contains2<Eq, Map, Check> { type Result; }

impl<Map, C, Ctx> Contains2<True, Map, C> for Ctx
{ type Result=Some<Map>; }

impl<Map, C> Contains2<False, Map, C> for EmptyCtx
{ type Result=None; }

impl<Check, First, T, Typ, Next, Eq, R>
Contains2<False, T, Check> for TypeCtx<First, Typ, Next> where
    Check: NatEq<First, Eq=Eq>,
    Next: Contains2<Eq, Typ, Check, Result=R>,
{ type Result=R; }

```

Figure 10: Membership function for contexts

if we're reached the end of the list, and one for checking if we've reached our target variable. TraitLang only allows one branch point and one return value per branch. To get around this, we create two functions, the first one passing the results of its check to the second as parameters. The second then branches once based on all the information. Even for this simple function, the code is difficult to read. We work through it below.

Figure 10 shows the code of the context membership function. The first two lines are the data structure, implemented like a linked list. It can be empty or contain the natural number id of a variable, a type, and the next node. The `Contains` function takes an id and returns an optional value, in the case of calling it on an empty context, it returns `None`. When called on a non-empty context, `Contains` checks for equality with the target, passes that result, along with the type, as parameters to `Contains2` called on the rest of the context and returns the result of `Contains2`. `Contains2` has enough information to chose one of the three end-points of the algorithm. If the prior equality check was true, it returns the prior type (called `map` in the code) regardless of the rest of the context. If the check was false and the rest of the context is empty, it returns `None`. If there is more context to process, it does an equality check on the next value and calls itself recursively the same way `Contains` did.

5.2 Type checking

Type checking starts by checking that the AST is well-formed. The code is in Figure 11. Most rules define the syntax that we're using as well-formed if its sub-syntax is well-formed. The exception is the `Lam` case, which requires a variable as its first item. There are different traits used for different parts of the syntax, like `WFNat` and `WFType`, to make sure they are used in the proper places. There is

```

trait WFNat {}
impl WFNat for Zero {}
impl<N: WFNat> WFNat for Succ<N> {}

trait WFType {}

struct Number;
impl WFType for Number {}

struct Arrow<T1, T2>(T1, T2);
impl<T1: WFType, T2: WFType> WFType for Arrow<T1, T2> {}

trait Expr {}

struct Num<N>(N);
impl<N: WFNat> Expr for Num<N> {}

struct Plus<N1, N2>(N1, N2);
impl<N1: Expr, N2: Expr> Expr for Plus<N1, N2> {}

struct Var<N>(N);
impl<N: WFNat> Expr for Var<N> {}

struct Lam<V, T, E>(V, T, E);
impl<N: WFNat, T: WFType, E: Expr> Expr for Lam<Var<N>, T, E> {}

struct App<E1, E2>(E1, E2);
impl<E1: Expr, E2: Expr> Expr for App<E1, E2> {}

```

Figure 11: Well-formedness checking logic

little complexity to the code, it mostly tags some constructions as appropriate.

Our type checking code in Figure 12 is among the most simple and elegant in this paper, because we are able to directly mirror the type checking rules. We use a trait called `Typed` parametrized by a context. Premises are found in the “where” clauses with the syntax form preceding them. The resulting type is an associated type, to make sure that there is only one type per value. Otherwise, the rules are direct translations of the typing rules for the lambda calculus. For example, the last rule, `App`, requires that the first expression (E_1) be an `Arrow` type (from T_1 to T_2) in the current context (ctx), and the second expression (E_2) be of the type at the front of the arrow (T_1), also in the current context. The type of the `App` expression is the type of the end of the arrow (T_2).

6 RUNNING OUR DSL

Running our code may require an initial conversion, but would otherwise be standard. When defining a struct in Rust, it generates a singleton constructor (parametrized as appropriate) with the same name. This is what we've been using in our AST nodes, while the struct itself is used in our implementation of traits. From the runtime perspective, each AST node is a different type, which can make coding up the evaluation difficult. A conversion to an AST using tagged variants of types (Rust's `enum`) would simplify the eval code.

7 RELATED WORK

A similar project is “turnstile” [1] for the Racket language. The authors similarly take advantage of a compile-time algorithm to do type checking. In their case, they use Racket's macro expander,

```

trait Typed<Ctx> { type T; }

impl<N,Ctx> Typed<Ctx> for Num<N> { type T=Number; }

impl<N1,N2,Ctx> Typed<Ctx> for Plus<N1,N2> where
    N1:Typed<Ctx,T=Number>,
    N2:Typed<Ctx,T=Number>,
{ type T=Number; }

impl<N,Ctx,T> Typed<Ctx> for Var<N> where
    Ctx:Contains<N,Result=Some<T>>
{ type T=T; }

impl<Ctx,N,T1,T2,E> Typed<Ctx> for Lam<Var<N>,T1,E> where
    E:Typed<TypeCtx<N,T1,Ctx>,T=T2>,
{ type T=Arrow<T1,T2>; }

impl<Ctx,E1,E2,T1,T2> Typed<Ctx> for App<E1,E2> where
    E1:Typed<Ctx,T=Arrow<T1,T2>>,
    E2:Typed<Ctx,T=T1>
{ type T=T2; }

```

Figure 12: Type checking rules for our DSL

adding typing annotations to the syntax objects it creates. Rust traits provide a declarative way to add meta-data to types, allowing much simpler use, and the ability to follow the typing rules more directly. On the other hand, Racket has more advanced capability in its macro system, allowing a layer of abstraction that lets the user follow typing rules as well. Racket also provides a mechanism for generating useful error messages.

8 CONCLUSION

We have shown the implementation of a DSL with a shallow embedding in Rust. The Rust parser, through the macro system, was used to parse it. The Rust compile-time algorithms were used to type check it. And we suggested a way for the Rust runtime system to run the code, since that is a far more common task. A full demo can be run and modified from <https://play.rust-lang.org/?gist=a0f5ec8999cb08de3842500a9aa959a7&version=stable&mode=debug&edition=2015>.

It is our hope that these explorations will inform further language design. Rust's traits were not originally intended to be used this way, as is obvious looking at error messages of some programs that fail to type check. We hope that meta-programming becomes more valuable in the future, and use cases like those demonstrated will highlight areas to work on.

REFERENCES

- [1] Stephen Chang, Alex Knauth, and Ben Greenman. Type systems as macros. In *Proceedings of the 44th ACM SIGPLAN Symposium on Principles of Programming Languages, POPL 2017, Paris, France, January 18-20, 2017*, pages 694–705, 2017.