

On Factorization as Bijection to Finite Sequences

Paul Tarau

Department of Computer Science and Engineering
University of North Texas
tarau@cse.unt.edu

Abstract—We describe mechanisms to emulate multiplicative properties of *prime numbers* by connecting them to computationally simpler representations involving bijections from natural numbers to *sequences* and *multisets* of natural numbers. Through a generic framework, we model interesting properties of a canonical representation of factorization and its simpler equivalents. Among their shared properties, we study the information lost through factorization, we establish a connection to Catalan families.

Keywords—sequence and multiset encodings and prime factorization, information loss through factorization, experimental mathematics and functional programming, declarative modeling of computational phenomena, factorization and Catalan families.

I. INTRODUCTION

The study of the prime factorization of natural numbers has been a driving force towards various abstraction layers materialized along the years as in research and application fields ranging from formal logic and foundation of mathematics to algebraic and analytic number theory and cryptography. It is not unusual that significant progress on prime numbers correlates with unexpected paradigm shifts, the prototypical example being Riemann’s paper [1] connecting primality and complex analysis, all evolving around the still unsolved *Riemann Hypothesis* [2], [3]

Prime numbers exhibit a number of fundamental properties of natural phenomena and human artifacts in an unusually pure form. For instance, *reversibility* is present as the ability to recover the operands of a product of distinct primes. This relates to the information-theoretical view of multiplication [4] and it suggests investigating connections between combinatorial properties of multisets and operations on multisets and multiplicative number theory.

At the same time, the growing number of conjectures [5] and the large number of still unsolved problems involving prime numbers [6] suggest that modeling the underlying generic mechanisms that connect, like factorization, a given natural number to a finite multiset of natural numbers, is worth exploring.

This paper is a sequel of [7] and [8], with focus on emulating properties of prime factorization, *seen as a bijection between \mathbb{N} and its finite sequences and multisets*. Towards this end we use two infinite families of such bijections, based on n -adic valuations [8] and a new size-proportionate encoding using bijective base- k numbers.

After encapsulating their commonalities as a Haskell type class, we study in this generic framework some interesting

properties of factorization and compare them with its simpler emulations.

The paper is organized as follows.

Section II introduces general mechanisms for defining bijections between \mathbb{N} and its finite sequences and multisets that will be used to emulate properties of prime factorization. Subsection II-A overviews the encoding of natural numbers in bijective base- k that will be used in subsection II-B to define a size-proportionate bijection from \mathbb{N} to finite lists of natural numbers. Subsection II-C describes a bijection connecting sequences and multisets.

In section III we introduce a generic framework encapsulating factorization as a bijection between \mathbb{N} and its finite sequences.

Section IV describes a connection between recursive application of our canonical factorization operation and Catalan families.

Section V studies the information loss through factorization for various instances of our framework, including the actual primes.

Section VI discusses some limitations of our approach and suggest some open problems.

Section VII overviews some related work and section VIII concludes the paper.

The paper is written as a *literate Haskell program*. The code extracted from it is also available as a standalone file at <http://www.cse.unt.edu/~tarau/research/2014/GPrimes.hs>.

As a convenience to the reviewers, a number of appendices summarize content from [8], making the paper and its related code self-contained.

Appendix A overviews the encoding of finite multisets with primes using the primality testing and factoring algorithms of Appendix B.

Appendix C describes a family of bijections between \mathbb{N} and sequences of natural numbers using n -adic valuations.

Appendix D provides algorithms for multiset operations.

II. BIJECTIONS BETWEEN NATURAL NUMBERS AND SEQUENCES AND MULTISETS OF NATURAL NUMBERS

We start by introducing a new mechanism for defining a family of bijections between \mathbb{N} and its finite sequences and multisets. We will use the resulting decomposition of a number to a multiset as an analogue of prime factorization which will model some interesting properties of the primes themselves.

A. Encoding numbers in bijective base- k

Conventional binary representation of a natural number n describes it as the value of a polynomial $P(x) = \sum_{i=0}^k a_i x^i$

with digits $a_i \in \{0, 1\}$ for $x=2$, i.e. $n = P(2)$. If one rewrites $P(x)$ using Horner's scheme as $Q(x) = a_0 + x(a_1 + x(a_2 + \dots))$, n can be represented as an iterated application of a function $f(a, x, v) = a + xv$ which can be further specialized as a sequence of applications of two functions $f_0(v) = 2v$ and $f_1(v) = 2v + 1$, corresponding to the encoding of n as the sequence of binary digits a_0, a_1, \dots, a_k . It is easy to see that this encoding is not bijective as a function to $\{0, 1\}^*$. For instance 1, 10, 100 etc. would all correspond to $n = 1$. The problem is that conventional numbering systems do not provide a *bijection* between arbitrary combinations of digits and natural numbers, given that leading 0s are ignored. As a fix, the functions $o(v) = 1 + 2v$ and $i(v) = 2 + 2v$ can be used instead, resulting in the so called *bijjective base-2* representation: $1 = o(0), 2 = i(0), 3 = o(o(0)), 4 = i(o(0)), 5 = o(i(0))$ etc. Note that this representation is *canonical* i.e. a unique sequence of symbols identifies each natural number. This mechanism is easily generalized for an arbitrary base k . Bijjective base- k numbers are referred in [9] pp. 90-92 as “m-adic” numbers.

An encoder for *numbers in bijjective base- k* that provides such a bijection is implemented as the function `toBijBase` which, like the decoder `fromBijBase`, works one digit a time, using the functions `getBijDigit` and `putBijDigit`. They are both parameterized by the base b of the numeration system.

```
type N = Integer

fromBijBase :: N → [N] → N
toBijBase :: N → N → [N]
```

The functions `toBijBase` and `fromBijBase` convert to and from bijjective base b .

```
toBijBase _ 0 = []
toBijBase b n | n > 0 = d : ds where
  (d, m) = getBijDigit b n
  ds = toBijBase b m

fromBijBase _ [] = 0
fromBijBase b (d:ds) = n where
  m = fromBijBase b ds
  n = putBijDigit b d m
```

The function `getBijDigit` extracts one bijjective base- b digit from natural number n . It also returns the “left-over” information content as the second component of an ordered pair.

```
getBijDigit :: N → N → (N, N)
getBijDigit b n | n > 0 =
  if d == 0 then (b-1, q-1) else (d-1, q) where
    (q, d) = quotRem n b
```

The function `putBijDigit` integrates a bijjective base- b digit d into a natural number m :

```
putBijDigit :: N → N → N → N
putBijDigit b d m | 0 ≤ d && d < b = 1 + d * b + m
```

The encoding/decoding to bijjective base- b works as follows:

```
*GPrimes> toBijBase 4 2013
[0, 2, 0, 2, 2, 0]
*GPrimes> fromBijBase 4 it
```

```
2013
```

```
*GPrimes> map (toBijBase 2) [0..7]
[[[]], [0], [1], [0, 0], [1, 0], [0, 1], [1, 1], [0, 0, 0]]
```

This encoding will be used for define a family of bijections between natural numbers and finite lists, in section II-B.

B. A size-proportionate bijection from \mathbb{N} to finite lists of natural numbers

Definition 1: We call *size-proportionate* a bijection f from \mathbb{N} to finite lists of natural numbers if the size of the bijective base-2 representation of $n \in \mathbb{N}$ is within a constant factor of the sum of the sizes of the bijective base-2 representations of the members of the list $f(n)$.

We will define a size-proportionate bijection between \mathbb{N} and finite lists of elements of \mathbb{N} (denoted $[N]$ from now on) by converting a number n to bijective base $k + 1$ and then using digit k as a separator. The function `nat2nats` implements this operation.

```
nat2nats _ 0 = []
nat2nats _ 1 = [0]
nat2nats k n | n > 0 = ns where
  n' = pred n
  k' = succ k
  xs = toBijBase k' n'
  nss = splitWith k xs
  ns = map (fromBijBase k) nss

splitWith sep xs = y : f ys where
  (y, ys) = break (==sep) xs

f [] = []
f (_:zs) = splitWith sep zs
```

The function `nats2nat k` reverses `nat2nats k` by intercalating the separator k between the base k representation of the numbers on a list ns and then converts the result to bijective base $k + 1$.

```
nats2nat _ [] = 0
nats2nat _ [0] = 1
nats2nat k ns = n where
  nss = map (toBijBase k) ns
  xs = intercalate [k] nss
  n' = fromBijBase (succ k) xs
  n = succ n'
```

Proposition 1: The bijection from \mathbb{N} to $[N]$ defined by `nat2nats` is size-proportionate.

Proof: It follows from the fact that digit k in the bijective base- $k + 1$ is used as a separator turning the bijective base- k representation of $n \in \mathbb{N}$ into a list. ■

The following example illustrates this property, which is easy to detect visually by noticing similarly-sized decimal representations.

```
*GPrimes> nats2nat 7
[2, 0, 1, 3, 2, 0, 1, 4, 9, 777, 888, 0, 999]
27662734980522719186436359235
*GPrimes> nat2nats 7 it
[2, 0, 1, 3, 2, 0, 1, 4, 9, 777, 888, 0, 999]
```

C. A bijection between finite multisets and lists of natural numbers

Multisets [10] are collections with repeated elements.

Non-decreasing sequences provide a canonical representation for multisets of natural numbers. While finite multisets and finite lists of elements of \mathbb{N} share a common representation $[\mathbb{N}]$, multisets are subject to the implicit constraint that their ordering is immaterial. This suggests that a multiset like $[4, 4, 1, 3, 3, 3]$ could be represented *canonically* as a sequence by first ordering it as $[1, 3, 3, 3, 4, 4]$ and then computing the differences between consecutive elements i.e., $[x_0, x_1 \dots x_i, x_{i+1} \dots] \rightarrow [x_0, x_1 - x_0, \dots, x_{i+1} - x_i \dots]$. This gives $[1, 2, 0, 0, 1, 0]$, with the first element 1 followed by the increments $[2, 0, 0, 1, 0]$, as implemented by `mset2list`:

```
mset2list xs = zipWith (-) (xs) (0:xs)
```

It is now clear that incremental sums of the numbers in such a sequence return the original multiset as implemented by `list2mset`:

```
list2mset ns = tail (scanl (+) 0 ns)
```

Note that a canonical representation (i.e., being sorted) is assumed for multisets.

The bijection between finite multisets and finite sequences in \mathbb{N} defined by the functions `mset2list` and `list2mset`, is illustrated by the following example:

```
*GPrimes> mset2list [1,3,3,3,4,4]
[1,2,0,0,1,0]
*GPrimes> list2mset [1,2,0,0,1,0]
[1,3,3,3,4,4]
```

One can also think about this bijection as a one-to-one way to encode *canonically represented multisets* as sequences.

The following example illustrates how putting together `nat2nats k` and `list2mset` maps a natural number (123456789) with a multiset $[10, 10, 10, 25, 54]$, the same way as factoring maps a natural number into the multiset of exponents of its prime factors.

```
*GPrimes> nat2nats 3 123456789
[0,0,0,3,2,1595]
*GPrimes> nat2nats 3 12345678
[10,0,0,15,29]
*GPrimes> list2mset it
[10,10,10,25,54]
*GPrimes> mset2list it
[10,0,0,15,29]
*GPrimes> nats2nat 3 it
12345678
```

We will now build a general framework for hosting such multiset decompositions of natural numbers via bijections between natural numbers and lists, like `nat2nats k`.

III. A GENERIC FRAMEWORK FOR STUDYING PRIMALITY THROUGH $\mathbb{N} \rightarrow [\mathbb{N}]$ BIJECTIONS

To facilitate the comparison of properties derived from bijections between prime factors and alternative bijections between \mathbb{N} and finite multisets and sequences of elements of \mathbb{N} , we will use a Haskell type class for sharing their common structure, abstracting away their basic component: *a bijection between natural numbers and finite sequences of natural numbers*. In particular, given that sequences of differences between their consecutive terms represent non-decreasingly ordered multisets canonically, an interesting instance of this

mechanism will be used to represent prime factorization in subsection III-C.

A. The type class defining our framework

Parameterized by the bijections `mix: $[\mathbb{N}] \rightarrow \mathbb{N}$` and `unmix: $\mathbb{N} \rightarrow [\mathbb{N}]$` , as well as the bookkeeping operations (`lift` and `unlift`), we design a generic framework that expresses isomorphisms between operations on natural numbers and sequences of natural numbers. We implement the framework as the Haskell type class `Composer`, whose instances will provide a playground for exploring emulations of interesting properties of factorization.

```
class Composer m where
  lift :: N -> m N
  unlift :: m N -> N

  mix :: [m N] -> m N
  unmix :: m N -> [m N]
```

The polymorphic type parameter `m` will be instantiated to Haskell data types providing alternative implementations of the framework.

B. The instances derived from the bijection `nat2nats`

The instances `A` and `B` will be used to encapsulate the “prime emulations” derived from the functions `nats2nat` and `nat2nats`.

```
newtype A a = A a deriving (Show, Read, Eq, Ord)
```

```
instance Composer A where
  mix = mixWith (nats2nat 2)

  unmix = unmixWith (nat2nats 2)

  lift n = A n
  unlift (A n) = n
```

```
newtype B a = B a deriving (Show, Read, Eq, Ord)
```

```
instance Composer B where
  mix = mixWith (nats2nat 3)

  unmix = unmixWith (nat2nats 3)

  lift n = B n
  unlift (B n) = n
```

C. The instance derived for primes

The `Composer` instance `P` describes the actual primes in our framework. Along the lines of [8], it is based on the bijections between \mathbb{N} and its finite sequences `plist2nat` and `nat2plist`, that will use the functions `nat2pmset` and `pmset2nat` that provide the view of natural numbers as multisets of their prime factors (see Appendix A).

```
newtype P a = P a deriving (Show, Read, Eq, Ord)
```

```
instance Composer P where
  mix = mixWith plist2nat
  unmix = unmixWith nat2plist

  lift n = P n
  unlift (P n) = n
```

The following example illustrates the factorization of 360, represented by the list associated to its predecessor 359, that are derived from the positions of the factors of 360 in the infinite stream of primes.

```
*GPrimes> unmix (P 359)
[P 0,P 0,P 0,P 1,P 0,P 1]
*GPrimes> mix it
P 359
```

So P_0 corresponds to prime 2 (in position 0 in the stream of primes), the next two P_0 values indicate no increase to the next factors (also 2). Next, P_1 is the increment to the prime 3's position, which stay unchanged once as marked by P_0 . Finally P_1 corresponds to the next prime 5 appearing once in $360 = 2 * 2 * 2 * 3 * 3 * 5$.

D. The instances derived from the bijection nat2xnats

The instance X encapsulates the prime emulation based on the dyadic valuation $\nu_2(n)$, as defined in [8] (see Appendix C).

```
newtype X a = X a deriving (Show, Read, Eq, Ord)
instance Composer X where
  lift n = X n
  unlift (X n) = n

  mix = mixWith (xnats2nat 2)
  unmix = unmixWith (nat2xnats 2)
```

A similar instance, Y based on $\nu_3(n)$ is defined as follows:

```
newtype Y a = Y a deriving (Show, Read, Eq, Ord)
instance Composer Y where
  lift n = Y n
  unlift (Y n) = n

  mix = mixWith (xnats2nat 3)
  unmix = unmixWith (nat2xnats 3)
```

Clearly, an infinite number of similar instances can be derived for each k -adic valuation ν_k .

E. Some shared operations

First we define two generic map functions converting between lifted and unlifted lists of natural numbers.

```
unlifts :: Composer m => [m N] -> [N]
unlifts = map unlift

lifts :: Composer m => [N] -> [m N]
lifts = map lift
```

Next we define combinators that lift *mixing* and *unmixing* functions (transforming between \mathbb{N} and $[\mathbb{N}]$) to our data types.

```
mixWith :: Composer m => ([N] -> [N]) -> ([m N] -> [m N])
mixWith f = lift . f . unlifts

unmixWith :: Composer m => (N -> [N]) -> (m N -> [m N])
unmixWith f = lifts . f . unlift
```

The combinators `liftFun` and `liftFuns` apply operations on \mathbb{N} and $[\mathbb{N}]$ to their lifted counterparts.

```
liftFun :: Composer m => (N -> N) -> m N -> m N
liftFun f = lift.f.unlift

liftFuns :: Composer m => ([N] -> [N]) -> [m N] -> [m N]
liftFuns f = lifts.f.unlifts
```

After defining a generic list to multiset operation (note that a multiset decomposition is based on elements of \mathbb{N}^+ rather than \mathbb{N})

```
to_xmset :: Composer m => m N -> [m N]
to_xmset = liftFuns (map succ . list2mset) . unmix
```

```
from_xmset :: Composer m => [m N] -> m N
from_xmset = mix . liftFuns (mset2list . map pred)
```

we can provide a generic “multiset prime” *recognizer* (`is_xprime`) and a *generator* (`xprimes_from`). They use the combinator `liftFun` to apply the successor and predecessor operations `succ` and `pred` to shift natural numbers between \mathbb{N} and \mathbb{N}^+ .

The predicate `is_xprime` recognizes singletons as “primes” in our generic framework.

```
is_xprime :: Composer m => m N -> Bool
is_xprime x = f xs where
  xs = to_xmset (liftFun pred x)
  f [p] = True
  f _ = False
```

The function `xprimes_from` generates the infinite stream of “primes” starting from x .

```
xprimes_from :: Composer m => m N -> [m N]
xprimes_from x = filter is_xprime
  (iterate (liftFun succ) x)
```

We can also define a generic bijection from *positions* in the stream of “primes” to the stream of “primes” (`from_xindices`), as well as a simple “factorization” operation, `to_xfactors`.

```
from_xindices :: Composer m => [m N] -> m N
from_xindices = (liftFun succ) . from_xmset
```

```
to_xindices :: Composer m => m N -> [m N]
to_xindices = to_xmset . (liftFun pred)
```

```
to_xfactors :: Composer m => m N -> [m N]
to_xfactors x = map i2f (to_xindices x) where
  ps = xprimes_from (lift 1)
  i2f i = ps `genericIndex` (pred (unlift i))
```

F. Examples comparing our instances based on actual and emulated factorizations

The following examples illustrate the flexibility of the “plug-in” mechanism that this framework provides, when working with actual and emulated factorizations.

```
*GPrimes> take 10 (xprimes_from (P 1))
[P 2,P 3,P 5,P 7,P 11,P 13,P 17,P 19,P 23,P 29]
*GPrimes> take 10 (xprimes_from (B 1))
[B 2,B 3,B 4,B 5,B 7,B 8,B 9,B 11,B 12,B 13]
*GPrimes> take 10 (xprimes_from (Y 1))
[Y 2,Y 4,Y 10,Y 28,Y 82,
 Y 244,Y 730,Y 2188,Y 6562,Y 19684]
```

One can also observe, that the following holds:

Proposition 2: There’s an infinite number of *multiset primes* in each member of the family X, Y, \dots derived from valuation ν_b and they are exactly the numbers of the form $b^n + 1$.

The generic equivalent of *factorization* that matches its expected behavior on the instance \mathbb{P} representing the usual primes, is illustrated by the following examples:

```
*GPrimes> to_xfactors (P 36)
[P 2,P 2,P 3,P 3]
*GPrimes> to_xfactors (A 23)
[A 2,A 2,A 3]
*GPrimes> to_xfactors (X 32)
[X 2,X 2,X 2,X 2]
*GPrimes> to_xfactors (X 33)
[X 33]
```

In [7] and [8] we provide a detailed comparison in terms of two specific instances of multisets: those corresponding here to data types \mathbf{X} , \mathbf{Y} and the primes, corresponding here to data type \mathbf{P} .

IV. A CATALAN CONNECTION

The recursive application of the `unmix` and `mix` functions can be used to unfold a natural number to a tree containing successive layers of canonical decomposition as sequences. The ordered rooted trees of type `CTree` (with empty leaves `C []`) will be used to host the result of this decomposition.

```
data CTree = C [CTree] deriving (Show,Read,Eq)
```

We define the unfolding operation `toTree` and its folding counterpart `fromTree` on a rooted ordered tree representation, a typical member of the *Catalan family* [11] of combinatorial structures.

```
toTree :: Composer m => m N -> CTree
toTree mn = C (map toTree (unmix mn))
```

```
fromTree :: Composer m => CTree -> m N
fromTree (C xs) = mix (map fromTree xs)
```

The function `morphTree` describes generically a transformation between two such tree views, provided by two instances of the type class `Composer`.

```
morphTree :: (Composer m', Composer m) => m' N -> m N
morphTree = fromTree . toTree
```

The function `tsize` computes the size of a tree.

```
tsize :: Composer m => (N -> m N) -> N -> N
tsize t n = ts (toTree (t n)) where
  ts (C []) = 1
  ts (C xs) = succ (sum (map ts xs))
```

After adding the usual instances

we can define the functions p , a , b , x , y corresponding, respectively, to target types P , A , B , X , Y .

```
p n = morphTree n :: P N
a n = morphTree n :: A N
b n = morphTree n :: B N
x n = morphTree n :: X N
y n = morphTree n :: Y N
```

The following examples illustrate the use of these functions.

[illegible]

```
P 2014
*GPrimes> toTree (X 2014)
C [C [C []],C [],C [],C [],
  C [C []],C [],C [],C [],C []]
```

In combination with `unlift`, the functions `a`, `b`, `p`, `x` and `y` can be used to define (infinite) permutations of \mathbb{N} .

```
*GPrimes> map (unlift.p.A) [0..15]
[0,1,2,4,3,10,6,5,30,16,9,8,24,7,12,126]
*GPrimes> map (unlift.p.X) [0..15]
[0,1,2,3,4,5,8,7,6,9,14,11,24,17,26,15]
```

Intuitively, these permutations of \mathbb{N} amplify the “imperfection” of the emulation of, in this case prime factorization, by the other instances of the class `Composer`, as shown in figures 1.

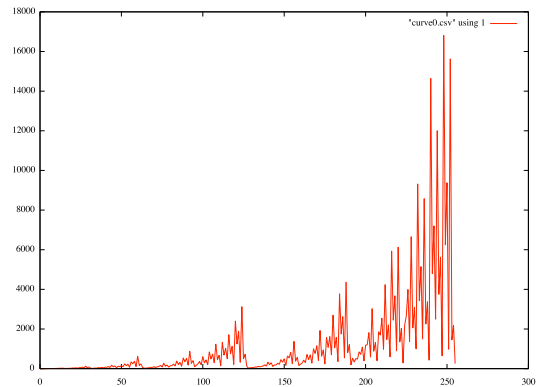


Fig. 1: Permutation of \mathbb{N} obtained by morphing between \mathbb{P} and \mathbb{X} views on $[0..2^8 - 1]$

V. MEASURING THE INFORMATION LOST THROUGH FACTORIZATION

In [4] a comprehensive information theoretic analysis of multiplication is presented, in terms of the information content of a number vs. its factors (and in particular its prime factors).

As prime factors contain together more bits than their product, multiplication results in general in *loss of information*. However, if one considers the *positions in the stream of primes* as a basis for canonical representations of factorization, we will see that this situation is *reversed*.

A. A multiplicative structure, generically

Multiplication corresponds to sorted concatenation of the multisets of factors, illustrated below:

```
*GPrimes> 12*15
180
*GPrimes> to_xfactors (P 12)
[P 2,P 2,P 3]
*GPrimes> to_xfactors (P 15)
[P 3,P 5]
*GPrimes> to_xfactors (P 180)
[P 2,P 2,P 3,P 3,P 5]
```

As shown in [8], one can also work with the the indices in the multiset of primes directly. We define multiplication along the lines of [8]. First we define a combinator `borrowBinOp` connecting to multiset operations (defined in Appendix D):


```
borrowBinOp f x y =
  from_xindices (f (to_xindices x) (to_xindices y))
```

Next we define generic multiplication, unit, gcd, lcm operations in term of their multiset counterparts.

```
multiply :: (Ord (m N), Composer m) =>
  (m N) -> (m N) -> (m N)
multiply = borrowBinOp msetSum
```

```
unit :: Composer m => m N
unit = lift 1
```

```
ggcd :: (Ord (m N), Composer m) => (m N) -> (m N) -> (m N)
ggcd = borrowBinOp msetInter
```

```
glcm :: (Ord (m N), Composer m) => (m N) -> (m N) -> (m N)
glcm = borrowBinOp msetUnion
```

One can observe that `multiply` works as expected on type `P` and it has similar basic properties on its alternatives.

```
*GPrimes> multiply (X 10) (X 5)
X 26
*GPrimes> multiply (X 5) (X 10)
X 26
*GPrimes> multiply (X 5) (X 1)
X 5
```

The following holds:

Proposition 3: The `multiply` operation defines a semigroup with unit `P 1` on data type `P` (corresponding to \mathbb{N}^+). The same structure is defined on data types `A`, `B`, `X` and `Y` and their generalizations for an arbitrary `k`.

Definition 2: We call *canonical factorization* the representation provided by our `unmix` bijection from \mathbb{N} to $[\mathbb{N}]$, for each of the instances of the type class `Composer`.

It follows from their bijection to arbitrary finite sequences in $[\mathbb{N}]$, that canonical factorization is information theoretically minimal.

In contrast to [4], where it is shown that information is lost by multiplication when considering the actual factors (and primes in particular), we will show that, when considering our *canonical factorization*, it is the other way around.

B. Computing the amount of information lost by factorization

After defining the representation size for bijective base-`b`

```
repsize b n = genericLength (toBijBase b n)
bitsize = repsize 2
```

we can proceed with the comparison between the information content corresponding to bijections between \mathbb{N} and $[\mathbb{N}]$, as well as their multiset counterparts.

```
infoLoss :: Composer m => (N -> m N) -> N -> N -> N
infoLoss t k n = f (unlifts.unmix.t) n where
  f u n = (repsize k n) - s where
    ns = u n
    s = sum (map (repsize k) ns)
```

The function `infoLoss` computes the difference between the representation sizes of a natural number and its expansion to a canonical list representation. In particular for $k = 2$ it computes the difference between their bitsizes in bijective base-2.

```
totalInfoLoss :: Composer m => (N -> m N) -> N -> N
totalInfoLoss t n = sum (map (infoLoss t 2) [0..2^n-1])
```

```
prefixSumLoss :: Composer m => (N -> m N) -> N -> [N]
prefixSumLoss t n =
  scanl1 (+) (map (infoLoss t 2) [0..2^n-1])
```

The functions `totalInfoLoss` and `prefixSumLoss` count respectively the total information loss and its prefix sums on the interval $[0..2^n - 1]$.

One can observe that that information loss is positive for each of them.

```
*GPrimes> map (infoLoss A 2) [0..15]
[0,1,0,1,2,0,0,2,1,1,2,2,2,3,0,1]
*GPrimes> map (infoLoss B 2) [0..15]
[0,1,0,1,0,2,0,1,1,2,0,0,0,2,0,1]
*GPrimes> map (infoLoss P 2) [0..15]
[0,1,0,2,1,1,0,3,2,2,1,2,1,1,1,4]
*GPrimes> map (infoLoss X 2) [0..15]
[0,1,0,2,1,1,1,3,1,2,1,2,2,2,2,4]
*GPrimes> map (infoLoss Y 2) [0..15]
[0,1,1,1,2,1,1,3,2,2,2,3,2,2,2,2]
```

Note however, that for larger bases (e.g., 10 and 20) that does not hold for every n :

```
*GPrimes> infoLoss P 10 104
-1
*GPrimes> infoLoss A 20 20
-1
*GPrimes> infoLoss X 10 10
-1
```

It is easy to prove that `infoLoss 2` is always positive in the case of the family of instances `A`, `B` ... as the list computed by `unmix` is obtained by removing a digit from their bijective base-2 representation.

It is also likely to be easy to prove that the same holds for instances `X`, `Y` ... where the left side of the bijection is exponentially larger than the right side.

The case for instance `P` is more intricate (but also more interesting), and one might first suspect that the the proof would involve deep properties of the primes.

We state this claim as:

Conjecture 1: The function `infoLoss 2` computing the difference of bitsizes of n and the sum of the bitsizes of the canonical factors of n , for instances `P`, `A`, `B`, ..., `X`, `Y`, ... is positive for all $n \in \mathbb{N}$.

After observing that `bitsize x` is the same as $b(x) = \lfloor \log_2(x+1) \rfloor$ we can restate this conjecture for instance `P` in terms of standard arithmetic notations (and prove it), as follows:

Proposition 4: For $n \in \mathbb{N}$, let $n+1 = p_0^{j_0} \dots p_i^{j_i}$, where p_0, \dots, p_i are the prime factors of $n+1$ in increasing order. Let $[q_0, \dots, q_m]$ be the list of the positions in the sequence of primes starting with 2 of these factors with their respective multiplicities. Let $[q_1 - q_0, \dots, q_m - q_{m-1}]$ be the list of their consecutive differences, clearly all ≥ 0 . Let $b(x) = \lfloor \log_2(x+1) \rfloor$.

Then $b(n) \geq \sum_{l=1}^m b(q_l - q_{l-1})$.

Proof: For $n = 0$ equality holds, the list on the right (factorization of 1) being empty. For $n > 0$ observe that $\sum_{l=1}^m b(q_l - q_{l-1}) = \sum_{l=1}^m \lfloor \log_2(1 + q_l - q_{l-1}) \rfloor \leq \lfloor \sum_{l=1}^m \log_2(1 + q_l - q_{l-1}) \rfloor \leq \lfloor \sum_{l=1}^m \log_2(q_l) \rfloor$

$$= \lfloor \log_2(\prod_{i=1}^m q_i) \rfloor \leq \lfloor \log_2(\prod_{i=0}^m p_i) \rfloor = \lfloor \log_2(n+1) \rfloor = b(n)$$

Observe that no “deep properties” of primes are involved, and the proof follows from a sequence of rather obvious inequalities.

Note that total information losses are comparable for all data types, P being closest to B.

```
*GPrimes> totalInfoLoss A 10
4379
*GPrimes> totalInfoLoss B 10
2987
*GPrimes> totalInfoLoss P 10
3965
*GPrimes> totalInfoLoss X 10
5447
*GPrimes> totalInfoLoss Y 10
5920
```

Figures 2, 3 and 4 illustrate the information lost by canonical factorization representations for instances P, A and X.

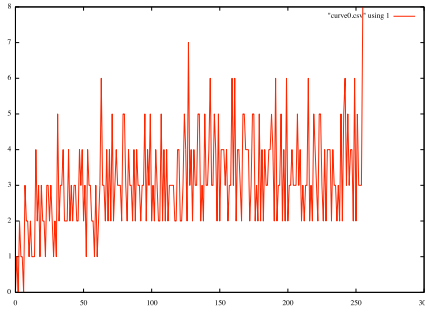


Fig. 2: Information loss for P on $[0..2^8 - 1]$

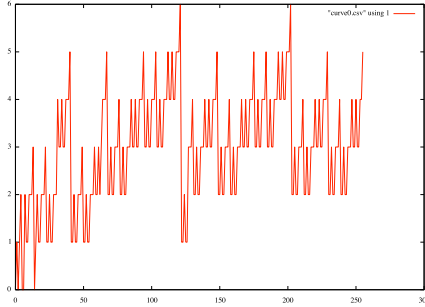


Fig. 3: Information loss for A on $[0..2^8 - 1]$

figure 5 compares prefix sums for the information lost by canonical factorization representations for type P and a instance matching it closely on a small initial segments of \mathbb{N} .

VI. DISCUSSION

Note that strictly speaking, the information content of a sequence would need to incorporate delimiting costs between its elements. In that sense one can see multiplication as implicitly adding a “separator” between factors, the same way as we explicitly do it with types A and B built with the size-proportionate encoding of section II. With implicit separation costs included, and given the reversibility of multiplication

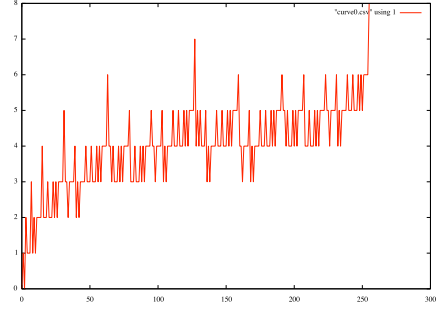


Fig. 4: Information loss for X on $[0..2^8 - 1]$

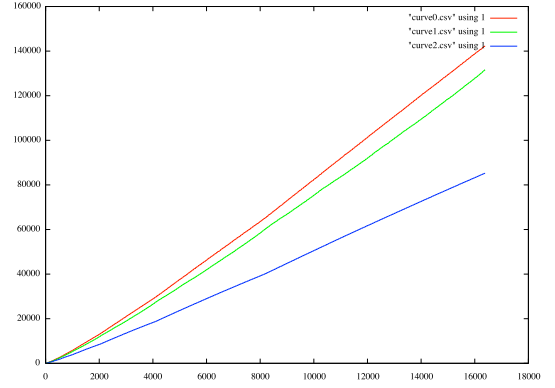


Fig. 5: Prefix sum of information loss for Y=red/upper, X=green/middle and P=blue/lower on $[1..2^{14}]$

through factoring, one can see the two sides as having exactly the same information content. It would be interesting in this sense to compare “efficiency” of a given encoding of a sequence seen as minimizing the separation costs in the corresponding aggregate.

The prevalent generalization of the unique decomposition of a natural number as a product of primes is provided by algebraic properties of commutative rings that have abstracted away the addition and multiplication operation. By contrast, our approach stays within elementary number theoretical concepts by using bijections between natural numbers and finite multisets of natural numbers. The resulting structure is a commutative semigroup with unit but additional properties of natural numbers are put at work that provide some interesting insight into the behavior of primes through the behavior of alternate instances sharing the same type class operations. For instance, all the resulting semigroups share the fact that they have an infinite number of generators. An important other property that our instances A, B, X, Y share with instance P corresponding to the actual primes, is that the bijection to multisets of natural numbers results in smaller numbers. Like in the case of primes, recursing over the sequence-equivalent of this decomposition maps these instances to Catalan families of combinatorial objects in a similar way.

An interesting open problem would be to discover a bijection similar to the ones the we have studied that would exhibit an additive structure. As it is possible to perform arithmetic operations directly with members of the Catalan family, as shown in [12] a venue would be to explore if the bijection to

Catalan families described in this paper (in section IV) could provide an answer to this problem.

VII. RELATED WORK

There's a huge amount of work on prime numbers and related aspects of multiplicative and additive number theory. Studies of the distribution of prime numbers and various probabilistic and information theoretic aspects also abound.

While we have not made use of any significantly advanced facts about prime numbers, the following references circumscribe the main topics to which our experiments can be connected [3], [6].

A recursive application of prime factorization similar to the one in section IV, used as a way to enumerate all trees, is given in [13]. In combinatorics, natural number encodings show up as *ranking* functions [14] that can be traced back to Gödel numberings [15], [16] associated to formulas. Together with their inverse *unranking* functions, they are also used in combinatorial generation algorithms for various data types [11].

In [7] and [8] a detailed comparison in terms of the specific instance of multisets corresponding here to data type X and primes, corresponding here to data type P , with focus on basic operations like product, gcd, lcm, etc., as well as in terms of more interesting concepts like the *rad* and *Mertens* functions.

Besides the generalized approach using Haskell's type classes, the most important new contributions of this paper cover the information loss through multiplication, the use of generalized n -adic valuation and size-proportionate bijections for multiset decompositions, their comparison with the decomposition provided by factoring and the correspondence between recursive application of these decomposition operations and Catalan families.

VIII. CONCLUSION

We have explored some computational analogies between multisets, natural number encodings and multiplicative properties of prime numbers in a framework for experimental mathematics implemented as a literate Haskell program.

We have lifted our Haskell implementations to a generic type class based model, along the lines of [17], which allows experimenting with instances parameterized by arbitrary bijections between \mathbb{N} and $[N]$.

A concept of canonical factorization, representing the decomposition of a number in a generic way has been introduced and its recursive application has provided a tree/DAG view uniquely associated to a natural number, through the corresponding member of the Catalan family of combinatorial structures.

Interesting correlations have been found between information losses due to factorization between primes and the other instances of our generic framework.

As an object of future work, of special interest in this direction are multiset decompositions of a natural number in $O(\log(\log(n)))$ factors, similar to the $\omega(x)$ and $\Omega(x)$ functions counting the distinct and non-distinct prime factors of x , to mimic more closely the distribution of primes.

REFERENCES

- [1] B. Riemann, "Ueber die Anzahl der Primzahlen unter einer gegebenen Grösse," *Monatsberichte der Berliner Akademie*, 1859.
- [2] G. L. Miller, "Riemann's Hypothesis and Tests for Primality," in *STOC*. ACM, 1975, pp. 234–239.
- [3] J. Derbyshire, *Prime Obsession: Bernhard Riemann and the Greatest Unsolved Problem in Mathematics*. New York: Penguin, 2004.
- [4] N. Pippenger, "The average amount of information lost in multiplication," *IEEE Transactions on Information Theory*, vol. 51, no. 2, pp. 684–687, 2005.
- [5] P. Cégielski, D. Richard, and M. Vsemirnov, "On the Additive Theory of Prime Numbers," *Fundam. Inform.*, vol. 81, no. 1-3, pp. 83–96, 2007.
- [6] R. Crandall and C. Pomerance, *Prime Numbers—a Computational Approach*, 2nd ed. New York: Springer, 2005.
- [7] P. Tarau, "Emulating Primality with Multiset representations of Natural Numbers," in *Proceedings of 8th International Colloquium on Theoretical Aspects of Computing, ICTAC 2011*, A. Cerone and P. Pihlajasaari, Eds. Johannesburg, South Africa: Springer, LNCS 6916, Sep. 2011, pp. 218–238.
- [8] —, "Towards a generic view of primality through multiset decompositions of natural numbers," *Theoretical Computer Science*, vol. 537, no. 0, pp. 105 – 124, 2014, theoretical Aspects of Computing (ICTAC 2011). [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0304397514003302>
- [9] A. Salomaa, *Formal Languages*. Academic Press, New York, 1973.
- [10] D. Singh, A. M. Ibrahim, T. Yohanna, and J. N. Singh, "An overview of the applications of multisets," *Novi Sad J. Math.*, vol. 52, no. 2, pp. 73–92, 2007.
- [11] D. L. Kreher and D. Stinson, *Combinatorial Algorithms: Generation, Enumeration, and Search*, ser. The CRC Press Series on Discrete Mathematics and its Applications. CRC PressINC, 1999.
- [12] P. Tarau, "Computing with Catalan Families," in *Proceedings of Language and Automata Theory and Applications, 8th International Conference, LATA 2014*, A.-H. Dediu, C. Martin-Vide, J.-L. Sierra, and B. Truthe, Eds. Madrid, Spain: Springer, LNCS, Mar. 2014, pp. 564–576.
- [13] D. W. Matula, "A natural rooted tree enumeration by prime factorization," *SIAM Review*, vol. 273, no. 10, pp. 267–323, 1968.
- [14] C. Martinez and X. Molinero, "Generic algorithms for the generation of combinatorial objects," in *MFCs*, ser. Lecture Notes in Computer Science, B. Rován and P. Vojtas, Eds., vol. 2747. Berlin Heidelberg: Springer, 2003, pp. 572–581.
- [15] K. Gödel, "Über formal unentscheidbare Sätze der Principia Mathematica und verwandter Systeme I," *Monatshefte für Mathematik und Physik*, vol. 38, pp. 173–198, 1931.
- [16] J. Hartmanis and T. P. Baker, "On Simple Goedel Numberings and Translations," in *ICALP*, ser. Lecture Notes in Computer Science, J. Loeckx, Ed., vol. 14. Berlin Heidelberg: Springer, 1974, pp. 301–316.
- [17] P. Tarau, "Declarative Modeling of Finite Mathematics," in *PPDP '10: Proceedings of the 12th international ACM SIGPLAN Symposium on Principles and Practice of Declarative Programming*. New York, NY, USA: ACM, 2010, pp. 131–142.

APPENDIX A ENCODING FINITE MULTISSETS WITH PRIMES

The mapping between finite multisets and primes described in this section goes back to Gödel's arithmetic encoding of formulas [15], [16].

A factorization of a natural number is uniquely described as a multiset of primes. Moreover, we can use the fact that each prime number is uniquely associated to its *position* in the infinite stream of primes, to obtain a bijection from multisets of natural numbers to natural numbers. It is provided by the function `to_prime_positions` defined in Appendix B. The function `nat2pmset` maps a natural number to the multiset of prime positions in its factorization.

```
nat2pmset 1 = []
nat2pmset n | n>1 = map succ (to_prime_positions n)
```

Clearly the following holds:

Proposition 5: p is prime if and only if its decomposition into a multiset given by `nat2pmset` is a singleton. The function `pmset2nat` (relying on `from_pos_in` and `primes` defined in Appendix B) maps back a multiset of positions of primes to the result of the product of the corresponding primes.

```
pmset2nat [] = 1
pmset2nat ns = from_prime_positions (map pred ns)
```

By using the bijection between lists and multisets we obtain:

```
nat2plist = mset2list . map pred . nat2pmset . succ
plist2nat = pred . pmset2nat . map succ . list2mset
```

The operations `nat2pmset` and `pmset2nat` form a bijection between \mathbb{N}^+ and $[\mathbb{N}^+]$ working as follows:

```
*GPrimes> nat2pmset 2014
[1,8,16]
*GPrimes> pmset2nat it
2014

*GPrimes> nat2plist 2014
[2,3,5]
*GPrimes> plist2nat it
2014
```

For instance, as the factorization of 2014 is $2 \cdot 19 \cdot 53$, the list `[1, 8, 16]` contains the *positions* of the factors, starting from position 1 of the factor 2, in the sequence of primes.

Note that the encoding of the prime factors of a number as a multiset of prime positions in the sequence of primes, rather than the multiset of the primes themselves, is significantly more compact, and its list encoding is even more so. In section V we will study in detail optimal information-theoretical encodings derived from factorization, as well as from alternative bijective representations.

APPENDIX B A SIMPLE ALGORITHM FOR PRIMALITY TESTING AND FACTORING

The following code implements a factorization function (`to_primes`), a primality test (`is_prime`) and a generator for the infinite stream of primes.

```
primes = 2 : filter is_prime [3,5..]
is_prime p = [p]==to_primes p

to_primes n | n>1 =
  to_factors n p ps where (p:ps) = primes

to_factors n p ps | p*p > n = [n]
to_factors n p ps | 0==n `mod` p =
  p : to_factors (n `div` p) p ps
to_factors n p ps@(hd:tl) = to_factors n hd tl

to_prime_positions n | n>1 =
  map (to_pos_in (h:ps)) qs where
    (h:ps)=genericTake n primes
    qs=to_factors n h ps
from_prime_positions ns = product
  (map (from_pos_in primes) ns)

to_pos_in xs x =
  fromIntegral i where Just i=elemIndex x xs
from_pos_in xs n = genericIndex xs n
```

APPENDIX C A BIJECTION BETWEEN \mathbb{N} AND $[\mathbb{N}]$ USING n -ADIC VALUATIONS

Along the lines of [7], [8] we will also describe a mechanism for deriving bijections between \mathbb{N} and $[\mathbb{N}]$ from one-solution Diophantine equations. Let's observe that

Proposition 6: $\forall z \in \mathbb{N} - \{0\}$ the exponential Diophantine equation

$$2^x(2y+1) = z \quad (1)$$

has exactly one solution $x, y \in \mathbb{N}$.

This follows immediately from the unicity of the decomposition of a natural number as a multiset of prime factors.

We will generalize this observation to obtain a *family of bijections* between \mathbb{N} and $[\mathbb{N}]$ by choosing an arbitrary base b instead of 2.

Definition 3: Given a number $n \in \mathbb{N}, n > 1$, the n -adic evaluation of a natural number m is the largest exponent k of n , such that n^k divides m . It is denoted $\nu_n(m)$.

Note that the solution x of the equation (1) is actually $\nu_2(z)$. This suggests deriving a similar diophantine equation for an arbitrary n -adic valuation.

We implement, for an arbitrary $b \in \mathbb{N}$ the functions `xcons` b and `xdecons` b , defined between $\mathbb{N} \times \mathbb{N}$ and \mathbb{N}^+ .

```
xcons :: N -> (N,N) -> N
xcons b (x,y') | b>1 = (b^x)*y where
  q=y' `div` (b-1)
  y=y'+q+1

xdecons :: N -> N -> (N,N)
xdecons b z | b>1 && z>0 = (x,y') where
  hd n = if n `mod` b > 0 then 0 else 1+hd (n `div` b)
  x = hd z
  y = z `div` (b^x)
  q = y `div` b
  y' = y-q-1
```

Using `xdecons` we define the head and tail projection functions `xhd` and `xtl`:

```
xhd, xtl :: N -> N -> N
xhd b = fst . xdecons b
xtl b = snd . xdecons b
```

Intuitively, their correctness follows from the fact that the quotient y , after dividing z with the largest possible power b^x is not divisible by b . We then “rebase” it to base $b - 1$, and make z correspond to a unique pair of natural numbers.

More precisely, non-divisibility of y by b , stated as $y = bq + m$, $b > m > 0$ can be rewritten as $y - q - 1 = bq - q + m - 1$, $b > m > 0$, or equivalently $y - q - 1 = (b - 1)q + (m - 1)$, $b > m > 0$ from where it follows that by setting $y' = y - q - 1$ and $m' = m - 1$, we map z to a pair (y', m') such that $y' = (b - 1)q + m'$ and $b - 1 > m' \geq 0$. So we can transform (y, m) such that $y = bq + m$ with $b > m > 0$, into a pair (y', m') , such that $y' = q(b - 1) + m'$ with $b - 1 > m' \geq 0$. Note that the transformation works also in the opposite direction with $y' = y - q - 1$ giving $y = y' + q + 1$, and with $m' = m - 1$ giving $m = m' + 1$.

The following examples illustrate the operations for base 3:

```
*GPrimes> xcons 3 (10,20)
1830519
*GPrimes> xhd 3 1830519
10
*GPrimes> xtl 3 1830519
20
```

Note that $\text{xhd } n \ x$ is the n -adic valuation of x , $\nu_n(x)$

Definition 4: We call $\text{xhd } n \ x$ the n -adic head of $x \in \mathbb{N}^+$, $\text{xtl } n \ x$ the n -adic tail of $x \in \mathbb{N}^+$ and $\text{xcons } n \ (x, y)$ the n -adic cons of $x, y \in \mathbb{N}$.

For each base $b > 1$ we obtain a pair of functions between natural numbers and lists of natural numbers in terms of xhd , xtl and xcons as follows:

```
nat2xnats :: N -> N -> [N]
nat2xnats _ 0 = []
nat2xnats k n | n > 0 = xhd k n : nat2xnats k (xtl k n)

xnats2nat :: N -> [N] -> N
xnats2nat _ [] = 0
xnats2nat k (x:xs) = xcons k (x, xnats2nat k xs)
```

Proposition 7: The function nat2xnats defines a bijection from \mathbb{N} to $[\mathbb{N}]$ and xnats2nat is its inverse.

Proof: The follows from the reversibility of the cons and decons functions. ■

The following example illustrates how they work:

```
*GPrimes> nat2xnats 3 2014
[0,0,1,2,0,2]
*GPrimes> xnats2nat 3 it
2014

*GPrimes> xnats2nat 2 [0,10,100]
5192296858534827628530496329222145
*GPrimes> nat2xnats 2 it
[0,10,100]
```

As the second example illustrates, this bijection is not size-proportionate as values on the $[\mathbb{N}]$ side result in exponentially larger representations on the \mathbb{N} side.

APPENDIX D MULTISET OPERATIONS

```
msetInter :: Ord a => [a] -> [a] -> [a]
msetDif :: Ord a => [a] -> [a] -> [a]
msetSum :: Ord a => [a] -> [a] -> [a]
```

```
msetSymDif :: Ord a => [a] -> [a] -> [a]
msetUnion :: Ord a => [a] -> [a] -> [a]
```

```
msetInter [] _ = []
msetInter _ [] = []
msetInter (x:xs) (y:ys) | x == y = x:msetInter xs ys
msetInter (x:xs) (y:ys) | x < y = msetInter xs (y:ys)
msetInter (x:xs) (y:ys) | x > y = msetInter (x:xs) ys
```

```
msetDif [] _ = []
msetDif xs [] = xs
msetDif (x:xs) (y:ys) | x == y = msetDif xs ys
msetDif (x:xs) (y:ys) | x < y = x:msetDif xs (y:ys)
msetDif (x:xs) (y:ys) | x > y = msetDif (x:xs) ys
```

```
msetSum xs ys = sort (xs ++ ys)
```

```
msetSymDif xs ys =
  msetSum (msetDif xs ys) (msetDif ys xs)
```

```
msetUnion xs ys =
  msetSum (msetSymDif xs ys) (msetInter xs ys)
```