# Power Management PCIe/NVMe Flash Cards

## Presenter: Bob Weisickle

CEO/Founder: OakGate Technololgy, Inc.

**OakGate Technology**

## Host

| Application |
| --- |
| Operating System |
| PCIe Storage Driver |
| HBA/RAID Card (SAS/SATA) |

**Interface (6G/s SATA or 12G/s SAS)**

| SAS/SATA Interface |
| --- |
| Flash Controller |
| Flash Memory |

**SSD**

**Interface (x4/x8 PCIe Gen-2/3)**

## Host

| Application |
| --- |
| Operating System |
| PCIe Storage Driver |
| Flash Controller |
| Flash Memory |

**SSD**

Improved:
o Complexity
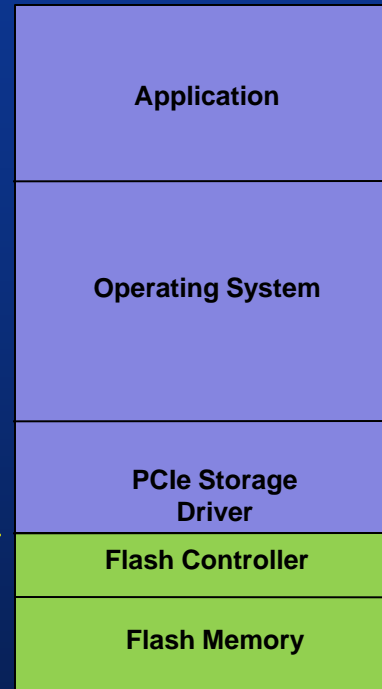o Latency
o Performance
o Reliability

# Performance/Validation Challenge

**Host**

## Requirements

- Access to PCIe Registers
- Queue Configuration
- Hot Discovery
- Error/Fault Injection
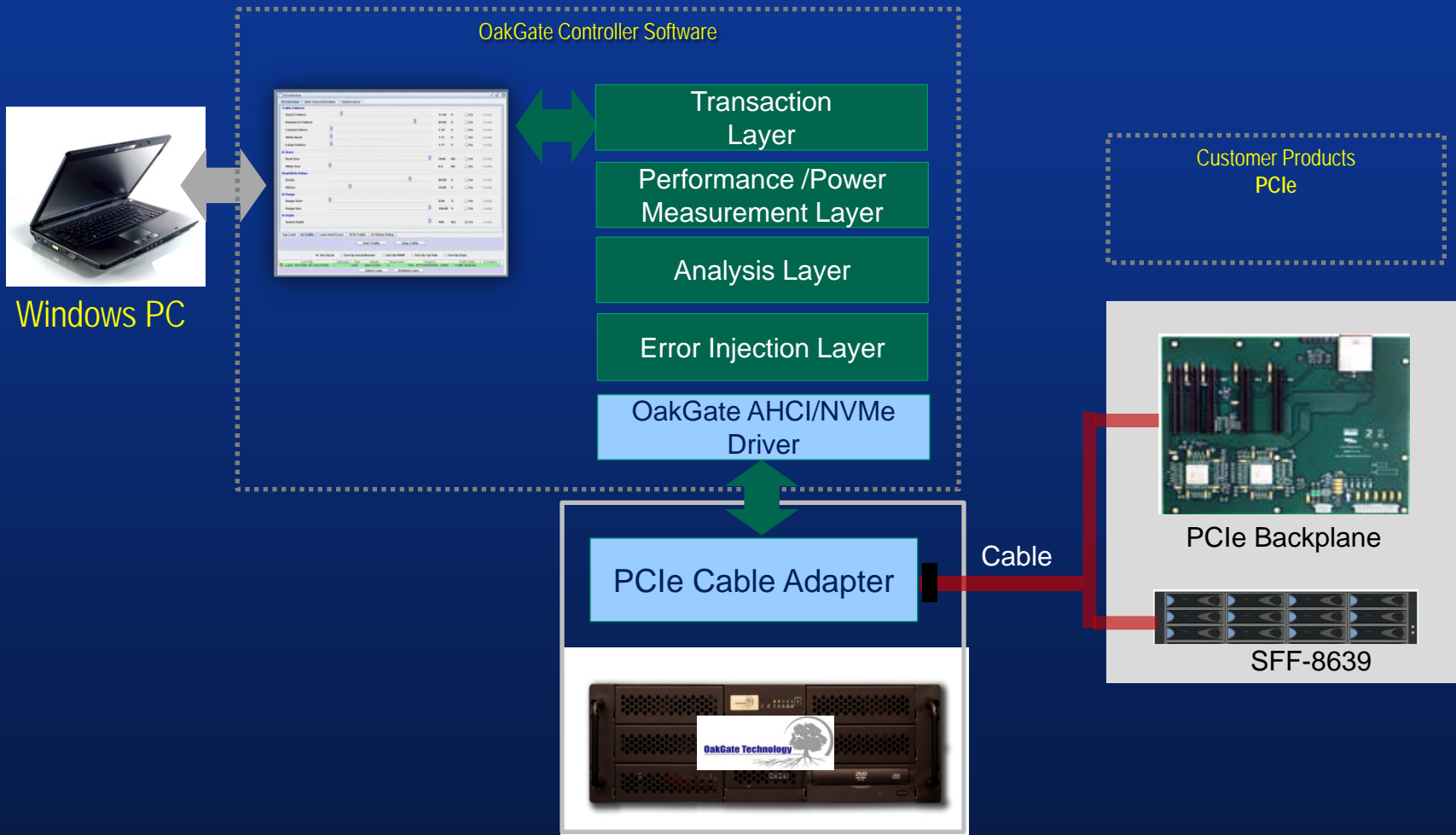- Power Loss/Recovery
- PCIe Resets
- NVMe Performance/Power

| |
|---|
| Application |
| Operating System |
| PCIe Storage Driver |
| Flash Controller |
| Flash Memory |

PCIe Interface

SSD

## Challenges

- Closely coupled to host
- Access via host driver
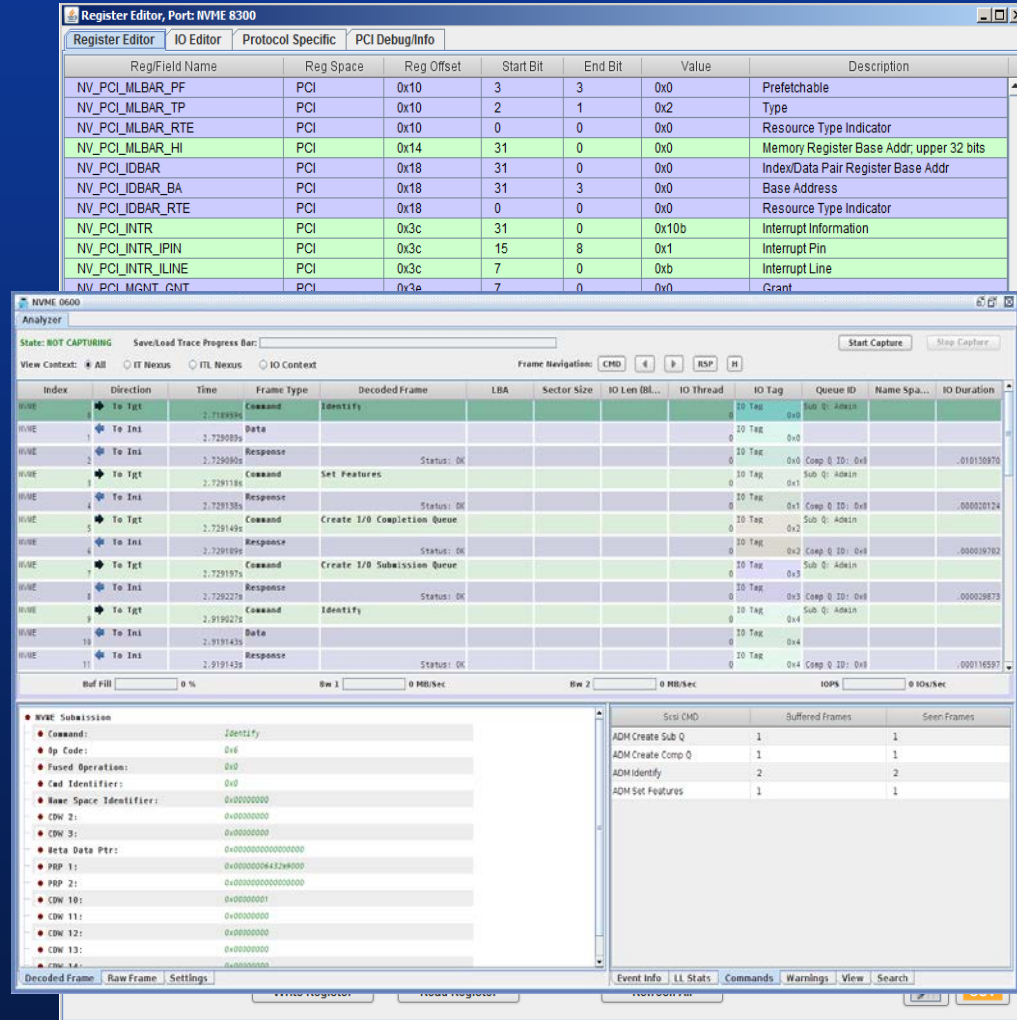- Shared environment
- Limited Hot Plug/Discovery
- Power Measurement

OakGate Technology

# Performance/Validation Environment



OakGate Controller Software

Transaction Layer

Performance /Power Measurement Layer

Analysis Layer

Error Injection Layer

OakGate AHCI/NVMe Driver

Windows PC

Customer Products
PCIe

PCIe Cable Adapter

Cable

PCIe Backplane

SFF-8639

Test System

# Base Level PCIe Validation

- ## Register Access
  - Bit Level verification
  - Error Injection (Bit/Register Level)
  - Base Address Visibility
  - Get/Set Features
- ## Command Analyzer
  - Posted commands w/ responses
  - Visibility to errors

# Queue Configuration

- **Multiple Sub/Comp Queues**

- **Definable Sub/Comp Queue Size**

- **Sub Queue IO Loading**

- **Multi-CPU Queue Distribution**

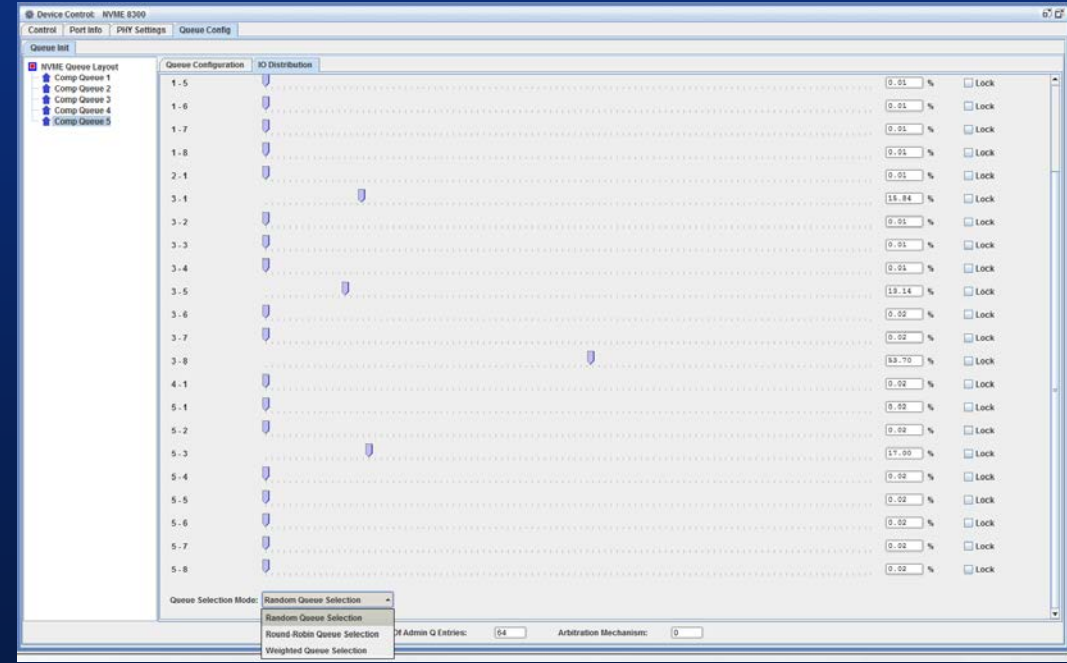- **Definable Interrupt Support**
  - MSI-X
  - MSI
  - Line

# Hot Power Loss/Recovery

- Host Independent Power Control
    - Ability to control power to each device
    - Measure power of each device

- Surprise Power Loss/Recovery
    - Detection of device being powered off
    - Auto re-discovery and drive state verification
    - No Host reboot required

- Auto Queue Management
    - Re-instantiate all Queues
    - Clean up of outstanding commands

- Auto Traffic Control
    - Continue with traffic workload
    - Check for data corruption

# Resets/Queues/Abort

- IO Aborts
- PCIe Resets
  - Function Resets
  - Hot Resets
  - Controller
  - Subsystem
- Auto Re-Discovery
- Queue Management & Recovery

# Data Retention/Corruption

- Validate Pre-existing Data
- Validate 'Acknowledged' Writes
- Evaluate 'Outstanding' Writes
  - Verify Atomic Write Functionality
  - Single LBA corruption
  - Visibility to Partial Writes

**Validate Data Integrity**

☑ Validate Data Integrity           [ Clear All Data Validation Tables in System ]

   ☑ Validate Written Data On Link Up    Nr Of Last Writes To Validate: 1000

     ☑ Validate LBA Atomicity
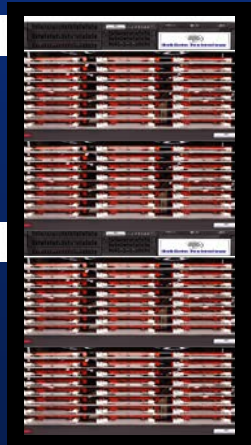
     ☑ Validate Outstanding IOs

> Warning: "Validate LBA Atomicity" and "Validate Outstanding IOs" should only be enabled to check IO's that have NOT been acknowledged. Enabling them could generate errors that may be proper behavior for the target.

# Power Assessment

- ▪ Real-time Power Measurement
  - Power vs Performance
  - Power vs Workload

# Infrastructure Scale

- Improved Space Utilization
- Better Power/DUT
- Enhanced Management

**_SAS/SATA_**
**44U Rack**
**96 SATA/SFF-8639 SSDs**
**OakGate**

**_PCIe_**
**44U Rack**
**256 PCIe SSDs**
**OakGate**

**_Thermal_**
**10.5x5.5 Box**
**1024 SATA SSDs**
**512 PCIe x4 SSDs**
**OakGate/Tanisys**

# Summary

- NVMe has momentum in the Market
    - Products exists from a several major suppliers
- Test/Validation Tools exist
    - With Power Control and Power Measurement
    - Detailed PCIe Register/Queue access
- PCIe Ecosystem is evolving
    - SFF-8639, M.2 etc becoming available
    - Dual Port Support
- NVMe Spec is evolving rapidly
    - Requires Validation Tools that stay current

# Thank You

Questions/Follow-up:
bob.weisickle@oakgatetech.com
www.oakgatetech.com