

From DSL to HPC Component-Based Runtime: A Multi-Stencil DSL Case Study

Julien Bigot, Hélène Coullon, Christian Perez

INRIA team Avalon
Maison de la simulation (CEA)

16th November 2015

Motivation

- + Domain Specific Languages

- ▶ Easy language for the user
- ▶ Implicit optimizations
- ▶ Implicit parallelization

- Domain Specific Languages

- ▶ Difficulties deported to the DSL designer
 - ▶ Low level high performance programming
 - ▶ Maintainability and portability
- ▶ As many DSLs as domains
 - ▶ Productivity and code-reuse
 - ▶ Separation of concerns

Motivation

Component models

- ▶ Divide an application into several functional black boxes
- ▶ Each component defines its connections with outer world
- ▶ Each component are built independently from the others
- ▶ Application = Assembly of components

+ Component models

- ▶ Maintainability through separation of concerns
- ▶ Code-reuse and productivity
- ▶ Dynamic assembly of components

Motivation

What if a DSL produces a component-based runtime?

- ▶ Is it efficient?
- ▶ Is it feasible?
- ▶ Does it improve issues of DSLs?
 - ▶ maintainability
 - ▶ portability
 - ▶ productivity

Let's take a useful example : *the Multi-Stencil Language* !

Table of contents

Multi-Stencil Language

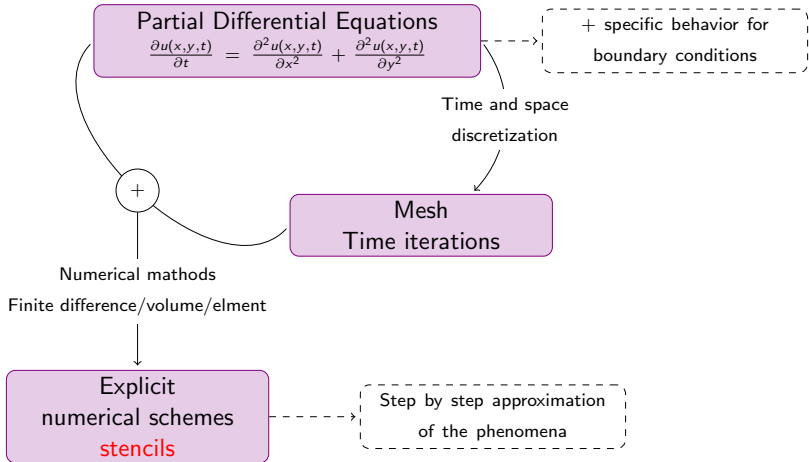
Overview

Compiler

Evaluation

Conclusion and perspectives

Numerical simulation = Multi-Stencil application



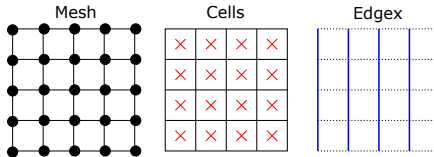
Time and Mesh

Time

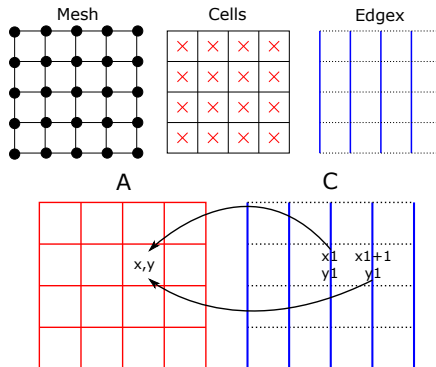
At each time iteration of the simulation are applied the *computations kernels* of the application.

Mesh

- ▶ A Mesh is a connected undirected graph $\mathcal{M} = (V, E)$ without bridges
- ▶ Mesh entities are a subset of $V \cup E$



Computation kernel : Examples



Compute A using $\{(A, local), (C, n2)\}$ on the computation domain called $dc2$ with the numerical expression

$$A(x, y) = A(x, y) + C(x1, y1) + C(x1 + 1, y1).$$

Multi-Stencil program

$$\mathcal{MSP}(T, \mathcal{M}, \mathcal{E}, \mathcal{D}, \Delta, \Gamma)$$

- ▶ T the set of time iterations to tun the simulation
- ▶ \mathcal{M} the mesh of the simulation
- ▶ \mathcal{E} the set of mesh entities
- ▶ \mathcal{D} the set of computation domains
- ▶ Δ the set of data
- ▶ Γ the set of computations

= the six sections of a Multi-Stencil Language program !

Multi-Stencil Language

MSL is not

- ▶ a new stencil optimizer/compiler
- ▶ a new distributed data structure

MSL is

- ▶ a high-level language for multi-stencil simulations
- ▶ agnostic from the type of mesh used (data structure)
- ▶ based on identifiers only

*MSL produces a ready-to-fill
component-based parallel pattern of the simulation !*

References

[illegible]

Original work

MSL to Component-based runtime

Ready-to-fill parallel pattern

- ▶ Data parallelism
 - ▶ External distributed data structure
 - ▶ Automatic detection of synchronizations
- ▶ Task parallelism (mid-grain)
 - ▶ Compile a static scheduling of computation kernels

The fine grain task parallelism is left to other languages :

- ▶ OpenMP in the kernels
- ▶ Kernels generated by stencil compilers (Pochoir, PATUS, Liszt etc.)

The diagram illustrates the system architecture, enclosed in a dashed box. The components and their connections are as follows:

- start**: A small circle representing the initial state, connected to the **Driver** component.
- Driver**: A rounded rectangle, highlighted with a red border. It is connected to the **DDS** component and the **Data** component.
- DDS** (\mathcal{M}, \mathcal{E}): A rounded rectangle representing the Distributed Data Store. It is connected to the **Driver** and the **Data** component.
- Data** (Δ, \mathcal{D}): A rounded rectangle representing the Data component. It is connected to the **Driver**, **DDS**, and **Computations** components.
- Time** (T): A rounded rectangle representing the Time component. It is connected to the **Driver** and the **Computations** component.
- Computations** (Γ): A dashed rounded rectangle representing the Computations component. It is connected to the **Data** and **Time** components.

Connections are marked with small circles (white or black) and lines. A black circle on a line indicates a specific state or event. A line ending in a black square and an asterisk (*) indicates an external connection or output.

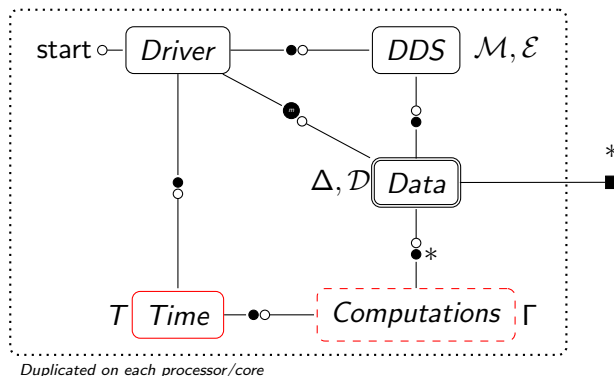
The diagram illustrates the system architecture, enclosed in a dotted box. The components and their connections are as follows:

- start**: A small circle representing the initial state, connected to the **Driver** component.
- Driver**: A rounded rectangle representing the main control unit. It has three outgoing connections:
 - A horizontal line to the **DDS** component.
 - A diagonal line to the **Data** component, passing through a small black circle and a small white circle.
 - A vertical line to the **Time** component, passing through a small black circle and a small white circle.
- DDS** (\mathcal{M}, \mathcal{E}): A rounded rectangle with a red border, representing the Distributed Data Store. It has a vertical connection to the **Data** component, passing through a small white circle and a small black circle.
- Data** (Δ, \mathcal{D}): A rounded rectangle with a red border, representing the Data component. It has a vertical connection to the **Computations** component, passing through a small white circle and a small black circle, with an asterisk (*) next to the black circle. It also has a horizontal connection to the right, passing through a small black circle and an asterisk (*).
- Time** (T): A rounded rectangle with a purple border, representing the Time component. It has a horizontal connection to the **Computations** component, passing through a small black circle and a small white circle.
- Computations** (Γ): A dashed rounded rectangle representing the Computation component.

The entire diagram is enclosed in a dotted box, and the text "Duplicated on each processor/core" is written below it.

MSL to Component-based runtime

$$\mathcal{MSP}(T, \mathcal{M}, \mathcal{E}, \mathcal{D}, \Delta, \Gamma)$$



Example

```

mesh: cart
mesh entities: cell, edgex, edgey
computation domains:
  allcell in cell
  alledgex in edgex
  alledgey in edgey
  part1edgex in edgex
  part2edgex in edgex
data:
  a, cell
  b, cell
  c, edgex
  d, edgex
  e, edgey

```

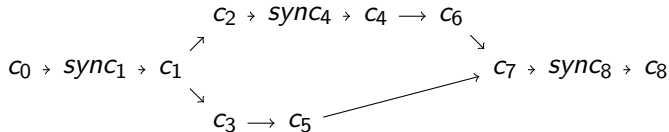
```

f, cell
g, edgey
h, edgex
i, cell
j, edgex
time: 500
computations:
  b[allcell]=c0(a)
  c[alledgex]=c1(b[n1])
  d[alledgex]=c2(c)
  e[alledgey]=c3(c)
  f[allcell]=c4(d[n1])
  g[alledgey]=c5(e)
  h[alledgex]=c6(f)
  i[allcell]=c7(g,h)
  j[partedgex]=c8(i[n1])

```


Mid-grain task parallelism

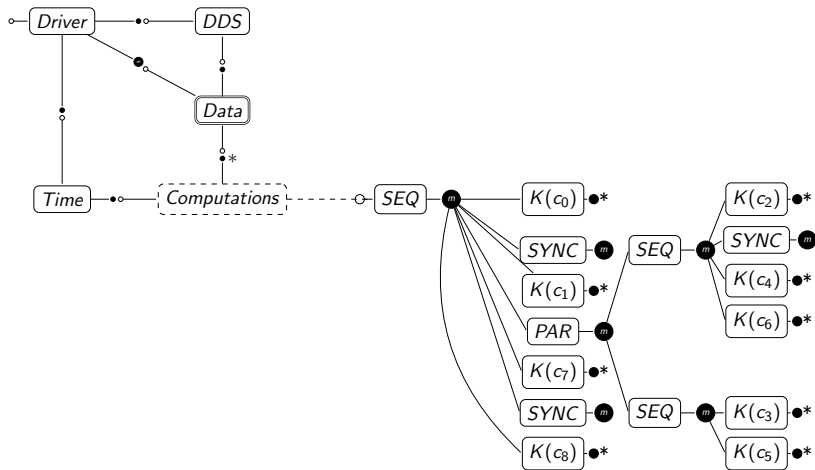
Build the read/write dependency graph.



Data parallelism + Task parallelism

- ▶ Use a dynamic scheduler
- ▶ Compute a static scheduling

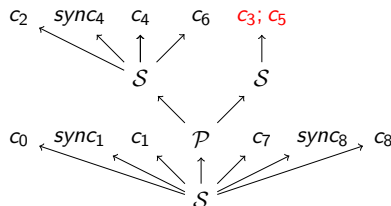
- ▶ We build a static scheduling of the multi-stencil application
- ▶ This static scheduling can be dumped to a component-based runtime



Back to Data Parallelism

Back to Data Parallelism only ?

Loop Fusion : Same computation domain under a sequence or a parallel node



Canonical form : $[c_0, sync_1, c_1, c_2, sync_4, c_4, c_6, \{c_3, c_5\}, c_7, sync_8, c_8]$

Implementation and evaluation

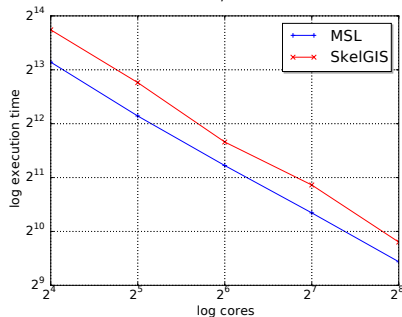
Implementation of MSL : Python, SkelGIS and L^2C

Shallow-water equations : 1 mesh, 3 mesh entities, 7 computation domains, 48 data, 98 computations (32 stencils, 66 local kernels)

Evaluation of the data parallelism

- ▶ Full SkelGIS implementation (DDS + specific interfaces to hide communications)
- ▶ MSL implementation which uses the SkelGIS DDS
- ▶ Cluster Edel of Grid'5000 : 2 8-cores Intel Xeon E5520, 24GB RAM / nodes

ions



- 26 / 28

Perspectives

- ▶ A comparison with an improved version of SkelGIS
- ▶ More evaluations on TGCC Curie (CEA) up to 32k cores
- ▶ An evaluation of the Data+Task parallelism (mixing the static scheduling and OpenMP)
- ▶ A new dump with a component responsible for the dynamic scheduling of the application (using OpenMP 4.x)
- ▶ A new dump to CPU+GPUs (using existing stencil compilers Pochoir, PATUS etc.)
- ▶ An alternative to the SkelGIS DDS using Global Arrays (PGAS)

Thank You !