



Technology Matters: A personal view

Robert Sharpe
23rd January 2014

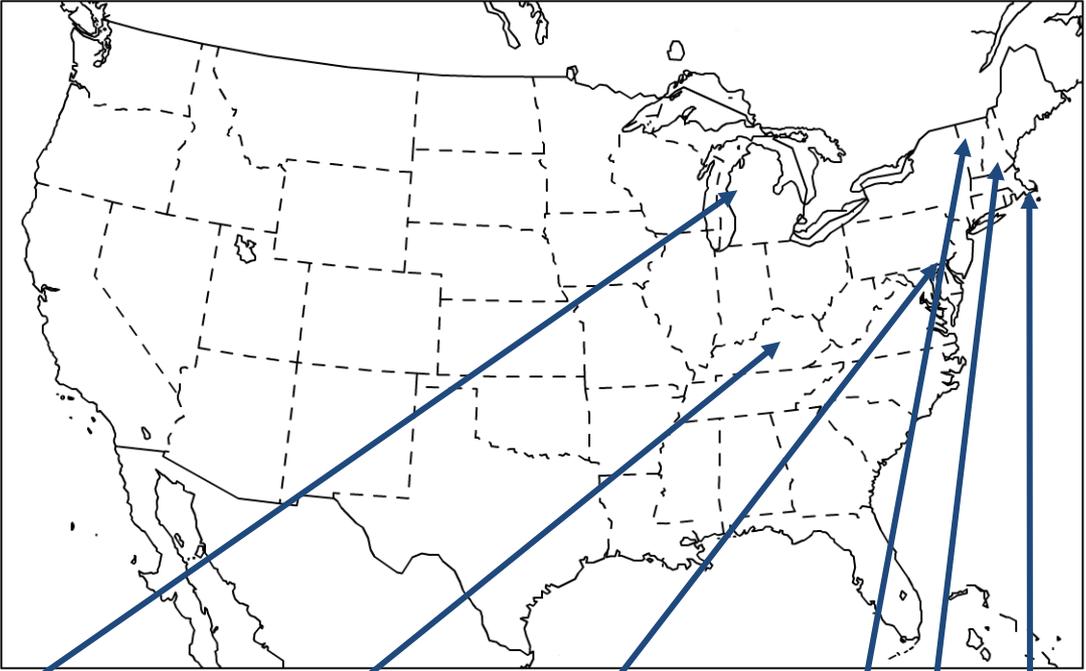
Agenda

- Quick introduction to Tessella
- What systems do I need for digital preservation?
- On-site & off-site systems
- How automatic can it be?:
 - Ingest example
- Scalability
- Conclusions

SDB: Digital Archiving Systems in 11 Countries across 4 Continents



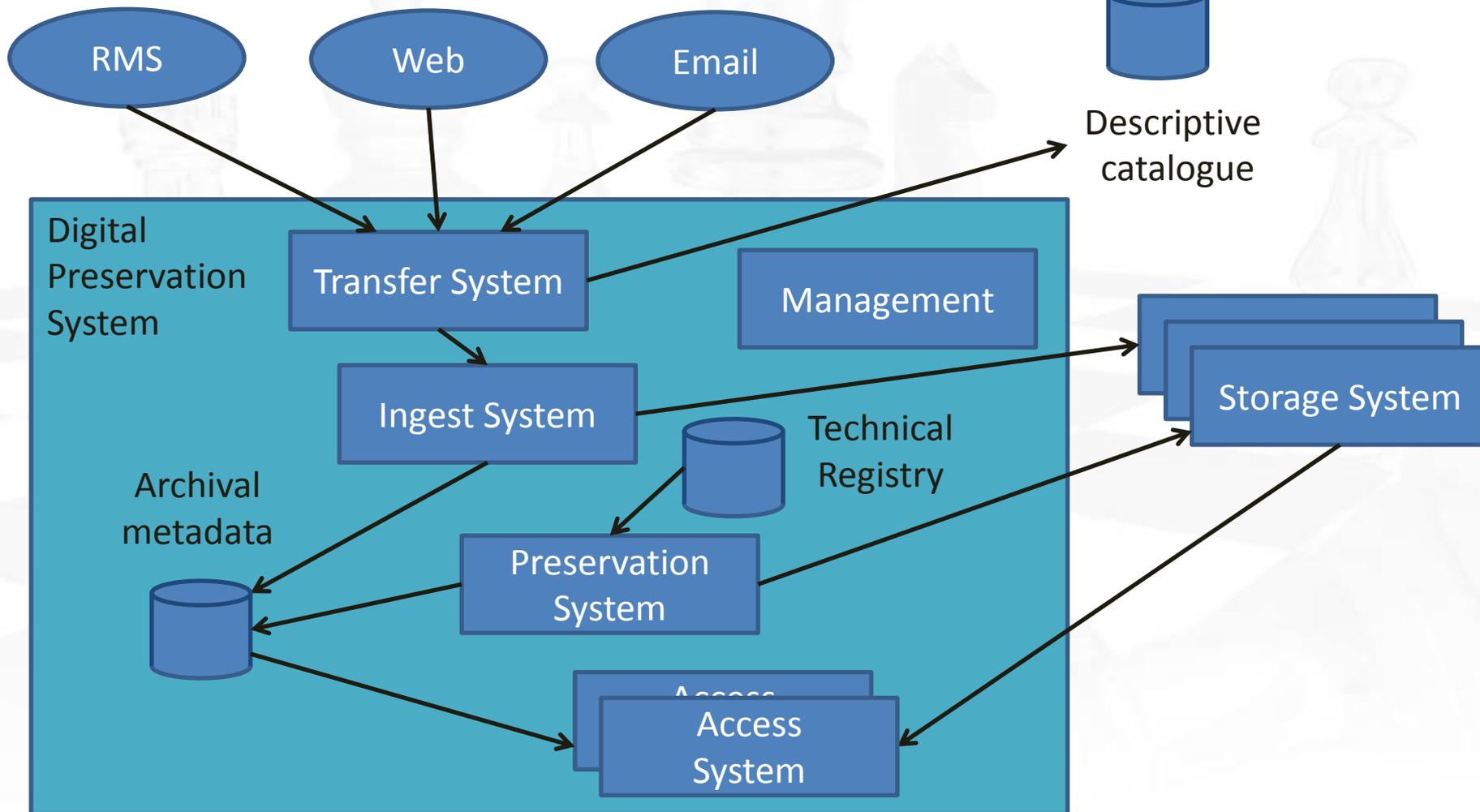
Preservica : Digital Preservation as a Service



Colby Sawyer College

What systems do I need?

External (source) systems



On-site Systems

- Pros:
 - Security concerns reduced
 - Have control
 - Easy (easier) to migrate to new system
 - Can customise (can considerably reduce manual overheads)
 - Bandwidth issues reduced
- Cons:
 - Need own hardware
 - Need own support
 - Costs more!
- Conclusion:
 - Good if have budget & bespoke needs
 - Good for large volumes

Off-site Systems (e.g., cloud)

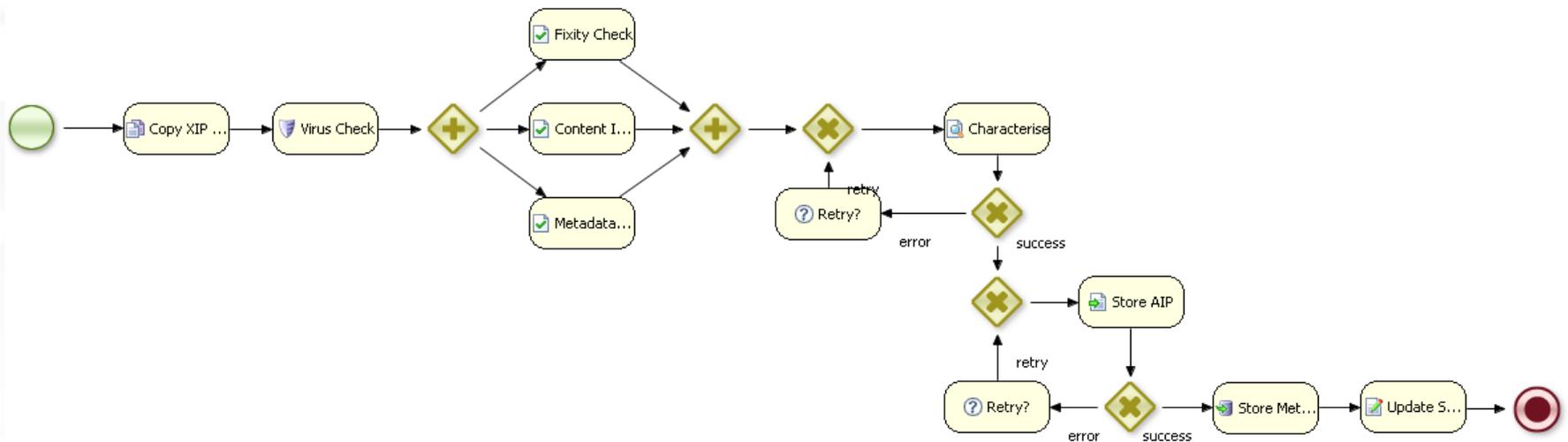
- Pros:
 - Cheaper
 - Low / zero up-front cost:
 - Don't buy hardware
 - Lower operational costs:
 - Shared support
 - Pay for what you need
- Cons:
 - Harder to customise
 - Possible security concerns
 - Bandwidth for large volumes
 - Harder to migrate
- Conclusion:
 - Good for low budgets
 - Limited ability to be bespoke

How automatic can it be?

- Golden rule:
 - Humans make judgements
 - Let software implement your judgements:
 - Will make less mistakes
 - Can be driven by machine-readable policy
- Sometimes lack of trust:
 - Good to test software
 - Once passed test, use it!
 - If issue occurs in production:
 - Fix it
 - Get your supplier to fix it

How automatic can it be? Ingest

- Capture human judgement as policy up front:
 - Decide what to keep?
 - Decide what to structure / catalogue?
 - Decide storage policy (how many copies to store)?
 - Decide which steps are necessary?



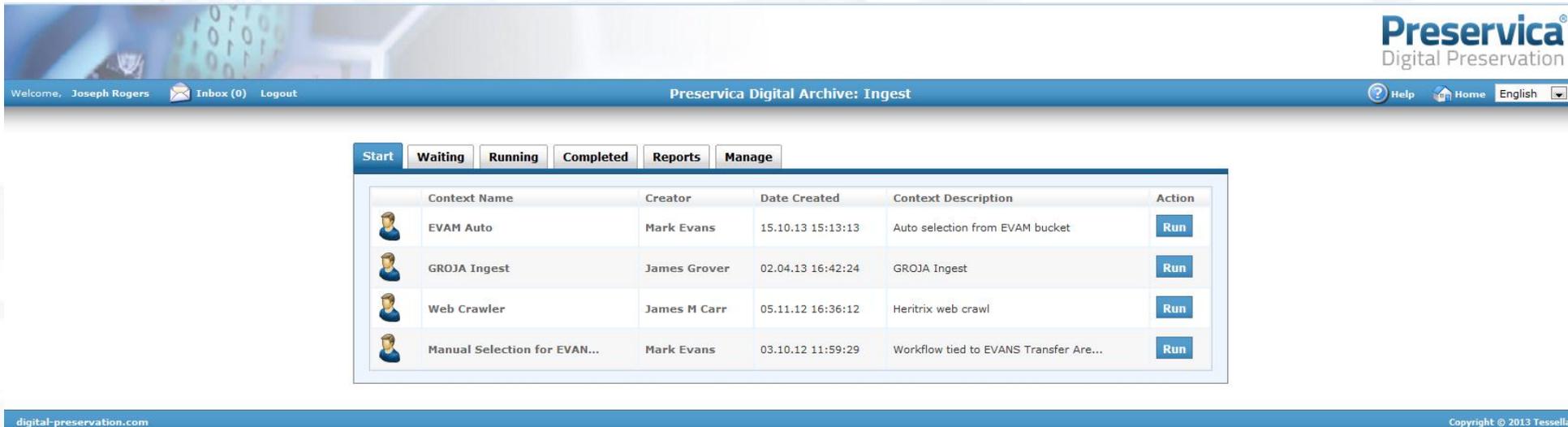
Real system example: Ingest

- In operation, let the software do its job:

The screenshot displays the Preservica Digital Archive Home page. The header includes the Preservica logo and the text 'Digital Preservation'. Below the header, a navigation bar shows 'Welcome, Mark Evans', 'Inbox (0)', 'Logout', 'Preservica Digital Archive: Home', 'Help', 'Home', and 'English'. The main content area features a navigation menu with the following items: Data Management, Preservation, Registry, Ingest (highlighted with a red circle), Storage, Administration, Search, Explorer, and Access. The footer contains the URL 'digital-preservation.com' and the copyright notice 'Copyright © 2013 Tessella'.

Real system example: Ingest

- Pick workflow to start:



The screenshot shows the Preservica Digital Archive: Ingest interface. The top navigation bar includes the Preservica logo, user information (Welcome, Joseph Rogers), an inbox icon with 0 items, and a logout link. The main content area features a table with workflow details and a 'Run' button for each entry.

Context Name	Creator	Date Created	Context Description	Action
EVAM Auto	Mark Evans	15.10.13 15:13:13	Auto selection from EVAM bucket	Run
GROJA Ingest	James Grover	02.04.13 16:42:24	GROJA Ingest	Run
Web Crawler	James M Carr	05.11.12 16:36:12	Heritrix web crawl	Run
Manual Selection for EVAN...	Mark Evans	03.10.12 11:59:29	Workflow tied to EVANS Transfer Are...	Run

- In fact even this is often automated:
 - Watch for arrival of complete SIPs

Real system example: Ingest

- Watch (if you want to):

Step Progress

State	Name	Progress	Started	Finished	Messages
	Select	<div style="width: 100%; height: 10px; background-color: green;"></div>	03.09.10 17:39:10	03.09.10 17:40:04	
	Copy XIP Package	<div style="width: 100%; height: 10px; background-color: green;"></div>	03.09.10 17:40:04	03.09.10 17:40:08	
	Fixity Check	<div style="width: 100%; height: 10px; background-color: green;"></div>	03.09.10 17:40:08	03.09.10 17:40:10	
	Metadata Integrity	<div style="width: 100%; height: 10px; background-color: green;"></div>	03.09.10 17:40:10	03.09.10 17:40:12	
	Content Integrity	<div style="width: 100%; height: 10px; background-color: green;"></div>	03.09.10 17:40:12	03.09.10 17:40:14	
	Characterise	<div style="width: 100%; height: 10px; background-color: green;"></div>	03.09.10 17:40:14	03.09.10 17:40:26	View
	Store Files	<div style="width: 100%; height: 10px; background-color: green;"></div>	03.09.10 17:40:26	03.09.10 17:40:36	
	Store Metadata	<div style="width: 100%; height: 10px; background-color: green;"></div>	03.09.10 17:40:36	03.09.10 17:40:38	
	Update Search Index	<div style="width: 100%; height: 10px; background-color: green;"></div>	03.09.10 17:40:38	03.09.10 17:40:42	

Real system example: Ingest

- Deal only with issues that the system can't:

Welcome, John R. Doe (Tenant : TESSELLA) Logout SDB Digital Archive: Ingest Home

Start Waiting Running **Completed** Reports Manage

Filter Workflows

Submission name	Collection Code	Top Level Record	Date Completed	Agency	Size	Files	Workflow Context
JPEG	JPEG	JPEG	03.09.10 17:40:42		12 KB	1	Manual Ingest
Manual Ingest			03.09.10 17:37:24			0	Manual Ingest

Help tessella.com Copyright © 2010 Tessella

How automatic can it be? Ingest

- Example issue: metadata impedance:
 - Source metadata:
 - Info in ERMS , e-mail system
 - Very little (e.g., web crawling)
 - Traditionally:
 - Translate to archival schemas (EAD etc.)
- Could manually map metadata:
 - As part of manual cataloguing
- Can automate:
 - Set up transform
- OR can bypass:
 - Embed original metadata
 - Use technology to view/edit/index/search without transform

Scalability

- Ingest: Series sequential steps
- Tool like DROID (format identification) typical time:
 - Small files: ~20s per 1000 files
 - Large files: ~ 8s per GB (c. 10TB per day)
- Large volumes:
 - Throughout more important than individual run speed
 - Need ability to run in parallel (multiple threads)
 - Automation important
 - Resilience important
- It can be done:
 - SDB ingests FamilySearch ingests at 50TB every day
 - Note doesn't need very expensive processing power:
 - 6 Application Servers @ c. \$5k each = \$30k
 - Ingest disk arrays and network Higher
 - Storage costs Dominant

Conclusions

- Try to minimise number of systems:
 - Will cost more in interfaces if you don't
- Choose system:
 - On-site / Off-site
- Archivists / Librarians / Curators are in charge:
 - Do what you are good at
 - Buy software / services to do the rest
- Automate everything that you can:
 - Use software that already does this
- Scale by system engineering:
 - Don't judge by speed of 1 thread on your desktop
- Lots of interesting issues to resolve:
 - But don't reinvent the wheel!