

SNIA Emerald Measurement Specification Training

By the SNIA Green TWG

Presented by: Patrick Stanko, Jim Espy, Herb
Tanzer, Dave Thiel



- Introduction
- SNIA Emerald Power Efficiency Measurement Specification Overview
 - ◆ Sections
 - ◆ Taxonomy
 - ◆ Measurement
 - ◆ Metrics
- Defining Product Family and Best Foot Forward
- Using the Best Foot Forward in SNIA Emerald™ Data Submissions
- SNIA Emerald™ Program and Data Submission Process

A break will be provided between one of the first three bullet points

➤ SNIA Green Storage Initiative (GSI)

- ◆ To conduct research on power and cooling issues confronting storage administrators
- ◆ Educate the vendor and user community about the importance of power efficiency in shared storage environments
- ◆ Leverage SNW and other SNIA and partner conference to focus attention on energy efficiency for networked storage infrastructures
- ◆ Provide input to the SNIA Green Storage TWG on requirements for green storage metrics and standards
- ◆ Provide external advocacy and support of SNIA Green TWG technical work

➤ SNIA Green Technical Working Group (TWG)

- ◆ Technical body working on green storage metrics and standards
- ◆ Gets direction from GSI
- ◆ Wrote the SNIA Emerald™ Power Efficiency Measurement Specification and related documents



Introduction (Continued)

➤ SNIA Emerald[™] Program

- ◆ Program run by SNIA GSI
- ◆ Open, public, web-based repository location for SNIA Emerald[™] Program Test Data Reports based on the SNIA Measurements Specification
- ◆ Easily identifiable program logo
- ◆ No SNIA membership required
- ◆ More information at the end of training

- Refer to the *SNIA Emerald™ Power Efficiency Measurement Specification v1.0.0 (herein known as Measurement Specification)* for definitive information. This training is based on the specification, but in the event any conflict, the specification takes precedence.
- This training is an overview of the measurement specification with suggestions on SUT configurations to measure
- The methods described in this training for configuring for the “best” results represent the combined expertise of the members of the Green Storage TWG but may not yield the optimal configuration for all products.

Emerald™ Measurement Specification Overview

➤ What does this measurement specification do

- ◆ Defines a proxy method to measure power efficiency of a storage system
- ◆ Covers Online (disk), Near Online (MAID), Removable Media Library (Tape/Optical Library), Virtual Media Library
- ◆ Supports measurement of block based storage
- ◆ Provides a storage taxonomy

➤ What it does not yet do (future revisions)

- ◆ Specify how to measure power efficiency of file system or object-based systems
- ◆ Define measurement procedures for adjunct products or interconnect elements

➤ Eight Sections

- ◆ Sections 1 - 4 cover Overview, References, Scope, Definitions, Symbols, Abbreviations, and Conventions
- ◆ Section 5 defines a storage taxonomy
- ◆ Section 6 provides an top level overview of capacity optimization techniques
- ◆ Section 7 describes the test procedure and requirements
 - › Online
 - › Near online
 - › Removable Media
 - › Virtual Media Library
- ◆ Section 8 names the metrics generated from the test procedure

- Need a taxonomy (product classification) to enable fair comparisons among storage products
- Similar green metrics may apply to all product categories, but different values establish particular best-in-class
- Unique considerations apply to special categories
- Clear taxonomy simplifies comparisons and aids regulatory efforts

Taxonomy Categories

- Six categories define broad storage market segments

Attribute	Category					
	Online	Near Online	Removable Media Library	Virtual Media Library	Adjunct Product	Interconnect Element
Access Pattern	Random/ Sequential	Random/ Sequential	Sequential	Sequential		
MaxTTFD (t)	t < 80 ms	t > 80 ms	t > 80 ms t < 5 min	t < 80 ms	t < 80 ms	t < 80 ms
User Accessible Data	Required	Required	Required	Required	Prohibited	Prohibited

Boxes with grey shading do not have classifications or measurement procedures defined

➤ Classifications are used to compare like products

- ◆ Group products that share common functionality or performance requirements

Category \ Level	Online	Near Online	Removable Media Library	Virtual Media Library
Consumer/ Component	Online 1	Near Online 1	Removable 1	Virtual 1
Low-end	Online 2	Near Online 2	Removable 2	Virtual 2
Mid-range	Online 3	Near Online 3	Removable 3	Virtual 3
	Online 4			
High-end	Online 5	Near Online 5	Removable 5	Virtual 5
Mainframe	Online 6	Near Online 6	Removable 6	Virtual 6

Boxes with grey shading do not have classifications or measurement procedures defined

Taxonomy Online Category

Attribute	Classification					
	Online 1	Online 2	Online 3	Online 4	Online 5	Online 6
Access Pattern	Random/ Sequential	Random/ Sequential	Random/ Sequential	Random/ Sequential	Random/ Sequential	Random/ Sequential
MaxTTFD (t)	t < 80 ms	t < 80 ms	t < 80 ms	t < 80 ms	t < 80 ms	t < 80 ms
User-Accessible Data	Required	Required	Required	Required	Required	Required
Connectivity	Not specified	Connected to single or multiple hosts	Network-connected	Network-connected	Network-connected	Network-connected
Consumer/ Component	Yes	No	No	No	No	No
Integrated Storage Controller	Optional	Optional	Required	Required	Required	Required
Storage Protection	Optional	Optional	Required	Required	Required	Required
No SPOF	Optional	Optional	Optional	Required	Required	Required
Non-Disruptive Serviceability	Optional	Optional	Optional	Optional	Required	Required
FBA/CKD Support	Optional	Optional	Optional	Optional	Optional	Required
Maximum Configuration	≥1	≥ 4	≥ 12	> 100	>400	>400



Taxonomy Near-Online Category



Advancing storage & information technology

Attribute	Classification					
	Near Online 1	Near Online 2	Near Online 3	Near Online 4	Near Online 5	Near Online 6
Access Pattern	Random/ Sequential	Random/ Sequential	Random/ Sequential		Random/ Sequential	Random/ Sequential
MaxTTFD (t)	t > 80 ms	t > 80 ms	t > 80 ms		t > 80 ms	t > 80 ms
User-accessible Data	Required	Required	Required		Required	Required
Connectivity	Not specified	Network connected	Network connected		Network connected	Network connected
Consumer/ Component	Yes	No	No		No	No
Integrated Storage Controller	Optional	Optional	Required		Required	Required
Storage Protection	Optional	Optional	Required		Required	Required
No SPOF	Optional	Optional	Optional		Optional	Required
Non-Disruptive Serviceability	Optional	Optional	Optional		Optional	Required
FBA/CKD Support	Optional	Optional	Optional		Optional	Optional
Maximum Configuration	≥ 1	≥ 4	≥ 12		> 100	> 1000



Green Storage Initiative

Green Storage TWG

Taxonomy Removable Media Library Category

Attribute	Classification					
	Removable 1	Removable 2	Removable 3	Removable 4	Removable 5	Removable 6
Access Pattern	Sequential	Sequential	Sequential		Sequential	Sequential
MaxTTFD (t)	80ms < t < 5m	80ms < t < 5m	80ms < t < 5m		80ms < t < 5m	80ms < t < 5m
User-Accessible Data	Required	Required	Required		Required	Required
No SPOF	Optional	Optional	Optional		Optional	Required
Robotics	Prohibited	Required	Required		Required	Required
No SPOF	Optional	Optional	Optional		Optional	Required
Non-disruptive Serviceability	Optional	Optional	Optional		Optional	Required
Maximum Drive Count	Not specified	4	≥ 5		≥ 25	≥ 25

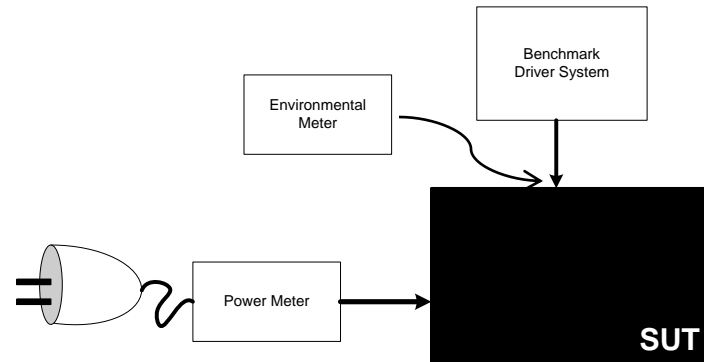
Taxonomy Virtual Media Library Category

Attribute	Classification					
	Virtual1	Virtual 2	Virtual 3	Virtual 4	Virtual 5	Virtual 6
Access Pattern	Sequential	Sequential	Sequential		Sequential	Sequential
MaxTTFD (t)	t < 80 ms	t < 80 ms	t < 80 ms		t < 80 ms	t < 80 ms
User-accessible Data	Required	Required	Required		Required	Required
Storage Protection	Optional	Optional	Required		Required	Required
No SPOF	Optional	Optional	Optional		Optional	Required
Non-Disruptive Serviceability	Optional	Optional	Optional		Optional	Required
Maximum Configuration	12	>12	> 48		> 96	> 96

- Section 7 defines the test/measurement procedure
 - ◆ Defines configuration guidelines and instrumentation requirements
 - ◆ A section for each of the six categories for test execution
- Basic measurement specification procedure
 - ◆ Four continuous test phases
 - › SUT Conditioning
 - › Active test
 - › Read Idle test
 - › Capacity Optimization test
 - ◆ Each category could have a different requirements for each of the test phases

➤ Basic configuration

- ◆ Not allowed to change configuration or tune parameters during test phases



➤ Environmental

- ◆ Climate controlled facility
- ◆ Class I ASHREA
- ◆ Measured at inlet of SUT

Test Phase Power and Measurement Requirements

➤ Input Voltage

NOMINAL INPUT VOLTAGE RANGE	Phases	AC
100-120 VAC RMS	1	47 – 63 Hz
200 – 240 VAC RMS	1	47 – 63 Hz
200 - 480 VAC RMS	3	47 – 63 Hz

➤ Power Meter

- ◆ Measured at a maximum interval of 5 seconds
- ◆ Input power accuracy of 1%
- ◆ Power measurement requirements

Power Consumption (p)	Minimum Accuracy
$p \leq 10 \text{ W}$	$\pm 0.01 \text{ W}$
$10 < p \leq 100 \text{ W}$	$\pm 0.1 \text{ W}$
$p > 100 \text{ W}$	$\pm 1.0 \text{ W}$

➤ Environmental Meter

- ◆ Temperature resolution of 0.1 degree
- ◆ Measured at a maximum interval of 1 minute
- ◆ Measured at inlet

➤ See appendices or SPEC website for recommended meters

Benchmark Driver Requirements

- See appendices for recommended benchmark drivers
 - ◆ VdBench
 - ◆ IOMeter
- Requirements
 - ◆ Uniform distributed random numbers over the range $[0, 2^{64} - 1]$
 - ◆ Reproduce the IO profiles for each taxonomy category and classification
 - ◆ Freely available open-source tools

- Example of scripts used for online system in the user guide download from sniaemerald.com website
 - ◆ Users can cut and copy from user guide
 - › Will have to modify to match system
 - › Instructions on how to modify are in user's guide
- Two set of scripts
 - ◆ First set is used to find IOPS required to meet the specified response time requirement
 - ◆ Second set runs the complete test sequence
- Based on VdBench
 - ◆ Open source benchmark driver
 - ◆ <http://Sourceforge.net/projects/vdbench>

VdBench Recommendations

- If you are not experienced with VDBENCH, you are urged to read the user guide included with the package
- If possible pre fill the disk with data. For small and some medium storage systems you can do this by running the conditioning part for a longer time; for medium sized to large systems this may be not possible to do in a reasonable time (i.e. less than 24 hours)
- Try as much as possible to do a single host run

- **Average response time**
 - ◆ Arithmetic mean of the system response time over a time interval
- **Average power**
 - ◆ Arithmetic mean of the system power over a time interval
- **Operations rate**
 - ◆ Rate of completed work over a time interval
 - ◆ Different for random and sequential work loads
- **Periodic Power Efficiency**
 - ◆ Ratio of operations rate over average power for the same time interval
- **Metric Stability**
 - ◆ 10 point weighted average of periodic power efficiency of one minute intervals through the complete measurement interval
 - ◆ Method to show when it is appropriate to start a measurement interval
- **Time interval defined in specification**
 - ◆ 1 minute
 - ◆ Measurement interval specified for each test phase

SUT Conditioning Test Phase

- Intended to provide a uniform initial condition for subsequent measurements
- Demonstrate the SUT's ability to process IO requests
- Assure that each storage device in the SUT is fully operational and capable of satisfying any supported request
- Achieve typical operational temperature
- Each taxonomy category will have different measurement interval requirement to demonstrate stability

IO Profile	IO Size (KiB)	Read/Write Percentage	IO Intensity	Transfer Alignment (KiB)	Access Pattern
Mixed Workload 1 (i=MW1)	8	70/30	100	8	Random
Mixed Workload 2 (i=MW2)	8	70/30	25	8	Random
Random Write (i=RW)	8	0/100	100	8	Random
Random Read (i=RR)	8	100/0	100	8	Random
Sequential Write (i=SW)	256	0/100	100	256	Sequential
Sequential Read (i=SR)	256	100/0	100	256	Sequential

- All or some of the IO profiles are used by the defined taxonomy categories
 - ◆ Drive enough IOs to reach the required response time or through-put specified in the measurement specification
 - ◆ The 25 IO intensity is 25% of the IO defined for MW1

SUT Active Test Phase

- ▶ **IO profiles used by taxonomy category**
 - ◆ Online and Near-Online use all six IO profiles
 - ◆ Removable Media and VML use only the sequential IO profiles
- ▶ **Run as an uninterrupted sequence of workloads**
 - ◆ Specification defines the order to be run for each taxonomy category

SUT Ready Idle Test Phase

- Defined as storage systems and components that are configured, powered up, connected to one or more hosts and capable of satisfying externally-initiated, application-level initiated IO requests within normal response time constraints, but no such IO requests are being submitted.
- Average power measured in the measurement window
- No external IO given by the host
- Can perform any IO within the taxonomy required response time interval

SUT Capacity Optimization Method Test Phase

➤ Heuristic tests

- ◆ Delta snapshots
- ◆ Thin provisioning
- ◆ Data de-duplication
- ◆ Parity RAID
- ◆ Compression

➤ Run after ready idle test phase

➤ C program generated by SNIA

- ◆ Download from sourceforge.net/projects/sniadeduptest
- ◆ Used for de-duplication and compression

➤ Taxonomy dependent

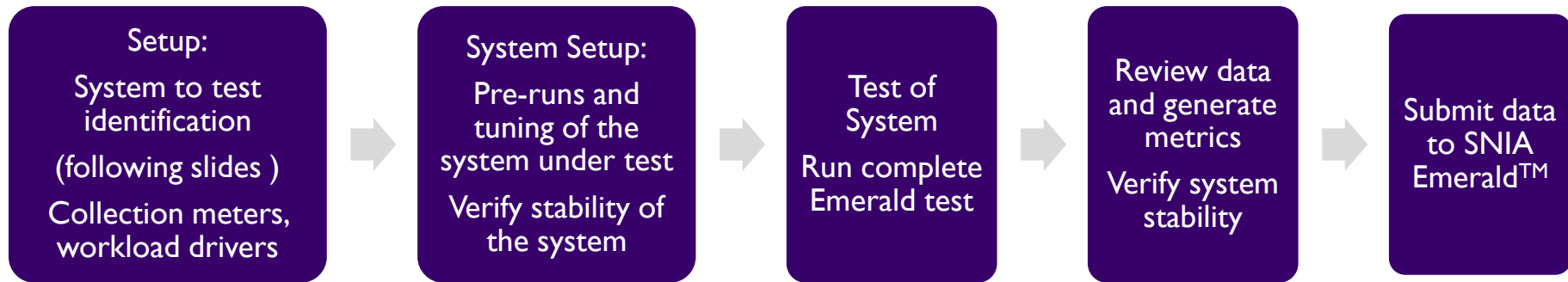
➤ Active (Primary)

- ◆ Ratio of operations rate over average power for the same measurement interval
 - EP_{MW1} (IOP/S/W) of the 70% mixed workload at maximum response time
 - EP_{MW2} (IOP/S/W) of the 25% of the IO used in MW1
 - EP_{RR1} (IOP/S/W) of the random read workload at maximum response time
 - EP_{RW1} (IOP/S/W) of the random write workload at maximum response time
 - EP_{SR1} (MiB/S/W) of the sequential read workload at maximum throughput
 - EP_{SW1} (MiB/S/W) of the sequential write workload at maximum throughput
- ◆ Number of active metrics generated dependent on the taxonomy category tested

IO Profile	IO Size (KiB)	Read/Write Percentage	IO Intensity	Transfer Alignment (KiB)	Access Pattern
Mixed Workload 1 (i=MW1)	8	70/30	100	8	Random
Mixed Workload 2 (i=MW2)	8	70/30	25	8	Random
Random Write (i=RW)	8	0/100	100	8	Random
Random Read (i=RR)	8	100/0	100	8	Random
Sequential Write (i=SW)	256	0/100	100	256	Sequential
Sequential Read (i=SR)	256	100/0	100	256	Sequential

- **Ready Idle (Primary)**
 - ◆ Ratio of raw capacity over average power measured in the defined measurement window
- **Capacity Optimization (Secondary)**
 - ◆ A yes/no for each Capacity Optimization Method tested
 - ◆ Do not have to test all COMs but if vendor declares to have a COM it must be tested and on during active test phase

Flow Needed for Valid Emerald Measurement



➤ General timeline

- ◆ Tune the system
- ◆ A day to run test
- ◆ A day to generate the required data and review it
- ◆ A few hours to submit the data

- Introduction
- SNIA Emerald Power Efficiency Measurement Specification Overview
 - ◆ Sections
 - ◆ Taxonomy
 - ◆ Measurement
 - ◆ Metrics
- Defining Product Family and Best Foot Forward
- Using the Best Foot Forward in SNIA Emerald[™] Data Submissions
- SNIA Emerald[™] Program and Data Submission Process

➤ Wide Spectrum of Storage-Oriented Products

- ◆ Created a taxonomy to narrow scope
- ◆ Categories: On-Line, Near-Line, etc.
- ◆ Classifications: Further granularity of each Category

➤ Still too Broad in Scope

- ◆ Vendors may have multiple products in a particular Category/Classification
- ◆ Each product may have many configuration variables

➤ Requirement/Challenge: Select Appropriate Test Configurations

- ◆ Comprehensive and usable results for customer
- ◆ Minimized, lower cost, but effective testing methods for vendor

Concept of Product and Family

➤ Product:

- ◆ Represents a fundamental performance capability space that separates it from any other potentially related products

➤ Product Family:

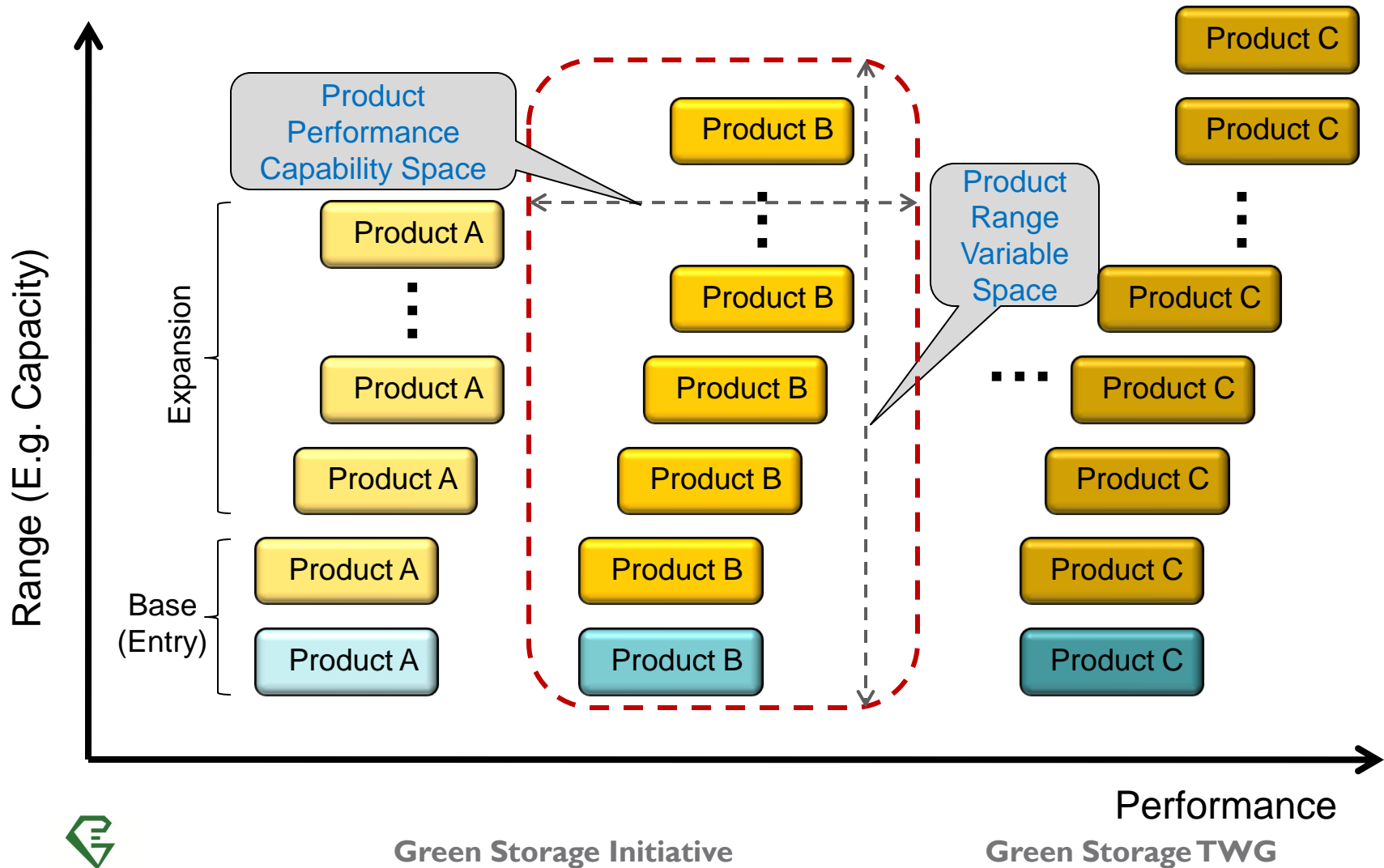
- ◆ Represents the full range space of configuration variables and options for a particular product.

➤ Term Usage:

- ◆ Terms *family* and *range* are used interchangeably and may include such aspects as number and type of storage device (spinning or solid state drive), cache size, availability levels, etc.

- **Vendor Aligns Product(s) with SNIA Taxonomy Category**
 - ◆ Hopefully straightforward – Taxonomy will adapt over time
- **Vendor Aligns Product(s) with Category Classification**
 - ◆ Will be some boundary gray areas - E.g. OL-3 or OL-4?
- **Vendor Further Defines Product/Family Configurations**
 - ◆ The really hard part...
- **Conceptual Representation**
 - ◆ Next slide depicts a possible product/family (range) differentiation
 - ◆ Believed applicable to most storage system architectures

Simplified Product/Family Representation



➤ Products Could be of Various Architectural Types

- ◆ Monolithic – Little or no scaling but may still have family aspects
- ◆ Scale-up – E.g. base controller + storage expansion
- ◆ Scale-out – E.g. base compute/storage + compute/storage expansion
- ◆ Others TBD

➤ Product Performance Typically Scales With Expansion

- ◆ Varying degrees
 - › Scale-up performance typically rolls off at varying degrees before max configuration
 - › Scale-out performance can be linear with increasing configurations
- ◆ Any inter-product performance overlap driven by vendor's market positioning

Family (Range) Discussion

➤ Range Variables

- ◆ Example on previous product/family depiction focuses on capacity
- ◆ Could involve other variables

➤ Range Variable Types

- ◆ Particular Items of highest potential energy consumption impact:
 - › Controller or related compute element – Typically defines performance aspect
 - › Cache – Also performance oriented - Not considered part of the user-addressable space
 - › Number and type of persistent storage devices – Defines user-addressable space
 - › RAS items – As necessary to meet reliability, availability, serviceability requirements
 - › Capacity optimization – Functionality (typically software) that more effectively utilizes physical storage space such as thin provisioning, compression and de-duplication
- ◆ Many other examples
 - › Power supplies, cooling, I/O, etc.

Approach to Range Variable Reduction

➤ Range Variable Reduction is Difficult

- ◆ Even with the 5 listed items still too many test cases
 - › Significant set-up and execution times
 - › Complex results sets
- ◆ Maximum system size testing is expensive and cumbersome to manage
- ◆ Need a simpler alternative...

➤ “Best Foot Forward” (aka Sweet Spot) - BFF

- ◆ Find proxy family configuration(s)
 - › Reasonably representative of the all range variables?
- ◆ Find test point(s) where Measurement Specification active metrics are best
 - › The “sweet spot”
- ◆ Suitable for any architecture
 - › E.g. scale-up, scale-out, hybrid, ...

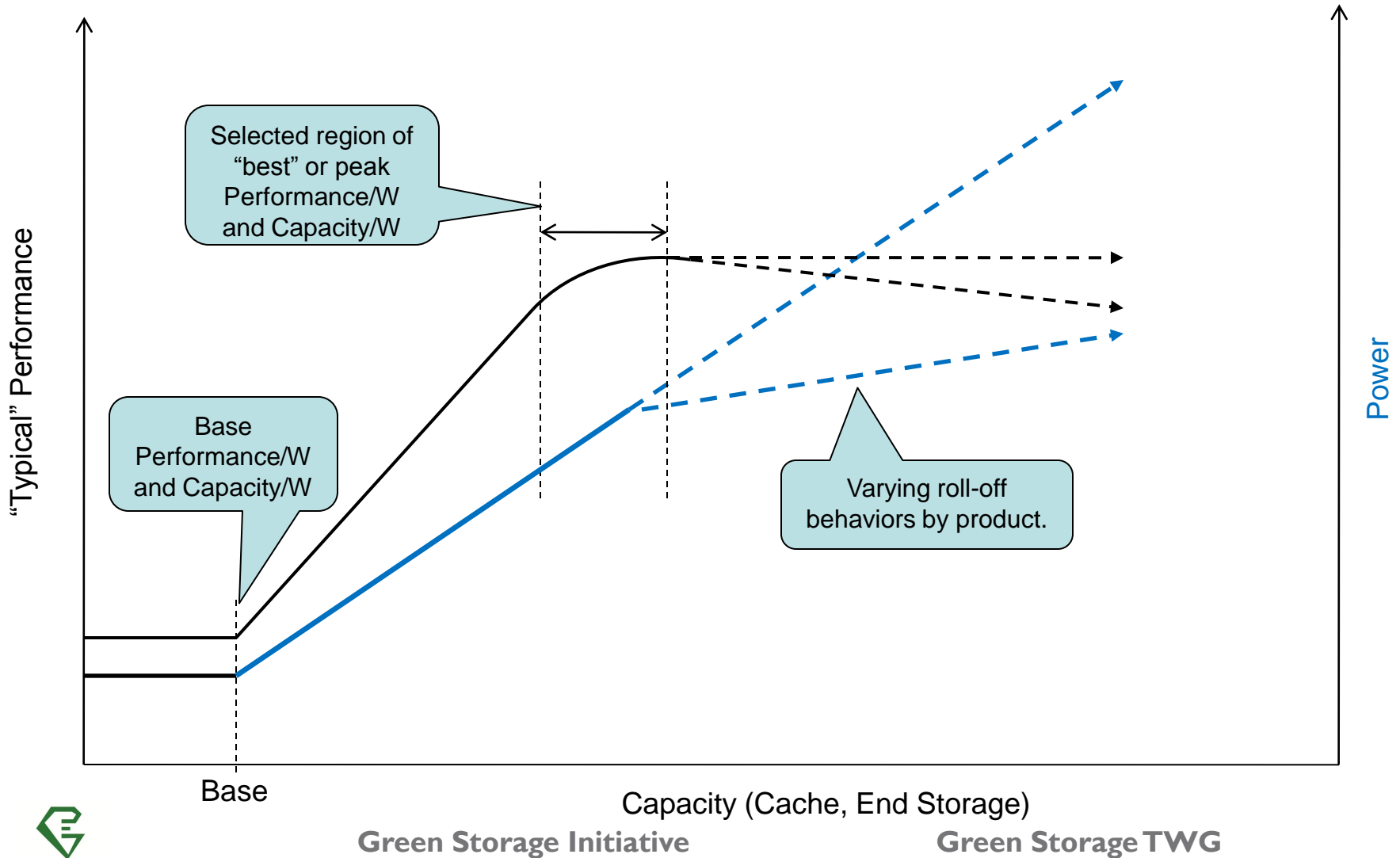
➤ BFF Looks Holistically at Storage System Product/Family

- ◆ Allows vendor to select and test one product/family configuration
 - › Or more if desired
- ◆ At operating points near the Measurement Spec metric peak values
 - › I.e. the “sweet spot”
- ◆ Results reasonably representative of the entire family
 - › Easier and less expensive for the vendor
 - › Simple and understandable results for the potential customer

➤ Scale Up Example on Following Slide

- ◆ Based on notion that Measurement Spec active metrics have peak values
- ◆ Peaks typically located at points well below maximum configurations

Best Foot Forward Approach Scale-Up System



➤ Previous Slide is a Rough Approximation

- ◆ Capacity increases are actually more stepwise
- ◆ Performance roll-off can vary by product
 - › Dashed lines attempt to show one (of possibly many) changes due to different storage technology tiers, e.g. scaling capacity w/large SATA drives
- ◆ Regardless, example depicts a smaller test configuration

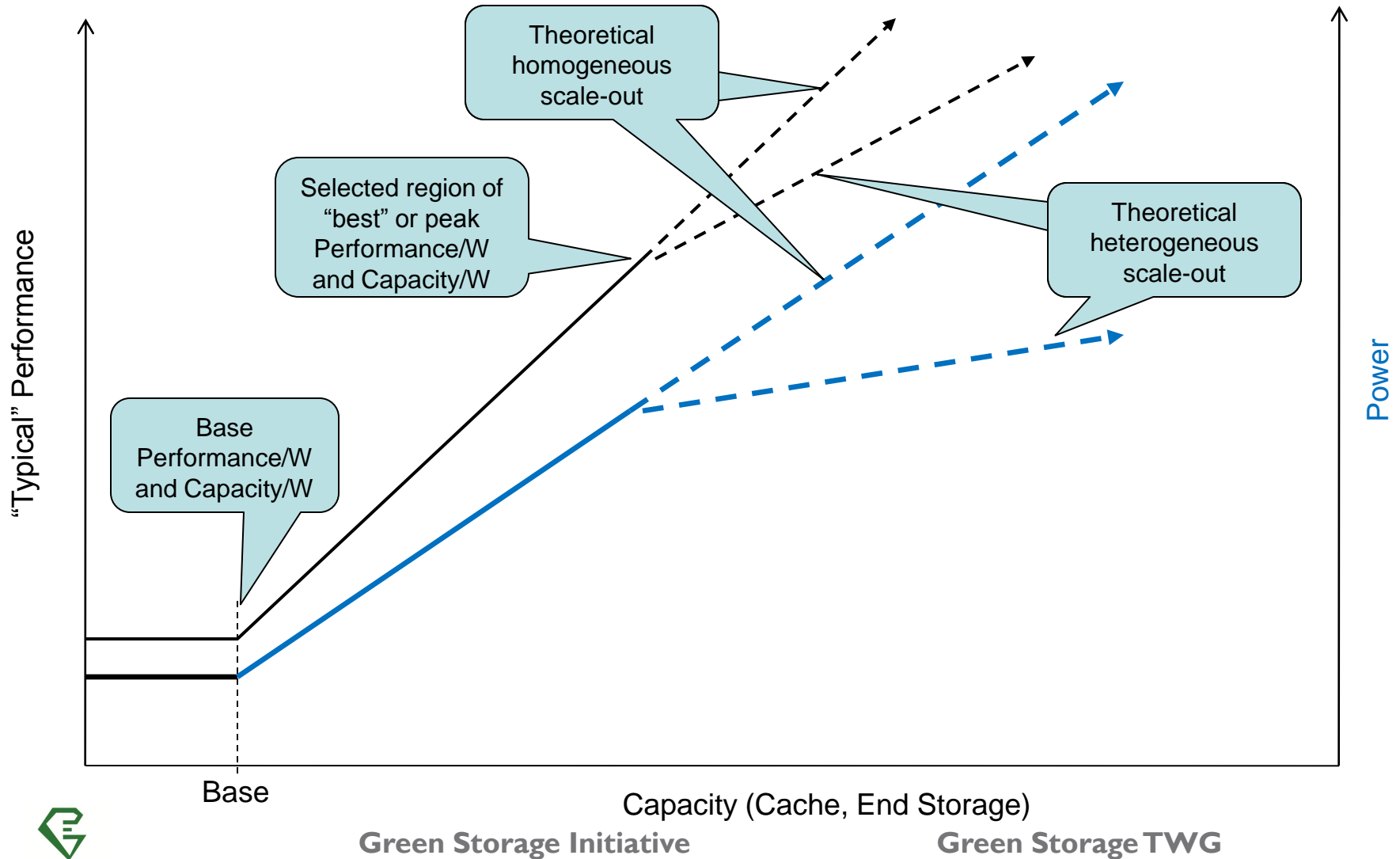
➤ What About Other Test Points?

- ◆ Could also test at base (entry point) but not required
- ◆ Key is no requirement to test beyond the peak point

➤ Scale Out Example on Following Slide

- ◆ What if there is no clearly discernable peak?

Best Foot Forward Approach Scale-Out System



Best Foot Forward Approach

➤ Again a Rough Approximation

- ◆ Capacity increases are actually more stepwise
- ◆ Dashed lines attempt to show one (of possibly many) changes due homogeneous vs heterogeneous scale-out configurations
- ◆ Can still select a smaller test configuration

- **Given Known Taxonomy Category and Classification**
 - ◆ Vendor determines one or more family representative configurations
 - ◆ Vendor locates Measurement Spec active metric peak points
 - ◆ Tests are performed on this reduced configuration (set)
 - › Note: For smaller systems, the BFF may in fact be the maximum configuration
- **Where is the Performance/W Peak?**
 - ◆ Depends on numerical increase of numerator vs denominator with capacity
 - ◆ If numerator initially increases more than denominator, a clear peak
 - ◆ Else it becomes harder – Just pick a point before it rolls off?
 - ◆ *Selection of the peak point(s) is the subject of the next slides*

- Introduction
- SNIA Emerald Power Efficiency Measurement Specification Overview
 - ◆ Sections
 - ◆ Taxonomy
 - ◆ Measurement
 - ◆ Metrics
- Defining Product Family and Best Foot Forward
- Using the Best Foot Forward in SNIA Emerald™ Data Submissions
- SNIA Emerald™ Program and Data Submission Process

Best Foot Forward aka Sweet-spot

- The benefit of Best Foot Forward (BFF) is to reduce the full range of variables of a product family to just a few test configurations; this reduced test set can be considered representative of the entire product family
- The BFF consists of the optimized configurations that will produce a set of peak power efficiency metrics of a product family for the different test phases
 - ◆ Random workloads [IOP/s/Watt]
 - ◆ Sequential throughput [MiB/s/Watt]
 - ◆ Idle Capacity [GB/Watt]

An Approach for Emerald Data Submission

- Start by aligning your product family within a taxonomy definition
- Baseline run to establish the test process
 - ◆ Start w/ available configuration; no particular “tuning” in affect
 - ◆ Identify any issues with conditioning, stability, response times, etc. (per the run rules), post-processing, reporting, etc.
- Consider all possible (and valid) product SKU’s to identify configurations that will give the peak power efficiency metrics
- Using estimator tools, identify the “best-foot- forward” or “sweet-spot” relative to each specific test profile
- Set-up, test, and measure the peak metric values for your first sweet-spot
 - ◆ Run through the complete sequence of test phases
 - ◆ Test validate and data correlate
- For each additional sweet-spot of interest, re-configure and re-test

Candidate SUT: A shipping Online-3 SAN

- ▶ Full redundancy except for single midplane in dual-controller enclosure
- ▶ Two controller performance points, with variable cache and front-end interfaces
- ▶ The lower product class can support 120xLFF or 250SFF and the higher product class can support 240xLFF or 450xSFF
 - ◆ 12 x LFF drive shelves
 - ◆ 25 x SFF drive shelves

Supported drives, 6Gb SAS

- ◆ SFF
 - > 146GB, 15K
 - > 300GB, 10K
 - > 450GB, 10K
 - > 600GB, 10K
 - > 500GB, 7.2K midline
 - > 200GB SSD*
 - > 400GB SSD*
- ◆ LFF
 - > 300GB, 15K
 - > 450GB, 15K
 - > 600GB, 15K
 - > 2TB, 7.2K midline

* Will characterize SSD's separate from spinning drives



Test Phase IO Profiles for Online & Near Online

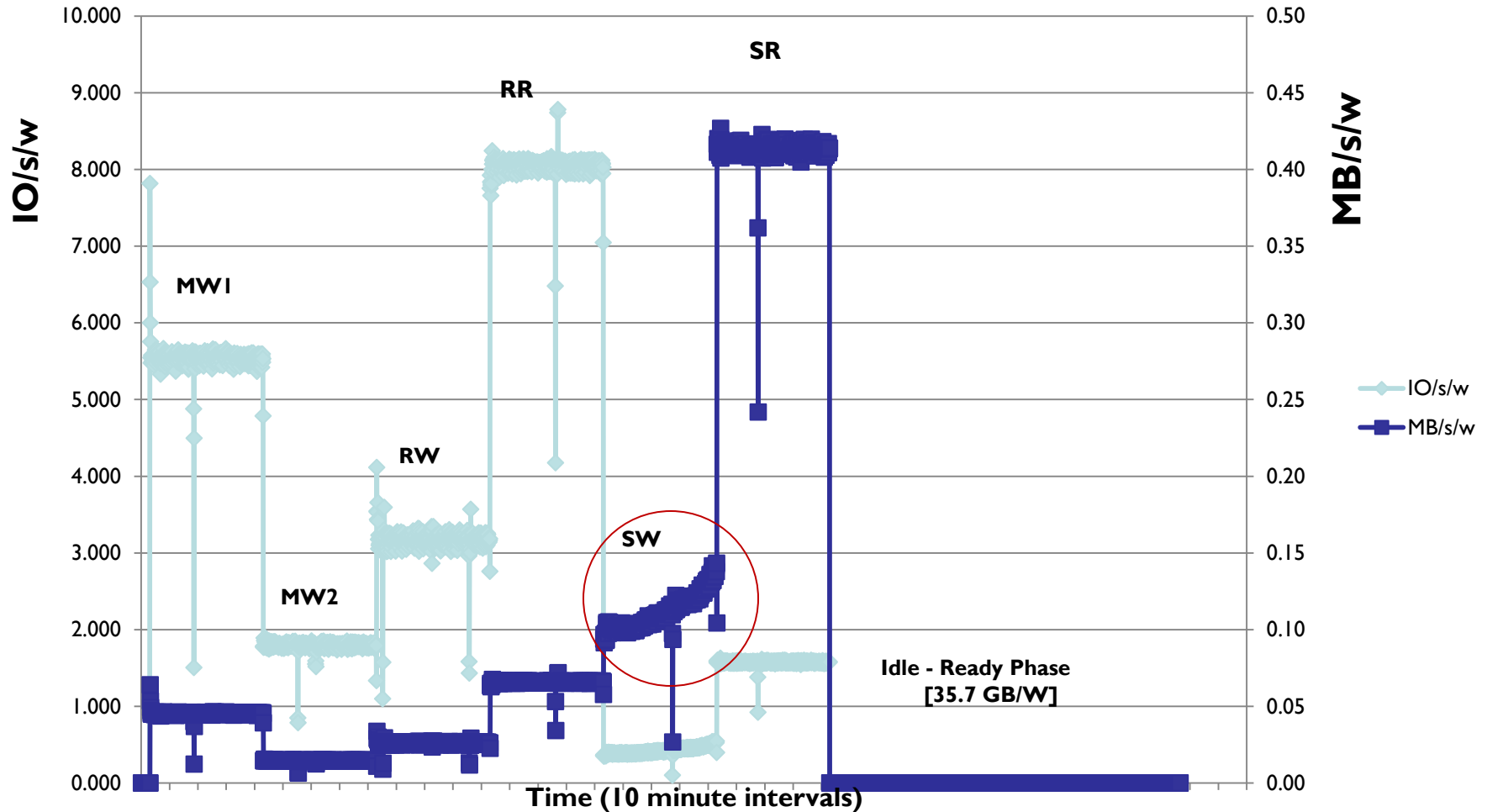
IO Profile	IO Size (KiB)	Read/Write Percentage	IO Intensity	Transfer Alignment (KiB)	Access Pattern
Mixed Workload 1 (i=MW1)	8	70/30	100	8	Random
Mixed Workload 2 (i=MW2)	8	70/30	25	8	Random
Random Write (i=RW)	8	0/100	100	8	Random
Random Read (i=RR)	8	100/0	100	8	Random
Sequential Write (i=SW)	256	0/100	100	256	Sequential
Sequential Read (i=SR)	256	100/0	100	256	Sequential

* The test phases are executed in an uninterrupted sequence. Each test phase shall last a minimum of 40 minutes, comprised of a minimum of 10 minutes to establish stability followed by 30 minutes as the measurement interval.

A 2-hour ready idle test follows the above active tests

Baseline Test Results for Candidate SUT

Starting available config: 370 SAS drives (no "tuning")
10 shelves x 25 SFF drives (10K-300GB)
10 shelves x 12 LFF drives (15K-600GB)



Baseline Test Observations

- ▶ Peak workload efficiency metric occurred during RR phase (~ 8 IOP/s/W)
- ▶ Peak throughput efficiency metric [MB/s/W] occurred during SR phase (~ 0.42 MB/s/W)
- ▶ Power consumed for any workload varies only ~12% (4100W at idle to 4600W during RR) → performance can have a bigger influence on the metric

Note: Easy to observe that SeqWrite (SW) measurement interval has not reached stability; stability is required in order to have a valid metric measurement

Finding the Best Foot Forward

- While there are 7 different Emerald test profiles (online); you may have from 1 to 7 possible different optimized configurations – vendor choice
 - ◆ 4 x Random [IOP/s/Watt]
 - ◆ 2 x Sequential [MiB/s/Watt]
 - ◆ 1 x Ready-Idle [raw capacity, GB/Watt]
- Recommended to use estimator tools that combine power and performance to predict the peak metrics
 - ◆ The alternative is educated derivations and potentially a lot of testing that is very labor and resource intensive
 - ◆ As long as the simulated results are reasonably accurate, the physical configuration selected for actual test to measure the peak value can be limited in range

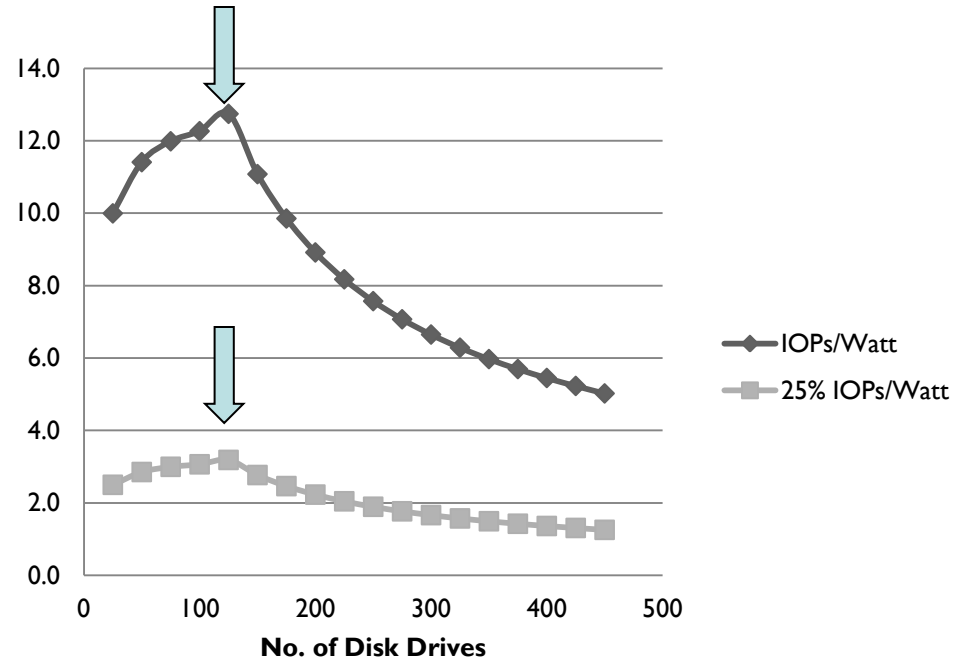
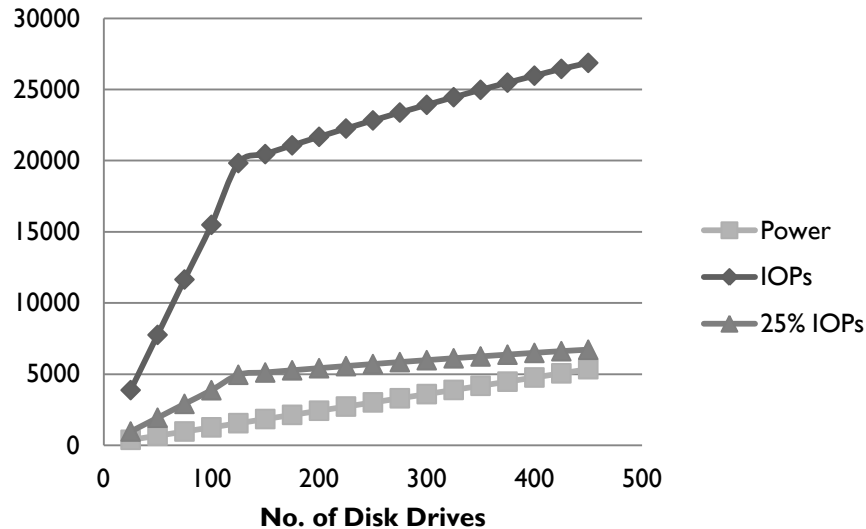
Predicted Peak Metrics for an Online-3 Test Candidate

Exercise #	Prediction basis
1 1.5	Mixed Workload, Random 70/30 R/W -- Granular level, single drives
2	Random Read (100/0 R/W) & Random Write (0/100 R/W)
3	Sequential Read (100/0 R/W) & Sequential Write (0/100 R/W)
4	RAID level
5	Ready Idle

Exercise I: Mixed Workload

8K Random 70/30 R/W

SFF 15K rpm, RAID 5

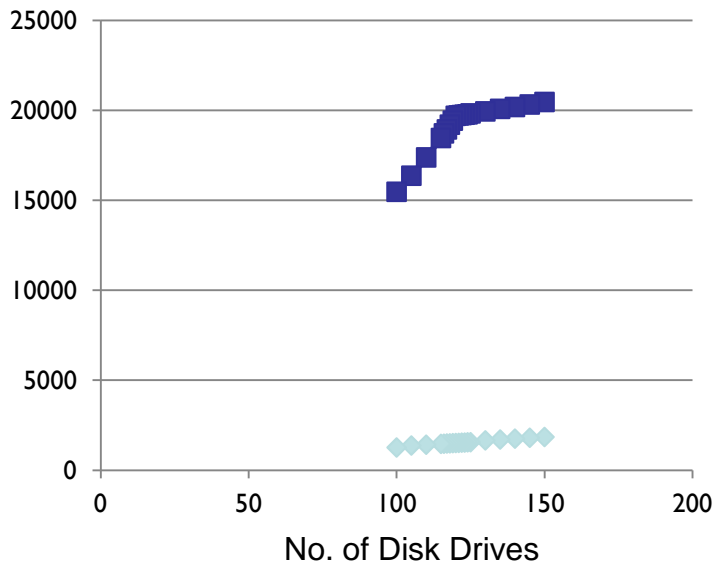


- Peak metric = 12.7 IOP/s/Watt at 125 drives
- Changing the read/write mix changed the metric but not the drive count
60/40 r/w = 11.5 IOP/s/W; 80/20 r/w = 14.9 IOP/s/W

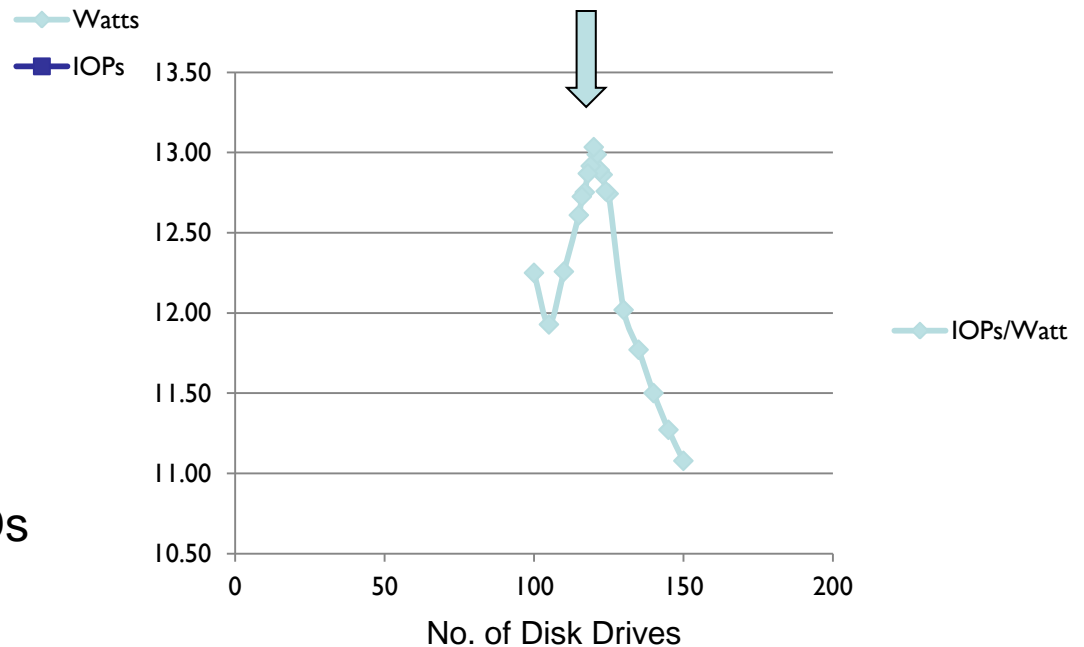
Note: Incrementing drive count by full JBOD

Exercise 1.5: Granular Drive Counts (Increment by Single HDDs)

8K Random, 70/30 R/W, SFF 15K rpm, RAID 5



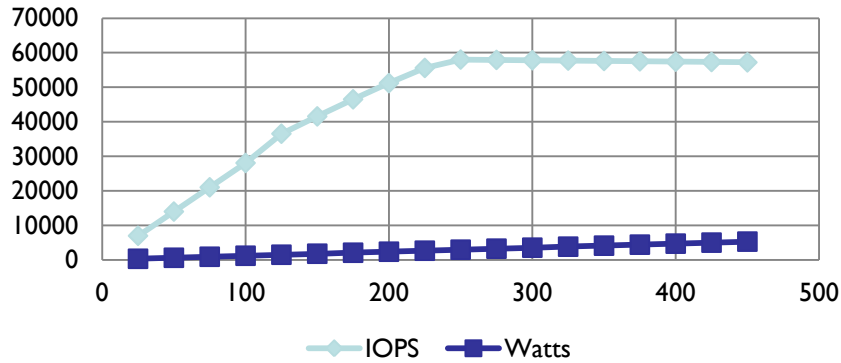
Peak Metric:
13.03 IOP/s/Watt & 120 HDDs



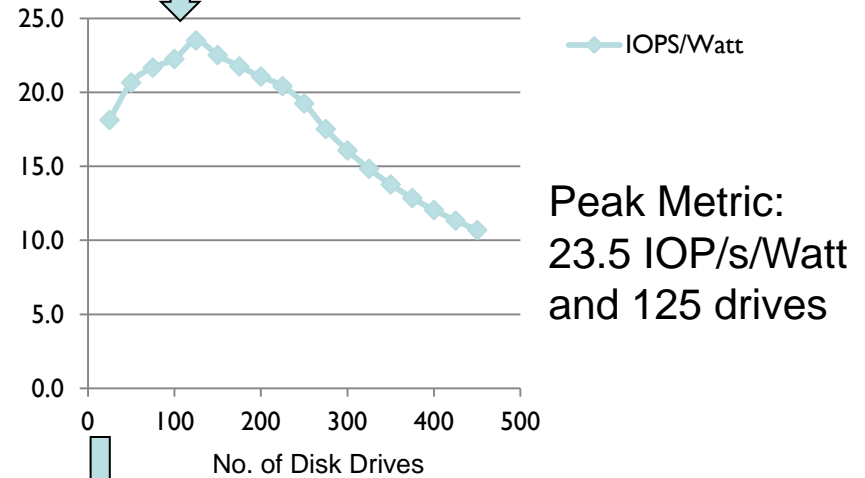
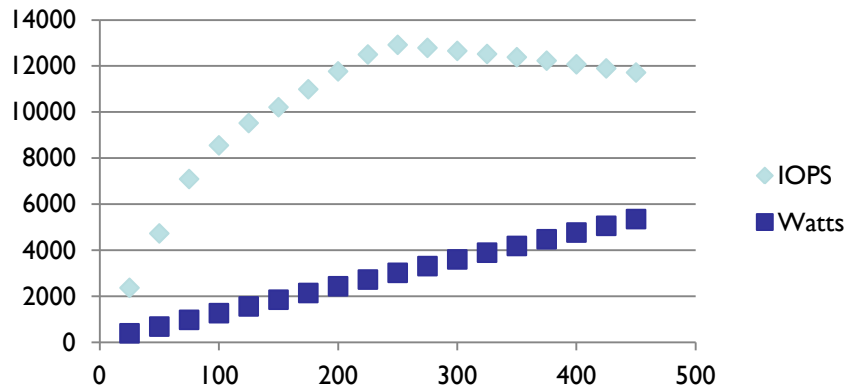
Exercise 2: 8K Random Read, Write

SFF 15K rpm, RAID 5

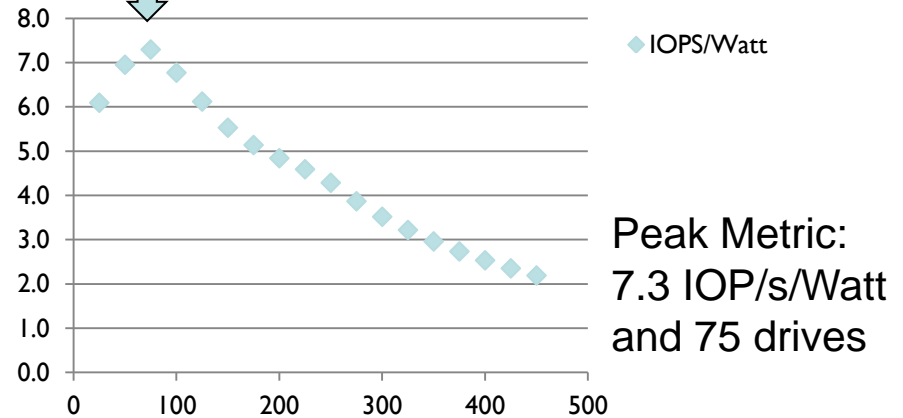
Random Read



Random Write



Peak Metric:
23.5 IOP/s/Watt
and 125 drives

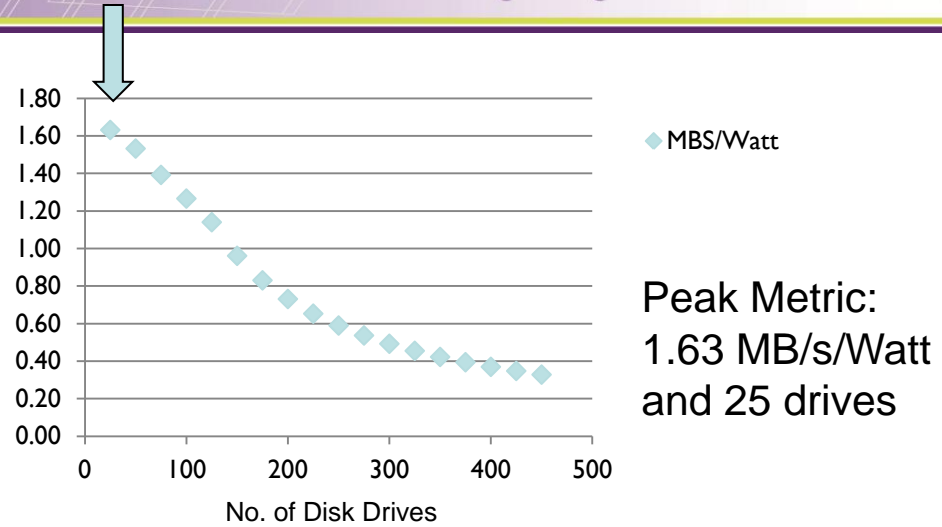
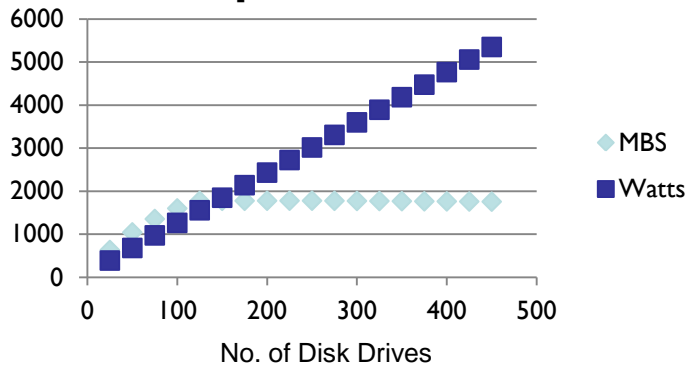


Peak Metric:
7.3 IOP/s/Watt
and 75 drives

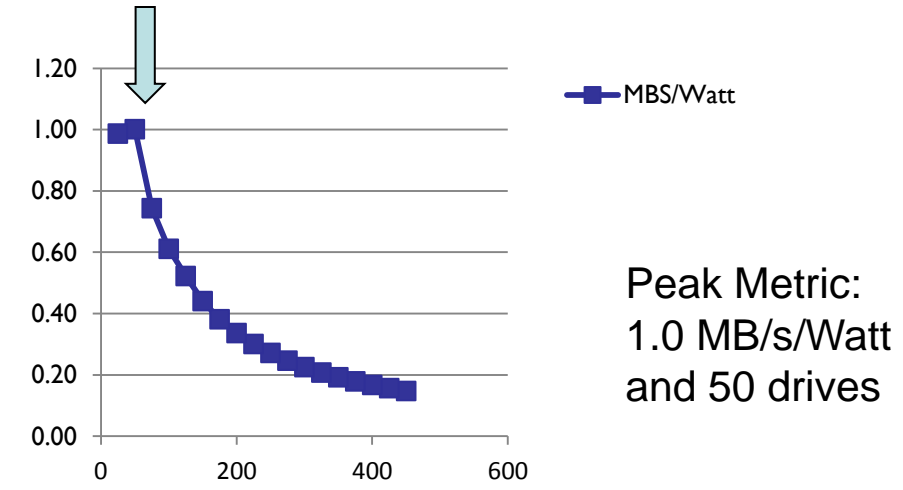
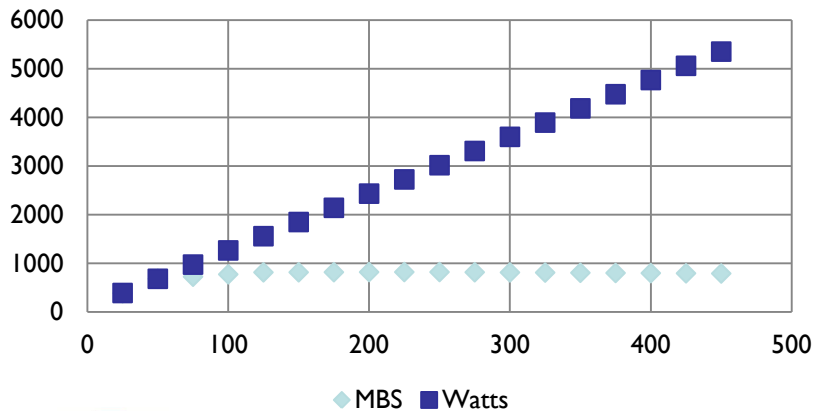
Exercise 3: 256KB Sequential Read, Write

SFF 15K rpm, RAID 5

Sequential Read



Sequential Write



Exercise 4: RAID level (SFF 15K rpm)

Peak Power Efficiency [IOP/s/Watt] or {MB/s/Watt}, # of HDDs

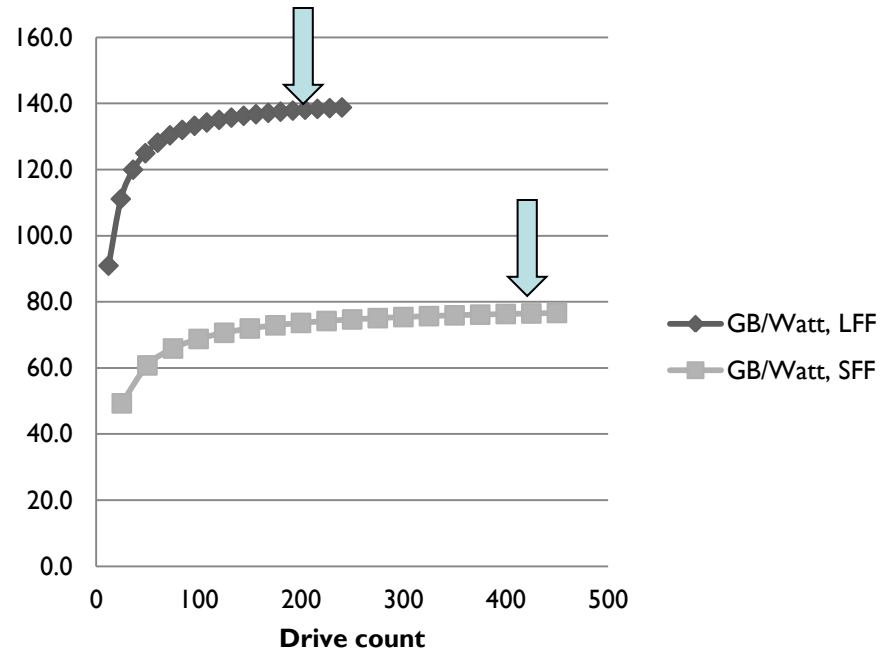
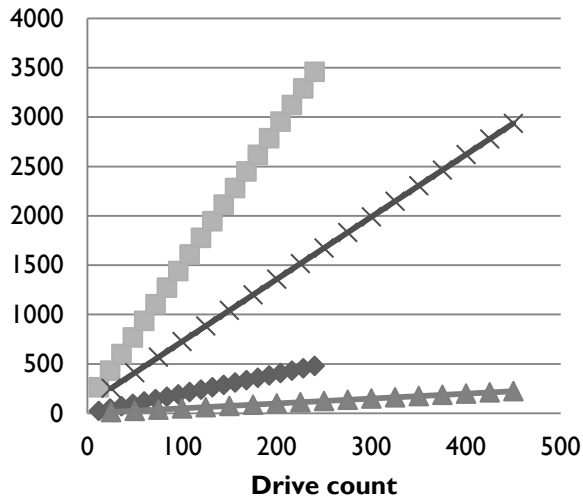
	8K Rand Mixed (70/30 R/W)	8K Rand Read	8K Rand Write	128K Seq Read	128K Seq Write
RAID 5 Distributed single parity	[12.7], 125	[23.5], 125	[7.3], 75	{1.63}, 25	{1.00}, 50
RAID 10 Blocks striped and mirrored	[18.6], 125	[23.2], 125	[13.3], 75	{2.0}, 25	{0.68}, 50

Notes:

- 1) The Online 3 category is required to have RAID protection
- 2) RAID 6 (double parity) will offer greater redundancy but poorer performance efficiency

Exercise 5: Ready-Idle

LFF 2TB 7.2K rpm and SFF 500GB 7.2K rpm drives at Ready-Idle



Peak Metrics:

LFF: 138.8 GB/W and 240 drives

SFF: 76.6 GB/W and 450 drives

General Observations for the Candidate SUT

- Active cases - the performance* reaches a roll off point relatively early (i.e., smaller drive count); then it either levels out or goes down slightly. The peak [Performance/Power] metric is reached at or before this performance roll off point
 - ◆ All peak predictions for Random are reached with the same drive type (15K, SFF) and close in drive count (125 or 75)
 - ◆ All peak predictions for Sequential reached with the same drive type (15K, SFF) and close in drive count (50 or 25)
- Ready-idle case - the peak metric levels but continues to slowly rise with drive count (as the controller electronics power is amortized over increasing numbers of drives)
 - ◆ Vendor can choose to test and submit a lower drive count configuration, and add to the notes a power calculator based projection for the largest drive count configuration

*Note: very dependent on specific Controller performance and bandwidth behavior

Sample Data Submission (Online-3 SUT)

Operational Power

Idle power test

Average watts	592.692 W
Raw capacity tested	7300 GB
EP _{RI}	12.317 GB/W

Standard idle metric

GB per Watt

Note: 1 GB = 10⁹ bytes; 1 GiB = 2³⁰ bytes
a GiB is about 7.4% larger than a GB

Active power tests

EP _{RR}	17.255	run length (minutes)	30	Average latency	12 ms
Small random reads	I/Os per second per Watt				
EP _{RW}	9.155	run length (minutes)	30	Average latency	7 ms
Small random writes	I/Os per second per Watt				
EP _{SR}	3.01	run length (minutes)	30	Average latency	6 ms
Large sequential reads	MiB per second, per Watt				
EP _{SW}	1.22	run length (minutes)	30	Average latency	9 ms
Large sequential writes	MiB per second, per Watt				
EP _{MW1}	10.501	run length (minutes)	30	Average latency	13 ms
Mixed workload 1	I/Os per second per Watt				
70% random, 30% sequential, I/O intensity = 100					
EP _{MW2}	3.486	run length (minutes)	30	Average latency	4 ms
Mixed workload 2	I/Os per second per Watt				
70% random, 30% sequential, I/O intensity = 25					

NOTE: power-related numbers are required to be reported to three significant digits



SNIA Emerald™

The SNIA Emerald Test Data Report

Disclosure for storage systems and products

Capacity Optimizations

	On during test?	Available in SUT?
Deduplication		NO
Compression		NO
Thin provisioning		YES
Parity RAID	YES	
Read-only delta snapshots		YES
Writeable delta snapshots		YES

Other mandatory disclosures, per spec

Test data provided is for a specific configuration that is tuned to achieve the best SeqRead and SeqWrite performance (the "sweet-spot"). The sweet spot data for alternate configurations that are tuned for the best Random and Idle metrics will be added in the near future.

Data Comparison Observations

	MWI (70/30 R/W)	Seq Write (0/100 R/W)	Seq Read (100/0 R/W)	Ready Idle
Baseline Test (no tuning, 250SFF+120LFF)	5.6 IOP/s/W	0.12 MB/s/W	0.42 MB/s/W	35.7 GB/W
Estimator Predicted BFF	13.03 (120 SFF, 15K)	1.0 MB/S/W (25 or 50 SFF, 15K)	1.63 MB/s/W (25 SFF) 1.54 MB/S/W (50 SFF)	138.8 GB/W (240 LFF, 2GB/7.2K) 76.6 GB/W (450 SFF, 500GB/7.2K)
Data Submission (50 SFF, 146GB, 15K)	10.5	1.22 MiB/s/W*	3.01 MiB/s/W*	12.317 GB/W

* Note: 1 MB = 10⁶ bytes; 1 MiB = 2²⁰ bytes
A MiB is 4.86% larger than a MB

- Introduction
- SNIA Emerald Power Efficiency Measurement Specification Overview
 - ◆ Sections
 - ◆ Taxonomy
 - ◆ Measurement
 - ◆ Metrics
- Defining Product Family and Best Foot Forward
- Using the Best Foot Forward in SNIA EmeraldTM Data Submissions
- **SNIA EmeraldTM Program and Data Submission Process**

➤ Purpose

- ◆ Provide open access to storage system power efficiency information using a well-defined testing procedure and additional information related to system power characteristics
- ◆ The report data can help IT professionals make storage platform selections as part of an overall Green IT and Sustainability objective

➤ Test procedure: SNIA Emerald™ Power Efficiency Measurement Specification

➤ Public access and submittal is through the sniaemerald.com web site

- ◆ No charge for access to test results, specifications or user guides
- ◆ Submission of results is for a modest fee, discounted or waived for SNIA/GSI members
- ◆ SNIA membership is not required to submit or to access test results
- ◆ Voluntary, non-exclusionary, low cost program for manufacturers - Options to self-measure or third party measurement

➤ Process

- ◆ Storage Vendors test their equipment and submit test results to the Emerald Program
- ◆ Emerald Program publishes results on the sniaemerald.com web site
- ◆ IT users (public) download results from the sniaemerald.com web site
- ◆ Vendor gains right to use the SNIA Emerald™ logo in conjunction with tested products

➤ Legal protections

- ◆ Terms of Use: conditions on use of test results agreed to by those downloading results
- ◆ Terms of Submission: agreed to by vendor submitting test results

➤ Sign up for the mailing list: sniaemerald.com

Why Should Storage Vendors Use the Emerald Program?

➤ SNIA Emerald Program seeks to

- ◆ Encourage storage vendors to build better products
- ◆ Stimulate the IT community to more rapidly deploy and operate multi-vendor storage technology efficiently

➤ SNIA Emerald Program

- ◆ Provides a level playing field for test sponsors
- ◆ Produces results that are powerful and yet simple to use
- ◆ Provides value for vendors as well as IT consumers and solution integrators
- ◆ Reports results in a manner that is easy to submit , audit and verify

Why Should Storage Consumers Use the Emerald Program?

- ▶ **SNIA Emerald Program seeks to**
 - ◆ Provide a collection of standard metrics and data that allows IT architects to objectively compare a range of possible storage solutions
- ▶ **SNIA Emerald Program**
 - ◆ Enables users to select the mode of storage usage that accomplishes their work objectives with the lowest overall energy consumption
 - ◆ Drives vendor companies to innovate and compete in the development of energy efficient products as measured by the standard yardsticks

Prepare Test Data Report

- Test Data Results may be submitted to the Emerald Program for publication on the sniaemerald.com web site
- Download Test Data Report template from sniaemerald.com using the “Documents and Downloads” menu item
- Complete the template with information on your product and your test results
 - ◆ Page 1 – basic information about the product and vendor company
 - ◆ Page 2 – report test metrics - exercise care to provide data in the appropriate units and precision
 - ◆ Pages 3-8 – additional product information - provide as much detail as you wish
- Note: after some administrative edits, this spreadsheet will be converted to a PDF verbatim and published – so make sure that the spreadsheet as submitted represents your product and company appropriately

Test Data Submission - I

- You must be registered and logged-in to the sniaemerald.com web site in order to submit a test result
 - ◆ If you are not registered as a user, register using the “create an account” menu item on the lower left side of sniaemerald.com
 - ◆ If you are not logged-in, log-in with you user name and password on the lower left side of sniaemerald.com
- Upload your completed Test Data Report spreadsheet to sniaemerald.com using the “Submit Test Data Results” menu item
- Fill in web form with basic information on your company and the tested product
 - ◆ Note: The “submitter” is the person making the submission; the “submitter company” is his/her company
 - ◆ Note: The “product vendor” is the company whose product was tested

Test Data Submission - 2

- Agree on behalf of your company to the legal “Terms of Submission” that are your companies agreement to the policies of the Emerald Program
 - ◆ Note: you may wish to review these terms in advance of doing the submission
- After processing by the Emerald Program, your results will be published on sniaemerald.com and will be available for download by anyone agreeing to the “Terms of Use”
 - ◆ If there are any issues with your submission, you will be contacted by the Emerald Program Director
 - ◆ Submissions are published in **Provisional** status for 2 months
 - ◆ If no serious issues arise regarding the submission within 2 months, it is advanced to **Accepted** status
- If you have problems or questions regarding the submission process or the Test Data Report template, contact the Emerald Program Director at emerald@snia.org

Thank You

- Questions
- Please fill out the online survey
 - ◆ <http://www.surveymonkey.com/s/EmeraldTraining>
- Copy of slides at
 - ◆ <http://www.sniaemerald.com>
- Email questions to:
 - ◆ Emerald Program
 - > Emerald@SNIA.org
 - ◆ Green Storage Initiative
 - > GSI@snia.org
 - ◆ Green TWG
 - > greentwg-chair@snia.org

- Storage Networking Industry Association
 - ◆ <http://www.snia.org>
- Green Storage Initiative
 - ◆ <http://www.snia.org/forums/green>
 - ◆ Green tutorials and white papers
- SNIA EmeraldTM Program (downloads)
 - ◆ <http://www.sniaemerald.com/>
 - ◆ Measurement Specification
 - ◆ User guide
 - ◆ Download test data reports
 - ◆ Submission process and requirements
 - ◆ Training materials

