# Optimizations for intensive signal processing applications on Systems-on-Chip

Calin Glitia

DaRT   lifL   Université Lille1 Sciences et Technologies   Cnrs   INRIA LILLE - NORD-EUROPE

September 6, 2010

## Detection systems



## Multimedia
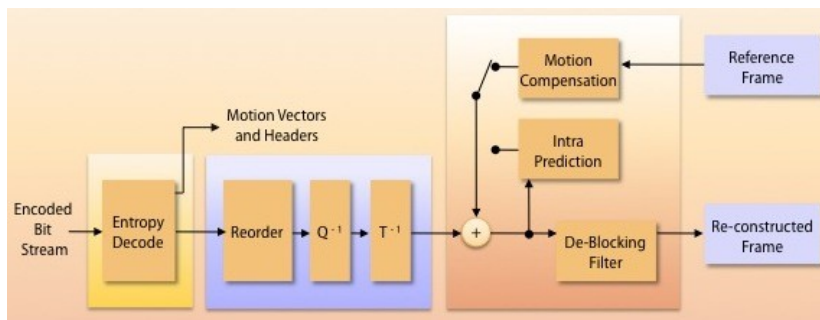
## Detection systems



## Multimedia



- Repetitive computations
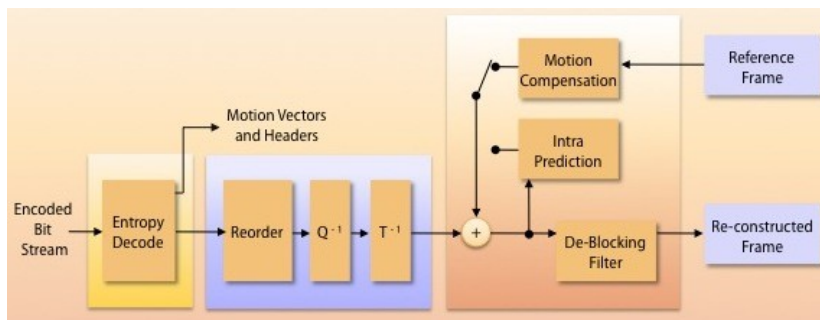- Considerable amount of data $\Rightarrow$ Multi-dimensional arrays

# Modular decomposition



## Data-flow oriented modeling

- **Logical parallelism**
  1. Task parallelism and pipeline
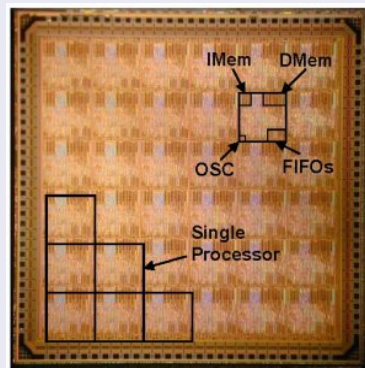  2. Data parallelism

# Modular decomposition



## Data-flow oriented modeling

- Logical parallelism
  1. Task parallelism and pipeline
  2. Data parallelism
- Complexity:
  - Elementary functions assemblage
  - Complex accesses at the data structures

Systems-on-Chip:
- Increase in the integration capacity
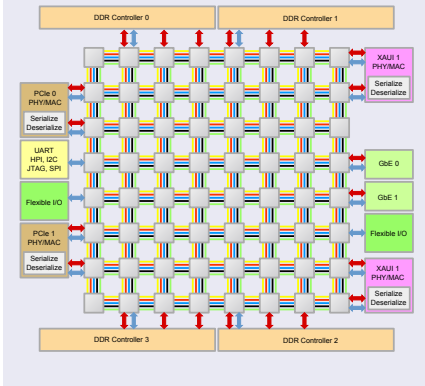- Multiprocessors

Multiprocessor SoC

# Execution platforms

Systems-on-Chip:
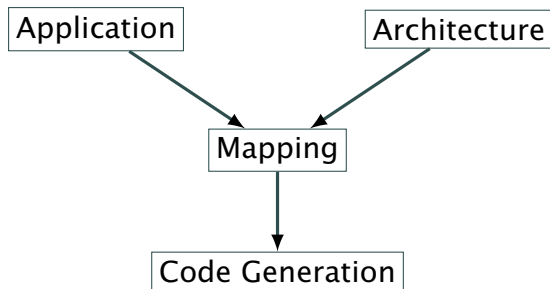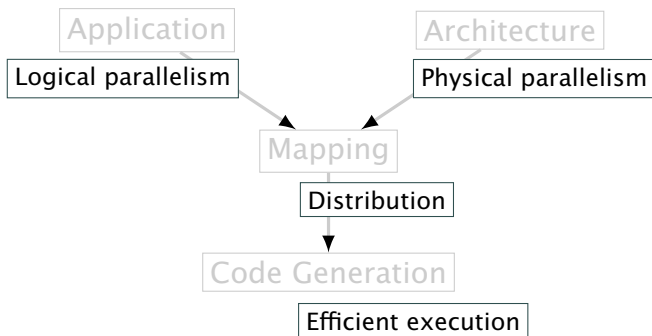- Increase in the integration capacity
- Multiprocessors

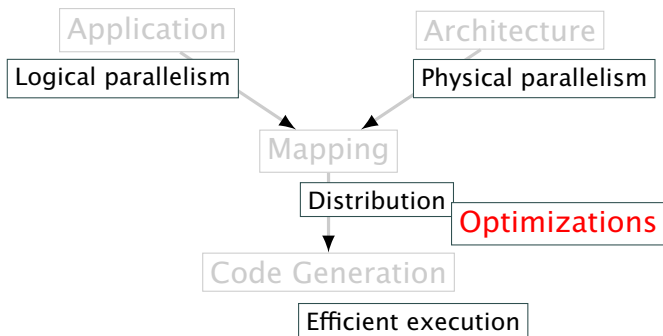Architecture models :
- Repetitive topologies
- Physical parallelism



Multiprocessor SoC

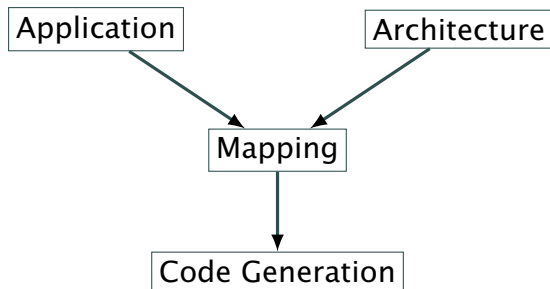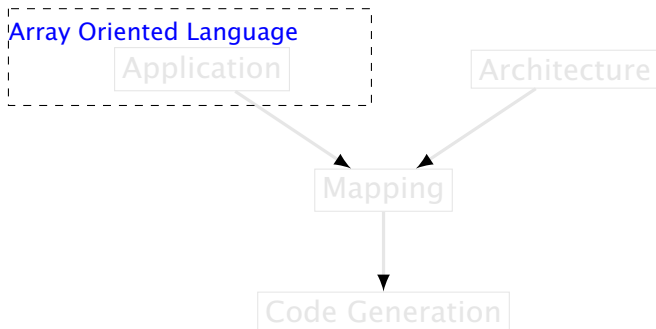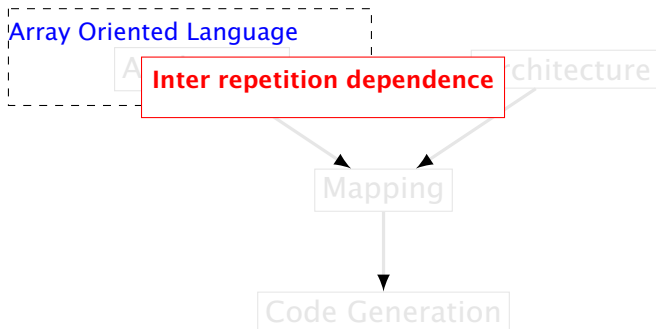Application — Logical parallelism

Architecture — Physical parallelism

Mapping — Distribution

Code Generation — Efficient execution

Array Oriented Language
Application

Architecture

Mapping

Code Generation

Modeling and Analysis of Real-Time Embedded Systems

Array Oriented Language

A... ...chitecture

**Inter repetition dependence**

Mapping

Code Generation

# Some of the existing languages

## Data-Flow

1 Synchronous Data Flow

## Data-Flow

1 Synchronous Data Flow
2 Extensions:
   - Cyclo-Static Data Flow
   - Multi-dimensional Synchronous Data Flow, Windowed Synchronous Data Flow
   - Boolean Data Flow, Dynamic Data Flow

# Some of the existing languages

## Data-Flow

1. Synchronous Data Flow
2. Extensions:
   - Cyclo-Static Data Flow
   - Multi-dimensional Synchronous Data Flow, Windowed Synchronous Data Flow
   - Boolean Data Flow, Dynamic Data Flow
3. **Array-OL**

# Some of the existing languages

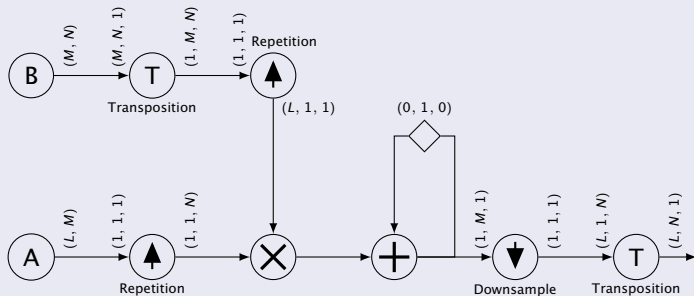## Data-Flow

1. Synchronous Data Flow
2. Extensions:
   - Cyclo-Static Data Flow
   - Multi-dimensional Synchronous Data Flow, Windowed Synchronous Data Flow
   - Boolean Data Flow, Dynamic Data Flow
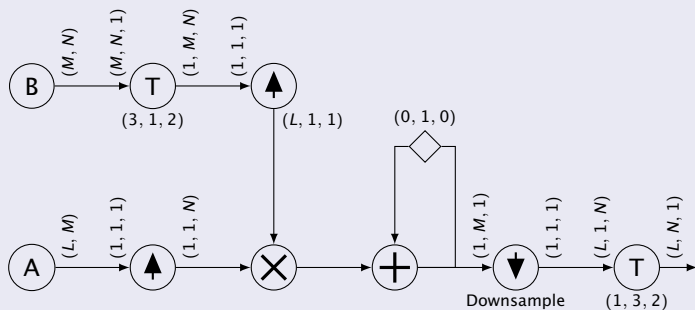3. **Array-OL**

## Functional languages

- Alpha: polyhedron, recurrence equations
- Sisal, Single Assignment C

## Matrix multiplication
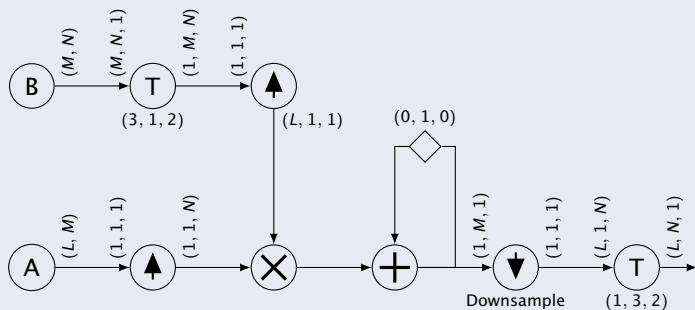
## Matrix multiplication



- Data-flow graph: actors consuming/producing MultiDimensional-tokens
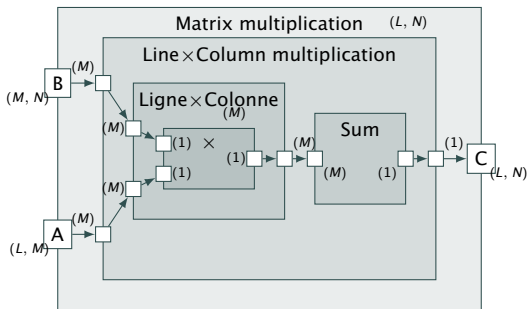- Static analysis/schedule

## Matrix multiplication



- Data-flow graph: actors consuming/producing MultiDimensional-tokens
- Static analysis/schedule
- Limitations: Multiple data consumptions
- Extensions

# Array Oriented Language – principles
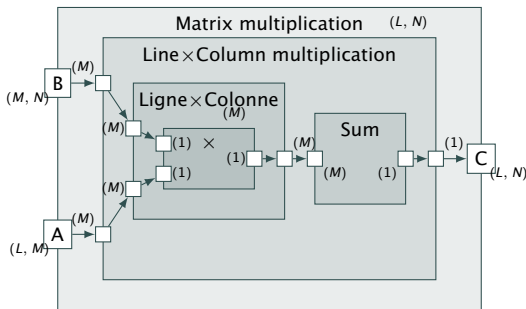
## Similarities with other data-flow languages

- Hierarchical decomposition – task parallelism
- Data-flow oriented formalism – multi-dimensional data structures

## Similarities with other data-flow languages

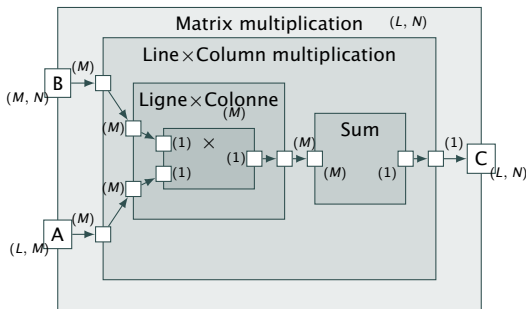- Hierarchical decomposition – task parallelism
- Data-flow oriented formalism – multi-dimensional data structures



- **Explicit data parallelism : data-parallel repetitions ("loops")**
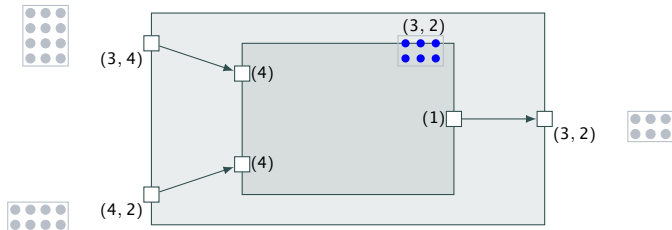- **Uniform paving by sub-arrays**

# Array Oriented Language – principles

## Similarities with other data-flow languages

- Hierarchical decomposition – task parallelism
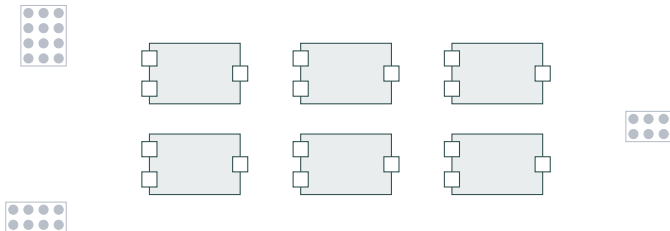- Data-flow oriented formalism – multi-dimensional data structures



- **Explicit data parallelism : data-parallel repetitions ("loops")**
- **Uniform paving by sub-arrays**
- **Space and time mixed as dimensions of the data structures**
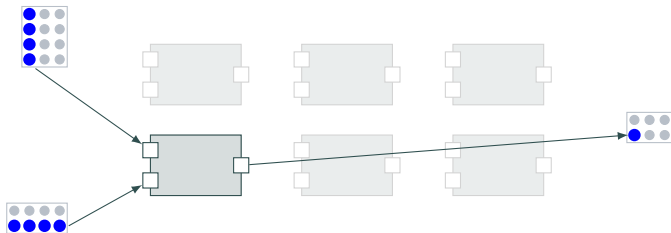
## Matrix multiplication: $N = 3, M = 4, L = 2$



Matrix multiplication

## Matrix multiplication: $N = 3, M = 4, L = 2$
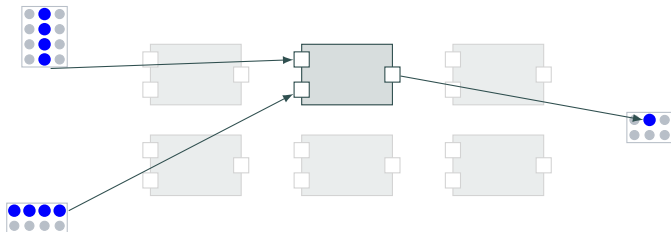


Line×Column multiplications

# Data-parallel repetitions

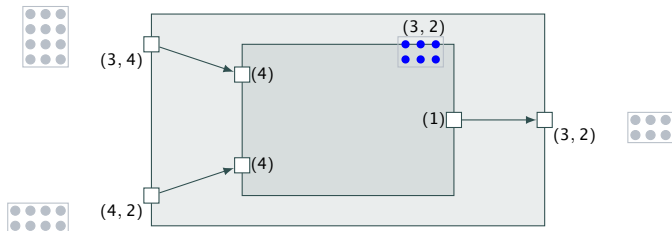## Matrix multiplication: $N = 3, M = 4, L = 2$



uniformly spaced input/output patterns
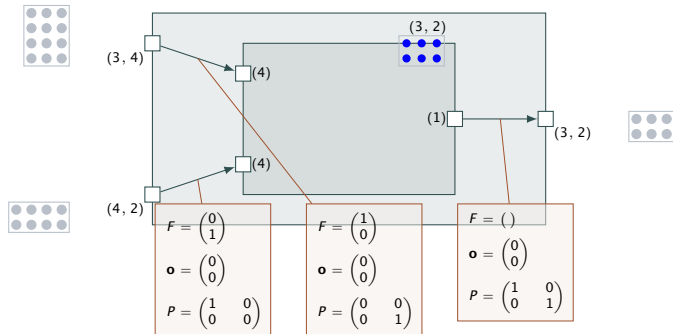
## Matrix multiplication: $N = 3, M = 4, L = 2$



DATA-PARALLEL instances

## Matrix multiplication: $N = 3, M = 4, L = 2$



Compact representation: repetition space
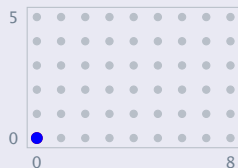
## Matrix multiplication: $N = 3, M = 4, L = 2$



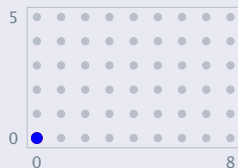Tilers – links between input and output patterns
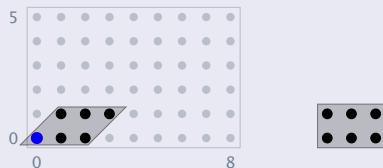
## Tiler

- **o**: origin of the reference pattern

## Tiler

- **o**: origin of the reference pattern

## Tiler

- **o**: origin of the reference pattern
- *F*: Fitting matrix – the shape of the tiles in the array

# Uniform sub-array accesses (patterns)

## Tiler

- **o**: origin of the reference pattern
- $F$: Fitting matrix – the shape of the tiles in the array
- $P$: Paving matrix – uniform spacing of the tiles



## Formal specification:

$$\mathbf{o} + (P\, F) \cdot \begin{pmatrix} \mathbf{r} \\ \mathbf{i} \end{pmatrix} \bmod \mathbf{s}_{\text{array}}$$

# Uniform sub-array accesses (patterns)

## Tiler

- **o**: origin of the reference pattern
- $F$: Fitting matrix – the shape of the tiles in the array
- $P$: Paving matrix – uniform spacing of the tiles



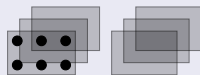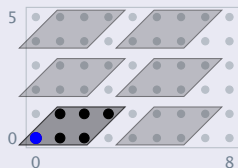## Formal specification:

$$\mathbf{o} + (P\ F) \cdot \begin{pmatrix} \mathbf{r} \\ \mathbf{i} \end{pmatrix} \bmod \mathbf{s}_{array}$$

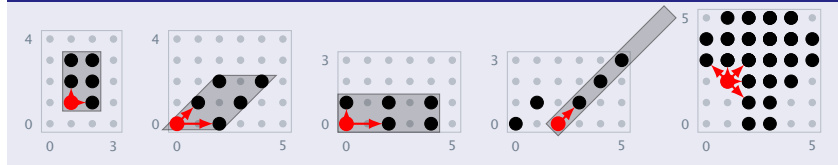# Common motif-based accesses

## Pattern examples

# Common motif-based accesses

## Pattern examples



## Paving example



$$F = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \qquad \mathbf{s}_{\text{pattern}} = \begin{pmatrix} 5 \\ 3 \end{pmatrix}$$

$$\mathbf{o} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \qquad \mathbf{s}_{\text{array}} = \begin{pmatrix} 10 \\ 5 \end{pmatrix}$$

$$P = \begin{pmatrix} 0 & 3 \\ 1 & 0 \end{pmatrix} \qquad \mathbf{s}_{\text{repetition}} = \begin{pmatrix} 3 \\ 3 \end{pmatrix}$$

## Specification

- Data-flow oriented visual formalism
- Express the regularity of computations/data accesses
- Exploit the parallelism

## Specification

- Data-flow oriented visual formalism
- Express the regularity of computations/data accesses
- Exploit the parallelism
- **Rules that allow static analysis**

## Specification

- Data-flow oriented visual formalism
- Express the regularity of computations/data accesses
- Exploit the parallelism
- **Rules that allow static analysis**

## Limitations

- Numerical values for the multidimensional spaces/accesses

## Specification

- Data-flow oriented visual formalism
- Express the regularity of computations/data accesses
- Exploit the parallelism
- **Rules that allow static analysis**

## Limitations

- Numerical values for the multidimensional spaces/accesses
- Cycles not allowed in the dependence graph

## Specification

- Data-flow oriented visual formalism
- Express the regularity of computations/data accesses
- Exploit the parallelism
- **Rules that allow static analysis**

## Limitations

- Numerical values for the multidimensional spaces/accesses
- Cycles not allowed in the dependence graph
- **Extension: inter-repetition dependences**

## State construction

- *Transfer data between different instances of the same repetition.*

## State construction

- *Transfer data between different instances of the same repetition.*
- Examples: Sum, **Integrate**
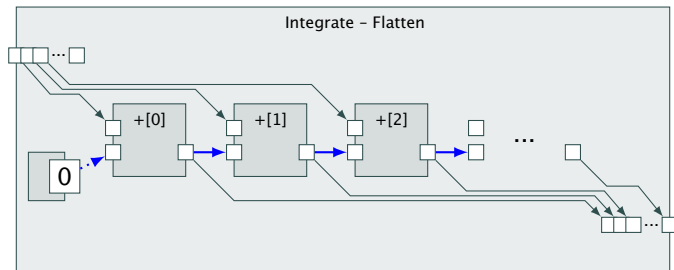
# Need for uniform dependences

## State construction

- *Transfer data between different instances of the same repetition.*
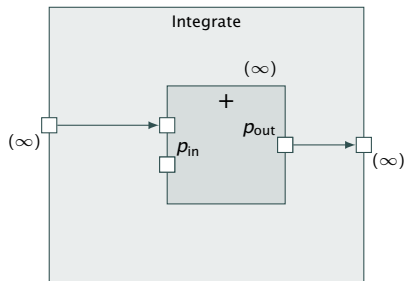- Examples: Sum, **Integrate**

## State construction

- *Transfer data between different instances of the same repetition.*
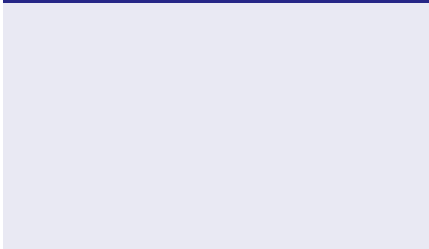- Examples: Sum, **Integrate**



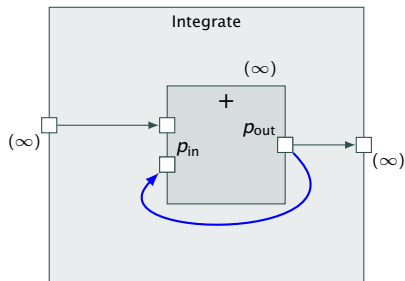Integrate – Flatten

- Uniform data dependences between instances of a repetition

Integrate

$(\infty)$

$+$

$p_{out}$

$(\infty)$

$p_{in}$

$(\infty)$

$(\infty)$

$\mathbf{d} = (1)$

## Inter-repetition dependence

1 Data dependence: $p_{out} \rightarrow p_{in}$
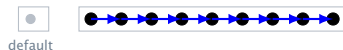
2 Dependence vector inside the repetition space

default

# Uniform dependences



## Inter-repetition dependence

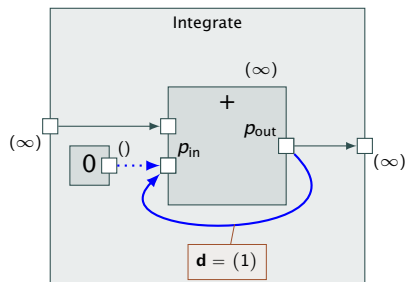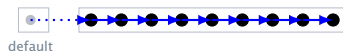1. Data dependence: $p_{out} \rightarrow p_{in}$

2. Dependence vector inside the repetition space

3. Initial values: default link for dependences that exit the repetition space



default

Integrate

$(\infty)$

$+$

$p_{out}$

$(\infty)$

$0$ $()$ $p_{in}$

$(\infty)$

$\mathbf{d} = (1)$

## Inter-repetition dependence

**1** Data dependence: $p_{out} \rightarrow p_{in}$

**2** Dependence vector inside the repetition space

**3** Initial values: default link for dependences that exit the repetition space
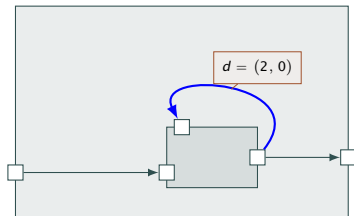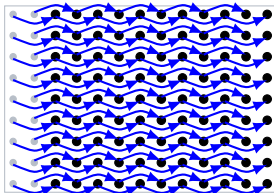


default

Calin Glitia, Philippe Dumont, and Pierre Boulet.

Array-OL with delays, a domain specific specification language for multidimensional intensive signal processing.
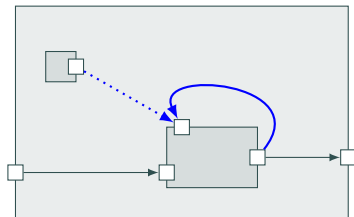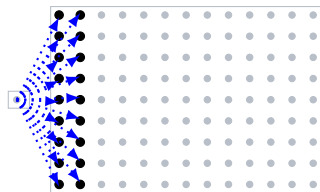
*Multidimensional Systems and Signal Processing*, 2009.
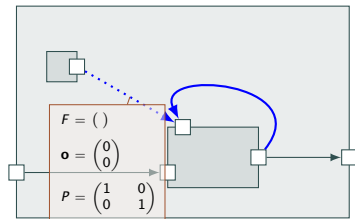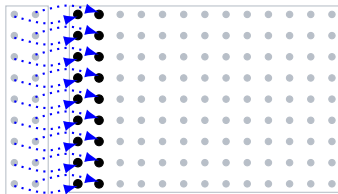
## Initial values – default link



$d = (2, 0)$

# Initial values

## Initial values – default link

- Same initial value

# Initial values

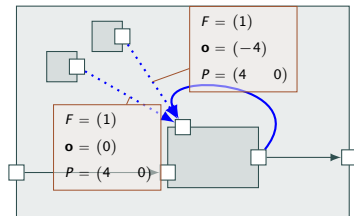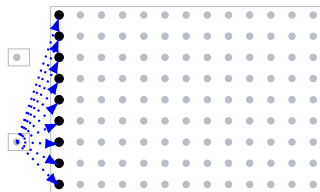## Initial values – default link

- Same initial value
- Different values – Tiler

## Initial values – default link

- Same initial value
- Different values – Tiler
- Different default links – Exclusive tilers



$F = (1)$
$\mathbf{o} = (-4)$
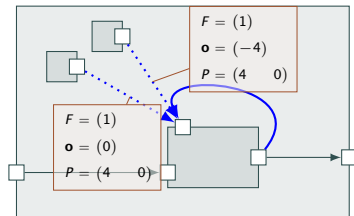$P = (4 \quad 0)$

$F = (1)$
$\mathbf{o} = (0)$
$P = (4 \quad 0)$

## Initial values – default link

- Same initial value
- Different values – Tiler
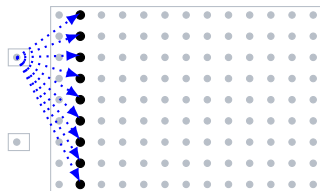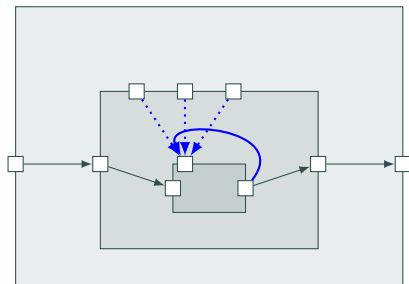- Different default links – Exclusive tilers



$F = (1)$
$\mathbf{o} = (-4)$
$P = (4 \quad 0)$

$F = (1)$
$\mathbf{o} = (0)$
$P = (4 \quad 0)$

# Complex dependences

## Dependence constructions:
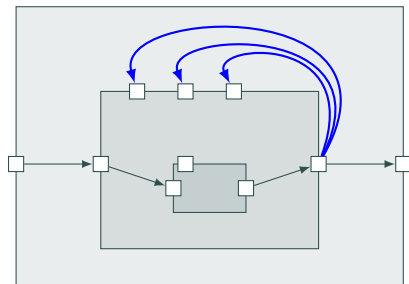
- **Multiple default links**

## Dependence constructions:

- Multiple default links
- Multiple dependences on a repetition space
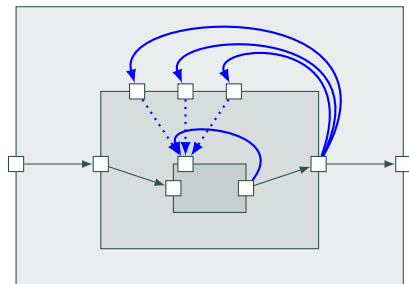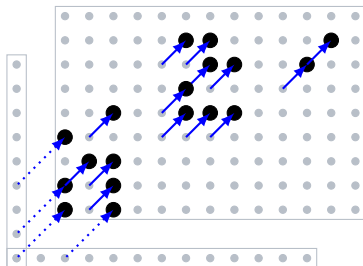
## Dependence constructions:

- Multiple default links
- Multiple dependences on a repetition space
- Dependences connected through the hierarchy



Dependences on the complete repetition space

- The repetition space is split in parallel hyper-planes

- The repetition space is split in parallel hyper-planes
- Pipeline execution following the distance vector

- The repetition space is split in parallel hyper-planes
- Pipeline execution following the distance vector

## Scheduling uniform loops

Alain Darte and Yves Robert.

Constructive methods for scheduling uniform loop nests.
*IEEE Trans. Parallel Distributed Systems*, 5(8):814–822, 1994.

# Modeling and Analysis of Real-Time Embedded Systems

## Profile UML – standard OMG

- Model Driven Engineering
- Co-design: application, architecture, mapping

## Repetitive Structure Modeling

- All the ARRAY-OL concepts are included
- Proposed by the DaRT team

## Physical connections between architecture components

- Compact expression

## Physical connections between architecture components

- Compact expression
- Cyclic uniform inter-connections

**1** Expression of state constructions

1 Expression of state constructions

2 Complex dependences through the hierarchy

1 Expression of state constructions



2 Complex dependences through the hierarchy



3 Parallelism – pipeline

1 Expression of state constructions



2 Complex dependences through the hierarchy



3 Parallelism – pipeline



4 Repeated inter-connected architectures

Modeling and Analysis of Real-Time Embedded Systems

Array Oriented Language

Inter repetition dependence

chitecture

**High-level refactoring**
  – data-parallel transformations
  – strategies

Mapping

Code Generation

## Logical space and time as mixed dimensions of multidimensional structure

- Specification: expresses the data dependences
  - between all the data elements that transits the system
- And a partial execution order
  - between all the execution of the tasks in of the system

## Logical space and time as mixed dimensions of multidimensional structure

- Specification: expresses the data dependences
    - between all the data elements that transits the system
- And a partial execution order
    - between all the execution of the tasks in of the system

## EFFICIENT execution

- Optimized code generation
- Projection of specification into physical space and time

# Execution

## Logical space and time as mixed dimensions of multidimensional structure

- Specification: expresses the data dependences
  - between all the data elements that transits the system
- And a partial execution order
  - between all the execution of the tasks in of the system

## EFFICIENT execution

- Optimized code generation
- Projection of specification into physical space and time

## Adapt a specification to the execution

- High-level refactoring
- Execution that reflects the specification

# Projection into space and time

## Multi-dimensional structures

repetition spaces                                    data structures

## Multi-dimensional structures

repetition spaces
$\Downarrow$
in space

data structures
$\Downarrow$
in time

## Multi-dimensional structures

repetition spaces
$\Downarrow$
in space

$\longleftrightarrow$
linked
(trade-off)

data structures
$\Downarrow$
in time

From a high-level specification to the execution

## Multi-dimensional structures

repetition spaces                           data structures
$\Downarrow$                                             $\Downarrow$
in space            $\longleftrightarrow$                 in time
linked
(trade-off)

## Take into account the execution constraints

- Data dependences
- Available resources

Horizontal Filter
$(240, 1080, \infty)$

Vertical Filter
$(720, 120, \infty)$

$(1920, 1080, \infty)$

$(720, 1080, \infty)$

$(720, 480, \infty)$

$(13)$

$(3)$

$(14)$

$(4)$

Horizontal Filter

$(240, 1080, \infty)$

$(1920, 1080, \infty)$

$(720, 1080, \infty)$

$(13)$

$(3)$

Vertical Filter

$(720, 120, \infty)$

$(720, 480, \infty)$

$(14)$

$(4)$

# Projection example



Horizontal Filter

$(240, 1080, \infty)$

$(1920, 1080, \infty)$

$(13)$

$(3)$

Vertical Filter

$(720, 120, \infty)$

$(720, 1080, \infty)$

$(720, 480, \infty)$

$(14)$

$(4)$

## Maximal parallelism

- Memory size
- Infinite data structures – Blocking points

# Projection example



## Pipeline

- Execution Order

The diagram shows a nested structure labeled $(240, 120, \infty)$ containing two filter blocks: "Horizontal Filtre" $(14)$ and "Vertical Filter" $(3)$. Input is labeled $(1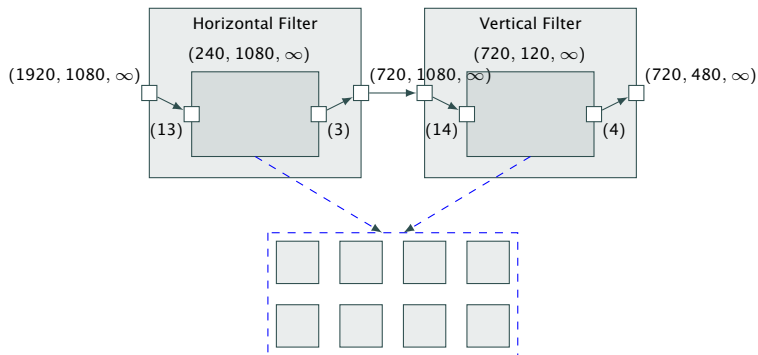920, 1080, \infty)$ with $(14, 13)$, internal labels $(13)$, $(3)$, $(3, 14)$, $(14)$, $(4)$, and output $(720, 480, \infty)$ with $(3, 4)$.

## Fusion of successive repetitions

- Minimize the arrays – macro-patterns
- Distribution of the common repetition
- Each processor its macro-patterns in memory

# Projection example



## Re-computations

- When intermediate values are consumed by multiple repetitions
- Trade-off
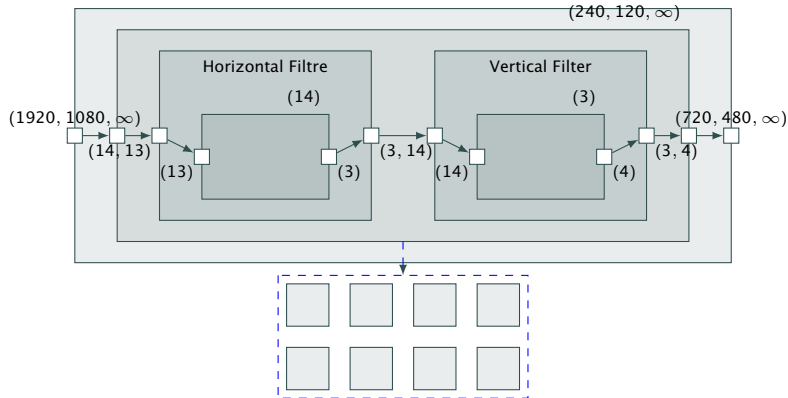  - Recompute values
  - Keep in memory – increase of memory size

## Adapt a specification to execution

- change the granularity of the repetitions
- array sizes reductions

## Adapt a specification to execution

- change the granularity of the repetitions
- array sizes reductions

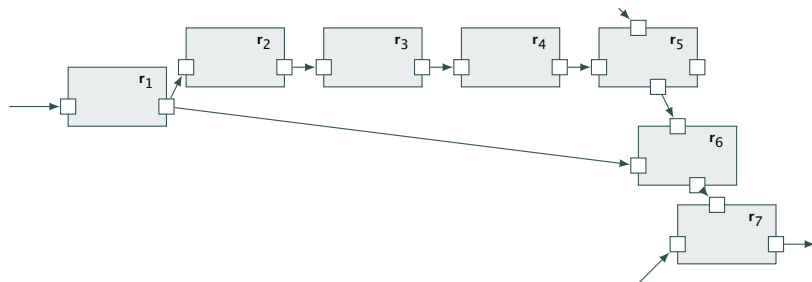## "High-level" loop transformations

- repetition = visual representation of data-parallel loop nest
- fusion, change paving, tiling, collapse, . . .
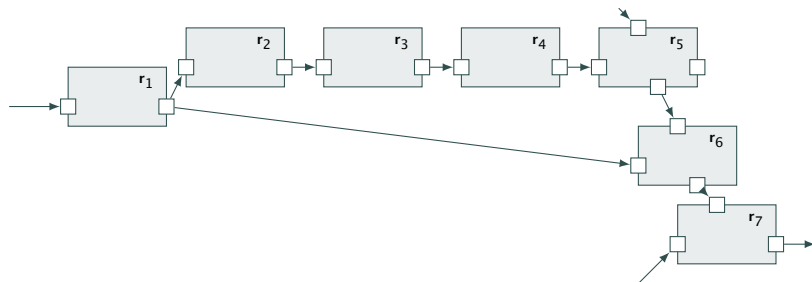
Calin Glitia and Pierre Boulet.

High level loop transformations for multidimensional signal processing embedded applications.

In *International Symposium on Systems, Architectures, Modeling, and Simulation (SAMOS VIII)*, Samos, Greece, July 2008.
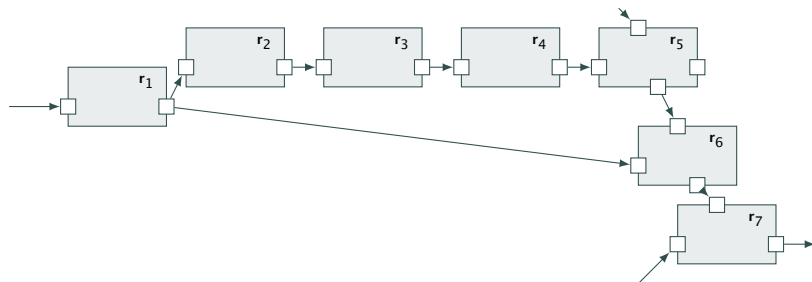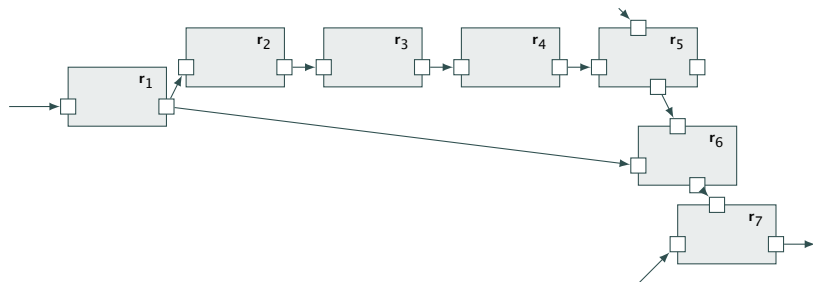
■ MAXIMAL reduction of the intermediate arrays

- MAXIMAL reduction of the intermediate arrays
- Fusion of **multiple repetitions**
  - Minimizes only the last intermediate array
  - Re-computations!

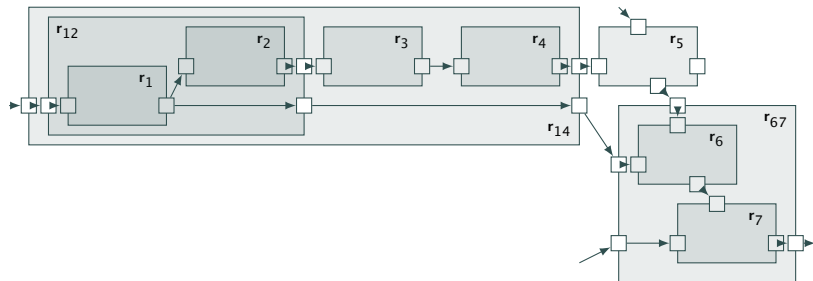# Optimization strategies – memory size reduction



- ■ MAXIMAL reduction of the intermediate arrays
- ■ Fusion of **multiple repetitions**
  - ■ Minimizes only the last intermediate array
  - ■ Re-computations!
- ■ Complete fusion?
  - ■ Too much re-computations
  - ■ Limited array reduction

- MAXIMAL reduction of the intermediate arrays

- Strategy that limits the re-computations
  - using result from complete fusion and two-by-two fusions
  - where re-computations are introduces and minimal achievable array reduction
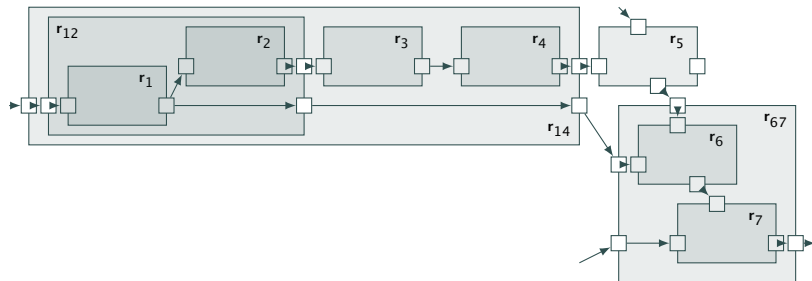
# Optimization strategies – memory size reduction



■ MAXIMAL reduction of the intermediate arrays

| Repetitions before fusion | Repetitions after fusion | | Re-computations (product) | Reduction factor of the output arrays |
|---|---|---|---|---|
| $8 \times 128 \times 96$ | | $119 \times \begin{pmatrix} 10 \times 8 \\ 1 \end{pmatrix}$ | 9.29 | 1228.8 |
| $119 \times 96$ | $96 \times$ | | 1 | 96 |
| $80 \times 80 \times 96$ | | $80 \times 80$ | 1 | 96 |
| $96$ | | $1$ | 1 | 1 |
| $128 \times 96 \times 80$ | $128 \times 96 \times 80$ | | 1 | 1 |
| $119 \times 128 \times 96$ | $128 \times 96 \times \begin{pmatrix} 119 \\ 1 \end{pmatrix}$ | | 1 | 12288 |
| $128 \times 96$ | | | 1 | 1 |

■ MAXIMAL reduction of the intermediate arrays

📄 Calin Glitia, Pierre Boulet, Éric Lenormand, and Michel Barreteau.

Repetitive model refactoring strategy for the design space exploration of intensive signal processing applications.

*Journal of Systems Architecture, Special Issue: Hardware/Software CoDesign.*

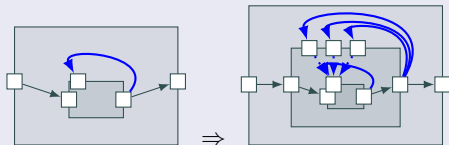# And the inter-repetition dependences?

## Why ?

- To allow the use of the refactoring tools on models with uniform dependences

# And the inter-repetition dependences?

## Why ?

- To allow the use of the refactoring tools on models with uniform dependences
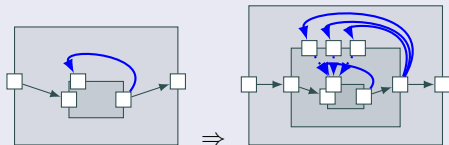


- Typically: $\Rightarrow$

# And the inter-repetition dependences?

## Why ?

■ To allow the use of the refactoring tools on models with uniform dependences
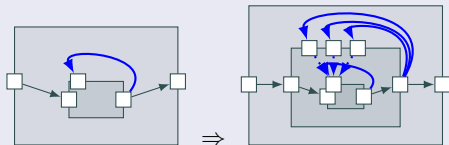


■ Typically: $\Rightarrow$

## Algorithm

■ The global accesses and dependences MUST remain unchanged
■ Automatically compute new dependences after a transformation

# And the inter-repetition dependences?

## Why ?

- To allow the use of the refactoring tools on models with uniform dependences
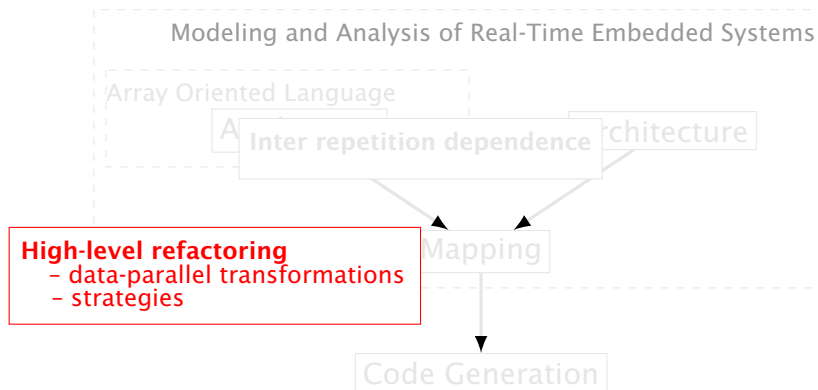


- Typically: $\Rightarrow$

## Algorithm

- The global accesses and dependences MUST remain unchanged
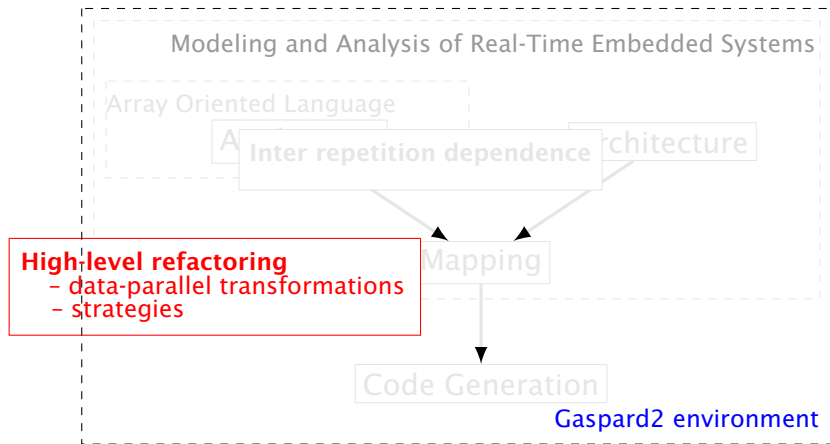- Automatically compute new dependences after a transformation

Calin Glitia and Pierre Boulet.

Interaction between inter-repetition dependences and high-level transformations in array-ol.
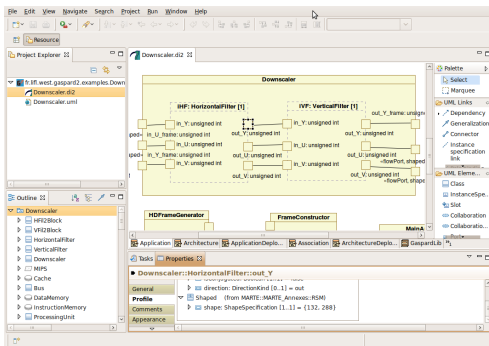
In *Conference on Design and Architectures for Signal and Image Processing (DASIP 2009)*, Sophia Antipolis, France, September 2009.

# Outline

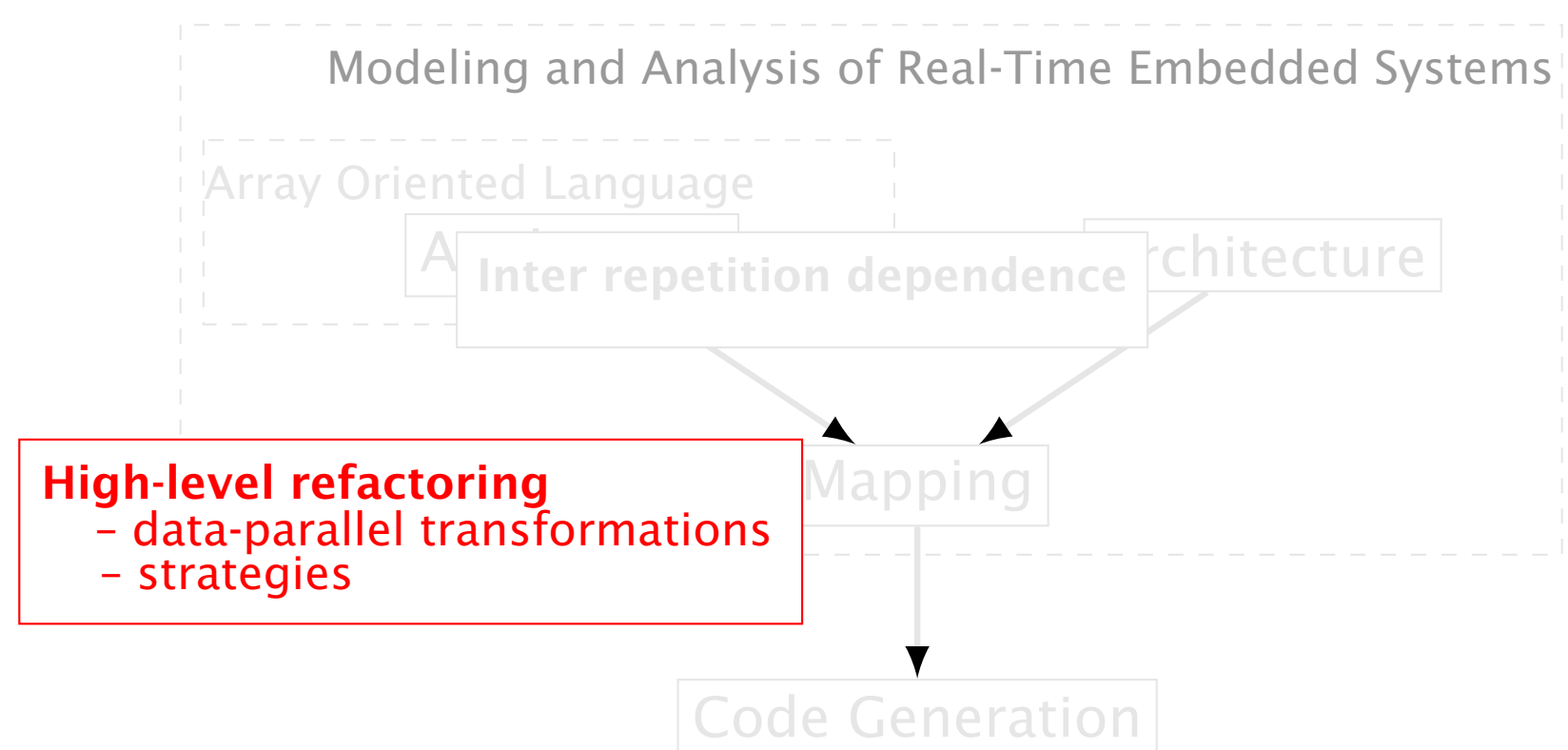

Modeling and Analysis of Real-Time Embedded Systems
Array Oriented Language
Inter repetition dependence
rchitecture
Mapping
Code Generation
High-level refactoring
– data-parallel transformations
– strategies

## SoC visual co-design

1. allows modeling, simulation and code generation of SoC
2. approach Model Driven Engineering
3. Subset of the MARTE UML profile

high-level specification
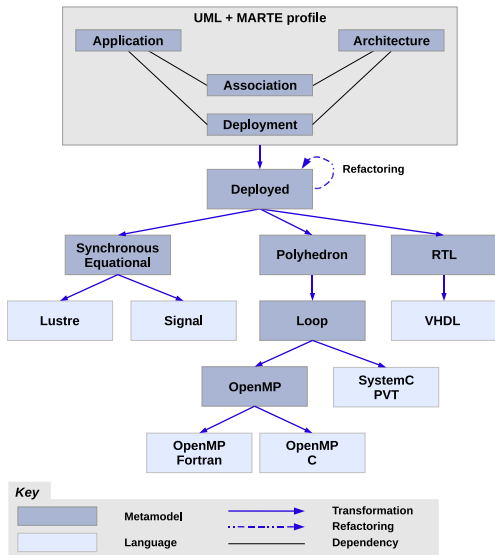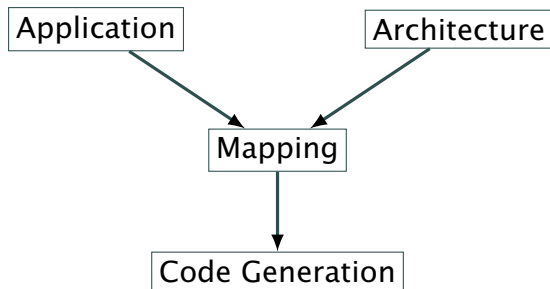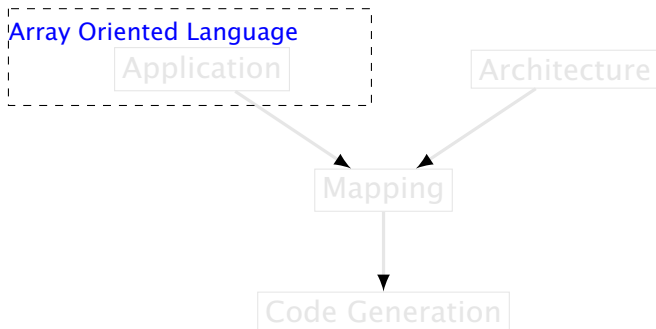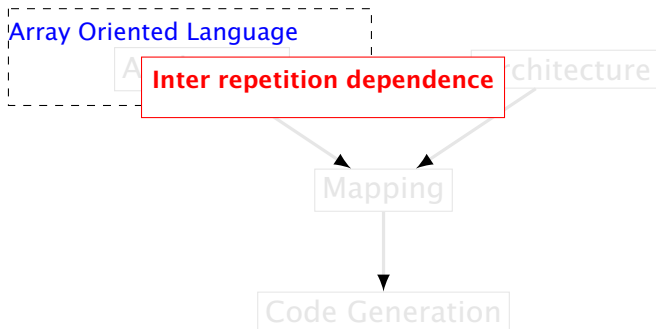
specializations

code generation

high-level specification

- inter-repetition dependences
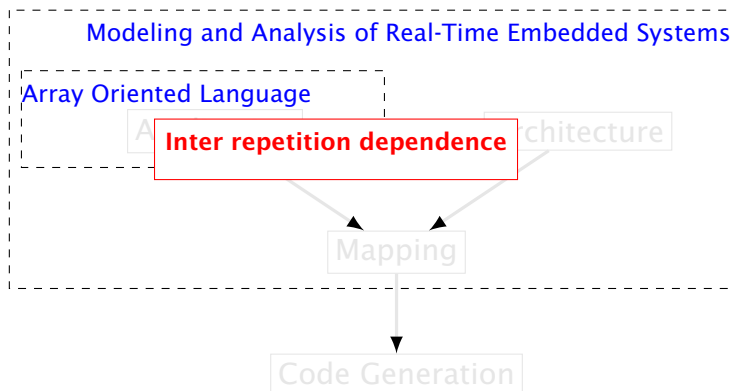- refactoring tools: implementation and integration
- MDE contributions to Gaspard2

specializations

code generation

Application        Architecture

Mapping

Code Generation

# Conclusion

Modeling and Analysis of Real-Time Embedded Systems

Array Oriented Language

**Inter repetition dependence** rchitecture

**High-level refactoring**
  – data-parallel transformations
  – strategies

Mapping

Code Generation

Modeling and Analysis of Real-Time Embedded Systems

Array Oriented Language

**Inter repetition dependence** rchitecture

**High level refactoring**
  – data-parallel transformations
  – strategies

Mapping

Code Generation

Gaspard2 environment

. . .