### Data Locality on Manycore Architectures

Duco van Amstel

Inria / Kalray - Grenoble, France

Ph.D defence - July 18th, 2016









# Introduction

## The Proverbial Memory Wall

Observations:

- Available computational power keeps growing exponentially
- Bandwidth of memory interfaces increases slower

Has been the case since the '90s and is called the Memory-Wall

### Definition

For a given executable / partial code :

Operational Intensity (OI) =  $\frac{\text{Number of instructions}}{\text{Number of memory operations}}$ 



## **Research** topics

### Working hypothesis

Local target memory of limited size communicating with distant memory of infinite size through memory operations

Idea: Improve data reuses and reduce IO operations of frequently executed code → Improve data locality of programs



### Overview



### Topics

**1** Modeling memory usage and IO costs

Memory usage & IO model (Chapter 1)

## Code representation

What we want:

- a compact representation
- 2 data reuses and no other dependencies
- straightforward evaluation of memory usage for any part of the representation

We refer to it as the memory-use graph.



Figure: Dataflow diagram







### Memory-use graph construction

We observe:

- inter-iteration edges
- intra-iteration edges

Some transformations need to be performed





### Memory-use graph construction

Steps:

1 Indicate the size of reused data on edges





Steps:

- Indicate the size of reused data on edges
- 2 Transform self-reuses into state data





Steps:

- Indicate the size of reused data on edges
- 2 Transform self-reuses into state data
- **3** Transform other inter-iteration reuses into **state and** intra-iteration reuse data





## Memory-use graph construction

Steps:

- Indicate the size of reused data on edges
- 2 Transform self-reuses into state data
- Transform other inter-iteration reuses into state and intra-iteration reuse data
- Add internal computation requirements to account for hidden intermediary values





```
for (int j = 1; j < N - 1; j + +) {
    /*S1*/A[j] = (A[j-2] + A[j-1] + A[j] + A[j+1] + A[j+2]) / 5;
    /*S2*/ B[j] = B[j] + A[j] * C[j];
    /*S3*/B[i] = B[i] - (B[i-1] - B[i]) * C[i-1];
  }
S2 :
                                                            S1
 mov [R15 + RSI*4 - 4], R8D ; Save C[j-1]
                                                                 C: 0
 mov R8D, [R15 + RSI*4] ;Load C[j]
 mov EBX, [R14 + RSI*4]; Load B[j]
 mov EAX, [R13 + RSI*4]; Load A[j]
                                                                S: 1
                                                           S2
 mul EAX, R8D ; Compute on RA
                                                                 C: 1
 add EBX, EAX ; Compute on RB
 mov [R14 + RSI*4], EBX ; Save B[j]
                                                                S: 1
                                                           S3
Registers used: 3
                                                                 C \cdot 1
```





## Evaluation of memory usage



Consider a partial execution:

$$\begin{array}{c} B_1 \rightarrow C_1 \rightarrow D_1 \rightarrow B_2 \rightarrow C_2 \rightarrow D_2 \rightarrow \\ B_3 \rightarrow C_3 \rightarrow D_3 \end{array}$$

Compute the size of data stored in memory in and between nodes:

The maximum is the **memory usage** of the partial execution, *i*.e 12

#### Figure: Iterated memory-use graph



## IO cost computation

For an arbitrary piece of code:

```
\textbf{memory usage} \leq \textbf{memory size}
```

Internal reuses do not require IO

Divide evaluated code in such pieces

- Only external reuses generate IO
- IO cost equal to cut between pieces



### Topics

2 Generalized tiling

Data locality



### Representing the iteration space



LB (Loop-body) = {S1,S2,S3,S4}



### Representing the iteration space



## Data locality optimizations

#### Code example:







## A missing piece







## A missing piece



## Comparison of tiling methods

#### Tile regularity

classical tiling



Schedule within a tile

4

Ordering of tiles

## Comparison of tiling methods



Schedule within a tile

Tile semi-regularity

generalized tiling



Defining a tiling solution means:

**1** Linearization / schedule the memory-use graph



ts
U.
ž
ē
at
st

				j	
	S1	S1	S1	S1	
	S3	S3	S3	S3	
	S2	S2	S2	S2	
ments	S4	S4	S4	S4	
state	S5	S5	S5	S5	
	,			🕤 KALRAY	l

Defining a tiling solution means:

- 1 Linearization / schedule the memory-use graph
- 2 Choosing points in the schedule where to place tile limits
- 3 Specifying a width for each tile





Defining a tiling solution means:

- 1 Linearization / schedule the memory-use graph
- 2 Choosing points in the schedule where to place tile limits
- 3 Specifying a width for each tile
- ... and there are multiple solutions





Defining a tiling solution means:

- 1 Linearization / schedule the memory-use graph
- 2 Choosing points in the schedule where to place tile limits
- 3 Specifying a width for each tile
- ... and there are multiple solutions





Objective Reducing the reach of reuses



Edges ordered by descending weight: green, red, orange, blue, brown



Objective Reducing the reach of reuses



Edges ordered by descending weight: green, red, orange, blue, brown Blue zones have been frozen



Objective Reducing the reach of reuses



Edges ordered by descending weight: green, red, orange, blue, brown Blue zones have been frozen



#### Objective Reducing the reach of reuses



Edges ordered by descending weight: green, red, orange, blue, brown Blue zones have been frozen







### Heavy-edge / Greedy

Perform heavy-edge scheduling and greedily compute locally optimal tiles





### Heavy-edge / Greedy

Perform heavy-edge scheduling and greedily compute locally optimal tiles





### Heavy-edge / Greedy

Perform heavy-edge scheduling and greedily compute locally optimal tiles

### Tile-aware / Heavy-edge

Contract edges one-by-one and create or expand tiles to include them; Roll back contraction if not possible





### Heavy-edge / Greedy

Perform heavy-edge scheduling and greedily compute locally optimal tiles

### Tile-aware / Heavy-edge

Contract edges one-by-one and create or expand tiles to include them; Roll back contraction if not possible

#### Tile-aware / Conservative

Contract edges one-by-one and create tiles to include them; Roll back contraction if not possible





## Topics

3 Applications in loop and dataflow optimization



#### Code generation

					J
	S1	S1	S1	S1	S1
	S2	S2	S2	S2	S2
nents	S3	S3	S3	S3	S3
	S4	S4	S4	S4	S4
staten	S5	S5	S5	S5	S5

Target tiling: Schedule S1, S3, S2, S4, S5 Tile 1 {S1, S3} - width 3 Tile 2 {S2, S4, S5} - width 2



#### Code generation

					$\xrightarrow{J}$
	S1	S1	S1	S1	S1
nents	S3	<b>S</b> 3	S3	<b>S</b> 3	S3
	S2	S2	S2	S2	S2
	S4	S4	S4	S4	S4
staten	↓ S5	S5	S5	S5	S5

Target tiling: Schedule S1, S3, S2, S4, S5 Tile 1 {S1, S3} - width 3 Tile 2 {S2, S4, S5} - width 2

Transformations to be applied:

Rescheduling



#### Code generation



Target tiling: Schedule S1, S3, S2, S4, S5 Tile 1 {S1, S3} - width 3 Tile 2 {S2, S4, S5} - width 2

- Rescheduling
- 2 Loop fission

#### Code generation



Target tiling: Schedule S1, S3, S2, S4, S5 Tile 1 {S1, S3} - width 3 Tile 2 {S2, S4, S5} - width 2

- Rescheduling
- 2 Loop fission
- 3 Loop unrolling



#### Code generation



Target tiling: Schedule S1, S3, S2, S4, S5 Tile 1 {S1, S3} - width 3 Tile 2 {S2, S4, S5} - width 2

- Rescheduling
- 2 Loop fission
- 3 Loop unrolling
- 4 Rescheduling



#### Code generation



Target tiling: Schedule S1, S3, S2, S4, S5 Tile 1 {S1, S3} - width 3 Tile 2 {S2, S4, S5} - width 2

Transformations to be applied:

- Rescheduling
- 2 Loop fission
- 3 Loop unrolling
- 4 Rescheduling

Prototyping in the Open64 compiler suggests the impact of generalized tiling on load/store operations could be significant.

# Implementation

LLVM



## Generalized dataflow tiling

#### Target

Static & cyclo-static dataflow languages

Prototyping was done in Streamlt:

- Usage of preprocessed StreamIt benchmarks
- Rescheduling of actors & execution scaling
- Simulation of cache-behaviour / evaluation of cache-misses



## Prototype results





## Topics

4 Performance debugging with convex graph partitioning



### Overview

### Performance debugging

Analysis of the execution of code in order to identify performance bottle-necks



### Convex partitioning



... is a convex partitioning



Non-convex



### Partition schemes









*k*-way partition scheme

### Max-distance criterion

**Idea:** For each vertex compute longest path originating at a source and terminating at a sink

- Source path ≤ Sink path
   → source partition
- 2 Source path > Sink path → sink partition





## Max-distance criterion

**Idea:** For each vertex compute longest path originating at a source and terminating at a sink

- Source path ≤ Sink path → source partition
- 2 Source path > Sink path → sink partition





## **Experimental results**



# Conclusion & perspectives

### Overview



### Contributions

### **Operational intensity**

Model to compute a close approximation of its value for arbitrary code

### Generalized tiling

Formalization of an optimization problem and proposal of several heuristics to compute solutions

### GraphUtilities library

Multi-threaded library for directed graph manipulation: state-of-the-art reachability queries, convex partitioning & associated metrics ( $\sim$  5k lines)



### Future work

- Stabilize / rewrite LLVM implementation of generalized tiling
- Pursue implementation of dataflow tiling
- Expand generalized tiling formalization for more flexibility
- Consider parallel execution of tiles
- Continue improving hierarachical partition schemes

• • • •







# **Backup slides**



## Generalized tiling in the polyhedral model?



Memory-use graph

	~				j
	S1	S1	S1	S1	S1
nents	S2	S2	S2	S2	S2
	S3	S3	S3	S3	S3
	S4	S4	S4	S4	S4
stater	S5	S5	S5	S5	S5



# Generalized tiling in the polyhedral model?



Need to create uniform reuse pattern for whole iteration space.

 $\rightarrow$  Duplicate all existing reuse edges for each point in iteration space?



# Generalized tiling in the polyhedral model?



Theoretically possible but... more false reuse edges then true ones. Very constrained optimization opportunities.





What if we consider another schedule?

$$\begin{array}{c} B_1 \rightarrow B_2 \rightarrow B_3 \rightarrow C_1 \rightarrow C_2 \rightarrow C_3 \rightarrow \\ D_1 \rightarrow B_3 \rightarrow D_3 \end{array}$$

Amount of data stored in and between nodes changes...

...and so does the resulting **memory usage**: 16

Figure: Iterated memory-use graph



#### 38 / 34

### Transition to $\Sigma$ -C

#### Context

Dataflow languages on a manycore processor

Illustration with the MPPA processor family by Kalray:

- 256 computational cores, 32 management cores
- 16 computation clusters of 17 cores each
- 4 IO clusters of 4 cores each
- Clusters each have own local memory
- Clusters linked through a Network-on-Chip

Distributed architecture

 $\rightarrow$  Dataflow programs must be mapped to ressources



### Communications over the Network-on-Chip



Contributions:

- Network traffic model
- Redesign of bandwidth limiter (patent filed)
- NoC simulator
   (~ 1.5k lines)

Remaining challenges:

- Map actors to clusters with size constraints
- Combine memory and network models



# Reachability querying

#### Goal

For a directed-acyclic graph and two of its vertices decide whether there exists a path from one to the other.



# Reachability querying

#### Goal

For a directed-acyclic graph and two of its vertices decide whether there exists a path from one to the other.

Existing approaches:

- Naïve depth-first search (prohibitive for many queries)
- Compute full transitive-closure (prohibitive for large graphs)
- Pre-compute information reducing search time (many variants)



# Reachability querying

#### Goal

For a directed-acyclic graph and two of its vertices decide whether there exists a path from one to the other.

#### Existing approaches:

- Naïve depth-first search (prohibitive for many queries)
- Compute full transitive-closure (prohibitive for large graphs)
- Pre-compute information reducing search time (many variants)

### Contribution

- Pre-indexation method based on reverse post-order traversals
- **2** Experimental results improving on state-of-the-art methods

Drawback: Index is very hard to maintain during graph partitioning

### Memory usage & IO cost



## Partition tree balance

#### Convexify partition trees

Benchmark	Min	Max	Std.Dev
adi	15	19	0.19
bicg	9	13	0.17
durbin	10	12	0.24
fdtd-2d	15	20	0.20
gemm	12	17	0.27
gesummv	11	13	0.13
jacobi-2d	15	20	0.26
ludcmp	11	16	0.15
mvt	12	14	0.15
seidel-2d	12	17	0.27
syr2k	13	17	0.24
syrk	12	15	0.12
trisolv	9	11	0.25
trmm	14	17	0.22

#### Max-distance partition trees

Benchmark	Min	Max	Std.Dev
adi	4	30	13.71
bicg	3	21	10.60
durbin	2	29	19.97
fdtd-2d	2	40	22.89
gemm	3	27	15.17
gesummv	4	23	11.06
jacobi-2d	2	35	18.74
ludcmp	2	42	23.71
mvt	4	27	13.56
seidel-2d	4	46	39.66
syr2k	3	23	14.05
syrk	2	23	12.58
trisolv	2	26	19.90
trmm	3	29	17.91

