# Seamlessly Scaling Applications with DAPHNE

COMPAS 2024, Nantes, France

---

Quentin GUILLOTEAU, Jonas H. Müller KORNDÖRFER, Florina M. CIORBA

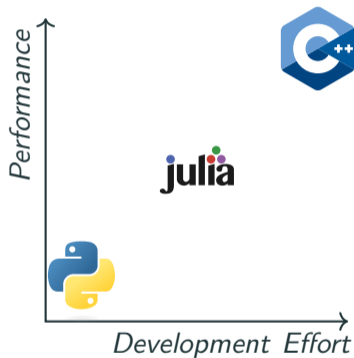2024-07-04

University of Basel, Switzerland
{Quentin.Guilloteau, Jonas.Korndorfer, Florina.Ciorba}@unibas.ch

## The "Two-Language Problem"

- **Fast** development (e.g. Python) …

- But reimplementation for **performance** (e.g. C++)

→ Julia as a "solution"

- **Trade-off**: Performance vs. Ease of Development

- **Fast** development (e.g. Python) …
- ~~Determined computation for **performance** (e.g. C+++)~~
- → Julia as a "solution"
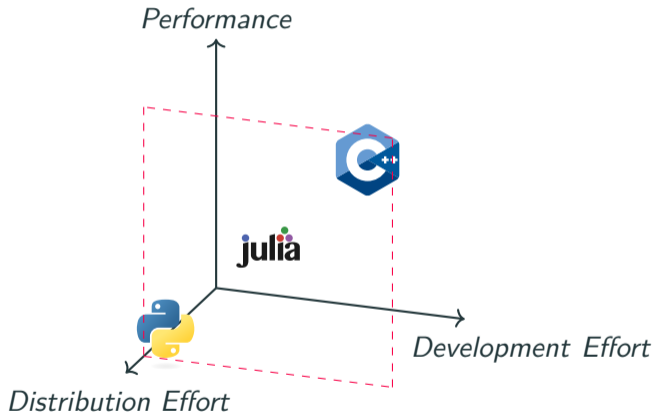- **Trade-off**: Performance vs. Ease of Development

↪ Only for a **single machine**!

Rewrite the application to:

- Distribute the computations
- Integrate communication library (e.g. **MPI**)

Requires:

- **More substantial effort**
- Error handling
- Additional expertise

Rewrite the application to:

- Distribute the computations
- Integrate communication library (e.g. **MPI**)

Requ

- **More substantial effort**
- Error handling
- Additional expertise

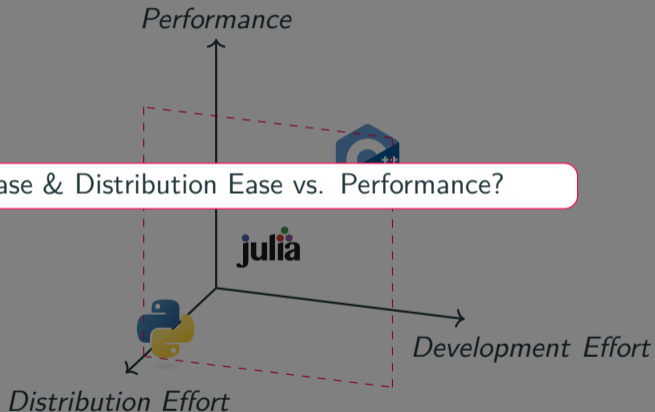*Performance*

↪ **Trade-off**: Development Ease & Distribution Ease vs. Performance?

julia

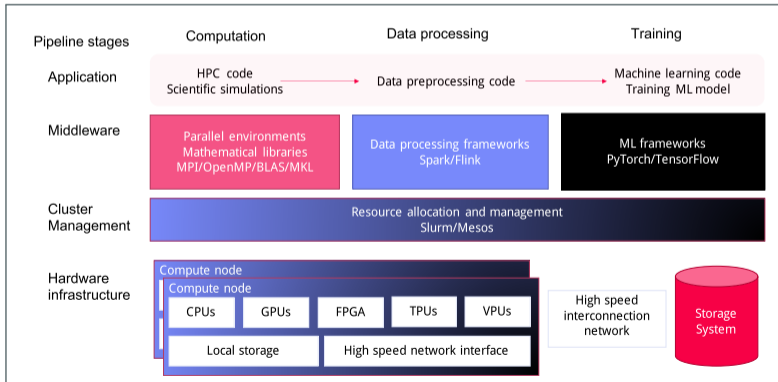*Development Effort*

*Distribution Effort*

# DAPHNE

**An Open and Extensible System Infrastructure for IDA Pipelines**

# DAPHNE – Motivation

- **I**ntegrated **D**ata **A**nalysis pipelines: HPC → DM → ML
- **Different** libraries, programming models, etc. at **every stage**!

- EU H2020 Project 🇪🇺 (2021-2024)
- Open-source: https://github.com/daphne-eu/daphne

**DAPHNE**



P. Damme et al., DAPHNE: An Open and Extensible System Infrastructure for Integrated Data Analysis Pipelines [CIDR 2022]

- EU H2020 Project 🇪🇺 (2021-2024)
- Open-source: https://github.com/daphne-eu/daphne

**DAPHNE**

Application

| DaphneLib (API) | Python API |

| DaphneDSL (Domain-specific Language) |

**Extensible Infrastructure**

Compilation

| **MLIR** | **DaphneIR** (MLIR Dialect) |
| MLIR-Based Compilation Chain | Optimization Passes |
| | New Runtime Abstractions for Data, Devices, Operations |
| | Hierarchical Scheduling |

**Multi-level Compilation/ Runtime**

Runtime System

**Scheduling!**

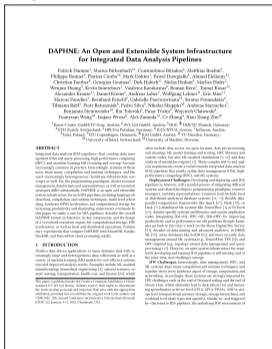| **Device Kernels** (CPU, GPU, FPGA, Storage) | **Vectorized Execution Engine** (Fused Op Pipelines) | **Sync/Async I/O Buffer/Memory Management** |

**Fine-grained Fusion and Parallelism**

Deployment

**Local (embedded) and Distributed Environments** (standalone, HPC, data lake, cloud, DB)

**Integration w/ Resource Mgmt & Prog. Models**



P. Damme et al., DAPHNE: An Open and Extensible System Infrastructure for Integrated Data Analysis Pipelines [CIDR 2022]
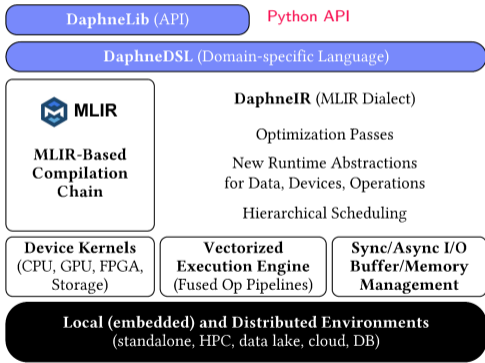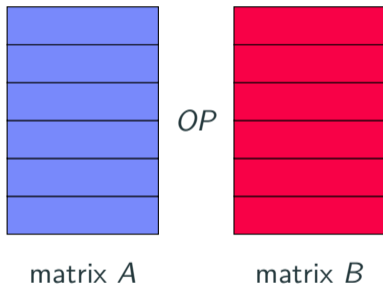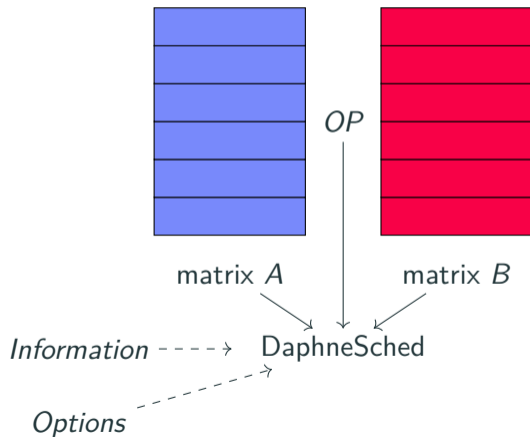
# DaphneSched

## Local and Distributed Versions

matrix $A$     $OP$     matrix $B$

matrix *A*  matrix *B*

Work Partitioner

New Task!

*OP*

matrix *A*     matrix *B*

Work Queue

Work Partitioner

matrix *A*        matrix *B*        Work Queue

Work Partitioner

New Task!

matrix *A*  matrix *B*  Work Queue

Work Partitioner

New Task!

OP

matrix A          matrix B

Work Queue

Work Partitioner

Work Queue

Worker1

Worker2

Work Queue

Work Response

**12 Work Partitioners**
STATIC (OpenMP static),

SS (OpenMP dynamic),

GSS (OpenMP guided),

TSS (in LLVM),

FAC2, TFSS, FISS, VISS, PLS,

MSTATIC, MFSC, PSS (in LB4OMP),

and AUTO

**3 Queue Layouts**
CENTRALIZED, PERGROUP*, PERCPU*

**4 Work Stealing* Strategies**
SEQ, SEQPRI, RND, RNDPRI



Local DaphneSched

**12 Work Partitioners**
STATIC (OpenMP static),

SS (OpenMP dynamic),

GSS (OpenMP guided),

TSS (in LLVM),

FAC2, TFSS, FISS, VISS, PLS,

MSTATIC, MFSC, PSS (in LB4OMP),

and AUTO

**3 Queue Layouts**
CENTRALIZED, PERGROUP*, PERCPU*

**4 Work Stealing* Strategies**
SEQ, SEQPRI, RND, RNDPRI



Connected Components in DaphneDSL, executed
with DaphneSched and **sparse** input matrix
(wikipedia-20070206) on a single node

## DaphneSched – Distributed Scheduler ⇒ Hierarchical Scheduling – NEW!

**Coordinator (MPI Rank 0)**

- Partitions work (data, ops) to the local DaphneSched instances
- Communicator Manager coordinates with local DaphneSched instances
- Different message types:
  - BM: Broadcast Message (Data)
  - CM: Compute Message (MLIR)
  - DM: Distribute Message (Data)
  - RM: Ready Message (Sync)

**Workers (MPI Ranks 1 .. P-1)**

- Listen incoming messages from coordinator
- Execute a local DaphneSched instance



Distributed DaphneSched

# Experimental Evaluation

## Design of Experiments

- Study the **effort and performance of scaling applications**

- Compare C++, Python, Julia, DaphneDSL

- **Connected Components** graph algorithm: broadcast, max, max
  (S. Beamer et al., GAP Benchmark suite)

- 2 Intel Broadwell E5-2640v4 CPUs, 10 cores each, 64GB of RAM

- 3 input matrices:

| Matrix | Size | Density (%) |
|---|---|---|
| amazon0601 | $403'394 \times 403'394$ | $2.08 \times 10^{-3}$ |
| wikipedia-20070206 | $3'566'907 \times 3'566'907$ | $0.354 \times 10^{-3}$ |
| ljournal-2008 | $5'363'260 \times 5'363'260$ | $0.275 \times 10^{-3}$ |

## Implementations – Connected Components (CC) Algorithm

| Language (abbreviation) | External dependencies | Lines of code per implementation | | |
|---|---|---|---|---|
| | | Sequential | Local Parallel | Distributed Parallel |
| C++ (cpp) | Eigen | $\simeq 25$ | $\simeq 25$ | $\simeq 120$ |
| Python (py) | Numpy, Scipy | $\simeq 10$ | $\simeq 10$ | $\simeq 100$ |
| Julia (jl) | MatrixMarket.jl | $\simeq 25$ | $\simeq 25$ | $\simeq 100$ |
| **DaphneDSL** (daph) | $\varnothing$ | $\simeq \mathbf{10}$ | $\simeq \mathbf{10}$ | $\simeq \mathbf{10}$ |

- CC: broadcast (dense vector to CSR matrix), row-wise max, vector-wise max
- cpp, jl: Broadcast implemented by hand
- MPI: encapsulates the local parallel version, between a scatter of CSR matrix and MPI_Allreduce with user-defined function

## Implementations – CC in DaphneDSL

```
1 G = readMatrix($f); // read sparse matrix from CLI argument
2 maxi = 100; // maximum number of iterations
3 start = now();
4
5 c = seq(1.0, as.f64(nrow(G)), 1.0); // initialization
6
7 for(iter in 1:maxi) {
8     c = max(aggMax(G * t(c), 0), c);
9 }
10
11 end=now();
12 print((end-start) / 1000000000.0);
```

Strong scaling on a *single node* (20 total cores) with threads (left) and MPI processes (right) for CC with `wikipedia-20070206` as input.

Error bars: 95% confidence intervals.

DaphneSched: `CENTRALIZED` queue w/ `STATIC`.

- DaphneDSL scales <u>with threads</u>
- The others scale <u>with processes</u>

Average execution time with 95% confidence intervals for CC with DaphneDSL and DaphneSched, for 12 scheduling schemes and 3 queue layouts, with 1 victim selection – SEQPRI, executed on 4 nodes and 1 MPI process / node.

Total degree of parallelism: $(4 - 1) \times 20 = 60$ DAPHNE workers.

**Take Away**: Local scheduling options have a non-negligible impact on distributed execution performance

Average execution time with 95% confidence intervals for CC with DaphneDSL and DaphneSched, for 12 scheduling schemes and 3 queue layouts, with 1 victim selection – SEQPRI, executed on 4 nodes and 1 MPI process / node.
Total degree of parallelism: $(4 - 1) \times 20 = 60$ DAPHNE workers.

Strong scaling: 1-9 compute nodes. Inside each node, work is parallelized using MPI for C++, Julia, and Python, and threads for DaphneDSL.

DaphneSched: used `CENTRALIZED + STATIC` (default) and `CENTRALIZED + AUTO`

⇒ highlight the impact of scheduling on performance.

- DAPHNE is outperformed by others on small inputs ☹    ⇒ Impact of scheduling
- DAPHNE outperforms others on larger inputs ☺    ✓ **No additional effort!**

**Take Away**: Seamlessly Scaling Applications with DAPHNE!

Strong scaling: 1-9 compute nodes. Inside each node, work is parallelized using MPI for C++, Julia, and Python, and threads for DaphneDSL.

DaphneSched: used `CENTRALIZED` + `STATIC` (default) and `CENTRALIZED` + `AUTO`
⇒ highlight the impact of scheduling on performance.

- DAPHNE is outperformed by others on small inputs ☹    ⇒ Impact of scheduling
- DAPHNE outperforms others on larger inputs ☺    ✓ **No additional effort!**

# Conclusion and Future Steps

**Conclusion**

- Distributing applications → Difficult, substantial effort, expertise
- **DAPHNE scales seamlessly** (without additional effort) ☺
- Best performance is not always guaranteed ☹
- Interesting trade-off: **Performance vs. Ease of Development**

**Future Steps**

- Dynamic partitioning by the coordinator in Distributed DaphneSched
- Communication/Stealing between the distributed workers
- Trade-off between colocating the coordinator with workers?
- Evaluation with full IDA pipelines

# Seamlessly Scaling Applications with DAPHNE

COMPAS 2024, Nantes, France

---

Quentin GUILLOTEAU, Jonas H. Müller KORNDÖRFER, Florina M. CIORBA

2024-07-04

University of Basel, Switzerland
{Quentin.Guilloteau, Jonas.Korndorfer, Florina.Ciorba}@unibas.ch

# Additional Slides

## Implementations – CC in DaphneDSL

```
1 G = readMatrix($f); // read sparse matrix from CLI argument
2 maxi = 100; // maximum number of iterations
3 start = now();
4
5 c = seq(1.0, as.f64(nrow(G)), 1.0); // initialization
6
7 for(iter in 1:maxi) {
8     c = max(aggMax(G * t(c), 0), c);
9 }
10
11 end=now();
12 print((end-start) / 1000000000.0);
```

## Implementations – CC in Python

```python
1 import sys
2 import time
3 from scipy.io import mmread
4 from scipy.sparse import csr_matrix, csr_array
5 import numpy as np
6
7 def cc(filename, maxi=100):
8     G = csr_matrix(mmread(filename))
9     n = G.shape[1]
10     start = time.time()
11     c = np.array([list(map(lambda i: float(i), range(1, n + 1, 1)))])
12     for iter in range(maxi):
13         x = G.multiply(c.transpose()).max(axis=0)
14         c = np.maximum(c, x.todense())
15     end = time.time()
16     print(end - start)
```

## Implementations – CC in Julia

```julia
using MatrixMarket
using SparseArrays
using SparseMatricesCSR

function G_broadcast_mult_c(G, c)
  cols = colvals(G)
  vals = nonzeros(G)
  m, n = size(G)
  maxs = zeros(n)
  for j = 1:m
     for i in nzrange(G, j)
        col = cols[i]
        val = vals[i]
        if val * c[j] > maxs[col]
          maxs[col] = val*c[j]
        end
     end
  end
```

## Implementations – CC in C++

```cpp
1  #include <iostream>
2  #include <Eigen/SparseCore>
3  #include <Eigen/Dense>
4  #include <Eigen/Sparse>
5  #include <Eigen/Core>
6  #include <unsupported/Eigen/SparseExtra>
7  #include <chrono>
8
9  typedef Eigen::SparseMatrix<double, Eigen::RowMajor> SpMatR;
10 typedef SpMatR::InnerIterator InIterMatR;
11
12 int main(int argc, char** argv) {
13   if (argc != 3) {
14     std::cout << "Usage: bin mat.mtx size" << std::endl;
15     return 1;
16   }
17   std::string filename = argv[1];
```