# Synthesizing MPI Implementations from Functional Data-Parallel Programs

Inferring data distributions using types

Tristan Aubrey-Jones
Bernd Fischer





# People need to program distributed memory architectures:

#### **GPUs**



Graphics/GPGPU CUDA/OpenCL

Memory levels: thread-local, block-shared, device-global.

#### Server farms/compute clusters



Big Data/HPC; MapReduce/HPF; Ethernet/Infiniband

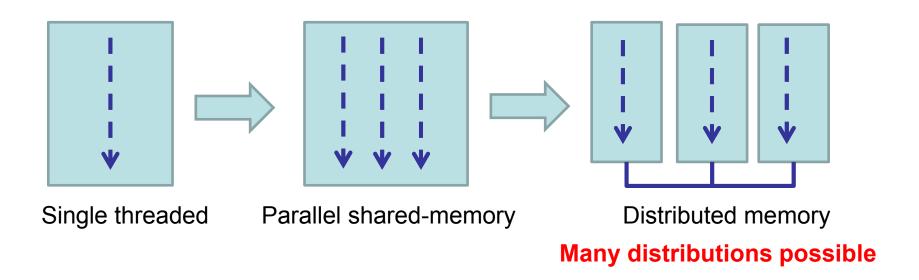
#### **Future many-core architectures?**



#### **Our Aim**



# To automatically generate distributed-memory cluster implementations (C++/MPI).



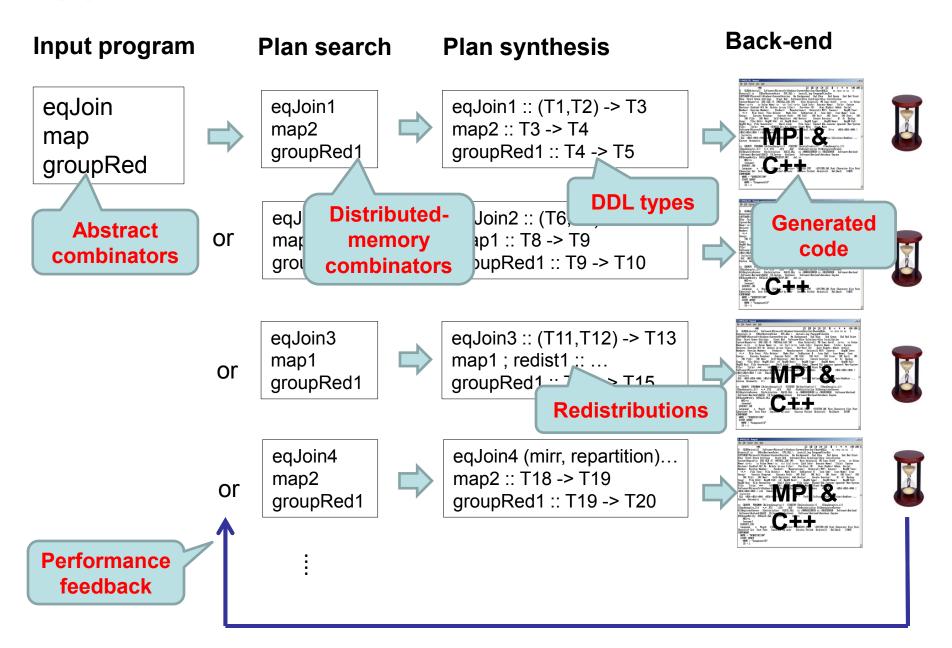
We want this **not just for arrays** (i.e., other collection types as well as disk backed collections).

# **Approach**

- We define Flocc, a **data-parallel DSL**, where parallelism is expressed via skeletons or **combinator functions** (HOFs).
- We use distributed data layout (DDL) types to carry data distribution information, and type inference to drive code generation.
- We develop a code generator that searches for well performing implementations, and generates C++/MPI implementations.

# **Approach**





#### What's new?

#### Last++ time:

- Functional DSL with data-parallel combinators
- Data distributions for distributed-memory implementations using dependent types
- Distributions for maps, arrays, and lists
- Synthesis of data distributions using type inference algorithm

#### Since then:

- MPI/C++ code generation
- Performance-feedbackbased data distribution search
- Automatic redistribution insertion
- Local data layouts
- Arrays with ghosting
- Type inference with E-unification
- (Thesis submitted)



What we presented last time.

**RE-CAP** 

# A Distributed Map type

The **DMap** type extends the basic **Map k v** type to symbolically describe how the Map should be distributed on the cluster

$$dt ::= \ldots \mid \mathsf{DMap}\ t_1\ t_2\ f\ d_1\ d_2\mid \ldots$$

- Key and value types t1 and t2
- dependent types! Partition function  $\mathbf{f}$ :  $(t1,t2) \rightarrow N$ 
  - takes key-value pairs to node coordinates
- Partition dimension identifier d1
  - specifies which nodes the coordinates map onto
- Mirror dimension identifier d2
  - specifies which nodes to mirror the partitions over

Also works for Arrays (DArr) and Lists (DList)



# Distribution types for group reduces

```
groupReduce2 : \forall k, k', v, v', d, m, \Pi(f,\_,\_,\_) : 
 ((k,v) → k', (k,v) → v', (v',v') → v', DMap k v f d m)
```

→ DMap k' v' fst d m

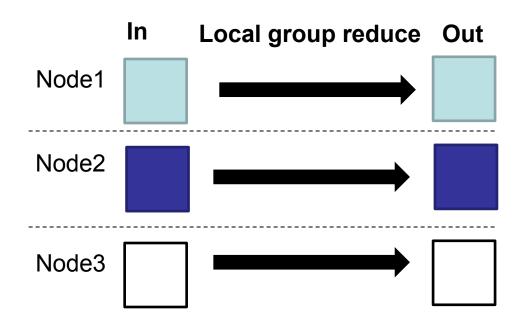
- Π binds concrete term in AST
- creates a reference to argument's AST when instantiated
- must match concrete reference when unifying (*rigid*)
- used to specify that a collection is distributed using a specific key projection function



# Distribution types for group reduces

```
groupReduce2 : \forall k, k', v, v', d, m, \Pi(f,\_,\_,\_) : 
 ((k,v) \rightarrow k', (k,v) \rightarrow v', (v',v') \rightarrow v', DMap k v f d m)
```

- → DMap k' v' fst d m
- Input must be partitioned using the groupReduce's key projection function f
- Result keys are already co-located

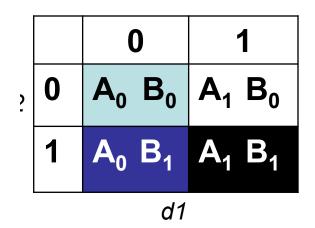


No inter-node communication necessary (but constrains input distribution)

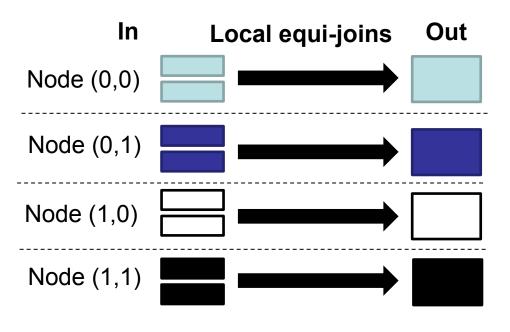
# Distribution types for joins

eqJoin3: ∀f, k', k1, k2, v1, v2, d1, d2, m,
 ((k1,v1) → k', (k2,v2) → k',
 DMap k1 v1 IftFun(f) d1 (d2,m),
 DMap k2 v2 rhtFun(f) d2, (d1,m))

 $\rightarrow$  DMap (k1,k2) (v1,v2) f (d1,d2) m



Left and right input partitioned and mirrored on orthogonal dims



- No inter-node communication
- Output can be partitioned by any f
- Can be more efficient than mirroring whole input

# **Deriving distribution plans**

#### Hidden from user

- Different distributed implementations of each combinator
- Enumerate different choices of combinator implementations
- Each combinator implementation has a DDL type
- Use type inference to check if a set of choices is sound and infer data distributions
- Backend templates for each combinator implementation

#### eqJoin1 / eqJoin2 / eqJoin3

#### Southampton

# Distributed matrix multiplication - #1

```
let R = eqJoin (((ai,aj),_) -> aj,
                    \((bi,bj),_) -> bi,
                    A, B) in
groupReduce (\(((ai,aj),(bi,bj)),_) -> (ai,bj),
                (\_,(av,bv)) \rightarrow mul(av,bv), add, R)
A: DMap (Int,Int) Float ((ai,aj), \_) \rightarrow ai d1 (d2,m)
B: DMap (Int,Int) Float ((bi,bj), ) \rightarrow bj d2 (d1,m)
R: DMap ((Int,Int),(Int,Int)) (Float,Float)
           (((ai,aj),(bi,bj)), ) \rightarrow (ai,bj)(d1,d2) m
C: eqJoin3: Vf, k', k1, k2, v1, v2, d1, d2, m,
      ((k1,v1) \rightarrow k', (k2,v2) \rightarrow k',
                                                     d2)
       DMap k1 v1 IftFun(f) d1 (d2,m),
       DMap k2 v2 rhtFun(f) d2, (d1,m))
                                                     pmputation.
    \rightarrow DMap (k1,k2) (v1,v2) f (d1,d2) m
```

### Distributed matrix multiplication - #1

R : DMap ((Int,Int),(Int,Int)) (Float,Float)  $(((ai,aj),(bi,bj)),\_) \rightarrow (ai,bj) (d1,d2) m$ 

C: DMap (Int,Int) Float fst (d1,d2) m

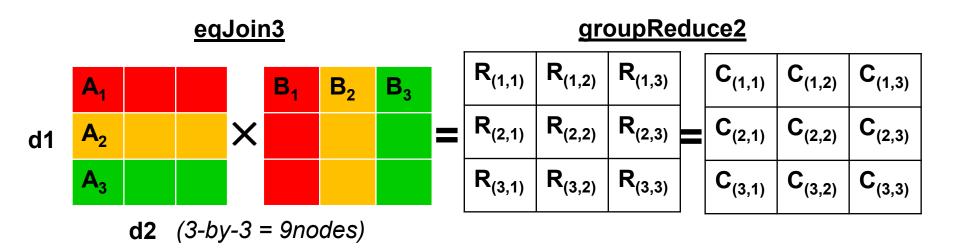
```
groupReduce2 : \forall k, k', v, v', d, m, \Pi(f,\_,\_,\_):
((k,v) \rightarrow k', (k,v) \rightarrow v', (v',v') \rightarrow v', DMap \ k \ v \ f \ d \ m)
\rightarrow DMap \ k' \ v' \ fst \ d \ m
```

# Distributed matrix multiplication - #1

- groupReduce2 in part by key Πf, out part by fst (out keys)
- eqJoin3 left part by lftFun(f), right rhtFun(f), out by any f

A: Partitioned by row along d1, mirrored along d2

B: Partitioned by column along d2, mirrored along d1



C: Partitioned by (row, column) along (d1, d2)

A common solution.

# Distributed matrix multiplication - #2

- eqJoin1 left part by Πf1, right by Πf2, aligned along d
- groupReduce1 in part by any f, out part by fst (out keys)

```
A: Partitioned by col along d

(aligned with A)
```

```
A: DMap (Int,Int) Float \((ai,aj),_)->aj d m
```

B: DMap (Int,Int) Float \((bi,bj),\_)->bi d m

R: DMap ((Int,Int),(Int,Int)) (Float,Float)

((\((ai,aj),\_) ->aj) · lft) d m

C: DMap (Int,Int) Float fst d m

C: Partitioned by (row, column) along d

Must exchange R during groupReduce1.



What's new since last time.

# **RECENT WORK**

#### Since last time...



- Local data layouts
- Distributed arrays with ghosting
- Automatic redistribution insertion
- E-unification in type inference
- MPI/C++ code generation
- Performance-feedback-based data distribution search
- Proof of concept (~25k Haskell loc)

#### Southampton

# **Local data layouts**

- Extra DDL type parameters for local layouts
  - Choice of data structure (or storage mode)
    - Sorted std::vector, hash-map, tree-map, value-stream etc...
  - Order of elements or key indexed by

```
groupReduce2 :: \Pi(f,\_,\_,\_) : ((i,v)->j, (i,v)->w, (w,w)->w, DMap \langle \ldots \rangle Stm f) -> DMap \langle \ldots \rangle Stm fst union :: (DMap \langle \ldots \rangle Stm fst, DMap \langle \ldots \rangle Iter fst) -> DMap \langle \ldots \rangle Stm fst diff1ii :: (DMap \langle \ldots \rangle Stm f, DMap \langle \ldots \rangle Hash fst) -> DMap \langle \ldots \rangle Stm f diff1iii :: (DMap \langle \ldots \rangle Stm f, DMap \langle \ldots \rangle Tree fst) -> DMap \langle \ldots \rangle Stm f diff1iv :: (DMap \langle \ldots \rangle Stm f, DMap \langle \ldots \rangle Vec fst) -> DMap \langle \ldots \rangle Stm f
```

#### Southampton

# **Extended array distributions**

- Supports
  - index offsets, axis-reversal (all HPF alignments)
  - cyclic, blocked, and block-cyclic distributions (all HPF dists)
  - ghosted regions/fringes
- Index transformer functions in DDL types
  - Take and return tuples of integer array indices
    - ▶ Block-sizes
    - ▶ Index directions
    - ▶ Index offsets
    - ▶ Index ghosted fringe sizes (left and right)
- Distribution for Jacobi 1D
  - DArr ... bs dir id id (+1) (+1) -> DArr ... bs dir id id id
- Relies on E-unification (later...)

#### **Automatic redistribution insertion**

- Data re-distribution and re-layout functions are type casts.
- For invalid Flocc plans (i.e., that don't type check) insert just enough redistributions (or re-layouts) to make type check.
- Means can synthesize a valid plan for any choice of combinator implementations.
- Finds implementations that benefit from redistributing data so more efficient combinator implementations can be used.

```
mirrMap :: DMap k v f d m -> DMap k v f d (m,m') repartMap :: DMap k v f1 d1 m -> DMap k v f2 d2 m saveVecMap :: DMap \langle \dots \rangle Stm f -> DMap \langle \dots \rangle Vec f sortVecMap :: DMap \langle \dots \rangle Vec f -> DMap \langle \dots \rangle Vec g readVecMap :: DMap \langle \dots \rangle Vec f -> DMap \langle \dots \rangle Iter f readIterMap :: DMap \langle \dots \rangle Iter f -> DMap \langle \dots \rangle Stm f
```

#### Southampton

#### E-unification

- Adding equational theories for projection, permutation, and index transformation functions to DDL type inference
- Use E-prover and "question" conjectures to return values for existentially qualified variables
- Allows improved array distributions and more flexible DDL types to be supported
- (Not integrated with current prototype)

# E-unification: projection functions

$$\beta \cdot \alpha \stackrel{def}{=} \langle \mathbf{x} \to (\beta(\alpha \mathbf{x})) \rangle$$

$$\alpha \otimes \beta \stackrel{def}{=} \langle (\mathbf{x}, \mathbf{y}) \to (\alpha \mathbf{x}, \beta \mathbf{y}) \rangle$$

$$id \stackrel{def}{=} \langle \mathbf{x} \to \mathbf{x} \rangle$$

$$\Delta \stackrel{def}{=} \langle \mathbf{x} \to (\mathbf{x}, \mathbf{x}) \rangle$$

$$\Pi_1 \stackrel{def}{=} \langle (\mathbf{x}, \mathbf{y}) \to \mathbf{x} \rangle$$

$$\Pi_2 \stackrel{def}{=} \langle (\mathbf{x}, \mathbf{y}) \to \mathbf{y} \rangle$$

$$f \cdot id = f$$

$$id \cdot f = f$$

$$\Pi_1 \cdot \Delta = id$$

$$\Pi_2 \cdot \Delta = id$$

$$(\Pi_1 \otimes \Pi_2) \cdot \Delta = id \otimes id$$

$$(\Pi_{2} \otimes \Pi_{1}) \cdot \Delta \cdot (\Pi_{2} \otimes \Pi_{1}) \cdot \Delta = \mathrm{id} \otimes \mathrm{id}$$

$$(f_{1} \otimes f_{2}) \cdot (\Pi_{2} \otimes \Pi_{1}) \cdot \Delta = (\Pi_{2} \otimes \Pi_{1}) \cdot \Delta \cdot (f_{2} \otimes f_{1})$$

$$\Pi_{1} \cdot (f_{1} \otimes f_{2}) = f_{1} \cdot \Pi_{1}$$

$$\Pi_{2} \cdot (f_{1} \otimes f_{2}) = f_{2} \cdot \Pi_{2}$$

$$\Pi_{1} \cdot (f_{1} \otimes f_{2}) \cdot \Delta = f_{1}$$

$$\Pi_{2} \cdot (f_{1} \otimes f_{2}) \cdot \Delta = f_{2}$$

$$\Delta \cdot f = (f \otimes f) \cdot \Delta$$

$$(f_{1} \cdot f_{2}) \cdot f_{3} = f_{1} \cdot (f_{2} \cdot f_{3})$$

$$(f_{1} \otimes f_{2}) \cdot (f_{3} \otimes f_{4}) = (f_{1} \cdot f_{3}) \otimes (f_{2} \cdot f_{4})$$

# E-unification: indexing functions

$$i^{-1-1} = i$$

$$1^{-1} = 1$$

$$+(i) \cdot -(i) = id$$

$$-(i) \cdot +(i) = id$$

$$\times (i) \cdot \times (i^{-1}) = id$$

$$\times (i^{-1}) \cdot \times (i) = id$$

$$g \cdot id = g$$

$$id \cdot g = g$$

$$id \cdot g = g$$

$$(g_1 \cdot g_2) \cdot g_3 = g_1 \cdot (g_2 \cdot g_3)$$

$$(g_1 \otimes g_2) \cdot (g_3 \otimes g_4) = (g_1 \cdot g_3) \otimes (g_2 \cdot g_4)$$

# E-unification: permutation functions

$$\alpha \odot \beta \stackrel{def}{=} \langle \mathbf{x} \to (\alpha \ \mathbf{x}, \beta \ \mathbf{x})$$
  
null  $\stackrel{def}{=} \langle - \rangle$ 

$$h \cdot id = h$$

$$id \cdot h = h$$

$$h \odot \text{null} = h$$

$$\text{null} \odot h = h$$

$$(h_1 \cdot h_2) \cdot h_3 = h_1 \cdot (h_2 \cdot h_3)$$

$$(h_1 \odot h_2) \odot h_3 = h_1 \odot (h_2 \odot h_3)$$

$$(h_1 \odot h_2) \cdot h_3 = ((h_1 \cdot h_3) \odot (h_2 \cdot h_3))$$

## **Code generation**

- Generates C++ and MPI from plans (i.e., Flocc programs with concrete combinator implementations and inferred DDL types)
- Transforms to DFG and uses expression templates to generate code.
- Currently supports mapand list-based combinator templates.
- Uses "stream" local storage mode to splice multiple consumers into a producer's loop body.

```
// BEGIN readList
   int v14 = 0;
   // BEGIN reduceList init
   double v17(v1);
   double v18:
   // END reduceList init
   std::vector<double >::iterator v15(v6->begin());
   std::vector<double >::iterator v12(v11->begin());
   std::vector<double >::iterator v13(v11->end()):
14 for (; v12 != v13 && v14 > 0; v12++) {
     // BEGIN zip consume:
     if (v15 < v6 -> end()) {
     // BEGIN mapList consume:
         double v16;
         v16 = (*v12) * (*v15);
         // BEGIN reduceList consume
         v18 = v17 + v16;
         // END reduceList consume
     // END mapList consume
     else if (v15 == v6 -> end()) v14 --;
     v15++:
     // END zip consume
31 }
33 // BEGIN reduceList fin
   cartComm.Allreduce(&v18, &v17, sizeof(double), MPI_PACKED, v103);
   if (cartComm != cartComm) {
     cartComm.Bcast(&v17, sizeof(double), MPI_PACKED, rootRank);
  // END reduceList fin
40 // END readList
```

#### Southampton

# **Code generation**

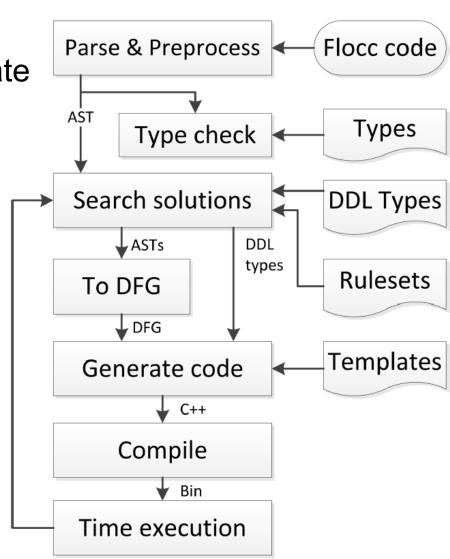
#### Performance comparable with hand-coded versions:

	PLINQ			Manual MPI		
Program	Speedup	Compiler	Data	Speedup	Compiler	Data
Dot product	4.96×	gcc -Ofast	2.2GB	0.99×	icc -O3	4.5GB
Simple linear regression	$137 \times$	gcc -Ofast	3GB	$0.89 \times$	icc -O3	3GB
				$0.61 \times$	gcc -O3	3GB
Standard deviation	$98.6 \times$	gcc -Ofast	3GB	$0.88 \times$	icc -O3	3GB
				$1.00 \times$	gcc - O3	3GB
Histogram	$31.5 \times$	gcc -Ofast	32MB	$0.73 \times$	icc -O3	8GB
Matrix multiply	$342 \times$	gcc -Ofast	1.4MB	$6.14 \times$	gcc -O3	140MB
				$0.49 \times$	icc -O3	140MB

- PLINQ comparisons run on quad-core (Intel Xeon W3520/2.67GHz) x64 desktop with 12GB RAM.
- C++/MPI comparisons run on <u>Iridis3&4</u>: a 3<sup>rd</sup> gen cluster with ~1000 Westmere compute nodes, each with two 6-core CPUs and 22GB RAM, over an InfiniBand network. Speedups compared to sequential, averaged over 1,2,3,4,8,9,16,32 nodes.

#### Performance-feedback-based search

- Tried different search algorithms to explore candidate implementations
- For each candidate we automatically insert redistributions to make it type check
- We evaluate each candidate by code generating, compiling, and running it on some test data
- Generates C++ using MPI



#### Performance-feedback-based search



- Tried 946 different combinations of search heuristics applied to 4 map-based example programs
- Heuristics composed of
  - Search algorithms

    - ▷ Depth/first exhaustive
  - Termination conditions
  - Runtime pruning
- Found
  - Genetic-searches successfully reduce search time
  - Fixed budget termination best
  - Fixed budget runtime pruning best
  - Need to enumerate different redistribution insertion variants

#### The Pros and Cons...

#### **Benefits**

- Multiple distributed collections:
   Maps, Arrays, Lists...
- Generates distributed algorithms fully automatically
- Performance feedback more accurate/flexible than cost metrics
- Finds algorithms including redistributions
- Synthesizes local layouts
- Can support in-memory and disk backed collections (e.g. for Big Data)

#### Limitations

- Current implementation mainly has list and map combinator backend templates
- Current implementation's redistribution insertion algorithm is slow

#### **Future work**



- Extend prototype implementation.
  - Array combinator backend templates
  - Faster redistribution insertion
- Integrate equational theories with implementations.
- Support more distributed memory architectures (GPUs).
- Retrofit into an existing functional language.
- Similar type inference for imperative languages?
- (pass PhD viva)

# **QUESTIONS?**