



**Organization:** OCLC Research

**Primary mentor:** Devon Smith, Lead Engineer, OCLC Research

Project Title	<b>Metadata Record Similarity: Identification and Clustering</b>
Description	Clustering within the WorldCat bibliographic utility en masse is often computationally infeasible since the data is so large and complex. This project will focus on finding ways to "fuzzy cluster" MARC (MACHINE readable cataloging) metadata records so that computationally expensive processes can be run on the resultant, smaller fuzzy clusters. In a very similar vein, when given an exemplar, quickly finding records that are similar to the exemplar will also be explored.
Problems/ Research Questions	How to quickly process large volumes of data, exclude "obviously" non-matching records, and include the most similar records.
Techniques	Any/all of string similarity metrics, match scoring, fingerprinting, locality sensitive hashing, Kolmogorov complexity. Most any approach that can be justified by some literature would be acceptable.
Tools/ Languages used	R or Python (libraries to be identified)
Data	Description: WorldCat Data Type: MARC (MACHINE readable cataloging) Data Size: Sample
Outcome	An algorithm and an implementation thereof that can be used to quickly identify likely match candidates or small fuzzy clusters for further, more computationally expensive processing. This outcome will feed into several ongoing research

	<p>initiatives at OCLC, furthering our goals of increasing the quality and utility of our data.</p>
Milestone Timeline	<ol style="list-style-type: none"> <li>1) Week 1: <ul style="list-style-type: none"> <li>• Establish environment</li> <li>• Preliminary data analysis and exploration</li> </ul> </li> <li>2) Month 1: <ul style="list-style-type: none"> <li>• Literature review</li> <li>• Experiment with similarity functions and clustering algorithms</li> <li>• Develop evaluation plan</li> </ul> </li> <li>3) Months 2-5: <ul style="list-style-type: none"> <li>• Iteratively develop, experiment, and evaluate a similarity function and a clustering algorithm</li> </ul> </li> <li>4) Month 6: <ul style="list-style-type: none"> <li>• Prepare and present results.</li> </ul> </li> </ol>
References	<ul style="list-style-type: none"> <li>• Mchine readable cataloging, see: <a href="https://www.loc.gov/marc/">https://www.loc.gov/marc/</a>, and explore the “Bibliographic” link, first link on left-hand, ToC (table of contents).</li> <li>• Worldcat, explore: <a href="https://www.worldcat.org/">https://www.worldcat.org/</a></li> </ul>