

TASK-BASED PROGRAMMING MODELS FOR SCALABLE ALGORITHMS

Application to matrix multiplication

E. Agullo, A. Buttari, A. Guermouche, J. Herrmann, A. Jego

February 10, 2022

Université de Toulouse - IRIT

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 a 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 a 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 a efficient, portable and maintainable way?

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

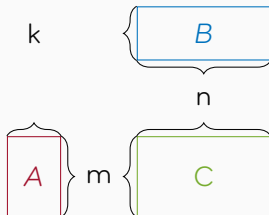
- **Algorithms** have to handle communications to achieve high scalability
 - Communication/computation overlap
 - Communication-avoiding algorithms
- **Software** should deliver a reliable abstraction to build mathematical software
 - parallel programming models
 - runtime systems

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

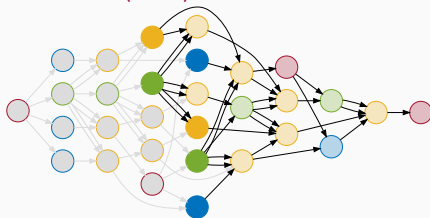
We focus on two aspects of this large question

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

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



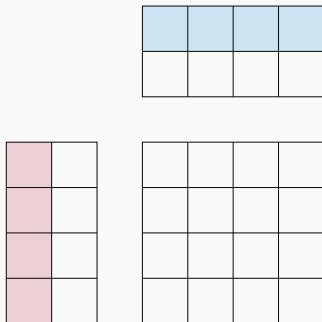
- **Software:** **task-based** parallel programming paradigm through the **sequential task flow (STF)** model



BACKGROUND

2D C-stationary

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



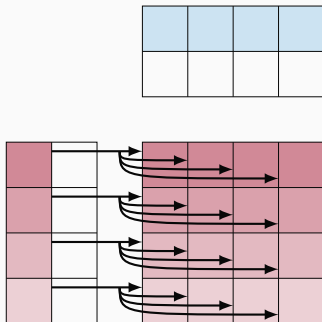
```

! 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

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



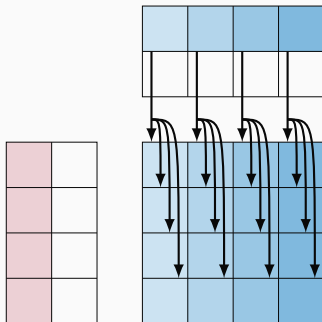
```

! 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

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



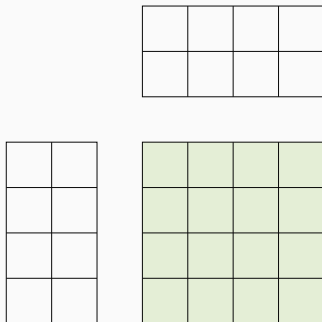
```

! 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

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



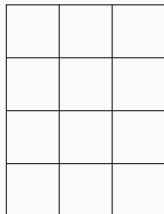
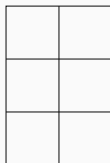
```

! 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

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



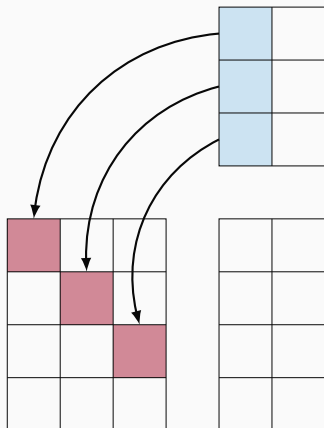
```

! 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

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



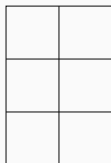
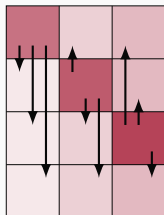
```

! I am r × 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

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

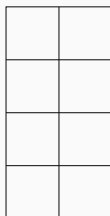
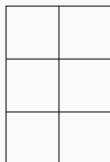
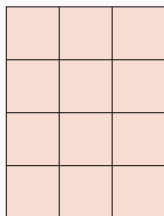


```

! 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

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



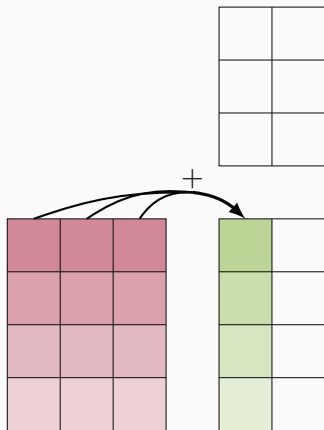
```

! I am r × 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

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



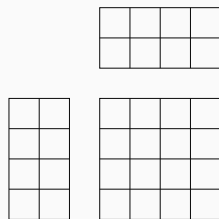
```

! 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

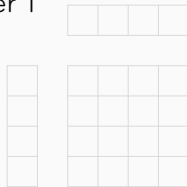
```



Layer 0



Layer 1



2.5D C-stationary

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

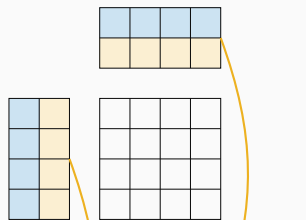
```

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

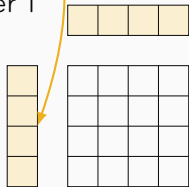
forall layer=1,h:
  call SUMMA_2DC(A[:, (layer : layer + 1) * k/h],
                B[(layer : layer + 1) * k/h, :], Ch)

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

Layer 0



Layer 1



2.5D C-stationary

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

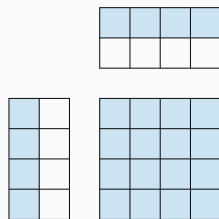
```

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

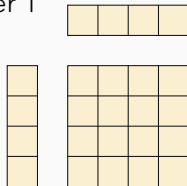
forall layer=1,h:
  call SUMMA_2DC(A[:, (layer : layer + 1) * k/h],
                B[(layer : layer + 1) * k/h, :], Ch)

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

Layer 0



Layer 1



2.5D C-stationary

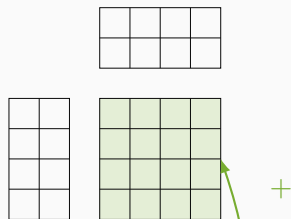
- **Data mapping**
 - 2D block-cyclic on 0^{th} layer
- **Task mapping**
 - Not trivially related to data mapping
- **Communications**
 - row-wise scatter+**broadcast** of A
 - column-wise scatter+**broadcast** of B
 - aisle-wise reduce of C

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

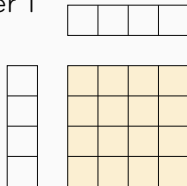
```
forall layer=1,h:
  call SUMMA_2DC(A[:, (layer : layer + 1) * k/h],
                B[(layer : layer + 1) * k/h, :], Ch) ←
```

```
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

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

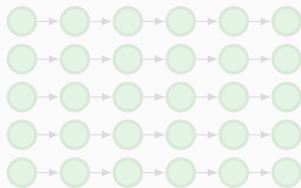
THE SEQUENTIAL TASK FLOW (STF) MODEL

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
```



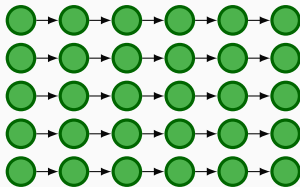
DAG of GEMM

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

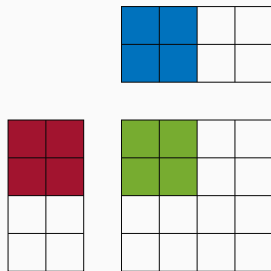
- Can we use the **STF model** to express a **scalable distributed-memory GEMM** algorithm ?
- How **efficient** is the resulting expression compared to reference libraries ?

DISTRIBUTED-MEMORY STF GEMM ALGORITHM

```

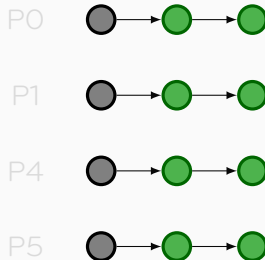
!
do i=1,m
  do j=1,n
    do l=1,k
      call insert_task(gemm,
                     Ai,l:R,
                     Bl,j:R,
                     Ci,j:RW)
    end do
  end do
end do

```



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, Favre, 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,j:R,
                      Bi,j:R,
                      Ci,j:RW)
    end do
  end do
end do
```

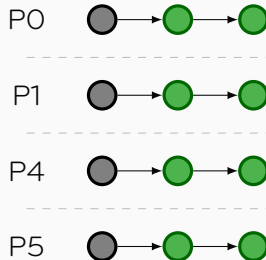
P0	P1	P2	P3
P4	P5	P6	P7

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, Favre, 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,j:R,
        Bi,j:R,
        Ci,j:RW)
    end do
  end do
end do
```

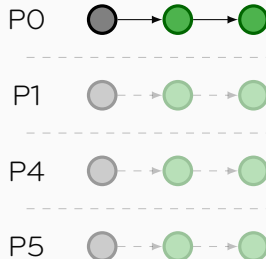
P0	P1	P2	P3
P4	P5	P6	P7

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

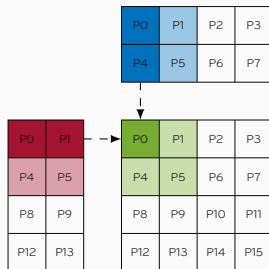
- M0-0 Data mapping
- M0-1 Task mapping **inferred** from DM (RW)
- M0-2 Point-to-point comms. **inferred** from TM



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

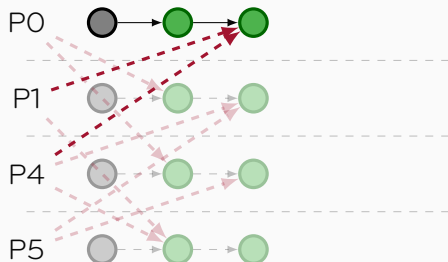
BASILINE STF MODEL

```
! 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,j:R,
                     Bi,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, Favege, Furmento, Pruvost, Sergent, and Thibault, 2017.

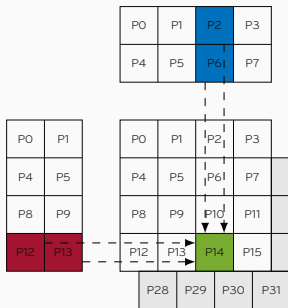
The baseline STF model **M0** does not allow the expression of the features of SUMMA algorithms

- **M0-1** : Task mapping is over-constrained by data mapping
- **M0-2** : Point-to-point communications patterns are inefficient
- **M0-3** : Reduction patterns are hard to make

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

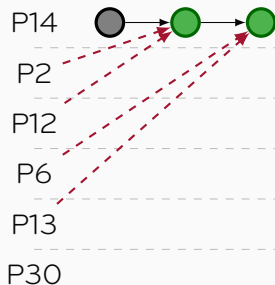
KEY 1 : TASK MAPPING

```
!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
```



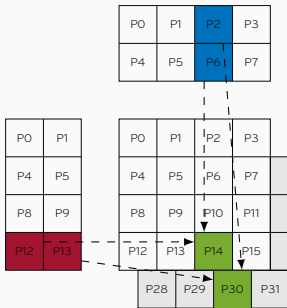
Advanced model MX

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



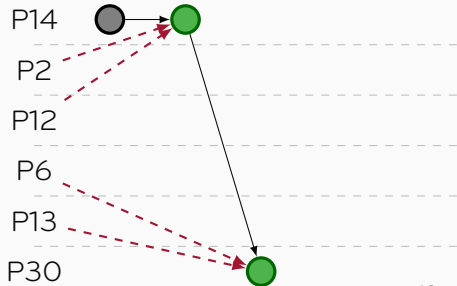
KEY 1 : TASK MAPPING

```
!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

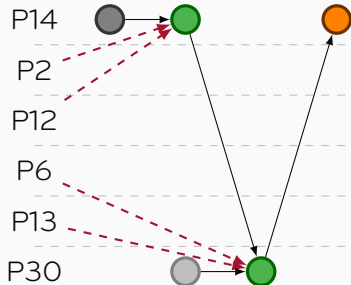
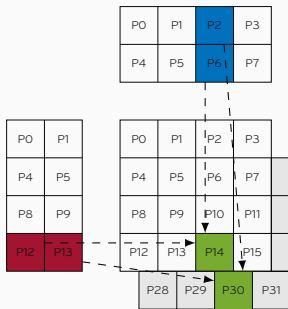


KEY 2 : REDUCTION PATTERN

```
!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

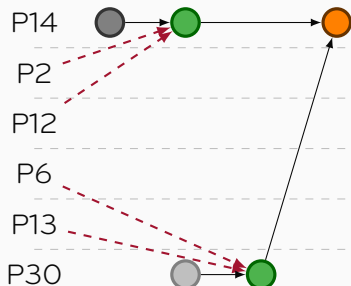
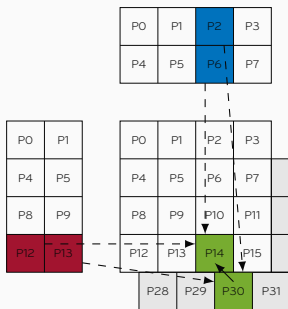


KEY 2 : REDUCTION PATTERN

```
!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,
                      $C_{i,j}$ :MPI_REDUX,
                     task_map[i,j,l])
    end do
  end do
end do
call reduce_to_initial_layer(C)
```

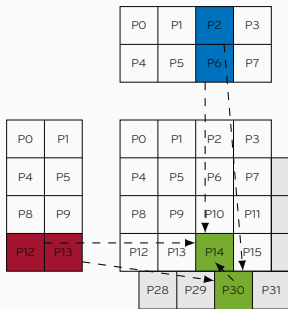
Advanced model MX

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



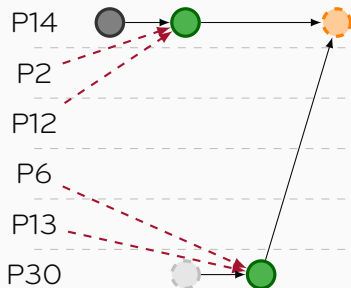
KEY 2 : REDUCTION PATTERN

```
!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,
                      $C_{i,j}$ :MPI_REDUX,
                     task_map[i,j,l])
    end do
  end do
end do
```



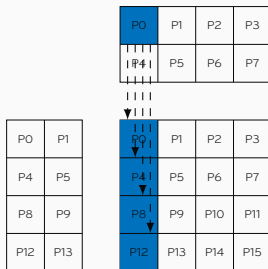
Advanced model MX

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



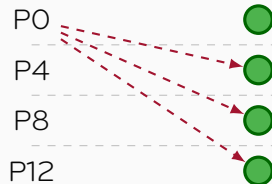
KEY 3 : COMMUNICATIONS

```
!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,
                       $C_{i,j}$ :MPI_REDUX,
                      task_map[i,j,l])
    end do
  end do
end do
```



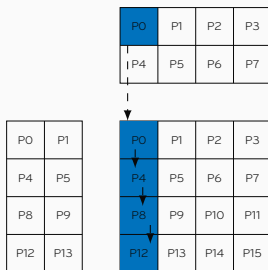
Advanced model MX

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



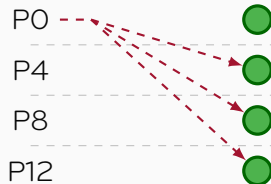
KEY 3 : COMMUNICATIONS

```
!data_map: A,B,C distributed on one layer
!task_map: executing node of  $A_{i,l}B_{l,j}$ 
!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,
         $A_{i,l}:\text{R}$ ,
         $B_{l,j}:\text{R}$ ,
         $C_{i,j}:\text{MPI\_REDUX}$ ,
        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^a
- MX-3 Implicit reduction pattern



^aDenis, 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,
                      $A_{i,l}$ :R,
                      $B_{l,j}$ :R,
                      $C_{i,j}$ :MPI_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 **Collective** communications **inferred** from task mapping
- MX-3 **Implicit** reduction pattern

- 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 owner of A computes
 - $\text{task_map}[i, j, l] = l \% c + c * (i \% r)$
 - $\text{access_mode} = \text{MPI_REDUX}$



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



The advanced model allows to cover all SUMMA variants within three nested loops.

- 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 owner of A computes
 - $\text{task_map}[i, j, l] = l \% c + c * (i \% r)$
 - $\text{access_mode} = \text{MPI_REDUX}$



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

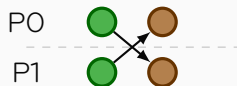


The advanced model allows to cover all SUMMA variants within three nested loops.

- 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 owner of A computes
 - $\text{task_map}[i, j, l] = l \% c + c * (i \% r)$
 - $\text{access_mode} = \text{MPI_REDUX}$



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



The advanced model allows to cover all SUMMA variants within three nested loops.

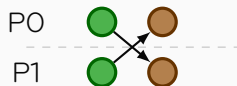
- 2D-C-Stat owner of C computes

- $\text{task_map}[i, j, l] = j \% c + c * (i \% r)$
- $\text{access_mode} = \text{RW.and.COMMUTE}$



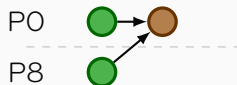
- 2D-A-Stat owner of A computes

- $\text{task_map}[i, j, l] = l \% c + c * (i \% r)$
- $\text{access_mode} = \text{MPI_REDUX}$



- 2.5D-C-Stat several nodes computes

- $\text{task_map}[i, j, l] = j \% c + c * (i \% r) + rc * \frac{l}{k/h}$
- $\text{access_mode} = \text{MPI_REDUX}$



The advanced model allows to cover all SUMMA variants within three nested loops.

EXPERIMENTAL RESULTS

- **StarPU** runtime
 - Support STF programming model
 - Developed by Storm Team at INRIA - Bordeaux
 - **NewMadeleine** communications backend (TaDaaM Team)

MX-0 Data mapping is provided.

MX-1 Explicit task mapping is provided.

MX-2 Dynamic collective communications are available in an experimental branch if compiled with newMadeleine.

MX-3 Reduction pattern exists in shared memory, **we have adapted it to distributed memory.**

- **StarPU** runtime
 - Support STF programming model
 - Developed by Storm Team at INRIA - Bordeaux
 - **NewMadeleine** communications backend (TaDaaM Team)

M0-0 Data mapping is provided.

MX-1 Explicit task mapping is provided.

MX-2 Dynamic collective communications are available in an experimental branch if compiled with newMadeleine.

MX-3 Reduction pattern exists in shared memory, **we have adapted it to distributed memory.**

- **StarPU** runtime
 - Support STF programming model
 - Developed by Storm Team at INRIA - Bordeaux
 - **NewMadeleine** communications backend (TaDaaM Team)

M0-0 Data mapping is provided.

MX-1 Explicit task mapping is provided.

MX-2 Dynamic collective communications are available in an experimental branch if compiled with newMadeleine.

MX-3 Reduction pattern exists in shared memory, **we have adapted it to distributed memory.**

- **StarPU** runtime
 - Support STF programming model
 - Developed by Storm Team at INRIA - Bordeaux
 - **NewMadeleine** communications backend (TaDaaM Team)

M0-0 Data mapping is provided.

MX-1 Explicit task mapping is provided.

MX-2 Dynamic collective communications are available in an experimental branch if compiled with newMadeleine.

MX-3 Reduction pattern exists in shared memory, **we have adapted it to distributed memory.**

- **StarPU** runtime
 - Support STF programming model
 - Developed by Storm Team at INRIA - Bordeaux
 - **NewMadeleine** communications backend (TaDaaM Team)

M0-0 Data mapping is provided.

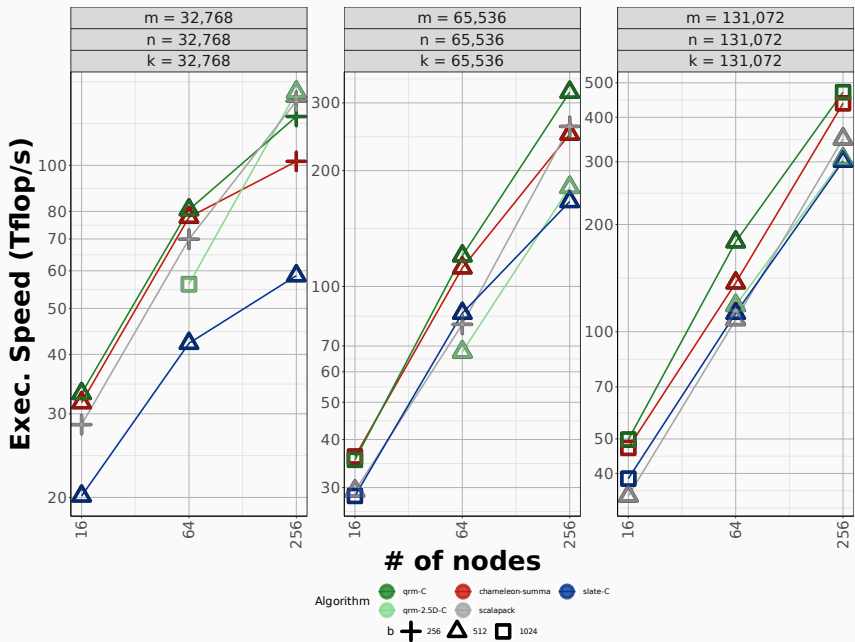
MX-1 Explicit task mapping is provided.

MX-2 Dynamic collective communications are available in an experimental branch if compiled with newMadeleine.

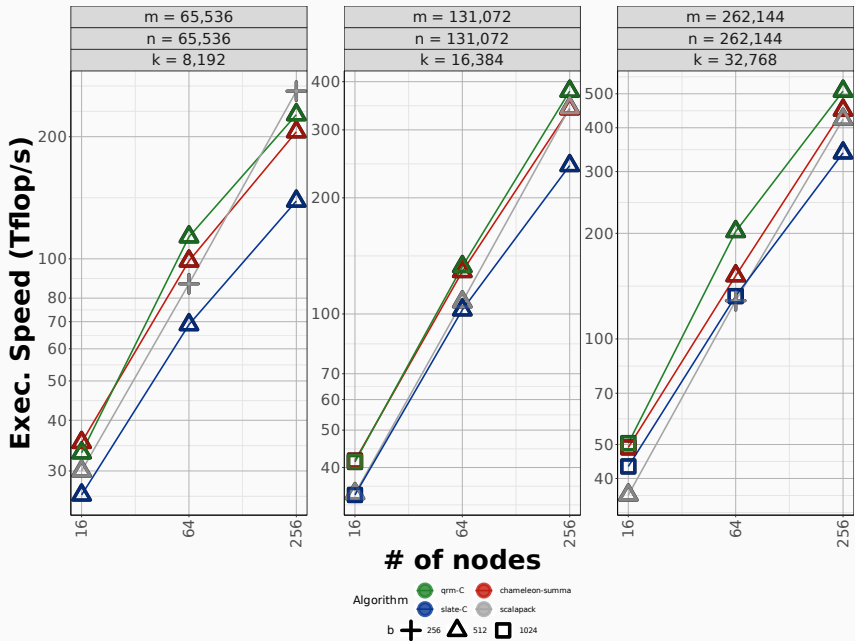
MX-3 Reduction pattern exists in shared memory, **we have adapted it to distributed memory.**

- **Hardware:** Très Grand Centre de Calcul (TGCC)
 - **Irène - Skylake** Infiniband EDR
 - Dual-processor Intel Skylake 8168 @ 2.7 GHz
 - 24 cores per processor, 192 GB DDR4 memory
 - **Irène - Rome** Infiniband HDR100
 - 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
 - **scalapack** – OpenMPI 4.0.5, 2 cores / MPI
 - Synchronous approach
 - 2D A,B,C-stationary hand-made routines
 - **slate** – OpenMPI 4.0.5, 1 MPI / node
 - Dynamic approach
 - 2D A,C-stationary based on OpenMP **task**
 - **chameleon** – nMad, 1 MPI / node
 - StarPU with explicit communications management
 - 2D C-stationary pipelined version
 - **qr_mumps ours** – nMad, 1 MPI / node
 - Based on StarPU with advanced STF model
 - 2.5D A,B,C-stationary version

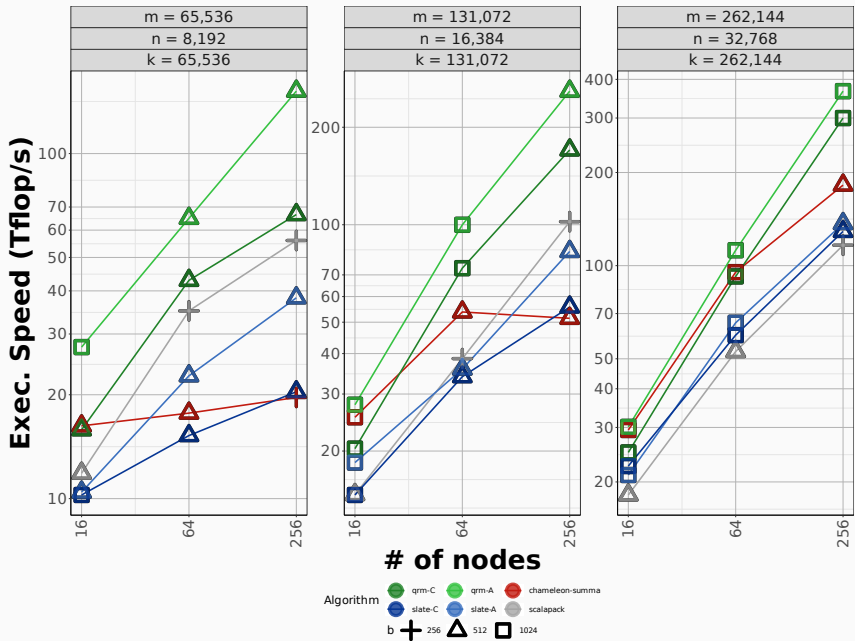
ROME NODES – SQUARE MATRICES, DGEMM



ROME NODES – LARGE C, DGEMM



SKYLAKE NODES – LARGE A, DGEMM



CONCLUSIONS

Contributions

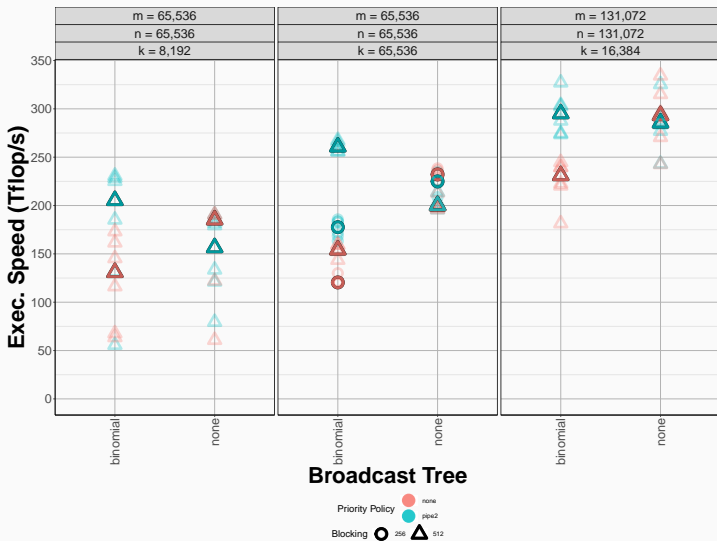
- **Show the STF programming model** allows for the expression of a complex distributed-memory algorithm such as SUMMA
- **Enhance StarPU programming interface** to cover this advanced model
- **Implement** our method and **compare it** with reference libraries, with interesting results

Future works

- Analyze and reliably **control memory consumption**
- **Analyze the impact of priorities** with communication patterns
- Use our approach to implement other numerical routines
 - **Dense LU, Cholesky factorization**
 - **Sparse factorization**

BACKUP SLIDES

DYNAMIC COLLECTIVE COMMUNICATIONS



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

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

```
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:MPI_REDUX
                       )
    end do
    call mpi_redux_tree(
      Ci,j)
  end do
end do; end do
```

2.5D-A-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:MPI_REDUX
                       )
    end do
    call mpi_redux_tree(
      Ci,j)
  end do
end do; end do
```

2.5D-B-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:MPI_REDUX
                       )
    end do
  end do
end do; end do
do i=1,m; do j=1,n
  call mpi_redux_tree(Ci,j)
end do; end do
```

2.5D-C-Stat

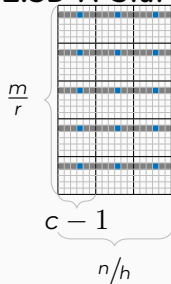
MEMORY CONSUMPTION CONCERN

Discarding owned blocks, Memory is consumed two ways

- by communication buffers
 - invalidated upon being used in all local contributions
- by reduction copies
 - freed after reduction pattern execution

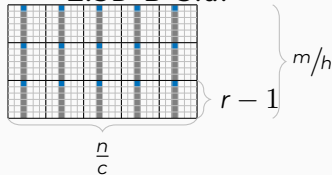
Temporary copies of C are limited

2.5D-A-Stat



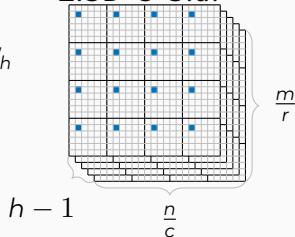
$$\frac{m}{r} n/h (c-1)$$

2.5D-B-Stat



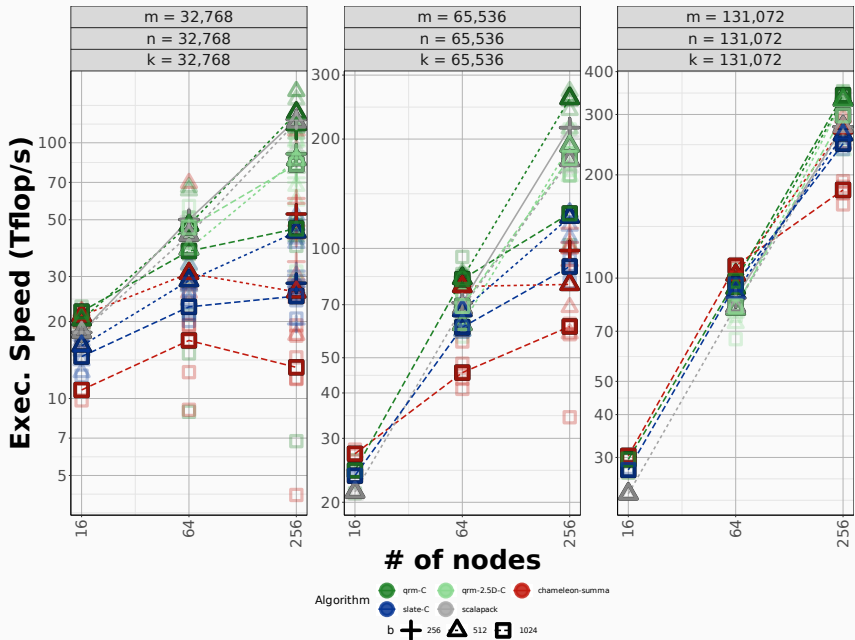
$$m/h \frac{n}{c} (r-1)$$

2.5D-C-Stat

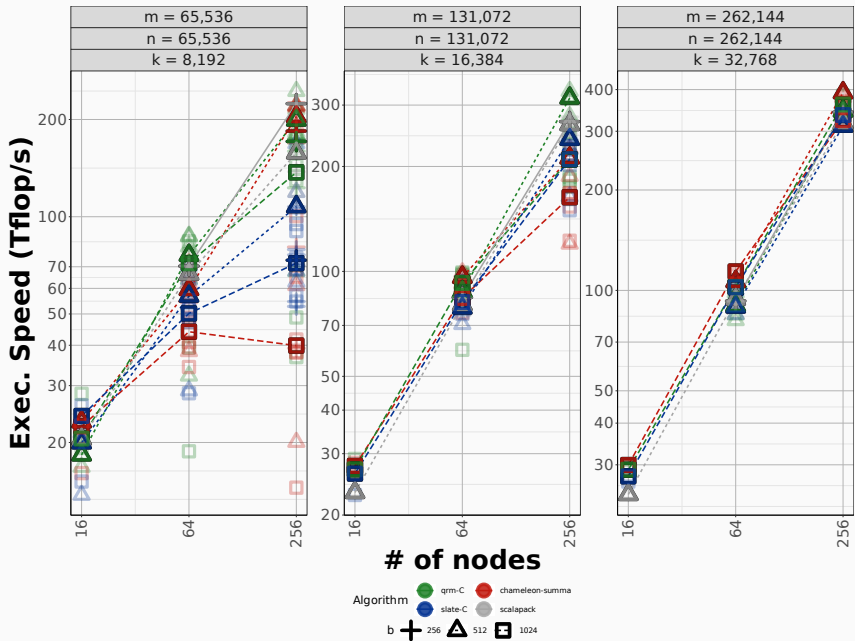


$$\frac{m}{r} \frac{n}{c} (h-1)$$

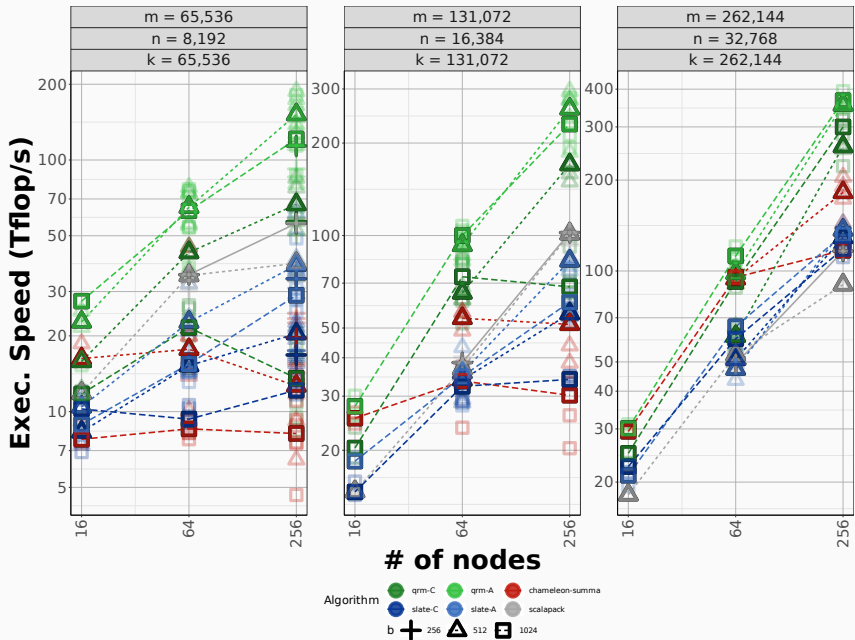
SKYLAKE NODES – SQUARE MATRICES, DGEMM (FULL)



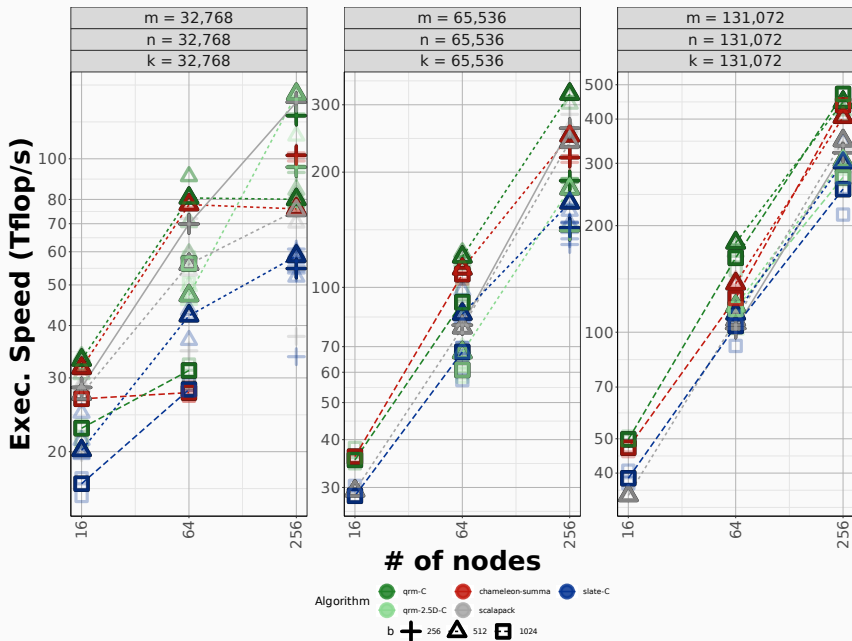
SKYLAKE NODES – LARGE C, DGEMM (FULL)



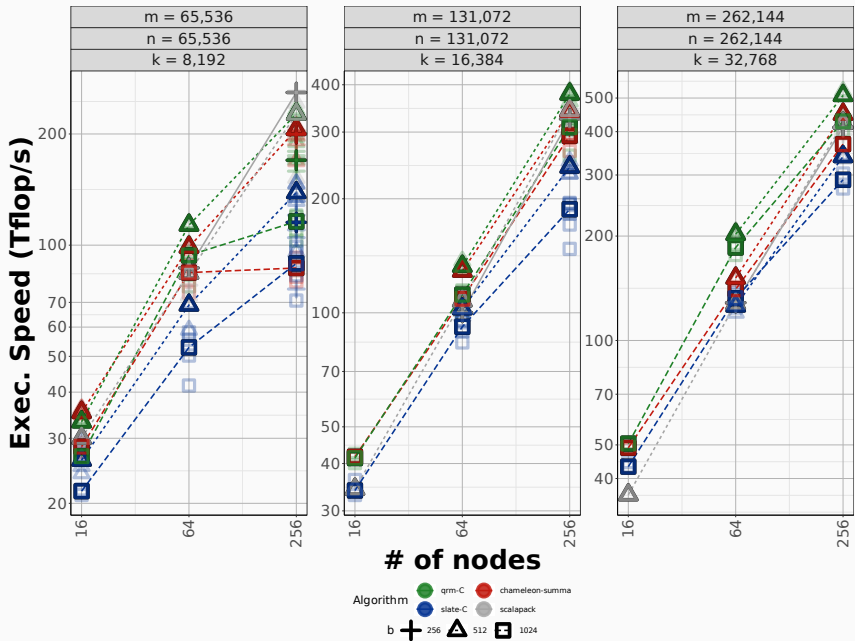
SKYLAKE NODES – LARGE A, DGEMM (FULL)



ROME NODES – SQUARE MATRICES, DGEMM (FULL)



ROME NODES – LARGE C, DGEMM (FULL)



ROME NODES – LARGE A, DGEMM (FULL)

