



**Utrecht
University**

Computing Science MSc Thesis

Generic Incremental Computation for Regular Datatypes in Haskell

Author

Jort van Gorkum (6142834)

Supervisors

Wouter Swierstra

Trevor McDonell

Faculty of Science

Department of Information and Computing Sciences

Programming Technology

July 6, 2022

Write abstract

When computing a result over a data type, a small change would lead to completely recomputing the result. Parts of the recomputation give the same results as the previous computation. By keeping track of the intermediate results and reusing the results when the input is the same, fewer computations have to be performed.

- Explain incremental computation and that is used by many things
- Use incremental computation by reusing the structure of the data structure
- Implementing in Haskell is time-consuming and error-prone
- Explain the use of generic version of incremental computation
- Explain that it is faster without writing additional code given the normal function.

Acknowledgements

Write acknowledgement

Contents

1	Introduction	5
1.1	Contributions	7
2	Specific Implementation	8
2.1	Merkle Tree (TreeH)	10
2.2	Zipper	11
2.2.1	Zipper TreeH	12
3	Datatype-Generic Programming	13
3.1	Introduction	13
3.2	Explicit recursion	15
3.3	Sums of Products	16
3.4	Mutually recursive datatypes	16
4	Generic Implementation	18
4.1	Regular	18
4.2	Generic Zipper	20
4.3	Memory Management	22
4.3.1	Cache Addition Policies	22
4.3.2	Cache Replacement Policies	22
4.4	Pattern Synonyms	23
5	Experiments	25
5.1	Method	25
5.2	Results	27
5.2.1	Execution Time	27
5.2.2	Memory Usage	28
5.2.3	Comparison Cache Addition Policies	29
6	Conclusion	33
6.1	Future work	33
A	Generic Programming	34
A.1	Functor instances for Pattern Functors	34
B	Results	35
B.1	Individual benchmark results - Worst Case	35
B.1.1	Linear trend line	35
B.1.2	Logarithmic trend line	37
B.2	Individual benchmark results - Average Case	38
B.2.1	Linear trend line	38

B.2.2	Logarithmic trend line	40
B.3	Individual benchmark results - Best Case	41
B.3.1	Linear trend line	41
B.3.2	Logarithmic trend line	43

Todo list

Write abstract	1
Write acknowledgement	1
Write more context	5
Add a problem statement	5
Add why my thesis is unique	5
Add the explanation of that we use the Zipper for updating the Tree	6
Check if hash and hash values are correctly used	8
Add the implementation HashMap[17]	9
Compare paper “A Lightweight Approach to Datatype-Generic Rewriting”	13
Use "representation types" instead of "pattern functor" for these datatypes	13
Add a description of what a universe is	16
Write about why Regular is chosen	18
Show how to implement this	22
Change from recursion call to recursion depth	22
Is this a good statement?	23
Show how to create these policies	23
Map over the results which removes the not valid values	23
Should I add this?	29
Move to appendix	30
Write future work	33

1

Introduction

Write more context

Add a problem statement

Add why my thesis is unique

Incremental computation is an approach to improve performance by computing only the result of the changed parts of the input, instead of computing the entire result. To determine what parts of the input has changed, we provide a *Zipper* which can be used by the user to efficiently update the previous input to a new input and keeps track of the changes to the input. As a result, we know which parts of the input stays the same. Therefore, we can reuse the results of a previously performed computation for the same input¹.

The task of implementing incremental computation for Haskell is explained in Chapter 2. However, the task for implementing incremental computation for multiple datatypes, becomes repetitive and error-prone. To prevent writing implementations for every datatype, we use datatype-generic programming. Datatype-generic programming is a technique that uses the structure of a datatype to define functions for a large class of datatypes, which is explained in more detail in Chapter 3. Then, using the Haskell generics library **regular**, we implement the generic version. In Chapter 4, the generic implementation is explained and additionally describe: complexity, different garbage collection strategies and different data structures for storing results.

To illustrate how the incremental computation functionality is used, an example is presented which implements the **max-path-sum** in a non-incremental and an incremental manner. The max-path-sum function traverses through a tree and returns the sum of the path with the highest value. For the example, we first define the definition of a tree and an example tree which is used as input for the max-path-sum. Then, the non-incremental function (**maxPathSum**) is implemented, by returning the value for the leaf; and for the node recursively calling the children and determine the highest value between the children and adding it to its own value.

```
data BinTree = Leaf Int
              | Node BinTree Int BinTree
```

¹The technique for reusing results based on input is called *memoization*

```

exampleTree :: BinTree
exampleTree = Node (Node (Leaf 8) 7 (Leaf 1)) 3 (Node (Leaf 5) 4 (Leaf 2))

-- The non-incremental implementation of max-path-sum
maxPathSum :: BinTree -> Int
maxPathSum (Leaf x)      = x
maxPathSum (Node l x r) = x + max (maxPathSum l) (maxPathSum r)

> maxPathSum exampleTree
18

> let exampleTree' = update (const (Leaf 6)) [Bttm] exampleTree

> maxPathSum exampleTree'
16

```

The incremental implementation, first adds hashes to the input which represent the internal data structure to the data structure (also known as a *merkle tree*). This is used to efficiently check what parts of the input has changed. Then, the initial computation is performed which gives the result and a map of all the intermediate results. Next, the tree gets updated, in this case the bottom left leaf (Leaf 8) is replaced with Leaf 6. Finally, the updated tree is recomputed reusing the previously computed results. As a result, the complete right side of the root node did not have to be recomputed, because it stayed the same.

Add the explanation of that we use the Zipper for updating the Tree

```

-- The incremental implementation of max-path-sum
incMaxPathSum :: BinTree -> Int
incMaxPathSum (Leaf_ x)      = x
incMaxPathSum (Node_ l x r) = x + max l r

-- Add hashes to data structure
> let merkleTree = merkle exampleTree

-- Initial computation
> let (y, m) = cataMerkle incMaxPathSum (merkleTree)
      (18, { "6dd": 18, "5df": 15, "fa0": 8, "8d0": 1, "f3b": 9, "84b": 5
            , "1ad": 2 })

-- Update Tree
> let merkleTree' = update (const (merkle (Leaf 6))) [Bttm] merkleTree

```

```

-- Incremental computation
> cataMerkleMap incMaxPathSum m (merkleTree')
(16, { "6dd": 18, "5df": 15, "fa0": 8, "bbd": 16, "91c": 13, "3af": 6
      , "8d0": 1, "f3b": 9, "84b": 5, "1ad": 2 })

```

1.1 Contributions

The main contributions of the Thesis are the following:

- We define an algorithm for incremental computation over recursive data structures. The algorithm uses hashes for comparing if the data structures are equal in a constant time and a Zipper to efficiently update the recursive data structure without rehashing the entire data structure.
- We use datatype-generic programming to write a generic version of the algorithm, to support a large class of datatypes, namely *regular datatypes*.
- We use pattern synonyms, to make the developer experience the same as implementing a non-incremental algorithm.
- We define cache addition policies and cache replacement policies to optimize the performance/memory usage for different use-cases.

2

Specific Implementation

```
data Tree a = Leaf a
            | Node (Leaf a) a (Leaf a)

sumTree :: Tree Int -> Int
sumTree (Leaf x)      = x
sumTree (Node l x r) = x + (sumTree l) + (sumTree r)
```

Computing a value of a data structure can easily be defined in Haskell, but every time there is a small change in the `Tree`, the entire `Tree` needs to be recomputed. This is inefficient, because most of the computations have already been performed in the previous computation.

To prevent recomputation of already computed values, the technique memoization is introduced. Memoization is a technique where the results of computational intensive tasks are stored and when the same input occurs, the result is reused.

The comparison of two values in Haskell is done with the `Eq` typeclass, which implements the equality operator `(==) :: a -> a -> Bool`. So, an example implementation of the `Eq` typeclass for the `Tree` datatype would be:

```
instance Eq a => Eq (Tree a) where
  Leaf x1      == Leaf x2      = x1 == x2
  Node l1 x1 r1 == Node l2 x2 r2 = x1 == x2 && l1 == l2 && r1 == r2
  _            == _            = False
```

The problem with using this implementation of the `Eq` typeclass for Memoization is that for every comparison of the `Tree` datatype the equality is computed. This is inefficient because the equality implementation has to traverse the complete `Tree` data structure to know if the `Tree`'s are equal.

Check if hash and hash values are correctly used

To efficiently compare the `Tree` datatypes, we need to represent the structure in a manner which does not lead to traversing to the complete `Tree` data structure. This can be accomplished using a `hash` function. A hash function is a process of transforming a data structure into an arbitrary fixed-size value, where the same input always generates the same output.

One disadvantage of using hashes is *hash collisions*. Hash collisions happen when two different

pieces of data have the same hash. This is because a hash function has a limited amount of bits to represent every possible combination of data. To calculate the chance of a hash collision occurring we use the formula $p = \epsilon^{\frac{-k(k-1)}{2N}}$ from *Hash Collision Probabilities*[14]. So, given a common hash function CRC-32[11], which has a digest size of 32bits, to get a 50% chance of a hash collision occurring in a collection, it needs $0.5 = \epsilon^{\frac{-k(k-1)}{2 \times 2^{32}}} \rightarrow k = 77163$ hash values.

Digest size (in bits)	Collection size
32	77163
64	5.06×10^9
128	2.17×10^{19}

Table 2.1: The collection size needed for a 50% chance of getting a hash collision

The hash function ultimately chosen in this paper is the *CityHash*[15]. CityHash is a non-cryptographic hash function which works well for hash tables[13]. It supports two different digest sizes: 64- and 128-bit. The digest size used in this paper will be the 64-bit variant. This digest size is chosen, because we think it gives a good balance between performance and the probability of getting a hash collision.

```
class Hashable a where
    hash :: a -> Digest

instance Show a => Hashable a where
    hash = cityHash64 . show

instance Hashable a => Hashable (Tree a) where
    hash (Leaf x)      = concatHash [hash "Leaf", hash x]
    hash (Node l x r) = concatHash [hash "Node", hash x, hash l, hash r]
```

The hashes can then be used to efficiently compare two *Tree* data structures, without having to traverse the entire *Tree* data structure.

Add the implementation *HashMap*[17]

To keep track of the intermediate results of the computation, we store the results in a *HashMap*[17]. A *HashMap* is an implementation of mapping a *hashable* key to a value. In our next example the *Digest* is the key and the value is the intermediate result.

```
sumTreeInc :: Tree Int -> (Int, HashMap Digest Int)
sumTreeInc l@(Leaf x)      = (x, insert (hash l) x empty)
sumTreeInc n@(Node l x r) = (y, insert (hash n) y (ml <> mr))
    where
        y = x + xl + xr
        (xl, ml) = sumTreeInc l
        (xr, mr) = sumTreeInc r
```

Then after the first computation over the entire `Tree`, we can recompute the `Tree` using the previously created `HashMap`. Thus, when we recompute the `Tree`, we first look in the `HashMap` if the computation has already been performed then return the result. Otherwise, compute the result and store it in the `HashMap`.

```
sumTreeIncMap :: HashMap Digest Int -> Tree Int -> (Int, HashMap Digest Int)
sumTreeIncMap m l@(Leaf x) = case lookup (hash l) m of
  Just x  -> (x, m)
  Nothing -> (x, insert (hash l) x empty)
sumTreeIncMap m n@(Node l x r) = case lookup (hash n) m of
  Just x  -> (x, m)
  Nothing -> (y, insert (hash n) y (ml <> mr))
  where
    y = x + xl + xr
    (xl, ml) = sumTreeIncMap m l
    (xr, mr) = sumTreeIncMap m r
```

Generating a hash for every computation over the data structure is time-consuming and unnecessary, because most of the `Tree` data structure stays the same. The work of Miraldo and Swierstra[8] inspired the use of the Merkle Tree. A Merkle Tree is a data structure which integrates the hashes within the data structure.

2.1 Merkle Tree (TreeH)

First we introduce a new datatype `TreeH`, which contains a `Digest` for every constructor in `Tree`. Then to convert the `Tree` datatype into the `TreeH` datatype, the structure of the `Tree` is hashed and stored into the datatype using the `merkle` function.

```
data TreeH a = LeafH Digest a
             | NodeH Digest (Leaf a) a (Leaf a)

merkle :: Tree Int -> TreeH Int
merkle l@(Leaf x) = LeafH (hash l) x
merkle (Node l x r) = NodeH h l' x r'
  where
    h = hash ["Node", x, getHash l', getHash r']
    l' = merkle l
    r' = merkle r
```

The precomputed hashes can then be used to easily create a `HashMap`, without computing the hashes every time the `sumTreeIncH` function is called.

```
sumTreeIncH :: TreeH Int -> (Int, HashMap Digest Int)
sumTreeIncH (LeafH h x)      = (x, insert h x empty)
sumTreeIncH (NodeH h l x r) = (y, insert h y (ml <> mr))
```

```

where
  y = x + xl + xr
  (xl, ml) = sumTreeInc l
  (xr, mr) = sumTreeInc r

```

The problem with this implementation is, that when the `Tree` datatype is updated, the entire `Tree` needs to be converted into a `TreeH`, which is linear in time. This can be done more efficiently, by only updating the hashes which are impacted by the changes. Which means that only the hashes of the change and the parents need to be updated.

The first intuition to fixing this would be using a pointer to the value that needs to be changed. But because Haskell is a functional programming language, there are no pointers. Luckily, there is a data structure which can be used to efficiently update the data structure, namely the Zipper[5].

2.2 Zipper

The Zipper is a technique for keeping track of how the data structure is being traversed through. The Zipper was first described by Huet[5] and is a solution for efficiently updating pure recursive data structures in a purely functional programming language (e.g., Haskell). This is accomplished by keeping track of the downward current subtree and the upward path, also known as the *location*.

To keep track of the upward path, we need to store the path we traverse to the current subtree. The traversed path is stored in the `Cxt` datatype. The `Cxt` datatype represents three options the path could be at: the `Top`, the path has traversed to the left (`L`), or the path has traversed to the right (`R`).

```

data Cxt a = Top
           | L (Cxt a) (Tree a) a
           | R (Cxt a) (Tree a) a

```

```

type Loc a = (Tree a, Cxt a)

```

```

enter :: Tree a -> Loc a
enter t = (t, Top)

```

Using the `Loc`, we can define multiple functions on how to traverse through the `Tree`. Then, when we get to the desired location in the `Tree`, we can call the `modify` function to change the `Tree` at the current location.

Eventually, when every value in the `Tree` has been changed, the entire `Tree` can then be rebuilt using the `Cxt`. By recursively calling the `up` function until the top is reached, the current subtree gets rebuilt. And when the top is reached, the entire tree is then returned.

```

left :: Loc a -> Loc a
left (Node l x r, c) = (l, L c r x)

```

```

right :: Loc a -> Loc a
right (Node l x r, c) = (r, R c l x)

up :: Loc a -> Loc a
up (t, L c r x) = (Node t x r, c)
up (t, R c l x) = (Node l x t, c)

modify :: (Tree a -> Tree a) -> Loc a -> Loc a
modify f (t, c) = (f t, c)

leave :: Loc a -> a
leave (t, Top) = t
leave l        = top (up l)

> leave $ modify (const (Leaf 4)) $ left $ enter (Node (Leaf 1) 2 (Leaf 3))
(Node (Leaf 4) 2 (Leaf 3))

```

2.2.1 Zipper TreeH

The implementation of the Zipper for the `TreeH` datatype is the same as for the `Tree` datatype. However, the `TreeH` also contains the hash of the current and underlying data structure. Therefore, when a value is modified in the `TreeH`, all the parent nodes of the modified value needs to be updated.

The `updateLoc` function modifies the value at the current location, then checks if the location has any parents. If the location has any parents, go up to that parent, update the hash of that parent and recursively update the parents hashes until we are at the top of the data structure. Otherwise, return the modified locations, because all the other hashes are not affected by the change.

```

updateLoc :: (TreeH a -> TreeH a) -> Loc a -> Loc a
updateLoc f l = if top l' then l' else updateParents (up l')
  where
    l' = modify f l
    updateParents :: Loc a -> Loc a
    updateParents (Loc x Top) = Loc (updateHash x) Top
    updateParents (Loc x cs)  = updateParents $ up (Loc (updateHash x) cs)

```

Then, the `update` function can be defined using the `updateLoc` function, by first traversing through the data structure with the given directions. Then modifying the location using the `updateLoc` function and then leave the location and the function results in the updated data structure.

```

update :: (TreeH a -> TreeH a) -> [Loc a -> Loc a] -> TreeH a -> TreeH a
update f dirs t = leave $ updateLoc f l'
  where
    l' = applyDirs dirs (enter t)

```

3

Datatype-Generic Programming

Compare paper “A Lightweight Approach to Datatype-Generic Rewriting”

The implementation in Chapter 2 is an efficient implementation for incrementally computing the summation over a `Tree` datatype. However, when we want to implement this functionality for a different datatype, a lot of code needs to be copied while the process remains the same. This results in poor maintainability, is error-prone and is in general boring work.

An example of reducing manual implementations for datatypes is the *deriving* mechanism in Haskell. The built-in classes of Haskell, such as `Show`, `Ord`, `Read`, can be derived for a large class of datatypes. However, deriving is not supported for custom classes. Therefore, we use *Datatype-Generic Programming*[3] to define functionality for a large class of datatypes.

In this chapter, we introduce Datatype-Generic Programming, also known as *generic programming* or *generics* in Haskell, as a technique that uses the structure of a datatype to define functions for a large class of datatypes. This prevents the need to write the previously defined functionality for every datatype.

3.1 Introduction

There are multiple generic programming libraries, however to demonstrate the workings of generic programming we will be using a single library as inspiration, named `regular`[7]. Here the generic representation of a datatype is called a *pattern functor*. A pattern functor is a stripped-down version of a data type, by only containing the constructor but not the recursive structure. The recursive structure is done explicitly by using a fixed-point operator.

First, the pattern functors defined in `regular` are 5 core pattern functors and 2 meta information pattern functors. The core pattern functors describe the datatypes. The meta information pattern functors only contain information (e.g., constructor name) but not any structural information.

Use "representation types" instead of "pattern functor" for these datatypes

```
data U r      = U                -- Empty constructors
data I r      = I r              -- Recursive call
```

```

data K a r      = K a                -- Constants
data (f :+: g) r = L (f r) | R (g r) -- Sums (Choice)
data (f **: g) r = (f r) **: (g r)   -- Products (Combine)

```

The conversion from regular datatypes into pattern functors is done by the `Regular` type class. The `Regular` type class has two functions. The `from` function converts the datatype into a pattern functor and the `to` function converts the pattern functor back into a datatype. In `regular`, the pattern functor is represented by a type family. Then using the `Regular` conversion to a pattern functor, we can write the `Tree` datatype from Chapter 2 as:

```

type family PF a :: * -> *

class Regular a where
  from :: a -> PF a a
  to   :: PF a a -> a

type instance PF (Tree a) = K a                -- Leaf
                               :+: (I :+: K a :+: I) -- Node

class Regular (Tree a) where
  from (Leaf x)      = L (K x)
  from (Node l x r) = R (I l :+: K x :+: I r)

  to (L (K x))          = Leaf x
  to (R (I l :+: K x :+: I r)) = Node l x r

```

To demonstrate the workings of generic programming, we are going to implement a simple generic function which determines the length of an arbitrary datatype. First, we define the length function within a type class. The type class is used, to define how to determine the length for every pattern functor `f`.

```

class GLength f where
  glength :: (a -> Int) -> f a -> Int

```

Writing instances for the empty constructor `U` and the constants `K` is simple because both pattern functors return zero. The `U` pattern functor returns zero, because it does not contain any children. The `K` pattern functor returns zero, because we do not count constants for the length.

```

instance GLength U where
  glength _ _ = 0

instance GLength (K a) where
  glength _ _ = 0

```

The instances for sums and products pattern functors are quite similar. The sums pattern functor recurses into the specified choice. The product pattern functor recurses in both constructors and

combines them.

```
instance (GLength f, GLength g) => GLength (f :+: g) where
  glength f (L x) = glength f x
  glength f (R x) = glength f x

instance (GLength f, GLength g) => GLength (f **: g) where
  glength f (x **: y) = glength f x + glength f y
```

The instance for the recursive call `I` needs an additional argument. Because, we do not know the type of `x`, so an additional function (`f :: a -> Int`) needs to be given which converts `x` into the length for that type.

```
instance GLength I where
  glength f (I x) = f x
```

Then using the `GLength` instances for all pattern functors, a function can be defined using the generic length function. By first, converting the datatype into a generic representation, then calling `glength` given recursively itself, and for every recursive call increase the length by one.

```
length :: (Regular a, GLength (PF a)) => a -> Int
length = 1 + glength length (from x)
```

```
> length [1, 2, 3]
3
> length (Node (Leaf 1) 2 (Leaf 3))
3
> length {"1": 1, "2": 2, "3": 3}
3
```

3.2 Explicit recursion

The previous implementation of the `length` function is implemented for a shallow representation. A shallow representation means that the recursion of the datatype is not explicitly marked. Therefore, we can only convert one layer of the value into a generic representation using the `from` function.

Alternatively, by marking the recursion of the datatype explicitly, also called the deep representation, the entire value can be converted into a generic representation in one go. To mark the recursion, a fixed-point operator (`Fix`) is introduced. Then, using the fixed-point operator we can define a `from` function that given the pattern functors have an instance of `Functor`¹, return a generic representation of the entire value.

```
data Fix f = In { unFix :: f (Fix f) }
```

¹The `Functor` instances for the pattern functors can be found in Section A.1


```

deepFrom :: (Regular a, Functor (PF a)) => a -> Fix (PF a)
deepFrom = In . fmap deepFrom . from

```

Subsequently, we can define a `cata` function which can use the explicitly marked recursion by applying a function at every level of the recursion. Then using the `cata` function we can define the same `length` function as in the previous section, but just in a single line. However, this deep representation does come at the cost that the implementation is less efficient than the shallow representation.

```

cata :: Functor f => (f a -> a) -> Fix f -> a
cata = f . fmap (cata f) . unFix

length' :: (Regular a, GLength (PF a), Functor (PF a), Foldable (PF a))
        => a -> Int
length' = cata ((1+) . sum) . deepFrom

```

3.3 Sums of Products

Add a description of what a universe is

A different way of describing datatypes in a generic representation, besides pattern functors, are *Sums of Products*[18] (SOP). SOP is a generic representation with additional constraints which more faithfully reflects the Haskell datatypes: each datatype is a single n-ary sum, where each component of the sum is a single n-ary product. The SOP universe is described using *codes* of kind `[[*]]`. The outer list describes an n-ary sum, representing the choice between constructors and each inner list an n-ary products, representing the constructor arguments. The code of kind `[[*]]` can then be interpreted to describe Haskell datatypes of kind `*`. To define a code, the tick mark ``` is used to lift the list to a type-level.

```
Code (Tree a) = `[ `[a], `[Tree a, a, Tree a]
```

The usage of SOP has a positive effect on expressing generic functions easily or at all. Additionally, the SOP completely divides the structural representation from the metadata. As a result, you do not have to deal with metadata while writing generic functions. However, the additional constraints on the generic representation makes the SOP universe size comparatively bigger than pattern functors. Therefore, it is more complex to extend the SOP than for pattern functors.

3.4 Mutually recursive datatypes

A large class of datatypes is supported by the previous section, namely *regular* datatypes. Regular datatypes are datatypes in which the recursion only goes into the same datatype. However, if we want to support the abstract syntax tree of many programming languages, we need to support datatypes which can recurse over different datatypes, namely mutually recursive datatypes.

```
data Tree a = Empty
            | Node (a, Forest a)

data Forest a = Nil
              | Cons (Tree a) (Forest a)
```

To support mutually recursive datatypes, we need to keep track of which recursive position points to which datatype. This is accomplished by using *indexed fixed-points*[19]. The indexed fixed-points works by creating a type family φ with n different types, where the types inside the family represent the indices for different kinds ($*_{\varphi}$). Using the limited set of kinds we can determine the type for the recursive positions. Thus, supporting mutually recursive datatypes is possible, but it adds a lot more complexity.

4

Generic Implementation

4.1 Regular

The **regular** generic programming library was chosen, because it has the smallest universe size compared to other libraries. Therefore, implementing the generic implementation is less complex. However, **regular** does not support *Sums of Products* and *mutually recursive datatypes*, but we expect that the results will be meaningful without supporting these features.

Write about why Regular is chosen

The first step of the incremental computation was computing the Merkle Tree. In other terms, we need to store the hash of the data structure inside the data structure. We accomplish this by defining a new type **Merkle** which is a fixed-point over the data structure where each of the recursive positions contains a hash (**K Digest**).

```
type Merkle f = Fix (f :: K Digest)
```

But, before the hash can be stored inside the data structure, the hash needs to be computed from the data structure. For this we need to know how to hash the generic datatypes. We introduce a typeclass named **Hashable** which defines a function **hash**, which converts the **f** datatype into a **Digest** (also known as a *hash value*).

```
class Hashable f where
  hash :: f (Merkle g) -> Digest
```

```
digest :: Show a => a -> Digest
digest = digestStr . show -- converts a string into a hash value
```

The **Hashable** instance of **U** is simple. The **digest** function is used to convert the constructor name **U** into a **Digest**. The **K** also uses the **digest** function to convert the constructor name into a **Digest**, but it also calls **digest** on the constant value of **K**. Therefore, the type of the value of **K** needs an instance for **Show**. Then both digests are combined into a single digest.

```
instance Hashable U where
  hash _ = digest "U"
```

```
instance (Show a) => Hashable (K a) where
  hash (K x) = digestConcat [digest "K", digest x]
```

The instances for `:+:`, `:*` and `C` are quite similar as the instance for the `K` datatype. However, the value inside the constructor are recursively called.

```
instance (Hashable f, Hashable g) => Hashable (f :+: g) where
  hash (L x) = digestConcat [digest "L", hash x]
  hash (R x) = digestConcat [digest "R", hash x]
```

```
instance (Hashable f, Hashable g) => Hashable (f **: g) where
  hash (x **: y) = digestConcat [digest "P", hash x, hash y]
```

```
instance (Hashable f) => Hashable (C c f) where
  hash (C x) = digestConcat [digest "C", hash x]
```

The `I` instance is different from the previous instances, because the recursive position is already converted into a Merkle Tree. Thus, we need to get the computed hash from the recursive position, digest the datatype name and combine the digests.

```
instance Hashable I where
  hash (I x) = digestConcat [digest "I", getDigest x]
  where
    getDigest :: Fix (f **: K Digest) -> Digest
    getDigest (In _ **: K h) = h
```

The `hash` implementation can then be used to define a function `merkleG` which converts from a shallow generic representation, to a generic representation where one layer of recursive positions contains a hash value.

Subsequently, we can define a function `merkle` which converts the entire generic representation, into a generic representation where every recursive position contains a hash value. We can define `merkle` using the same implementation as in Section 3.2, but we add a step where after all the children are recursively called, the `merkleG` function is applied.

```
merkleG :: Hashable f => f (Merkle g) -> (f **: K Digest) (Merkle g)
merkleG f = f **: K (hash f)
```

```
merkle :: (Regular a, Hashable (PF a), Functor (PF a))
=> a -> Merkle (PF a)
merkle = In . merkleG . fmap merkle . from
```

The `Merkle` representation can then be used to define a function `cataMerkleState` which given a function `alg :: (f a -> a)` which converts the generic representation `f a` into a value of type `a` and the `Merkle f` data structure, and returns a `State of (HashMap Digest a) a`. The `cataMerkleState` function starts with retrieving the `State`, which keeps track of the intermediate

results and stores them into a `HashMap Digest a`. Then, given the hash value of the recursive position, we look into the `HashMap` if the value has been computed. If the value has been computed, then return the value. Otherwise, recursively compute all the children, apply the given function `alg`, insert the new value into the `HashMap` and return the computed value.

```
cataMerkleState :: (Functor f, Traversable f)
                => (f a -> a) -> Merkle f -> State (HashMap Digest a) a
cataMerkleState alg (In (x :: K h))
  = do m <- get
      case lookup h m of
        Just a  -> return a
        Nothing -> do y <- mapM (cataMerkleState alg) x
                       let r = alg y
                       modify (insert h r) >> return r
```

The `cataMerkleState` function can be used, but to execute the function we first need to give the function a `HashMap Digest a`. To simplify the use of `cataMerkleState`, we define a function `cataMerkle`, which executes the `cataMerkleState` with an empty `HashMap` and returns the final computed result and the final state as the result.

```
cataMerkle :: (Functor f, Traversable f)
            => (f a -> a) -> Merkle f -> (a, HashMap Digest a)
cataMerkle alg t = runState (cataMerkleState alg t) empty
```

Finally, we have all the necessary functionality defined to write a function over the generic representation and automatically generate all the intermediate results and the final result. The example below computes the sum over the generic representation of `Tree` by adding all the values of the leaf and nodes.

```
cataSum :: Merkle (PF (Tree Int)) -> (Int, HashMap Digest Int)
cataSum = cataMerkle
  (\case
    L (C (K x))          -> x
    R (C (I l :: K x :: I r)) -> l + x + r
  )
```

```
> cataSum $ merkle $ Node (Leaf 1) 2 (Leaf 3)
(6, {"931090e5": 1, "7d1ef1c9": 3, "ba811ed5": 6})
```

4.2 Generic Zipper

For the implementation of a generic zipper, we need to define (A) a datatype which keeps track of the location inside the data structure, (B) a context which keeps track of the locations which have been traversed through (C) and functions which facilitate traversing through the data structure.

The implementation of the general zipper is based on the paper “Generic representations of tree transformations” by Bransen and Magalhaes[2].

The context (`Ctx`) is implemented using a type family¹[6]. Type families can be defined in two manners, standalone or associated with a type class. For the explanation we use the standalone definition because the code is less clumped, making it easier to explain. However, the actual implementation uses type synonym families, which makes it clearer how the type should be used and has better error messages.

The standalone definition uses a `data family`. Then, we want to write for every representation type an instance on how to represent the representation type in the context.

```
data family Ctx (f :: * -> *) :: * -> *
```

The `K` and `U` representation-types do not have a datatype, because these representation-types cannot be traversed through.

```
data instance Ctx (K a) r
data instance Ctx U      r
```

The sum representation-type does get traversed, but only has one choice by either traversing the left `CL` or the right `CR` side.

```
data instance Ctx (f :+: g) r = CL (Ctx f r) | CR (Ctx g r)
```

The product representation-type does have a choice between two traversals, either traverse through the left side and store the right side or traverse through the right side and store the left side.

```
data instance Ctx (f :*: g) r = C1 (Ctx f r) (g r) | C2 (f r) (Ctx g r)
```

The recursive representation-type does not recursively go into a new `Ctx`, but into the recursive type `r`, thus we only need to define datatype which indicates that it is a recursive position.

```
data instance Ctx I r = CId
```

The `Zipper` type class can now be defined using the `Ctx`. There are 6 primary functions, which we need to build navigating functions. The `cmap` function works like the `fmap`, but over contexts. The `fill` function fills the hole in a context with a given value, which is used to reconstruct the data structure. The final 4 functions (`first`, `last`, `next` and `prev`) are the *primary* navigation operations used to build *interface* functions such as `left`, `right`, `up`, `down`, etc.

```
class Functor f => Zipper f where
  cmap      :: (a -> b) -> Ctx f a -> Ctx f b
  fill      :: Ctx f a -> a -> f a
  first, last :: f a -> Maybe (a, Ctx f a)
  next, prev  :: Ctx f a -> a -> Maybe (a, Ctx f a)
```

Finally, we can define the location `Loc` using the `Ctx` and the `Zipper`. The location takes a datatype with the constraints that it has an instance for `Regular` and a pattern functor which works with

¹“Type families are to vanilla data types what type class methods are to regular functions.”[4]

the `Zipper` type class, a list of contexts and returns a location datatype `Loc`.

```
data Loc :: * -> * where
  Loc :: (Regular a, Zipper (PF a)) => a -> [Ctx (PF a) a] -> Loc a
```

To use the `Merkle` type, we need to fulfill the previous constraints. Thus, we need to define a type instance for the pattern functor and an instance for the `Regular` typeclass. The pattern functor type instance is the same definition as the `Merkle` type but without the fixed-point. The `Regular` instance for the `Merkle` type is folding/unfolding the fixed-point. Moreover, the `Merkle` type also needs to update the digests of its parents when a value is changed. This is accomplished in the same manner as in Section 2.2.1.

```
type instance PF (Merkle f) = f :: K Digest
instance Regular (Merkle f) where
  from = out
  to   = In
```

4.3 Memory Management

In the optimal case, each intermediate value can be stored. Unfortunately, we are limited by the amount of memory there is available on the machine. Therefore, the amount of intermediate values needs to be limited. There are two suggestions of limiting the amount of intermediate values with *Cache Addition Policies* and *Cache Replacement Policies*. Cache addition policies are policies which indicate if an intermediate value can be added to the cache. And the cache replacement policies are policies which remove intermediate values from the cache based on metrics (e.g., recency, frequency, computational cost, etc.).

4.3.1 Cache Addition Policies

Show how to implement this

Change from recursion call to recursion depth

An example of a cache addition policy would be to determine the amount of recursive calls a data structure has. So, given a node in a tree we then know how many times we have to go into recursion before reaching a leaf. Using that information we can define a filter for the cache, where only elements with a recursion depth of $> i$ can be added to the cache. This filters out a large amount of intermediate values and could potentially lead to a speed-up, because the lookup in the `HashMap` could take longer than the calculation itself. However, this does not prevent the memory from filling-up. Instead, intermediate values need to be removed when there is no memory left. This is where cache replacement policies are needed.

4.3.2 Cache Replacement Policies

Cache replacement policies can be based on multiple metrics.

- Recency: when added in the HashMap
- Frequency: how many times the value is looked up in the HashMap.
- Computational cost: how many recursive calls are needed
- Iterations: how many iterations the object stays in the HashMap.

The cache replacement policy can be based on a single metric of a combination of multiple metrics. Unfortunately, we cannot determine the overall best cache replacement policy, because the best cache replacement policy is application specific as stated in “Selective memoization”². As a result, this paper will only describe possible policies which can be used by developers, but not show any results.

Is this a good statement?

Show how to create these policies

Map over the results which removes the not valid values

Possible cache replacement policies:

- Random replacement: remove random elements from the cache.
- Recency-based policies: remove elements based on when the element was added.
- Frequency-based policies: remove elements based on the amount of lookups.
- Computational cost policies: remove elements based on the amount of time it takes to compute the value.
- Combination of policies mentioned above. For example, combine the frequency based policy with the computation cost policy.

4.4 Pattern Synonyms

The developer experience using `cataMerkle` is difficult, because the developer needs to know the pattern functor of its datatype to define a function and the function definitions are quite verbose. To make the use of `cataMerkle` easier, we introduce *pattern synonyms*[12].

Pattern synonyms add an abstraction over patterns, which allows the user to move additional logic from guards and case expressions into patterns. For example, the pattern functor of the `Tree` datatype can be represented using a `pattern`.

```
pattern Leaf_ :: a -> PF (Tree a) r
pattern Leaf_ x <- L (C (K x)) where
  Leaf_ x = L (C (K x))
```

²“In general the replacement policy must be application-specific, because, for any fixed policy, there are programs whose performance is made worse by that choice.”[1]


```

pattern Node_ :: r -> a -> r -> PF (Tree a) r
pattern Node_ l x r <- R (C (I l :: K x :: I r)) where
  Node_ l x r = R (C (I l :: K x :: I r))

```

The previously defined `patterns` can then be used to define the `cataSum` as the original datatype `Tree`, but the constructor names leading with an additional underscore. However, writing all the patterns for all pattern functors is an arduous task. Luckily, we can use `TemplateHaskell` to generate the pattern synonyms³.

```

cataSum :: Merkle (PF (Tree Int)) -> (Int, HashMap Digest Int)
cataSum = cataMerkle
  (\case
    Leaf_ x      -> x
    Node_ l x r -> l + x + r
  )

```

³An example of using `TemplateHaskell` to generate pattern synonyms can be found at *Generics-MRSOP-TH*

5

Experiments

5.1 Method

We conducted experiments over three type of functions: the `Cata Sum`, `Generic Cata Sum` and `Incremental Cata Sum`. The `Cata Sum` is the simple function which traverses through the entire tree and sums all the values. The `Generic Cata Sum` is the initial `Incremental Cata Sum`, which starts with an empty `HashMap`. And the `Incremental Cata Sum`, which already has a `HashMap` filled with intermediate results and keeps track over multiple iterations.

```
cataSum :: Tree Int -> Int
cataSum (Leaf x)      = x
cataSum (Node l x r) = x + cataSum l + cataSum r

genericCataSum :: Merkle (PF (Tree Int)) -> (Int, HashMap Digest Int)
genericCataSum = cataMerkle
  (\case
    Leaf_ x      -> x
    Node_ l x r -> l + x + r
  )

incCataSum :: HashMap Digest Int
            -> Merkle (PF (Tree Int)) -> (Int, HashMap Digest Int)
incCataSum = cataMerkleMap
  (\case
    Leaf_ x      -> x
    Node_ l x r -> l + x + r
  )
```

The experiments will be benchmarks executed with the Haskell package `criterion`[10]. `Criterion` performs the benchmarks multiple times to get an average result. The benchmarks will track two metrics: the execution time and the memory usage. The execution time will be in seconds and the memory usage will be the max-bytes used. The results gather for execution time comes from the `criterion` package, however the memory usage will not come from the package. This is because

criterion only keeps track of the memory allocation and not usage. Therefore, we measure the memory usage with the GHC profiler[16] and profile every benchmark individually to know its memory usage.

To test how well the three functions perform, we perform three types of updates, multiple times. These benchmarks will be based on the three type of cases: worst, average and best. (1) The worst case updates the lowest left leaf to a new leaf. (2) The average case updates a node in the middle with a new leaf. (3) And the best case replaces the left child of the root-node with a new leaf.

```
worstCase :: Merkle (PF (Tree Int)) -> Merkle (PF (Tree Int))
worstCase = update (const (Leaf i)) [Btm]

averageCase :: Merkle (PF (Tree Int)) -> Merkle (PF (Tree Int))
averageCase = update (const (Leaf i)) (replicate n Dwn)
  where
    n = round (logBase 2.0 (fromIntegral n) / 2.0)

bestCase :: Merkle (PF (Tree Int)) -> Merkle (PF (Tree Int))
bestCase = update (const (Leaf i)) [Dwn]
```

5.2 Results

5.2.1 Execution Time

For all the benchmarks, the Incremental Cata Sum is faster when the tree contains more than 10^3 nodes. However, for every benchmark the execution time is better or worse depending on the type of update that is performed. The best case scenario has the biggest difference, then the average case and the closest performance difference is the worst case. The execution time for Cata Sum and Generic Cata Sum seems to be linear with the amount of nodes and the Incremental Cata Sum is constant/logarithmic. The discrepancy in the range between $10^2 - 10^3$ is probably because the amount of nodes is too low to get stable results.

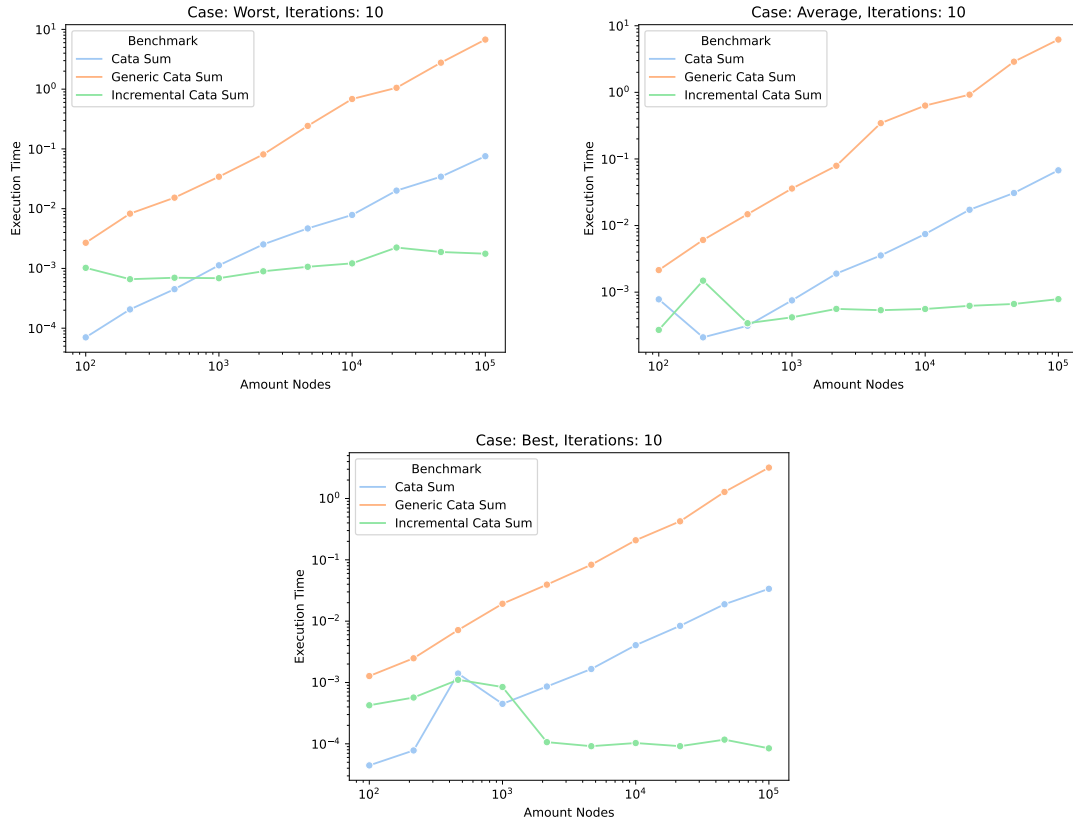


Figure 5.1: The execution time over 10 executions for the Worst, Average and Best case.

5.2.2 Memory Usage

The memory usage results are all linear with the amount of nodes. The Cata Sum uses the least amount of memory, then the Incremental Cata Sum and the most used bytes is the Generic Cata Sum. We think that the Incremental Cata Sum uses less memory than the Generic Cata Sum, because the Incremental Cata Sum needs to load-in fewer parts of the data structure than the Generic Cata Sum. Also, the same discrepancy occurs in the range between $10^2 - 10^3$ as with the execution time.

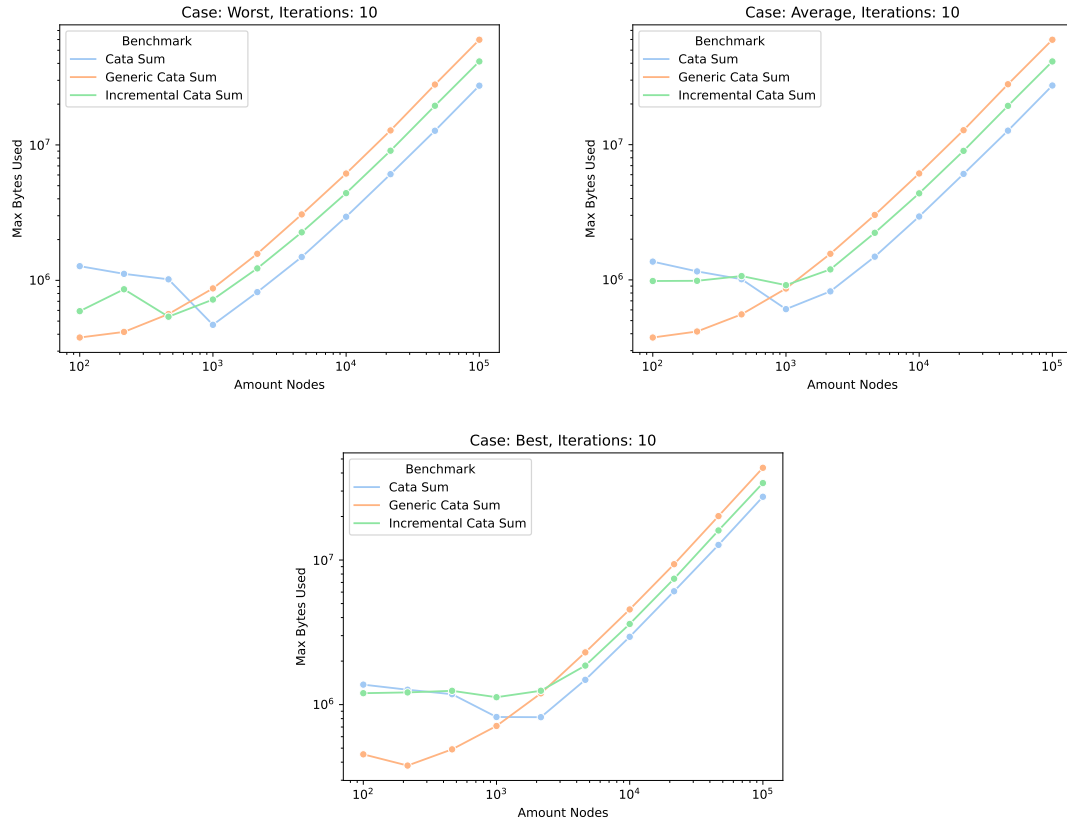


Figure 5.2: The max-bytes-used over 10 executions for the Worst, Average and Best case.

Should I add this?

5.2.3 Comparison Cache Addition Policies

Minimum of 5 recursive calls

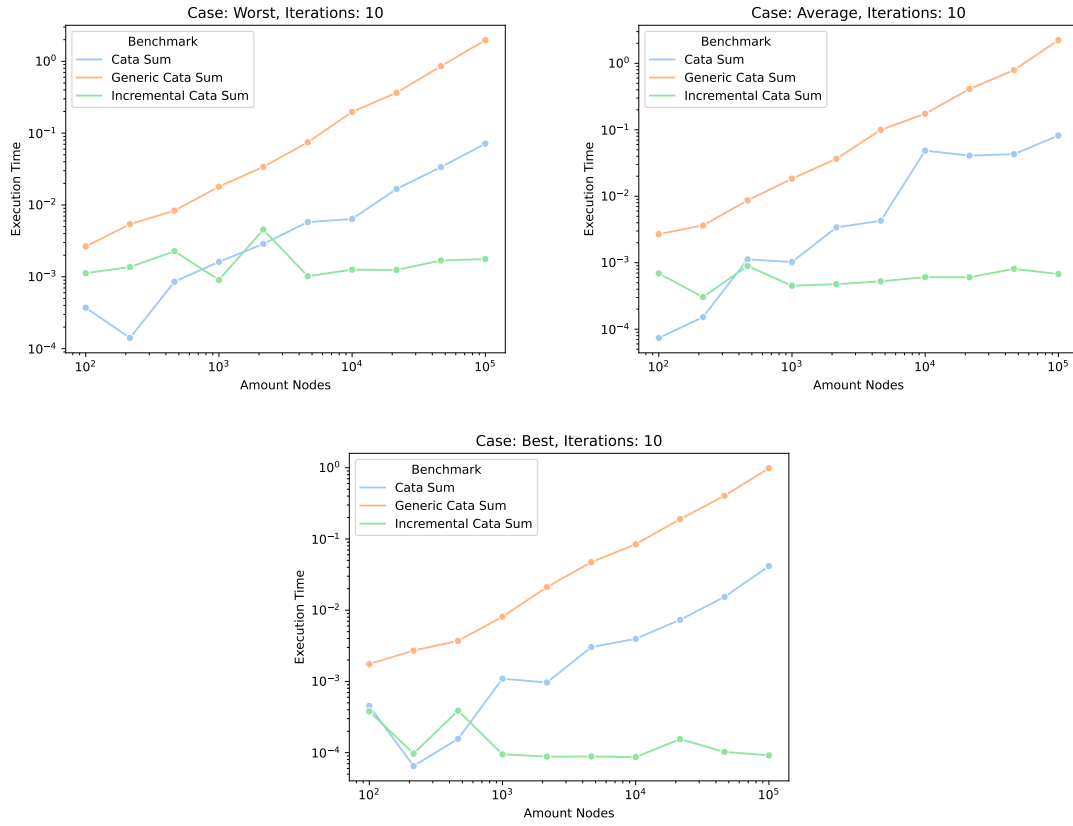


Figure 5.3: The execution time over 10 executions for the Worst, Average and Best case, where the minimum amount of recursive calls is 5.

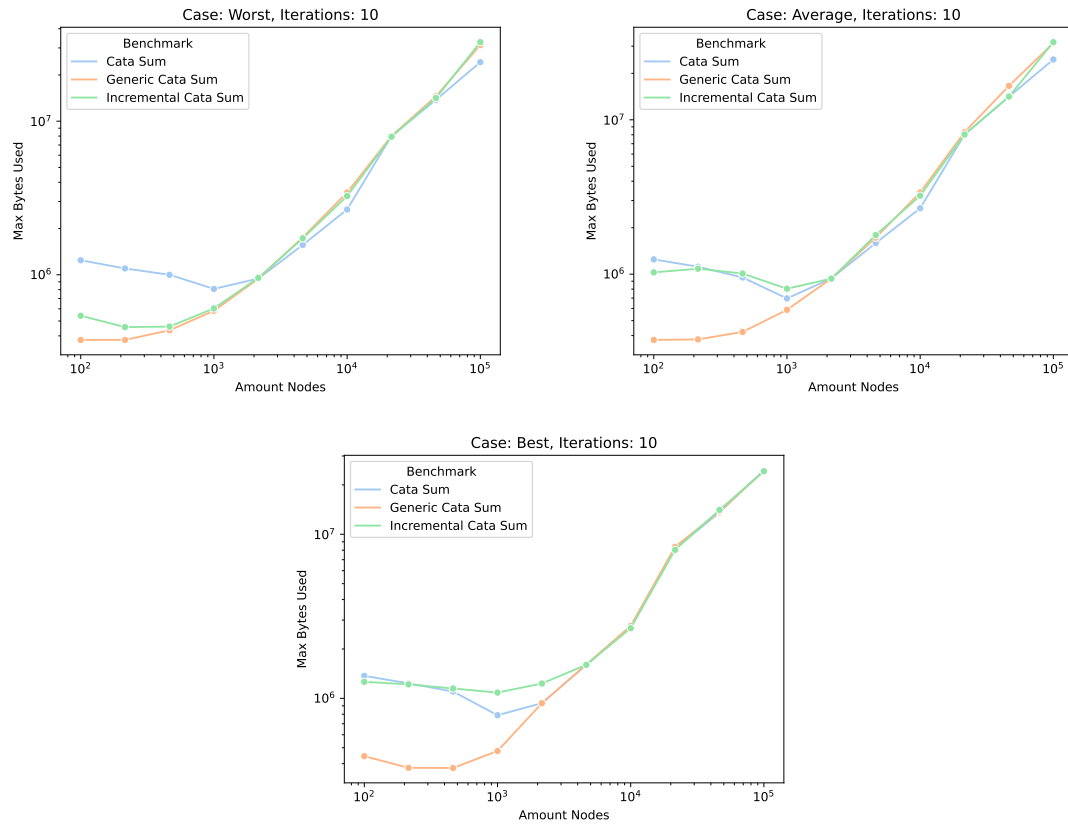


Figure 5.4: The max-bytes-used over 10 executions for the Worst, Average and Best case, where the minimum amount of recursive calls is 5.

Move to appendix

Minimum of 10 recursive calls

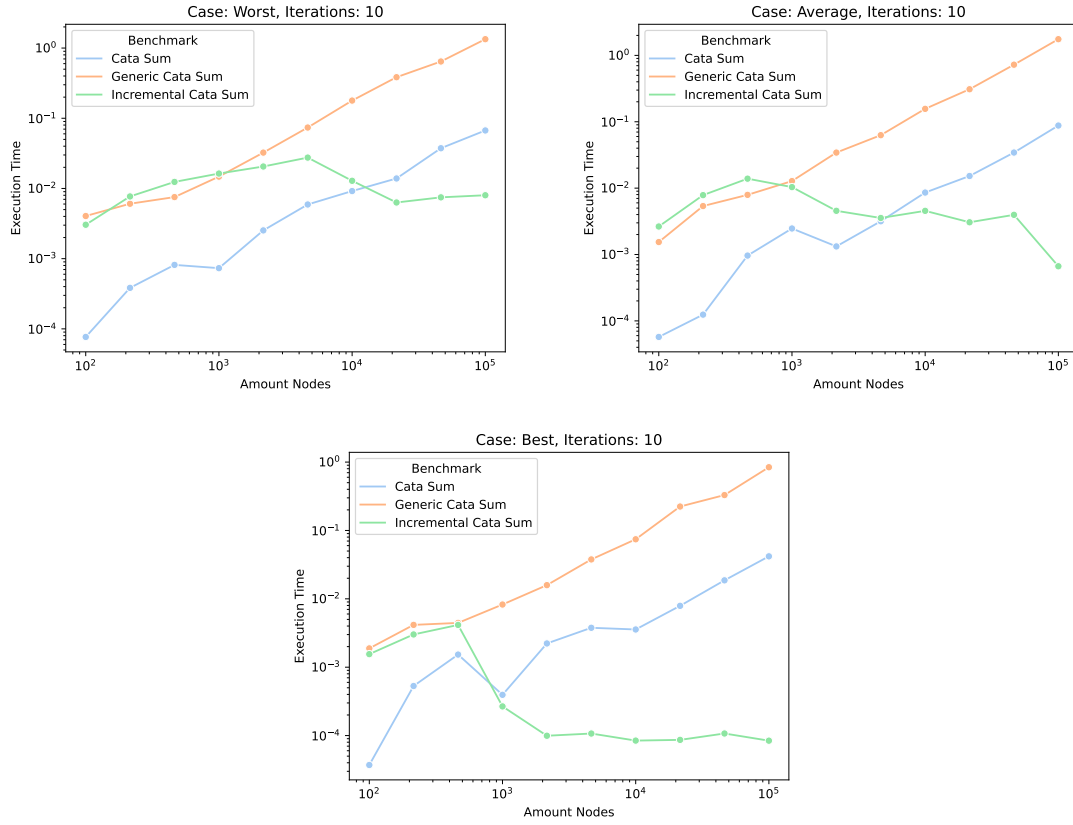


Figure 5.5: The execution time over 10 executions for the Worst, Average and Best case, where the minimum amount of recursive calls is 10.

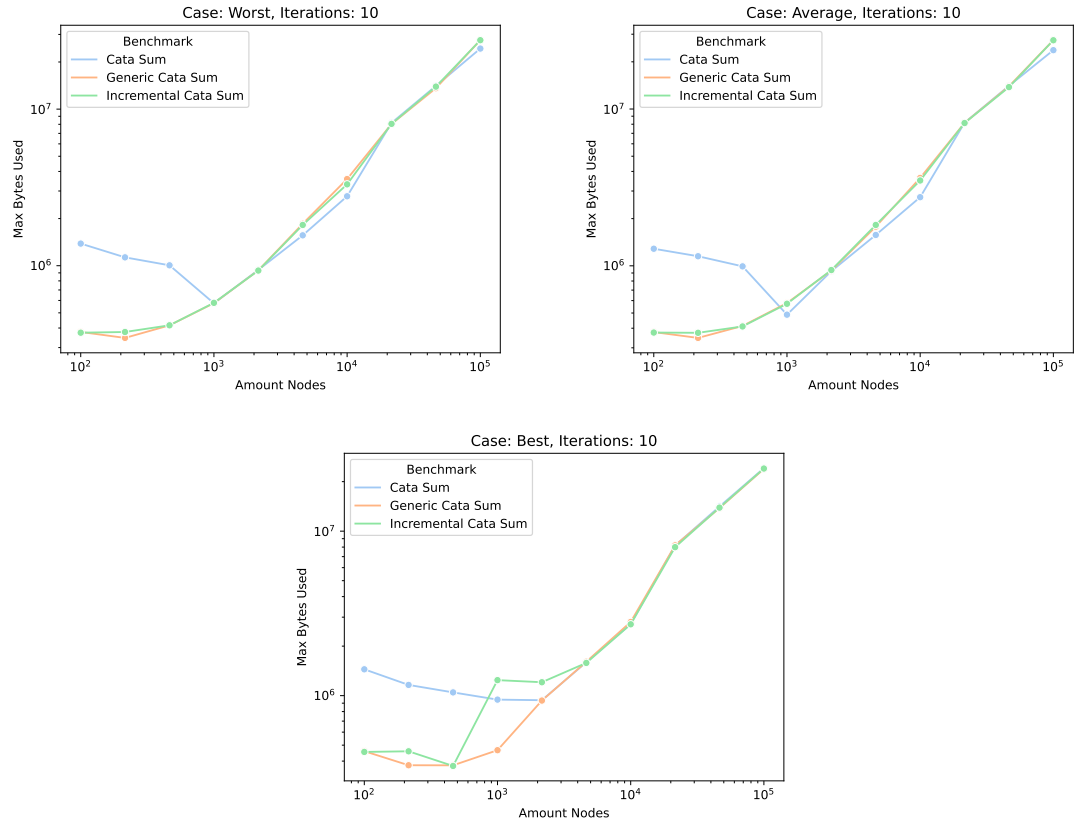


Figure 5.6: The max-bytes-used over 10 executions for the Worst, Average and Best case, where the minimum amount of recursive calls is 10.

6

Conclusion

Setup of conclusion:

- Start with the goal/contributions of the paper.
- Show that the results are successful
- Conclude with a one-liner.

6.1 Future work

Write future work

- Only results for a sum over data structure computation (no real-world results / only synthetic results)
- Only support regular datatypes (could support a larger class of datatypes, for example, mutually recursive datatypes)
- Does not use Sums-of-Products
-



Generic Programming

A.1 Functor instances for Pattern Functors

```
instance Functor I where
  fmap f (I r) = I (f r)
```

```
instance Functor (K a) where
  fmap _ (K a) = K a
```

```
instance Functor U where
  fmap _ U = U
```

```
instance (Functor f, Functor g) => Functor (f :+: g) where
  fmap f (L x) = L (fmap f x)
  fmap f (R y) = R (fmap f y)
```

```
instance (Functor f, Functor g) => Functor (f **: g) where
  fmap f (x **: y) = fmap f x **: fmap f y
```

B

Results

B.1 Individual benchmark results - Worst Case

B.1.1 Linear trend line

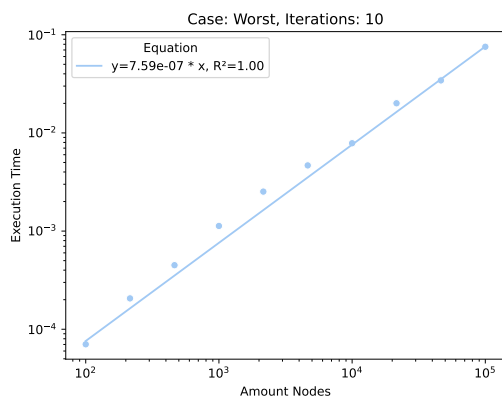


Figure B.1: Execution Time - Cata Sum

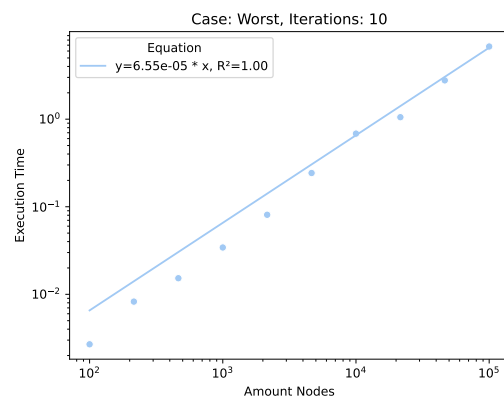


Figure B.2: Execution Time - Generic Cata Sum

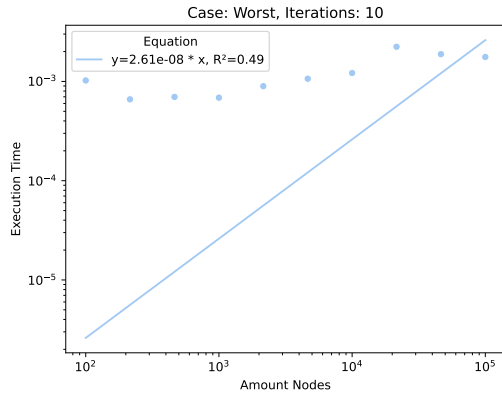


Figure B.3: Execution Time -
Incremental Cata Sum

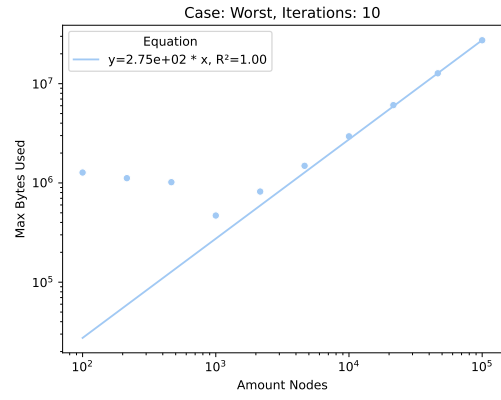


Figure B.4: Memory Usage - Cata
Sum

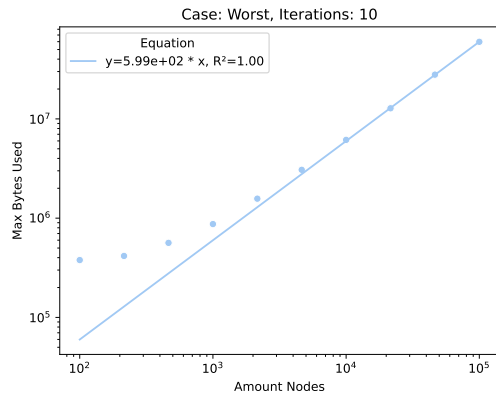


Figure B.5: Memory Usage -
Generic Cata Sum

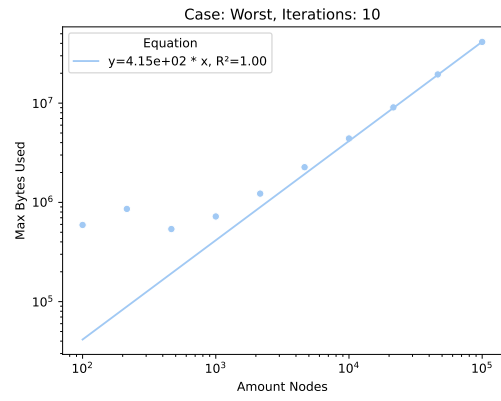


Figure B.6: Memory Usage -
Incremental Cata Sum

B.1.2 Logarithmic trend line

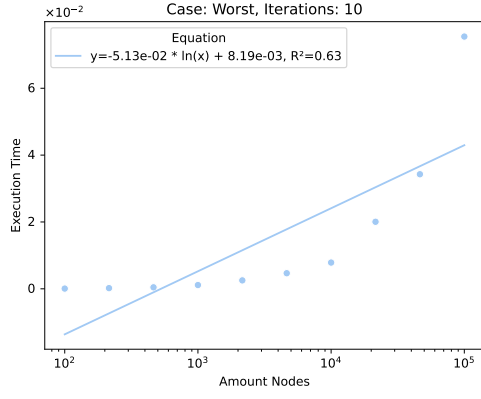


Figure B.7: Execution Time - Cata Sum

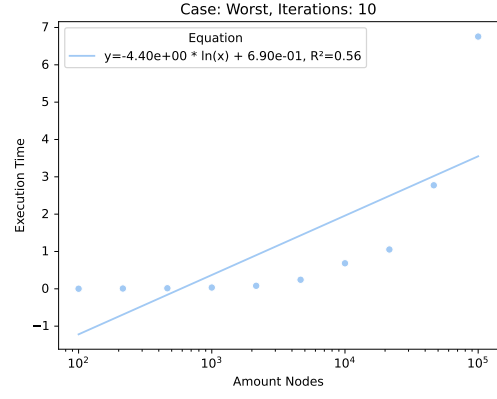


Figure B.8: Execution Time - Generic Cata Sum

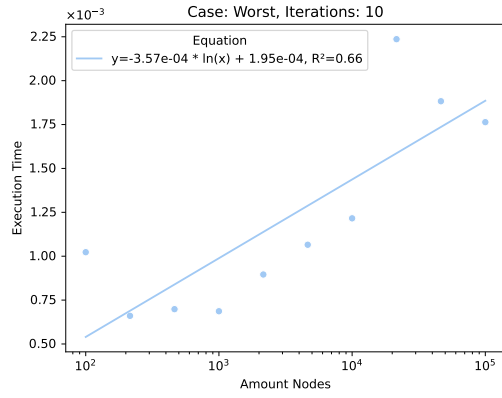


Figure B.9: Execution Time - Incremental Cata Sum

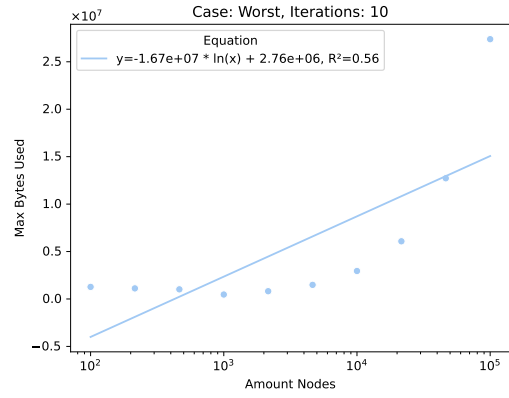


Figure B.10: Memory Usage - Cata Sum

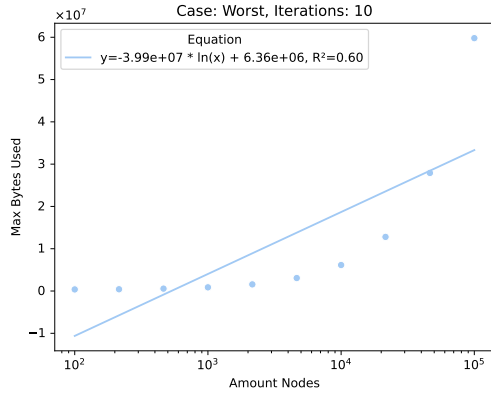


Figure B.11: Memory Usage -
Generic Cata Sum

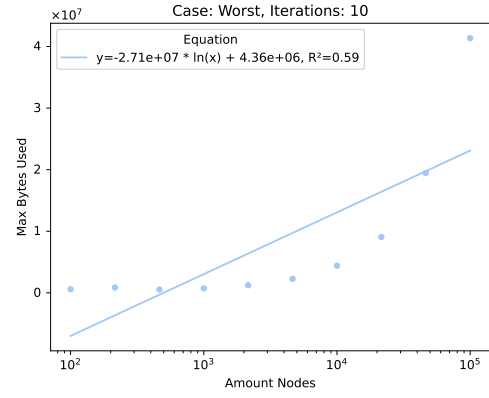


Figure B.12: Memory Usage -
Incremental Cata Sum

B.2 Individual benchmark results - Average Case

B.2.1 Linear trend line

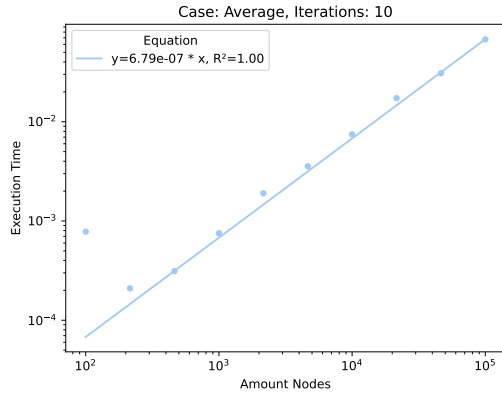


Figure B.13: Execution Time -
Cata Sum

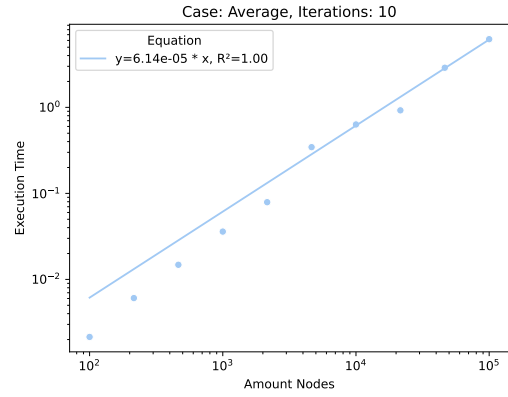


Figure B.14: Execution Time -
Generic Cata Sum

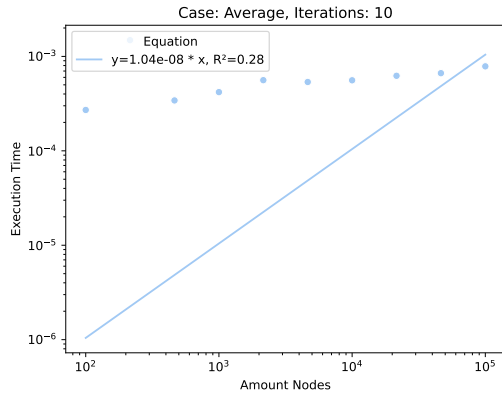


Figure B.15: Execution Time -
Incremental Cata Sum

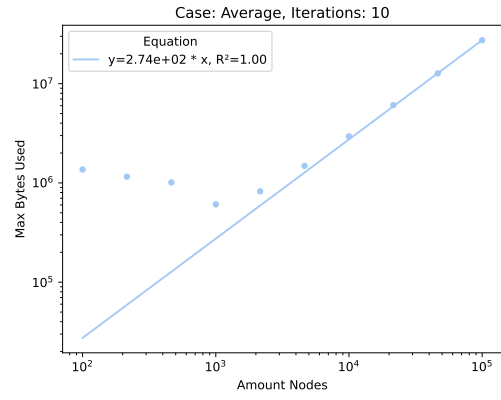


Figure B.16: Memory Usage - Cata
Sum

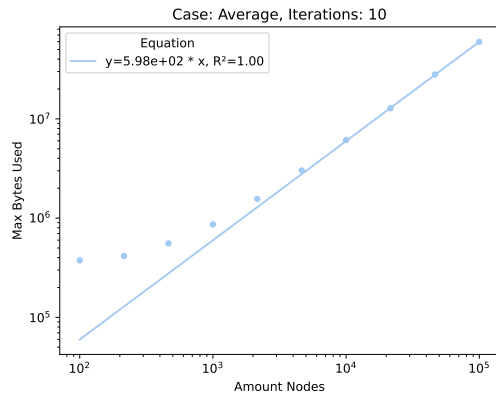


Figure B.17: Memory Usage -
Generic Cata Sum

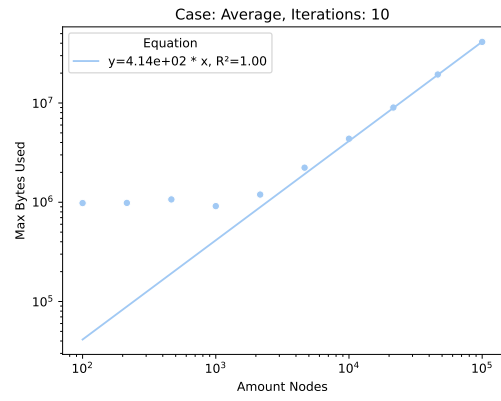


Figure B.18: Memory Usage -
Incremental Cata Sum

B.2.2 Logarithmic trend line

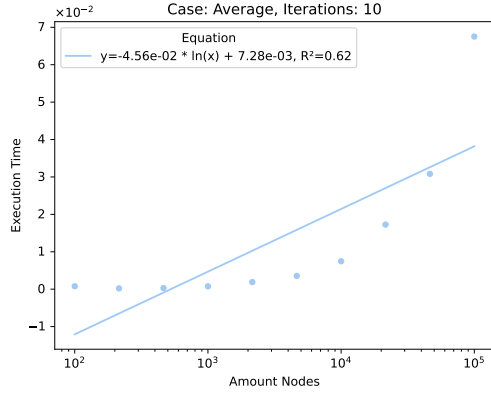


Figure B.19: Execution Time - Cata Sum

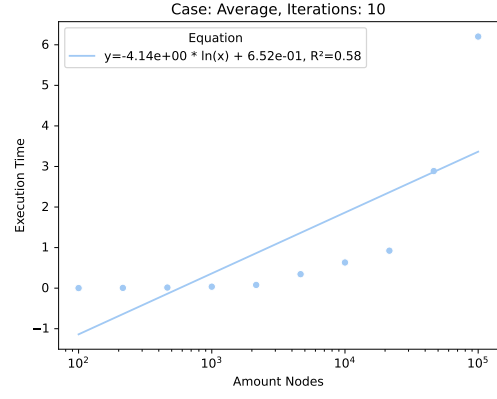


Figure B.20: Execution Time - Generic Cata Sum

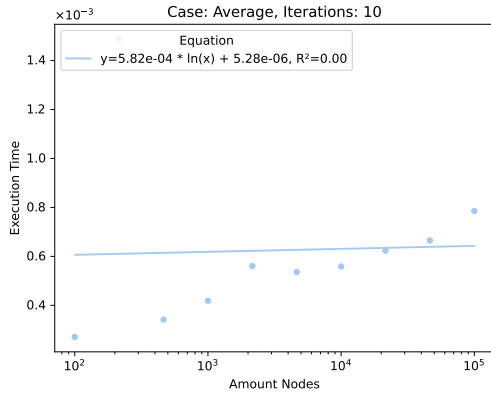


Figure B.21: Execution Time - Incremental Cata Sum

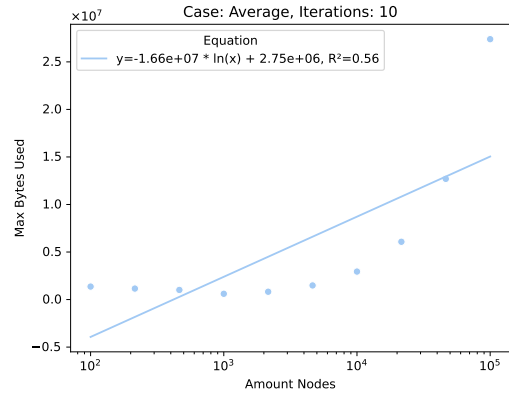


Figure B.22: Memory Usage - Cata Sum

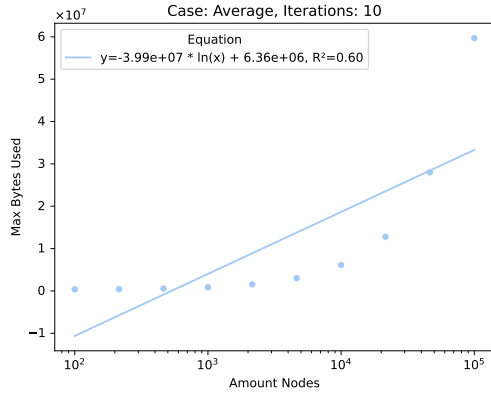


Figure B.23: Memory Usage -
Generic Cata Sum

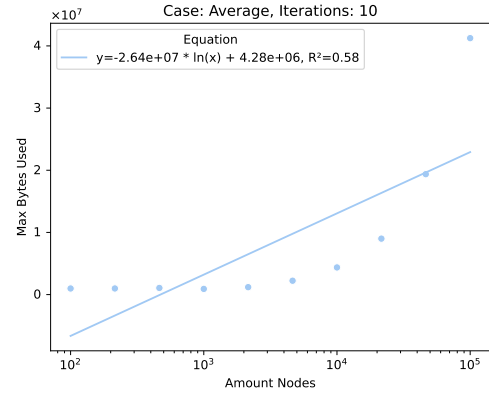


Figure B.24: Memory Usage -
Incremental Cata Sum

B.3 Individual benchmark results - Best Case

B.3.1 Linear trend line

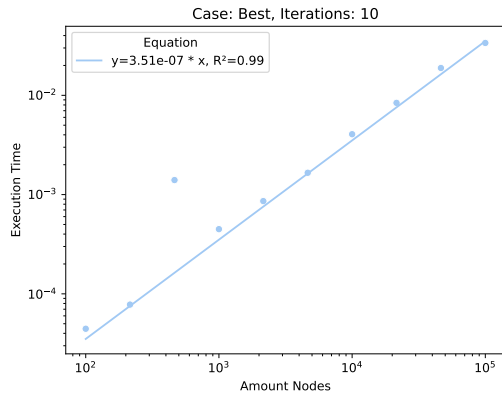


Figure B.25: Execution Time -
Cata Sum

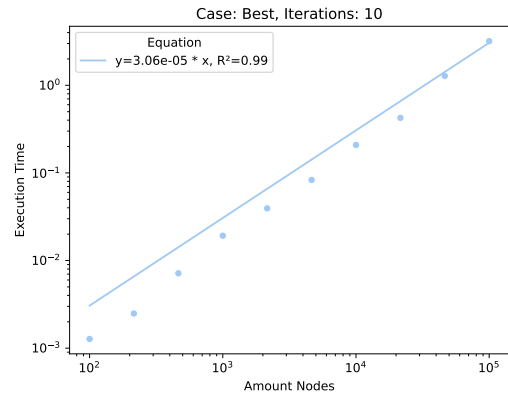


Figure B.26: Execution Time -
Generic Cata Sum

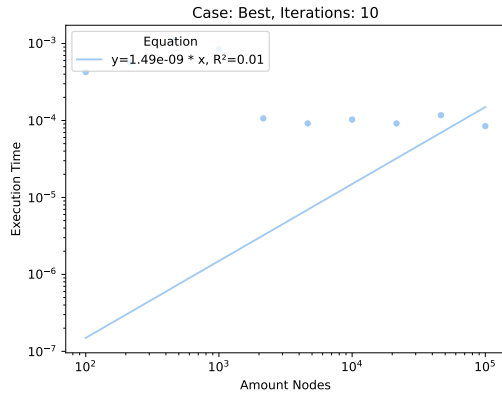


Figure B.27: Execution Time -
Incremental Cata Sum

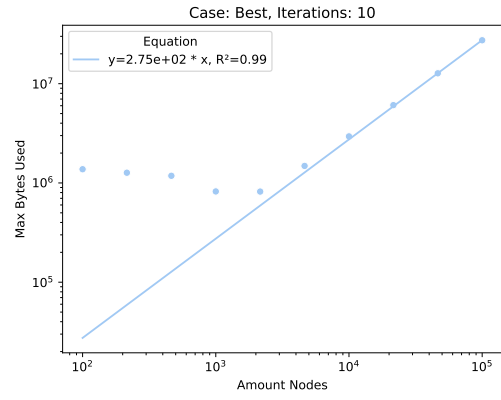


Figure B.28: Memory Usage - Cata
Sum

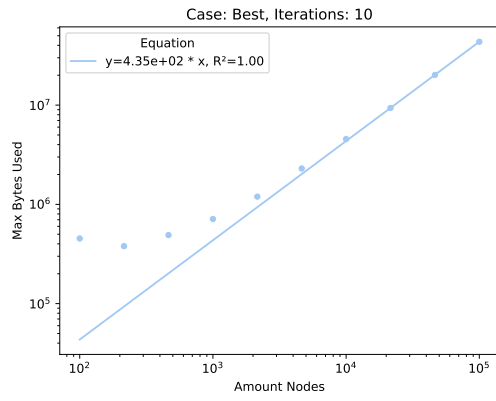


Figure B.29: Memory Usage -
Generic Cata Sum

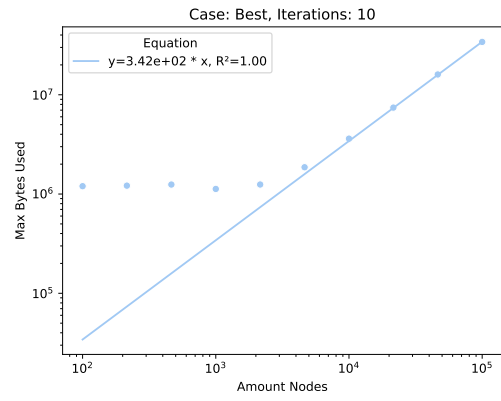


Figure B.30: Memory Usage -
Incremental Cata Sum

B.3.2 Logarithmic trend line

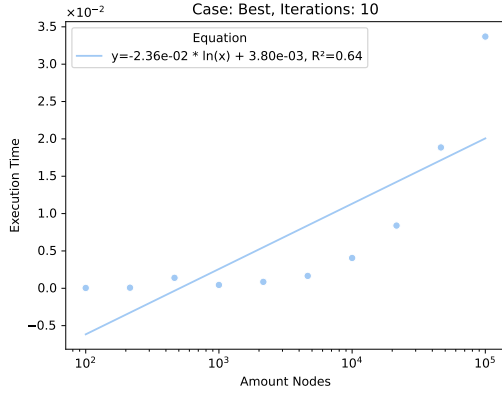


Figure B.31: Execution Time - Cata Sum

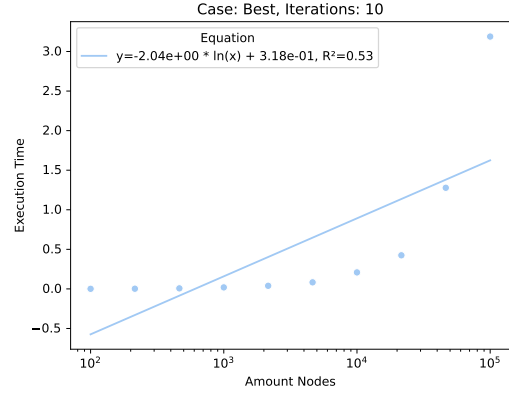


Figure B.32: Execution Time - Generic Cata Sum

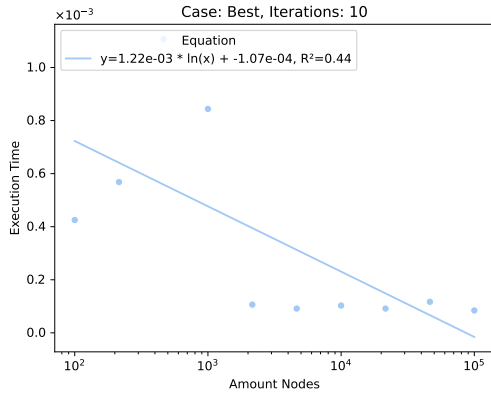


Figure B.33: Execution Time - Incremental Cata Sum

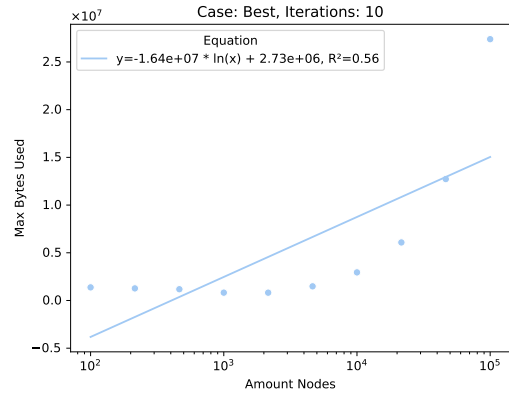


Figure B.34: Memory Usage - Cata Sum

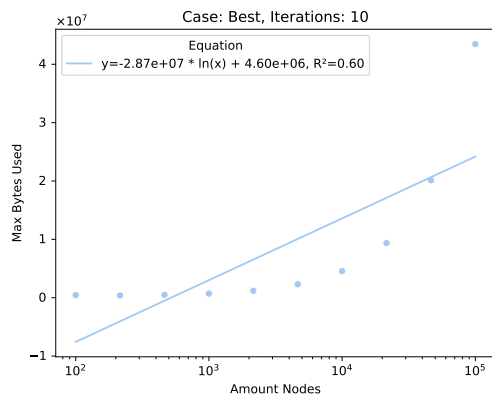


Figure B.35: Memory Usage -
Generic Cata Sum

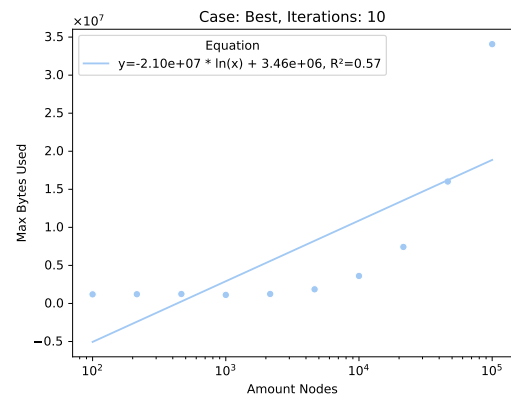


Figure B.36: Memory Usage -
Incremental Cata Sum

Bibliography

- [1] Umut A Acar, Guy E Blelloch, and Robert Harper. “Selective memoization”. In: *ACM SIGPLAN Notices* 38.1 (2003), pp. 14–25.
- [2] Jeroen Bransen and José Pedro Magalhaes. “Generic representations of tree transformations”. In: *Proceedings of the 9th ACM SIGPLAN workshop on Generic programming*. 2013, pp. 73–84.
- [3] Jeremy Gibbons. “Datatype-generic programming”. In: *International Spring School on Datatype-Generic Programming*. Springer. 2006, pp. 1–71.
- [4] HaskellWiki. *GHC/Type families*. 2021. URL: https://wiki.haskell.org/GHC/Type_families (visited on May 25, 2022).
- [5] Gérard Huet. “The zipper”. In: *Journal of functional programming* 7.5 (1997), pp. 549–554.
- [6] Oleg Kiselyov, Simon Peyton Jones, and Chung-chieh Shan. “Fun with type functions”. In: *Reflections on the Work of CAR Hoare*. Springer, 2010, pp. 301–331.
- [7] Jose Pedro Magalhaes. *Generic programming library for regular datatypes*. URL: <https://hackage.haskell.org/package/regular> (visited on Apr. 28, 2022).
- [8] Victor Cacciari Miraldo and Wouter Swierstra. “An efficient algorithm for type-safe structural diffing”. In: *Proceedings of the ACM on Programming Languages* 3.ICFP (2019), pp. 1–29.
- [9] Thomas van Noort et al. “A Lightweight Approach to Datatype-Generic Rewriting”. In: (2008).
- [10] Bryan O’Sullivan. *Robust, reliable performance measurement and analysis*. URL: <https://hackage.haskell.org/package/criterion> (visited on June 14, 2022).
- [11] William Wesley Peterson and Daniel T Brown. “Cyclic codes for error detection”. In: *Proceedings of the IRE* 49.1 (1961), pp. 228–235.
- [12] Matthew Pickering et al. “Pattern synonyms”. In: *Proceedings of the 9th International Symposium on Haskell*. 2016, pp. 80–91.
- [13] Geoff Pike, Jyrki Alakuijala, and Software Engineering Team. *Introducing CityHash*. Apr. 11, 2011. URL: <https://opensource.googleblog.com/2011/04/introducing-cityhash.html> (visited on June 14, 2022).
- [14] Jeff Preshing. *Hash Collision Probabilities*. 2011. URL: <https://preshing.com/20110504/hash-collision-probabilities> (visited on May 3, 2022).
- [15] Austin Seipp. *Bindings to CityHash*. URL: <https://hackage.haskell.org/package/cityhash> (visited on June 15, 2022).
- [16] GHC Team. *Running a compiled program*. URL: https://downloads.haskell.org/~ghc/6.12.1/docs/html/users_guide/runtime-control.html (visited on June 14, 2022).
- [17] Johan Tibell. *Efficient hashing-based container types*. URL: <https://hackage.haskell.org/package/unordered-containers-0.2.19.1/docs/Data-HashMap-Strict.html> (visited on June 15, 2022).

- [18] Edsko de Vries and Andres Löh. “True sums of products”. In: *Proceedings of the 10th ACM SIGPLAN workshop on Generic programming*. 2014, pp. 83–94.
- [19] Alexey Rodriguez Yakushev et al. “Generic programming with fixed points for mutually recursive datatypes”. In: *ACM Sigplan Notices* 44.9 (2009), pp. 233–244.