

TASK-BASED PROGRAMMING MODELS FOR SCALABLE ALGORITHMS

Application to matrix multiplication

E. Agullo, A. Buttari, A. Guermouche, J. Herrmann, A. Jégou
June 17, 2022



INTRODUCTION

Architectures are becoming increasingly **large, heterogeneous** and **complex**. This has two effects:

- **Algorithms** have to handle communications to achieve high scalability
 - Communication/computation overlap
 - Communication-avoiding algorithms
- **Software** should deliver a reliable abstraction to build mathematical software
 - parallel programming models
 - runtime systems

Are currently available parallel programming models and runtimes expressive enough to implement scalable algorithms in an efficient, portable and maintainable way?

Architectures are becoming increasingly **large, heterogeneous** and **complex**. This has two effects:

- **Algorithms** have to handle communications to achieve high scalability
 - Communication/computation overlap
 - Communication-avoiding algorithms
- **Software** should deliver a reliable abstraction to build mathematical software
 - parallel programming models
 - runtime systems

Are currently available parallel programming models and runtimes expressive enough to implement scalable algorithms in an efficient, portable and maintainable way?

Architectures are becoming increasingly **large, heterogeneous** and **complex**. This has two effects:

- **Algorithms** have to handle communications to achieve high scalability
 - Communication/computation overlap
 - Communication-avoiding algorithms
- **Software** should deliver a reliable abstraction to build mathematical software
 - parallel programming models
 - runtime systems

Are currently available parallel programming models and runtimes expressive enough to implement scalable algorithms in an efficient, portable and maintainable way?

Architectures are becoming increasingly **large, heterogeneous** and **complex**. This has two effects:

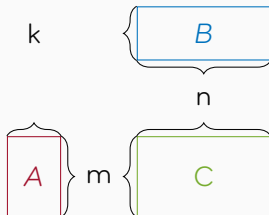
- **Algorithms** have to handle communications to achieve high scalability
 - Communication/computation overlap
 - Communication-avoiding algorithms
- **Software** should deliver a reliable abstraction to build mathematical software
 - parallel programming models
 - runtime systems

Are currently available parallel programming models and runtimes expressive enough to implement scalable algorithms in an efficient, portable and maintainable way?

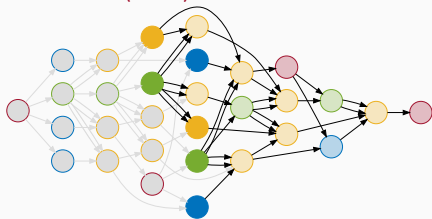
We take the following examples to answer this large question

- **Algorithms:** scalable dense matrix multiplication (GEMM) algorithms such as **SUMMA** and **2.5D SUMMA**

$$C = \alpha AB + \beta C$$



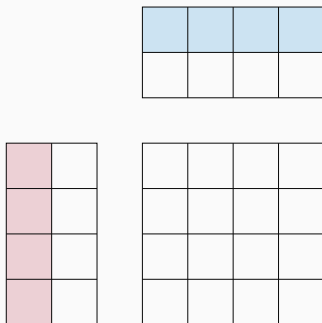
- **Software:** **task-based** parallel programming paradigm through the **sequential task flow (STF)** model



BACKGROUND

2D C-stationary SUMMA

- **Data mapping**
 - 2D block-cyclic
- **Workload mapping**
 - owner of C computes
- **Communications**
 - row-wise broadcast of A
 - column-wise broadcast of B



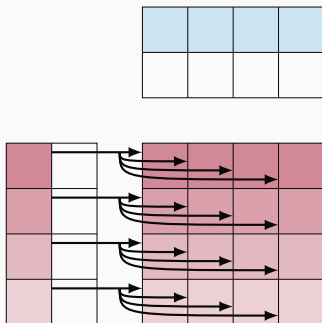
```

! I am r x c
do l=1, k
  forall i, i%p==r :
    bcast(Ai,l, r×*)
  forall j, j%q==c :
    bcast(Bl,j, *×c)
  forall i, i%p==r :
    forall j, j%q==c :
      call gemm(Ai,l, Bl,j, Ci,j)
end do

```

2D C-stationary SUMMA

- **Data mapping**
 - 2D block-cyclic
- **Workload mapping**
 - owner of C computes
- **Communications**
 - row-wise **broadcast** of A
 - column-wise broadcast of B



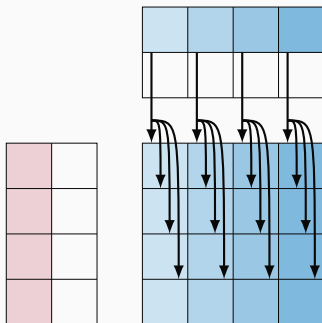
```

! I am r x c
do l=1, k
  forall i, i%p==r :
    bcast(Ai,l, r×*)
  forall j, j%q==c :
    bcast(Bl,j, *×c)
  forall i, i%p==r :
    forall j, j%q==c :
      call gemm(Ai,l, Bl,j, Ci,j)
end do

```

2D C-stationary SUMMA

- **Data mapping**
 - 2D block-cyclic
- **Workload mapping**
 - owner of C computes
- **Communications**
 - row-wise broadcast of A
 - column-wise **broadcast** of B



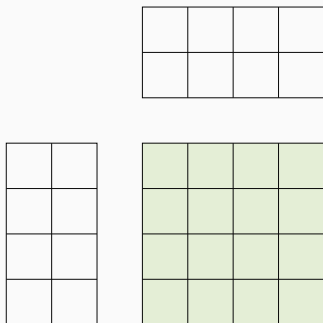
```

! I am r x c
do l=1, k
  forall i, i%p==r :
    bcast(Ai,l, r×*)
  forall j, j%q==c :
    bcast(Bl,j, *×c) ←
  forall i, i%p==r :
    forall j, j%q==c :
      call gemm(Ai,l, Bl,j, Ci,j)
end do

```

2D C-stationary SUMMA

- **Data mapping**
 - 2D block-cyclic
- **Workload mapping**
 - owner of C **computes**
- **Communications**
 - row-wise broadcast of A
 - column-wise broadcast of B



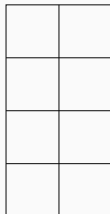
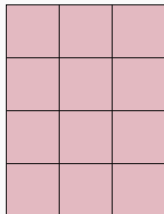
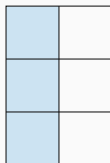
```

! I am r x c
do l=1, k
  forall i, i%p==r :
    bcast(Ai,l, r×*)
  forall j, j%q==c :
    bcast(Bl,j, *×c)
  forall i, i%p==r :
    forall j, j%q==c :
      call gemm(Ai,l, Bl,j, Ci,j) ←
end do

```

2D A-stationary SUMMA

- **Data mapping**
 - 2D block-cyclic
- **Workload mapping**
 - owner of A computes
- **Communications**
 - column-wise scatter+broadcast of B
 - row-wise reduction of C



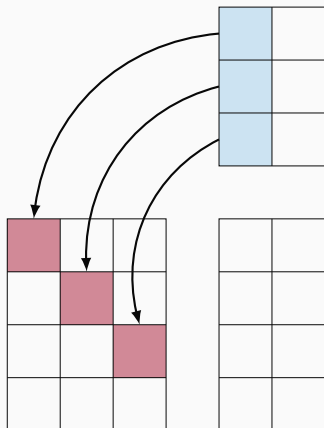
```

! I am r x c
do j=1, n
  forall 1, 1%p==r.and.1%q==c :
    recv(Bl,j, 1%p x j%q)
  forall 1, 1%p==c :
    bcast(Bl,j, *x c)
  forall i, i%p==r :
    forall 1, j%q==c :
      call gemm(Ai,l, Bl,j, Ci,i)
      reduce(Ci,i, i%r x j%q)
end do

```

2D A-stationary SUMMA

- **Data mapping**
 - 2D block-cyclic
- **Workload mapping**
 - owner of A computes
- **Communications**
 - column-wise **scatter**+broadcast of B
 - row-wise reduction of C



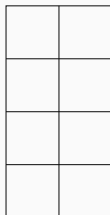
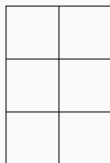
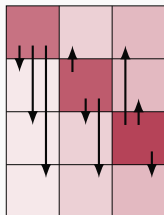
```

! I am r x c
do j=1, n
  forall 1, 1%p==r.and.1%q==c :
    recv(Bi,j, 1%p×j%q) ←
  forall 1, 1%p==c :
    bcast(Bi,j, *×c)
  forall i, i%p==r :
    forall 1, j%q==c :
      call gemm(Ai,i, Bi,j, Ci,i)
      reduce(Ci,i, i%r×j%q)
end do

```

2D A-stationary SUMMA

- **Data mapping**
 - 2D block-cyclic
- **Worload mapping**
 - owner of A computes
- **Communications**
 - column-wise scatter+**broadcast** of B
 - row-wise reduction of C



```

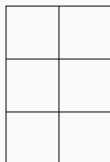
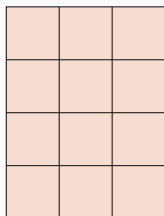
! I am r x c
do j=1, n
  forall 1, 1%p==r.and.1%q==c :
    recv(Bl,j, 1%p×j%q)
  forall 1, 1%p==c :
    bcast(Bl,j, *×c)
  forall i, i%p==r :
    forall 1, j%q==c :
      call gemm(Ai,l, Bl,j, Ci,i)
      reduce(Ci,i, i%r×j%q)
end do

```

←

2D A-stationary SUMMA

- **Data mapping**
 - 2D block-cyclic
- **Workload mapping**
 - owner of A **computes**
- **Communications**
 - column-wise scatter+broadcast of B
 - row-wise reduction of C



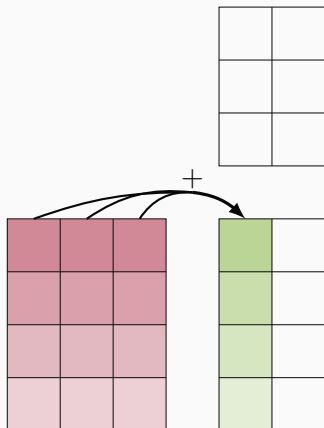
```

! I am r x c
do j=1, n
  forall 1, 1%p==r.and.1%q==c :
    recv(Bi,j, 1%p×j%q)
  forall 1, 1%p==c :
    bcast(Bi,j, *×c)
  forall i, i%p==r :
    forall 1, j%q==c :
      call gemm(Ai,i, Bi,j, Ci,i) ←
    reduce(Ci,i, i%r×j%q)
end do

```


2D A-stationary SUMMA

- **Data mapping**
 - 2D block-cyclic
- **Workload mapping**
 - owner of A computes
- **Communications**
 - column-wise scatter+broadcast of B
 - row-wise **reduction** of C



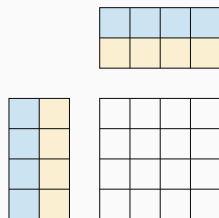
```

! I am r x c
do j=1, n
  forall 1, 1%p==r.and.1%q==c :
    recv(Bl,j, 1%p x j%q)
  forall 1, 1%p==c :
    bcast(Bl,j, *x c)
  forall i, i%p==r :
    forall 1, j%q==c :
      call gemm(Ai,l, Bl,j, Ci,i)
      reduce(Ci,i, i%r x j%q)
end do

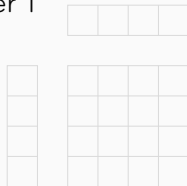
```



Layer 0



Layer 1



2.5D C-stationary SUMMA

- **Data mapping**
 - 2D block-cyclic on 0^{th} layer
- **Workload mapping**
 - Not trivially related to data mapping
- **Communications**
 - row-wise scatter+broadcast of A
 - column-wise scatter+broadcast of B
 - aisle-wise reduction of C

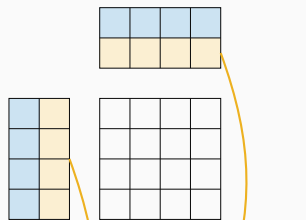
```

! I am r x c x h
scatter Ai,j on layer i/h
scatter Bi,j on layer i/h

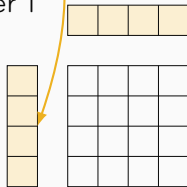
forall layer=1,h:
  call SUMMA_2DC(A[:, (layer : layer + 1) * k/h],
                B[(layer : layer + 1) * k/h, :], Ch)

forall i, i%p==r :
  forall j, j%q==c :
    reduce(Chi,i, r x c x 0)
  
```

Layer 0



Layer 1



2.5D C-stationary SUMMA

- **Data mapping**
 - 2D block-cyclic on 0^{th} layer
- **Workload mapping**
 - Not trivially related to data mapping
- **Communications**
 - row-wise **scatter**+broadcast of A
 - column-wise **scatter**+broadcast of B
 - aisle-wise reduction of C

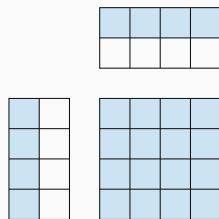
```

! I am r × c × h
scatter Ai,j on layer i/h
scatter Bi,j on layer i/h

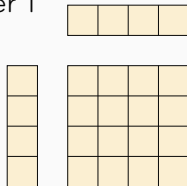
forall layer=1,h:
  call SUMMA_2DC(A[:, (layer : layer + 1) * k/h],
                B[(layer : layer + 1) * k/h, :], Ch)

forall i, i%p==r :
  forall j, j%q==c :
    reduce(Chi,i, r × c × 0)
  
```

Layer 0



Layer 1



2.5D C-stationary SUMMA

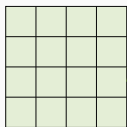
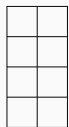
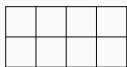
- **Data mapping**
 - 2D block-cyclic on 0^{th} layer
- **Workload mapping**
 - Not trivially related to data mapping
- **Communications**
 - row-wise scatter+**broadcast** of A
 - column-wise scatter+**broadcast** of B
 - aisle-wise reduction of C

```
! I am r x c x h
scatter Ai,j on layer i/h
scatter Bi,j on layer i/h
```

```
forall layer=1,h:
  call SUMMA_2DC(A[:, (layer : layer + 1) * k/h],
                B[(layer : layer + 1) * k/h, :], Ch) ←
```

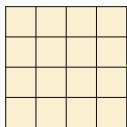
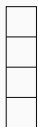
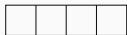
```
forall i, i%p==r :
  forall j, j%q==c :
    reduce(Ci,ih, r x c x 0)
```

Layer 0



+

Layer 1



2.5D C-stationary SUMMA

- **Data mapping**
 - 2D block-cyclic on 0^{th} layer
- **Workload mapping**
 - Not trivially related to data mapping
- **Communications**
 - row-wise scatter+broadcast of A
 - column-wise scatter+broadcast of B
 - aisle-wise **reduction** of C

```

! I am r × c × h
scatter Ai,j on layer i/h
scatter Bi,j on layer i/h

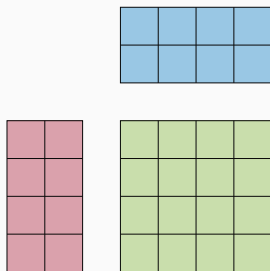
forall layer=1,h:
  call SUMMA_2DC(A[:, (layer : layer + 1) * k/h],
                B[(layer : layer + 1) * k/h, :], Ch)

forall i, i%p==r :
  forall j, j%q==c :
    reduce(Chi,i, r × c × 0)
  
```

THE SEQUENTIAL TASK FLOW (STF) MODEL

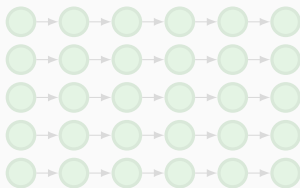
Sequential GEMM

```
do i=1,m
  do j=1,n
    do l=1,k
      call gemm(
        Ai,l,
        Bl,j,
        Ci,j )
    end do
  end do
end do
```



Shared-memory STF GEMM

```
do i=1,m
  do j=1,n
    do l=1,k
      call insert_task( gemm,
        Ai,l:R,
        Bl,j:R,
        Ci,j:RW)
    end do
  end do
end do
```

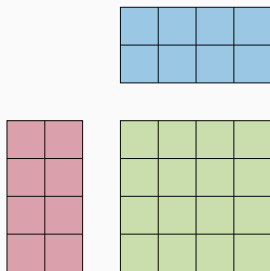


DAG of GEMM

THE SEQUENTIAL TASK FLOW (STF) MODEL

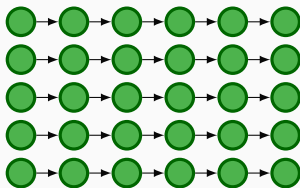
Sequential GEMM

```
do i=1,m
  do j=1,n
    do l=1,k
      call gemm(
        Al,l,
        Bl,j,
        Ci,j )
    end do
  end do
end do
```



Shared-memory STF GEMM

```
do i=1,m
  do j=1,n
    do l=1,k
      call insert_task( gemm,
        Al,l:R,
        Bl,j:R,
        Ci,j:RW)
    end do
  end do
end do
```



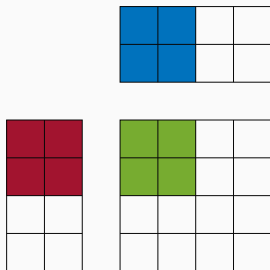
DAG of GEMM

- Can we use the **STF model** to express a **scalable distributed-memory GEMM** algorithm ?
- How **efficient** is the resulting expression compared to reference libraries ?

DISTRIBUTED-MEMORY STF GEMM ALGORITHM

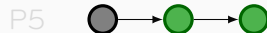
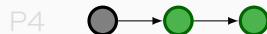
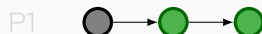
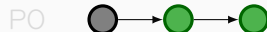
```

!
do i=1,m
  do j=1,n
    do l=1,k
      call insert_task(gemm,
                     Ai,l:R,
                     Bl,j:R,
                     Ci,j:RW)
    end do
  end do
end do
    
```



Baseline model (M0)^a

- M0-0 Data mapping
- M0-1 Task mapping **inferred** from DM (RW)
- M0-2 Point-to-point comms. **inferred** from TM



^aAgullo, Aumage, Faverge, Furmento, Pruvost, Sergent, and Thibault, 2017.

BASILINE STF MODEL

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

P0	P1	P2	P3
P4	P5	P6	P7

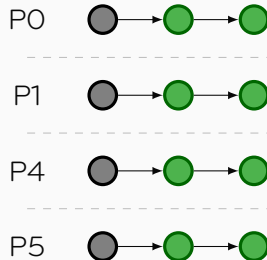
P0	P1
P4	P5
P8	P9
P12	P13

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

Baseline model (M0)^a

- M0-0 Data mapping
- M0-1 Task mapping **inferred** from DM (RW)
- M0-2 Point-to-point comms. **inferred** from TM

^aAgullo, Aumage, Faverge, Furmento, Pruvost, Sergent, and Thibault, 2017.



BASILINE STF MODEL

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

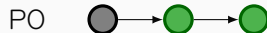
P0	P1	P2	P3
P4	P5	P6	P7

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)^a

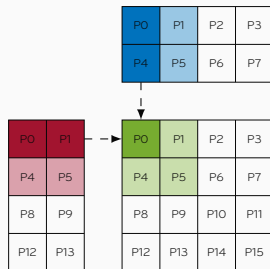
- M0-0 Data mapping
- M0-1 Task mapping **inferred** from DM (RW)
- M0-2 Point-to-point comms. **inferred** from TM



^aAgullo, Aumage, Faverge, Furmento, Pruvost, Sergent, and Thibault, 2017.

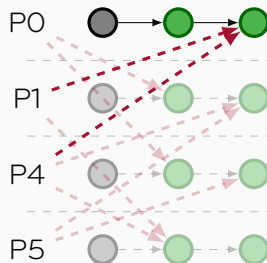
BASILINE STF MODEL

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



Baseline model (M0)^a

- M0-0 Data mapping
- M0-1 Task mapping **inferred** from DM (RW)
- M0-2 Point-to-point comms. **inferred** from TM



^aAgullo, Aumage, Faverge, Furmento, Pruvost, Sergent, and Thibault, 2017.

The baseline STF model M_0 does not allow the expression of SUMMA algorithms

- Task mapping is bound to data mapping (M_0-1)
- Reduction patterns are hard to make
- Point-to-point communications patterns are inefficient (M_0-2)

We have identified key functionalities required in a STF model M_X , extended from M_0 , to allow the expression of SUMMA algorithms.

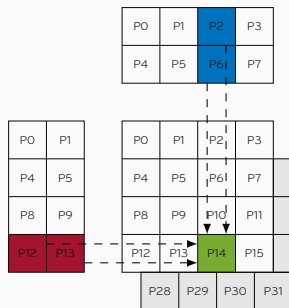
The baseline STF model **M0** does not allow the expression of SUMMA algorithms

- Task mapping is bound to data mapping (**M0-1**)
- Reduction patterns are hard to make
- Point-to-point communications patterns are inefficient (**M0-2**)

We have identified key functionalities required in a STF model **MX**, extended from **M0**, to allow the expression of SUMMA algorithms.

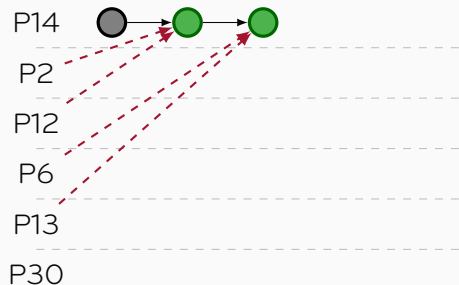
KEY 1 – TASK MAPPING

```
!data_map: A,B,C distributed on one layer
!
!
do i=1,m
  do j=1,n
    do l=1,k
      call insert_task(gemm,
                      Ai,l:R,
                      Bl,j:R,
                      Ci,j:RW)
    end do
  end do
end do
```



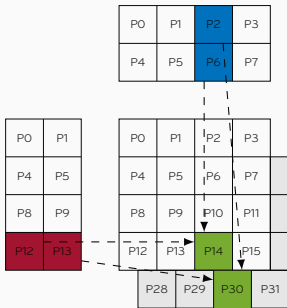
Advanced model MX

- M0-0 Data mapping
- M0-1 Implicit task mapping
- M0-2 Comms. inferred



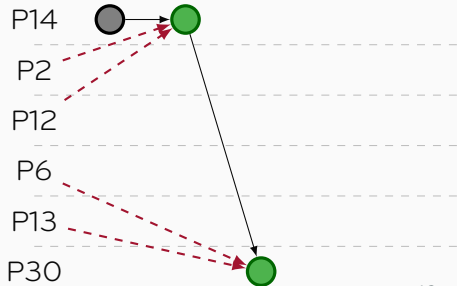
KEY 1 – TASK MAPPING

```
!data_map: A,B,C distributed on one layer
!task_map: executing node of  $A_{i,l}B_{l,j}$ 
!
do i=1,m
  do j=1,n
    do l=1,k
      call insert_task(gemm,
         $A_{i,l}$ :R,
         $B_{l,j}$ :R,
         $C_{i,j}$ :RW,
        task_map[i,j,l])
    end do
  end do
end do
```



Advanced model MX

- M0-0 Data mapping
- MX-1 **Explicit** task mapping
- M0-2 Comms. inferred

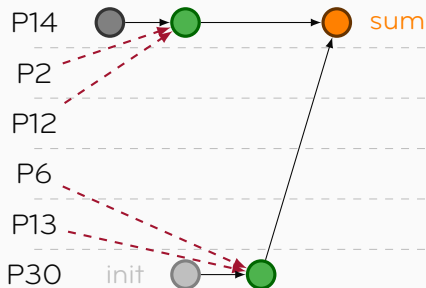
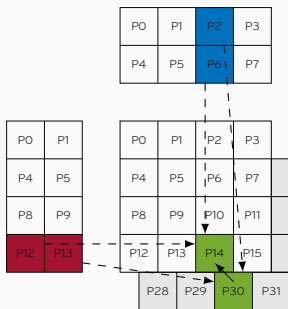


KEY 2 – REDUCTION PATTERN

```
!data_map: A,B,C distributed on one layer
!task_map: executing node of  $A_{i,l}B_{l,j}$ 
!
call instantiate_local(C)
do i=1,m
  do j=1,n
    do l=1,k
      call insert_task(gemm,
                      $A_{i,l}$ :R,
                      $B_{l,j}$ :R,
                      $C_{i,j}$ :RW,
                     task_map[i,j,l])
    end do
  end do
end do
call reduce_to_initial_layer(C)
```

Advanced model MX

- M0-0 Data mapping
- MX-1 Explicit task mapping
- M0-2 Comms. inferred
- M0-3 Explicit reduction pattern

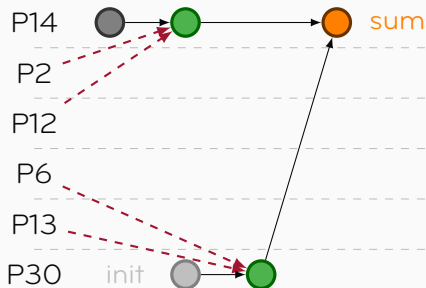
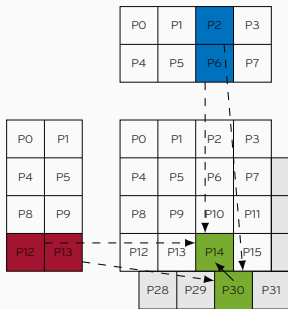


KEY 2 – REDUCTION PATTERN

```
!data_map: A,B,C distributed on one layer
!task_map: executing node of  $A_{i,l}B_{l,j}$ 
!
call instantiate_local(C)
do i=1,m
  do j=1,n
    do l=1,k
      call insert_task(gemm,
                      $A_{i,l}$ :R,
                      $B_{l,j}$ :R,
                      $C_{i,j}$ :RANK_REDUX,
                     task_map[i,j,l])
    end do
  end do
end do
call reduce_to_initial_layer(C)
```

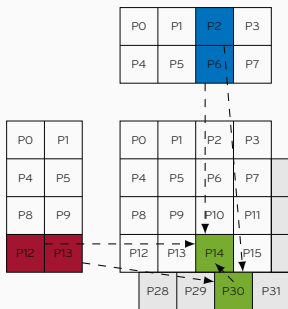
Advanced model MX

- M0-0 Data mapping
- MX-1 Explicit task mapping
- M0-2 Comms. inferred
- M0-3 Explicit reduction pattern



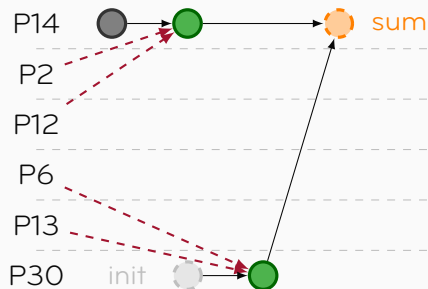
KEY 2 – REDUCTION PATTERN

```
!data_map: A,B,C distributed on one layer
!task_map: executing node of  $A_{i,l}B_{l,j}$ 
!
call bind_methods(C, init, sum)
do i=1,m
  do j=1,n
    do l=1,k
      call insert_task(gemm,
                      $A_{i,l}$ :R,
                      $B_{l,j}$ :R,
                      $C_{i,j}$ :RANK_REDUX,
                     task_map[i,j,l])
    end do
  end do
end do
```



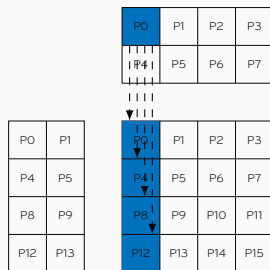
Advanced model MX

- M0-0 Data mapping
- MX-1 Explicit task mapping
- M0-2 Comms. inferred
- MX-3 **Implicit** reduction pattern



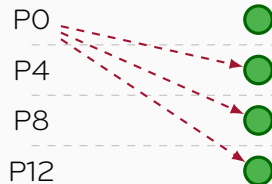
KEY 3 – COMMUNICATIONS

```
!data_map: A,B,C distributed on one layer
!task_map: executing node of  $A_{i,l}B_{l,j}$ 
!
call bind_methods(C, init, sum)
do i=1,m
  do j=1,n
    do l=1,k
      call insert_task(gemm,
                      $A_{i,l}$ :R,
                      $B_{l,j}$ :R,
                      $C_{i,j}$ :RANK_REDUX,
                     task_map[i,j,l])
    end do
  end do
end do
```



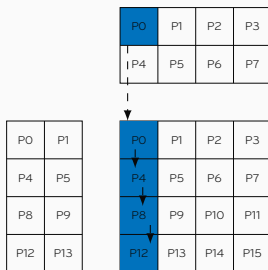
Advanced model MX

- M0-0 Data mapping
- MX-1 Explicit task mapping
- M0-2 Comms. inferred
- MX-3 Implicit reduction pattern



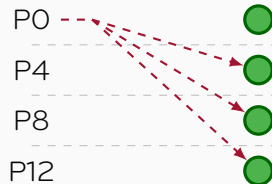
```

!data_map: A,B,C distributed on one layer
!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 insert_task(gemm,
                      $A_{i,l}$ :R,
                      $B_{l,j}$ :R,
                      $C_{i,j}$ :RANK_REDUX,
                     task_map[i,j,l])
    end do
  end do
end do
    
```



Advanced model MX

- M0-0 Data mapping
- MX-1 Explicit task mapping
- MX-2 **Coll. comms.** inferred ¹
- MX-3 Implicit reduction pattern



¹Denis, Jeannot, Swartvagher, and Thibault, 2020.

```
!data_map: A,B,C distributed on one layer
!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 insert_task(gemm,
                      $A_{i,l}:\text{R}$ ,
                      $B_{l,j}:\text{R}$ ,
                      $C_{i,j}:\text{RANK\_REDUX}$ ,
                     task_map[i,j,l])
    end do
  end do
end do
```

- `task_map` can be set to cover any SUMMA variant
- Agnostic of:
 - the data distribution
 - the scheduling
 - the processing units

Advanced model MX

- M0-0 Data mapping
- MX-1 **Explicit** task mapping
- MX-2 **Collective** communications **inferred** from task mapping
- MX-3 **Implicit** reduction pattern

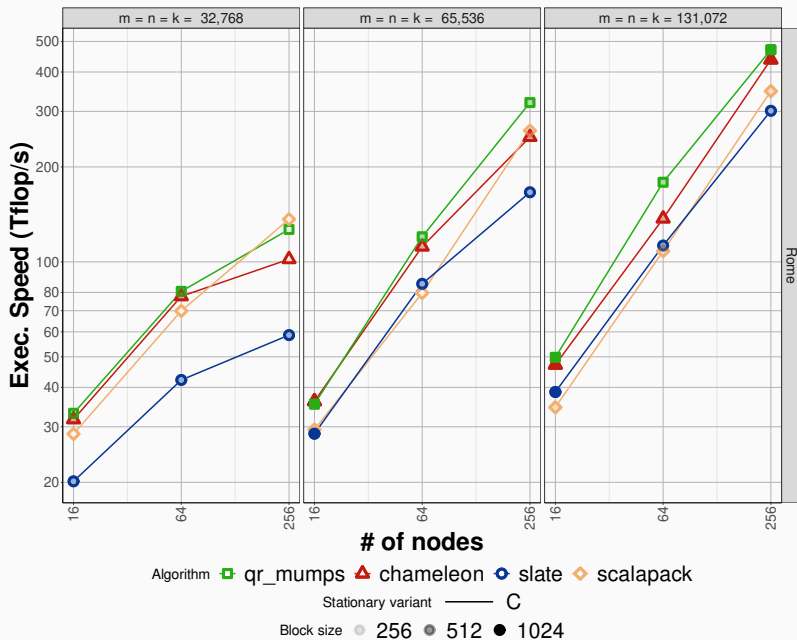
EXPERIMENTAL RESULTS

- **Software:** GNU 9.3.0, MKL 21.3.0
 - **qr_mumps ours** – newMadeleine, 1 MPI / node
 - StarPU with advanced model **MX** + seq. BLAS
 - 2.5D A,B,C-stationary version
 - **scalapack** – OpenMPI 4.0.5, 2 cores / MPI
 - Synchronous MPI + seq. BLAS
 - 2D A,B,C-stationary
 - **slate** – OpenMPI 4.0.5, 1 MPI / node
 - MPI + OpenMP with **task** + seq. BLAS
 - 2D A,C-stationary with lookahead
 - **chameleon** – newMadeleine, 1 MPI / node
 - StarPU with explicit comms. management + seq. BLAS
 - 2D C-stationary pipelined with lookahead
- **Hardware:** Très Grand Centre de Calcul (TGCC)
 - **Irène - Skylake** Infiniband EDR, 4x links
 - Dual-processor Intel Skylake 8168 @ 2.7 GHz
 - 24 cores per processor, 192 GB DDR4 memory
 - **Irène - Rome** Infiniband HDR100, 2x links
 - Dual-processor AMD Rome (Epyc) @ 2.6 GHz
 - 64 cores per processor, 256 GB DDR4 memory

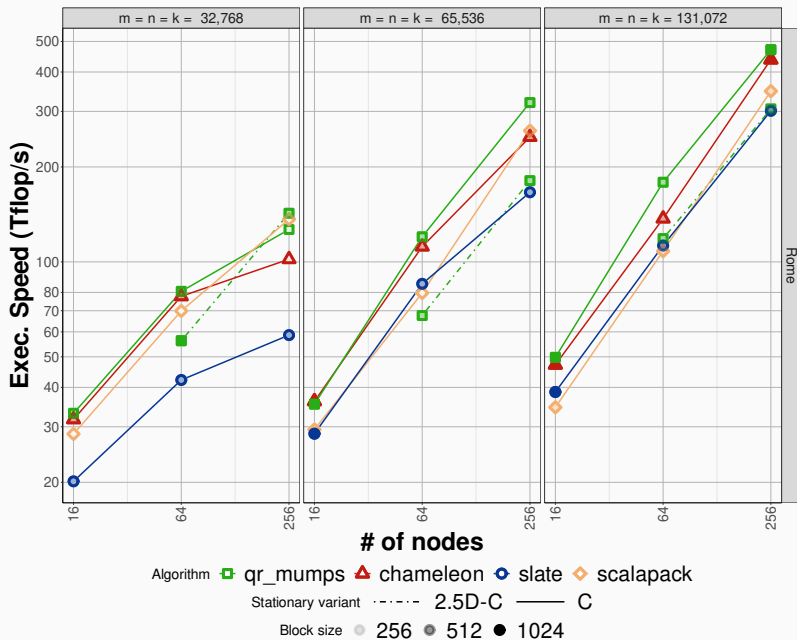
- **Software:** GNU 9.3.0, MKL 21.3.0
 - **qr_mumps ours** – newMadeleine, 1 MPI / node
 - StarPU with advanced model **MX** + seq. BLAS
 - 2.5D A,B,C-stationary version
 - **scalapack** – OpenMPI 4.0.5, 2 cores / MPI
 - Synchronous MPI + seq. BLAS
 - 2D A,B,C-stationary
 - **slate** – OpenMPI 4.0.5, 1 MPI / node
 - MPI + OpenMP with **task** + seq. BLAS
 - 2D A,C-stationary with lookahead
 - **chameleon** – newMadeleine, 1 MPI / node
 - StarPU with explicit comms. management + seq. BLAS
 - 2D C-stationary pipelined with lookahead
- **Hardware:** Très Grand Centre de Calcul (TGCC)
 - **Irène - Skylake** Infiniband EDR, 4x links
 - Dual-processor Intel Skylake 8168 @ 2.7 GHz
 - 24 cores per processor, 192 GB DDR4 memory
 - **Irène - Rome** Infiniband HDR100, 2x links
 - Dual-processor AMD Rome (Epyc) @ 2.6 GHz
 - 64 cores per processor, 256 GB DDR4 memory

- **Software:** GNU 9.3.0, MKL 21.3.0
 - **qr_mumps ours** – newMadeleine, 1 MPI / node
 - StarPU with advanced model **MX** + seq. BLAS
 - 2.5D A,B,C-stationary version
 - **scalapack** – OpenMPI 4.0.5, 2 cores / MPI
 - Synchronous MPI + seq. BLAS
 - 2D A,B,C-stationary
 - **slate** – OpenMPI 4.0.5, 1 MPI / node
 - MPI + OpenMP with **task** + seq. BLAS
 - 2D A,C-stationary with lookahead
 - **chameleon** – newMadeleine, 1 MPI / node
 - StarPU with explicit comms. management + seq. BLAS
 - 2D C-stationary pipelined with lookahead
- **Hardware:** Très Grand Centre de Calcul (TGCC)
 - **Irène - Skylake** Infiniband EDR, 4x links
 - Dual-processor Intel Skylake 8168 @ 2.7 GHz
 - 24 cores per processor, 192 GB DDR4 memory
 - **Irène - Rome** Infiniband HDR100, 2x links
 - Dual-processor AMD Rome (Epyc) @ 2.6 GHz
 - 64 cores per processor, 256 GB DDR4 memory

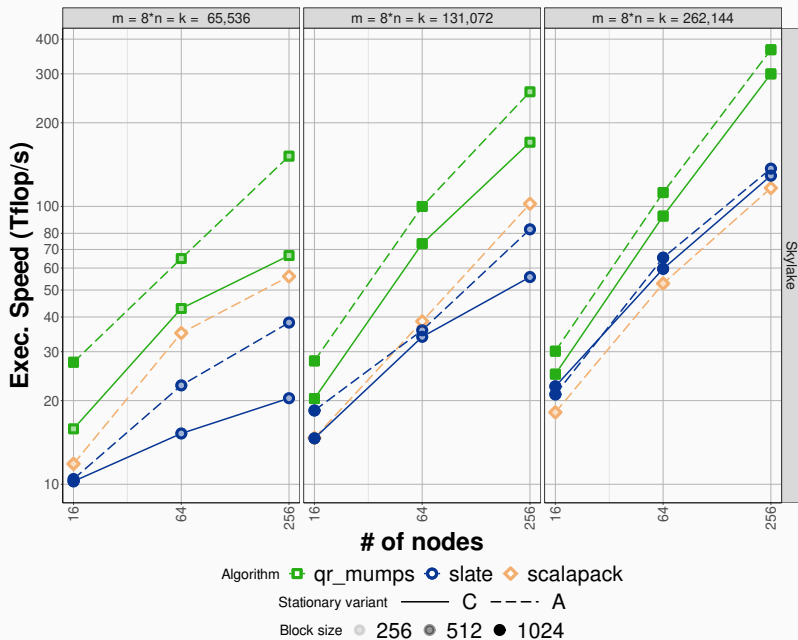
ROME (UP TO 32,768 CORES) – SQUARE MATRICES, DGEMM



ROME (UP TO 32,768 CORES) – SQUARE MATRICES, DGEMM



SKYLAKE (UP TO 12,288 CORES) – LARGE A, DGEMM



CONCLUSIONS

Contributions

- Show **the STF programming model** can be suitably **extended** to **express** a complex distributed-memory algorithm
 - **Portable**, easily **maintainable**
 - Independent of matrices' distribution, workflow's scheduling
 - **Competitive** with reference libraries

Technical report: <https://hal.inria.fr/hal-03588491/>

On-going work

- Studying SYMM in the context of a scientific application
 - In collaboration with O. Beaumont, O. Coulaud, L. Eyraus-Dubois, A. Franc, R. Peressoni, F. Pruvost, M. Vérité
- Studying the scheduling of communication patterns in GEMM
 - In collaboration with O. Beaumont, A. Denis, L. Eyraus-Dubois, E. Jeannot, M. Vérité, P. Swartvagher

We want to address more distributed-memory scalable algorithms

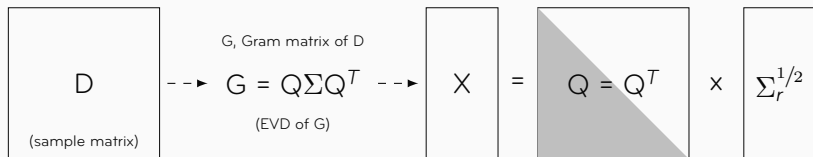
- Dense **2.5D** Cholesky, LU **factorizations**
- Sparse factorizations

However, we need further **extensions** to the current programming model

- **Replication of computation** as required by **state-of-the-art 2.5D** algorithms
- **All-reduce pattern** to ensure **stability** through tournament pivoting
- **Memory-awareness** to ensure an algorithm will **not exceed** a user-defined **memory constraint**

ON-GOING WORK

Our versatile approach is of interest in the **Fast Methods for Randomized numerical algebra (FMR)** data mining library.

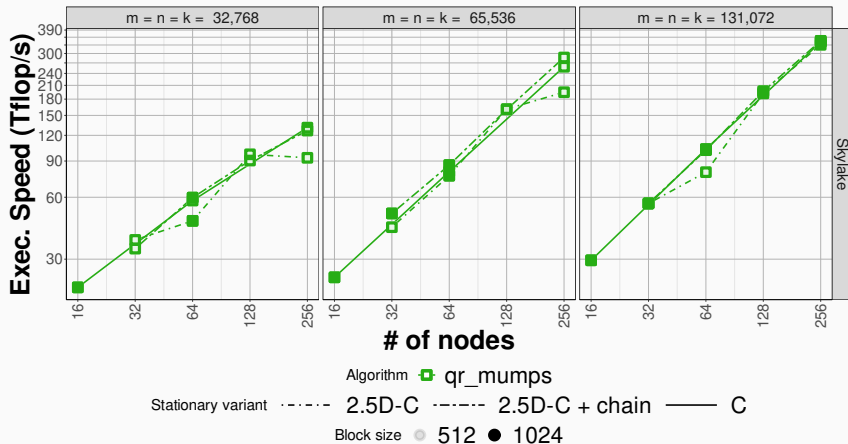


- **A-stationary algorithm:** involved matrices are disproportionate
- **distribution-agnosticism:** to deploy *symmetric distributions*

⇒ In collaboration with:

Olivier Beaumont, Olivier Coulaud, Lionel Eyraud-Dubois, Alain Franc, Romain Peressoni, Florient Pruvost, Mathieu V erit e

Experiments questioning the scheduling of communications



⇒ In collaboration with:

Olivier Beaumont, Alexandre Denis, Lionel Eyraus-Dubois,
Emmanuel Jeannot, Mathieu Vérité, Philippe Swartvagher

2.5D factorization algorithms have **more complex** critical path.

They require more advanced key functionalities:

1. Accumulation and reduction of **updates over layers**
2. **Replication** of factorization tasks
3. **All-reduce pattern** to implement tournament pivoting

The programming model currently lacks support for 2. and 3.

2.5D factorization algorithms have **more complex** critical path.

They require more advanced key functionalities:

1. Accumulation and reduction of **updates over layers**
2. **Replication** of factorization tasks
3. **All-reduce pattern** to implement tournament pivoting

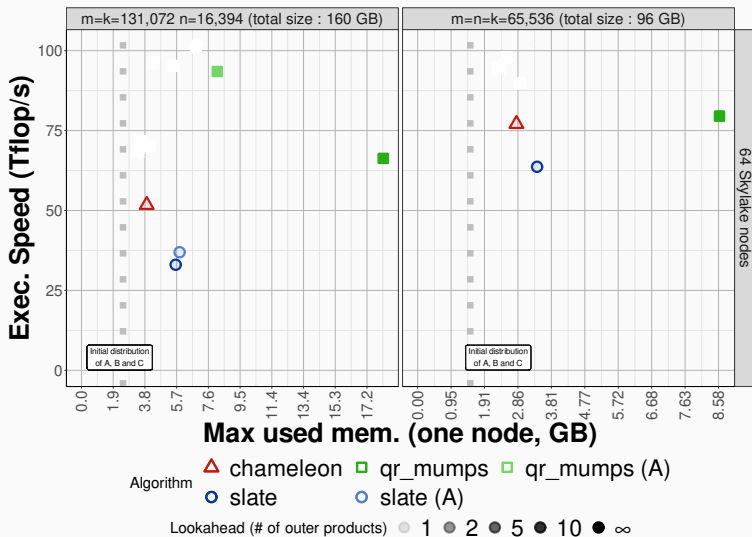
The programming model currently lacks support for 2. and 3.

BACKUP SLIDES

MEMORY CONSUMPTION CONCERN - EXPERIMENT

Submitting the **full DAG** allows **greedy** memory consumption

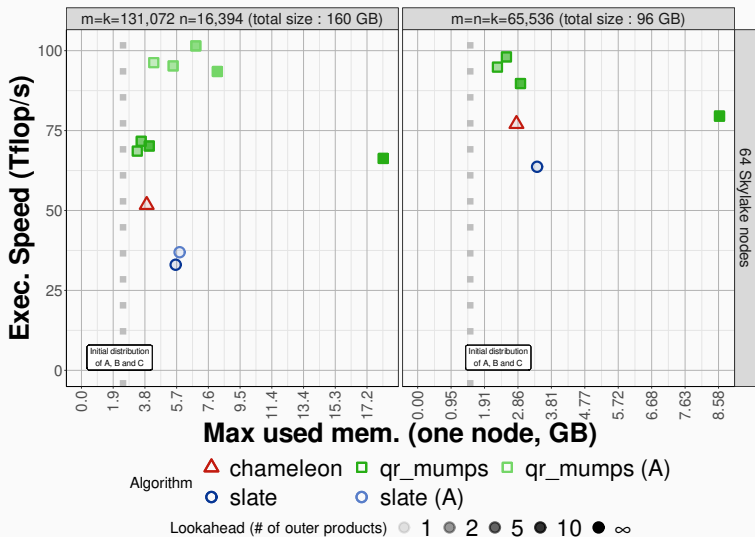
Task window mechanisms can help limit memory consumption



MEMORY CONSUMPTION CONCERN - EXPERIMENT

Submitting the **full DAG** allows **greedy** memory consumption

Task window mechanisms can help limit memory consumption



Some common operations are done before switching on the `stat` parameter:

- Binding C to the `init` and `sum` codelets
- (submitting) the scaling of C by β

2.5D-A-Stat

```
do jh=1,n/h; do layer=1,h
  j = jh * layer
  do i=1,m
    do l=1,k
      Ail => transa ? Al,i : Ai,l
      Blj => transb ? Bj,l : Bl,j
      Cij => Ci,j
      work_node = Ail%owner
                  + r*c*layer
      call insert_task( gemm,
        work_node:ON_NODE,
        Ail:R,Blj:R,
        Cij:RANK_REDUX )
      end do
      call mpi_redux_tree(Ci,j)
    end do
  end do; end do
```

2.5D-B-Stat

```
do ih=1,m/h; do layer=1,h
  i = ih * layer
  do j=1,n
    do l=1,k
      Ail => transa ? Al,i : Ai,l
      Blj => transb ? Bj,l : Bl,j
      Cij => Ci,j
      work_node = Blj%owner
                  + r*c*layer
      call insert_task( gemm,
        work_node:ON_NODE,
        Ail:R,Blj:R,
        Cij:RANK_REDUX )
      end do
      call mpi_redux_tree(Ci,j)
    end do
  end do; end do
```

2.5D-C-Stat

```
do lh=1,k/h; do layer=1,h
  l = lh * layer
  do i=1,m
    do j=1,n
      Ail => transa ? Al,i : Ai,l
      Blj => transb ? Bj,l : Bl,j
      Cij => Ci,j
      work_node = Cij%owner
                  + r*c*layer
      call insert_task( gemm,
        work_node:ON_NODE,
        Ail:R,Blj:R,
        Cij:RANK_REDUX )
      end do
    end do; end do
  do i=1,m; do j=1,n
    call mpi_redux_tree(Ci,j)
  end do; end do
```

STARPU – ADDITIONAL FEATURES ?

Some additional features to StarPU have been considered.

- Automatically insert reduction tasks

- Further simplify factorization implementations

```
call starpu_mpi_redux(Ai,j)  
call insert_task(potrf, Ai,j:RW)
```



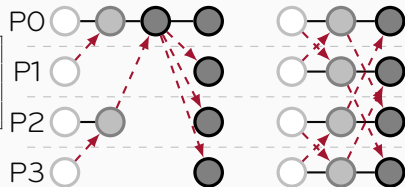
```
! implicit reduction  
call insert_task(potrf, Ai,j:RW)
```

- Allow tasks replication across nodes

- Required by state-of-the-art algorithms
- Each resulting tile can be sent to different subsets of nodes

```
call insert_task(potrf, Ai,j:RW,  
                ON_NODES:task_map[l,l,:])  
call insert_task(trsm, Ai,j:R, Ai,j:RW,  
                FROM_NODE:task_map[l,l,i%h])
```

- Allow all-reduce pattern
 - Shorten the critical path as replicated data is already available



- 2D-C-Stat owner of C computes
 - $\text{task_map}[i, j, l] = j \% c + c * (i \% r)$
 - $\text{access_mode} = \text{RW}$



- 2D-A-Stat owners of A compute
 - $\text{task_map}[i, j, l] = l \% c + c * (i \% r)$
 - $\text{access_mode} = \text{RANK_REDUX}$



- 2.5D-C-Stat several nodes compute
 - $\text{task_map}[i, j, l] = j \% c + c * (i \% r) + rc * \frac{l}{k/h}$
 - $\text{access_mode} = \text{RANK_REDUX}$



The advanced model allows to **cover all variants** of SUMMA algorithms with three nested loops.

In a **distribution-agnostic** fashion.

- 2D-C-Stat owner of C computes
 - $\text{task_map}[i, j, l] = j \% c + c * (i \% r)$
 - $\text{access_mode} = \text{RW. and. COMMUTE}$



- 2D-A-Stat owners of A compute
 - $\text{task_map}[i, j, l] = l \% c + c * (i \% r)$
 - $\text{access_mode} = \text{RANK_REDUX}$



- 2.5D-C-Stat several nodes compute
 - $\text{task_map}[i, j, l] = j \% c + c * (i \% r) + rc * \frac{l}{k/h}$
 - $\text{access_mode} = \text{RANK_REDUX}$



The advanced model allows to **cover all variants** of SUMMA algorithms with three nested loops.

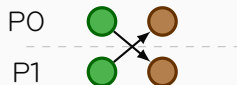
In a **distribution-agnostic** fashion.

IDEAL EXPRESSION - TASK MAPPINGS

- 2D-C-Stat owner of C computes
 - $\text{task_map}[i, j, l] = j \% c + c * (i \% r)$
 - $\text{access_mode} = \text{RW. and. COMMUTE}$



- 2D-A-Stat owners of A compute
 - $\text{task_map}[i, j, l] = l \% c + c * (i \% r)$
 - $\text{access_mode} = \text{RANK_REDUX}$



- 2.5D-C-Stat several nodes compute
 - $\text{task_map}[i, j, l] = j \% c + c * (i \% r) + rc * \frac{l}{k/h}$
 - $\text{access_mode} = \text{RANK_REDUX}$



The advanced model allows to **cover all variants** of SUMMA algorithms with three nested loops.

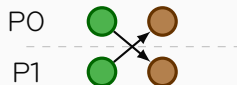
In a **distribution-agnostic** fashion.

IDEAL EXPRESSION - TASK MAPPINGS

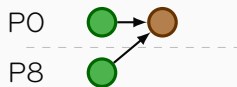
- 2D-C-Stat owner of C computes
 - $\text{task_map}[i, j, l] = j \% c + c * (i \% r)$
 - $\text{access_mode} = \text{RW. and. COMMUTE}$



- 2D-A-Stat owners of A compute
 - $\text{task_map}[i, j, l] = l \% c + c * (i \% r)$
 - $\text{access_mode} = \text{RANK_REDUX}$



- 2.5D-C-Stat several nodes compute
 - $\text{task_map}[i, j, l] = j \% c + c * (i \% r) + rc * \frac{l}{k/h}$
 - $\text{access_mode} = \text{RANK_REDUX}$



The advanced model allows to **cover all variants** of SUMMA algorithms with three nested loops.

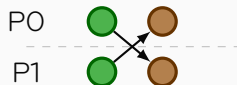
In a **distribution-agnostic** fashion.

IDEAL EXPRESSION - TASK MAPPINGS

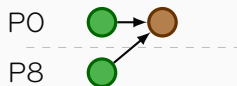
- 2D-C-Stat owner of C computes
 - $\text{task_map}[i, j, l] = C_{i,j}.\text{owner_rank}$
 - $\text{access_mode} = \text{RW. and. COMMUTE}$



- 2D-A-Stat owners of A compute
 - $\text{task_map}[i, j, l] = A_{i,j}.\text{owner_rank}$
 - $\text{access_mode} = \text{RANK_REDUX}$

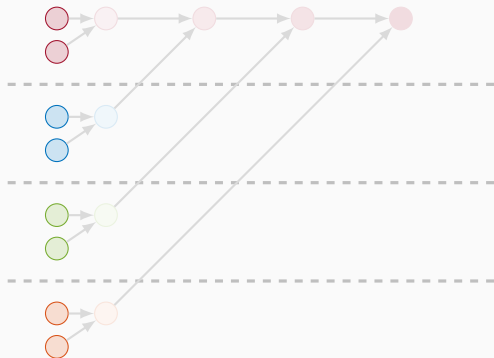


- 2.5D-C-Stat several nodes compute
 - $\text{task_map}[i, j, l] = C_{i,j}.\text{owner_rank} + rc * \frac{l}{k/h}$
 - $\text{access_mode} = \text{RANK_REDUX}$



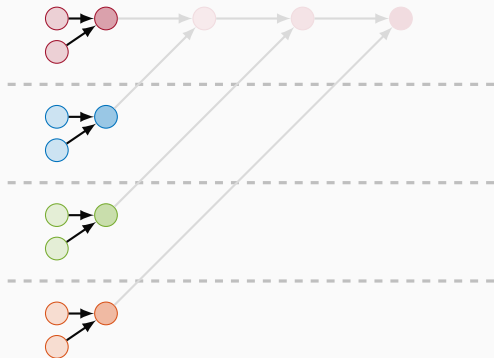
The advanced model allows to **cover all variants** of SUMMA algorithms with three nested loops. In a **distribution-agnostic** fashion.

- **STARPU_REDUX**
 - One contribution per **core**
 - Intra-node reduction through binary tree
 - Remain available in StarPU (master)
- `starpu_mpi_redux_data`
 - **All** nodes in `MPI_COMM_World` contribute through a **flat** tree
 - **Blocking** method

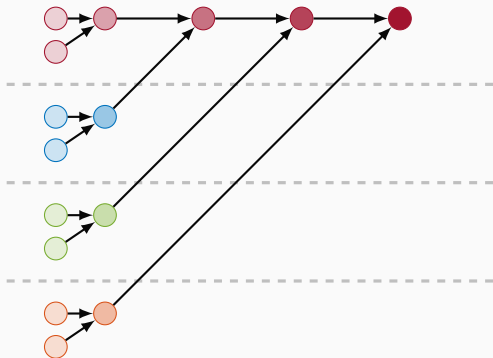


FEATURES OF STARPU – 1.3

- **STARPU_REDUX**
 - One contribution per **core**
 - Intra-node reduction through binary tree
 - Remain available in StarPU (master)
- `starpu_mpi_redux_data`
 - **All** nodes in `MPI_COMM_World` contribute through a **flat** tree
 - **Blocking** method

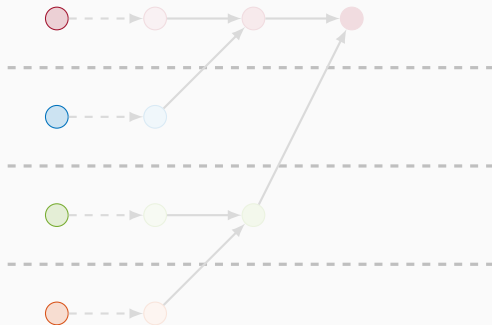


- **STARPU_REDUX**
 - One contribution per **core**
 - Intra-node reduction through binary tree
 - Remain available in StarPU (master)
- **starpu_mpi_redux_data**
 - **All** nodes in MPI_COMM_World contribute through a **flat** tree
 - **Blocking** method



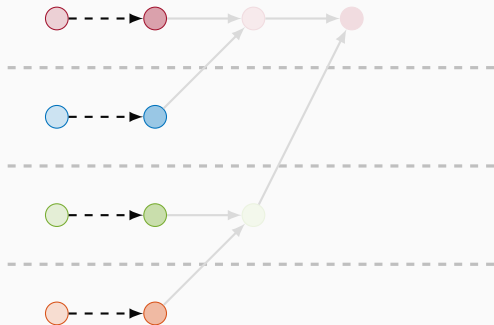
FEATURES OF STARPU – OUR CONTRIBUTION

- STARPU_MPI_REDUX (**Feature!** issue 3)
 - One contribution per **contributing node**
 - Interpreted as STARPU_RW|STARPU_COMMUTE
- starpu_mpi_redux_data_tree (**Enhancement!** merge req. 21)
 - All **contributing** nodes join in the reduction through a **n-ary** tree
 - **Non-blocking** method



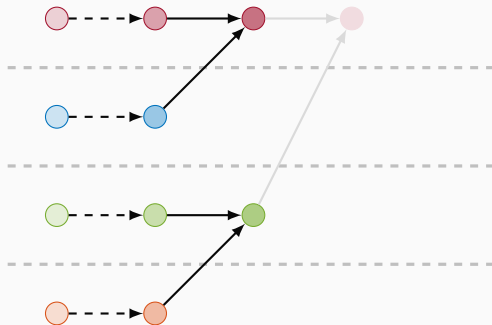
FEATURES OF STARPU – OUR CONTRIBUTION

- STARPU_MPI_REDUX (**Feature!** issue 3)
 - One contribution per **contributing node**
 - Interpreted as STARPU_RW|STARPU_COMMUTE
- starpu_mpi_redux_data_tree (**Enhancement!** merge req. 21)
 - All **contributing** nodes join in the reduction through a **n-ary** tree
 - **Non-blocking** method



FEATURES OF STARPU – OUR CONTRIBUTION

- STARPU_MPI_REDUX (**Feature!** issue 3)
 - One contribution per **contributing node**
 - Interpreted as STARPU_RW|STARPU_COMMUTE
- starpu_mpi_redux_data_tree (**Enhancement!** merge req. 21)
 - All **contributing** nodes join in the reduction through a **n-ary** tree
 - **Non-blocking** method



FEATURES OF STARPU – OUR CONTRIBUTION

- STARPU_MPI_REDUX (**Feature!** issue 3)
 - One contribution per **contributing node**
 - Interpreted as STARPU_RW|STARPU_COMMUTE
- starpu_mpi_redux_data_tree (**Enhancement!** merge req. 21)
 - All **contributing** nodes join in the reduction through a **n-ary** tree
 - **Non-blocking** method

