

data origin



metadata
(plain text)

paper source
(LaTeX/DOCX)

data sets

unarXive
S2ORC_{LaTeX}
arXMLiv



intermediate
(JATS XML)

PMC-OAS



distribution
(PDF)

CORE
S2ORC



distribution
(various)

MAG
OpenAlex
crossref



title
authors
abstract
etc.

references

