# NEXTGenIO Performance Tools for In-Memory I/O

holger.brunst@tu-dresden.de
ZIH, Technische Universität Dresden

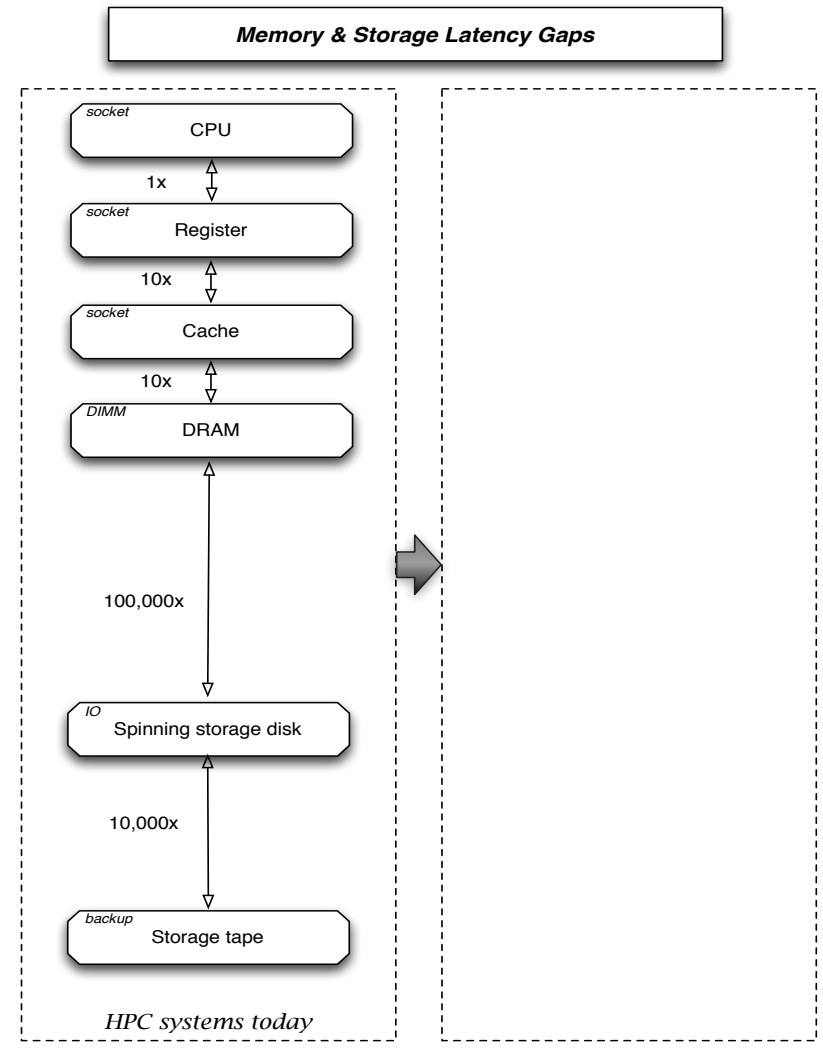UIOP Workshop, Hamburg 22nd-23rd March 2017

# Credits

Intro slides by Adrian Jackson (EPCC)

# A new hierarchy

- New non-volatile memory technology is going to change the memory hierarchy we have

- What does that mean for applications, particularly scientific simulations?

- I/O performance is one of the critical components for scaling up HPC applications and enabling HPDA applications at scale



Memory & Storage Latency Gaps

HPC systems today

# NEXTGenIO summary

## Project

- Research & Innovation Action
- 36 month duration
- €8.1 million
- Approx. 50% committed to hardware development

## Partners

- EPCC
- INTEL
- FUJITSU
- BSC
- TUD
- ALLINEA
- ECMWF
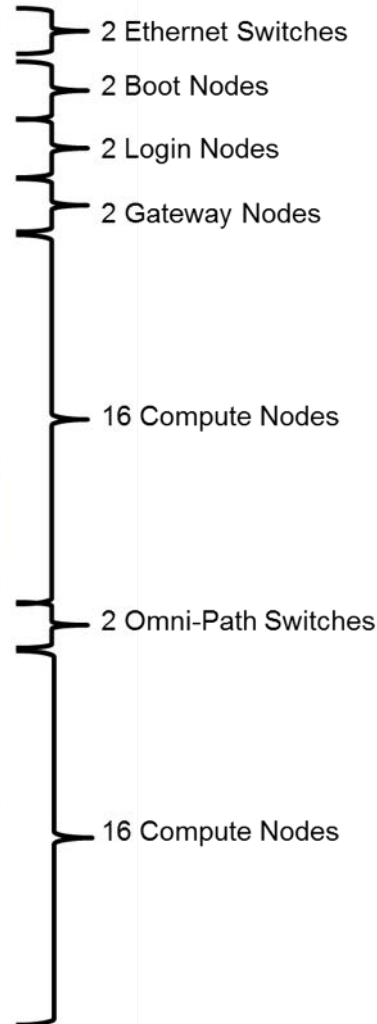- ARCTUR

# NEXTGenIO objectives

- Develop a new server architecture using next generation processor and memory advances
  - Based on Intel® Xeon and Intel DIMM based on 3D XPoint™ memory technology
- Investigate the best ways of utilising these technologies in HPC
  - Develop the systemware to support their use, particularly at scale
- Model different I/O workloads and use this understanding in a co-design process
  - Representative of real HPC centre workloads

# Hardware

- The project is developing a new HPC platform with focus on I/O performance
  - System and motherboard designed & built by Fujitsu

- The *Complete Compute Nodes* are based on
  - Intel™ CPUs - 2 sockets per node
  - Intel™ DIMMs
  - Intel™ OmniPath

# Prototype



2 Ethernet Switches
2 Boot Nodes
2 Login Nodes
2 Gateway Nodes

16 Compute Nodes

2 Omni-Path Switches

16 Compute Nodes

Note: final configuration may differ

# Intel™ DIMMs

- Non-volatile RAM
  - 3D XPoint technology
- Much larger capacity than DRAM
- Slower than DRAM by a small factor, but significantly faster than SSDs ™
- 12 DIMM slots per socket
  - Populated by combination of DDR4 and Intel™ DIMMs
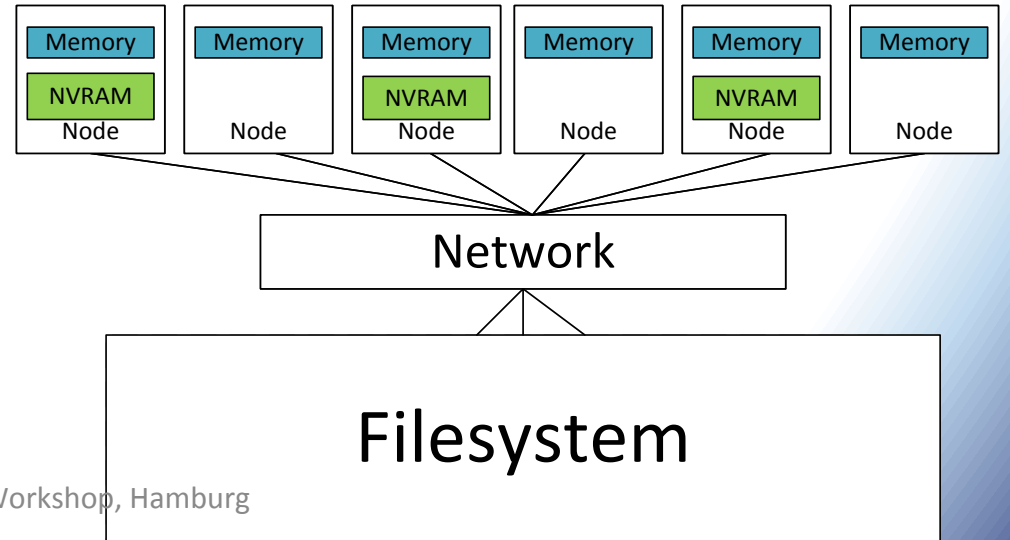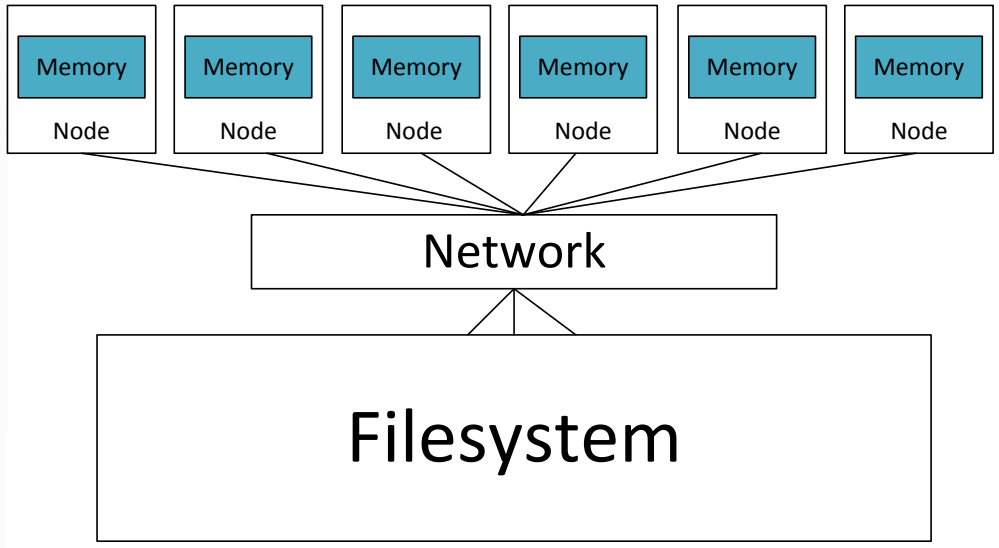
# Intel™ DIMMs – Usage Models

- The "memory" usage model allows for the extension of the main memory with Intel™ DIMMs
  - The data is volatile like normal DRAM based main memory

- The "storage" usage model which supports the use of Intel™ DIMMs like a classic block device
  - E.g. like a very fast SSD

- The "application direct" usage model maps persistent storage from the Intel™ DIMMs directly into the main memory address space
  - Direct CPU load/store instructions for persistent main memory regions

# Remote access

- Complete compute nodes and network hardware will support remote access to NVDIMMs from other CCNs
  - Using RDMA between nodes will allow data in NVDIMMs to be shared between CCNs if required by the applications using them
- Systemware will support remote access and use it for data partitioning and replication.
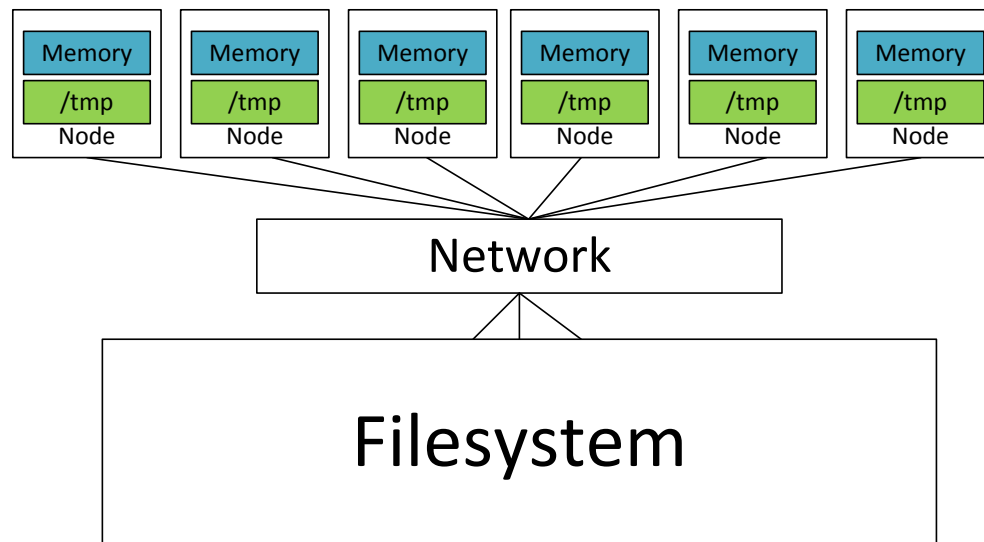
# Exploiting distributed storage

# Using distributed storage

- Without changing applications
  - Large memory space/in-memory database etc…
  - Local filesystem

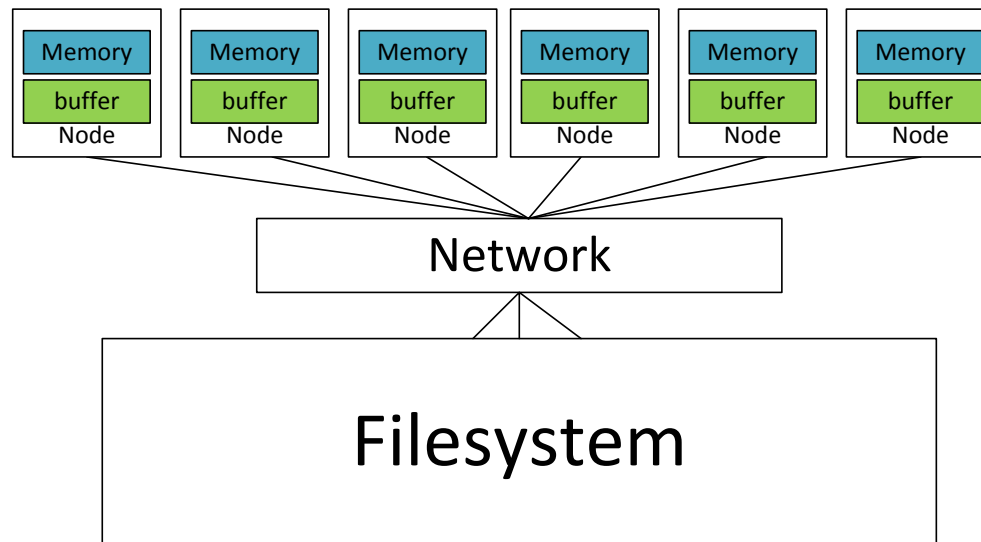| Memory | Memory | Memory | Memory | Memory | Memory |
|:------:|:------:|:------:|:------:|:------:|:------:|
| /tmp | /tmp | /tmp | /tmp | /tmp | /tmp |
| Node | Node | Node | Node | Node | Node |

Network

Filesystem

- Users manage data themselves
- No global data access/namespace, large number of files
- Still require global filesystem for persistence

# Using distributed storage

- Without changing applications
  - Filesystem buffer

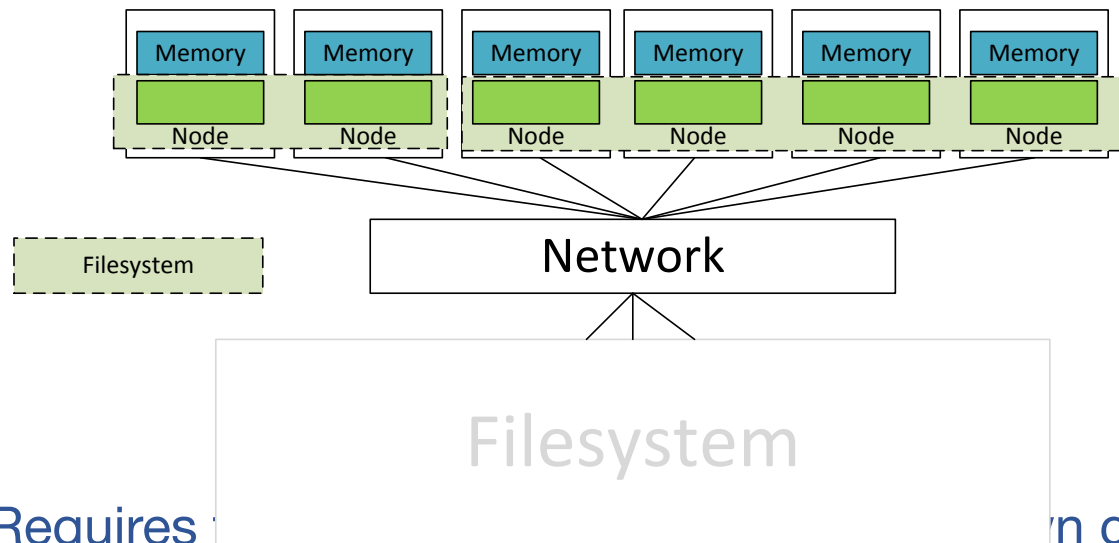| Memory | Memory | Memory | Memory | Memory | Memory |
|--------|--------|--------|--------|--------|--------|
| buffer | buffer | buffer | buffer | buffer | buffer |
| Node | Node | Node | Node | Node | Node |

**Network**

**Filesystem**

- Pre-load data into NVRAM from filesystem
- Use NVRAM for I/O and write data back to filesystem at the end
- Requires systemware to preload and postmove data
- Uses filesystem as namespace manager

# Using distributed storage
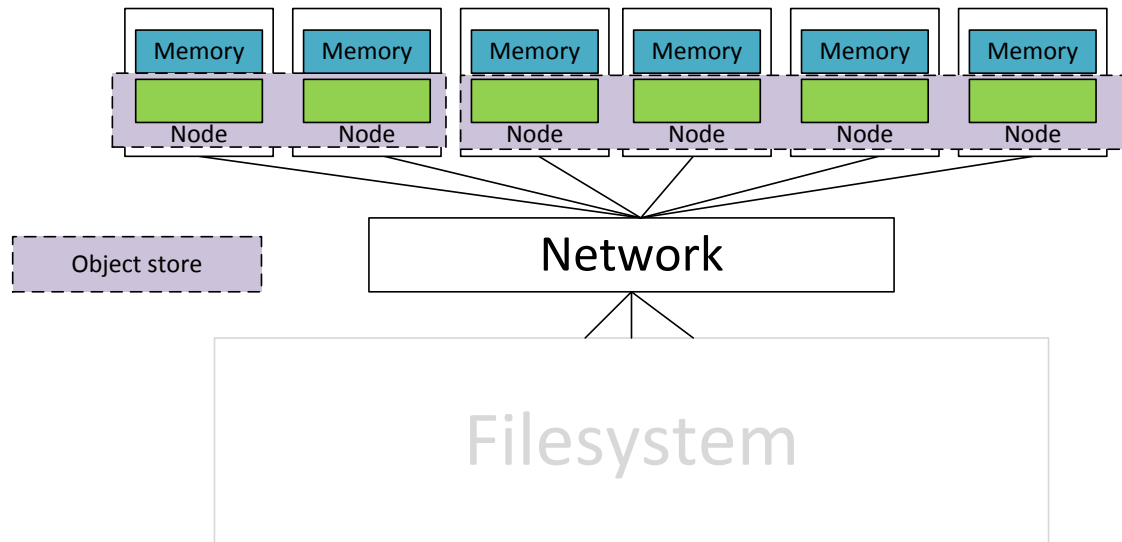
- Without changing applications
  - Global filesystem



- Requires ~~~~~~~~~~~~~~~~~~~~~~~~~n global filesystems for individual jobs
- Requires filesystem that works across nodes
- Requires functionality to preload and postmove filesystems
- Need to be able to support multiple filesystems across system

# Using distributed storage

- With changes to applications
  - Object store



- Needs same functionality as global filesystem
- Removes need for POSIX, or POSIX-like functionality
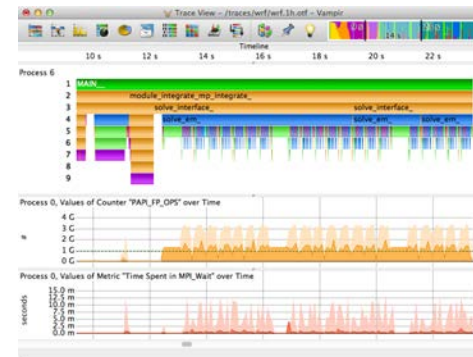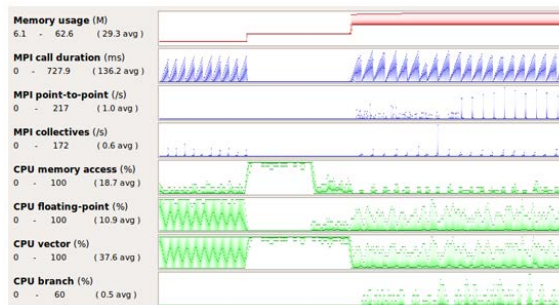
# The Challenge of distributed storage

- Enabling all the use cases in  multi-user, multi-job environment is the real challenge
    - Heterogeneous scheduling mix
    - Different requirements on the NVRAM
    - Scheduling across these resources
    - Enabling sharing of nodes
    - etc….
- Enabling applications to do more I/O
    - Large numbers of our applications don't heavily use I/O at the moment
    - What can we enable if I/O is significantly cheaper
- NEXTGenIO is tackling these
    - Job scheduler
    - Data scheduler
    - Data movers

# Tools effort

- Performance analysis tools need to understand new memory hierarchy and its impact on applications
    - TUD's Vampir & Allinea's MAP
- At the same time, tools themselves can exploit NVRAM to rapidly store sampling/tracing data

# IO workload simulation by ECMWF

- Need to quantify improvements in job runtime and throughput
  - Measure and understand current bottlenecks
- Create a workload simulator and generator
  - Simulator can be used to derive system configuration options
  - Generator can be used to create scaled down version of data centre workload

# Performance Features

- Monitoring NVRAM resources
- Recording File I/O operations
- Providing memory access statistics
- Adding system topology information
- Merging perf. data (workflow support)

- Hence, extensions are required in:
  - Data formats
  - Measurement infrastructure
  - Visualization

# Functional extensions in measurement infrastructure

# Monitoring extensions

1. File I/O operations
2. Hardware/software counters
   1. Non-Volatile Memory Library (NVML) provides information on NVRAM allocation
   2. Counters enable memory statistics on NVRAM
3. System topology information related to NVRAM

**Figure 1 : Metric Architecture with the Memory Access Metric extension. Differences between TUD and Allinea**

# NVM Library Wrapper

- Enter and leave events
- Additional information being recorded
  - Requested memory size
  - Usable memory size
  - High Water Mark metric for the utilization of memory pool over its entire execution
  - Size and number of elements available in the persistent array
- Byte access activities to/from NVRAM remain out of scope (e.g. memory mapped files)
- NVRAM health status
  - Intel management API access (system configuration)
  - S.M.A.R.T. (on node level)
  - Intel NVDIM_API

# Memory Access Statistics

- Memory access hotspots for using DRAM and NVRAM?
  - Where? When? Type of memory?

- Metric collection needs to be extended
  1. DRAM local access
  2. DRAM remote access (on a different socket)
  3. NVRAM local access
  4. NVRAM remote access (on a different socket)

# Access to PMU using perf

- Architectural independent counters
  - May introduce some overhead
    - MEM_TRANS_RETIRED.LOAD_LATENCY
    - MEM_TRANS_RETIRED.PRECISE_STORE
    - Guess: It will also work for NVRAM?

- Architectural dependent counters
  - Counter for DRAM
    - MEM_LOAD_UOPS_L3_MISS_RETIRED.REMOTE_DRAM
    - MEM_LOAD_UOPS_L3_MISS_RETIRED.LOCAL_DRAM
    - MEM_LOAD_UOPS_*.REMOTE_NVRAM ?
    - MEM_LOAD_UOPS_*.LOCAL_NVRAM ?

# Information of system topology in compute node

- Following information will be incorporated in the system topology of each compute node:
    1. Using procfs:
        - Total size of NVRAM
        - Total size of DRAM
        - Total processors
    2. Using sysfs for:
        - Namespaces for used NVDIMMs
    3. Using libnuma / numactl:
        - Total number of NUMA nodes
        - Memory size of each NUMA node
        - Processing Cores allocated to each NUMA node
        - NUMA distances
    4. SLURM API will be used to provide meta information of following records:
        - Job memory size of NVRAM
        - Job memory size of DRAM
        - Number of compute nodes
        - Number of compute processors

**Figure 2 : Extension of Score-P infrastructure to incorporate system topology in OTF2 Trace**

# Functional Extensions in Visualization (Vampir)

# Display stacked I/O layers

- I/O layers
  - POSIX
  - MPI-I/O
  - HDF5
  - NetCDF
  - PNetCDF
  - Adios
  - (Lustre)

- Data of interest
  - File Open/Create Operations
  - File Close Operations
  - Data Transfer operations

# I/O operations over time



Individual I/O Operation

I/O Runtime Contribution

# I/O data rate over time



I/O Data Rate of single thread

# I/O summaries with totals

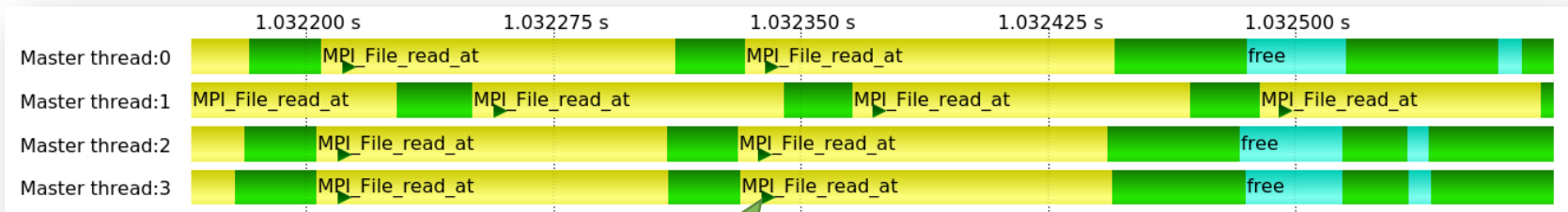All Processes, Aggregated I/O Transaction Size per Operation Type

|  5.0 MiB | 3.8 MiB | 2.5 MiB | 1.2 MiB | 0B | |
|---|---|---|---|---|---|
| 6 MiB | | | | | Sum |
| | | 3 MiB | | | WRITE |
| | | 3 MiB | | | READ |

**Other Metrics:**
- IOPS
- I/O Time
- I/O Size
- I/O Bandwidth

All Processes, Average I/O Bandwidth per Operation Type

| 1,920 MiB/s | 1,280 MiB/s | 640 MiB/s | 0B/s | |
|---|---|---|---|---|
| 2.315 GiB/s | | | | READ |
| | | 153.514 MiB/s | | WRITE |

# I/O summaries per file



All Processes, Aggregated I/O Transaction Time per File Name

| Time | File Name |
|---|---|
| 5.232 s | /N/dc2/scratch/wardmod/input/fineGrid/frate.wheat.nc |
| 3.628 s | /N/dc2/scratch/wardmod/input/fineGrid/frate.corn.nc |
| 2.503 s | /N/dc2/scratch/wardmod/input/fineGrid/frate.soy.nc |
| 2.027 s | /N/dc2/scratch/wardmod/inp...dx/daily/Intr_relh_1948.nc |
| 1.986 s | /N/dc2/scratch/wardmod/inp...dx/daily/Intr_rads_1948.nc |
| 1.359 s | /N/dc2/scratch/wardmod/in...x/daily/Intr_wspd_1948.nc |
| 1.182 s | /N/dc2/scratch/wardmod/i...C_60year_climo_mm_day.nc |
| 0.942 s | /tmp/mpi/proc6/output/hourly/1948_hourly4thmb.nc |
| 0.798 s | /N/dc2/scratch/wardmod/inp...dx/daily/Intr_tmin_1948.nc |
| 0.768 s | /N/dc2/scratch/wardmod/inp...dx/daily/Intr_tmax_1948.nc |
| 0.704 s | /N/dc2/scratch/wardmod/i...ntr_RADS_60year_climo.nc |
| 0.694 s | /N/dc2/scratch/wardmod/i...ntr_WSPD_60year_climo.nc |
| 0.59 s | /N/dc2/scratch/wardmod/i...ntr_TMAX_60year_climo.nc |

# Missing: timeline view per file

# I/O operations per thread or per file
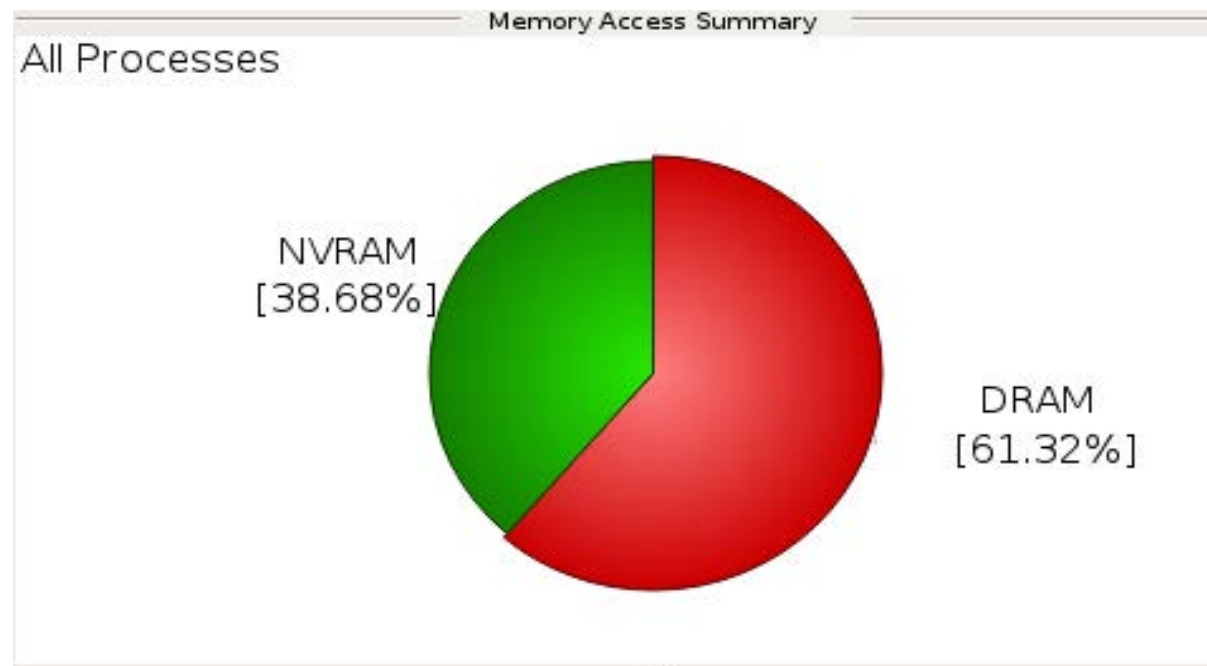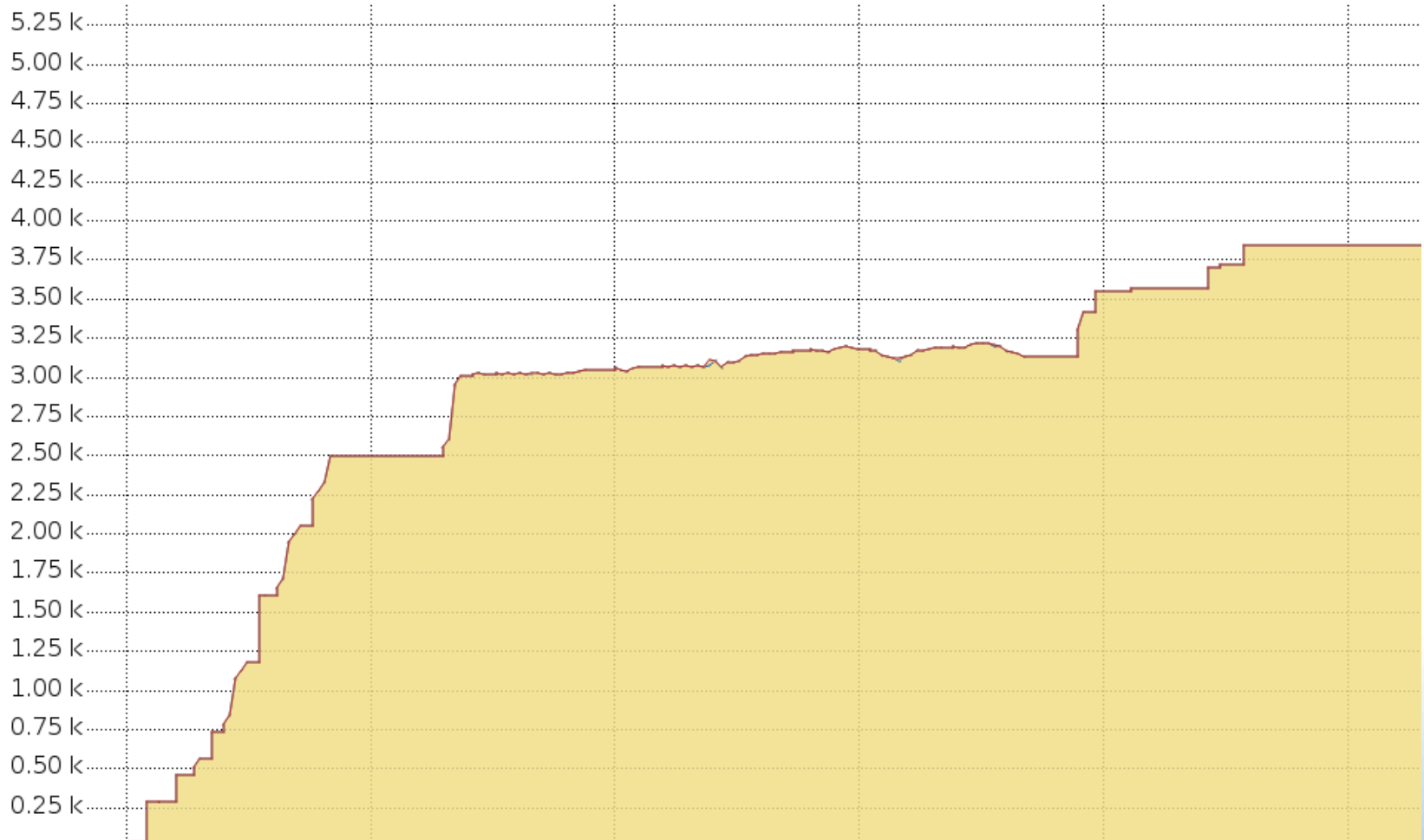
# I/O comparison example

# Visualization of memory access statistics

- Currently, Vampir has many different views presenting counter metrics

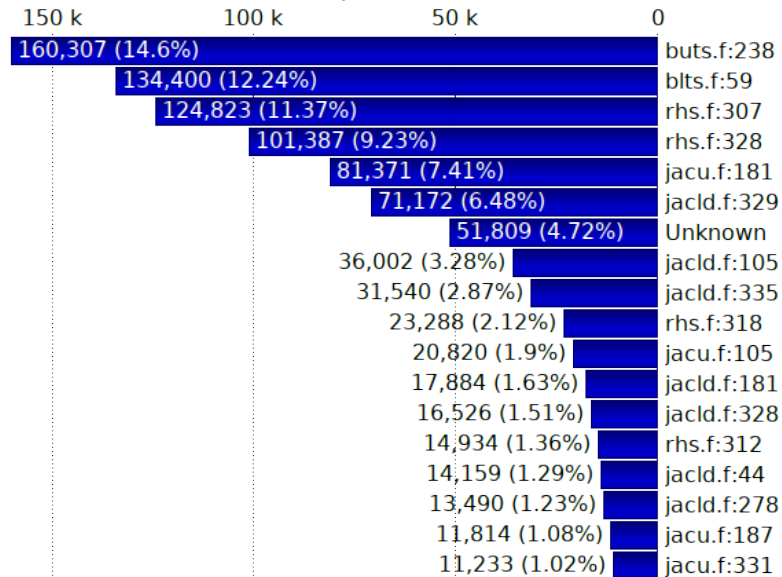- Future: Absolute number of memory accesses performed by an application for different types of memory



Memory Access Summary

All Processes

NVRAM [38.68%]

DRAM [61.32%]

# NVRAM allocation over time

# NVRAM accesses over time?

Counters?

# Future Stats



All Processes, Number of Hits per Source Code Location

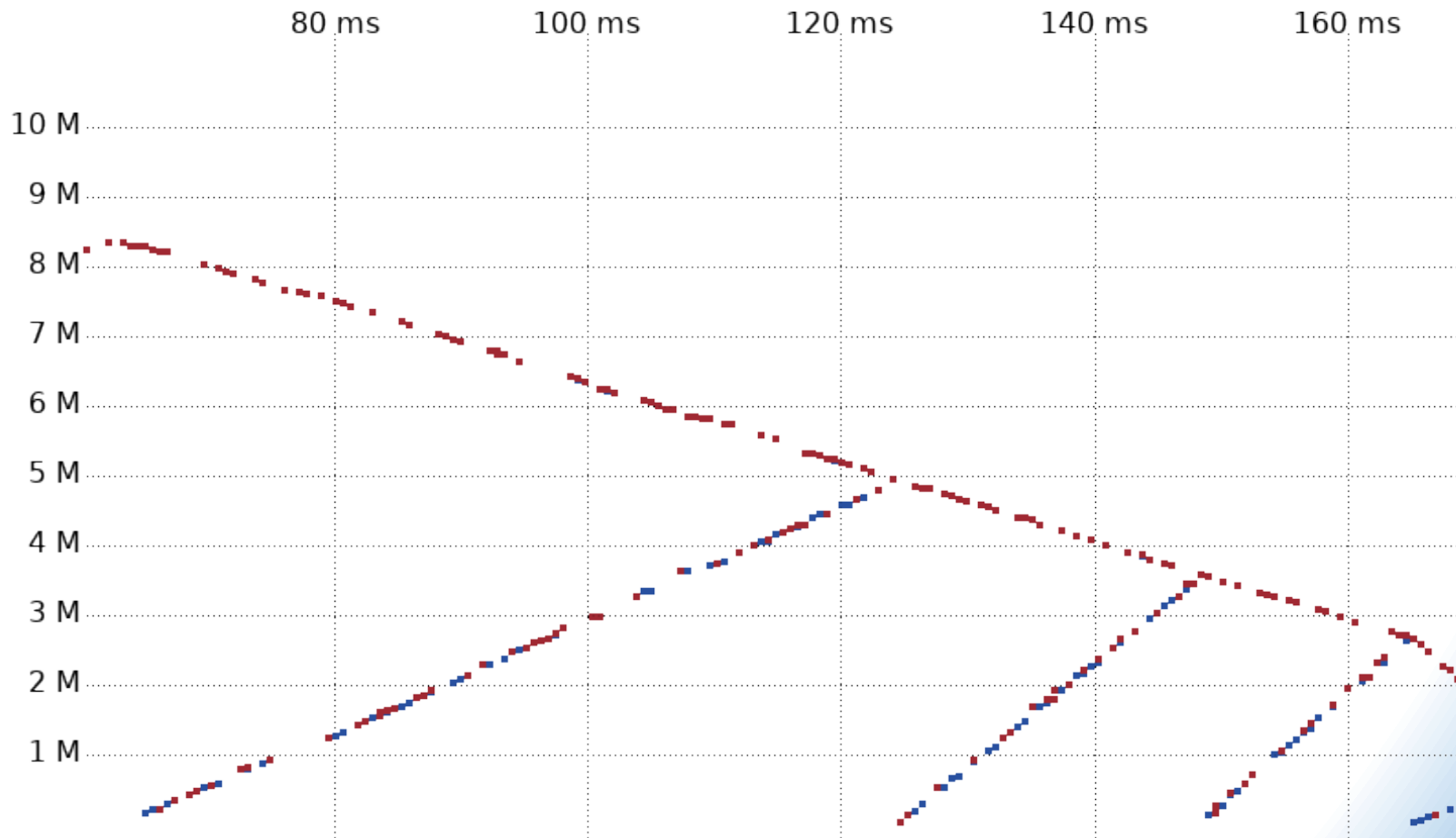| 150 k | 100 k | 50 k | 0 | |
|---|---|---|---|---|
| 160,307 (14.6%) | | | | buts.f:238 |
| 134,400 (12.24%) | | | | blts.f:59 |
| 124,823 (11.37%) | | | | rhs.f:307 |
| 101,387 (9.23%) | | | | rhs.f:328 |
| 81,371 (7.41%) | | | | jacu.f:181 |
| 71,172 (6.48%) | | | | jacld.f:329 |
| 51,809 (4.72%) | | | | Unknown |
| 36,002 (3.28%) | | | | jacld.f:105 |
| 31,540 (2.87%) | | | | jacld.f:335 |
| 23,288 (2.12%) | | | | rhs.f:318 |
| 20,820 (1.9%) | | | | jacu.f:105 |
| 17,884 (1.63%) | | | | jacld.f:181 |
| 16,526 (1.51%) | | | | jacld.f:328 |
| 14,934 (1.36%) | | | | rhs.f:312 |
| 14,159 (1.29%) | | | | jacld.f:44 |
| 13,490 (1.23%) | | | | jacld.f:278 |
| 11,814 (1.08%) | | | | jacu.f:187 |
| 11,233 (1.02%) | | | | jacu.f:331 |

All Processes, Number of Hits per Source Code Location

| 3 k | 2 k | 1 k | 0 | |
|---|---|---|---|---|
| 3,173 (14.81%) | | | | buts.f:238 |
| 2,707 (12.63%) | | | | blts.f:59 |
| 2,411 (11.25%) | | | | rhs.f:307 |
| 1,884 (8.79%) | | | | rhs.f:328 |
| 1,611 (7.52%) | | | | jacu.f:181 |
| 1,417 (6.61%) | | | | jacld.f:329 |
| 990 (4.62%) | | | | Unknown |
| 663 (3.09%) | | | | jacld.f:105 |
| 643 (3%) | | | | jacld.f:335 |
| 411 (1.92%) | | | | rhs.f:318 |
| 396 (1.85%) | | | | jacu.f:105 |
| 351 (1.64%) | | | | jacld.f:181 |
| 324 (1.51%) | | | | jacld.f:328 |
| 302 (1.41%) | | | | jacld.f:278 |
| 281 (1.31%) | | | | rhs.f:312 |
| 267 (1.25%) | | | | jacld.f:44 |
| 254 (1.19%) | | | | jacu.f:187 |
| 230 (1.07%) | | | | jacu.f:331 |

# Future Stats



hread, Values of Metric "long int* vector" over Time

# Summary

- NEXTGenIO developing a ***full*** hardware and software solution

- Requirements capture and first architectural designs completed
  - Hardware under development
  - Systemware under development

- Potential to both significantly reduce I/O costs and enable new usage models for HPC and HPDA
  - Proper convergence between HPC and HPDA