



Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

Top-Down Inductive Synthesis with Higher-Order Functions

Master Thesis

Alexandra Maximova

August 15, 2016

Advisors: Prof. Dr. Martin Vechev, Dimitar Dimitrov

Department of Computer Science, ETH Zürich

Abstract

TODO: write me :)

Contents

Contents	iii
1 Introduction	1
1.1 About program synthesis in general	1
1.2 Problem definition	1
1.3 Contributions	2
2 Related Work	3
3 Top down type driven synthesis	9
3.1 Problem	9
3.1.1 The 'replicate' example	9
3.1.2 Lessons learned	12
3.2 Calculus	13
3.2.1 Terms and Types	13
3.2.2 Encodings	14
3.2.3 Evaluation semantics	15
3.2.4 Type checking	16
3.2.5 Type unification	17
3.3 Search	18
3.3.1 Search space	18
3.3.2 Exploration	19
3.4 Cost functions	20
3.5 Black list	21
3.6 Templates	21
3.6.1 Successor rules	23
4 Implementation	25
5 Evaluation	27
5.1 Experimental set up	27

CONTENTS

5.1.1	Evaluation on benchmarks	27
5.1.2	Automatic black list generation	29
5.2	Factors affecting runtime	29
5.2.1	Size of the solution	29
5.2.2	Number of components	30
5.2.3	Examples	30
5.2.4	Blacklist	31
5.2.5	Templates	32
5.2.6	Cost functions	32
5.2.7	Stack vs Queue expansion	33
5.3	Synthesised solutions	34
5.4	Automatic black list	35
5.5	Comparison to related work	36
6	Conclusions	39
6.1	Conclusions	39
6.2	Future Work	39
	Bibliography	41

Chapter 1

Introduction

1.1 About program synthesis in general

put in some context, link ??

- What is program synthesis
- Problem too general, needs to be restricted. In Chapter 2 we will see how other restrict the search space. Restrict to components.
- Motivate why it is still interesting to have such a synthesiser (thin interfaces and so on?)
- List of contributions
- Explain the structure of the thesis. (maybe merge with the contributions?)

for motivation you could write something about the extremely restricted list library in ocaml :)

Consider you are writing code in a functional programming language with a smaller choice of library functions than you are used to. For example (true story), if you are using OCaml and you are surprised that `replicate` is missing even in the more complete core library [cite jane street core](#), you could spend a couple of minutes writing your recursive version of `replicate`. Or you can synthesize it with TAMANDU in less than one second from other components. **TODO:** rewrite without 'you' and maybe don't mention ocaml, core and all that thing?.

1.2 Problem definition

have many components, put them together into a program, no lambdas, no if-then-else, no recursion

1.3 Contributions

Evaluation, exploring the baseline algorithm, exploring the search space

Chapter 2

Related Work

Here discuss only synthesis of functional programs. Input-output examples are an intuitive and easy way to specify programs. Types are good to prune the search space, but simple types are not enough to specify a program. So the research went in two directions: on one side, there was the need to do something to the I/O-examples to make the search more efficient, on the other side more complex and expressive types that can act as a specification, for example the *refinement types* from [7]. Lately, there was the idea that one can automatically generate those ty[ypes from I/O-examples, to have both the advantages of an intuitive and easy specification and the research done on type systems, liquid types and something like that.

Let's start. My 5 papers, organize some information about them.

- Synquid (Nadia) put the replicate example if you can (should not be that difficult). Not that "user friendly". Lambda, conditionals, recursion, components. Synthesis starting from a partial program possible. Fast, good for sorting. It is possible to download the tool (put link), try it online and there is even an emacs-mode for it.
 1. What is the specification? The refinement type of the desired program, that is a type decorated with logical predicates. Can bring the replicate example.
 2. What is the target language? Haskell, I think. I saw lambda abstractions and recursion. Can generate and use higher-order components. You can specify own datatypes and own components. Pattern matching, structural recursion, ability to use and generate polymorphic higher-order functions, reasoning about universal and recursive properties of data structures.
The Synquid language has lambda expressions, pattern matching, conditional and fixpoint.

3. What can it do well and fast? How fast? Most advanced benchmark: generate various sorting algorithms for data structures with complex invariants (search trees, heaps, balanced trees). This is out of scope of my synthesis procedure. For the "easy" benchmarks (what I have) everything under 0.4 s. For more complicated stuff like sorting and tree insertion under 5 seconds. For most complicated stuff (red-black trees) under 20 s.
 4. What is the difficulty? What can they not generate (or take a long time to generate)? How much time do they need? It is not always easy for an inexperienced user to specify a program and to find the right 'measure' function for a datatype.
They need up to 20 s to generate balance over red-black trees. They usually don't give more than 5 components and 6 measure functions over the needed data types.
The specification is quite big compared to the synthesised program, half of the programs the specification is one third of the generated program or more.
 5. What do they do? In [7] SYNQUID is proposed. Refinement types (types decorated with logical predicates) are used to prune the search space. SMT-solvers are used to satisfy the logical predicates. The key is the new procedure for type inference (called modular refinement type reconstruction), which thank to its modularity scales better than other existing inference procedures for refinement types. Programs can therefore be type checked even before they are put together.
- Lambda square (Feser). How would you specify replicate? Give an idea of how it is generated. User friendly, as only I/O-examples need to be specified, type is reconstructed automatically.
 1. What is the specification? I/O-examples
 2. What is the target language? Their synthesis algorithm targets a functional programming language that permits higher-order functions, combinators like map, fold and filter, pattern-matching, recursion, and a flexible set of primitive operators and constants. Actually, lambda calculus with algebraic data types and recursion.
 3. What can it do well and fast? How fast? Generating programs over nested structures like trees of lists, lists of lists and lists of trees. They generate the half of the benchmarks in under 0.43 s.

-
4. What is the difficulty? What can they not generate (or take a long time)? How much time do they need? Only 7 higher-order combinators and not much more first-order components.
They use relatively many examples (between 3 and 12, mostly 4 to 6).
They need up to 320 s to generate droplast. Other benchmarks that take more than 100 s to synthesise are removing duplicates from a list, inserting a tree under each leaf of another tree and dropping the smallest number in a list of lists.
 5. What do they do? Inductive generalization, deduction, enumerative search. First, generate *hypotheses* (that is, programs with free variables like $\lambda x. \text{map } ?f \ x$ where $?f$ is a placeholder for an unknown program that needs to be generated) in a type-aware manner (the type is inferred from the I/O-examples). Then deduction based on the semantics of the higher-order combinators is used either to quickly refute a hypothesis or infer new I/O-examples to guide the synthesis of missing functions. Best-first enumerative search is used to enumerate candidate programs to fill in the missing parts of hypotheses.
The tool proposed in [3] is called λ^2 and generates its output in λ -calculus with algebraic types and recursion. The user specifies the desired program providing input-output examples. No particular knowledge is required from the user, as was demonstrated using random input-output examples. The examples are inductively generalized in a type-aware manner to a set of hypotheses (programs that possibly have free variables). The key idea are the hard-coded deduction rules used to prune the search space depending on the semantics of some of the higher-order combinators (map, fold, filter and a few others). Deduction is also used to infer new input-output examples in order to generate the programs needed to fill in the holes in the hypotheses. This tool is able to synthesize programs manipulating recursive data structures like lists, trees and nested data structures such as lists of lists and trees of lists. The tool is able to synthesise all benchmark programs in under 7 minutes.
- Escher (Kincaid). User friendly, as untyped and only I/O-examples need to be specified. Powerful (recursion, special if-then-else, components) Goal graph.
 1. What is the specification? I/O-examples. Oracle simulating interaction with the user.
 2. What is the target language? Recursion, special treatment for con-

ditionals, components. That's really all they have: either apply a component to its arguments (including component self referring to the program being generated for recursion), a constant, an input variable or a conditional. Programming language is untyped. Components may be higher-order.

3. What can they do well and fast? How fast? Synthesise recursive programs (tail recursive, divide-and-conquer, mutually recursive). In the evaluation, they give 23 components to all benchmarks, none of them is higher-order. Can generate all benchmarks in under 11 s, most of them in under 1 s.
4. What is the difficulty? What can they not generate (or take a long time)? How much time do they need? The user must give a *closed* example set, that is for each example, there must be another example that can be used for a recursive call.
5. What do they do? Forward search: inductive enumeration, add a new component to an already synthesised program. Conditional inference: use the goal graph to see if you can join to synthesised programs by a conditional statement. A heuristic guides the alternation between forward search and conditional inference. Programs are associated with value vectors. Search space is also pruned based on *observational equivalence*, that is programs with equivalent value vectors are treated as equivalent programs. Can be formalised as a non-deterministic transition system over configurations (triples consisting of a set of synthesised programs, a goal graph and a list of input-output examples) with six transition rules: an initial rule to start the search, a terminate rule to terminate the search and four synthesis rules (one for forward search, two for conditional split and one to ask for new input-output examples to evaluate a recursive call). The recursive calls are answered by the oracle. Rule scheduling is added to turn this into a practical system (different heuristics).
In [2] ESCHER is presented, an inductive synthesis algorithm that learns a recursive procedure from input-output examples provided by the user. The user must provide a "closed" set of examples, otherwise recursion cannot be handled properly. The target language is untyped, first-order and purely functional. The algorithm is parametrized by components that can be instantiated differently to suit different domains. The approach combines enumerative search and conditional inference. The key idea is to use a special data structure, a *goal graph*, to infer conditional branches

instead of treating if-then-else as a component. Observational equivalence is also used to prune the search space. Programs with the same value vectors (output of the program when applied to the inputs of the input-output examples) are considered equivalent and only one of them is synthesized. An implementation of the tool was tested on a benchmark consisting of recursive programs (including tail-recursive, divide-and-conquer and mutually recursive programs) drawn from functional programming assignments and standard list and tree manipulation programs. For all examples the same fixed set of components was used. The tool is able to synthesize all of them quickly. There is very little information on how many input-output examples were needed to synthesize the benchmarks and how difficult it is for a non-experienced user to come up with a "closed" set of examples.

- Myth (Osera). Like mine, requires type and I/O-examples. Refinement tree.
 1. What is the specification? A type signature, the components and a list of input-output examples.
 2. What is the target language? pattern matching, recursion, higher-order functions in typed programming languages. Can synthesise higher-order functions, programs using higher-order functions and work with large algebraic datatypes. ML-like language with algebraic data types, match, top-level function definitions and explicitly recursive functions.
 3. What can they do well and fast? How fast? Recursive programs with pattern matching. It's very fast, many programs are synthesised in around 0.1 s. It can also generate larger programs (75 AST nodes) in reasonable time (3 s for calculating the set of free variables in an untyped lambda-calculus).
 4. What is the difficulty? What can they not generate (or take a long time)? How much time do they need? To generate recursive functions, they also need a closed set of examples, so that a recursive call to the function being synthesised can be answered by an input-output example. They also require relatively many examples. They use a relatively large context, but they do not say how big it is. Lacks support richer types like products and polymorphic types.
 5. What do they do? The tool in [5] is called MYTH and uses not only type information but also input-output examples to restrict

2. RELATED WORK

the search space. The special data structure used to hold this information is the *refinement tree*. This system can synthesize higher-order functions, programs that use higher order functions and work with large algebraic data types.

There is an ML-like type system that incorporates input-output examples. Two pieces: a *refinement tree* and an enumerative search. Two major operations: refine the goal type and the examples (push them down in the refinement tree) and guess a term of the right type that matches the examples (for one of the nodes of the refinement tree).

Chapter 3

Top down type driven synthesis

In this chapter we will formally define our type-directed top-down synthesis procedure. We will start with an intuitive description of the problem and move on to the formal definitions of the target programming language and the search space. Finally, we will define the synthesis procedure and present some enhancements.

3.1 Problem

3.1.1 The 'replicate' example

Recall the example from Chapter 1 where we wanted to synthesise `replicate`. As our synthesis procedure is type-driven and example-based, the user specifies a program by providing its type along with a few I/O-examples. Let us specify `replicate` as follows.

```
replicate :: ∀X. Int → X → List X
replicate 3 1 = [1,1,1]
replicate 2 [] = [[] ,[]]
```

We also need a library of components, from which we are going to compose our program. Let us assume we have the standard list combinators `map` and `foldr` with `enumTo`, the function that returns a list from 1 up to its argument, and `const`, that always returns its first argument, along with the list constructors `cons` and `[]` and the integer constructors `succ` and `0`. Even with so few library components, the search space is quite big.

The goal is to put together components from the library in order to get a list. More concretely, we fix `X` to be a fixed input type variable, we fix `n` to be an integer and `x` to be a fixed input variable of type `X`. Now the goal looks like

```
replicate n x = ?p
?p :: List X
```

, where $?p$ is a *hole*, a new fresh variable whose type is known.

Which components can we use in order to get something of type `List X`? We cannot use `enumTo`, because it only produces a list of integers, but all other possibilities are open. We could either fold with an interesting function, or map some function, or even use `const` with a smart first argument. But the first and easiest thing that has the type `List X` is `[]`. That is, our first program looks as follows.

```
replicate n x = []
?p = [] :: List X
```

This is a *closed* program, that is, a program without holes that can be evaluated on the input-output examples. Alas, it does not satisfy any of them. Therefore we must try the other components. We have following possibilities to fill in the hole $?p$.

```
replicate n x = cons ?x ?xs
?p = cons ?x ?xs :: List X
?x :: X
?xs :: List X
```

```
replicate n x = map ?f ?xs
?p = map ?f ?xs :: List X
?f :: ?Y → X
?xs :: List ?Y
```

Here $?Y$ is a fresh type variable that will be instantiated with something later. As of now we have no idea about the type of the argument of $?f$, we only know that it has to match the type of the elements of $?xs$.

```
replicate n x = foldr ?f ?init ?xs
?p = foldr ?f ?init ?xs :: List X
?f :: ?Y → List X → List X
?init :: List X
?xs :: List ?Y
```

```
replicate n x = const ?xs ?s
?p = const ?xs ?s :: List X
?xs :: List X
?s :: ?Y
```

Now we take the most promising program and try to fill in one of its *holes* (the fresh variables starting with $?$). Let us decide that the first one is the most promising. We now have two holes to fill in, a function that can take something and return an `X` and a list of something. Note that it has to be

possible for all possible instantiations of X and note that we have only one value of type X , namely x , the second argument to `replicate`.

What are the possibilities to instantiate $?f$? Obviously, we cannot use `map` or `enumTo`, because they return lists. But all other possibilities are still open. We have to add following programs to our pool of possible solutions.

```
replicate n x = map (foldr ?g ?init) ?xs
?p = map ?f ?xs :: List X
?f = foldr ?g ?init :: List ?Z → X
?g :: ?Z → X → X
?init :: X
?xs :: List (List ?Z)
```

Note that we instantiated $?Y$ with `List ?Z`, because `foldr` takes a list as its last argument.

```
replicate n x = map (const ?x) ?xs
?f = const ?x :: ?Y → X
?x :: X
?xs :: List ?Y
```

As you noticed, we maintain a frontier of programs with holes and expand one of the holes of the most promising program. Let us do it again. From the 6 programs generated so far we choose the most promising one. Let us decide that it is the last one. It has two holes to fill in. For the first hole, $?x$, we have only one possibility. As this hole must be of type X , the only thing we can take is the second argument of `replicate`, x .

```
replicate n x = map (const x) ?xs
?p = map ?f ?xs :: List X
?f = const ?x :: ?Y → X
?x = x :: X
?xs :: List ?Y
```

Let us directly decide that this is the most promising program so far. Later, in Section 3.4, we will define cost functions of programs and define the most promising program to be the one with the smallest cost. For now we just choose the one that will lead us to the solution.

Like in the beginning, we have to generate a list. However, since this time the type of the elements is not fixed, we cannot rule out `enumTo`. Therefore we have a lot of possibilities, starting with `replicate n x = map (const x) []`, where we instantiate $?xs$ with `[]`, and ending with

```
replicate n x = map (const x) (enumTo ?n)
?p = map ?f ?xs :: List X
```

```
?f = const ?x :: ?Y → X
?x = x :: X
?xs = enumTo ?n :: List Int
?n :: Int
```

First, we are going to evaluate the closed program `replicate n x = map (const x) []`. This program does not satisfy any of the input-output examples too.

The next step is to expand one of the holes of the most promising program, let us decide that it is the last one. That is, we are looking for an integer. An integer can either be the constructor 0, or the constructor `succ` applied to some integer hole, the first argument of `replicate`, `n`, or `const` with some clever first argument applied to something.

Notice that among the new programs there are two closed programs, `replicate n x = map (const x) (enumTo 0)` and `replicate n x = map (const x) (enumTo n)`. Evaluation shows that only the second one satisfies the I/O-examples.

3.1.2 Concepts

This example shows some important concepts that are defined formally in the next sections of this chapter.

- Hole: unknown part of a program that can be instantiated with some other programs. Only its type is known.
- Closed program: a program without holes that can be evaluated on the input-output examples. The terms and the types of our calculus are formally defined in Section 3.2.1.
- Type-aware expansion of holes: we expand holes based on their type. In Section 3.3.1 we can find the rules according to which a program is expanded.
- Best first search: a frontier of programs with holes is maintained, one hole of the most promising is expanded in every iteration. The best first search algorithm is defined in Section 3.3.2 and a notion of most promising program is presented in Section 3.4.
- Superfluous program: a program that is equivalent to a simpler program. For example, a human programmer would not instantiate a hole with `const ?x ?y`, because he could have written just `?x` instead, which is always a shorter and preferable program. Section 3.5 shows one way to rule out superfluous programs.

3.2 Calculus

In this section we will introduce three different calculus.

- System F as in the excellent book of Pierce [6]
- The system we derive from System F and in which the components are defined (includes holes, input variables and recursion for terms and types)
- A subset of the second system, used for generation (only application of components, holes or input variables is allowed).

The exposition will closely follow Pierce [6].

Let's start with the standart System F. Our calculus is based on System F. A subset of our calculus is used for generation.

- generate programs from restricted calculus
- calculus itself is still powerful, because we use it to define the library components as well.

3.2.1 Terms and Types

- similar to Pierce
- except holes, components and free variables
- recursion
- recursive types that can take parameters
- the search space is not the whole language but only a subset of it

Real System F

$$t ::= x \mid \lambda x : T. t \mid t t \mid \Lambda X. t \mid t [T] \quad (\text{terms})$$

$$T ::= X \mid T \rightarrow T \mid \forall X. T \quad (\text{types})$$

$$\Gamma ::= \emptyset \mid \Gamma \cup \{x : T\} \mid \Gamma \cup \{X\} \quad (\text{variable bindings})$$

We use an extension of System F featuring holes $?x$, input variables i as well as named library components c and named types C . The use of the names enables recursion in the definition of library components and types. Named types also support type parameters. The number of type parameters supported by a named type is denoted as K in its definition.

Question: where is the output? Can you define input and output pairs? Our system

$$t ::= x \mid \lambda x : T. t \mid t t \mid \Lambda X. t \mid t [T] \mid c \mid ?x \mid i \quad (\text{terms})$$

$$T ::= X \mid T \rightarrow T \mid \forall X. T \mid ?X \mid I \mid C [T, \dots, T] \quad (\text{types})$$

$$\Gamma ::= \emptyset \mid \Gamma \cup \{x : T\} \mid \Gamma \cup \{X\} \quad (\text{variable bindings})$$

$$\Xi ::= \emptyset \mid \Xi \cup \{?x : T\} \mid \Xi \cup \{?X\} \quad (\text{hole bindings})$$

$$\Phi ::= \emptyset \mid \Phi \cup \{i = t : T\} \mid \Phi \cup \{I = T : K\} \quad (\text{input variable bindings})$$

$$\Delta ::= \emptyset \mid \Delta \cup \{c = t : T\} \mid \Delta \cup \{C = T : K\} \quad (\text{library components})$$

Question: do we need the definitions of library components for synthesis? We use them only for evaluation, but we always evaluate programs during synthesis. Same for the input variables. Actually, for each I/O-example we get the pair (Φ, o) , where Φ instantiates all input variables and input types of a program and o is the expected output. Subset of our system that builds the search space

$$t ::= t t \mid t [T] \mid c \mid ?x \mid i \quad (\text{terms})$$

$$T ::= X \mid T \rightarrow T \mid \forall X. T \mid ?X \mid I \mid C [T, \dots, T] \quad (\text{types})$$

$$\Xi ::= \emptyset \mid \Xi \cup \{?x : T\} \mid \Xi \cup \{?X\} \quad (\text{hole bindings})$$

$$\Phi ::= \emptyset \mid \Phi \cup \{i = t : T\} \mid \Phi \cup \{I = T : K\} \quad (\text{input variable bindings})$$

$$\Delta ::= \emptyset \mid \Delta \cup \{c = t : T\} \mid \Delta \cup \{C = T : K\} \quad (\text{library components})$$

A program is the quadriplet $\{\Xi, \Phi, \Delta \vdash t :: T\}$. A term is called *closed* if it does not contain holes. A program is closed, if Ξ is empty and t and T do not contain holes.

3.2.2 Encodings

Note that in the definition of types we do not see familiar types such as booleans, integers, lists or trees. All these types can be encoded in System F using either the Church's or the Scott's encoding [1]. We opted for the Scott's encoding because it's more efficient in our case. Scott's booleans coincide with Church's booleans and are encoded as follows.

$$\begin{aligned} \text{Bool} &= \forall R. R \rightarrow R \rightarrow R \\ \text{true} &= \Lambda R. \lambda x_1 : R. \lambda x_2 : R. x_1 \\ &\quad : \text{Bool} \\ \text{false} &= \Lambda R. \lambda x_1 : R. \lambda x_2 : R. x_2 \\ &\quad : \text{Bool} \end{aligned}$$

```

if-then-else =  $\Lambda X. \lambda b:\text{Bool}. \lambda t:X. \lambda f:X. b [X] t f$ 
              :  $\forall X. \text{Bool} \rightarrow X \rightarrow X \rightarrow X$ 

```

Scott's integers differ from Church's integers as they are more suitable for pattern matching. *because they don't unwrap the whole integer every time.*

```

Int =  $\forall R. R \rightarrow (\text{Int} \rightarrow R) \rightarrow R$ 
zero =  $\Lambda R. \lambda z:R. \lambda s:\text{Int} \rightarrow R. z$ 
      : Int
succ =  $\lambda n:\text{Int}. \Lambda R. \lambda z:R. \lambda s:\text{Int} \rightarrow R. s n$ 
      : Int  $\rightarrow$  Int
case =  $\Lambda R. \lambda n:\text{Int}. \lambda a:R. \lambda f:\text{Int} \rightarrow R. n [R] a f$ 
      :  $\forall R. \text{Int} \rightarrow R \rightarrow (\text{Int} \rightarrow R) \rightarrow R$ 

```

Analogously, Scott's lists are a recursive type and naturally support pattern matching.

```

List X =  $\forall R. R \rightarrow (X \rightarrow \text{List } X \rightarrow R) \rightarrow R$ 
nil =  $\Lambda X. \Lambda R. \lambda n:R. \lambda c:X \rightarrow \text{List } X \rightarrow R. n$ 
     :  $\forall X. \text{List } X$ 
con =  $\Lambda X. \lambda x:X. \lambda xs:\text{List } X. \Lambda R. \lambda n:R. \lambda c:X \rightarrow \text{List } X$ 
      $\rightarrow R. c x xs$ 
     :  $\forall X. X \rightarrow \text{List } X \rightarrow \text{List } X$ 
case =  $\Lambda X. \Lambda Y. \lambda l:\text{List } X. \lambda n:Y. \lambda c:X \rightarrow \text{List } X \rightarrow Y.$ 
      $l [Y] n c$ 
     :  $\forall X. \forall Y. \text{List } X \rightarrow Y \rightarrow (X \rightarrow \text{List } X \rightarrow Y) \rightarrow Y$ 

```

3.2.3 Evaluation semantics

In this section we present the evaluation semantics of our extended System F, that is the second calculus introduced in Section 3.2.1. The evaluation semantics is a standard eager evaluation and we refer to the excellent book of Benjamin Pierce about type systems [6] for an introduction to the evaluation semantics of System F and evaluation rules in general.

The evaluation judgement $\Phi, \Delta \vdash t \longrightarrow t'$ means that the term t evaluates in one step to the term t' under the free variable bindings library Φ , that contains concrete instantiations for the input variables, and the component library Δ , that contains the definitions of the library components. Before listing the evaluation rules, let us define *value* v to be a term to which no evaluation rule apply.

$$\frac{c = t : T \in \Delta}{\Phi, \Delta \vdash c \longrightarrow t} \text{E-Lib}$$

$$\frac{i = t : T \in \Phi}{\Phi, \Delta \vdash i \longrightarrow t} \text{E-Inp}$$

$$\frac{\Phi, \Delta \vdash t_1 \longrightarrow t'_1}{\Phi, \Delta \vdash t_1 t_2 \longrightarrow t'_1 t_2} \text{E-App1}$$

$$\frac{\Phi, \Delta \vdash t_2 \longrightarrow t'_2}{\Phi, \Delta \vdash v_1 t_2 \longrightarrow v_1 t'_2} \text{E-App2}$$

$$\frac{}{\Phi, \Delta \vdash (\lambda x : T_{11}. t_{12}) v_2 \longrightarrow [x \mapsto v_2] t_{12}} \text{E-AppAbs}$$

$$\frac{}{\Phi, \Delta \vdash (\Lambda X. t_2) [T_2] \longrightarrow [X \mapsto T_2] t_2} \text{E-APPABS}$$

Rules E-Lib and E-Inp load the definitions of library components or input variables from the respective library. E-App1 and E-App2 evaluate the left hand side, respectively the right hand side, of an application. E-AppAbs and E-APPABS get rid of an abstraction and substitute the argument into the body.

Note that E-App2 applies only if the left hand side of the application cannot be evaluated further and that E-AppAbs applies only when the argument of the lambda abstraction is a value, determining the order of evaluation.

3.2.4 Type checking

In this sections we will present the typing rules of our extended System F, that is the second calculus presented in Section 3.2.1. The typing judgement $\Gamma, \Xi, \Phi, \Delta \vdash t : T$ means the term t has type T in the contexts Γ and Ξ , binding respectively variables and holes, and Φ and Δ , containing signatures and definitions of respectively input variables and library components. The typing judgement is similar to the typing judgement of System F presented in the book about type systems of Benjamin Pierce [6] and we refer to the book for more details.

$$\frac{x : T \in \Gamma}{\Gamma, \Xi, \Phi, \Delta \vdash x : T} \text{T-Var}$$

$$\frac{?x : T \in \Xi}{\Gamma, \Xi, \Phi, \Delta \vdash ?x : T} \text{T-Hol}$$

$$\frac{i = t : T \in \Phi}{\Gamma, \Xi, \Phi, \Delta \vdash i : T} \text{T-Inp}$$

$$\frac{c = t : T \in \Delta}{\Gamma, \Xi, \Phi, \Delta \vdash c : T} \text{T-Lib}$$

$$\frac{\Gamma \cup \{x : T_1\}, \Xi, \Phi, \Delta \vdash t_2 : T_2}{\Gamma, \Xi, \Phi, \Delta \vdash \lambda x : T_1. t_2 : T_1 \rightarrow T_2} \text{T-Abs}$$

$$\frac{\Gamma, \Xi, \Phi, \Delta \vdash t_1 : T_1 \rightarrow T_2 \quad \Gamma, \Xi, \Phi, \Delta \vdash t_2 : T_1}{\Gamma, \Xi, \Phi, \Delta \vdash t_1 t_2 : T_2} \text{T-App}$$

$$\frac{\Gamma \cup \{X\}, \Xi, \Phi, \Delta \vdash t_2 : T_2}{\Gamma, \Xi, \Phi, \Delta \vdash \Lambda X. t_2 : \forall X. T_2} \text{T-ABS}$$

$$\frac{\Gamma, \Xi, \Phi, \Delta \vdash t_1 : \forall X. T_{12}}{\Gamma, \Xi, \Phi, \Delta \vdash t_1 [T_2] : [X \mapsto T_2] T_{12}} \text{T-APP}$$

The rules T-Var, T-Hol, T-Inp and T-Lib load the signature of variables, holes, input variables or library components from the respective context. T-Abs types a lambda abstraction as an arrow type from the type of its variable to the type of its body. T-Abs gives a type abstraction as a universal type in accordance to the type of its body. An application typechecks only if the left hand side has an arrow type and the type of the argument is equal to the type of the argument of the left hand side. Analogously, a type application typechecks only if the left hand side has a universal type.

3.2.5 Type unification

TODO: Write this

Since unification on System F types is undecidable [4], we decided to do something simpler. Explain what. Boom, fertig, party.

The unification algorithm is based on the unification algorithm for typed lambda calculus from the book of Pierce **TODO: cite the book properly!!!** and slightly modified to fit our needs.

How did you do type unification? Unification of universal types is not implemented. If you need to unify to universal types, you should transform all bound variables into holes and remove the universal quantifier.

A set of constraints is a set of types that should be equal under a substitution. The unification algorithm is supposed to output a substitution σ so that $\sigma(S) = \sigma(T)$ for every constraint $S = T$ in \mathcal{C} .

Input: Set of constraints $\mathcal{C} = \{T_{11} = T_{12}, T_{21} = T_{22}, \dots\}$
Output: Substitution σ so that $\sigma(T_{i1}) = \sigma(T_{i2})$ for every constraint $T_{i1} = T_{i2}$ in \mathcal{C}

Function *unify*(\mathcal{C}) **is**

```

  if  $\mathcal{C} = \emptyset$  then []
  else
    let  $\{T_1 = T_2\} \cup \mathcal{C}' = \mathcal{C}$  in
    if  $T_1 = T_2$  then
      | unify( $\mathcal{C}'$ )
    else if  $T_1 = ?X$  and  $?X$  does not occur in  $T_2$  then
      | unify( $[X \mapsto T_2]\mathcal{C}' \circ [X \mapsto T_2]$ )
    else if  $T_2 = ?X$  and  $?X$  does not occur in  $T_1$  then
      | unify( $[X \mapsto T_1]\mathcal{C}' \circ [X \mapsto T_1]$ )
    else if  $T_1 = T_{11} \rightarrow T_{12}$  and  $T_2 = T_{21} \rightarrow T_{22}$  then
      | unify( $\mathcal{C}' \cup \{T_{11} = T_{21}, T_{12} = T_{22}\}$ )
    else if  $T_1 = C [T_{11}, T_{12}, \dots, T_{1k}]$  and  $T_2 = C [T_{21}, T_{22}, \dots, T_{2k}]$ 
      then
      | unify( $\mathcal{C}' \cup \{T_{11} = T_{21}, T_{12} = T_{22}, \dots, T_{1k} = T_{2k}\}$ )
    else
      | fail
    end
  end
end

```

Algorithm 1: Type unification

3.3 Search

how do we explore the search space (best-first search).

3.3.1 Search space

Use only third system presented in Section System F.

Formal problem definition. Given a library Δ , a goal type T and a list of I/O-examples $[(\Phi_1, o_1), \dots, (\Phi_N, o_N)]$, find a closed term t so that $\emptyset, \emptyset, \Phi_1, \Delta \vdash t : T$ (that is, t has the goal type under an empty variable binding context and an empty hole binding context) and t satisfies all I/O-examples, that is $\Phi_n, \Delta \vdash t \longrightarrow^* \vdash t'$ and $\Phi_n, \Delta \vdash o_n \longrightarrow^* \vdash t'$ for all $n = 1, \dots, N$.

We see the search space as a graph of programs with holes (see third syntax presented in Section System F) where there is an edge between two terms t_1 and t_2 if and only if the judgement *derive* defined below $\Xi, \Phi, \Delta \vdash t_1 :: T_1 \Rightarrow \Xi', \Phi, \Delta \vdash t_2 :: T_2$ holds between the two.

To express the rules in a more compact form, we introduce *evaluation contexts*.

An evaluation context is an expression with exactly one syntactic hole $[]$ in which we can plug in any term. For example, if we have the context \mathcal{E} we can place the term t into its hole and denote this new term by $\mathcal{E}[t]$.

A hole $?x$ can be turned into a library component c from the context Δ or an input variable i from the context Φ . The procedure $\text{fresh}(T)$ transforms universally quantified type variables into fresh type variables $?X$ not used in Ξ . The notation $\sigma(\Delta)$ denotes the application of the substitution σ to all types contained in the context Δ .

$$\frac{c : T_c \in \Delta \quad \sigma \text{ unifies } T \text{ with } \text{fresh}(T_c)}{\Xi, \Phi, \Delta \vdash ?x :: T \Rightarrow \sigma(\Xi \setminus \{?x : T\}), \Phi, \Delta \vdash c :: \sigma(T)} \text{D-VarLib}$$

$$\frac{i : T_i \in \Phi \quad \sigma \text{ unifies } T \text{ with } \text{fresh}(T_i)}{\Xi, \Phi, \Delta \vdash ?x :: T \Rightarrow \sigma(\Xi \setminus \{?x : T\}), \Phi, \Delta \vdash i :: \sigma(T)} \text{D-VarInp}$$

A hole can also be turned into a function application of two new active holes.

$$\frac{?X \text{ is a fresh type variable} \quad \Xi' = \Xi \setminus \{?x : T\} \cup \{?x_1 : ?X \rightarrow T, ?x_2 : ?X, ?X\}}{\Xi, \Phi, \Delta \vdash ?x :: T \Rightarrow \Xi', \Phi, \Delta \vdash ?x_1 ?x_2 :: T} \text{D-VarApp}$$

In all other cases we just choose a hole and expand it according to the three rules above.

$$\frac{\Xi, \Phi, \Delta \vdash ?x :: T_1 \Rightarrow \Xi', \Phi, \Delta \vdash t'_1 :: T'_1}{\Xi, \Phi, \Delta \vdash t[?x] :: T \Rightarrow \Xi', \Phi, \Delta \vdash t[t_1] :: [T_1 \mapsto T'_1]T} \text{D-App}$$

Note that the types of all successor programs unify with the types of their ancestors. Thus, the search is type directed. Only programs of the right type are generated.

3.3.2 Exploration

Write also about stack vs queue of open holes

We explore the search graph using a best first search. We can play with the algorithm in two points (marked in blue): first, we can define which hole to expand first, second, we can choose the compare function of the priority queue. Different approaches to define the compare function are discussed in the next session. Concerning the order of expansion of the holes, we tried two strategies: the first one was maintaining a stack of holes (Ξ implemented as a stack), which leads to the expansion of the deepest hole first, the second one was maintaining a queue of holes (Ξ implemented as a queue), which leads to the expansion of the oldest hole first.

Input: goal type T , library components Δ , list of input-output examples $[(\Phi_1, o_1), \dots, (\Phi_N, o_N)]$
Output: closed program $\{\Xi, \Phi_1, \Delta \vdash t :: T\}$ that satisfies all I/O-examples

```

queue ← PriorityQueue.empty compare
queue ← PriorityQueue.push queue  $\{\Xi, \Phi_1, \Delta \vdash ?x :: T\}$ 
while not ((PriorityQueue.top queue) satisfies all I/O-examples) do
  | successors ← successor (PriorityQueue.top queue)
  | queue ← PriorityQueue.pop queue
  | for all  $s$  in successors do
    | queue ← PriorityQueue.push queue  $s$ 
  | end
end
return PriorityQueue.top queue

```

Algorithm 2: Best first search

3.4 Cost functions

The compare function in the best first search algorithm is cost p_1 - cost p_2 . There are different possibilities to implement this cost function. We will present four alternatives.

number of nodes The first cost function is based only on the number of nodes of the term. Longer and more complicated terms are disadvantaged.

Input: term t
Output: weighted number of nodes in t

```

nof-nodes( $c$ ) = 1
nof-nodes( $?x$ ) = 2
nof-nodes( $i$ ) = 0
nof-nodes( $x$ ) = 1
nof-nodes( $\lambda x : T. t$ ) = 1 + nof-nodes( $t$ )
nof-nodes( $t_1 t_2$ ) = 1 + nof-nodes( $t_1$ ) + nof-nodes( $t_2$ )
nof-nodes( $\Lambda X. t$ ) = 1 + nof-nodes( $t$ )
nof-nodes( $t [T]$ ) = 1 + nof-nodes( $t$ )

```

Algorithm 3: Cost function based on the number of nodes

number of nodes and types The second cost function also adds a factor based on the types appearing in the term, thus penalizing terms with type application and very complicated types.

Input: term t

Output: sum of weighted number of nodes in term t and weighted number of nodes in the types appearing in t

$$\text{nof-nodes-type}(X) = 1$$

$$\text{nof-nodes-type}(\text{?}X) = 0$$

$$\text{nof-nodes-type}(I) = 0$$

$$\text{nof-nodes-type}(C [T_1, \dots, T_k]) = 0$$

$$\text{nof-nodes-type}(T_1 \rightarrow T_2) = 3 + \text{nof-nodes-type}(T_1) + \text{nof-nodes-type}(T_2)$$

$$\text{nof-nodes-type}(\forall X. T) = 1 + \text{nof-nodes-type}(T)$$

$$\text{nof-nodes-term}(c) = 1$$

$$\text{nof-nodes-term}(\text{?}x) = 2$$

$$\text{nof-nodes-term}(i) = 0$$

$$\text{nof-nodes-term}(x) = 1$$

$$\text{nof-nodes-term}(\lambda x : T. t) = 1 + \text{nof-nodes-term}(t) + \text{nof-nodes-type}(T)$$

$$\text{nof-nodes-term}(t_1 t_2) = 1 + \text{nof-nodes-term}(t_1) + \text{nof-nodes-term}(t_2)$$

$$\text{nof-nodes-term}(\wedge X. t) = 1 + \text{nof-nodes-term}(t)$$

$$\text{nof-nodes-term}(t [T]) = 1 + \text{nof-nodes-term}(t) + \text{nof-nodes-type}(T)$$

$$\text{nof-nodes-and-types}(t) = \text{nof-nodes-term}(t)$$

Algorithm 4: Cost function based on the number of nodes and types

no same component In the third function we additionally penalize terms using the same component more than once.

length of the string The simplest and most imprecise method to take both the number of nodes and the complexity of the types appearing in the term is taking the length of the string of that term.

3.5 Black list

automatic generation of black list discussed in evaluation A black list is a list of terms. Programs containing a black term as a subterm are not allowed to have successors. Thus, the algorithm above is modified as follows.

One could use the synthesis algorithm presented in Section 3.3.2 to automatically synthesise black lists. For example, one could synthesise many programs corresponding to the identity function or to the empty list or to any other term, and add all those programs but one to the black list.

3.6 Templates

Top-down type-driven synthesis.

Input: term t

Output: sum of the weighted number of nodes in term t , the weighted number of nodes in the types appearing in t and the weighted number of library components appearing more than once in t .

$$\text{nof-nodes-type}(X) = 10$$

$$\text{nof-nodes-type}(\text{?}X) = 3$$

$$\text{nof-nodes-type}(I) = 0$$

$$\text{nof-nodes-type}(C [T_1, \dots, T_k]) =$$

$$4 + \text{nof-nodes-type}(T_1) + \dots + \text{nof-nodes-type}(T_k)$$

$$\text{nof-nodes-type}(T_1 \rightarrow T_2) = 5 + \text{nof-nodes-type}(T_1) + \text{nof-nodes-type}(T_2)$$

$$\text{nof-nodes-type}(\forall X. T) = 10 + \text{nof-nodes-type}(T)$$

$$\text{nof-nodes-term}(c) = 3$$

$$\text{nof-nodes-term}(\text{?}x) = 2$$

$$\text{nof-nodes-term}(i) = 0$$

$$\text{nof-nodes-term}(x) = 10$$

$$\text{nof-nodes-term}(\lambda x : T. t) = 10 + \text{nof-nodes-term}(t) + \text{nof-nodes-type}(T)$$

$$\text{nof-nodes-term}(t_1 \ t_2) = 6 + \text{nof-nodes-term}(t_1) + \text{nof-nodes-term}(t_2)$$

$$\text{nof-nodes-term}(\Lambda X. t) = 10 + \text{nof-nodes-term}(t)$$

$$\text{nof-nodes-term}(t [T]) = 5 + \text{nof-nodes-term}(t) + \text{nof-nodes-type}(T)$$

$$\text{count}(t) = \sum_{c_i \text{ appears in } t} (\text{occurrences of } c_i \text{ in } t) - 1$$

$$\text{no-same-component}(t) = \text{nof-nodes-term}(t) + 3 \text{ count}(t)$$

Algorithm 5: Cost function based on the number of nodes and types penalizing the use of a library component more than once

A template is a program with holes. We are interested in templates where all higher-order components are fixed and there are holes for the first-order components. The search space is thus similar to the search space described in 3.3.1, with the exception that Δ contains only the higher-order components. One of the new things are *closed holes* $\text{?}x$. Those are holes that are supposed to be filled in later with first-order components.

The idea behind the templates is that once the higher-order components are fixed, it should be easy and fast to find a first-order assignment to get the right program. So we could do a limited search from a template and if we do not find a program satisfying all of the I/O-examples we can move quickly to the next template.

We additionally restrict the space by requiring a template to have no more than M higher-order components and no more than P closed holes.

Input: goal type T , library components Δ , list of input-output examples $[(\Phi_1, o_1), \dots, (\Phi_N, o_N)]$, black list $[b_1, \dots, b_M]$

queue \leftarrow PriorityQueue.empty compare

queue \leftarrow PriorityQueue.push queue $\{\Xi, \Phi_1, \Delta \vdash ?x :: T\}$

while not ((PriorityQueue.top queue) satisfies all I/O-examples) **do**

if not ((PriorityQueue.top queue) contains subterm from black list) **then**

 successors \leftarrow successor (PriorityQueue.top queue)

 queue \leftarrow PriorityQueue.pop queue

for all s in successors **do**

 queue \leftarrow PriorityQueue.push queue s

end

else

 queue \leftarrow PriorityQueue.pop queue

end

end

Output: PriorityQueue.top queue

Algorithm 6: Best first search with black list

3.6.1 Successor rules

The successor rules are very similar to the ones defined in 3.3.1, apart from little modifications. That is, now we have a successor rule to close a hole, and we can not instantiate a hole with an input variable any more, because that is supposed to be done in the next step. All the rules are modified to take into account the restriction on the number of components and the number of closed holes. In order to do this, we need to pass along m , the number of higher-order components in the term whose subterms we are traversing.

So we can *close* a hole.

$$\frac{\begin{array}{l} |\Xi| \leq P \text{ and } m \leq M \\ T \text{ is a type a first-order component can have} \end{array}}{\Xi, \Phi, \Delta, m \vdash ?x :: T \Rightarrow \Xi, \Phi, \Delta, m \vdash \underline{?x} :: T} \text{G-VarClose}$$

We can instantiate a hole with a (higher-order) library component.

$$\frac{\begin{array}{l} |\Xi| \leq P \text{ and } m < M \\ c = t_c : T_c \in \Delta \\ \sigma \text{ unifies } T \text{ with fresh}(T_c) \end{array}}{\Xi, \Phi, \Delta, m \vdash ?x :: T \Rightarrow \sigma(\Xi \setminus \{?x : T\}), \Phi, \Delta, m + 1 \vdash c :: \sigma(T)} \text{G-VarLib}$$

We can instantiate a hole with a function application of two fresh holes.

$$\begin{array}{c}
|\Xi| < P \text{ and } m \leq M \\
?X \text{ is a fresh type hole, } ?x_1 \text{ and } ?x_2 \text{ are fresh term holes} \\
\frac{\Xi' = \Xi \setminus \{?x : T\} \cup \{?x_1 : ?X \rightarrow T, ?x_2 : ?X, ?X\}}{\Xi, \Phi, \Delta \vdash ?x :: T \Rightarrow \Xi', \Phi, \Delta \vdash ?x_1 ?x_2 :: T} \text{G-VarApp}
\end{array}$$

We can expand one of the holes of the program according to one of the three rules above.

$$\frac{\Xi, \Phi, \Delta, m \vdash ?x :: T_1 \Rightarrow \Xi', \Phi, \Delta, m' \vdash t'_1 :: T'_1}{\Xi, \Phi, \Delta, m \vdash t[?x] :: T \Rightarrow \Xi', \Phi, \Delta, m' \vdash t[t_1] :: [T_1 \mapsto T'_1]T} \text{G-App}$$

Chapter 4

Implementation

What could I talk about in this chapter?

- Programming language and compiler version
- put the type definitions and explain them (What are Fun, FUN and BuiltinFun) (built-in integers for speed)
- Library syntax and the type-checking when added to the library?
- eager evaluation, describe evaluator
- Table of implemented components

Chapter 5

Evaluation

The main goal of this chapter is to give some insights about the factors that affect performance and compare different variants of the synthesis procedure described in Chapter 3. The chapter also puts our synthesis procedure in relation with the related work discussed in Chapter 2.

5.1 Experimental set up

This section presents the set up of the two experiments we are going to discuss in the rest of the chapter.

The goal of the first experiment is to assess the quality and the performance of the synthesiser on standard benchmarks. The detailed set up is described in Section 5.1.1 and the results are discussed in Section 5.2.

In the second experiment the synthesiser is used to automatically generate a *black list* that can be successively be used to prune the search space. We refer back to Section 3.5 for a description of pruning based on black lists. Section 5.1.2 describes how we used the synthesis procedure to generate a black list and Section 5.4 reviews the quality of the generated black list.

5.1.1 Evaluation on benchmarks

TODO: Insert giant table. Refer to the right table, make sure you describe only columns that are really there

Make and insert mentioned figures: manual black list

We evaluated 9 variants of our synthesis procedure, crossing the 3 exploration strategies with 3 of the cost functions described in Chapter 3. The three exploration strategies we evaluated are the following.

plain implements the basic synthesis procedure based on best first search described in Section 3.3.2.

blacklist implements the pruning of the search space based on a manually compiled black list provided in Table ?? . We refer to Section 3.5 for more details.

template implements the double best first search introduced in Section 3.6. As you probably recall, the procedure first looks for a *template* featuring at most *nof_comp* higher-order components and at most *nof_hol* holes and as soon as such a template is found the procedure falls back on the *plain* variant up to a certain depth using only the first-order components.

For each exploration strategy, we instantiated the cost function with 3 of the cost functions described in Section 3.4, that is with *nof-nodes*, *nof-nodes-simple-type* and *no-same-component*. We refer back to the corresponding section for more details.

We exercised the 9 different variants of our synthesis procedure on a benchmark of 23 programs over lists, mostly taken from related work or standard functional programming assignments. Table ?? and Table ?? list the benchmarks along with the specification size, that is the number of input-output examples we used to generate them, and the number of library components we provided to the synthesiser. The other columns report for each variant the running time, the ratio to the minimum running time for that benchmark and the size of the generated solution expressed in number of nodes.

All experiments were run on an Intel quad core 3.2 GHz with 16 GB RAM. Since the code is sequential, the performance could not benefit from the number of cores. The performance numbers are averages from 1 to 3 different executions all sharing the same specification, that is the goal type, the given examples and the set of components do not change between different executions.

In Table ?? all benchmarks except for *nth* share the same set of components, from which we took out the benchmark to synthesise, if it was one of the components. In Table ??, in order to meet the need of all benchmarks, we used four different sets of 19 components.

Programs are enumerated only up to a timeout based on the number of programs that have been analysed so far. For the exploration strategies **plain** and **blacklist** the execution had been stopped after examining 2500000 programs (with or without holes). The exploration strategy **template** was restricted to generate templates with at most 2 higher-order components and at most 5 holes, the depth of the first-order search was limited to 10 calls to the *plain* procedure. For the cost function *nof-nodes* this corresponds to circa 4 min.

5.1.2 Automatic black list generation

We also used our system to generate an automatic black list based on the identity function. Using the **plain** exploration strategy and the *nof-nodes* cost function, we generated 100 programs representing the identity functions over integers, 100 over lists of integers and 100 over lists of lists of integers. The examples were also generated automatically with the **plain** exploration strategy and the *nof-nodes* cost function providing only constructors as library components.

We chose not to generate the polymorphic identity function. As during pruning we are ignoring types, holes and input variables, the programs that would have been generated for the polymorphic identity function are also generated for the identity over any specific type. Table ?? summarizes 30 typical programs.

5.2 Factors affecting runtime

Two variants of our synthesis procedure were able to synthesis all 23 benchmarks in the presence of 37 library components, all variants synthesised at least 15 benchmark programs within the time limit. 78% of the benchmarks were synthesised within 1 s using **blacklist** as the exploration strategy and *nof-nodes-simple-type* as the cost function.

The search space is of exponential nature and depends on many factors: most notably the number of library components and the size of the solution to be synthesised. In the remainder of this section we look at these and other factors and their influence on the runtime.

5.2.1 Size of the solution

TODO: If you have time, redo the graphics in latex, or at least find a way to move trend line in the legend Figure 5.1 shows that the average running time for all nine variants of the synthesis procedure depends exponentially on the number of nodes. This goes along with the intuition that a bigger program is more difficult to synthesise. For example, if we have n possibilities to generate a program consisting of one node, that is $?x$ where we have n possibilities to instantiate the hole $?x$, then we will have n^2 possibilities to generate a program with three nodes, that is $?x_1 ?x_2$ where we have n possibilities to instantiate the hole.

TODO: Delete this question or do something with it. We cannot answer it, we do not have enough data. More surprisingly, the individual results show that there must be some other factor influencing the runtime. Take, for example, `enumFromTo`, `stutter` and `nth`. All three of them have a solution with exactly 13 nodes, but their runtimes differ at least by an order of magnitude.

What does make `nth` generate in less than 1s, stutter a hundred times slower and `enumFromTo` to time out in most of the cases?

In our simple intuitive explanation of the exponential dependency of the synthesis time with the size of the solution we completely ignored the contribution of types to search space pruning.

5.2.2 Number of components

TODO: Write this

Discuss that `enumFromTo` can be generated with 19 but not with 37 components even if we provide `enumTo`. That providing only the components it needs, `enumFromTo` can be synthesised pretty quickly also without providing `enumTo`.

How number of components affects synthesis time

In particular, not only the number but also the number of components of the same type. If you have more functions with the same type you will need much more time to find the program you are looking for (more possible successors).

5.2.3 Examples

Another factor that greatly impacts on performance is the choice and the number of provided input-output examples. As our procedure evaluates every closed program it synthesises on at least the first input-output example, we must make sure that the first input-output example is

- a. small enough, so that also undesirable programs like `enumTo (prod (enumTo (prod xs)))` do not get stuck or run out of memory trying to construct a list with 479001600 elements, which happens already for the at first sight innocent input `[2,2,3]`.
- b. expressive enough to rule out many programs, so that there is no need to fall back on the other, often bigger, input-output examples.

Clearly, using as few and as small input-output examples as possible has a positive effect on performance. On the other side, too few and too general input-output examples can lead to the synthesis of the wrong program, that is a program that satisfies all provided input-output examples but that does not generalise in the expected way. This was especially a problem with `enumFromTo` and `member`. **TODO: Put a concrete example with `enumFromTo`**
`con Int _0 (drop Int _0)enumTo _1))` well, this is a real solution
`filter Int (b_leq _0)(con Int b_zero (enumTo _1))`
 On input `enumFromTo 1 3 = [1,2,3]` and `enumFromTo 2 4 = [2,3,4]` we get the program `con Int _0 (con Int (b_succ _0)(con Int _1 (nil Int)))`

for (1,2) and (1,3) we get `enumTo _1` or `b_foldNatNat (List Int)(con Int)(nil Int)_1` depending on the components we give.

for (1,2) and (2,3) we get `con Int _0 (con Int _1 (nil Int))`

for (1,2) and (2,4) we get `con Int _0 (map Int Int (b_add _0)(enumTo _0))`

5.2.4 Blacklist

TODO: Make and insert mentioned figures: the manual black list

The search space abounds of superfluous programs that are equivalent to smaller ones. In Section 3.5 we introduced a way to leverage this inconvenience: Pruning based on black lists. This approach allows us not to explore further programs that will surely lead to a solution bigger than the optimal one, like `append [X] (nil [X])?xs`, or not lead to a solution at all, like `(head [?X1 → ?X2 → X3] (nil [?X1 → ?X2 → X3])) ?x1 ?x2`.

A longer black list allows to prune more superfluous programs and sinks considerably the number of programs our synthesis procedure needs to consider before finding a solution. However, black list pruning is extremely expensive in our implementation. Each element of the black list is matched against every subtree of every program with holes that is generated. That is, there is a trade-off between the length of the black list and the gain in performance that we can get.

Figure ?? shows the black list we used to evaluate the benchmarks. We manually compiled it combining unwanted patterns often seen in the search space with some carefully chosen automatically generated identity functions. We also added some extremely increasing functions like `foldNat [Int] (foldNat [Int] (mul n)1)1 m` that represented a problem for our evaluator.

TODO: If you have time, try to figure out what mathematical function it corresponds to.

In Table ?? we see that the runtime profits the most from the introducing of black list pruning when we use the cost function *nof-nodes*. The running time drops less significantly if we use other cost functions. A possible explanation of this behaviour could be the fact that other cost functions give a higher cost to those programs that are filtered with our black list.

We could also empirically see that pruning using black lists is very helpful in the presence of "useless" functions that increase the branching factor, for example `flip`, `const` or `uncurry`. Forbidding a fully applied `flip`, `const` or `uncurry` has a comparable effect on performance to taking those components out of the library. However, since we are not taking those components out of the library, we are still able to synthesis functions that need them.

5.2.5 Templates

TODO: Write this

Explain why templates perform so poorly (hm... because they don't do BFS, that is they have to explore every branch to the "end"?)

Short section explaining that the templates you generate are not the templates you expect and why I found one example, where templates help! For `dropmax` I had a run out of memory exception with plain enumeration and I could synthesise it in 7 seconds with templates! (Ok, I modified templates to put in only really higher-order components and close every hole, no matter the type). More resistant to the choice of examples.

5.2.6 Cost functions

TODO: Write this

How cost functions influence the search space (`(head (^26 → ^27 → Int) (nil (^26 → ^27 → Int)))?2 ?3` has too complicated types appearing in the term, so it will have a higher cost in some of the cost functions. `foldnat foldnat mul (mul 3 3)3` takes a lot of time to evaluate, it will have a higher cost in no-same-component).

- Idea behind `nof-nodes`: smaller programs generalize better to the examples
- Idea behind `nof-nodes-simple-types`: programs with smaller types are less "useless" (example with `head nil`). Favor the use of input variables. But now a lot of "simple" but difficult to evaluate programs are synthesised, like the `enumTo prod enumTo prod`.
- Idea behind `no-same-component`: there are components that are usually used only once in a program, like `foldr`, `foldNat`. Another thing is that we penalize also complicated types of the form `List (List (^5 → List ^4))` that had no additional cost in `nof-nodes-simple-types`. But now types overweight the number of nodes and simple programs like `foldr add zero _0` weight more than programs like `add mul sub prod enumTo...`
- Idea behind `no-same-component-bigger-constants`: make all constants bigger so that we can differentiate more between the constants and so that nodes count more than types.
- Idea behind `no-same-component-even-bigger-constants`: see above. (well, then you failed. your constants for types are pretty much as big as those for terms...). Why it's bad? See above.

5.2.7 Stack vs Queue expansion

As already mentioned in Section 3.3.1, we have two open questions in our best first search:

- a. what program to expand next
- b. which hole of this program to expand first

In the previous section we addressed the first question with different cost functions. In this section we focus on the second one.

Among all possible heuristics to determine which hole of the least-cost program to expand next, we chose to discuss two. In the first one the open holes of a program are organised in a stack, as opposed to the second one, where the open holes are kept in a queue.

Organising the open holes of a program into a stack leads to the expansion of the holes from left to right. To give some intuition we provide a derivation of `mapAdd` featuring the stack of open holes on the right of the program, where `xs` represents the input list and `n` the amount of the increment.

```
(?x0, [x0]) →
(?x1 ?x2, [?x1, ?x2]) →
(?x3 ?x4 ?x2, [?x3, ?x4, ?x2]) →
(map [Int] ?x4 ?x2, [?x4, ?x2]) →
(map [Int] (?x5 ?x6) ?x2, [?x5, ?x6, ?x2]) →
(map [Int] (add ?x6) ?x2, [?x6, ?x2]) →
(map [Int] (add n) ?x2, [?x2]) →
(map [Int] (add n) xs, [])
```

Left-to-right expansion often lead to faster synthesis, because leftmost holes have usually more constraints on their type. Consider the program `?x3 ?x4 ?x2` from the derivation of `.` We know more about `?x3` than about `?x2`: The first one must be a function that takes two arguments of some type and returns a list of integers, whereas the second hole could be anything. Furthermore, the instantiation of `?x3` with `map [Int]` imposes some constraints on the types of `?x4` and `?x2`.

Keeping the open holes of a program in a queue leads to the expansion of the hole with the smallest depth first. This could be useful to control the depth of a program, but in practice it has a substantial drawback. Consider again the derivation of `mapAdd`. The first three steps are the same, but in the program `?x3 ?x4 ?x2` we would now try to expand the hole `?x2`, that we have absolutely no information about. Every library component and every input variable are valid instantiations of this hole. Thus, this expanding strategy leads to a higher branching factor and explores many useless programs like `?x3 ?x4 map` and `map (?x5 foldr) xs`.

We used the stack-based expansion strategy throughout all runtime evaluations of the benchmarks.

5.3 Synthesised solutions

TODO: Rewrite programs in normal style

Most of the synthesised solutions are precisely the ones we would have written by hand. Interestingly, for some programs different variants of the synthesis procedure found two different valid programs of the same size. For example, for `replicate` we find both `map Int &0 (const &0 Int _0)(enumTo _1)` and `b_foldNat (List &0)(con &0 _0)(nil &0)_1`.

Sometimes, as for example for `length`, two solutions of different sizes are found: `foldr &0 Int (const (Int → Int)&0 b_succ)b_zero _0` and `sum (map &0 Int (const Int &0 (b_succ b_zero))_0)`. The reason is that the cost function *nof-nodes-simple-type* assigns a higher cost to the first solution, as it contains an arrow type.

For the few benchmarks that can use `foldl` and `foldr` interchangeably, like `concat`, `maximum` and `sum`, different variants find different programs. The different variants of the synthesis procedure do not show clear preference for the one or the other. They can use `foldr` for one of such programs and `foldl` for the other.

Interesting is the case of `multfirst` and `multlast`, where we find `map &0 &0 (const &0 &0 (head &0 _0))_0` and `replicate &0 (length &0 _0)(head &0 _0)` for `multfirst` and analogous programs for `multlast`. The different synthesis procedures show a clear preference for the one or the other. For example, all synthesis procedures using the **templates** exploration strategy seem to prefer the use of the higher-order `map` to the first-order `replicate`. This has to do with the depth of the first-order search. The search from a particular template (in this case the template with no higher-order functions) times out before reaching programs with three components. On the other hand, the search starting from the template `map ^0 &0 (const ^0 &0 ?1)?2` succeeds within the timeout.

The preference of the **template** exploration strategy for solutions containing higher-order functions leads to unexpectedly large programs. For example, one of the solutions for `factorial` is `prod (foldr (List Int)(List Int)(foldl (List Int)Int (const (List Int)Int))(enumTo _0)(nil (List Int)))` instead of the much simpler `prod (enumTo _0)`. Note that since we are folding over an empty list, the two programs are completely equivalent.

There is a tendency to represent the constant integer 1 as `prod (nil Int)` instead of `succ zero`. And even if we forbid with a black list the use of `enumFromTo (succ zero)_` and `prod nil`, the synthesis procedure with

the **blacklist** exploration strategy still finds a way to express `enumTo` using `enumFromTo`: It simply falls back to `enumFromTo (b_div _0 _0)_0`. Of course, if we take the component `enumFromTo` out of the library we can generate the desired program `b_foldNatNat (List Int)(con Int)(nil Int)_0`.

Small programs are not always efficient. For example, for `enumFromTo` we find the solution `con Int _0 (b_foldNat (List Int)(tail Int)(enumTo _1)_0)`, that corresponds to generating a list from 1 to the second input and then dropping the first part of the list. The more efficient solution `con Int _0 (map Int Int (add _0)(enumTo (sub _1 _0)))` is larger and thus it is generated only if we take `foldNat` out of the library.

Sometimes the solution found by the synthesiser suggested other benchmarks we could try to synthesise. For example, after examining the aforementioned solution for `enumFromTo` we realized that `drop` can be implemented as `b_foldNat (List &0)(tail &0)_1 _0`. Analogously, a “wrong” solution for `member` turned out to be a clever implementation of `isEven`, namely `foldNat not true _0`. Folding over an integer is similar to recursion over that integer. In this case the base case of the recursion is `true` and in the inductive case we negate `isEven (n-1)`.

Another unexpectedly clever solution is due to the absence of a polymorphic equality function. We only had equality over integers, thus the benchmark `member` is not polymorphic but has the type `Int → List Int → Bool`. This allowed the synthesiser to generate, along with the expected solution `not (is_nil Int (filter Int (b_eq _0)_1))`, a special version that makes use of the fact that the product of a list is 0 if the list contains at least one 0 and that two numbers are equal if their difference is 0: `b_is_zero (prod (map Int Int (b_sub _0)_1))`. This works in our implementation because our built-in integers can have negative values.

5.4 Automatic black list

We were able to automatically synthesise 300 programs corresponding to the identity function for three different types using automatically generated input-output examples in under 2s. Because of their incremental nature, we need 8 automatically synthesised input-output examples as opposed to the 2-3 manual ones that would have been enough. Below that number there is no list of lists of integers that contains something but `nil` or no list of length two. This implies that synthesis using automatically synthesised input-output examples is slower than synthesis using manual ones.

For the evaluation of the benchmarks we preferred compiling a manual black list mainly because of three reasons:

1. All programs in the automatically generated black list correspond to the identity function.

2. Many automatically generated black list programs are unnecessary. For example, `append (append nil nil) _` and `concat (append nil _)` are not needed if the black list already contains `append nil _`.
3. The automatically generated programs are all closed programs and as such they are too concrete. For example, instead of `foldNatNat max _ zero`, `foldNatNat const _ zero` and `foldNatNat drop _ zero` we could just have the one program `foldNatNat _ _ zero` that generalises all the programs with the idea that folding over the integer 0 is the same as taking the initial value, no matter which function is used for folding.

The first two points can be addressed with small modifications to the experimental set up: generate `nil`, `zero`, `undefined` and other constants as well for the first and prune the black list after or during generation for the second. The third point is way more complex. Partial evaluation of programs with holes could help to some extent, but at the end it's about the ability to abstract and generalise over programs.

5.5 Comparison to related work

TODO: Put the table in LaTeX

TODO: Write this Comparison to Feser, to Nadia, to Escher and to Myth. Say that you cannot really compare the numbers because they weren't run on the same machine. Say the thing with the size of the specification and the number of components. Say you have components capturing (tail?) recursive behaviour (fold over lists and fold over integers).

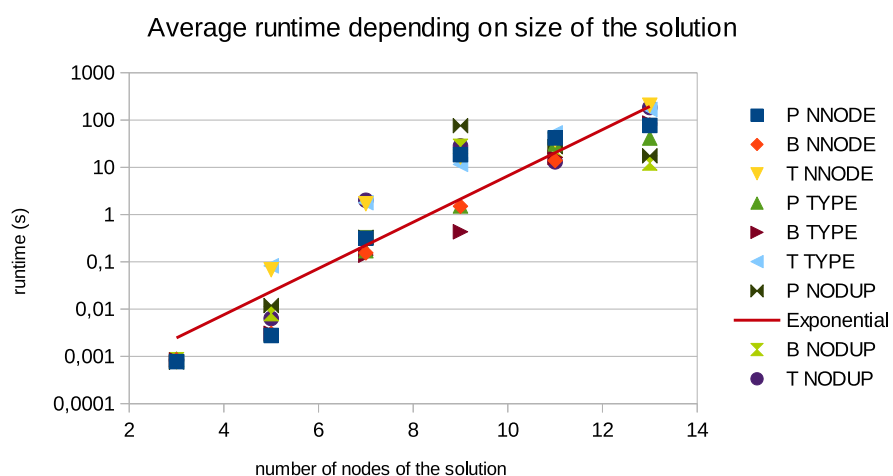


Figure 5.1: Average running time of the variants of the synthesis procedure depending on the number of nodes of the solution.

Conclusions

6.1 Conclusions

The baseline is not that bad. Gathered some data about the search space. Incremental development (synthesise `enumTo` before synthesising `enumFromTo` and use it as a component).

6.2 Future Work

Templates done well, augmented examples?

Bibliography

- [1] M Abadi, L Cardelli, and G Plotkin. Types for the scott numerals, 1993.
- [2] Aws Albarghouthi, Sumit Gulwani, and Zachary Kincaid. Recursive program synthesis. In *Proceedings of the 25th International Conference on Computer Aided Verification, CAV'13*, pages 934–950, Berlin, Heidelberg, 2013. Springer-Verlag.
- [3] John K. Feser, Swarat Chaudhuri, and Isil Dillig. Synthesizing data structure transformations from input-output examples. In *Proceedings of the 36th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2015*, pages 229–239, New York, NY, USA, 2015. ACM.
- [4] Gérard P. Huet. Unification in typed lambda calculus. In *Proceedings of the Symposium on Lambda-Calculus and Computer Science Theory*, pages 192–212, London, UK, UK, 1975. Springer-Verlag.
- [5] Peter-Michael Osera and Steve Zdancewic. Type-and-example-directed program synthesis. In *Proceedings of the 36th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2015*, pages 619–630, New York, NY, USA, 2015. ACM.
- [6] B.C. Pierce. *Types and Programming Languages*. MIT Press, 2002.
- [7] Nadia Polikarpova and Armando Solar-Lezama. Program synthesis from polymorphic refinement types. *CoRR*, abs/1510.08419, 2015.



Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

Declaration of originality

The signed declaration of originality is a component of every semester paper, Bachelor's thesis, Master's thesis and any other degree paper undertaken during the course of studies, including the respective electronic versions.

Lecturers may also require a declaration of originality for other written papers compiled for their courses.

I hereby confirm that I am the sole author of the written work here enclosed and that I have compiled it in my own words. Parts excepted are corrections of form and content by the supervisor.

Title of work (in block letters):

Authored by (in block letters):

For papers written by groups the names of all authors are required.

Name(s):

First name(s):

With my signature I confirm that

- I have committed none of the forms of plagiarism described in the '[Citation etiquette](#)' information sheet.
- I have documented all methods, data and processes truthfully.
- I have not manipulated any data.
- I have mentioned all persons who were significant facilitators of the work.

I am aware that the work may be screened electronically for plagiarism.

Place, date

Signature(s)

For papers written by groups the names of all authors are required. Their signatures collectively guarantee the entire content of the written paper.