

# Data-Only Flattening for Nested Data Parallelism

Lars Bergstrom  
John Reppy  
Stephen Rosen  
Adam Shaw

University of Chicago  
{larsberg,jhr,sirosen,ams}@cs.uchicago.edu

Matthew Fluet  
Rochester Institute of Technology  
mtf@cs.rit.edu

Mike Rainey  
Max Planck Institute for  
Software Systems  
mrainey@mpi-sws.org

## Abstract

Data parallelism has proven to be an effective technique for high-level programming of a certain class of parallel applications, but it is not well suited to irregular parallel computations. Blelloch and others proposed *nested data parallelism* (NDP) as a language mechanism for programming irregular parallel applications in a declarative data-parallel style. The key to this approach is a compiler transformation that *flattens* the NDP computation and data structures into a form that can be executed efficiently on a wide-vector SIMD architecture. Unfortunately, this technique is ill suited to execution on today's multicore machines. We present a new technique, called *data-only flattening*, for the compilation of NDP, which is suitable for multicore architectures. Data-only flattening transforms nested data structures in order to expose programs to various optimizations while leaving control structures intact. We present a formal semantics of data-only flattening in a core language with a rewriting system. We demonstrate the effectiveness of this technique in the Parallel ML implementation and we report encouraging experimental results across various benchmark applications.

**Categories and Subject Descriptors** D.3.0 [Programming Languages]: General; D.3.2 [Programming Languages]: Language Classifications – Applicative (Functional) Programming, Concurrent, distributed, and parallel languages; D.3.4 [Programming Languages]: Processors – Compilers, Optimization

**Keywords** multicore, NESL, nested data parallelism, compilers

## 1. Introduction

Data-parallel computations are ones in which a function is applied to the elements of a collection (e.g., set or sequence) in parallel. Data parallelism is an effective technique to take advantage of parallel hardware and is especially suited to large-scale parallelism [10], but most languages that support data parallelism limit that support to *flat data parallelism* (FDP), where the computation being mapped over the collection does not contain nested data parallel computation. While FDP is very effective for many regular-parallel applications, it is not well-suited for irregular parallel ap-

plications. To address this weakness, Blelloch and others proposed *nested data parallelism* (NDP) [4, 5, 8, 20].

The basic operation in both flat and nested data parallelism is the *parallel map operation*, which applies a function to the elements of a collection in parallel. What distinguishes NDP from FDP is that elements of the collection may themselves be collections, and the mapped computation may itself involve parallel maps over the nested collections. Because the nested collections in an NDP computation may vary in size, it is difficult to ensure a balanced partitioning of work across multiple processors and it is difficult to execute the parallel computation with Single-Instruction-Multiple-Data (SIMD) architectures.

Blelloch addressed these challenges with an approach that he called *flattening*. Flattening (also called *vectorisation*) is a technique for converting irregular nested computations into regular computations on flat arrays. This approach, which was invented for first-order NDP by Sabot and Blelloch [2, 5] and extended to full-featured higher-order functional languages by Keller, Chakravarty, and others [7, 8, 14, 16], transforms both the data representations and the code so that the computation can be executed by a SIMD machine. An alternative approach, which is used in the Manticore system, is to execute the parallel map operations as fork-join parallelism and to rely on efficient work-stealing techniques to handle load balancing [1, 23].

In this paper, we introduce a new approach to implementing NDP constructs that is based on the idea of flattening the nested data representations, but not vectorising the code [25]. In our prototype implementation, we build on the Manticore compiler for PML [11], a parallel dialect of Standard ML [17].

This paper makes the following contributions:

1. We introduce a novel approach to implementing NDP that is well suited to execution on MIMD architectures, such as modern multicore processors.
2. We provide a formalization of our approach using a core calculus.
3. We present empirical evidence that our approach improves the performance of code that executes over irregular data while also preserving performance on regular data.

## 2. Full flattening for NDP

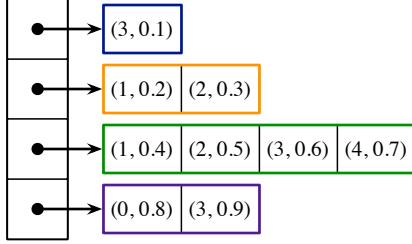
As a motivating example, consider the NDP code to implement sparse-matrix times dense-vector multiplication. We represent a dense vector as a parallel array of floating-point values:

```
type dense_vec = float parray
```

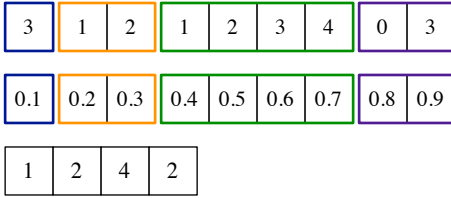
and a sparse vector is a parallel array of index-value pairs:

$$\begin{bmatrix} 0.0 & 0.0 & 0.0 & 0.1 & 0.0 \\ 0.0 & 0.2 & 0.3 & 0.0 & 0.0 \\ 0.0 & 0.4 & 0.5 & 0.6 & 0.7 \\ 0.8 & 0.0 & 0.0 & 0.9 & 0.0 \end{bmatrix}$$

(a) a matrix



(b) its sparse representation



(c) its flat representation

**Figure 1.** Representations of a sparse matrix

```
type sparse_vec = (int * float) parray
```

In this representation, we record only the non-zero entries of the vector, paired with their indices. A sparse matrix is a parallel array of sparse vectors:

```
type sparse_mat = sparse_vec parray
```

Figure 1(a) and (b) illustrate a matrix and its sparse representation.

The `sumP` operator is a parallel reduction that computes the sum of the elements of a parallel array of floats:

```
val sumP : float parray -> float
```

We define multiplying a dense vector by a sparse-matrix as the following high-level NDP function:

```
fun smvm (sm : sparse_mat, v : dense_vec) =
  [| sumP [| x * v!i | (i,x) in sv |] | sv in sm |]
```

The technical challenge is to find the meeting point between this elegant declarative code and an efficient implementation that can exploit powerful multicore architectures.

As one solution to this challenge, Blelloch and Sabot introduced the NESL language and the *flattening transformation* to compile nested-data-parallel programs for wide-vector parallel machines [2, 5]. NESL is an ML-like language with nested-data-parallel features. NESL provides only scalars and associated operators, sequences, simple datatypes, conditionals, let bindings, top-level function definitions, and a parallel *apply-to-each* construct. Blelloch’s flattening transformation then substantially changes both the data structures and code. This transformation is especially effective in treating irregular parallelism: certain operations like the parallel *segmented sum* operation, where sums are computed over an array of arrays of numbers, complete in the same number of steps regardless of the irregularity of the shapes of the segments.

```
datatype 'a rope
  = Leaf of 'a seq
  | Cat of 'a rope * 'a rope
```

```
datatype shape
  = Lf of int * int
  | Nd of shape rope
```

```
datatype 'a farray
  = FArray of 'a rope * shape
```

**Figure 2.** Datatypes `rope`, `shape` and `farray`.

The flattening transformation yields flattened arrays consisting of two components: one or more flat data vectors (more than one in the case of unzipped tuples), containing the elements of the nested array in left-to-right order, and one or more *segment-descriptors*. In NESL, a segment-descriptor is always a flat vector of integers, and a flattened array carries with it one segment-descriptor for each level of nesting. Figure 1(c) illustrates the flat representation of the example sparse matrix.

## 2.1 Challenges for full flattening

Full flattening, while a successful innovation, can produce inefficient code in common cases, including conditionals and certain regular nested parallel programs.

When full flattening is applied to conditionals, the resulting code generates many intermediate vectors, as we demonstrate here in a simplified example. Here is a function that replaces the zeroes in a vector with ones:

```
fun g (xs : dense_vec) =
  [| if x=0.0 then 1.0 else x | x in xs |]
```

The transformed code (which uses standard `split` and `combine` operations [2]) builds intermediate vectors to handle partitioning the data elements, performs the appropriate conditional work on each partition of the data, and then combines them:

```
fun g_full_flat (xs : dense_vec) = let
  val flags = [| x=0.0 | x in xs |]
  val (zs, nzs) = split (xs, flags)
  val zs' = [| 1.0 | z in zs |]
in
  combine (zs', nzs, flags)
end
```

This approach makes sense for SIMD architectures with large penalties for failing to keep the vector registers full, but on multicore machines those benefits can be overwhelmed by the large number of memory operations. There are many more extended examples of this kind of transformation in the literature [2, 13]. Our approach, by contrast, will transform `g` to operate on a flattened dense vector, but will not otherwise transform the code; to the point, it will not necessitate generating any intermediates.

With respect to regular nested parallel programs, problems similar to dense matrix multiplication suffer from a polynomial space increase under full flattening, which is a known serious problem [26] also due to excessive data copying. The present work describes a system designed to address these problems with traditional full flattening by avoiding extra data copies, both due to splitting of conditional branches and duplication of vectors within nested parallel applications.

## 3. Data-Only Flattening

We present data-only flattening in the context of a broader system for *hybrid flattening*. Hybrid flattening is a program representation that allows both flat and nested representations of parallel arrays and which has coercions for transforming between representations.

This representation allows flexibility in choosing when it is profitable to use a flat representation vs. a nested representation. Hybrid flattening does not itself express or embody a particular transformation policy.

Whereas Manticore without flattening compiles nested parallel arrays to nested ropes [6], with flattening it compiles nested parallel arrays to flattened arrays. Flattened arrays, like nested arrays as compiled by NESL, consist of two pieces: a flat data vector, and a value representing the structure of the nested array called a *shape tree*. By means of standard unzipping transformations, nested arrays of tuples are compiled to tuples of flattened arrays. In our implementation, flattened arrays are represented by the polymorphic `farray` datatype. To represent the flat data vector part of flattened arrays, we use ropes, exploiting Manticore’s existing rope infrastructure to compute in parallel with them. Figure 2 presents the PML datatype definitions for `rope`, `shape`, for shape trees, and `farray`. Note that ropes and shape trees are internal representations; the programmer uses them only indirectly.

Shape trees are our adaptation of segment descriptors in the NESL tradition. A *shape* is an  $n$ -ary tree whose leaves store integer pairs. Each leaf contains the starting index and the index of the element following the segment of data in an `farray`. The shape `Shape.Lf (i, i+n)` describes a length- $n$  segment starting at  $i$  and ending at the last position before  $i+n$ .<sup>1</sup>

A simple, flat parallel array of integers such as

```
[ | 1, 2, 3 | ]
```

has the following `farray` representation.<sup>2</sup>

```
FArray (Rope.Leaf [1,2,3], Shape.Lf (0,3))
```

The data in the original sequence appears here in the original order as a `Rope.Leaf` and the accompanying `shape` indicates that the flattened array’s only segment begins at position 0 and ends at position 2.

Nested parallel arrays are translated as follows. Consider the following nested array:

```
[ | [ | 1, 2 | ], [ | ], [ | 3, 4, 5, 6 | ] | ]
```

Its flattened array representation is the following:

```
FArray (Rope.Leaf [1,2,3,4,5,6],
        Shape.Nd [Shape.Lf (0,2),
                  Shape.Lf (2,2),
                  Shape.Lf (2,6)])
```

The flat data appears in order in a `Rope.Leaf`. The `shape` is a `Nd` with three leaves: this means that the parallel array consists of three subsequences. The leaves tell us that the first sequence begins at position 0 and ends at 1, the second sequence is empty at position 2, and the third sequence begins at position 2 and ends at 5. This representation naturally scales to any nesting depth.

Literal values can be flattened at transformation time directly, with the representation change incurring no runtime cost. In general, however, the compiler must cope with arbitrarily nested arrays whose dimensions are known only at runtime. In such cases, the compiler needs to arrange for flattening to take place at runtime. We handle this issue, by inserting coercion operators that perform runtime flattening. When nested arrays are transformed into flattened arrays, all operations applied to those array values must be correspondingly transformed. Our approach to this problem is to provide a core group of type-indexed families of array operators, each of which performs its operation at every array type in its family.

<sup>1</sup> In practice, this choice of bounds is a convenient convention that guards against common fencepost errors.

<sup>2</sup> For brevity, we present the sequence in the `Rope.Leaf` node and the sequence of leaves in the `Shape.Nd` with list syntax.

$\tau$	$::=$ $g$ <span style="float: right;"><i>ground types</i></span> $ $ $(\tau, \tau)$ <span style="float: right;"><i>pairs</i></span> $ $ $\tau \rightarrow \tau$ <span style="float: right;"><i>functions</i></span> $ $ $[\tau]$ <span style="float: right;"><i>parallel arrays</i></span> $ $ $\{\tau; \nu\}$ <span style="float: right;"><i>flattened parallel arrays</i></span>
$\nu$	$::=$ $lf$ <span style="float: right;"><i>structure of flat arrays</i></span> $ $ $nd(\nu)$ <span style="float: right;"><i>structure of nested arrays</i></span>
$g$	$::=$ $int$ $ $ $bool$

Figure 3. Flatland: types.

$t$	$::=$ $e^\tau$	
$e$	$::=$ $b$ <span style="float: right;"><i>ground terms</i></span> $ $ $x$ <span style="float: right;"><i>variables</i></span> $ $ <b>if</b> $t$ <b>then</b> $t$ <b>else</b> $t$ <span style="float: right;"><i>conditionals</i></span> $ $ <b>let</b> $x = t$ <b>in</b> $t$ <span style="float: right;"><i>let expressions</i></span> $ $ <b>fun</b> $f x^\tau = t$ <b>in</b> $t$ <span style="float: right;"><i>function expressions</i></span> $ $ $t \circ t$ <span style="float: right;"><i>function composition</i></span> $ $ $t t$ <span style="float: right;"><i>application</i></span> $ $ $(t, t)$ <span style="float: right;"><i>pairs</i></span> $ $ $\pi_i t$ <span style="float: right;"><i>projection (<math>i \in \{1, 2\}</math>)</i></span> $ $ $[t, \dots, t]$ <span style="float: right;"><i>arrays</i></span> $ $ $\{t, \dots, t; s\}$ <span style="float: right;"><i>flattened arrays</i></span> $ $ $t !_\tau t$ <span style="float: right;"><i>array subscript</i></span> $ $ <b>map</b> $_{(\tau, \tau, \tau, \tau)}(t, t)$ <span style="float: right;"><i>array map</i></span> $ $ <b>flt</b> $_{(\tau, \tau)}(t, t)$ <span style="float: right;"><i>array filter</i></span> $ $ <b>red</b> $_{(\tau, \tau)}(t, t, t)$ <span style="float: right;"><i>array reduction</i></span> $ $ $\tau \triangleright \tau$ <span style="float: right;"><i>type coercions</i></span>	
$s$	$::=$ <b>lf</b> $(t, t)$ <span style="float: right;"><i>leaves</i></span> $ $ <b>nd</b> $[s, \dots, s]$ <span style="float: right;"><i>nodes</i></span>	
$b$	$::=$ $true$ $ $ $false$ $ $ $0$ $ $ $1$ $ $ $\dots$ $ $ $not$ $ $ $+$ $ $ $\dots$	

Figure 4. Flatland: terms.

This group contains parallel array subscripting and parallel maps, filters, and reductions over parallel arrays. All operations on parallel arrays are either members of this core group or are built from members of this group. As such, transformation of the type-indexed operators matching the transformations of data structures is sufficient to preserve the program’s behavior.

## 4. Formalization

The system presented here, *Flatland*, consists of a model language and a variety of rewriting systems and judgments. Its model language is an explicitly-typed, monomorphic, strict, pure functional language. Flattened and non-flattened terms commingle in Flatland: there is no inherent distinction between source language and target language.

Figure 3 presents Flatland’s types, ranged over by the metavariable  $\tau$ . We use subscript indices ( $\tau_i$ ) and overbars ( $\bar{\tau}$ ) to distinguish types from one another. The type language consists of ground types *int* and *bool*, pairs, functions, parallel arrays, and flattened parallel arrays. Flattened-parallel types include *shape types* as subcomponents. Shape types, ranged over by the metavariable  $\nu$  (for “nesting”), record an array’s nesting depth.

$$\frac{\Gamma(x) = \tau}{\Gamma \vdash x^\tau \text{ ok}} \quad \frac{BE(b) = \tau \quad (BE = \text{basis env})}{\Gamma \vdash b^\tau \text{ ok}}$$

$$\frac{\Gamma \vdash e_1^{\text{bool}} \text{ ok} \quad \Gamma \vdash e_2^\tau \text{ ok} \quad \Gamma \vdash e_3^\tau \text{ ok}}{\Gamma \vdash (\text{if } e_1^{\text{bool}} \text{ then } e_2^\tau \text{ else } e_3^\tau)^\tau \text{ ok}}$$

$$\frac{\Gamma \vdash e_1^\tau \text{ ok} \quad \Gamma \vdash e_2^{\text{int}} \text{ ok} \quad \tau' = (! \tau)}{\Gamma \vdash (e_1^\tau !_\tau e_2^{\text{int}})^{\tau'} \text{ ok}}$$

$$\frac{\Gamma \vdash e_1^{\tau_1 \rightarrow \tau_2} \text{ ok} \quad \Gamma \vdash e_2^{\tau_3} \text{ ok} \quad \tau_3 \mathbf{A} \tau_1 \quad \tau_4 \mathbf{A} \tau_2}{\Gamma \vdash (\mathbf{map}_{(\tau_1, \tau_2, \tau_3, \tau_4)} (e_1^{\tau_1 \rightarrow \tau_2}, e_2^{\tau_3}))^{\tau_4} \text{ ok}}$$

$$\frac{\Gamma \vdash e_1^{\tau_1 \rightarrow \text{bool}} \text{ ok} \quad \Gamma \vdash e_2^{\tau_2} \text{ ok} \quad \tau_2 \mathbf{A} \tau_1}{\Gamma \vdash (\mathbf{filt}_{(\tau_1, \tau_2)} (e_1^{\tau_1 \rightarrow \text{bool}}, e_2^{\tau_2}))^{\tau_2} \text{ ok}}$$

$$\frac{\Gamma \vdash e_1^{(\tau_1, \tau_1) \rightarrow \tau_1} \text{ ok} \quad \Gamma \vdash e_2^{\tau_1} \text{ ok} \quad \Gamma \vdash e_3^{\tau_2} \text{ ok} \quad \tau_2 \mathbf{A} \tau_1}{\Gamma \vdash (\mathbf{red}_{(\tau_1, \tau_2)} (e_1^{(\tau_1, \tau_1) \rightarrow \tau_1}, e_2^{\tau_1}, e_3^{\tau_2}))^{\tau_1} \text{ ok}}$$

$$\frac{\vdash \tau \triangleright \bar{\tau} \text{ ok}}{\Gamma \vdash (\tau \triangleright \bar{\tau})^{\tau \rightarrow \bar{\tau}} \text{ ok}} \quad \dots$$

**Figure 5.** Well-typedness of terms (selected rules).

$$\overline{\vdash \tau \triangleright \tau \text{ ok}}$$

$$\overline{\vdash [\tau] \triangleright \{\tau; lf\} \text{ ok}} \quad \overline{\vdash \{\tau; lf\} \triangleright [\tau] \text{ ok}}$$

$$\overline{\vdash [(\tau_1, \tau_2)] \triangleright ([\tau_1], [\tau_2]) \text{ ok}} \quad \overline{\vdash ([\tau_1], [\tau_2]) \triangleright [(\tau_1, \tau_2)] \text{ ok}}$$

$$\frac{\vdash \tau_1 \triangleright \tau_2 \text{ ok} \quad \vdash \tau_2 \triangleright \tau_3 \text{ ok}}{\vdash \tau_1 \triangleright \tau_3 \text{ ok}} \quad \dots$$

**Figure 6.** Well-formedness of coercions (selected rules).

Figure 4 contains Flatland’s term language. Every term  $t$  includes an explicit type as a superscript. The metavariable  $b$  ranges over constants, and  $x$  ranges over variables. As in the type language, parallel-array types are written with square brackets, and flattened-array types with curly braces. Every flattened array includes a shape tree. Since the language is monomorphic, array operators — in Flatland these are subscript, map, filter, and reduce — do not have polymorphic implementations. Instead, we assume there is a type-indexed family for each one.

Flattening and unflattening operators are represented in our system by *coercions*. For the coercion that transforms values of type  $\tau_1$  into values of type  $\tau_2$ , we write  $\tau_1 \triangleright \tau_2$ . Except for identity coercions, these are potentially expensive representation-changing operations; during compilation, we try to eliminate as many of them as possible.

Well-formedness of coercions, so that nonsensical coercions like  $\text{int} \triangleright (\text{int}, \text{int})$  are rejected, is given in Figure 6. Well-typedness of terms is in Figure 5, and the rules for calculating shape types appear in Figure 7.<sup>3</sup> For array type  $\tau$ , we use the notation

<sup>3</sup>Due to space limitations, only selected rules are presented; the full rule sets appear in [25].

$$\frac{\Gamma \vdash e_1^{\text{int}} \text{ ok} \quad \Gamma \vdash e_2^{\text{int}} \text{ ok}}{\Gamma \vdash \mathbf{lf}(e_1^{\text{int}}, e_2^{\text{int}}) : \text{lf}}$$

$$\frac{\Gamma \vdash s_1 : \nu \quad \dots \quad \Gamma \vdash s_n : \nu}{\Gamma \vdash \mathbf{nd}[s_1, \dots, s_n] : \text{nd}(\nu)}$$

**Figure 7.** Calculation of shape types.

$$\begin{aligned} ![\tau] &= \tau \\ !\{\tau; lf\} &= \tau \\ !\{\tau; \text{nd}(\nu)\} &= \{\tau; \nu\} \\ !(\tau_1, \tau_2) &= (!\tau_1, !\tau_2) \end{aligned}$$

**Figure 8.** Array-type subscripting.

$$\overline{\tau \mathbf{L} \tau}$$

$$\frac{\tau_1 \mathbf{L} \tau_2}{\tau_2 \mathbf{L} \tau_1} \quad \frac{\tau_1 \mathbf{L} \tau_2 \quad \tau_2 \mathbf{L} \tau_3}{\tau_1 \mathbf{L} \tau_3}$$

$$\overline{[(\tau_1, \tau_2)] \mathbf{L} ([\tau_1], [\tau_2])}$$

$$\overline{\{\{\tau; \nu\}; lf\} \mathbf{L} \{\tau; \text{nd}(\nu)\}}$$

$$\overline{\{\{\tau; \nu\}; \text{nd}(\nu')\} \mathbf{L} \{\{\tau; \text{nd}(\nu)\}; \nu'\}}$$

$$\frac{\tau_1 \mathbf{L} \tau_2}{[\tau_1] \mathbf{L} [\tau_2]} \quad \frac{\tau_1 \mathbf{L} \tau_2}{[\tau_1] \mathbf{L} \{\tau_2; lf\}}$$

$$\frac{\tau_1 \mathbf{L} \tau_2}{\{\tau_1; \nu\} \mathbf{L} \{\tau_2; \nu\}}$$

$$\frac{\tau_1 \mathbf{L} \tau'_1 \quad \tau_2 \mathbf{L} \tau'_2}{(\tau_1, \tau_2) \mathbf{L} (\tau'_1, \tau'_2)} \quad \frac{\tau_1 \mathbf{L} \tau'_1 \quad \tau_2 \mathbf{L} \tau'_2}{\tau_1 \rightarrow \tau_2 \mathbf{L} \tau'_1 \rightarrow \tau'_2}$$

$$\frac{\tau_1 \mathbf{L} [\tau_2]}{\tau_1 \mathbf{A} \tau_2}$$

**Figure 9.** Definition of  $\tau_1 \mathbf{A} \tau_2$  and its auxiliary relation  $\tau_1 \mathbf{L} \tau_2$ .

(!  $\tau$ ) to mean the type of the element selected by subscript out of a value of type  $\tau$ . The return type of a particular subscript operator is calculated from its domain. Figure 8 gives the definition of (!  $\tau$ ). If (!  $\tau$ ) cannot be computed from these rules, it is undefined.

In order to understand the typing of the type-indexed operators, one needs to understand the relation  $\mathbf{A}$ . For a type  $\tau$ , we cannot definitively write  $[\tau]$  to mean “array of  $\tau$ ,” because  $\{\tau; lf\}$  is also an array of  $\tau$ , albeit in a different representation. Furthermore, if  $\tau$  is a pair type  $(\tau_1, \tau_2)$ , then  $[(\tau_1, \tau_2)]$  and  $([\tau_1], [\tau_2])$  (and more) are also arrays of  $\tau$ , and so on. Thus we appeal to the relation  $\tau_1 \mathbf{A} \tau_2$  for “ $\tau_1$  is an array of  $\tau_2$ .” For arrays of  $(\text{int}, \text{int})$ , for example, all

$$\begin{array}{c}
\frac{}{e^\tau \mapsto (\tau \triangleright \tau) e^\tau} \quad \frac{}{(\tau \triangleright \tau) e^\tau \mapsto e^\tau} \\
\frac{\vdash \tau_1 \triangleright \tau_2 \text{ ok} \quad \vdash \tau_2 \triangleright \tau_3 \text{ ok}}{\tau_1 \triangleright \tau_3 \mapsto (\tau_2 \triangleright \tau_3) \circ (\tau_1 \triangleright \tau_2)} \quad \frac{}{(\tau_2 \triangleright \tau_3) \circ (\tau_1 \triangleright \tau_2) \mapsto \tau_1 \triangleright \tau_3} \\
\frac{}{(\text{let } x = t_1 \text{ in } (\bar{\tau} \triangleright \tau) t_2)^\tau \mapsto ((\bar{\tau} \triangleright \tau) (\text{let } x = t_1 \text{ in } t_2))^\tau} \\
\frac{\vdash \bar{\tau} \triangleright \tau \text{ ok}}{x^\tau \mapsto ((\bar{\tau} \triangleright \tau) \bar{x}^\tau)^\tau} \quad \dots
\end{array}$$

**Figure 10.** Coercion distribution (selected rules), by which coercions are introduced, eliminated and otherwise manipulated.

of the following are established by **A**:

$$\begin{array}{l}
[(int, int)] \quad \mathbf{A} \quad (int, int) \\
([int], [int]) \quad \mathbf{A} \quad (int, int) \\
\{(int, int); lf\} \quad \mathbf{A} \quad (int, int) \\
\{(int; lf), \{int; lf\}\} \quad \mathbf{A} \quad (int, int)
\end{array}$$

These types correspond to arrays of integer pairs and all equivalent representations under well-formed coercions. The relation **A** is defined in Figure 9. **A** is defined by defining an auxiliary relation **L** that verifies that two types are both representations of the same scalar or array type. We use the name **L** for the relation because it tells us that one type is “on the same level as” another. All arrays of type  $\tau$  can be coerced between one another. Formally, if  $\tau_1 \mathbf{A} \tau$ , then  $\tau_2 \mathbf{A} \tau \Leftrightarrow \vdash \tau_1 \triangleright \tau_2 \text{ ok}$  (proven elsewhere [25].)

We specify array operators by writing their type indices in a subscript.  $!_\tau$  is the operator that selects elements from an array type  $\tau$ .  $\mathbf{map}_{(\tau_1, \tau_2, \tau_3, \tau_4)}$  takes two arguments, a function of type  $\tau_1 \rightarrow \tau_2$  and a term of array type  $\tau_3$ , and produces a term of array type  $\tau_4$ .  $\mathbf{filt}$  and  $\mathbf{red}$  carry similar type subscripts. The typing rules of  $\mathbf{map}$ ,  $\mathbf{filt}$  and  $\mathbf{red}$  all appeal to **A**.

To perform a flattening step, we insert a coercion. A standard step in flattening transformations is to unzip arrays of pairs of scalars—that is, to reshape an array of pairs into a pair of arrays. The coercion that unzips an array of integer pairs is written  $[(int, int)] \triangleright ([int], [int])$ . Its inverse coercion  $([int], [int]) \triangleright [(int, int)]$  is also part of the language of coercions. If it is ever the case that a pair of inverse coercions, like these two, are successively applied to a value, they may be rewritten to an identity coercion, which can in turn be removed from the program.

Well-formed programs are guaranteed to remain well-formed under any legal transformation in our rewriting system. The top-level type of the whole program remains fixed under transformation, but the types of the subexpressions within a program may change. The well-formedness guarantee is maintained by introducing inverse coercions for every coercion introduced into the program, to maintain the stability of the types in the program. For example, every time a coercion is applied to the value to which a variable  $x$  is bound, the inverse coercion is introduced at every use of coerced  $x$ . This ensures that no clients of  $x$  are put in a position to compute with a term of the wrong type, post-coercion. Type-indexed operators such as  $\mathbf{map}$  and  $\mathbf{filt}$  have the ability to absorb type coercions, since, for each one, there are many implementations from which to choose, and this helps reduce the number of type coercions in the transformed program.

To provide intuition about how this machinery works, here is an example of *Flatland* in action. Consider the following program:

$$\text{let } ns = [1, 2, 3] \text{ in } (ns !_{[int]} 0)$$

$$\begin{array}{c}
\frac{\vdash \tau \triangleright \bar{\tau} \text{ ok}}{(t_1 !_\tau t_2)^{(! \tau)} \mapsto ((! \bar{\tau} \triangleright ! \tau) (((\tau \triangleright \bar{\tau}) t_1) !_{\bar{\tau}} t_2)^{(! \bar{\tau})})^{(! \tau)}} \\
\frac{\bar{\tau}_2 \mathbf{A} \tau_1}{(\mathbf{filt}_{(\tau_1, \tau_2)} (t_1, t_2))^{\tau_2} \mapsto ((\bar{\tau}_2 \triangleright \tau_2) (\mathbf{filt}_{(\tau_1, \bar{\tau}_2)} (t_1, (\tau_2 \triangleright \bar{\tau}_2) t_2)))^{\tau_2}} \\
\frac{\bar{\tau}_3 \mathbf{A} \tau_1 \quad \bar{\tau}_4 \mathbf{A} \tau_2}{(\mathbf{map}_{(\tau_1, \tau_2, \tau_3, \tau_4)} (t_1, t_2))^{\tau_4} \mapsto ((\bar{\tau}_4 \triangleright \tau_4) (\mathbf{map}_{(\tau_1, \tau_2, \bar{\tau}_3, \bar{\tau}_4)} (t_1, (\tau_3 \triangleright \bar{\tau}_3) t_2)))^{\tau_4}} \\
\frac{\bar{\tau}_2 \mathbf{A} \tau_1}{(\mathbf{red}_{(\tau_1, \tau_2)} (t_1, t_2, t_3))^{\tau_1} \mapsto (\mathbf{red}_{(\tau_1, \bar{\tau}_2)} (t_1, t_2, (\tau_2 \triangleright \bar{\tau}_2) t_3))^{\tau_1}}
\end{array}$$

**Figure 11.** Coercion introductions with type-indexed operators. Note how the indexing types change to accommodate the coercions introduced.

$$\begin{array}{c}
\frac{\vdash \tau_1 \triangleright \bar{\tau}_1 \text{ ok} \quad \bar{x} \text{ fresh} \quad S = [x^{\tau_1} / ((\bar{\tau}_1 \triangleright \tau_1) \bar{x}^{\bar{\tau}_1})^{\tau_1}]}{(\text{let } x = t_1 \text{ in } t_2)^{\tau_0} \mapsto (\text{let } \bar{x} = (\tau_1 \triangleright \bar{\tau}_1) t_1 \text{ in } S t_2)^{\tau_0}} \\
\frac{\vdash \tau_0 \triangleright \bar{\tau}_0 \text{ ok} \quad \bar{f}, \bar{x} \text{ fresh} \quad S = [f^{\tau_0 \rightarrow \tau_1} / (\bar{f}^{\bar{\tau}_0 \rightarrow \tau_1} \circ (\tau_0 \triangleright \bar{\tau}_0))^{\tau_0 \rightarrow \tau_1}]}{S' = S[x^{\tau_0} / ((\bar{\tau}_0 \triangleright \tau_0) \bar{x}^{\bar{\tau}_0})^{\tau_0}]} \\
\frac{}{(\mathbf{fun } f x^{\tau_0} = t_1 \text{ in } t_2)^{\tau_2} \mapsto (\mathbf{fun } \bar{f} \bar{x}^{\bar{\tau}_0} = S' t_1 \text{ in } S t_2)^{\tau_2}}
\end{array}$$

**Figure 12.** Coercion propagation rules. Coercions introduced at binding sites necessitate substitutions accordingly.

$$\begin{array}{l}
\text{let } ns = [1, 2, 3] \text{ in } (ns !_{[int]} 0) \\
\mapsto \text{let } \bar{ns} = \downarrow [1, 2, 3] \text{ in } (\uparrow \bar{ns}) !_{[int]} 0 \\
\mapsto \text{let } \bar{ns} = \downarrow [1, 2, 3] \text{ in } (int \triangleright int) ((\downarrow (\uparrow \bar{ns})) !_{\{int; lf\}} 0) \\
\mapsto \text{let } \bar{ns} = \downarrow [1, 2, 3] \text{ in } (\downarrow (\uparrow \bar{ns})) !_{\{int; lf\}} 0 \\
\mapsto \text{let } \bar{ns} = \downarrow [1, 2, 3] \text{ in } ((\downarrow \circ \uparrow) \bar{ns}) !_{\{int; lf\}} 0 \\
\mapsto \text{let } \bar{ns} = \downarrow [1, 2, 3] \text{ in } ((\{int; lf\} \triangleright \{int; lf\}) \bar{ns}) !_{\{int; lf\}} 0 \\
\mapsto \text{let } \bar{ns} = \downarrow [1, 2, 3] \text{ in } (\bar{ns} !_{\{int; lf\}} 0)
\end{array}$$

**Figure 13.** Transforming a program by  $\mapsto$ .

For brevity, we let

$$\downarrow = [int] \triangleright \{int; lf\}$$

and

$$\uparrow = \{int; lf\} \triangleright [int]$$

Think of  $\downarrow$  as “flatten” and  $\uparrow$  as “unflatten.” The step-by-step transformation of this let-expression appears in Figure 13. The parallel array  $[1, 2, 3]$  is transformed to its flattened array equivalent,  $(\downarrow [1, 2, 3])$ , which evaluates to  $\{1, 2, 3; \mathbf{lf}(0, 3)\}$ . Furthermore, by rewriting, we exchange one type-indexed subscript operator for another, thereby eliminating coercion operations. The coercions  $\downarrow$  and  $\uparrow$  are inverse coercions. Note by the fourth rule in Figure 10 we have

$$([int] \triangleright \{int; lf\}) \circ (\{int; lf\} \triangleright [int]) \mapsto \{int; lf\} \triangleright \{int; lf\}$$

$$\begin{aligned}
\mathbf{F}[g] &= g \\
\mathbf{F}[\{g\}] &= \{g; lf\} \\
\mathbf{F}[\tau_1 \rightarrow \tau_2] &= \mathbf{F}[\tau_1] \rightarrow \mathbf{F}[\tau_2] \\
\mathbf{F}[(\tau_1, \tau_2)] &= (\mathbf{F}[\tau_1], \mathbf{F}[\tau_2]) \\
\mathbf{F}[[\tau_1 \rightarrow \tau_2]] &= \{\mathbf{F}[\tau_1 \rightarrow \tau_2]; lf\} \\
\mathbf{F}[[\tau_1, \tau_2]] &= (\mathbf{F}[[\tau_1]], \mathbf{F}[[\tau_2]]) \\
\mathbf{F}[[[\tau]]] &= \mathbf{N}[(\mathbf{F}[[\tau]])] \\
\mathbf{N}[(\tau_1, \tau_2)] &= (\mathbf{N}[\tau_1], \mathbf{N}[\tau_2]) \\
\mathbf{N}[\{\tau; \nu\}] &= \{\tau; nd(\nu)\}
\end{aligned}$$

**Figure 14.** Type flattening.

$$\frac{\{\} \vdash e^\tau \searrow (\bar{\tau} \triangleright \tau) \diamond \bar{e}^{\bar{\tau}}}{e^\tau \Downarrow ((\bar{\tau} \triangleright \tau) \bar{e}^{\bar{\tau}})^\tau}$$

**Figure 15.** Top-level data-only flattening.

so the composition of  $\downarrow$  and  $\uparrow$  is mutually annihilating. Post-transformation, the representation of the array bound to  $ns$  is coerced exactly once, to the differently-typed fresh variable  $\bar{ns}$ .

#### 4.1 Formal data-only flattening

We now present data-only flattening in the context of the *Flatland* system, as a transformation from a source language to a target language. A *source type* is a type that is neither a flattened-array type, nor contains any flattened-array types, generated by the grammar

$$\tau ::= g \mid \tau \rightarrow \tau \mid (\tau, \tau) \mid [\tau]$$

We define *source programs* as a term  $e^\tau$  for source type  $\tau$ , all of whose subterms have source types and contain no coercions. By this definition, we have made it illegal to write down flattened arrays anywhere in a source program. The transformation will introduce all flattened-array values and all coercions.

The target language is defined in terms of *flat types*. A type  $\tau$  is *flat* if

- it is a ground type  $g$ ,
- it is a function type  $\tau_1 \rightarrow \tau_2$  and  $\tau_1$  and  $\tau_2$  are flat,
- it is a pair type  $(\tau_1, \tau_2)$  and  $\tau_1$  and  $\tau_2$  are flat, or
- it is an array type  $\{\tau; \nu\}$  and  $\tau$  is a ground type or a flat function type.

If a type is not flat, we say it is *nonflat*. Note that source types and nonflat types are not the same. For example,  $[(int, bool)]$  is a source type, and the related type  $\{\{int; lf\}, \{bool; lf\}\}$  is a flat type. But the related type  $\{(int, bool); lf\}$  is neither a source type nor a flat type; it is disqualified as a source type since it is a flattened-array type and disqualified as a flat type since it includes a pair inside an array. A *target program* is an expression whose outermost type is a source type, yet all of whose subexpressions have flat types. The restriction on its outermost type is a consequence of the type-preservation property of *Flatland*'s rewriting. The restriction is enforced by the application of one last “unflattening” coercion to the transformed program at the top level. Within the program, all subexpressions are flattened.

Flattening of whole programs is written as a type-preserving relation  $\Downarrow$ . Figure 15 gives the sole judgment for  $\Downarrow$ , which immediately delegates its work to an auxiliary relation  $\searrow$ . Whole-program

$$\begin{aligned}
&\frac{}{\Delta \vdash b^\tau \searrow (\tau \triangleright \tau) \diamond b^\tau} \\
&\frac{\Delta(x^\tau) = \bar{x}^{\bar{\tau}}}{\Delta \vdash x^\tau \searrow (\bar{\tau} \triangleright \tau) \diamond \bar{x}^{\bar{\tau}}} \\
&\frac{\Delta \vdash e_1^{bool} \searrow (bool \triangleright bool) \diamond \bar{e}_1^{bool} \quad \Delta \vdash e_2^{\bar{\tau}} \searrow (\bar{\tau} \triangleright \tau) \diamond \bar{e}_2^{\bar{\tau}} \quad \Delta \vdash e_3^{\bar{\tau}} \searrow (\bar{\tau} \triangleright \tau) \diamond \bar{e}_3^{\bar{\tau}}}{\Delta \vdash (\text{if } e_1^{bool} \text{ then } e_2^{\bar{\tau}} \text{ else } e_3^{\bar{\tau}})^\tau \searrow (\bar{\tau} \triangleright \tau) \diamond (\text{if } \bar{e}_1^{bool} \text{ then } \bar{e}_2^{\bar{\tau}} \text{ else } \bar{e}_3^{\bar{\tau}})^\tau} \\
&\frac{\Delta \vdash e_1^{\tau_1} \searrow (\bar{\tau}_1 \triangleright \tau_1) \diamond \bar{e}_1^{\bar{\tau}_1} \quad \Delta' = \Delta[x^{\tau_1} \mapsto \bar{x}^{\bar{\tau}_1}], \bar{x} \text{ fresh} \quad \Delta' \vdash e_2^{\tau_2} \searrow (\bar{\tau}_2 \triangleright \tau_2) \diamond \bar{e}_2^{\bar{\tau}_2}}{\Delta \vdash (\text{let } x = e_1^{\tau_1} \text{ in } e_2^{\tau_2})^{\tau_2} \searrow (\bar{\tau}_2 \triangleright \tau_2) \diamond (\text{let } \bar{x} = \bar{e}_1^{\bar{\tau}_1} \text{ in } \bar{e}_2^{\bar{\tau}_2})^{\bar{\tau}_2}} \\
&\frac{\Delta' = \Delta[f^{\tau_0 \rightarrow \tau_1} \mapsto \bar{f}^{\bar{\tau}_0 \rightarrow \bar{\tau}_1}], \bar{f} \text{ fresh} \quad \Delta'' = \Delta'[x^{\tau_0} \mapsto \bar{x}^{\bar{\tau}_0}], \bar{x} \text{ fresh} \quad \Delta'' \vdash e_1^{\tau_1} \searrow (\bar{\tau}_1 \triangleright \tau_1) \diamond \bar{e}_1^{\bar{\tau}_1} \quad \Delta'' \vdash e_2^{\tau_2} \searrow (\bar{\tau}_2 \triangleright \tau_2) \diamond \bar{e}_2^{\bar{\tau}_2}}{\Delta \vdash (\text{fun } f \text{ } x^{\tau_0} = e_1^{\tau_1} \text{ in } e_2^{\tau_2})^{\tau_2} \searrow (\bar{\tau}_2 \triangleright \tau_2) \diamond (\text{fun } \bar{f} \bar{x}^{\bar{\tau}_0} = \bar{e}_1^{\bar{\tau}_1} \text{ in } \bar{e}_2^{\bar{\tau}_2})^{\bar{\tau}_2}} \\
&\frac{\Delta \vdash e_1^{\tau_1} \searrow (\bar{\tau}_1 \triangleright \tau_1) \diamond \bar{e}_1^{\bar{\tau}_1} \quad \Delta \vdash e_2^{\tau_2} \searrow (\bar{\tau}_2 \triangleright \tau_2) \diamond \bar{e}_2^{\bar{\tau}_2}}{\Delta \vdash (e_1^{\tau_1}, e_2^{\tau_2})^{(\tau_1, \tau_2)} \searrow ((\bar{\tau}_1, \bar{\tau}_2) \triangleright (\tau_1, \tau_2)) \diamond (\bar{e}_1^{\bar{\tau}_1}, \bar{e}_2^{\bar{\tau}_2})^{(\bar{\tau}_1, \bar{\tau}_2)}} \\
&\frac{\Delta \vdash e^{(\tau_1, \tau_2)} \searrow ((\bar{\tau}_1, \bar{\tau}_2) \triangleright (\tau_1, \tau_2)) \diamond \bar{e}^{(\bar{\tau}_1, \bar{\tau}_2)}}{\Delta \vdash (\pi_1 e^{(\tau_1, \tau_2)})^{\tau_1} \searrow (\bar{\tau}_1 \triangleright \tau_1) \diamond (\pi_1 \bar{e}^{(\bar{\tau}_1, \bar{\tau}_2)})^{\bar{\tau}_1}} \quad (\pi_2 \text{ sim.})
\end{aligned}$$

**Figure 16.** Data-only flattening, group 1.

flattening consists of transforming a program  $e^\tau$  of source type  $\tau$  to a program  $\bar{e}^{\bar{\tau}}$  of flat type  $\bar{\tau}$  and then coercing the transformed program to original type  $\tau$  at the top level. Note that in the cases where  $\tau$  is, for example, a ground type (see Figure 13), the outermost coercion is an identity coercion and has no effect.

The auxiliary relation of data-only flattening is given in Figures 16 and 17. The syntax of  $\searrow$  is as follows:

$$\Delta \vdash e^\tau \searrow (\bar{\tau} \triangleright \tau) \diamond \bar{e}^{\bar{\tau}}$$

$\Delta$  is a finite map from variable terms to variable terms; it is used to implement propagations through let-expressions and functions. On the right-hand side of the relation, a diamond ( $\diamond$ ) is used to construct a pair out of a coercion and a transformed expression. The relation produces an unflattening coercion, along with the transformed expression, for use in one of the following ways. If the expression transformed is the whole program, the  $\Downarrow$  relation applies the coercion to preserve the program's original type (as per the rule in Figure 15). If the expression is not the whole program, the accompanying coercion is used as a building block for further coercions as program transformation proceeds outward.

The important work in data-only flattening takes place at array terms; see the second rule in Figure 17 (which in turn appeals to

$$\begin{array}{c}
\frac{\Delta \vdash e_1^{\tau_1 \rightarrow \tau_2} \searrow (\bar{\tau}_1 \rightarrow \bar{\tau}_2 \triangleright \tau_1 \rightarrow \tau_2) \diamond \bar{e}_1^{\bar{\tau}_1 \rightarrow \bar{\tau}_2} \quad \Delta \vdash e_2^{\tau_1} \searrow (\bar{\tau}_1 \triangleright \tau_1) \diamond \bar{e}_2^{\bar{\tau}_1}}{\Delta \vdash (e_1^{\tau_1 \rightarrow \tau_2} e_2^{\tau_1})^{\tau_2} \searrow (\bar{\tau}_2 \triangleright \tau_2) \diamond (\bar{e}_1^{\bar{\tau}_1 \rightarrow \bar{\tau}_2} \bar{e}_2^{\bar{\tau}_1})^{\bar{\tau}_2}} \\
\\
\frac{\mathbf{F}[[\tau]] = \bar{\tau}}{\Delta \vdash e^{[\tau]} \searrow (\bar{\tau} \triangleright [\tau]) \diamond (([\tau] \triangleright \bar{\tau}) e^{[\tau]})^{\bar{\tau}}} \\
\\
\frac{\Delta \vdash e_1^{\tau_2 \rightarrow \tau_3} \searrow (\bar{\tau}_2 \rightarrow \bar{\tau}_3 \triangleright \tau_2 \rightarrow \tau_3) \diamond \bar{e}_1^{\bar{\tau}_2 \rightarrow \bar{\tau}_3} \quad \Delta \vdash e_2^{\tau_1 \rightarrow \tau_2} \searrow (\bar{\tau}_1 \rightarrow \bar{\tau}_2 \triangleright \tau_1 \rightarrow \tau_2) \diamond \bar{e}_2^{\bar{\tau}_1 \rightarrow \bar{\tau}_2}}{\Delta \vdash (e_1 \circ e_2)^{\tau_1 \rightarrow \tau_3} \searrow (\bar{\tau}_1 \rightarrow \bar{\tau}_3 \triangleright \tau_1 \rightarrow \tau_3) \diamond (\bar{e}_1 \circ \bar{e}_2)^{\bar{\tau}_1 \rightarrow \bar{\tau}_3}} \\
\\
\frac{\Delta \vdash e_1^{\tau} \searrow (\bar{\tau} \triangleright \tau) \diamond \bar{e}_1^{\bar{\tau}} \quad \Delta \vdash e_2^{int} \searrow (int \triangleright int) \diamond \bar{e}_2^{int}}{\Delta \vdash (e_1^{\tau} !_{\tau} e_2^{int})^{(! \tau)} \searrow ((! \bar{\tau}) \triangleright (! \tau)) \diamond (\bar{e}_1^{\bar{\tau}} !_{\bar{\tau}} \bar{e}_2^{int})^{(! \bar{\tau})}} \\
\\
\frac{\Delta \vdash e_1^{\tau_1 \rightarrow \tau_2} \searrow (\bar{\tau}_1 \rightarrow \bar{\tau}_2 \triangleright \tau_1 \rightarrow \tau_2) \diamond \bar{e}_1^{\bar{\tau}_1 \rightarrow \bar{\tau}_2} \quad \Delta \vdash e_2^{\tau_3} \searrow (\bar{\tau}_3 \triangleright \tau_3) \diamond \bar{e}_2^{\bar{\tau}_3}}{\mathbf{F}[[\tau_4]] = \bar{\tau}_4} \\
\frac{\Delta \vdash (\mathbf{map}_{(\tau_1, \tau_2, \tau_3, \tau_4)} (e_1^{\tau_1 \rightarrow \tau_2}, e_2^{\tau_3}))^{\tau_4} \searrow (\bar{\tau}_4 \triangleright \tau_4) \diamond (\mathbf{map}_{(\bar{\tau}_1, \bar{\tau}_2, \bar{\tau}_3, \bar{\tau}_4)} (\bar{e}_1^{\bar{\tau}_1 \rightarrow \bar{\tau}_2}, \bar{e}_2^{\bar{\tau}_3}))^{\bar{\tau}_4}}{\Delta \vdash e_1^{\tau_1 \rightarrow bool} \searrow (\bar{\tau}_1 \rightarrow bool \triangleright \tau_1 \rightarrow bool) \diamond \bar{e}_1^{\bar{\tau}_1 \rightarrow bool}} \\
\frac{\Delta \vdash e_2^{\tau_2} \searrow (\bar{\tau}_2 \triangleright \tau_2) \diamond \bar{e}_2^{\bar{\tau}_2}}{\Delta \vdash (\mathbf{filt}_{(\tau_1, \tau_2)} (e_1^{\tau_1 \rightarrow bool}, e_2^{\tau_2}))^{\tau_2} \searrow (\bar{\tau}_2 \triangleright \tau_2) \diamond (\mathbf{filt}_{(\bar{\tau}_1, \bar{\tau}_2)} (\bar{e}_1^{\bar{\tau}_1 \rightarrow bool}, \bar{e}_2^{\bar{\tau}_2}))^{\bar{\tau}_2}} \\
\\
\frac{\Delta \vdash e_1^{(\tau_1, \tau_1) \rightarrow \tau_1} \searrow ((\bar{\tau}_1, \bar{\tau}_1) \rightarrow \bar{\tau}_1 \triangleright (\tau_1, \tau_1) \rightarrow \tau_1) \diamond \bar{e}_1^{(\bar{\tau}_1, \bar{\tau}_1) \rightarrow \bar{\tau}_1} \quad \Delta \vdash e_2^{\tau_1} \searrow (\bar{\tau}_1 \triangleright \tau_1) \diamond \bar{e}_2^{\bar{\tau}_1} \quad \Delta \vdash e_3^{\tau_2} \searrow (\bar{\tau}_2 \triangleright \tau_2) \diamond \bar{e}_3^{\bar{\tau}_2}}{\Delta \vdash (\mathbf{red}_{(\tau_1, \tau_2)} (e_1^{(\tau_1, \tau_1) \rightarrow \tau_1}, e_2^{\tau_1}, e_3^{\tau_2}))^{\tau_2} \searrow (\bar{\tau}_2 \triangleright \tau_2) \diamond (\mathbf{red}_{(\bar{\tau}_1, \bar{\tau}_2)} (\bar{e}_1^{(\bar{\tau}_1, \bar{\tau}_1) \rightarrow \bar{\tau}_1}, \bar{e}_2^{\bar{\tau}_1}, \bar{e}_3^{\bar{\tau}_2}))^{\bar{\tau}_2}}
\end{array}$$

Figure 17. Data-only flattening, group 2.

type flattening in Figure 14.) This rule introduces coercions from parallel arrays to flattened arrays. The variable rule (second rule, Figure 16) substitutes typed variables for their flattened replacements per the map carried by  $\Delta$ . The array-operator rules exchange operators indexed by source type to operators indexed by the corresponding flat types. The other rules are administrative, recursively propagating transformations through expressions. There is exactly one rule for every distinct syntactic form, so the rules describes both a semantic specification and an algorithm.

## 5. Implementation

Data-only flattening in PML is accomplished in three successive phases: an *abstract flattening* phase, whereby abstract flattening operations — symbolic values that stand in for actual implementations — are inserted throughout the code; a *fusion* phase where canceling coercions (adjacent coercions that undo one another’s work) are eliminated; and a *concrete flattening* phase, where symbolic flattening operations are replaced by monomorphic code. Due to space limitations, the detailed operations of these phases is not presented here and may be found elsewhere [25].

## 5.1 Optimizations

Flattened PML programs are amenable to various optimizations that cannot be applied to non-flattened ones. In this section, we discuss several such optimizations, each of which is responsible, in part, for the performance improvements we report in our benchmark results.

*Monomorphization.* Monomorphization is an optimization whereby a polymorphic data structure containing uniformly-represented (*i.e.*, boxed) elements is transformed to a representation containing raw (unboxed) elements in their place. The flattening transformation, by virtue of unzipping arrays of tuples, exposes more opportunities for monomorphization than otherwise. In PML, arrays of double pairs, for example, become pairs of double arrays, which in turn become farrays, each containing a specialized rope of doubles as its flat data vector. Monomorphization is well known to be valuable even outside the context of nested data parallel compilation. MLton [18], an optimizing whole-program SML compiler, performs monomorphization to generate better-performing sequential code. PML stands to benefit from monomorphization even without flattening (PML currently does no monomorphization unless flattening is enabled), although it will never be the case that, without unzipping tuples, non-flattened PML will have as many opportunities to do it.

*Tab flattening.* Nested parallel comprehensions over ranges have *regular* structure: at each dimension, the length of every array is fixed a constant. The regularity of such structures can be exploited by the *tab flattening* optimization, which performs simple integer arithmetic operations to collapse multidimensional tabulations into linear ones.

Every one-dimensional parallel comprehension of scalars is trivially regular:

```
val xs = [| Double.fromInt i | i in [| 0 to 9 |] |]
```

The straightforward, and inefficient, implementation of this parallel comprehension is to translate it to a map over the parallel array containing the integers from 0 to 9.

```
PArray.map Double.fromInt [| 0 to 9 |]
```

This naïve translation entails building an ephemeral data structure that is immediately computed with and discarded. To save the cost associated with this intermediate structure, the compiler rewrites parallel comprehensions over ranges as tabulations:

```
PArray.tabulate (10, Double.fromInt)
```

Tabulating over integer intervals requires no intermediate data structures, and realizes a performance improvement over the build-and-map strategy outlined above.

Nested parallel comprehensions naturally give rise to nested tabulations. The computation of `xss` in this excerpt

```
val xss = [| [| (i*10)+j | j in [| 0 to 9 |] |] |
            | i in [| 0 to 9 |] |]
```

can be naturally expressed by a tabulate within a tabulate as follows:

```
PArray.tabulate (10, fn i =>
  PArray.tabulate (10, fn j =>
    (i*10) + j))
```

This translation is already better than using maps with ephemeral structures, but the shape of our flattened array representations allows us to use tab flattening to improve on nested tabulations. Recall our evaluation of `xss` results in an farray containing a flat data vector and a shape tree. We name the result `xssF` and sketch it as follows:

```

val xssF = FArray (Rope.Leaf [0, 1, 2, ..., 99],
                  Shape.Nd [Shape.Lf (0,10),
                          ...,
                          Shape.Lf (90,100)])

```

We can generate the flat data vector of `xssF` in one tabulation, over a single counter representing the total number of elements in the nested array, by performing the appropriate index arithmetic on the counter:

```

let fun f k = let
  val (i, j) = (k div 10, k mod 10)
  in
    (i*10) + j
  end
in
  PArray.tabulate (10*10, f)
end

```

The shape tree in rectangular cases has a simple regular structure as well, and be computed from the dimensions of a regular array in a straightforward way. Tab flattening operation scales to any number of dimensions for regular nested arrays.

*Segmented reductions.* NESL’s fast segmented operations are an important element of NESL’s ability to perform well on irregular nested data parallel programs, and an important one for PML to emulate. NESL’s segmented sum operation, for example, is able to compute the sums of a nested array of numbers in a fixed number of steps regardless of the irregularity of the array’s structure. This operation is critical to the performance of sparse-matrix/vector multiplication (see below). Here is an example of an irregular sum computation in PML:

```

let val nss = [| [| 1, 2 |],
                [| |],
                [| 3, 4, 5, 6 |] |]
in
  [| sum ns | ns in nss |]
end

```

If this parallel comprehension is rewritten such that the `sum` operation is simply mapped over the array-valued elements of `nss`, the irregularity of the structure of `nss`, if there is wide variation in the lengths of its elements, is bound to affect load balancing adversely. As such, the compiler replaces nested reductions like these

```

val ss = [| PArray.reduce add 0.0 ns | ns in nss |]
with “segmented reductions”
val ss = PArray.segreduce add 0.0 nss

```

The non-flattened implementation of `PArray.segreduce oper ident` is simply `PArray.map (PArray.reduce oper ident)`. As part of the flattening transformation, the non-flattened `segreduce` implementation is replaced by a faster `segreduce` written to exploit the shape of `farray` data structures.

## 6. Evaluation

These benchmarks compare the performance and scalability of Manticore’s current implementation of parallel arrays [1] against data-only flattening. The baseline performance is from the sequential version of the benchmark, which runs on a single processor and eliminates all overhead from the parallel language constructs and associated runtime features. The sequential version does not take advantage of any of the data transformations provided by data-only flattening, relying on the default polymorphic array representation.

These benchmarks were selected to show that we gain performance on code that executes over irregular data, as in SMVM, while preserving most of the performance of benchmarks over

regular data. Flattening happens at compile time in the following sense: nested arrays such as those in the main expressions of `mandelbrot` and `raytracer` are created flat in the first place, not constructed nested and then flattened.

### 6.1 Experimental method

Our benchmark machine is a Dell PowerEdge R815 server, outfitted with 48 cores and 128 GB physical memory. This machine runs x86\_64 Ubuntu Linux 10.04.2 LTS, kernel version 2.6.32-42. The 48 cores are provided by four 12 core AMD Opteron 6172 “Magny Cours” processors. Each core operates at 2.1 GHz and has 64 KB each of instruction and data L1 cache and 512 KB of L2 cache. There are two 6 MB L3 caches per processor, each of which is shared by six cores, for a total of 48 MB of L3 cache.

We ran each experiment configuration 30 times, and we report the average performance results in our graphs and tables. Times are reported in seconds.

### 6.2 Mandelbrot

We compute the Mandelbrot set by means of a function `elt` which consumes a pair of integers and produces an integer. The argument to `elt` represents a location in the complex plane. Its return value is the number of iterations required, according to the standard iterating Mandelbrot set membership test, for a given point to diverge outside the set (by having a modulus greater than 2). A point is a member of the Mandelbrot set if it fails to diverge before reaching a fixed upper limit of iterations (we use 1000). We execute this simple function in parallel over a  $2048 \times 2048$  range using the following PML code:

```

fun mandelbrot n = let
  val rng = [| 0 to (n-1) |]
  in
    [| [| elt (i, j) | j in rng |] | i in rng |]
  end

```

Figure 18(a) shows PML speedups, with and without flattening, against the sequential baseline. Due to the relatively small amount of computation at each element, the benefits of the data-only flattening transformation provide only a 5% speedup at 48 cores. These benefits come from the reduced amount of memory traffic when using a monomorphic array representation, avoiding an extra allocation per result element and associated garbage collector pressure.

### 6.3 Raytracer

Our ray tracing benchmark computes the image of a scene graph consisting of a group of overlapping spheres with transparency and reflection. The code is translated from a parallel program in the implicitly-parallel language Id90 [19]. It is a brute-force implementation and does not use any acceleration data structures. The Raytracer benchmark renders a  $2048 \times 2048$  image in parallel as a two-dimensional sequence.

Similar to the Mandelbrot benchmark, we write the body of the main function as a nested parallel comprehension:

```

fun raytracer n = let
  val ns = [| 0 to (n-1) |]
  in
    [| [| trace (i, j) | j in ns |] | i in ns |]
  end

```

In all versions of the program, the Manticore compiler transforms this nested parallel comprehension into a two-dimensional tabulation over fixed ranges, avoiding creation of an intermediate array of values to iterate over. In the data-only version of this program, that two-dimensional tabulation is instead performed as a



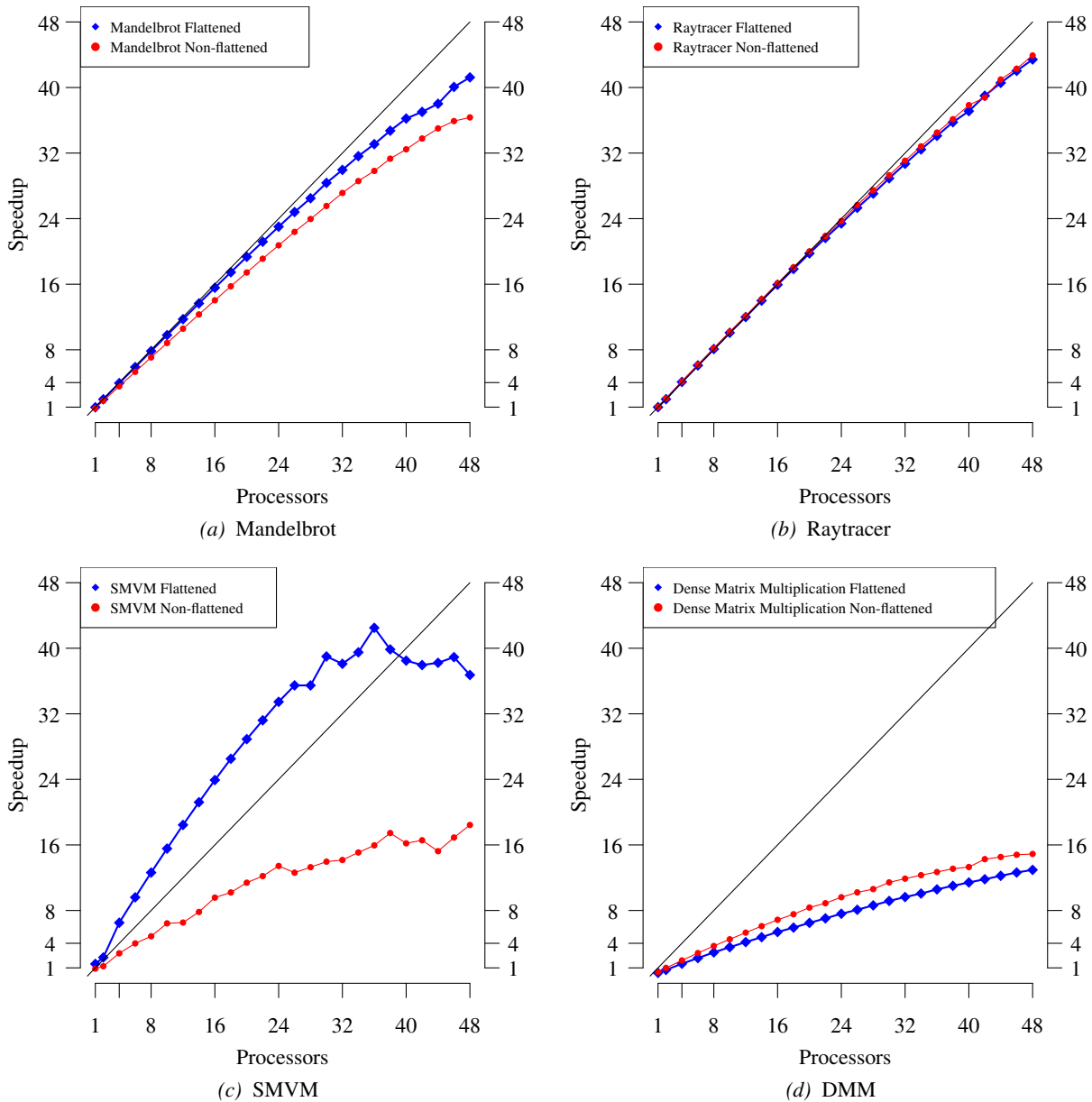


Figure 18. Comparison of flattened and non-flattened performance.

one-dimensional tabulation with adjusted indices. As is shown in Figure 18(b), this transformation slightly reduces the overhead required by the work-stealing scheduler to evenly balance the remaining work, since work on the one-dimensional data structure requires only splitting the index in half instead of the default finger-tree splitting required in the two-dimensional version [1].

#### 6.4 Sparse-Matrix Vector Multiplication

Among our benchmarks, sparse-matrix vector multiplication profits most from the data-only flattening transformation. Sparse-matrix vector multiplication is expressed concisely in the following PML code:

```

fun smvm (sm, v) = let
  fun add (a,b) = a+b
  in
    [| PArray.reduce add 0.0
      [| x*(v!i) | (i,x) in sv |]
     | sv in sm |]
  end

```

The inner parallel array expression computes the dot-product of the sparse matrix and the vector. The outer parallel array then computes the sum of each of those resulting dot-products. The summation is implemented on top of `PArray.reduce`, which allows the runtime to process the reduction in parallel.

In both versions of the benchmark, the Manticore compiler automatically optimizes the parallel reduction over multiple values

into a single segmented reduction. Segmented operators have been shown by Blelloch and others to result in more balanced chunks of parallel work, across a variety of platforms [3, 24]. All versions are transformed into the following PML code:

```

fun smvm (sm, v) = let
  val prods = products (sm, v)
  val sums = segsum prods
in
  sums
end

```

The flattened version of `smvm` uses monomorphic vectors for both the intermediate representation of the dot-products and the final result of the segmented reduction. Further, the implementation of segmented reduction over monomorphic vectors takes advantage of the layout, performing the segmented reduction with far less overhead than the sequential and non-flattened versions.

Figure 18(c) gives the speedups of PML over its sequential baseline, both with and without flattening. The sparse matrix is  $10,000 \times 10,000$ , with a random number of entries between 100 and 500 in each row. Flattened `smvm` is substantially faster than non-flattened `smvm` for all numbers of processors up to 48, and furthermore has much better performance with respect to the sequential baseline. The super-linear speedups are due to the relatively small amount of work performed on each element compared to the improvement due to the representation change. Above 36 processors, our performance improvement flattens due to having insufficient data to take advantage of the processors. Unfortunately, limitations in the Manticore runtime currently prevent us from further increasing the size of the data.

## 6.5 Dense Matrix Matrix Multiplication

The dense matrix multiplication (DMM) benchmark is a dense-matrix by dense-matrix multiplication in which each matrix is  $600 \times 600$ . As mentioned in Section 2, this benchmark has traditionally had extremely poor performance under flattening. As shown in Figure 18(d), our approach does still result in a slowdown in performance due to the creation of some intermediate arrays (resulting in a factor of 3 increase in memory usage). This penalty is roughly 13% and could be reduced through the introduction of additional fusion operations to avoid those intermediate arrays.

## 6.6 Conclusion

These benchmarks demonstrate that the data-only flattening transformation significantly improves a benchmark with irregular data (SMVM), does not experience the polynomial blowup typical to full flattening on DMM, and does not dramatically change the performance of other programs.

## 7. Related work

The incremental extension of NESL’s foundation to a more feature-rich platform has ultimately taken the form of Data Parallel Haskell [9]. Chakravarty *et al.* first present the language Nepal in 2001 [8], characterized as a version of Haskell including nested data parallelism. Nepal is succeeded by Chakravarty *et al.*’s Data Parallel Haskell [9], bringing together Nepal-style nested data parallelism with Haskell as implemented in the Glasgow Haskell Compiler (GHC) [12]. In 2008, Peyton Jones *et al.* give a thorough overview of the Data Parallel Haskell language and its compilation [22], including an updated account of how Data Parallel Haskell uses the flattening transformation in their implementation. The especially germane question is this: what happens after a system has compiled a Data Parallel Haskell program in the manner of NESL, and yet executes on an SMP rather than a SIMD

machine? Their answer to this technical problem is to adapt the `split` and `join` mechanisms originally presented in Keller’s dissertation [13] to implement NESL-style operations across the multiple processing elements on a multicore computer. Parallel computations are `split` across processing elements and subsequently `joined` on completion. To eliminate unnecessary synchronization points, GHC uses rewrite rules [21] to erase successive applications of `split` and `join` (per the identity equivalence rule originally given by Keller). (In addition to this, various advanced fusion techniques are employed to streamline the resulting post-flattened program, an overview of which is given in their paper.)

Though Data Parallel Haskell represents broad advances in NESL-style flattening in numerous ways, and although it has been adapted to run on multicore machines, its compilation strategy continues to reflect the SIMD orientation of its predecessors, though there has been recent work on vectorisation in DPH to avoid excessive flattening [15]. This work differs from that line of research by never performing the full vectorisation transformation on the code, though both approaches flatten nested data.

## Acknowledgments

This material is based upon work supported by the National Science Foundation under Grants CCF-0811389 and CCF-1010568, and upon work performed in part while John Reppy was serving at the National Science Foundation. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of these organizations or the U.S. Government.

## References

- [1] L. Bergstrom, M. Fluet, M. Rainey, J. Reppy, and A. Shaw. Lazy tree splitting. In *ICFP ’10*, pages 93–104, New York, NY, September 2010. ACM.
- [2] G. E. Blelloch. *Vector models for data-parallel computing*. MIT Press, Cambridge, MA, USA, 1990.
- [3] G. E. Blelloch. Prefix sums and their applications. Technical Report CMU-CS-90-190, School of Computer Science, Carnegie Mellon University, Nov. 1990.
- [4] G. E. Blelloch. Programming parallel algorithms. *CACM*, 39(3):85–97, Mar. 1996.
- [5] G. E. Blelloch and G. W. Sabot. Compiling collection-oriented languages onto massively parallel computers. *JPDC*, 8(2):119–134, 1990.
- [6] H.-J. Boehm, R. Atkinson, and M. Plass. Ropes: an alternative to strings. *SP&E*, 25(12):1315–1330, Dec. 1995. ISSN 0038-0644.
- [7] M. M. T. Chakravarty and G. Keller. More types for nested data parallel programming. In *ICFP ’00*, pages 94–105, New York, NY, Sept. 2000. ACM.
- [8] M. M. T. Chakravarty, G. Keller, R. Leshchinskiy, and W. Pfannenstiel. Nepal – nested data parallelism in Haskell. In *Euro-Par ’01*, volume 2150 of *LNCS*, pages 524–534, New York, NY, Aug. 2001. Springer-Verlag.
- [9] M. M. T. Chakravarty, R. Leshchinskiy, S. Peyton Jones, G. Keller, and S. Marlow. Data Parallel Haskell: A status report. In *DAMP ’07*, pages 10–18, New York, NY, Jan. 2007. ACM.
- [10] J. Dean and S. Ghemawat. MapReduce: Simplified data processing on large clusters. In *OSDI ’04*, pages 137–150, Berkeley, CA, Dec. 2004. USENIX Association.
- [11] M. Fluet, N. Ford, M. Rainey, J. Reppy, A. Shaw, and Y. Xiao. Status Report: The Manticore Project. In *ML ’07*, pages 15–24, New York, NY, Oct. 2007. ACM.
- [12] GHC. The Glasgow Haskell Compiler. Available from <http://www.haskell.org/ghc>.

- [13] G. Keller. *Transformation-based Implementation of Nested Data Parallelism for Distributed Memory Machines*. PhD thesis, Technische Universität Berlin, Berlin, Germany, 1999.
- [14] G. Keller and M. M. T. Chakravarty. Flattening trees. In *Euro-Par '98: Proceedings of the 4th International Euro-Par Conference on Parallel Processing*, pages 709–719, London, UK, 1998. Springer-Verlag.
- [15] G. Keller, M. M. T. Chakravarty, R. Leshchinskiy, B. Lippmeier, and S. Peyton Jones. Vectorisation Avoidance. In *HASKELL '12*, New York, NY, Sept. 2012. ACM. Forthcoming.
- [16] R. Leshchinskiy, M. M. T. Chakravarty, and G. Keller. Higher order flattening. In V. Alexandrov, D. van Albada, P. Sloot, and J. Dongarra, editors, *ICCS '06*, number 3992 in LNCS, pages 920–928, New York, NY, May 2006. Springer-Verlag.
- [17] R. Milner, M. Tofte, R. Harper, and D. MacQueen. *The Definition of Standard ML (Revised)*. The MIT Press, Cambridge, MA, 1997.
- [18] MLton. The MLton Standard ML compiler. Available at <http://mlton.org>.
- [19] R. S. Nikhil. *ID Language Reference Manual*. Laboratory for Computer Science, MIT, Cambridge, MA, July 1991.
- [20] D. W. Palmer, J. F. Prins, and S. Westfold. Work-efficient nested data-parallelism. In *FoMPP5*, pages 186–193, Los Alamitos, CA, 1995. IEEE Computer Society Press.
- [21] S. Peyton Jones, A. Tolmach, and T. Hoare. Playing by the rules: Rewriting as a practical optimization technique in GHC. In *Proceedings of the 2001 Haskell Workshop*, pages 203–233, Sept. 2001.
- [22] S. Peyton Jones, R. Leshchinskiy, G. Keller, and M. M. T. Chakravarty. Harnessing the Multicores: Nested Data Parallelism in Haskell. In *APLAS '08*, pages 138–138, New York, NY, Dec. 2008. Springer-Verlag.
- [23] M. Rainey. *Effective Scheduling Techniques for High-Level Parallel Programming Languages*. PhD thesis, University of Chicago, Aug. 2010. Available from <http://manticore.cs.uchicago.edu>.
- [24] S. Sengupta, M. Harris, Y. Zhang, and J. D. Owens. Scan primitives for GPU computing. In *GH '07*, pages 97–106, Aire-la-Ville, Switzerland, Aug. 2007. Eurographics Association.
- [25] A. Shaw. *Implementation techniques for nested-data-parallel languages*. PhD thesis, University of Chicago, Aug. 2011. Available from <http://manticore.cs.uchicago.edu>.
- [26] D. Spoonhower, G. E. Blelloch, R. Harper, and P. B. Gibbons. Space profiling for parallel functional programs. In *ICFP '08*, pages 253–264, New York, NY, Sept. 2008. ACM.