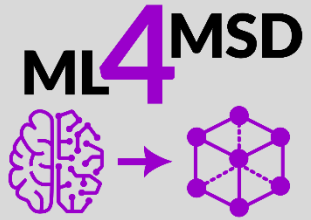


ME 5374-ST



# Machine Learning for Materials Science and Discovery

Fall 2025

Asst. Prof. Peter Schindler

## Lecture 9 – Data-Types and Databases in Materials Science

- Data-Types in Materials Science
- Atomistic File Types
- Application Programming Interfaces (APIs)
- Materials Databases
- FAIR Data Principles

# Types of Materials Data



Text

Scientific Literature

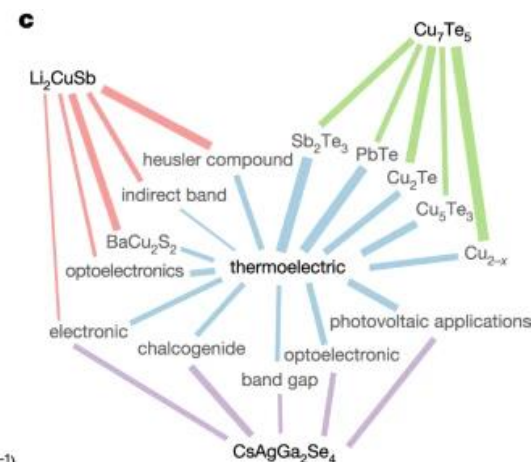
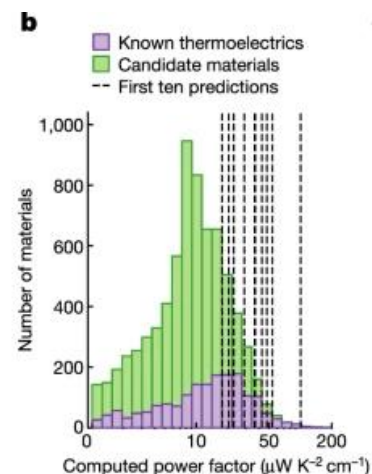
NLP and LLMs

**a**

Cosine similarity to 'thermoelectric'

1.	$\text{Bi}_2\text{Te}_3$	✓
2.	$\text{MgAgSb}$	✓
3.	$\text{PbTe}$	✓
...		✓
326.	$\text{Li}_2\text{CuSb}$	?
...		✓
328.	$\text{In}_4\text{Te}_3$	✓
...		✓
345.	$\text{Cu}_3\text{Nb}_2\text{O}_8$	?
...		✓

✓ Known thermoelectrics  
? Predictions



Tshitoyan, V., et al. *Nature* **571**, 95–98 (2019)

# Types of Materials Data

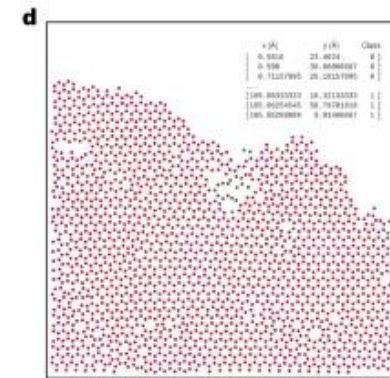
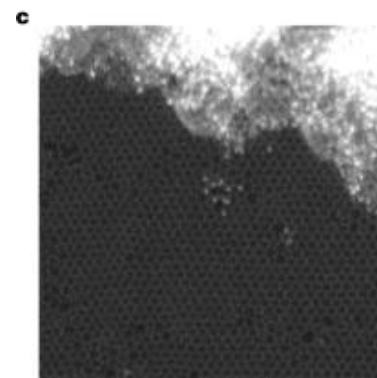
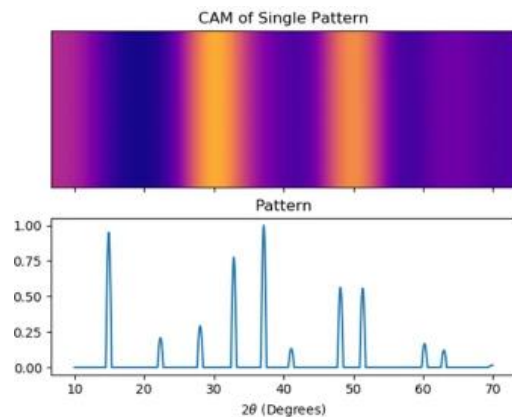
Experimental

Spectra,  
Images

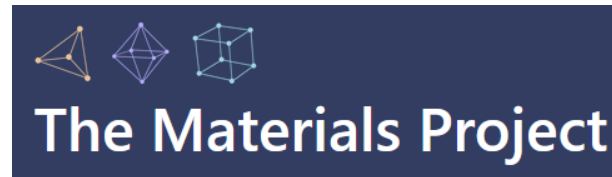
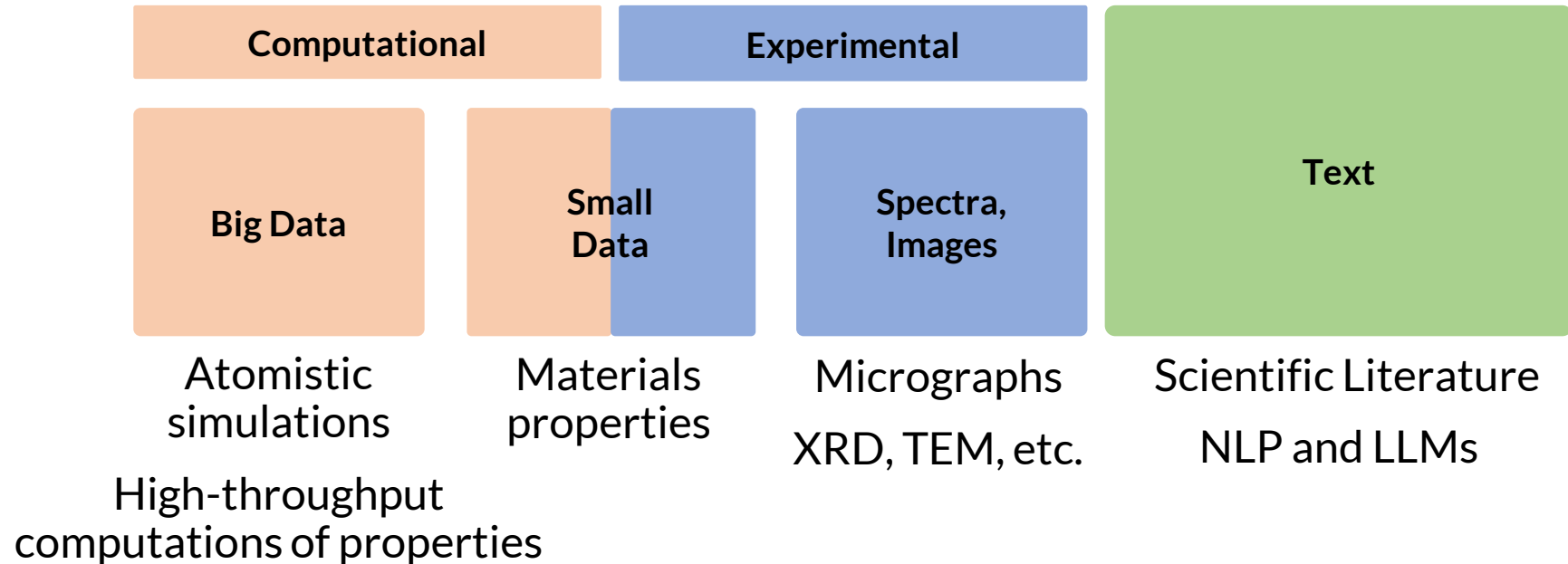
Text

Micrographs  
XRD, TEM, etc.

Scientific Literature  
NLP and LLMs



# Types of Materials Data



# ML Paradigms in Materials Science

## Model-centric AI

How change the model/architecture to improve performance?



## Data-centric AI

How systematically change data (x/y) to improve performance?

Big Data



"Good Data"

## Shallow ML + feature engineering

Computational

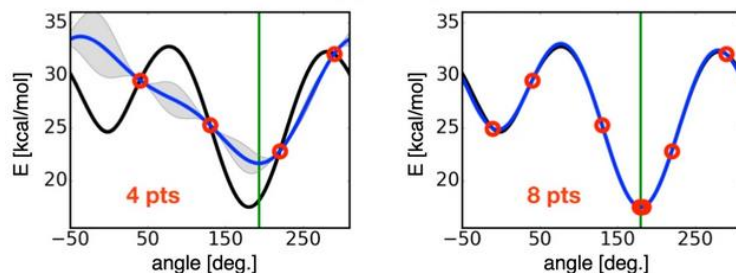
Experimental

Big Data

Small Data

Spectra, Images

## Active Learning



Computational

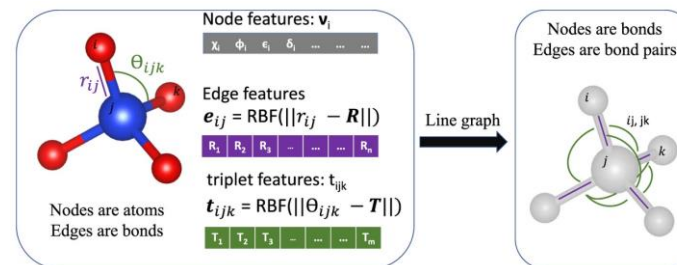
Experimental

Big Data

Small Data

Spectra, Images

## Deep Learning



Computational

Experimental

Big Data

Small Data

Spectra, Images

## Transfer Learning

Computational

Experimental

Big Data

Small Data

Spectra, Images

# Crystal Information File (CIF)

```
# generated using pymatgen
data_Mo2NO
_symmetry_space_group_name_H-M 'P 1'
_cell_length_a 3.10000000
_cell_length_b 6.20000000
_cell_length_c 3.10000000
_cell_angle_alpha 90.00000000
_cell_angle_beta 90.00000000
_cell_angle_gamma 90.00000000
_symmetry_Int_Tables_number 1
_chemical_formula_structural Mo2NO
_chemical_formula_sum 'Mo2 N1 O1'
_cell_volume 59.58200000
_cell_formula_units_Z 1
loop_
_symmetry_equiv_pos_site_id
_symmetry_equiv_pos_as_xyz
1 'x, y, z'
loop_
_atom_site_type_symbol
_atom_site_label
_atom_site_symmetry_multiplicity
_atom_site_fract_x
_atom_site_fract_y
_atom_site_fract_z
_atom_site_occupancy
Mo Mo0 1 0.00000000 0.00000000 0.00000000 1
N Mo 1 0.00000000 0.50000000 0.00000000 1
Mo Mo2 1 0.50000000 0.70000000 0.70000000 1
O O3 1 0.00000000 0.25000000 0.25000000 1
```

← No explicit symmetry

Lattice

Symmetry/Equivalent Positions

Atomic sites (species,  
occupation, coordinates,...)

```
_space_group_IT_number 29
_symmetry_space_group_name_Hall 'P 2c -2ac'
_symmetry_space_group_name_H-M 'P c a 21'
_cell_angle_alpha 90
_cell_angle_beta 90
_cell_angle_gamma 90
_cell_length_a 5.5833
_cell_length_b 5.5892
_cell_length_c 5.5812
_cell_formula_units_Z 4
_cell_volume 174.168
_database_code_amcsd 0005243
_exptl_crystal_density_diffn 6.328
_cod_original_formula_sum 'Co As S'
_cod_database_code 9004218
loop_
_space_group_symop_operation_xyz
x,y,z
1/2+x,-y,z
1/2-x,y,1/2+z
-x,-y,1/2+z
-
```

More details on syntax:

<https://www.iucr.org/resources/cif/spec/version1.1/cifsyntax>

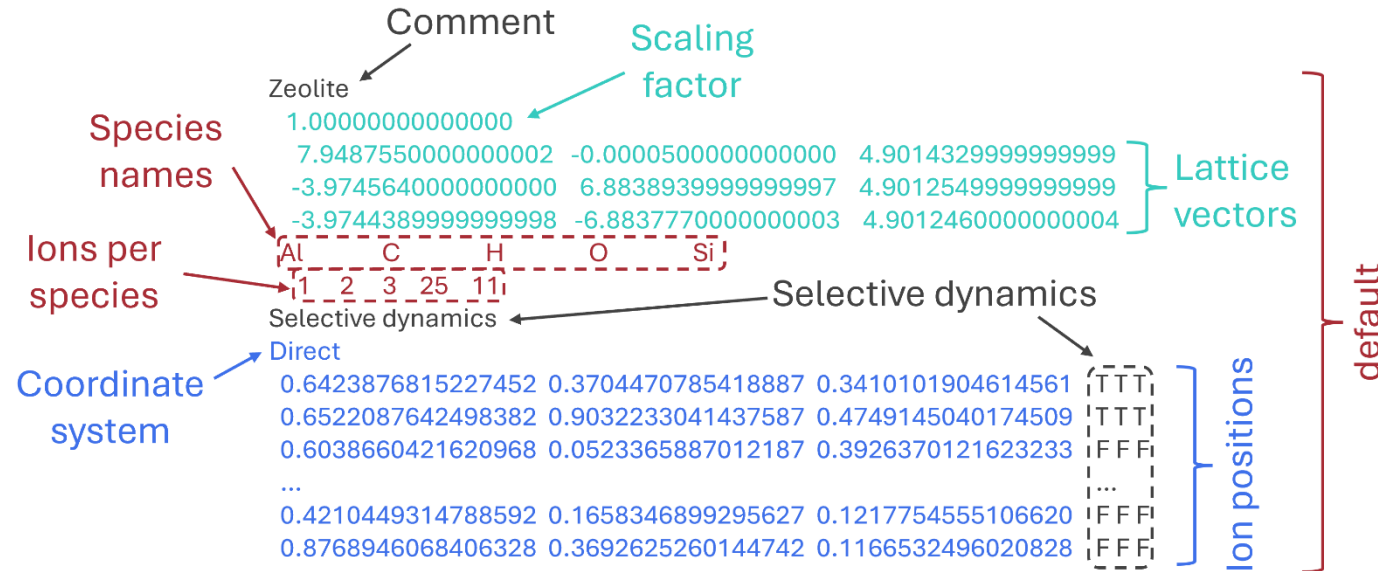
# XYZ File

```
<number of atoms>  
comment line  
<element> <X> <Y> <Z>  
...
```

- Stores Cartesian coordinates (in Angstrom) and atom types of a molecule
- Extended XYZ format allows definition of lattice (and other information) in the comment line

```
8  
Lattice="5.44 0.0 0.0 0.0 5.44 0.0 0.0 0.0 5.44" Properties=species:S:1:pos:R:3 Time=0.0  
Si      0.00000000    0.00000000    0.00000000  
Si      1.36000000    1.36000000    1.36000000  
Si      2.72000000    2.72000000    0.00000000  
Si      4.08000000    4.08000000    1.36000000  
Si      2.72000000    0.00000000    2.72000000  
Si      4.08000000    1.36000000    4.08000000  
Si      0.00000000    2.72000000    2.72000000  
Si      1.36000000    4.08000000    4.08000000
```

# POSCAR File (VASP)



- Crystal structure input file for VASP (Density Functional Theory)
- Coordinate systems: “Direct” (=Fractional) or “Cartesian”
- More details: <https://www.vasp.at/wiki/index.php/POSCAR>



# Application Programming Interfaces (APIs)

API is a set of defined functions, procedures, methods, or classes which enable a **structured way of exchanging data** between programs. *It facilitates...*

- uploading, examining, and downloading of data (without manual steps)
- selectively querying (fetch data that adhere to constraints)
- keeping track of versions
- standardization of interface

*Representational State Transfer* (REST API or RESTful interface) is the most commonly used framework that interfaces via stateless endpoints through request URLs.

(see here for details: <https://aws.amazon.com/what-is/restful-api/> )

Request URL

```
https://api.materialsproject.org/materials/summary/?formula=SiO2&deprecated=false&_per_page=100&_skip=0&_limit=100&_all_fields=false&license=BY-C
```

<https://api.materialsproject.org/docs>

**OPTIMADE**



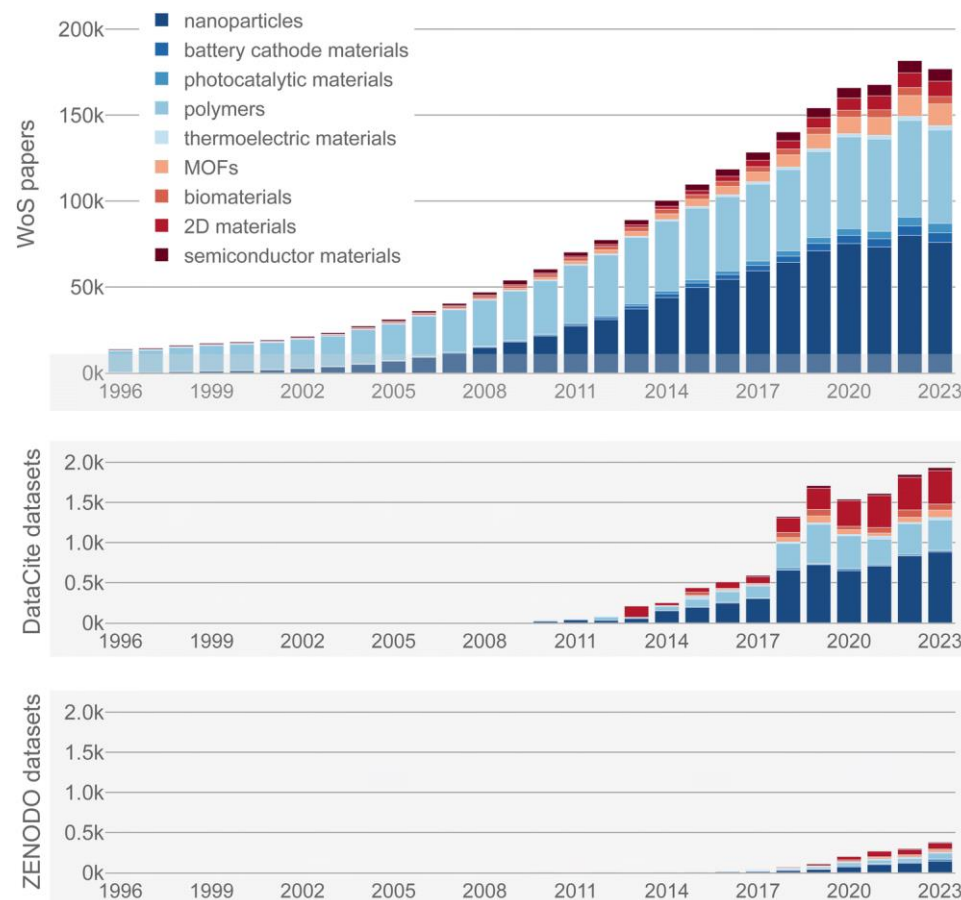
OPTIMADE  
Open Databases Integration  
for Materials Design

API standard developed for materials science and adopted by many large-scale databases

# Data Sources

Source type	Access methods	Examples	Access type	Drawbacks
<b>Papers &amp; textbooks</b>	Hand extract (copy/paste) Download PDF/LaTeX/XML Text-mine with NLP/LLMs	Journal articles, Landolt-Börnstein ( <a href="https://materials.springer.com/bookshelf">https://materials.springer.com/bookshelf</a> )	Manual (unless text-mining with NLP/LLMs)	Time-consuming; error-prone; parsing PDFs difficult; coverage limited
<b>Public websites (no API)</b>	Manual download, copy/paste Web crawling (legal grayzone)	University pages, project web pages, vendor datasheets	Manual	Data formats inconsistent; scraping may be legally/ethically restricted; fragile if site changes
<b>Public repositories (no API)</b>	Manual download of files or bulk repositories	Figshare, Zenodo, NIMS ( <a href="https://mits.nims.go.jp/">https://mits.nims.go.jp/</a> )	Manual	No programmatic access; updates require re-download; can be very large
<b>Public databases (with API)</b>	Programmatic access (REST/GraphQL APIs)	Materials Project, PubChem, OQMD	API	Rate limits (sometimes) Learning curve for API queries
<b>Commercial databases</b>	API + web portals; licensed access	Proprietary vendor/industry databases ICSD, Pauling File	API (licensed)	Expensive; restrictive licenses; not open for redistribution
<b>Automated experiments (own data)</b>	Instrument/robot control, high-throughput workflows	Autonomous labs; Laboratory Information Management Systems (LIMS)-managed workflows	Manual (experiment) + API (instrument control)	Very costly; needs infrastructure; time-intensive

# Potential Value of Text-Mining



# Recap: Database File Formats

Format	Structure	Human-readable	Loading speed	Strengths	Limitations
<b>CSV</b> (Comma-Separated Values)	Rows and columns separated by delimiters (flat, tabular)	✓ Yes	⚡ Medium	Simple, lightweight, works well with spreadsheets	No nested structures, no data types (all text)
<b>JSON</b> (JavaScript Object Notation)	Key-value pairs, lists, nested objects (hierarchical)	✓ Yes	🐢 Slow	Handles complex data, preserves data types, widely used in APIs (Application Programming Interfaces)	More verbose, harder to edit manually than CSV
<b>HDF5</b> (Hierarchical Data Format)	Binary format with hierarchical groups and datasets	✗ No	⚡ Fast	Efficient for very large datasets, supports metadata	Requires special libraries (e.g., h5py, PyTables)
<b>Pickle</b>	Python-specific serialized objects	✗ No	⚡ Fast	Can store almost any Python object easily	Not portable outside Python, unsafe if source is untrusted
<b>YAML</b> (YAML Ain't Markup Language)	Human-readable key-value and nested structure	✓ Yes	🐢 Slow	More flexible and readable than JSON, allows comments	Less standardized than JSON

# Examples: Property-Focused Materials Databases (mostly computational)

name	structure information	mechanical properties	thermal properties	electronic properties	API <sup>a</sup>	data license	refs
Materials Project	Y	Y	Y	Y	Y	CC BY 4.0	85
Open Quantum Materials Database	Y	N	Y	Y	Y	CC BY 4.0	86
AFLOW for Materials Discovery	Y	Y	Y	Y	Y	<i>b</i>	87
Novel Materials Discovery (NOMAD)	Y	Y	Y	Y	Y	CC BY 4.0	88
Open Materials Database	Y	N	Y	Y	Y	CC BY 4.0	89
Citrine Informatics	Y	Y	Y	Y	Y	CC BY	90
Materials Platform for Data Science (MPDS)	Y	Y	Y	Y	Y	CC BY 4.0	91
AiiDA/Materials Cloud	Y	Y	Y	Y	Y	Varies	92, 93
NREL MatDB	Y	N	Y	Y	N	Own license	94
NIST TRC Alloy Data	N	N	Y	N	On request	Free	95
NIST TRC ThermoData	N	N	Y	N	N	NIST SRD	96
NIST JARVIS-DFT/-ML Database	Y	Y	Y	Y	Y	Public domain	97, 98
MatWeb	N	Y	Y	N	N	Paid	99
Total Materia	N	Y	Y	N	N	Paid	100
Ansys Granta (MaterialUniverse repository)	N	Y	Y	N	N	Paid	101
MATDAT	N	Y	Y	N	N	Paid	102

*Machine Learning for Materials Scientists: An Introductory Guide toward Best Practices,*

A. Y.-T. Wang, R. J. Murdock, S. K. Kauwe, A. O. Oliynyk, A. Gurlo, J. Brgoch, K. A. Persson, and T. D. Sparks, *Chemistry of Materials* 2020 32 (12), 4954-4965

# Examples: Structure-Focused Materials Databases (mostly experimental)

name	no. records <sup>a</sup>	API	Data license	ref
Cambridge Structural Database (CSD)	1,055,780	Y	Paid	103
Inorganic Crystal Structure Database (ICSD)	216,302	N	Paid	104
Pearson's Crystal Data (PCD)	335,000	N	Paid	105
International Centre for Diffraction Data (ICDD)	1,004,568	N	Paid	106
Crystallography Open Database (COD)	455,714	Y	Open-access	107
Pauling File	357,612	Y	Paid	108
CrystMet database	160,000	N	Paid	109

<sup>a</sup>Note: values for number of records were updated as of the submission date (May 2020).

*Machine Learning for Materials Scientists: An Introductory Guide toward Best Practices,*

A. Y.-T. Wang, R. J. Murdock, S. K. Kauwe, A. O. Oliynyk, A. Gurlo, J. Brgoch, K. A. Persson, and T. D. Sparks, *Chemistry of Materials* 2020 32 (12), 4954-4965

# Compilation of Databases and Tools in Materials Informatics

<https://github.com/blaiszik/awesome-matchem-datasets>

<https://github.com/tilde-lab/awesome-materials-informatics>

The screenshot shows the GitHub repository 'awesome-matchem-datasets' by user 'blaiszik'. It is a public repository with 6 branches and 0 tags. The commit history shows a recent update to README.md by 'blaiszik' 4 days ago (commit 4164bbc) and 91 total commits. The file list includes .Jlms, LICENSE, README.md, and matchem-datasets.png. The README is selected, showing the title 'Awesome Materials & Chemistry Datasets' and a description: 'A curated list of the most useful datasets in materials science and chemistry for training machine learning and AI foundation models. This includes experimental, computational, and literature-mined datasets—prioritizing open-access resources and community contributions.'

The screenshot shows the GitHub repository 'awesome-materials-informatics' by user 'blokhin'. It is a public repository with 2 branches and 1 tag. The commit history shows a recent update to README.md by 'blokhin' 4 months ago (commit 8c00f57) and 150 total commits. The file list includes .github/workflows, CITATION.cff, CONTRIBUTING.md, README.md, and \_config.yml. The README is selected, showing the title 'Awesome Materials Informatics' with an 'awesome' badge and a DOI '10.5281/zenodo.7693349'. The description states: 'The novel discipline of materials informatics is a junction of materials, computer, and data sciences. It aims to unite the nowadays competing physics- and data-intensive efforts for the most impactful applied science, that transformed our society in the 20th century.'

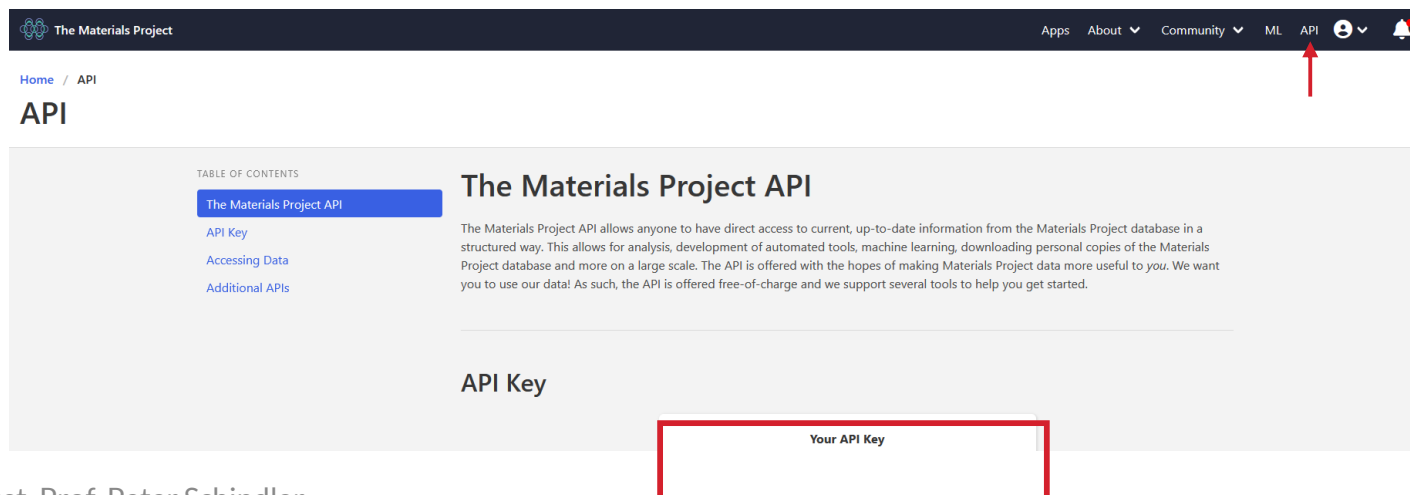
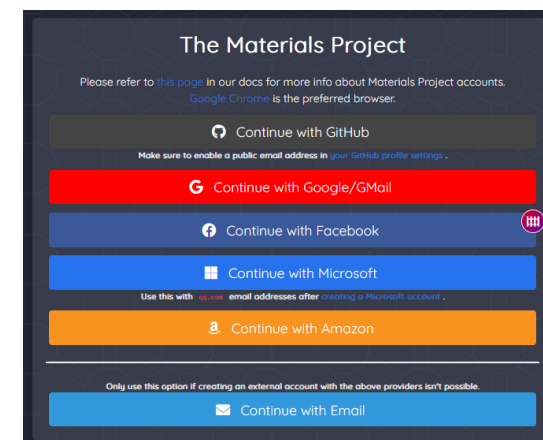
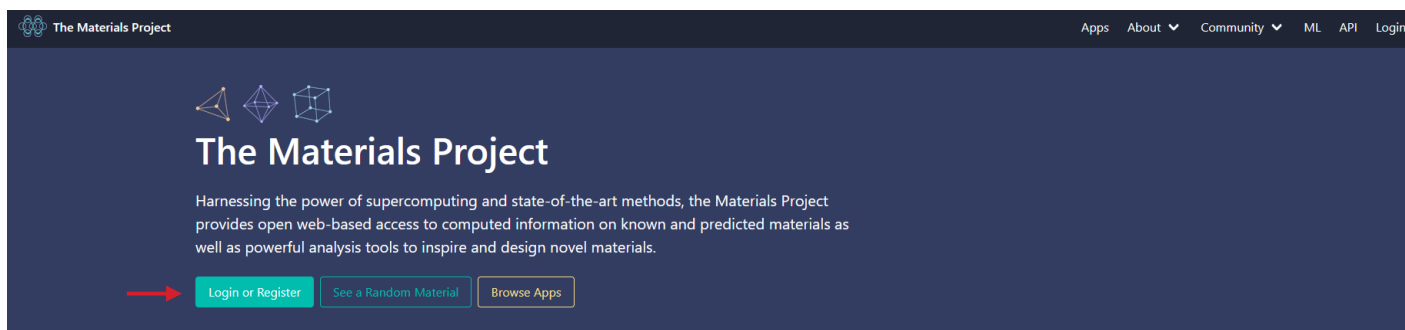
# FAIR Data Principles

Letter / Term	Purpose	Example Solution
<b>F – Findable</b>	Ensure data can be located easily by humans/machines	Assign persistent identifiers (DOIs) Provide searchable metadata in repositories
<b>A – Accessible</b>	Make data retrievable via standard protocols	Open web protocols (HTTP/HTTPS) or APIs, Ensure data can be accessed with authentication if needed (e.g., PubChem API)
<b>I – Interoperable</b>	Enable data integration across platforms and tools	Store crystallographic data in CIF format, adopt common ontologies (e.g., OPTIMADE for materials APIs)
<b>R – Reusable</b>	Facilitate long-term use and reproducibility	Provide detailed metadata (methods, units), apply open licenses (e.g., CC-BY), deposit in repositories like Zenodo or Figshare

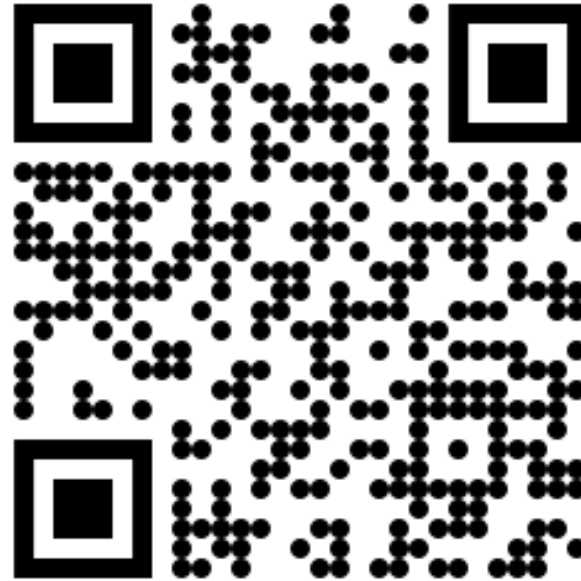


# Materials Project: Obtaining API Key

Go to: <https://materialsproject.org>



# Lecture Feedback



Please, scan the QR code and take a minute to let me know how the lecture was and mention any **feedback/questions**

This form is **anonymous!**