# Format Migrations at Harvard Library
## General Framework and Plan Development

### Project Description
This project through the National Digital Stewardship Resident program involved designing a format migration framework for obsolete digital formats in the digital repository. The format migration framework strategizes how to manage and execute migration projects and documents the general process for preparing for and performing a format migration, including but not limited to the steps that need to be taken, the decisions that need to be made, key stakeholders to include, the types of research and testing that needs to be done, migration artifacts that should be preserved, and templates to facilitate this process. The framework was developed by working through several real use cases with the Library's Preservation Services staff: Kodak PhotoCD (the primary case study for this poster), RealAudio, and SMIL audio playlists.

## Overall Migration Framework

The most desired outcome from this project was a framework that could be used to develop a migration plan for any format in the repository. These are the essential components of the framework.

| Component | Activity |
|---|---|
| Project Start-Up: What Must Be in Place | Define stakeholder groups and inform of project initialization. Refer to framework and any relevant past projects to serve as guidance for setting up the project. Refer to general data management/migration plans for a sense of how to engage constituents and when (DataCave, NSW). |
| | Identify other parallel library projects that might impact the timing or functionality of the migration plan (e.g. changes to database, batch ingest processing tools, metadata schemas, etc.) |
| | Create a protected environment for necessary testing tools, ensure that all platform requirements are satisfied in order to adequately analyze the content (e.g. text editors, image viewers, etc.) Additionally, have a secure way to pull content from the repository and for accessing unique applications or software (obsolete or in some cases Harvard-specific) that reduces risks of virus or malware. |

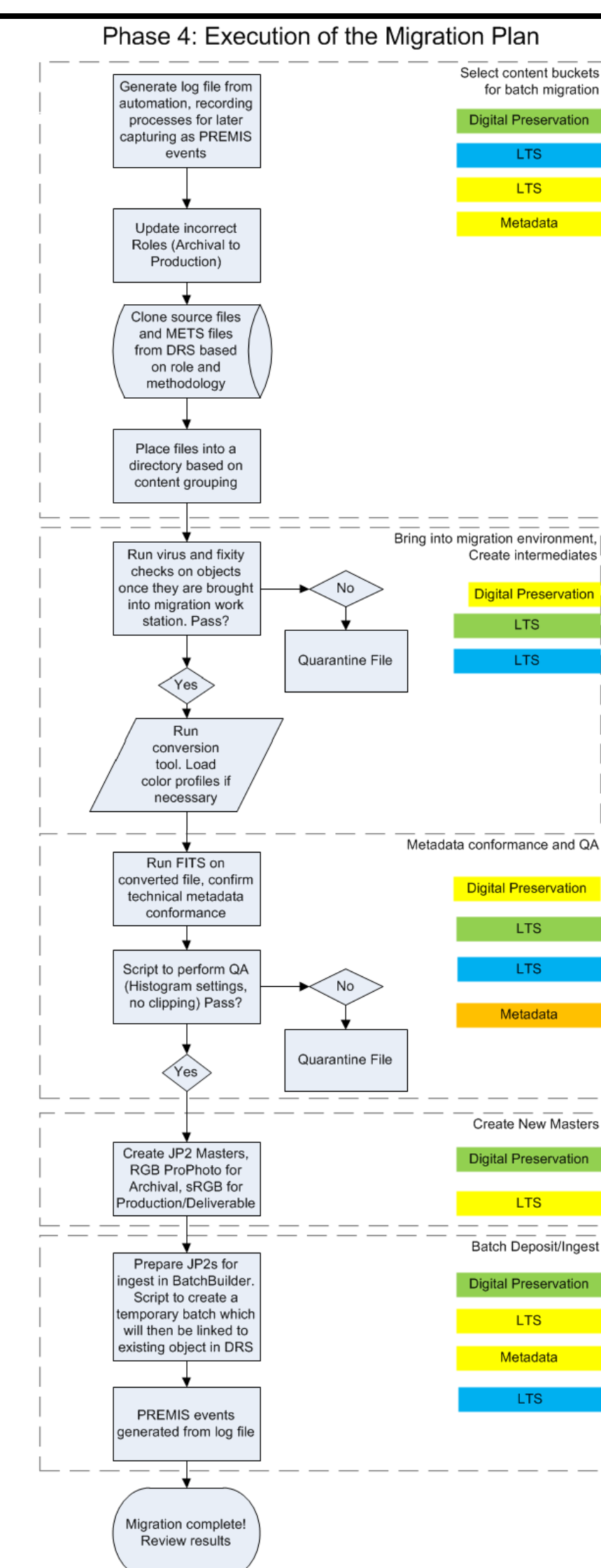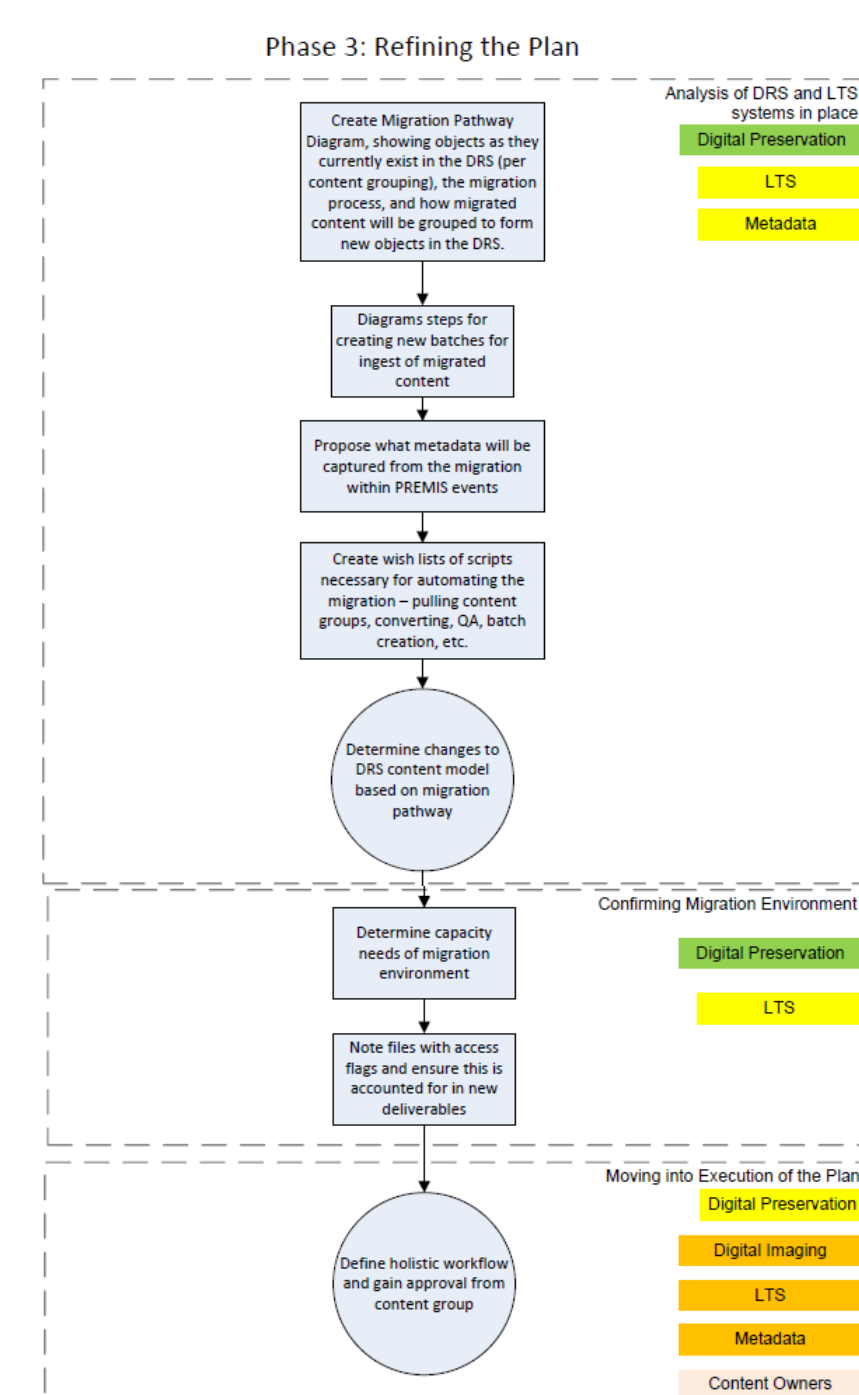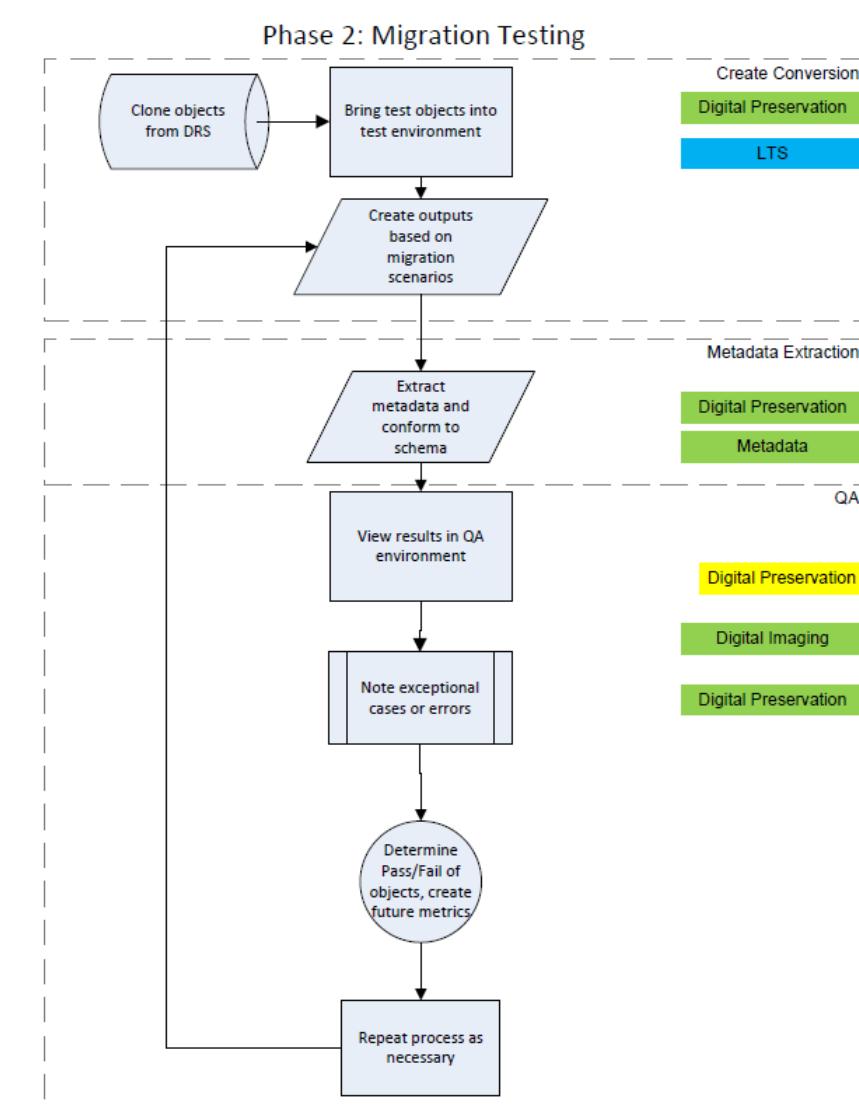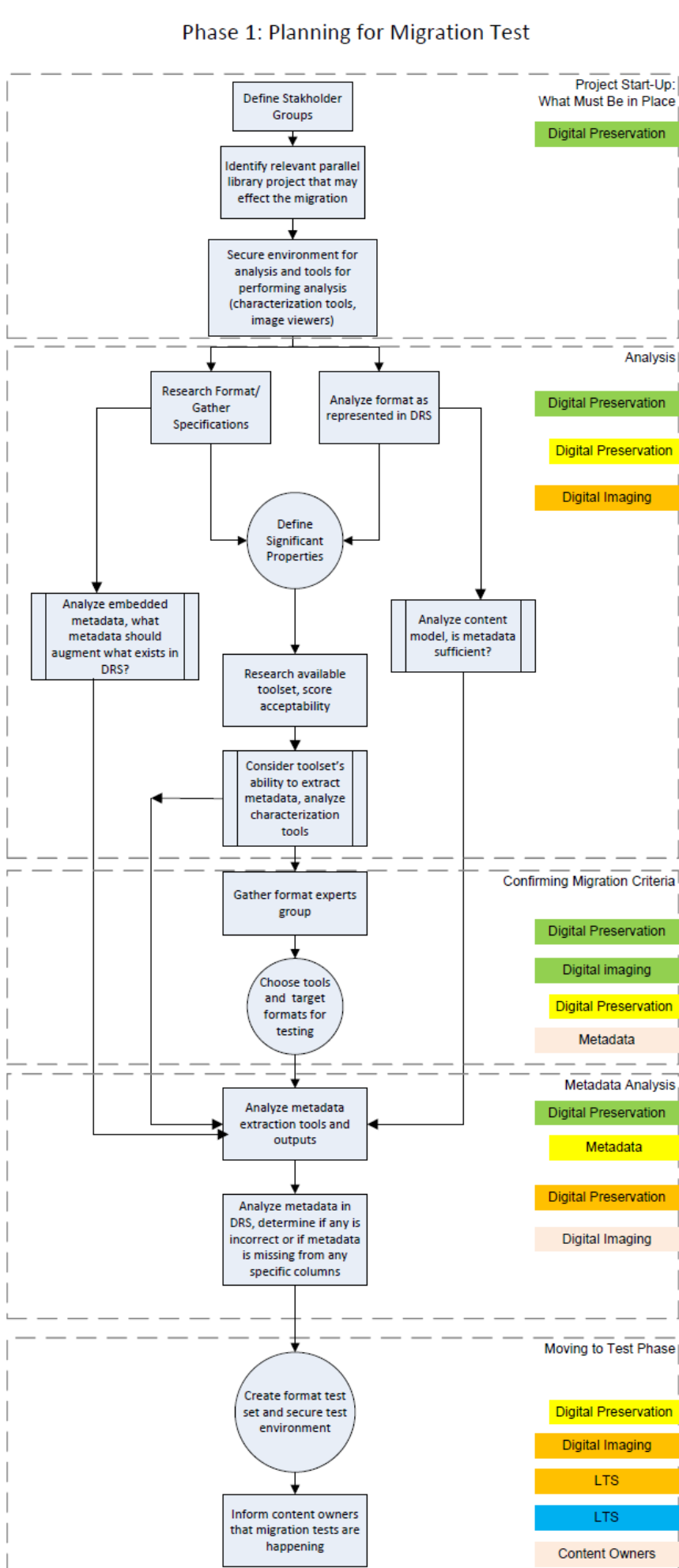| Component | Activity |
|---|---|
| Analysis | Research format specifications and additional literature around recommended tools, target format, and processes for successful migration. Distill essential significant properties. |
| | Analyze content in the DRS grouped by shared technical characteristics. Distill into a plan for performing tests based on definitive groupings of content within a format (useful metadata to consider might be Roles, Methodology, Relationships, and specific technical metadata). |
| | Additional questions to potentially address in analysis: <br> • How were the materials originally digitized <br> • What was the chain of custody (donor deposits, Harvard owners, vendors, ingest of digital material, etc. <br> • Any specific purpose or intent behind the objects to determine the appropriate preservation actions |
| Confirming Migration Criteria | Gather format expert groups and proposal for defining significant properties and migration criteria, confirm goals for successfully migrating content |
| Metadata Analysis | Assess existing metadata for content and identify any modifications that need to occur. Select appropriate characterization tools and compare against FITS output. Begin to assess any augmentation that needs to occur with the current schema and any migration-related metadata that should be added, for example a PREMIS event. |
| Moving into Test Phase | Determine the necessary OS, hardware, applications and other system requirements for the task. Secure the test environment and develop a strategy for creating a test set of content based on grouping, a variety of custom settings to apply for comparison, and metrics for performing quality assurance. |

## Project Workflow Example: Kodak Photo CD

The framework is used to guide the creation of a format-specific workflow


Phase 1: Planning for Migration Test

## Deliverables and Metrics: Kodak Photo CD
Key deliverables are identified at each step of the framework to serve as evidence of the decisions made along the process and to gain consensus across stakeholder groups.

### Stakeholder Involvement – Define Essential Roles Across the Project
- **Project Management** (Digital Preservation)
- **Technical Guidance/Format Expertise** (those who understand the format best – Imaging Services)
- **Analysis, Requirements and Specifications** (Digital Preservation, though some documents may originate from other departments)
- **Quality Assurance/Plan Approval** (Digital Preservation/Imaging Services)
- **Systems Conformance/Technical Infrastructure** (Library Technology Services)
- **Content Ownership** (curators or collection managers, involvement is generally just to be informed of major decisions)

**Analysis:** Distill information from format/tool specifications and files within the Digital Repository Service (DRS) to build test criteria

DRS – Content Groupings
Create essential "buckets" of content based on shared technical characteristics and object structures

Specifications – Significant Properties

Assessment of tools
Scoring acceptability of available conversion tools based on significant properties



YCC color space

Image Pac Compression

Scene Balance Algorithms for adjustments in light/color

### Phase 1: Planning for Migration Test

---

| Component | Activity |
|---|---|
| Create Conversion | Copy sample files from DRS based on defined content buckets, perform conversions with a variety of custom settings and target formats/wrappers/codecs. |
| Metadata Extraction | Extract metadata from both source file and target file, compare outputs and determine significant attributes to be captured. Compare with FITS output to determine that the tool can satisfactorily parse the file. |
| QA | Perform quality assurance with format experts on converted files based on defined metrics |

Phase 2: Migration Testing

Acquire PhotoCD-specific film terms, import into pcdMagic, generate several different file format outputs for comparison.

**Create several outputs**
- KODAK Photo CD 4050 E-6 V3.4
- KODAK Photo CD 4050 K-14 V3.4
- KODAK Photo CD Color Negative V3.0
- KODAK Photo CD Universal E-6 V3.2
- KODAK Photo CD Universal K-14 V3.2

Export To...    JPEG... / TIFF... / DNG...

Use Photoshop to check color balance (if image contains color bars). Compare gamut of histograms to find best range.

QA results
Kodak Color Negative
Kodak Universal E6

…Or channels balanced

Noting subtle differences – though seemingly similar, different film terms produce slightly different results.

### Phase 2: Testing

---

| Component | Activity |
|---|---|
| Moving into Executing the Plan | Based on test, define migration pathway with confirmed custom settings relative to content groupings. Diagram how migrated content will be added to objects. |
| | Define project roles based on the Migration Workflow Diagram and RACI model assignments. Circulate this to stakeholder groups and gain approval. |
| Analysis of DRS and LTS systems in place | Diagram modifications to content model, file-to-file relationships, and ingest processes in order to accommodate migrated content. |
| | Create and approve a plan with Library Technology Services for batch ingests, design method to pull files from DRS based on content groupings and bring into migration environment. |

| Component | Activity |
|---|---|
| Confirming Migration Environment | In confirming that the technical environment and processes for connecting them are in place, make sure the following questions are answered: <br> • What is the environment for the migration? What is the total filesize of everything that will be pulled, staged, generated, and re-packaged? <br> • What services are needed to transfer material and will any tools need to be built for this process? <br> • Does the chosen tool interface well with automated scripts such that migration can be performed based on designated content groups? <br> • How will QA be performed en masse? <br> • How will new files and their relationships be handled through batch ingest processes? <br> • Are there other external dependencies which are required to render the file (streaming service, application, etc.) or for it to function as an object within the DRS? <br> • Do any exceptions need to be made in separating content based on the Access Flag? This will be especially relevant for deliverable content. |

Phase 3: Refining the Plan

Migration Pathway and Batch Ingest diagram for Kodak PhotoCD

Diagram migration pathway of the content, paying particular attention to relationships between files (masters, derivatives, related metadata files) and object structures within the repository:

- How content looks now
- What happens to the content at each step of the migration
- How content will be structured such that it complies with repository ingest policies

Begin to think about what tools and services are needed to automate this process.

### Phase 3: Refining the Plan

---

| Component | Activity |
|---|---|
| Schedule Migration | Set timelines relative to existing projects and departmental workloads. Anticipate setbacks that could occur as well as deadlines that need to be met. |
| Custom development | Develop scripts for automating the process <br><br> **Prepare Migration Environment** <br> • Locate grouping of objects based on DRS criteria (e.g. Production Masters in a given format, methodology statements, relationship to a list METS record IDs, etc. <br> • Move grouping of objects to migration environment (e.g. async), mark status of object to avoid offline issues in DRS <br> • Prepare log file to record migration events, later for recording within PREMIS <br> • Determine process for recording file's relationship to its object such that new content can be joined with the object during ingest <br><br> **Perform Migration** <br> • Interface with conversion tool to convert files based on custom settings and outputs |

| Component | Activity |
|---|---|
| Custom development | QA <br> • QA results based on pass/fail metrics (e.g. no color clipping in Photoshop, audio file successfully plays in multiple browsers through SDS, etc.) <br><br> Ingest <br> • Create batch submissions with a particular focus on generating new metadata (from migration), creating new relationships with files in the DRS ('HAS_SOURCE', Deprecating replaced file), and running file fixity and characterization validation. <br><br> Reassign URNs for replacing deliverables (or create new ones), may also need to deactivate current URNs. Also ensure that new deliverables match the access flags of the deliverables that they are replacing. <br><br> Create batch deposit with batch control files designed for content buckets relative to their intended relationships structure within a content model |

Phase 4: Execution of the Migration Plan

Harvard's BatchBuilder tool for managing ingests into the DRS

**Necessary scripts and software augmentation to enable automatic workflow:**
- Clone files from DRS based on role and methodology, METS file to follow the file to later link with existing objects upon ingest.
- Batch convert in **pcdMagic** based on content groupings
- Run **FITS** on converted files to confirm technical metadata conformance
- **Photoshop** API scripting for QA (histogram checking) and conversion to target format (JP2, ProPhoto RGB for Archival Masters, sRGB for Production
- Devise strategy within **BatchBuilder** (ingest processing software) to create temporary batch ingest which is later linked with existing objects
- Reassign URNs, OSNs, and Access Flag from old deliverables to new, change status of old masters to "Deprecated"

Workflow Legend

Roles:
- Responsible
- Accountable
- Consulted
- Informed
- Technical Support

### Phase 4: Executing the Plan

---

| Component | Activity |
|---|---|
| Verify results Post-Ingest | Verify that migrations were successful (batch QA reports, ingest reports) |
| Encapsulate project Deliverables/Artifacts | Ensure that all deliverables/artifacts and project documentation is complete. Determine long-term disposition of the documents (what is deposited in DRS and what is kept in separate project folders in share drive). Confirm plan to deaccession any content. |

Joey Heinen
National Digital Stewardship Resident
Harvard Library
joeygheinen.jh@gmail.com
917-684-9495

HARVARD LIBRARY

NDSR BOSTON

### Phase 5: Project Wrap-Up