# Understanding and Improving I/O Performance of HPC systems
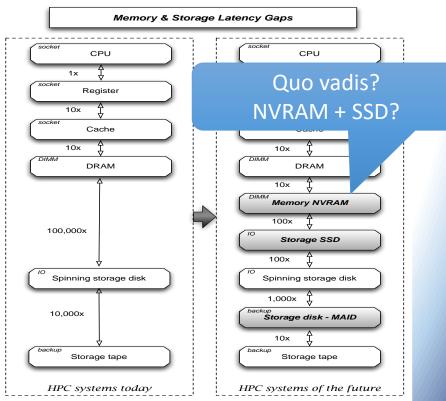
Muhammad Sarim Zafar (sarim.zafar@tu-dresden.de)

ZIH, Technische Universität Dresden

ISC17, Frankfurt, June 18th 2017

# New Members in Memory Hierarchy

- New memory technology

- Changes the memory hierarchy we have

- Impact on applications e.g. simulations?

- I/O performance is one of the critical components for scaling applications



Memory & Storage Latency Gaps

Quo vadis?
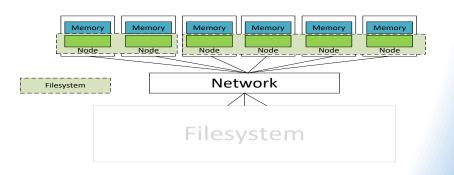NVRAM + SSD?

# Three Usage Models

- The "memory" usage model
  - Extension of the main memory
  - Data is *volatile* like normal main memory

- The "storage" usage model
  - Classic *persistent* block device
  - Like a very fast SSD

- The "application direct" usage model
  - Maps *persistent* storage into address space
  - Direct CPU load/store instructions
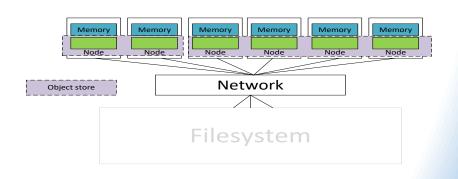
# Using Distributed Storage

- Distributed global file system
  - No changes to apps
  - Support for NVRAM

- Required functionality
  - Preload and post load filesystems
  - Create and share files amongst the jobs
  - Works across nodes
  - Support multiple filesystems across cluster

- I/O Performance
  - Sum of memory layers available
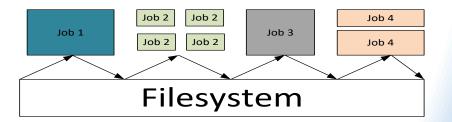
# Using an Object Store

- Needs changes at the application layer
  - Needs same functionality as global filesystem
- I/O Performance
  - Mapping to objects
  - Different type of abstraction (APIs etc.)
  - Different kind of Instrumentation

# On the larger perspective

- Towards Workflows
  - Workflows
    - Producer-consumer model
  - Resident data sets
    - Sharing preloaded data across a range of jobs
    - Data analytic workflows
      - E.g. Total IO operations performed by the workflow
      - E.g. Number of IO operations performed on DRAM and NVRAM
  - I/O Performance
    - Data merging/integration between the jobs
    - Remove file system from intermediate stages

# In Context of: Vampir & Score-P

# New Tool Key Objectives

- Analysis tools need to
  - Reveal performance interdependencies in I/O and memory hierarchy
  - Support workflow visualization(Workload modelling)
  - Exploit NVRAM to store data for themselves

# Exploiting I/O Layers

- I/O layers
  - E.g. Distributed FS
    - Client side – Easy to access
    - Server side – Complex
  - Kernel – IO : rusage etc
  - POSIX – fopen etc
  - MPI-I/O
  - HDF5
  - NetCDF
  - PNetCDF

# I/O performance data

- Event data AND aggregated numbers
  - Open/Create/Close operations (meta data)
  - Data transfer operations

- Client side instrumentation
  - procfs via PAPI-Components counters, Score-P (user space)

- Server side instrumentation (FS dependent)
  - Data base (custom made, system space)
  - Mapping to nodes/applications

# In the presence of NVMs…

- Information
    - Memory size (requested, usable)
    - High Water Mark metric
    - Size and number of elements in memory
- Individual NVRAM load/stores
    - Remain out of scope (e.g. memory mapped files)
- NVRAM health status
    - Not measurable at high frequencies

# I/O Operations over Time



Individual I/O Operation

I/O Runtime Contribution

All Processes, Accumulated Exclusive Time per Function Group

| | | |
|---|---|---|
| 87.42% | | MPI |
| 11.44% | | MPI-IO |
| 0.44% | | MPI-Collectives |
| 0.4% | | POSIX_IO_API |
| 0.26% | | Application |
| 0.02% | | MEMORY_ALLOCATE |
| 0.01% | | MEMORY_DEALLOCATE |

# 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 |

Scale marks: 4.5 s, 3.0 s, 1.5 s, 0s

Aggregated data for specific resource

# I/O Operations per File



I/O Events of File "/lustre/scratch2/cherold/ior_test/testFile" over Time

| | 0.860 s | 0.862 s | 0.864 s | 0.866 s | 0.868 s |

Master thread:0 — WRITE ... WRITE WRITE WRITE
Master thread:1 — WRITE ... WRITE ... WRITE WRITE
Master thread:2 — WRITE WRITE ... WRITE WRITE
Master thread:3 — WRITE ... WRITE ... WRITE WRITE

Focus on specific resource

I/O Events over Time

| | 1.09 s | 1.12 s | 1.15 s | 1.18 s | 1.21 s | 1.24 s |

/scratch/cherold/ior_test/testFile
all
stderr
stdin
stdout
wp3.ior
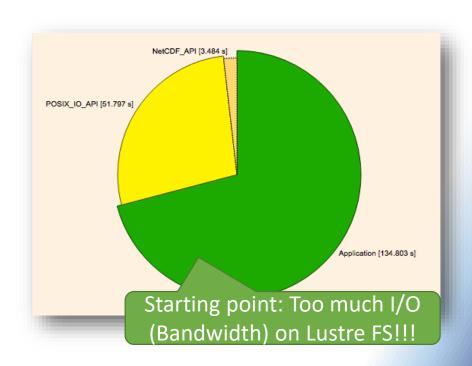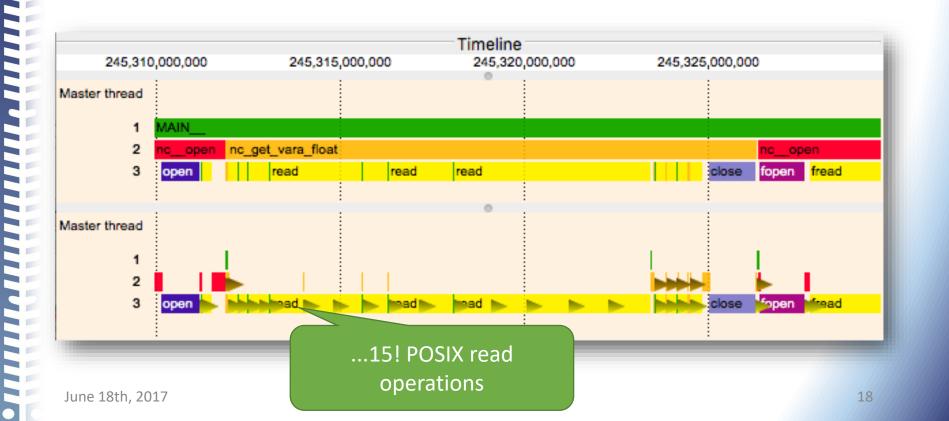
Show all resources

# Example

- Bringing the Lustre system I/O down
  - with a single (serial) application
- Higher I/O demand than IOR benchmark
- Why?



Starting point: Too much I/O (Bandwidth) on Lustre FS!!!

# Details Make a Difference



...15! POSIX read operations

# Coarse Grained Time Series Reveal Some Clue, but...



Single Thread Lustre I/O Bandwidth over time

# Impact of Lustre Prefetching

# Future Stats



All Processes, Number of Hits per Source Code Location

| | | |
|---|---|---|
| 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

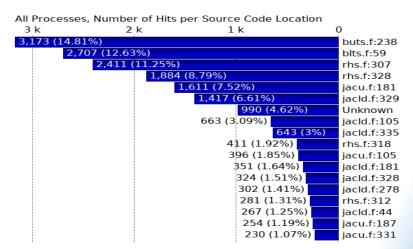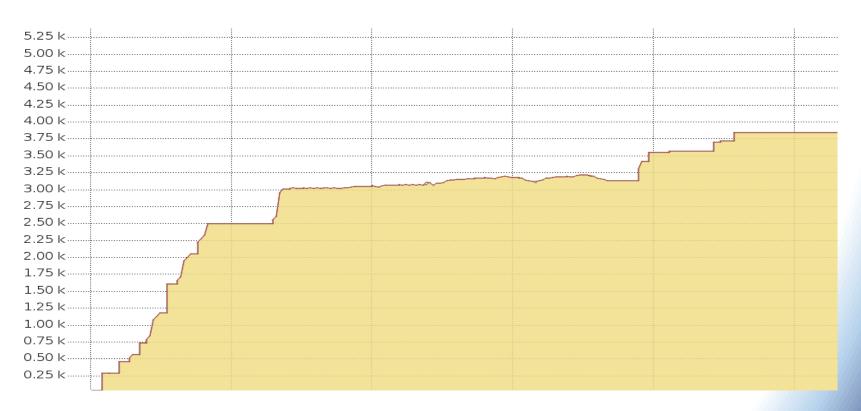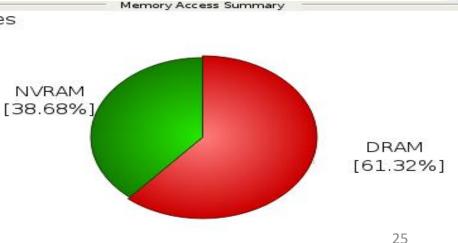| | | |
|---|---|---|
| 3,173 (14.81%) | | buts.f:238 |
| 2,707 (12.63%) | | blts.f:59 |
| 2,411 (11.25%) | | rhs.f:307 |
| 1,884 (8.79%) | | rhs.f:328 |
| 1,611 (7.52%) | | jacu.f:181 |
| 1,417 (6.61%) | | jacld.f:329 |
| 990 (4.62%) | | Unknown |
| 663 (3.09%) | | jacld.f:105 |
| 643 (3%) | | jacld.f:335 |
| 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 |

# NVRAM allocation over time

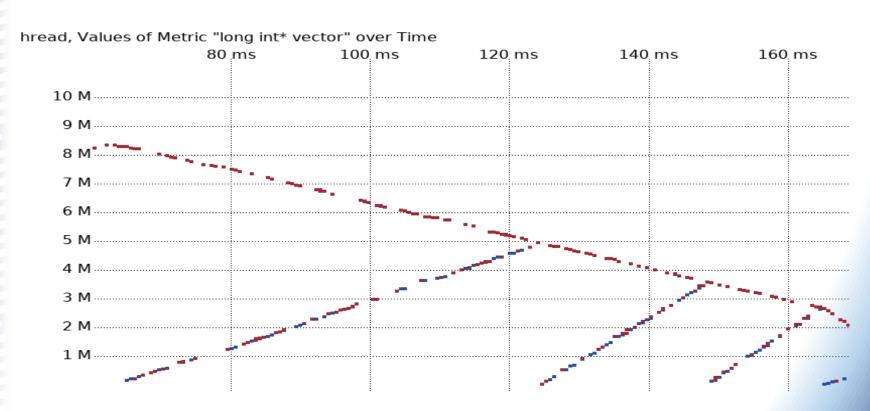# Visualization of memory access statistics

- 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%]

# Future Stats



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

# Summary

- NEXTGenIO developing a full hardware **and** software solution
  - Solution includes Lustre
- Performance focus
  - Consider complete I/O stack
  - Incorporate new I/O paradigms
  - Study implications of NVRAM

- Reduce I/O costs
- New usage models for HPC and HPDA

# Hands-on Slides

# Objective

- Familiarize with the usage oge Score-P/Vampir tool suite

- Execises are based on a small portable IO benchmarking application.

# First stage: Instrument Benchio

```
# Connect
% ssh –X userid@abc.edu
```

```
# make separate dir
% mkdir scorep-vampir

# get the sources
% wget
http://www.archer.ac.uk/training/course-material/2016/09/160929_AdvMPI_EPCC/benchio.tar

# untar
```

# Preparation

```
# Load modules
% module load cray-netcdf-hdf5parallel/4.4.0
% module load cray-hdf5-parallel/1.8.16
% module use /work/d131/d131/shared/vampir/modules
% module load vampir/io-dev
% module load scorep/cray
```

# Building the benchmark

```
% vi Makefile
MF=      Makefile
# ARCHER
FC= ftn
FFLAGS=
LFLAGS=
# Ubuntu
#FC=      mpif90
#FFLAGS=-I/usr/include
#LFLAGS=-L/usr/lib -lnetcdff -L/usr/lib/x86_64-linux-gnu/ -lhdf5_fortran
EXE=     benchio
SRC= \
         benchio.f90 \
         mpiio.f90 \
         netcdf.f90 \
         hdf5.f90 \
         benchclock.f90
```

Clean and make, without instrumentation

Modify specification of compiler / linker

Specify the MPI paradigm for CRAY

Once finised, clean and make

```
% vi          M   ile
MF=          M   ile
# ARCHER
FC= scorep --mpp=mpi ftn
FFLAGS=
LFLAGS=
# Ubuntu
#FC=     mpif90
#FFLAGS=-I/usr/include
#LFLAGS=-L/usr/lib -lnetcdff -L/usr/lib/x86_64-linux-gnu/ -lhdf5_fortran
EXE=     benchio
SRC= \
         benchio.f90 \
         mpiio.f90 \
         netcdf.f90 \
         hdf5.f90 \
         benchclock.f90
```

# Job submission

```
# By default, only profiling is enabled. Enable the tracing in
# job script benchio.pbs
#!/bin/bash --login

#PBS -A d131
#PBS -N benchio
#PBS -l select=3
#PBS -l walltime=00:05:00
#PBS -j oe

# change directory to where the job was submitted from

export SCOREP_ENABLE_TRACING=true

…
```

```
# Visualize Trace
% cd scorep-<traceid>
% vampir traces.otf2
```