

## **TACC Technical Report IMP-20**

# **A mathematical formalization of data parallel operations**

Victor Eijkhout\*

May 13, 2016

This technical report is a preprint of a paper intended for publication in a journal or proceedings. Since changes may be made before publication, this preprint is made available with the understanding that anyone wanting to cite or reproduce it ascertains that no published version in journal or proceedings exists.

Permission to copy this report is granted for electronic viewing and single-copy printing. Permissible uses are research and browsing. Specifically prohibited are *sales* of any copy, whether electronic or hardcopy, for any purpose. Also prohibited is copying, excerpting or extensive quoting of any report in another work without the written permission of one of the report's authors.

The University of Texas at Austin and the Texas Advanced Computing Center make no warranty, express or implied, nor assume any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed.

\* [eijkhout@tacc.utexas.edu](mailto:eijkhout@tacc.utexas.edu), Texas Advanced Computing Center, The University of Texas at Austin

## **Abstract**

We give a mathematical treatment of generalized data parallel operations, showing that formulating an algorithm in terms of data parallel operations allows for automatic derivation of a dataflow formulation. To this end we give a formal definition of the concept of ‘distribution’, and we introduce the concept of a ‘signature function’. The former is strictly a description of parallel data, independent of any algorithm considerations, whereas the latter is strictly an algorithm property, not involving any mention of parallel execution. Formally, our result is then that given distributions and signature functions, any task dependencies (whether realized as synchronization or as message passing) can be systematically derived.

The following IMP reports are available or under construction:

- IMP-00** The IMP Elevator Pitch
- IMP-01** IMP Distribution Theory
- IMP-02** The deep theory of the Integrative Model
- IMP-03** The type system of the Integrative Model
- IMP-04** Task execution in the Integrative Model
- IMP-05** Processors in the Integrative Model
- IMP-06** Definition of a ‘communication avoiding’ compiler in the Integrative Model (under construction)
- IMP-07** Associative messaging in the Integrative Model (under construction)
- IMP-08** Resilience in the Integrative Model (under construction)
- IMP-09** Tree codes in the Integrative Model
- IMP-10** Thoughts on models for parallelism
- IMP-11** A gentle introduction to the Integrative Model for Parallelism
- IMP-12** K-means clustering in the Integrative Model
- IMP-13** Sparse Operations in the Integrative Model for Parallelism
- IMP-14** 1.5D All-pairs Methods in the Integrative Model for Parallelism (under construction)
- IMP-15** Collectives in the Integrative Model for Parallelism
- IMP-16** Processor-local code generation (under construction)
- IMP-17** The CG method in the Integrative Model for Parallelism (under construction)
- IMP-18** A tutorial introduction to IMP software (under construction)
- IMP-19** Report on NSF EAGER 1451204.
- IMP-20** A mathematical formalization of data parallel operations
- IMP-21** Adaptive mesh refinement (under construction)
- IMP-22** Implementing LULESH in IMP (under construction)
- IMP-23** Distributed computing theory in IMP (under construction)

## 1 Introduction

Much theoretical work has been done about parallel and concurrent programming. However, almost without exception this work starts from an implicit assumption of independent tasks that interact. Any overall behaviour of the parallel assembly of processes is at best an emergent property. This does not do justice to the nature of the way parallelism often occurs in scientific computing: there an essentially sequential program is executed (‘matrix times vector, inner product of the output with another vector, scale vector by the inner product value’, et cetera) of operations that are, in a generalized sense, data parallel. The parallelism is solely due to the fact that the objects are distributed.

The problem with coding in this manner is that, in any but the most regular applications, the synchronizations and communications between the underlying processes are hard to derive by a compiler or middleware layer. Thus, systems based on such ‘sequential semantics’ have had limited success in scientific computing.

In this paper we give a mathematical foundation for algorithms that can be formulated in terms of generalized data parallelism. We give a non-standard definition of the concept of ‘distribution’ and we define the ‘signature function’ of a data parallel operation. Taken together, these then allow for task synchronizations and communications to be formally derived.

In particular, we show that we can derive a dataflow formulation of the algorithm from the sequential semantics description. In a practical programming system this dataflow can then be realized in terms of task dependencies or message passing, or hybrid combinations of these.

## 2 Motivating example

We consider a simple data parallel example, and show how it leads to the basic distribution concepts of Integrative Model for Parallelism (IMP): the three-point operation

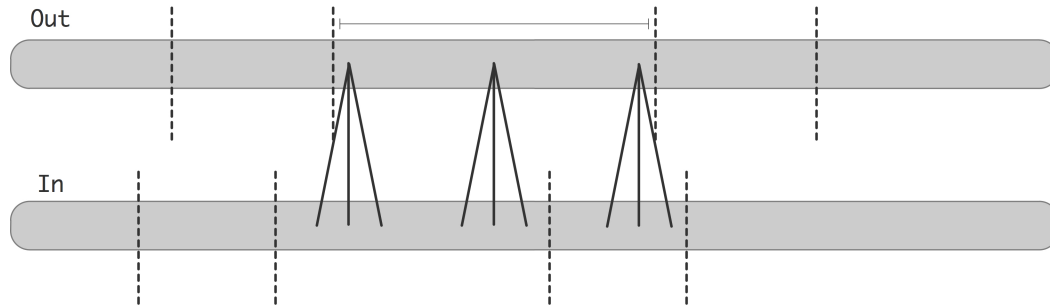
$$\forall_i: y_i = f(x_i, x_{i-1}, x_{i+1})$$

which describes for instance the 1D heat equation

$$y_i = 2x_i - x_{i-1} - x_{i+1}.$$

(Stencil operations are much studied; see e.g., [7] and the polyhedral model, e.g., [1]. However, we claim far greater generality for our model.) We illustrate this graphically

by depicting the input and output vectors, stored distributed over the processors by contiguous blocks, and the three-point combining operation:

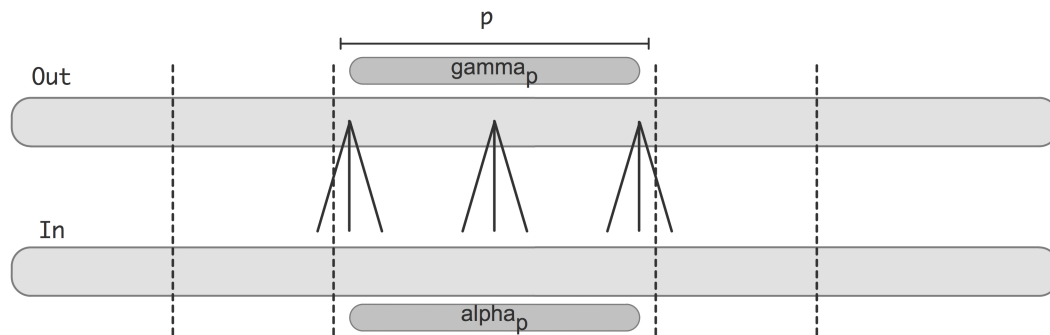


The distribution indicated by vertical dotted lines we call the  $\alpha$ -distribution for the input, and the  $\gamma$ -distribution for the output. These distributions are mathematically given as an assignment from processors to sets of indices:

$$\alpha: p \mapsto [i_{p,\min}, \dots, i_{p,\max}].$$

The traditional concept of distributions in parallel programming systems is that of an assignment of data indices to a processor, reflecting that each index ‘lives on’ one processor, or that that processor is responsible for computing that index of the output. We turn this upside down: we define a distribution as a mapping from processors to indices. This means that an index can ‘belong’ to more than one processor. (The utility of this for redundant computing is obvious. However, it will also seen to be crucial for our general framework.)

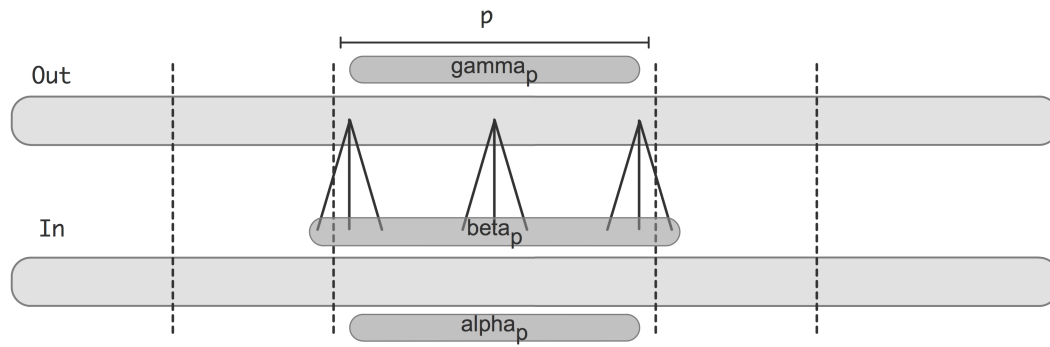
For purposes of exposition we will now equate the input  $\alpha$ -distribution and the output  $\gamma$ -distribution, although that will not be necessary in general.



This picture shows how, for the three-point operation, some of the output elements on processor  $p$  need inputs that are not present on  $p$ . For instance, the computation of  $y_i$  for  $i_{p,\min}$  takes an element from processor  $p - 1$ . This gives rise to what we call the  $\beta$ -distribution:

$\beta(p)$  is the set of indices that processor  $p$  needs to compute the indices in  $\gamma(p)$ .

The next illustration depicts the different distributions for one particular process:



Observe that the  $\beta$ -distribution, unlike the  $\alpha$  and  $\gamma$  ones, is not disjoint: certain elements live on more than one processing element. It is also, unlike the  $\alpha$  and  $\gamma$  distributions, not specified by the programmer: it is derived from the  $\gamma$ -distribution by applying the shift operations of the stencil. That is,

The  $\beta$ -distribution brings together properties of the algorithm and of the data distribution.

We will formalize this derivation below.

## 2.1 Deriving the dataflow formulation

This gives us all the ingredients for reasoning about parallelism. Defining a *kernel* as a mapping from one distributed data set to another, and a *task* as a kernel on one particular process(or), all data dependence of a task results from transforming data from  $\alpha$  to  $\beta$ -distribution. By analyzing the relation between these two we derive at dependencies between processors or tasks: each processor  $p$  depends on some predecessors  $q_i$ , and this set of predecessors can be derived from the  $\alpha, \beta$  distributions:  $q_i$  is a predecessor if

$$\alpha(q_i) \cap \beta(p) \neq \emptyset.$$

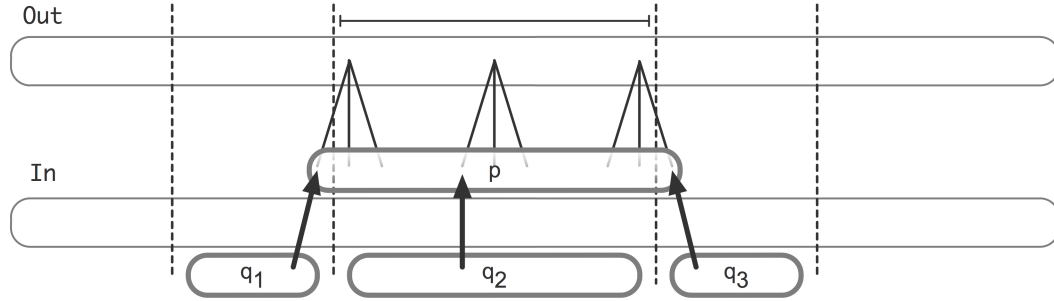


Figure 1 illustrates this: the left DAG is the sequential program of a heat equation evolution; the right DAG is the dataflow representation derived when this sequential program is run on six processors.

In message passing, these dataflow dependences obviously corresponds to actual messages: for each process  $p$ , the processes  $q$  that have elements in  $\beta(p)$  send data to  $p$ . (If  $p = q$ , of course at most a copy is called for.) Interestingly, this story has an interpretation in tasks on shared memory too. If we identify the  $\alpha$ -distribution on the input with tasks that produce this input, then the  $\beta$ -distribution describes what input-producing tasks a task  $p$  is dependent on. In this case, the transformation from  $\alpha$  to  $\beta$ -distribution gives rise to a *dataflow* formulation of the algorithm.

## 2.2 Programming the model

In our motivating example we showed how the concept of ‘ $\beta$ -distribution’ arises, and the role it plays combining properties of the data distributions and of the algorithm’s data dependencies. This distribution generalizes concepts such as the ‘halo region’ in distributed stencil calculations, but its applicability extends to all of (scientific) parallel computing. For instance, for collectives we can define a  $\beta$ -distribution, which is seen to equal the  $\gamma$ -distribution.

It remains to be argued that the  $\beta$  distribution can actually be used as the basis for a software system. To show this, we associate with the function  $f$  that we are computing an expression of the algorithm (not the parallel!) data dependencies, called the ‘signature function’, denoted  $\sigma_f$ . For instance for the computation of  $y_i = f(x_i, x_{i-1}, x_{i+1})$ , the signature function is

$$\sigma_f(i) = \{i, i-1, i+1\}.$$

With this, we state (without proof; for which see section 3.3 and [2]) that

$$\beta = \sigma_f(\gamma).$$

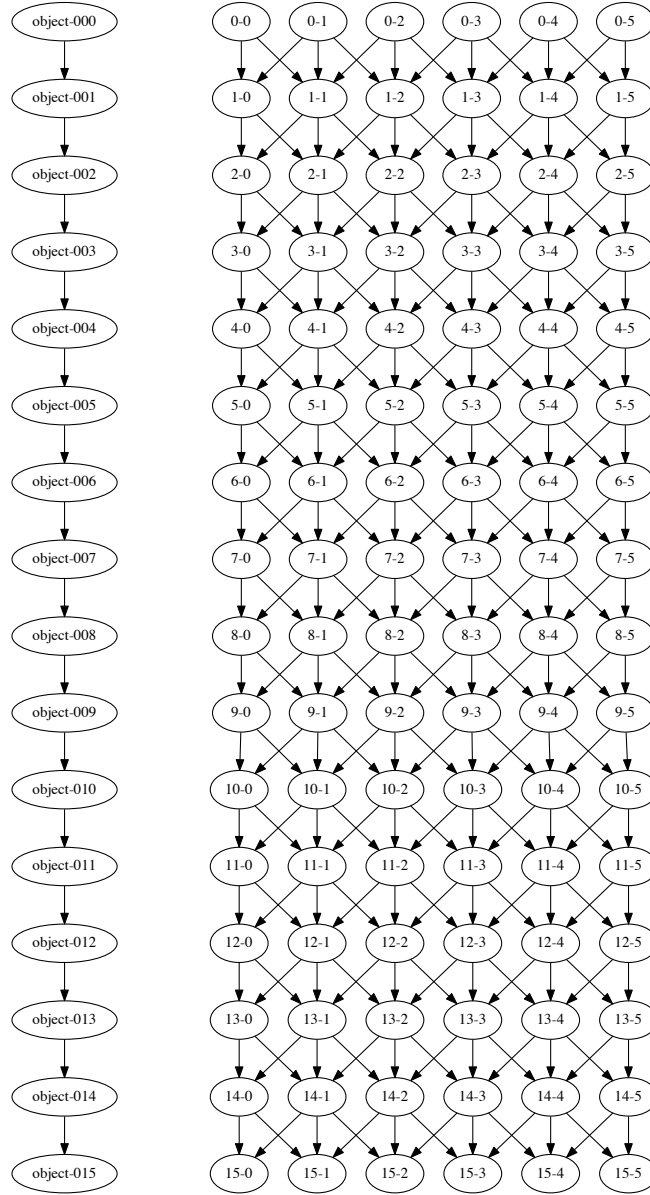


Figure 1: Kernel (left) and task relations (right) for the one-dimensional heat equation, executing 15 steps on 6 processors.

It follows that, if the programmer can specify the data dependencies of the algorithm, a compiler/runtime system can derive the  $\beta$  distribution, and from it, task dependencies and messages for parallel execution.

Specifying the signature function is quite feasible, but the precise implementation depends on the context. For instance, for regular applications we can adopt a syntax similar to stencil compilers such as the Pochoir compiler [7]. For sparse matrix applications the signature function is isomorphic to the adjacency graph; for collective operations,  $\beta = \gamma$  often holds; et cetera.

### 3 Formal definition

#### 3.1 Data parallel computation

The Integrative Model for Parallelism (IMP) is a theory of data parallel functions. By this we mean functions where each element of a distributed output object is computed from one or more elements of one or more distributed input objects.

- Without loss of generality we limit ourselves to a single input object.
- Since all output elements can be computed independently of each other, we call this a ‘data parallel’ function. In our context this has no connotations of SIMD or synchronization; it merely expresses independence.

Formally, a data parallel computation is the use of a function with a single output to compute the elements of a distributed object:

$$\text{Func} \equiv \text{Real}^k \rightarrow \text{Real}$$

where  $k$  is some integer.

Since we will mostly talk about indices rather than data, we define  $\text{Ind} \equiv 2^N$  and we describe the structure of the data parallel computation through a ‘signature function’:

$$\text{Signature} \equiv N \rightarrow \text{Ind}.$$

In our motivating example, where we computed  $y_i = f(x_{i-1}, x_i, x_{i+1})$ , our signature function was

$$\sigma_f \equiv i \mapsto \{i-1, i, i+1\}.$$

- The signature function can be compactly rendered in cases of a stencil computation.



- In general it describes the bi-partite graph of data dependencies. Thus, for sparse computations it is isomorphic to the sparse matrix, and can be specified as such.
- In certain cases, the signature function can be most compactly be rendered as a function recipe. For instance, for 1D multigrid restriction it would be given as  $\sigma(i) = \{2i, 2i + 1\}$ .
- For collectives such as an ‘allreduce’, the signature function expresses that the output is a function of all inputs:  $\forall_i: \sigma(i) = N$ .

### 3.2 Distributions

We now formally define distributions as mappings from processors to sets of indices:

$$\text{Distr} \equiv \text{Proc} \rightarrow \text{Ind}.$$

Traditionally, distributions are considered as mappings from data elements to processors, which assumes a model where a data element lives uniquely on one processor. By turning this definition around we have an elegant way of describing:

- Overlapping distributions such as halo data, where data has been gathered on a processor for local operations. Traditionally, this is considered a copy of data ‘owned’ by another processor.
- Rootless collectives: rather than stating that all processors receive an identical copy of a result, we consider them to actually own the same item.
- Redundant execution. There are various reasons for having operations executed redundantly on more than one processor. This can for instance happen in the top levels of a coarsening multilevel method, or in redundant computation for resilience.

We now bring together the concepts of signature function and distribution:

1. We can extend the signature function concept, defined above as mapping integers to sets of integers, to a mapping from sets to sets: with the obvious definition that, for  $\sigma \in \text{Signature}, S \in \text{Ind}$ :

$$\sigma(S) = \{\sigma(i) : i \in S\}.$$

In our motivating example,

$$\sigma([i_{\min}, i_{\max}]) = [i_{\min} - 1, i_{\max} + 1].$$

2. We then extend this to distributions with the definition that for  $\sigma \in \text{Signature}$  and  $u \in \text{Distr}$

$$\sigma(u) \equiv p \mapsto \sigma(u(p)) \quad \text{where} \quad \sigma(u(p)) = \{\sigma(i) \mid i \in u(p)\}$$

We now have the tools for our grand result.

### 3.3 Definition and use of $\beta$ -distribution

Consider a data parallel operation  $y = f(x)$  where  $y$  has distribution  $\gamma$ , and  $x$  has distribution  $\alpha$ . We call a local operation to be one where every processor has all the elements of  $x$  needed to compute its part of  $y$ . By the above overloading mechanism, we find that the total needed input on processor  $p$  is  $\sigma(\gamma(p))$ .

This leads us to define a *local operation* formally as:

**Definition 1** We call a kernel  $y = f(x)$  a local operation if  $x$  has distribution  $\alpha$ ,  $y$  has distribution  $\gamma$ , and

$$\alpha \supset \sigma_f(\gamma).$$

That is, for a local operation every processor already owns all the elements it needs for its part of the computation.

Next, we call  $\sigma_f(\gamma)$  the ‘ $\beta$ -distribution’ of a function  $f$ :

**Definition 2** If  $\gamma$  is the output distribution of a computation  $f$ , we define the  $\beta$ -distribution as

$$\beta = \sigma_f(\gamma).$$

Clearly, if  $\alpha \supset \beta$ , each processor has all its needed inputs, and the computation can proceed locally. However, this is often not the case, and considering the difference between  $\beta$  and  $\alpha$  gives us the description of the task/process communication:

**Corollary 1** If  $\alpha$  is the input distribution of a data parallel operation, and  $\beta$  as above, then processor  $q$  is a predecessor of processor  $p$  if

$$\alpha(q) \cap \beta(p) \neq \emptyset.$$

*Proof.* The set  $\beta(p)$  describes all indices needed by processor  $p$ ; if the stored elements in  $q$  overlap with this, the computation on  $q$  that produces these is a predecessor of the subsequent computation on  $p$ .

This predecessor relation takes a specific form depending on the parallelism mode. For instance, in message passing it takes form of an actual message from  $q$  to  $p$ ; in a Directed Acyclic Graph (DAG) model such as OpenMP tasking it becomes a ‘task wait’ operation.

**Remark 1** *In the context of Partial Differential Equation (PDE) based applications, our  $\beta$ -distribution corresponds loosely to the ‘halo’ region. The process of constructing the  $\beta$ -distribution is implicitly part of such packages as PETSc [6], where the communication resulting from it is constructed in the `MatAssembly` call. Our work takes this ad-hoc calculation, and shows that it can formally be seen to underlie a large part of scientific parallel computing.*

## 4 Practical importance of this theory

The above discussion considered operations that can be described as ‘generalized data parallel’. From such operations one can construct many scientific algorithms. For instance, in a multigrid method a red-black smoother is data parallel, as are the restriction and prolongation operators.

In the IMP model these are termed ‘kernels’, and each kernel gives rise to one layer of task dependencies; see section 2.1. Taking together the dependencies for the single kernels then gives us a complete task graph for a parallel execution; the edges in this graph can be interpreted as MPI messages or strict dependencies in a DAG execution model.

Demonstration software along these lines has been built, showing performance comparable to hand-coded software; see [3].

## 5 Summary and further reading

In this report we have motivated and defined our mathematical underpinning of a concept that is sometimes known as ‘halo region’ and shown how it is a basic tool for parallel computation. The story of the IMP model is further developed in our reports series. We particularly draw attention to:

- [2] goes in more detail on the concept of distributions, giving many examples.
- [5] gives a type system of IMP, showing how our computer code is a direct implementation of the math; [4] gives a tutorial in the use of this code.
- [3] reports on how our work is progressing and shows results from some proof-of-concept implementations of algorithms.

## Acknowledgement

This work was supported by NSF EAGER grant 1451204. The code for this project is available at <https://bitbucket.org/VictorEijkhout/imp-code>.

## References

- [1] Roshan Dathathri, Chandan Reddy, Thejas Ramashekar, and Uday Bondhugula. Generating efficient data movement code for heterogeneous architectures with distributed-memory. In *Proceedings of the 22nd international conference on Parallel architectures and compilation techniques*, PACT '13, pages 375–386, Piscataway, NJ, USA, 2013. IEEE Press.
- [2] Victor Eijkhout. IMP distribution theory. Technical Report IMP-01, Integrative Programming Lab, Texas Advanced Computing Center, The University of Texas at Austin, 2014.
- [3] Victor Eijkhout. Report on NSF EAGER 1451204. Technical Report IMP-19, Integrative Programming Lab, Texas Advanced Computing Center, The University of Texas at Austin, 2014.
- [4] Victor Eijkhout. A tutorial introduction to IMP software (under construction). Technical Report IMP-18, Integrative Programming Lab, Texas Advanced Computing Center, The University of Texas at Austin, 2014.
- [5] Victor Eijkhout. The type system of the integrative model. Technical Report IMP-03, Integrative Programming Lab, Texas Advanced Computing Center, The University of Texas at Austin, 2014.
- [6] W. D. Gropp and B. F. Smith. Scalable, extensible, and portable numerical libraries. In *Proceedings of the Scalable Parallel Libraries Conference, IEEE 1994*, pages 87–93.
- [7] Yuan Tang, Rezaul Alam Chowdhury, Bradley C. Kuszmaul, Chi-Keung Luk, and Charles E. Leiserson. The pochoir stencil compiler. In *Proceedings of the 23rd ACM symposium on Parallelism in algorithms and architectures*, SPAA '11, pages 117–128, New York, NY, USA, 2011. ACM.