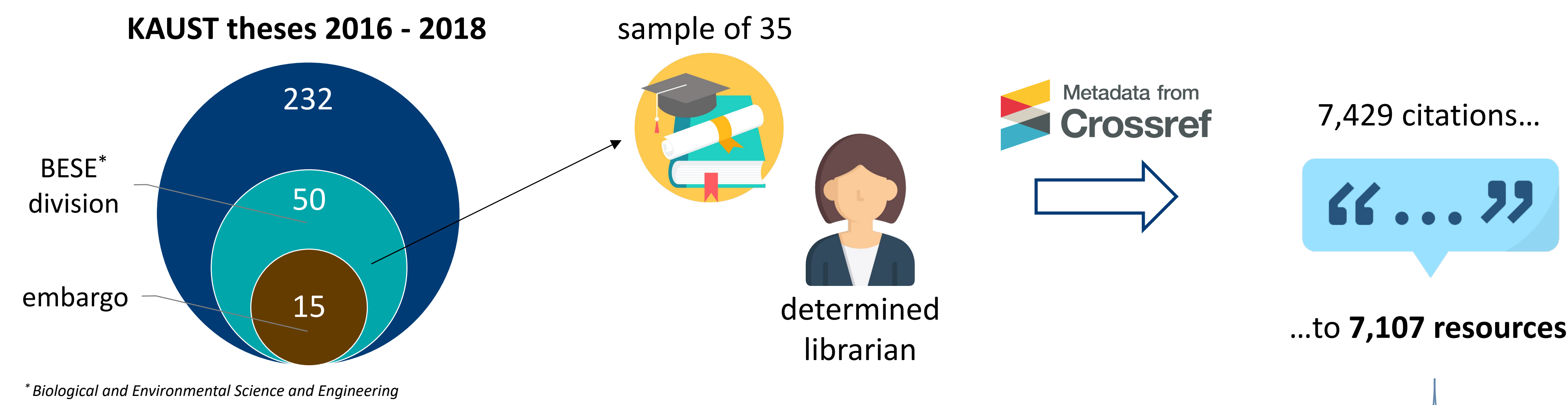


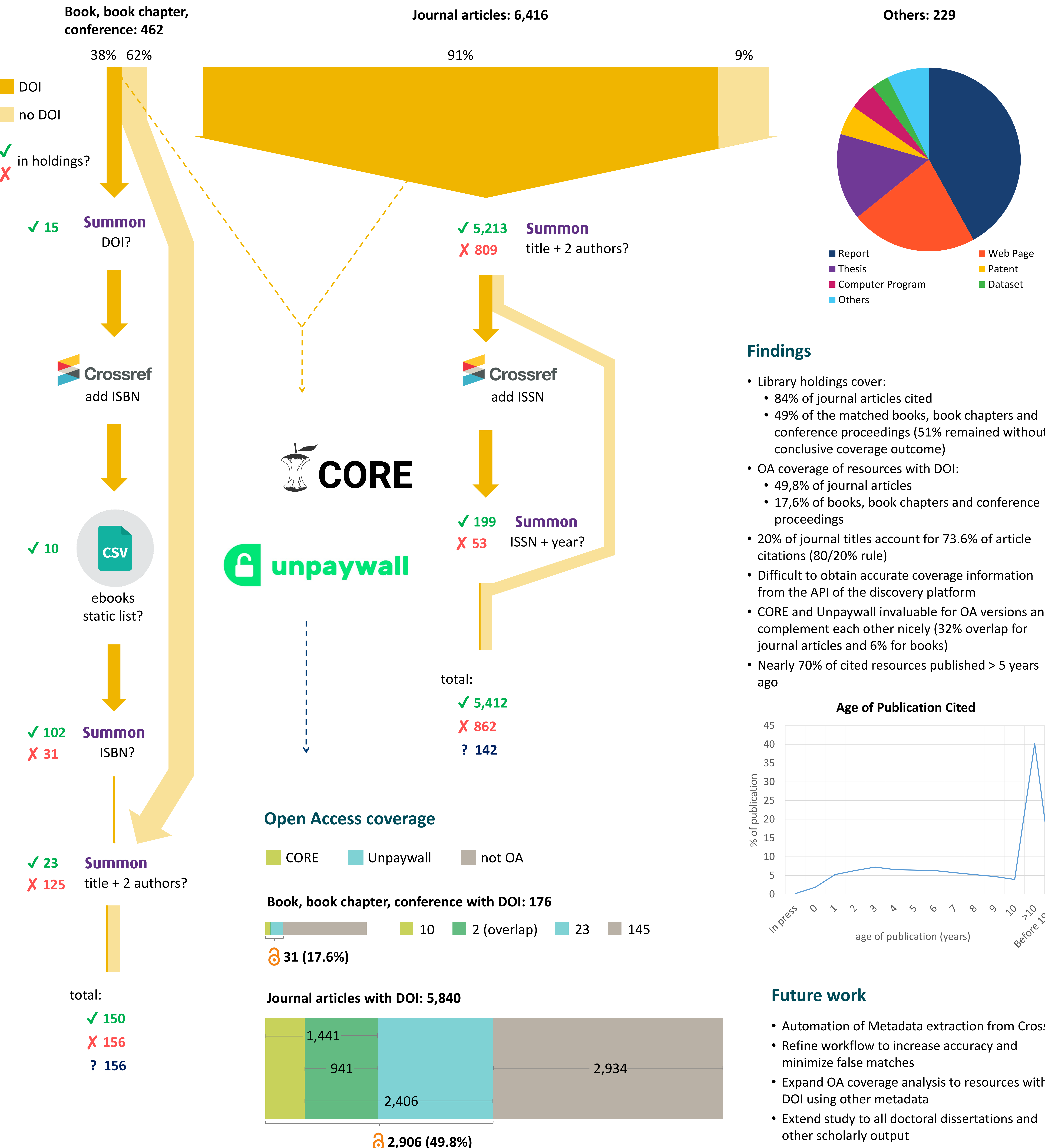
# Data mining of Citations in Theses: a workflow for automated analysis of Open Access and library holdings coverage

Jose I. Martín Cuesta (Systems Specialist), Lee Yen Han (Subject Specialist)



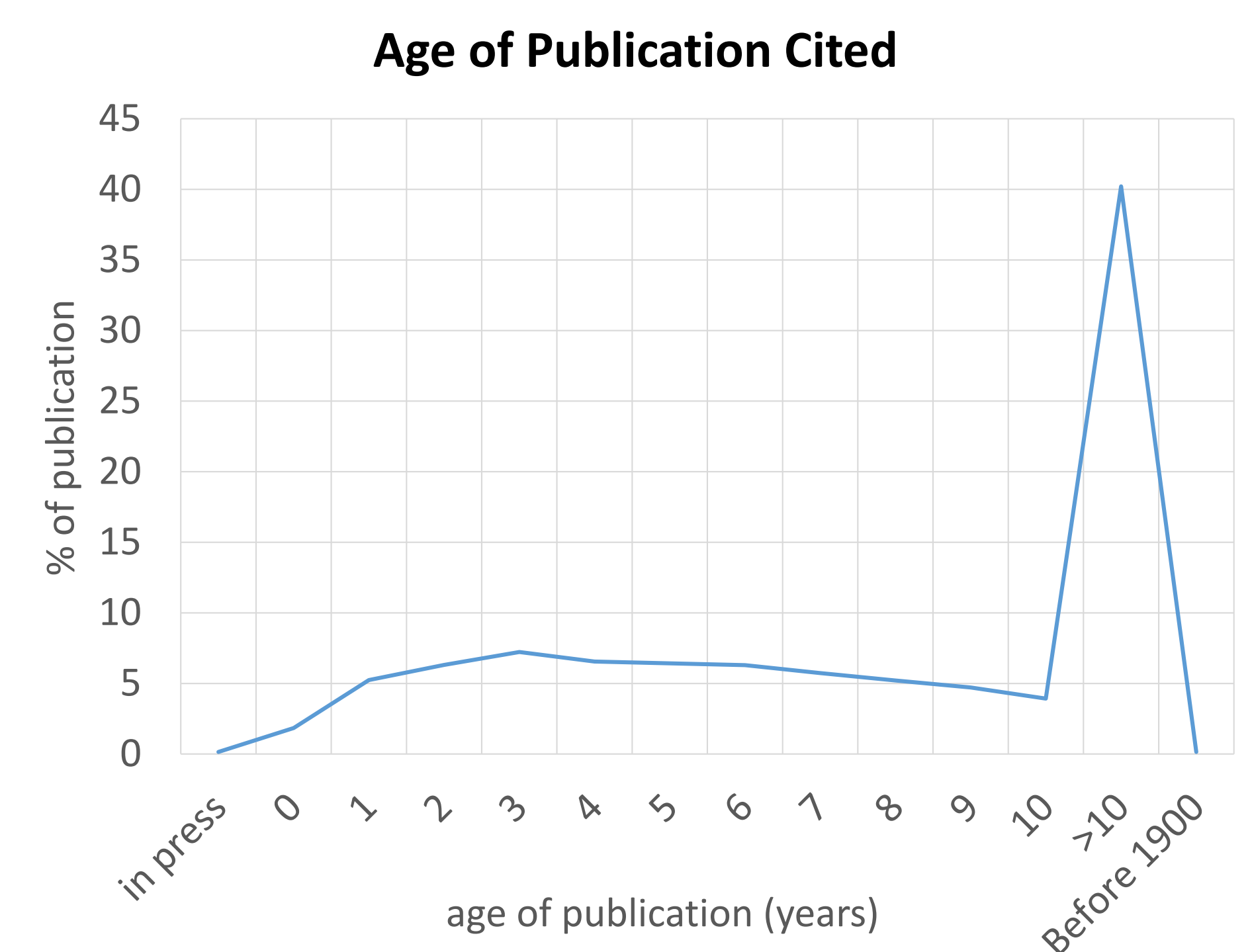
## Objectives

- Assess library collection coverage of resources used in doctoral dissertations
- How many of those resources are Open Access?
- Complete and compare with other resource usage metrics



## Findings

- Library holdings cover:
  - 84% of journal articles cited
  - 49% of the matched books, book chapters and conference proceedings (51% remained without conclusive coverage outcome)
- OA coverage of resources with DOI:
  - 49,8% of journal articles
  - 17,6% of books, book chapters and conference proceedings
- 20% of journal titles account for 73.6% of article citations (80/20% rule)
- Difficult to obtain accurate coverage information from the API of the discovery platform
- CORE and Unpaywall invaluable for OA versions and complement each other nicely (32% overlap for journal articles and 6% for books)
- Nearly 70% of cited resources published > 5 years ago



## Future work

- Automation of Metadata extraction from Crossref
- Refine workflow to increase accuracy and minimize false matches
- Expand OA coverage analysis to resources without DOI using other metadata
- Extend study to all doctoral dissertations and other scholarly output