Heterogeneous platform "HybriLIT"

HLIT-

60





The International IT-School "Machine Learning, Parallel and Hybrid Computations & Big Data Analytics" in frames of the International Conference "Mathematical Modeling and Computational Physics, 2019" Stará Lesná, High Tatra Mountains, Slovakia 01–04 July 2019

Heterogeneous platform HybriLIT



The commissioned supercomputer "Govorun" named after Nikolai Nikolayevich Govorun, with the development of information whom technologies at JINR has been connected since 1966, became a natural elaboration of the heterogeneous cluster HybriLIT being the MICC component. The supercomputer "Govorun" is a joint project of the Bogoliubov Laboratory of Theoretical Physics and the Laboratory of Information Technologies supported by the JINR Directorate. The project is aimed at the significant speed-up of complex theoretical and experimental studies in the field of nuclear physics and condensed matter physics underway at JINR including the development of the NICA megaproject computing.

05

0

HI F

<u>ה</u>

0

0

U

60

D

0)

The platform consists of two elements: the education and testing polygon and the supercomputer "Govorun", combined by the unified software and information environment.



The supercomputer "Govorun" components



GPU-component based on NVIDIA DGX-1 Volta. CPU-component based on the newest Intel architectures: Intel Xeon Phi and Intel Skylake processors The supercomputer "Govorun" is designed to carry out resource-intensive and massive parallel calculations for the solution of a wide range of challenges facing JINR, which becomes possible due to the heterogeneity (presence of different types of computing accelerators) of the supercomputer computing architecture.

At present, the supercomputer is used to solve problems that require massive parallel calculations in various fields of nuclear physics and high energy physics, particularly in lattice quantum chromodynamics to study the properties of hadronic matter at high energy density and baryon charge and in the presence of supramaximal electromagnetic fields, for mathematical modeling of interactions of antiprotons with protons and nuclei using DPM, FTF and UrQMD+SMM generators developed at JINR and being of interest for the NICA-MPD experiment, as well as for solving applied problems.



with store/retrieve data. This will add 230 TFlops to "Govorun" system, almost doubling the CPU part performance.

The GPU-component of the supercomputer "Govorun"



The GPU-component consists of 5 NVIDIA DGX-1 servers. Each server has 8 GPU NVIDIA Tesla V100 based on the latest architecture NVIDIA Volta. Moreover, one server NVIDIA DGX-1 has 40960 cores CUDA, which are equivalent to 800 high-performance central processors. A whole number of novel technologies are used in DGX-1, including the NVLink 2.0 wire with the bandwidth up to 300 Gb/s.

The GPU-component gives a users of the supercomputer a possibility to allow as massively parallel computation for general-purpose tasks using such technologies as CUDA and OpenCL, as well as use applications already adapted for this architecture. Also, GPU-component allow to use machine learning and deep learning algorithms for solving applied problems by neural network approach.





Caffe



h







Unified software-hardware environment

The unified software and information environment of the HybriLIT platform allows users to use the education and testing polygon is aimed at exploring the possibilities of novel computing architectures, IT-solutions, to develop and debug their applications, furthermore, carry out calculations on the supercomputer, which allows them to effectively use the supercomputer resources.

Software and information environment of the HybriLIT



The unified software and information environment including the unified system level (the operation system, the job scheduler, file systems and software) as well as a set of services allowing users to quickly get the answers to their questions, jointly develop parallel applications, information about receive conferences, seminars and meetings dedicated to parallel programming technologies.

0 D \bigcirc U 60 D 0)



Implementation of the neural network approach, methods

To provide all the possibilities both for developing mathematical models and algorithms and carrying out resource-intensive calculations including graphics accelerators, which significantly reduce the calculation time, an ecosystem for tasks of machine learning and deep learning (ML/DL) and data analysis has been created and is actively developing for HybriLIT platform users.



The created ecosystem has two components:

• the computing component is aimed at carrying out resource-intensive, massive parallel tasks of neural network training using NVIDIA graphics accelerators;

• the component for the development of models and algorithms on the *JupyterHub* basis, i.e. a multi-user platform for working with *Jupyter Notebook* (known as *IPython* with the possibility to work in a web browser), including several libraries and frameworks.

Educational program on the HybriLIT platform

The HybriLIT platform is used as a polygon for training students, post-graduate students and young scientists in the field of HPC. Training is conducted both by young staff members of the Institute and its Member States and by students of Dubna University. The main tool in the educational process is an education and testing polygon which is the basic platform for basic courses on "Computer System Architecture", "High Performance Computing Technologies" and "Mathematical Models in Physics"; the number of students is about 200 people per academic year.

 \bigcirc

60

D

d

0000



10



In the last decade novel computational facilities and technologies has become available: MPI – OpenMP – CUDA – OpenCL...



It is not easy to follow modern trends. Modification of the existing codes or developments of new ones ?

Parallel technologies: levels of parallelism



How to control hybrid hardware: MPI – OpenMP – CUDA - OpenCL ...



Remote access to the cluster

For Windows users:

T

0

gen

etel

In order to connect from the Windows OS, it is necessary to use a special program – SSH-client, e.g. PuTTY.

- 1. To install PuTTY, it is necessary to copy **putty.exe** file to your computer and launch it.
- 2. Enter **hydra.jinr.ru** in the field **Host Name (or IP address)**
- Enter the preferable name for the current connection in the field **Saved Sessions**
- 3. To support remote graphical interface, choose **Connection > SSH > X11** and mark it in the field **Enable X11 forwarding**
- 4. Check if the field Local ports accept connection from other hosts is marked in **Connection > SSH > Tunnels**
- 5. After the chosen settings go back to the **Sessions** and click **Save** in order to save the current settings.
- 6. Click *Open* to connect HybriLIT and enter login and password.

😵 PuTTY Configuration X	🕵 PuTTY Configuration X	🕵 PuTTY Configuration X	PuTTY Configuration
Basic options for your PuTTY session Category: Session Logging Terminal Methods Specify the destination you want to connect to Host Name (or IP address) Point Point Point Provide Appearance Behaviour Translation Selection Colours Connection Data Provy Teinet Riogin Bestail Close window on exit: Aways Never Only on clean exit	Category: Bell Options controlling SSH X11 forwarding Features Window Authorization Selection Colours Connection Data Proxy Teinett Riogin SSH Mathematication Particle for local display Browse Browse Browse	Category: Options controlling SSH pot forwarding Peatures Pot forwarding Pot of forwarding Pot forwarding Pot of forwarding	Category: Session Basic options for your PuTTY session Specify the destination you want to connect to Hot Name (or IP address) Point Peatures Window - Appearance Behaviour Translation - Colours Connection Data - Proxy - Teinet - Riogin - SSH - Kex - Cipher - Aduth
About Open Cancel	About Open Cancel	About Open Cancel	About Open Cancel

Remote access to the cluster for Linux users:

Launch the terminal, type in the command line:

\$ ssh username@hydra.jinr.ru

where **username** is your login

13





5 basic steps to carry out computations on the platform

User login to the cluster:

ssh username@hydra.jinr.ru



3. Prepared script:
3. Prepared script:
GPU-Script.sh
#!/bin/sh
Example for one-GPU applications
#SBATCH -p gpu
#
Number of GPUs
#SBATCH --gres=gpu:1

./test1



The list of main commands in the work with the modules:

Module commands	Description
module avail	list available modules
module add <name module="" of="" the=""></name>	load chosen module
module list	lists currently loaded modules
module rm <name module="" of="" the=""></name>	unloads module foo and removes all changes
	that it made in the environment
module clear	unloads all modules

The list of main SLURM commands:

SLURM commands	Description	
sinfo	reports the state of partitions and nodes	
	managed by SLURM	
squeue	reports the state of jobs	
<pre>sbatch <name of="" script="" the=""></name></pre>	submit a job script for later execution (the script	
	typically contains one or more srun commands to	
	launch parallel tasks)	
scancel	cancel a pending or running job	





OpenMP parallel programming technology: example

Task: find the sum of two arrays C = A + B, the elements of which are determined by some functions.



Compilation and launch of OpenMP-applications

- \$ module add intel/v2019.0.018
- \$ icc -qopenmp solver_Openmp.c -mkl -o test1
- \$ sbatch script_omp.sh

```
script_omp.sh
    #!/bin/sh
1.
2.
    # Example with 5 cores for OpenMP
    #SBATCH -p tut
3.
    # Number of cores per task
4.
5.
    #SBATCH - c 5
   #Set OMP NUM THREADS to the same# value as -c
6.
7.
    if [ -n "$SLURM_CPUS_PER_TASK" ]; then
      omp_threads=$SLURM_CPUS_PER_TASK
      else omp threads=1
8.
9.
    fi
    export OMP_NUM_THREADS=$omp_threads
10.
    export OMP DISPLAY ENV=TRUE
11.
12.
        ./test1
13. # End of submit file
```





S

0

U en

2

MPI parallel programming technology: example







MPI parallel programming technology

rank = 0

#include <stdio.h>
#include <mpi.h>

int main(int argc, char* argv[]){
 int numprocs, rank;

MPI_Init(&argc,&argv); MPI_Comm_size(MPI_COMM_WORLD,&numprocs); MPI_Comm_rank(MPI_COMM_WORLD,&rank); printf("Numprocs is %d; Hello from %d\n", numprocs, rank);

MPI_Finalize();

Numprocs is 6; Hello from 0 Numprocs is 6; Hello from 1 Numprocs is 6; Hello from 5 Numprocs is 6; Hello from 3 Numprocs is 6; Hello from 4 Numprocs is 6; Hello from 2

rank = 0 🛛 📉	rank = 1	rank = numprocs - 1
<pre>#include <stdio.h> #include <mpi b=""></mpi></stdio.h></pre>	<pre>#include <stdio.h> #include <mpi h=""></mpi></stdio.h></pre>	<pre>#include <stdio.h> #include <mpi h=""></mpi></stdio.h></pre>
<pre>int main(int argc, char* argv[]){ int numprocs, rank;</pre>	<pre>int main(int argc, char* argv[]){ int numprocs, rank;</pre>	<pre>int main(int argc, char* argv[]){ int numprocs, rank;</pre>
<pre>MPI_Init(&argc,&argv); MPI_Comm_size(MPI_COMM_WORLD,</pre>	<pre>MPI_Init(&argc,&argv); MPI_Comm_size(MPI_COMM_WORLD, &numprocs); MPI_Comm_rank(MPI_COMM_WORLD, &rank); printf("Numprocs is %d; Hello from %d\n", numprocs, rank); MPI_Finalize(); }</pre>	<pre>MPI_Init(&argc,&argv); MPI_Comm_size(MPI_COMM_WORLD, &numprocs); MPI_Comm_rank(MPI_COMM_WORLD, &rank); printf("Numprocs is %d; Hello from %d\n", numprocs, rank) MPI_Finalize(); }</pre>



MPI parallel programming technology

The 6 most important MPI commands:



gather

reduction

Point-to-Point Operation Send count element with type Process Process type to process dest dest tag Buffer - Program address space that MPI_Send(buffer, count, type, dest, tag, comm) references the data that is to be sent Process k Receive count Process Process element with type type tag from process source source **Buffer** - Program address space that references the MPI Recv(buffer, count, type, source, tag, comm, status) data that is to be received. **Collective Communication MPI_Barrier()** – barrier synchronization **MPI Bcast()** – broadcasting **MPI_Gather** ()– gathering of data from all process to the master process MPI_Scatter() - scattering of data from master process to all processes in group MPI_Reduce(), MPI_AllReduce() - global operations (*MPI_MIN, MPI_MAX, MPI_SUM, MPI_PROD* etc)





Intel: compilers, tools for developing, debugging and profiling parallel applications, mathematical library ...

Intel Inspector Memory and Threading Checking

Intel Math Kernel Library Optimized Routines for Science, Engineering, and Financial Intel Aplication Performance Snapshot Tool for application profiling

Intel Data Analytics Acceleration Library Optimized for Data Analytics & Machine Learning Intel Advisor Vectorization Optimization and Thread Prototyping

Intel® Integrated Performance Primitives Image, Signal, and Compression Routines

Intel Threading Building Blocks Task-Based Parallel C++ Template Library



Source: https://software.intel.com/en-us/articles/intel-parallel-studio-xe-release-notes

Intel Math Kernel Library. LAPACK : Solving Systems of Linear Equations

Computes the solution to the system of linear

lapack_int LAPACKE_dgtsv (

and multiple right-hand sides.

int matrix_layout ,

lapack_int n ,

double * dL ,

double * d ,

double * du ,

double * b ,

lapack_int Ldb);

lapack_int nrhs ,

equations with a tridiagonal coefficient matrix A

?gtsv

$$\mathbf{A} \cdot Y = B$$



Compilation

24

- \$ module add intel/v2019.0.018
- \$ icc tridiagonal_solver.c -mkl -qopenmp -o test1

https://software.intel.com/en-us/mkl-developer-reference-c-lapack-linear-equation-driver-routines

60



CUDA programming model

Definitions: *device* = GPU; *host* = CPU; *kernel* = function that runs on the device

Parallel parts of an application are executed on the device as *kernels*

- One *kernel* is executed at a time
- Several threads execute each *kernel*

kernel <<< dim3 dG, dim3 dB, nSgMem, nStream >>> (args)

- dim3 **dG** number of blocks dimension (grid)
- dim3 **dB** dimension of blocks
- **nShMem** the amount of shared memory to be allocated for each block
- **nStream** the stream number

dim3 is an integer vector type that can be used in CUDA code.

```
dim3 dG(512); // 512x1x1
```

```
dim3 dB(1024, 1024); // 1024x1024x1
```

Function Type Qualifiers

Declaration specifier	Callable from	Executed on
global	host	device
host	host	host
device	device	device

CUDA application: GPU vectors sum

$$c_i = a_i + b_i, i = 0, ..., N - 1$$

Serial code	CUDA code
<pre>void sum(int N, int *a, int *b,</pre>	global void kernel(int N, int *a, int *b, int *c)
<pre>int tid = 0; while(tid < N){</pre>	<pre>int tid = threadIdx.x + blockIdx.x *</pre>
<pre>c[tid] = a[tid] + b[tid]; tid += 1; }</pre>	<pre>while(tid < N){ c[tid] = a[tid] + b[tid]; tid += blockDim.x * gridDim.x;</pre>
}	<pre>} }</pre>





Compilation and execution CUDA applications



Add module \$ module add cuda/v10.1-1

Compilation

\$ nvcc -arch=compute_35 deviceInfo.cu -o test1

1.	#!/bin/sh	script_CUDA.sl
2.	#SBATCH -p tut	
3.	# Number of GPU per node	
4.	#SBATCHgres=gpu:1	
5.	<pre># Set time of work (in minutes)</pre>	
6.	#SBATCH -t 60	
7.	<pre># Connect module for CUDA</pre>	
8.	module add cuda/v10.1-1	
9.	./test1	
10.	<pre># End of submit file</pre>	2

Run in batch mode **\$ sbatch script_CUDA**

Ecosystem for ML/DL, data analysis tasks (basic libraries and tools)



scikit-learn – machine learning library for Python programming language.



NumPy the is fundamental package for scientific computing with Python. It contains among other things: a powerful Ndimensional array object sophisticated (broadcasting) functions tools for integrating C/C++ and Fortran code useful linear algebra, Fourier transform, random and number capabilities



SciPy – is a Pythonbased ecosystem of opensoftware for source mathematics, science, and engineering. SciPy contains modules for optimization, linear algebra, integration, special interpolation, functions, FFT, signal and processing, ODE image solvers and other tasks common in science and engineering.



matplotlib – is a Python 2D plotting library which produces publication quality figures



pandas – Flexible and powerful data analysis / manipulation library for Python, providing labeled data structures similar to R data.

TensorFlow – is a free and opensource a symbolic math library, and is also used for machine learning applications such as neural networks. It is used for both research and production at Google.





Ecosystem for ML/DL tasks





Ecosystem for tasks of machine learning, deep learning and data analysis Work within the ML/DL ecosystem To get started you need to: 1. Log in with your HybriLIT account in GitLab: https://gitlab-hybrilit.jinr.ru/ 2. Enter the components (the authorization is done via GitLab): Gitl ab **Component for carrying out resource-intensive Development component** (without graphics accelerators) calculations (with graphics accelerators NVIDIA) **NVIDIA** https://jhub.jinr.ru https://jhub2.jinr.ru Jupyter Notebook After the authorization the **Jupyter Notebook** interactive environment opens: jupyter Logout Control Pane Create: Folder: Files Running Clusters **File Python 3** Select items to perform actions on them Upload New -Notebool h / 2019-IT-SH 0 -Name 4 Python 3 C Other

> Text File Folder Terminal

60

30

2019-IT-SH.ipynb

Home directories are available for users, the

file systems NFS, ZFS, EOS are reinstalled.



SLURM submission script

Submission script for MPI application (Intel realization) #!/bin/sh #Partition for run; #SBATCH -p tut #Number of MPI tasks; #SBATCH -n 8 #Set time of work (in minutes); #SBATCH -t 60 #Connect module for Intel MPI module add intel/v2019.0.018 #Submit a job for execution; mpirun ./exec_mpi

Submission script for MPI application (GNU realization)
#!/bin/sh
#Partition for run;
#SBATCH -p tut
#Number of MPI tasks;
#SBATCH -n 8
#Set time of work (in minutes);
#SBATCH -t 60
#Connect module for OpenMPI
module add openmpi/v3.1.3
#Submit a job for execution;
mpirun ./exec_mpi

Submission script for OpenMP application
#!/bin/sh
#Partition for run;
#SBATCH -p tut
#Number of cores per task;
#SBATCH -c 5
#Number of MPI tasks;
#SBATCH -n 1
#Set OMP_NUM_THREADS to the same value as -c;
export OMP_NUM_THREADS=5
#Set time of work (in minutes);
#SBATCH -t 60
#Submit a job for execution;
./exec_openmp

Submission script for CUDA and application #!/bin/sh #Partition for run; #SBATCH -p tut #Number of GPU per node; #SBATCH -gres=gpu:1 #Set time of work (in minutes); #SBATCH -t 60 #Connect module for CUDA module add cuda/v10.1-1 #Submit a job for execution; ./exec_cuda

JOINT INSTITUTE FOR NUCLEAR RESEARCH http://jinr.ru 141980, Dubna, Moscow region, Joliot-Curie 6

LABORATORY OF INFORMATION TECHNOLOGIES

http://lit.jinr.ru Phone: (+7 49621) 64-019 Fax: (+7 49621) 65-145

MULTIFUNCTIONAL INFORMATION AND COMPUTING COMPLEX https://micc.jinr.ru

HETEROGENEOUS PLATFORM HYBRILIT IN LIT JINR http://hlit.jinr.ru

