



HAL
open science

Automatic construction of an accessible linked dataset from scientific literature for superconducting materials discovery

Luca Foppiano, Sae Dieb, Akira Suzuki, Pedro Baptista de Castro, Suguru Iwasaki, Yan Meng, Kensei Terashima, Yoshihiko Takano, Masashi Ishii

► To cite this version:

Luca Foppiano, Sae Dieb, Akira Suzuki, Pedro Baptista de Castro, Suguru Iwasaki, et al.. Automatic construction of an accessible linked dataset from scientific literature for superconducting materials discovery. FAIR-DI: Conference on a FAIR Data Infrastructure for Materials Genomics, Jun 2020, Virtual meeting, Germany. hal-02870900

HAL Id: hal-02870900

<https://hal.science/hal-02870900>

Submitted on 17 Jun 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Automatic construction of an accessible linked dataset from scientific literature for superconducting materials discovery

Luca Foppiano¹, Sae Dieb¹, Akira Suzuki¹, Pedro Baptista de Castro², Suguru Iwasaki², Yan Meng², Kensei Terashima², Yoshihiko Takano², Masashi Ishii¹

¹ Materials Database Group, MaDIS, NIMS, Japan

² Nano Frontier Superconducting Materials Group, MANA, NIMS, Japan

Introduction and Motivation

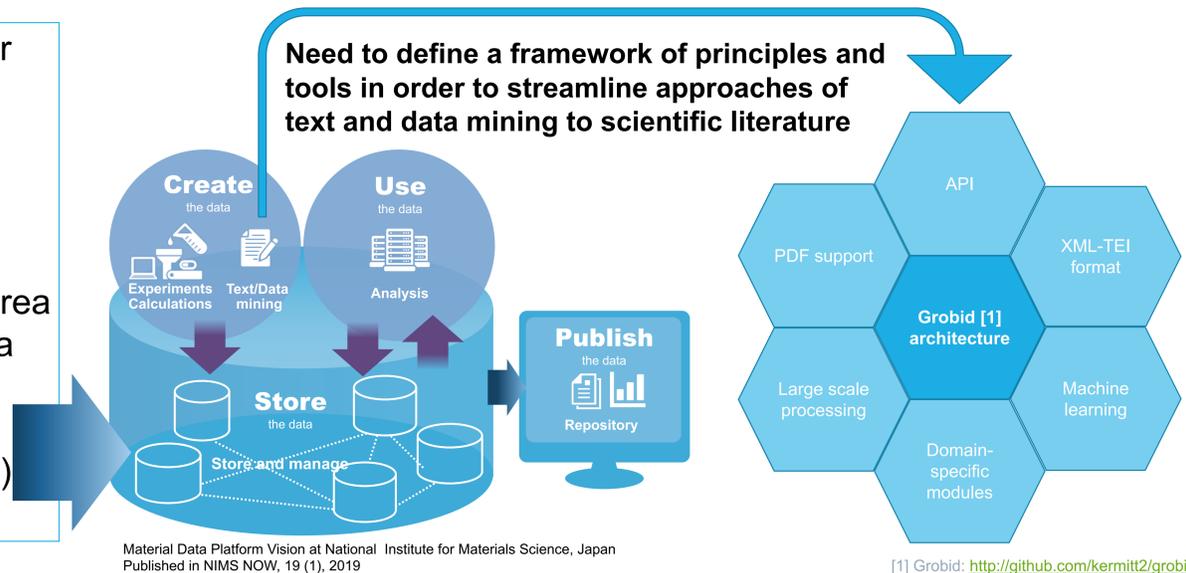
Structured data is key ingredient in machine learning for materials research.

Scientific publications provide large scale unstructured data, however:

- Text-mining is challenging
- Solutions can be reused within the same scientific area
- Machine learning requires availability of training data

Needs of FAIR infrastructure for data and tools:

- Findable, Accessible, Interoperable Reusable (FAIR) data principles



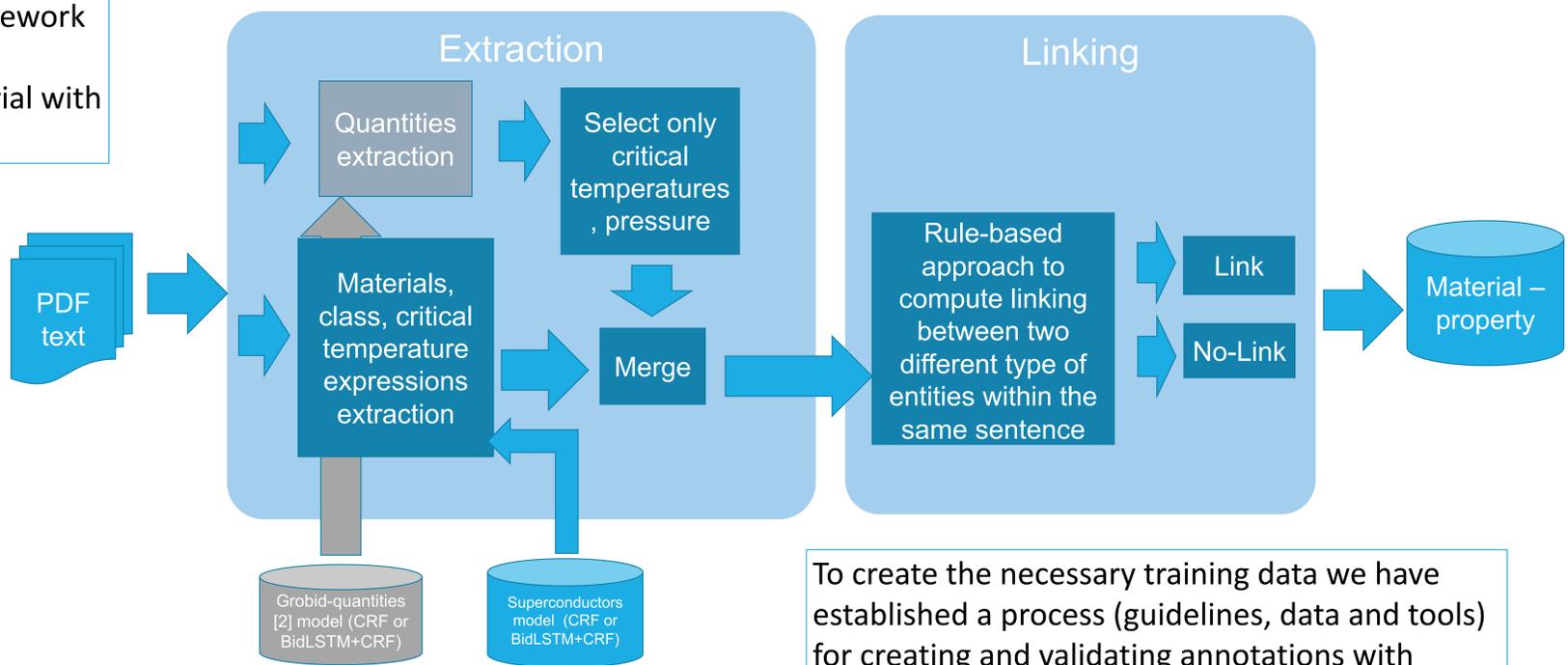
Material Data Platform Vision at National Institute for Materials Science, Japan
Published in NIMS NOW, 19 (1), 2019

[1] Grobid: <http://github.com/kermitt2/grobid>

Superconductors case study

Implementation of our framework within the superconductors domain for extracting material with their respective properties.

Our tool combines ML-based entity "Extraction" through sequence labelling and heuristic "Linking" of the identified materials with their respective physical properties.

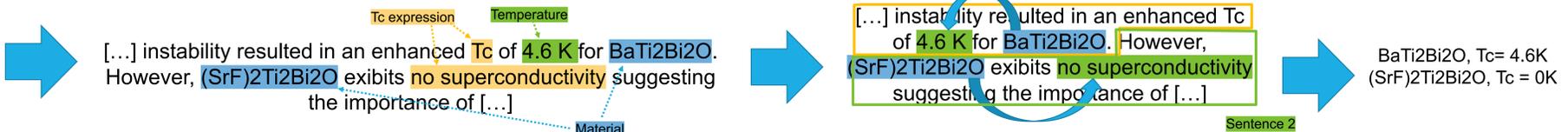


[2] Grobid-quantities is a module of Grobid focusing on Measurement extraction: <http://github.com/kermitt2/grobid-quantities>

[3] excerpt from article "Synthesis and Physical Properties of the New Oxybismuthides BaTi₂Bi₂O and (SrF)₂Ti₂Bi₂O with a d1 Square Net, Takeshi et al., 2012, JPS"

Example

[...] instability resulted in an enhanced T_c of 4.6 K for BaTi₂Bi₂O. However, (SrF)₂Ti₂Bi₂O exhibits no superconductivity suggesting the importance of [...]



Evaluation

Extraction

Corpus of 114 PDF papers evaluated using 10-Fold cross validation

Label	Precision	Recall	F1-Score
Class	81.66	72.36	76.64
Material	81.89	80.09	80.96
Measurement method	73.57	71.24	72.26
Critical pressure	55.27	34.4	41.54
Critical temperature expression	82.3	77.61	79.83
Critical temperature value	74.64	61.8	67.56
All (micro avg)	80.41	75.89	78.07

End to end Evaluation (Extraction + Linking)

Corpus of 500 PDF papers from American Institute of Physics (AIP), American Physical Society (APS) and Institute of Physics (IOP). Results calculated using manual evaluation

Total	Correct	Wrong	Precision	Recall	F1-score
597	441	156	73.86	66.33	69.90

Future work

- Linking improvement with probabilistic approach
- Increase the training data and extract more information (e.g. magnetic field)

Contact information: FOPPIANO.Luca@nims.go.jp