

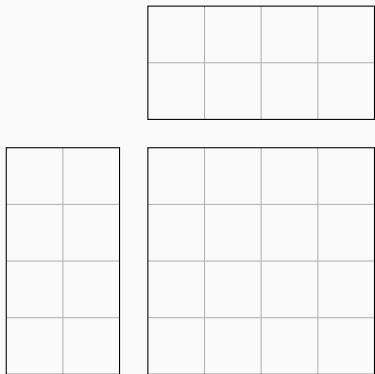
TASK-BASED PROGRAMMING MODELS FOR SCALABLE DISTRIBUTED MEMORY PARALLEL ALGORITHMS

E. Agullo, A. Buttari, A. Guermouche, J. Herrmann, A. Jego
December 4, 2020

Université de Toulouse - IRIT

GENERAL MATRIX MULTIPLICATION ALGORITHMS

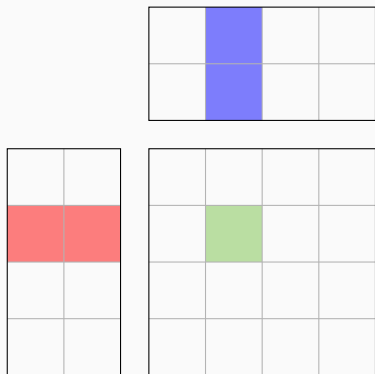
SEQUENTIAL GEMM



```
do i=1,m
  do j=1,n
    do l=1,k
      call gemm( A[i,l], B[l,j], C[i,j])
    end do
  end do
end do
```

$$\forall i,j, C_{ij} += \sum A_{i,k} B_{k,j}$$

SEQUENTIAL GEMM



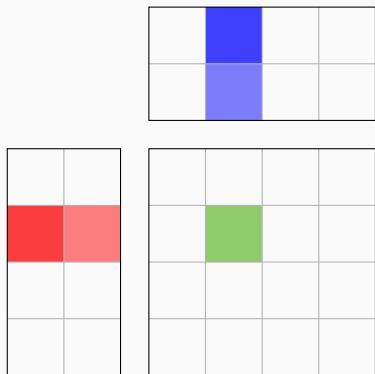
As an example

$$C_{2,2} =$$

```
do i=1,m
  do j=1,n
    do l=1,k
      call gemm( A[i,l], B[l,j], C[i,j])
    end do
  end do
end do
```

$$\forall i,j, C_{ij} += \sum A_{i,k} B_{k,j}$$

SEQUENTIAL GEMM



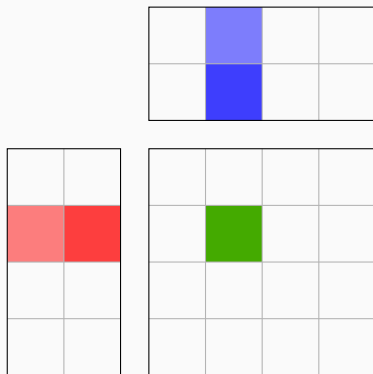
```
do i=1,m
  do j=1,n
    do l=1,k
      call gemm( A[i,l], B[l,j], C[i,j])
    end do
  end do
end do
```

$$\forall i,j, C_{i,j} += \sum A_{i,k} B_{k,j}$$

As an example

$$C_{2,2} = A_{2,1} * B_{1,2}$$

SEQUENTIAL GEMM



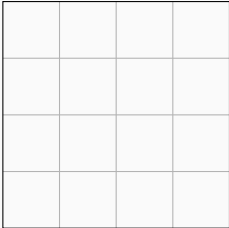
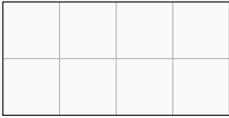
```
do i=1,m
  do j=1,n
    do l=1,k
      call gemm( A[i,l], B[l,j], C[i,j])
    end do
  end do
end do
```

$$\forall i,j, C_{i,j} += \sum A_{i,k} B_{k,j}$$

As an example

$$C_{2,2} = A_{2,1} * B_{1,2} + A_{2,2} * B_{2,2}$$

Ingredients



P0	P1	P2	P3
P4	P5	P6	P7

P0	P1
P4	P5
P8	P9
P12	P13

P0	P1	P2	P3
P4	P5	P6	P7
P8	P9	P10	P11
P12	P13	P14	P15

Ingredients

- Data Mapping - here, 2D Block Cyclic

GEMM : DISTRIBUTED MEMORY (BASIC)

P0	P1	P2	P3
P4	P5	P6	P7

P0	P1
P4	P5
P8	P9
P12	P13

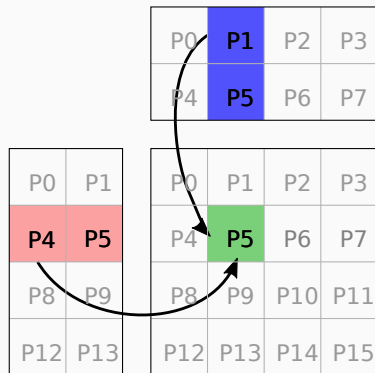
P0	P1	P2	P3
P4	P5	P6	P7
P8	P9	P10	P11
P12	P13	P14	P15

Ingredients

- Data Mapping - here, 2D Block Cyclic
- Task Mapping (= Data Mapping)

```
do k=1,2
  call Recv(A[i,k])
  call Recv(B[k,j])
  call gemm(A[i,k], B[k,j], C[i,j])
end do
```

GEMM : DISTRIBUTED MEMORY (BASIC)

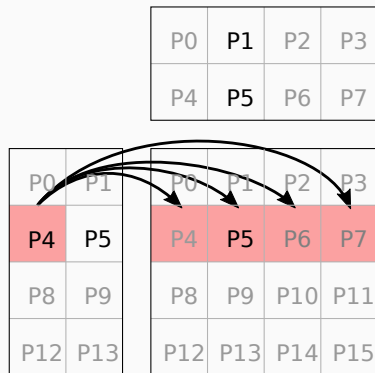


Ingredients

- Data Mapping - here, 2D Block Cyclic
- Task Mapping (= Data Mapping)
- Communicat°

```
do k=1,2
  call Send(A[i,j]) on my row
  call Send(B[i,j]) on my column
  call gemm(A[i,k],B[k,j],C[i,j])
end do
```

GEMM : DISTRIBUTED MEMORY (SUMMA)



Sender and receivers match color.

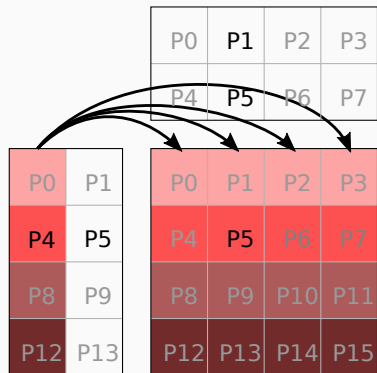
Sending block of A from P4

Ingredients

- Data Mapping - here, 2D Block Cyclic
- Task Mapping (= Data Mapping)
- Communicat^o patterns

```
do k=1,2
  call Bcast(A[i,j]) on my row
  call Bcast(B[i,j]) on my column
  call gemm(A[i,k],B[k,j],C[i,j])
end do
```

GEMM : DISTRIBUTED MEMORY (SUMMA)



Sender and receivers match color.

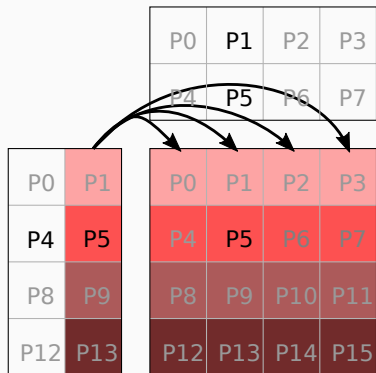
Sending blocks of A from **the column** of P4

Ingredients

- Data Mapping - here, 2D Block Cyclic
- Task Mapping (= Data Mapping)
- Communicat° patterns

```
do k=1,2
  call Bcast(A[i,j]) on my row
  call Bcast(B[i,j]) on my column
  call gemm(A[i,k],B[k,j],C[i,j])
end do
```

GEMM : DISTRIBUTED MEMORY (SUMMA)



Sender and receivers match color.

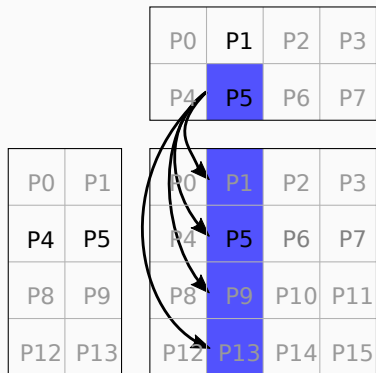
Sending blocks of A from **the column** of P5

Ingredients

- Data Mapping - here, 2D Block Cyclic
- Task Mapping (= Data Mapping)
- Communicat° patterns

```
do k=1,2
  call Bcast(A[i,j]) on my row
  call Bcast(B[i,j]) on my column
  call gemm(A[i,k],B[k,j],C[i,j])
end do
```

GEMM : DISTRIBUTED MEMORY (SUMMA)



Sender and receivers match color.

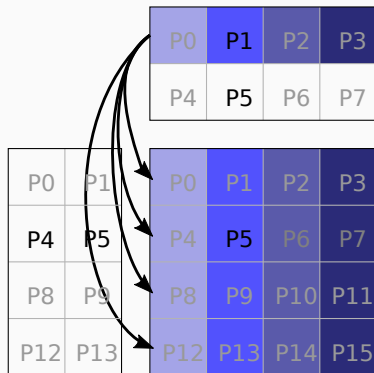
Sending block of B from P5

Ingredients

- Data Mapping - here, 2D Block Cyclic
- Task Mapping (= Data Mapping)
- Communicat^o patterns

```
do k=1,2
  call Bcast(A[i,j]) on my row
  call Bcast(B[i,j]) on my column
  call gemm(A[i,k],B[k,j],C[i,j])
end do
```

GEMM : DISTRIBUTED MEMORY (SUMMA)



Sender and receivers match color.

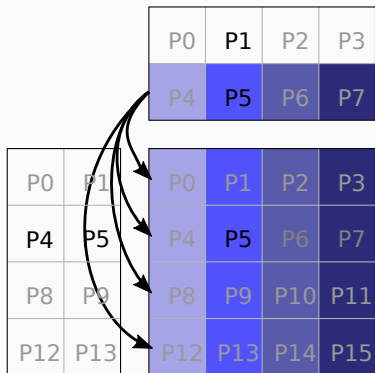
Sending blocks of B from **the row** of P1

Ingredients

- Data Mapping - here, 2D Block Cyclic
- Task Mapping (= Data Mapping)
- Communicat^o patterns

```
do k=1,2
  call Bcast(A[i,j]) on my row
  call Bcast(B[i,j]) on my column
  call gemm(A[i,k],B[k,j],C[i,j])
end do
```

GEMM : DISTRIBUTED MEMORY (SUMMA)



Ingredients

- Data Mapping - here, 2D Block Cyclic
- Task Mapping (= Data Mapping)
- Communicat^o patterns

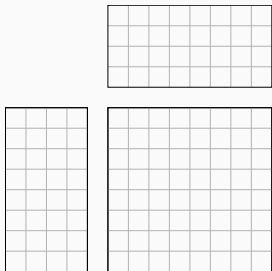
```
do k=1,2
  call Bcast(A[i,j]) on my row
  call Bcast(B[i,j]) on my column
  call gemm(A[i,k],B[k,j],C[i,j])
end do
```

Sender and receivers match color.

Sending blocks of B from **the row** of P5

SUMMA ALGORITHMS

m, m, k 2D product on a $\sqrt{p} \times \sqrt{p}$ grid



Communications critical path (*latency* α and *bandwidth* β)

SUMMA ALGORITHMS

m, m, k 2D product on a $\sqrt{p} \times \sqrt{p}$ grid

0	1	2	3	4	5	6	7
8							
16							
24							

0	1	2	3
8			
16			
24			
32			
40			
48			
56	57	58	59

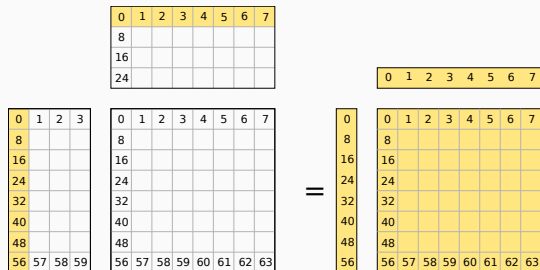
0	1	2	3	4	5	6	7
8							
16							
24							
32							
40							
48							
56	57	58	59	60	61	62	63

Communications critical path (*latency* α and *bandwidth* β)

2D: $\left(\quad \right) (\alpha + \beta b^2)$

SUMMA ALGORITHMS

m, m, k 2D product on a $\sqrt{p} \times \sqrt{p}$ grid

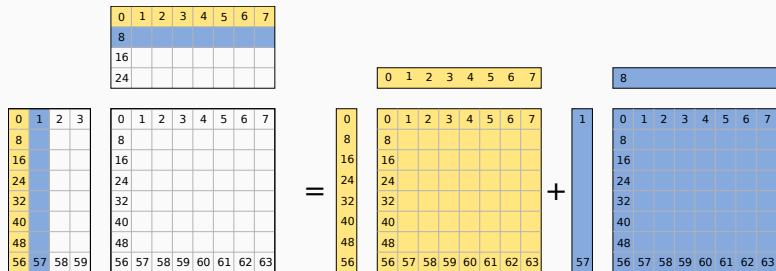


Communications critical path (*latency α and bandwidth β*)

$$2D: \left(\frac{2m}{b\sqrt{p}} \log(\sqrt{p}) \right) (\alpha + \beta b^2)$$

SUMMA ALGORITHMS

m, m, k 2D product on a $\sqrt{p} \times \sqrt{p}$ grid

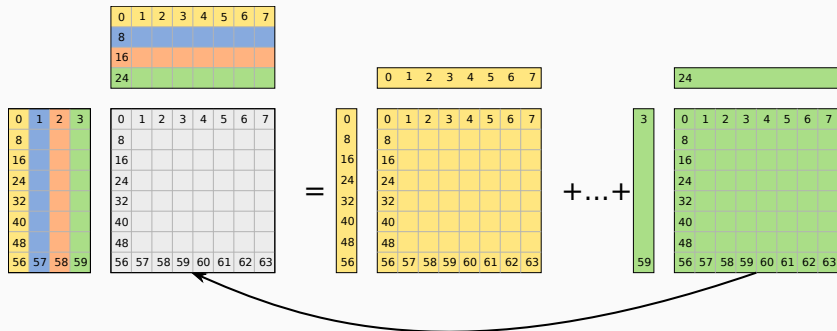


Communications critical path (*latency* α and *bandwidth* β)

$$2D: \left(\frac{2m}{b\sqrt{p}} \log(\sqrt{p}) \right) (\alpha + \beta b^2)$$

SUMMA ALGORITHMS

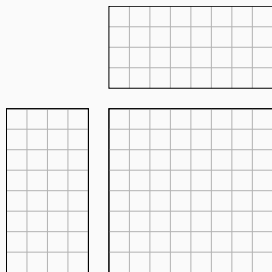
m, m, k 2D product on a $\sqrt{p} \times \sqrt{p}$ grid



Communications critical path (*latency* α and *bandwidth* β)

$$2D: \frac{k}{b} \left(\frac{2m}{b\sqrt{p}} \log(\sqrt{p}) \right) (\alpha + \beta b^2)$$

m, m, k 3D product on a $\sqrt[3]{p} \times \sqrt[3]{p} \times \sqrt[3]{p}$ grid



Communications critical path (*latency* α and *bandwidth* β)

$$2D: \frac{k}{b} \left(\frac{2m}{b\sqrt{p}} \log(\sqrt{p}) \right) (\alpha + \beta b^2)$$

m, m, k 3D product on a $\sqrt[3]{p} \times \sqrt[3]{p} \times \sqrt[3]{p}$ grid

0	1	2	3	0	1	2	3
4							
8							
12							

0	1	2	3
4			
8			
12			
0			
4			
8			
12	13	14	15

0	1	2	3	0	1	2	3
4							
8							
12							
0							
4							
8							
12	13	14	15	12	13	14	15

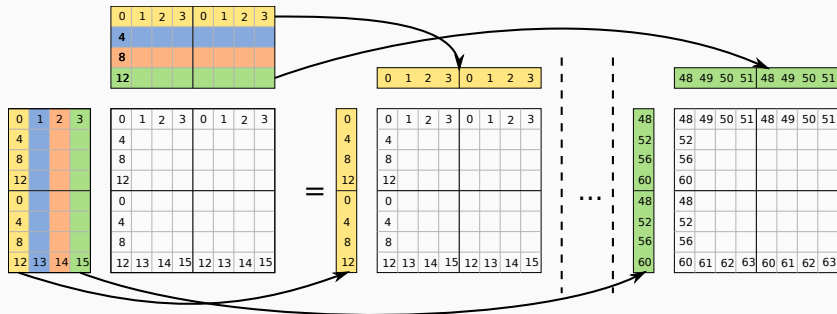
Communications critical path (*latency* α and *bandwidth* β)

$$2D: \frac{k}{b} \left(\frac{2m}{b\sqrt{p}} \log(\sqrt{p}) \right) (\alpha + \beta b^2)$$

$$3D: \left(\right) (\alpha + \beta b^2)$$

SUMMA ALGORITHMS

m, m, k 3D product on a $\sqrt[3]{p} \times \sqrt[3]{p} \times \sqrt[3]{p}$ grid



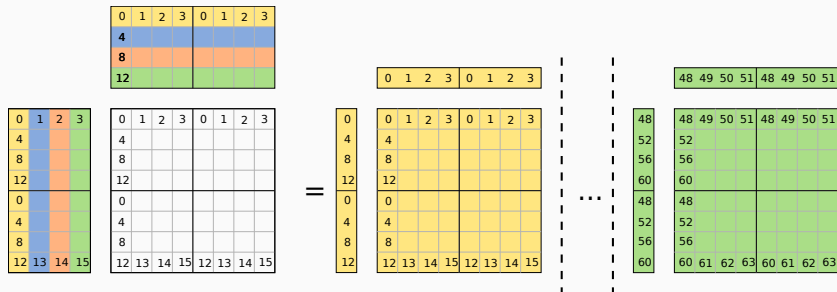
Communications critical path (*latency* α and *bandwidth* β)

$$2D: \frac{k}{b} \left(\frac{2m}{b\sqrt{p}} \log(\sqrt{p}) \right) (\alpha + \beta b^2)$$

$$3D: \left(\frac{2mk}{b^2 \sqrt[3]{p^2}} \right) (\alpha + \beta b^2)$$

SUMMA ALGORITHMS

m, m, k 3D product on a $\sqrt[3]{p} \times \sqrt[3]{p} \times \sqrt[3]{p}$ grid



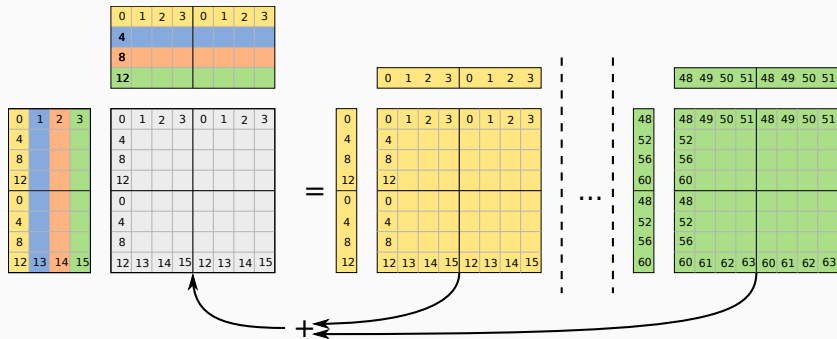
Communications critical path (*latency* α and *bandwidth* β)

$$2D: \frac{k}{b} \left(\frac{2m}{b\sqrt{p}} \log(\sqrt{p}) \right) (\alpha + \beta b^2)$$

$$3D: \left(\frac{2mk}{b^2 \sqrt[3]{p^2}} + \frac{k}{b\sqrt[3]{p}} \left(\frac{2m}{b\sqrt[3]{p}} \log(\sqrt[3]{p}) \right) \right) (\alpha + \beta b^2)$$

SUMMA ALGORITHMS

m, m, k 3D product on a $\sqrt[3]{p} \times \sqrt[3]{p} \times \sqrt[3]{p}$ grid



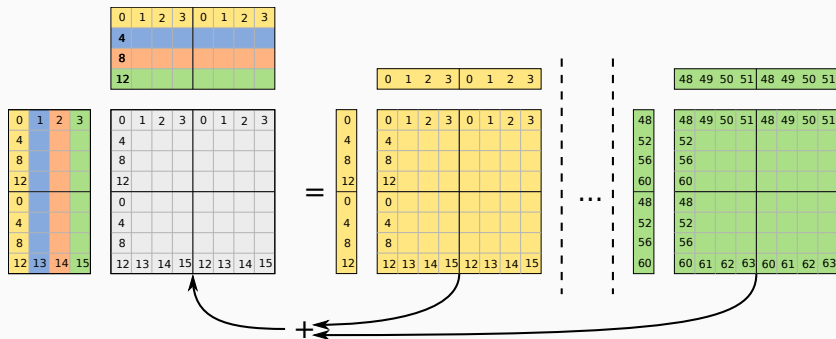
Communications critical path (*latency* α and *bandwidth* β)

$$2D: \frac{k}{b} \left(\frac{2m}{b\sqrt{p}} \log(\sqrt{p}) \right) (\alpha + \beta b^2)$$

$$3D: \left(\frac{2mk}{b^2 \sqrt[3]{p^2}} + \frac{k}{b\sqrt[3]{p}} \left(\frac{2m}{b\sqrt[3]{p}} \log(\sqrt[3]{p}) \right) + \frac{m^2}{b^2 \sqrt[3]{p^2}} \log(\sqrt[3]{p}) \right) (\alpha + \beta b^2)$$

SUMMA ALGORITHMS

m, m, k 3D product on a $\sqrt[3]{p} \times \sqrt[3]{p} \times \sqrt[3]{p}$ grid



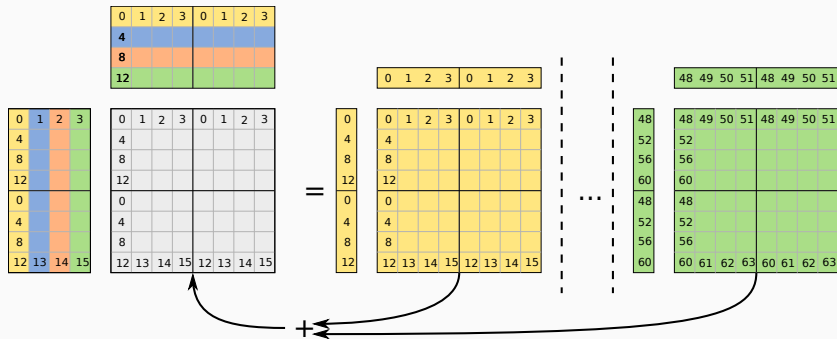
Communications critical path (*latency α and bandwidth β*)

$$2D: 4 \left(\frac{2 \times 8}{8} \log(8) \right) (\alpha + \beta b^2)$$

$$3D: \left(\frac{2 \times 8 \times 4}{4^2} + \frac{4}{4} \left(\frac{2 \times 8}{4} \log(4) \right) + \frac{8 \times 8}{4^2} \log(4) \right) (\alpha + \beta b^2)$$

SUMMA ALGORITHMS

m, m, k 3D product on a $\sqrt[3]{p} \times \sqrt[3]{p} \times \sqrt[3]{p}$ grid



Communications critical path (*latency* α and *bandwidth* β)

$$2D: (24\log(2)) (\alpha + \beta b^2)$$

$$3D: (4 + 16\log(2)) (\alpha + \beta b^2)$$

SUMMA is implemented in reknown libraries

SUMMA is implemented in reknown libraries

- Multiple SUMMA variants are implemented in **Sca1aPACK**
- A variant using pipeline is implemented in **Chameleon** using **StarPU**

The **PARSEC** runtime has helped with distributed **GEMM** algorithms

SUMMA is implemented in reknown libraries

- Multiple SUMMA variants are implemented in **ScaLaPACK**
- A variant using pipeline is implemented in **Chameleon** using **StarPU**

The **PARSEC** runtime has helped with distributed **GEMM** algorithms

These implementations use either

- **explicit communication management**
- **specific data structures** required by their programming models

SUMMA

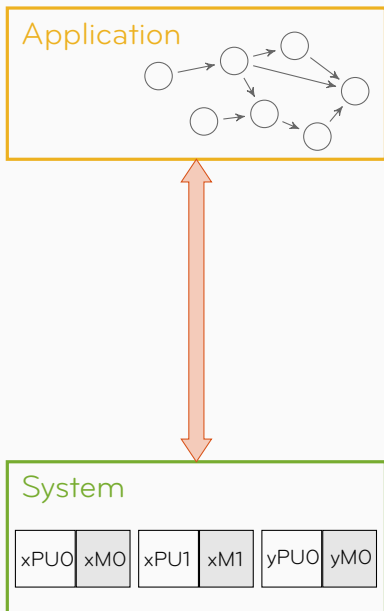
SUMMA

- Data Mapping
- Task Mapping to describe 2.5D decomposition
- Collective communications

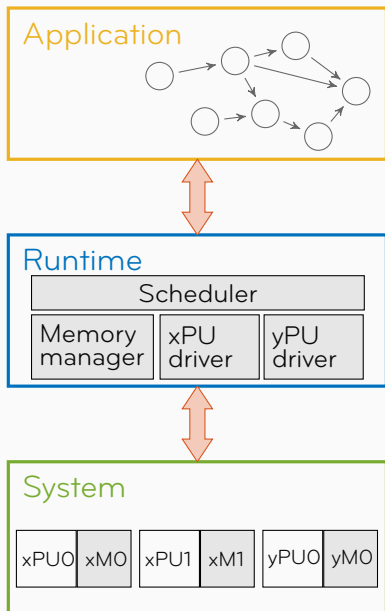
SUMMA

- Data Mapping
- Task Mapping to describe 2.5D decomposition
- Collective communications

Can we feature **complex** communications and reduction **patterns** using **programming models** offered by runtime systems ?



- The classical approach is based on a mixture of technologies (e.g., MPI+OpenMP+CUDA) which
 - requires a big programming effort
 - is difficult to maintain and update
 - is prone to (performance) portability issues



- The classical approach is based on a mixture of technologies (e.g., MPI+OpenMP+CUDA) which
 - requires a big programming effort
 - is difficult to maintain and update
 - is prone to (performance) portability issues
- **runtimes** provide an abstraction layer that hides the architecture details:
 - portable programming interface
 - handling of data
 - support for multiple architectures
 - customizable scheduling policies
 - ...

TASK-BASED MATRIX MULTIPLICATION WITH STF

```
do i=1,m
  do j=1,n
    do l=1,k
      call gemm(
                A[i,l],
                B[l,j],
                C[i,j] )
    end do
  end do
end do
```

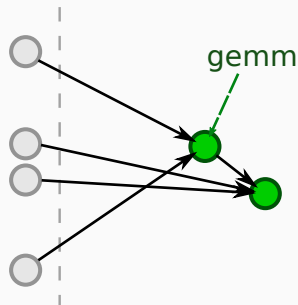
P0	P1	P2	P3
P4	P5	P6	P7

P0	P1
P4	P5
P8	P9
P12	P13

P0	P1	P2	P3
P4	P5	P6	P7
P8	P9	P10	P11
P12	P13	P14	P15

Directed Acyclic Graph (DAG) :

- **Node** : Elementary task (gemm)
- **Edge** : Sequential dependency



DAG for $C_{2,2}$

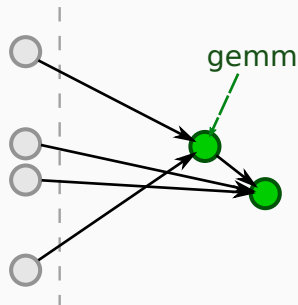
TASK-BASED MATRIX MULTIPLICATION WITH STF

```
do i=1,m
  do j=1,n
    do l=1,k
      call gemm(
                A[i,l],
                B[l,j],
                c[i,j] )
    end do
  end do
end do
```

Directed Acyclic Graph (DAG) :

- **Node** : Elementary task (`gemm`)
- **Edge** : Sequential dependency

Sequential Task Flow : how to describe the DAG to the runtime



DAG for $C_{2,2}$

TASK-BASED MATRIX MULTIPLICATION WITH STF

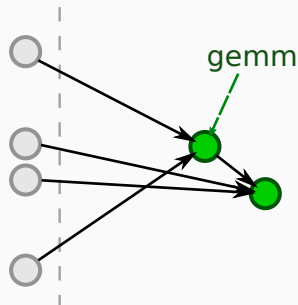
```
do i=1,m
  do j=1,n
    do l=1,k
      call submit( gemm,
                  A[i,l],
                  B[l,j],
                  c[i,j])
    end do
  end do
end do
```

Sequential Task Flow : how to describe the DAG to the runtime

- Submit tasks **sequentially**

Directed Acyclic Graph (DAG) :

- **Node** : Elementary task (`gemm`)
- **Edge** : Sequential dependency



DAG for $C_{2,2}$

TASK-BASED MATRIX MULTIPLICATION WITH STF

```
do i=1,m
  do j=1,n
    do l=1,k
      call submit( gemm,
                  A[i,l]:R,
                  B[l,j]:R,
                  C[i,j]:RW)
    end do
  end do
end do
```

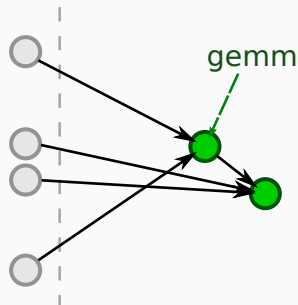
This works well in **shared memory** !

Directed Acyclic Graph (DAG) :

- **Node** : Elementary task (gemm)
- **Edge** : Sequential dependency

Sequential Task Flow : how to describe the DAG to the runtime

- Submit tasks **sequentially**
- Access modes Read, Write, ReadWrite

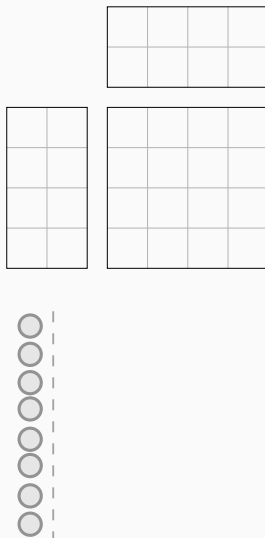


DAG for $C_{2,2}$


```

!
do i=1,m
  do j=1,n
    do l=1,k
      call submit( gemm,
                  A[i,l]:R,
                  B[l,j]:R,
                  C[i,j]:RW)
    end do
  end do
end do
    
```

Baseline model (M0)



```

! data_map: A,B,C distributed on a grid
do i=1,m
  do j=1,n
    do l=1,k
      call submit( gemm,
                  A[i,l]:R,
                  B[l,j]:R,
                  C[i,j]:RW)
    end do
  end do
end do
    
```

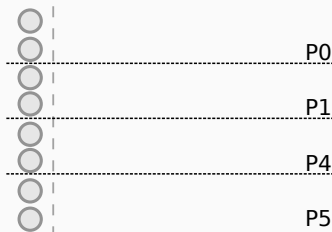
Baseline model (M0)

- M0-0 *Data Mapping*

P0	P1	P2	P3
P4	P5	P6	P7

P0	P1
P4	P5
P8	P9
P12	P13

P0	P1	P2	P3
P4	P5	P6	P7
P8	P9	P10	P11
P12	P13	P14	P15



```

! data_map: A,B,C distributed on a grid
do i=1,m
  do j=1,n
    do l=1,k
      call submit( gemm,
                   A[i,l]:R,
                   B[l,j]:R,
                   C[i,j]:RW)
    end do
  end do
end do
    
```

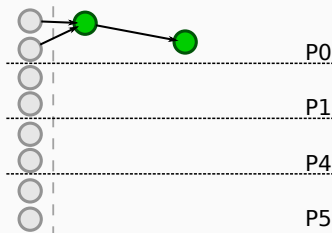
P0	P1	P2	P3
P4	P5	P6	P7

P0	P1
P4	P5
P8	P9
P12	P13

P0	P1	P2	P3
P4	P5	P6	P7
P8	P9	P10	P11
P12	P13	P14	P15

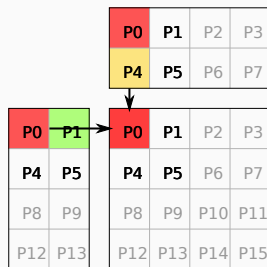
Baseline model (M0)

- M0-0 *Data Mapping*
- M0-1 *Task Mapping inferred from DM (RW)*



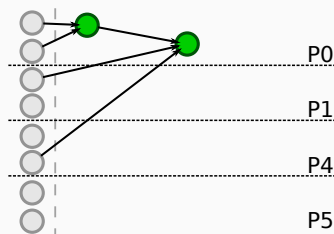
```

! data_map: A,B,C distributed on a grid
do i=1,m
  do j=1,n
    do l=1,k
      call submit( gemm,
                   A[i,l]:R,
                   B[l,j]:R,
                   C[i,j]:RW)
    end do
  end do
end do
    
```



Baseline model (M0)

- M0-0 *Data Mapping*
- M0-1 *Task Mapping inferred* from DM (RW)
- M0-2 *Point-to-Point Comm. inferred* from TM

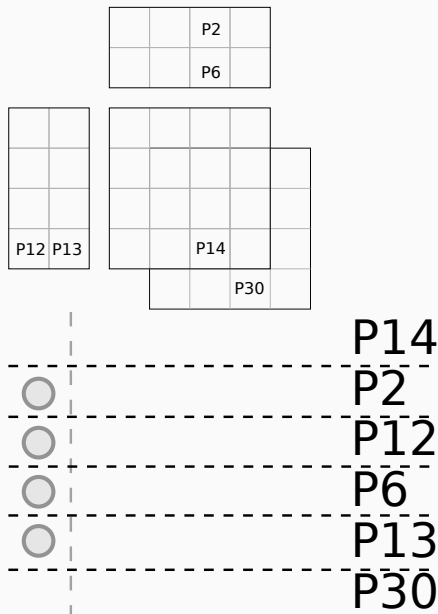


LIMITS ON IMPLEMENTING A TASK-BASED 2.5D GEMM ALGORITHM

LIMIT L1 : HOW TO DISTRIBUTE A DOT PRODUCT ?

```
! data_map: A,B,C distributed on a grid
do i=1,m
  do j=1,n
    do l=1,k
      call submit( gemm,
                  A[i,l]:R,
                  B[l,j]:R,
                  C[i,j]:RW)
    end do
  end do
end do
```

Limit L1

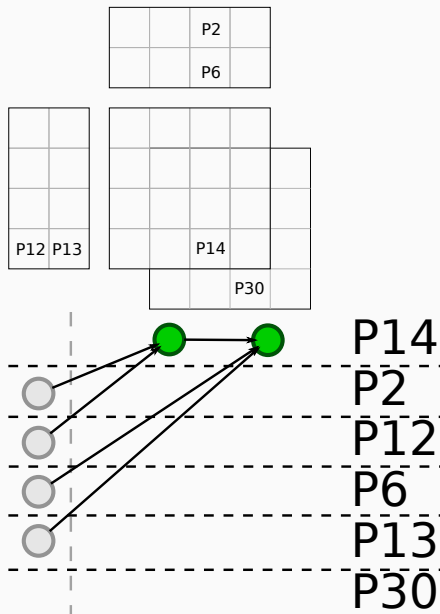


LIMIT L1 : HOW TO DISTRIBUTE A DOT PRODUCT ?

```
! data_map: A,B,C distributed on a grid
do i=1,m
  do j=1,n
    do l=1,k
      call submit( gemm,
                  A[i,l]:R,
                  B[l,j]:R,
                  C[i,j]:RW)
    end do
  end do
end do
```

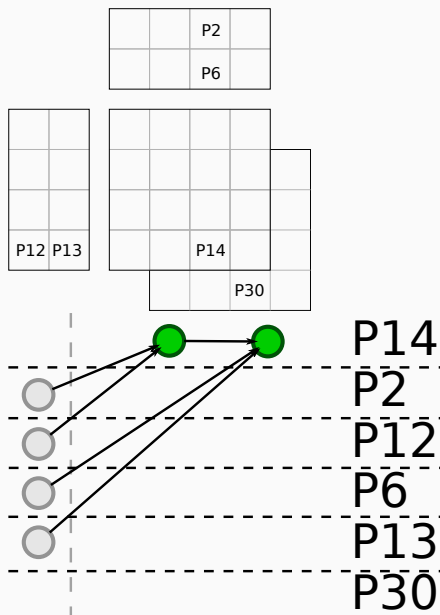
Limit L1

M0-1 : the Data Mapping forces
P14 to execute the tasks



LIMIT L2 : HOW COMPACT IS IT ?

```
! data_map: A,B,C distributed on a grid
do i=1,m
  do j=1,n
    do l=1,k
      call submit( gemm,
                  A[i,l]:R,
                  B[l,j]:R,
                  C[i,j]:RW)
    end do
  end do
end do
```



Limit L2

New tasks end up in the DAG

- `init` tasks
- `sum` tasks

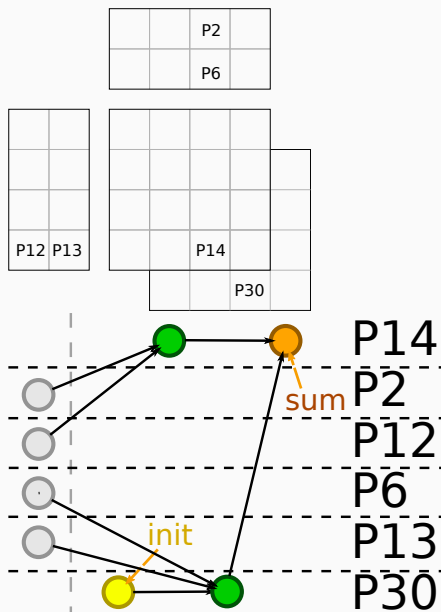
LIMIT L2 : HOW COMPACT IS IT ?

```
! data_map: A,B,C distributed on a grid
instantiate C_copy on every layer
do i=1,m
  do j=1,n
    do l=1,k
      call submit( gemm,
                  A[i,l]:R,
                  B[l,j]:R,
                  C_copy[i,j]:RW)
    end do
  end do
end do
do l=1,h
  do i=1,m; do j=1,n
    call submit(sum,
                C_copy[i,j]:R,
                C[i,j]:RW )
  end do
end do
```

Limit L2

New tasks end up in the DAG

- `init` tasks
- `sum` tasks



LIMIT L3 : COLLECTIVE COMMUNICATIONS ?

```
! data_map: A,B,C distributed on a grid
do i=1,m
  do j=1,n
    do l=1,k
      call submit( gemm,
                   A[i,l]:R,
                   B[l,j]:R,
                   C[i,j]:RW)
    end do
  end do
end do
```

P0	P1	P2	P3
P4	P5	P6	P7

P0	P1
P4	P5
P8	P9
P12	P13

P0	P1	P2	P3
P4	P5	P6	P7
P8	P9	P10	P11
P12	P13	P14	P15

Limit L3

LIMIT L3 : COLLECTIVE COMMUNICATIONS ?

```
! data_map: A,B,C distributed on a grid
do i=1,m
  do j=1,n
    do l=1,k
      call submit( gemm,
                  A[i,l]:R,
                  B[l,j]:R,
                  C[i,j]:RW)
    end do
  end do
end do
```

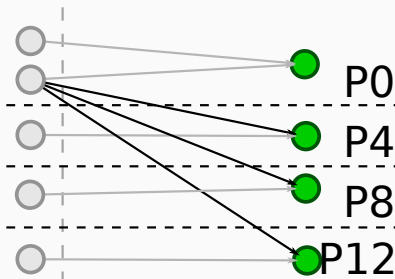
Limit L3

M0-2 : *Point-to-point*

Communications

P0	P1	P2	P3
P4	P5	P6	P7

P0	P1	P2	P3
P4	P5	P6	P7
P8	P9	P10	P11
P12	P13	P14	P15



Limits on M0

- **L1** Over-constrained Task Mapping
- **L2** Tedious reduction pattern
- **L3** Point-to-point Communications)

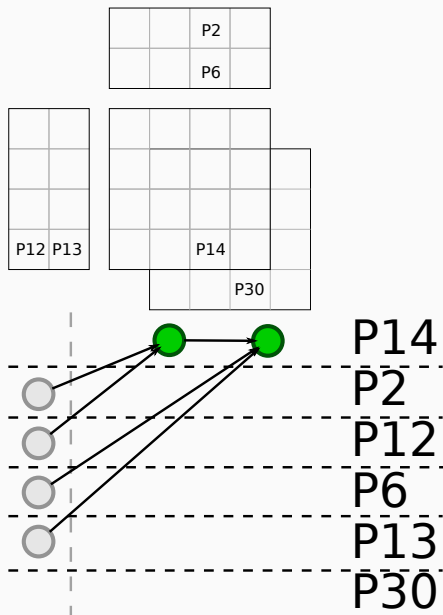
EXTENDING THE BASELINE MODEL

DISTRIBUTING A DOT PRODUCT (L1)

```
!data_map: A,B,C distributed on a grid
!  
do i=1,m  
  do j=1,n  
    do l=1,k  
      call submit(gemm,  
                 A[i,1]:R,  
                 B[1,j]:R,  
                 C[i,j]:RW)  
    end do  
  end do  
end do
```

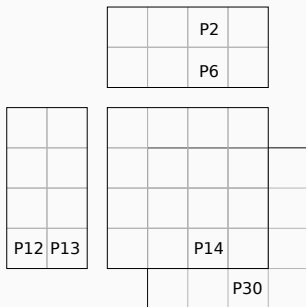
Extended model MX

- M0-0 Data Mapping
- M0-1 *Task Mapping inferred from Data Mapping*
- M0-2 Communications inferred from Task Mapping



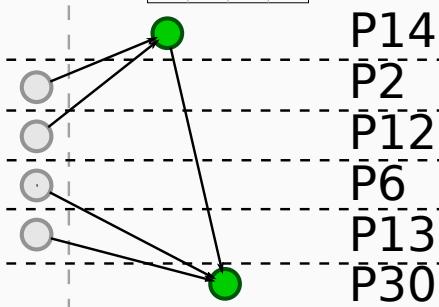
DISTRIBUTING A DOT PRODUCT (L1)

```
!data_map: A,B,C distributed on a grid
!task_map: executing node of A[i,1]B[1,j]
do i=1,m
  do j=1,n
    do l=1,k
      call submit(gemm,
                  A[i,1]:R,
                  B[1,j]:R,
                  C[i,j]:RW,
                  task_map[i,j,l])
    end do
  end do
end do
! reduction tasks submission
```



Extended model MX

- M0-0 Data Mapping
- MX-1 **Explicit** Task Mapping
- M0-2 Communications inferred from Task Mapping

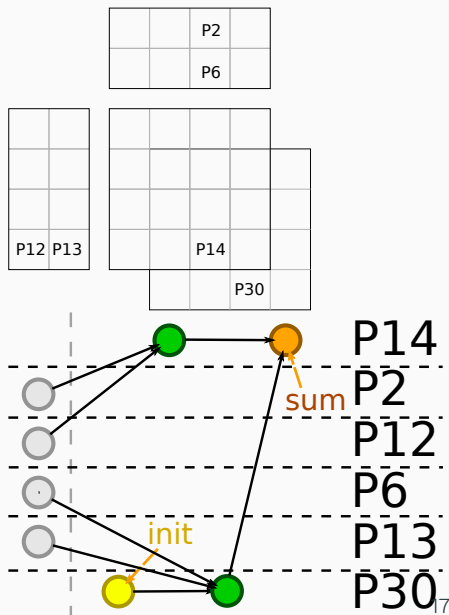


DELEGATING THE REDUCTION PATTERN (L2)

```
!data_map: A,B,C distributed on one layer
!task_map: executing node of A[i,1]B[1,j]
instantiate C_copy on every layer
do i=1,m
  do j=1,n
    do l=1,k
      call submit(gemm,
                 A[i,1]:R,
                 B[1,j]:R,
                 C_copy[i,j]:RW,
                 task_map[i,j,l])
    end do
  end do
end do
reduce C using C_copy
```

Extended model MX

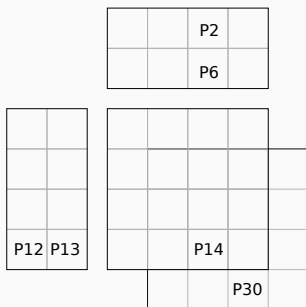
- M0-0 Data Mapping
- MX-1 Explicit Task Mapping
- M0-2 Communications inferred from Task Mapping
- MX-3



DELEGATING THE REDUCTION PATTERN (L2)

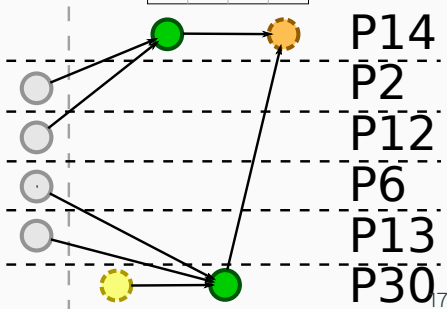
```
!data_map: A,B,C distributed on a grid  
!task_map: executing node of A[i,1]B[1,j]
```

```
do i=1,m  
  do j=1,n  
    do l=1,k  
      call submit( gemm,  
                  A[i,1]:R,  
                  B[1,j]:R,  
                  C[i,j]:MPI_REDUX,  
                  task_map[i,j,l])  
    end do  
  end do  
end do
```



Extended model MX

- M0-0 Data Mapping
- MX-1 Explicit Task Mapping
- M0-2 Communications inferred from Task Mapping
- MX-3 **Implicit** reduction pattern

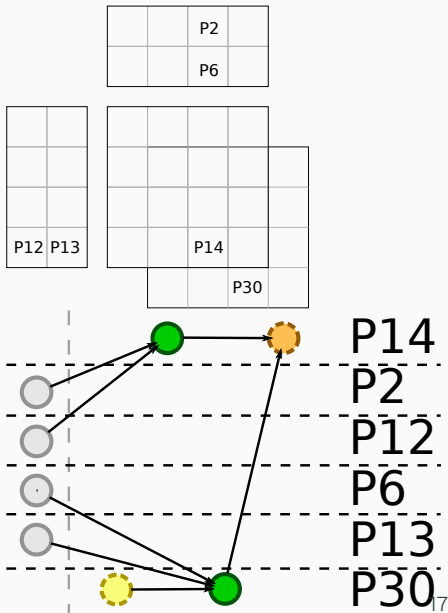


DELEGATING THE REDUCTION PATTERN (L2)

```
!data_map: A,B,C distributed on a grid
!task_map: executing node of A[i,1]B[1,j]
call bind_methods(C, init, sum)
do i=1,m
  do j=1,n
    do l=1,k
      call submit( gemm,
                  A[i,1]:R,
                  B[1,j]:R,
                  C[i,j]:MPI_REDUX,
                  task_map[i,j,l])
    end do
  end do
end do
```

Extended model MX

- M0-0 Data Mapping
- MX-1 Explicit Task Mapping
- M0-2 Communications inferred from Task Mapping
- MX-3 **Implicit** reduction pattern



DETECTING COLLECTIVE COMMUNICATIONS (L3)

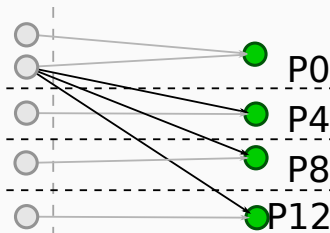
```
!data_map: A,B,C distributed on a grid
!task_map: executing node of A[i,l]B[l,j]
!dyn_tree: type of broadcast tree
call bind_methods(C, init, sum)
do i=1,m
  do j=1,n
    do l=1,k
      call submit( gemm,
                  A[i,l]:R,
                  B[l,j]:R,
                  C[i,j]:MPI_REDUX,
                  task_map[i,j,l])
    end do
  end do
end do
```

P0	P1	P2	P3
P4	P5	P6	P7

P0	P1	P2	P3
P4	P5	P6	P7
P8	P9	P10	P11
P12	P13	P14	P15

Extended model MX

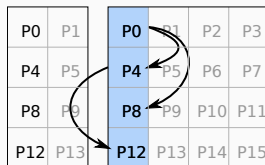
- M0-0 Data Mapping
- MX-1 Explicit Task Mapping
- MX-2 Comms. from TM
- MX-3 Implicit reduction pattern



DETECTING COLLECTIVE COMMUNICATIONS (L3)

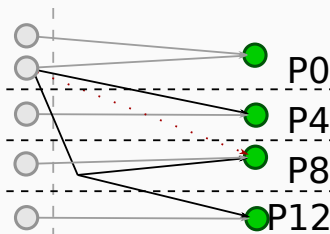
```
!data_map: A,B,C distributed on a grid
!task_map: executing node of A[i,l]B[l,j]
!dyn_tree: type of broadcast tree
call bind_methods(C, init, sum)
do i=1,m
  do j=1,n
    do l=1,k
      call submit( gemm,
                  A[i,l]:R,
                  B[l,j]:R,
                  C[i,j]:MPI_REDUX,
                  task_map[i,j,l])
    end do
  end do
end do
```

P0	P1	P2	P3
P4	P5	P6	P7



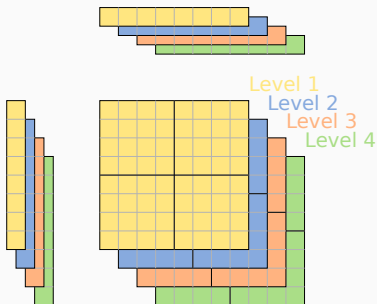
Extended model MX

- M0-0 Data Mapping
- MX-1 Explicit Task Mapping
- MX-2 **Coll.** Comms. from TM
- MX-3 Implicit reduction pattern



IDEAL EXPRESSION

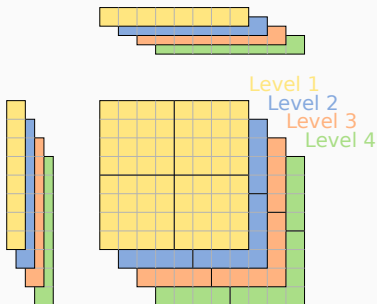
```
!data_map: A,B,C distributed on a grid
!  
!  
do i=1,m  
  do j=1,n  
    do l=1,k  
      call submit( gemm,  
                  A[i,l]:R,  
                  B[l,j]:R,  
                  C[i,j]:RW)  
    end do  
  end do  
end do
```



Extended model (MX)

- M0-0 *Data Mapping*

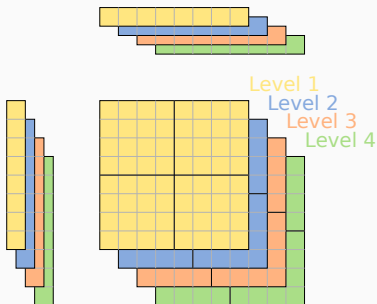
```
!data_map: A,B,C distributed on a grid
!  
!  
do i=1,m  
  do j=1,n  
    do l=1,k  
      call submit( gemm,  
                  A[i,l]:R,  
                  B[l,j]:R,  
                  C[i,j]:RW)  
    end do  
  end do  
end do
```



Extended model (MX)

- M0-0 *Data Mapping*
- M0-1 *Task Mapping* **inferred** from Data Mapping (RW)

```
!data_map: A,B,C distributed on a grid
!task_map: executing node of A[i,l]B[l,j]
!
do i=1,m
  do j=1,n
    do l=1,k
      call submit( gemm,
                  A[i,l]:R,
                  B[l,j]:R,
                  C[i,j]:RW,
                  task_map[i,j,l])
    end do
  end do
end do
```

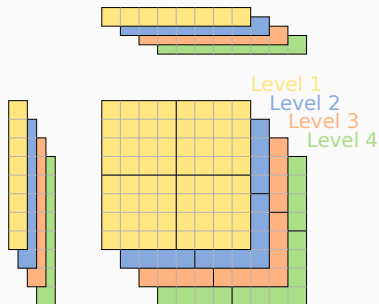


Extended model (MX)

- M0-0 *Data Mapping*
- MX-1 **Explicit** *Task Mapping*


```

!data_map: A,B,C distributed on a grid
!task_map: executing node of A[i,l]B[l,j]
!
do i=1,m
  do j=1,n
    do l=1,k
      call submit( gemm,
                  A[i,l]:R,
                  B[l,j]:R,
                  C[i,j]:RW,
                  task_map[i,j,l])
    end do
  end do
end do
    
```

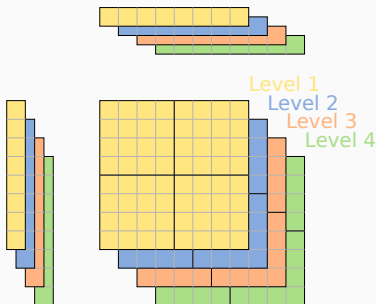


Extended model (MX)

- M0-0 *Data Mapping*
- MX-1 **Explicit** *Task Mapping*
- M0-2 Communicat^o **inferred** from Task Mapping

```
!data_map: A,B,C distributed on a grid
!task_map: executing node of A[i,l]B[l,j]
!dyn_tree: type of diffusion tree

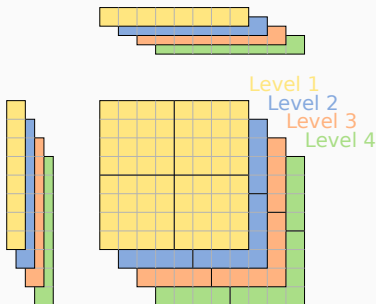
do i=1,m
  do j=1,n
    do l=1,k
      call submit( gemm,
                  A[i,l]:R,
                  B[l,j]:R,
                  C[i,j]:RW,
                  task_map[i,j,l])
    end do
  end do
end do
```



Extended model (MX)

- M0-0 *Data Mapping*
- MX-1 **Explicit** *Task Mapping*
- MX-2 *Collective Communicat^o* **inferred** from Task Mapping

```
!data_map: A,B,C distributed on a grid
!task_map: executing node of A[i,l]B[l,j]
!dyn_tree: type of broadcast tree
call bind_methods(C, init, red)
do i=1,m
  do j=1,n
    do l=1,k
      call submit( gemm,
                  A[i,l]:R,
                  B[l,j]:R,
                  C[i,j]:MPI_REDUX,
                  task_map[i,j,l])
    end do
  end do
end do
```



Extended model (MX)

- M0-0 *Data Mapping*
- MX-1 **Explicit** *Task Mapping*
- MX-2 *Collective Communicat^o* **inferred** from Task Mapping
- MX-3 **implicit** reduction pattern

ACTUAL IMPLEMENTATION AND RESULTS

- StarPU
 - STF-based *runtime*
 - Developed by Storm Team in INRIA - Bordeaux

- StarPU
 - STF-based *runtime*
 - Developed by Storm Team in INRIA - Bordeaux
- **L1** StarPU provides the API

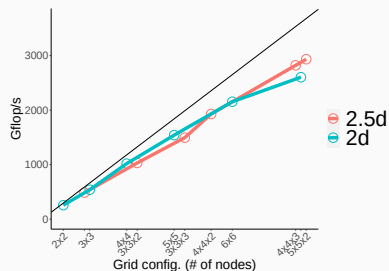
- StarPU
 - STF-based *runtime*
 - Developed by Storm Team in INRIA - Bordeaux
- **L1** StarPU provides the API
- **L2** StarPU provides the API *in shared memory*
 - Adaptations are required for distributed computing

- StarPU
 - STF-based *runtime*
 - Developed by Storm Team in INRIA - Bordeaux
- NewMadeleine
 - Communications engine based on MPI standards
 - Developed by TaDaaM Team in INRIA- Bordeaux
- **L1** StarPU provides the API
- **L2** StarPU provides the API *in shared memory*
 - Adaptations are required for distributed computing
- **L3** NewMadeleine & StarPU provide the detection
 - Ongoing research (Philippe Swartvagher's PhD. thesis)

PRELIMINARY RESULTS

Ran on **CALMIP** Olympe with *Intel Gold Skylake CPUs*

- 1 core used per node
 - Strong scaling using sq. mat. of size 65k
- ≈ 80% peak performance



PRELIMINARY RESULTS

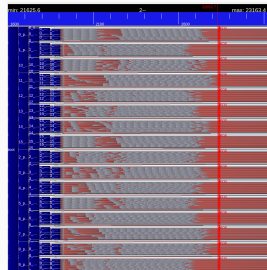
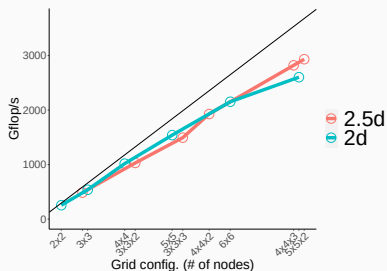
Ran on **CALMIP** Olympe with *Intel Gold Skylake CPUs*

- 1 core used per node
 - Strong scaling using sq. mat. of size 65k
- ≈ 80% peak performance

- 16 cores used per node
 - sq. mat. of size 16k
 - limited to 2D
- ≈ 55% peak performance

Trace legend

- grey : gemm tasks
- red : idle time



CONCLUSION

Contributions

- Study of extensions to the STF model to **compactly** write **distributed** algorithms
 - Model M0 → Model MX
- Design and implementation of a prototype using complex patterns
 - Exchanges on StarPU (see gitlab issue#2, issue#3)

Contributions

- Study of extensions to the STF model to **compactly** write **distributed** algorithms
 - Model M0 → Model MX
- Design and implementation of a prototype using complex patterns
 - Exchanges on StarPU (see gitlab issue#2, issue#3)

Future Works

- Broaden the scope of the experiments
- Further extend the model
 - Implement other algorithms (i.e. factorizations)
 - Adapt to sparse algebra

Thank you

APPENDIX

Behaviour

Test Program

With $2n$ nodes (n pairs)

ADJUSTING THE REDUCTION

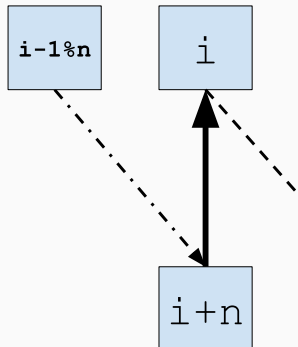
Test Program

With $2n$ nodes (n pairs)

```
a_i = 0, c_i = my_rank
```

Behaviour

```
a_i=0 ; c_i=my_rank
```



ADJUSTING THE REDUCTION

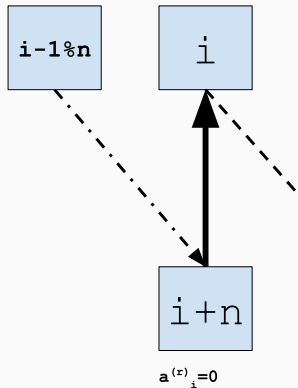
Test Program

With $2n$ nodes (n pairs)

```
a_i = 0, c_i = my_rank  
a_i^(r) = 0
```

Behaviour

$a_i=0$; $c_i=my_rank$



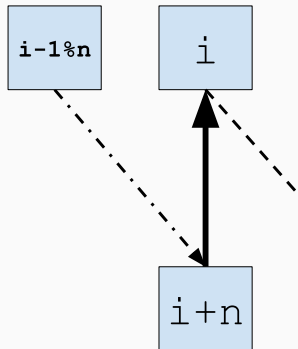
Test Program

With $2n$ nodes (n pairs)

```
a_i = 0, c_i = my_rank  
a_i^(r) = 0  
a_i^(r) += c_(i-1%n)
```

Behaviour

$a_i=0$; $c_i=my_rank$



$a_i^{(x)}=0$

$a_i^{(x)} += c_{i-1\%n}$

Test Program

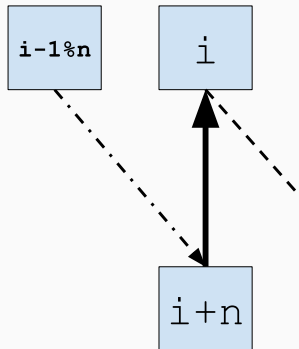
With $2n$ nodes (n pairs)

```
a_i = 0, c_i = my_rank  
a_i^(r) = 0  
a_i^(r) += c_(i-1%n)  
a_i += a_i^(r)
```

Behaviour

$a_i = 0 ; c_i = \text{my_rank}$

$a_i += a_i^{(r)}$



$a_i^{(r)} = 0$

$a_i^{(r)} += c_{i-1\%n}$

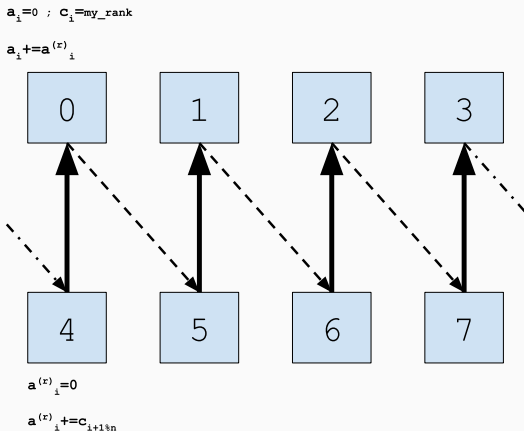
Test Program

With $2n$ nodes (n pairs)

```

a_i = 0, c_i = my_rank
a_i^(r) = 0
a_i^(r) += c_(i-1%n)
a_i += a_i^(r)

```

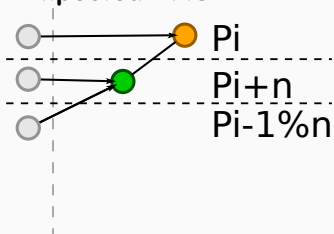
Behaviour ($n=4$)

Test Program

With $2n$ nodes (n pairs)

```
a_i = 0, c_i = my_rank  
a_i^(r) = 0  
a_i^(r) += c_(i-1%n)  
a_i += a_i^(r)
```

Expected DAG



Test Program

With $2n$ nodes (n pairs)

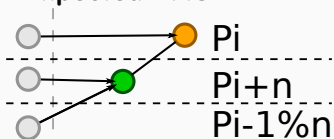
```

a_i = 0, c_i = my_rank
a_i^(r) = 0
a_i^(r) += c_(i-1%n)
a_i += a_i^(r)
    
```

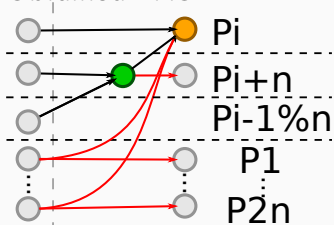
Issues

- All nodes partake in the reduction

Expected DAG



Obtained DAG



Test Program

With $2n$ nodes (n pairs)

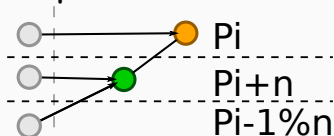
```

a_i = 0, c_i = my_rank
a_i^(r) = 0
a_i^(r) += c_(i-1%n)
a_i += a_i^(r)
    
```

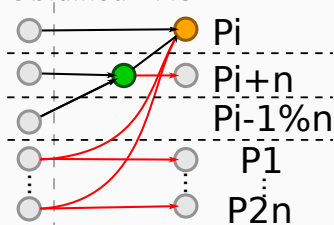
Issues

- All nodes partake in the reduction
- ⇒ Many instances w/
REDUX access

Expected DAG



Obtained DAG



Test Program

With $2n$ nodes (n pairs)

```
a_i = 0, c_i = my_rank  
a_i^(r) = 0  
a_i^(r) += c_(i-1%n)  
a_i += a_i^(r)
```

Issues

- All nodes partake in the reduction

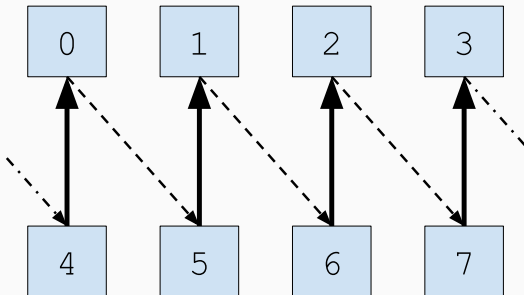
⇒ Many instances w/
REDUX access

Solution

Splitting MPI_COMM_WORLD

$a_i = 0 ; c_i = \text{my_rank}$

$a_i += a_i^{(x)}$



$a_i^{(x)} = 0$

$a_i^{(x)} += c_{i+1\%n}$

Test Program

With $2n$ nodes (n pairs)

```
a_i = 0, c_i = my_rank
a_i^(r) = 0
a_i^(r) += c_(i-1%n)
a_i += a_i^(r)
```

Issues

- All nodes partake in the reduction

⇒ Many instances w/
REDUX access

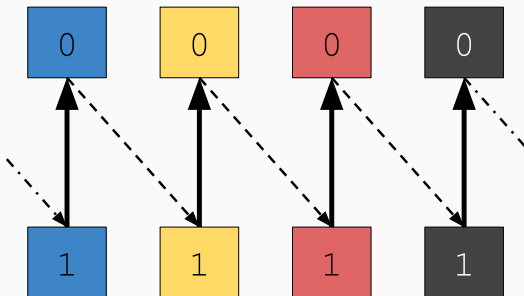
Solution

- One communicator per pair to *reduce*

Splitting MPI_COMM_WORLD

$a_i = 0 ; c_i = \text{my_rank}$

$a_i += a_i^{(x)}$

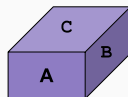


$a_i^{(x)} = 0$

$a_i^{(x)} += c_{i+1\%n}$

Tensor representation

one element \leftrightarrow one task

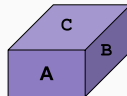


Reminder

- $C_{ij+} = A_{il} * B_{lj}$ on *task_{ijl}*

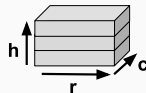
Tensor representation

one element \leftrightarrow one task



Grid of nodes

h layers of $r \times c$ nodes



Reminder

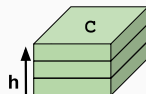
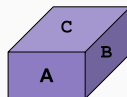
- $C_{ij} = A_{il} * B_{lj}$ on task ijl
- X_{ij} on $j \% c + c * (i \% r)$ (2D Block Cyclic distrib.)

Task_map SUMMA-C

- $task_{ijl} :=$
 $j\%c + c * (i\%r) + r * c * (l\%h)$

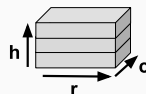
Tensor representation

one element \leftrightarrow one task



Grid of nodes

h layers of $r * c$ nodes



Reminder

- $C_{ij} = A_{il} * B_{lj}$ on $task_{ijl}$
- X_{ij} on $j\%c + c * (i\%r)$ (2D Block Cyclic distrib.)

Task_map SUMMA-C

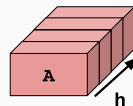
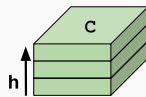
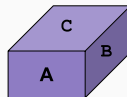
- $task_{ijl} :=$
 $j\%c + c * (i\%r) + r * c * (l\%h)$

Task_map SUMMA-A

- $task_{ijl} :=$
 $l\%c + c * (i\%r) + r * c * (j\%h)$

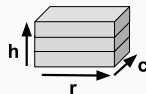
Tensor representation

one element \leftrightarrow one task



Grid of nodes

h layers of $r * c$ nodes



Reminder

- $C_{ij} += A_{il} * B_{lj}$ on $task_{ijl}$
- X_{ij} on $j\%c + c * (i\%r)$ (2D Block Cyclic distrib.)

Task_map SUMMA-C

- $task_{ijl} :=$
 $j\%c + c * (i\%r) + r * c * (l\%h)$

Task_map SUMMA-A

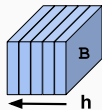
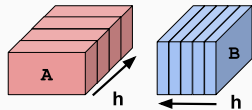
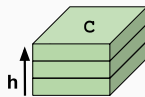
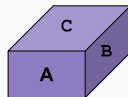
- $task_{ijl} :=$
 $l\%c + c * (i\%r) + r * c * (j\%h)$

Task_map SUMMA-B

- $task_{ijl} :=$
 $j\%c + c * (l\%r) + r * c * (i\%h)$

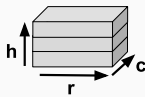
Tensor representation

one element \leftrightarrow one task



Grid of nodes

h layers of $r * c$ nodes



Reminder

- $C_{ij} = A_{il} * B_{lj}$ on $task_{ijl}$
- X_{ij} on $j\%c + c * (i\%r)$ (2D Block Cyclic distrib.)