

Challenges and solutions to long-term records retention and access

Dr Matthew Addis
Arkivum Ltd
SAG Edinburgh 2013

Contents

- Retention of electronic documents and records
- What are the challenges?
- How can they be overcome?

“Security is of key importance to our business, Arkivum’s A-Stor Pharma service allows us to store our encrypted data for the long term in a cost efficient way that is entirely scalable and reduces pressure on our internal IT infrastructure.”

Dan Watkins, Oxford Fertility Unit

WHY KEEP ELECTRONIC DOCUMENTS AND RECORDS?

Why keep electronic documents and records?

- Regulation and compliance
- Intellectual Property
- Reuse
- Save money
- Better use of in-house resources

GxP Regulations

- FDA 21 CFR Part 11
- EudraLex Volume 4 Annex 11
- OECD GLP



MHRA

- MHRA GCP guide
- Medicines for Human Use Clinical Trials Act (SI 2004/1031)



STATUTORY INSTRUMENTS	
2004 No. 1031	
MEDICINES	
The Medicines for Human Use (Clinical Trials) Regulations 2004	
<i>Made</i> - - - -	<i>31st March 2004</i>
<i>Laid before Parliament</i>	<i>1st April 2004</i>
<i>Coming into force</i> - -	<i>1st May 2004</i>
<p>The Secretary of State, being a Minister designated(a) for the purposes of section 2(2) of the European Communities Act 1972(b) in relation to medicinal products, in exercise of the powers conferred by the said section 2(2), and of all other powers enabling him in that behalf, hereby makes the following Regulations:</p>	

Research

MRC ethics series

Good research practice:
Principles and guidelines



- MRC GRP guide



Medical Research Council
July 2012

Records Management

ISO 15489-1:2001

Information and documentation -- Records management -- Part 1: General

ISO 23081-1:2006

Information and documentation -- Records management processes -- Metadata for records -- Part 1: Principles



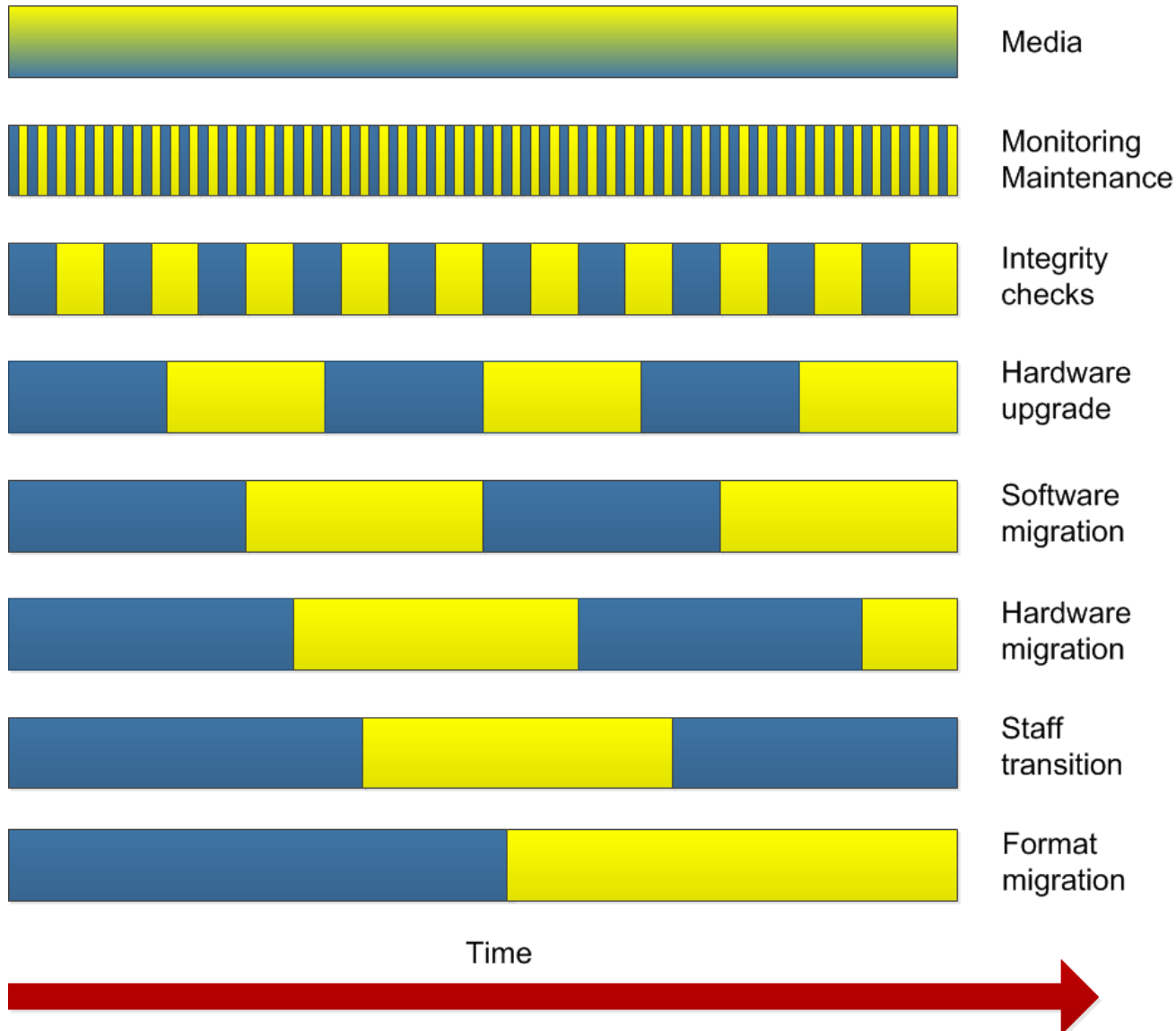
Modular Requirements
for Records Systems

Common Requirements

- Integrity
- Authenticity / reliability
- Access control / confidentiality
- Usability / ready access / readability
- Responsibility / named individuals
- SoPs, validation, audit trails
- ***Risk management***

WHAT ARE THE CHALLENGES?

20 years of keeping content alive



Digital Preservation

*“Digital information lasts forever -
or five years, whichever comes first.”*

Jeff Rothenberg

Perpetual Change

- Change costs money
- Change takes time
- Change introduces risk
- Change requires validation
- Change needs planning and management

Change and Risk

- Do nothing
- Do the wrong thing
- Delayed action
- In-house
- Use a service provider

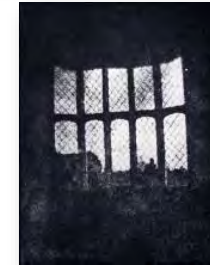


Risks

- Technical obsolescence, e.g. formats and apps
- Hardware failures, e.g. digital storage systems
- Loss of staff, e.g. skilled archive and IT staff
- Insufficient budget, e.g. storage too expensive
- Accidental loss, e.g. human error
- Selection, e.g. don't retain what you should
- Stakeholders, e.g. retention no longer a priority
- Underestimation of resources or effort
- Fire, flood, meteors, aliens...



Risk Identifier:	R55	
Risk Name:	Loss of integrity of information	
Risk Description:	Repository is incapable of demonstrating that the integrity of information has been maintained since its receipt, and that what is stored corresponds exactly with what was originally received.	
Is this Risk Relevant?:	<ul style="list-style-type: none"> Does repository commit to preservation of information integrity? 	
Example Risk Manifestation(s):	<ul style="list-style-type: none"> Records documenting government expenditure have been subjected to unauthorised or unanticipated changes, rendering them no longer representative of originally deposited content 	
Nature of Risk:	Physical environment	
	Personnel, management and administration procedures	
	Operations and service delivery	X
	Hardware, software or communications equipment and facilities	
Owner:	Preservation	
Escalation Owner:	Preservation	
Stakeholders:	Management; financiers; staff; depositors; users; producers	
Mitigation strategy(ies):	<p>Avoidance strategies:</p> <ul style="list-style-type: none"> Ensure policies and procedures are conceived with due consideration of integrity requirements Maintain and review policies and procedures to ensure adequate recording and comparison of checksums to demonstrate that archived information has suffered no loss of integrity since its deposit or receipt Ensure software and hardware systems and preservation strategies are capable of preserving information integrity <p>In the event of risk's execution:</p> <ul style="list-style-type: none"> Invoke treatment strategies to alleviate loss of reputation or trust 	
Risk Relationships:	→R01 [contagious] →R02 [contagious]	
Risk Probability:	4	
Risk Potential Impact:	3	
Risk Severity:	12	



30,000 years

2,200 years

1,200 years

200 years

20 years

Each new generation of technology:

**1000x times denser
lasts 1/10th as long**



Adobe Reader



Adobe Reader could not open 'ArkivumPresentation.pdf' because it is either not a supported file type or because the file has been damaged (for example, it was sent as an email attachment and wasn't correctly decoded).

OK

It's not getting any better!

- 1000 times more HDD capacity over last 15 years
- Only 10 times lower Bit Error Rates (BER)
- HDD BER = 10^{-14}
- 1 TB = 10^{13} bits
- 10% chance of an error when reading all of a HDD
- *Within a few years, you are more likely than not to get a read error when copying a HDD*

No storage media is 100% safe

- SSD
- CD, DVD
- Data Tape

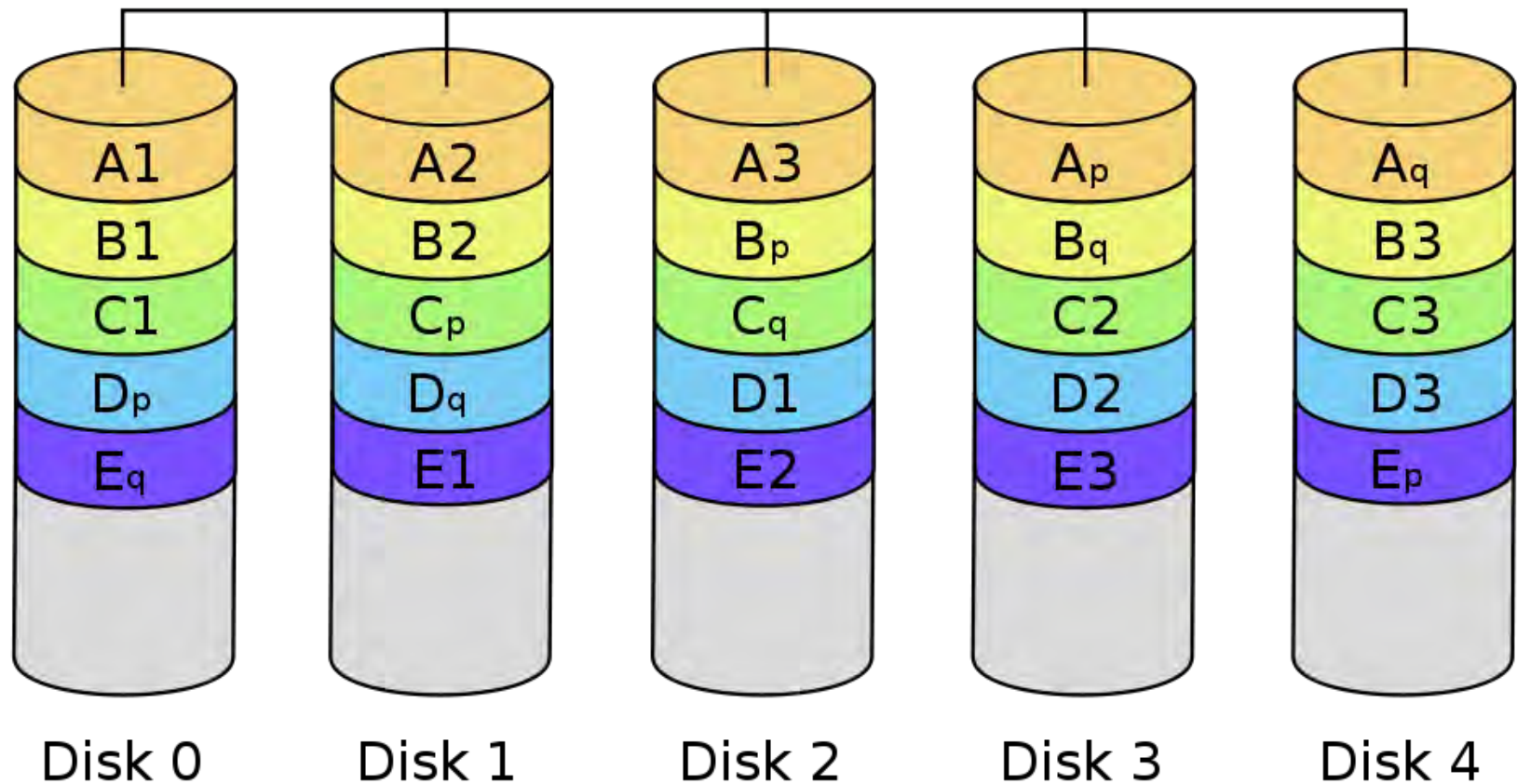


Examples of damaged discs.

Failure	Seen?	Devices exhibiting that failure
Bit Corruption	Y	SSD#11, SSD#12, SSD#15
Flying Writes	N	-
Shorn Writes	Y	SSD#5, SSD#14, SSD#15
Unserializable Writes	Y	SSD#2, SSD#4, SSD#7, SSD#8, SSD#9, SSD#11, SSD#12, SSD#13, HDD#1
Metadata Corruption	Y	SSD#3
Dead Device	Y	SSD#1
None	Y	SSD#6, SSD#10, HDD#2

The IT Industry knows this already

RAID 6

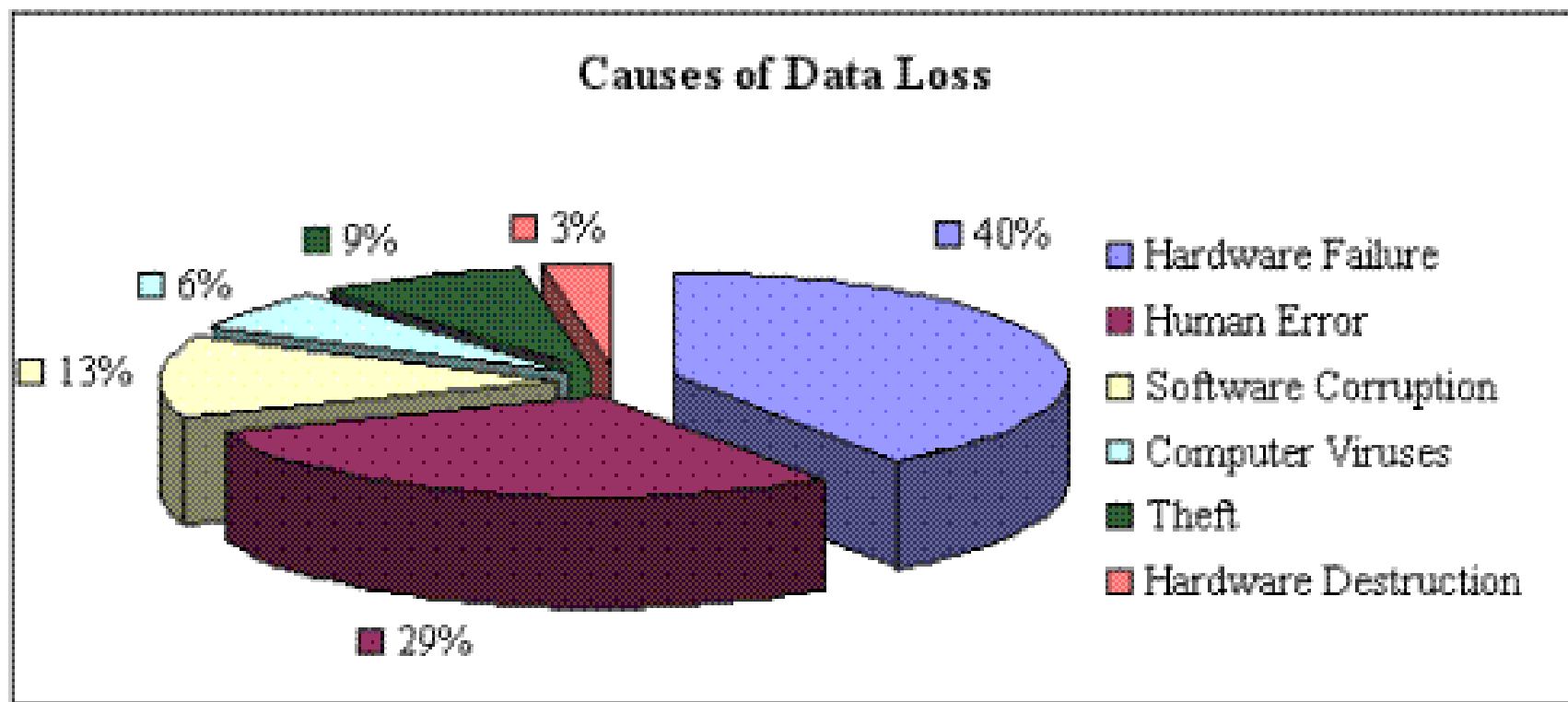


Systems bring their own problems

“Disk failures are not always a dominant factor of storage subsystem failures, and a reliability study for storage subsystems cannot only focus on disk failures. Resilient mechanisms should target all failure types”

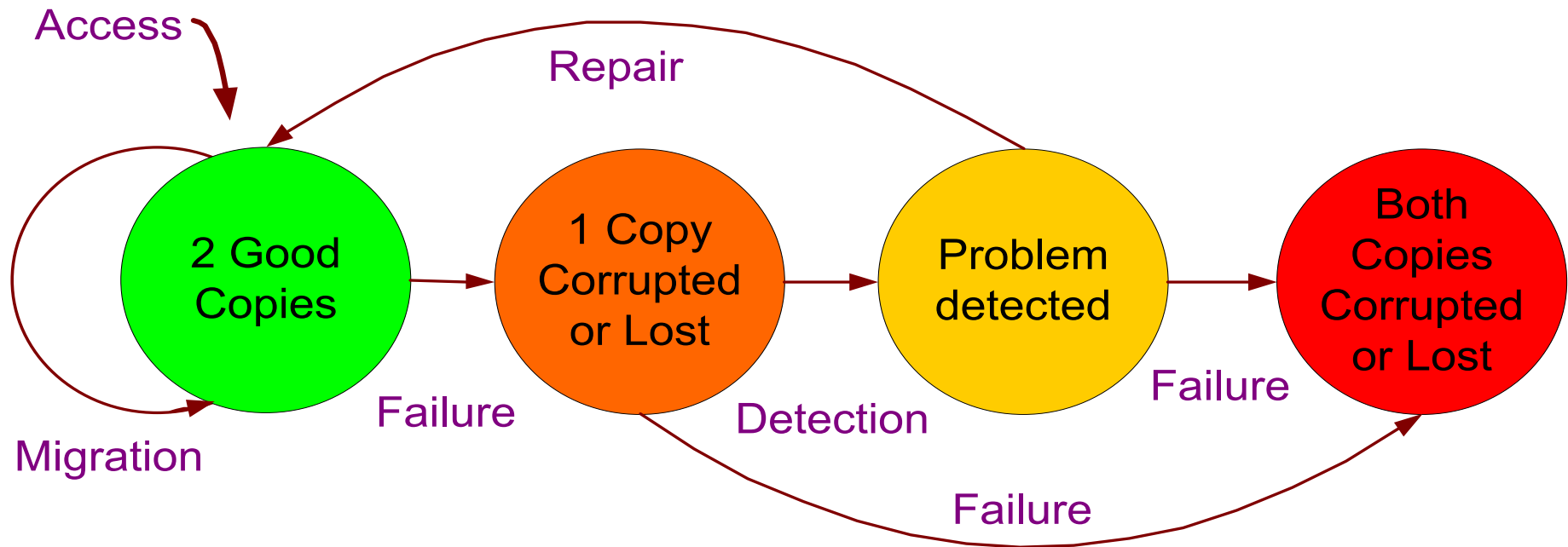
NetApp study of 1.8M HDD in 155,000 systems

People cause data loss too!



Drive risk down

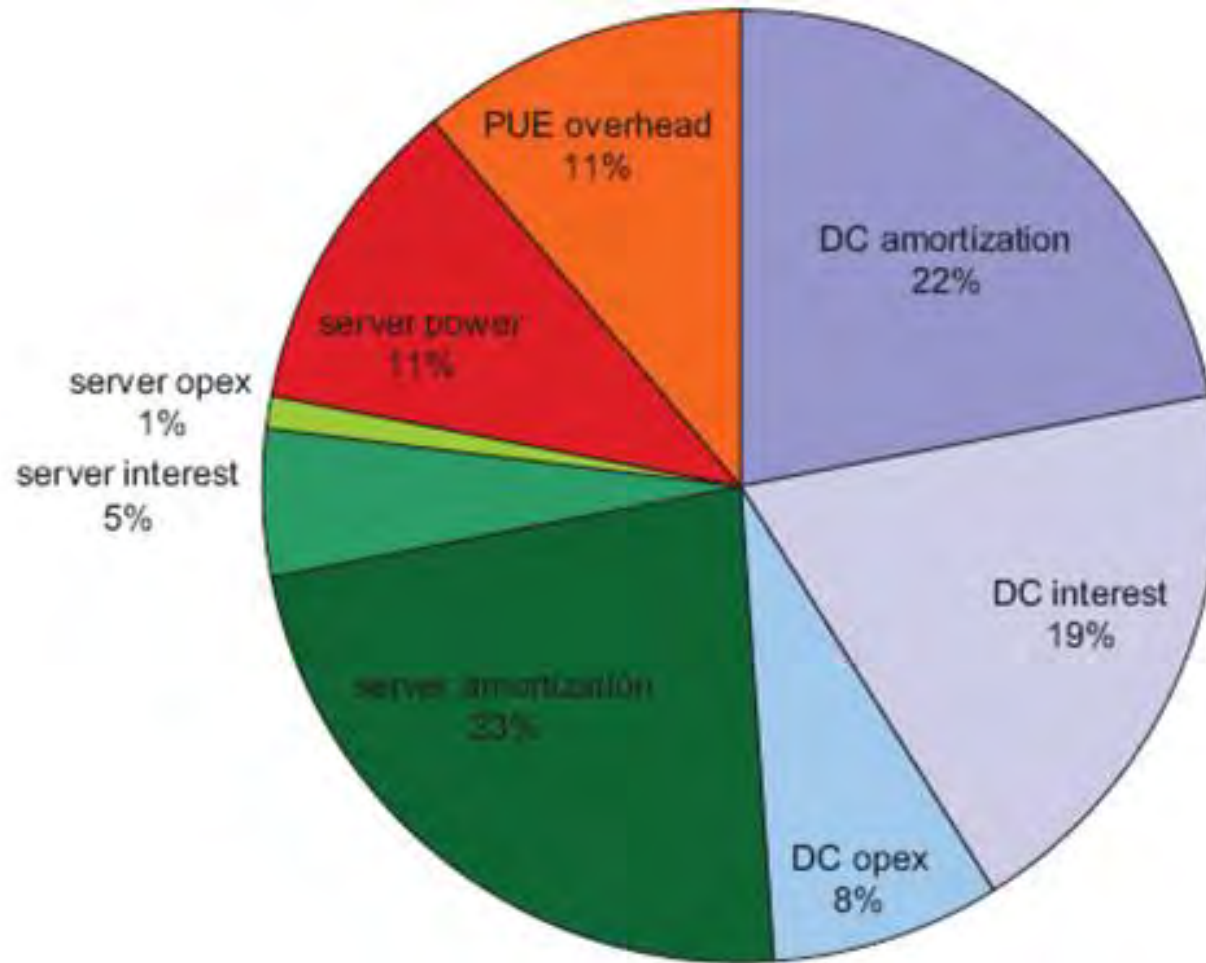
- Diversity keeps things safe
- Active management of data integrity
- Migration to address obsolescence



Why not make more copies?



Storage TCO



Google

Data

- Documents: eTMF, eCTD
- Supporting data
 - LIMS
 - NGS
 - MRI, CT, 3D ultrasound, 2D x-ray
 - Sensors and monitoring
 - Large cohort studies
- Large *numbers* of files and/or large volumes

Something has to give

- >100% CAGR in data volumes / records
- 20% CAGR in HDD capacity per £
- 2% CAGR in IT Budgets

IHS iSuppli

ComputerEconomics

→ Store less

→ Compress

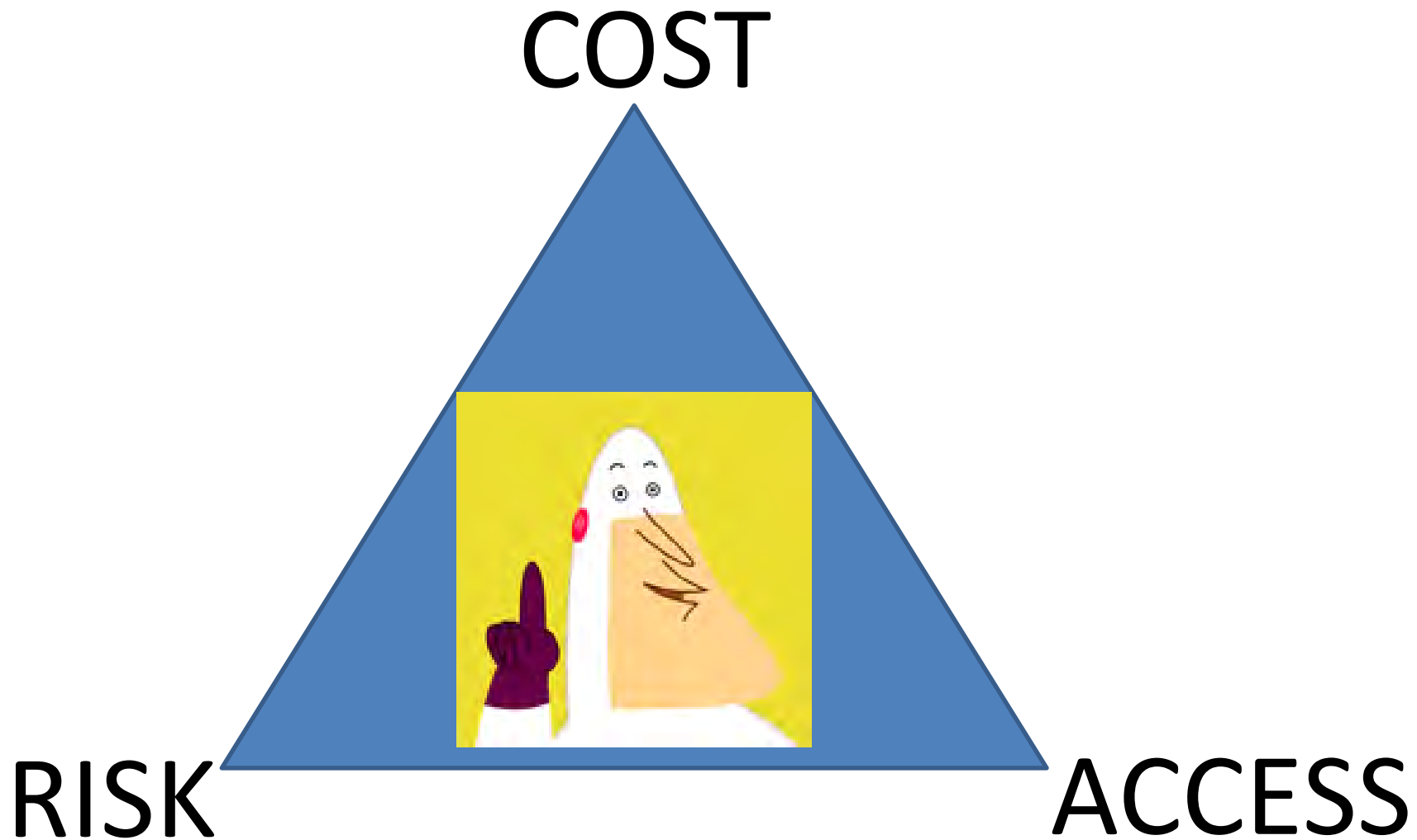
→ Archive outside of EDM systems

MAKING CHOICES

Cost v.s. safety v.s. access

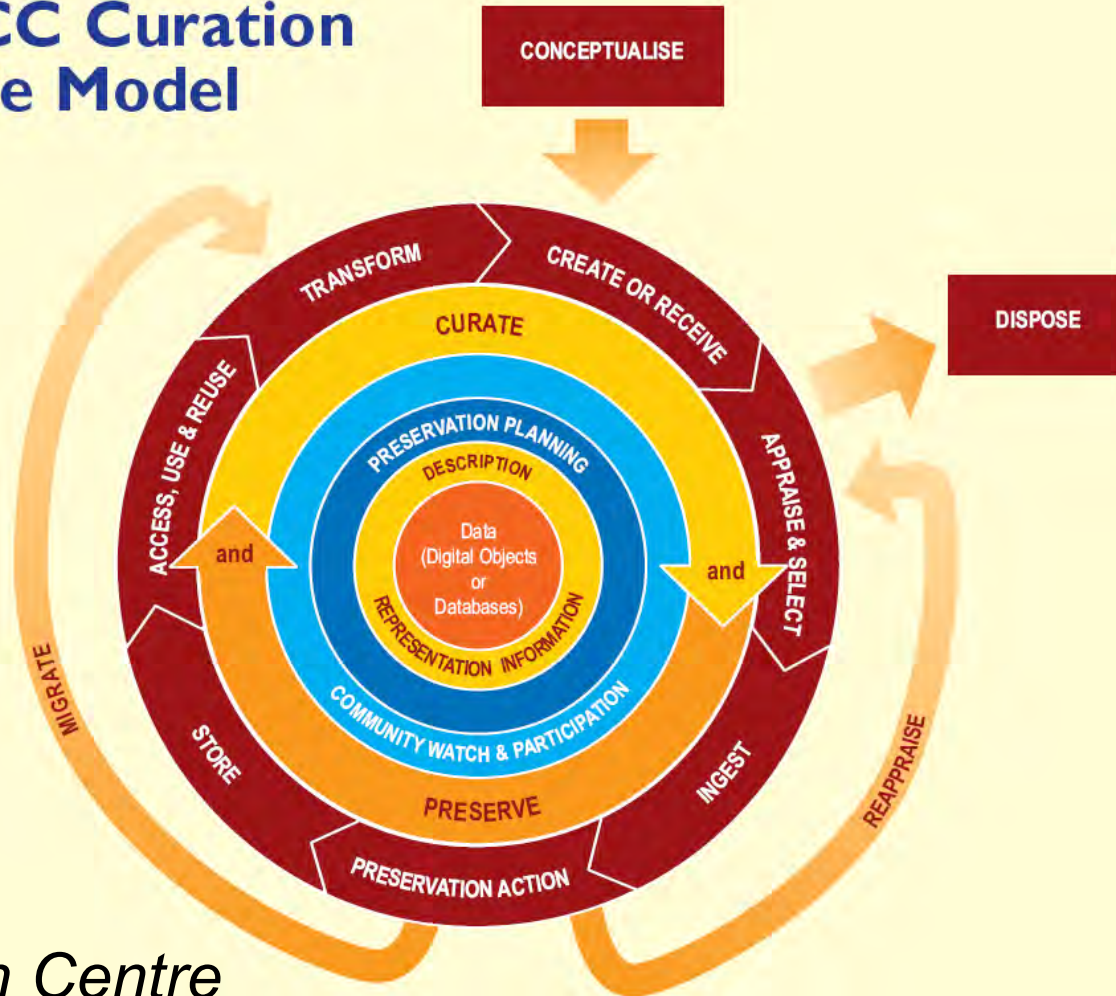
	Data tape on shelves	HDD in servers	Storage as a Service
Storage Cost	Low (media, shelves, climate control)	High (servers, power, cooling, maintenance)	High (fully managed service)
Access Cost	High (people retrieve and load media)	Low (internal network, automated)	High (bandwidth, charges for i/o)
Latent Failures	Low (data tape is reliable)	Med (‘bit rot’)	Low (replication and monitoring)
Access Failures	Medium (people handle tapes)	Low/Medium (depends on system)	Low (automated checks)

How do you make a choice?



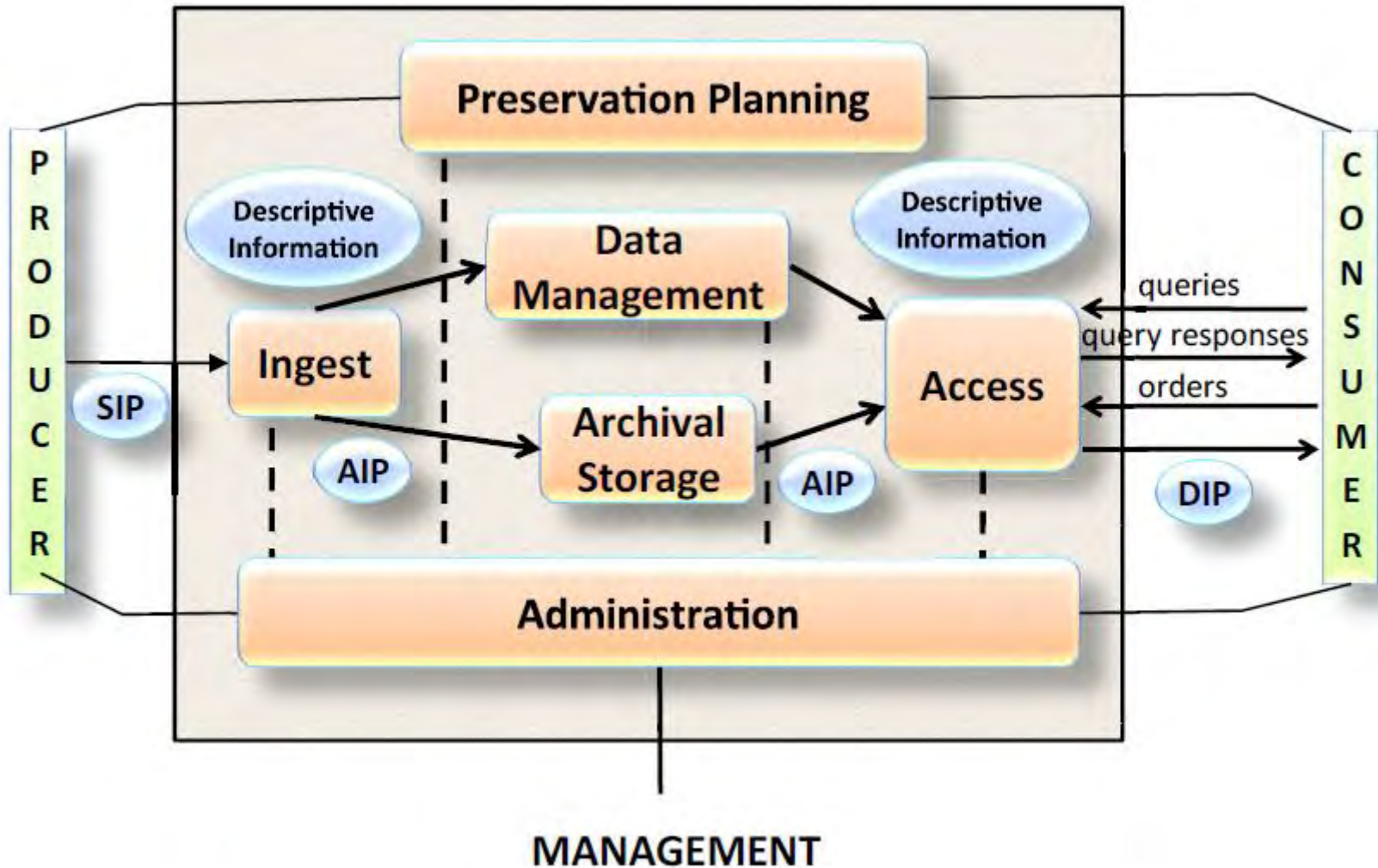
Digital Content Lifecycle

The DCC Curation Lifecycle Model



Digital Curation Centre

ISO 14721 (OAIS) 2012



Metadata standards



ISO 16363 (TDR) 2012



The Consultative Committee for Space Data Systems

Recommendation for Space Data System Practices

**AUDIT AND
CERTIFICATION OF
TRUSTWORTHY DIGITAL
REPOSITORIES**

DRAMBORA interactive

Digital Repository Audit Method
Based on Risk Assessment

A SIMPLE STRATEGY

Simple strategy

Descriptive
Technical
Preservation
Rights
Metadata Standards

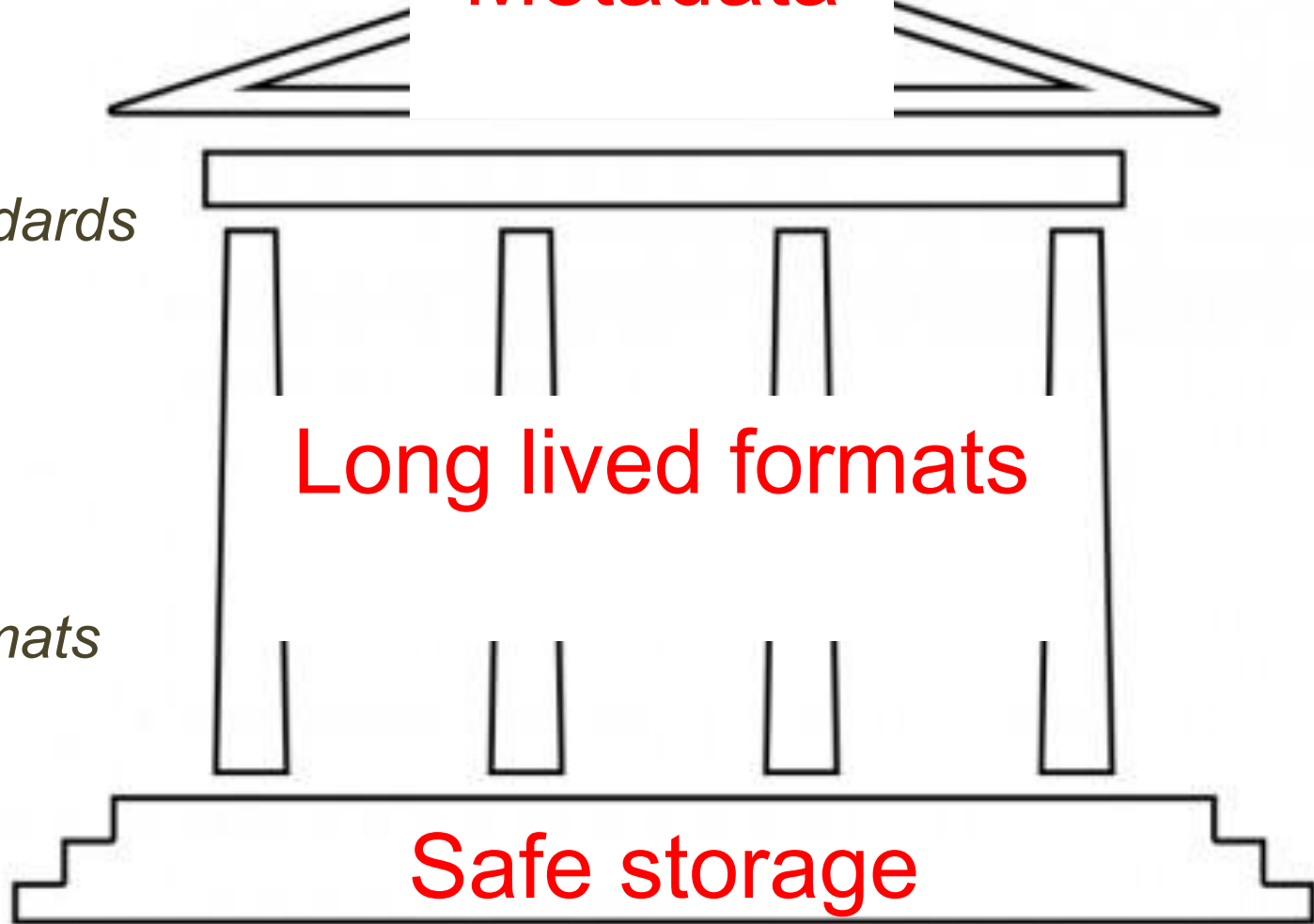
Metadata

Creation
Preservation
Access
Canonical Formats

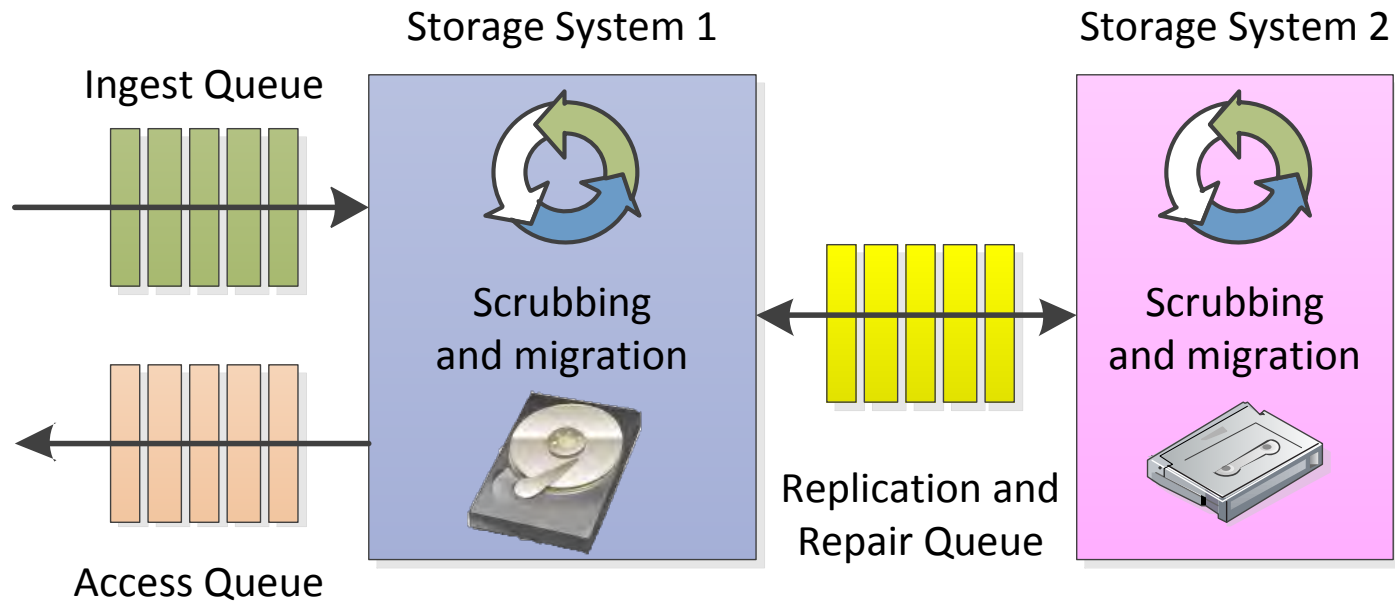
Long lived formats

*Optimise
for archive*

Safe storage



Active archiving



- Preservation best practice (diversity, intervention)
 - Multiple copies in different locations
 - Different technologies and different people
 - Active management: migration, integrity

Burnline

- Metadata and data on a file system
- Open standards and formats
- Drive down costs and risks



HOW ARKIVUM DOES IT

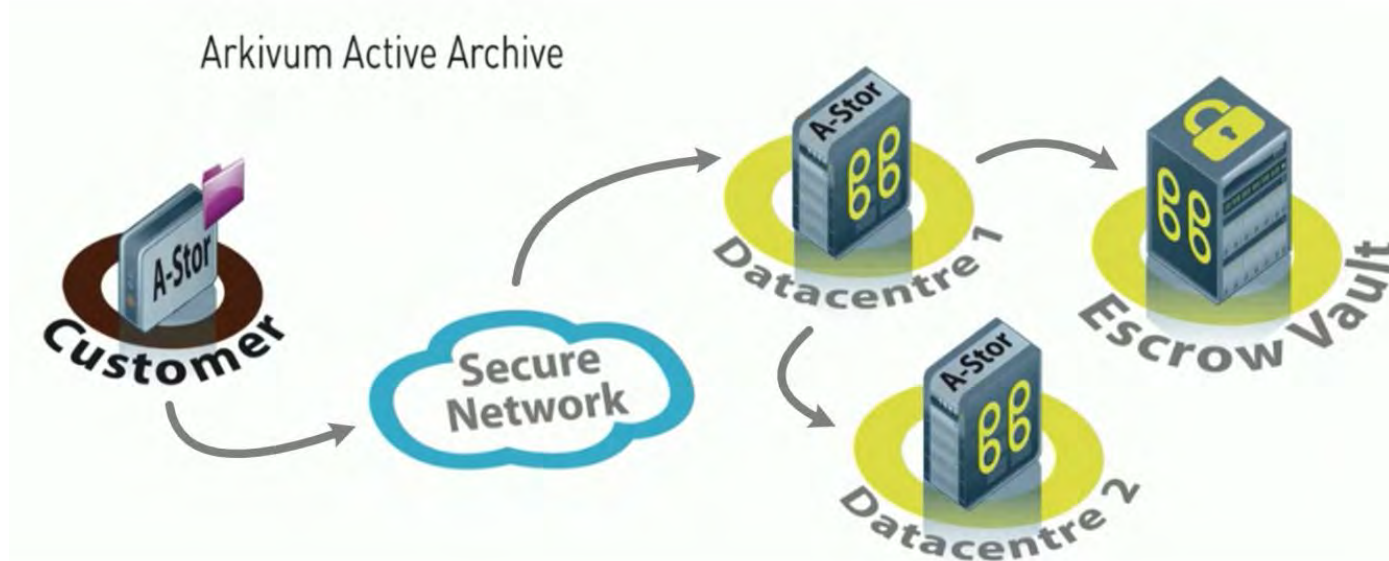
Arkivum

- Online data archiving as a service
- Spin-out of University of Southampton
- Decade of know-how working with archives
- Safe, secure, accessible data storage
- Designed from ground-up for retention and access

ARKIVUM
ASSURED ARCHIVING

100% data integrity guarantee
Keep your data safe & secure forever

Approach



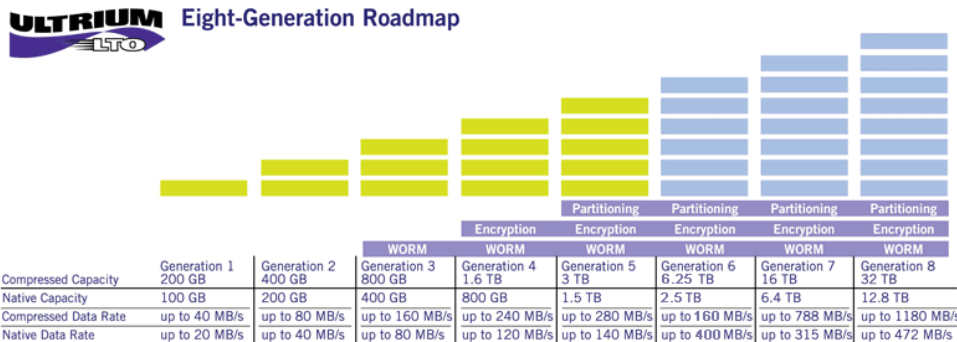
- Best practice: 3 copies, 3 locations, online and offline
- Active archiving: integrity, obsolescence
- Risk management: ISO27001, Drambora, TDR

100% data integrity guarantee

- All data is returned 'bit perfect'
- No restriction on time
- No restriction on volume
- Included in the SLA
- Worldwide insurance backed: £5M per loss event

Data escrow

- Copy of data offline at third-party escrow site
- LTFS on LTO tape with open source tools
- Customer has access to escrow copy:
 - If we fail to provide the service
 - If the customer decides to leave

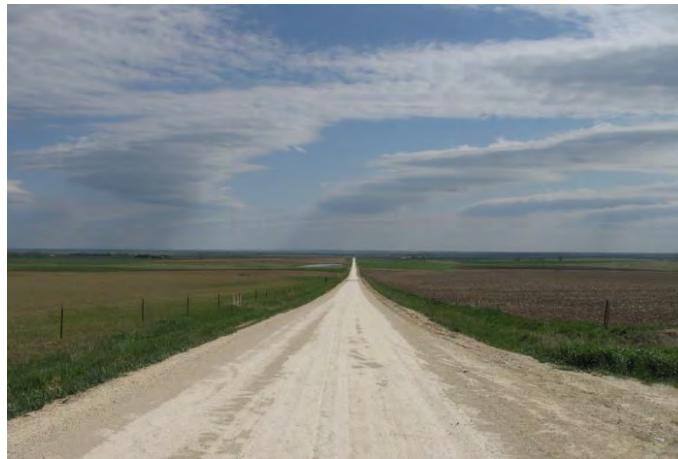


Note: Compressed capacities for generations 1-5 assume 2:1 compression. Compressed capacities for generations 6-8 assume 2.5:1 compression (achieved with larger compression history buffer).
Source: The LTO Program. The LTO Ultrium roadmap is subject to change without notice and represents goals and objectives only.
Linear Tape-Open, LTO, the LTO logo, Ultrium, and the Ultrium logo are registered trademarks of HP, IBM and Quantum in the US and other countries.



Pricing

- PAYG or Paid-Up for 5,10 or 25 year
- No ingress or egress charges
- Migrations and refreshes included
- Audits and certification included
- Escrow copy included, no exit costs



Thank you

- www.arkivum.com
- matthew.addis@arkivum.com

“Arkivum has helped us to create a robust archiving solution that will allow us to focus our budget on the business rather than yet more storage”

“Archived documents can then be seamlessly accessed from within the document management system, in the same way current documents are”.