European Research Council
Established by the European Commission

*Slide of the seminar*

# NewTurb Cluster

New cluster for Turbulence simulations

# NewTurb Cluster

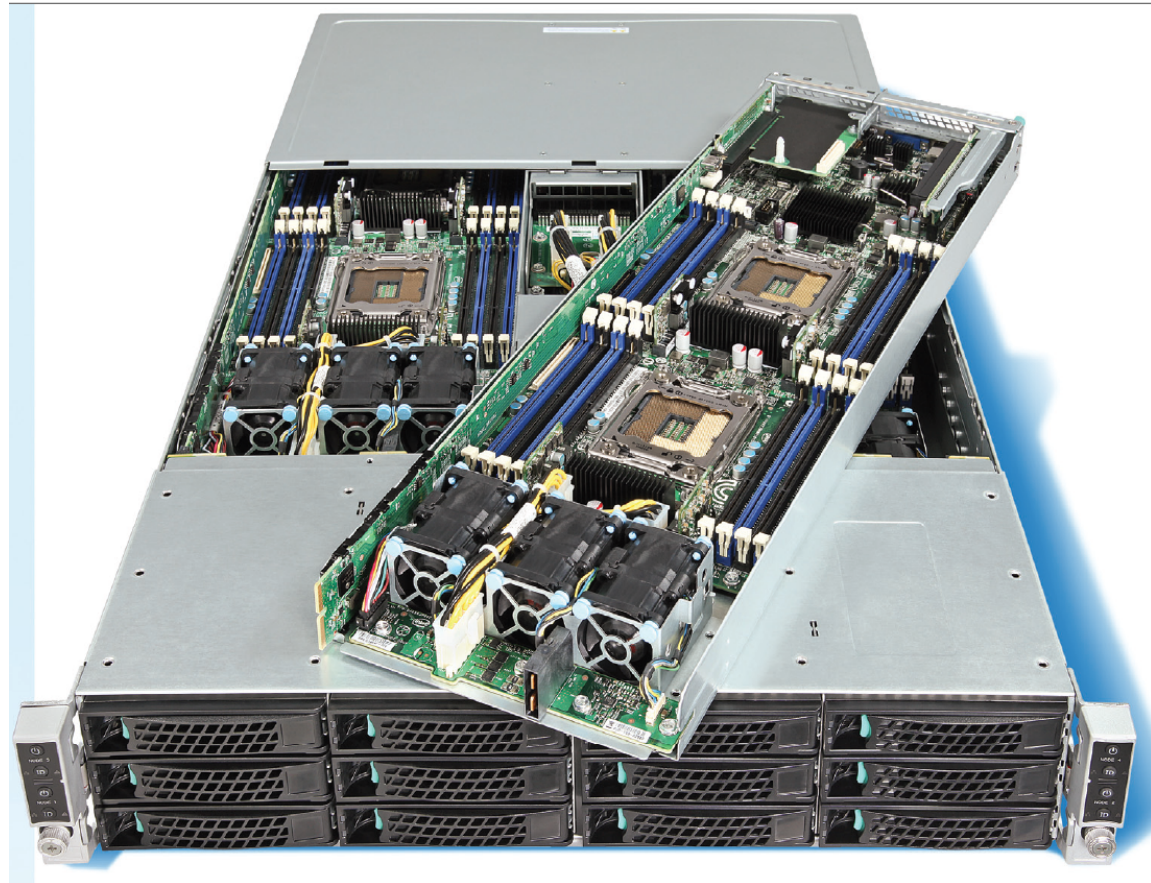New cluster for Turbulence simulations

# Outline

- Hardware presentation
- Installed software
- Batch usage
- Visualization
- Code development
- Next

# Hardware

- Cluster is composed by 3 elements:
  - Frontend
    - Two 12-cores 2.4GHz intel Xeon server
    - 256 GB RAM 1.6GHz DDR3 RAM
    - Two 56 Gbit/s QDR Infiniband
    - 1 Nvidia Quadro K4000 GPU

  - 16 computing nodes, each with
    - Two 12-cores 2.4GHz Intel Xeon chip, for a total of 24 cores
    - 256 GB 1.6GHz DDR3 RAM
    - One 56Gbit/s QDR Infiniband

  - Storage server
    - 64 GB RAM
    - 250 TB disk in RAID5 configuration
    - Two 56 Gbit/s QDR Infiniband

# 4 nodes in 2U

- Highly compact nodes:

  - 4 nodes in a 2U server

  - Energy efficient

- Maintenance

  - Hot swappable node

  - Hot swappable Hds

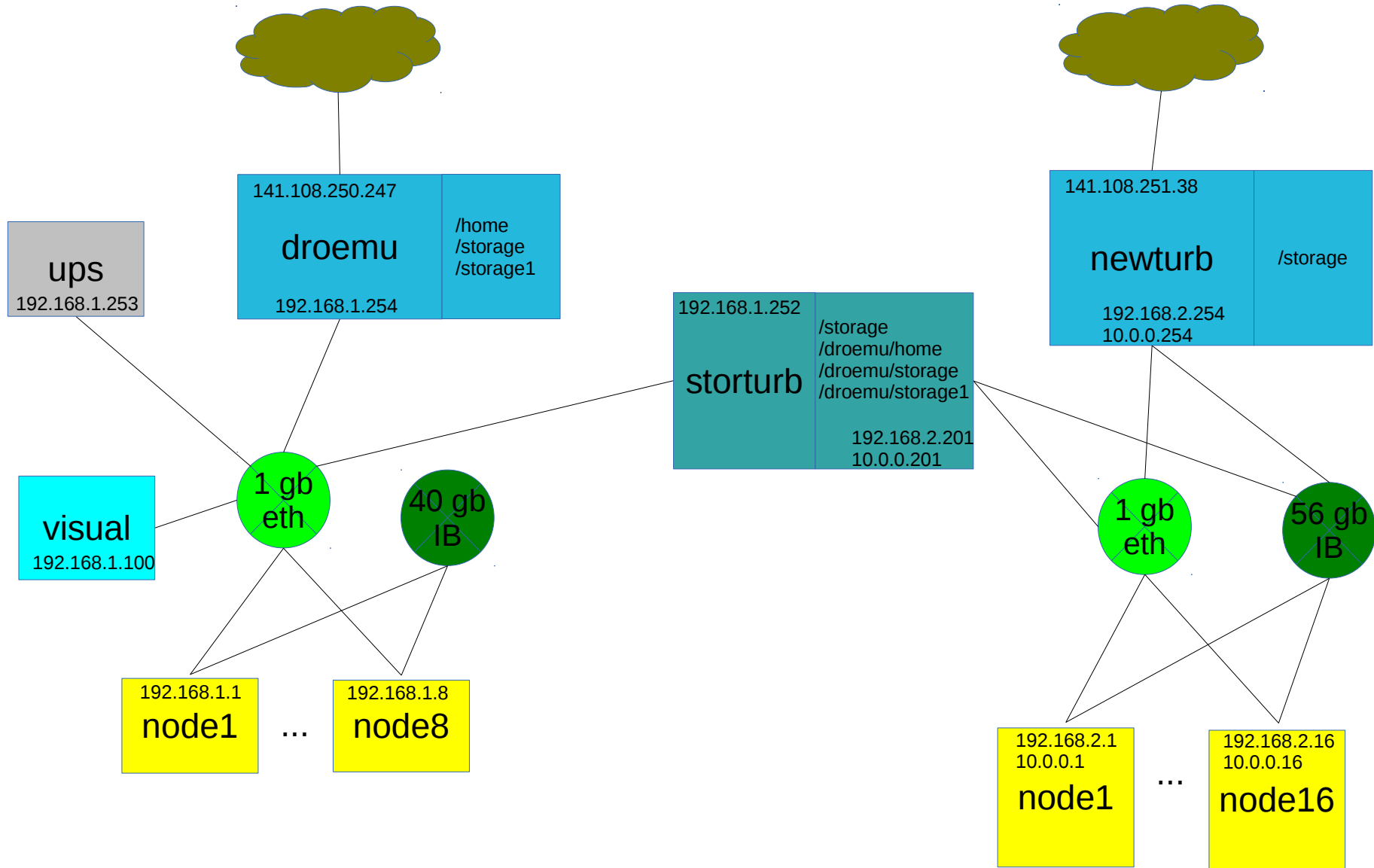- We have 4 of these

# Storage Server

- 72 x 4TB Hds

  – 288 TB raw space

- 64 GB Ram

- Dual Intel XEON 8 cores @ 2.60GHz

- Dual Infiniband FDR 56 Gb/s

- 3 HW Raid controller

# Frontend: newturb

- 256 GB Ram

- Same 2x12core Intel CPU

- Two Infiniband FDR 56Gb/s

- Quadro K4000 GPU

ups
192.168.1.253

141.108.250.247
droemu
/home
/storage
/storage1
192.168.1.254

141.108.251.38
newturb
/storage
192.168.2.254
10.0.0.254

192.168.1.252
storturb
/storage
/droemu/home
/droemu/storage
/droemu/storage1
192.168.2.201
10.0.0.201

visual
192.168.1.100

1 gb
eth

40 gb
IB

1 gb
eth

56 gb
IB

192.168.1.1
node1

...

192.168.1.8
node8

192.168.2.1
10.0.0.1
node1

...

192.168.2.16
10.0.0.16
node16

# Infiniband benchmark: bandwidth

| #bytes | #iterations | BW peak MB/s | BW avg MB/s |
|---|---|---|---|
| 2 | 10 | 1 | 1 |
| 16 | 10 | 12 | 11 |
| 64 | 10 | 49 | 45 |
| 256 | 10 | 152 | 152 |
| 512 | 10 | 347 | 312 |
| 1k | 10 | 883 | 880 |
| 4k | 10 | 3120 | 3119 |
| 16k | 10 | 4637 | 4636 |
| 32k | 10 | 5221 | 5221 |
| 128k | 10 | 5780 | 5780 |
| 1M | 10 | 5967 | 5967 |
| 2M | 10 | 5976 | 5976 |
| 4M | 10 | 5940 | 5940 |
| 8M | 10 | 5992 | 5992 |

# Infiniband benchmark: latency

| #bytes | #iterations | t_min<br>usec | | t_max<br>usec | t_typical<br>usec |
|---|---|---|---|---|---|
| 2 | 10.0 | 0.9 | | 9.7 | 1.2 |
| 16 | 10.0 | 0.9 | | 2.8 | 1.0 |
| 64 | 10.0 | 1.0 | | 2.7 | 1.0 |
| 256 | 10.0 | 1.4 | | 8.4 | 1.4 |
| 512 | 10.0 | 1.5 | | 4.5 | 1.6 |
| 1k | 10.0 | 1.8 | | 9.1 | 1.9 |
| 4k | 10.0 | 3.1 | | 4.7 | 3.2 |
| 16k | 10.0 | 5.0 | | 9.4 | 5.3 |
| 32k | 10.0 | 7.8 | | 11.8 | 7.8 |
| 128k | 10.0 | 23.5 | | 27.8 | 23.8 |
| 1M | 10.0 | 168.7 | | 174.2 | 169.4 |
| 2M | 10.0 | 334.9 | | 340.2 | 335.3 |
| 4M | 10.0 | 666.1 | | 671.6 | 667.4 |
| 8M | 10.0 | 1328.6 | | 1333.3 | 1330.8 |

# Software

- Operating system is CentOS v6.5 64bit
  - Community edition of Redhat OS: very well supported
- Libraries
  - OpenMPI v1.8.1
  - HDF v1.8.13
  - fftw v2, v3
  - P3DFFT
  - GSL
- Tools
  - CMake
  - Paraview
  - CUDA v6

# Users

- Each user logins on newturb.roma2.infn.it
- His/her home is a LOCAL directory
  - NOT shared with the whole cluster (SPEED hack)
- Place for job preparation:
  - /storage/<<USER>>
  - SHARED with the cluster
    - NFS v4
    - Infiniband link 56Gb/s (hw link speed, not user speed)

# Users: filesystems layout

- /storage is the new 250TB area
  - Shared via Infiniband

- /droemu area visible in read-only
  - Gigabit link

- Symlinks for setting up simulations

# Code

- Compile
- Link
- ldd
- valgrind
- efence
- GDB

# Code: Valgrind

- Simulate execution on a virtual PC
- Detects memory errors, memory leaks,...
- Simulate the execution
  - Slow
- Needs some help with system libs
  - Unless you want to contribute to opensource world!

# Code: Valgrind

```c
#include <stdlib.h>
int main()
{
    int p,t,b[10];
    if(p==5)        ERROR
        t=p+1;
        b[p]=100;  ERROR
    return 0;
}
```

```
                                                    fabio@pcEuhit:~
File  Edit  View  Search  Terminal  Help
==3821==     at 0x40047C: main (in /scratch/fab/presStuff/a.out)
==3821==
==3821== Use of uninitialised value of size 8
==3821==     at 0x40048C: main (in /scratch/fab/presStuff/a.out)
==3821==
==3821==
==3821== HEAP SUMMARY:
==3821==     in use at exit: 0 bytes in 0 blocks
==3821==   total heap usage: 0 allocs, 0 frees, 0 bytes allocated
==3821==
==3821== All heap blocks were freed -- no leaks are possible
==3821==
==3821== For counts of detected and suppressed errors, rerun with: -v
==3821== Use --track-origins=yes to see where uninitialised values come from
==3821== ERROR SUMMARY: 2 errors from 2 contexts (suppressed: 6 from 6)
-bash-4.1$ cc -g valg.c
-bash-4.1$ valgrind --tool=memcheck --leak-check=full ./a.out
==3890== Memcheck, a memory error detector
==3890== Copyright (C) 2002-2012, and GNU GPL'd, by Julian Seward et al.
==3890== Using Valgrind-3.8.1 and LibVEX; rerun with -h for copyright info
==3890== Command: ./a.out
==3890==
==3890== Conditional jump or move depends on uninitialised value(s)
==3890==     at 0x40047C: main (valg.c:5)
==3890==
==3890== Use of uninitialised value of size 8
==3890==     at 0x40048C: main (valg.c:7)
==3890==
==3890==
==3890== HEAP SUMMARY:
==3890==     in use at exit: 0 bytes in 0 blocks
==3890==   total heap usage: 0 allocs, 0 frees, 0 bytes allocated
==3890==
==3890== All heap blocks were freed -- no leaks are possible
==3890==
==3890== For counts of detected and suppressed errors, rerun with: -v
==3890== Use --track-origins=yes to see where uninitialised values come from
==3890== ERROR SUMMARY: 2 errors from 2 contexts (suppressed: 6 from 6)
-bash-4.1$
```

# Code: ElectricFence

- Catch memory errors as they occur
  - Run-Time Bound checking
  - HW based: fast but uses RAM
- One guard-area for each allocation
  - 4kb plus your allocation
- Two trips for extra security
  - Band after mem
  - Band before mem

# Code: ElectricFence

```c
#include <stdio.h>
#include <stdlib.h>
        int main (void)
        {
            int i;
            int *a = (int*) malloc( 9*sizeof(int));

            for ( i=0; i<=9; ++i){
                a[i] = i;
                printf ("%d\n ", a[i]);
            }

            free(a);
            return 0;
        }
```

```
fabio@pcEuhit:~

File   Edit   View   Search   Terminal   Tabs   Help

fabio@pcEuhit:~                                         ✖  root@newt
-bash-4.1$ cc efe.c; ./a.out
0
 1
 2
 3
 4
 5
 6
 7
 8
 9
-bash-4.1$
-bash-4.1$
-bash-4.1$ cc efe.c -lefence; ./a.out

  Electric Fence 2.2.2 Copyright (C) 1987-1999 Bruce Perens <bruce@perens.com>
0
 1
 2
 3
 4
 5
 6
 7
 8
Segmentation fault (core dumped)
-bash-4.1$ ▮
```

# Code: ElectricFence + GDB



- Debugger helps!

# Code: gdb

- Compile with
  - -g -O0
- Run the program
  - command r
- Step
  - command s
- Breakpoint
  - command b
- Pops in case of errors

# Code: IDE

- Eclipse

- Integrates with GDB

# Code: ldd

- Shows the shared libs that will be used

# Batch system

- Torque v4.2.7
  - Defines the queue
  - The server that accepts jobs


- Maui Scheduler v3.3.1
  - Selects job to run
  - The server that shows which job is running

# Batch system: 2

- 4 Queues defined:
  - 1 "routing" queue
    - Named route
    - Routes job to the final destination queue based on requested resources

  - 3 "execution" queue
    - Reg_256
      - Highest priority: reserved to jobs of 256 processes

    - Reg_64
      - Middle prio: for jobs of 64 processes

    - Batch
      - Lowest prio: until 64 processes

# Batch system: 3

- Job script
  - Text shell script with PBS keywords...
    - #PBS -l nodes=xx:ppn=yy,walltime=zz:zz:zz
    - nodes: number of computing nodes you want
    - ppn: number of processes in each node
    - walltime: maximum time allowed
  - ...and commands
    - mpiexec -np zz -mca btl openib,self -hostfile $PBS_NODEFILE exec-file
      - np: real number of process created

- Submit a job:
  - qsub jobScript.pbs

- Check the queue:
  - showq
  - pbstop

# Batch system: 4

- The Job is running...
  - PIC

- StdOut and StdErr will be given back at the END
  - JobScript.oxxx
  - JobScript.exxx
  - Xxx is the jobID queuing number

- Job preparation:
  - Create a directory fo each job under /storage/USER/jobxxx
  - qsub from here
  - Files will not interfere with each other
  - "Easy" debugging

# Batch system: 5

- Scheduler policy:
  - 256 cores -> Prio 1256
  - 64 cores -> Prio 164
  - 1-63 cores -> Prio 1-63

- Backfill active:
  - Small jobs can be scheduled if they don't prevent big ones from running
  - to fill the machine meanwhile

- Queue time into account
- Soon will have a factor for Fair share usage

# Queue: showq

| ACTIVE | JOBS | | | | |
|---|---|---|---|---|---|
| JOBNAME | USERNAME | STATE | PROC | REMAINING | STARTTIME |
| 372 | malapaka | Running | 64 | 14:10:00 | 06:52:26 AM |
| 373 | malapaka | Running | 64 | 17:38:50 | 10:21:16 AM |
| 380 | sahoo | Running | 64 | 22:43:14 | 03:25:40 PM |
| 3 Active Jobs 192 of 384 Processors Active (50.00%) | | | | | |
| 8 of 16 Nodes Active (50.00%) | | | | | |
| IDLE | JOBS | | | | |
| JOBNAME | USERNAME | STATE | PROC | WCLIMIT | QUEUETIME |
| 0 | Idle | Jobs | | | |
| BLOCKED | JOBS | | | | |
| JOBNAME | USERNAME | STATE | PROC | WCLIMIT | QUEUETIME |
| 374 | malapaka | Hold | 64 | 23:00:00 | 21:31:34 |
| Total Jobs:4 Active Jobs:3 Idle Jobs:0 Blocked Jobs: 1 | | | | | |

# Queue: pbstop

# Visualization

- Paraview
  - Can handle VTK, HDF5 files

  - Lots of filters

  - Can be used with a remote machine that has the data
    - No need to move data

# Visualization:2

- VTK files are native
  - Serial and parallel
  - ASCII and BINARY


- HDF5 files need a wrapper
  - XDMF, xml based describes underlying data
  - Can handle scalar, 3D, time-varying,....

# Visualization: HDF5 files

- Example HDF5 file

```
HDF5 "cb_outP24_000001.h5" {
GROUP "/" {
  GROUP "PS3D" {
    DATASET "b" {
      DATATYPE  H5T_IEEE_F64LE
      DATASPACE  SIMPLE { ( 256,
256, 129 ) / ( 256, 256, 129 ) }
    }
  }
}
}
```

- Wrapper file

```
<?xml version="1.0" ?> <!DOCTYPE Xdmf SYSTEM "Xdmf.dtd" []>
<Xdmf>
<Domain>
<Grid Name="my_Grid" GridType="Uniform">
<Topology TopologyType="3DCoRectMesh" Dimensions="256 256 129" />

<Geometry GeometryType="Origin_DxDyDz">
 <DataItem Dimensions="3" NumberType="Integer" Format="XML">0 0 0</DataItem>
<DataItem Dimensions="3" NumberType="Integer" Format="XML"> 1 1 1</DataItem>
</Geometry>

 <Attribute Name="b" AttributeType="Scalar" Center="Node">
<DataItem Dimensions="256 256 129" NumberType="Double" Precision="8"
Format="HDF">
cb_outP24_000001.h5:/PS3D/b
</DataItem>
</Attribute>
  </Grid>
 </Domain>
</Xdmf>
```
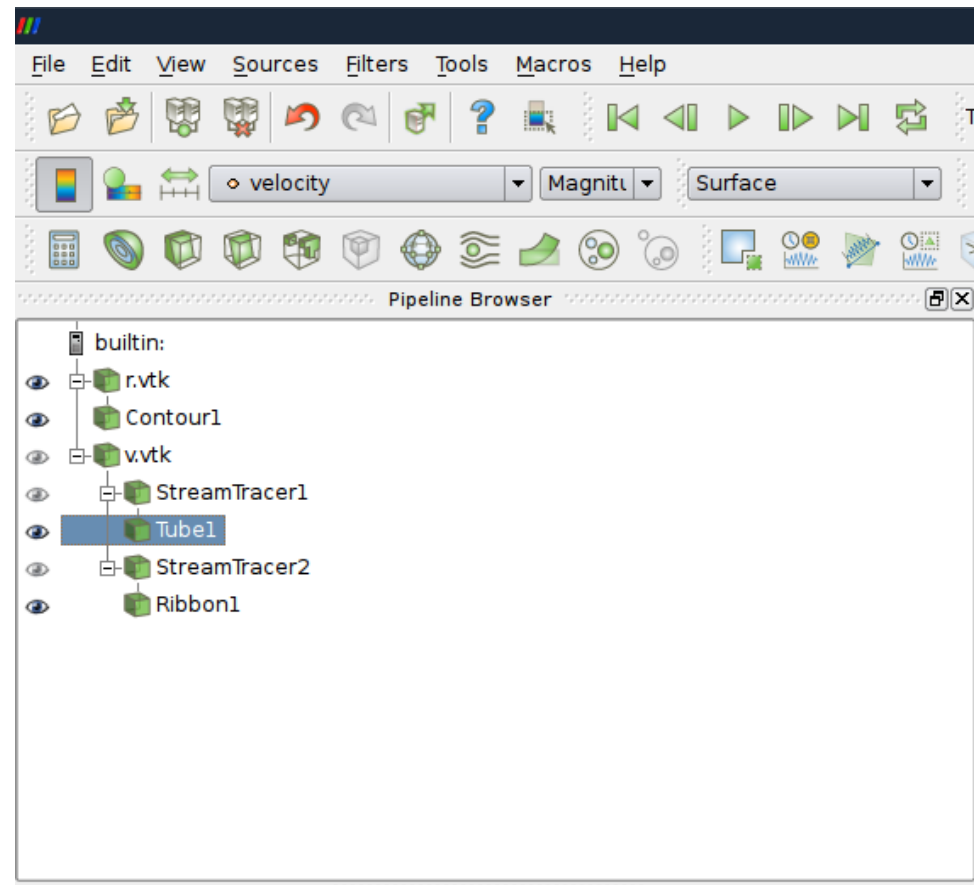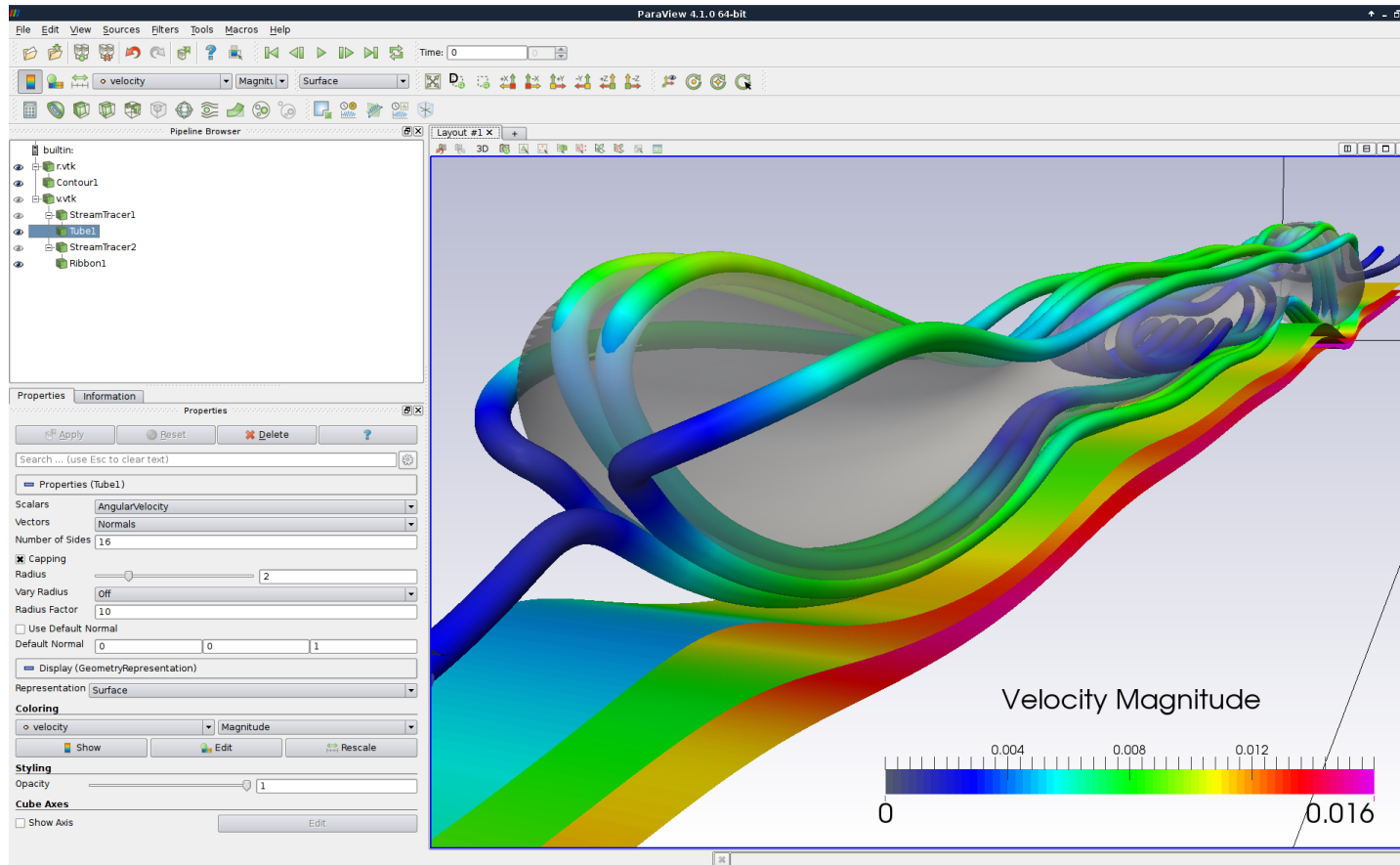
# Paraview: Filters

- Load data
- Manipulate with filters
  - Isosurfaces
  - Streamlines
  - Cutting
  - Projecting
- Pipelining metaphore
  - 1$^{st}$ filter output is 2$^{nd}$ filter input
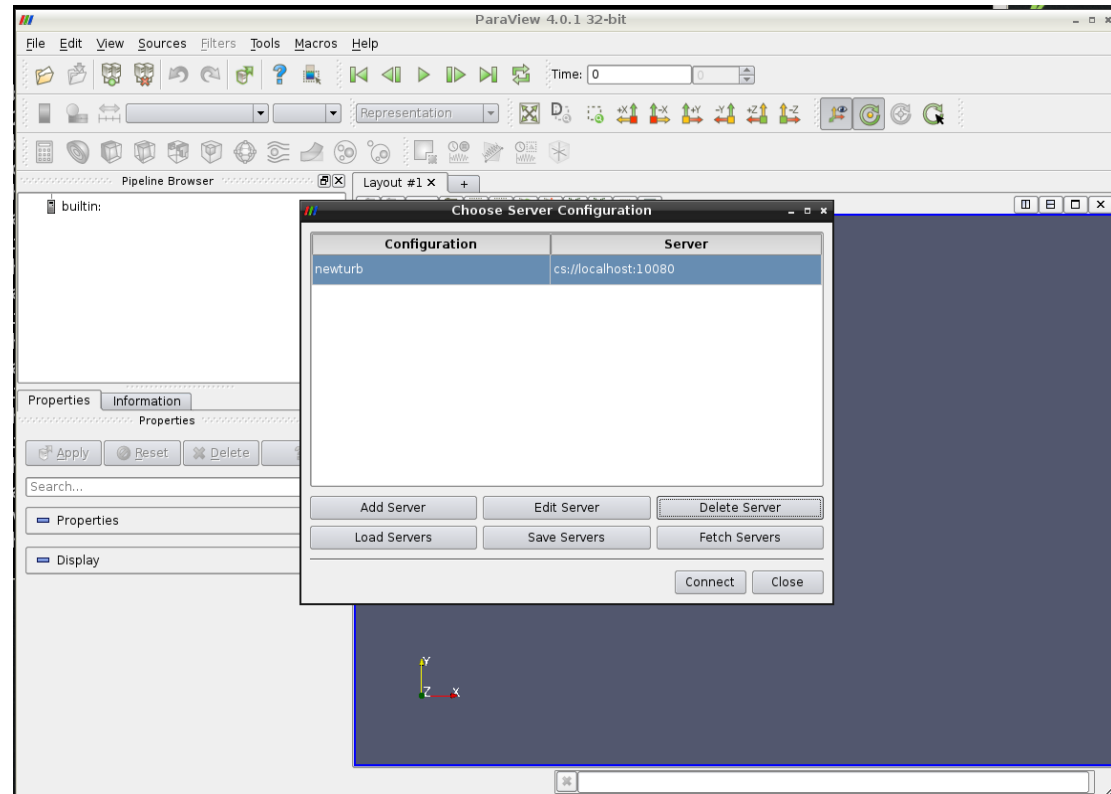
# Paraview: demo

# Paraview:example

# Paraview: client+server

- Version match
  - v4.0.1

- Load Config file
  - First time only

- File → Connect

- You're on newturb!

# Quadro K4000 GPU



## SPECIFICATIONS

| | |
|---|---|
| GPU Memory | 3GB GDDR5 |
| Memory Interface | 192-bit |
| Memory Bandwidth | 134.0GB/s |
| CUDA Cores | 768 |
| System Interface | PCI Express 2.0 x16 |
| Max Power Consumption | 80W |
| Thermal Solution | Ultra-quiet active fansink |
| Form Factor | 4.376" H × 9.50" L, Single Slot, Full Height |
| Display Connectors | DVI-I DL + 2x DP1.2 |
| Max Simultaneous Displays | 3 direct, 4 DP1.2, 2 Win XP |
| Max DP 1.2 Resolution | 3840 × 2160 at 60Hz |
| Max DVI-I DL Resolution | 2560 × 1600 at 60Hz |
| Max DVI-I SL Resolution | 1920 × 1200 at 60Hz |
| Max VGA Resolution | 2048 × 1536 at 85Hz |
| Graphics APIs | Shader Model 5.0, OpenGL 4.4, DirectX 11 |
| Compute APIs | CUDA, DirectCompute, OpenCL |

- 768 Kepler Cuda cores
- 800 Mhz, 2 FLOP/cyc
- ~1250TFlop/s SP

<p style="color:red">BUT</p>

- 3 GB RAM
- HOST <-> GPU BW

# NEXT

- Other nodes are arriving
  - Cluster will soon have little more then 512 procs
    - 2.5 TFLOP/s DP

  - RAM will be upgraded to 5.5 TB