Osaka University/NANOPROGRAM June 24, 2024

Introduction to Materials Informatics

Tamio Oguchi Center for Spintronics Research Network Osaka University

アウトライン



まとめ

Paradigm of Science



Kepler's laws of planetary motion



Tycho Brahe 1546-1601

> **Johannes Kepler** 1571-1630

Galileo Galilei 1564-1642 Telescope [1609]



https://en.wikipedia.org/wiki/

Isaac Newton 1642-1727

Newton's law of universal gravitation and law of motion

Observations \rightarrow Laws \rightarrow **Principles** \rightarrow **Applications**

The 5th Solvay Conference 1927.10 "Electrons and photons"



https://en.wikipedia.org/wiki/Solvay_Conference

"Most of Physics and all of Chemistry are solved." Paul Dirac 1902-1984



"You can not understand it, until you know how to calculate it." J. C. Slater 1900-1976



Growing Power of Supercomputer

2021

Fugaku

Peak Performance (FLOPS)



Cray-1

- First Supercomputer in the World
 - Installed in Los Alamos National Lab in 1976
 - Peak Performance:100MFLOPS
 - Main Memory: 8MB



Seymour Cray



Los Alamos National Laboratory



説

第三の物理学としての計算物理学

A. J. Freeman

訳 大 西 楢 平*

本稿は,著者が山田財団の援助により2カ月間日本に滞在した機会に行った 第37回日本物理学会年会(1982年10月)の特別講演から,計算物理学に関する一 般的考察の部分を新たに展開,加筆したものである.米国における計算物理学 の進展の現状と物理学の諸分野とのかかわりあいを多くの例をあげて紹介して いる.複雑な系をとり扱う物理学としての計算物理学,その理論科学的,実験 科学的,及び応用科学的側面,また計算物理学者の特性などが議論されている.

日本物理学会誌 1983

12

Paradigm of Science





Number of publications per year (1975–2014) on topics ("density functional" or "DFT"), according to the Web of Science Core Collection (February 2015)

Jones (2015) 14

The Nobel Prize in Chemistry 1998



Walter Kohn



John A. Pople

The Nobel Prize in Chemistry 1998 was divided equally between Walter Kohn "for his development of the densityfunctional theory" and John A. Pople "for his development of computational methods in quantum chemistry". **Current Issues in Materials Science and Engineering**

- Complexity in Materials
 - Multi-scale, multi-physics, hierarchal structures in time and space, leading to the importance of top-down and bottom-up connections in many contexts
 - Strong requirements not only to get high performance for cost and efficiency with sustainability for environment and energy problems, and but also to accelerate R&D
- We are confronted with limitations in the accuracy and efficiency of traditional research strategies such as experimental, theoretical, and computational approaches.
- Importance of "Data-Intensive Scientific Discovery" as the fourth paradigm of science has been rapidly rising quite recently in materials science & engineering, as called

★ MATERIALS INFORMATICS

SCIENCE

The New York Times

A Deluge of Data Shapes a New Era in Computing

DEC. 14, 2009 In a speech given just a few weeks before he was lost at sea off the California coast in January 2007, Jim Gray, a database software pioneer and a <u>Microsoft</u> researcher, sketched out an argument that computing was fundamentally transforming the practice of science.

Dr. Gray called the shift a <u>"fourth paradigm." The first three</u> paradigms were experimental, theoretical and, more recently, computational science. He explained this paradigm as an evolving era in which an "exaflood" of observational data was threatening to overwhelm scientists. The only way to cope with it, he argued, was a new generation of scientific computing tools to manage, visualize and analyze the data flood.

.....

Now, as a testimony to his passion and vision, colleagues at Microsoft Research, the company's laboratory that is focused on science and computer science, have published a tribute to Dr. Gray's perspective in <u>"The Fourth Paradigm: Data-Intensive Scientific</u> <u>Discovery.</u>" It is a collection of essays written by Microsoft's scientists and outside scientists, some of whose research is being financed by the software publisher.

Paradigm of Science



"The Fourth Paradigm: Data-Intensive Scientific Discovery" in Materials Science

MATERIALS INFORMATICS



Materials Genome Initiative

https://www.mgi.gov



Materials Informatics



Materials Research



Materials Research



Big Data



Sources: McKinsey Global Institute, Twitter, Cisco, Gartner, EMC, SAS, IBM, MEPTEC, QAS

http://www.ibmbigdatahub.com/infographic/four-vs-big-data

Materials Database

- MatNavi
 http://mits.nims.go.jp
 - AtomWork: inorganic material DB *Free*



Structure: 82,000, Property: 55,000, Phase diagram: 15,000



★ AtomWork-Adv Fee-Based

Structure: 273,830, Property: 298,021, Phase diagram: 40,301

Limited number of data for particular purpose !

Computational Materials Design with Machine Learning



Computational Materials Database



https://materialsproject.org

OQMD: The Open Quantum Materials Database

http://oqmd.org



https://nomad-coe.eu



http://aflowlib.org

Role of First-Principles Calculation

- Most of material properties are governed by the electronic states described by the quantum theory.
- For concrete materials, many of properties can be predicted from first principles, offering another route to solve the direct problem.



- Electron theory may provide new axes x (descriptors) that P spans.
- First-principles calculations can give some valuable pieces of information at less cost, even for non-existing and unstable materials.

Direct and Virtual Screening





Machine Learning



Machine Learning and Deep Learning

Machine Learning $\underbrace{\mathsf{Machine Learning}}_{\mathsf{Input}} \xrightarrow{\mathsf{Car}}_{\mathsf{Not Car}} \xrightarrow{\mathsf{Car}}_{\mathsf{Output}}$ $\underbrace{\mathsf{Deep Learning}}_{\mathsf{Input}} \xrightarrow{\mathsf{Car}}_{\mathsf{Output}} \xrightarrow{\mathsf{Car}}_{\mathsf{Output}}$



Descriptor for Structural Stability in Binary Compounds

• Empirical 2D descriptors by van Vechten and Phillips



• Atomic number: not a good descriptor



L.M. Ghiringhelli et al. PRL (2015).

Conditions of Good Descriptor

- A descriptor *di* uniquely characterizes the material *i* as well as property-relevant elementary processes.
- Materials that are very different (similar) should be characterized by very different (similar) descriptor values.
- The determination of the descriptor must not involve calculations as intensive as those needed for the evaluation of the property to be predicted.
- The dimension of the descriptor should be as low as possible (for a certain accuracy request).

L.M. Ghiringhelli et al. PRL (2015).

• 説明変数(explanatory variable)、特徴量(feature)

Sparse Modeling by LASSO and Exhaustive Search



EA(B)

highest occupied KS level of B lowest unoccupied KS level of B

$$\frac{r_s(\mathbf{A}) - r_p(\mathbf{B})|}{\exp(r_s(\mathbf{A}))}$$

 $r_p(A)^2$

$$\frac{|r_p(\mathbf{B}) - r_s(\mathbf{B})|}{\exp(\underline{r_d(\mathbf{A})})}$$

density maximum radius of A-d orbital

$$\Delta E_{\rm AB} = E(\rm AB_{\rm RS}) - E(\rm AB_{\rm ZB})$$



L.M. Ghiringhelli et al. PRL (2015).

First-Principles Calculations



Computational Materials Design with Machine Learning



Crystal Structure Search for Materials Discovery and Design

Acc. Chem. Res. 1994, 27, 309-314

Are Crystal Structures Predictable?

ANGELO GAVEZZOTTI^{*}

Dipartimento di Chimica Fisica ed Elettrochimica, Università di Milano, Milano, Italy

Received May 16, 1994

"No": by just writing down this concise statement, in what would be the first one-word paper in the chemical literature, one could safely summarize the present state of affairs, earn an honorarium from the American Chemical Society, and do a reasonably good service to his or her own reputation. In the mainstream of academic tradition, one could then concede a "maybe", or even a conditional "yes", thus making a good point for discussion; and then, in the mainstream of publication policy tradition, proceed eventually to have his or her papers rejected by referees taking the opposite stand.
Crystal structure prediction





Properties

- First-principles calculations provide a powerful tool to predict properties for realistic materials with the information of crystal structure.
- Therefore, crystal structure prediction is indispensable for materials discovery without prior knowledge on structure.



Issues in crystal structure prediction

★ Huge Structure Space

- **1.** How to search the structure space globally
 - Local optimization by first-principles calculations or classical MD simulations
- 2. How to represent structure space (descriptor)
 - Easy and efficient representation of structure vector



1. Global structure search

Difficult due to many local minima

 Many computationally demanding structure optimization processes are required.



Development of an efficient global search algorithm is highly desired.

Structure search algorithm

• Random search algorithm

 $\checkmark\,$ Random selection of lattice parameters and atomic positions

• Evolution-type algorithm

- ✓ Evolutionary algorithm (USPEX)
- ✓ **Particle swarm optimization (CALYPSO)**

Learning-type algorithm

 \checkmark Bayesian optimization

→ A sequential design strategy for global <u>optimization</u> of black-box functions that <u>doesn't require derivatives</u>.
[Wikipedia]

★ The target functions may not be analytic nor continuous.

Key features for structure prediction

Exploration

 to search the structure space as globally as possible

Exploitation

 to search the structure space without missing important minima

Bayesian optimization (BO)



BO enables to balance trade-off between exploration and exploitation of the search algorithm.

2. Descriptor for crystal structure

Fingerprint: based on radial distribution function

A. R. Oganov and M. Valle, J. Chem. Phys. 130, 104504 (2009).

$$F_{AB}(R) = \sum_{A_{i},\text{cell}} \sum_{B_{j}} \frac{\delta(R - R_{ij})}{4\pi R_{ij}^{2} \frac{N_{A}N_{B}}{V}\Delta} - 1 \qquad F_{AB}(0) = -1$$
$$F_{AB}(\infty) = 0$$

- > *i*-th atom of type A within the unit cell
- > *j*-th atom of type B within the distance R_{max} (~5 Å)
- \succ discretized over bins of width \triangle (50 discretized R points)





Flowchart of Bayesian optimization



47

Lattice constants



Code development

CrySPY

https://github.com/Tomoki-YAMASHITA/CrySPY distributed under the MIT License

• Algorithm

- / Random search
- ✓ Bayesian optimization
- ✓ Evolutionary algorithm

Space group by find_wy

H. Kino, https://github.com/nims-hrkn/find_wy

Bayesian optimization library, COMBO

T. Ueno *et al.*, Materials Discovery 4, 18 (2016). https://github.com/tsudalab/combo

/ LAQA: Look Ahead based on Quadratic Approximation A fine-grained method to reduce local optimization steps

Local optimization

Interfaced with

- √ VASP
- ✓ Quantum Espresso
- √ soiap
- ✓ LAMMPS

Test simulations

- Both random search and Bayesian optimization are applied to two known systems
 - ▶ Na₈Cl₈ (16 atoms/cell)

Rocksalt structure



Y₂Co₁7 (19 atoms/cell)

Th₂Zn₁₇-type structure (R-3m) Ferromagnetic



Test simulation: Na₈Cl₈



Test simulation: Na₈Cl₈

To investigate the statistical efficiency, 200 BO simulations were carried out.



Frequency distribution of number of trials required to find the rocksalt structure

★ BO reduces trials by 31% compared with RS.

Test simulation: Na₈Cl₈



Learned correlation between energy and structure distance

800 structure data

One of the good candidates in the global-minimum valley is selected within 70 trials, so BO works well in the exploitation phase.

Test simulation: Y₂Co₁₇





[Conventional] Hexagonal

19 atoms/cell

[Primitive]

Rhombohedral

Test simulation: Y₂Co₁₇



★ BO reduces trials by 39% compared with RS.

Test simulation: Y₂Co₁₇



Learned correlation between energy and structure distance

> It's quite difficult to search for the global minimum, but BO works well in the exploration phase.

: 4 good candidates



Summary

A tool for crystal structure search

- ➤ Random search
- > Bayesian optimization (BO)
- Test simulations for Na₈Cl₈ and Y₂Co₁₇

Bayesian optimization is highly efficient and significantly reduce the number of searching trials required to find the global minimum structure.

Phys. Rev. Materials <u>2</u>, 013803 (2018). Phys. Rev. Materials <u>4</u>, 033801 (2020). Sci. Technol. Adv. Mater.: Methods <u>1</u>, 87 (2021). Sci. Technol. Adv. Mater.: Methods <u>2</u>, 67 (2022). Crystal Structure Map for Materials Classification and Modeling **Current Hot Issue in Materials Research**

- Materials Discovery Assisted by Data Science Approach
 - Today, a wealth of materials data are being accumulated, experimentally and computationally.
 - ✓ Where and how is each existing materials system placed for design and discovery ?

"Materials Map"

Мар



Q: How to define such a low-dimensional axis system and the coordinate of materials targets ?

Structure as Axes of the Map

- Most fundamental attribute in condensed matter systems, especially crystalline materials
- Governing the electronic states as under the Born-Oppenheimer approximation.

 $\mathcal{H}_{ ext{electron}} = \mathcal{H}_{ ext{electron}} \left[\{ \boldsymbol{R}_n \} \right]$

★ Our understanding of materials properties is often based on the correlation with structure.

Crystal Structure Map

✓ Coordinates classify huge related types of structures.

✓ Axes provide structure features for modeling properties.

Crystal Structure DB



http://www.crystallography.net/cod/

https://materialsproject.org

What's CIF ?

- Crystallographic Information Framework
- Established by International Union of Crystallography (IUCr) – Dictionary: CoreCIF
- Widely used in structure database such as COD, ICSD, Materials Project, MatNavi, etc.
- Containing all information about crystal structure necessary for electronic structure calculations
- Conversion apps from CIF to VASP, QE, etc.: cif2cell*, C-Tools^{\$}, cifconv[#]

*: T. Björkman, Comput. Phys. Commun., 182, 1183 (2011). \$: https://sourceforge.net/projects/c-tools/ #: F. Izumi and K. Miyazaki, Ceramics 54, 473 (2019).

de facto standard for crystal structure format in experiment and theory/computation

Procedure of Structure Classification for a Set of CIFs

- Step 1: To represent a (often highly dimensional) vector as structure feature
 - Simple, tractable, objective, and translationally and rotationally invariant enough to apply to general structures
- Step 2: To define distance between feature vectors for measuring similarity
 - ✓ Degree of similarity

- d^(1,2)
- $d^{(i,j)} < d_{th}$: equivalent within a given *threshold*
- $d^{(i,j)} \sim d_{th}$: close to equivalent but minor dissimilarity
- $d^{(1,2)} < d^{(1,3)}$ means that $d^{(1,2)}$ is closer in similarity than $d^{(1,3)}$.
- Step 3: To map the position of each structure in a lowdimensional space so that the distances calculated from it well approximate the original ones
 - Classification (clustering) of crystal structures for further analysis

Step 1: Feature Vector for Crystal Structure

✓ Coordination Number (CN)

simple feature of local structure around each atomic site no unique definition of neighbors in general structures

✓ Cluster Expansion

useful for considering different configurations like alloy ordering

✓ Radial Distribution Function (RDF)

complete pair-wise information very simple measurable in EXAFS



http://titan.nusr.nagoya-u.ac.jp/Tabuchi/BL5S1/lib/exe/ fetch.php?media=tabuchi:gairon-20181210-v2.pdf

F-Fingerprint

Oganov and Valle, J. Chem. Phys. 130, 104504 (2009)



Step 2: Distance as a Measure of Similarity

- Distance between two *n*-D vectors X and Y
- 1. Euclidean distance (ED)

$$d(\boldsymbol{X}, \boldsymbol{Y}) = \left[\sum_{i=1}^{n} (X_i - Y_i)^2\right]^{1/2}$$



2. Pearson correlation coefficient (PCC)

$$r(\boldsymbol{X}, \boldsymbol{Y}) = \frac{\sum_{i=1}^{n} (X_i - \bar{X})(Y_i - \bar{Y})}{\left(\sum_{i=1}^{n} (X_i - \bar{X})^2 \sum_{i=1}^{n} (Y_i - \bar{Y})^2\right)^{1/2}}$$
$$-1 \le r \le 1$$

→ Cosine distance (CD)

$$d_{\text{cosine}} = \frac{1}{2}(1-r) \qquad 0 \le d_{\text{cosine}} \le 1$$

When X and Y are normalized $d = 2(1 - r) = 4d_{\text{cosine}}$

https://en.wikipedia.org/wiki/Cosine_similarity

Step 3: Mapping Structure

W. S. Torgerson, Psychometrika 17, 401 (1952)

- Multidimensional Scaling
 - ✓ For a given pair-wise (Euclidean/non-Euclidean) distance matrix D of N points, we seek N-point coordinates X in a low (k ≪ n) dimensional space so that pair-wise Euclidean distance matrix D^x calculated using X is the closest approximation to D.
 - ✓ Often called "Dimension Reduction" for data mapping
 - Equivalent to "Principal Component Analysis (PCA)" for Euclidean distance
 - ✓ Implemented in MATLAB, R, Scikit-learn, ...

Multidimensional Scaling (MDS)

W. S. Torgerson, Psychometrika 17, 401 (1952)

• Algorithm squared distance 1. $B = XX^{T} = -\frac{1}{2}JD^{(2)}J, J = I - \frac{1}{N}11^{T}$ \leftarrow double centering 2. $B = QAQ^{T}$ eigenvalue decomposition I: identity matrix 1: identity vector 3. $X_{+} = Q_{+}A_{+}^{1/2}$ first k positive-eigenvalue part only $X: (N \times d)$ $(N \times k)$ $k \ll d:$ dimension reduction d-D coordinates centered

What dimensionality k one should choose?

The larger eigenvalues, the more important coordinates, because the sum of the eigenvalues in Λ_+ should approximate the sum of all eigenvalues in Λ , so that small negative eigenvalues cancel out small positive eigenvalues.

R. Sibson, J. R. Statist. Soc. B 41, 217 (1979)

$$\mathbf{B} \approx \mathbf{B}_{+} = \mathbf{X}_{+} \mathbf{X}_{+}^{T} = (\mathbf{Q}_{+} \mathbf{\Lambda}_{+}^{1/2}) (\mathbf{Q}_{+} \mathbf{\Lambda}_{+}^{1/2})^{T}$$
69

MDS Example: Kansai-Area Prefecture Center Distance

	Hyogo	Wakayama	Osaka	Nara	Shiga	Kyoto
Hyogo	0	134	85	116	118	60
Wakayama	134	0	68	66	145	141
Osaka	85	68	0	32	83	75
Nara	116	66	32	0	79	95
Shiga	118	145	83	79	0	63
Kyoto	60	141	75	95	63	0



MDS and Principal Component Analysis (PCA)

 New coordinates are span in the axes given as eigenvectors of principal component on original *n*-D space

$$\mathbf{B} = \hat{\mathbf{X}}\hat{\mathbf{X}}^{T} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^{T} \approx \mathbf{Q}_{+}\mathbf{\Lambda}_{+}\mathbf{Q}_{+}^{T} \equiv \mathbf{X}_{+}\mathbf{X}_{+}^{T}$$

$$\hat{\mathbf{X}} = \mathbf{X}_{-}\mathbf{I}\mathbf{X}$$

$$\mathbf{X}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}_{-}\mathbf{U}$$

$$\begin{split} \mathbf{u}_{+} &= \hat{\mathbf{X}}^{T} \mathbf{Q}_{+} \mathbf{\Lambda}_{+}^{-\frac{1}{2}} = \hat{\mathbf{X}}^{T} \mathbf{X}_{+} \mathbf{\Lambda}_{+}^{-1} = \mathbf{X}^{T} \mathbf{X}_{+} \mathbf{\Lambda}_{+}^{-1} &: \textit{d-D eigenvectors of PC} \\ \hline (d \times \textit{k}) & \mathbf{Y}_{+} = \hat{\mathbf{Y}} \mathbf{u}_{+} &: \textit{projection on the PC space} \\ \hat{\mathbf{Z}} &= \mathbf{Z}_{+} \mathbf{u}_{+}^{T} &: \textit{inverse projection} \end{split}$$

C. M. Bishop, Pattern Recognition and Machine Learning (Springer-Verlag, New York, 2006)

XAT

Generation of Structure with CIF



*: H. T. Stokes, D. M. Hatch, and B. J. Campbell, ISOTROPY Software Suite, iso.byu.edu. H. T. Stokes and D. M. Hatch, J. Appl. Cryst. 38, 237 (2005).

Only three inequivalent structures



- "cif2esc" converts crystal structure information from CIF files to input data of several DFT codes for electronic structure calculations and to fingerprint data for structure analysis.
- "diagnosis" classifies crystal structure according to fingerprint distance as its similarity computed by "distance" and "cmds" extracts principal features by using dimension reduction.
- "atls" generates list of possible atom-type combinations from given atom compositions, passing to input for the "findsym*" application with lattice information and atom coordinates to get CIF files.

*: H. T. Stokes, D. M. Hatch, and B. J. Campbell, ISOTROPY Software Suite, iso.byu.edu. H. T. Stokes and D. M. Hatch, J. Appl. Cryst. 38, 237 (2005).

Demonstration: Al₂O₃ polymorph

Crystal Structure of Al₂O₃

https://materialsproject.org


Crystal Structure of Al₂O₃

Materials Project	SG	HOF [meV/atom]	Gap [eV]	Comments
1 mp-1143	R-3c	U (U)	5.599 (6.044)	most stable
2 mp-1245063	P1	+218 (+204)	3.137 (3.326)	amorphous model
3 mp-1938	Pbcn	+92(+94)	5.047 (5.313)	
4 mp-2254	Pna21	+18(+17)	4.592 (4.830)	
5 mp-32591	Cm	+63(+57)	2.854 (3.075)	
6 mp-638765	P-1	+683 (+809)	1.119 (1.136)	
7 mp-642363	Cmcm	+266 (+281)	4.009 (4.249)	C ₄₄ <0, C ₅₅ <0
8 mp-754531	P21/c	+98 (+77)	4.141 (4.374)	
9 mp-755483	R3m	+93 (+91)	3.863 (4.098)	
10 mp-776475	la-3	+44 (+30)	4.982 (5.216)	In ₂ O ₃ stable phase
11 mp-985587	P321	+208 (+210)	4.686 (4.807)	artificial film [excluded]
12 mp-7048	C2/m	+17(+10)	4.239 (4.605)	Ga ₂ O ₃ stable phase
13 mp-754401	Cmc21	+74 (+46)	4.021 (4.235)	
14 mp-754624	R-3	+72 (+66)	5.386 (5.797)	
15 mp-755175	P-31c	+79 (+72)	5.227 (5.494)	
a-pristine	P1	+179	3.725	amorphous model*

data with parentheses: https://materialsproject.org Memide at al. The Stat ISAD Aut

*: H. Momida *et al.*, The 81st JSAP Autumn Meeting 2020, 11p-Z07-13, 11 Sept. 2020.

Where is the position of amorphous structure?

Structure Map of Al₂O₃ w/o P321



 Amorphous model structues (P1) can be distinguished from P-1 by the PC C4.

Structure Map of Al₂O₃ w/o P321



78

Summary

- Features of crystal structure are extracted by dimension reduction from the fingerprint distances.
- The principal components C1 & C2 may be related to the coordination numbers of O@AI and O@O, while C3 & C4 might be so to AI@AI.
- The features can be used as descriptors to model the total energies and energy gaps by regression analysis.



TO: Sci. Technol. Adv. Mater.: Methods <u>4</u>, 2355860 (2024).

計算機によって物性科学はどう変わるか

…インプットからアウトプットまでの膨大な計算の労力から解放されることによって、我々は何をインプットすべきかという創造的な部分と、何をアウトプットから読み取るかという判断力・想像力を要する部分とに専念することができる。…

豊沢 豊 日本物理学会誌1985