

**COPAN**  
SYSTEMS

# MAID for Active Archive Data

**Aloke Guha**

**CTO, COPAN Systems**

**aloke.guha@copansys.com**

**Panel on Emerging Technologies**

**NASA/IEEE MSST 2004**

**12th NASA Goddard/21st IEEE Conference on  
Mass Storage Systems & Technologies**

**The Inn and Conference Center**

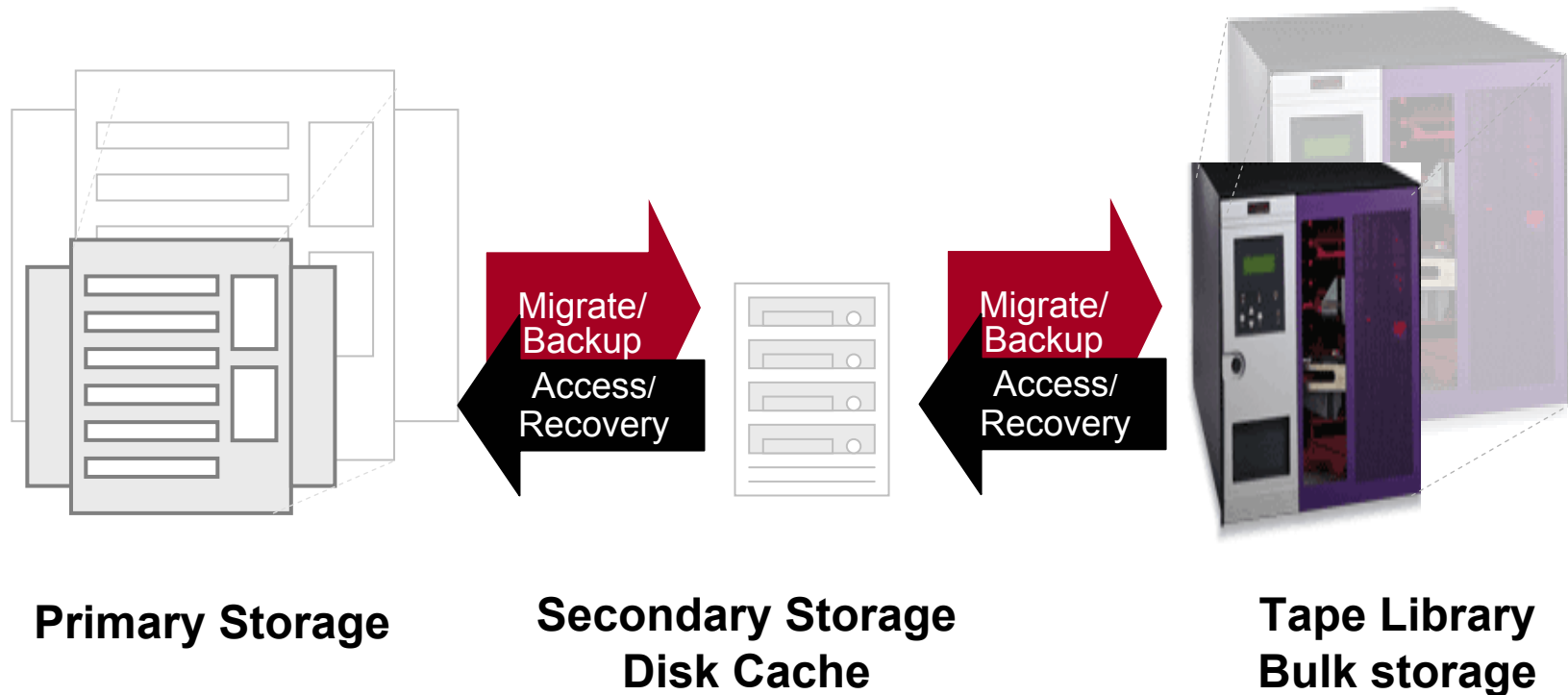
**University of Maryland University College**

**Adelphi MD USA**

**April 13-16, 2004**

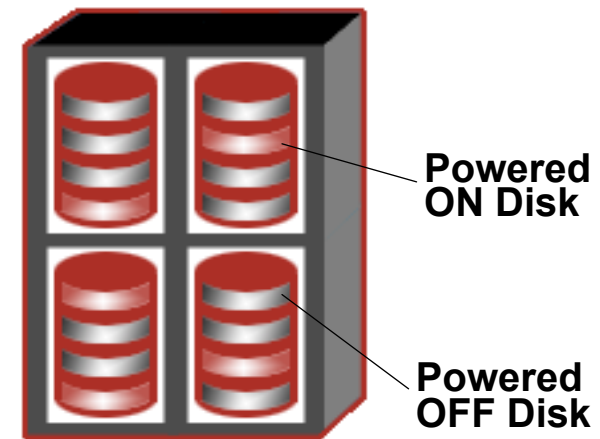


# Today's Hierarchy: Go Slow or Pay Up



# MAID . . . . Power-Managed Disk

- **Large number of power-managed drives**
  - More than 50% drives powered OFF
  - Power-cycling by policy for application
- **Benefits**
  - Scale
  - Cost
  - Service life
  - Energy
- **Ref: Colarelli and Grunwald, FAST 2002, SC 2002**
  - Tradeoff of disk power vs. performance
  - Virtualization of ON (cache) and OFF drives, RAID-0
  - Caching not beneficial for archive workload



# MAID Applicability

Content Type	Write	Update	Read	Technology	Metric
<b>Dynamic</b>	Many	Yes	Many	Disk	IOPs
<b>Active Archive</b>	Once	No	0 to n	MAID	Bandwidth
<b>Deep Archive</b>	Once	No	Rarely	Tape	Cost

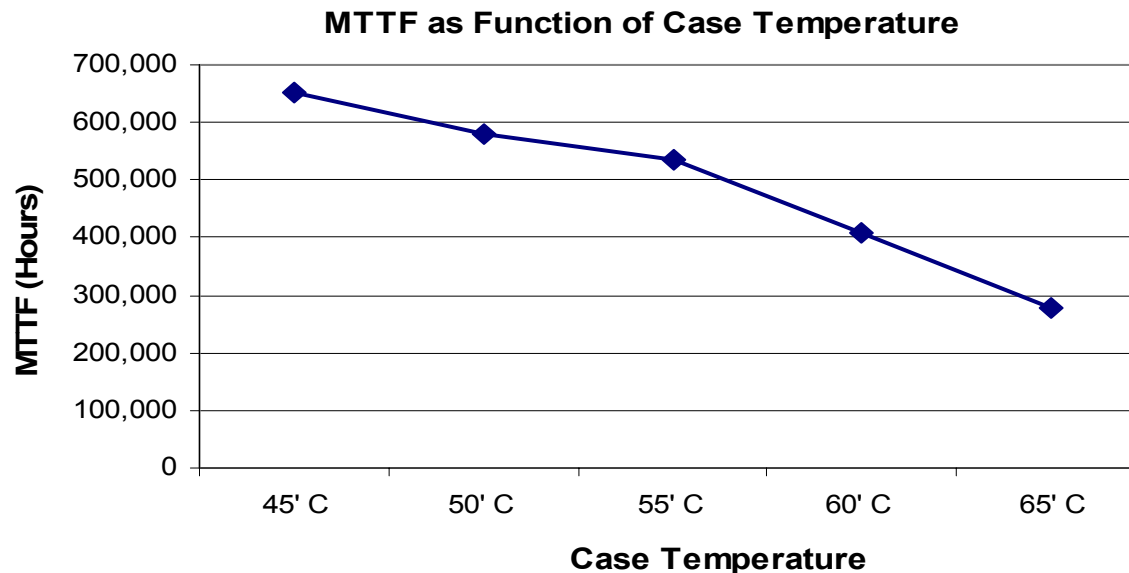
# Scale: Storage Capacity

- **Benefit**

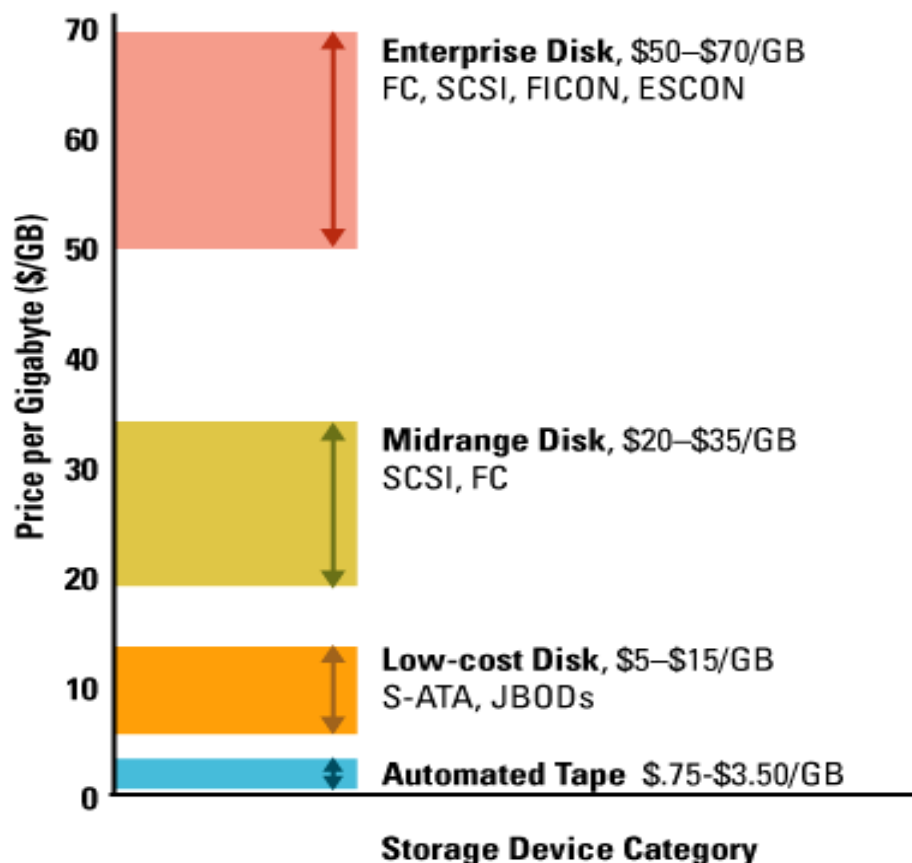
- Large Number of Drives/Single System Footprint
- 0(1000) drives  $\Rightarrow$  250TB - 400TB

- **Needs**

- High-Density Interconnect Architecture
- Manage environmental conditions



## Disk and Tape Pricing Guidelines



- > The price per gigabyte decreases as the ratio of cartridges to drives increases, diverging from disk costs
- > Disk prices are for working subsystems
- > Automated tape prices include drives, media and library
- > Tape cartridge capacity growing faster than disk drive capacity
- > Automated tape nominally about 1/15 to 1/20 the price of disk for Unix, Linux, Win2K; 1/25 or less for mainframes.

Source: Horison Information Strategies

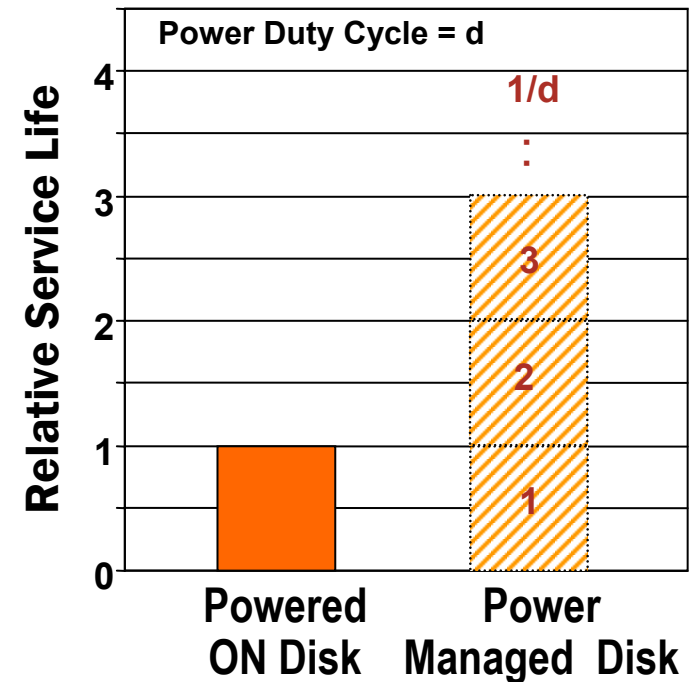
# Understanding Storage Costs

- **Non-media storage cost far exceeds media \$/GB**
- **Cost Efficiency:** 
$$\frac{\text{Media Cost}}{\text{Storage System Cost}}$$
- **Disk Example**
  - **250GB SATA ~\$1/GB vs Storage: \$5-\$15/GB**
  - **Cost Efficiency: 0.07 – 0.2**
- **Tape Example**
  - **200GB LT02 ≤ \$0.5/GB vs Storage: \$0.75-\$3.5/GB\***
  - **Cost Efficiency: 0.14 – 0.67**
- **Traditional disk systems 3x or more cost of tape\***
- **MAID levels playing field between Disk and Tape!**

\*Native uncompressed capacity: cost/GB depends on ratio of cartridges to drives, typ. 20:1 – 80:1  
Assume same compression applied to data on disk or tape

# Reliability: Drive Service Life

- **Effective drive service life**
  - **MTBF  $\propto$  1/Annual Failure Rate**
  - **AFR  $\propto$  Power On Hours**
- **Service life  $\propto$  1/(Power Duty Cycle)**
- **Data Rel.  $\propto$  1/(Power Duty Cycle)<sup>2</sup>**
- **Needs**
  - **Manage start stops  $\leq$  50K**
  - **Data protection and integrity**



Power duty cycle = # of powered-ON drives/# of powered-OFF drives



# Managing Start Stops

- **Bandwidth/capacity limits SS to 3% of max<sup>1</sup>**
- **Archive #mounts/volume limits SS to <5 % of max<sup>1, 2</sup>**

	Industry	Volumes Used	Daily #Mounts/Volumes Used		
			Average	Max	Median
1	Telco	373	0.0	0.1	0.0
2	Telco	1,015	0.1	0.4	0.0
3	Telco	688	0.1	0.5	0.0
4	Telco	1,189	0.0	0.0	0.0
...					
33	Utility	278	2.6	4.9	3.5
34	Govt	3,393	0.4	0.7	0.5
35	Govt	84	0.1	0.4	0.1
...					
	<b>Average</b>	<b>1,122</b>	<b>0.6</b>	<b>1.1</b>	<b>0.6</b>

## Ave # Start-Stops over 5 yr. ops

	Capacity (TB)		
<b>Bandwidth (TB/hr)</b>	<b>150</b>	<b>200</b>	<b>250</b>
<b>2</b>	<b>584</b>	<b>438</b>	<b>350</b>
<b>3</b>	<b>876</b>	<b>657</b>	<b>526</b>
<b>4</b>	<b>1168</b>	<b>876</b>	<b>701</b>

**Typical Specified Limit: 50K**

<sup>1</sup>Over 5-year period

<sup>2</sup>Source: SW Vendor - data from 43 archives on tape: Volumes\_Used excludes tape volumes not allocated in ATL

# Performance Implications

- **Bandwidth increases as power duty cycle**
- **Access time depends drive state, power duty cycle**
- **Needs**
  - **Limit duty cycle, else increase device failure rates**
  - **Optimize overall architecture for performance wrt cost**

<b>Device Type</b>	<b>Load (secs)</b>	<b>First Byte (secs)</b>	<b>File Access (secs)</b>	<b>Unload (secs)</b>	<b>Total/File* (secs)</b>
<b>Single Drive</b>	<b>&lt;10</b>	<b>0.1s</b>	<b>10-12</b>	<b>0.1</b>	<b>12</b>
<b>RAID(n)</b>	<b>&lt;10*n</b>	<b>4s</b>	<b>14-15</b>	<b>0.1</b>	<b>≥15</b>
<b>Tape Drive</b>	<b>18</b>	<b>41</b>	<b>59</b>	<b>18</b>	<b>77</b>

\* Transfer time depends on size of RAID set

# Optimized MAID Fills the Gap in Hierarchy

