# Solid State Disks Testing with PROOF
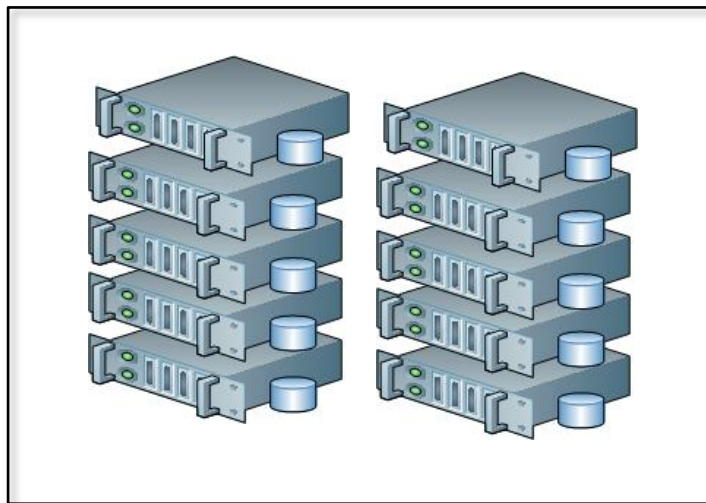
**Sergey Panitkin, Robert Petkus, Ofer Rind**

**BNL**

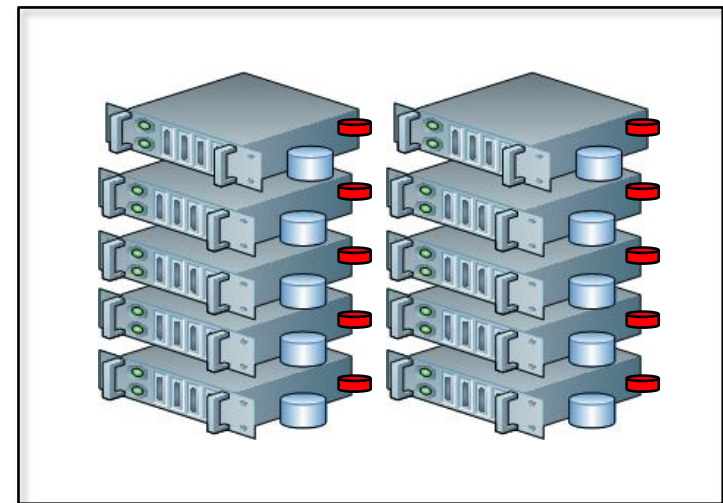# Current BNL PROOF Farm Configuration

"Old farm"-Production

- 10 nodes – 4 GB RAM each
- 40 cores: 1.8 GHz Opterons
- 20 TB of HDD space (10x4x500 GB)

"New Farm" – test site

- 10 nodes - 16 GB RAM each
- 80 cores: 2.0 GHz Kentsfields
- 5 TB of HDD space (10x500 GB)
- 640 GB SSD space (10x64 GB)

# New Solid State Disks@ BNL

- Model: Mtron  MSP-SATA7035064
- Capacity 64 GB
- Average access time ~0.1 ms (typical HD ~10ms)
- Sustained read ~120MB/s
- Sustained write ~80 MB/s
- IOPS (Sequential/ Random)  81,000/18,000
- Write endurance  >140 years @ 50GB write per day
- MTBF  1,000,000 hours
- 7-bit Error Correction Code

# Test configuration

- 1+1 or 1+8 nodes PROOF farm configurations

- 2x4 core Kentsfield CPUs per node, 16 GB RAM per node

- All default settings in software and OS

- Different configuration of SSD and HDD hardware depending on tests

- Root 5.18.00 – latest production version

- "PROOF Bench" suit of benchmark scripts to simulate analysis in root. Part of root distribution.

    - http://root.cern.ch/twiki/bin/view/ROOT/ProofBench

    - Data simulate HEP events ~1k per event

    - Single  ~3+ GB file per PROOF worker in this tests

- Reboot before every test to avoid memory caching effects

- This set of tests emulates interactive, command prompt root session

    - Plot one variable, scan ~10E7 events, ala D3PD analysis

- Looking at read performance of I/O subsystem

Typical test session in root

# SSD vs HDD



CPU limited

Read rate vs number of PROOF workers per node

- SSD 1 disk
- HDD 1 disk

Single variable scan. Single node

Rate, MB/s

Number of workers
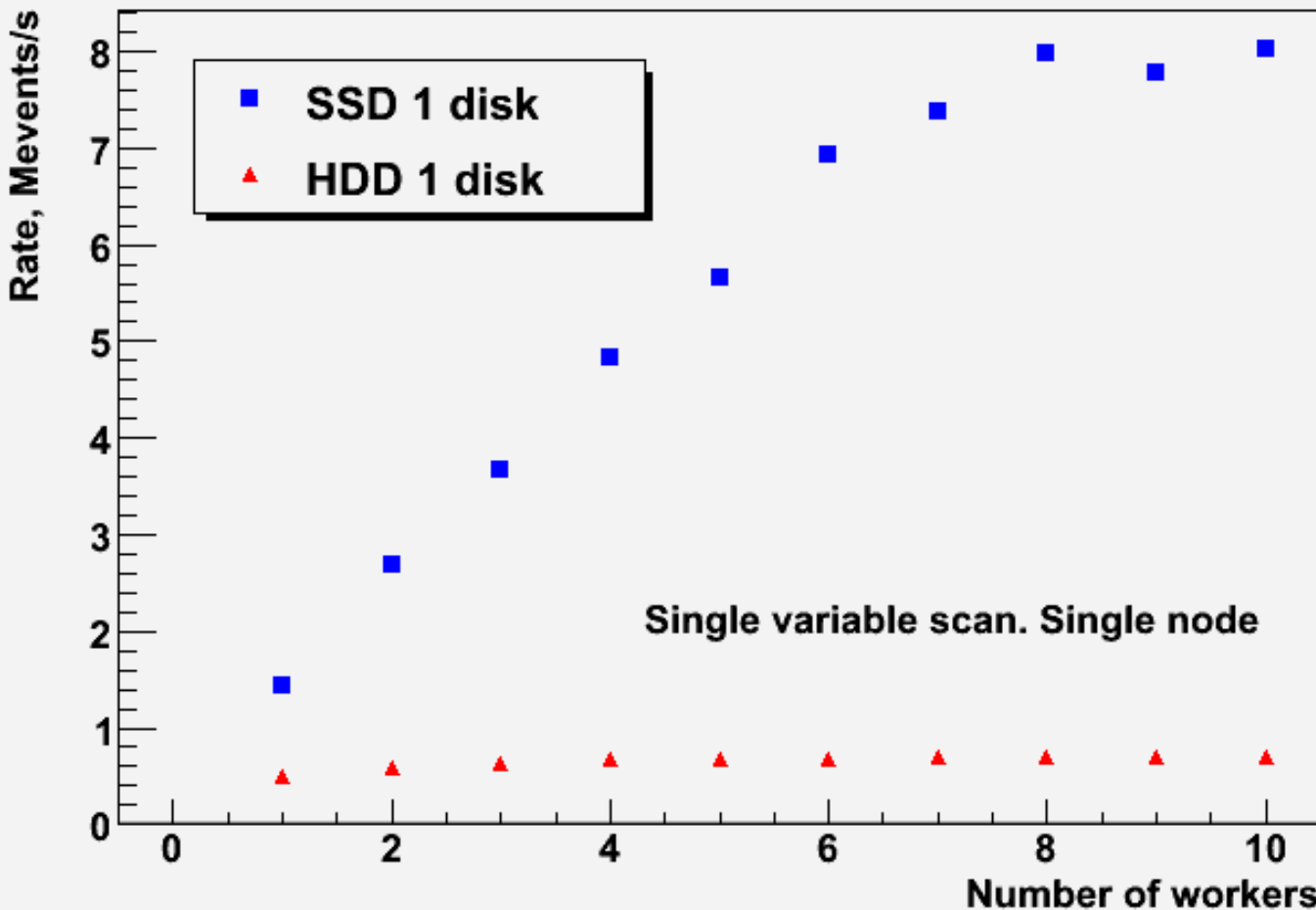
- ➢ SSD holds clear speed advantage
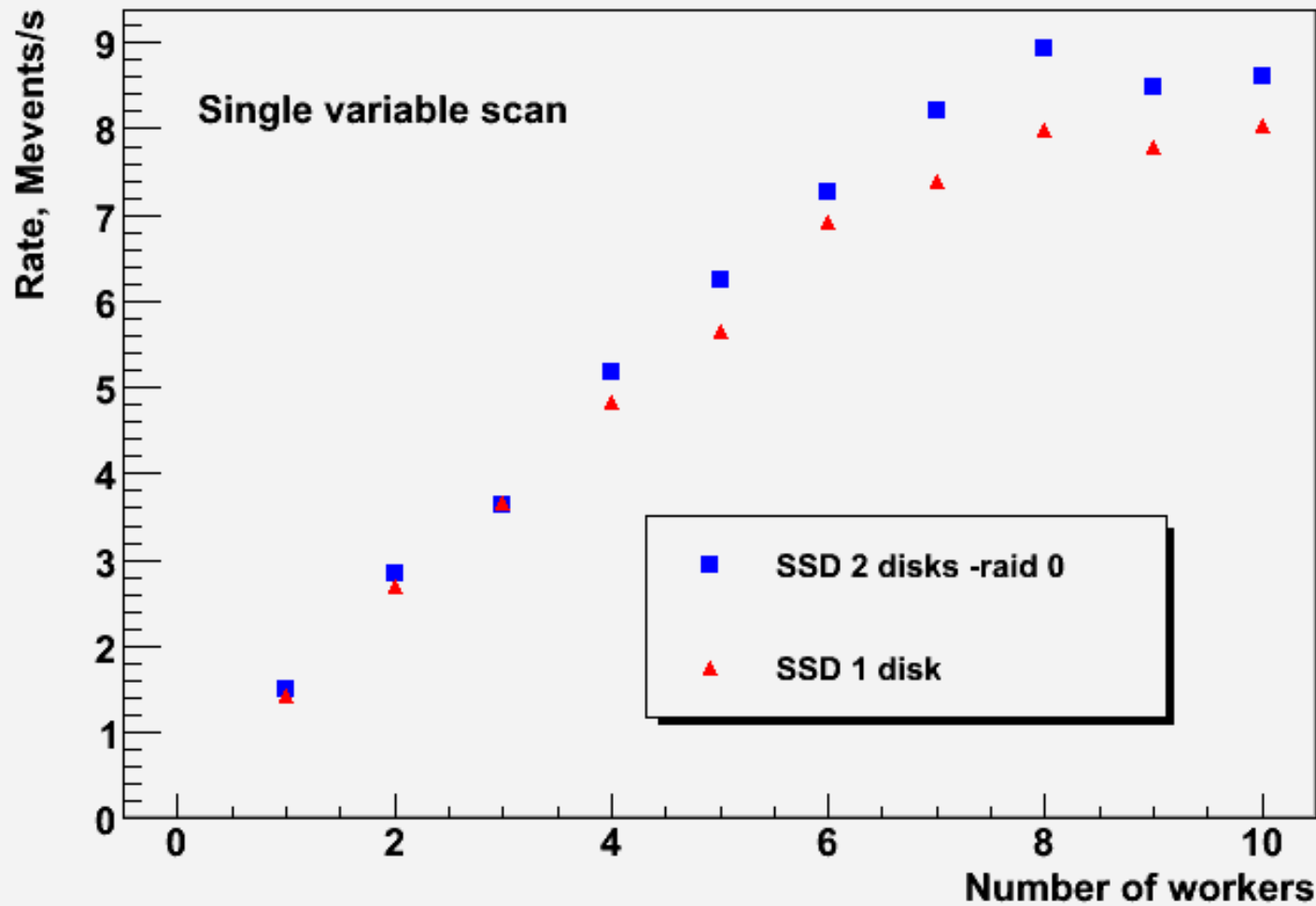- ➢ ~ 10 times faster in concurrent read scenario

**Analysis rate vs number of PROOF workers per node**

With 1 worker :    5.3M events,   15.8 MB read out of ~3 GB of data on disk
With 8 workers: 42.5M events,  126.5 MB read out of ~24 GB of data

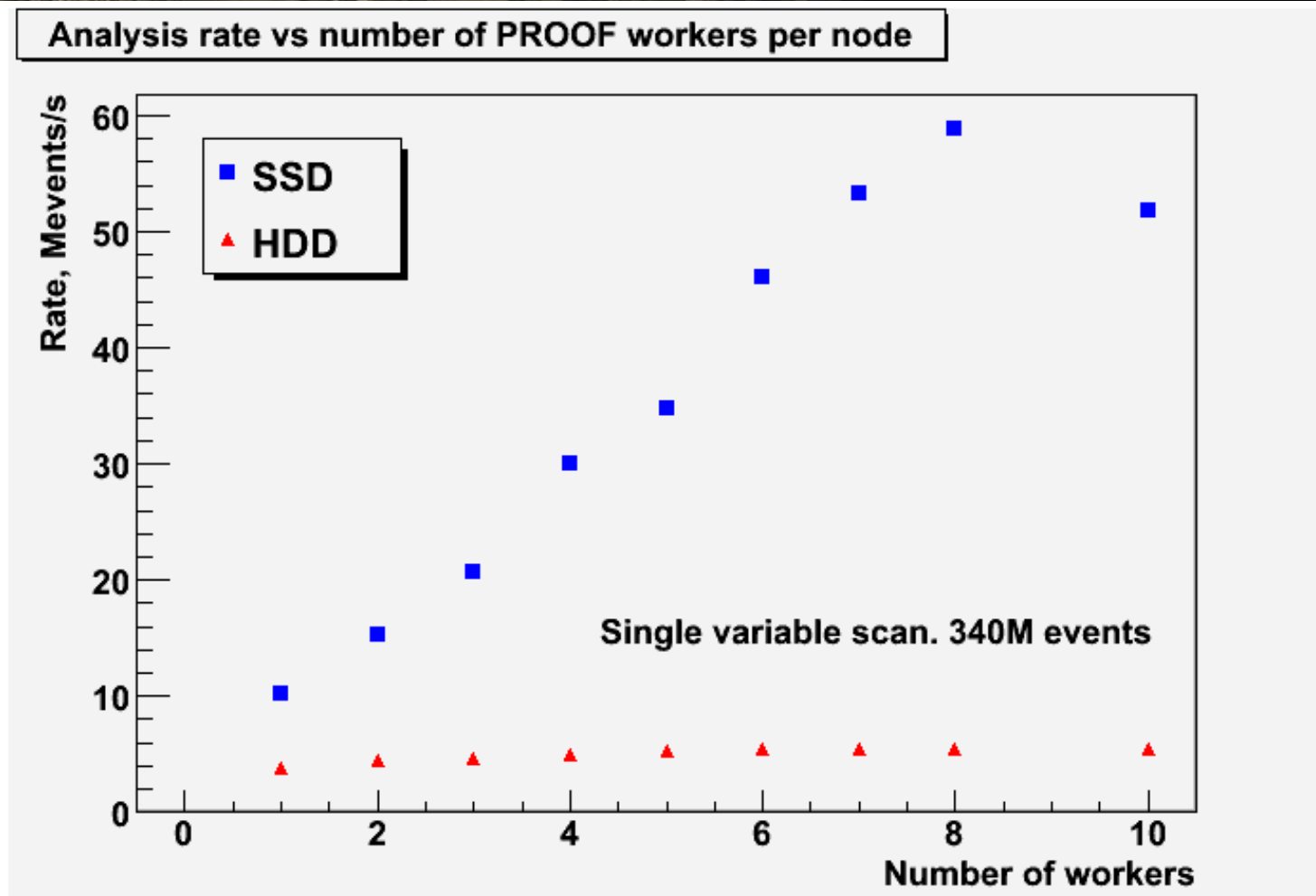# SSD vs HDD. 8 node farm



Analysis rate vs number of PROOF workers per node

Single variable scan. 340M events

Aggregate (8 node farm) analysis rate as a function of number of workers per node

Almost linear scaling with number of nodes

Sergey Panitkin

# Summary

- SSD offer significant performance advantage in concurrent analysis environment

- ~x10 better read performance than HD in our test

- More results will be shown tomorrow

- More tests are planned

    - ARA on AODs and DPDs

    - Different hardware configurations