VIII Международная молодежная научно-практическая школа "Высокопроизводительные вычисления на GRID системах"

# Будущее вычислений: от мобильных платформ до экзафлопсных суп<mark>ерко</mark>мпьютерных систем

#### Чл.-корр. РАН, профессор, Заместитель директора НИВЦ МГУ, Заведующий кафедрой Суперкомпьютеров и квантовой информатики ВМК МГУ

Вл.В.Воеводин

voevodin@parallel.ru

6-е февраля, 2017, САФУ, Архангельск

VIII International Youth Scientific School "High-Performance Computing Using GRID Systems"

# Future Computing: from Mobile Platforms to Exascale Supercomputing Systems

Prof. Vladimir Voevodin Deputy Director, Research Computing Center, MSU Head of Department on Supercomputers and Quantum Informatics, CMC, MSU

voevodin@parallel.ru

February 6<sup>th</sup>, 2017, NARFU, Arkhangelsk

#### http://msu.mnc.ru

# Lomonosov Moscow State University (established in 1755)

41 faculties 350+ departments 5 major research institutes

45 000+ students, 2500 full doctors (Dr.Sci.), 6000 PhDs, 1000+ full professors, 5000 researchers.

# MSU Computing Center, 1956



# MSU Computing Center, 1959



# Computing History of Moscow State University (from 1956 up to now)



Computing History of Moscow State University (from 1956 up to now)

Today the key statement in "computing area" is "Building Up Parallel Computing World". Why ?

There are a number of reasons...

# Computer World



Supercomputers

Servers...

PCs, Laptops...

Tablets, Smartphones...



And a second s

A DESCRIPTION OF THE OWNER OF THE

# Software Development Process

#### (We always make a compromise...)



# Another reason – diversity of architectures



# Software Development Process

#### (We always make a compromise...)



# Portability... Is It Always Difficult?



# *Portability... Is It Always Difficult?*



#### **Complexity of Supercomputing Centers** (reasons to talk about "Building Up Parallel Computing World")



# **Complexity of Supercomputing Centers**

# Is it difficult to control few components ? A few ?..



# A few? Info on MSU "Lomonosov" Supercomputer : (1.7 Pflops, 6000 computing nodes, 12K CPUs, 2K GPUs...)



Current trend: all these numbers grow extremely fast !

It's impossible to predict/describe a state of supercomputers... We have almost lost control over supercomputers...

#### Structure of algorithms and applications (one more reason to talk about "Building Up Parallel Computing World")



### Efficiency of Supercomputing Centers (another reason to talk about "Building Up Parallel Computing World")



Average performance (one core) of "Chebyshev" supercomputer for 3 days

Today the key statement in "computing area" is "Building Up Parallel Computing World". Why ?

> There are a number of reasons... There are a lot of questions...

How well do we know properties, features of Parallel Algorithms ? Should we think about Parallel Algorithms? Yes... Unfortunately, Yes...

How well do we know static and dynamic characteristics of Parallel Programs ? Should we think about properties of Parallel Programs? Yes... Unfortunately, Yes...

How well do we know architectures of parallel computers ? Should we think about architectures? Yes... Unfortunately, Yes...

# Have we ever met these questions before ? Yes, often...

teledert

#### Main stages to solve problems (Supercomputing co-design chain)



# Main stages to solve problems: Methods and Algorithms (GAUSS elimination method)



#### Main stages to solve problems (Ubiquitous parallelism)



#### Main stages to solve problems (Where supercomputing co-design begins)

#### How to ensure efficiency of this supercomputing co-design chain ?



### Supercomputing Co-Design Technologies and Tools (A practical approach)



#### Supercomputing Co-Design Technologies and Tools (A practical approach)



### Supercomputing Co-Design Technologies and Tools (JobDigest of an application)



#### Supercomputing Co-Design Technologies and Tools (analysis at different levels)



#### Supercomputing Co-Design Technologies and Tools (How well do we control a state of supercomputers?)



#### What could be a reason of this situation?

- Hardware failure?

Yes, it could be ...

Yes ... Yes ...

Yes ....

- Software failure?
- Error in the code?
- Algorithmic problem?

Large numbers in supercomputers: cores, processors, accelerators, nodes, HW&SW components, files, indexes, users, projects, processes, threads, running and queued jobs...

F & INATOOPHI

We don't know for sure the current state of supercomputer's components ...

*What is now?* We hope that a component works until we get an evidence that it has failed.

What do we need?

# Our expectations = Reality

**FA**-INATOOPHI

We need a guarantee: if something goes wrong inside a supercomputer we shall be notified immediately.

We want a system behaves in a way we expect it should behave.

If discrepancy occurs between our expectations and supercomputer behavior (i.e. reality) we need to know immediately about it.

How can we do that? Supercomputer is huge, we can't control it to a full extent any more.

But... supercomputer can do it itself (instead of us), if we explain what "our expectations" are.

### Supercomputing Co-Design Technologies and Tools (A practical approach)





Supercomputers should be autonomous in self-control. (They become more dynamic, more sophisticated, more and more parallel)

... The larger supercomputers, the more autonomous they should be...
### Supercomputing Co-Design Technologies and Tools (OctoTron: predictability and autonomous life of supercomputers)

A guarantee of "our expectations = reality", how this can be done?

- a formal model of supercomputers (model is a graph),
- a set of formal rules,
- a set of reactions,

Autonomous life and control of MSU supercomputers:

-"Chebyshev" supercomputer, 60 Tflops, 625 CPUs: 10 228 nodes, 24 698 edges, 205 044 attributes, 160 rules, 100 reactions;
- "Lomonosov" supercomputer, 1.7 Pflops, 12 000 CPUs, 2 000 GPU:

116 000 nodes, 332 000 edges, 2 400 000 attributes,...

Initial deployment, Detection of faults, critical and emergency situations, Turning off minimum amount of hardware, Self diagnostics, Previous accidents, etc. are done according to a model and rules.

Current trend: many decisions about control over HW&SW of supercomputers must be taken automatically. Algorithms and programs... Are they important in the supercomputing co-design chain?

# What is a good parallel program / algorithm?



Have we ever met these questions before ? Yes, often...

Six generations of computer architectures – – six battles for parallel applications, for high performance, good scalability and efficiency...



Vector computers

Mid 70-s.

**Features**: pipelined functional units, vector instructions, vector registers.

**Programming**: vectorization of the innermost loops.

Cray-1 supercomputer



Cray X-MP supercomputer



Cray Y-MP supercomputer

### Vector-parallel computers

80-s.

**Features**: pipelined functional units, vector instructions, vector registers. 2-32 processors, shared memory.

**Programming**: vectorization of the innermost loops, parallelization of outer loops.



Cray T3D supercomputer



Intel Paragon XPS140 supercomputer

Massive-parallel computers

90-s.

**Features**: thousands of processors, distributed memory.

**Programming**: explicit message passing and data distribution, MPI.



**DEC AlphaServer** 



Sun StarFire supercomputer

Shared-memory computers

Mid 90-s.

**Features**: tens/hundreds of processors, shared memory.

**Programming**: single address space, local and shared variables, OpenMP.



MSU supercomputer "Chebyshev"



Clusters of SPM-nodes

2000-s.

**Features**: tens/hundreds thousands of nodes, distributed memory. Tens of processors/cores, shared memory.

**Programming**: hybrid parallel programming technology MPI + OpenMP.

"K" supercomputer



MSU supercomputer "Lomonosov"



Mid 2000-s.

**Features**: tens/hundreds thousands of nodes, distributed memory. Tens of processors/cores, shared memory. Several accelerators (GPUs, Phi's).

Programming: MPI+OpenMP+OpenCL/CUDA;



Tianhe-2 supercomputer

### Generations of Parallel Computer Architectures (or How often have we have to rewrite applications completely?)

Parallel programming paradigms (from the 70s up to now): 70s - Loop Vectorization (innermost) 80s - Loop Parallelization (outer) + Vectorization (innermost) 90s - MPI mid 90s - OpenMP mid 2000s - MPI+OpenMP 2010s - CUDA, OpenCL, MPI+OpenMP+accelerators (GPU, Xeon Phi)

Changes in computer architectures do not change algorithms! But...

For each generation of a new computing platform we have to:

- Analyze algorithms to find a way to match better features and properties of the platform ;

- Express the properties of algorithms we found to obtain efficient Implementation for the platform.

### Changes in computer architectures do not change algorithms! (Algorithms remain the same)

Are these figures different?











### Generations of Parallel Computer Architectures (or How often have we have to rewrite applications completely?)

Parallel programming paradigms (from the 70s up to now): 70s - Loop Vectorization (innermost) 80s - Loop Parallelization (outer) + Vectorization (innermost) 90s - MPI mid 90s - OpenMP mid 2000s - MPI+OpenMP 2010s - CUDA, OpenCL, MPI+OpenMP+accelerators once and for all

Changes in computer architectures <u>do not change algorithms</u>! But For each generation of a new computing platform we have to: - Analyze algorithms to find a way to match better features and properties of the platform ;

- Express the properties of algorithms we found to obtain efficient Implementation for the platform. What are key properties of an algorithm we need to analyze and describe to obtain an efficient implementation in the future?

What properties are important?

Unified (complete) description of an algorithm: What do we need to take into account ?

### **Description of Algorithms** (What should be included in this description?)

Information Graph **Determinacy** Computational kernel Locality of computations Macrostructure Scalability Performance Data locality Mathematical description Communication profile **Properties and Features** Efficiency Serial Complexity Computational intensity Resource of Parallelism Input / Output data

### Description

For positive definite Hermitian matrices (symmetric matrices in the real case), we use the decomposition  $A = LL^*$ , where L is the lower triangular matrix  $\mathbf{B}$ , or the decomposition  $A = U^*U$ , where U is the upper triangular matrix  $\mathbf{B}$ . These forms of the Cholesky decomposition are equivalent in the sense of the amount of arithmetic operations and are different in the sense of data representation. The essence of this decomposition consists in the implementation of formulas obtained uniquely for the elements of the matrix L from the above equality. The Cholesky decomposition is widely used due to the following features.

Mathematical Description	Remarks on the Algorithm
Input data: a symmetric positive definite matrix $A$ whose elements are denoted by $a_{ij}$ ).	The Cholesky decomposition allows one to
Output data: the lower triangular matrix $L$ whose elements are denoted by $l_{ij}$ ).	use the so-called accumulation mode due to
The Cholesky algorithm can be represented in the form	the fact that the significant part of
1 /	computation involves dot product
$l_{11} = \sqrt{a_{11}},$	operations. Hence, these dot products can
$l_{j1} = \frac{a_{j1}}{l_{11}},  j \in [2, n],$	be accumulated in double precision for
	additional accuracy. In this mode, the
i-1	Cholesky method has the least equivalent
$l_{ii} = \int a_{ii} - \sum l_i^2 \int i \in [2, n],$	perturbation. During the process of
$\sim_{ii}$ $\bigvee \sim_{ii}$ $\sum_{n=1}^{in}$ $\sim_{ip}$ , $v \in [-, v]$ ,	decomposition, no growth of the matrix
p-1	elements can occur, since the matrix is
$l = \left( \sum_{i=1}^{i-1} l_i \right) / l_i = i \in [0, n-1]$	symmetric and positive definite. Thus, the
$\iota_{ji} = \left( a_{ji} - \sum_{n=1} \iota_{ip} \iota_{jp} \right) / \iota_{ii},  i \in [2, n-1], j \in [i+1, n].$	Cholesky algorithm is unconditionally stable.
$\mu = 1$	



#### Additional Info

There exist block versions of this algorithm;





### 1 Properties and structure of the algorithm

### 1.1 General description

**The Cholesky decomposition algorithm** was first proposed by Andre-Louis Cholesky (October 15, 1875 - August 31, 1918) at the end of the First World War shortly before he was killed in battle. He was a French military officer and mathematician. The idea of this algorithm was published in 1924 by his fellow officer and, later, was used by Banachiewicz in 1938 [7]. In the Russian mathematical literature, the Cholesky decomposition is also known as the square-root method [1-3] due to the square root operations used in this decomposition and not used in Gaussian elimination.



Summary

Originally, the Cholesky decomposition was used only for dense real symmetric positive definite matrices. At present, the application of this decomposition is much wider. For example, it can also be employed for the case of Hermitian matrices. In order to increase the computing performance, its block versions are often applied.

In the case of sparse matrices, the Cholesky decomposition is also widely used as the main stage of a direct method for solving linear systems. In order to reduce the memory requirements and the profile of the matrix, special reordering strategies are applied to minimize the number of arithmetic operations. A number of reordering strategies are used to identify the independent matrix blocks for parallel computing systems.



\* Scalability, performance, efficiency were measured on MSU "Lomonosov" Supercomputer



\* Dynamic Properties were obtained on MSU "Lomonosov" Supercomputer

# It is very useful information about the algorithm, we desperately need it.

But... CHALLENGES are everywhere...



- How to show dependency of the graph on a problem size ?





Do i = 1, n  
Do j = 1, n  

$$A(i,j) = 0$$
  
Do k = 1, n  
 $A(i,j) = A(i,j) + B(i,k)*C(k,j)$ 

How to formalize the information structure? Polyhedra, inequalities, vector functions...



$$\begin{cases} 1 \le i \le n \\ 1 \le j \le n \\ k = 1 \\ \begin{cases} i_1 = i \\ j_1 = j \\ u_3 (1) \end{cases} \end{cases}$$

2



## Potential parallelism: how to extract, describe, show...? (challenges of the algorithm description)



### Potential parallelism: how to extract, describe, show...? (challenges of the algorithm description)



(resource of "mathematical" parallelism)

## Potential parallelism: how to extract, describe, show...? (challenges of the algorithm description)



### Data locality: a number of open questions (challenges of the algorithm description)

How to evaluate spatial and temporal data locality of a program ?

How to compare spatial and temporal data locality of programs ?



Data locality: a number of open questions (challenges of the algorithm description)

How to evaluate spatial and temporal data locality of a program ?

How to compare spatial and temporal data locality of programs ?

Can we predict data locality in future implementations by using information from algorithms only ?

What does it mean "data locality of algorithm" ? What does it mean "algorithm with good/bad locality" ? **There are no data structures in algorithms**, locality can't be applied to algorithms! At the same time algorithms form the basis of programs...

## **Properties and Structures of Algorithms...**

Can we represent and describe them?

Can we analyze algorithms **once and for all** 

Changes in computer architectures <u>do not change algorithms</u>!

For each generation of a new computing platform we have to:

 Analyze algorithms to find a way to match better features and properties of the platform ;

- Express the properties of algorithms we found to obtain efficient Implementation for the platform.

## **Properties and Structures of Algorithms...**

# Yes, we can! AlgoWiki

Etcledellalde

http://AlgoWiki-Project.org

Can we analyze algorithms **once and for all** 

 Analyze algorithms to find a way to match better features and properties of the platform ;

### Supercomputing Co-Design Technologies and Tools (A practical approach)


#### **Description of Algorithms** (What should be included in this description?)

Information Graph Determinacy Computational kernel Algorithms: Theoretical Part (machine-independent properties, "Once and for all") Resource of Parlelism Input / Output data

D

### <u>AlgoWiki</u>



#### http://AlgoWiki-Project.org

### Building Up Intelligible Parallel Computing World

## A final remark...

### Building Up Parallel Computing World

Problems and Challenges







Parallel Computing Technologies







Parallel Computing Education



VIII International Youth Scientific School "High-Performance Computing Using GRID Systems"

# Thank you!

Prof. Vladimir Voevodin Deputy Director, Research Computing Center, MSU Head of Department on Supercomputers and Quantum Informatics, CMC, MSU

voevodin@parallel.ru

February 6<sup>th</sup>, 2017, NARFU, Arkhangelsk

VIII International Youth Scientific School "High-Performance Computing Using GRID Systems"

### Благодарю за внимание!

Prof. Vladimir Voevodin Deputy Director, Research Computing Center, MSU Head of Department on Supercomputers and Quantum Informatics, CMC, MSU

voevodin@parallel.ru

February 6<sup>th</sup>, 2017, NARFU, Arkhangelsk