

Functional Programs as Compressed Data

Naoki Kobayashi

Tohoku University
koba@ecei.tohoku.ac.jp

Kazutaka Matsuda

Tohoku University
kztk@kb.ecei.tohoku.ac.jp

Ayumi Shinohara

Tohoku University
ayumi@ecei.tohoku.ac.jp

1

Abstract

We propose an application of programming language techniques to lossless data compression, where tree data are compressed as functional programs that generate them. This “functional programs as compressed data” approach has several advantages. First, it follows from the standard argument of Kolmogorov complexity that the size of compressed data can be optimal up to an additive constant. Secondly, a compression algorithm is clean: it is just a sequence of β -expansions for λ -terms. Thirdly, one can use program verification and transformation techniques (higher-order model checking, in particular) to apply certain operations on data without decompression. In the paper, we present algorithms for data compression and manipulation based on the approach, and prove their correctness. We also report preliminary experiments on prototype data compression/transformation systems.

1. Introduction

Data compression plays an important role in today’s information processing technologies. Its advantages are not limited to the decrease of data size, which enables more data to be stored in a device. Recent computer systems have a large memory hierarchy, from CPU registers to several levels of cache memory, main memory, hard disk, etc., so that decreasing the data size enables more data to be stored in a higher level of the hierarchy, leading to more efficient computation. Some data compression schemes allow various operations to be performed without decompression in time polynomial in the size of the compressed data, so that one can sometimes achieve super-polynomial speed-up by compressing data. Data compression can also be applied to knowledge discovery [10].

In this paper, we are interested in the (lossless) compression of string/tree data as *functional programs*. The idea of “programs as compressed data” can be traced back at least to Kolmogorov complexity [20, 21], where the complexity of data is defined as the size of the smallest program that generates the data. The use of the λ -calculus in the context of Kolmogorov complexity has also been studied before [37]. Despite the generality and potential of the “functional programs as compressed data” approach, however, it did not seem to have attracted enough attention, especially in the programming language community.

The goal of the present paper is to show that we can use programming language techniques, program verification/transformation techniques in particular, to strengthen the “functional programs as compressed data” approach, so that the approach becomes not only of theoretical interest but potentially of practical interest. The approach has the following advantages.

¹This is an extended version of the paper that appeared in Proceedings of ACM SIGPLAN 2012 Workshop on Partial Evaluation and Program Manipulation (PEPM’12).

1. Generality and optimality: In principle, it subsumes arbitrary compression schemes. Imagine some compression scheme and suppose that w is the compressed form of data v in the scheme. Let f be a functional program for decompression. Then, v can be expressed as a (closed) functional program $f w$. This is larger than w only by a constant, i.e. the size of the program f . This is actually the same as the argument for Kolmogorov complexity. We use a functional language (or more precisely, the λ -calculus) instead of a universal Turing machine, but it is easy to observe that the size of (a certain binary representation of) a λ -term representing the original data can be optimal with respect to Kolmogorov complexity, up to an additive constant.

We can also *naturally* mimic popular compression schemes used in practice. For example, consider the run-length coding. The string “abaabaababbbb” can be compressed as [3, “aba”, 4, “b”], meaning that the string consists of 3 repetitions of “aba” and 4 repetitions of “b”. This can be expressed as:

`(repeat 3 "aba" (repeat 4 "b" ""))`

where `repeat` is a function that takes a non-negative integer n and strings s_1 and s_2 , and returns the string $s_1^n s_2$. For another example, consider grammar-based compression, where strings or trees are expressed as (a restricted form of) context-free (tree) grammars [4, 13, 22]. The grammar-based compression has recently been studied actively, and used in practice for compression of XML data [4]. For instance, consider the tree shown in Figure 1 (which has been taken from in [4]). It can be compressed as the following tree grammar:

$$S = B(B(A)) \quad A = c(a, a) \quad B(y) = c(A, d(A, y))$$

Using the λ -calculus, we can express it by:

`let A = c a a in let B = $\lambda y. c A (d A y)$ in B(B(A))`

where the sharing of the tree context $\lambda y. c(A, d(A, y))$ is naturally expressed by the λ -term. The data compression by a common pattern extraction then corresponds to an inverse β -reduction step. The previous grammar-based compression uses context-free grammars and their variants, while the λ -calculus has at least the same expressive power as higher-order grammars. Thus, as far as data compression is concerned, our approach can be considered a higher-order extension of the grammar-based compression and can achieve a theoretically higher compression ratio.

2. Data manipulation without decompression: Besides the compression ratio and the efficiency of the compression/decompression algorithms, an important criterion is what operations can be directly applied to compressed data without decompression. In fact, the main strength of the grammar-based approach [4, 13, 23, 24, 29] is that a large set of operations, such as pattern matching and string replacement, can be performed without decompression. It is particularly important when the size of original data is too large to fit into memory, but the size of the compressed data is small enough. As we show in the present paper, the “functional programs as compressed data” approach also enjoys such a property, by using program verification and transformation techniques. For example, consider a

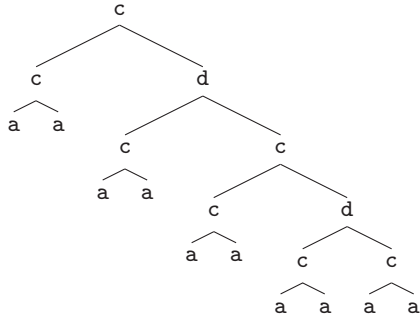


Figure 1. A tree

query to ask whether a given tree T matches a given pattern P . It can be rephrased as the query of whether the given tree T is accepted by a tree automaton M_P . Then, the problem of answering the query without decompression is: “Given a functional program p , is the tree generated by p accepted by M_P ?”. If p is a simply-typed program, then this is just an instance of higher-order model checking problems [15, 16, 26].

Pattern matching should often return not just a yes/no-answer, but extra information such as the position of the first match and the number of occurrences. Such operations can be expressed by transducers. Thus, the problem of performing such operations without decompression can be formalized as follows (see also the diagram in Figure 2):

“Given a tree transducer f and a functional program p that generates a tree T , construct a program q that generates $f(T)$.”

Thanks to the “functional programs as compressed data” approach, the construction of the program q is trivial: $q = \hat{f}(p)$, where \hat{f} is a representation of transducer f as a functional program. Of course, $\hat{f}(p)$ may not be an ideal representation, both in terms of the size of the program and the efficiency for further transformations. Fortunately, p is a tree generator and \hat{f} is a consumer, so that we can apply the standard fusion transformation [8] to simplify $\hat{f}(p)$. An alternative, more sophisticated approach is, as discussed later, to extend a higher-order model checking algorithm to directly construct q .

3. Applications to knowledge and program discovery: This is a more speculative advantage. It is folklore that compressed data contains the essence of the data, hence knowledge can be discovered by compressing data to the extreme [10]. As we already discussed, the use of functional programs allows us to compress data to the limit (up to a constant factor), so that we may be able to extract knowledge, represented in the form of a program, by compressing data. In fact, consider the following Church numeral representation of 9: $\lambda s.\lambda z.s(s(s(s(s(s(s(s(z))))))))$. Our prototype compressor for λ -terms produces:

$$(\lambda n.\lambda f.n(nf))(\lambda s.\lambda x.s(s(s(x)))).$$

The part $\lambda s.\lambda x.s(s(s(x)))$ is the Church numeral 3, and the part $\lambda n.\lambda f.n(nf)$ is a square function for Church numerals. Thus, the equation $3^2 = 9$ and the square function have been automatically discovered by compression.

In the rest of this paper, we first introduce the λ -calculus as the language for expressing compressed data, and discuss the relationship with Kolmogorov complexity in Section 2. We then describe an algorithm for compressing trees as λ -terms in Section 3. In Section 4, we extend and apply program verification/transformation techniques to achieve processing of compressed trees (represented

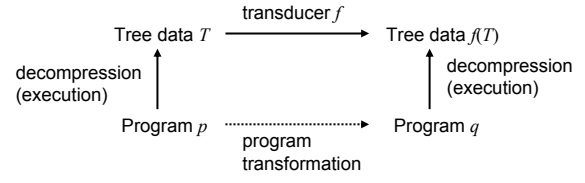


Figure 2. Applying transducers without decompression

in the form of λ -terms) without decompression. Section 5 reports preliminary experiments on data compression and processing. Section 6 discusses related work and Section 7 concludes.

The main contributions of this paper are: (i) Showing that *typed* λ -calculus with intersection types provides an optimal compression size up to an additive constant (Section 2.2). (ii) Developing an algorithm to compress trees as λ -terms (Section 3). (iii) Showing that higher-order model checking can be used to answer pattern match queries without decompression (Section 4.1). (iv) An extension of higher-order model checking to manipulate compressed data without decompression (Section 4.2). (v) Implementation and experiments on the algorithms for data compression and data manipulations without decompression (Section 5).

2. λ -Calculus as a Data Compression Language

2.1 Syntax

We use the λ -calculus for describing tree data and tree-generating programs. To represent a tree, we assume a *ranked alphabet* (i.e., a mapping from a finite set of symbols to non-negative integers) Σ . We write a, b, \dots for elements of the domain of Σ and call them *terminal symbols* (or just symbols). They are used as tree constructors below.

The set $Terms_\Sigma$ of λ -terms, ranged over by M , is defined by:

$$M ::= x \mid a \mid \lambda x.M \mid M_1 M_2.$$

Here, the meta-variables x and a range over variables and symbols (a, b, \dots) respectively. Note that, if symbols are considered free variables, this is exactly the syntax of the λ -calculus. As usual, λx is a binder for the variable x , and we identify terms up to the α -conversion. We also use the standard convention that the application $M_1 M_2$ is left-associative, and binds tighter than lambda-abstractions, so that $\lambda x.a x x$ means $\lambda x.((a x) x)$. We sometimes write $\text{let } x = M_1 \text{ in } M_2$ for $(\lambda x.M_2)M_1$.

The *size* of M , written $\#M$, is defined by:

$$\begin{aligned} \#x &= \#a = 1 & \#(\lambda x.M) &= \#M + 1 \\ \#(M_1 M_2) &= \#M_1 + \#M_2 + 1 \end{aligned}$$

The set of Σ -labeled trees, written \mathcal{T}_Σ , is the least subset of λ -terms closed under the rule:

$$\forall M_1, \dots, M_n \in \mathcal{T}_\Sigma. \Sigma(a) = n \Rightarrow a M_1 \dots M_n \in \mathcal{T}_\Sigma.$$

(Note here that n may be 0, which constitutes the base case.) We often use the meta-variable T to denote an element of \mathcal{T}_Σ .

If M has a β -normal form, we write $\llbracket M \rrbracket$ for it. In the present paper, we are interested in the case $\llbracket M \rrbracket$ is a tree (i.e. an element of \mathcal{T}_Σ). When $\llbracket M \rrbracket$ is a tree T , we often call M a *program* that generates T , or a program for T in short. The goal of our data compression is, given a tree T , to find a small program for T .

Example 2.1. Let T be $a^9(c)$, i.e., $\underbrace{a(a(\dots(a(c))\dots))}_9$.

It is generated by the following program M :

$$(\lambda n.n(n a c))(\lambda s.\lambda x.s(s(s(x)))).$$

Note that $\#T = 19 > \#M = 18$.

In general, the size of a program M for T can be hyper-exponentially smaller than the size of T . For example, consider the tree:

$$T := \underbrace{a^{2^{2^{\dots^2}}}}_n c$$

It is generated by: $M := (\lambda f. \underbrace{f f \dots f}_n a c)(\lambda g. \lambda x. g(g(x)))$.

Example 2.2. Consider the following term, which generates a unary tree $a^{57}(c)$.

```
let twice = λf.λx.f(f(x)) in
let b0 = λn.λs.λz.n (twice s) z in
let b1 = λn.λs.λz.n (twice s) (s z) in
let zero = λs.λz.z in
b1(b0(b0(b1(b1(b1(zero)))))) a c
```

The part $b_1(b_0(b_0(b_1(b_1(b_1(zero))))))$ corresponds to the binary representation 111001 of 57, with the least significant bit first. The last line can be replaced with:

```
let thrice = λf.λx.f(f(f(x))) in b1(twice b0(thrice b1(zero))) a c
```

The part $b_1(twice b_0(thrice b_1(zero)))$ then corresponds to the runlength coding of the binary representation 111001. \square

2.2 Typing

We considered the untyped λ -calculus above, but we can actually assume that any program that generates a tree is well-typed in the intersection type system given below. The assumption that programs are well-typed is important for the program transformations discussed in Section 4. The use of intersection types is important for guaranteeing that we do not lose any expressive power for expressing finite trees: see Theorem 2.1 below.

The set of (*intersection*) types, ranged over by τ , is given by:

$$\tau ::= \circ \mid \tau_1 \wedge \dots \wedge \tau_k \rightarrow \tau$$

Here, k may be 0, in which case we write $\top \rightarrow \tau$ for $\tau_1 \wedge \dots \wedge \tau_k \rightarrow \tau$. Intuitively, \circ describes a tree, and $\tau_1 \wedge \dots \wedge \tau_k \rightarrow \tau$ describes a function that takes an element having all of the types τ_1, \dots, τ_k , and returns an element of type τ . We sometimes write $\bigwedge_{i \in \{1, \dots, k\}} \tau_i \rightarrow \tau$ for $\tau_1 \wedge \dots \wedge \tau_k \rightarrow \tau$.

A *type environment* is a finite set of type bindings of the form $x : \tau$. Unlike ordinary type environments, we allow multiple occurrences of the same variable, like $\{x : \circ \rightarrow \circ, x : (\circ \rightarrow \circ) \rightarrow (\circ \rightarrow \circ)\}$. We often omit $\{\}$ and just write $x_1 : \tau_1, \dots, x_n : \tau_n$ for $\{x_1 : \tau_1, \dots, x_n : \tau_n\}$.

The type judgment relation is of the form $\Gamma \vdash M : \tau$ where Γ is a type environment. It is inductively defined by the following typing rules:

$$\frac{}{\Gamma, x : \tau \vdash x : \tau} \quad \frac{\Sigma(a) = k}{\Gamma \vdash a : \underbrace{\circ \rightarrow \dots \rightarrow \circ}_k \rightarrow \circ}$$

$$\frac{\Gamma, x : \tau_1, \dots, x : \tau_n \vdash M : \tau \quad x \text{ does not occur in } \Gamma}{\Gamma \vdash \lambda x. M : \tau_1 \wedge \dots \wedge \tau_n \rightarrow \tau}$$

$$\frac{\Gamma \vdash M_1 : \tau_1 \wedge \dots \wedge \tau_n \rightarrow \tau \quad \forall i \in \{1, \dots, n\}. \Gamma \vdash M_2 : \tau_i}{\Gamma \vdash M_1 M_2 : \tau}$$

It follows from the standard argument for intersection types [2, 39] that any program that generates a tree is well-typed. (Note that even a term like $(\lambda x.c)\Omega$ that contains the divergent term Ω is typed, by assigning $\top \rightarrow \circ$ to $\lambda x.c$.)

Theorem 2.1. *Let M be a λ -term. Then, $\emptyset \vdash M : \circ$ if and only if $\llbracket M \rrbracket$ is a tree.*

Proof. The “only if” part follows from the facts that (i) $\emptyset \vdash M : \circ$ implies that M has a β -normal form, (ii) if M is in β -normal form and $\emptyset \vdash M : \circ$, then M is a tree, and (iii) typing is preserved by β -reductions. The “if” part follows from the facts that (iv) if M is a tree, then $\emptyset \vdash M : \circ$, and (v) typing is preserved by β -expansions (i.e. the inverse of β -reductions). \square

As we consider only programs representing (finite) trees, thanks to the theorem above, we can safely assume that all the programs in consideration are well-typed (in the intersection type system above) in the rest of this paper. By using the standard argument for Kolmogorov complexity, we can show that our representation of data as (implicitly-)typed λ -terms can be optimal up to an additive constant with respect to Kolmogorov complexity [20, 21]: see Appendix A for more details. A reader not familiar with Kolmogorov complexity may also wish to consult [20, 21].

Example 2.3. Consider the term:

$$(\lambda \text{twice}. \text{twice twice } a c) \lambda f. \lambda x. f(f(x)).$$

The part $\lambda \text{twice}. \dots$ is given the following intersection type:

$$\begin{aligned} & (((\circ \rightarrow \circ) \rightarrow \circ \rightarrow \circ) \\ & \wedge (((\circ \rightarrow \circ) \rightarrow \circ \rightarrow \circ) \rightarrow (\circ \rightarrow \circ) \rightarrow \circ \rightarrow \circ)) \rightarrow \circ \quad \square \end{aligned}$$

2.3 Relationship with Grammar-based Compression

Grammar-based compression schemes, in which a string or a tree is expressed as a grammar that generates it, have been actively studied recently [4, 13, 22]. Our compression scheme using the λ -calculus can *naturally* mimic grammar-based compression schemes. For example, consider the compression scheme using context-free grammars (with the restriction of cycle-freeness) or straight-line programs. A string s is expressed as a grammar of the following form:

$$X_1 = e_1, X_2 = e_2, \dots, X_n = e_n,$$

where e_i is either a terminal symbol a , or $X_j X_k$ where $1 \leq j, k < i$. It can be expressed as

$$\begin{aligned} & \text{let } X_1 = \lambda y. e_1^{(y)} \text{ in let } X_2 = \lambda y. e_2^{(y)} \text{ in } \dots \\ & \text{let } X_n = \lambda y. e_n^{(y)} \text{ in } X_n(e), \end{aligned}$$

where $e^{(y)}$ is defined by: $a^{(y)} = a(y)$ and $(X_j X_k)^{(y)} = X_j(X_k(y))$. It generates s in the form of a linear tree, with e as an end-marker. We can also express various extensions of straight-line programs, such as context-free tree grammars [4] and collage systems [13] as λ -terms.

Example 2.4. Fibonacci words² are variations of Fibonacci numbers, obtained by replacing the addition $+$ with the string concatenation, and the first and second elements with b and a . The n -th word is expressed by the following straight-line program:

$$X_0 = b, X_1 = a, X_2 = X_1 X_0, \dots, X_n = X_{n-1} X_{n-2}.$$

It is encoded as:

$$\begin{aligned} & \text{let } X_0 = b \text{ in let } X_1 = a \text{ in let } X_2 = \lambda x. X_1(X_0(x)) \text{ in } \dots \\ & \text{let } X_n = \lambda x. X_{n-1}(X_{n-2}(x)) \text{ in } X_n(e) \end{aligned}$$

For $n = 2^m$, we have a more compact encoding:

$$\begin{aligned} & \text{let } \text{concat} = \lambda x. \lambda y. \lambda z. x(y(z)) \text{ in} \\ & \text{let } g = \lambda k. \lambda x. \lambda y. k y (\underbrace{\text{concat } y x}_m) \text{ in} \\ & \quad \text{twice}(\dots(\text{twice}(g))\dots) (\lambda x. \lambda y. x) b a e \end{aligned}$$

²Fibonacci words and its generalization called *Sturmian words* have been studied in a field called *Stringology* [6].

```

compressTerm(M) =
  let M1 = compressAsTree(M) in
  let M2 = simplify(M1) in
  if #M2 ≥ #M then M else compressTerm(M2)

```

Figure 3. Compression Algorithm for λ -terms

A similar encoding is also possible for an arbitrary number n by using b_0 and b_1 in Example 2.2. \square

3. Compression as β -Expansions

This section discusses a compression algorithm which, given a tree T , finds a small λ -term M such that $\llbracket M \rrbracket = T$. A naive method would be to enumerate all the λ -terms M smaller than T , and check whether $\llbracket M \rrbracket = T$. There are however obvious difficulties. First, the number of terms smaller than T is exponential in the size of T . Secondly, and even worse, it is in general undecidable to check whether $\llbracket M \rrbracket = T$. Even if we restrict M to the simply-typed λ -terms, the number of reductions required to compute $\llbracket M \rrbracket$ is in general non-elementary [33].

Here, we instead suggest reusing existing algorithms for (context-free) grammar-based tree compression [4, 22], by regarding a λ -term as a term tree (identified up to α -conversion) as follows.

$$x^\# = x \quad a^\# = a \quad (\lambda x.M)^\# = \lambda x \quad (MN)^\# = \begin{array}{c} @ \\ | \\ M^\# \quad N^\# \end{array}$$

As we have seen in Section 2.3, compressed data in the form of a context-free grammar can be easily translated to a λ -term. Thus, a grammar-based tree compression algorithm can be regarded as an algorithm for compression of λ -terms. By repeatedly applying such an algorithm to an initial tree T , we can obtain a small λ -term M such that $\llbracket M \rrbracket = T$. (There is, however, no guarantee that the resulting term is the smallest such M .) Note that the repeated applications are possible because the input and output languages for the compression algorithm are the same: the λ -calculus.

Figure 3 shows our algorithm, parametrized by two auxiliary algorithms: *compressAsTree* and *simplify*. Given a λ -term M (or a tree as a special case), we just invoke a tree compression algorithm to obtain compressed data in the form of λ -term M_1 . It is then simplified by using properties of λ -terms (such as the η -equality). We repeat these steps until the size of a term can no longer be reduced. (In the actual implementation, *compressAsTree* returns multiple candidates, which are inspected for further compression in a breadth-first manner. The termination condition $\#M_2 \geq \#M$ is also relaxed to deal with the case where the term size does not monotonically decrease: See Section 5.)

Because of the repeated applications of *compressAsTree*, we can actually use the following very simple algorithm for *compressAsTree*, which just finds and extracts a common tree context, rather than more sophisticated algorithms [4, 22]. Let us define a context with (up to) k -holes by:

$$C ::= []_1 \mid \dots \mid []_k \mid x \mid a \mid MC \mid CM \mid \lambda x.C$$

We write $C[M_1, \dots, M_k]$ for the term obtained by replacing each $[]_i$ in C with M_i . Note that ignoring binders, a context is just a tree context with up to k holes. Then, *compressAsTree* just needs to find (non-deterministically) contexts C_0, C_1, C_2, C_3 and terms $M_1, \dots, M_k, N_1, \dots, N_k$ such that

- (i) $M = C_0[C_1[C_2[M_1, \dots, M_k], C_2[N_1, \dots, N_k]]]$ or (ii) $M = C_0[C_1[C_2[M_1, \dots, M_k]]] \wedge M_i = C_3[C_2[N_1, \dots, N_k]]$; and

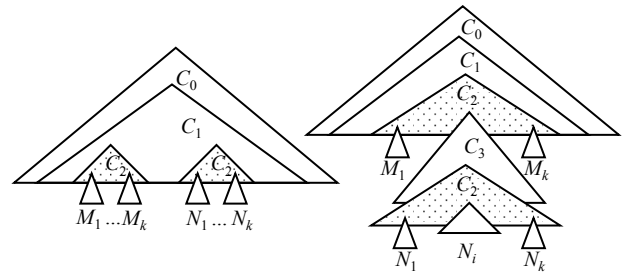


Figure 4. Cases where the common context C_2 occurs horizontally (left, case (i)), and vertically (right, case (ii))

- the free variables in $M_1, \dots, M_k, N_1, \dots, N_k$ are not bound in C_2 , and the free variables in C_2 are not bound in C_1 .

Here, (i) and (ii) are the cases where the common context C_2 occurs horizontally and vertically, respectively. The output is:

$$C_0[(\lambda f.C_1[f M_1 \dots M_k, f N_1 \dots N_k])(\lambda \tilde{x}.C_2[\tilde{x}])]$$

in case (i), and

$$C_0[(\lambda f.C_1[f M_1 \dots M_{i-1} M'_i M_{i+1} \dots M_k])(\lambda \tilde{x}.C_2[\tilde{x}])]$$

where $M'_i = C_3[f N_1 \dots N_k]$ in case (ii), and \tilde{x} denotes the sequence x_1, \dots, x_k . This transformation is a restricted form of β -expansion step: $C[\llbracket N/x \rrbracket M] \rightarrow C[(\lambda x.M)N]$, applicable only when M contains two occurrences of x .

The sub-procedure *compressAsTree* above is highly non-deterministic in the choice of contexts. In our prototype implementation, we pick every pair (M', M'') of subterms of M and find the maximum common context C_2 such that $M' = C_2[M_1, \dots, M_k]$ and $M'' = C_2[N_1, \dots, N_k]$. For splitting the enclosing context into C_0 and C_1 , we choose the largest C_1 that satisfies the condition on bound variables. The resulting procedure is still non-deterministic in the choice of the pairs (M', M'') , and our implementation applies the depth-first search. See Section 5 for more details.

For the simplification procedure *simplify*, we apply the following rules until no more rules become applicable.

- η -conversion: $\lambda x.M x \rightarrow M$
- β -reduction when the argument is a variable: $(\lambda x.M)y \rightarrow [y/x]M$
- β -reduction for linear functions: $(\lambda x.M)N \rightarrow [N/x]M$ where x occurs at most once in M .

Example 3.1. Consider a tree $\mathbf{a}^9(c)$. Let $C_0, C_1, C_2, C_3, M_1, N_1$ be:

$$C_0 = []_1 \quad C_1 = []_1 \quad C_2 = \mathbf{a}^3[]_1 \quad C_3 = []_1 \\ M_1 = \mathbf{a}^6(c) \quad N_1 = \mathbf{a}^3(c)$$

Then, case (ii) applies and the following term is obtained:

$$(\lambda f.f(f(\mathbf{a}^3(c))))\lambda x.\mathbf{a}^3(x)$$

Next, we again extract the common context $\mathbf{a}^3[]_1$, and obtain

$$(\lambda g.(\lambda f.f(f(g(c))))\lambda x.gx)(\lambda x.\mathbf{a}^3(x))$$

By using the η -equality $\lambda x.gx = g$, we get:

$$(\lambda g.(\lambda f.f(f(g(c))))g)(\lambda x.\mathbf{a}^3(x)).$$

As a part of the simplification procedure, we also β -reduce terms of the form $(\lambda x.M)y$ and obtain: $(\lambda g.g(g(g(c))))(\lambda x.\mathbf{a}^3(x))$. In the third iteration, we can extract the common context $[]_1([]_1([]_1[]_2))$

and obtain $(\lambda h.(\lambda g.h g c)(\lambda x.h a x))(\lambda f.\lambda x.f(f f x))$. By simplifying the term (by η -conversion and β -reduction for the linear function $\lambda g.h g c$), we obtain:

$$(\lambda h.(h (h a c)))(\lambda f.\lambda x.f(f f x)). \quad \square$$

Example 3.2. Recall the tree in Figure 1. By extracting the first two occurrences of the common context (with zero holes) $c a a$, we obtain: $(\lambda x.c x (d x (c (c a a) (d (c a a) (c a a))))(c a a)$. By further extracting the common context $c a a$ repeatedly (and applying *simplify*), we get $(\lambda C.c C (d C (c C (d C C)))(c a a)$. By extracting the common context a , we obtain

$$(\lambda C.c C (d C (c C (d C C)))(\lambda A.(c A A))a).$$

This corresponds to the DAG representation in Figure 1 of [23] and also to the regular grammar representation in Figure 1 of [4]. By extracting the common context $\lambda y.c C (d C y)$, the term is further transformed to:

$$(\lambda C.(\lambda B.B(B(C)))(\lambda y.c C (d C y)))(\lambda A.(c A A))a).$$

This corresponds to the sharing graph representation in Figure 1 of [23] and to the CFG representation in Figure 1 of [4]. \square

Relationship with CFG-based Tree Compression Algorithms

As demonstrated in Example 3.2, context-free grammar-based tree compression algorithms [4, 22] can be mimicked by our compression method based on λ -calculus. In fact, they may be viewed as a controlled and restricted form of our compression algorithm. For example, for efficient compression, Busatto et al. [4] impose restrictions on the number of holes and the size of common contexts, and also introduce certain priorities among subterms from which common contexts are searched. (There is also another difference that Busatto's algorithm finds more than two occurrences of a common context at once, but it can be mimicked by repeated applications of *compressAsTree* above.)

A more fundamental restriction of the previous approaches is that they [4] extract only common tree contexts with first-order types, of the form $\circ \rightarrow \dots \rightarrow \circ \rightarrow \circ$. Because of this difference, our compression algorithm based on the λ -calculus is more powerful than ordinary grammar-based compression algorithms. For example, the compression discussed in Example 3.1 is not possible with CFG-based compression: note that the final term cannot be expressed by a context-free tree grammar.

Limitations

The compression algorithm sketched above is not complete: there is a λ -term M and a tree T such that $\llbracket M \rrbracket = T$ but M cannot be obtained from T by the algorithm. For example, consider the following tree (represented as a term) T :

$$\text{br } (a (b (x (y (c (d e)))))) (a (b (z (c (d e))))))$$

It can be expressed by the following term M :

$$\text{let } f = \lambda g.a (b (g (c (d e)))) \text{ in } \text{br } (f \lambda u.x(y u)) (f z),$$

but M cannot be obtained by our algorithm. Note that $\lambda g.a (b (g (c d e)))$ is a higher-order context, having a tree context g as a hole. To enable the above compression, it suffices to β -expand T to:

$$\text{br } (a (b ((\lambda u.x(y u)) (c (d e)))))) (a (b (z (c (d e))))))$$

before applying our algorithm. Such pre-processing is however non-trivial in general.

For another example, consider the following tree (represented as a term) T :

$$\text{br } (a_1 (a_2 \dots (a_n e) \dots)) (a_n \dots (a_2 (a_1 (e))) \dots),$$

which consists of a linear tree representing the sequence a_1, \dots, a_n and its reverse.

The following term M generates T , but the common pattern h cannot be found by our algorithm.

$$\begin{aligned} &\text{let } h = \lambda a.\lambda k.\lambda x.\lambda y.k (a x) (\lambda z.y(a(z))) \text{ in} \\ &\text{let } id = \lambda z.z \text{ in} \\ &h a_n (\dots (h a_2 (h a_1 (\lambda x.\lambda y.\text{br } x (y e)))) \dots) e id \end{aligned}$$

Actually, finding terms h, k, x, y such that

$$\llbracket h a_n (\dots (h a_2 (h a_1 k)) x y) \rrbracket = T$$

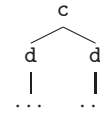
is an instance of the higher-order matching problem [34]. Thus, higher-order matching algorithms may be applicable to our data compression scheme, but it is left for future work.

4. Processing of Compressed Data

This section discusses how to process compressed data *without decompression*.

4.1 Pattern Matching as Higher-Order Model Checking

We first discuss the problem of answering whether $\llbracket M \rrbracket$ matches P , given a program M and a regular tree pattern P . For instance, we may wish to check whether some path from the root of the tree $\llbracket M \rrbracket$ contains ab as a subpath, or check whether $\llbracket M \rrbracket$ contains a subtree of the shape:



Such a pattern matching problem can be formalized as an acceptance problem for tree automata [5].

Below we write $\text{dom}(f)$ for the domain of a map f .

Definition 4.1 (tree automata). A (top-down, non-deterministic) tree automaton \mathcal{A} is a quadruple (Σ, Q, q_I, Δ) , where Σ is a ranked alphabet, Q is a set of states, q_I is the initial state, and $\Delta (\subseteq Q \times \text{dom}(\Sigma) \times Q^*)$ is a transition function, such that $(q, a, q_1 \dots q_n) \in \Delta$ implies $\Sigma(a) = n$. The reduction relation $S_1 \rightarrow S_2$ on subsets of $Q \times \mathcal{T}_\Sigma$ is defined by: $S \cup \{(q, a T_1 \dots T_n)\} \rightarrow S \cup \{(q_1, T_1), \dots, (q_n, T_n)\}$ if $(q, a, q_1 \dots q_n) \in \Delta$. A tree T is accepted by \mathcal{A} if $\{q_I, T\} \rightarrow^* \emptyset$. We write $\mathcal{L}(\mathcal{A})$ for the set of trees accepted by \mathcal{A} .

A tree automaton $\mathcal{A} = (\Sigma, Q, q_I, \Delta)$ is *deterministic* if, for each pair $(q, a) \in Q \times \text{dom}(\Sigma)$, there is at most one (q_1, \dots, q_n) such that $(q, a, q_1 \dots q_n) \in \Delta$.

Example 4.1. Let $\Sigma_1 = \{a \mapsto 1, b \mapsto 1, e \mapsto 0\}$. Consider the automaton $\mathcal{A} = (\Sigma_1, \{q_0, q_1, q_f\}, q_0, \Delta)$ where Δ is given by:

$$\Delta = \{(q_0, a, q_1), (q_0, b, q_0), (q_1, a, q_1), (q_1, b, q_f), (q_f, a, q_f), (q_f, b, q_f), (q_f, e, \epsilon)\}$$

Then, a Σ_1 -labeled tree T contains a subtree of the form $a(b(\dots))$ if and only if T is accepted by \mathcal{A} . \square

Example 4.2. Let $\Sigma_2 = \{b \mapsto 1, c \mapsto 2, d \mapsto 1, e \mapsto 0\}$. Consider the automaton $\mathcal{A}_2 = (\Sigma_2, \{q_0, q_1, q_f\}, q_0, \Delta_2)$ where Δ_2 is given by:

$$\Delta_2 = \{(q_0, b, q_0), (q_0, c, q_1 q_1), (q_0, c, q_0 q_f), (q_0, c, q_f q_0), (q_0, d, q_0), (q_1, d, q_f), (q_f, b, q_f), (q_f, c, q_f q_f), (q_f, d, q_f), (q_f, e, \epsilon)\}$$

Then, a Σ_2 -labeled tree T contains a subtree of the form $c(d \dots)(d \dots)$ if and only if T is accepted by \mathcal{A}_2 . \square

The goal here is, given a (well-typed) program M and an automaton \mathcal{A} , to check whether $\llbracket M \rrbracket \in \mathcal{L}(\mathcal{A})$ holds. A simple decision algorithm is to decompress M (i.e. fully β -reduce M) to a tree $T (= \llbracket M \rrbracket)$ and run the automaton \mathcal{A} for T . This is however

inefficient if T is large or the reduction sequence of M is long. Instead, we use the type-based technique for model checking higher-order recursion schemes [16, 38], to reduce $\llbracket M \rrbracket \stackrel{?}{\in} \mathcal{L}(\mathcal{A})$ to a type-checking problem for M . Because of subtle differences between higher-order recursion schemes [14, 26] and the language considered here (see Remark 4.1), we give a direct construction of the type system below.

Let $\mathcal{A} = (\Sigma, Q, q_I, \Delta)$ be a tree automaton. The set \mathbf{RTy}_Q of *refinement intersection types*, ranged over by θ , is given by:

$$\theta ::= q(\in Q) \mid \theta_1 \wedge \dots \wedge \theta_k \rightarrow \theta$$

Here, k may be 0. In $\theta_1 \wedge \dots \wedge \theta_k \rightarrow \theta$, we require that $\theta_i \neq \theta_j$ if $i \neq j$. Intuitively, q describes the set of trees accepted by M from the state q (i.e., accepted by (Σ, Q, q, Δ)). The type $\theta_1 \wedge \dots \wedge \theta_k \rightarrow \theta$ describes a function that takes an element of types $\theta_1, \dots, \theta_k$, and returns an element of type θ .

A *refinement type environment* Ψ is a finite set of type bindings of the form $x : \theta$, where multiple occurrences of the same variable are allowed as in the intersection type system in Section 2.2. We write $\text{dom}(\Psi)$ for the set $\{x \mid x : \theta \in \Psi\}$ of variables. The type judgment relation $\Psi \vdash_{\mathcal{A}} M : \theta$ is defined by:

$$\frac{}{\Psi, x : \theta \vdash_{\mathcal{A}} x : \theta} \quad \frac{(q, a, q_1 \dots q_k) \in \Delta}{\Psi \vdash_{\mathcal{A}} a : q_1 \rightarrow \dots \rightarrow q_k \rightarrow q}$$

$$\frac{\Psi, x : \theta_1, \dots, x : \theta_n \vdash_{\mathcal{A}} M : \theta \quad x \text{ does not occur in } \Psi}{\Psi \vdash_{\mathcal{A}} \lambda x. M : \theta_1 \wedge \dots \wedge \theta_n \rightarrow \theta}$$

$$\frac{\Psi \vdash_{\mathcal{A}} M_1 : \theta_1 \wedge \dots \wedge \theta_n \rightarrow \theta \quad \forall i \in \{1, \dots, n\}. \Psi \vdash_{\mathcal{A}} M_2 : \theta_i}{\Psi \vdash_{\mathcal{A}} M_1 M_2 : \theta}$$

Note that these typing rules are the same as those for the intersection type system in Section 2.2 except the rule for constants. The refinement type system is sound and complete for the problem in consideration.

Theorem 4.1. *Let M be a program and $\mathcal{A} = (\Sigma, Q, q_I, \Delta)$ be a tree automaton. Then, $\llbracket M \rrbracket \in \mathcal{L}(\mathcal{A})$ if and only if $\emptyset \vdash_{\mathcal{A}} M : q_I$.*

Proof. This follows from the facts that (i) for any tree T , $T \in \mathcal{L}(\mathcal{A})$ if and only if $\emptyset \vdash_{\mathcal{A}} T : q_I$, and that (ii) the typing is preserved by β -reduction and its inverse. \square

Suppose that a derivation for $\emptyset \vdash M : \circ$ is given. The result of Tsukada and Kobayashi ([38], Theorem 5) implies that to check whether $\emptyset \vdash_{\mathcal{A}} M : q_I$ holds, we just need to generate a *finite* set of candidates of derivation trees for $\emptyset \vdash_{\mathcal{A}} M : q_I$, and check whether one of them is valid. To state it more formally, we need to introduce some terminologies. The refinement relation $\theta :: \tau$ is defined by:

$$\frac{q \in Q \quad \theta :: \tau \quad \forall j \in \{1, \dots, m\}. \exists i \in \{1, \dots, k\}. \theta_j :: \tau_i}{q :: \circ \quad (\theta_1 \wedge \dots \wedge \theta_m \rightarrow \theta) :: (\tau_1 \wedge \dots \wedge \tau_k \rightarrow \tau)}$$

We extend the refinement relation to the relation on type environments by:

$$\Psi :: \Gamma \Leftrightarrow \forall x : \theta \in \Psi. \exists \tau. (x : \tau \in \Gamma \wedge \theta :: \tau).$$

Let π and π' be derivation trees for $\Psi \vdash_{\mathcal{A}} M : \theta$ and $\Gamma \vdash M : \tau$ respectively. π is a *refinement* of π' , if for each node labeled by $\Psi_1 \vdash_{\mathcal{A}} M_1 : \theta_1$ in π , there exists a corresponding node labeled by $\Gamma_1 \vdash M_1 : \tau_1$ in π' such that $\Psi_1 :: \Gamma_1$ and $\theta_1 :: \tau_1$. The following is the result of Tsukada and Kobayashi ([38], Theorem 5), rephrased for the language of this paper.

Theorem 4.2 ([38]). *If there are derivation trees π and π' respectively for $\emptyset \vdash M : \circ$ and $\emptyset \vdash_{\mathcal{A}} M : q_I$, then there exists a derivation tree π'' for $\emptyset \vdash_{\mathcal{A}} M : q_I$ such that π'' is a refinement of π .*

Let us define the *type size* of a judgment $\Gamma \vdash M : \tau$ by:

$$\begin{aligned} \#(\Gamma \vdash M : \tau) &= \#\Gamma + \#\tau \\ \#(x_1 : \tau_1, \dots, x_n : \tau_n) &= \#\tau_1 + \dots + \#\tau_n \\ \#\circ &= 1 \quad \#(\tau_1 \wedge \dots \wedge \tau_k \rightarrow \tau) = \#\tau_1 + \dots + \#\tau_n + \#\tau + 1 \end{aligned}$$

Define the *type width* of a derivation tree for $\Gamma \vdash M : \tau$ as the largest type size of a node of the derivation.

The following theorem follows immediately from Theorem 4.2 above.

Theorem 4.3. *Given an automaton \mathcal{A} and a type derivation tree for $\emptyset \vdash M : \circ$, $\llbracket M \rrbracket \stackrel{?}{\in} \mathcal{L}(\mathcal{A})$ can be decided in time linear in the size of M , under the assumption that the size of \mathcal{A} and the type width of derivation trees are bounded by a constant.*

Proof. Due to Theorems 4.1 and 4.2, it suffices to check whether there exists a derivation tree π for $\emptyset \vdash_{\mathcal{A}} M : q_I$ that is a refinement of the derivation tree π' for $\emptyset \vdash M : \circ$. Since the type width of π' is bounded by a constant, for each subterm N of M , the number of possible judgments that can occur in π is also bounded by a constant (although the constant can be huge). Thus, based on the refinement typing rules, we can enumerate all the valid judgments for N in time linear in the size of N : for example, to enumerate the typing for $M_1 M_2$, first enumerate valid typings for M_1 and M_2 and combine them by using the application rule. Thus, valid typings for M can also be enumerated in time linear in the size of M , and then it suffices to just check whether $\emptyset \vdash_{\mathcal{A}} M : q_I$ is among the valid judgments. \square

The fixed-parameter linear-time algorithm in the proof above is impractical due to the huge constant factor. We can instead use Kobayashi's fixed-parameter linear-time algorithm for higher-order model checking [17].³

Remark 4.1. Higher-order model checking [16, 26] usually refers to model checking of the tree generated by a *higher-order recursion scheme*, which can be considered a functional program. The only differences between the language of higher-order recursion schemes and our language are: (i) recursion schemes can be used to describe infinite trees, while our language is only used for describing finite trees, and (ii) recursion schemes must be simply-typed, but our language allows intersection types. Actually, our language can be considered a restriction of the extension of recursion schemes considered by Tsukada and Kobayashi [38].

4.2 Data Processing as Program Transformation

In the previous subsection, we considered pattern match queries to answer just yes or no. In practice, it is often required to provide extra information, such as the position of the first match and the number of matching positions. Computation of such extra information can be expressed as tree transducers [5].

Definition 4.2 (tree transducers). A (top-down, non-deterministic) tree transducer \mathcal{X} is a quadruple (Σ, Q, q_I, Θ) , where Σ is a ranked alphabet, Q is a set of states, q_I is the initial state, and $\Theta (\subseteq Q \times \text{dom}(\Sigma) \times Q^* \times \text{Terms}_{\Sigma})$ satisfies: if $(q, a, q_1 \dots q_n, M) \in \Theta$, then $\Sigma(a) = n$ and $\emptyset \vdash M : \underbrace{\circ \rightarrow \dots \rightarrow \circ}_n \rightarrow \circ$. The transduction

relation $(q, T) \xrightarrow{\mathcal{X}} M$ is defined inductively by the rule:

$$\frac{(q_i, T_i) \xrightarrow{\mathcal{X}} M_i \text{ for each } i \in \{1, \dots, n\} \quad (q, a, q_1 \dots q_n, M) \in \Theta}{(q, a T_1 \dots T_n) \xrightarrow{\mathcal{X}} M M_1 \dots M_n}$$

³Kobayashi's algorithm [17] is for the *simply-typed* λ -calculus with recursion. We can however adapt it for our language with intersection types.

We write $\mathcal{X}(T)$ for the set of trees $\{\llbracket M \rrbracket \mid (q_I, T) \longrightarrow_{\mathcal{X}} M\}$, and call an element of $\mathcal{X}(T)$ an *output* of the transducer \mathcal{X} for T .

A transducer $\mathcal{X} = (\Sigma, Q, q_I, \Theta)$ is *deterministic* if, for each pair $(q, a) \in Q \times \text{dom}(\Sigma)$, there is at most one (q_1, \dots, q_n, M) such that $(q, a, q_1 \dots q_n, M) \in \Theta$.

For a deterministic transducer, $\mathcal{X}(T)$ is empty or a singleton set. When $\mathcal{X}(T)$ is singleton, by abuse of notation, we sometimes write $\mathcal{X}(T)$ for the element of $\mathcal{X}(T)$.

Example 4.3. Let $\Sigma_2 = \{\mathbf{a} \mapsto 1, \mathbf{b} \mapsto 1, \mathbf{s} \mapsto 1, \mathbf{e} \mapsto 0\}$. Consider the transducer $\mathcal{X} = (\Sigma_2, \{q_0, q_1, q_f\}, q_0, \Theta)$ where Θ is given by:

$$\Theta = \{(q_0, \mathbf{b}, q_1, \lambda x.x), (q_0, \mathbf{a}, q_0, \mathbf{s}), (q_1, \mathbf{b}, q_1, \mathbf{s}), (q_f, \mathbf{e}, \epsilon, \mathbf{e}), (q_1, \mathbf{a}, q_f, \lambda x.\mathbf{e}), (q_f, \mathbf{a}, q_f, \lambda x.\mathbf{e}), (q_f, \mathbf{b}, q_f, \lambda x.\mathbf{e})\}$$

Given a Σ_2 -labeled tree T without \mathbf{s} , \mathcal{X} returns the depth of the first occurrence of a subterm of the form $\mathbf{b}(\mathbf{a}(\dots))$ in unary representation. For example, for $T = \mathbf{a}(\mathbf{b}(\mathbf{a}(\mathbf{a}(\mathbf{b}(\mathbf{e}))))))$, $\mathcal{X}(T) = \{\mathbf{s}(\mathbf{e})\}$. \square

Example 4.4. Let $\Sigma_3 = \{\mathbf{a} \mapsto 1, \mathbf{b} \mapsto 1, \mathbf{e} \mapsto 0, \mathbf{br} \mapsto 2\}$. Consider the transducer $\mathcal{X} = (\Sigma_3, \{q_0, q_{\text{odd}}, q_{\text{even}}\}, q_0, \Theta)$ where Θ is given by:

$$\Theta = \{(q_0, \mathbf{br}, q_{\text{odd}} q_{\text{even}}, \lambda x.\lambda y.\mathbf{br} y x), (q_0, \mathbf{br}, q_{\text{odd}} q_{\text{odd}}, \mathbf{br}), (q_0, \mathbf{br}, q_{\text{even}} q_{\text{even}}, \mathbf{br}), (q_0, \mathbf{br}, q_{\text{even}} q_{\text{odd}}, \mathbf{br}), (q_{\text{odd}}, \mathbf{a}, q_{\text{even}}, \mathbf{a}), (q_{\text{even}}, \mathbf{a}, q_{\text{odd}}, \mathbf{a}), (q_{\text{even}}, \mathbf{e}, \epsilon, \mathbf{e})\}$$

It takes a tree of the form $\mathbf{br}(\mathbf{a}^m(\mathbf{e}))(\mathbf{a}^n(\mathbf{e}))$ as an input, and swaps the subtrees $\mathbf{a}^m(\mathbf{e})$ and $\mathbf{a}^n(\mathbf{e})$ only if m is odd and n is even. It is a non-deterministic transducer, but outputs a singleton set. \square

The goal of this subsection is, given a program M and a tree transducer \mathcal{X} , to construct a program N that produces an element of $\mathcal{X}(\llbracket M \rrbracket)$. The construction should satisfy the following properties.

1. It should be reasonably efficient (which also implies N is not too large). In particular, it should be often faster than actually constructing $\llbracket M \rrbracket$ and then applying the transducer.
2. It should be easy to apply further operations (such as pattern matching as discussed in the previous section) on N .

A naive approach to construct N would be to express the transducer \mathcal{X} as a program f , and let N be $f(M)$.⁴ This approach obviously does not satisfy the second criterion, however.

We discuss an approach based on an extension of higher-order model checking below. An alternative approach, based on fusion transformation, is discussed in Appendix C.

4.2.1 Model Checking Approach

We can extend the higher-order model checking discussed in Section 4.1 to compute the output of a transducer (without decomposition). Let $\mathcal{X} = (\Sigma, Q, q_I, \Theta)$ be a tree transducer. We shall define a type-directed, non-deterministic transformation relation $\Psi \vdash_{\mathcal{X}} M : \theta \Longrightarrow N$, where θ is a refinement type introduced in Section 4.1. Intuitively, it means that if the value of M is traversed by transducer \mathcal{X} as specified by θ , then the output of the transducer is (the tree or function on trees represented by) N . As a special case, if M represents a tree and if $\Psi \vdash_{\mathcal{X}} M : q_I \Longrightarrow N$, then N is an output of \mathcal{X} , i.e., $\llbracket N \rrbracket \in \mathcal{X}(\llbracket M \rrbracket)$. The relation $\Psi \vdash_{\mathcal{X}} M : \theta \Longrightarrow N$ is inductively defined by the following rules.

⁴Strictly speaking, as our language does not have destructors for tree constructors $a_1, \dots, a_n \in \text{dom}(\Sigma)$, we need to transform M into $M' a_1 \dots a_n$ where M' is a pure λ -term, and then transform it into $f M' a_1 \dots a_n$.

$$\frac{}{\Psi, x : \theta \vdash_{\mathcal{X}} x : \theta \Longrightarrow x_{\theta}} \quad (\text{TR-VAR})$$

$$\frac{(q, a, q_1 \dots q_k, N_{q,a,q_1 \dots q_k}) \in \Theta}{\Psi \vdash_{\mathcal{X}} a : q_1 \rightarrow \dots \rightarrow q_k \rightarrow q \Longrightarrow N_{q,a,q_1 \dots q_k}} \quad (\text{TR-CONST})$$

$$\frac{\Psi, x : \theta_1, \dots, x : \theta_n \vdash_{\mathcal{X}} M : \theta \Longrightarrow N \quad x \text{ does not occur in } \Psi}{\Psi \vdash_{\mathcal{X}} \lambda x.M : \theta_1 \wedge \dots \wedge \theta_n \rightarrow \theta \Longrightarrow \lambda x_{\theta_1} \dots \lambda x_{\theta_n}.N} \quad (\text{TR-ABS})$$

$$\frac{\Psi \vdash_{\mathcal{X}} M_1 : \theta_1 \wedge \dots \wedge \theta_n \rightarrow \theta \Longrightarrow N_1 \quad \forall i \in \{1, \dots, n\}. \Psi \vdash_{\mathcal{X}} M_2 : \theta_i \Longrightarrow N_{2,i}}{\Psi \vdash_{\mathcal{X}} M_1 M_2 : \theta \Longrightarrow N_1 N_{2,1} \dots N_{2,n}} \quad (\text{TR-APP})$$

Basically, the transformation works in a compositional manner. Note that if we remove the part “ $\Longrightarrow N$ ”, the rules above are essentially the same as refinement typing rules. In rule TR-CONST, the transformation for constants is determined by the transducer. In rule TR-APP, if the argument M_2 in the original program should have multiple refinement types $\theta_1, \dots, \theta_n$, we separately translate the argument M_2 for each type, and duplicate the argument. Thus, in rule TR-ABS for functions, the function $\lambda x.M$ of type $\theta_1 \wedge \dots \wedge \theta_n \rightarrow \theta$ is transformed into a function that takes multiple arguments.

Example 4.5. Consider the following program to compute $(\mathbf{ab})^2 \mathbf{e}$:

let $\textit{twice} = \lambda f.\lambda z.f(f(z))$ **in** $\textit{twice}(\lambda z.\mathbf{a}(\mathbf{b}(z))) \mathbf{e}$

Let us consider the transducer \mathcal{X} given in Example 4.3. Let ρ and Ψ be $(q_1 \rightarrow q_0) \wedge (q_f \rightarrow q_1) \rightarrow q_f \rightarrow q_0$ and $\textit{twice} : \rho$ respectively. Then, we have:

$$\begin{aligned} \Psi \vdash_{\mathcal{X}} \textit{twice} : \rho &\Longrightarrow \textit{twice}_{\rho} \\ \Psi \vdash_{\mathcal{X}} \lambda z.\mathbf{a}(\mathbf{b}(z)) : q_1 \rightarrow q_0 &\Longrightarrow \lambda z_{q_1}.\mathbf{s}((\lambda x.x)z_{q_1}) \\ \Psi \vdash_{\mathcal{X}} \lambda z.\mathbf{a}(\mathbf{b}(z)) : q_f \rightarrow q_1 &\Longrightarrow \lambda z_{q_f}.\lambda x.\mathbf{e}((\lambda x.\mathbf{e})z_{q_f}) \\ \Psi \vdash_{\mathcal{X}} \mathbf{e} : q_f &\Longrightarrow \mathbf{e} \end{aligned}$$

Thus, we get:

$$\Psi \vdash_{\mathcal{X}} \textit{twice}(\lambda z.\mathbf{a}(\mathbf{b}(z))) \mathbf{e} : q_0 \Longrightarrow \textit{twice}_{\rho}(\lambda z_{q_1}.\mathbf{s}((\lambda x.x)z_{q_1}))(\lambda z_{q_f}.\lambda x.\mathbf{e}((\lambda x.\mathbf{e})z_{q_f})) \mathbf{e}$$

The body of \textit{twice} is transformed as follows.

$$\begin{aligned} \emptyset \vdash_{\mathcal{X}} \lambda f.\lambda z.f(f(z)) : \rho &\Longrightarrow \\ \lambda f_{q_1 \rightarrow q_0}.\lambda f_{q_f \rightarrow q_1}.\lambda z_{q_f}.f_{q_1 \rightarrow q_0}(f_{q_f \rightarrow q_1} z_{q_f}) & \end{aligned}$$

Thus, after some obvious simplifications (such as $(\lambda x.\mathbf{e})M =_{\beta} \mathbf{e}$), we obtain the following program.

let $\textit{twice} = \lambda f_1.\lambda f_2.\lambda z.f_1(f_2(z))$ **in** $\textit{twice}(\lambda z.\mathbf{s} z)(\lambda z.\mathbf{e} z) \mathbf{e}$

By evaluating it, we get $(\mathbf{s} \mathbf{e})$, which is the output of \mathcal{X} applied to $(\mathbf{ab})^2 \mathbf{e}$. \square

The following theorem guarantees the correctness of the transformation. A proof is given in Appendix B.

Theorem 4.4. *Let \mathcal{X} be a tree transducer. If $\emptyset \vdash_{\mathcal{X}} M : q_I \Longrightarrow N$, then $\llbracket N \rrbracket \in \mathcal{X}(\llbracket M \rrbracket)$. Conversely, if $\mathcal{X}(\llbracket M \rrbracket)$ is not empty, then there exists N such that $\emptyset \vdash_{\mathcal{X}} M : q_I \Longrightarrow N$ and $\llbracket N \rrbracket \in \mathcal{X}(\llbracket M \rrbracket)$.*

Remark 4.2. The second part of the theorem above does not guarantee that every element of $\mathcal{X}(\llbracket M \rrbracket)$ is obtained by the transformation.

Remark 4.3. The transformation above is also applicable to an extension of transducers called *high-level transducers* [7].

Algorithm. Suppose that a derivation tree for $\emptyset \vdash M : \circ$ and a transducer \mathcal{X} are given. Thanks to the above theorem, we can reuse the algorithm presented in Section 4.1, to decide whether $\mathcal{X}(\llbracket M \rrbracket)$ is non-empty, and if so, output an element of $\mathcal{X}(\llbracket M \rrbracket)$: Let $\mathcal{A}_{\mathcal{X}}$ be an associated automaton (Σ, Q, q_I, Δ) , where $\Delta = \{(q, a, q_1 \cdots q_n) \mid (q, a, q_1 \cdots q_n, M') \in \Theta\}$. Given a program M that generates a tree, we can first check whether $\emptyset \vdash_{\mathcal{A}_{\mathcal{X}}} M : q_I$ holds. If it does not hold, then $\mathcal{X}(\llbracket M \rrbracket)$ is empty, so we are done. Otherwise, we have a derivation tree for $\emptyset \vdash_{\mathcal{A}_{\mathcal{X}}} M : q_I$, from which we can construct a derivation tree for the program transformation relation: $\emptyset \vdash_{\mathcal{X}} M : q_I \Longrightarrow N$, and output N .

By Theorem 4.3, the above algorithm runs in time linear in the size of M , under the assumption that the size of \mathcal{X} and the type width of the derivation tree for $\emptyset \vdash M : \circ$ is bounded by a constant (though the constant factor can be huge as in higher-order model checking).

5. Implementation and Experiments

We have implemented the following two prototype systems, which can be tested at <http://www.kb.ecei.tohoku.ac.jp/~koba/compress/>.

1. A data compression system based on the algorithm described in Section 3: It takes a tree as an input, and outputs a λ -term that generates the tree. It is based on the algorithm described in Section 3, but it has a few parameters to adjust heuristics: D, N, W . D is the depth of the search of the algorithm of Figure 3. The system first applies *compressAsTree* up to depth D , and returns up to W smallest terms. The system then repeats this up to N times. (Thus, the total search depth is $N \times D$, but some candidates are dropped due to the width parameter W .)

2. A system to manipulate compressed data: It takes a program M in the form of a higher-order recursion scheme [26] and an automaton \mathcal{A} (or a transducer \mathcal{X} , resp.) as input, and answers whether $\llbracket M \rrbracket$ is accepted by \mathcal{A} (or outputs a program that generates $\mathcal{X}(\llbracket M \rrbracket)$). We have implemented a new version of a higher-order model checker based on a refinement of Kobayashi's linear-time algorithm [17] (as the previous model checkers [15, 17] are not fast enough for our purpose), and then added a feature to produce the output of a transducer based on the transformation given in Section 4.

5.1 Compression

We report the experimental on the data compression system. The main purposes of the experiments are: (i) to check whether interesting patterns can be obtained (to confirm the third advantage discussed in Section 1), and (ii) to check whether there is an advantage in terms of the compression ratio. As the current implementation is naive, scalability is a secondary issue at the moment.

5.1.1 Knowledge/Program Discovery

Natural Number The first experiment is for (unary) trees of the form $a^n(\mathbf{e})$. For $n = 9$ (with parameters $N = 3, D = 1, W = 4$), the output was:

```
let thrice =  $\lambda f. \lambda x. f(f(f(x)))$  in thrice(thrice a)e.
```

For $n = 16$ (with $N = 10, D = 1, W = 4$), the output was:

```
let twice =  $\lambda f. \lambda x. f(f(x))$  in (twice twice) twice a e.
```

Here, we have renamed variables with common names such as *twice*. Thus, common functions such as *twice* and *thrice* have been automatically discovered. The part *thrice(thrice a)* also corresponds to the square function for Church numerals, and (*twice twice*) *twice* corresponds to the exponential $2^{2^2} = 16$. This indicates that our algorithm can achieve hyper-exponential

compression ratio. (In fact, by running our algorithm by hand, we get $65536 = 2^{2^{2^2}}$.)

Thue-Morse Sequence Thue-Morse Sequence (A010060 in <http://oeis.org/>) t_n is the 0-1 sequence generated by:

$$t_0 = 0 \quad t_n = t_{n-1} s_{n-1}$$

where s_i is the sequence obtained from t_i by interchanging 0s and 1s. For example, $t_3 = 01101001$ and $t_4 = 0110100110010110$.

We have encoded a 0-1-sequence into a unary tree consisting of **a** (for 0), **b** (for 1), and **e** (for the end of the sequence): for example, 011 was represented by **a(b(b e))**. For the 10th sequence t_{10} (with $N = 20, D = 1, W = 4$), the output was:

```
let rep =  $\lambda x. \lambda y. \lambda z. x(y(y(xz)))$  in
let step =  $\lambda f. \lambda a. \lambda b. rep(fab)(fba)$  in
let iter = step(step(step rep)) in
let  $t_8 = iter \mathbf{a} \mathbf{b}$  in let  $s_8 = iter \mathbf{b} \mathbf{a}$  in
 $t_8(s_8(s_8(t_8 \mathbf{e})))$ 
```

This is an interesting encoding of the Thue-Morse Sequence. It is known that $t_n = t_{n-2} s_{n-2} t_{n-2}$ holds for all $n \geq 2$. The above encoding uses this recurrence equation (which has been somehow discovered automatically from only the 10th sequence, not from the definition of Thue-Morse Sequence!), and represents t_{10} as $t_8 s_8 s_8 t_8$. Using the above equation, t_8 and s_8 were represented by (*step*³ *rep*) **a b** and (*step*³ *rep*) **b a** respectively.

As for the compression ratio, the length of n -th Thue-Morse Sequence is $O(2^n)$, while the size of the above representation is $O(n)$. For a larger k , the part *step* ^{k} *rep* (in *iter* above) can further be compressed as in the compression of natural numbers discussed above; thus the hyper-exponential compression ratio is achieved by our algorithm.

Fibonacci Word For the 7th Fibonacci word **abaababaabaababaababa** (with $N = 10, D = 1, W = 4$), one of the outputs was:

```
let  $f_2 = \lambda y. \mathbf{a}(by)$  in let  $f_3 = \lambda y. f_2(ay)$  in
let  $f_4 = \lambda y. f_3(f_2y)$  in let  $f_5 = \lambda y. f_4(f_3y)$  in
 $f_5(f_4(f_5 \mathbf{e}))$ 
```

This is almost the definition of Fibonacci word; the last line is equivalent to **let $f_6 = \lambda y. f_5(f_4y)$ in $f_6(f_5 \mathbf{e})$** . (Note again that we have not given the definition of Fibonacci word; we have only given the specific instance.) The system could not, however, find a more compact representation such as the one in Example 2.4. This is probably due to the limitation discussed at the end of Section 3, that our compression algorithm is not powerful enough to extract some higher-order patterns.

L-system Consider an instance of L-systems, defined by [25]:

$$F_0 = \mathbf{f} \quad F_{n+1} = F_n[+F_n]F_n[-F_n]F_n$$

where “[”, “]”, “+”, “-” and **f** are terminal symbols. Given the unary tree representation of the sequence F_3 (which is given in Figure 6 of [25]), our system (with $N = 50, D = 1, W = 4$) output the following program in 38 seconds:

```
let step =  $\lambda g. \lambda z. g(\mathbf{let} h = \lambda z. g([\mathbf{g}z]) \mathbf{in} [(+(h([-h z]))))])$ 
in step(step(step(f))) e.
```

The function *step* is equivalent to: $\lambda g. \lambda z. g[+g]g[-g]g z$, where the applications are treated as right-associative here to avoid too many parentheses. The above output is exactly (a compressed form of) the definition of F_3 . The output of Sequitur [25] is more complex.

5.1.2 Compression of Real-Life Data

In this subsection, we report experiments to check whether our compression scheme is also effective (in terms of the compression

ratio and whether interesting patterns are discovered) for real-life data. We have tested two kinds of data: XML data and natural language sentences.

To evaluate the effectiveness of compression, we compared the outputs of our compression system with those of grammar-based tree compression tool TreeRePair [22] (<http://code.google.com/p/treerepair/>). The sizes of compressed data were similar for the experiments below, and we could not observe a clear advantage over TreeRePair.⁵ Potential reasons for the lack of clear advantages are: (i) Natural language sentences may not have much redundancy, (ii) the data in the experiments are too small to benefit from extraction of higher-order structures, and (iii) our compression algorithm is not good enough to extract complex higher-order structures (recall Section 3).

Our system could, however, find some interesting higher-order patterns, as reported below.

XML Data of Wikipedia We tested the structure of the first 11280 lines (1% in lines) of enwik8, which is the target data of Hutter Prize⁶, a compression competition. As in the experiments for tree compression in [4, 22], we removed PCData and attributes and used the binary-tree encoding.

In the output of our system, the higher-order pattern *twice* was extracted to compress the repetition such as `<namespace/> . . . <namespace/>` appearing in the tree. In addition, the following flipped version of the combinator $B (\lambda f. \lambda g. \lambda x. f (g x))$ was extracted.

$$Q = \lambda f. \lambda g. \lambda x. g (f x)$$

With these functions, we can compress a term like

```
let h1 = λx.f1(g x) in . . . let hn = λx.fn(g x) in hn
```

with $2n$ edges into

```
let Q = λf. λg. λx. g (f x) in let g' = Q g in
let h1 = g' f1 in . . . let hn = g' fn in hn
```

with $n + 3$ edges.

English Text We examined a part of (simple) English text extracted from the article of “Jupiter” in Simple-English version of Wikipedia <http://simple.wikipedia.org/wiki/Jupiter>. The text had 1017 words including punctuations; e.g., “Jupiter’s” is considered as 3 words “Jupiter”, an apostrophe, and “s”. The text was encoded as a unary tree, whose node is labelled by a word instead of a character.

In addition to frequently-appearing phrases such as “million miles away”, “in 1979.”, and “km/h”, interesting higher-order patterns were extracted, such as:

```
let s = λy.APOS (“s” y) in let possessive = Q s in . . .
```

The pattern $Q s = \lambda n. \lambda y. n (\text{APOS} (“s” y))$ expresses the possessive form “A’s B”. The combinators B and Q were also extracted.

TreeRePair compressed the XML corresponding to the text to a grammar with 915 edges. Here we encoded a text to a horizontal sequence of elements enclosed with root element, e.g., “This is good.” is encoded to `<root><This/><is/><good/><PERIOD></root>`. When we fed the output of TreeRePair to our system, the output had 892 edges. This size reduction was mainly caused by introduction of the combinators B and Q .

⁵ We omit size data here, since for a fair comparison, we need to encode both λ -terms and grammars into bit strings. Note however that if we first use TreeRePair before applying our system, the result cannot be worse than that of TreeRePair.

⁶ <http://prize.hutter1.net/>

English-French Translation We have also tested our system to compress a sequence of pairs of an English sentence and its French translation. We have taken simple English sentences from a textbook used in an English course of a kindergarten and prepared their French translations by using Google translation (<http://translate.google.com/>). Given an input containing:

```
pair(I(like(him(period))))(Je(le(aime(bien(period))))))
```

and

```
pair(I(like(her(period))))(Je(la(aime(bien(period)))))),
```

our system produced the following output:

```
let xE = λz.pair (I(like(z(period)))) in
let xF = λz.(Je(z(aime(bien(period)))))) in
. . . (xE him (xF (le))) . . . (xE her (xF (la))) . . .
```

Thus, the correspondences like “him” vs “le”, “her” vs “la”, and “I like xxx” vs “Je xxx aime bien” have been discovered. Other word or phrase-phrase correspondences that have been found from a sequence of 14 pairs of sentences include: “plays with a . . .” vs “joue avec un . . .”, “friend” vs “ami”, “ball” vs “ballon”, etc.

5.2 Data Processing

We have applied various pattern match queries and transformations to Fibonacci words, to check the scalability of our system with respect to the size of compressed data. Table 1 shows the results. The 2^m -th Fibonacci words (for $m = 4, 6, 8, 10, 12, 14$) were represented by using the encoding of Example 2.4. The size of the representation of the n -th Fibonacci word is $O(\log n)$ (or $O(m)$). The queries and transformations are: Q_1 : contains aa, Q_2 : contains no bb, Q_3 : contains no aaa, T_1 : the first occurrence of aab, T_2 : count the number of ab, T_3 : replace ab with bb, TQ_3 : T_3 followed by query “contains bbb?”. In the row TQ_3 , the times do not contain those for applying T_3 (which are shown in the row T_3). All the experiments were conducted on a machine with Intel(R) Xeon(R) CPU with 3Ghz and 8GB memory.

Our system based on higher-order model checking could quickly answer pattern match queries or apply transformations. The increase of the time with respect to m varies depending on the query or transformation, but an exponential slowdown was not observed for any of the queries and transformations. Note that the length of n -th Fibonacci word is greater than 1.6^{n-1} , so that it is impossible to actually (no matter whether eagerly or lazily) construct the word and then apply a pattern match query or transformation. Even with the grammar-based compression based on context-free grammars, the size of the representation of n -th Fibonacci word is $O(n)$; thus our approach (which runs in time $O(\log n)$) is exponentially faster than the grammar-based approach for this experiment. Our system was relatively slower for TQ_3 . This is probably because the transformation T_3 increased the arity of functions, which had a bad effect on model checking. It may be possible to reduce this problem by post-processing the output of the transformation using other program transformation techniques.

To check the scalability of our system with respect to the size of patterns, we have also tested our system for the following larger queries on 1024th Fibonacci words: Q_4 : contains abaababaab, and Q_5 : contains no abababaaba. The times to answer the queries were 0.51 and 0.44 seconds respectively.

It should be noted however that the above result is an extreme case that shows the advantage of our approach. When the effect of compression is small as in the experiments in Section 5.1.2, the advantage of compressing data and manipulating them without decompression can be easily offset by the inefficiency of the current higher-order model checker.

	m=4	m=6	m=8	m=10	m=12	m=14
Q_1	0.12	0.12	0.12	0.26	0.27	0.27
Q_2	0.03	0.03	0.04	0.12	0.12	0.12
Q_3	0.14	0.14	0.14	0.14	0.14	0.14
T_1	0.13	0.13	0.13	0.27	0.27	0.27
T_2	0.04	0.04	0.12	0.12	0.13	0.13
T_3	0.04	0.04	0.12	0.12	0.13	0.13
TQ_3	0.47	0.62	1.32	1.53	1.84	2.09

Table 1. Times for processing queries and transformations on 2^m th Fibonacci words, measured in seconds.

6. Related Work

The idea of compressing strings or tree data as functional programs is probably not new; in fact, Tromp [37] studied Kolmogorov complexity in the setting of λ -calculus. We are, however, not aware of any serious previous studies of the approach that propose data compression/manipulation algorithms with a similar capability.

In the context of higher-order model checking, Broadbent et al. ([3], Corollary 3) showed that if t is the tree generated by an order- n higher-order recursion scheme and \mathcal{I} is a well-formed MSO-interpretation, $\mathcal{I}(t)$ can be generated by an order- $(n+1)$ recursion scheme. As a higher-order recursion scheme can be viewed as a simply-typed λ -term (with recursion) and a transducer can be expressed as a MSO-interpretation, this gives another procedure for the data transformation discussed in Section 4.2 (for the case where the program M is simply-typed). Their transformation is however indirect and quite complex, as it goes via collapsible higher-order pushdown automata [9]. Their transformation also increases the order of the program, as opposed to our transformation given in Section 4.2.1. Thus, we think their transformation is mainly of theoretical interest (indeed, it has never been implemented).

As discussed in Section 2.3, our approach to compress data as functional programs can be regarded as a generalization of grammar-based compression [1, 4, 22, 25, 28, 30, 31]. Since the problem of computing the smallest CFG that exactly generates w is known to be NP-hard [35], various heuristic compression algorithms have been proposed, including Re-pair [19, 22]. Processing of compressed data without decompression has been a hot topic of studies in the grammar-based compression, and our result in Section 4 can be considered a generalization of it to higher-order grammars. In the context of CFG-based compression, however, more operations can be performed without decompression, including the equivalence checking (“given two compressed strings, do they represent the same string?”) [27] and compressed pattern matching (“given a compressed string and a *compressed* pattern, does the string match the pattern?”) [12]. It is left for future work to investigate whether those operations extend to higher-order grammars.

Our experiment to discover knowledge from the compression of English-French translation, discussed in Section 5.1, appears to be related to studies of example-based machine translation [32], in particular, automatic extraction of translation templates from a bilingual corpus [11]. Nevill-Manning and Witten [25] also report inference of hierarchical (not higher-order, in the sense of the present paper) structures by grammar-based compression.

We have not discussed the issue of how to compactly represent a λ -term (obtained by our compression algorithm) as a bit string. Tromp [37] gave a few schemes for representing untyped λ -terms as bit strings. Vytiniotis and Kennedy [40] introduced a game-based method for representing simply-typed λ -terms as bit-strings.

7. Conclusion

We have studied the approach of compressing data as functional programs, and shown that programming language techniques can be used for compressing and manipulating data. In particular, we have extended a higher-order model checking algorithm to transform compressed data without decompression. The prototype compression and transformation systems have been implemented and interesting experimental results have been obtained. Larger experiments are left for future work.

Acknowledgments We would like to thank Tatsunari Nakajima, Kunihiko Sadakane, Kazuya Yaguchi, and anonymous referees for discussions and comments. This work was partially supported by Kakenhi 23220001.

References

- [1] A. Apostolico and S. Lonardi. Some theory and proactive of greedy off-line textual substitution. In *Data Compression Conference 1998 (DCC98)*, pages 119–128, 1998.
- [2] H. Barendregt, M. Coppo, and M. Dezani-Ciancaglini. A filter lambda model and the completeness of type assignment. *J. Symb. Log.*, 48(4):931–940, 1983.
- [3] C. H. Broadbent, A. Carayol, C.-H. L. Ong, and O. Serre. Recursion schemes and logical reflection. In *Proceedings of LICS 2010*, pages 120–129. IEEE Computer Society Press, 2010.
- [4] G. Busatto, M. Lohrey, and S. Maneth. Efficient memory representation of XML document trees. *Inf. Syst.*, 33(4-5):456–474, 2008.
- [5] H. Comon, M. Dauchet, R. Gilleron, C. Löding, F. Jacquemard, D. Lugiez, S. Tison, and M. Tommasi. Tree automata techniques and applications. Available on: <http://www.grappa.univ-lille3.fr/tata>, 2007. release October, 12th 2007.
- [6] M. Crochemore and W. Rytter. *Jewels of Stringology*. World Scientific, 2002.
- [7] J. Engelfriet and H. Vogler. High level tree transducers and iterated pushdown tree transducers. *Acta Inf.*, 26(1/2):131–192, 1988.
- [8] A. J. Gill, J. Launchbury, and S. L. Peyton-Jones. A short cut to deforestation. In *FPCA*, pages 223–232, 1993.
- [9] M. Hague, A. Murawski, C.-H. L. Ong, and O. Serre. Collapsible pushdown automata and recursion schemes. In *Proceedings of LICS 2008*, pages 452–461. IEEE Computer Society, 2008.
- [10] M. Hutter. *Universal Artificial Intelligence: Sequential Decisions based on Algorithmic Probability*. Springer-Verlag, Berlin, 2004.
- [11] H. Kaji, Y. Kida, and Y. Morimoto. Learning translation templates from bilingual text. In *COLING*, pages 672–678, 1992.
- [12] M. Karpinski, W. Rytter, and A. Shinohara. An efficient pattern matching algorithm for strings with short description. *Nordic Journal on Computing*, 4(2):129–144, 1997.
- [13] T. Kida, T. Matsumoto, Y. Shibata, M. Takeda, A. Shinohara, and S. Arikawa. Collage system: a unifying framework for compressed pattern matching. *Theor. Comput. Sci.*, 1(298):253–272, 2003.
- [14] T. Knapik, D. Niwinski, and P. Uzczyzyn. Higher-order pushdown trees are easy. In *Proceedings of FOSSACS 2002*, volume 2303 of LNCS, pages 205–222. Springer-Verlag, 2002.
- [15] N. Kobayashi. Model-checking higher-order functions. In *Proceedings of PPDP 2009*, pages 25–36. ACM Press, 2009.
- [16] N. Kobayashi. Types and higher-order recursion schemes for verification of higher-order programs. In *Proc. of POPL*, pages 416–428, 2009.
- [17] N. Kobayashi. A practical linear time algorithm for trivial automata model checking of higher-order recursion schemes. In *Proceedings of FOSSACS 2011*, volume 6604 of LNCS, pages 260–274. Springer-Verlag, 2011.
- [18] N. Kobayashi and C.-H. L. Ong. Complexity of model checking recursion schemes for fragments of the modal mu-calculus. In *Proceed-*

- ings of *ICALP 2009*, volume 5556 of *LNCS*, pages 223–234. Springer-Verlag, 2009.
- [19] N. J. Larsson and A. Moffat. Offline dictionary-based compression. In *Proc. Data Compression Conference '99 (DCC'99)*, page 296. IEEE Computer Society, 1999.
- [20] M. Li and P. M. B. Vitányi. Kolmogorov complexity and its applications. In *Handbook of Theoretical Computer Science, Volume A: Algorithms and Complexity*, pages 187–254. The MIT Press, 1990.
- [21] M. Li and P. M. B. Vitányi. *An introduction to Kolmogorov complexity and its applications (3rd ed.)*. Texts in computer science. Springer-Verlag, 2009.
- [22] M. Lohrey, S. Maneth, and R. Mennicke. Tree structure compression with repair. In *2011 Data Compression Conference (DCC 2011)*, pages 353–362. IEEE Computer Society, 2011.
- [23] S. Maneth and G. Busatto. Tree transducers and tree compressions. In *Proceedings of FOSSACS 2004*, volume 2987 of *LNCS*, pages 363–377. Springer-Verlag, 2004.
- [24] W. Matsubara, S. Inenaga, A. Ishino, A. Shinohara, T. Nakamura, and K. Hashimoto. Efficient algorithms to compute compressed longest common substrings and compressed palindromes. *Theor. Comput. Sci.*, 410(8-10):900–913, 2009.
- [25] C. G. Nevill-Manning and I. H. Witten. Compression and explanation using hierarchical grammars. *Comput. J.*, 40(2/3):103–116, 1997.
- [26] C.-H. L. Ong. On model-checking trees generated by higher-order recursion schemes. In *LICS 2006*, pages 81–90. IEEE Computer Society Press, 2006.
- [27] W. Plandowski. Testing equivalence of morphisms on context-free languages. In *ESA '94*, volume 855 of *LNCS*, pages 460–470. Springer-Verlag, 1994.
- [28] W. Rytter. Application of Lempel-Ziv factorization to the approximation of grammar-based compression. *Theor. Comput. Sci.*, 302(1-3):211–222, 2003.
- [29] W. Rytter. Grammar compression, LZ-encodings, and string algorithms with implicit input. In *ICALP'04*, volume 3142 of *LNCS*, pages 15–27. Springer-Verlag, 2004.
- [30] H. Sakamoto. A fully linear-time approximation algorithm for grammar-based compression. *Journal of Discrete Algorithms*, 3(2-4):416–430, 2005.
- [31] H. Sakamoto, S. Maruyama, T. Kida, and S. Shimozone. A space-saving approximation algorithm for grammar-based compression. *IE-ICE Trans. on Information and Systems*, E92-D(2):158–165, 2009.
- [32] H. L. Somers. Review article: Example-based machine translation. *Machine Translation*, 14(2):113–157, 1999.
- [33] R. Statman. The typed lambda-calculus is not elementary recursive. *Theor. Comput. Sci.*, 9:73–81, 1979.
- [34] C. Stirling. Decidability of higher-order matching. *Logical Methods in Computer Science*, 5(3), 2009.
- [35] J. Storer. NP-completeness results concerning data compression. Technical Report 234, Department of Electrical Engineering and Computer Science, Princeton University, Princeton, N.J., 1997.
- [36] A. Takano and E. Meijer. Shortcut deforestation in calculational form. In *Proceedings of the seventh international conference on Functional programming languages and computer architecture, FPCA '95*, pages 306–313, New York, NY, USA, 1995. ACM.
- [37] J. Tromp. Binary lambda calculus and combinatory logic. In C. S. Claude, editor, *Randomness And Complexity, from Leibniz To Chaitin*, pages 237–260. World Scientific Publishing Company, 2008.
- [38] T. Tsukada and N. Kobayashi. Untyped recursion schemes and infinite intersection types. In *Proceedings of FOSSACS 2010*, volume 6014 of *LNCS*, pages 343–357. Springer-Verlag, 2010.
- [39] S. van Bakel. Intersection type assignment systems. *Theor. Comput. Sci.*, 151(2):385–435, 1995.
- [40] D. Vytiniotis and A. Kennedy. Functional pearl: every bit counts. In *Proceedings of ICFP 2010*, pages 15–26, 2010.

Appendix

A. Relationship with Kolmogorov Complexity

As already discussed in Section 1, the “functional programs as compressed data” approach provides a universal compression scheme, in the sense that any compressed data can be expressed in the form of (typed) λ -terms. As sketched below, our representation of compressed data in the form of λ -terms is optimal with respect to Kolmogorov complexity, up to an additive constant [20, 21].

Let U be a universal Turing machine. Kolmogorov complexity [20, 21] of a binary string (an element of $\{0, 1\}^*$) v , written $K(v)$, is defined by:

$$K(v) := \min\{|w| \mid w \in \{0, 1\}^*, U(w) = v\}.$$

Here, $|w|$ is the length of w . Since the λ -calculus is Turing complete, there exists a λ -term U_λ such that $U(w) = v$ if and only if $\llbracket U_\lambda(\hat{w}) \rrbracket = \hat{v}$. Here, \hat{w} is an encoding of a binary string w into a λ -term such that the size of \hat{w} is linear in $|w|$; One can, for example, use the encoding of [37]. Define $K_{\lambda, \#}(v)$ by:

$$K_{\lambda, \#}(v) := \min\{\#M \mid \llbracket M \rrbracket = \hat{v}\}.$$

Then, by the existence of U_λ above, there exist c_1 and c_2 such that

$$\forall v \in \{0, 1\}^*. K_{\lambda, \#}(v) \leq c_1 K(v) + c_2.$$

Here, c_1 is a constant such that $\forall w \in \{0, 1\}^*. \# \hat{w} \leq c_1 |w|$, which depends on the encoding $\hat{\cdot}$.

For a more precise comparison of the size of the λ -term $U_\lambda(\hat{w})$ with Kolmogorov complexity, we need to encode the λ -term into a bit string. Given a λ -term of the form $M \hat{w}$ (where w is a bit string), let us write $B_{M \hat{w}}$ for the following bit-string:

Part A	0	Part B	w
--------	---	--------	-----

where Part A is the unary encoding of the length of Part B and Part B is some encoding of M into a bit string (such as Tromp’s one [37]).

Let $K_{\lambda, \#'}(v)$ be:

$$K_{\lambda, \#'}(v) := \min\{|B_{M \hat{w}}| \mid \llbracket M \hat{w} \rrbracket = \hat{v}\}.$$

Then, $K_{\lambda, \#'}(v) \leq |B_{U_\lambda \hat{w}}| = K(v) + 1 + 2c$, where w is a bit string such that $U(w) = v$ with $|w| = K(v)$, and c is the length of the encoding of U_λ into a bit string. As c is independent of v , the result implies that our “compressed data as functional programs” approach achieve an optimal compression size up to an additive constant.

B. Proof of Theorem 4.4

Lemma B.1. *If $\emptyset \vdash_{\mathcal{X}} M : q \implies N$, then $\emptyset \vdash M : \circ$ and $\emptyset \vdash N : \circ$.*

Proof. $\emptyset \vdash M : \circ$ follows immediately from the facts that $\emptyset \vdash_{\mathcal{X}} M : q \implies N$ implies $\emptyset \vdash_{\mathcal{A}_{\mathcal{X}}} M : q$ and that the type system for $\Gamma \vdash M : \tau$ is actually the special case for the refinement type system where the underlying automaton is $\mathcal{A} = (\Sigma, \{\circ\}, \circ, \Delta)$ with $\Delta = \{(\circ, a, \circ^{\Sigma(a)}) \mid a \in \text{dom}(\Sigma)\}$.

To show $\emptyset \vdash N : \circ$, we prepare the following generalized lemma:

$$\text{If } \Psi \vdash_{\mathcal{X}} M : \theta \implies N, \text{ then } \Psi^b \vdash_{\mathcal{A}_{\mathcal{X}}} N : \theta^b. \quad (**)$$

Here, the translation \cdot^b is defined by:

$$\begin{aligned} \circ^b &= \circ \\ (\theta_1 \wedge \dots \wedge \theta_k \rightarrow \theta)^b &= \theta_1^b \rightarrow \dots \rightarrow \theta_k^b \rightarrow \theta^b \\ (x_1 : \theta_1, \dots, x_k : \theta_k)^b &= x_{1, \theta_1} : \theta_1^b, \dots, x_{k, \theta_k} : \theta_k^b \end{aligned}$$

(**) follows by straightforward induction on the derivation of $\Psi \vdash_{\mathcal{X}} M : \theta \implies N$. \square

Lemma B.2. *Suppose that $\Psi, x : \theta_1, \dots, x : \theta_k \vdash_{\mathcal{X}} M : \theta \implies N$ and $\Psi \vdash_{\mathcal{X}} M_0 : \theta_i \implies N_i$ for each $i \in \{1, \dots, k\}$, with $x \notin \text{dom}(\Psi)$. Then $\Psi \vdash_{\mathcal{X}} [M_0/x]M_1 : \theta_1 \implies [N_1/x_{\theta_1}, \dots, N_k/x_{\theta_k}]N$.*

Proof. The derivation for

$\Psi \vdash_{\mathcal{X}} [M_0/x]M_1 : \theta_1 \implies [N_1/x_{\theta_1}, \dots, N_k/x_{\theta_k}]N$ can be obtained from that for $\Psi, x : \theta_1, \dots, x : \theta_k \vdash_{\mathcal{X}} M : \theta \implies N$, by replacing each leaf $\Psi, \Psi', x : \theta_1, \dots, x : \theta_k \vdash_{\mathcal{X}} x : \theta_i \implies x_{\theta_i}$ of the derivation with $\Psi, \Psi' \vdash_{\mathcal{X}} M_0 : \theta_i \implies N_i$. \square

Lemma B.3. *If $\Psi \vdash_{\mathcal{X}} M : \theta \implies N$ and $M \longrightarrow_{\beta} M'$, then there exists N' such that $N \longrightarrow_{\beta}^* N'$ and $\Psi \vdash_{\mathcal{X}} M' : \theta \implies N'$.*

Proof. This follows by induction on the derivation of $M \longrightarrow_{\beta} M'$. Since the induction steps are trivial, we show only the base case, where $M = (\lambda x. M_1)M_2$ and $M' = [M_2/x]M_1$. Suppose $\Psi \vdash_{\mathcal{X}} M : \theta \implies N$. Then, we have:

$$\begin{aligned} N &= (\lambda x_{\theta_1}, \dots, x_{\theta_k}. N_1) N_{2,1} \dots N_{2,k} \\ \Psi \vdash_{\mathcal{X}} M_2 : \theta_i &\implies N_{2,i} \text{ for each } i \in \{1, \dots, k\} \\ \Psi, x : \theta_1, \dots, x : \theta_k \vdash_{\mathcal{X}} M_1 : \theta &\implies N_1 \end{aligned}$$

Let N' be $[N_{2,1}/x_{\theta_1}, \dots, N_{2,k}/x_{\theta_k}]N_1$. By Lemma B.2, we have $\Psi \vdash_{\mathcal{X}} M' : \theta \implies N'$. Furthermore, $N \longrightarrow_{\beta}^* N'$ holds as required. \square

Lemma B.4. *If $\Psi \vdash_{\mathcal{X}} M' : \theta \implies N'$ and $M \longrightarrow_{\beta} M'$, then there exists N such that $\Psi \vdash_{\mathcal{X}} M : \theta \implies N$.*

Proof. The proof proceeds by induction on the derivation of $M \longrightarrow_{\beta} M'$. As the induction steps are straightforward, we discuss only the base case, where $M = (\lambda x. M_1)M_2$ and $M' = [M_2/x]M_1$. Suppose $\Psi \vdash_{\mathcal{X}} M' : \theta \implies N'$. From the derivation of $\Psi \vdash_{\mathcal{X}} M' : \theta \implies N'$, we can construct a derivation for $\Psi, x : \theta_1, \dots, x : \theta_k \vdash_{\mathcal{X}} M_1 \implies N_1$ by replacing subderivations of the form $\Psi, \Psi' \vdash_{\mathcal{X}} M_2 : \theta_i \implies N_{2,i,j}$ with those for $\Psi, \Psi', x : \theta_1, \dots, x : \theta_k \vdash_{\mathcal{X}} x \implies x_{\theta_i}$. For each $i \in \{1, \dots, k\}$, choose one such $N_{2,i,j}$ and let it be $N_{2,i}$. Let N be $(\lambda x_{\theta_1}, \dots, x_{\theta_k}. N_1)N_{2,1} \dots N_{2,k}$. Then, we have $\Psi \vdash_{\mathcal{X}} M : \theta \implies N$ as required. \square

Lemma B.5. *Let T be a Σ -labeled tree. Then, $\emptyset \vdash_{\mathcal{X}} T : q \implies N$ if and only if $(q, T) \longrightarrow_{\mathcal{X}} N$.*

Proof. This follows by induction on the structure of T . Suppose $T = a T_1 \dots T_k$ with $\Sigma(a) = k$. If $\emptyset \vdash_{\mathcal{X}} T : q \implies N$, then we have:

$$\begin{aligned} N &= M_{q,a,q_1 \dots q_k} N_1 \dots N_k \\ (q, a, q_1 \dots q_k, M_{q,a,q_1 \dots q_k}) &\in \Theta_{\mathcal{X}} \\ \emptyset \vdash_{\mathcal{X}} T_i : q_i &\implies N_i \text{ for each } i \in \{1, \dots, k\} \end{aligned}$$

By the induction hypothesis, we get $(q_i, T_i) \longrightarrow_{\mathcal{X}} N_i$. Thus, we have $(q, T) \longrightarrow_{\mathcal{X}} N$.

Conversely, suppose $(q, T) \longrightarrow_{\mathcal{X}} N$. Then, we have:

$$\begin{aligned} N &= M_{q,a,q_1 \dots q_k} N_1 \dots N_k \\ (q, a, q_1 \dots q_k, M_{q,a,q_1 \dots q_k}) &\in \Theta_{\mathcal{X}} \\ (q_i, T_i) &\longrightarrow_{\mathcal{X}} N_i \end{aligned}$$

By the induction hypothesis, we get $\emptyset \vdash_{\mathcal{X}} T_i : q_i \implies N_i$ for each $i \in \{1, \dots, k\}$. Thus, we have $\emptyset \vdash_{\mathcal{X}} T : q \implies N$ as required. \square

Proof of Theorem 4.4. Suppose $\emptyset \vdash_{\mathcal{X}} M : q_I \Longrightarrow N$. By Lemma B.1 and Theorem 2.1, there exists $T (= \llbracket M \rrbracket)$ such that $M \xrightarrow{\beta^*} T$. By Lemma B.3, there exists N' such that $N \xrightarrow{\beta^*} N'$ and $\emptyset \vdash_{\mathcal{X}} T : q_I \Longrightarrow N'$. By Lemma B.5, we have $(q_I, T) \xrightarrow{\mathcal{X}} N'$, which implies $\llbracket N' \rrbracket \in \mathcal{X}(\llbracket M \rrbracket)$. By the condition $N \xrightarrow{\beta^*} N'$, we have $\llbracket N \rrbracket = \llbracket N' \rrbracket \in \mathcal{X}(\llbracket M \rrbracket)$, as required.

Conversely, suppose that $T \in \mathcal{X}(\llbracket M \rrbracket)$. By the definition of $\mathcal{X}(\llbracket M \rrbracket)$, there exists N' such that $(q_I, \llbracket M \rrbracket) \xrightarrow{\mathcal{X}} N'$. By Lemma B.5, we have $\emptyset \vdash_{\mathcal{X}} \llbracket M \rrbracket : q_I \Longrightarrow N'$. By Lemma B.4, there exists N such that $\emptyset \vdash_{\mathcal{X}} M : q_I \Longrightarrow N$. By the first part of this theorem, we have $\llbracket N \rrbracket \in \mathcal{X}(\llbracket M \rrbracket)$ as required. \square

C. Fusion Approach

We discuss another approach, based on the idea of shortcut fusion [8, 36]. Recall that the goal was to construct a program N that produces $\mathcal{X}(\llbracket M \rrbracket)$. Here, we can regard M as a tree generator, and transducer \mathcal{X} as a tree consumer. Thus, by using shortcut fusion, we can construct a program N that computes $\mathcal{X}(\llbracket M \rrbracket)$ without constructing the intermediate data $\llbracket M \rrbracket$. For the sake of simplicity, we assume below that the transducer $\mathcal{X} = (\Sigma, Q, q_0, \Theta)$ is deterministic and total, i.e., for each $(q, a) \in Q \times \text{dom}(\Sigma)$, there is exactly one (q_1, \dots, q_k, M) such that $(q, a, q_1 \dots q_k, M) \in \Theta$.

A (deterministic and total) transducer $\mathcal{X} = (\Sigma, Q, q_0, \Theta)$ (where $Q = \{q_0, \dots, q_{n-1}\}$) can be viewed as the following homomorphism h from \mathfrak{o} (i.e., the set of Σ -labeled trees) to $Q \rightarrow \mathfrak{o}$:

$$h(a \ x_1 \ \dots \ x_k) = f_a(h \ x_1) \ \dots \ (h \ x_k) \quad (\Sigma(a) = k)$$

Here, f_a is given by:

$$\begin{aligned} f_a \ g_1 \ \dots \ g_k \ q = \\ \text{case } q \text{ of} \\ \quad q_0 \quad \Rightarrow M_0(g_1 \ q_{0,1}) \ \dots \ (g_k \ q_{0,k}) \\ \quad \vdots \\ \quad q_{n-1} \Rightarrow M_{n-1}(g_1 \ q_{n-1,1}) \ \dots \ (g_k \ q_{n-1,k}) \end{aligned}$$

where $(q_i, a, q_{i,1} \dots q_{i,k}, M_i) \in \Theta$. By using Church encoding, q_i can be encoded as the function $\lambda q_0, \dots, q_{n-1}.q_i$. (Here, $\lambda x_0, \dots, x_{n-1}.M$ is an abbreviated form of $\lambda x_0. \dots \lambda x_{n-1}.M$.) Thus, f_a above becomes:

$$\begin{aligned} \lambda g_1, \dots, g_k. \lambda q. \\ \quad q(M_0(g_1 \ q_{0,1}) \ \dots \ (g_k \ q'_{0,k})) \ \dots \\ \quad (M_{n-1}(g_1 \ q'_{n-1,1}) \ \dots \ (g_k \ q'_{n-1,k})), \end{aligned}$$

of type $\underbrace{\tau_h \rightarrow \dots \rightarrow \tau_h}_k \rightarrow \tau_h$, where $\tau_h = \underbrace{(\mathfrak{o} \rightarrow \dots \rightarrow \mathfrak{o})}_n$

$\mathfrak{o} \rightarrow \mathfrak{o}$ and $q'_{i,j} = \lambda q_0, \dots, q_{n-1}.q_{i,j}$. It is easy to verify that, for every Σ -labelled tree T , $\llbracket h(T) \lambda q_0, \dots, q_{n-1}.q_0 \rrbracket = \mathcal{X}(T)$ holds.

Now, extend h to the homomorphism on λ -terms by:

$$\begin{aligned} h(x) = x \quad h(a) = f_a \\ h(M_1 M_2) = h(M_1) h(M_2) \quad h(\lambda x. M) = \lambda x. h(M) \end{aligned}$$

h just replaces each tree constructor a with f_a . Then, $N = h(M) \lambda q_0, \dots, q_{n-1}.q_0$ is what we want.

Theorem C.1. *Suppose that M is a program of type \mathfrak{o} . Let $\mathcal{X} = (\Sigma, Q, q_0, \Theta)$ be a deterministic and total transducer, and let h be the homomorphism as defined above. Then, we have:*

$$\llbracket h(M) \lambda q_0, \dots, q_{n-1}.q_0 \rrbracket = \mathcal{X}(\llbracket M \rrbracket).$$

Proof. Let $T = \llbracket M \rrbracket$, i.e., $M \xrightarrow{\beta^*} T$. As M does not have tree destructors, we have $h(M) \xrightarrow{\beta^*} h(T)$. Thus, we get $N = h(M) \lambda q_0, \dots, q_{n-1}.q_0 \xrightarrow{\beta^*} h(T) \lambda q_0, \dots, q_{n-1}.q_0 \xrightarrow{\beta^*} \mathcal{X}(T)$ as required. \square

We used syntactic reasoning in the above proof. Alternatively, we can use semantic techniques [8, 36].

We have assumed above that \mathcal{X} is deterministic and total. To remove the assumption, it suffices to extend h to a homomorphism from \mathfrak{o} to $Q \rightarrow 2^{\mathfrak{o}}$.

Example C.1. Recall Example 4.5. With the fusion-based approach, we get the following program as the output of transformation:

```
let q'_0 = λq_0, q_1, q_f. q_0 in
let q'_1 = λq_0, q_1, q_f. q_1 in
let q'_f = λq_0, q_1, q_f. q_f in
let f_a = λg. λq. q (s(g q'_0)) ((λx. e)(g q'_f)) ((λx. e)(g q'_f)) in
let f_b = λg. λq. q ((λx. x)(g q'_1)) (s(g q'_1)) ((λx. e)(g q'_f)) in
let f_e = λq. q ⊥ e in
let twice = λf. λz. f(f(z)) in twice(λz. f_a(f_b(z))) f_e q'_0
```

Here, \perp is a special tree constructor denoting an undefined tree, introduced to make \mathcal{X} total. \square

Remark C.1. There are trade-offs between the two approaches (model checking and fusion approaches).

- As is clear from Examples 4.5 and C.1, the model checking approach reduces the program more aggressively than the fusion approach, which just postpones the computation involved in the transducer.
- The model checking approach runs in time linear in the program size only under the assumption that the type width for the type derivation of the program is fixed. The constant factor can be huge, although it does not always show up according to experiments. The fusion approach runs in time linear in the program size unconditionally.
- The fusion approach raises the order of the program (where the order of a program is the largest order of the type of a function; the order of a type is defined by: $order(\mathfrak{o}) = 0$, $order(\tau_1 \wedge \dots \wedge \tau_k \rightarrow \tau) = \max(order(\tau_1) + 1, \dots, order(\tau_k) + 1, order(\tau))$) by two. This can have a very bad effect on further pattern match queries or data manipulations based on higher-order model checking (described in Section 4.1). The model checking approach does not raise the order, although it may raise the arity of functions (e.g., in Example 4.5, the function *twice* takes three arguments after the transformation). Note that for higher-order model checking, the order of programs is the most important factor that affects the worst-case complexity [18, 26].

Because of the first and last points, we think the model checking approach is preferable, and the fusion approach (or the other naive approaches discussed at the beginning of Section 4.1) should be used only when the model checking approach is too slow.