

Web Archives

A Storage & Content Management Use Case

Laura Graham

Library of Congress Web Archives Team

The Tools & Formats

- ❖ **Heritrix: Produces the Content:**
 - ❖ Developed by Internet Archive
 - ❖ Open Source Java App
- ❖ **Output File Format:**
 - ❖ ARCs
 - ❖ Warcs: ISO 28500:2009
- ❖ **Wayback “Viewer” Provides the Access:**
 - ❖ Developed by Internet Archive
 - ❖ Open Source Java App

Content: A “Plain Vanilla” Use Case

- ❖ **File & Content is fixed**

- ❖ What Heritrix writes = what we transfer = what we store = what we make accessible to the public
 - ❖ No master or derivatives
 - ❖ WARC files themselves may remain compressed

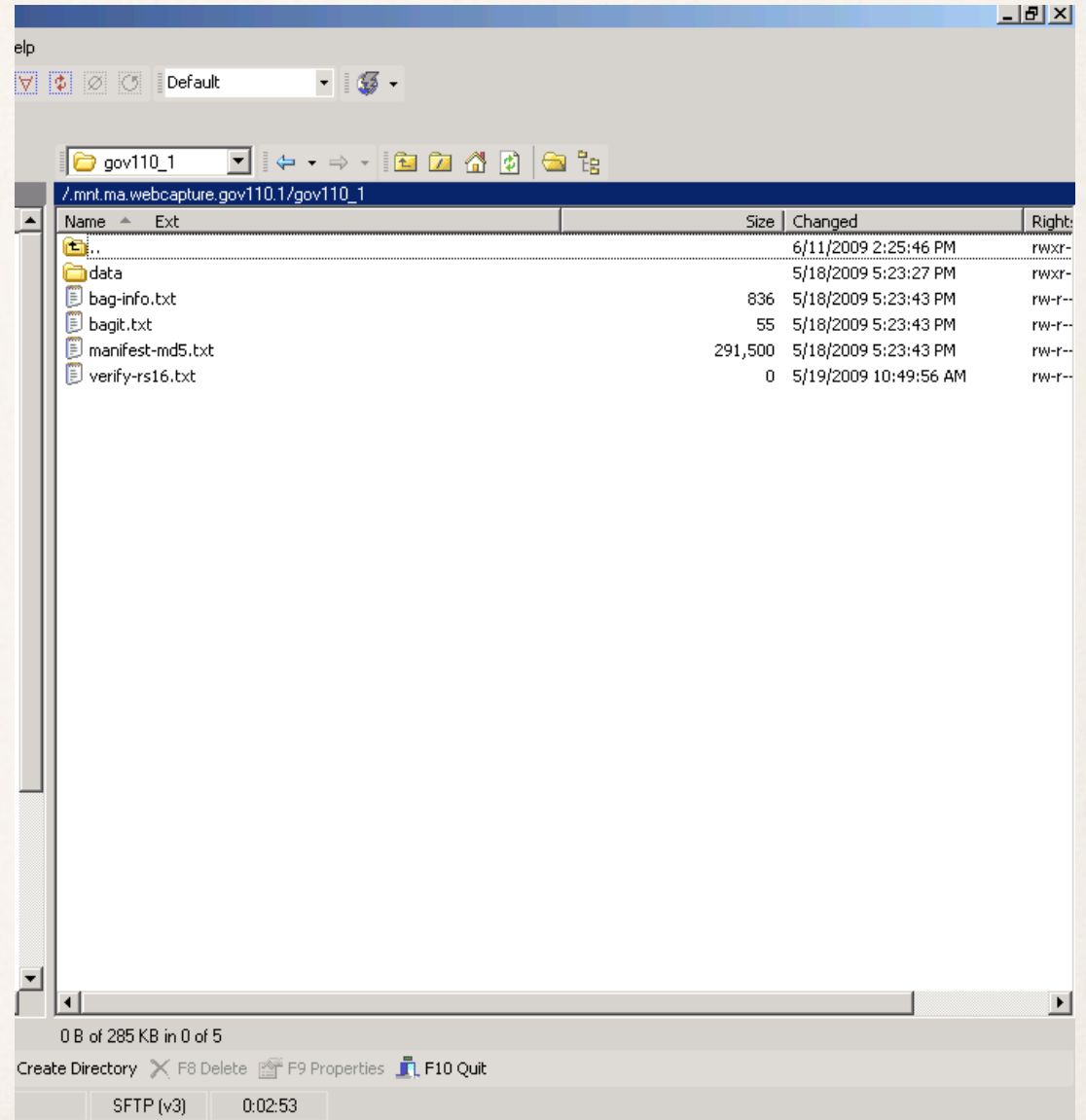
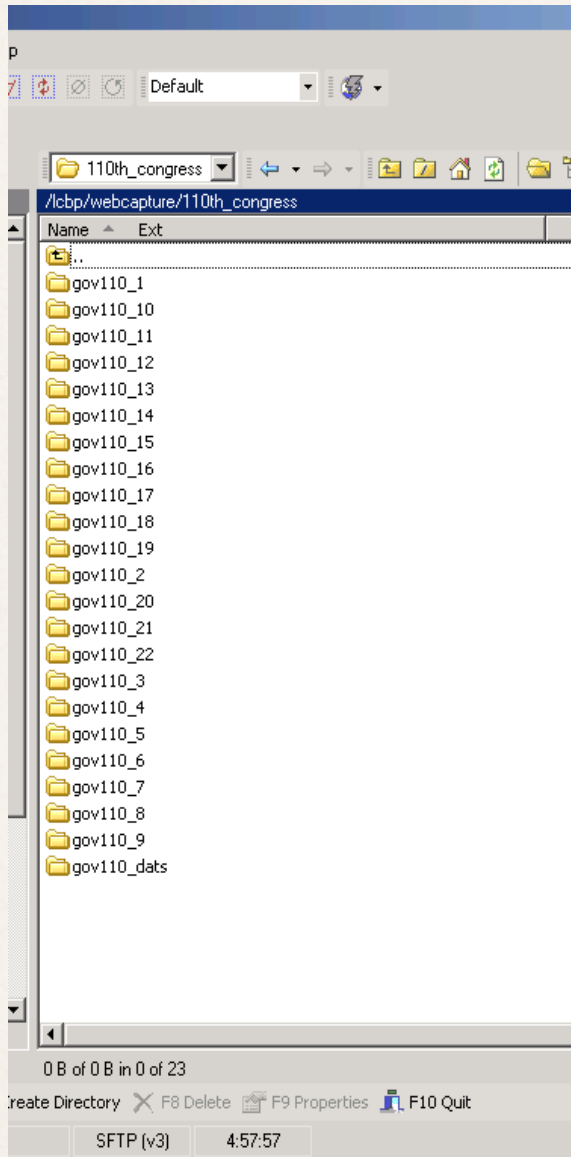
- ❖ **Organization**

- ❖ No items ... just files in buckets
 - ❖ Collection
 - ❖ Crawl Frequency

- ❖ **Bagit File Packaging Format**

- ❖ Bags persist from transfer to storage to access

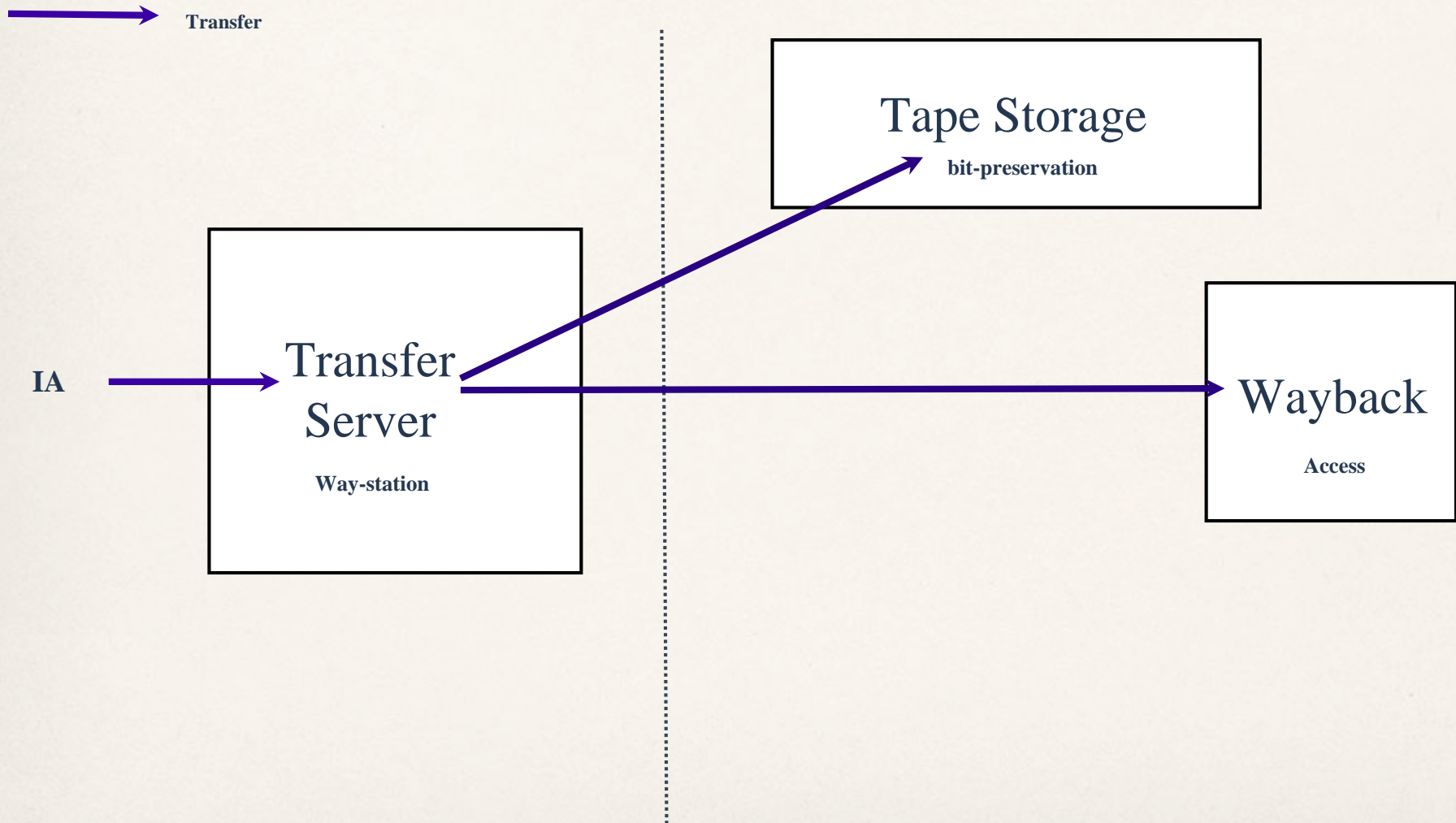
Bags and more bags ... !



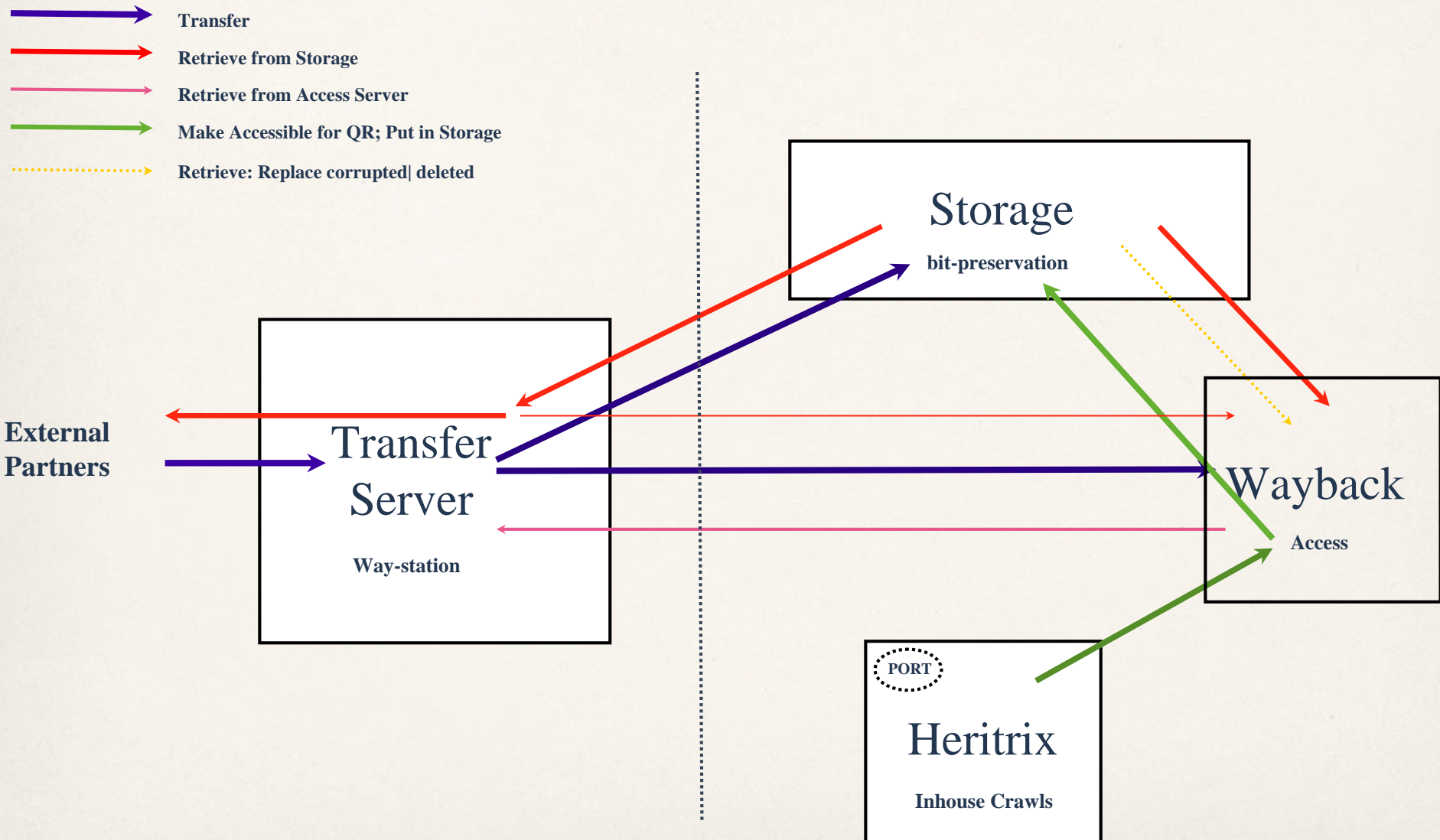
Scale of Content (Internet Archive)

Year	GBs	Change % - +	Transferred
2001	136		0
2002	10221	7415	0
2003	762	-92	0
2004	7771	919	0
2005	16364	110	0
2006	21557	31	0
2007	19003	-11	9391
2008	12863	-32	16569
2009	32875	155	38746
	121552		64706

Starting out in 2008...



...and here we are in 2009



Requirements (aka wish list) ...

- ❖ Easy & efficient ‘interoperability’ between machines
 - ❖ automation / elimination of “manual labor” steps
- ❖ Universal Tracking
 - ❖ workflow status & data attributes of everything everywhere ... not just in storage

Example: Copying & Verifying

- ❖ Copying Bags:

- ❖ Where possible, we'd like *not* to have to move stuff so much...



- ❖ Verifying Bags:

- ❖ Straightforward file fixity
- ❖ We want to know our content is safe
 - ❖ ...but we'd like to verify-on-the fly just about everywhere we go

Example: Repackaging

- ❖ Automated ‘Repackaging’ of Bags “en route” to target server
 - ❖ Workflows/Machines have different Requirements
 - ❖ Size:
 - ❖ Transfer: 1+ TB is efficient
 - ❖ Storage: ≥ 650 -800 GBs good for management
 - ❖ Access Server: 300 GBs upper limit
 - ❖ Source bag \rightarrow Destination bags
 - ❖ 1 to many
 - ❖ Many to 1
 - ❖ and more...

In Sum: Use Case = Interaction

- ❖ 2008: Web Archives began with a simple, one-directional workflow, and a focus on the content in the “storage box”
- ❖ 2009: Storage box now part of a complex of interactions, workflows, and machines
 - ❖ The User “lives” in that “complex”
 - ❖ Tools & utilities to manage the relationship between storage & its larger environment