# Compiler Optimization for Irregular Memory Access Patterns in PGAS Programs

Thomas B. Rolinger (UMD/LPS)

Christopher D. Krieger (LPS)

Alan Sussman (UMD)

Contact: tbrolin@cs.umd.edu

**LCPC 2022** 





#### Motivation

- Why target irregular memory accesses?
- Why target Partitioned Global Address Space (PGAS) programs?
- Why are compiler optimizations needed?

Irregular memory access patterns are prevalent in:

- graph analytics
- sparse linear algebra
- PDE solvers
- some machine learning applications

In these applications, performance is often memory-bound

• irregular memory accesses are **performance bottlenecks** 

Irregular memory access patterns are prevalent in:

- graph analytics
- sparse linear algebra
- PDE solvers
- some machine learning applications

In these applications, performance is often memory-bound

data is large, sparse and unstructured

irregular memory accesses are performance bottlenecks

Irregular memory access patterns are prevalent in:

- graph analytics
- sparse linear algebra
- PDE solvers
- some machine learning applications

In these applications, performance is often memory-bound

• irregular memory accesses are **performance bottlenecks** 

### Partitioned Global Address Space (PGAS)

- Presents a view of a distributed memory system that resembles a single shared address space
- Each node "owns" a partition of the address space, so there is still the concept of **local vs remote memory**
- Provides one-sided communication (puts/gets), allowing for languages/libraries that implement the PGAS model to hide data distribution/communication details





Distributed Memory e.g. MPI

Shared Memory e.g. OpenMP





PGAS

```
1 #pragma omp parallel for
2 for(int i = 0; i < numrows; i++) {
3 double accum = 0.0;
4 for(int k = rowPtr[i]; k < rowPtr[i+1]; k++) {
5 accum += values[k] * x[col_idx[k]];
6 }
7 b[i] = accum;
8 }
```



```
1 #pragma omp parallel for
2 for(int i = 0; i < numrows; i++) {
3 double accum = 0.0;
4 for(int k = rowPtr[i]; k < rowPtr[i+1]; k++)
5 accum += values[k] * x[col_idx[k]];
6 }
7 b[i] = accum;
8 }
```

#### Sparse Matrix Vector Multiply (SpMV) MPI+OpenMP distributed- and shared-memory

```
int *recvCounts = (int*)malloc(sizeof(int)*numRanks);
2 int *displs = (int*)malloc(sizeof(int)*numRanks);
3 MPI_Allgather(&localNumRows, 1, MPI_INT, recvCounts, 1, MPI_INT, MPI_COMM_WORLD);
4 recv_counts[rank] = localNumRows;
5 int currOffset = 0;
6 for (int i = 0; i < numRanks; i++) {
    displs[i] = currOffset;
    currOffset += recvCounts[i];
8
9 }
10
11 MPI_Allgatherv(locX, localNumRows, MPI_DOUBLE, x, recvCounts, displs, MPI_DOUBLE, MPI_COMM_WORLD);
12 #pragma omp parallel for
13 for (int i = 0; i < \text{localNumRows}; i++) {
    double accum = 0.0;
14
    for (int k = rowPtr[i]; k < rowPtr[i+1]; k++) {</pre>
15
      accum += values[k] * x[col_idx[k]];
16
17
    locB[i] = accum;
18
19 }
```

```
1 #pragma omp parallel for
<sup>2</sup> for(int i = 0; i < numrows; i++) {
    double accum = 0.0;
3
    for(int k = rowPtr[i]; k < rowPtr[i+1]; k++)</pre>
4
      accum += values[k] * x[col_idx[k]];
5
6
    b[i] = accum;
7
8 }
```

9

16 17

18

#### Sparse Matrix Vector Multiply (SpMV) **MPI+OpenMP distributed- and shared-memory**

int \*recvCounts = (int\*)malloc(sizeof(int)\*numRanks); 2 int \*displs = (int\*)malloc(sizeof(int)\*numRanks); 3 MPI\_Allgather(&localNumRows, 1, MPI\_INT, recvCounts, 1, MPI\_INT, MPI\_COMM\_WORLD); 4 recv\_counts[rank] = localNumRows: 5 int currOffset = 0; 6 for (int i = 0; i < numRanks; i++) { displs[i] = currOffset; currOffset += recvCounts[i]; 1 MPI\_Allgatherv(locX, localNumRows, MPI\_DOUBLE, x, recvCounts, displs, MPI\_DOUBLE, MPI\_COMM\_WORLD) 12 #pragma omp parallel 101 13 for (int i = 0; i < localNumRows; i++) { double accum = 0.0; for (int k = rowPtr[i]; k < rowPtr[i+1]; k++)</pre> accum += values[k] \* x[col\_idx[k]]; locB[i] = accum; 19 } **MPI\_Allgatherv**  $\rightarrow$  full replication of x on every process, since we don't know which remote values we need

#### Sparse Matrix Vector Multiply (SpMV) MPI+OpenMP distributed- and shared-memory

```
1 #pragma omp parallel for
                                                                                    int *recvCounts = (int*)malloc(sizeof(int)*numRanks);
                                                                                   2 int *displs = (int*)malloc(sizeof(int)*numRanks);
<sup>2</sup> for(int i = 0; i < numrows; i++) {
                                                                                   3 MPI_Allgather(&localNumRows, 1, MPI_INT, recvCounts, 1, MPI_INT, MPI_COMM_WORLD);
                                                                                    4 recv_counts[rank] = localNumRows:
      double accum = 0.0;
3
                                                                                    5 int currOffset = 0;
                                                                                   6 for (int i = 0; i < numRanks; i++) {
      for(int k = rowPtr[i]; k < rowPtr[i+1]; k++)</pre>
4
                                                                                       displs[i] = currOffset;
         accum += values[k] * x[col idx[k]];
                                                                                       currOffset += recvCounts[i];
                                                                                    Q
5
                                                                                   9
6
                                                                                   10
                                                                                   11 MPI_Allgatherv(locX, localNumRows, MPI_DOUBLE, x, recvCounts, displs, MPI_DOUBLE, MPI_COMM_WORLD);
      b[i] = accum;
7
                                                                                   12 #pragma omp parallel for
8 }
                                                                                   13 for (int i = 0; i < localNumRows; i++) {
                                                                                       double accum = 0.0;
                                                                                       for (int k = rowPtr[i]; k < rowPtr[i+1]; k++) {</pre>
                                                                                         accum += values[k] * x[col_idx[k]];
                                                                                   16
                                                                                   17
                                                                                       locB[i] = accum;
                                                                                   18
```

19 }

```
Sparse Matrix Vector Multiply (SpMV)
```

Chapel + Partitioned Global Address Space (PGAS)

```
1 forall row in Rows {
2   const id = row.id;
3   var accum : real = 0;
4   for k in row.columnOffset {
5     accum += values[k] * x[col_idx[k]];
6   }
7   b[id] = accum;
8 }
```

#### Sparse Matrix Vector Multiply (SpMV) MPI+OpenMP distributed- and shared-memory









#### High productivity does not always lead to high performance



#### High productivity does not always lead to high performance



#### Just a simple transformation...

```
1 if doOptimization {
    if doInspector(x) || doInspector(col idx) {
       inspectorPreamble(x);
       forall loc in Locales do on loc {
         ref replD = x.replArray.replD;
         ref arrDom = x.replArray.arrDom;
         const ref indices = Rows.localSubdomain();
         for i in indices {
           ref row = Rows[i]:
           const id = row.id;
10
           var accum : real = 0;
11
           for k in row.columnOffset {
12
             inspectAccess(replD, arrDom, col idx[k]);
13
14
15
16
       sortReplIndices(x);
17
       inspectorOff(x);
18
19
       inspectorOff(col idx);
20
     executorPreamble(x):
21
     forall row in Rows with (ref arrDom = x.replArray.arrDom.
22
                              ref replArr = x.replArray.replArr) {
23
       const id = row.id;
24
       var accum : real = 0;
25
       for k in row.columnOffset {
26
        accum += values[k] * executeAccess(arrDom, replArr, x, col idx[k]);
27
28
       b[id] = accum:
29
30
31 else {
     forall row in Rows {
32
       const id = row.id;
33
       var accum : real = 0:
34
       for k in row.columnOffset {
35
        accum += values[k] * x[col idx[k]];
36
37
38
       b[id] = accum:
39
40
```

```
1 if doOptimization {
    if doInspector(x) || doInspector(col idx) {
      inspectorPreamble(x);
       forall loc in Locales do on loc {
        ref replD = x.replArray.replD;
        ref arrDom = x.replArray.arrDom;
        const ref indices = Rows.localSubdomain();
         for i in indices {
          ref row = Rows[i]:
          const id = row.id:
10
          var accum : real = 0;
11
          for k in row.columnOffset {
12
             inspectAccess(replD, arrDom, col idx[k]);
13
14
15
16
      sortReplIndices(x);
17
      inspectorOff(x);
18
      inspectorOff(col idx);
19
20
     executorPreamble(x);
21
     forall row in Rows with (ref arrDom = x.replArray.arrDom,
22
                              ref replArr = x.replArray.replArr) {
23
      const id = row.id:
24
      var accum : real = 0:
25
       for k in row.columnOffset {
26
        accum += values[k] * executeAccess(arrDom, replArr, x, col idx[k]);
27
28
       b[id] = accum;
29
30
31 else
     forall row in Rows
32
      const id = row.id;
33
      var accum : real = 0;
34
      for k in row.columnOffset {
35
        accum += values[k] * x[col idx[k]];
36
37
      b[id] = accum;
38
39
40
```

#### NAS-CG (Conjugate Gradient) Problem Size D (73 million non-zeros)



#### Manual optimizations can drastically improve performance

**NAS-CG (Conjugate Gradient) Problem Size D (73 million non-zeros)** ref replD = x.replArray.replD; ref arrDom = x.replArray.arrDom; for i in In this talk: for k i Automatic optimization of irregular memory access patterns 16 18 new code Preserve the high productivity of the PGAS model while forall row in achieving better performance for k in ro accum += values[k] \* executeAccess(arrDom, replArr, x, col idx[k]); 32 16 # nodes forall row in Rows { **Manual optimizations can** var accum : real = 0; for k in row.columnOffset { accum += values[k] \* x[col idx[k]]; drastically improve performance 39 40

### Outline

- Brief introduction Chapel
- Optimization: selective data replication
- Implementation within compiler:
  - Code transformations
  - Static analysis
- Performance evaluation:
  - NAS-CG
  - PageRank
- Future work and conclusions

- Chapel Parallel Programming Language
  - high-level language that implements the PGAS model
  - designed for productive parallel computing at scale
- Terminology
  - task: set of computations that can be executed in parallel
  - locale: machine resources on which tasks execute (i.e., a node in a cluster)

- Chapel Parallel Programming Language
  - high-level language that implements the PGAS model
  - designed for productive parallel computing at scale
- Terminology
  - task: set of computations that can be executed in parallel
  - locale: machine resources on which tasks execute (i.e., a node in a cluster)

```
1 var D : domain(1) = {0..5};
```

- 2 var data : [D] int;
- 3 data[0] = 1;
- domain **D** with indices 0 through 5
- array data defined over D

- Chapel Parallel Programming Language
  - high-level language that implements the PGAS model
  - designed for productive parallel computing at scale
- Terminology
  - task: set of computations that can be executed in parallel
  - locale: machine resources on which tasks execute (i.e., a node in a cluster)

```
1 var D = newBlockDom({0..#16});
2 var A : [D] int;
3 // or
4 var A = newBlockArr({0..#16}, int);
```

- block distributed domain D
- access to remote array elements looks no different than accessing local elements

- Chapel Parallel Programming Language
  - high-level language that implements the PGAS model
  - designed for productive parallel computing at scale
- Terminology
  - task: set of computations that can be executed in parallel
  - locale: machine resources on which tasks execute (i.e., a node in a cluster)



- forall loop provides data parallelism
- Rows is a block distributed array
- task assigned the iteration for a given row will execute on the locale where row is
  - *locale affinity* defined implicitly by the loop iterand 24

- We focus on accesses of the form **A[B[i]]** in **forall** loops
  - A is a distributed array and the values in **B** are not known until runtime
  - more complex access patterns are supported (see paper for details)

- We focus on accesses of the form **A[B[i]]** in **forall** loops
  - A is a distributed array and the values in **B** are not known until runtime
  - more complex access patterns are supported (see paper for details)
- Goal: replicate remotely accessed elements of A so they can be accessed locally in the forall
  - **inspector:** runtime analysis that determines remote accesses
  - executor: optimized version of the forall that redirects remote accesses to the replicated copies
  - both generated by the **compiler** without user intervention

- We focus on accesses of the form **A[B[i]]** in **forall** loops
  - A is a distributed array and the values in **B** are not known until runtime
  - more complex access patterns are supported (see paper for details)
- Goal: replicate remotely accessed elements of A so they can be accessed locally in the forall
  - **inspector:** runtime analysis that determines remote accesses
  - **executor:** optimized version of the **forall** that redirects remote accesses to the replicated copies
  - both generated by the **compiler** without user intervention
- Requirements:
  - access-of-interest is **read-only** in the **forall** loop

The inspector-executor technique has been around for a long time

We **leverage** the prior work on inspector-executors and **adapt** them for **Chapel+PGAS** 

- we work within Chapel's compiler, before it lowers the code to LLVM
- requires different approaches to statically reason about domains/arrays, forall loops and implicit communication

• Requirements:

а

access-of-interest is read-only in the forall loop

### **Compiler: Code Transformations**

```
1 forall i in B.domain {
2 C[i] += A[B[i]];
3 }
```



```
1 if doOptimization {
     if doInspector(A, B) {
2
       inspectorPreamble(A);
3
       forall i in inspectorIter(B.domain) {
4
         inspectAccess(A, B[i]);
5
       }
6
       inspectorOff(A, B);
7
    }
8
    executorPreamble(A);
9
     forall i in B.domain {
10
       C[i] += executeAccess(A, B[i]);
11
    }
12
13 }
14 else {
    forall i in B.domain {
15
       C[i] += A[B[i]];
16
     }
17
18 }
                                       29
```

# Compiler: Code Transformations (cont.)

**Inspector:** performs memory access analysis

 inspector should only be performed (1) the first time we encounter the loop and (2) anytime the access pattern A[B[i]] could have changed

**Executor:** executes the loop but redirects remote accesses to the replicated copies

- **executorPreamble()** initializes replicated elements of **A** with values from original array
  - we only replicate an element once, regardless of how many times it is accessed

1	<pre>if doOptimization {</pre>
2	<pre>if doInspector(A, B) {</pre>
3	<pre>inspectorPreamble(A);</pre>
4	<pre>forall i in inspectorIter(B.domain) {</pre>
5	<pre>inspectAccess(A, B[i]);</pre>
6	}
7	<pre>inspectorOff(A, B);</pre>
8	}
9	<pre>executorPreamble(A);</pre>
10	<pre>forall i in B.domain {</pre>
11	<pre>C[i] += executeAccess(A, B[i]);</pre>
12	}
13	}
14	else {
15	<pre>forall i in B.domain {</pre>
16	C[i] += A[B[i]];
17	}
18	<b>}</b> 30

### **Compiler: Static Analysis**

- Everything just discussed for code transformations only holds if the optimization **CAN** and **SHOULD** be applied:
  - optimization needs to maintain correct program results
  - optimization should improve program performance
- The goal of our static analysis is to achieve the above, all without requiring user intervention
  - no code annotations, hints, pragmas, etc.
  - just a **compiler flag** that is turned on when compiling the program

- When could code transformations produce **incorrect behavior**?
  - if the inspector is not executed when needed, the replicated copies will be out of date
  - → need to be able to statically determine when the inspector should execute

- **1. forall** must iterate over a distributed array or domain
- 2. i from A[B[i]] must be *yielded* by the loop that contains A[B[i]] and that loop must iterate over a domain/array
- **3. A**, **B**, and their **domains** cannot be modified in the **forall**

- (1) ensures the optimization can reason about which locales the task execute on
  - 1 forall i in B.domain {
    2 C[i] += A[B[i]];
    3 }

Controls locale affinity of the loop. If B.domain changes, we know that the locale affinity could change.

- **1. forall** must iterate over a distributed array or domain
- 2. i from A[B[i]] must be *yielded* by the loop that contains A[B[i]] and that loop must iterate over a domain/array
- **3. A**, **B**, and their **domains** cannot be modified in the **forall**

1 forall i in B.domain {
2 C[i] = A[B[i]];
3 }

• (3) ensures the optimization can reason about the values of i, specifically when they would change (i.e., when the array/domain changes)

- **1. forall** must iterate over a distributed array or domain
- 2. i from A[B[i]] must be *yielded* by the loop that contains A[B[i]] and that loop must iterate over a domain/array
- **3. A**, **B**, and their **domains** cannot be modified in the **forall**



#### Assume candidate access A[B[i]]

- **1. forall** must iterate over a distributed array or domain
- 2. i from A[B[i]] must be *yielded* by the loop that contains A[B[i]] and that loop must iterate over a domain/array
- **3. A**, **B**, and their **domains** cannot be modified in the **forall**

 (4) ensures that the values in B analyzed by the inspector will be the same as those used in the executor and original loop

Assume candidate access A[B[i]]

		The optimization statically identifies	ributed array or
1	forall	modifications to the arrays/domains and	d by the loop that
2	C[i]	inserts code to <b>turn on flags at runtime</b> to	p must iterate over a
3	۲ B[1]	indicate that the inspector needs to be	at he modified in
4	5	executed	iot be moulled m

• (4) ensures that the values in B analyzed by the inspector will be the same as those used in the executor and original loop

- When could code transformations produce **poor performance**?
  - if the **inspector** is executed each time the **forall** is executed

- I. forall must be nested in an outer serial loop
- **II. B** nor its **domain** can be modified within the outer loop that the **forall** is nested
- **III.** A's domain cannot be modified within the outer loop that the forall is nested in

- When could code transformations produce **poor performance**?
  - if the **inspector** is executed each time the **forall** is executed

- I. forall must be nested in an outer serial loop
- II. B nor its domain can be modified within the outer loop that the forall is nested
- III. A's domain cannot be modified within the outer loop that the forall is nested in

- (I) ensures that the **forall** is likely to be executed multiple times
  - 1 for step in 0..#numSteps {

- When could code transformations produce **poor performance**?
  - if the **inspector** is executed each time the **forall** is executed

#### Assume candidate access A[B[i]]

- I. forall must be nested in an outer serial loop
- **II. B** nor its **domain** can be modified within the outer loop that the **forall** is nested
- III. A's domain cannot be modified within the outer loop that the forall is nested in

• (II) ensures the inspector will not execute each time the forall is performed (if B or its domain is modified, the inspector must be rerun)

- When could code transformations produce **poor performance**?
  - if the **inspector** is executed each time the **forall** is executed

- I. forall must be nested in an outer serial loop
- II. B nor its domain can be modified within the outer loop that the forall is nested
- **III. A's domain** cannot be modified within the outer loop that the **forall** is nested in
- (III) is the same as (II) but for A's domain; changes to A's values do not change the access pattern A[B[i]] but changes to the domain can change where elements of A are located

- Additional analyses (see paper for details):
  - Non-affine expression analysis to identify candidate accesses beyond A[B[i]]
  - Interprocedural analysis to track modifications to arrays/domains across function calls
  - Alias analysis to track modifications to arrays/domains across to any aliases created
  - Call path analysis to identify invalid call paths to the function containing the forall and "turning off" the optimization at runtime for the invalid paths

### **Performance Evaluation**

- Applications:
  - NAS-CG (conjugate gradient)
  - PageRank (iterative SpMV-like operations)
- Systems:
  - FDR Infiniband, 20 cores per node, 512 GB of memory per node
  - Cray XC, Aries interconnect, 44 cores per node, 128 GB of memory per node
- Experiments:
  - measured runtime speed-ups achieved by optimization relative to the original Chapel code
  - includes any overhead incurred by the inspector

Name	Rows	Non-zeros	Density (%)	# of SpMVs	<
С	150k	39M	0.17	1950	
D	150k	73M	0.32	2600	
E	9M	6.6B	0.008	2600	
F	54M	55B	0.002	2600	

NAS-CG Data sets

- Each outer iteration performs 26 SpMVs
- Problem size C performs 75 outer iterations
- Other sizes perform 100 outer iterations
- No changes to A[B[i]] access pattern during CG

NAS-CG Data sets

Name	Rows	Non-zeros	Density (%)	# of SpMVs
C	150k	39M	0.17	1950
D	150k	73M	0.32	2600
E	9M	6.6B	0.008	2600
F	54M	55B	0.002	2600

#### NAS-CG Problem Size E Optimization Runtime Speed-ups



#### **NAS-CG Optimization Speed-ups**

		Cray XC				Infiniband		
Locales	C	D	E	F	C	D	E	F
2	3.2	2.8	_	_	8.9	6	357	_
4	3.6	3.4	17.5	_	15.8	10.4	345	_
8	5.7	6.2	36.7	_	115	127	364	_
16	8.6	11	22.5	_	238	330	258	270
32	6.4	8.4	34	52.3	160	240	195	165
64	4.1	4.9	16.7	25.4	NA	NA	NA	NA
geomean	5	5.5	24.1	36.4	57.3	57.5	296	211

#### Take-aways:

- "—" means not enough memory, "NA" means not enough nodes
- **high degree of data reuse** in the kernel, so the optimization performs very well
- **inspector overhead is small** due to many iterations w/o the access pattern changing
  - optimization provides larger gains on Infiniband
    - higher latency for small messages than Aries

PageRank Data sets

Name		<b>E</b>	Density (%)	Iterations	]←
webbase-2001	118M	992M	7.1e – 6	33	]
sk-2005	51M	1.9B	7.5e – 5	40	

- PageRank runs until a convergence threshold is met
- # of iterations depends on graph structure

#### PageRank: webbase-2001

		0		
Name	$ \mathbf{V} $	<b> E</b>	Density (%)	Iterations
webbase-2001	118M	992M	7.1e – 6	33
sk-2005	51M	1.9B	7.5e – 5	40

PageRank Data sets

#### **Optimization Runtime Speed-ups** speed-up over baselines 15 12 8.6 8.6 10 5.2 5



□ Cray XC ■ Infiniband

	Cray X	С	Infiniba		
Locales	webbase-2001	sk-2005	webbase-2001	sk-2005	Take-aways
2	0.88	1.2	5.2	2	• smaller
4	0.98	1.6	8.6	7.1	
8	0.97	1.3	12	6	than NA
16	0.94	1.7	9.6	5.4	• bea
32	1.3	1.4	4.5	4.2	<ul> <li>speed-u</li> </ul>
64	1.2	2.1	NA	NA	data reu
geomean	1.04	1.5	7.3	4.5	neverth

0

- speed-ups overall due to fewer iterations S-CG and less data reuse
  - cause of both the algorithm and the graphs
- ps on the Cray can be negative when the ise is low (webbase-2001)
- nevertheless, still significant speed-ups overall

#### PageRank Optimization Speed-ups

#### Future Work

- Limitations exists for this type of data replication
  - forall must execute multiple times without the memory access pattern changing
  - could use a lot of memory for the replication
  - currently limited to read-only data
- Future work: additional optimizations
  - adaptive prefetching for Chapel's remote cache
  - remote data aggregation for remote writes
  - end goal is a single framework that can apply all these optimizations automatically, deciding which one to apply considering the specific scenario

#### NAS-CG (Conjugate Gradient) Problem Size D (73 million non-zeros)



Optimization **preserves** original **high productivity** of code but achieves significantly **better performance**:

• runtimes reduced from **hours/days to minutes** 

Acknowledgements: Chapel team for providing access to the Cray system and compiler support