

Lazy wheel sieves and spirals of primes

COLIN RUNCIMAN

Department of Computer Science, University of York,

Heslington, York, YO1 5DD, UK

e-mail: colin@cs.york.ac.uk

1 Summary

The popular method of enumerating the primes is the *Sieve of Eratosthenes*. It can be programmed very neatly in a lazy functional language, but runs rather slowly. A little-known alternative method is the *Wheel Sieve*, originally formulated as a fast imperative algorithm for obtaining all primes up to a given limit, assuming destructive access to a bit-array. This article describes functional variants of the wheel sieve that enumerate all primes as a lazy list.

2 A standard solution

Few readers of this journal will be unfamiliar with the following program to enumerate the primes using *The Sieve of Eratosthenes*.

```
primes = sieve [2..]
sieve (p:xs) = p : sieve [x | x <- xs, x mod p > 0]
```

This little program, or something very like it, has been part of the stock-in-trade of lazy list-processing for over twenty years (Turner, 1975). For all its dainty appearance, however, the program makes brutal demands on the reduction machine. Suppose we compute p_k , the k^{th} prime. Sample reduction counts[†] are:

p_{100} : 36,558 p_{1000} : 3,124,142 p_{5000} : 75,951,324

Why is the computation so costly? There are two main reasons for inefficiency in *generate-and-test* programs: the generator may be too simple, producing comparatively few successful candidates among many failures; the test may be too complex, performing unnecessarily extensive checks. Both reasons apply here. The generator [2..] offers *all* integers > 1 as candidate primes. The multi-stage test makes p_{k+1} an item in p_k *intermediate lists* between the initial list of candidates and the final list of primes — `sieve`'s argument at recursion depth d being the list of numbers $> p_d$ not divisible by any of $p_1 \cdots p_d$.

[†] Reduction counts in this article are those reported by Mark Jones' HUGS interpreter, computing p_k as the value of `primes!!(k-1)` for some definition of `primes`.

```

primes = sieve [2..]

sieve (p:xs) = p : [x | x <- xs, noFactorIn primes squares x]

noFactorIn (p:ps) (q:qs) x =
  q > x || x mod p > 0 && noFactorIn ps qs x

squares = [p*p | p <- primes]

```

Fig. 1. A circular program with fewer intermediate lists and fewer tests for divisors than the standard solution using a recursive sieve.

3 A circular program with a cheaper test

The technique of *circular programming* (Bird, 1984) can be used to avoid the intermediate lists of the recursive sieve. By reformulating `sieve` so that it refers directly to the list of `primes` it is producing, it need no longer be recursive. Instead of checking for a single prime factor in the list comprehension of each recursive call, we check for all possible prime factors in one qualifier.

To implement this scheme, we need some way of cutting short the search for factors in `primes`. Otherwise the computation will fall into an unproductive *black hole* of self-reference. We use an elementary fact: if x is composite, it must have a prime factor p with $p \leq \sqrt{x}$; so if we reach a prime larger than \sqrt{x} without passing a factor of x , we may conclude that x is prime. The revised program is shown in Figure 1. The addition of the test `q > x` not only avoids the black hole; it also avoids many of the tests for divisors carried out by the recursive sieve. This double gain is reflected in reduction counts less than 35% of the original even when computing only as far as p_{100} , and less than 4% when computing p_{5000} .

$$p_{100}: 12,395 (< 35\%) \quad p_{1000}: 267,152 (< 10\%) \quad p_{5000}: 2,325,932 (< 4\%)$$

4 Wheels: a series of generators

Imagine a wheel of unit circumference, with a spike at one point on its rim. Roll this wheel along a tape. After n revolutions, there are n regularly-spaced holes through the tape, one unit apart. This wheel is equivalent to the generator we have used so far. Numbering the first hole 2, it generates 2, 3, 4, 5, ...

But this wheel is only the smallest, W_0 , in an infinite series of wheels, W_k for $k = 0, 1, 2, 3, \dots$. Following (Pritchard, 1982), let $\Pi_k = p_1.p_2.\dots.p_k$ the product of the first k primes. Then W_k is a wheel of circumference Π_k , with spikes positioned at exactly those points x units round the circumference where $x \bmod p_n > 0$ for $n = 1 \dots k$. Because $\Pi_k \bmod p_j = 0$ for all $j \leq k$, no matter how far W_k is rolled, beyond p_k the numbers spiked are exactly those without prime divisors $\leq p_k$.

5 Computing wheels from primes ...

A wheel can be represented by a construction including its circumference and a list of spike positions.

```
data Wheel = Wheel Int [Int]
```

For example, W_0 is represented by `Wheel 1 [1]`. The full infinite series of wheels can be defined as in Figure 2. W_{k+1} is generated by rolling W_k around a rim of circumference Π_{k+1} , excluding multiples of p_{k+1} from the spike positions obtained.

```
wheels =
  Wheel 1 [1] :
  zipWith nextSize wheels primes

nextSize (Wheel s ns) p =
  Wheel (s*p) [n' | o <- [0,s..(p-1)*s],
                n <- ns,
                n' <- [n+o], n' mod p > 0]
```

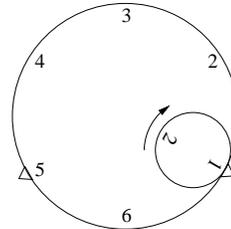


Fig. 2. Defining an infinite series of wheels of increasing size. The diagram shows how W_1 is used to generate W_2 , illustrating the application `nextSize (Wheel 2 [1]) 3`.

... and primes from wheels

Our `primes` programs so far have only used W_0 as generator. What options are there for the use of larger wheels?

Fixed wheel solutions: We could simply use W_k for some fixed k . But this has two disadvantages. First, we have to make some special arrangement to compute the first k primes (or else write them into the program explicitly) as they are needed to build W_k . Secondly, whatever k we choose, it will be a poor choice in some circumstances: if few primes are needed, a small wheel is best; if many are needed, a large one is best.

Variable wheel solutions: The alternative is a circular program that computes both `wheels` and `primes` according to demand. The larger the primes to be computed, the larger the wheel used. But this raises a key question. *When exactly should the wheel be changed?* At least two different answers can be given:

1. Change when it is most convenient to do so. The program will be simple, and the change-over will be cheap. Specifically, roll W_k exactly $(p_{k+1} - 1)$ times before changing to W_{k+1} . Then the change occurs at exactly the point where W_{k+1} would have started its second revolution had it already been in use.
2. Change just often enough to avoid *all* residual sieving. This maximises the benefits of using a series of wheels, but at the cost of a more complex program. Specifically, roll W_k only until it reaches p_{k+1}^2 , then change to W_{k+1} which must take up the computation part way through a revolution.

We shall consider each alternative, in order.

6 Lazy wheel sieve Mark I

First, a solution changing wheels at a point that allows the new wheel to start a complete revolution. Figure 3 gives the new definitions of `primes` and `sieve`.

```
primes = sieve wheels primes squares

sieve (Wheel s ns : ws) ps qs =
  [ n' | o <- s : [2*s,3*s..(head ps-1)*s],
    n <- ns,
    n' <- [n+o | s <= 2 || noFactorIn ps qs n'] ]
  ++ sieve ws (tail ps) (tail qs)
```

Fig. 3. The Mark I wheel-sieve.

The comprehension in the `sieve` computation is similar to that in `nextSize`: both roll a given wheel a number of times determined by a given prime. But circularity makes the `sieve` definition more delicate, and in several places we must tread carefully to avoid a black hole. Rather than pattern-match eagerly against the sequence of primes `ps` presented as `sieve`'s second argument, `head` and `tail` are applied when needed. We further avoid dependence on `ps` by separating the first of the offsets `o` from a range computation for the rest that is strict in its upper limit, and by searching for factors in `ps` only when `s > 2`.

The performance of the Mark I Wheel Sieve compares favourably with the circular sieve (which uses only the equivalent of W_0). Sample reduction counts are:

p_{100} : 4,920 ($\sim 40\%$) p_{1000} : 131,713 ($\sim 50\%$) p_{5000} : 1,445,789 ($\sim 60\%$)

7 Lazy wheel sieve Mark II

In the alternative wheel-changing solution, each wheel is changed just in time to avoid residual sieving. This method is more awkward to program because the *old wheel must stop rolling in mid revolution*. List comprehensions no longer provide a convenient way to express wheel-rolling. Instead of a comprehension such as

```
roll ns s p =
  [n' | o <- 0,s,..(p-1)*s,
    n <- ns,
    n' <- [n+o], c n']
```

we resort to explicit recursion and the use of `foldr`:

```
roll ns s p = roll' (p-1) 0
  where
    roll' 0 _ = []
    roll' t o = foldr (turn o) (roll' (t-1) (o+s)) ns
    turn o n xs =
      let n' = o+n in
      if c n' then n':xs else xs
```

To express a wheel change in mid-revolution, we can redefine the `turn` auxiliary so that it conditionally substitutes for `xs` a recursive call using a different wheel.

Another tricky aspect of the Mark II wheel-sieve is that the *new wheel must start rolling in mid revolution*. The simple list-of-spikes representation would force an ungainly and expensive `dropWhile (< p*p)`. This motivates a revised representation of wheels: split the spikes of W_k into *two* lists, the first containing all spikes $< p_k^2$ and the second those $> p_k^2$.

```
data Wheel = Wheel Int [Int] [Int]
```

A change of representation for wheels entails a new definition of `nextSize`. Because of the split spike-list we use `foldr` in place of a comprehension. As there is no longer any residual sieving, the `sieve` function is renamed `spiral`: it too uses the `foldr` technique, and the major recursive call representing a wheel-change is made within the `turn` auxiliary. Figure 4 shows a full program for the Mark II sieve.

The application of `dropWhile` to the recursive `spiral` computation may be surprising. Didn't we change the representation of wheels to make this unnecessary? Indeed, the revised representation makes it easy to skip values $< p_k^2$ from the *first revolution* of a new wheel W_k ; but for very small values of k the circumference Π_k is less than p_k^2 , so it is necessary to skip values *beyond* the first revolution. After

```
wheels =
  Wheel 1 [1] [] :
  zipWith3 nextSize wheels primes squares

nextSize (Wheel s ms ns) p q =
  Wheel (s*p) ms' ns'
  where
    (xs,ns') = span (<=q) (foldr (turn o) (roll (p-1) s) ns)
    ms' = foldr (turn 0) xs ms
    roll 0 _ = []
    roll t o =
      foldr (turn o) (foldr (turn o) (roll (t-1) (o+s)) ns) ms
    turn o n rs =
      let n' = o+n in [n' | n' mod p > 0] ++ rs

primes = spiral wheels primes squares

spiral (Wheel s ms ns : ws) ps qs =
  foldr (turn 0) (roll s) ns
  where
    roll o =
      foldr (turn o) (foldr (turn o) (roll (o+s)) ns) ms
    turn o n rs =
      let n' = o+n in
        if n'==2 || n' < head qs then n':rs
        else dropWhile (<n') (spiral ws (tail ps) (tail qs))
```

Fig. 4. The Mark II wheel sieve.

three wheel changes, `dropWhile` acts as an identity but for the small overhead of one comparison at each subsequent change.

Measured by reduction counts, Mark II out-paces Mark I, but the margin between the two is very slight until several hundred `primes` have been computed. Even by p_{1000} the margin is only 12%.

$$p_{100}: 4,841 (\sim 98\%) \quad p_{1000}: 116,646 (\sim 88\%) \quad p_{5000}: 1,245,756 (\sim 86\%)$$

8 Giant wheels and lazy spirals

The similar performance of the two wheel-sieve variants is more striking when one considers the sizes of wheels involved. Suppose we evaluate `primes` as far as p_{5000} (= 48,611). The largest wheel used by the Mark I wheel-sieve is W_6 since

$$\Pi_6 = 30,030 < p_{5000} < 510,510 = \Pi_7.$$

Drawn to the scale of Figure 2, W_6 would comfortably encircle a large building – quite a big wheel, though its representation does not involve numbers beyond the scope of single-word machine arithmetic. The Mark II wheel-sieve, however, reaches W_{47} since

$$p_{47} = 211 < \sqrt{p_{5000}} < 223 = p_{48}.$$

Now Π_{47} , the circumference of W_{47} , is an 86-digit number. On the same scale as before, the dimensions of W_{47} far exceed those of the visible universe! We must not construct any more of such a huge wheel than is strictly necessary. Lazy evaluation is essential. Continuing with the illustration of computing `primes` up to p_{5000} , of the three components in the representation of W_{47} : the circumference Π_{47} is *never evaluated*; the spikes less than $p_{47}^2 = 44,521$ are *never evaluated*; only a *fragment* of the rim of W_{47} beyond p_{47}^2 is constructed, containing just 373 spikes.

Figure 5 shows a way of tracing the enumeration of primes by a wheel-sieve program. Drawing the wheels concentrically, the computation can be traced as a spiral. Beginning on the rim of W_0 , the spiral orbits the centre once for each completed revolution of a wheel, and its increasing radius is equal to that of W_k at the point where the program switches to W_k as a new wheel. Each generated prime is marked at the appropriate point on the spiral, positioned on a radius passing through the spike that generated it.

9 Final remarks

There is some scope for improving the wheel-sieve programs given here. For example, during the construction of W_{k+1} using W_k , each candidate spike is tested for non-divisibility by p_{k+1} . But the failing candidates are exactly the products of p_{k+1} and a spike in W_k — a fact exploited in the array-based algorithm of (Pritchard, 1982). (That algorithm also changes wheels according the principle used for our Mark II program, so it has hardly a vestige of sieving, despite its name.)

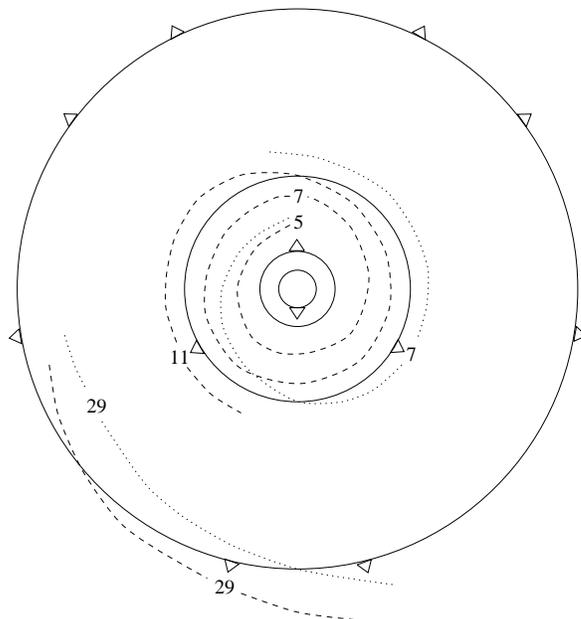


Fig. 5. Fragments early in the spirals of primes resulting from the two alternative wheel-changing rules: the dotted spiral traces Mark I, and the dashed spiral Mark II. (W_3 is shown *half-size* to resolve a conflict of scale.)

Reduction counts are handy and machine-independent. But they give only a rough indication of comparative costs. Among the factors they ignore, for example, the Mark II program uses more memory than Mark I.

One could of course write wheel-sieve variants for the indefinite enumeration of primes in an imperative language, but I'd rather not. On the way to the functional solutions the worst one meets is the odd black hole, which a helpful compiler-writer may arrange to be pin-pointed in the source program. All the most intricate problems of scheduling and memory management are solved for free as part of the paradigm.

It is hard to beat the cute simplicity of the two-line sieve with which we began. But even after improvement, it cannot match the wheel-sieve programs for speed. Besides, the spirals have an elegance of their own.

Acknowledgements

My thanks to Richard Bird and John Hughes for helpful comments.

References

- Bird, R. S. (1984). Using circular programs to eliminate multiple traversals of data. *Acta Informatica*, **21**(3), 239–50.
- Pritchard, P. (1982). Explaining the wheel sieve. *Acta Informatica*, **17**, 477–485.
- Turner, D. A. (1975). *SASL language manual*. Tech. rept. CS/75/1. Department of Computational Science, University of St. Andrews.