

Development and Verification of Parallel Algorithms in the Data Field Model

Björn Lisper¹ and Jonas Holmerin²

¹ Dept. of Computer Engineering, Mälardalen University, P.O. Box 883, S-721 23 Västerås, SWEDEN, bjorn.lisper@mdh.se

² Department of Numerical Analysis and Computing Science, Royal Institute of Technology, S-100 44 Stockholm, SWEDEN, joho@nada.kth.se

Abstract. Data fields are partial functions provided with explicit domain information. They provide a very general, formal model for collections of data. Algorithms computing data collections can be described in this formalism at various levels of abstraction: in particular, explicit data distributions are easy to model. Parallel versions of algorithms can then be formally verified against algorithm specifications in the model. Functions computing data fields can be directly programmed in the language Data Field Haskell. In this paper we give a brief introduction to the data field model. We then describe Data Field Haskell and make a small case study of how an algorithm and a parallel version of it both can be specified in the language. We then verify the correctness of the parallel version in the data field model.

1 Introduction

Many computing applications require indexed data structures. In many applications the indexing capability provides an important part of the model. On the other hand, the memory space of a parallel computer architecture is also indexed. Thus, indexed structures can describe data close to both the problem domain and architectural domain.

The canonical indexed data structure is the array. However, in particular when dealing with sparse, distributed applications, other, more dynamic indexed data structures are needed. Their low-level representations can be intricate. Thus, it is hard to design algorithms using them from scratch, in particular since the algorithms themselves may be complex, and to port algorithms to different architectures. Therefore, it is desirable to develop such algorithms on a high level first, where implementation details are hidden.

The data field model¹ is a formal model where indexed data structures are modeled as partial functions supplied with explicit information about their domains. A small formal language can be used to define data fields mathematically, and known proof techniques can be used to prove properties about data fields defined within the language. In particular, it is interesting to prove correspondences between different data fields, since this can be used to prove the correctness of refinement steps taken in a design process where abstract specifications are successively refined into distributed implementations.

¹ “Field” should be understood as in physics, as an entity which is a function of space and possibly time.

However, it is also helpful to validate such steps experimentally, in particular since they sometimes are only approximately correct and can yield, for instance, different numerical properties. The language Data Field Haskell can be used for this. It provides an instance of data fields, which is general and flexible enough to specify data fields both on a very high level of abstraction and on a level with explicitly parallel data distributions.

The rest of this paper is organized as follows. Sect. 2 gives a brief description of the underlying data field model. Sect. 3 describes Data Field Haskell. Sect. 4 provides an example of how a parallel program can be developed from a specification in Data Field Haskell. In Sect. 5 we verify formally that the parallel program indeed is correct w.r.t. the original specification. Sect. 6 provides an account for related work. In Sect. 7, finally, the story is wrapped up. The limited space does not allow a complete description of Data Field Haskell here – see [10] for the details.

Various versions of the data field model have been described elsewhere [7, 12–14]. The contribution of this paper is a description of a concrete implementation, an example of how it can be used in parallel program design, and a demonstration how the model can be used for formal verification of parallel programs.

2 The Data Field Model

2.1 Partial Functions

The concept of data fields is based on the more abstract model of indexed data structures as functions with finite domain [7, 12]. An array with range $[1..n]$, for instance, can be seen as a function from $\{1, \dots, n\}$, but we could also model “irregular” indexed structures as functions with non-contiguous, possibly non-numerical domains. Calls to a function outside its domain return an error value “*”, with algebraic properties similar to the divergent element \perp .

We define partial functions in a variation of the metalanguage for continuous functions in [23], extended with the constant *. Within this language, we can now define most types of collection-oriented operations [20] as higher order functions operating on partial functions [7, 14]. The language, however, lacks operations which extract the domain of a function or any information pertaining to it, like the size (“information operations” as defined in [20]). The lack of domain information makes it impossible to define collective operations such as *reduction*, since they need to know which elements to include. A partial function f over an enumerable cpo models a data structure which contains elements $f(x)$ exactly for those fully defined² x where $f(x) \neq *$. We thus define, for all partial functions f , its *domain* $dom(f) = \{x \mid f(x) \neq * \text{ and } x \text{ fully defined}\}$. $dom(f)$ is not computable in general, but we can still use it in abstract algorithm specifications.

Now assume, for any enumerable cpo C under consideration and any finite $C' \subseteq C$, with fully defined elements only, a bijection $i_{C'}: \{0, \dots, |C'| \perp 1\} \rightarrow C'$. Equipped with these functions, and the set operation $|\cdot|$ (cardinality), we define:

$$size(f) = |dom(f)| \quad enum_f = i_{dom(f)}$$

² x is fully defined, or *maximal*, in the cpo $\langle C, \sqsubseteq \rangle$ if $\forall y. x \sqsubseteq y \implies x = y$. Unless the data structure has a lazy lookup function, a lookup is successful only if the index is fully defined.

We can now define reduction over nonempty, partial functions with finite *size* w.r.t. a binary operation \oplus :

$$red(\oplus, f) = red'(\oplus, f, size(f) \perp 1)$$

where

$$\begin{aligned} red'(\oplus, f, 0) &= f(enum_f(0)) \\ red'(\oplus, f, n + 1) &= red'(\oplus, f, n) \oplus f(enum_f(n + 1)) \end{aligned}$$

If \oplus has an identity element e , then we define $red(\oplus, \lambda x.*) = e$ for the empty partial function $\lambda x.*$. Note the analogy with reduction over lists in the Bird-Meertens formalism [1]: we indeed have

$$red(\oplus, f) = \oplus/[f(enum_f(0)), \dots, f(enum_f(size(f) \perp 1))]$$

We will develop this analogy further in Sect. 5, where we will use the partial function model to verify the correctness of a parallelised algorithm.

2.2 Data Fields

The partial function model is simple and powerful, but in order to actually implement it explicit information about the function domains is needed. Thus we consider entities (f, b) – the data fields – where f is a function and b is a *bound*, a set representation which bounds the domain of the corresponding function. We require that the following operations are defined for bounds:

- For each bound b an interpretation $\llbracket b \rrbracket$ as a predicate (which in turn defines a set $\{\{b\}\} = \{x \mid \llbracket b \rrbracket(x) = true \text{ and } x \text{ fully defined}\}$).
- A predicate classifying each bound as either *finite* or *infinite*, depending on whether its set is surely finite or possibly infinite.
- For every bound b defining a finite set, $size(b)$ which yields the size of the set and $enum(b)$ which is a function enumerating its elements.
- Binary operations \sqcap, \sqcup on bounds (“intersection”, “union”).
- The bounds *all* and *nothing* representing the universal and empty set, respectively.

These operations are chosen to support the operations on partial functions which require the domain of the functions. They must fulfil certain properties, see [14]. We define *data field application* (“lookup”) viz.:

$$\begin{aligned} (f, b) ! x &= if(\llbracket b \rrbracket(x), f(x), *) \quad x \text{ fully defined} \\ (f, b) ! x &= \perp \quad \text{otherwise} \end{aligned}$$

Every data field d then defines a partial function $\llbracket d \rrbracket_{\alpha \rightarrow \beta}$ through $\llbracket d \rrbracket_{\alpha \rightarrow \beta}(x) = d ! x$ for all x . It immediately follows that $dom(\llbracket (f, b) \rrbracket_{\alpha \rightarrow \beta}) \subseteq \{\{b\}\}$. Thus, the bound of a data field always bounds the domain of its function. The bound does not have to be tight – we can have $(f, b) ! x = *$ even for some $x \in \{\{b\}\}$.

The theory of data fields also defines *φ -abstraction*, a syntax for convenient definition of data fields which parallels λ -abstraction for functions. The idea is that $\llbracket \varphi x.t \rrbracket_{\alpha \rightarrow \beta}$ should be the same function as $\lambda x.t$, except possibly for some pathological cases. Different definitions of $\varphi x.t$ fulfilling this are possible: some are given in [14]. For these definitions it is possible to prove the following result, which is stated more precisely in [14]:

Theorem 1. *If the bound of $\varphi x.t$ is fully defined, then $\varphi x.t ! y = \lambda x.t y$ for all fully defined y .*

Theorem 1 holds under the condition that \perp and $*$ are identified. Even if not, it still holds for total functions (i.e., which return a fully defined result given fully defined arguments). For data fields which define total functions, the result thus means that formal reasoning about data fields can be carried out in the more abstract model of partial functions. We will use this in Sect. 5.

3 Data Field Haskell

Data Field Haskell is a Haskell dialect where the arrays have been replaced by an instance of data fields, a variation of the *sparse/dense arrays* of [13, 14]. Our implementation of Data Field Haskell is based on the NHC compiler [18] for Haskell v. 1.3. The implementation is sequential and we have not implemented any advanced optimizations. In the following the reader is assumed to have a working knowledge of Haskell.

Data Field Haskell has data types `Datafield a b` for datafields and `Bounds a` for the corresponding bounds. Table 1 lists the most important functions to form and operate on data fields and bounds, and gives some algebraic laws for them.

Data Field Haskell has a rich variety of finite and infinite bounds: *dense bounds*, i.e., traditional array bounds, *sparse finite bounds*, which represent general finite sets, *predicate bounds*, which are classified as infinite, *universe* which represents the universal set, and *empty* which represents the empty set. *Product bounds* represent Cartesian products and generalise multidimensional array bounds. Some two-dimensional bounds are illustrated in Fig. 1.

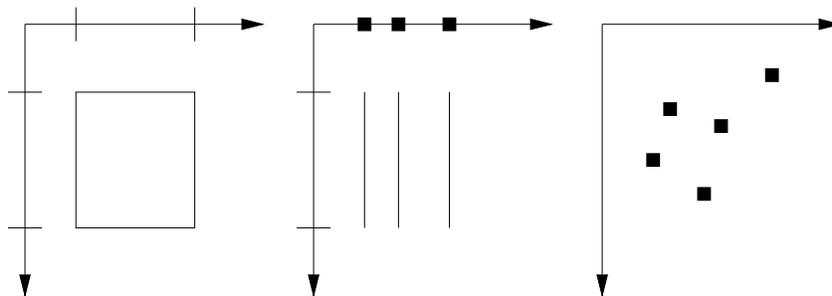


Fig. 1. Some two-dimensional bounds: two product bounds, and a sparse two-dimensional bound.

The most basic operations on bounds are `join` and `meet`, which provide \sqcup and \sqcap . The kind of bound computed depends on the arguments as shown in Table 1.

3.1 Forall-abstraction

Data Field Haskell provides φ -abstraction, with the following syntax (described in the metasyntax of the Haskell report [17]):

$$\text{forall } \text{apat}_1 \dots \text{apat}_n \rightarrow \text{exp}$$

```

datafield    defines data field from function and bound
!           data field indexing
outofBounds the out-of-bounds error value *
outofBounds the out-of-bounds error value *
isoutofBounds test for outofBounds
<:>        forms dense bound from pair of index tuples
sparse      forms finite sparse bound from list
predicate   forms predicate bound from predicate
<*>       forms product bound from two bounds
prod_n      forms n-dimensional product bound
join, meet   $\sqcup$  and  $\sqcap$  on bounds
<\>       explicit restriction of data field with bound
finite      tests bound for finiteness
enumerate   list of elements of the set defined by a finite bound in a given order
lowerBound  first enumerated element in finite bound
upperBound  last enumerated element in finite bound
size        the number of elements in a finite bound
inBounds    checks if an element belongs to the set defined by a bound
foldlDf     folds (reduces) finite data field w.r.t. binary operation
hstrictTab  evaluates all elements in finite data field fully

```

```

(bx1 <*> by1) 'op' (bx2 <*> by2) = (bx1 'op' bx2) <*> (by1 'op' by2)
  universe 'meet' x = x 'meet' universe = x
  empty 'join' x = x 'join' empty = x
  universe 'join' x = universe
  x 'join' universe = universe, x fully defined
  empty 'meet' x = empty
  x 'meet' empty = empty, x fully defined
bounds (datafield f b) = b
(datafield f b1) <\> b2 = datafield f (b2 'meet' b1)

```

meet	<i>E U S D P ×</i>	join	<i>E U S D P ×</i>
<i>E</i>	<i>E E E E E E</i>	<i>E</i>	<i>E U S D P ×</i>
<i>U</i>	<i>U S D P ×</i>	<i>U</i>	<i>U U U U U</i>
<i>S</i>	<i>S S S S</i>	<i>S</i>	<i>S S P S/P</i>
<i>D</i>	<i>D S ×</i>	<i>D</i>	<i>D P ×</i>
<i>P</i>	<i>P P</i>	<i>P</i>	<i>P P</i>
×	×	×	×

Table 1. Selected operations on data fields and bounds, some algebraic laws, and tables for result “types” of `join` and `meet` as a function of the argument “types”. “op” is any of `join`, `meet`. *E* = empty, *U* = universe, *S* = sparse, *D* = dense, *P* = predicate, × = product bound. “*S/P*” in the table for `join` means that the result is `sparse` if the product bound is finite, and a `predicate` otherwise. (Backquotes around binary functions, like in ‘`meet`’, turn them into infix operators in Haskell.)

Thus, the syntax is analogous to λ -abstraction in Haskell. The semantics is, with some minor deviations, the same as for φ -abstraction in [14]. It can be thought of as an implicitly parallel, functional `forall` statement where first `b` is computed and then, if needed, `\x -> t` is computed for all `x` in `b` in any (possibly parallel) order. The rules for computing `b` generalise existing array language principles for computing implicitly given bounds in the following cases: elementwise applied operations, selection of row/column from multidimensional array, and translation with constant offset. The exact rules are given in [10]: some representative examples are shown in Table 2.

```

        bounds (forall x -> 17) = universe
    bounds (forall x -> outofBounds) = empty
        bounds (forall x -> a!x + b!x) =
            (bounds a) 'meet' (bounds b)
    bounds (forall (x,y) -> a!x + b!y) =
            (bounds a) <*> (bounds b)
        bounds (forall x -> if a!x then b!x else c!x) =
    (bounds a) 'meet' ((bounds b) 'join' (bounds c))
            bounds (forall x -> a!(1,x)) = b2
            bounds (forall x -> a!(x,x)) = b1 'meet' b2
        bounds (forall (x,y) -> a!(y,x)) = b2 <*> b1
    bounds (forall x -> (datafield f (1 <:> 5))!(x+1)) = 0 <:> 4

```

Table 2. Some examples of bounds for `forall`-abstraction. Here, `bounds a = b1 <*> b2`.

4 An Example

We now give a simple example how a prototype for parallel code, which could be taken as a starting point for a real implementation, can be stepwise derived in Data Field Haskell from a specification. Our example is Jacobi’s algorithm, which is a classical, iterative method to solve linear systems of equations. Iterative methods of this kind are interesting to apply in particular to very large sparse systems, where direct methods can be prohibitively expensive. Thus, interesting issues in the development of parallel algorithms for Jacobi’s method are the ability to handle sparse structures at a suitable level of abstraction, and also how to develop methods for load balancing which are both reasonably efficient and give a reasonably good result. Notably, it should be possible to easily plug in different methods for load balancing since it is common to have matrices with a certain structure, for which specialised methods for load balancing may exist.

Jacobi’s method solves the equation $Ax = b$, where A is an $n \times n$ -matrix and x, b are n -vectors, by computing iterates $x^{(k)}$ according to

$$x_i^{(k+1)} = \frac{b_i - \sum_{j=1, j \neq i}^n a_{ij} x_j^{(k)}}{a_{ii}}, \quad i = 1, \dots, n \quad (1)$$

until two successive iterates satisfy some convergence criterion, for instance that $\|x^{(k+1)} - x^{(k)}\| < \epsilon$. Here, $\|\cdot\|$ is a vector norm: in this example, we use the maximum norm $\|y\| = \max_{i=1}^n |y_i|$.

```

sumDf = foldlDf (+) 0

iter f x eps conv =
  let xnew = f x
  in if (conv xnew x eps) then xnew else iter f xnew eps conv

jacobi a b eps x = iter (jacobi_iter a b) x eps conv

jacobi_iter a b x =
  forall i -> ((b!i) - sumDf ((forall j -> a!(i,j)*(x!j))
    <\> predicate (\j -> j /= i)))/a!(i,i)

conv x y eps = maxnorm (forall i -> (x!i - y!i)) < eps

maxnorm x = foldl1Df max (forall i -> abs (x!i))

```

Fig. 2. Data Field Haskell specification of Jacobi's algorithm.

A specification in Data Field Haskell of Jacobi's algorithm, according to (1), is given in Fig. 2. Here, `sumDf` sums all the elements in a finite data field (`foldlDf` is analogous to the Haskell function `foldl` on lists), and `iter` is a tail-recursive skeleton for general iterative algorithms with a convergence check. `jacobi` iterates `jacobi_iter` repeatedly until convergence according to `conv` is reached. Note the close resemblance of the definition of `jacobi_iter` and (1) – this exemplifies the ability to program with `forall`-abstraction in a style very close to mathematical notation. `conv`, finally, uses `foldl1Df` which is analogous to `foldl1` over lists.

Note that the bounds of `a` do not show up explicitly in the code: especially, the sums over `forall j -> a!(i,j)*(x!j)` have their limits given implicitly by the bounds of this data field expression, for the different values of `i`. These bounds are the `meet` of the bounds for `forall j -> a!(i,j)` and `x`. If `bounds a` equals $(1,1) \langle : \rangle (n,n)$ (representing a dense $n \times n$ -matrix) and if `bounds x` = $1 \langle : \rangle n$ (a dense n -vector) then the `meet` is $1 \langle : \rangle n$, and the sum will thus be performed over all $j \in \{1, \dots, n\}$. If, however, `a` has a sparse bound, then `forall j -> a!(i,j)` will have a sparse bound as well³ and the `meet` will also be sparse: thus, only the defined elements will be summed. The algorithm specification in Fig. 2 therefore defines both a dense and a sparse algorithm depending on the kind of data field!

We now manually refine this specification into a parallel algorithm. The idea is to distribute the rows of `a` and the elements of `b` according to a partitioning of the interval $\{1, \dots, n\}$ as indicated in Fig. 3. The inner products can then be kept local, which keeps the communication down, and this design decision leads to a parallel algorithm where each processor computes a segment of `xnew`. These are then assembled and broadcast for the new parallel iteration. The convergence test can be performed as a parallel reduction. If `a` is sparse, then it can pay off to perform a nontrivial load balancing, which however can be made statically since the structure of `a` does not

³ In the current implementation of Data Field Haskell this bound becomes wider than necessary. We expect to rectify this in later releases.

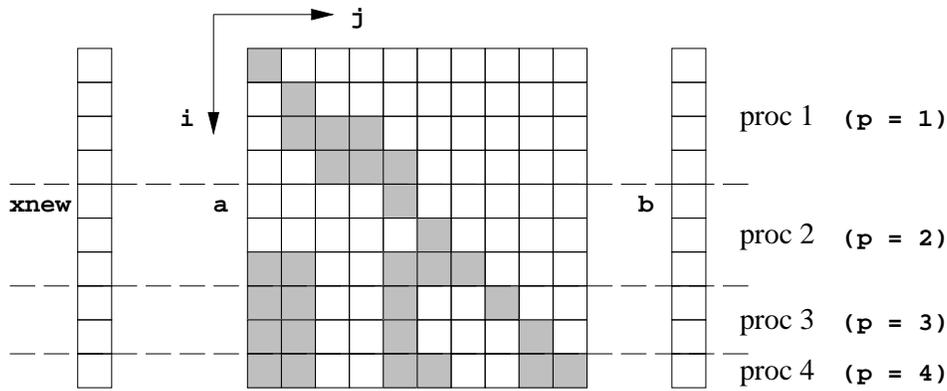


Fig. 3. Distribution of data in the parallel Jacobi algorithm.

change during the course of the algorithm. Our algorithm takes a static load balancing as a parameter.

```

p_jacobi ad bd eps xd = iter (p_jacobi_iter ad bd) xd eps p_conv

p_jacobi_iter ad bd xd = forall p -> jacobi_iter (ad!p) (bd!p) (prUnion xd)

p_conv xd yd eps =
  maxnorm (forall p ->
    maxnorm ((forall i -> ((xd!p)!i - (yd!p)!i))
      <\>(bounds (xd!p) 'meet' bounds (yd!p)))) < eps

prunion d1 d2 =
  forall x -> if (inBounds x (bounds d1)) then d1!x else d2!x

emptyfield = datafield (\x -> outofBounds) empty

prUnion = foldlDf prunion emptyfield

```

Fig. 4. Data Field Haskell program for parallelised simulation.

The distributed program is shown in Fig. 4. It makes heavy use of nested data fields, where the outermost data fields have indices (p in Fig. 4) which can be seen as processor addresses. `p_jacobi` uses the same basic iteration pattern as the original specification, but now the entities are distributed: `ad` is a nested datafield of two-dimensional data fields, and `bd` and `xd` are nested one-dimensional data fields. Each iteration step is performed by `p_jacobi_iter`: it assembles a global array from the distributed array `xd`, which is then used by all processors for computing local updates of its segments of `xd` using its parts of `ad` and `bd`. `p_conv`, finally, computes the maximum norm of `xd` in parallel, by first computing local maximal differences for each segment of old and new `xd`, and then computing a global maximum of these.

The global array is assembled through `prUnion`, which “flattens” a distributed one-dimensional data field of one-dimensional data fields. It is a reduction over the binary operation `prunion` on data fields. This operation computes the “union” of two

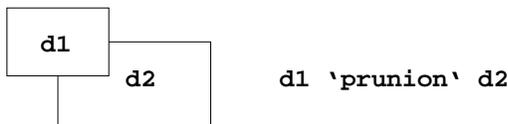


Fig. 5. prunion.

data fields giving priority to its first argument, see Fig. 5. If the data fields have disjoint bounds, then the operation acts somewhat like concatenation of lists. This similarity is not a mere coincidence, see Sect. 5.

The parallel algorithm above works on distributed data fields. The functions in Fig. 6 compute distributions. `divide a k lo hi` computes a data field of bounds which represents a k -partitioning of the interval $\{lo, \dots, hi\}$. It partitions the rows of `a` as to make the numbers of elements of `a` in each partition as equal as possible, thus obtaining a static load balancing of the work. It is a divide-and-conquer algorithm which itself is parallel. In each step, a point `midpoint` is selected between `lo` and `hi` such that the number of elements in rows `lo...midpoint` is as close to the number of elements in rows `midpoint + 1...hi` as possible. The same procedure is then recursively applied until k intervals have been formed.

`run_p_jacobi` uses the computed distribution to distribute initial data and start the parallel algorithm. It uses the help functions `matrix_curry`, which turns a two-dimensional data field into a data field of data fields, and `distr_1` and `distr_2` which convert one- and two-dimensional data fields, respectively, into nested, partitioned data fields according to a partitioning specified by a data field of bounds.

5 A Formal Verification of the Correctness of the Distributed Jacobi Algorithm

We now consider how to verify formally that `p_jacobi` implements `jacobi` correctly. Following Theorem 1, we use the more abstract model of partial functions, with \perp and $*$ identified, and rather than verifying directly that the implementation relation holds between the data fields we will define partial function versions `jacobi` and `p_jacobi` of them and prove that `p_jacobi` implements `jacobi`.

Definition 1. $f: A \rightarrow B$ is implemented by $g: C \rightarrow (D \rightarrow B)$ under (ϕ_0, ϕ_1) if $\forall x \in A. f(x) = g(\phi_0 x)(\phi_1 x)$.

Definition 1 is tailored for implementations by curried functions, which model nested data fields. Now consider the case when f and g are recursively defined, say we have a simultaneous recursive definition $(f, g) = F(f, g)$ where F is a continuous function. The condition on f and g in Definition 1 is an *inclusive* (or *admissible*) predicate, which means we can use *Scott's fixed-point induction* to prove it [23].

Proposition 1. Define (f_*, g_*) as the least fixed-point to $(f, g) = F(f, g)$, and let $(f_0, g_0) = (\perp, \perp)$, $(f_i, g_i) = F(f_{i-1}, g_{i-1})$ for $i > 0$. Then f_* is implemented by g_* under (ϕ_0, ϕ_1) if $\forall i > 0. (\forall x. (f_{i-1}(x) = g_{i-1}(\phi_0 x)(\phi_1 x)) \implies \forall x. (f_i(x) = g_i(\phi_0 x)(\phi_1 x)))$.

Proof. Scott's fixed-point induction principle for inclusive predicates P states that for any nondecreasing chain $\{d_i\}_{i=0}^{\infty}$, $P(\bigsqcup_{i=0}^{\infty} d_i)$ holds if: $P(d_0)$ holds, and $\forall i >$

```

divide a 1 lo hi = datafield (\i -> (lo<:>hi)) (1<:>1)

divide a k lo hi =
  let mid = midsearch a lo hi
  in (divide a (k 'div' 2) lo mid) 'concat'
     (divide a (k - (k 'div' 2)) (mid+1) hi)

midsearch a lo hi =
  midsearch_loc a lo hi lo hi (midpoint lo hi)
  where
    midsearch_loc a lo hi loloc hiloc mid
      | mid == hiloc || mid == loloc = adjust a lo hi mid
      | slice_size a lo mid > slice_size a (mid+1) hi
        = midsearch_loc a lo hi loloc mid (midpoint loloc mid)
      | otherwise
        = midsearch_loc a lo hi mid hiloc (midpoint mid hiloc)

midpoint m n = m + ((n-m) 'div' 2)

h_slice a b = a <\> (b<*>universe)

slice_size a lo hi = size (bounds (h_slice a (lo<:>hi)))

adjust a lo hi mid =
  if abs ((slice_size a lo (mid+1))-(slice_size a (mid+2) hi)) <
    abs ((slice_size a lo mid)-(slice_size a (mid+1) hi))
  then mid+1 else mid

concat d1 d2 = prunion d1 (forall x -> d2!(x-upperBound (bounds d1)))

run_p_jacobi a b eps x no_procs =
  let b_ac = bounds (matrix_curry a)
      dst  = divide a no_procs (lowerBound b_ac) (upperBound b_ac)
  in p_jacobi (distr_2 a dst) (distr_1 b dst) eps (distr_1 x dst)

matrix_curry a = forall i -> forall j -> a!(i,j)

distr_1 a b = forall p -> a <\> (b!p)

distr_2 a b = forall p -> h_slice a (b!p)

```

Fig. 6. Code for computing a distribution of data fields, for distributing initial data, and for starting the parallel algorithm.

$0.P(d_{i-1}) \implies P(d_i)$. In our case, $\{(f_i, g_i)\}_{i=0}^\infty$ is a nondecreasing chain where $(f_0, g_0) = (\perp, \perp)$ and $\bigsqcup_{i=0}^\infty (f_i, g_i) = (f_*, g_*)$. We trivially have $\forall x. \perp(x) = \perp(\phi_0 x)(\phi_1 x)$ which means the base case holds. What remains to prove is the induction step, which is the stated condition in the proposition.

We need some more results in order to prove that *p_jacobi* implements *jacobi*. *p_jacobi* uses data fields which are partitioned and distributed over processors. At certain points, they are reassembled from their parts using *prUnion*. The correctness of this reassembly depends on whether the distribution really is a true partitioning. We define partial function versions of *prunion* and *prUnion*:

$$\begin{aligned} f \text{ pr } g &= \lambda x. \text{if } (x \in \text{dom}(f), f(x), g(x)) \\ \text{Pr}(f) &= \text{red}(\text{pr}, f) \end{aligned}$$

It is easy to see that *pr* is associative. Partitioning and distributing a data field corresponds to a particular kind of implementation:

Definition 2. $g: C \rightarrow (A \rightarrow B)$ is a partitioning of $f: A \rightarrow B$ under ϕ if f is implemented by g under (ϕ, id) and if $\phi(\text{dom}(f)) = \text{dom}(g)$.

Proposition 2. If g is a partitioning of f under some ϕ , then $\text{Pr}(g) = f$.

Proof. In general, it holds that $(\text{Pr}(g))(i) = g(p)(i)$ for some $p \in \text{dom}(g)$. Furthermore, $\text{dom}(\text{Pr}(g)) = \bigcup_{p \in \text{dom}(g)} \text{dom}(g(p))$. We have $g(\phi(i))(i) = f(i)$ for all $i \in \text{dom}(f)$, and $\text{dom}(g(\phi(i))) = \{j \mid \phi(j) = \phi(i)\}$: thus, $\{\text{dom}(g(\phi(i))) \mid i \in \text{dom}(f)\}$ is a (set) partitioning of $\text{dom}(\text{Pr}(g))$. Therefore, $(\text{Pr}(g))(i) = g(\phi(i))(i)$ for all $i \in \text{dom}(f)$, and thus $(\text{Pr}(g))(i) = f(i)$. We also have, for all $i \in \text{dom}(f)$, $\text{dom}(g(\phi(i))) = \{j \mid j \in \text{dom}(f) \wedge \phi(j) = \phi(i)\}$. Since $\phi(\text{dom}(f)) = \text{dom}(g)$ every $p \in \text{dom}(g)$ is an image of some $i \in \text{dom}(f)$: thus, $\text{dom}(\text{Pr}(g)) = \bigcup_{i \in \text{dom}(f)} \text{dom}(g(\phi(i))) = \text{dom}(f)$. $\text{Pr}(g) = f$ follows.

Second, we need to prove that under certain conditions the result of a nested reduction, over a distributed, partitioned data field, equals a straight reduction over the original data field. The latter result corresponds to the BMF “reduce promotion” law $\oplus/\oplus/l = \oplus/+/l$ for nested lists l and associative operators \oplus [1]. This law holds since concatenation of lists imposes a certain ordering on the elements of respective lists. $\text{Pr}(f)$ will behave as $+/l$ if $\text{dom}(f(p)) \cap \text{dom}(f(q)) = \emptyset$ for $p \neq q$, and if p being enumerated before q by enum_f implies that each $x \in \text{dom}(f(p))$ is enumerated before each $y \in \text{dom}(f(q))$ by $\text{enum}_{\text{Pr}(f)}$.

Definition 3. $f: A \rightarrow (B \rightarrow C)$ has compatible enumerations if:

- $\text{dom}(f(x)) \cap \text{dom}(f(y)) = \emptyset$ when $x \neq y$,
- $\text{enum}_{\text{Pr}(f)}(n) = \text{enum}_{f(\text{enum}_f(m(n)))}(n \perp s(n))$, where $s(n), m(n)$ are given by $s(n) = \sum_{i < m(n)} \text{size}(f(\text{enum}_f(i)))$ and $s(n) \leq n < s(n) + \text{size}(f(\text{enum}_f(m(n))))$.

See Fig. 7. We can now state a law for curried partial functions which corresponds to the BMF law for nested lists above:

Theorem 2. If f has compatible enumerations and \oplus is associative, then $\text{red}(\oplus, \lambda x. \text{red}(\oplus, f(x))) = \text{red}(\oplus, \text{Pr}(f))$.

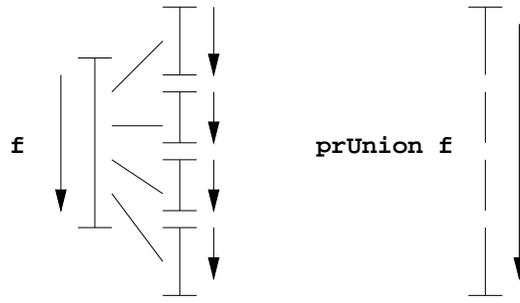


Fig. 7. Compatible enumerations.

Proof. (Sketch.) We first prove $red(\oplus, fprg) = red(\oplus, f) \oplus red(\oplus, g)$ provided $dom(f) \cap dom(g) = \emptyset$, $enum_{fprg}(n) = enum_f(n)$ when $0 \leq n < size(f)$, and $enum_{fprg}(n) = enum_g(n \perp size(f))$ for $size(f) \leq n < size(fprg)$. Then the result is proved by induction over the enumeration of f , using the definitions of Pr and red .

Finally two simple results: the first was originally stated in [7], the second follows directly from the definition of red :

Proposition 3. *If the n -ary function f is total and strict in all arguments, and if all g_i are total, $1 \leq i \leq n$, then $dom(\lambda x.f(g_1(x), \dots, g_n(x))) = dom(g_1) \cap \dots \cap dom(g_n)$.*

Proposition 4. *If f and \oplus are total, then $red(\oplus, f)$ yields a fully defined cpo element.*

$$f \setminus b = \lambda x.if(b(x), f(x), *) \quad (\text{explicit restriction})$$

$$sum(f) = red(+, f)$$

$$iter(f, x, \epsilon, conv) = let\ x_{new} = f(x)\ in\ if(conv(x_{new}, x, \epsilon), x_{new}, iter(f, x_{new}, \epsilon, conv))$$

$$jacobi(a, b, \epsilon, x) = iter(\lambda x.jacobi_iter(a, b, x), x, \epsilon, conv)$$

$$jacobi_iter(a, b, x) = \lambda i.(b(i) - sum(\lambda j.a(i, j) * x(j)) \setminus \lambda j.j \neq i) / a(i, i)$$

$$conv(x, y, \epsilon) = red(max, \lambda i. |x(i) - y(i)|) < \epsilon$$

$$p_jacobi(a', b', \epsilon, x') = iter(\lambda x'.p_jacobi_iter(a', b', x'), x', \epsilon, p_conv)$$

$$p_jacobi_iter(a', b', x') = \lambda p.jacobi_iter(a'(p), b'(p), Pr(x'))$$

$$p_conv(x', y', \epsilon) = red(max, \lambda p.red(max, \lambda i. |x'(p)(i) - y'(p)(i)|)) < \epsilon$$

Fig. 8. Definitions of $jacobi$ and p_jacobi .

Fig. 8 gives the definitions of $jacobi$ and p_jacobi : cf. Figs. 2 and 4. We now prove that p_jacobi does provide a partitioning of $jacobi$, given that the inputs to the functions do.

Theorem 3. *If: a, b, x are total; $\text{dom}(b) = \text{dom}(x) \subseteq \{1, \dots, n\}$; $a(i, i) \neq 0, *$ for $1 \leq i \leq n$; a' is a partitioning of a under ϕ ; b' and x' are partitionings of b, x , respectively, under ϕ' , where $\forall j. \phi(i, j) = \phi'(i)$; and x' has compatible enumerations, then $p\text{-jacobi}(a', b', \epsilon, x')$ is a partitioning of $\text{jacobi}(a, b, \epsilon, x)$ under ϕ' .*

Proof. Since jacobi and $p\text{-jacobi}$ are nonrecursively defined in terms of iter we can perform the fixed-point induction over iter . We have

$$\begin{aligned} \text{iter}_n(f, x, \epsilon, \text{conv}) = \\ \text{let } x_{\text{new}} = f(x) \text{ in if}(\text{conv}(x_{\text{new}}, x, \epsilon), x_{\text{new}}, \text{iter}_{n-1}(f, x_{\text{new}}, \epsilon, \text{conv})) \end{aligned}$$

Let $x_{\text{new}} = \text{jacobi_iter}(a, b, x)$ and $x'_{\text{new}} = p\text{-jacobi_iter}(a', b', x')$. From the definition of iter_n above and the definitions of jacobi and $p\text{-jacobi}$, we see that the result holds if: $\text{conv}(x_{\text{new}}, x, \epsilon) = p\text{-conv}(x'_{\text{new}}, x', \epsilon)$, x'_{new} is a partitioning of x_{new} under ϕ' , and $\text{iter}_{n-1}(\lambda x'. p\text{-jacobi_iter}(a', b', x'), x'_{\text{new}}, \epsilon, p\text{-conv})$ is a partitioning of $\text{iter}_{n-1}(\lambda x. \text{jacobi_iter}(a, b, x), x_{\text{new}}, \epsilon, \text{conv})$ under ϕ' . In the calls to iter_{n-1} , only x_{new} and x'_{new} have changed. Thus, by the fixed-point induction hypothesis, the partitioning relation between these calls holds if: x_{new} is total, x'_{new} is a partitioning of x_{new} under ϕ' , and if x'_{new} has compatible enumerations. Note that this induction hypothesis is more complex than the one in Proposition 1: this is since we also need to prove the result about conv and $p\text{-conv}$, which are reductions. (However, it is easy to see that the base case still holds, so it suffices to prove the inductive step.) We now prove the required conditions in turn:

- x_{new} total: we have $x_{\text{new}} = \text{jacobi_iter}(a, b, x) = \lambda i. (b(i) \perp \text{sum}(\lambda j. a(i, j) * x(j)) \setminus \lambda j. j \neq i) / a(i, i)$. Since all operations involved are total, and x, a , and b are total, this function is total.
- x'_{new} a partitioning of x_{new} under ϕ' : we first prove that $\text{dom}(x_{\text{new}}) = \text{dom}(x)$ and $\text{dom}(x'_{\text{new}}) = \text{dom}(x')$, which then implies $\phi'(\text{dom}(x_{\text{new}})) = \text{dom}(x'_{\text{new}})$.
 - $\text{dom}(x_{\text{new}}) = \text{dom}(x)$: we have $\text{dom}(x) = \text{dom}(b)$. We now prove $\text{dom}(b) = \text{dom}(x_{\text{new}})$ which yields the result. Reconsider the definition of x_{new} above. We have $\text{dom}(b) \subseteq \{1, \dots, n\}$, $\text{dom}(\lambda i. \text{sum}(\dots)) = U$ (the universal set), and $\text{dom}(\lambda i. a(i, i)) \supseteq \{1, \dots, n\}$. Since the arithmetic operations are strict and total (also division when $a(i, i) \neq 0$), Proposition 3 yields $\text{dom}(x_{\text{new}}) = \text{dom}(b) \cap U \cap \text{dom}(\lambda i. a(i, i)) = \text{dom}(b)$.
 - $\text{dom}(x'_{\text{new}}) = \text{dom}(x')$: since both b' and x' are partitionings of b and x under ϕ' , we have $\text{dom}(b') = \phi'(\text{dom}(b)) = \phi'(\text{dom}(x)) = \text{dom}(x')$. Furthermore, $x'_{\text{new}} = \lambda p. \text{jacobi_iter}(a'(p), b'(p), Pr(x'))$. Similar to above, we can check that $\text{jacobi_iter}(a'(p), b'(p), Pr(x'))$ is defined exactly when $b'(p)$ is defined, which implies $\text{dom}(x'_{\text{new}}) = \text{dom}(b')$. Thus, $\text{dom}(x'_{\text{new}}) = \text{dom}(x')$.

It remains to show that $x_{\text{new}}(i) = x'_{\text{new}}(\phi'(i))(i)$ for all i . We have $x'_{\text{new}} = \lambda p. \text{jacobi_iter}(a'(p), b'(p), Pr(x'))$. By Proposition 2 we have $Pr(x') = x$, thus $x'_{\text{new}} = \lambda p. \text{jacobi_iter}(a'(p), b'(p), x)$. For any i we now have (since $\forall j. \phi(i, j) = \phi'(i)$):

$$\begin{aligned} x'_{\text{new}}(\phi'(i))(i) &= \text{jacobi_iter}(a'(\phi(i, j)), b'(\phi'(i)), x)(i) \\ &= \dots b'(\phi'(i))(i) \dots a'(\phi(i, j))(i, j) \dots x(j) \dots a'(\phi(i, i))(i, i) \\ &= \dots b(i) \dots a(i, j) \dots x(j) \dots a(i, i) \\ &= \text{jacobi_iter}(a, b, x)(i) \\ &= x_{\text{new}}(i) \end{aligned}$$

- x'_{new} has compatible enumerations: by the above, it holds that $dom(x'_{new}) = dom(x')$ and $\forall p. dom(x'_{new}(p)) = dom(x'(p))$ (the latter follows from the fact that x'_{new} partitions x_{new} and x' partitions x under the same ϕ' , with $dom(x_{new}) = dom(x)$). Since x' has compatible enumerations, and compatible enumerations is a property only of the domains, the result follows.

It remains to prove $conv(x_{new}, x, \epsilon) = p_conv(x'_{new}, x', \epsilon)$. Let $y = \lambda i. |x_{new}(i) \perp x(i)|$ and $y' = \lambda p. \lambda i. |x'_{new}(p)(i) \perp x'(p)(i)|$. If $red(\max, \lambda p. red(\max, y'(p))) = red(\max, y)$, then the equality holds. It is easy to see (through Proposition 3), that y' , $y'(p)$, and y have the same domains as x' , $x'(p)$ and x , respectively. Thus, y' has compatible enumerations, and by Theorem 2 $red(\max, \lambda p. red(\max, y'(p))) = red(\max, Pr(y'))$ since \max is associative. It is also easy to see that y' is a partitioning of y under ϕ' . Proposition 2 then yields $Pr(y') = y$ which proves the equality.

6 Related Work

An excellent survey of collection-oriented languages up to around 1990 is found in [20]. The computation of bounds for `forall`-abstraction yields the implicit intersection rule of FIDIL [19]. The arrays in FIDIL resemble data fields also in other respects, for instance they can have a wider variety of shapes than traditional array bounds.

Examples of functional data parallel and array languages are Connection Machine Lisp [22], Id [3], Sisal [4], NESL [2], Data Parallel Haskell [9], and pH [15]. These languages are intended for direct parallel implementation whereas Data Field Haskell targets the specification phase. Haskell itself [17] has also been suggested for data parallel programming [16]. FISH [11] is an imperative array language, which shares some features with Data Field Haskell such as advanced polymorphism. It is, however restricted to regular arrays and certain recursion patterns which enables the generation of good code but makes it less suitable for specification of sparse or dynamic algorithms. A survey of the research in parallel functional programming is found in [8].

High-level specification and formal derivation of parallel programs has been considered in the BSP model [21] and in the Bird-Meertens Formalism [6]. Data Field Haskell could serve as a flexible format supporting these methods. Transformation of algorithm specifications in ML into programs for SIMD computers was considered in [5].

7 Conclusions and Further Research

We have presented the data field model and Data Field Haskell, a Haskell dialect which provides an instance of data fields. A possible application of Data Field Haskell is as a rapid prototyping tool for the early specification phase of parallel algorithms. We exemplified with the specification and initial development of a simple iterative linear equation solver, which could provide a starting point for the development of real production code. The initial specification was very close to the mathematical formulation, and the parallel version, once conceived, was straightforward to express with nested data fields. We made a formal verification that the parallel algorithm implements the specification through a partitioning of the data. A formal verification of this kind would probably be very hard to do for parallel programs written in current production languages.

The current implementation of Data Field Haskell could be improved in many ways, and it could certainly be given a parallel implementation. An elegant way to introduce explicit parallelism would be to use the type system of Haskell: “processor id” types could be introduced, with a 1-1-mapping to a given parallel architecture, and a parallelising compiler could then directly distribute the computation of data fields indexed with such types. An *efficient* parallel implementation would, however, require an efficient implementation of Haskell itself! Therefore, at least in a short perspective, we believe more in Data Field Haskell as a vehicle for specification of successively more refined parallel algorithms, which can serve as blueprints for actual parallel implementations.

We derived an interesting counterpart to a BMF law for reduction over curried partial functions. This line of investigation could certainly be pursued further, and also be extended to data fields. A difference is that lists always are ordered, while data fields only need to be ordered for certain operations. For instance, it is possible to formulate a version of Theorem 2 where the binary operation also is assumed commutative: then the requirement of compatible enumerations for the partial function can be dropped. There is clearly a rich variety of laws to explore for the subsequent use in program transformations and algorithm synthesis. The “added value” of such a theory over BMF would primarily be a greater ability to work with less ordered structures than lists, like indexed multidimensional structures. We believe this is beneficial for parallel algorithm design in application areas like scientific computing, where indexed structures prevail and the mapping phase to processors is important.

References

1. R. S. Bird. An introduction to the theory of lists. In M. Broy, editor, *Logic of Programming and Calculi of Discrete Design*, volume 36 of *NATO ASI Series F*, pages 3–42. Springer-Verlag, Berlin, 1987.
2. G. E. Blelloch. Programming parallel algorithms. *Comm. ACM*, 39(3), Mar. 1996.
3. K. Ekanadham. A perspective on Id. In B. K. Szymanski, editor, *Parallel Functional Languages and Compilers*, chapter 6, pages 197–253. Addison-Wesley, 1991.
4. J. T. Feo, D. C. Cann, and R. R. Oldehoeft. A report on the Sisal language project. *J. Parallel Distrib. Comput.*, 10:349–366, 1990.
5. S. Fitzpatrick, T. J. Harmer, A. Stewart, M. Clint, and J. M. Boyle. The automated transformation of abstract specifications of numerical algorithms into efficient array processor implementations. *Science of Computer Programming*, 28(1):1–41, 1997.
6. S. Gorlatch. Extracting and implementing list homomorphisms in parallel program development. *Science of Computer Programming*, 33(1):1–27, 1999.
7. P. Hammarlund and B. Lisper. On the relation between functional and data parallel programming languages. In *Proc. Sixth Conference on Functional Programming Languages and Computer Architecture*, pages 210–222. ACM Press, June 1993.
8. K. Hammond and G. Michaelson, editors. *Research Directions in Parallel Functional Programming*. Springer-Verlag, 1999.

9. J. M. D. Hill. Data Parallel Haskell: Mixing old and new glue. Tech. Rep. 611, Queen Mary and Westfield College, Dec. 1992.
10. J. Holmerin. Implementing data fields in Haskell. Technical Report TRITA-IT R 99:04, Dept. of Teleinformatics, KTH, Stockholm, Nov. 1999. <ftp://ftp.it.kth.se/Reports/paradis/DFH-report.ps.gz>.
11. C. B. Jay and P. A. Steckler. The functional imperative: shape! In C. Hankin, editor, *Proc. 7th European Symposium on Programming*, volume 1381 of *Lecture Notes in Comput. Sci.*, pages 139–53, Lisbon, Portugal, Mar. 1998. Springer-Verlag.
12. B. Lisper. Data parallelism and functional programming. In G.-R. Perrin and A. Darte, editors, *The Data Parallel Programming Model: Foundations, HPF Realization, and Scientific Applications*, Vol. 1132 of *Lecture Notes in Comput. Sci.*, pages 220–251, Les Ménuires, France, Mar. 1996. Springer-Verlag.
13. B. Lisper. Data fields. In *Proc. Workshop on Generic Programming*, Marstrand, Sweden, June 1998. <http://wsinwp01.win.tue.nl:1234/WGPPProceedings/>.
14. B. Lisper and P. Hammarlund. The data field model. Submitted. Preliminary version available as Tech. Rep. TRITA-IT R 99:02, Dept. of Teleinformatics, KTH, Stockholm, 2000.
15. R. S. Nikhil, Arvind, J. E. Hicks, S. Aditya, L. Augustsson, J.-W. Maessen, and Y. Zhou. pH language reference manual, version 1.0. Technical Report CSG-Memo-369, Massachusetts Institute of Technology, Laboratory for Computer Science, Jan. 1995.
16. J. T. O'Donnell. Data parallelism. In Hammond and Michaelson [8], chapter 7, pages 191–206.
17. J. Peterson, K. Hammond, L. Augustsson, B. Boutel, W. Burton, J. Fasel, A. D. Gordon, J. Hughes, P. Hudak, T. Johnsson, M. Jones, E. Meijer, S. L. Peyton Jones, A. Reid, and P. Wadler. Report on the programming language Haskell: A non-strict purely functional language, version 1.4, Apr. 1997. <http://www.haskell.org/definition/>.
18. N. Røjemo. *Garbage Collection, and Memory Efficiency, in Lazy Functional Languages*. PhD thesis, Department of Computing Science, Chalmers University of Technology, Gothenburg, Sweden, 1995.
19. L. Semenzato and P. Hilfinger. Arrays in FIDIL. In L. M. R. Mullin, M. Jenkins, G. Hains, R. Bernecky, and G. Gao, editors, *Arrays, Functional Languages, and Parallel Systems*, chapter 10, pages 155–169. Kluwer Academic Publishers, Boston, 1991.
20. J. M. Sipelstein and G. E. Blelloch. Collection-oriented languages. *Proc. IEEE*, 79(4):504–523, Apr. 1991.
21. D. B. Skillicorn. Building BSP programs using the refinement calculus. External Technical Report TR96-400, Dept. of Computing and Information Science, Queen's University, Kingston, Ontario, Oct. 1996.
22. G. L. Steele and W. D. Hillis. Connection Machine LISP: Fine grained parallel symbolic programming. In *Proc. 1986 ACM Conference on LISP and Functional Programming*, pages 279–297, Cambridge, MA, 1986. ACM.
23. G. Winskel. *The Formal Semantics of Programming Languages – An Introduction*. MIT Press, 1993.