# An algebra of scans

Ralf Hinze

Institut für Informatik III, Universität Bonn
Römerstraße 164, 53117 Bonn, Germany
`ralf@informatik.uni-bonn.de`
`http://www.informatik.uni-bonn.de/~ralf/`

**Abstract.** A parallel prefix circuit takes $n$ inputs $x_1$, $x_2$, ..., $x_n$ and produces the $n$ outputs $x_1$, $x_1 \circ x_2$, ..., $x_1 \circ x_2 \circ \cdots \circ x_n$, where '$\circ$' is an arbitrary associative binary operation. Parallel prefix circuits and their counterparts in software, parallel prefix computations or scans, have numerous applications ranging from fast integer addition over parallel sorting to convex hull problems. A parallel prefix circuit can be implemented in a variety of ways taking into account constraints on size, depth, or fanout. Traditionally, implementations are either defined graphically or by enumerating the underlying graph. Both approaches have their pros and cons. A figure if well drawn conveys the possibly recursive structure of the scan but it is not amenable to formal manipulation. A description in form of a graph while rigorous obscures the structure of a scan and is equally hard to manipulate. In this paper we show that parallel prefix circuits enjoy a very pleasant algebra. Using only two basic building blocks and four combinators all standard designs can be described succinctly and rigorously. The rules of the algebra allow us to prove the circuits correct and to derive circuit designs in a systematic manner.

> LORD DARLINGTON. ... [*Sees a fan lying on the table.*] *And what a wonderful fan! May I look at it?*
> LADY WINDERMERE. *Do. Pretty, isn't it! It's got my name on it, and everything. I have only just seen it myself. It's my husband's birthday present to me. You know to-day is my birthday?*
> — Oscar Wilde, *Lady Windermere's Fan*

## 1   Introduction

A *parallel prefix computation* determines the sums of all prefixes of a given sequence of elements. The term sum has to be understood in a broad sense: parallel prefix computations are not confined to addition, any associative operation can be used. Functional programmers know parallel prefix computations as *scans*, a term which originates from the language APL [1]. We will use both terms synonymously.

Parallel prefix computations have numerous applications; the most well-known is probably the carry-lookahead adder [2], a *parallel prefix circuit*. Other

applications include the maximum segment sum problem, parallel sorting, solving recurrences, and convex hull problems, see [3].

A parallel prefix computation seems to be inherently sequential. However, it can be made to run in logarithmic time on a parallel architecture or in hardware. In fact, scans can be implemented in a variety of ways taking into account constraints on measures such as size, depth, or fan-out.

A particular implementation can be modelled as a directed acyclic oriented graph and this is what most papers on the subject actually do. The structure is a graph as opposed to a tree because subcomputations can be shared. Actually, it is an ordered graph, that is, the inputs of a node are ordered, because the underlying binary operation is not necessarily commutative. Here is an example graph.



The edges are directed downwards; a node of in-degree two, an *operation node*, represents the sum of its two inputs; a node of in-degree one and out-degree greater than one, a *duplication node*, distributes its input to its outputs.

Different implementations can be compared with respect to several different measures: the *size* is the number of operation nodes, the *depth* is the maximum number of operation nodes on any path, and the *fan-out* is the maximal out-degree of an operation node. In the example above the size is 74, the depth is 5, and the fan-out is 17. If implemented in hardware, the size and the fan-out determine the required chip area, the depth influences the speed. Other factors include regularity of layout and interconnection.

It is not too hard—but perhaps slightly boring—to convince oneself that the above circuit is correct: given $n = 32$ inputs $x_1, x_2, \ldots, x_n$ it produces the $n$ outputs $x_1, x_1 \circ x_2, \ldots, x_1 \circ x_2 \circ \cdots \circ x_n$, where '$\circ$' is the underlying binary operation. The 'picture as proof' technique works reasonably well for a parallel prefix circuit of a small fixed width. However, an implementation usually defines a family of circuits, one for each number of inputs. In this case, the graphical approach is not an option, especially, when it comes to proving correctness. Some papers define a family of graphs by numbering the nodes and enumerating the edges, see, for instance, [4]. While this certainly counts as a rigorous definition it is way too concrete: an explicit graph representation obscures the structure of the design and is hard to manipulate formally.

In this paper we show that parallel prefix circuits enjoy a pleasant algebra. Using only two basic building blocks and four combinators all standard designs can be described succinctly and rigorously. The rules of the algebra allow us to prove the circuits correct and to derive new designs in a systematic manner.

The rest of the paper is structured as follows. Section 2 motivates the basic combinators and their associated laws. Section 3 introduces two scan combinators: horizontal and vertical composition of scans. Using these combinators

various recursive constructions can be defined and proven correct, see Section 4. Section 5 discusses more sophisticated designs: minimum depth circuits that have the minimal number of operation nodes. Section 6 then considers size-optimal circuits with bounded fan-out. Finally, Section 7 reviews related work and Section 8 concludes.

## 2  Basic combinators

This section defines the algebra of scans. Throughout the paper we employ the programming language Haskell [5] as the meta language. In particular, Haskell's class system is put to good use: classes allow us to define algebras and instances allow us to define associated models.

### 2.1  Monoids

The binary operation underlying a scan must be associative. Without loss of generality we assume that it also has a neutral element so that we have a monoidal structure.

```
class Monoid α where
  ε    ::  α
  (∘)  ::  α → α → α
```

Each instance of *Monoid* must satisfy the following laws.

$$
\begin{aligned}
\varepsilon \circ x &= x \\
x \circ \varepsilon &= x \\
x \circ (y \circ z) &= (x \circ y) \circ z
\end{aligned}
$$

For example, the parallel prefix circuit that computes carries in a carry-lookahead adder is based on the following monoid.

```
data KPG  =  K | P | G
instance Monoid KPG where
  ε          =  P
  K ∘ f    =  K
  P ∘ f    =  f
  G ∘ f    =  G
```

The elements of the type $KPG$ represent *carry propagation functions*: $K$ kills a carry ($\lambda c \to 0$), $P$ propagates a carry ($\lambda c \to c$), and $G$ generates a carry ($\lambda c \to 1$). The operation '∘' implements function composition, which is associative and has the identity, $P$, as its neutral element.

## 2.2    The algebra of fans and scans

Reconsidering the example graph of the introduction we note that a parallel prefix circuit can be seen as a composition of *fans*. Here are fans of different widths in isolation.



A fan adds its first input—counting from left to right—to each of its remaining inputs. It consists of a duplication node and $n - 1$ operation nodes. A scan is constructed by arranging fans horizontally and vertically. As an example, the following scan consists of three fans: a 3-fan placed below two 2-fans.



Placing two circuits side by side is called *parallel* or *horizontal composition*, denoted '×'.



Placing two circuits on top of each other is called *serial* or *vertical composition*, denoted '⅋'. We require that the two circuits have the same width.



Horizontal and vertical composition, however, are not sufficient as combining forms as the following circuit demonstrates (which occurs as a subcircuit in the introductory example).



At first sight, it seems that a more general fan combinator is needed. The fans in the middle part are not contiguous: the first input is only propagated to each second remaining input, the other inputs are wired through. However, a moment's reflection reveals that the middle part is really the previous circuit stretched by a factor of two. This observation motivates the introduction of a stretch combinator: generalizing from a single stretch factor, the combinator '≻−' takes a list of widths and stretches a given circuit accordingly.

The inputs of the resulting circuit are grouped according to the given widths. In the example above, we have four groups, each of width 2. The *last* input of each group is connected to the argument circuit; the other inputs are wired through.

To summarize, the example parallel prefix circuit is denoted by the following algebraic expression ($fan_i$ represents a fan of width $i$ and $id_i$ represents the identity circuit of width $i$).

$$fan_2 \times fan_2 \times fan_2 \times fan_2 \;\fatsemi$$
$$[2, 2, 2, 2] \succ (fan_2 \times fan_2 \;\fatsemi\; id_1 \times fan_3) \;\fatsemi$$
$$id_1 \times fan_2 \times fan_2 \times fan_2 \times id_1$$

The following class declaration defines the algebra of fans and scans. Note that the class *Circuit* abstracts over a *type constructor* $\gamma$ which in turn is parameterized by the underlying monoid. The type variable $\gamma$ serves as a placeholder for the carrier of the algebra.

**type** $Width$ $=$ $Nat$
**type** $Width^+$ $=$ $Nat^+$
**class** $Circuit\ \gamma$ **where**

| | | |
|---|---|---|
| $fan$ | :: | $(Monoid\ \alpha) \Rightarrow Width \to \gamma\ \alpha$ |
| $id$ | :: | $Width \to \gamma\ \alpha$ |
| $(\fatsemi)$ | :: | $\gamma\ \alpha \to \gamma\ \alpha \to \gamma\ \alpha$ |
| $(\times)$ | :: | $\gamma\ \alpha \to \gamma\ \alpha \to \gamma\ \alpha$ |
| $(\succ)$ | :: | $[\,Width^+\,] \to \gamma\ \alpha \to \gamma\ \alpha$ |
| $(\prec)$ | :: | $\gamma\ \alpha \to [\,Width^+\,] \to \gamma\ \alpha$ |
| $\lvert \cdot \rvert$ | :: | $\gamma\ \alpha \to Width$ |

The above class declaration makes explicit that only fans rely on the underlying monoidal structure; the remaining combinators can be seen as glue. The class additionally introduces a second stretch combinator '$\prec$' which is similar to '$\succ$' except that it connects the *first* input of each group to its argument circuit. The following pictures illustrate the difference between the two combinators.

$$[2, 3, 1] \;\succ\; fan_3 \quad = \quad \text{(diagram)} \qquad fan_3 \prec [2, 3, 1] \quad = \quad \text{(diagram)}$$

We shall see that '$\succ$' is useful for combining scans, while '$\prec$' is a natural choice for combining fans. The list argument of the stretch combinators must contain positive widths ($Nat^+$ is the type of naturals excluding zero).

The width of a circuit, say $f$, is denoted $\lvert f \rvert$. Being able to query the width of a circuit is important as some combinators are subject to width constraints: $f \fatsemi g$ is only defined if $\lvert f \rvert = \lvert g \rvert$, $f \prec x$ and $x \succ f$ require that $\lvert f \rvert = \#x$ where $\#x$ denotes the length of the list $x$. In particular, $f \prec [\,]$ is only valid if $\lvert f \rvert = 0$. We lift the width combinator to lists of circuits abbreviating $[\,\lvert f \rvert \mid f \leftarrow fs\,]$ by $\lvert fs \rvert$.

To save parentheses we agree that '$\prec$' and '$\succ$' bind more tightly than '$\times$', which in turn takes precedence over '$\mathbin{\raise0.3ex\hbox{$\scriptscriptstyle\circ\atop\circ$}}$'.

> **infixr** 1  $\mathbin{\raise0.3ex\hbox{$\scriptscriptstyle\circ\atop\circ$}}$
> **infixr** 2  $\times$
> **infix** 4  $\succ\!\!-, \prec\!\!-$

The fixity declarations furthermore ensure that the combinators bind less tightly than Haskell's list concatenation '$+\!\!\!+$'. As an example, $f \times g \prec\!\!- x +\!\!\!+ y \mathbin{\raise0.3ex\hbox{$\scriptscriptstyle\circ\atop\circ$}} h$ abbreviates $(f \times (g \prec\!\!- (x +\!\!\!+ y))) \mathbin{\raise0.3ex\hbox{$\scriptscriptstyle\circ\atop\circ$}} h$.

The following derived combinators will prove useful in the sequel.

> $par$   ::  $(Circuit\ \gamma) \Rightarrow [\gamma\ \alpha] \rightarrow \gamma\ \alpha$
> $par$   =  $foldr\ (\times)\ id_0$
> $seq$   ::  $(Circuit\ \gamma) \Rightarrow Width \rightarrow [\gamma\ \alpha] \rightarrow \gamma\ \alpha$
> $seq_n$  =  $foldr\ (\mathbin{\raise0.3ex\hbox{$\scriptscriptstyle\circ\atop\circ$}})\ id_n$

The combinator $par$ generalizes '$\times$' and places a list of circuits side by side. Likewise, $seq$ generalizes '$\mathbin{\raise0.3ex\hbox{$\scriptscriptstyle\circ\atop\circ$}}$' and places a list of circuits above each other.

> **infix** 4  $\succ, \prec$
> $(\succ)$    ::  $(Circuit\ \gamma) \Rightarrow [\gamma\ \alpha] \rightarrow \gamma\ \alpha \rightarrow \gamma\ \alpha$
> $fs \succ f$  =  $par\ fs \mathbin{\raise0.3ex\hbox{$\scriptscriptstyle\circ\atop\circ$}} |fs| \succ\!\!- f$
> $(\prec)$    ::  $(Circuit\ \gamma) \Rightarrow \gamma\ \alpha \rightarrow [\gamma\ \alpha] \rightarrow \gamma\ \alpha$
> $f \prec fs$  =  $f \prec\!\!- |fs| \mathbin{\raise0.3ex\hbox{$\scriptscriptstyle\circ\atop\circ$}} par\ fs$

The combinators '$\succ$' and '$\prec$' are convenient variants of '$\succ\!\!-$' and '$\prec\!\!-$': the expression $f \prec [f_1, \ldots, f_n]$ connects the $i$-th output of $f$ to the first input of $f_i$ while $[f_1, \ldots, f_n] \succ f$ connects the last output of $f_i$ to the $i$-th input of $f$. Thus, '$\succ$' is similar to the composition of an $n$-ary function with $n$ argument functions.

In Haskell, we can model circuits as list processing functions of type $[\alpha] \rightarrow [\alpha]$ where $\alpha$ is the underlying monoid. Serial composition is then simply forward function composition; parallel composition satisfies $(f \times g)\ (x +\!\!\!+ y) = f\ x +\!\!\!+ g\ y$ where '$+\!\!\!+$' denotes list concatenation and $|f| = \#x$, $|g| = \#y$. Figure 1 displays the complete instance declaration, which can be seen as *the standard model of Circuit*. Put differently, the intended semantics of the combinators is given by the list processing functions in Figure 1. Some remarks are in order. The expression $\Sigma x$ denotes the sum of the elements of the list $x$. The function *group* that is used in the definition of '$\prec\!\!-$' and '$\succ\!\!-$' takes a list of lengths and partitions its second argument accordingly. The expression $[e \mid a \leftarrow x \mid b \leftarrow y]$ is a *parallel list comprehension* and abbreviates $[e \mid (a, b) \leftarrow zip\ x\ y]$.

The algebraic laws each instance of the class *Circuit* has to satisfy are listed in Figure 2. The reader is invited to convince themself that the instance of Figure 1 is indeed a model in that sense. The list is not complete though: Figure 2 includes only the *structural laws*, rules that do not involve fans. The properties of fans will be discussed in separate paragraph below. Most of the laws except, perhaps, those concerned with '$\prec\!\!-$' and '$\succ\!\!-$' are straightforward: '$\mathbin{\raise0.3ex\hbox{$\scriptscriptstyle\circ\atop\circ$}}$' is associative with

**data** $Trans\ \alpha$ = $Trans\{\,width :: Width, apply :: [\alpha] \to [\alpha]\,\}$

**instance** $Circuit\ Trans$ **where**

$\quad fan_n \quad = \quad Trans\ n\ (\lambda u \to \textbf{case}\ u\ \textbf{of}$
$\qquad\qquad\qquad\qquad\qquad [\,] \to [\,]$
$\qquad\qquad\qquad\qquad\qquad a : as \to a : [\,a \circ b \mid b \leftarrow as\,])$

$\quad id_n \quad\ \ = \quad Trans\ n\ (\lambda u \to u)$

$\quad f \ \mathring{,}\ g \quad = \quad Trans\ |f|\ (\lambda u \to apply\ g\ (apply\ f\ u))$

$\quad f \times g \quad = \quad Trans\ (|f| + |g|)\ (\lambda u \to \textbf{let}\ (y, z) = splitAt\ |f|\ u$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \textbf{in}\ apply\ f\ y \mathbin{+\!\!+} apply\ g\ z)$

$\quad x \succ\!\!- f \quad = \quad Trans\ (\Sigma x)\ (\lambda u \to \textbf{let}\ ys = group\ x\ u$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad as = apply\ f\ [\,last\ y \mid y \leftarrow ys\,]$
$\qquad\qquad\qquad\qquad\qquad\qquad \textbf{in}\ concat\ [\,init\ y \mathbin{+\!\!+} [\,a\,] \mid y \leftarrow ys \mid a \leftarrow as\,])$

$\quad f -\!\!\prec x \quad = \quad Trans\ (\Sigma x)\ (\lambda u \to \textbf{let}\ ys = group\ x\ u$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad as = apply\ f\ [\,head\ y \mid y \leftarrow ys\,]$
$\qquad\qquad\qquad\qquad\qquad\qquad \textbf{in}\ concat\ [\,[\,a\,] \mathbin{+\!\!+} tail\ y \mid y \leftarrow ys \mid a \leftarrow as\,])$

$\quad |f| \qquad = \quad width\ f$

$group \qquad\qquad\ \ :: \quad [Int] \to [\alpha] \to [[\alpha]]$
$group\ [\,]\ as \qquad = \quad [\,]$
$group\ (i : x)\ as \quad = \quad bs : group\ x\ cs$
$\quad \textbf{where}\ (bs, cs) \quad = \quad splitAt\ i\ as$

**Fig. 1.** The standard model of the scan algebra

$$id_{|f|} \ \mathring{,}\ f \ = \ f$$
$$f \ \mathring{,}\ id_{|f|} \ = \ f$$
$$f \ \mathring{,}\ (g \ \mathring{,}\ h) \ = \ (f \ \mathring{,}\ g) \ \mathring{,}\ h$$

$$id_0 \times f \ = \ f$$
$$f \times id_0 \ = \ f$$
$$f \times (g \times h) \ = \ (f \times g) \times h$$
$$id_m \times id_n \ = \ id_{m+n}$$
$$(f \times g) \ \mathring{,}\ (f' \times g') \ = \ (f \ \mathring{,}\ f') \times (g \ \mathring{,}\ g')$$

$$|id_n| \ = \ n$$
$$|f \ \mathring{,}\ g| \ = \ |f| = |g|$$
$$|f \times g| \ = \ |f| + |g|$$
$$|fan_n| \ = \ n$$
$$|f -\!\!\prec x| \ = \ \Sigma x$$
$$|x \succ\!\!- f| \ = \ \Sigma x$$

$$id_{\#x} -\!\!\prec x \ = \ id_{\Sigma x}$$
$$f -\!\!\prec replicate\ |f|\ 1 \ = \ f$$
$$(f \ \mathring{,}\ g) -\!\!\prec x \ = \ (f -\!\!\prec x) \ \mathring{,}\ (g -\!\!\prec x)$$
$$(f \times g) -\!\!\prec (x \mathbin{+\!\!+} y) \ = \ (f -\!\!\prec x) \times (g -\!\!\prec y)$$
$$(f -\!\!\prec x) -\!\!\prec y \ = \ f -\!\!\prec [\,\Sigma z \mid z \leftarrow group\ x\ y\,]$$
$$id_{i-1} \times (f -\!\!\prec y \mathbin{+\!\!+} [k]) \ = \ ([i] \mathbin{+\!\!+} y \succ\!\!- f) \times id_{k-1}$$

**Fig. 2.** The structural laws of the scan algebra

$id_n$ as its neutral element; '$\times$' is associative with $id_0$ as its neutral element; '$\times$' preserves identity and vertical composition. Most of the laws are subject to width constraints: $(f \times g) \mathbin{\fatsemi} (f' \times g') = (f \mathbin{\fatsemi} f') \times (g \mathbin{\fatsemi} g')$, for instance, is only valid if $|f| = |f'|$ and $|g| = |g'|$. Use of these laws in subsequent proofs will be signalled by the hint *composition*.

Figure 2 only lists the laws for '$\prec$'; its companion combinator '$\succ$' satisfies analogous properties. The equations show that '$\prec$' preserves identity and composition (*replicate n a* constructs a list containing exactly $n$ copies of $a$). The second but last law in the right column demonstrates that nested occurrences of stretch combinators can be flattened. The last equation, termed *flip law*, shows that '$\prec$' can be defined in terms of '$\succ$' and vice versa. Recall that '$\succ$' connects last inputs and '$\prec$' connects first inputs. So strictly, only one stretch combinator is necessary. It is, however, convenient to have both at our disposal. Use of these laws will be signalled by the hint *stretching*.

As a warm-up in scan calculations, let us derive two simple consequences, which we need later on.

$$f \prec x + [j + k] \qquad = \quad (f \prec x + [j]) \times id_k \tag{1}$$
$$(f \times id_{\#y-1}) \prec x + y \quad = \quad f \prec x + [\Sigma y] \tag{2}$$

The rules allow us to push the identity, $id_n$, in and out of a stretch. To prove (1) we argue

$$f \prec x + [j + k]$$
$$= \quad \{ \text{ flip law } \}$$
$$([1] + x \succ f) \times id_{j+k-1}$$
$$= \quad \{ \text{ composition } \}$$
$$([1] + x \succ f) \times id_{j-1} \times id_k$$
$$= \quad \{ \text{ flip law } \}$$
$$(f \prec x + [j]) \times id_k$$

Property (2) is equally easy to show.

$$(f \times id_{\#y-1}) \prec x + y$$
$$= \quad \{ \text{ stretching } \}$$
$$(f \prec x + [head\ y]) \times (id_{\#y-1} \prec tail\ y)$$
$$= \quad \{ \text{ stretching } \}$$
$$(f \prec x + [head\ y]) \times id_{\Sigma(tail\ y)}$$
$$= \quad \{ \text{ derived stretch law (1) } \}$$
$$f \prec x + [\Sigma y]$$

Let us now turn to the axioms involving fans. Fans of width less than two are equal to the identities.

$$fan_0 \quad = \quad id_0$$
$$fan_1 \quad = \quad id_1$$

As an aside, this implies that the identity combinator, $id_n$, can be defined as a horizontal composition of fans.

$$id_n \;=\; \underbrace{fan_1 \times \cdots \times fan_1}_{n \text{ times}}$$

The first non-trivial fan law, equation (3) below, allows the designer of scans to trade depth for fan-out. Here is an instance of the law.



The circuit on the left has a depth of 2 and a fan-out of 5 while the circuit on the right has depth 1 and fan-out 8. The *first fan law* generalizes from the example.

$$fan_{1+n} \prec [fan_m \prec fs] + gs \;=\; fan_{m+n} \prec fs + gs \tag{3}$$

Interestingly, this rule is still structural as it does not rely on any properties of the underlying operator. Only the very last law, equation (4) below, employs the associativity of '∘'. Before we discuss the rule let us first take a look at some examples.



Both circuits have the same depth but the circuit on the right has fewer operation nodes. The left circuit consists of a big fan below a layer of smaller fans. The big fan adds its first input to each of the intermediate values; the same effect is achieved on the right by broadcasting the first input to each of the smaller fans. Here is the smallest instance of this optimization.



The left circuit, $id_1 \times fan_2 \,\fatsemi\, fan_3 = id_2 \prec [id_1, fan_2] \,\fatsemi\, fan_3$, maps the inputs $x_1$, $x_2$, $x_3$ to the outputs $x_1$, $x_1 \circ x_2$, $x_1 \circ (x_2 \circ x_3)$, while the right circuit, $fan_2 \times id_1 \,\fatsemi\, id_1 \times fan_2 = fan_2 \prec [fan_1, fan_2]$, maps $x_1$, $x_2$, $x_3$ to $x_1$, $x_1 \circ x_2$, $(x_1 \circ x_2) \circ x_3$. Clearly, the outputs are equal if and only if '∘' is associative. However, the first circuit consists of three operation nodes while the second requires only two. The *second fan law* captures this optimization.

$$\begin{aligned} & id_{1+\#x} \prec [id_i] + [fan_j \mid j \leftarrow y] \,\fatsemi\, fan_{i+\Sigma x} \\ &= \; fan_{1+\#x} \prec [fan_i] + [fan_j \mid j \leftarrow y] \end{aligned} \tag{4}$$

The size of the circuit of the right-hand side is always at most the size of the circuit on the left-hand side. Unless all the 'small' circuits are trivial, the depth of both circuits is the same. Thus, the second fan law is *the* central rule when it comes to optimizing scans.

In the sequel we will also need the following derived law, which is essentially a binary version of the second fan law.

$$id_m \times fan_{n+1} \mathbin{\text{\textsemicolon}} fan_{m+n+1} \;\;=\;\; fan_{1+m} \times id_n \mathbin{\text{\textsemicolon}} id_m \times fan_{n+1} \tag{5}$$

We argue as follows.

$$id_m \times fan_{n+1} \mathbin{\text{\textsemicolon}} fan_{m+n+1}$$
$=$    { second fan law }
$$fan_2 \prec [fan_m, fan_{n+1}]$$
$=$    { stretching }
$$fan_2 \prec [fan_m \prec replicate\ m\ id_1, fan_{n+1}]$$
$=$    { first fan law }
$$fan_{1+m} \prec replicate\ m\ id_1 \mathbin{+\!\!+} [fan_{n+1}]$$
$=$    { definition of '$\prec$' }
$$fan_{1+m} \prec\!\!\prec replicate\ m\ 1 \mathbin{+\!\!+} [n+1] \mathbin{\text{\textsemicolon}} par\ (replicate\ m\ id_1 \mathbin{+\!\!+} [fan_{n+1}])$$
$=$    { derived stretch law (1) }
$$(fan_{1+m} \prec\!\!\prec replicate\ m\ 1 \mathbin{+\!\!+} [1]) \times id_n \mathbin{\text{\textsemicolon}} par\ (replicate\ m\ id_1 \mathbin{+\!\!+} [fan_{n+1}])$$
$=$    { stretching }
$$fan_{1+m} \times id_n \mathbin{\text{\textsemicolon}} par\ (replicate\ m\ id_1 \mathbin{+\!\!+} [fan_{n+1}])$$
$=$    { composition }
$$fan_{1+m} \times id_n \mathbin{\text{\textsemicolon}} id_m \times fan_{n+1}$$

## 3   Serial and parallel scan combinators

The combinators we have seen so far are the basic building blocks of scans. The blocks can be composed in a multitude of ways, the resulting circuits not necessarily implementing parallel prefix circuits. By contrast, the combining forms introduced in this section take scans to scans, they are *scan combinators*.

Before we proceed, we should first make precise what we mean by 'scan' in our framework. Scans are, like fans, parameterized by the width of the circuit. We specify

$$
\begin{aligned}
scan_0 &\;\;=\;\; id_0 \\
scan_{n+1} &\;\;=\;\; succ\ scan_n
\end{aligned}
$$

where *succ* is given by

$$
\begin{aligned}
succ &\;\; :: \;\; (Circuit\ \gamma, Monoid\ \alpha) \Rightarrow \gamma\ \alpha \rightarrow \gamma\ \alpha \\
succ\ f &\;\;=\;\; id_1 \times f \mathbin{\text{\textsemicolon}} fan_{|f|+1}
\end{aligned}
$$

Whenever we introduce a new implementation of scans in the sequel, we will show using the laws of the algebra that the family of circuits is equal to $scan_n$.

The first scan combinator implements the *serial or vertical composition of scans*: the last output of the first circuit is fed into the first input of the second circuit.

**infixr** 3 $\backslash\!\backslash$
$(\backslash\!\backslash) \quad :: \quad (Circuit\ \gamma) \Rightarrow \gamma\ \alpha \to \gamma\ \alpha \to \gamma\ \alpha$
$f \backslash\!\backslash g \;\;=\;\; f \times id_{|g|-1} \,\fatsemi\, id_{|f|-1} \times g$

Because of the overlap the width of the resulting circuit is one less than the sum of the widths of the two arguments: $|f \backslash\!\backslash g| = |f| + |g| - 1$. The depth does not necessarily increase, as the following example illustrates.



The rightmost operation node of the first circuit is placed upon the uppermost leftmost duplication node of the second circuit.

Serial composition of scans is associative with $id_1$ as its neutral element.

$$id_1 \backslash\!\backslash f \quad\;\; = \;\; f$$
$$f \backslash\!\backslash id_1 \quad\;\; = \;\; f$$
$$f \backslash\!\backslash (g \backslash\!\backslash h) \;\; = \;\; (f \backslash\!\backslash g) \backslash\!\backslash h$$

The first two laws are straightforward to show; the proof of associativity is quite instructive: it reveals that $f \backslash\!\backslash (g \backslash\!\backslash h)$ and $(f \backslash\!\backslash g) \backslash\!\backslash h$ are even structurally equivalent, that is, they can be rewritten into each other using only structural rules.

$$\begin{aligned}
&\quad f \backslash\!\backslash (g \backslash\!\backslash h) \\
=&\quad \{\text{ definition of `}\backslash\!\backslash\text{' }\} \\
&\quad f \times id_{|g|+|h|-1} \,\fatsemi\, id_{|f|-1} \times (g \backslash\!\backslash h) \\
=&\quad \{\text{ definition of `}\backslash\!\backslash\text{' }\} \\
&\quad f \times id_{|g|+|h|-1} \,\fatsemi\, id_{|f|-1} \times (g \times id_{|h|-1} \,\fatsemi\, id_{|g|-1} \times h) \\
=&\quad \{\text{ composition }\} \\
&\quad f \times id_{|g|+|h|-1} \,\fatsemi\, id_{|f|-1} \times g \times id_{|h|-1} \,\fatsemi\, id_{|f|+|g|-2} \times h \\
=&\quad \{\text{ composition }\} \\
&\quad (f \times id_{|g|-1} \,\fatsemi\, id_{|f|-1} \times g) \times id_{|h|-1} \,\fatsemi\, id_{|f|+|g|-2} \times h \\
=&\quad \{\text{ definition of `}\backslash\!\backslash\text{' }\} \\
&\quad (f \backslash\!\backslash g) \times id_{|h|-1} \,\fatsemi\, id_{|f|+|g|-2} \times h \\
=&\quad \{\text{ definition of `}\backslash\!\backslash\text{' }\} \\
&\quad (f \backslash\!\backslash g) \backslash\!\backslash h
\end{aligned}$$

Serial composition interacts nicely with stretching. Let $\#x = |f| - 1$ and $\#y = |g|$, then

$$(f \backslash\!\backslash g) \prec x \mathbin{+\!\!+} y \;\; = \;\; (f \prec x \mathbin{+\!\!+} [1]) \backslash\!\backslash (g \prec y) \tag{6}$$

The proof builds upon the derived stretch laws.

$$(f \setminus\!\!\setminus g) \prec x \mathbin{+\mkern-8mu+} y$$
$$= \quad \{ \text{ definition of `}\setminus\!\!\setminus\text{' } \}$$
$$(f \times id_{|g|-1} \mathbin{\fatsemi} id_{|f|-1} \times g) \prec x \mathbin{+\mkern-8mu+} y$$
$$= \quad \{ \text{ stretching } \}$$
$$(f \times id_{|g|-1}) \prec x \mathbin{+\mkern-8mu+} y \mathbin{\fatsemi} (id_{|f|-1} \times g) \prec x \mathbin{+\mkern-8mu+} y$$
$$= \quad \{ \text{ derived stretch laws (1) and (2) } \}$$
$$(f \prec x \mathbin{+\mkern-8mu+} [1]) \times id_{\Sigma y - 1} \mathbin{\fatsemi} (id_{|f|-1} \times g) \prec x \mathbin{+\mkern-8mu+} y$$
$$= \quad \{ \text{ stretching } \}$$
$$(f \prec x \mathbin{+\mkern-8mu+} [1]) \times id_{\Sigma y - 1} \mathbin{\fatsemi} id_{\Sigma x} \times (g \prec y)$$
$$= \quad \{ \text{ definition of `}\setminus\!\!\setminus\text{' } \}$$
$$(f \prec x \mathbin{+\mkern-8mu+} [1]) \setminus\!\!\setminus (g \prec y)$$

The second scan combinator is the *parallel or horizontal composition of scans*: both circuits are placed side by side, an additional fan adds the last output of the left circuit to each output of the right circuit.

**infixl** 3 $\rrbracket$
$(\rrbracket)$ :: $(Circuit\ \gamma, Monoid\ \alpha) \Rightarrow \gamma\ \alpha \rightarrow \gamma\ \alpha \rightarrow \gamma\ \alpha$
$f \rrbracket g = f \times g \mathbin{\fatsemi} id_{|f|-1} \times fan_{|g|+1}$

The widths sum up: $|f \rrbracket g| = |f| + |g|$. Because of the additional fan the depth increases by one. Here is an example application of `$\rrbracket$'.



Before we turn to the algebraic properties of `$\rrbracket$', let us first note that the parallel composition of scans is really a serial composition in disguise.

$$f \rrbracket g = f \setminus\!\!\setminus succ\ g$$

The proof is straightforward.

$$f \rrbracket g$$
$$= \quad \{ \text{ definition of `}\rrbracket\text{' } \}$$
$$f \times g \mathbin{\fatsemi} id_{|f|-1} \times fan_{|g|+1}$$
$$= \quad \{ \text{ composition } \}$$
$$f \times id_{|g|} \mathbin{\fatsemi} id_{|f|} \times g \mathbin{\fatsemi} id_{|f|-1} \times fan_{|g|+1}$$
$$= \quad \{ \text{ composition } \}$$
$$f \times id_{|g|} \mathbin{\fatsemi} id_{|f|-1} \times (id_1 \times g \mathbin{\fatsemi} fan_{|g|+1})$$
$$= \quad \{ \text{ definition of `}\setminus\!\!\setminus\text{' } \}$$

$$f \, \backslash\!\backslash \, (id_1 \times g \,\fatsemi\, fan_{|g|+1})$$
$$= \quad \{ \text{ definition of } succ \}$$
$$f \, \backslash\!\backslash \, succ \; g$$

Parallel composition is associative, as well, and has $id_0$ as its right unit. It does not possess a left unit though as $id_0 \, [] \, f$ is undefined (the first argument must have a positive width).

$$f \, [] \, id_0 \quad\;\; = \; f$$
$$f \, [] \, (g \, [] \, h) \;=\; (f \, [] \, g) \, [] \, h$$

As opposed to serial composition, the circuits $f \, [] \, (g \, [] \, h)$ and $(f \, [] \, g) \, [] \, h$ are *not* structurally equivalent: the latter circuit has fewer operation nodes. The proof rests upon the above characterization of parallel composition.

$$f \, [] \, (g \, [] \, h)$$
$$= \quad \{ \text{ characterization of '[]' } \}$$
$$f \, \backslash\!\backslash \, succ \; (g \, \backslash\!\backslash \, succ \; h)$$
$$= \quad \{ \text{ see below } \}$$
$$f \, \backslash\!\backslash \, succ \; g \, \backslash\!\backslash \, succ \; h$$
$$= \quad \{ \text{ characterization of '[]' } \}$$
$$(f \, [] \, g) \, [] \, h$$

The second step is justified by the following calculations.

$$succ \; (f \, \backslash\!\backslash \, succ \; g)$$
$$= \quad \{ \text{ definition of } succ \text{ and '}\backslash\!\backslash\text{' } \}$$
$$id_1 \times f \times id_{|g|} \,\fatsemi\, id_{|f|+1} \times g \,\fatsemi\, id_{|f|} \times fan_{|g|+1} \,\fatsemi\, fan_{|f|+|g|+1}$$
$$= \quad \{ \text{ derived fan law (5) } \}$$
$$id_1 \times f \times id_{|g|} \,\fatsemi\, id_{|f|+1} \times g \,\fatsemi\, fan_{|f|+1} \times id_{|g|} \,\fatsemi\, id_{|f|} \times fan_{|g|+1}$$
$$= \quad \{ \text{ composition } \}$$
$$id_1 \times f \times id_{|g|} \,\fatsemi\, fan_{|f|+1} \times id_{|g|} \,\fatsemi\, id_{|f|+1} \times g \,\fatsemi\, id_{|f|} \times fan_{|g|+1}$$
$$= \quad \{ \text{ definition of } succ \text{ and '}\backslash\!\backslash\text{' } \}$$
$$succ \; f \, \backslash\!\backslash \, succ \; g$$

Since the proof relies on the second fan law, $succ \; f \, \backslash\!\backslash \, succ \; g$ has fewer nodes than $succ \; (f \, \backslash\!\backslash \, succ \; g)$.

Let us finally record the fact that $succ$, '$\backslash\!\backslash$' and '$[]$' are *scan combinators*.

$$succ \; scan_n \quad\quad = \;\; scan_{n+1}$$
$$scan_{m+1} \, \backslash\!\backslash \, scan_n \;=\;\; scan_{m+n}$$
$$scan_m \, [] \, scan_n \quad = \;\; scan_{m+n}$$

The first law holds by definition. The third equation implies the second and the third equation can be shown by a straightforward induction over $m$.

## 4    Simple scans

It is high time to look at some implementations of parallel prefix circuits. We have already encountered one of the most straightforward implementations, a simple nest of fans, which serves as the specification.

$$
\begin{array}{rcl}
scan_0 & = & id_0 \\
scan_{n+1} & = & succ\ scan_n
\end{array}
$$

Here is an example circuit of width 8.



The circuit $scan_n$ is, in fact, the worst possible implementation as it has maximum depth and the maximal number of operation nodes, namely, $n*(n-1)/2$ among all scans of the same width. Since $succ\ f = id_1 \diagdown succ\ f = id_1 \parallel f$, we can alternatively define $scan_n$ as a parallel composition of trivial circuits.

$$
scan_{n+1} \;\; = \;\; id_1 \parallel scan_n
$$

Now, if we bracket the parallel composition differently, we obtain the *serial scan*, whose correctness is immediate.

$$
\begin{array}{rcl}
ser_0 & = & id_0 \\
ser_1 & = & id_1 \\
ser_{n+1} & = & ser_n \parallel id_1
\end{array}
$$

The graphical representation illustrates why $ser_n$ is called serial scan.



The serial scan has maximum depth, but the least number of operation nodes, namely, $n-1$ among all scans of the same width. In a sequential language $ser_n$ is the implementation of choice; it corresponds, for instance, to Haskell's *scanl* operation. Using $f \parallel id_1 = f \diagdown succ\ id_1 = f \diagdown fan_2$ we can rewrite the definition of $ser_n$ to emphasize its serial nature.

$$
ser_{n+1} \;\; = \;\; ser_n \diagdown fan_2
$$

Now, if we balance the parallel composition more evenly, we obtain parallel prefix circuits of minimum depth.

$$
\begin{aligned}
&rec_n \\
&\quad | \ n \leqslant 1 \qquad = \quad id_n \\
&\quad | \ otherwise \quad = \quad rec_{\lceil n/2 \rceil} \ [\![ \ rec_{\lfloor n/2 \rfloor}
\end{aligned}
$$

Here is a minimum-depth circuit of width 32.



Note that the tree of operation nodes that computes the last output is fully balanced, which explains why the depth is minimal. If the width is not a power of two, then $rec_n$ constructs a slightly skewed tree, known as a Braun tree [6]. Since '$[\![$' is associative, we can, of course, realize arbitrary tree shapes; other choices include *left-complete trees* or *quasi left-complete trees* [7]. For your amusement, here is a Fibonacci-tree of width 34



defined in the obvious way.

$$
\begin{aligned}
fib_0 \quad &= \quad id_0 \\
fib_1 \quad &= \quad id_1 \\
fib_{n+2} \quad &= \quad fib_{n+1} \ [\![ \ fib_n
\end{aligned}
$$

## 5    Depth-optimal scans

### 5.1    Brent-Kung circuits

The $rec_n$ family of circuits implements a simple divide-and-conquer scheme. A different recursive decomposition was devised by Brent and Kung [8]. As an example, here is a Brent-Kung circuit of width 32.

The inputs are 'paired' using a layer of 2-fans. Every second output is then fed into a Brent-Kung circuit of half the width; the other inputs are wired through. A final layer of 2-fans, shifted by one position, distributes the results of the nested Brent-Kung circuit to the wired-through signals. Every recursive step halves the number of inputs and increases the depth by *two*. Consequently, Brent-Kung circuits have logarithmic but not minimum depth. On the other hand, they use fewer operation nodes than the $rec_n$ circuits and furthermore they have only a fan-out of 2!

Turning to the algebraic description, we note that the first layer of 2-fans can be generalized to a layer of *scans* of arbitrary, not necessarily equal widths.

$$
\begin{array}{lll}
(\rhd) & :: & (Circuit\ \gamma, Monoid\ \alpha) \Rightarrow [\gamma\ \alpha] \rightarrow \gamma\ \alpha \rightarrow \gamma\ \alpha \\
[\,]\rhd g & = & g \\
(f:fs) \rhd g & = & (f:fs) \succ g \fatsemi id_{|f|-1} \times par\ gs \\
\quad \textbf{where}\ gs & = & [fan_{|f|} \mid f \leftarrow fs] \plusplus [id_1]
\end{array}
$$

Each scan, except the first one, is complemented by a *fan* in the final layer, shifted one position to the left. The operator '$\rhd$' is also a *scan combinator*; it takes a list of scans and a scan to a resulting scan.

$$
[scan_i \mid i \leftarrow x] \rhd scan_{\#x} \quad = \quad scan_{\Sigma x} \tag{7}
$$

The Brent-Kung circuit is given by the following definition.

$$
\begin{array}{ll}
bk_n & \\
\quad \mid n \leqslant 1 & = \quad id_n \\
\quad \mid otherwise & = \quad (replicate\ \lfloor n/2 \rfloor\ fan_2 \plusplus [id_1 \mid odd\ n]) \rhd bk_{\lceil n/2 \rceil}
\end{array}
$$

The nested scan has width $\lceil n/2 \rceil$: if the number of inputs is odd, then the nested scan additionally takes the last input. As an aside to non-Haskell experts, the idiom $[e \mid b]$ is a trivial list comprehension that evaluates to $[\,]$ if $b$ is *False* and to $[e]$ if $b$ is *True*. Furthermore note, that $bk_n$ is a so-called *restricted* parallel prefix circuit, whose last output has minimum depth.

The Brent-Kung decomposition is based on the binary number system. Since the operator '$\rhd$' works for arbitrary scans, it is not hard to generalize the decomposition to an arbitrary base.

$$
\begin{array}{lll}
gbk\ b\ n \\
\quad |\ n \leqslant b & = & ser_n \\
\quad |\ r == 0 & = & (replicate\ d\ ser_b) \rhd gbk\ b\ d \\
\quad |\ otherwise & = & (replicate\ d\ ser_b \mathbin{+\!\!+} [ser_r]) \rhd gbk\ b\ (d+1) \\
\quad \textbf{where}\ (d, r) & = & divMod\ n\ b
\end{array}
$$

The definition of $gbk$ uses serial scans as 'base' circuits. This is, of course, an arbitrary choice; any scan will do. Here is a base-3 circuit of width 27.



This circuit has size 46 and depth 8, while its binary cousin has size 47 and depth 10.

Let us turn to the proof that '$\rhd$' is a scan combinator. Property (7) can be proven by induction over the length of $x$. We confine ourselves to showing the induction step. Let $k = \#ss = \#fs$, $i = |s|$, $j = |head\ ss|$, $n = j + \Sigma|fs|$ and finally $fs = [fan_{|s|} \mid s \leftarrow tail\ ss] \mathbin{+\!\!+} [fan_1]$, then

$$
\begin{array}{ll}
& s : ss \rhd scan_{k+1} \\
= & \quad \{\ \text{definition of '}\rhd\text{'}\ \} \\
& s : ss \succ scan_{k+1} \mathbin{\fatsemi} id_{i-1} \times par\ (fan_j : fs) \\
= & \quad \{\ \text{property of scan}\ \} \\
& s : ss \succ (id_1 \mathbin{[\![} scan_k) \mathbin{\fatsemi} id_{i-1} \times par\ (fan_j : fs) \\
= & \quad \{\ \text{definition of '}[\![\text{'}\ \} \\
& s : ss \succ (id_1 \times scan_k \mathbin{\fatsemi} fan_{k+1}) \mathbin{\fatsemi} id_{i-1} \times par\ (fan_j : fs) \\
= & \quad \{\ \text{stretching}\ \} \\
& s : ss \succ (id_1 \times scan_k \mathbin{\fatsemi} fan_{k+1}) \mathbin{\fatsemi} id_{i-1} \times (id_{k+1} \prec fan_j : fs) \\
= & \quad \{\ \text{shift law (8), see below}\ \} \\
& s : ss \succ (id_1 \times scan_k) \mathbin{\fatsemi} id_{i-1} \times (fan_{k+1} \prec fan_j : fs) \\
= & \quad \{\ \text{fan law}\ \} \\
& s : ss \succ (id_1 \times scan_k) \mathbin{\fatsemi} id_{i-1} \times (id_{k+1} \prec id_j : fs \mathbin{\fatsemi} fan_n) \\
= & \quad \{\ \text{stretching}\ \} \\
& s : ss \succ (id_1 \times scan_k) \mathbin{\fatsemi} id_{i-1} \times (par\ (id_j : fs) \mathbin{\fatsemi} fan_n)
\end{array}
$$

$\qquad = \quad \{ \text{ composition } \}$

$\qquad s \times (ss \succ scan_k) \mathbin{\text{\textfractionsolidus}} id_{i-1} \times (par\ (id_j : fs) \mathbin{\text{\textfractionsolidus}} fan_n)$

$\qquad = \quad \{ \text{ composition } \}$

$\qquad s \times (ss \succ scan_k) \mathbin{\text{\textfractionsolidus}} id_{i-1} \times par\ (id_j : fs) \mathbin{\text{\textfractionsolidus}} id_{i-1} \times fan_n$

$\qquad = \quad \{ \text{ composition } \}$

$\qquad s \times (ss \succ scan_k \mathbin{\text{\textfractionsolidus}} par\ (id_{j-1} : fs)) \mathbin{\text{\textfractionsolidus}} id_{i-1} \times fan_n$

$\qquad = \quad \{ \text{ definition of '}[\,]\text{' } \}$

$\qquad s \,[\,]\, (ss \succ scan_k \mathbin{\text{\textfractionsolidus}} par\ (id_{j-1} : fs))$

$\qquad = \quad \{ \text{ definition of } par \}$

$\qquad s \,[\,]\, (ss \succ scan_k \mathbin{\text{\textfractionsolidus}} id_{j-1} \times par\ fs)$

$\qquad = \quad \{ \text{ definition of '}\rhd\text{' } \}$

$\qquad s \,[\,]\, (ss \rhd scan_k)$

The *shift law*, used in the fourth step, is a combination of the flip law and the laws for stretching.

$$(fs \succ (l \mathbin{\text{\textfractionsolidus}} m)) \times f \mathbin{\text{\textfractionsolidus}} g \times (r \prec gs) \;=\; (fs \succ l) \times f \mathbin{\text{\textfractionsolidus}} g \times ((m \mathbin{\text{\textfractionsolidus}} r) \prec gs) \quad (8)$$

We reason as follows.

$\qquad (fs \succ (l \mathbin{\text{\textfractionsolidus}} m)) \times f \mathbin{\text{\textfractionsolidus}} g \times (r \prec gs)$

$\qquad = \quad \{ \text{ composition } \}$

$\qquad par\ fs \times f \mathbin{\text{\textfractionsolidus}} (|fs| \succ\!\!- (l \mathbin{\text{\textfractionsolidus}} m)) \times id_{|f|} \mathbin{\text{\textfractionsolidus}} id_{|g|} \times (r \prec\!\!- |gs|) \mathbin{\text{\textfractionsolidus}} g \times par\ gs$

$\qquad = \quad \{ \text{ flip law } \}$

$\qquad par\ fs \times f \mathbin{\text{\textfractionsolidus}} id_{|g|} \times ((l \mathbin{\text{\textfractionsolidus}} m) \prec\!\!- |gs|) \mathbin{\text{\textfractionsolidus}} id_{|g|} \times (r \prec\!\!- |gs|) \mathbin{\text{\textfractionsolidus}} g \times par\ gs$

$\qquad = \quad \{ \text{ stretching } \}$

$\qquad par\ fs \times f \mathbin{\text{\textfractionsolidus}} id_{|g|} \times (l \prec\!\!- |gs|) \mathbin{\text{\textfractionsolidus}} id_{|g|} \times ((m \mathbin{\text{\textfractionsolidus}} r) \prec\!\!- |gs|) \mathbin{\text{\textfractionsolidus}} g \times par\ gs$

$\qquad = \quad \{ \text{ flip law } \}$

$\qquad par\ fs \times f \mathbin{\text{\textfractionsolidus}} (|fs| \succ\!\!- l) \times id_{|f|} \mathbin{\text{\textfractionsolidus}} id_{|g|} \times ((m \mathbin{\text{\textfractionsolidus}} r) \prec\!\!- |gs|) \mathbin{\text{\textfractionsolidus}} g \times par\ gs$

$\qquad = \quad \{ \text{ composition } \}$

$\qquad (fs \succ l) \times f \mathbin{\text{\textfractionsolidus}} g \times ((m \mathbin{\text{\textfractionsolidus}} r) \prec gs)$

### 5.2   Ladner-Fischer circuits

Can we combine the good properties of *rec* and *bk*—*rec* has minimum depth, while *bk* gets away with fewer operation nodes? Yes, we can! Reconsider the circuit $rec_{32}$ in Section 4 and note that the left part does not occupy the bottom level. The idea, which is due to Ladner and Fischer [9], is to use the Brent-Kung decomposition for the left part—recall that it increases the depth by two—and the 'usual' decomposition for the right part. The following combinator captures one step of the Brent-Kung scheme.

$$double \quad :: \quad (Circuit\ \gamma, Monoid\ \alpha) \Rightarrow (Width \to \gamma\ \alpha) \to (Width \to \gamma\ \alpha)$$
$$double\ s\ n \quad = \quad (replicate\ \lfloor n/2 \rfloor\ fan_2 \mathbin{+\!\!+} [id_1 \mid odd\ n]) \rhd s\ \lceil n/2 \rceil$$

Using *double* we can define a *depth-optimal* parallel prefix circuit that has the minimal number of operation nodes among all minimum-depth circuits [10].

$$
\begin{array}{ll}
opt\ n & \\
\quad \mid n \leqslant 1 & = \quad id_n \\
\quad \mid otherwise & = \quad double\ opt\ \lceil n/2 \rceil\ [\![]\!]\ opt\ \lfloor n/2 \rfloor
\end{array}
$$

The following example circuit of width 32 illustrates that all layers are nicely exploited.



The size of the circuit is 74. By contrast, $rec_{32}$ consists of 80 operation nodes.

The *double* combinator allows the scan designer to trade depth for size. The *Ladner-Fischer circuit*, defined below, generalizes *opt* introducing the notion of *extra depth*: the first argument of *lf* specifies the extra depth that the designer is willing to accept in return for a smaller size.

$$
\begin{array}{ll}
lf\ k\ 0 & = \quad id_0 \\
lf\ k\ 1 & = \quad id_1 \\
lf\ 0\ n & = \quad lf\ 1\ \lceil n/2 \rceil\ [\![]\!]\ lf\ 0\ \lfloor n/2 \rfloor \\
lf\ (k+1)\ n & = \quad double\ (lf\ k)\ n
\end{array}
$$

It is not hard to see that *lf* 0 specializes to *opt* and *lf* $\infty$ specializes to *bk*. In a sense, Ladner-Fischer mediates between the two recursive decompositions.

## 6    Size-optimal scans

### 6.1    Lin-Hsiao circuits

The '$\triangleright$' combinator constructs a slightly asymmetric circuit: not every scan has a corresponding fan. Circuits with a more symmetric design were recently introduced by Lin and Hsiao [4]. As an example, here is one of their circuits of width 25, called $wl_6$.

Every scan in the upper part is complemented by a corresponding fan in the lower part. The two parts are joined by a '$\lrcorner\ulcorner$'-like shape (turned 90° degrees clockwise) that connects the first input to the last output. The '$\lrcorner\ulcorner$' combinator is easy to derive.

$$
\begin{aligned}
& scan_{n+1} \\
= \quad & \{ \text{ property of } scan \} \\
& id_1 \parallel scan_n \\
= \quad & \{ \text{ definition of `}\parallel\text{' } \} \\
& id_1 \times scan_n \mathbin{\fatsemi} fan_{n+1} \\
= \quad & \{ \text{ stretching } \} \\
& [id_1, scan_n] \succ id_2 \mathbin{\fatsemi} fan_{n+1} \\
= \quad & \{ \text{ fan laws } \} \\
& [id_1, scan_n] \succ id_2 \mathbin{\fatsemi} fan_2 \prec [fan_n, id_1]
\end{aligned}
$$

Thus, we define

$$
\begin{aligned}
(\lrcorner\ulcorner) \quad &:: \quad (Circuit\ \gamma, Monoid\ \alpha) \Rightarrow \gamma\ \alpha \to \gamma\ \alpha \to \gamma\ \alpha \\
f \lrcorner\ulcorner g \quad &= \quad [id_1, f] \succ id_2 \mathbin{\fatsemi} fan_2 \prec [g, id_1]
\end{aligned}
$$

We have $|f \lrcorner\ulcorner g| = |f| + 1 = |g| + 1$. The derivation above implies that $scan_n \lrcorner\ulcorner fan_n = scan_{n+1}$. The '$\lrcorner\ulcorner$' combinator constructs a so-called *zig-zag circuit* whose height difference is one. The *height difference* is the length of the path from the first input to the last output. The low height difference of one renders zig-zag circuits attractive for serial composition. This is utilized in [4] to construct size-optimal circuits. A *size-optimal circuit* has the minimal number of operation nodes among all circuits of a *fixed* given depth.

Perhaps surprisingly, a serial composition of two zig-zag circuits can again be written as a zig-zag circuit.

$$
(l_1 \lrcorner\ulcorner u_1) \big\backslash\!\big\backslash (l_2 \lrcorner\ulcorner u_2) \quad = \quad ([l_1, l_2] \succ scan_2) \lrcorner\ulcorner (fan_2 \prec [u_1, u_2]) \tag{9}
$$

To justify this we argue ($i_1 = |l_1| = |u_1|$ and $i_2 = |l_2| = |u_2|$)

$$
\begin{aligned}
& (l_1 \lrcorner\ulcorner u_1) \big\backslash\!\big\backslash (l_2 \lrcorner\ulcorner u_2) \\
= \quad & \{ \text{ definition of `}\big\backslash\!\big\backslash\text{' } \} \\
& (l_1 \lrcorner\ulcorner u_1) \times id_{i_2} \mathbin{\fatsemi} id_{i_1} \times (l_2 \lrcorner\ulcorner u_2) \\
= \quad & \{ \text{ definition of `}\lrcorner\ulcorner\text{' } \} \\
& (id_1 \times l_1 \mathbin{\fatsemi} fan_2 \prec [u_1, id_1]) \times id_{i_2} \mathbin{\fatsemi} id_{i_1} \times (id_1 \times l_2 \mathbin{\fatsemi} fan_2 \prec [u_2, id_1]) \\
= \quad & \{ \text{ composition } \} \\
& id_1 \times l_1 \times l_2 \mathbin{\fatsemi} (fan_2 \prec [i_1, 1]) \times id_{i_2} \mathbin{\fatsemi} id_{i_1} \times (fan_2 \prec [i_2, 1]) \mathbin{\fatsemi} u_1 \times u_2 \times id_1 \\
= \quad & \{ \text{ derived stretch law (6) } \} \\
& id_1 \times l_1 \times l_2 \mathbin{\fatsemi} (fan_2 \big\backslash\!\big\backslash fan_2) \prec [i_1, i_2, 1] \mathbin{\fatsemi} u_1 \times u_2 \times id_1 \\
= \quad & \{ fan_2 \big\backslash\!\big\backslash fan_2 = scan_3 = id_1 \parallel scan_2 = id_1 \times scan_2 \mathbin{\fatsemi} fan_3 \}
\end{aligned}
$$

$$id_1 \times l_1 \times l_2 \; \mathbf{\mathring{;}} \; (id_1 \times scan_2 \; \mathbf{\mathring{;}} \; fan_3) \prec\!\!\!\prec [i_1, i_2, 1] \; \mathbf{\mathring{;}} \; u_1 \times u_2 \times id_1$$

$=$    { stretching }

$$id_1 \times l_1 \times l_2 \; \mathbf{\mathring{;}} \; [1, i_1, i_2] \succ\!\!\!- (id_1 \times scan_2) \; \mathbf{\mathring{;}} \; fan_3 \prec\!\!\!\prec [i_1, i_2, 1] \; \mathbf{\mathring{;}} \; u_1 \times u_2 \times id_1$$

$=$    { composition }

$$[id_1, l_1, l_2] \succ (id_1 \times scan_2) \; \mathbf{\mathring{;}} \; fan_3 \prec [u_1, u_2, id_1]$$

$=$    { composition and fan law }

$$id_1 \times ([l_1, l_2] \succ scan_2) \; \mathbf{\mathring{;}} \; fan_2 \prec [fan_2 \prec [u_1, u_2], id_1]$$

$=$    { definition of '$\!\!\_\!\!\lceil$' }

$$([l_1, l_2] \succ scan_2) \; \_\!\!\lceil \; (fan_2 \prec [u_1, u_2])$$

The property can even be generalized to an $n$-fold composition.

$$(l_1 \; \_\!\!\lceil \; u_1) \, \big\backslash \cdots \big\backslash \, (l_n \; \_\!\!\lceil \; u_n) \;\; = \;\; ([l_1, \ldots, l_n] \succ scan_n) \; \_\!\!\lceil \; (fan_n \prec [u_1, \ldots, u_n])$$

The proof of this property proceeds by a simple induction. We only show the induction step.

$$(l_1 \; \_\!\!\lceil \; u_1) \, \big\backslash \cdots \big\backslash \, (l_n \; \_\!\!\lceil \; u_n) \, \big\backslash \, (l_{n+1} \; \_\!\!\lceil \; u_{n+1})$$

$=$    { ex hypothesi }

$$(([l_1, \ldots, l_n] \succ scan_n) \; \_\!\!\lceil \; (fan_n \prec [u_1, \ldots, u_n])) \, \big\backslash \, (l_{n+1} \; \_\!\!\lceil \; u_{n+1})$$

$=$    { see above }

$$([[l_1, \ldots, l_n] \succ scan_n, l_{n+1}] \succ scan_2) \; \_\!\!\lceil \; (fan_2 \prec [fan_n \prec [u_1, \ldots, u_n], u_{n+1}])$$

$=$    { scan law (10), see below }

$$([l_1, \ldots, l_n, l_{n+1}] \succ scan_{n+1}) \; \_\!\!\lceil \; (fan_2 \prec [fan_n \prec [u_1, \ldots, u_n], u_{n+1}])$$

$=$    { fan law }

$$([l_1, \ldots, l_n, l_{n+1}] \succ scan_{n+1}) \; \_\!\!\lceil \; (fan_{n+1} \prec [u_1, \ldots, u_n, u_{n+1}])$$

The scan law used in the third step is analogous to the first fan law.

$$[fs \succ scan_m] \mathbin{+\!\!+} gs \succ scan_{1+n} \;\; = \;\; fs \mathbin{+\!\!+} gs \succ scan_{m+n} \tag{10}$$

The proof is left as an exercise to the reader.

   To summarize, '$\_\!\!\lceil$' combines a tree of scans with a corresponding tree of fans to a scan. The combinator allows us to shape a scan after an arbitrary tree structure. This makes it easy, for instance, to take constraints on the fan-out into account—the fan-out corresponds directly to the degree of a tree. As an example, let us define the Lin-Hsiao circuit $wl$ shown above. The following Haskell data declaration introduces a suitable tree type and its associated fold operation.

$$
\begin{array}{lll}
\textbf{data } \textit{Tree } \alpha & = & \textit{Leaf } \alpha \mid \textit{Node } [\textit{Tree } \alpha] \\
\textit{fold} & :: & (\alpha \to \beta) \to ([\beta] \to \beta) \to (\textit{Tree } \alpha \to \beta) \\
\textit{fold leaf node } (\textit{Leaf } a) & = & \textit{leaf } a \\
\textit{fold leaf node } (\textit{Node ts}) & = & \textit{node } [\textit{fold leaf node } t \mid t \leftarrow ts]
\end{array}
$$

The scan tree and the fan tree of a zig-zag circuit can be implemented as two simple folds.

$$
\begin{aligned}
\textit{zig-zag} \quad &::\quad (\textit{Circuit } \gamma, \textit{Monoid } \alpha) \Rightarrow \textit{Tree Width} \rightarrow \gamma\ \alpha \\
\textit{zig-zag t} \quad &=\quad \textit{fold ser s-node t} \smallfrown \textit{fold fan f-node t} \\
\textit{s-node ts} \quad &=\quad ts \succ ser_{\#ts} \\
\textit{f-node ts} \quad &=\quad fan_{\#ts} \prec ts
\end{aligned}
$$

The 'base' circuit of the *s-node* can be any scan. The same is true of the *f-node*—recall that the first fan law allows us to rewrite a single fan as a nest of fans.

Now, the tree underlying the *wl* circuit is given by the following definition (note that the argument does *not* correspond to the width).

$$
\begin{aligned}
\textit{wl-tree}_5 \quad &=\quad \textit{Node}\,[\textit{Leaf } 4, \textit{Leaf } 4, \textit{Leaf } 4] \\
\textit{wl-tree}_{n+1} \quad &=\quad \textit{Node}\,[\textit{wl-tree}_n, \textit{wl-tree}_n]
\end{aligned}
$$

The circuit is then simply defined as the composition of *zig-zag* and *wl-tree*.

$$
\textit{wl}_n \quad = \quad \textit{zig-zag wl-tree}_n
$$

Lin and Hsiao show that a slightly optimized version of $wl_n$—using the first fan law the two 2-fans in the center are merged into a 3-fan— is size-optimal [4].

## 6.2   Brent-Kung, revisited

Interestingly, the Brent-Kung circuit can be seen as a zig-zag circuit in disguise, or rather, as a serial composition of zig-zag circuits. Reconsider the example graph given in Section 5.1 and note that the right part has the characteristic shape of a zig-zag circuit: the tree in the upper part is mirrored in the lower part, in fact, they can be mapped onto each other through a 180° rotation (this is because binary fans and binary scans are equal).

The tree shape underlying a Brent-Kung circuit is that of a Braun tree.

$$
\begin{aligned}
\textit{braun}_n & \\
\mid n \leqslant 2 \quad &=\quad \textit{Leaf } n \\
\mid \textit{otherwise} \quad &=\quad \textit{Node}\,[\textit{braun}_{\lceil n/2 \rceil}, \textit{braun}_{\lfloor n/2 \rfloor}]
\end{aligned}
$$

Here is the alternative definition of Brent-Kung as a serial composition of zig-zag circuits.

$$
\begin{aligned}
\textit{bk}'_n & \\
\mid n \leqslant 2 \quad &=\quad \textit{ser}_n \\
\mid \textit{otherwise} \quad &=\quad \textit{bk}'_{d+r} \,\big\backslash\!\!\backslash\, \textit{zig-zag braun}_d \\
\textbf{where } (d, r) \quad &=\quad \textit{divMod n } 2
\end{aligned}
$$

The graphical representation reveals that this variant is more condensed: every fan is placed at the topmost possible position.

## 7   Related work

Parallel prefix computations are nearly as old as the history of computers. One of the first implementations of fast integer addition using carry-lookahead was described by Weinberger and Smith [11]. However, the operation of the circuit seemed to rely on the particularities of carry propagation. It was only 20 years later that Ladner and Fischer formulated the abstract problem of prefix computation and showed that carry computation is an instance of this class [9]. In fact, they showed the more general result that any finite-state transducer can be simulated in logarithmic time using a parallel prefix circuit.

As an aside, the idea underlying this proof is particularly appealing to functional programmers as it relies on currying. Let $\phi :: (X, A) \to A$ be an arbitrary binary operation not necessarily associative. To compute the value of

$$\phi\,(x_1, \phi\,(x_2, \ldots\ \phi\,(x_n, a)\ \ldots))$$

and all of the intermediate results we rewrite the expression into a form suitable for a prefix computation

$$(\mathit{curry}\ \phi\ x_1 \cdot \mathit{curry}\ \phi\ x_2 \cdot \cdots \cdot \mathit{curry}\ \phi\ x_n)\ a$$

The underlying binary operation is then simply function composition. An implementation in hardware additionally requires that the elements of $A \to A$ can be finitely represented (see Section 2.1).

Fich later proved that the Ladner-Fischer family of scans is depth-optimal [10]. Furthermore, he improved the design for an extra depth of one. Since then various other families have been proposed taking into account restrictions on depth and, in particular, on fan-out. Lin and Hsiao, for instance, describe a family of size-optimal scans with a fan-out of 4 and a small depth. One main ingredient is the circuit $wl$ introduced in Section 6.1. The construction is given as an algorithm that transforms an explicit graph representing $wl_n$ into a graph representing $wl_{n+1}$. The transformation essentially implements the rule

$$(l_1 \mathbin{\raise1pt\hbox{$\ulcorner$}} u_1) \mathbin{\big\backslash} (l_2 \mathbin{\raise1pt\hbox{$\ulcorner$}} u_2) \;\; = \;\; ([\,l_1, l_2\,] \succ \mathit{scan}_2) \mathbin{\raise1pt\hbox{$\ulcorner$}} (\mathit{fan}_2 \prec [\,u_1, u_2\,])$$

However, since the graph representation is too concrete, the algorithm is hard to understand and even harder to prove correct.

There are a few papers that deal with the derivation of parallel prefix circuits. Misra [12] calculates the Brent-Kung circuit via the data structure of *powerlists*.[1] Since powerlists capture the recursive decomposition of Brent-Kung, the approach while elegant is not easily applicable to other implementations of scans. In a recent pearl, O'Donnell and Rünger [13] derive the recursive implementation using the digital circuit description language Hydra. The resulting specification contains all the necessary information to simulate or fabricate a circuit.

The parallel prefix computation also serves as a building block of parallel programming. We have already noted in the introduction that many algorithms can be conveniently expressed in terms of scans [3]. Besides encouraging well-structured programming this coarse-grained approach to parallelism allows for various program optimizations. Gorlach and Lengauer [14], for instance, show that a composition of two scans can be transformed into a single scan. The scan function itself is an instance of a so-called *list homomorphism*. For this class of functions, parallel programs can be derived in a systematic manner [15]. Applying the approach of [15] to scan yields the optimal hypercube algorithm. This algorithm can be seen as a *clocked circuit*. Consequently, there is no direct correspondence to any of the algorithms given here, which are purely combinatorial.

## 8 Conclusion

This paper shows that parallel prefix circuits enjoy a surprisingly rich algebra. The algebraic approach has several benefits: it allows us to specify scans in a readable and concise way, to prove them correct, and to derive new designs. In the process of preparing the paper the algebra of scans has undergone several redesigns. We hope that the final version presented here will stand the test of time.

### Acknowledgements

---

[1] Misra actually claims to derive the Ladner-Fischer scheme. However, the function presented in the paper implements Brent-Kung—recall in this respect that *lf* $\infty$ specializes to *bk*.

# References

1. Iverson, K.E.: APL: A Programming Language. John Wiley & Sons (1962)
2. Cormen, T.H., Leiserson, C.E., Rivest, R.L.: Introduction to Algorithms. First edn. The MIT Press, Cambridge, Massachusetts (1990)
3. Blelloch, G.: Prefix sums and their applications. Technical Report CMU-CS-90-190, School of Computer Science, Carnegie Mellon University, Pittsburgh, PA (1990) Also appears in *Synthesis of Parallel Algorithms*, ed. John E. Reif, Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1993.
4. Lin, Y.C., Hsiao, J.W.: A new approach to constructing optimal prefix circuits with small depth. In: Proceedings of the International Symposium on Parallel Architectures, Algorithms and Networks (ISPAN'02). (2002)
5. Peyton Jones, S.: Haskell 98 Language and Libraries. Cambridge University Press (2003)
6. Braun, W., Rem, M.: A logarithmic implementation of flexible arrays. Memorandum MR83/4, Eindhoven University of Technology (1983)
7. Hinze, R.: Constructing red-black trees. In Okasaki, C., ed.: Proceedings of the Workshop on Algorithmic Aspects of Advanced Programming Languages, WAAAPL'99, Paris, France. (1999) 89–99 The proceedings appeared as a technical report of Columbia University, CUCS-023-99, also available from `http://www.cs.columbia.edu/~cdo/waaapl.html`.
8. Brent, R., Kung, H.: The chip complexity of binary arithmetic. In: Twelfth Annual ACM Symposium on Theory of Computing (STOC '80), New York, ACM Press (1980) 190–200
9. Ladner, R.E., Fischer, M.J.: Parallel prefix computation. Journal of the ACM **27** (1980) 831–838
10. Fich, F.E.: New bounds for parallel prefix circuits. In ACM, ed.: Proceedings of the fifteenth annual ACM Symposium on Theory of Computing, Boston, Massachusetts, April 25–27, 1983, New York, NY, USA, ACM Press (1983) 100–109
11. Weinberger, A., Smith, J.: A one-microsecond adder using one-megacycle circuitry. IRE Transactions on Electronic Computers **EC-5** (1956)
12. Misra, J.: Powerlist: A structure for parallel recursion. ACM Transactions on Programming Languages and Systems **16** (1994) 1737–1767
13. O'Donnell, J.T., Rünger, G.: Derivation of a logarithmic time carry lookahead addition circuit. Journal of Functional Programming, Special Issue on Functional Pearls (2004) to appear.
14. Gorlatch, S., Lengauer, C.: (De)Composition rules for parallel scan and reduction. In: Proc. 3rd Int. Working Conf. on Massively Parallel Programming Models (MPPM'97), IEEE Computer Society Press, pages 23–32 (1998)
15. Gorlatch, S.: Extracting and implementing list homomorphisms in parallel program development. Science of Computer Programming **33** (1999) 1–27