



Profiling the DPLA

A Project Introduction

Content

1. Introduction to the DPLA
2. DPLA Website and API
3. Project Goals and Objectives
4. Interesting Questions
5. Designing Our Work
 - a. System Architecture
 - b. Goals for Python Code
 - c. Project Map
 - d. Output Structure
6. Implementation: Challenges and Methods
 - a. Improving the Code
 - b. Problems Encountered
 - c. Outcome
7. Resources to Continue this Project

THE DIGITAL PUBLIC LIBRARY OF AMERICA (DPLA)

- ❖ Is an online portal for digital cultural heritage material
- ❖ Accepts material from a group of about 30 Providers (also called Hubs)
- ❖ These Providers get material from individual institutions, like libraries, archives and museums, known as Data Providers.
- ❖ An institution can be both a Provider and a Data Provider
- ❖ Data Providers can add metadata to items
- ❖ Providers (Hubs) can add metadata to items
- ❖ The DPLA adds metadata to items
- ❖ There is no way to identify which information came from a the Data Provider vs. the Provider
- ❖ The basic unit of information is an item



Exhibitions

View all »



Explore by Place



Explore by Date

Timeline »

1946 1947 1948 1949

News

New Exhibition! Battle on the Balkans: US Presidential Elections Oct 18

A Wealth of Knowledge

explore 14,438,423 items from libraries, archives, and museums

Apps

The DPLA is a platform. Developers make apps that use the library's data in many different ways. Here are just a few. [App Library](#) »

DPLA API Results:
Displayed in Postman (Below)

DPLA Homepage (above)

```

4616     }
4617     ],
4618     "facets": {
4619     "sourceResource.collection.title": {
4620     "type": "terms",
4621     "missing": 4415285,
4622     "total": 12315907,
4623     "other": 6527411,
4624     "terms": [
4625     {
4626     "term": "Records of the National Aeronautics and Space Administration",
4627     "count": 638062
4628     },
4629     {
4630     "term": "Botany",
4631     "count": 468836
4632     },
4633     {
4634     "term": "DPLA: Include in Digital Public Library of America",
4635     "count": 456748
4636     },
4637     {
4638     "term": "Flowering plants and ferns",
4639     "count": 404686
4640     },
4641     {
4642     "term": "Anthropology",
4643     "count": 285233
4644     },
4645     {
4646     "term": "Records of the Office of the Secretary of Defense",
4647     "count": 269704
4648     },
4649     {
4650     "term": "War Department Collection of Confederate Records",
4651     "count": 229825
4652     },
4653     {
4654     "term": "Collections Online",

```

Goal and Objectives

What We Hope to Do:

Extract the data available from the DPLA. Compile the data in a new manner, which will allow us to answer interesting questions about how metadata attributes are being used... or not being used.

What We Hope to Accomplish:

Provide insights to the usage of metadata attribute fields that will allow institutions providing the data to leverage their metadata creation in a manner that will result in the most exposure of their items.

In order to meet our goals and objectives we identified some interesting questions we wanted to answer about the usage of metadata in the DPLA. These questions guided the features we choose to implement in our code. These questions all helped us understanding of how the “collection” attribute was being used by Providers

Interesting Questions

How many providers has the DPLA collected from? And who are they?

How many collections does each provider contribute to?

How many items has each provider contributed?

What's the distribution of the number of collections referred to by an item?

How many items have the same number of collections?

How many items does each collection share with other collections?

How many collections are referred to by more than one data Provider?

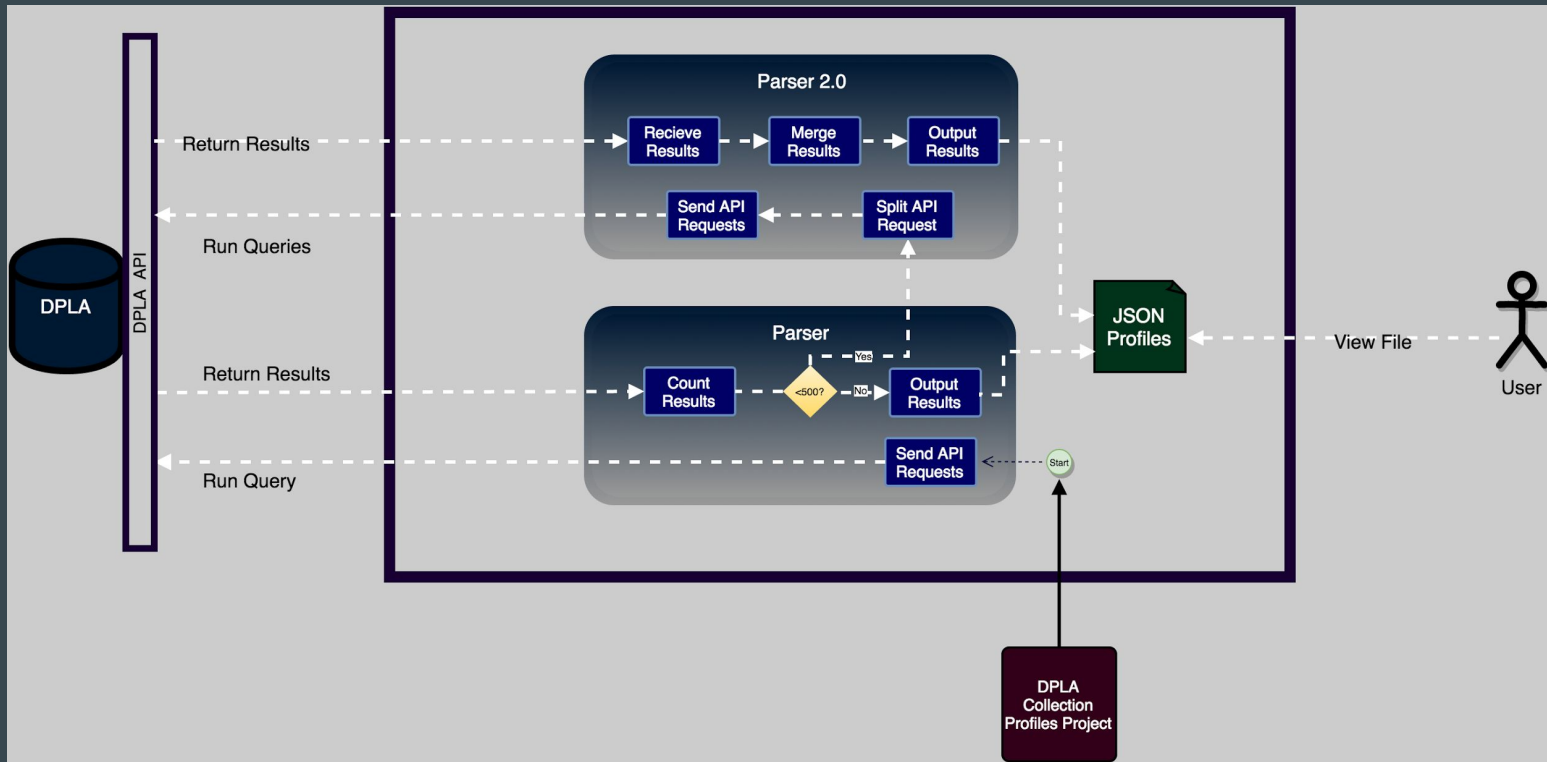
What collections share the most items?

How many items are not part of a collection?

Which providers / data providers are giving items with no collections?

Documenting Our Work

System Architecture Diagram



Goals for the Python Code

Features of the Code:

For Each Provider:

- Count of items contributed
- Count of collections used
- Count of items with no collection
- Count of items in each collection
- Count of dataProviders
- Item details: id, collection & dataProvider

The Code Should:

- Be modular and expandable
- Be efficient (Have no redundancy)
- Be clean and readable
- Use variables with meaning

Project Map

CCD Data Object and Property Labels

Variable Label	Property	Usage/Description	Level	Type
CCD_OBJ_LBL_adminDetails	adminDetails	holds the analysisDate and analysisTime properties	File	class
CCD_PROP_LBL_analysisDate	analysisDate	day file was generated	File	
CCD_PROP_LBL_analysisTime	analysisTime	time file was generated	File	
CCD_OBJ_LBL_dplaData	"dplaData"		DPLA	class
CCD_OBJ_LBL_providerData	"providerData"	object holding the providerId and the providerName	Provider	object
CCD_PROP_LBL_providerId	"@id"	URI for the provider page the DPLA API	Provider	string
CCD_PROP_LBL_providerName	"name"	Human-readable version of provider name	Provider	string
CCD_PROP_LBL_itemCountByName	"itemCountByName"	Count of total DPLA items by Provider name	Provider	numeric
CCD_PROP_LBL_itemCountById	"itemCountById"	Count of total DPLA items by Provider id	Provider	numeric
CCD_PROP_LBL_itemsInCollections	"itemsInCollections"	Number of items in collections by Provider	Provider	numeric
CCD_PROP_LBL_collectionCount	"collectionCount"	Count of Collections by Provider	Provider	numeric
CCD_PROP_LBL_dataProviderCount	"dataProviderCount"	Count of Data Providers by Provider	Provider	numeric
CCD_PROP_LBL_providerCountByName	"providerCountByName"	Count of providers based on provider.name field	Provider	numeric
CCD_PROP_LBL_providerCountById	"providerCountById"	Count of providers based on provider.id field	Provider	numeric
CCD_PROP_LBL_missingItemsByName	"missingItemsByName"	Missing items by Provider Name	Provider	
CCD_PROP_LBL_missingItemsById	"missingItemsById"	Missing Items by Provider Id	Provider	
CCD_PROP_LBL_itemCount	"itemCount"	Count of total number of items	DPLA	numeric
CCD_PROP_LBL_dataProviders	"dataProviders"	List of Data Providers by Provider	Provider	
CCD_PROP_LBL_provider	"provider"	Name of Provider	Provider	string
CCD_PROP_LBL_providerItemCount	"providerItemCount"	Number of items contributed by a provider	Provider	numeric

The Project Map identifies all the variables and file names used in the code and how those variables are used/ defined.

Provider Level Collection Details JSON Object Structure

```
"-providerColl"  
  /object  
    providerItemCount  
    noCollectionCount  
    provider  
    providerCollectionCount  
  /Collection  
    CollectionTitle  
    itemCount  
    dataProviderCount  
    providerDataProviderCount
```

Item Level Item Details JSON Object Structure

```
"-providerItem.json"  
  /object  
    /itemsInCollections  
      itemID : collectionTitle  
    provider
```

Provider Level Data Providers Details JSON Object Structure

```
"-providerDP"  
  /object  
    dataProviderCount  
    noDataProviderItemCount  
    provider  
  /dataProviders  
    itemCount  
    dataProvider
```

Implementation: Challenges and Methods

Improving the Code

- ❖ Previously the code...
 - Was long and clunky
 - Had duplicate functions in different files
- ❖ It was refactored to:
 - Increase modularity
 - Reduce redundancy
 - Add new functions

```
> def profile_dpla(base_dpla_filename):  
  
    # return 1 if the collection volume field>0, if filed=0 mean there is no usage of this field  
> def get_usage(collection_volume_field):  
  
    # get the volume for the filed in the location  
> def get_item_volume(field, location):  
  
> def have_or_not(a):  
    Begin code execution here  
    # ""  
  
    # setup argument parser  
    parser = argparse.ArgumentParser(description='please insert parameters')  
> parser.add_argument("-f", "--filename",
```

Problems We Encountered

- ❖ Default limit of 50 facets returned per query.
- ❖ Facet_size limitation is 2000
- ❖ Duplication of items in query results

The Code Now:

Is clean and readable

Is modular and extensible

Identifies all the desired features

Generates output with the data necessary to answer all the interesting questions.

Provider Level Details: Output File

```
▼ object {2}
  totalItems : 17620
  ▼ collectionCount {9}
    Chinese Rubbings Collection : 1104
    Studies in Scarlet: Marriage and Sexuality in the US and UK, 1815-1914 : 332
    Medieval Manuscripts at Houghton Library : 226
    Daguerreotypes at Harvard : 3362
    Latin American Pamphlet Digital Collection : 5726
    Dying Speeches and Bloody Murders: Crime Broad­sides, 1707-1891 : 502
    Colonial North American Project at Harvard : 4737
    Emily Dickinson Archive : 1531
    The Artemas Ward House and Its Collections : 100
```

Resources to Continue this Project:

The `DPLA_Profiles_Documentation_and_Resources_Guide`, located in UT Box has a complete inventory of all project files, locations and content. Data needed to continue the project is located in 3 places:

Github: `yanxian0924/DPLA_metadata` (publicly available)

UT Box Folder: `DPLA Profile Creation Documentation` (need permission from Unmil)

Class Server: `saab.ischool.utexas.edu/export/home/u09/unmil/ComputationalCollectionDescriptions`
(need permission from Unmil)