

**Organization:** Digital Public Library of America (DPLA)

**Project title:** DPLA Resources and Vocabulary Enrichment for Analytics

**Primary mentor:** Gretchen Gueguen

**Supporting mentors:** Emily Gore, Michael Della Bitta, Audrey Altman



<b>Description (2 project options)</b>	DPLA's data has more than 3 million unique subject headings, with only a portion coming from controlled vocabularies. Issues arise when records use slight variations of terms, or synonyms for the same concept. <b>Project 1:</b> Develop and test an effective method for analyzing record content and matching content, including keywords, with relevant controlled term from a defined list (LCSH or potentially Wikidata). Methods used will create a consistent vocabulary to aid users, and be replicable as a data source as it is re-ingested into DPLA. <b>Project 2:</b> Extends previous work focusing on the format terms in DPLA, to analyze this data and a vocabulary drawing from a select set of 30 terms from the <i>Art and Architecture Thesaurus</i> (a linked data ontology supported by the Getty Research Institute). The mechanics of this project are similar to project 1, however, the scope of the target vocabulary has already been identified.
<b>Problems</b>	Address challenges that exist when records use slight variations of terms, or synonyms for the same concept.
<b>Techniques</b>	Term extraction/matching, using NLP and potentially machine learning
<b>Data</b>	DPLA metadata records
<b>Outcome</b>	A method and an approach for identifying term variations and applying terminology more consistent terminology to support analytics