

EBU Storage Workshop 21-22 Nov 2011

Media Application - File System - Storage Behaviour

Relevance of Specifications

Luc Andries - CTO CandIT-media (a VRT-media-lab spinoff)

Setting the Scene

Relevance of Specifications

- **Production Media Storage**

(Editing, Transcoding, Rewrapping, ..., ~ File-based Production – **???Media Cloud ???**)



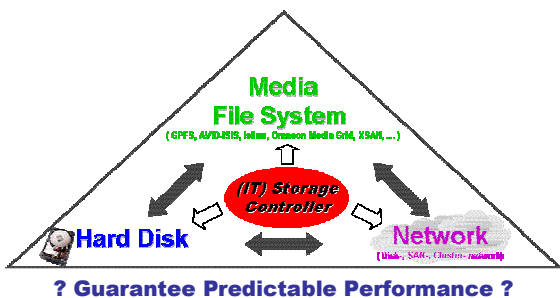
- **Guaranteed Predictable Performance**

Setting the Scene

Relevance of Specifications

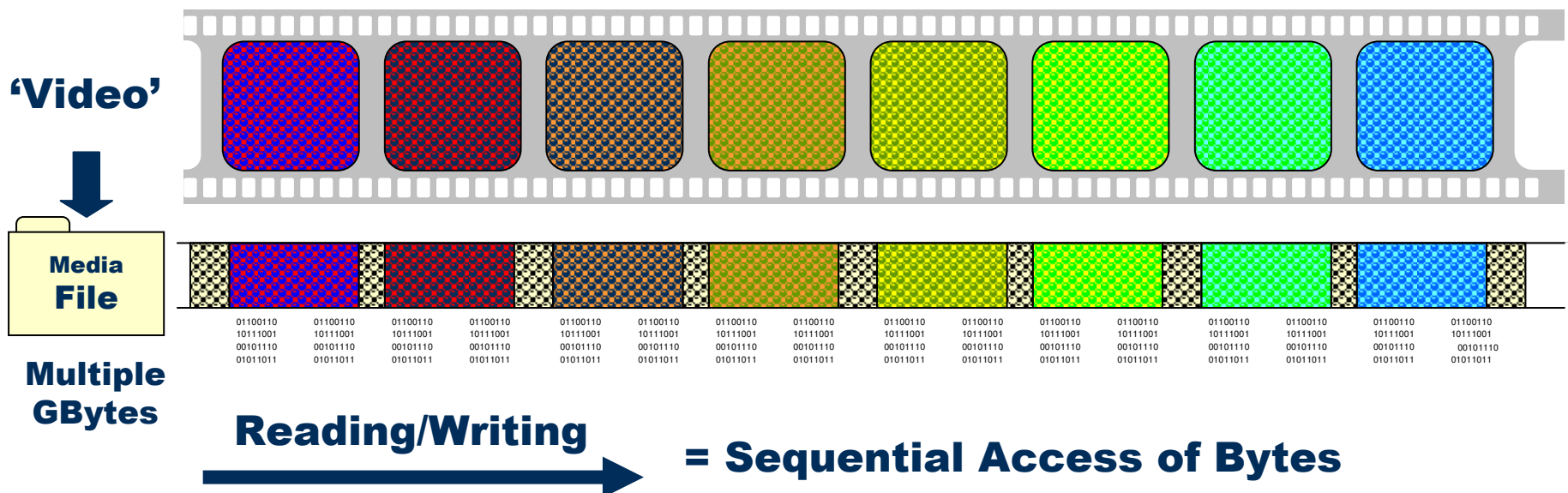
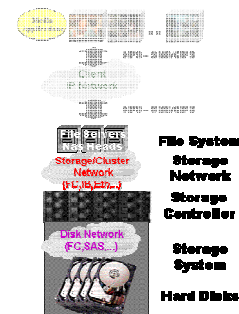
- **Media Application Behaviour**
- **Media File System Behaviour**
- **Media Storage Behaviour**

(Media) Application Behaviour



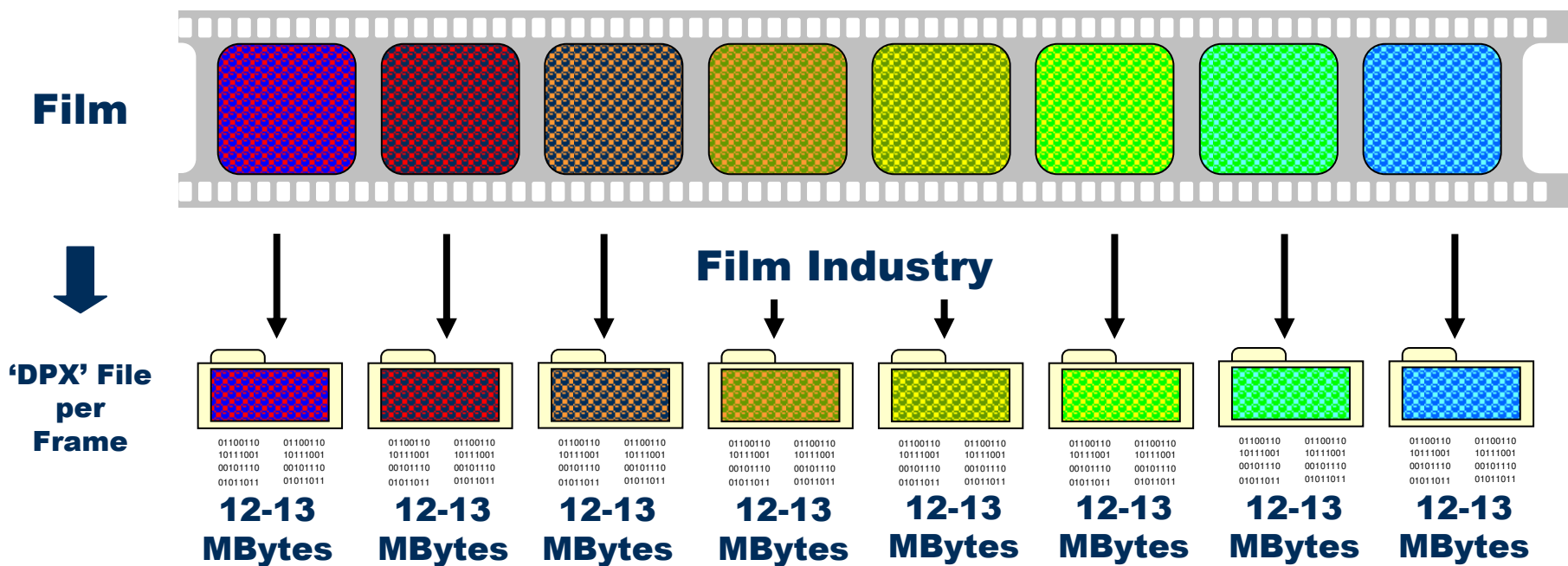
Common View of Media-File Access

“Sustained Sequential Access” of Large File(s)



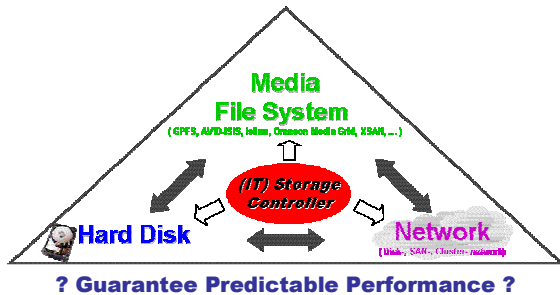
? How are these Bytes stored on Media Storage ?

Broadcasting vs Film Industry



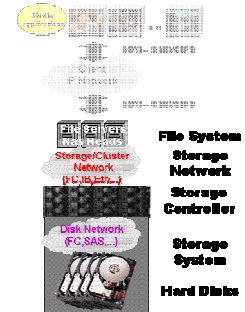
**“Sustained *Sequential* Access”
of ‘Large’ (?) Files**

(Media) Application Behaviour



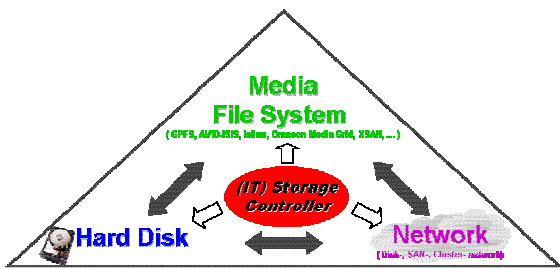
Media Application Behaviour

~

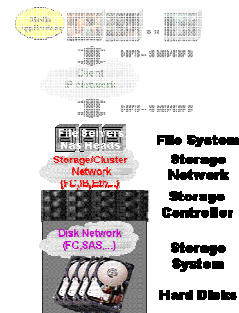


- Access 'Large' Files
- At High Throughput
- In a Multi-User/Multi-Process Context
- At Sustained Load

(Media) Storage Specification

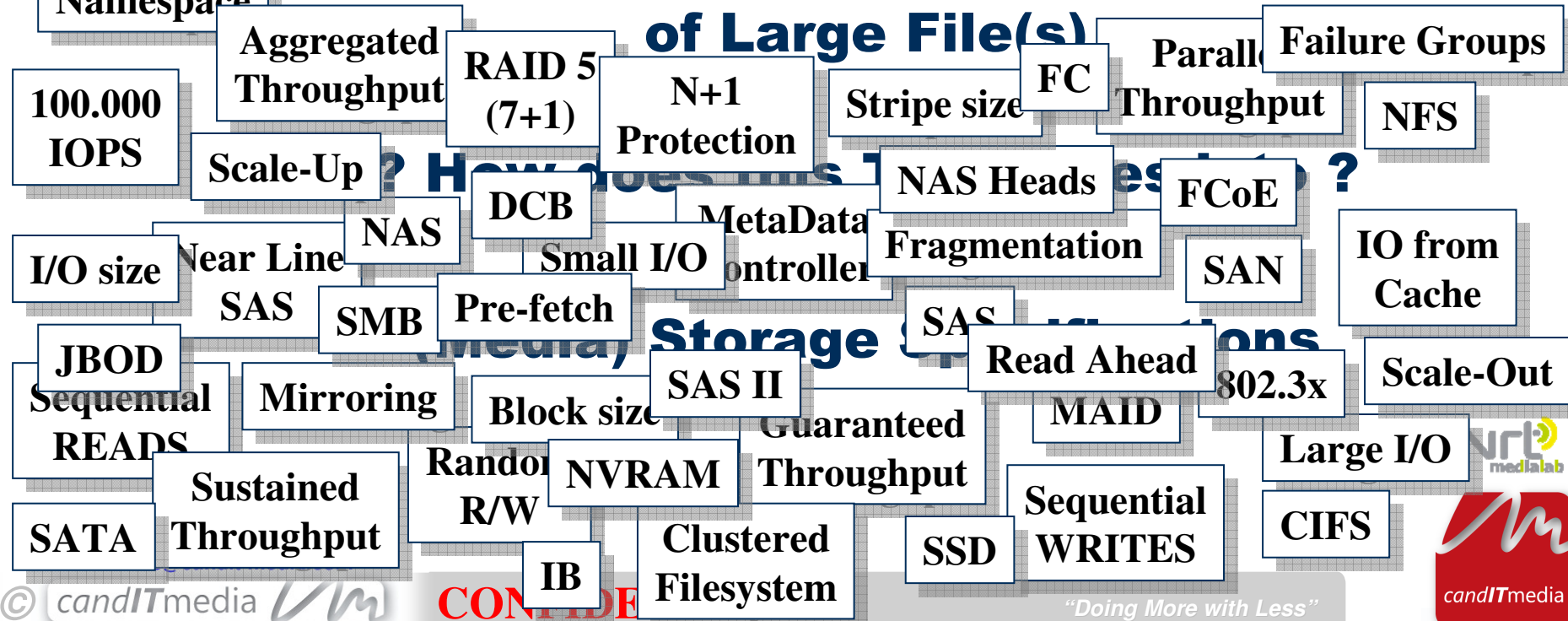


Common Perception of Media-File Access

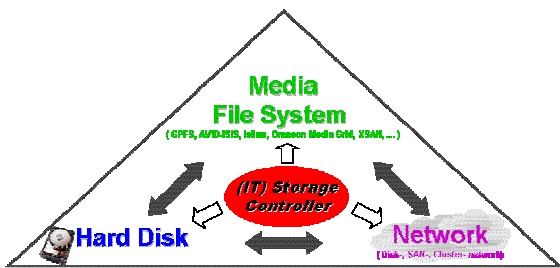


? Guarantee Predictable Performance ?

"Sustained Sequential Access" of Large File(s)

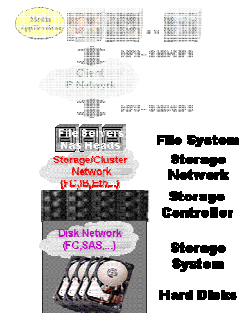


(Media) Storage Specification

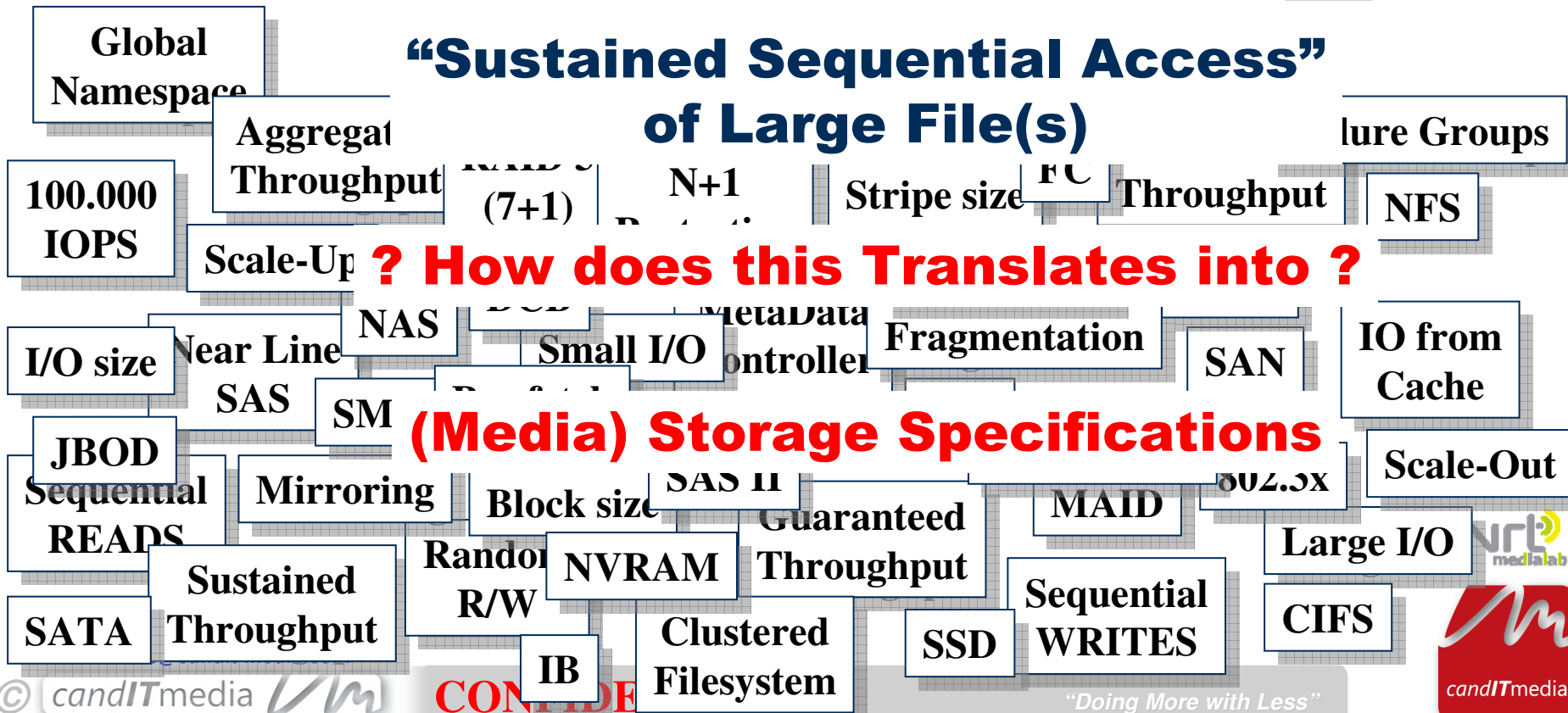


? Guarantee Predictable Performance ?

Common Perception of Media-File Access



"Sustained Sequential Access" of Large File(s)



(Media) Storage Specifications

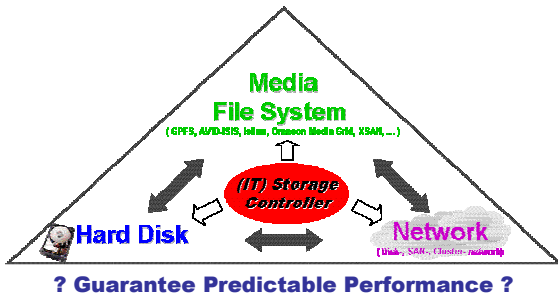
threaded architecture allows performance to linearly scale with advances in underlying hardware. This same architecture allows the SFA10000 to do what no other RAID controller has been able to do to date: to perform in the extreme range of both throughput and IOPS. The SFA10000 delivers random over 1 million burst IOPS up to and 300,000 random to disk. Sequential block throughput performance is 10GB/s for both reads and writes. Designed to house the most scalable unstructured file data, the system supports up to up to 1,200 drives raw storage while enabling a combination of SAS, SATA or SSD drives.

...onal primary file storage, providing over 1.4 million NFS ops and 85 gigabyte/sec (GBps) of aggregate throughput, all from a single file system. To accomplish these extraordinary speeds, the S200 combines the power of ultra-high performance solid state drives (SSD) and 10,000 RPM 2.5-inch Serial Attached SCSI (SAS) drive technology. It also features for dual 1GbE and dual 10GbE front-end networking, dual quad core Intel CPUs, a high-performance InfiniBand back-end, and up to 13.8 TB of globally coherent cache. In addition, the Isilon S200 leverages enterprise SSD technology to accelerate namespace-intensive metadata operations, and, also, provides an ability to place the mission critical, latency-sensitive data on SSDs in a SmartPools™ environment.

parallel architecture and a hardware accelerate sustained and predictable performance, innovative performance, quality of service, data protection, accelerate your high-concurrent and sequential-throughput applications. deploy and manage fewer systems to meet today's changing data patterns and the explosive growth in capacity and performance requirements. Significantly lowering both acquisition costs and total cost of ownership. ideal for performance driven, real-time streaming applications to 6GB/s of sustained throughput; and to consolidate archives, efficiently storing up to 3.6 petabytes in ju

Effortless real time performance and efficient bandwidth delivery set ISIS 7000 apart; With each ISIS 7000 Engine providing up to 400 MB/s of useable bandwidth, predictable, linear performance is sustained even as client counts increase, right up to 4.8 GB/s.

(Media) Application Behaviour



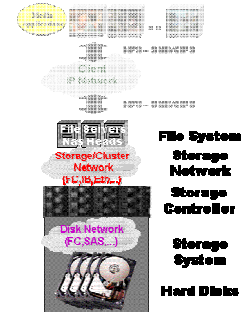
Media Application Behaviour

~

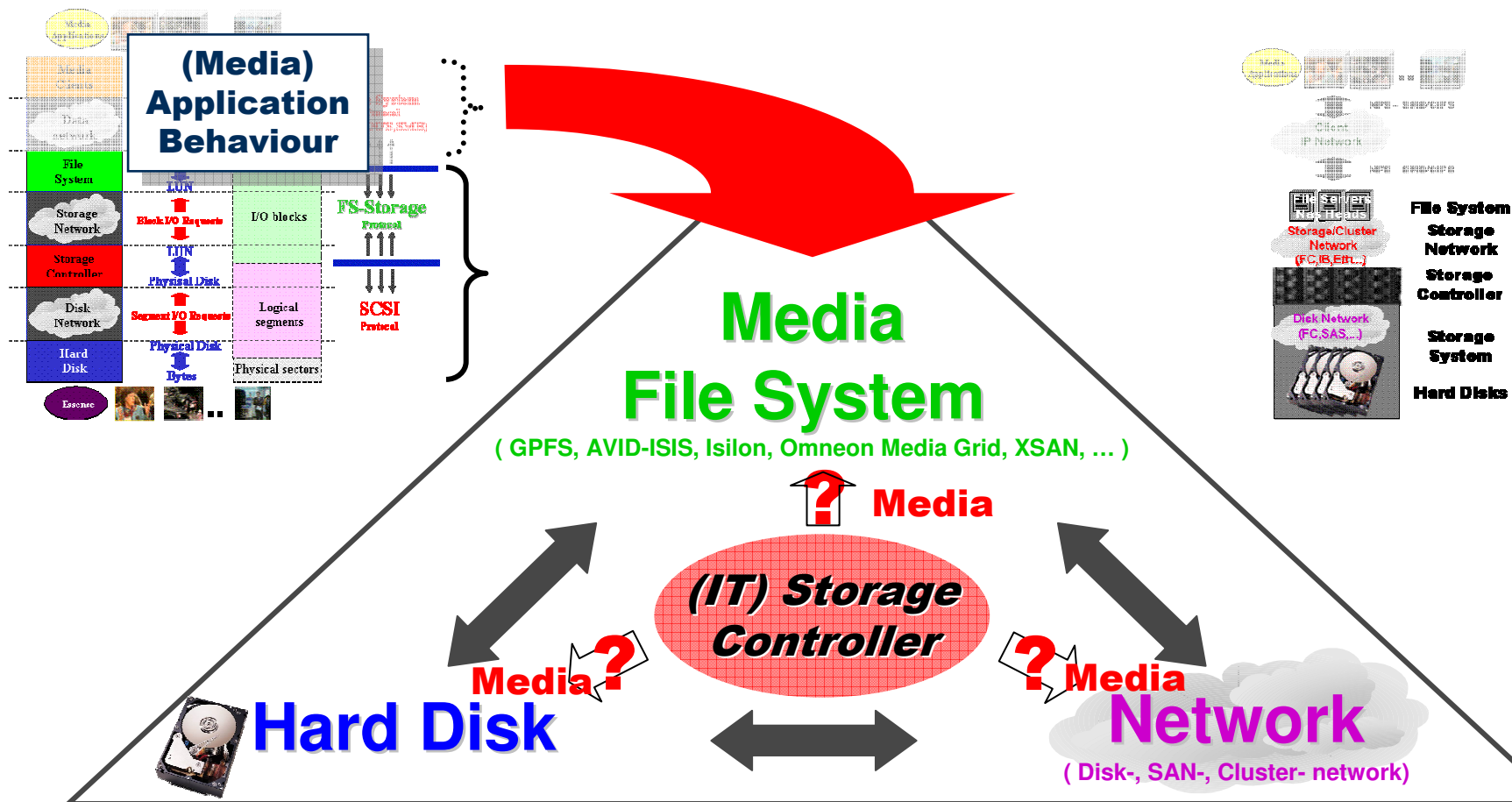
Sequential Access

Versus

- Access 'Large' Files
- At High Throughput
- In a Multi-User/Multi-Process Context
- At Sustained Load

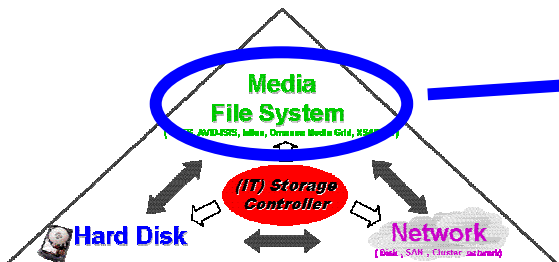


The Basic Triangle



? Guarantee Predictable Performance ?

(Media) File System Behaviour

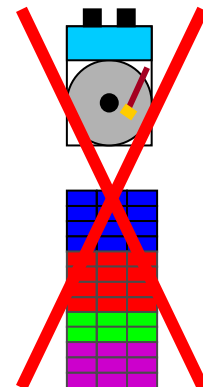
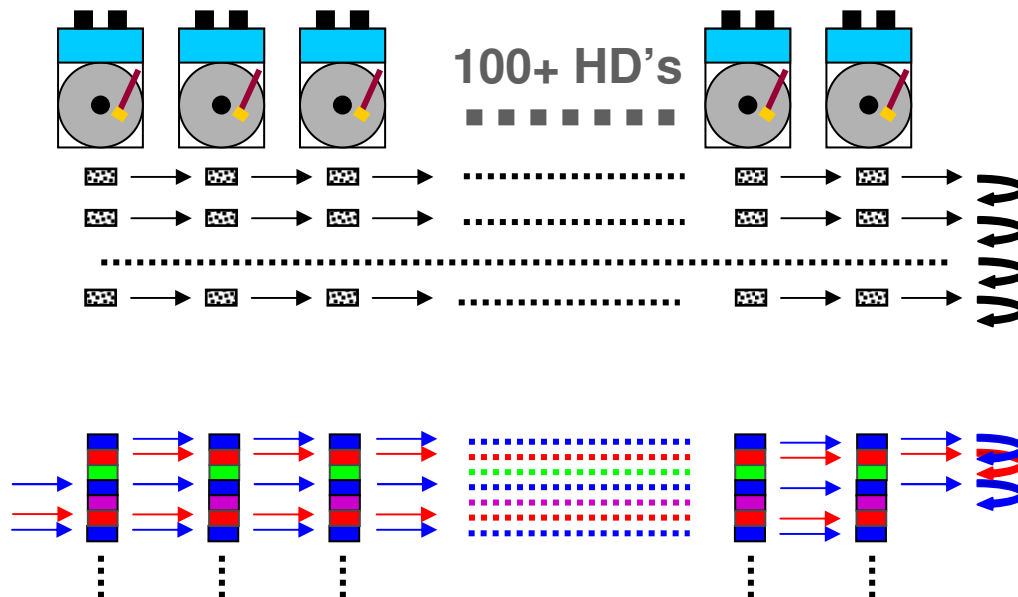


? Guarantee Predictable Performance ?

// Striping:

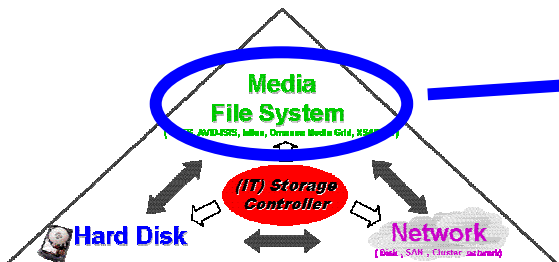
- Access 'Large' Files
- At High Throughput
- At Sustained Load
- In a Multi-User/Multi-Process Context

Process 1: Media File
Process 2: Media File
Process 3: Media File
Process 4: Media File



Media File System **Write** Behaviour Leads to
Natural Fragmentation of Disks

(Media) File System Behaviour



Media File System
Write Behaviour Leads to
Natural Fragmentation of Disks

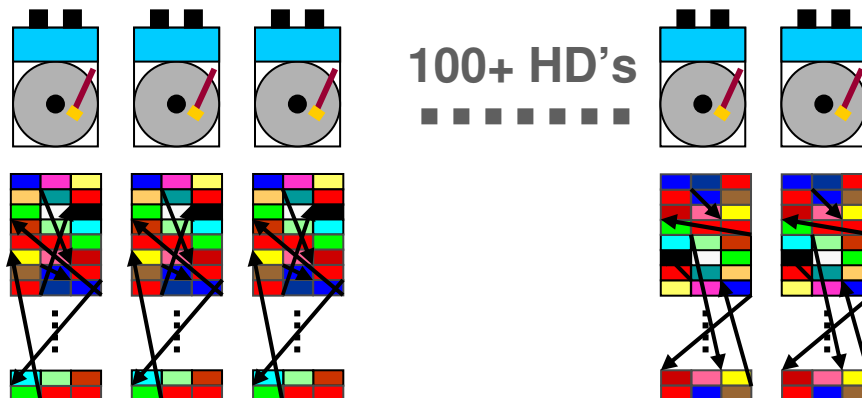
- Access 'Large' Files
- At High Throughput
- In a Multi-User/Multi-Process Context
- At Sustained Load

? Guarantee Predictable Performance ?

// Striping:

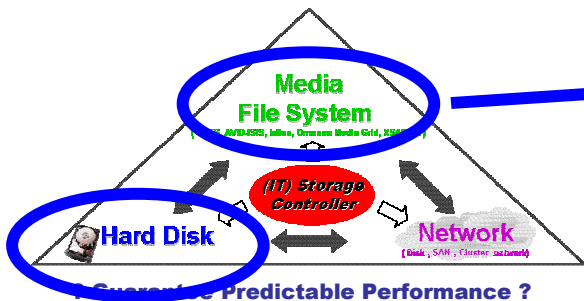
Multiple Reading Processes:

Process 1: Media File
Process 2: Media File
Process 3: Media File
...
Process n: Media File



Natural Fragmentation Leads to a
Random Access
Media File System Read Behaviour
on each Individual Disk

(Media) File System Behaviour



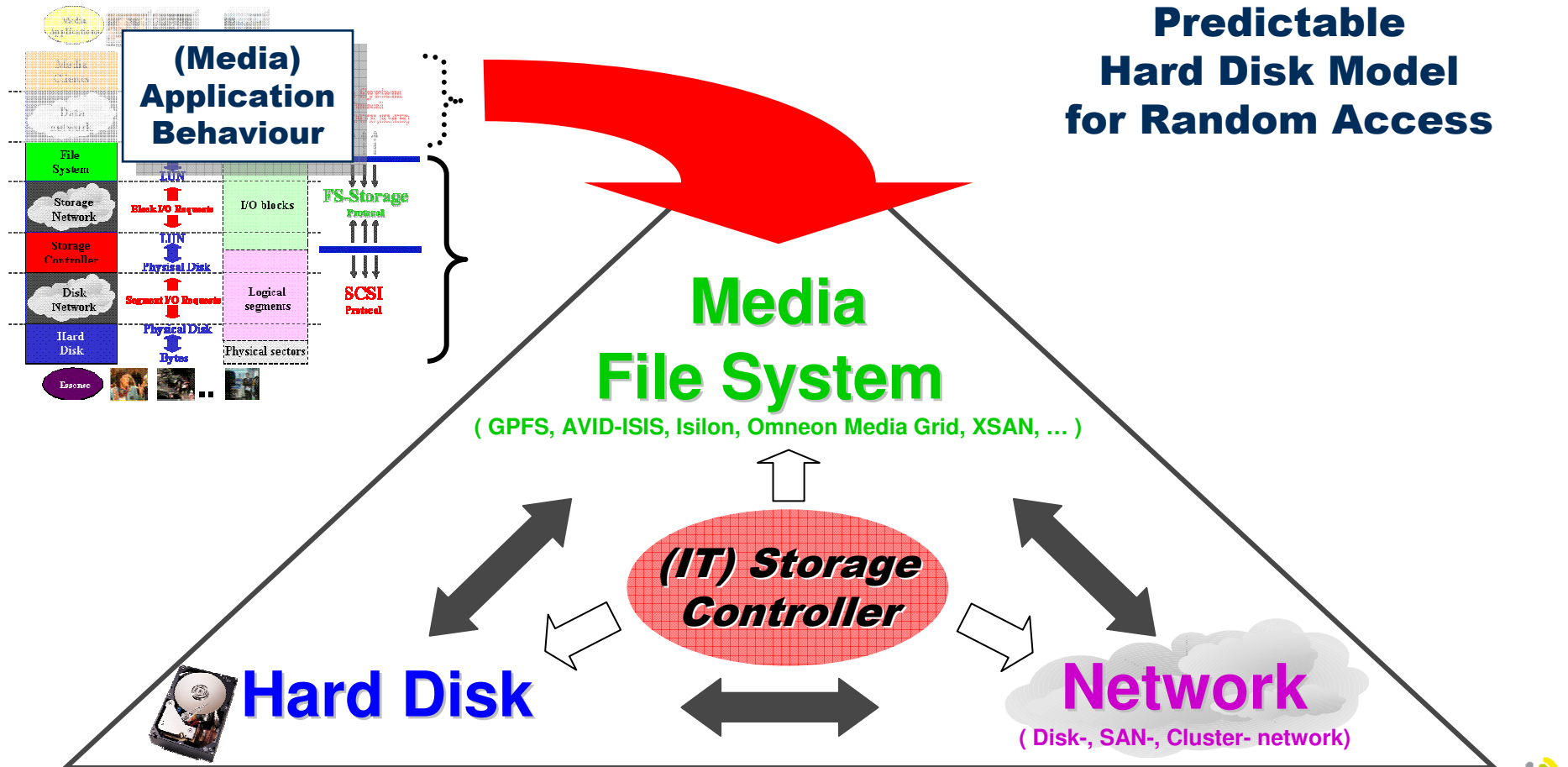
Media File System
Write Behaviour Leads to
Natural Fragmentation of Disks

- Access 'Large' Files
- At High Throughput
- In a Multi-User/Multi-Process Context
- At Sustained Load

Natural Fragmentation Leads to a
Random Access
Media File System **Read** Behaviour
on each Individual Disk

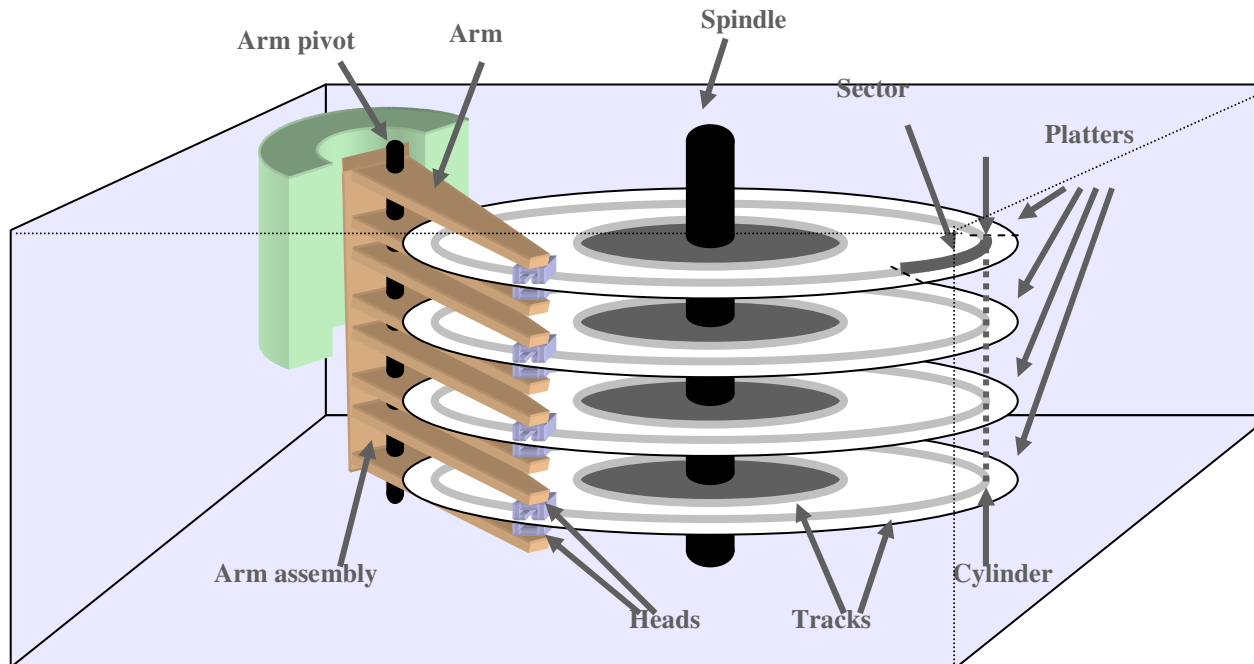
Predictable Media Storage Performance
Requires a
Predictable Hard Disk Model
for **Random Access**

A Hard Disk Model



A Hard Disk Model

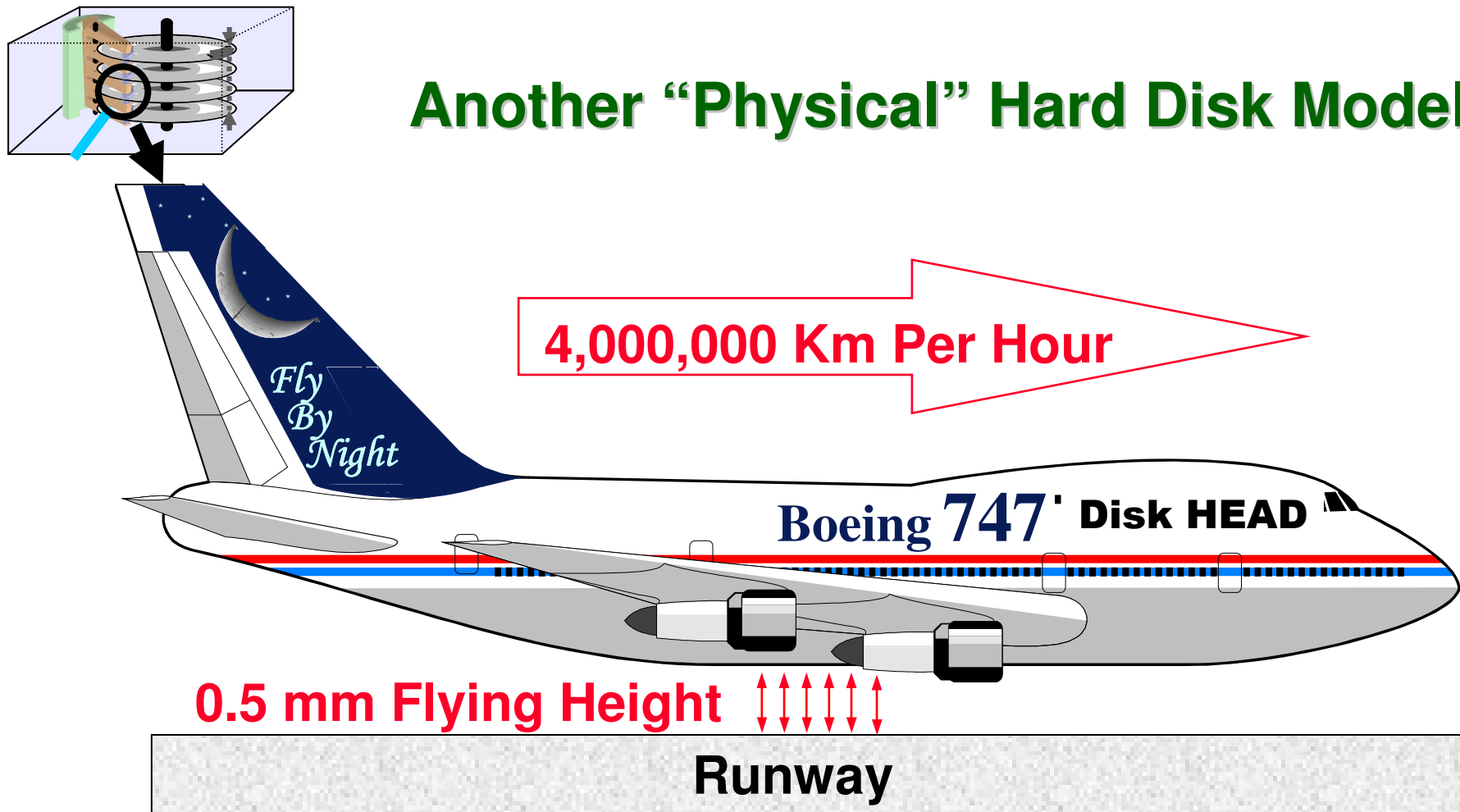
A Physical Hard Disk Model



Don't do this at Home !!!

A 'Physical' Disk Model

Another "Physical" Hard Disk Model



Sideways movement @ 80 g

A Hard Disk Model

SAS vs SATA



Different name comes from **Interface Technology**

SCSI (Small Computer System Interface)
FC(-AL) (Fibre Channel Arbitrated Loop)
SAS (serial attached SCSI)



ATA (Advanced Technology Attachment)
IDE (Integrated Disk Electronics)
SATA (serial ATA)

A Hard Disk Model

SAS vs SATA



Different specifications come from **Target Application Environment**

High End Systems / Reliability
In Groups / Racks / Shelves
Random Access / Performance
Heavy Duty 24h/24h 7d/7d



Low Cost (PC's)
Single Disk in PC
Personal Use
8h/d 5d/w

A Hard Disk Model

SAS vs SATA



Different specifications come from **Disk Technology**

2.5" Platter diameter
Spinning speed 15 Krpm
Bearings at both shaft ends
Separate processors (Servo - I/O)
Stiffness / Mechanical rigidity
Stronger actuators / Motors
Filters / O-ring sealings / Dessicants



3.7" Platter diameter
Spinning speed 7.2 Krpm
Bearings at single end
Single processor
Less robust

A Hard Disk Model

SAS vs SATA



Differences result in **different**

Data Rate

Reliability
MTBF

Data Density

Capacity

Applications
Environment

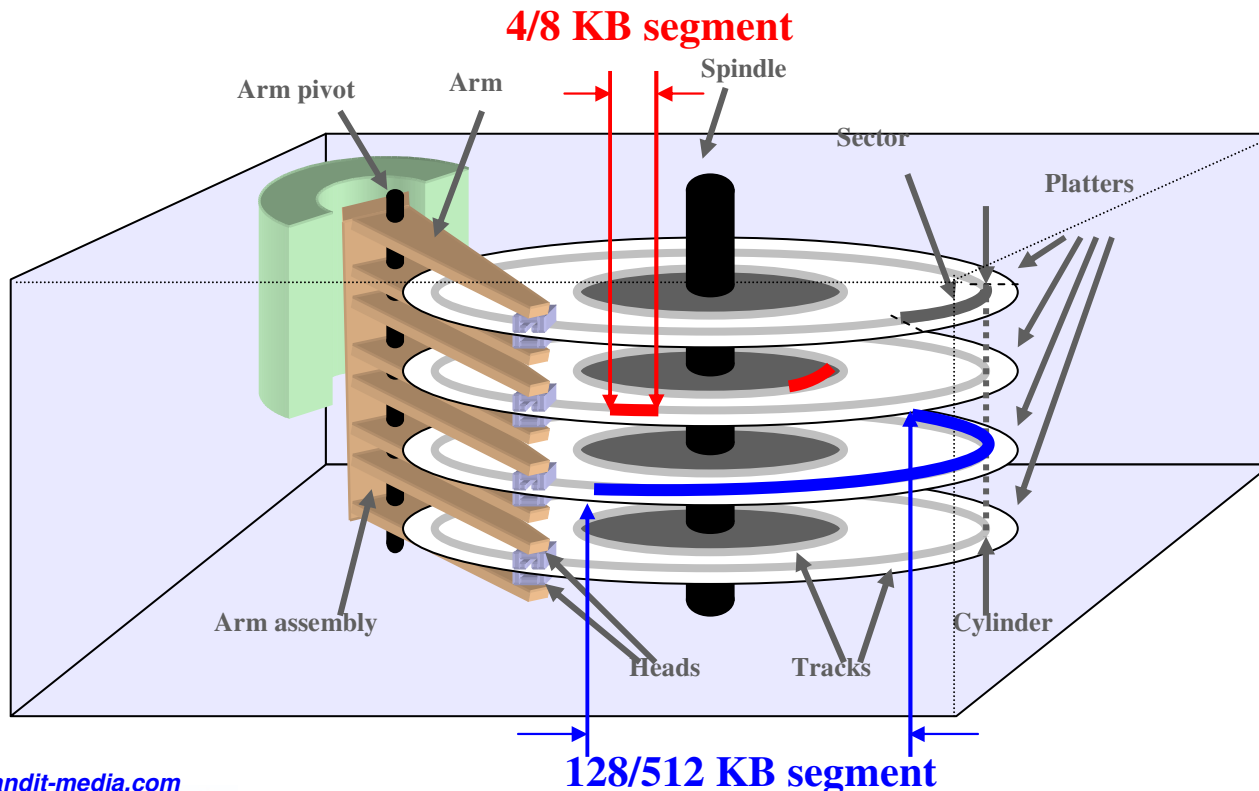
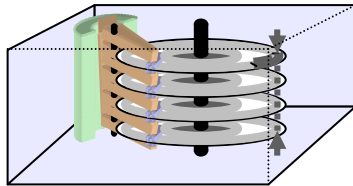
Cost

Random Access
Performance

A 'Physical' Disk Model

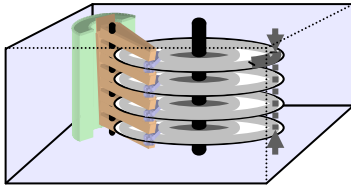
Relevant Parameters

- Parameters:
- Lateral seek time
 - Rotational speed
 - **Segment size**
- } = Head positioning
- = Reading/Writing time



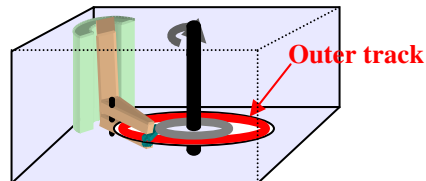
A 'Media' Hard Disk Model

Random Access Behaviour

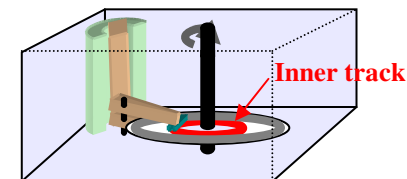


- Parameters:
- Lateral seek time
 - Rotational speed
 - **Segment size**
- } = Head positioning
- = Reading/Writing time

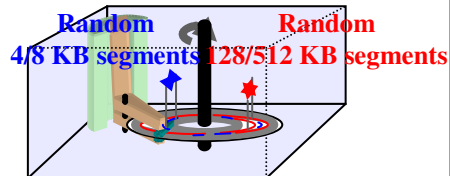
R/W on the Outer Track



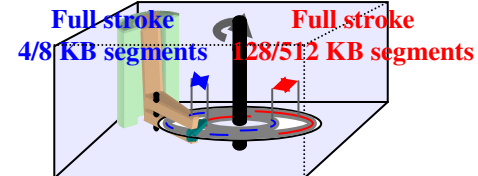
R/W on the Inner Track



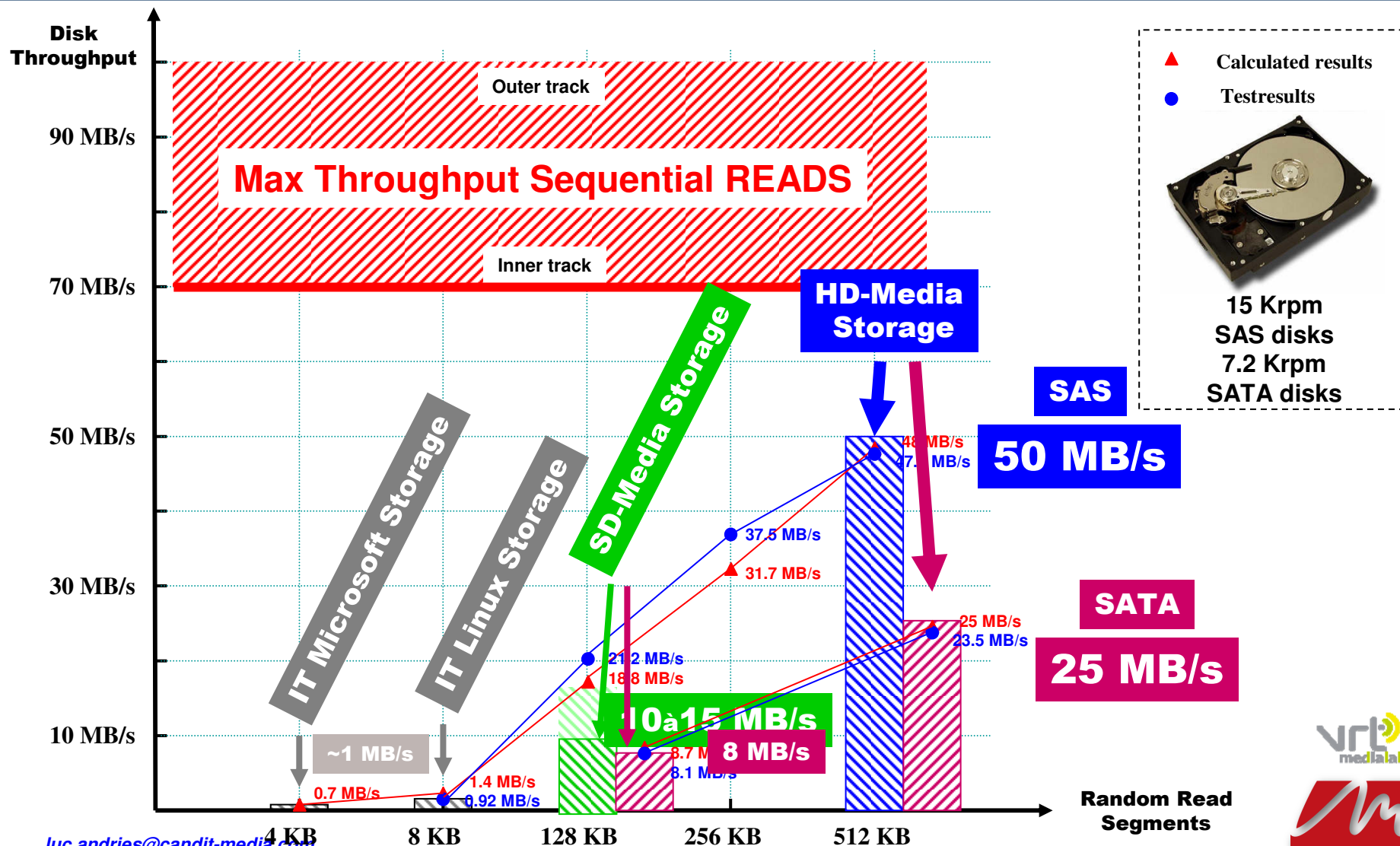
Random R/W
4/8KB – 128/512 KB Segments



Full Stroke R/W
4/8KB – 128/512 KB Segments

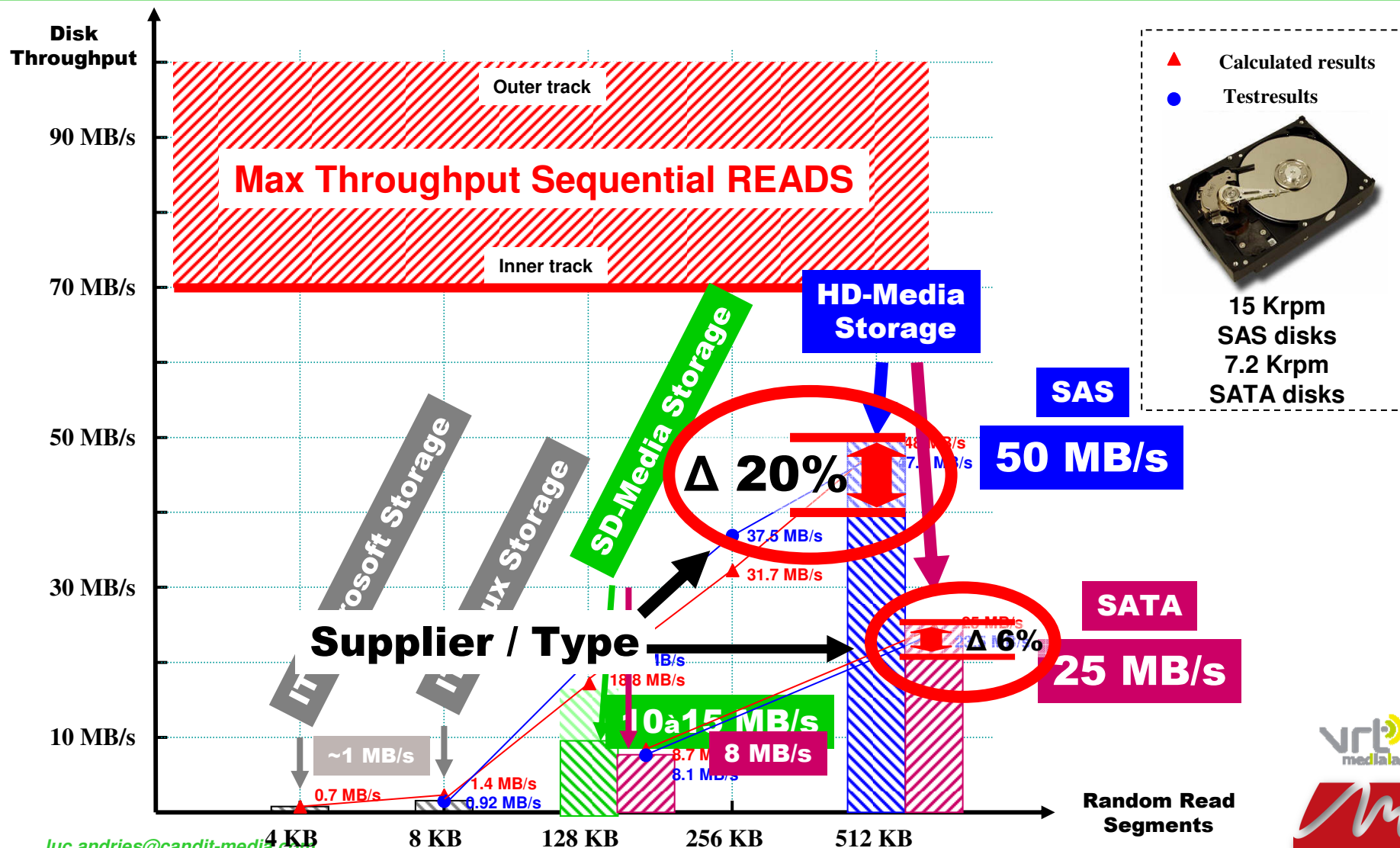


Some Hard Figures

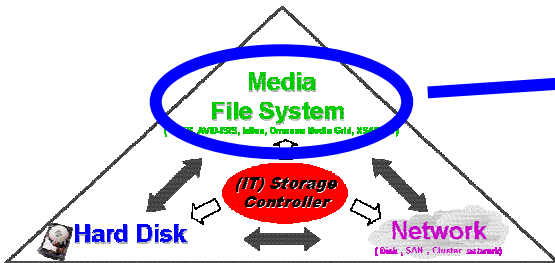


luc.andries@candit-media.com

Some Remarkable Figures



(Media) File System Behaviour



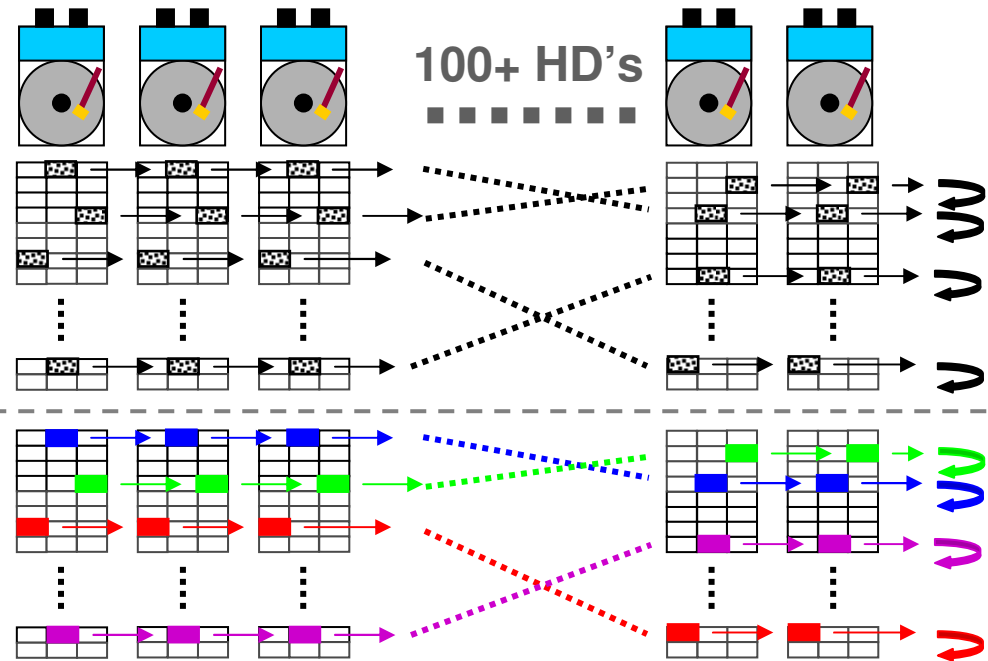
Media File System
should Force **Scattered**
Write Behaviour

- Access 'Large' Files
- At High Throughput
- In a Multi-User/Multi-Process Context
- At Sustained Load

? Guarantee Predictable Performance ?

// Striping:

- Access 'Large' Files
- At High Throughput
- At Sustained Load
- **Force scattered Writing Pattern**



- In a Multi-User/Multi-Process Context

Process 1: Media File

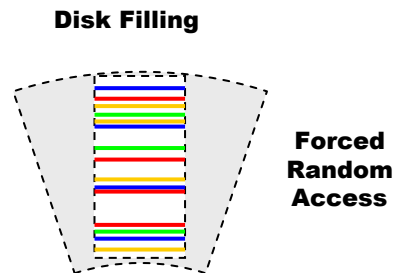
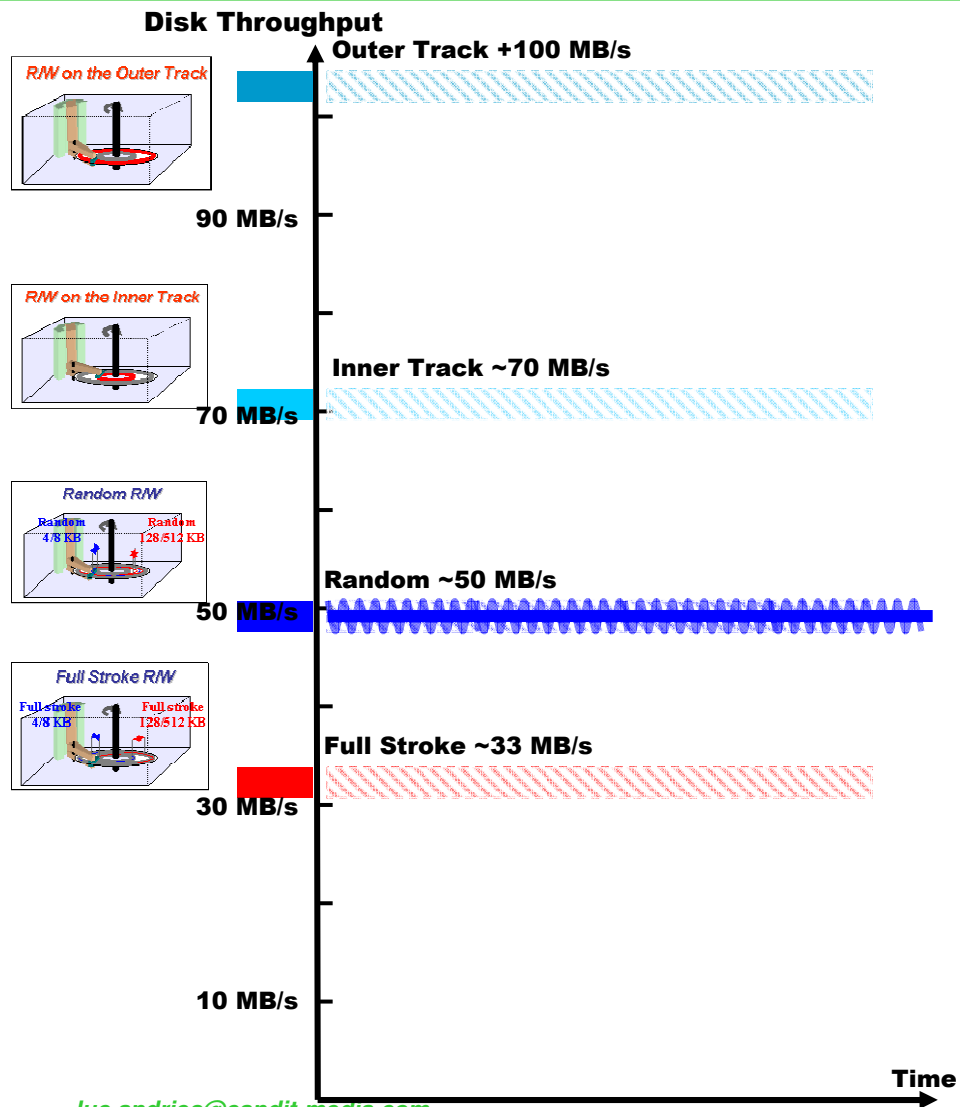
Process 2: Media File

Process 3: Media File

Process 4: Media File

Media File System has to Force Scattered Writing Pattern
to get Predictable Disk Reading/Writing Behaviour

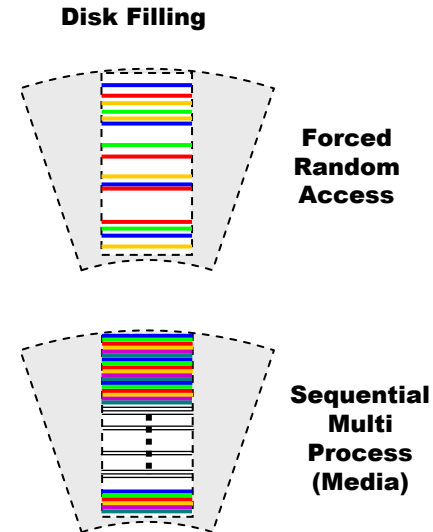
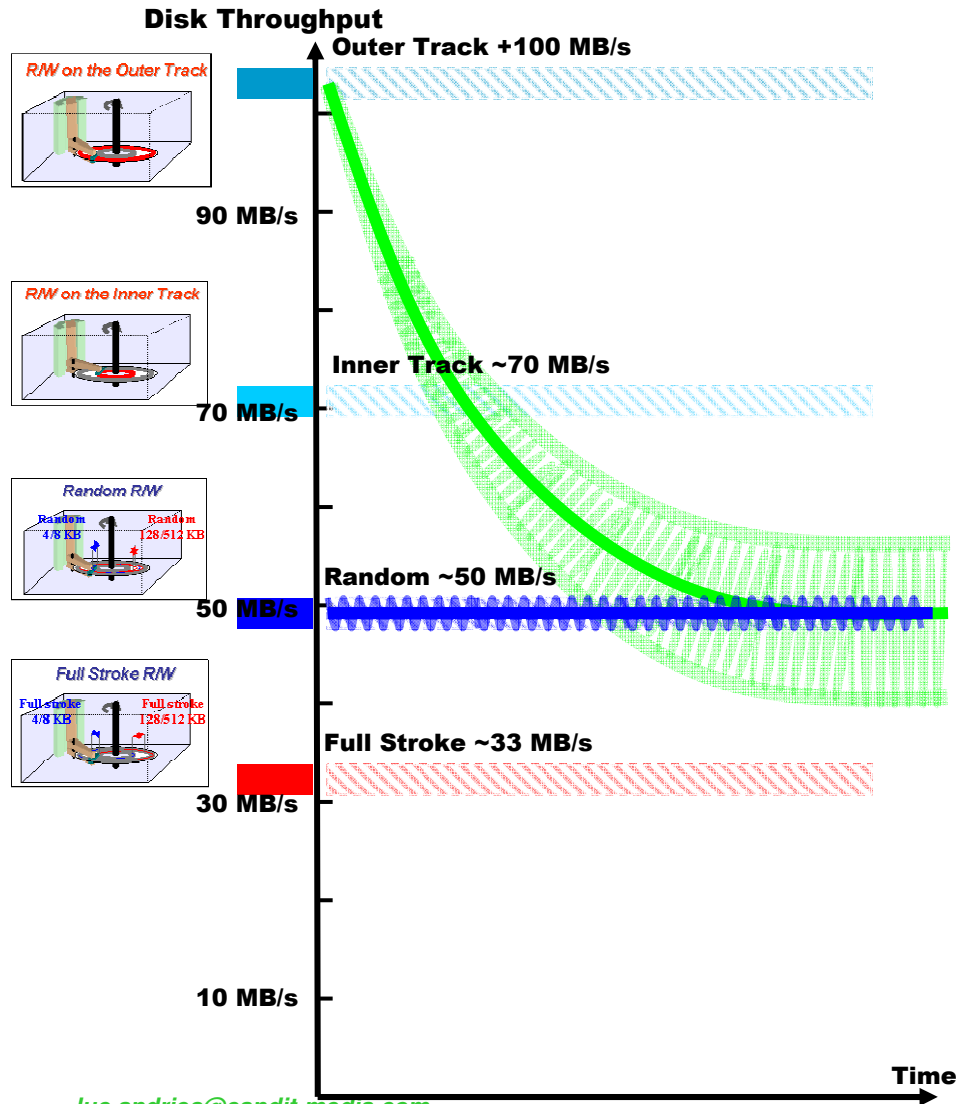
? Random Access = Worst Case ?



Relevance of Specifications

EBU Media Storage Workshop
21-22 Nov 2011

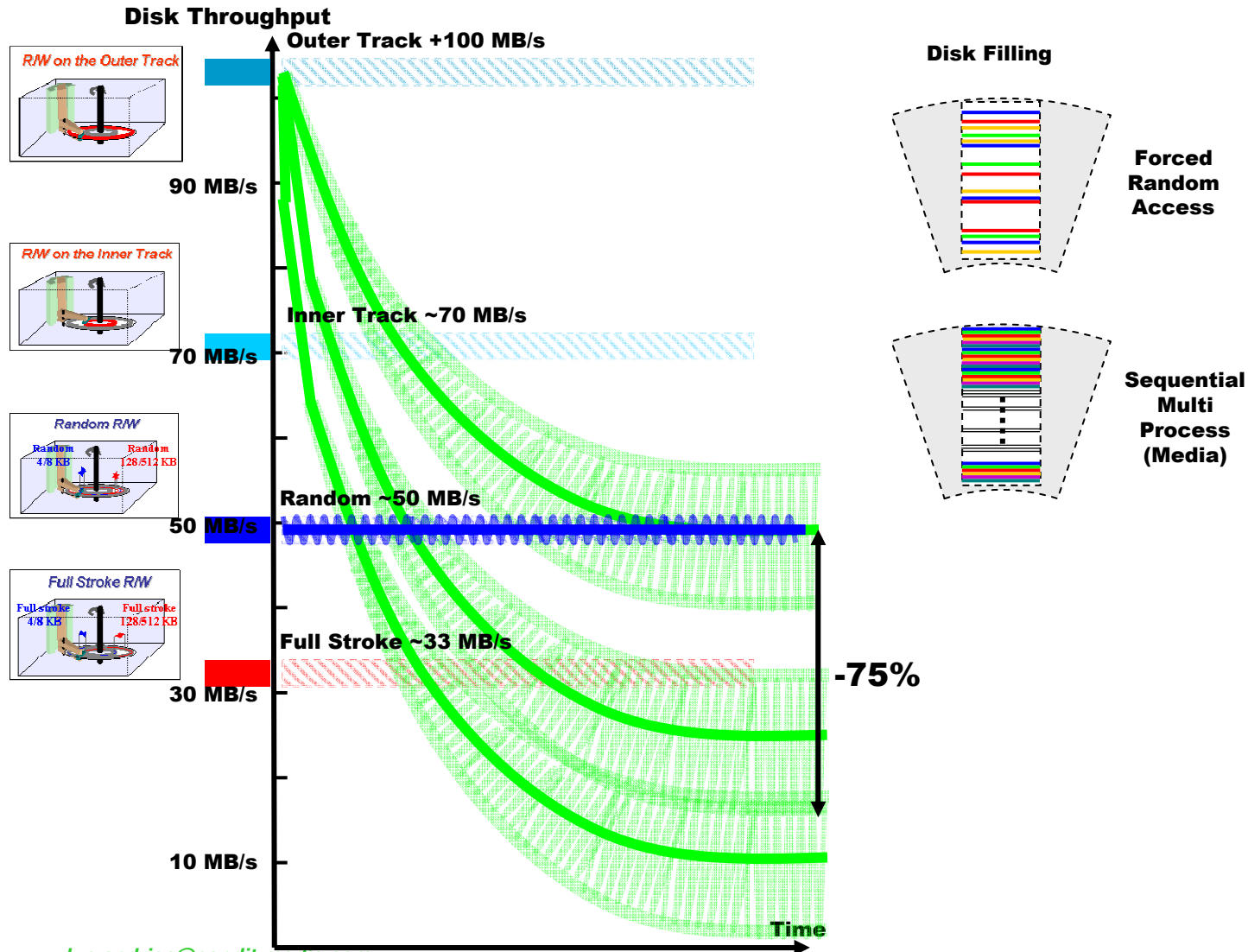
? Random Access = Worst Case ?



Relevance of Specifications

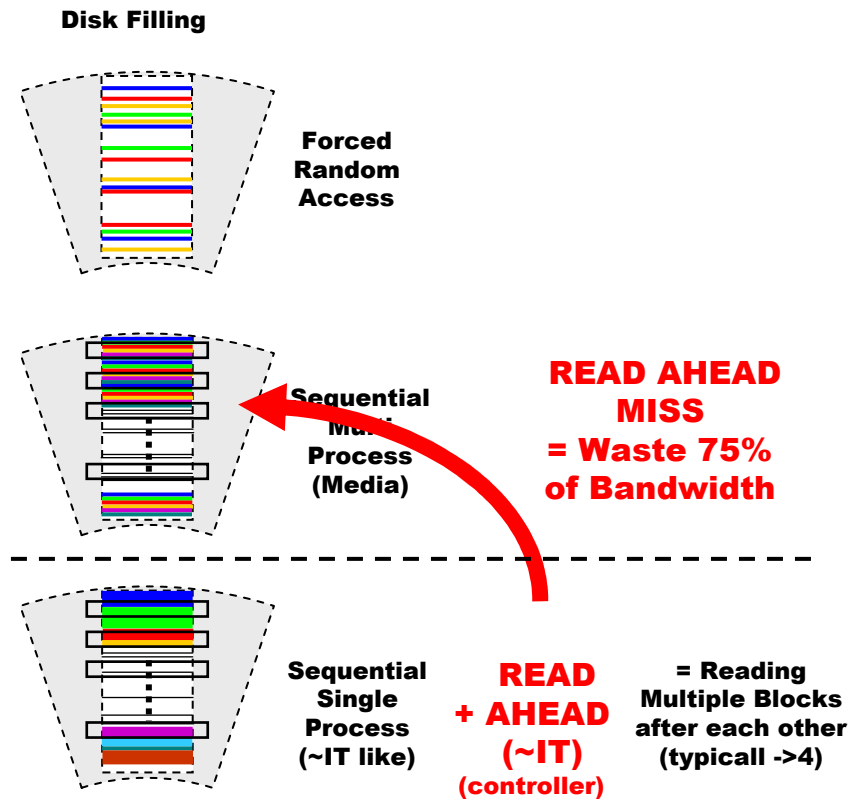
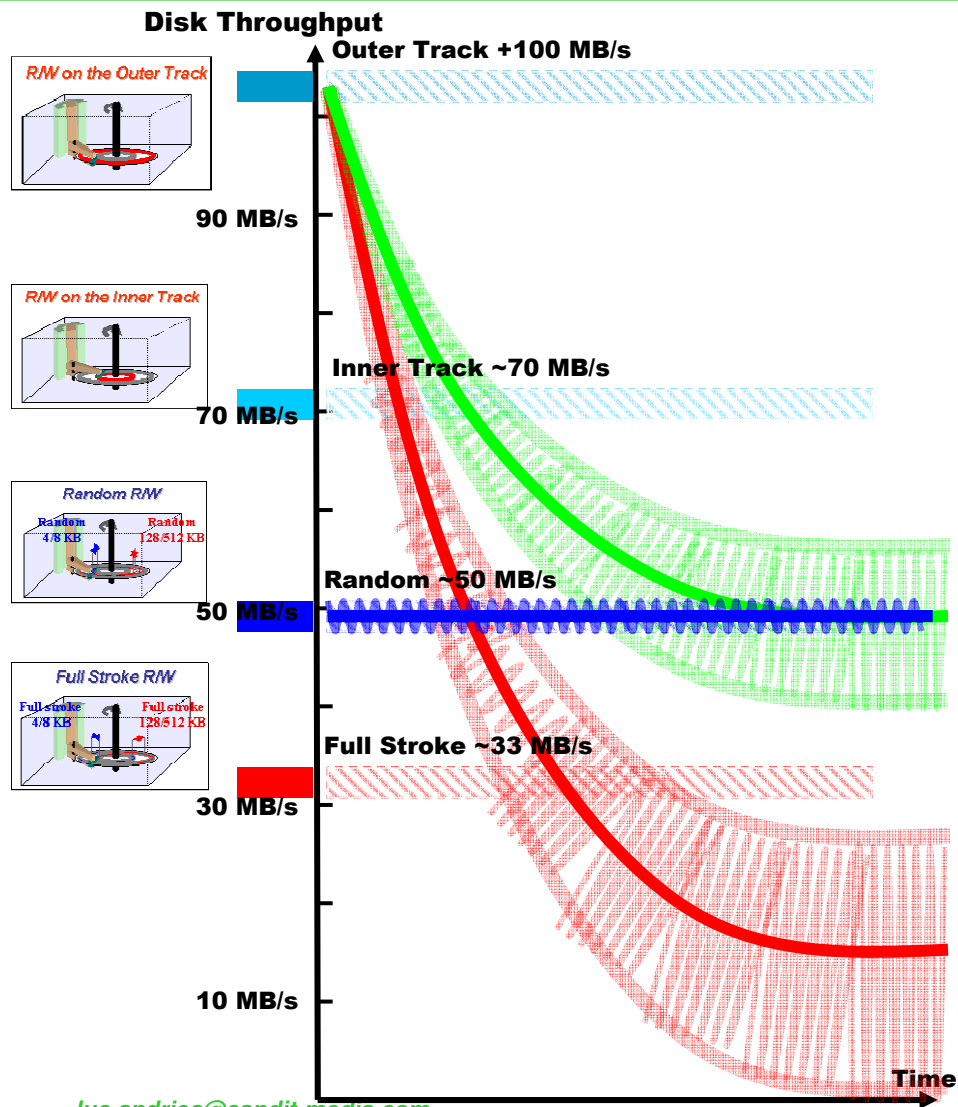
EBU Media Storage Workshop
21-22 Nov 2011

? Random Access = Worst Case ?

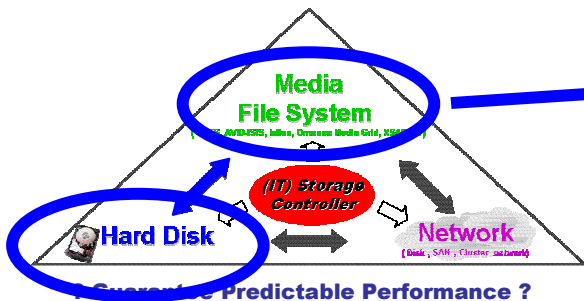


luc.andries@candit-media.com

? Random Access = Worst Case ?



(Media) File System Behaviour



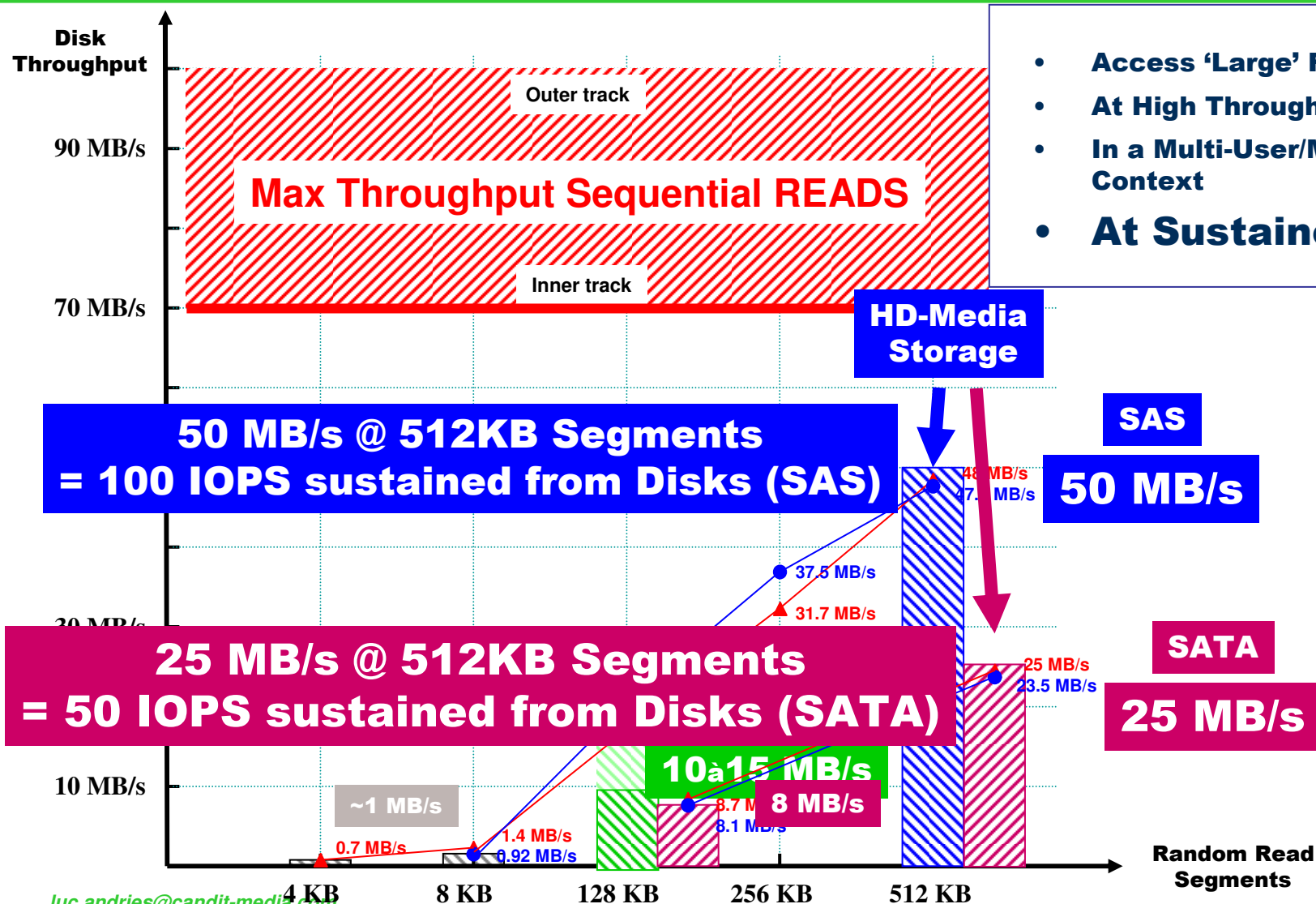
Media File System
should Force **Scattered**
Write Behaviour

- Access 'Large' Files
- At High Throughput
- In a Multi-User/Multi-Process Context
- At Sustained Load

Forced Fragmentation Guarantees
Consistent Random Access
Media File System Behaviour
on each Individual Disk

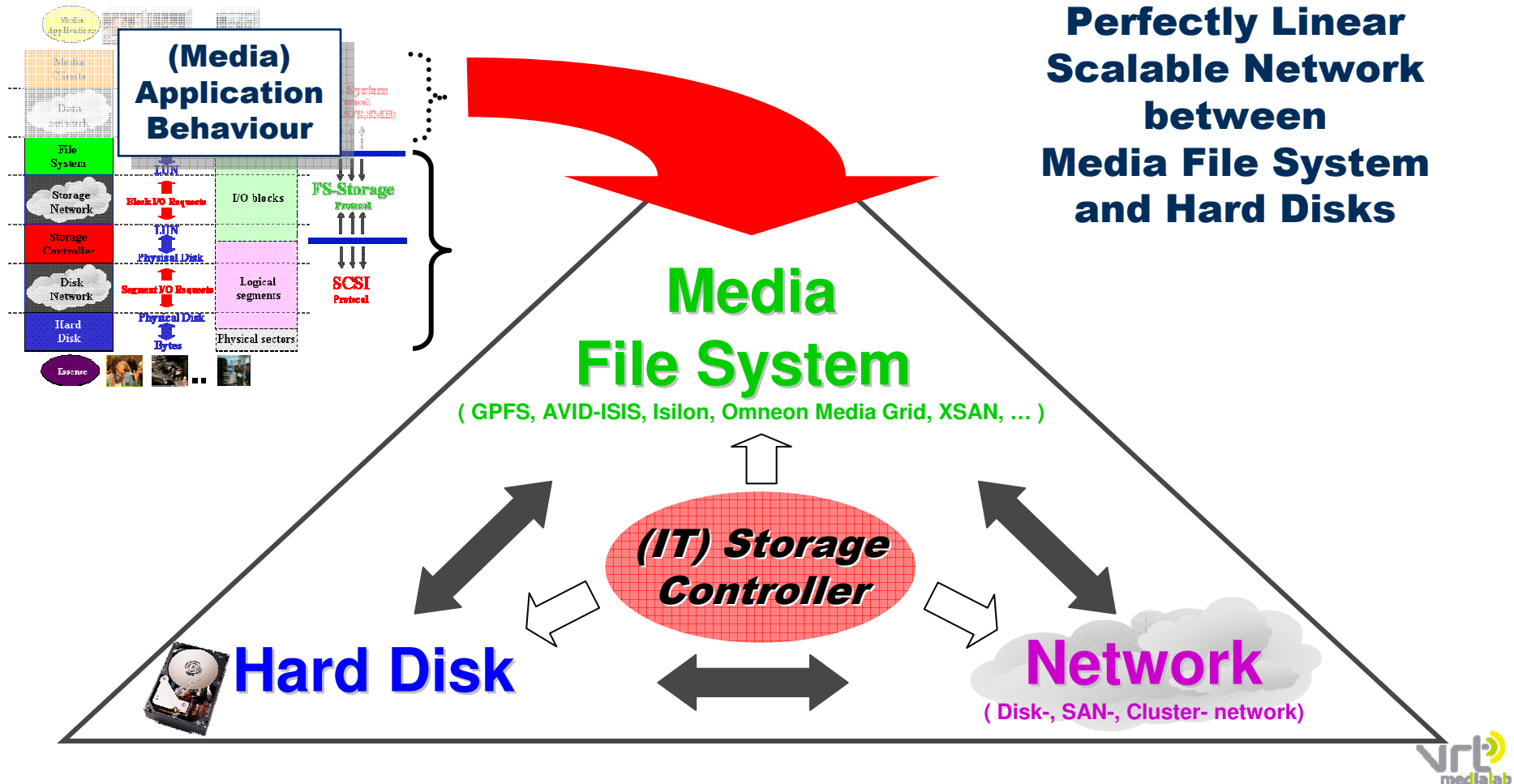
Guarantees
Predictable
Hard Disk Behaviour

Relevance of IOPS (I/O per Second)



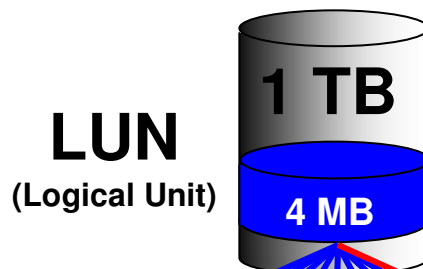
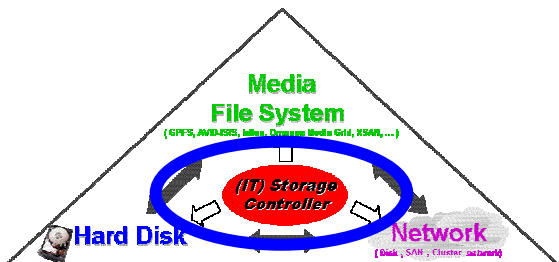
- Access 'Large' Files
- At High Throughput
- In a Multi-User/Multi-Process Context
- At Sustained Load

Media Storage Network



? Guarantee Predictable Performance ?

'Media' Storage Controller



**RAID 5 (8+1)
protection**

4 MB
I/O block



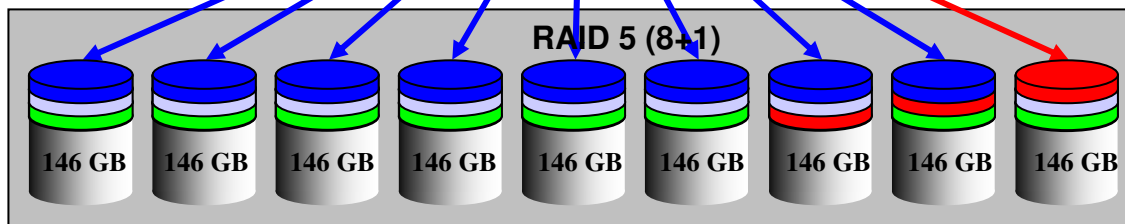
8x512 KB



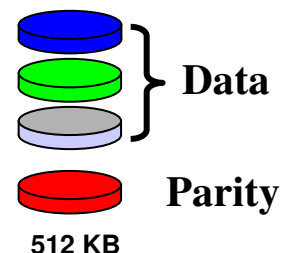
**Disk
Segments**

**Large stripe depth
(4MB) !!!**

**Disk efficiency
=> Large
segment size
(512 KB)**

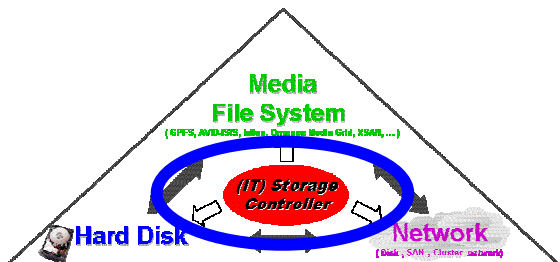


Physical disks

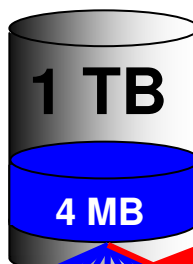


1 I/O (4MB) striped over 8 disks (+ parity)

'Media' Storage Controller



LUN
(Logical Unit)



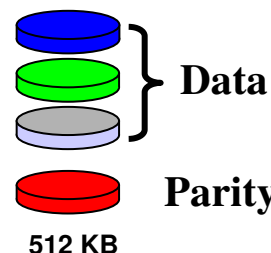
RAID 5 (8+1)
prot

**Large stripe depth
(4MB) !!!**

**Disk efficiency
=> Large
segment**

8x512 KB

**Disk
Segments**

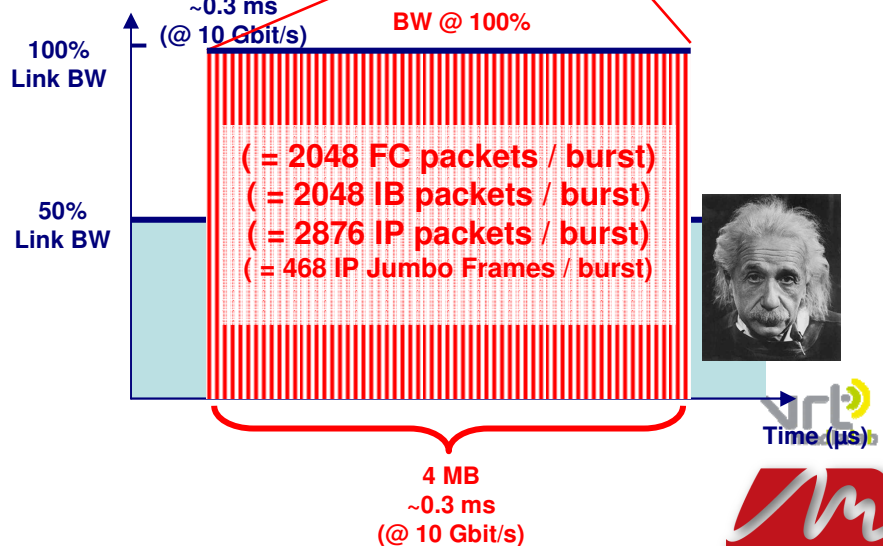
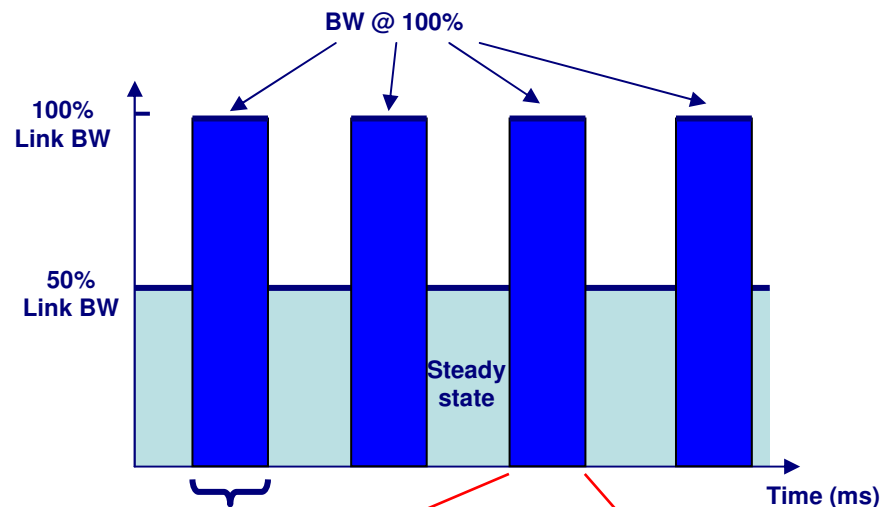
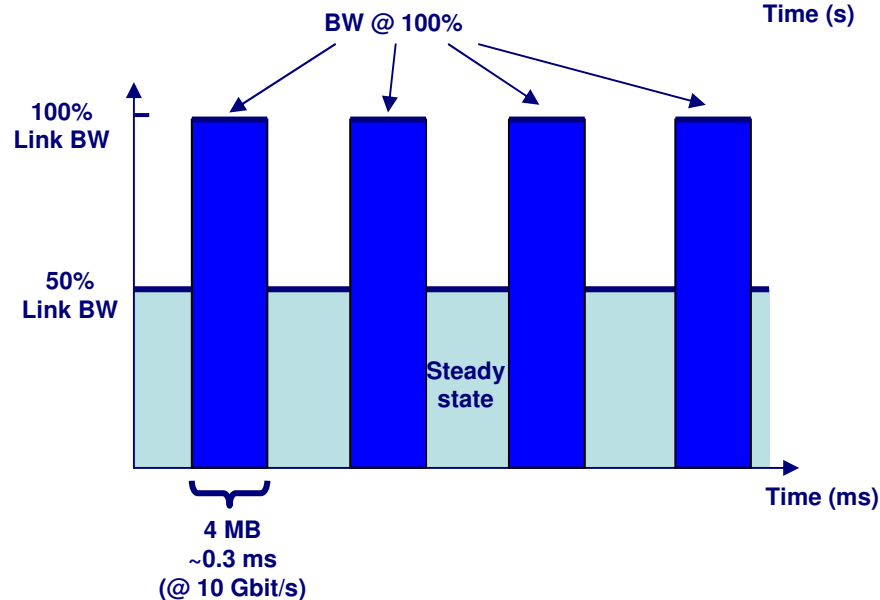
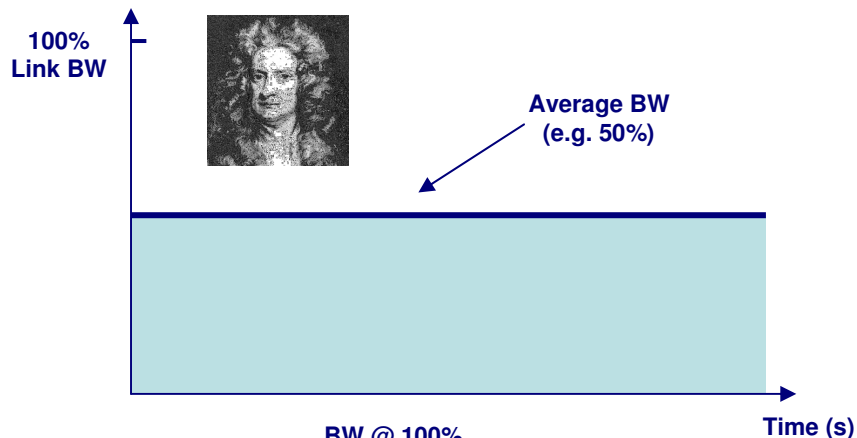


Physical disks

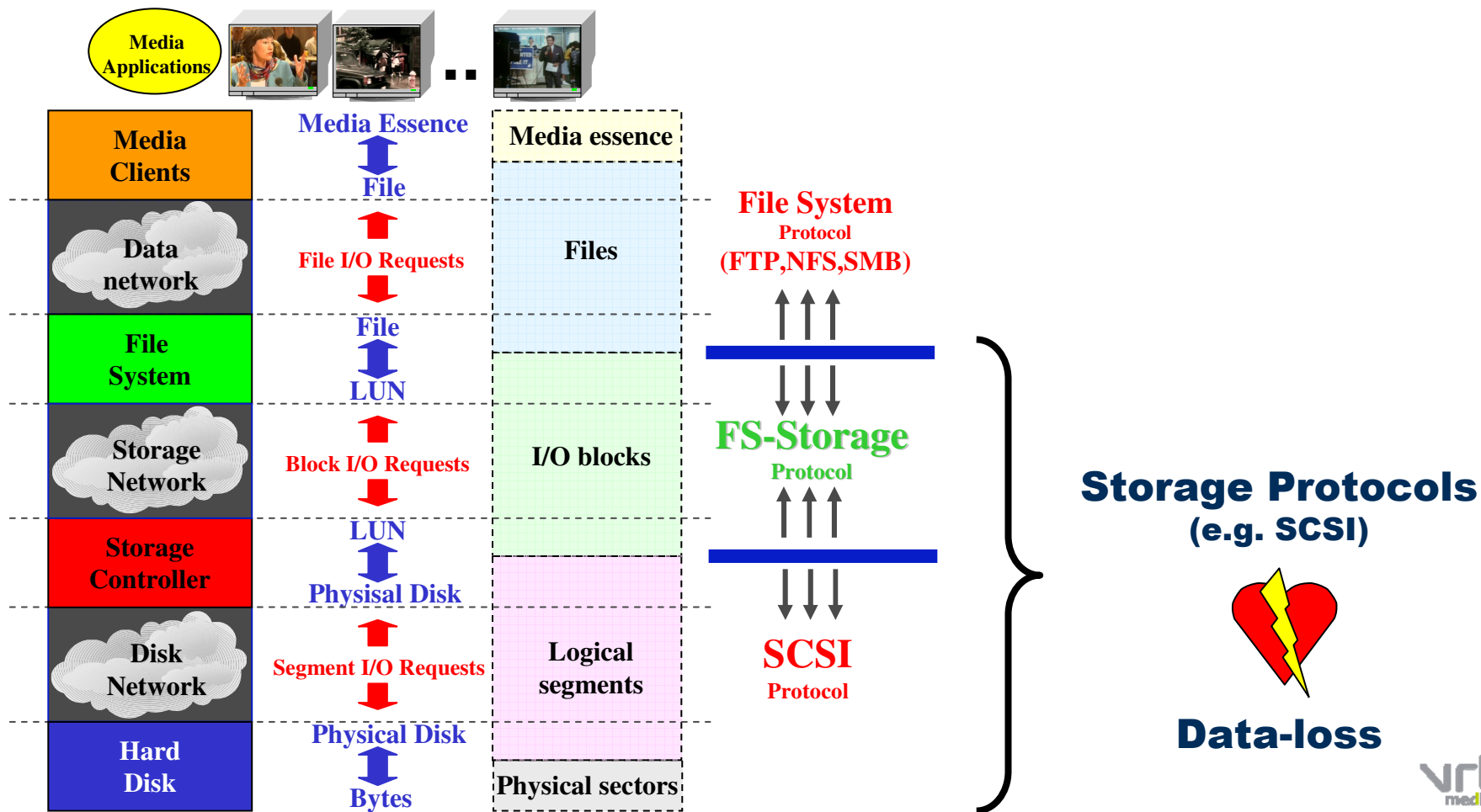


1 I/O (4MB) striped over 8 disks (+ parity)

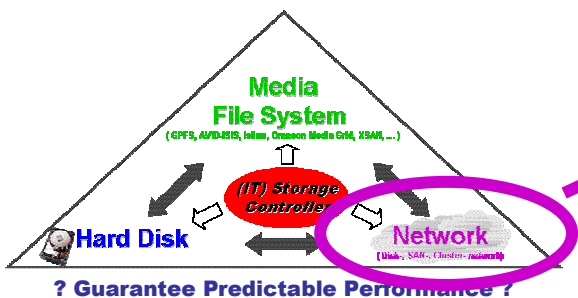
'Media' Storage Traffic Behaviour



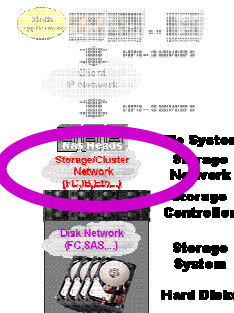
Media Storage Network



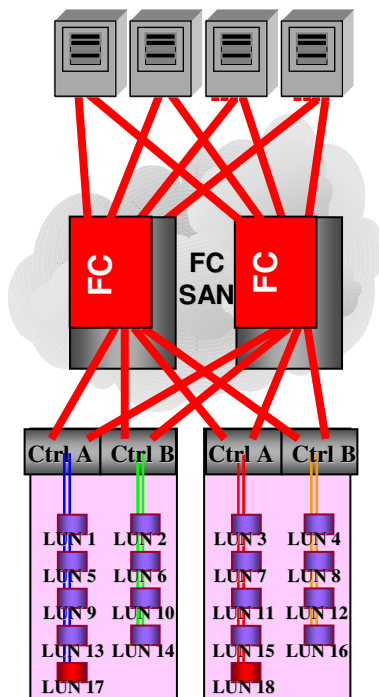
Media Storage Network



**Storage Networks
Require
Lossless
Network Technology**



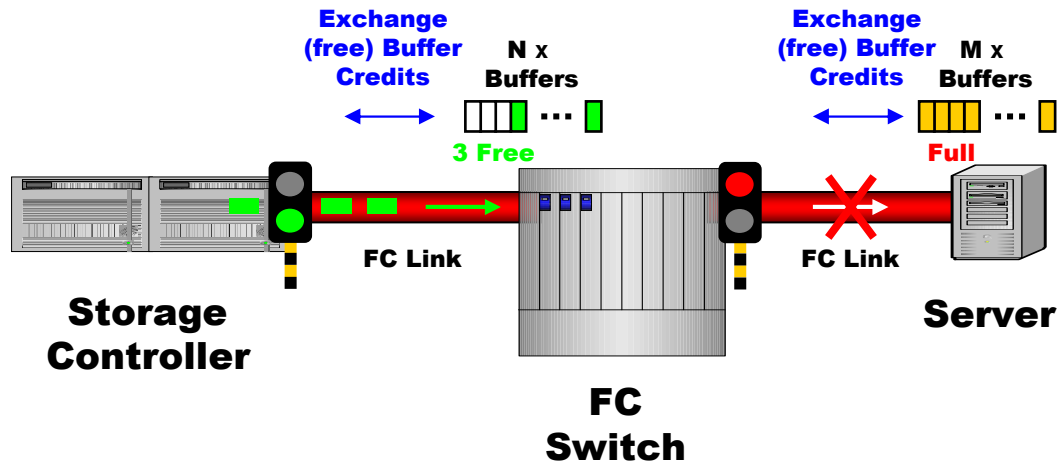
**Lossless
Technology
by Design**



**Fibre Channel (FC)
State of the Art
(IT) Storage
Network**

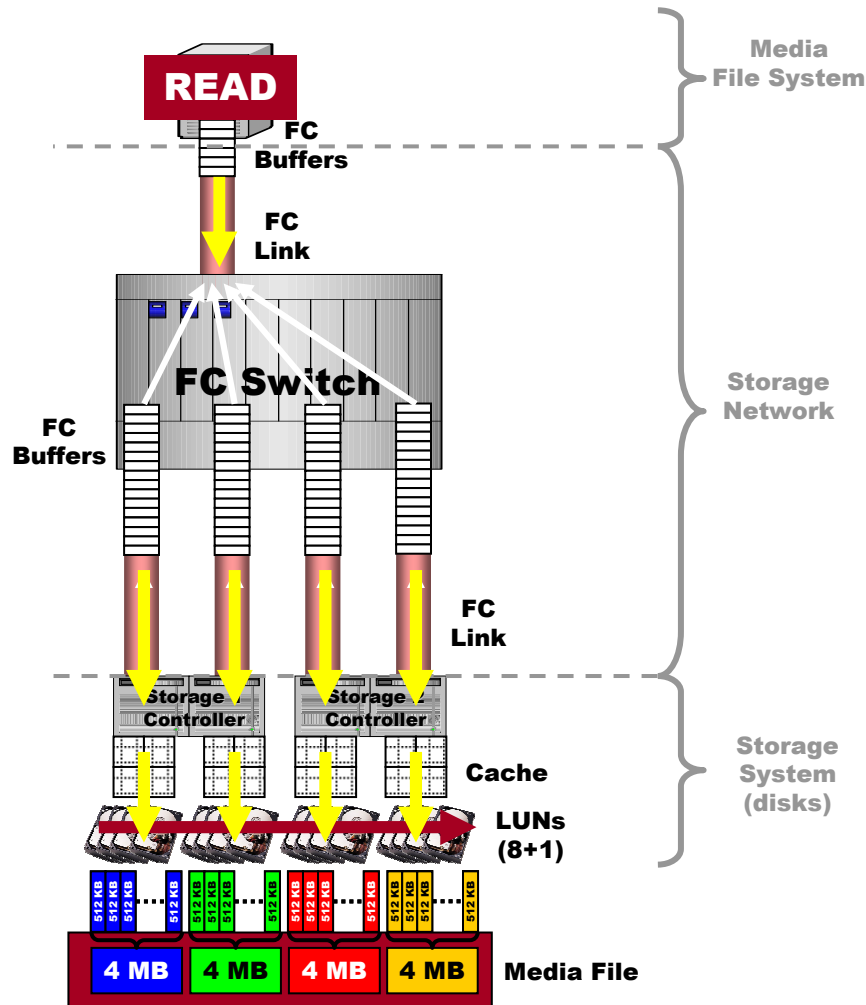
Lossless Fibre Channel (FC)

Lossless by Exchange of Available Buffers (Buffer_to_Buffer Credits)



**Packets are only
Transmitted
if Free Buffers are Available
at the (Link)-Destination**

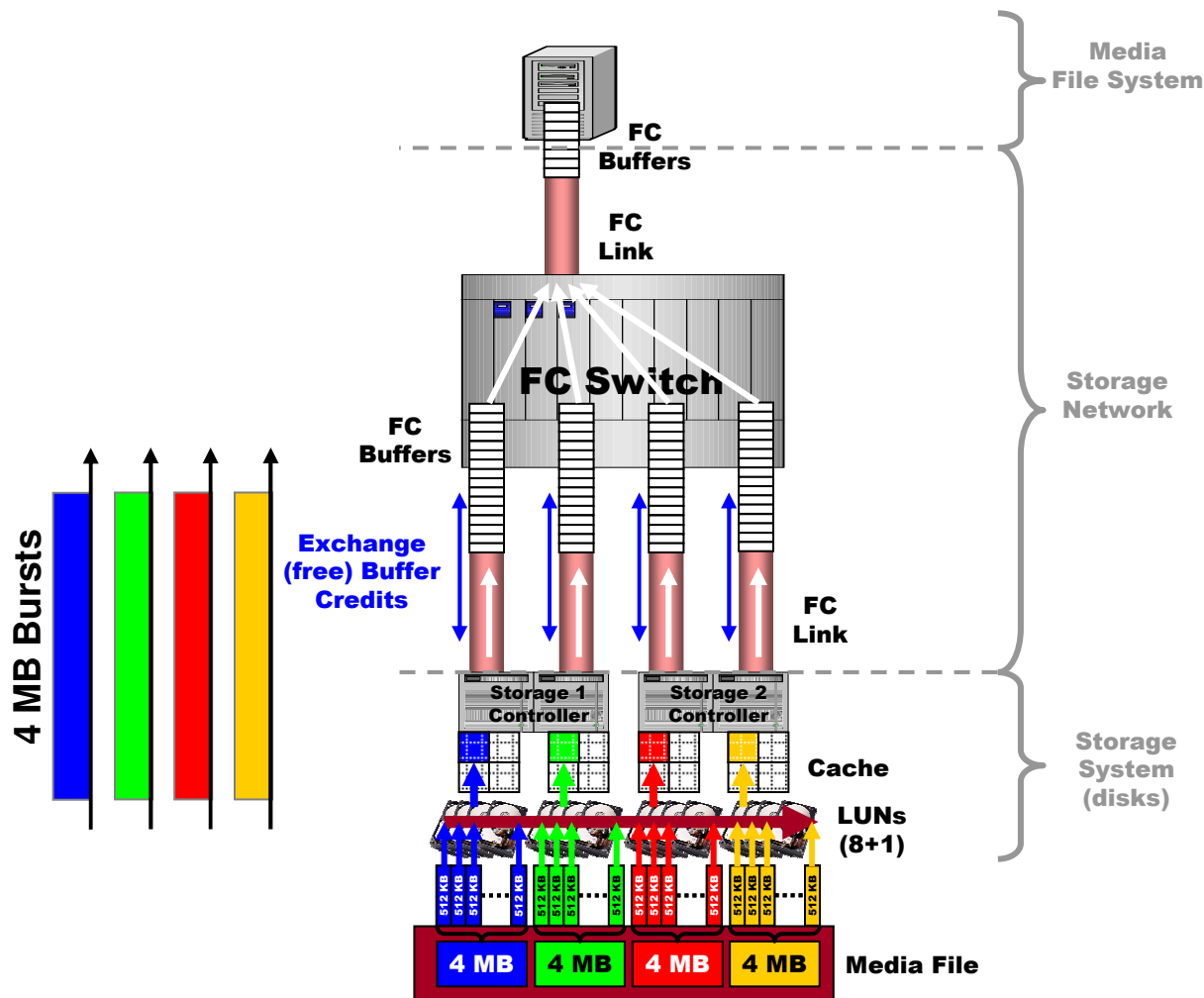
FC Storage Network in Media



Reading a Media File

- Media File Striping
- Large Block Size/Disk
- RAID 5 (8+1) -> 4MB
- File System READ

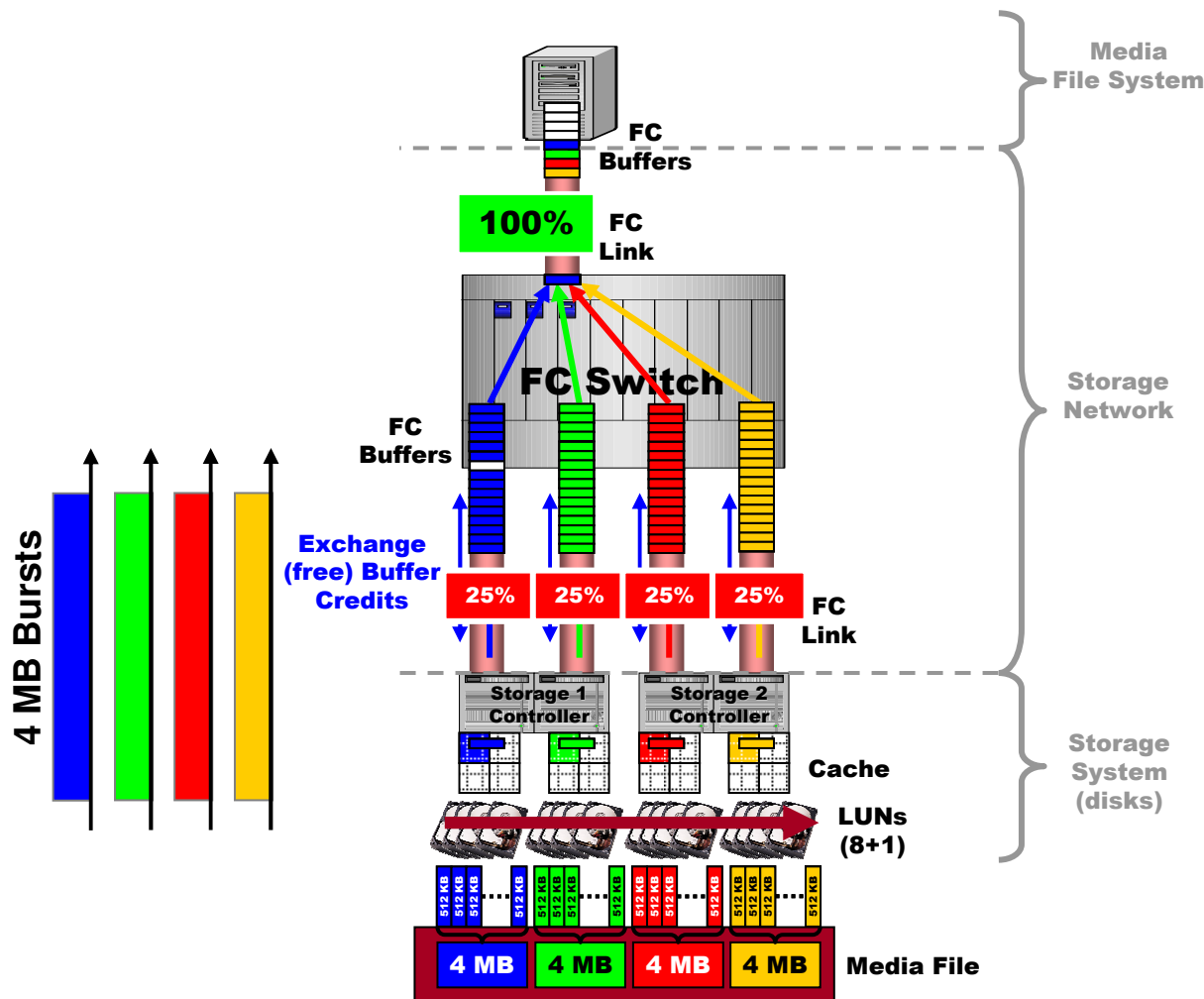
FC Storage Network in Media



Reading a Media File

- **Media File Striping**
- **Large Block Size/Disk**
- **RAID 5 (8+1) -> 4MB**
- **File System READ**
- **READ from Disk -> Cache**
- **4 MB Traffic Bursts**
- **Buffer_to_Buffer Credits**

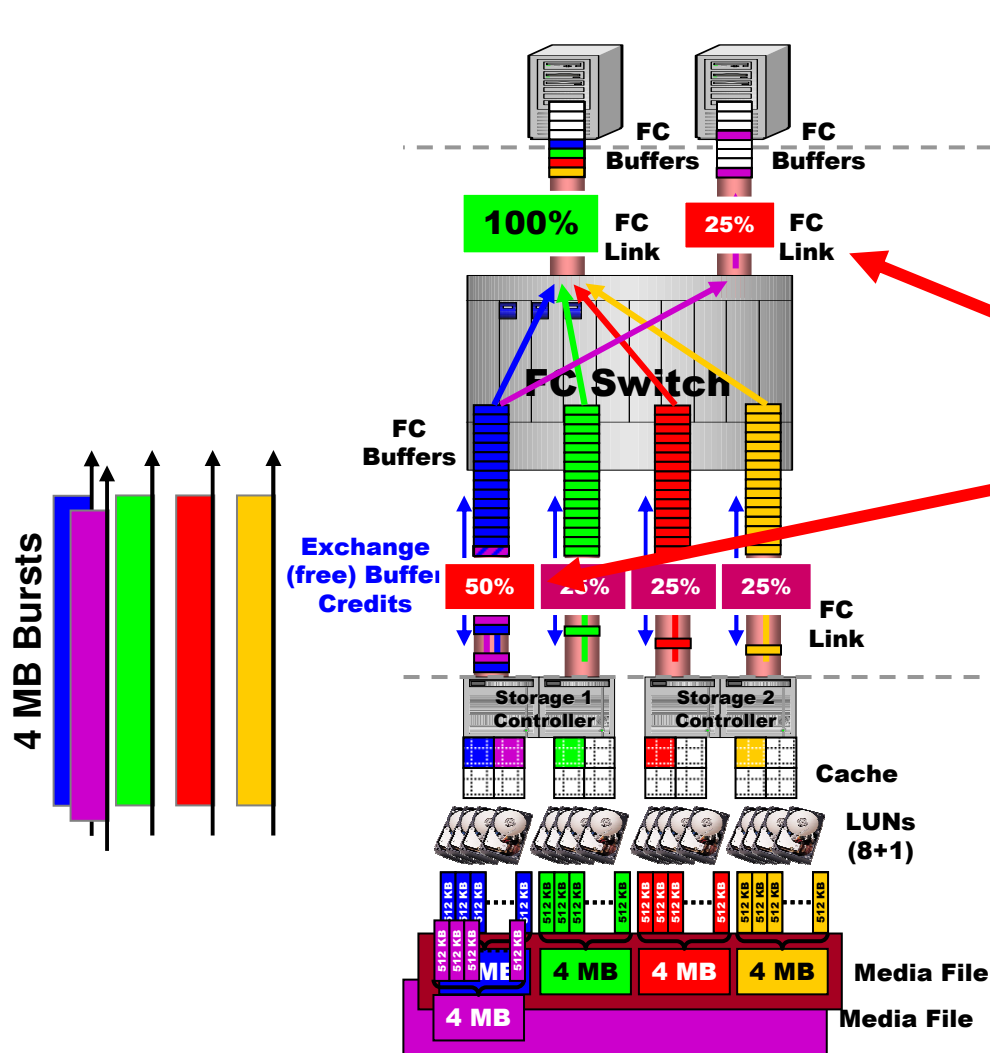
FC Storage Network in Media



Reading a Media File

- Media File Striping
- Large Block Size/Disk
- RAID 5 (8+1) -> 4MB
- File System READ
- READ from Disk -> Cache
- 4 MB Traffic Bursts
- Buffer_to_Buffer Credits
- Data Transfers (4:1)
- Flow Control Kicks in
- No Packet Loss!!!

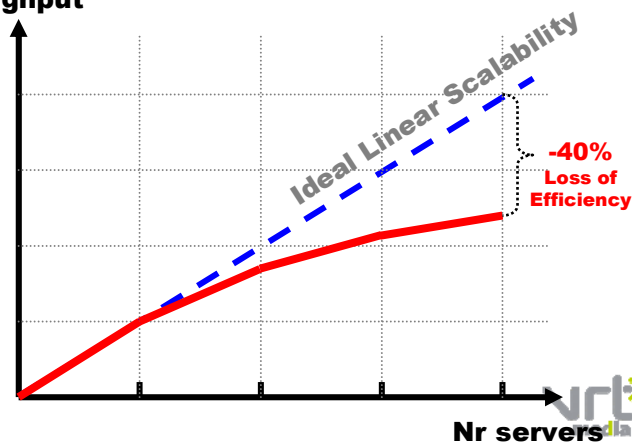
Media Traffic Interference



Reading Multiple Media Files

- 2nd Server
- Reading a 2nd File
- Sharing the 'Blue' Link
- 4 MB Traffic Bursts
- Link Level Flow Control Throttles both Flows
- Traffic Interference

Throughput

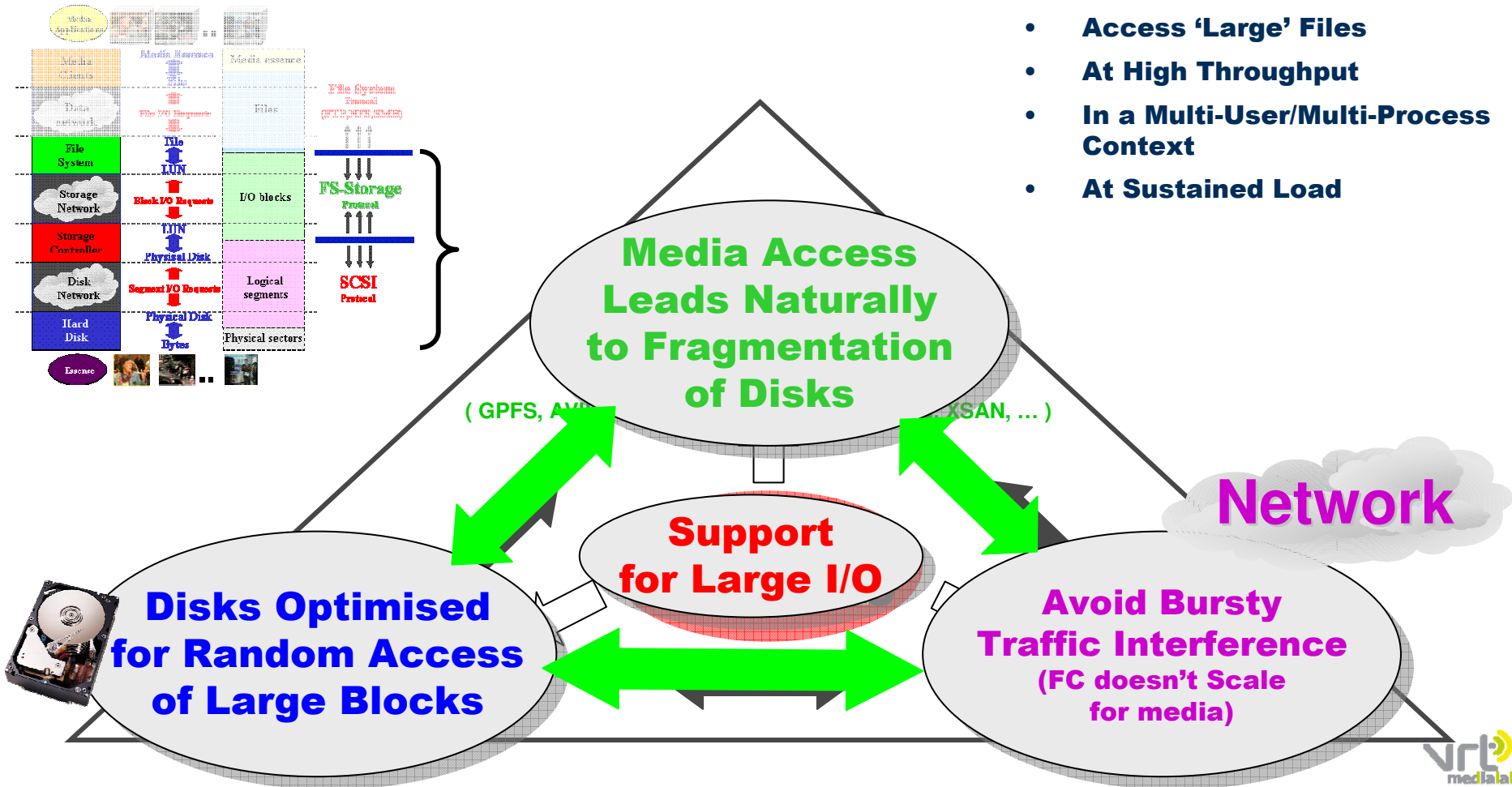


**FC Does NOT Scale
for Media**

Some Conclusions

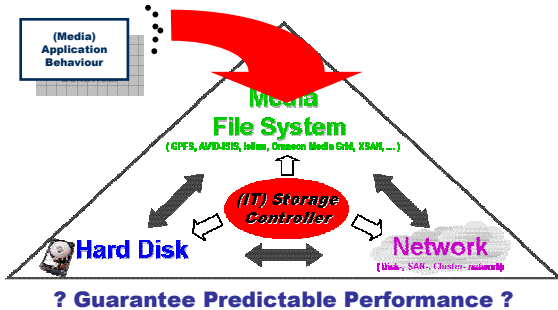
Relevance of Specifications
EBU Media Storage Workshop
21-22 Nov 2011

- Access 'Large' Files
- At High Throughput
- In a Multi-User/Multi-Process Context
- At Sustained Load



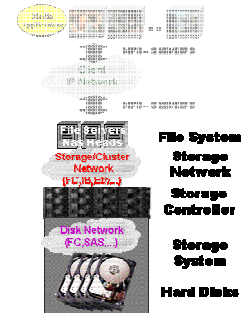
?! Guaranteed Predictable Performance ?!

(Media) Application Behaviour



Media Application Behaviour

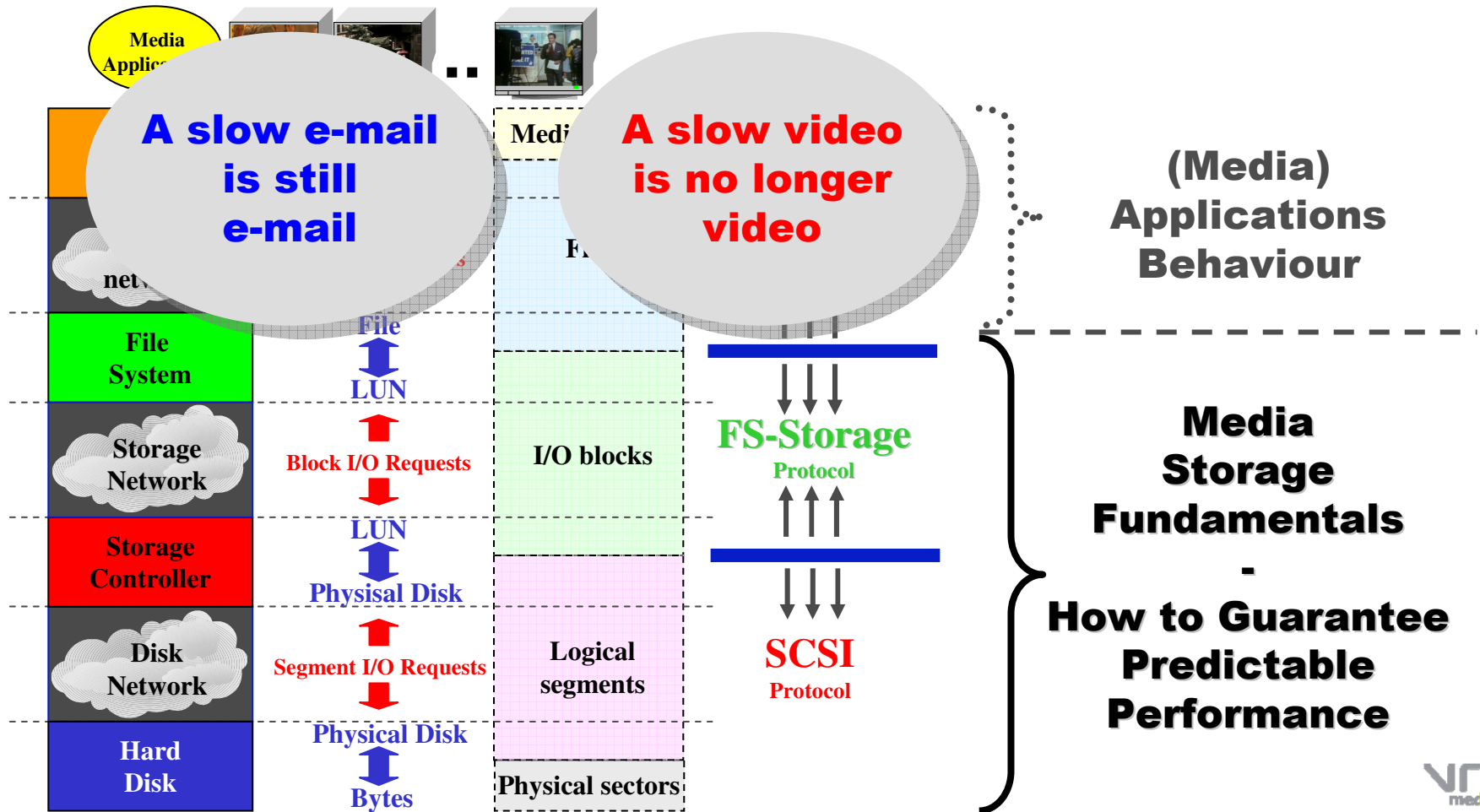
~



- Depends upon Application
- Depends on File system Protocol used
- Depends on Network Protocol used
- Depends on File System Intelligence
- Depends on the Storage
- Depends on the Network

EBU Networks Seminar
28-29 June 2011

Final Conclusion



Thank You !

