

# Performance Analysis and Optimization of the Tiled Cholesky Factorization on NUMA Machines

Emmanuel Jeannot

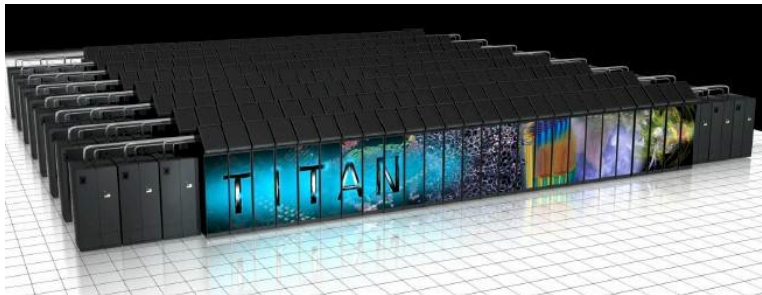
Inria Bordeaux Sud-Ouest, France

September 24, 2013

# Context

## High Performance Computing

- Treating very large problem
- on very large machine
- as fast as possible



Oak-ridge Titan supercomputer. 200 cabinets, 18 688 nodes, AMD 16-cores + NVIDIA TESLA, 20+ Pflops

# NUMA machines

## Hierarchical with several levels

- machine
- node
- socket
- caches
- cores

But one address space!



# Today's talk

Obtaining peak performance is not easy on modern NUMA machine

## This study

- What can be done statically
- What needs to be done dynamically
- Interaction between static and dynamic optimization
- Importance of the data allocation

# Outline

- 1 Introduction
- 2 The Cholesky factorization
- 3 Static analysis of compact DAG
- 4 Implementation
- 5 Conclusion

# Outline

- 1 Introduction
- 2 The Cholesky factorization**
- 3 Static analysis of compact DAG
- 4 Implementation
- 5 Conclusion

## Algorithm

The tiled version. A square matrix, symmetric definite positive decomposed in squares tiles. Compute  $L$ , such that  $A = LL^T$

```

for k = 0...T - 1 do
  A[k][k] ← DPOTRF(A[k][k])
  for m = k + 1...T - 1 do
    | A[m][k] ← DTRSM(A[k][k], A[m][k])
  end
  for n = k + 1...T - 1 do
    | A[n][n] ← DSYRK(A[n][k], A[n][n])
    | for m = n + 1...T - 1 do
      | | A[m][n] ← DGEMM(A[m][k], A[n][k], A[m][n])
    end
  end
end
end

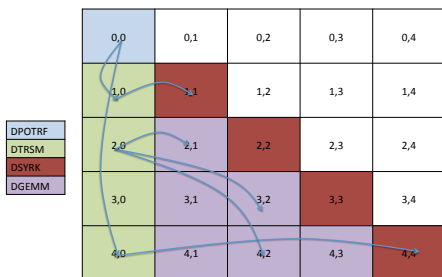
```

	0,0	0,1	0,2	0,3	0,4
	1,0	1,1	1,2	1,3	1,4
	2,0	2,1	2,2	2,3	2,4
	3,0	3,1	3,2	3,3	3,4
	4,0	4,1	4,2	4,3	4,4

DPOTRF
DTRSM
DSYRK
DGEMM

## Dependencies analysis within an iteration

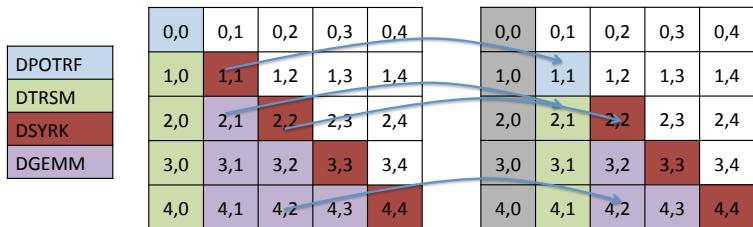


If  $T \geq 1$ :

- $\{\text{DPOTRF}(k) | 0 \leq k \leq T-1\} \rightarrow \{\text{DTRSM}(k; j) | k+1 \leq j \leq T-1\} : \{A[k][k]\}$
- $\{\text{DTRSM}(k; n) | 0 \leq k \leq T-1; k+1 \leq n \leq T-1\} \rightarrow \{\text{DSYRK}(k; n)\} : \{A[n][k]\}$
- $\{\text{DTRSM}(k; n) | 0 \leq k \leq T-1; k+1 \leq n \leq T-1\} \rightarrow \{\text{DGEMM}(k; j; n) | n+1 \leq j \leq T-1\} : \{A[n][k]\}$
- $\{\text{DTRSM}(k; n) | 0 \leq k \leq T-1; k+1 \leq n \leq T-1\} \rightarrow \{\text{DGEMM}(k; n; j) | k+1 \leq j \leq n-1\} : \{A[n][k]\}$



## Dependencies analysis between iterations

If  $T \geq 1$ :

- $\{\text{DSYRK}(k; m) | 0 \leq k \leq T-1; m = k+1\} \rightarrow \{\text{DPOTRF}(k+1)\} : \{A[m][m]\}$
- $\{\text{DSYRK}(k; m) | 0 \leq k \leq T-1; k+2 \leq m \leq T-1\} \rightarrow \{\text{DSYRK}(k+1; m)\} : \{A[m][m]\}$
- $\{\text{DGEMM}(k; m; n) | 0 \leq k \leq T-2; k+2 \leq n \leq T-1; m = k+1\} \rightarrow \{\text{DTRSM}(k+1; n)\} : \{A[n][m]\}$
- $\{\text{DGEMM}(k; m; n) | 0 \leq k \leq T-2; k+2 \leq n \leq T-1; k+2 \leq m \leq n-1\} \rightarrow \{\text{DGEMM}(k+1; m; n)\} : \{A[n][m]\}$

Dependencies analysis  $\implies$  parameterized task graph

```

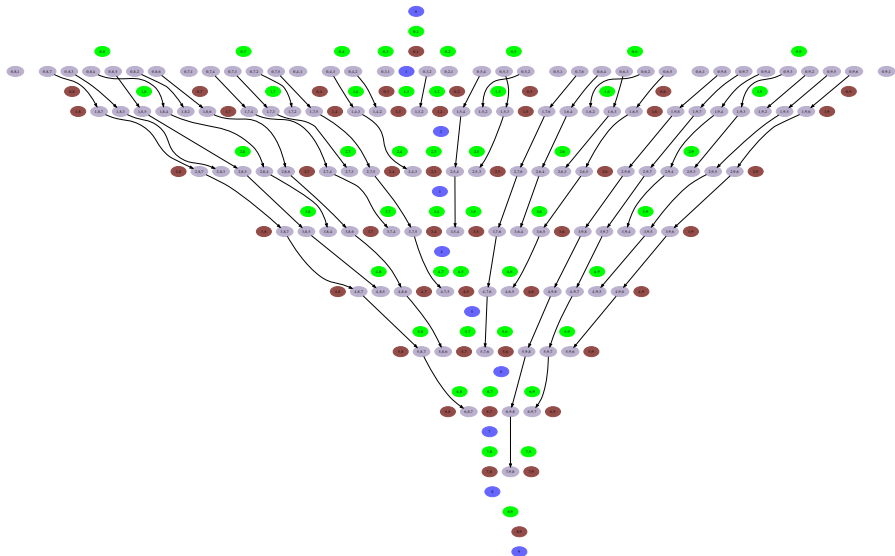
for  $k = 0 \dots T - 1$  do
   $A[k][k] \leftarrow \text{DPOTRF}(A[k][k])$ 
  for  $n = k + 1 \dots T - 1$  do
     $A[n][k] \leftarrow \text{DTRSM}(A[k][k], A[n][k])$ 
  end
  for  $m = k + 1 \dots T - 1$  do
     $A[m][m] \leftarrow \text{DSYRK}(A[m][k], A[m][m])$ 
    for  $n = m + 1 \dots T - 1$  do
       $A[n][m] \leftarrow \text{DGEMM}(A[n][k], A[m][k], A[n][m])$ 
    end
  end
end

```

If  $T \geq 1$ :

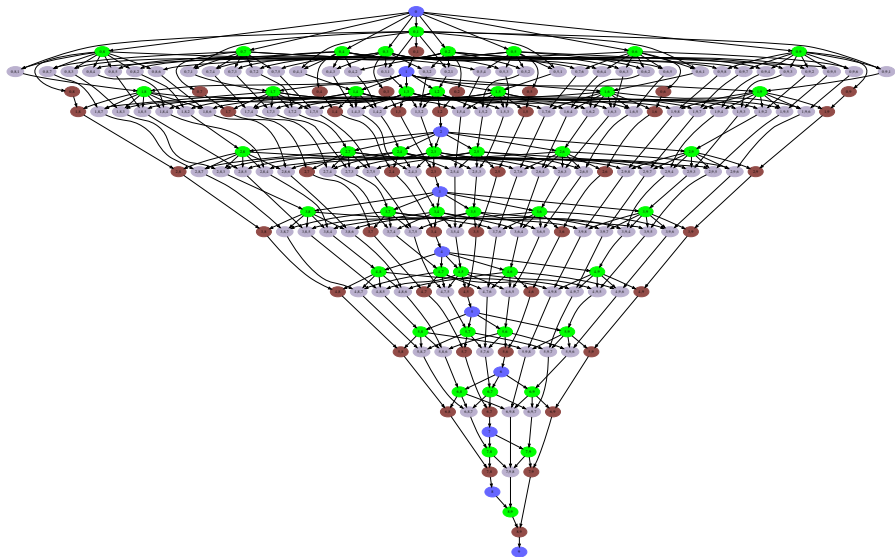
- $\{\text{DPOTRF}(k) \mid 0 \leq k \leq T - 1\} \rightarrow \{\text{DTRSM}(k; j) \mid k + 1 \leq j \leq T - 1\} : \{A[k][k]\}$
- $\{\text{DTRSM}(k; n) \mid 0 \leq k \leq T - 1; k + 1 \leq n \leq T - 1\} \rightarrow \{\text{DSYRK}(k; n)\} : \{A[n][k]\}$
- $\{\text{DTRSM}(k; n) \mid 0 \leq k \leq T - 1; k + 1 \leq n \leq T - 1\} \rightarrow \{\text{DGEMM}(k; j; n) \mid n + 1 \leq j \leq T - 1\} : \{A[n][k]\}$
- $\{\text{DTRSM}(k; n) \mid 0 \leq k \leq T - 1; k + 1 \leq n \leq T - 1\} \rightarrow \{\text{DGEMM}(k; n; j) \mid k + 1 \leq j \leq n - 1\} : \{A[n][k]\}$
- $\{\text{DSYRK}(k; m) \mid 0 \leq k \leq T - 1; m = k + 1\} \rightarrow \{\text{DPOTRF}(k + 1)\} : \{A[m][m]\}$
- $\{\text{DSYRK}(k; m) \mid 0 \leq k \leq T - 1; k + 2 \leq m \leq T - 1\} \rightarrow \{\text{DSYRK}(k + 1; m)\} : \{A[m][m]\}$
- $\{\text{DGEMM}(k; m; n) \mid 0 \leq k \leq T - 2; k + 2 \leq n \leq T - 1; m = k + 1\} \rightarrow \{\text{DTRSM}(k + 1; n)\} : \{A[n][m]\}$
- $\{\text{DGEMM}(k; m; n) \mid 0 \leq k \leq T - 2; k + 2 \leq n \leq T - 1; k + 2 \leq m \leq n - 1\} \rightarrow \{\text{DGEMM}(k + 1; m; n)\} : \{A[n][m]\}$

Rule 8,  $T = 10$ 

$$\{\text{DGEMM}(k; m; n) \mid 0 \leq k \leq T - 2; k + 2 \leq n \leq T - 1; k + 2 \leq m \leq n - 1\} \rightarrow \{\text{DGEMM}(k + 1; m; n)\} : \{A[n][m]\}$$


Expansion of the rules,  $T=10$ 

All rules



# Outline

- 1 Introduction
- 2 The Cholesky factorization
- 3 Static analysis of compact DAG**
- 4 Implementation
- 5 Conclusion

# Task/Data Allocation

## Parameterized Task Graph (PTG)

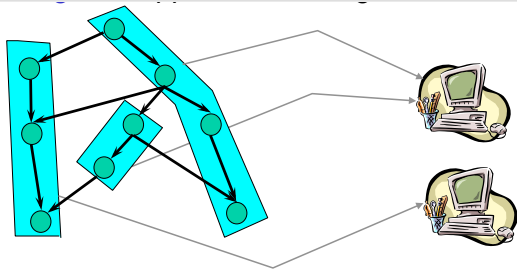
- Set of rules:

$$T_a(\vec{u}) \rightarrow T_b(\vec{v}) : D(\vec{y})|P$$

- $\forall \vec{u} \vec{v}$  and  $\vec{y}$  in  $P$  task  $T_a(\vec{u})$  sends data  $D(\vec{y})$  to task  $T_b(\vec{v})$ .
- A set of parameter. If we instantiate the parameter: we obtain the full DAG
- Symbolic allocation find a function that maps task/data to resources

# Symbolic allocation

- Using the *PTG* determine a **function** that gives the (virtual) **processor** where a task has to be **executed** and where data needs to be **stored**
- Principle** : put tasks that communicate on the same processor
- Problem** : suppress all communications  $\implies$  no more parallelism
- Trade-off** between suppressing communication and getting parallelism



# Overview of SMA

## Symbolic Mapping and Allocation (SMA) algorithm

3 steps:

- Select point to point (**bijection**) communication rules
- Suppress **conflicting** rules
- Perform mapping by expressing constraints: compute **symbolic allocation**



# Overview of SMA

## Symbolic Mapping and Allocation (SMA) algorithm

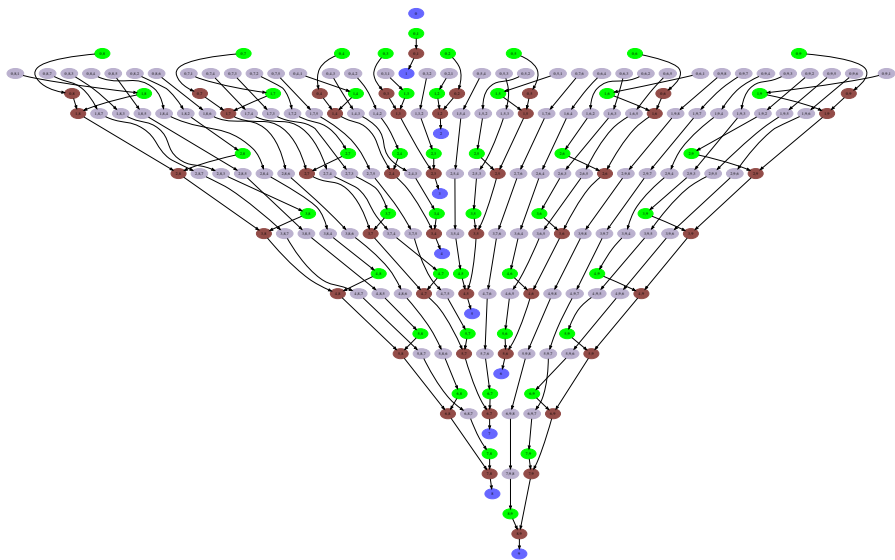
3 steps:

- Select point to point (**bijective**) communication rules
- Suppress **conflicting** rules
- Perform mapping by expressing constraints: compute **symbolic allocation**

5 bijection rules:

- $\{ \text{DPOTRF}(k) \mid 0 \leq k \leq T-1 \} \rightarrow \{ \text{DTRSM}(k; j) \mid k+1 \leq j \leq T-1 \} : \{ A[k][k] \}$
- $\{ \text{DTRSM}(k; n) \mid 0 \leq k \leq T-1; k+1 \leq n \leq T-1 \} \rightarrow \{ \text{DSYRK}(k; n) \} : \{ A[n][k] \}$
- $\{ \text{DTRSM}(k; n) \mid 0 \leq k \leq T-1; k+1 \leq n \leq T-1 \} \rightarrow \{ \text{DGEMM}(k; j; n) \mid n+1 \leq j \leq T-1 \} : \{ A[n][k] \}$
- $\{ \text{DTRSM}(k; n) \mid 0 \leq k \leq T-1; k+1 \leq n \leq T-1 \} \rightarrow \{ \text{DGEMM}(k; n; j) \mid k+1 \leq j \leq n-1 \} : \{ A[n][k] \}$
- $\{ \text{DSYRK}(k; m) \mid 0 \leq k \leq T-1; m = k+1 \} \rightarrow \{ \text{DPOTRF}(k+1) \} : \{ A[m][m] \}$
- $\{ \text{DSYRK}(k; m) \mid 0 \leq k \leq T-1; k+2 \leq m \leq T-1 \} \rightarrow \{ \text{DSYRK}(k+1; m) \} : \{ A[m][m] \}$
- $\{ \text{DGEMM}(k; m; n) \mid 0 \leq k \leq T-2; k+2 \leq n \leq T-1; m = k+1 \} \rightarrow \{ \text{DTRSM}(k+1; n) \} : \{ A[n][m] \}$
- $\{ \text{DGEMM}(k; m; n) \mid 0 \leq k \leq T-2; k+2 \leq n \leq T-1; k+2 \leq m \leq n-1 \} \rightarrow \{ \text{DGEMM}(k+1; m; n) \} : \{ A[n][m] \}$

## All bijection rules



# Conflicting Rules?

## Two types of conflicts

- fork conflict** : if two rules describe the same sending task instance for different receiving tasks instances
- join conflict** : if two rules describe the same receiving task instance for different sending tasks instances

## 1 fork conflict

- c  $\{DTRSM(k; n) | 0 \leq k \leq T - 1; k + 1 \leq n \leq T - 1\} \rightarrow \{DSYRK(k; n)\} : \{A[n][k]\}$
- $\{DSYRK(k; m) | 0 \leq k \leq T - 1; m = k + 1\} \rightarrow \{DPOTRF(k + 1)\} : \{A[m][m]\}$
- c  $\{DSYRK(k; m) | 0 \leq k \leq T - 1; k + 2 \leq m \leq T - 1\} \rightarrow \{DSYRK(k + 1; m)\} : \{A[m][m]\}$
- $\{DGEMM(k; m; n) | 0 \leq k \leq T - 2; k + 2 \leq n \leq T - 1; m = k + 1\} \rightarrow \{DTRSM(k + 1; n)\} : \{A[n][m]\}$
- $\{DGEMM(k; m; n) | 0 \leq k \leq T - 2; k + 2 \leq n \leq T - 1; k + 2 \leq m \leq n - 1\} \rightarrow \{DGEMM(k + 1; m; n)\} : \{A[n][m]\}$

To ensure as much parallelism as possible

Only non conflicting rules need to be clustered (selected)

## Identification

Find a **functions** that gives the **cluster number** of task  $T_a(\vec{u})$  for any valid value of  $\vec{u}$  and the mapping of tile  $A[i][j]$

- 1 Let  $\mathcal{Z}$  a set of non conflicting rules
- 2  $\vec{\rho}$  the vector of parameters
- 3 Each rule:  $T_a(\vec{u}) \rightarrow T_b(\vec{v}) : D(\vec{y})|P$
- 4 Affine [fea94] data mapping function  $\mu(D, \vec{y}) = \alpha_D \vec{y} + \beta_D + \gamma_D \vec{\rho}$
- 5 Affine clustering function  $\kappa(T_a, \vec{u}) = \alpha_a \vec{u} + \beta_a + \gamma_a \vec{\rho}$
- 6 Build system of equations for all rules in  $\mathcal{Z}$
- 7 Solve the system

# Example

## Data mapping (owner-compute rule)

- From the code, we have:  $T_a(\vec{u})$  writes  $D(\vec{y})$
- Therefore  $\kappa(T_a, \vec{u}) = \mu(D, \vec{y})$
- Hence,  $\vec{\alpha}_a \vec{u} + \vec{\beta}_a + \vec{\gamma}_a \vec{p} = \vec{\alpha}_D \vec{y} + \vec{\beta}_D + \vec{\gamma}_D \vec{p}$

## Clustering

- Each rule clustered rule:  $T_a(\vec{u}) \longrightarrow T_b(\vec{v}) : D(\vec{y})|P$
- $\kappa(T_a, \vec{u}) = \kappa(T_b, \vec{v})$
- Therefore:  $\vec{\alpha}_a \vec{u} + \vec{\beta}_a + \vec{\gamma}_a \vec{p} = \vec{\alpha}_b \vec{v} + \vec{\beta}_b + \vec{\gamma}_b \vec{p}$

## Cholesky Case

If we select the first conflicting rule:

- $\mu(A[i][j]) = 0$
- $\kappa(\text{DPOTRF}(k)) = k$
- $\kappa(\text{DTRSM}(k, n)) = k + n$
- $\kappa(\text{DSYRK}(k, m)) = m$
- $\kappa(\text{DGEMM}(k, m, n)) = m + n$

If we select the second conflicting rule:

- $\mu(A[i][j]) = i + j$
- $\kappa(\text{DPOTRF}(k)) = 2 \times k$
- $\kappa(\text{DSYRK}(k, n)) = 2 \times n$
- $\kappa(\text{DTRSM}(k, n)) = k + n$
- $\kappa(\text{DGEMM}(k, m, n)) = m + n$



# Outline

- 1 Introduction
- 2 The Cholesky factorization
- 3 Static analysis of compact DAG
- 4 Implementation**
- 5 Conclusion



# Memory node?

**Solution:** put tiles  $A[i][j]$  on cluster  $i + j$

In practice: tile  $A[i][j]$  is put on memory:

$$i + j \pmod G$$

## G?

- NUMA machine: memory hierarchy
- Which granularity is the most efficient?
  - $G = 1$  (shared memory)?
  - $G =$  number of memory bank?
  - $G =$  number of core (cache)?

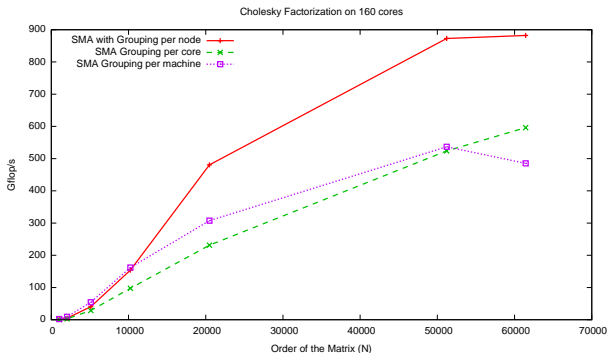
# Experiment conditions

- 160 cores NUMA machine composed of 20 nodes of one 8 cores socket Intel Nehalem Eagleton (E7-8837) at 2.67GHz
- icc 11.1-075
- MKL 11.1-075 (comparison and compute kernels)
- Plasma 1.4.5 + Quark runtime system

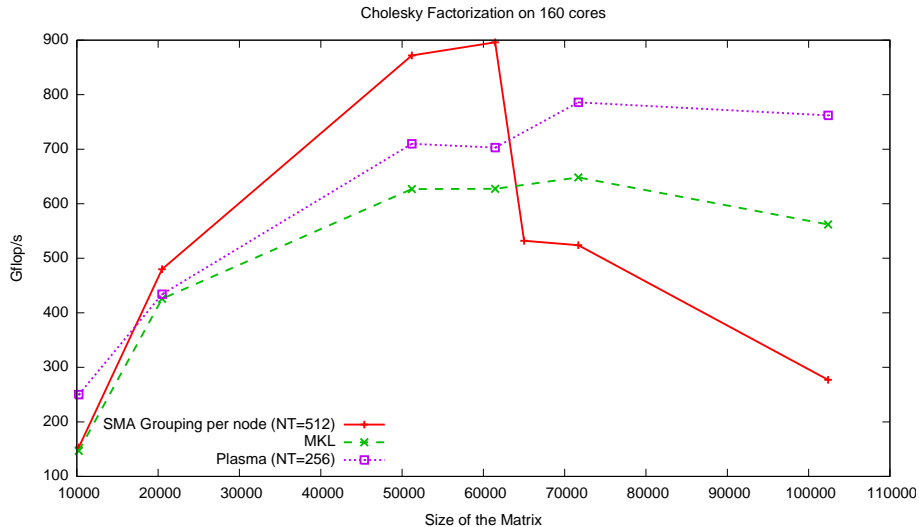
$$i + j \bmod G$$

G?

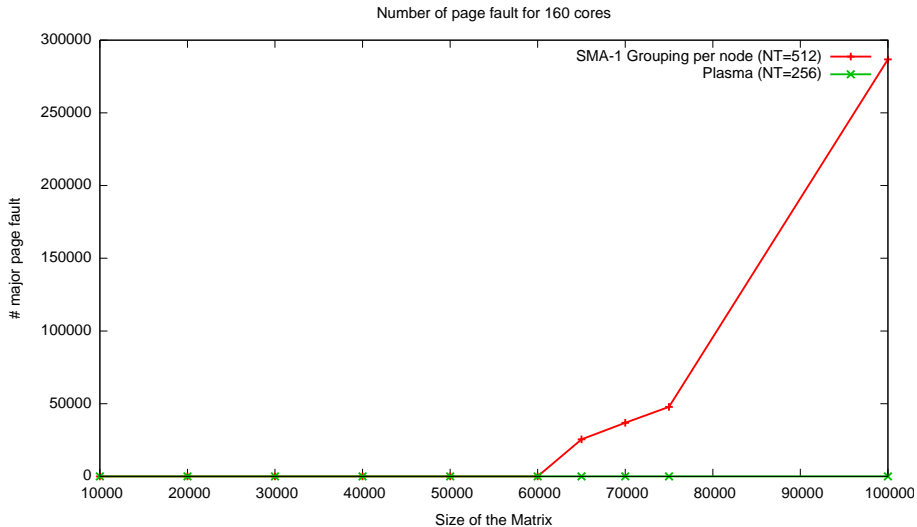
- Best compromise:
  - number of memory bank
- bind threads to memory node (hwloc)
- bind tile to memory node (`numa_alloc_onnode()`)



## Pure performance



# Page fault



# Using swap on a 600 Gb machine?

## Memory usage

- $N=64000$ : more than 30 Gb to store the matrix
- 20 nodes: 30 Gb/node

# Using swap on a 600 Gb machine?

## Memory usage

- $N=64000$ : more than 30 Gb to store the matrix
- 20 nodes: 30 Gb/node

## Memory allocation

- With `malloc` pages are put on the first node that writes it
- Sequential allocation: everything is put on the same node
- $N > 64000$  not enough memory on a node

# Using swap on a 600 Gb machine?

## Memory usage

- $N=64000$ : more than 30 Gb to store the matrix
- 20 nodes: 30 Gb/node

## Memory allocation

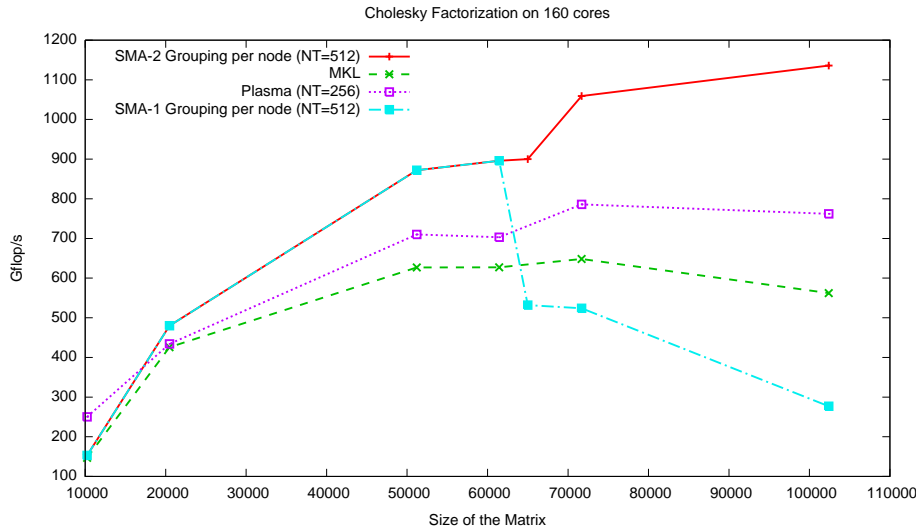
- With `malloc` pages are put on the first node that writes it
- Sequential allocation: everything is put on the same node
- $N > 64000$  not enough memory on a node

## Solutions

- directly store data in tile layout
- have a multithreaded I/O for creating the matrix in lapack format
- allocate pages of matrix in round robin fashion onto the nodes (`numa_alloc_interleaved`)



## Final result



# Comparison with Block-Cyclic Mapping

a <sub>11</sub>	a <sub>12</sub>	a <sub>13</sub>	a <sub>14</sub>	a <sub>15</sub>	a <sub>16</sub>	a <sub>17</sub>	a <sub>18</sub>	a <sub>19</sub>
a <sub>21</sub>	a <sub>22</sub>	a <sub>23</sub>	a <sub>24</sub>	a <sub>25</sub>	a <sub>26</sub>	a <sub>27</sub>	a <sub>28</sub>	a <sub>29</sub>
a <sub>31</sub>	a <sub>32</sub>	a <sub>33</sub>	a <sub>34</sub>	a <sub>35</sub>	a <sub>36</sub>	a <sub>37</sub>	a <sub>38</sub>	a <sub>39</sub>
a <sub>41</sub>	a <sub>42</sub>	a <sub>43</sub>	a <sub>44</sub>	a <sub>45</sub>	a <sub>46</sub>	a <sub>47</sub>	a <sub>48</sub>	a <sub>49</sub>
a <sub>51</sub>	a <sub>52</sub>	a <sub>53</sub>	a <sub>54</sub>	a <sub>55</sub>	a <sub>56</sub>	a <sub>57</sub>	a <sub>58</sub>	a <sub>59</sub>
a <sub>61</sub>	a <sub>62</sub>	a <sub>63</sub>	a <sub>64</sub>	a <sub>65</sub>	a <sub>66</sub>	a <sub>67</sub>	a <sub>68</sub>	a <sub>69</sub>
a <sub>71</sub>	a <sub>72</sub>	a <sub>73</sub>	a <sub>74</sub>	a <sub>75</sub>	a <sub>76</sub>	a <sub>77</sub>	a <sub>78</sub>	a <sub>79</sub>
a <sub>81</sub>	a <sub>82</sub>	a <sub>83</sub>	a <sub>84</sub>	a <sub>85</sub>	a <sub>86</sub>	a <sub>87</sub>	a <sub>88</sub>	a <sub>89</sub>
a <sub>91</sub>	a <sub>92</sub>	a <sub>93</sub>	a <sub>94</sub>	a <sub>95</sub>	a <sub>96</sub>	a <sub>97</sub>	a <sub>98</sub>	a <sub>99</sub>

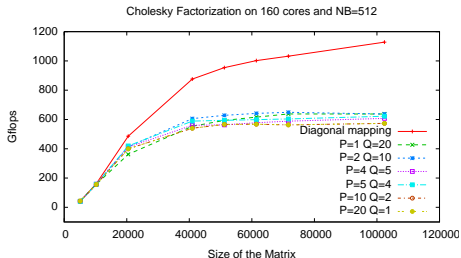
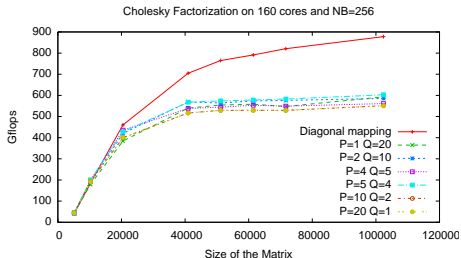
Global View

	0			1			2		
	a <sub>11</sub>	a <sub>12</sub>	a <sub>17</sub>	a <sub>18</sub>	a <sub>13</sub>	a <sub>14</sub>	a <sub>19</sub>	a <sub>15</sub>	a <sub>16</sub>
	a <sub>21</sub>	a <sub>22</sub>	a <sub>27</sub>	a <sub>28</sub>	a <sub>23</sub>	a <sub>24</sub>	a <sub>29</sub>	a <sub>25</sub>	a <sub>26</sub>
0	a <sub>51</sub>	a <sub>52</sub>	a <sub>57</sub>	a <sub>58</sub>	a <sub>53</sub>	a <sub>54</sub>	a <sub>59</sub>	a <sub>55</sub>	a <sub>56</sub>
	a <sub>61</sub>	a <sub>62</sub>	a <sub>67</sub>	a <sub>68</sub>	a <sub>63</sub>	a <sub>64</sub>	a <sub>69</sub>	a <sub>65</sub>	a <sub>66</sub>
	a <sub>91</sub>	a <sub>92</sub>	a <sub>97</sub>	a <sub>98</sub>	a <sub>93</sub>	a <sub>94</sub>	a <sub>99</sub>	a <sub>95</sub>	a <sub>96</sub>
	a <sub>31</sub>	a <sub>32</sub>	a <sub>37</sub>	a <sub>38</sub>	a <sub>33</sub>	a <sub>34</sub>	a <sub>39</sub>	a <sub>35</sub>	a <sub>36</sub>
	a <sub>41</sub>	a <sub>42</sub>	a <sub>47</sub>	a <sub>48</sub>	a <sub>43</sub>	a <sub>44</sub>	a <sub>49</sub>	a <sub>45</sub>	a <sub>46</sub>
1	a <sub>71</sub>	a <sub>72</sub>	a <sub>77</sub>	a <sub>78</sub>	a <sub>73</sub>	a <sub>74</sub>	a <sub>79</sub>	a <sub>75</sub>	a <sub>76</sub>
	a <sub>81</sub>	a <sub>82</sub>	a <sub>87</sub>	a <sub>88</sub>	a <sub>83</sub>	a <sub>84</sub>	a <sub>89</sub>	a <sub>85</sub>	a <sub>86</sub>

Local (distributed) View

P=2, Q=3

# Comparison with Block-Cyclic Mapping



# Conclusion

## Shared memory machine

Shared memory machine gives the illusion of a flat address space but:

- data allocation
- data movement

have a huge impact on the performance

# Conclusion

## Shared memory machine

Shared memory machine gives the illusion of a flat address space but:

- data allocation
- data movement

have a huge impact on the performance

## Scientific challenges

Low-level system commands can help to solve:

- data allocation
- data (re)distribution

# Outline

- 1 Introduction
- 2 The Cholesky factorization
- 3 Static analysis of compact DAG
- 4 Implementation
- 5 Conclusion**

# Conclusion

## Wrap-up

- SMA can be applied to new tiled algorithms: results are promising
- Memory is not flat. Data allocation have a huge impact on performance
- Many thing can be compiled statically
- But a good runtime system is necessary to achieve efficiency!

Acknowledgement: Emmanuel Agullo, George Bosilca, Michel Cosnard, Brice Goglin, *Morse associated team*, Tao Yang, . . .