# Data analysis and quality assessment by fingerprinting

MARTIN KUBAN

HU BERLIN & IRIS ADLERSHOF

06.08.2023

# Introduction

NOMAD Repository:
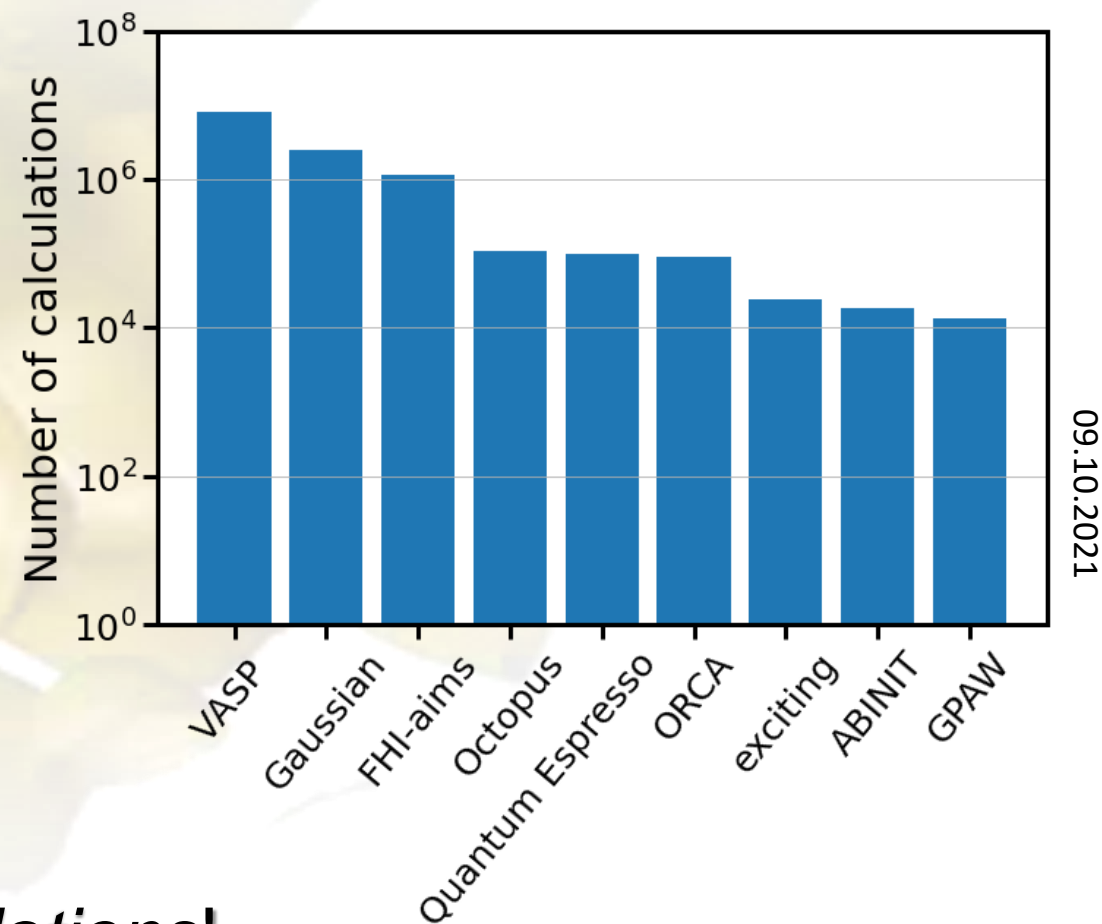
> 100 million calculations

> 40 codes
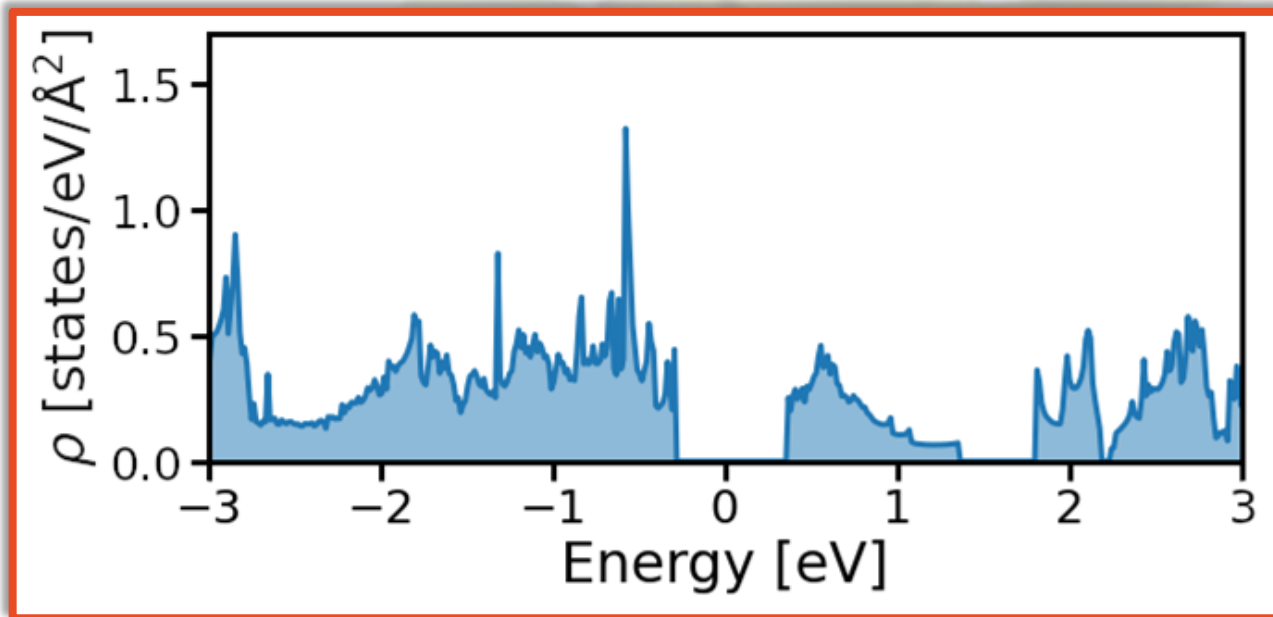
Interoperability?
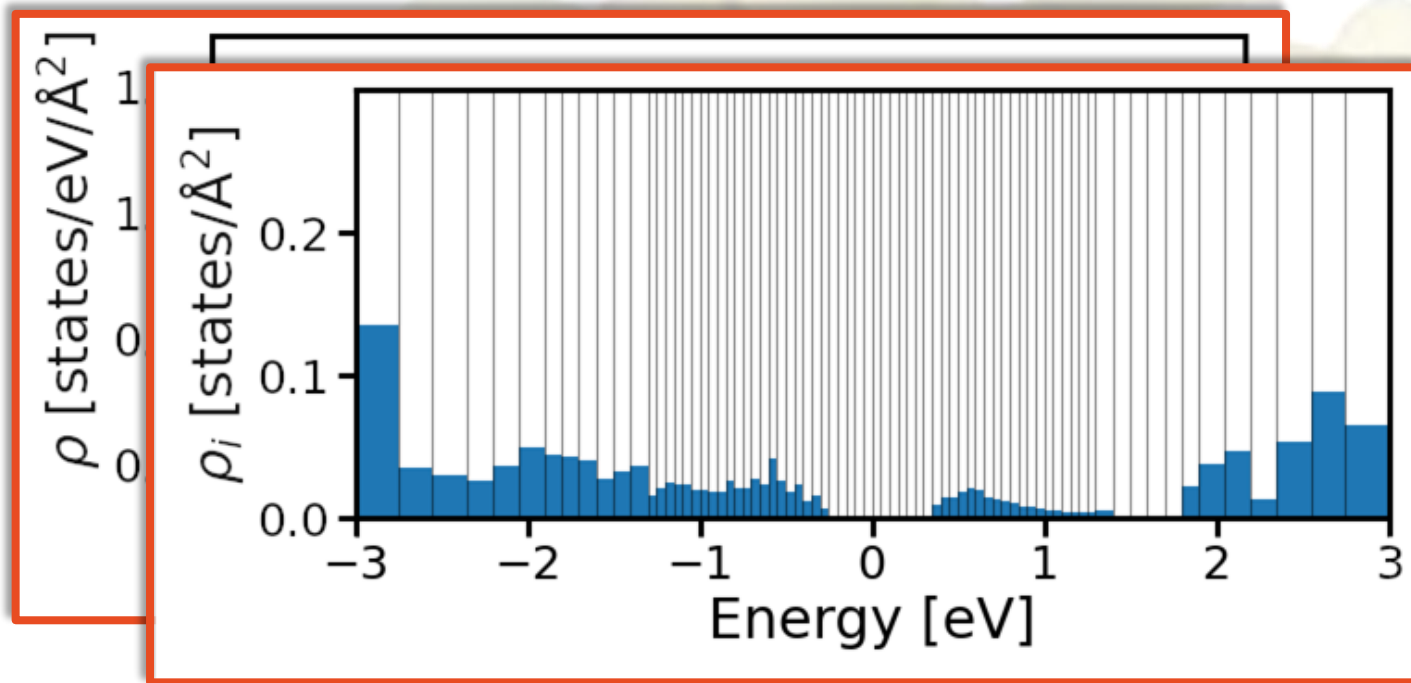
Findability?

Analyze data using *similarity relations*!



09.10.2021

# Spectrum fingerprints



M. Kuban *et al. Sci Data* **9**, 646 (2022).

Inspired by O. Isayev, et al., Chem. Mater. 27, 735 (2015)

# Spectrum fingerprints



M. Kuban *et al. Sci Data* **9**, 646 (2022).

Inspired by O. Isayev, et al., Chem. Mater. 27, 735 (2015)

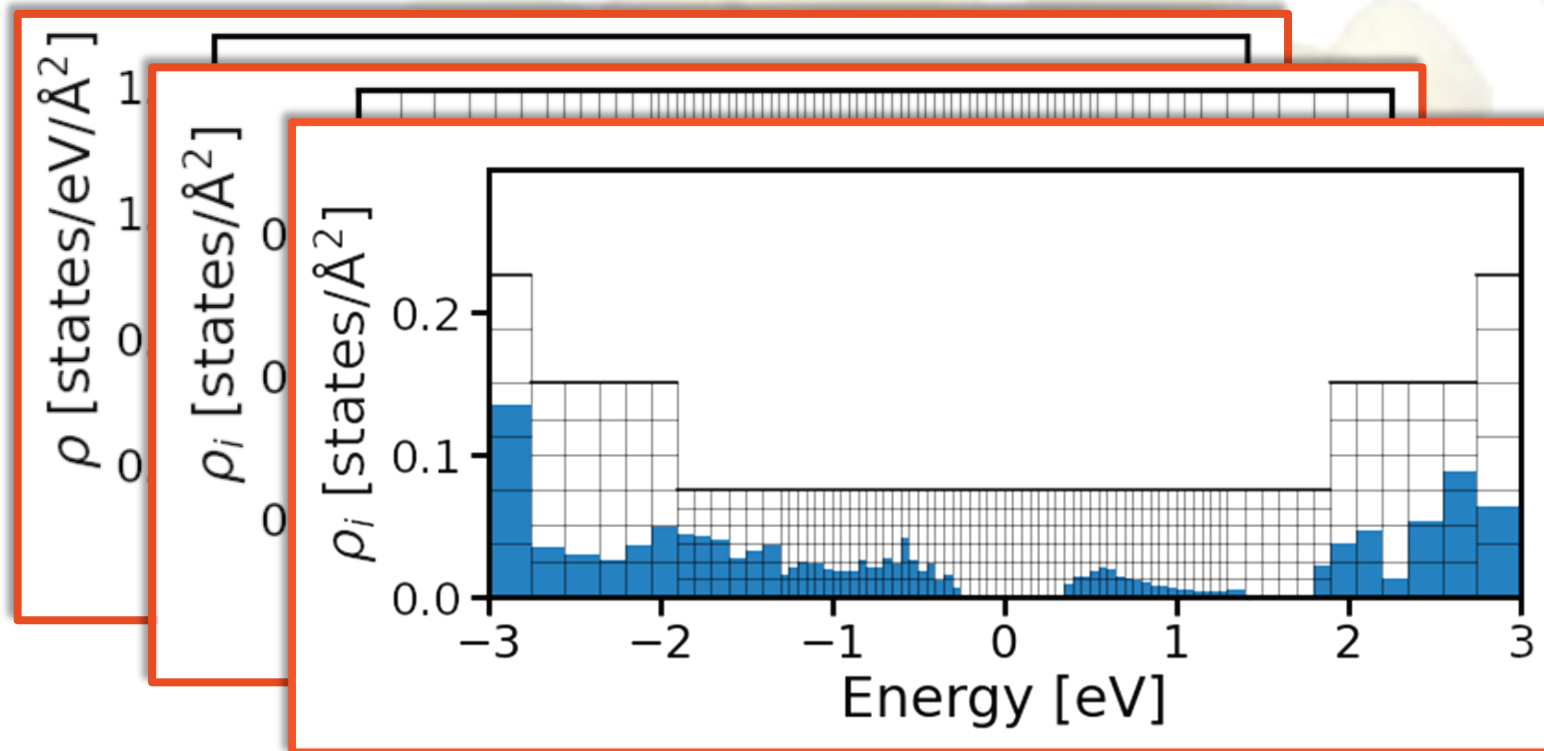# Spectrum fingerprints



M. Kuban *et al. Sci Data* **9**, 646 (2022).

Inspired by O. Isayev, et al., Chem. Mater. 27, 735 (2015)

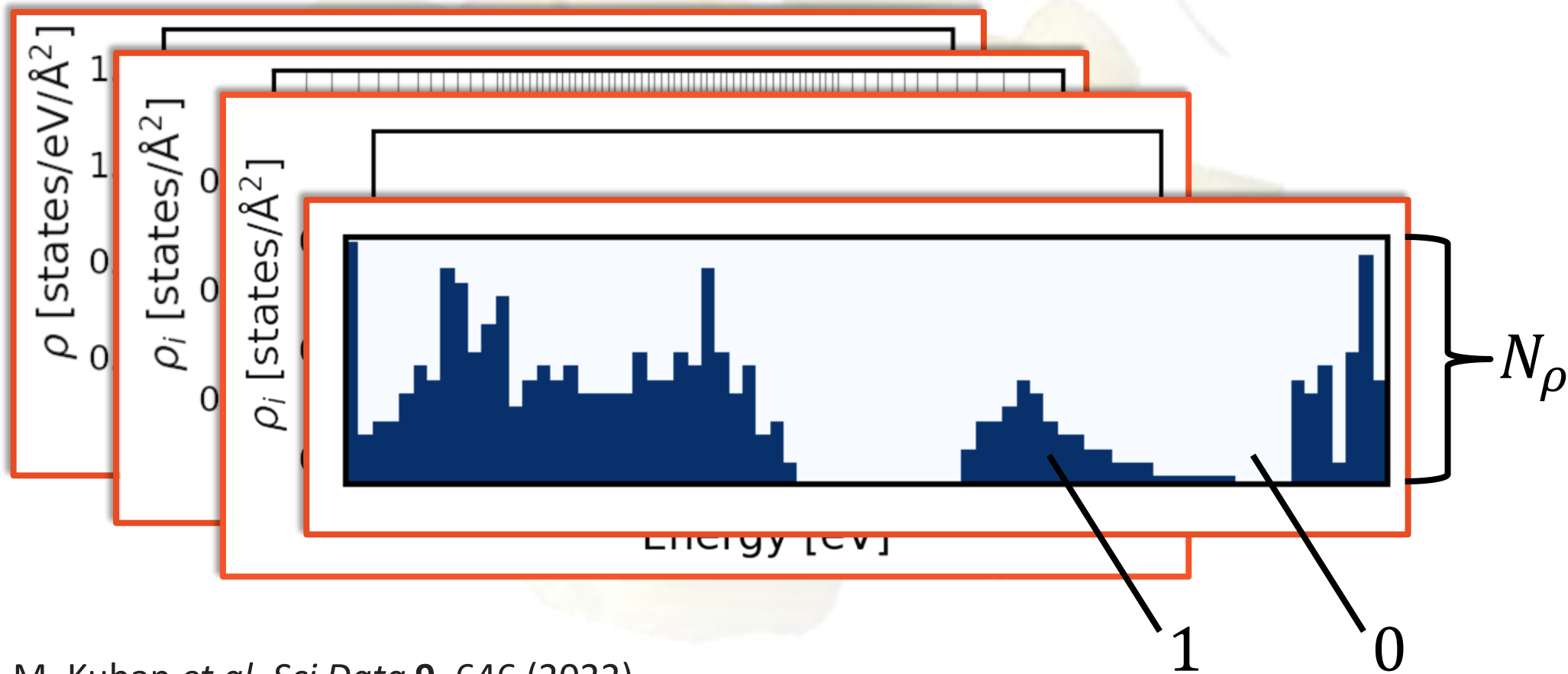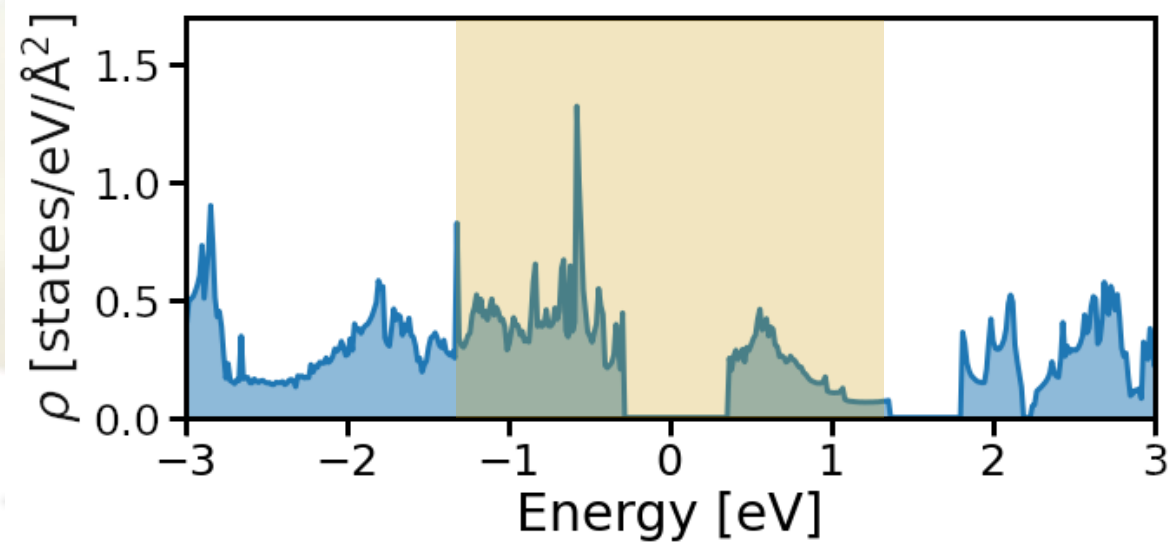# Spectrum fingerprints



M. Kuban *et al. Sci Data* **9**, 646 (2022).

Inspired by O. Isayev, et al., Chem. Mater. 27, 735 (2015)

# Fingerprint feature region



M. Kuban *et al.* Sci Data **9**, 646 (2022).

# Similarity metric

Tanimoto coefficient:

$$\mathrm{Tc}(\boldsymbol{A}, \boldsymbol{B}) = \frac{\boldsymbol{A} \cdot \boldsymbol{B}}{\boldsymbol{A}^2 + \boldsymbol{B}^2 - \boldsymbol{A} \cdot \boldsymbol{B}}$$

**Interpretable**: Intersection divided by union

**Metric**: For binary-valued descriptors

**Computationally cheap**: Can be described by bitwise operations

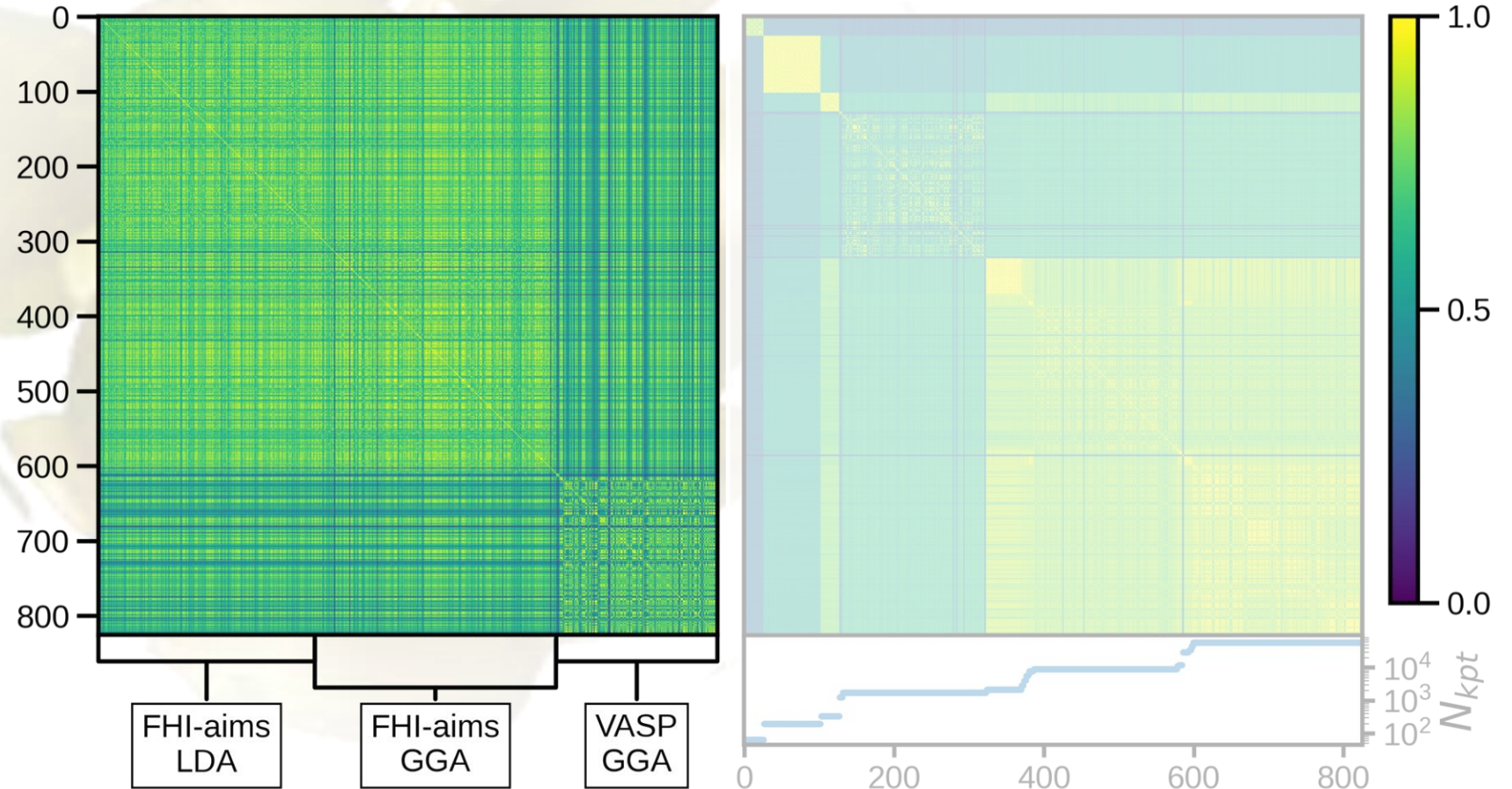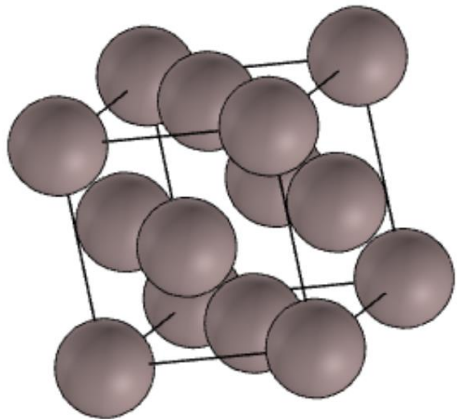Peter Willet et al., *J. Chem. Inf. Comput. Sci.* **38**, 983–996 (1998)

# Similarity metric



M. Kuban *et al.* Sci Data **9**, 646 (2022).

# Data quality assessment

# Code and computational parameters

**NOMAD**

fcc Al



FHI-aims LDA | FHI-aims GGA | VASP GGA

M. Kuban *et al.* MRS Bulletin **47**, 991–999 (2022).

# Code and computational parameters



fcc Al

M. Kuban *et al.* MRS Bulletin **47**, 991–999 (2022).

# Code and computational parameters



fcc Al

M. Kuban *et al.* MRS Bulletin **47**, 991–999 (2022).

# Code and computational parameters



fcc Al

M. Kuban *et al.* MRS Bulletin **47**, 991–999 (2022).

# Code and computational parameters



NOMAD

fcc Al

M. Kuban *et al.* MRS Bulletin **47**, 991–999 (2022).

# Code and computational parameters



fcc Al

FHI-aims LDA | FHI-aims GGA | VASP GGA

M. Kuban *et al.* MRS Bulletin **47**, 991–999 (2022).

# Code and computational parameters



fcc Al

M. Kuban *et al.* MRS Bulletin **47**, 991–999 (2022).

# Code and computational parameters



fcc Al

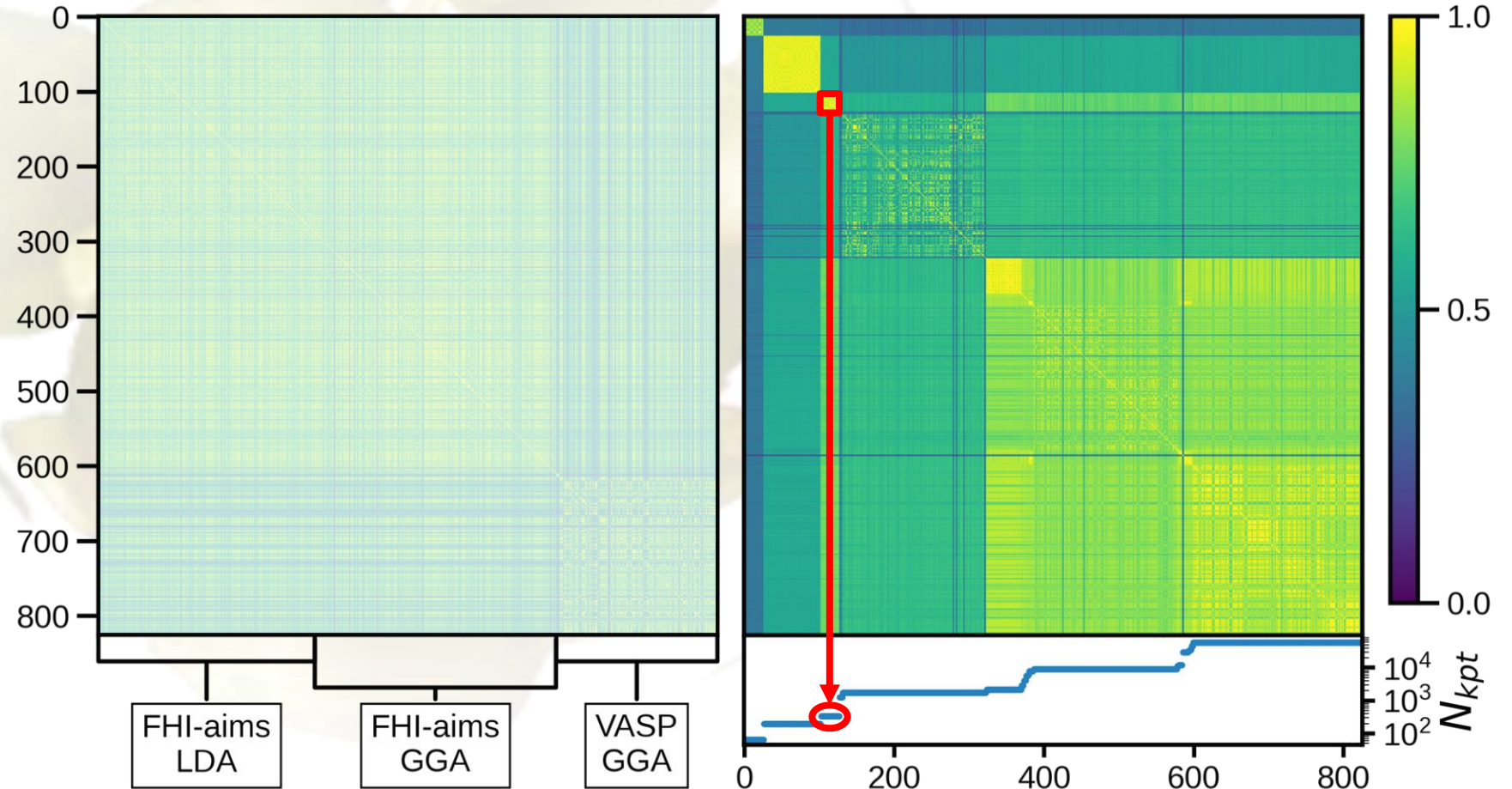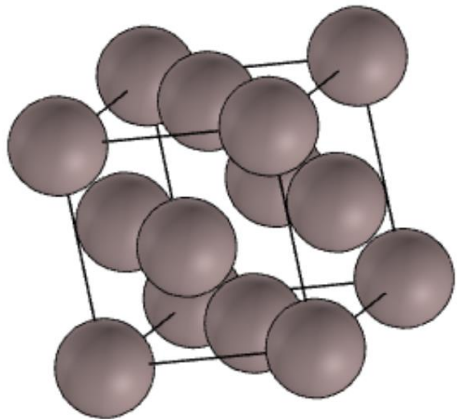M. Kuban *et al.* MRS Bulletin **47**, 991–999 (2022).
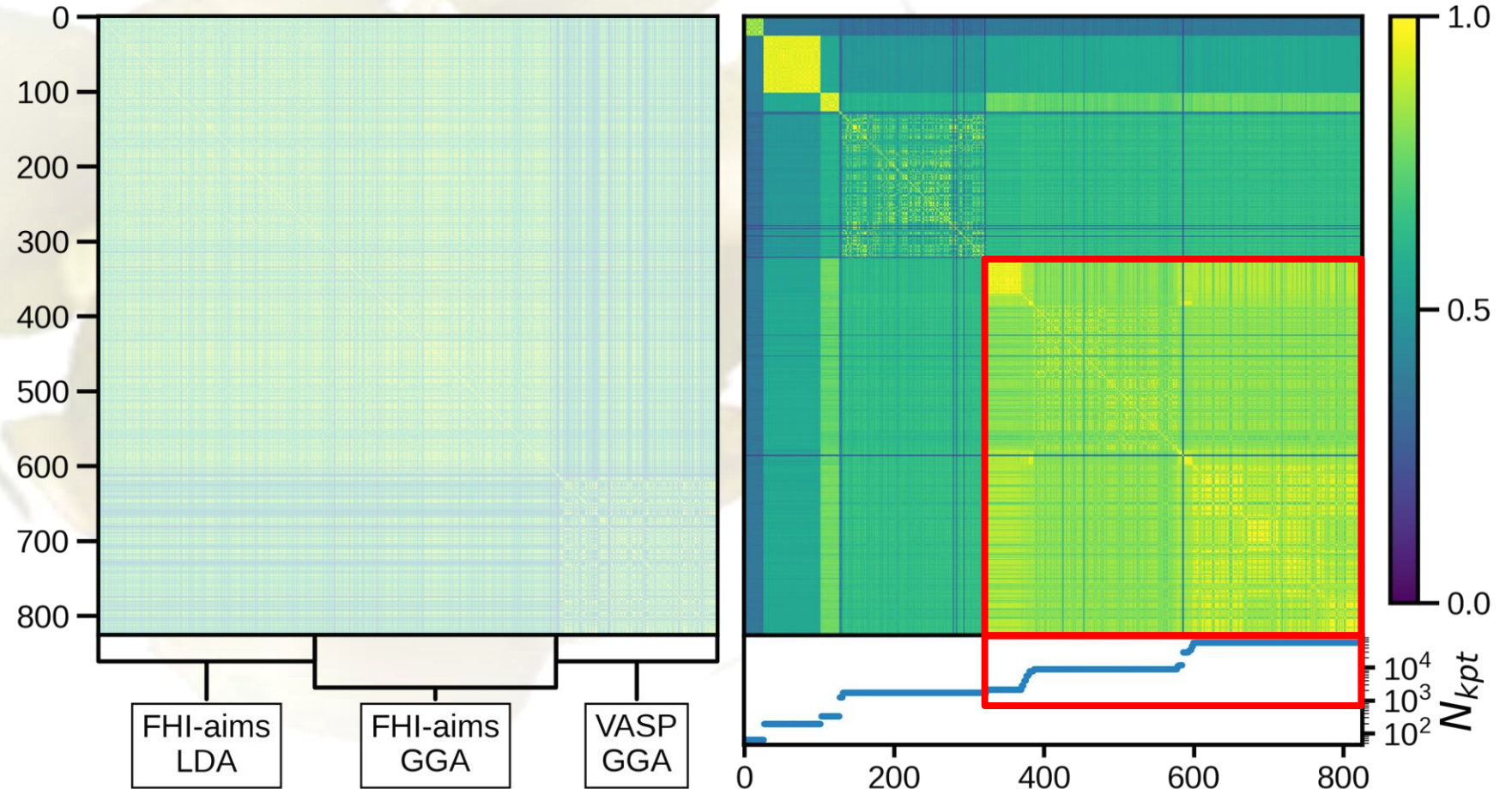
# Code and computational parameters



fcc Al

FHI-aims LDA | FHI-aims GGA | VASP GGA

M. Kuban *et al.* MRS Bulletin **47**, 991–999 (2022).

# Methodology

# Methodology

| PBE vs. | $S_{total}$ | $S_{valence}$ | $S_{conduction}$ |
|---------|-------------|---------------|------------------|
| PBE + SOC | 0.71 | 0.75 | 0.67 |
| HSE | 0.69 | 0.73 | 0.45 |



M. Kuban *et al.* MRS Bulletin **47**, 991–999 (2022).

# Methodology

| PBE vs. | $S_{total}$ | $S_{valence}$ | $S_{conduction}$ |
|---------|-------------|---------------|------------------|
| PBE + SOC | 0.71 | 0.75 | 0.67 |
| HSE | 0.69 | 0.73 | 0.45 |



M. Kuban *et al.* MRS Bulletin **47**, 991–999 (2022).

# Methodology



| PBE vs. | $S_{total}$ | $S_{valence}$ | $S_{conduction}$ |
|---------|-------------|---------------|------------------|
| PBE + SOC | 0.71 | 0.75 | 0.67 |
| HSE | 0.69 | 0.73 | 0.45 |



M. Kuban *et al.* MRS Bulletin **47**, 991–999 (2022).

# Excitonic spectra

$$\phi^{\lambda}(\mathbf{r}_e, \mathbf{r}_h) = \sum A_{vck}^{\lambda} \psi_{vk}^*(\mathbf{r}_h) \psi_{ck}(\mathbf{r}_e)$$



$\Psi_{\lambda=1}(\mathbf{r}_h, \mathbf{r}_e)$

M. Kuban *et al.* MRS Bulletin **47**, 991–999 (2022).

# Exploring data spaces
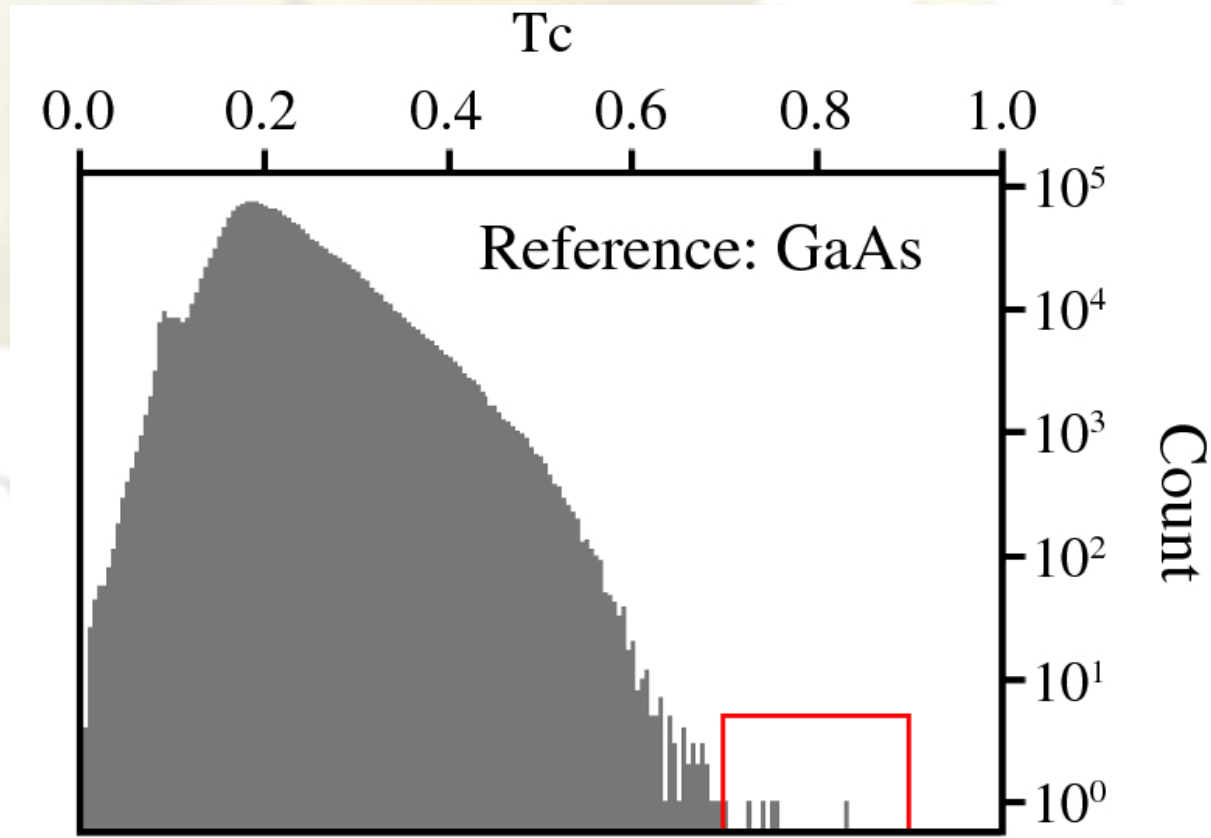
# Finding similar materials

For any material, which other materials are most similar to them?

1.8 million materials

Calculate pairwise similarities

**Search most similar materials**
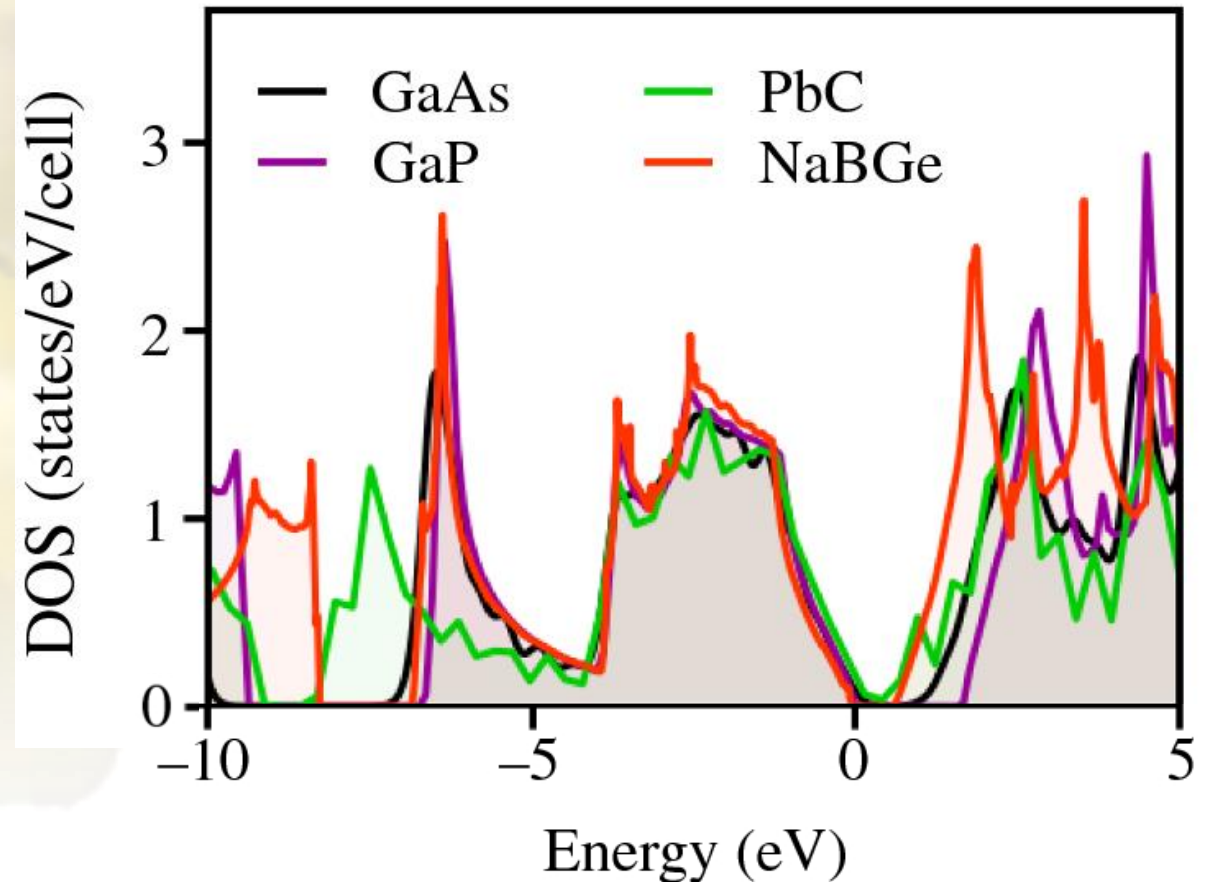
# Finding similar materials

M. Kuban *et al.* MRS Bulletin **47**, 991–999 (2022).

# Finding similar materials

**Similarity of GaAs to:**

GaP:      0.83
PbC:      0.75
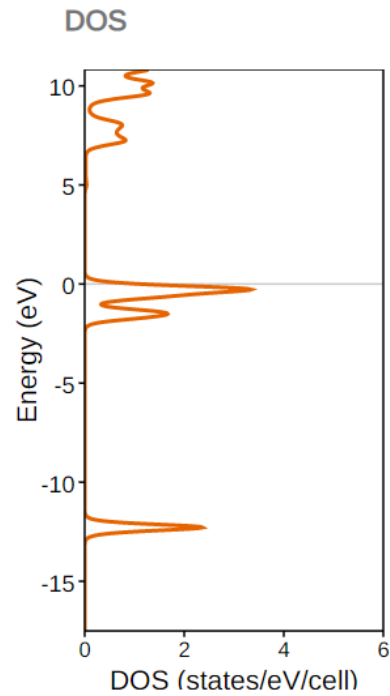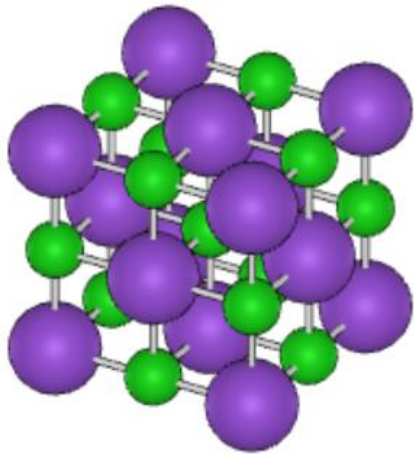NaBGe: 0.74



M. Kuban *et al.* MRS Bulletin **47**, 991–999 (2022).

# Finding similar materials



NaCl - space group 225

DOS

From calculation **zz6kLYmE**
(GGA - VASP)

Similar materials ∨

Similar materials ∧

| Formula (space group) | : Tc |
|---|---|
| $Br_2KLi$ (65) | : 0.565 |
| BrClSr (156) | : 0.555 |
| $Cl_2KLi$ (225) | : 0.541 |
| $BrLi_3Se$ (221) | : 0.539 |
| $ClLi_3Se$ (221) | : 0.537 |

# Finding similar materials

# Exploration of the C2DB

Computational 2D Materials Database [1,2]

High-throughput database

Atomically thin systems

4047 structures, 63 different elements
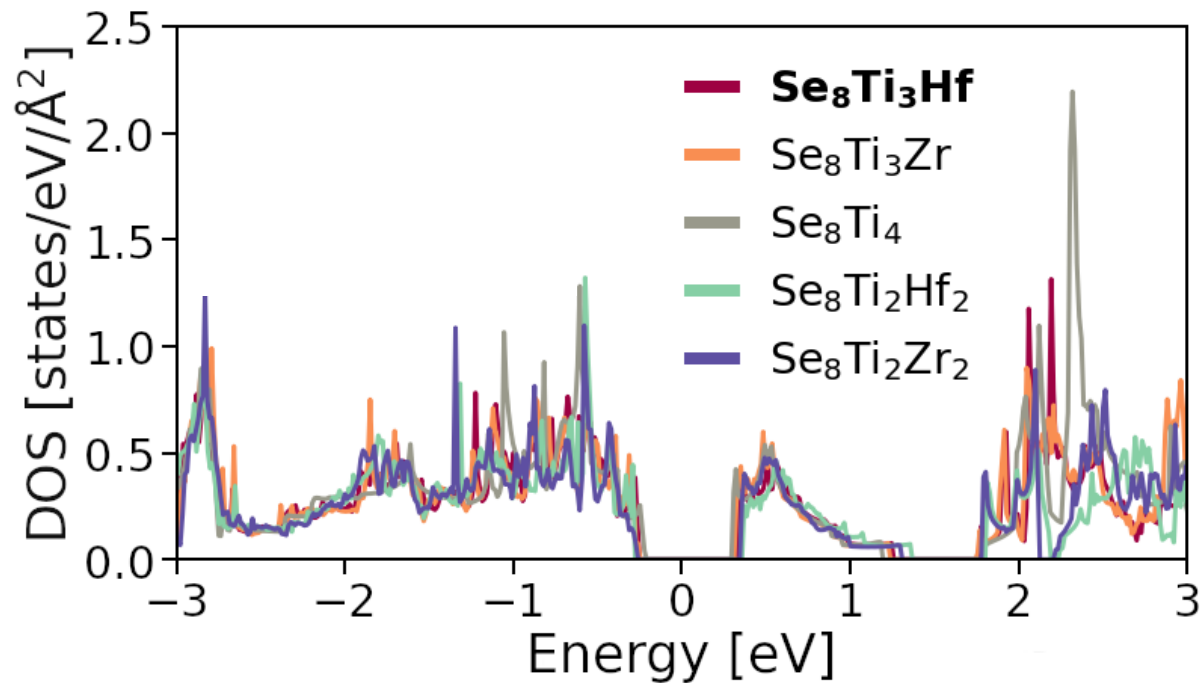
Projected DOS for 3491 structures


The Computational 2D Materials Database

Clustering by using a threshold based algorithm

[1] Sten Haastrup  *et al.*, 2D Materials 5, 042002 (2018)
[2] M. N. Gjerding *et al.*, 2D Materials 8, 044002 (2021)

# Cluster analysis



M. Kuban *et al.* Sci Data **9**, 646 (2022).

# Cluster analysis



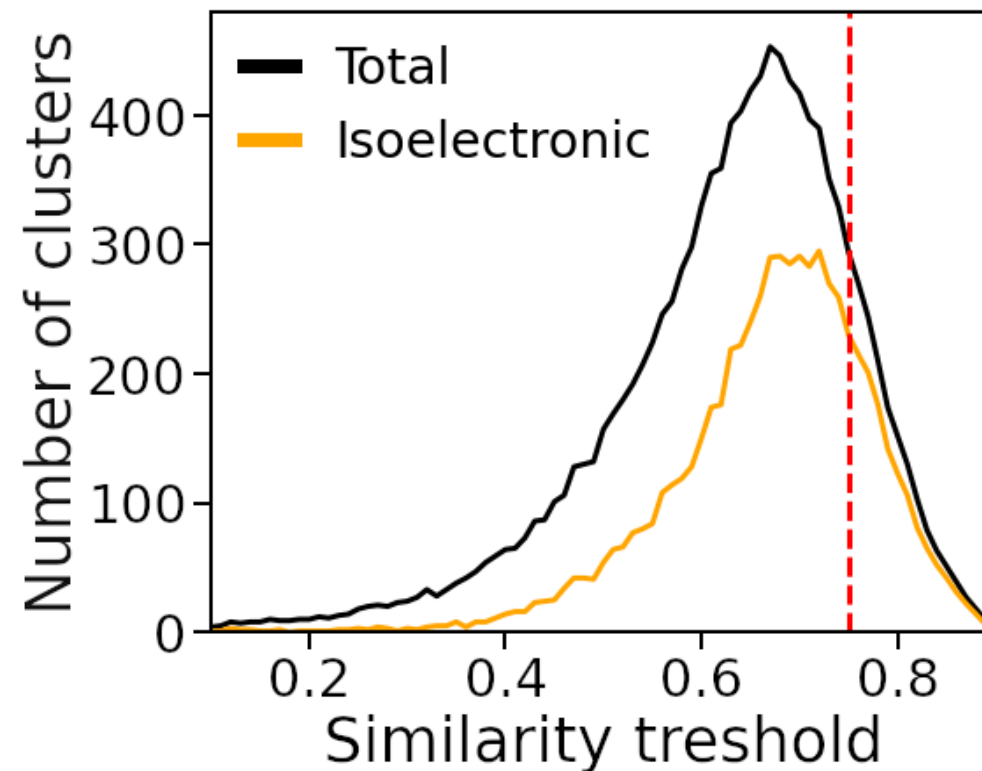M. Kuban *et al.* Sci Data **9**, 646 (2022).

# Isoelectronic clusters
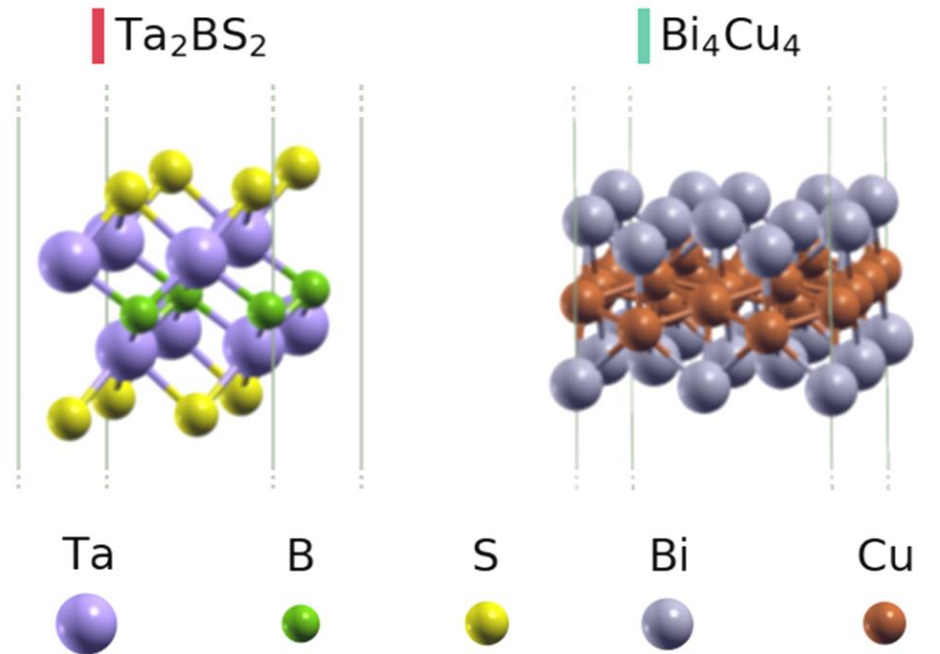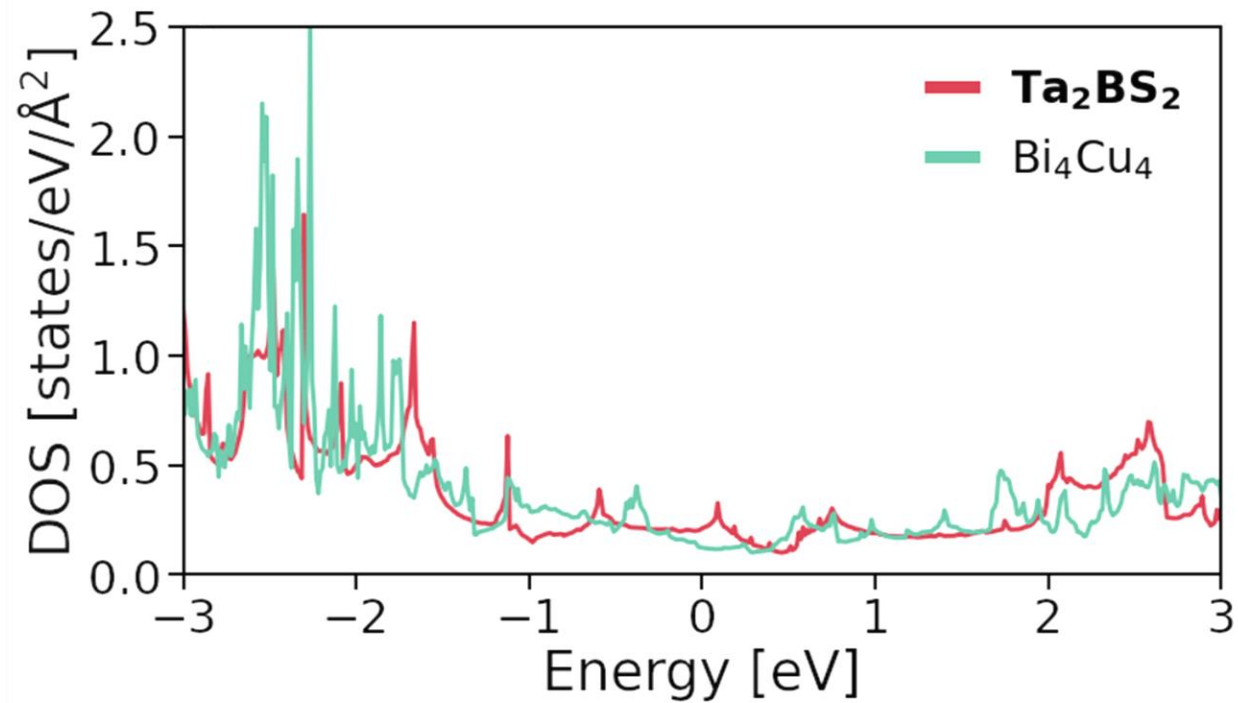
Descriptor for isoelectronic clusters:

$$\bar{c}_m = \frac{1}{N_{Atoms}} \sum_{i=1}^{N_{Atoms}} c_i$$

230 isoelectronic clusters (78%)



M. Kuban *et al.* Sci Data **9**, 646 (2022).

# Outliers



M. Kuban *et al.* Sci Data **9**, 646 (2022).

# Conclusions

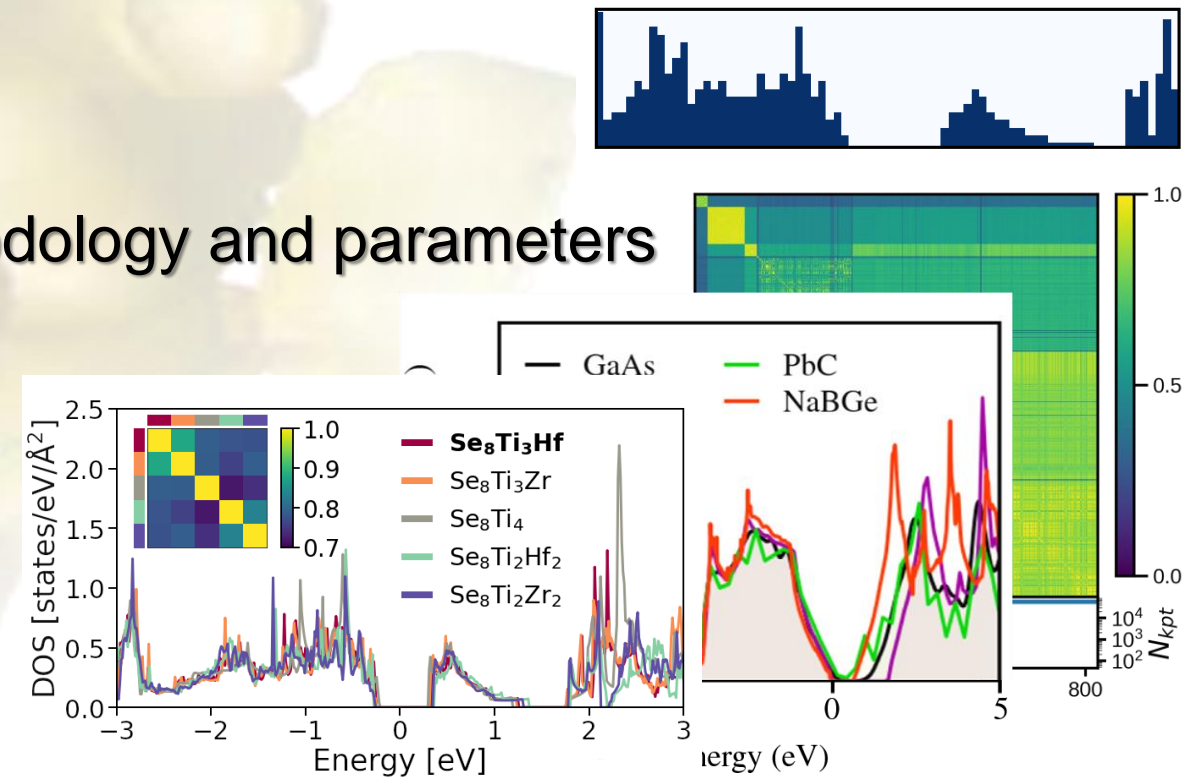**Spectral fingerprint to quantitatively evaluate the similarity of spectra**

**Quality assessment:**

Measuring the impact of methodology and parameters

**Data analytics:**

Finding similar materials

Unsupervised learning

# Acknowledgements

# Further reading

Kuban, M., Rigamonti, S., Scheidgen, M., Draxl, C., Density-of-states similarity descriptor for unsupervised learning from materials data. *Sci Data* **9**, 646 (2022)

Kuban, M., Gabaj, Š., Aggoune, W. *et al.* Similarity of materials and data-quality assessment by fingerprinting. *MRS Bulletin* **47**, 991–999 (2022).

NOMAD AI Toolkit Tutorial:

`https://nomad-lab.eu/aitoolkit/tutorial-dos-similarity`