

NOVEL MATERIALS DISCOVERY



MAX-PLANCK-GESELLSCHAFT

HPC for Computationally and Data-Intensive Problems

Luca M. Ghiringhelli
Fritz-Haber-Institut der MPG, Berlin

SuperMUC-NG
Next-Gen Science Symposium
Munich, November 22, 2018

Data-driven materials science

The Big Picture

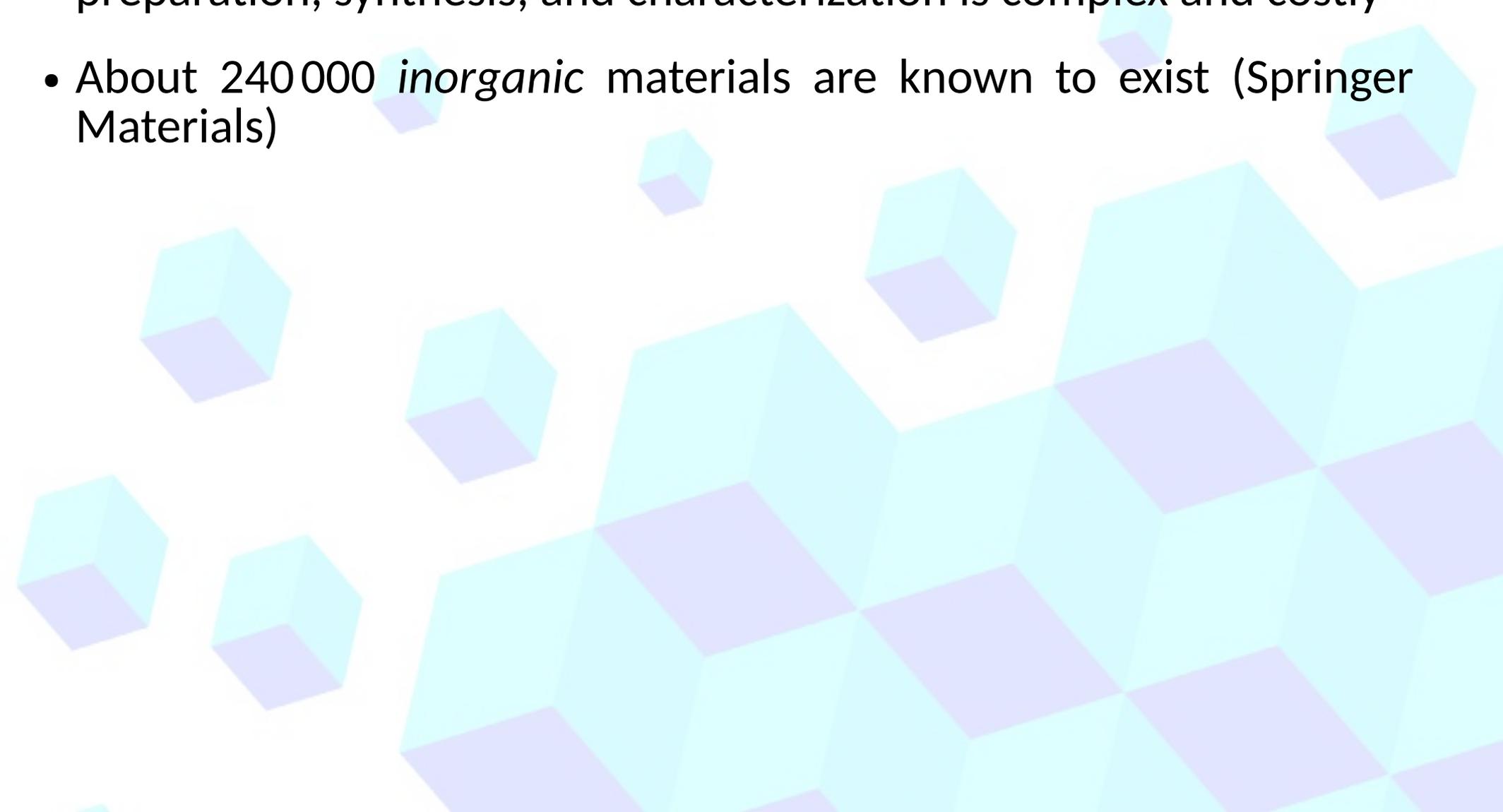
- Design of new materials: preparation, synthesis, and characterization is complex and costly



Data-driven materials science

The Big Picture

- Design of new materials: preparation, synthesis, and characterization is complex and costly
- About 240 000 *inorganic* materials are known to exist (Springer Materials)



Data-driven materials science

The Big Picture

- Design of new materials: preparation, synthesis, and characterization is complex and costly
- About 240 000 *inorganic* materials are known to exist (Springer Materials)
- Basic properties determined for very few of them

Data-driven materials science

The Big Picture

- Design of new materials: preparation, synthesis, and characterization is complex and costly
- About 240 000 *inorganic* materials are known to exist (Springer Materials)
- Basic properties determined for very few of them
- Number of possible materials: practically infinite

Data-driven materials science

The Big Picture

- Design of new materials: preparation, synthesis, and characterization is complex and costly
 - About 240 000 *inorganic* materials are known to exist (Springer Materials)
 - Basic properties determined for very few of them
 - Number of possible materials: practically infinite
- ⇒ New materials with superior properties exist but not yet known

Data-driven materials science

The Big Picture

- Design of new materials: preparation, synthesis, and characterization is complex and costly
 - About 240 000 *inorganic* materials are known to exist (Springer Materials)
 - Basic properties determined for very few of them
 - Number of possible materials: practically infinite
- ⇒ New materials with superior properties exist but not yet known
- Data analytics tools will help to identify trends and anomalies in data and guide discovery of new materials



NOMAD

NOVEL MATERIALS DISCOVERY

The NOMAD Laboratory

<https://nomad-coe.eu>



NOMAD

The NOMAD Laboratory

A European Centre of Excellence

[HOME](#) [PROJECT](#) [INDUSTRY](#) [TEAM](#) [RELATED PROJECTS](#) [NEWS](#) [PRESS KIT](#) [CONTACT US](#)

Enter Search... 



NOMAD REPOSITORY



THE ARCHIVE



ENCYCLOPEDIA



BIG-DATA ANALYTICS



ADVANCED GRAPHICS



HPC INFRASTRUCTURE



OUTREACH

The Novel Materials Discovery (NOMAD) Laboratory maintains the largest Repository, for input and output files of all important computational materials science codes.

From its open-access data, it builds several *Big-Data Services* helping to advance materials science and engineering.

To learn more, click on the buttons above. You can also watch our 3-minute summary on the *NOMAD Laboratory CoE* at [YouTube](#) (or at [YOUKU](#) in China).

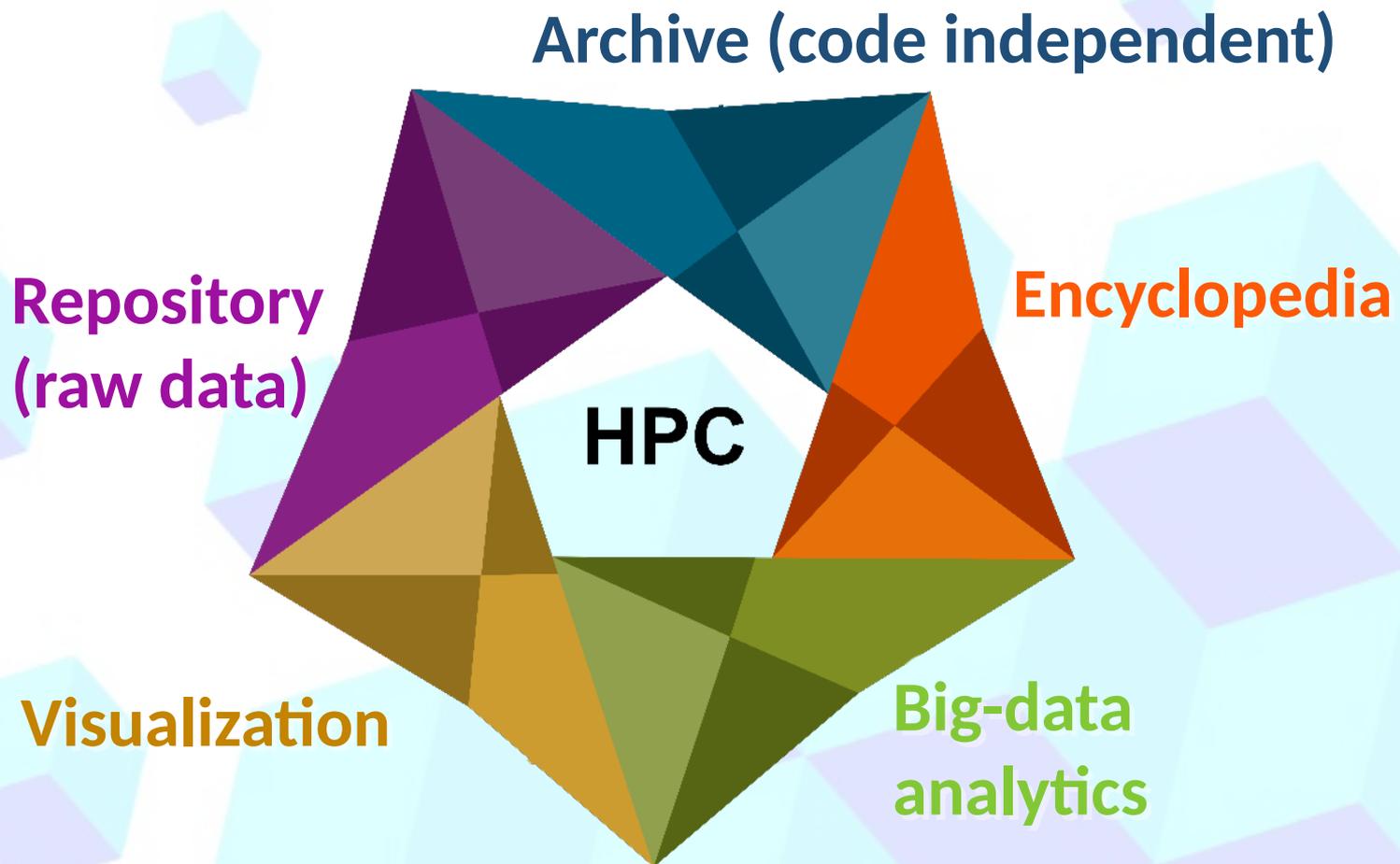
NOMAD Scope and Overview

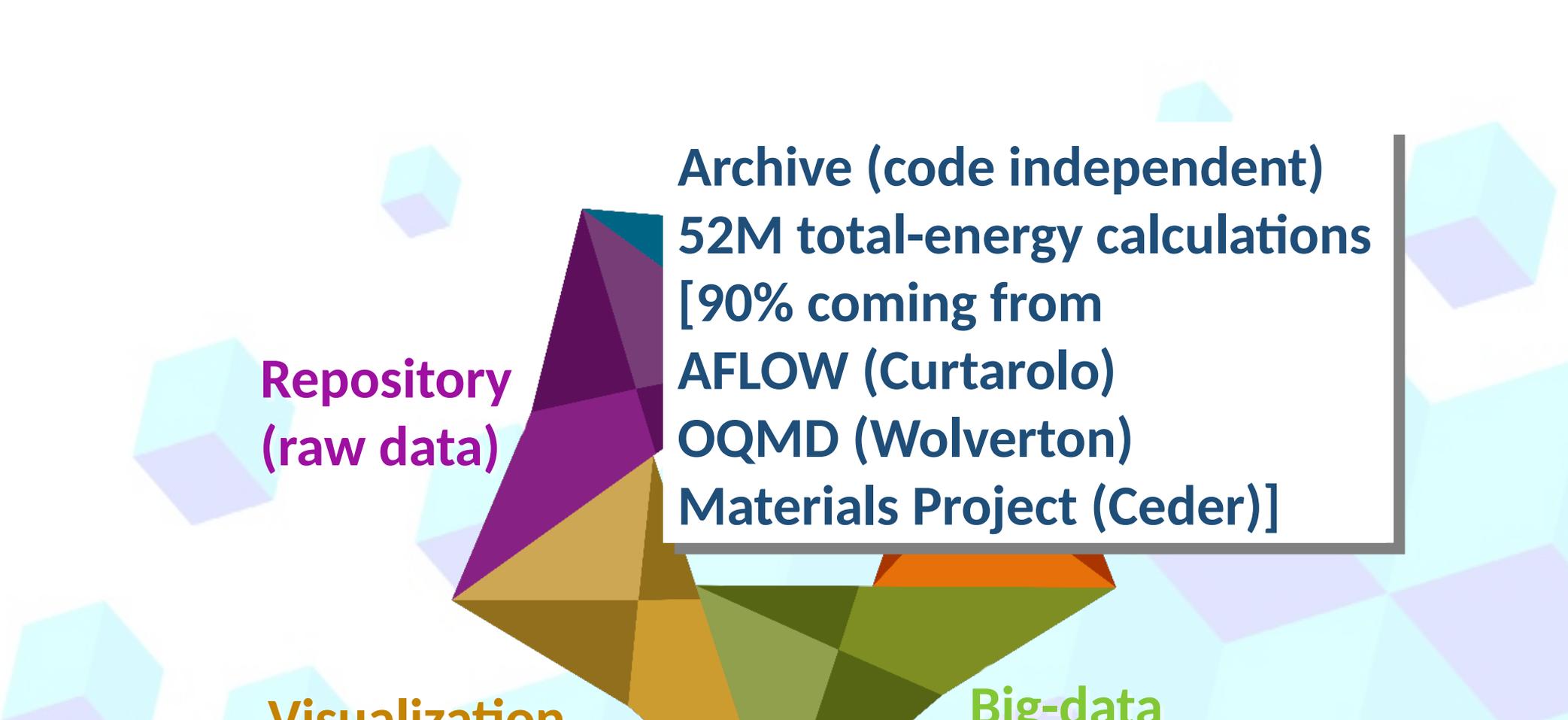
NOMAD Success Stories

Finding local patterns and structure
in big-data of materials-science
remains a challenge



New data mining tools must be developed
to help uncover hidden relations
in materials-science data





**Repository
(raw data)**

**Archive (code independent)
52M total-energy calculations
[90% coming from
AFLOW (Curtarolo)
OQMD (Wolverton)
Materials Project (Ceder)]**

Visualization

**Big-data
analytics**

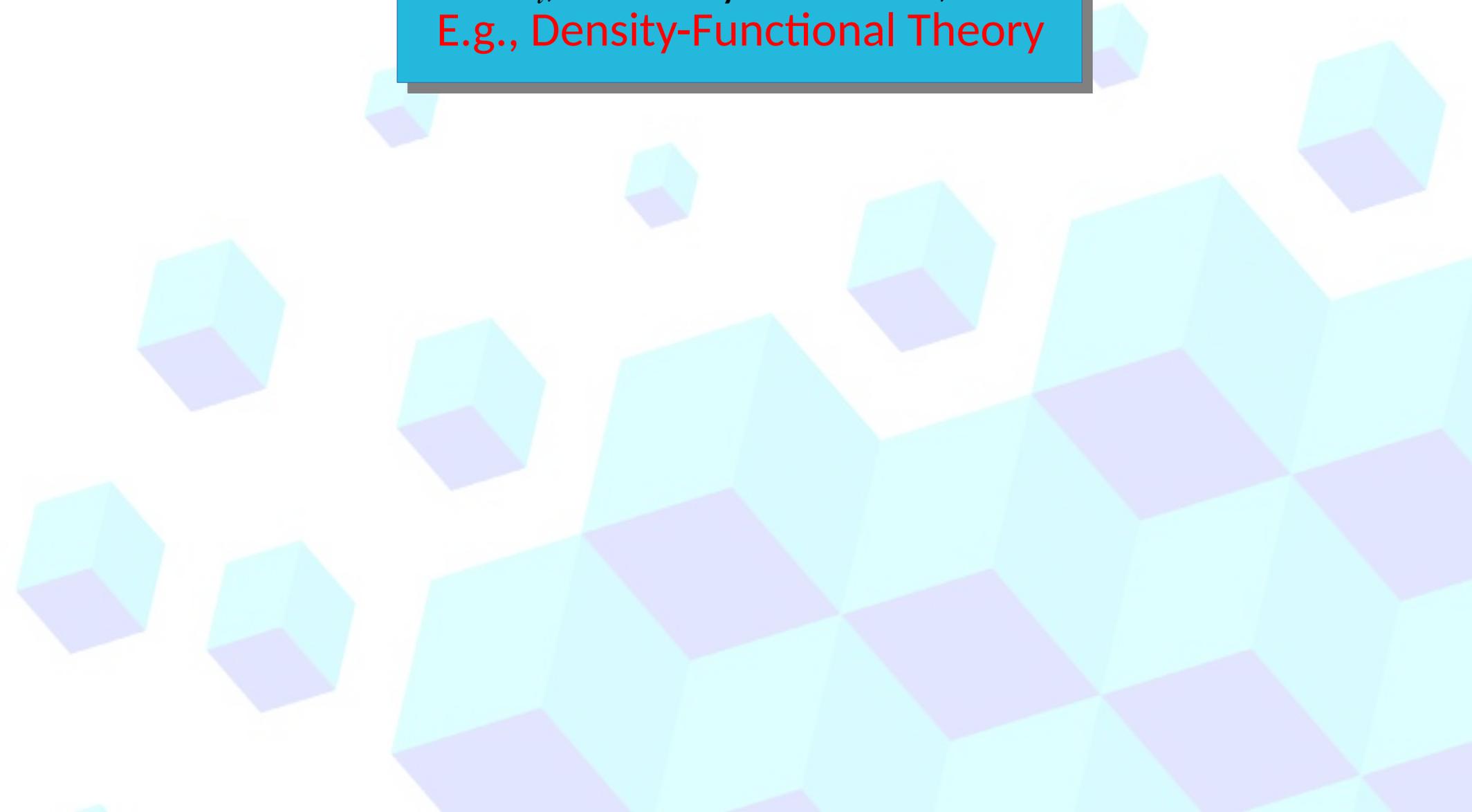
Data analytics: an ideal flow chart

Training set

Calculate properties and functions

P_i , for many *materials*, i

E.g., Density-Functional Theory



Data analytics: an ideal flow chart

Training set

Calculate properties and functions

P_i , for many *materials*, i

E.g., Density-Functional Theory

Descriptor

Find the appropriate descriptor d_i , build a table:

i	d_i	P_i
-----	-------	-------

Data analytics: an ideal flow chart

Training set

Calculate properties and functions

P_i , for many *materials*, i

E.g., Density-Functional Theory

Descriptor

Find the appropriate descriptor d_i , build a table:

i	d_i	P_i
-----	-------	-------

Learning

Find the function $P(d)$ for the table.

Build a **chart** for the property

Statistical learning

Data analytics: an ideal flow chart

Training set

Calculate properties and functions P_i , for many *materials*, i
E.g., Density-Functional Theory

Fast Prediction

Calculate properties and functions for new values of d (new materials)

Descriptor

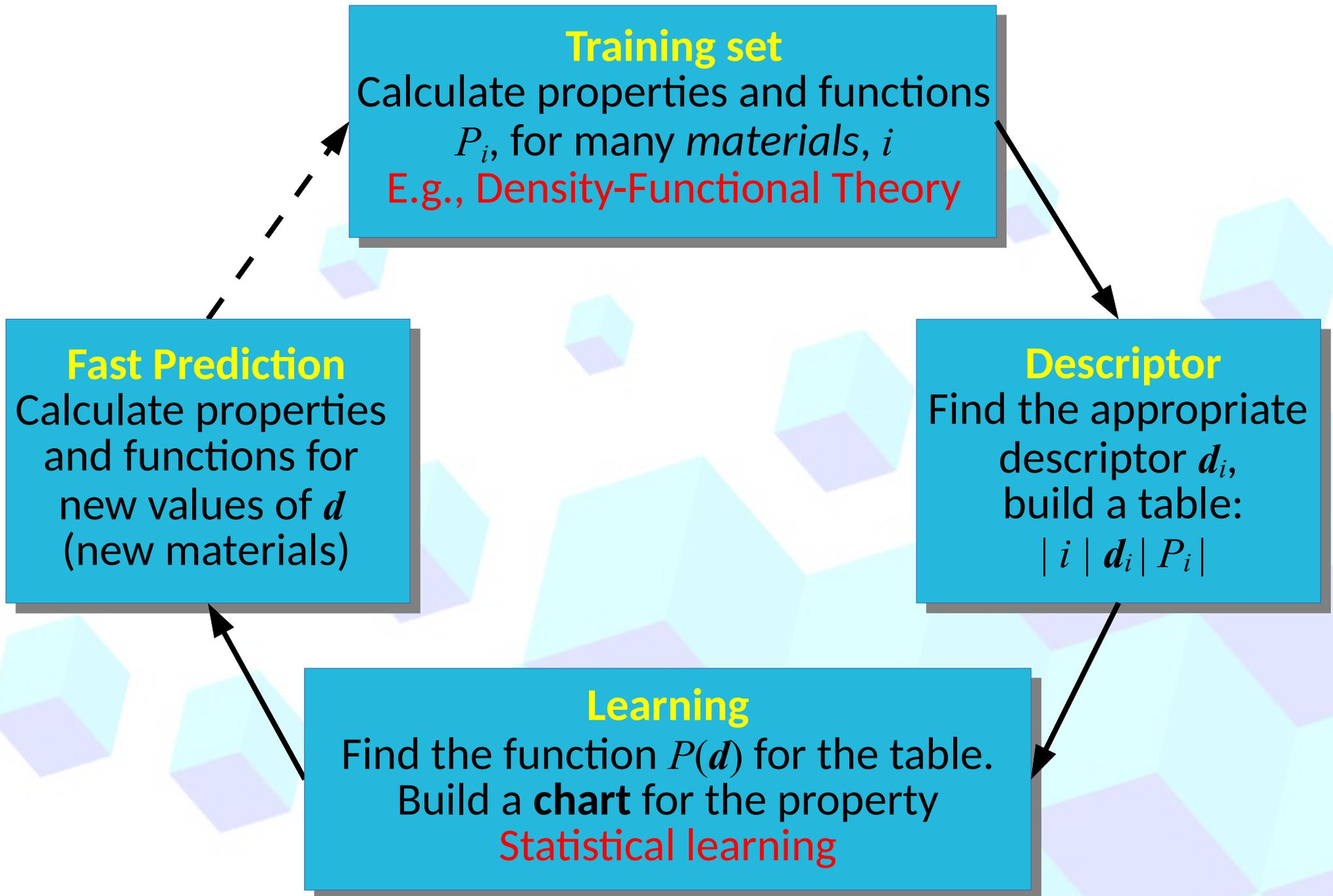
Find the appropriate descriptor d_i , build a table:

i	d_i	P_i
-----	-------	-------

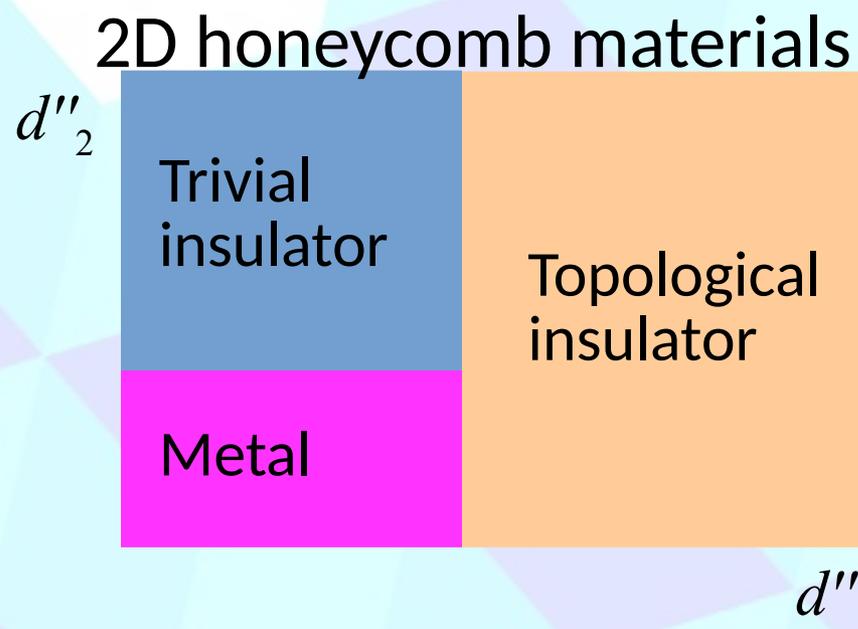
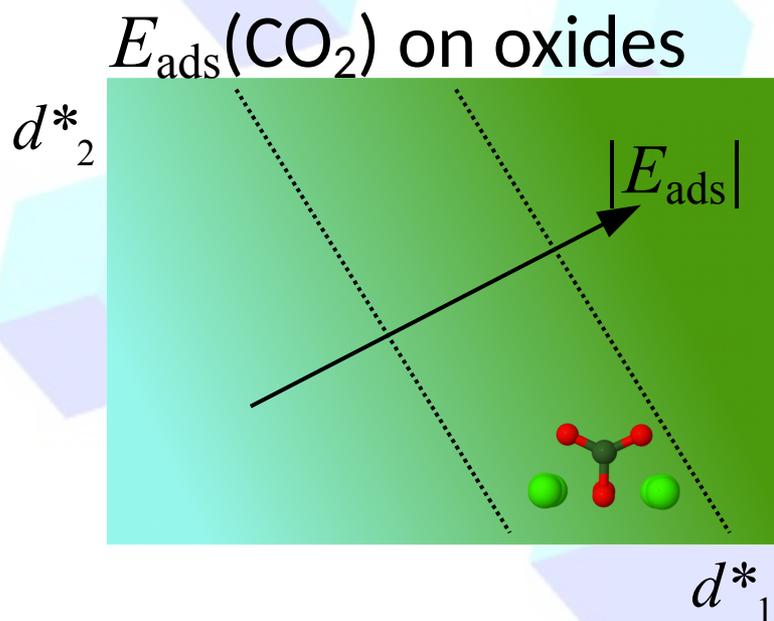
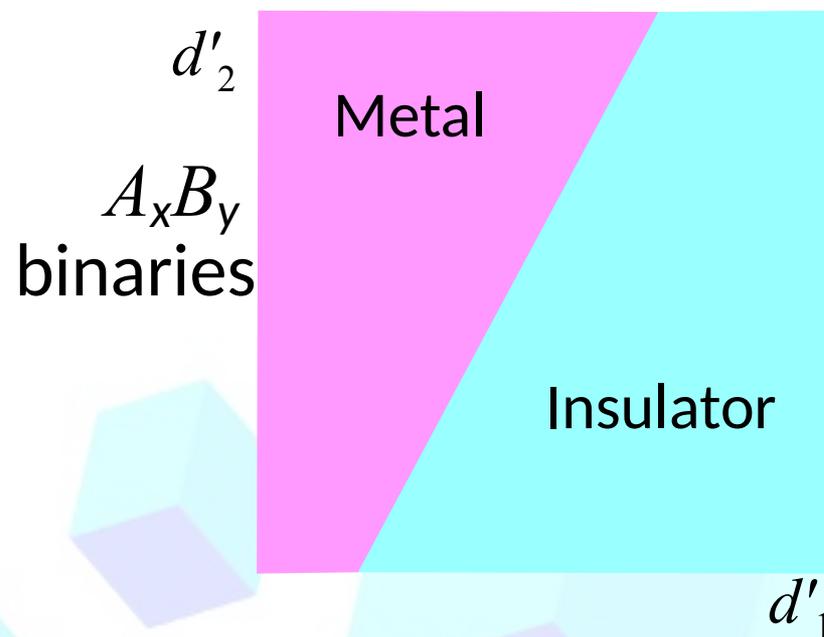
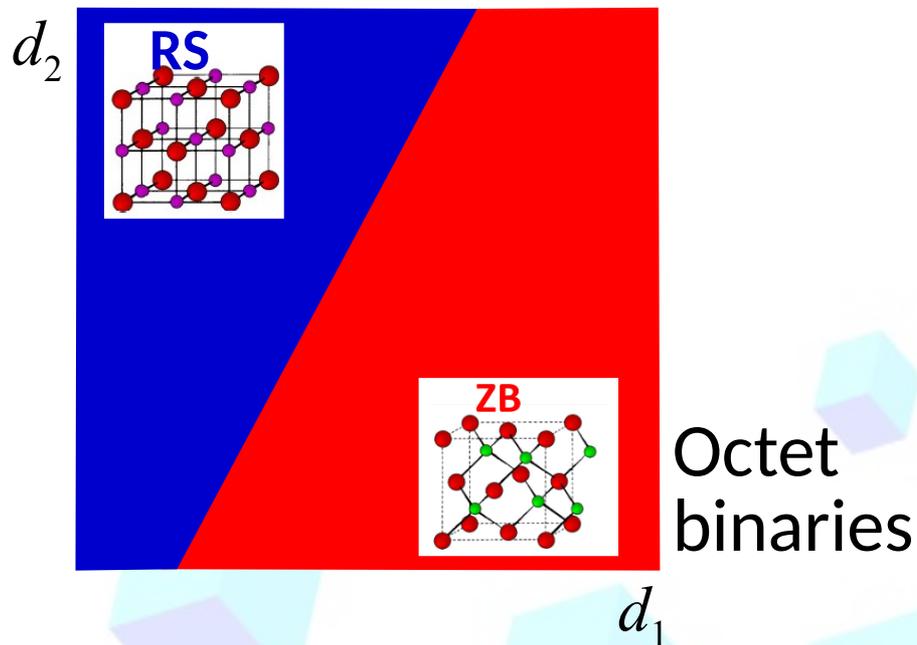
Learning

Find the function $P(d)$ for the table.
Build a **chart** for the property
Statistical learning

Data analytics: an ideal flow chart

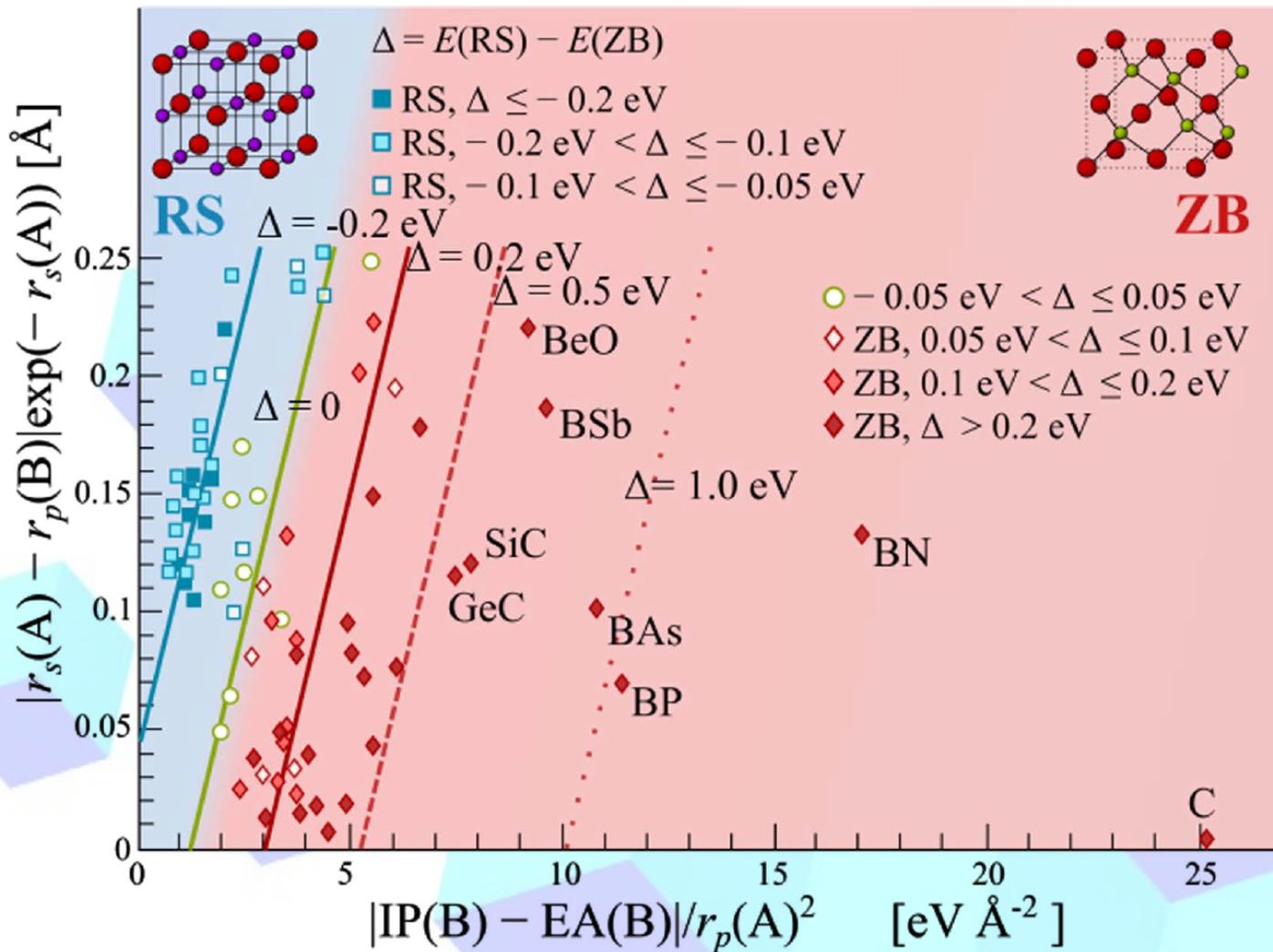


Learning/discovering maps of materials properties. A quantum many-body problem



Compressed sensing: the quest for descriptors and predictive models

Structure map with compressed-sensing algorithm, starting from 7 atomic features



Data analytics with compressed sensing

Compressed sensing

Aim: finding descriptors and learning predictive models

Ansatz:

$$\mathbf{P} = c_1 \mathbf{d}_1 + c_2 \mathbf{d}_2 + \dots + c_n \mathbf{d}_n$$

\mathbf{P} : property of interest

$\mathbf{d}_1, \dots, \mathbf{d}_n$: features, i.e., (nonlinear) functions of *primary features* (EA, IP, ...)

c_1, \dots, c_n : unknown coefficients => as few as possible are *nonzero*

Data analytics with compressed sensing

Compressed sensing

Aim: finding descriptors and learning predictive models

Ansatz:

$$P = c_1 d_1 + c_2 d_2 + \dots + c_n d_n$$

P : property of interest

d_1, \dots, d_n : features, i.e., (nonlinear) functions of *primary features* (EA, IP, ...)

c_1, \dots, c_n : unknown coefficients => as few as possible are *nonzero*

d_i : iterative construction with features and operators (+, ×, /, ², ...)

$$\frac{IP(B) - EA(B)}{r_p(A)^2}, \frac{|r_s(A) - r_p(B)|}{\exp(r_s(A))}$$

Data analytics with compressed sensing

Compressed sensing

Aim: finding descriptors and learning predictive models

Ansatz:

$$\mathbf{P} = c_1 \mathbf{d}_1 + c_2 \mathbf{d}_2 + \dots + c_n \mathbf{d}_n$$

\mathbf{P} : property of interest

$\mathbf{d}_1, \dots, \mathbf{d}_n$: features, i.e., (nonlinear) functions of *primary features* (EA, IP, ...)

c_1, \dots, c_n : unknown coefficients => as few as possible are *nonzero*

$$\operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^M} \|\mathbf{P} - \mathbf{D}\mathbf{c}\|_2^2 + \lambda \|\mathbf{c}\|_0$$

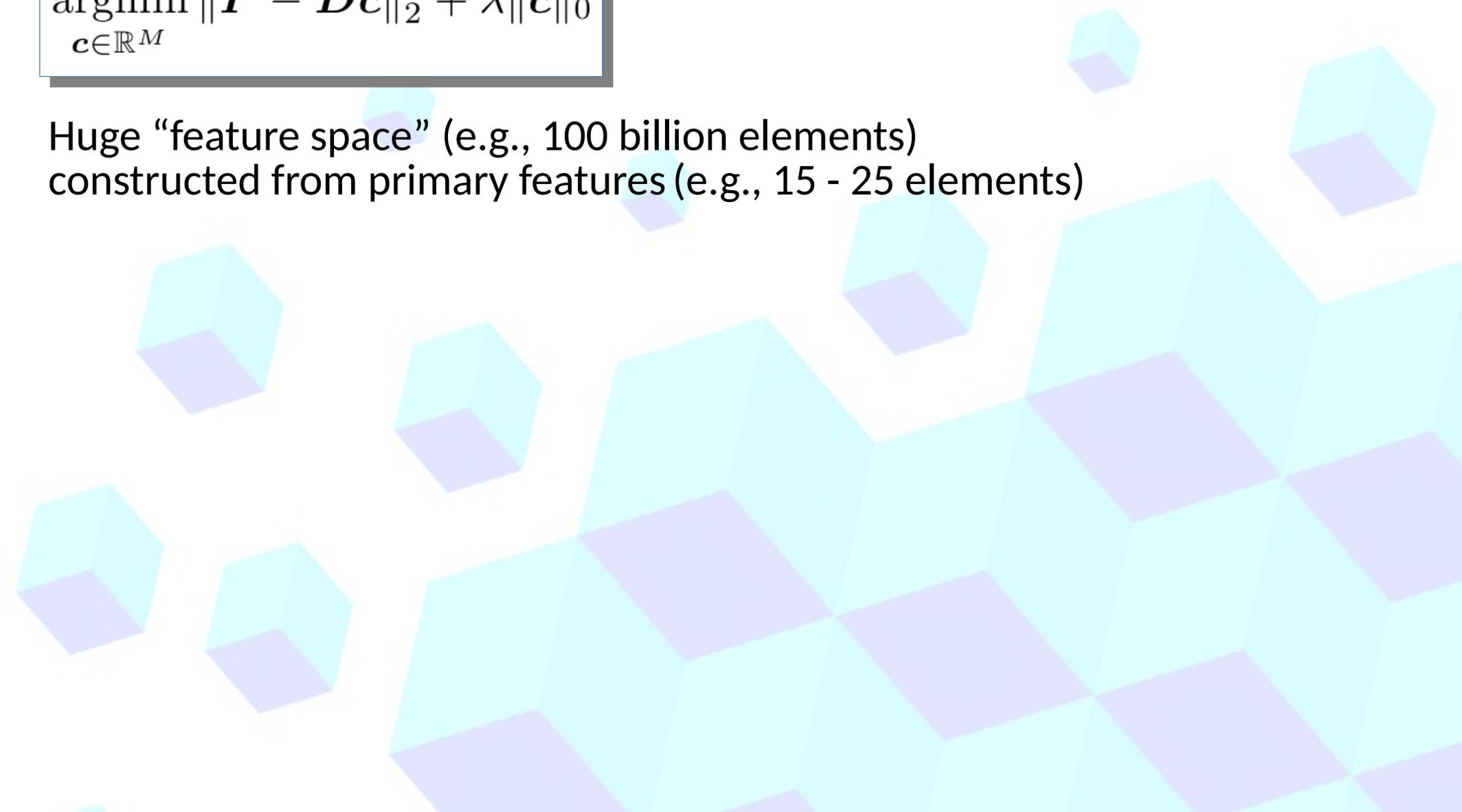
\mathbf{d}_i : iterative construction with features and operators (+, ×, /, ², ...)

$$\frac{\text{IP}(\text{B}) - \text{EA}(\text{B})}{r_p(\text{A})^2}, \frac{|r_s(\text{A}) - r_p(\text{B})|}{\exp(r_s(\text{A}))}$$

SISSO: sure independence screening plus sparsifying operator

$$\operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^M} \|\mathbf{P} - \mathbf{D}\mathbf{c}\|_2^2 + \lambda \|\mathbf{c}\|_0$$

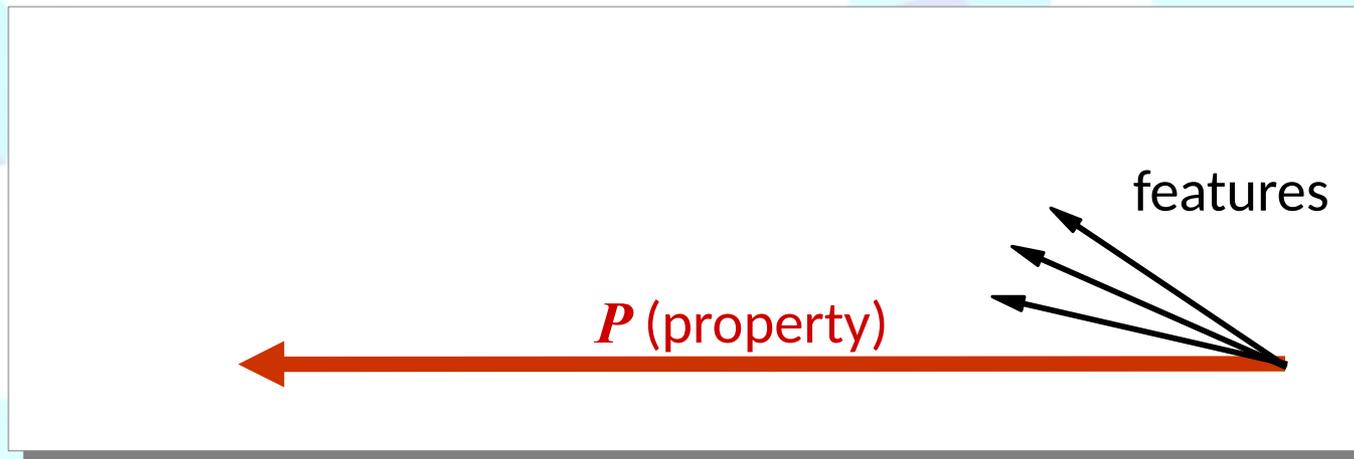
Huge “feature space” (e.g., 100 billion elements)
constructed from primary features (e.g., 15 - 25 elements)



SISSO: sure independence screening plus sparsifying operator

$$\operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^M} \|\mathbf{P} - \mathbf{D}\mathbf{c}\|_2^2 + \lambda \|\mathbf{c}\|_0$$

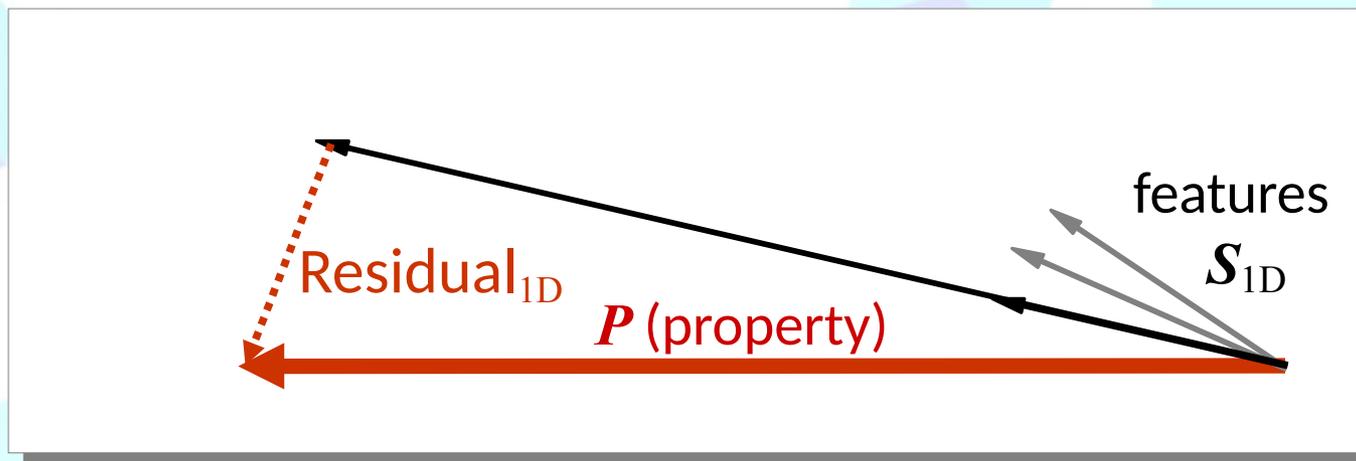
Huge “feature space” (e.g., 100 billion elements)
constructed from primary features (e.g., 15 - 25 elements)



SISSO: sure independence screening plus sparsifying operator

$$\operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^M} \|\mathbf{P} - \mathbf{D}\mathbf{c}\|_2^2 + \lambda \|\mathbf{c}\|_0$$

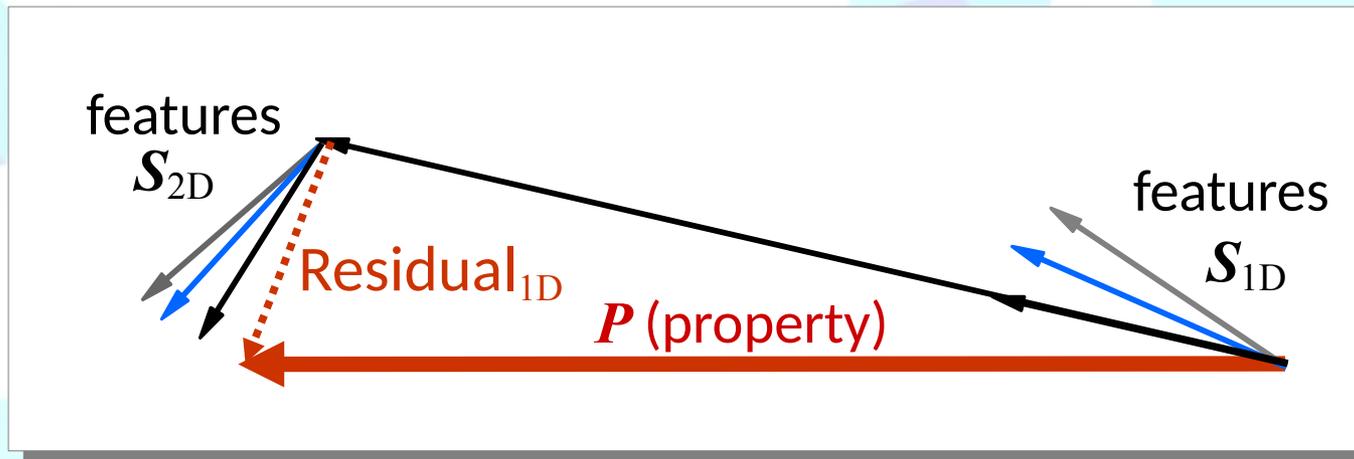
Huge “feature space” (e.g., 100 billion elements)
constructed from primary features (e.g., 15 - 25 elements)



SISSO: sure independence screening plus sparsifying operator

$$\operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^M} \|\mathbf{P} - \mathbf{D}\mathbf{c}\|_2^2 + \lambda \|\mathbf{c}\|_0$$

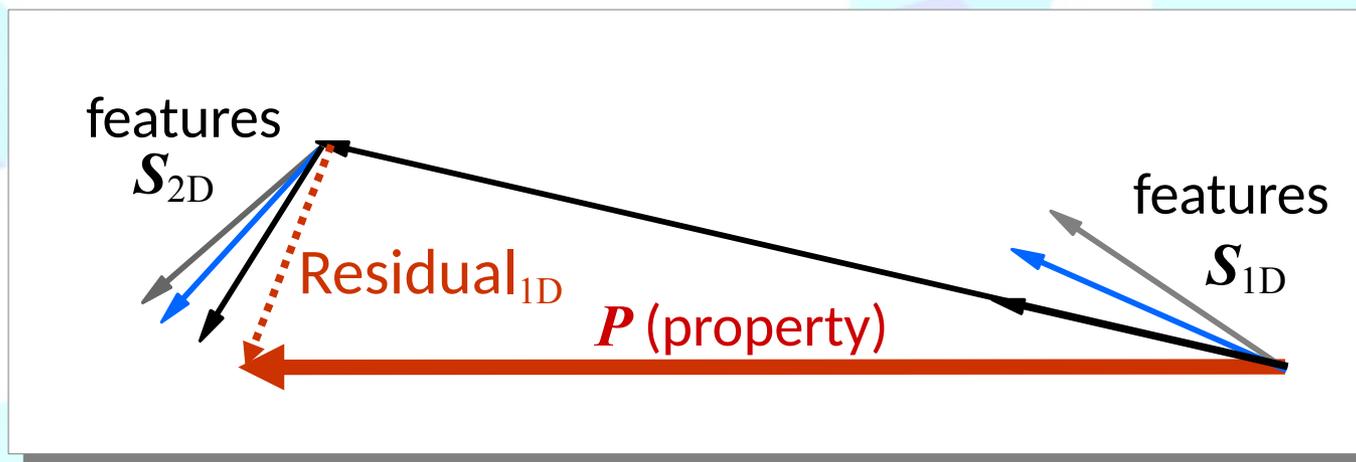
Huge “feature space” (e.g., 100 billion elements)
constructed from primary features (e.g., 15 - 25 elements)



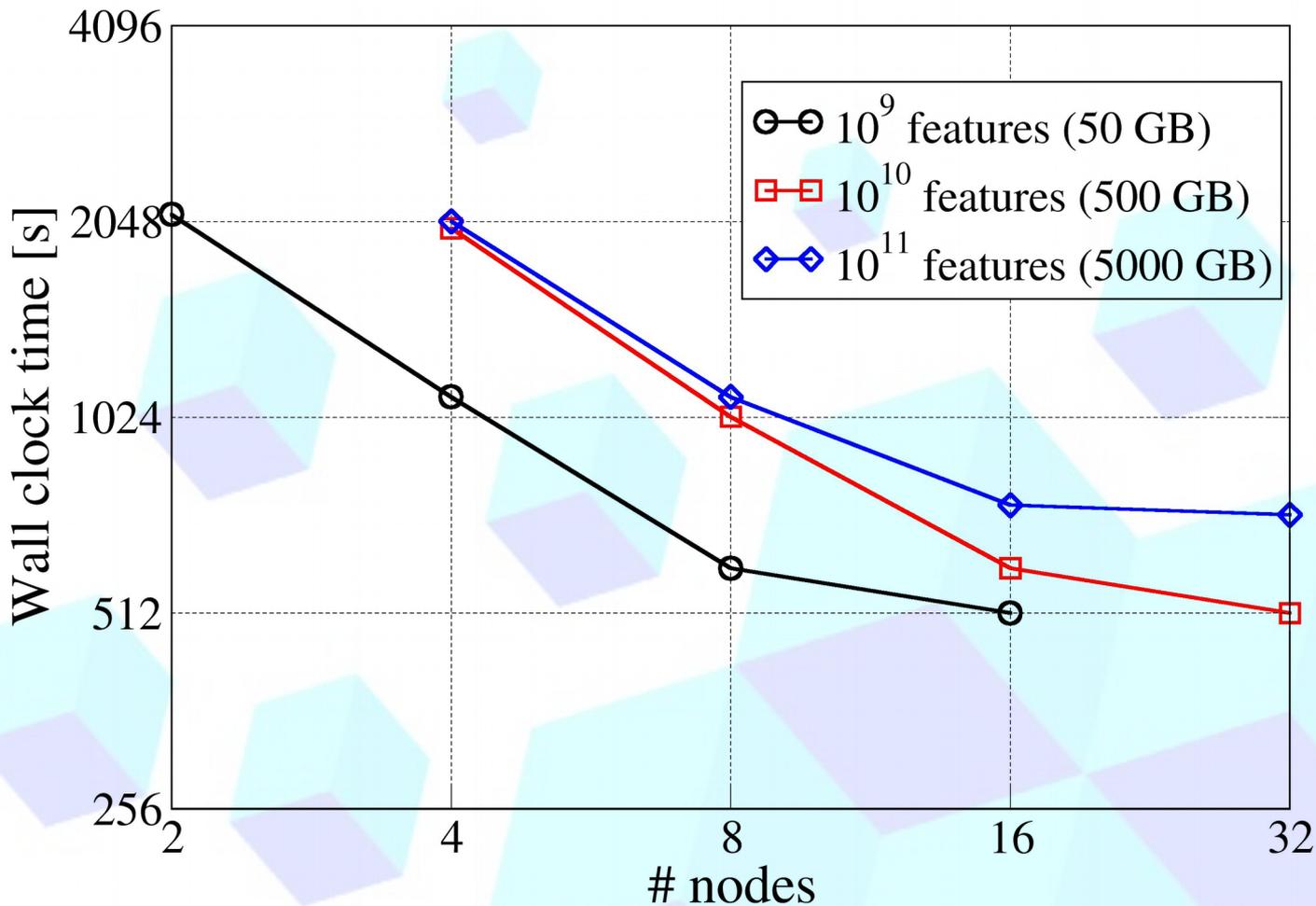
SISSO: sure independence screening plus sparsifying operator

Embarrassingly parallel

- + SIS: independent scalar products of features on property or residual)
- partial ranking
- + SO: independent least square regression
- partial ranking
- + outer parallelization for cross validation
- smart strategies needed for matrix storage

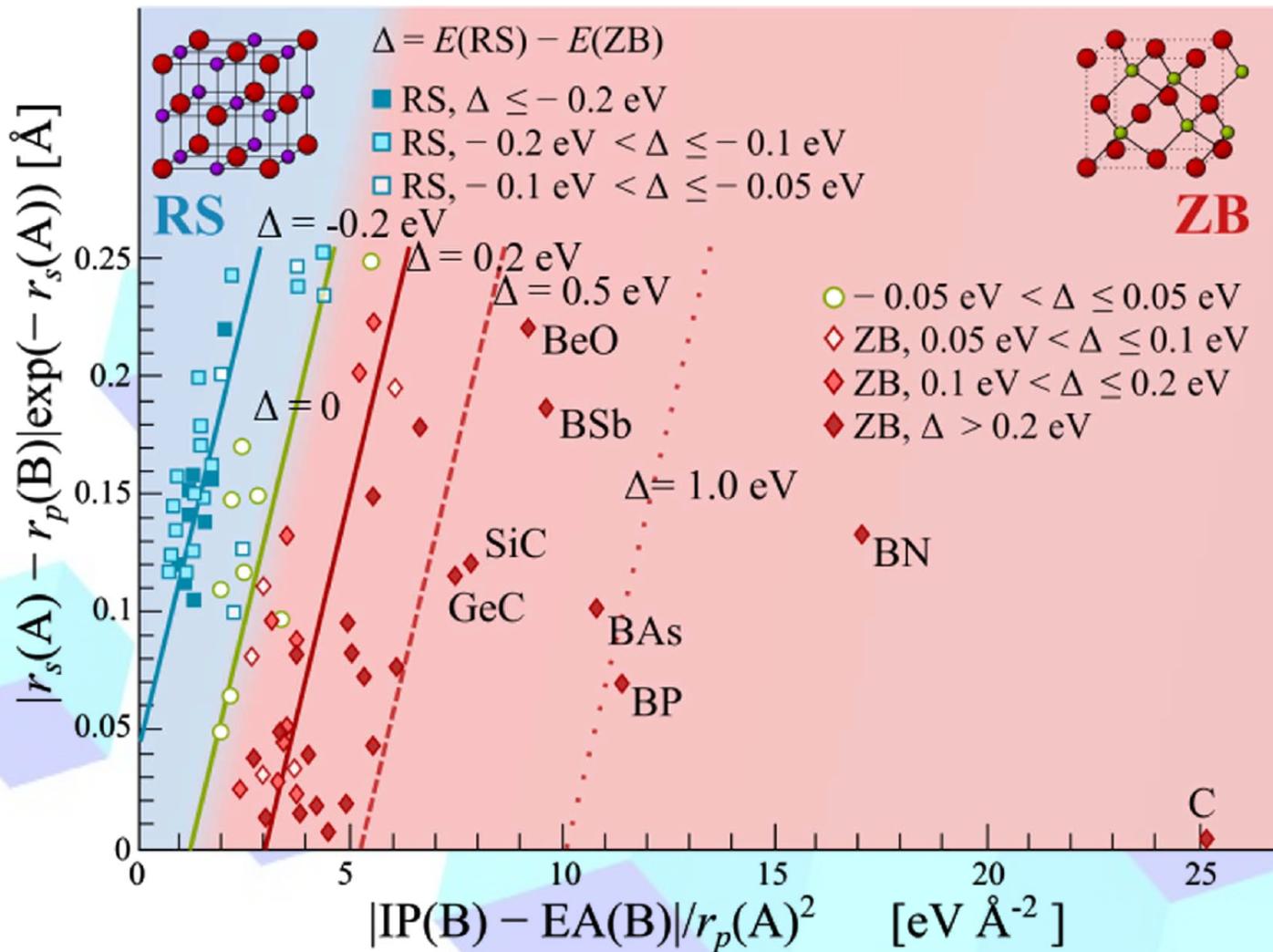


SISSO: sure independence screening plus sparsifying operator



Compressed sensing: the quest for descriptors and predictive models

Structure map with compressed-sensing algorithm, starting from 7 atomic features



Charts/maps of materials

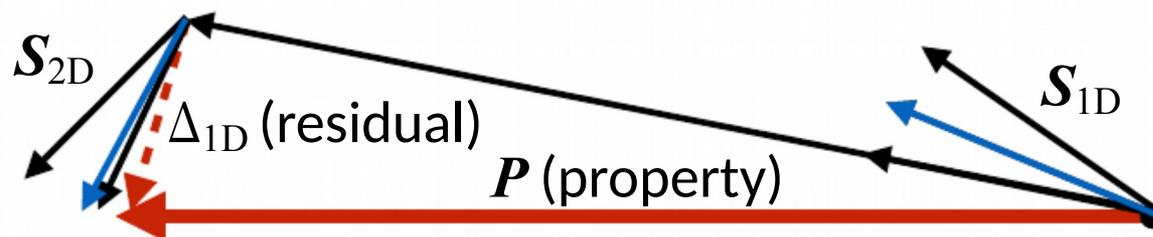
$$\operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^M} (\|\mathbf{P} - \mathbf{D}\mathbf{c}\|_2^2 + \lambda \|\mathbf{c}\|_0)$$

New cost function to be minimized:
overlap of convex domains

$$\hat{\mathbf{c}} \equiv \operatorname{arg\,min}_{\mathbf{c}} \left(\sum_{i=1}^{M-1} \sum_{j=i+1}^M O_{ij} + \lambda \|\mathbf{c}\|_0 \right)$$

Number of data points in
the overlap region, as
function of selected \mathbf{d}

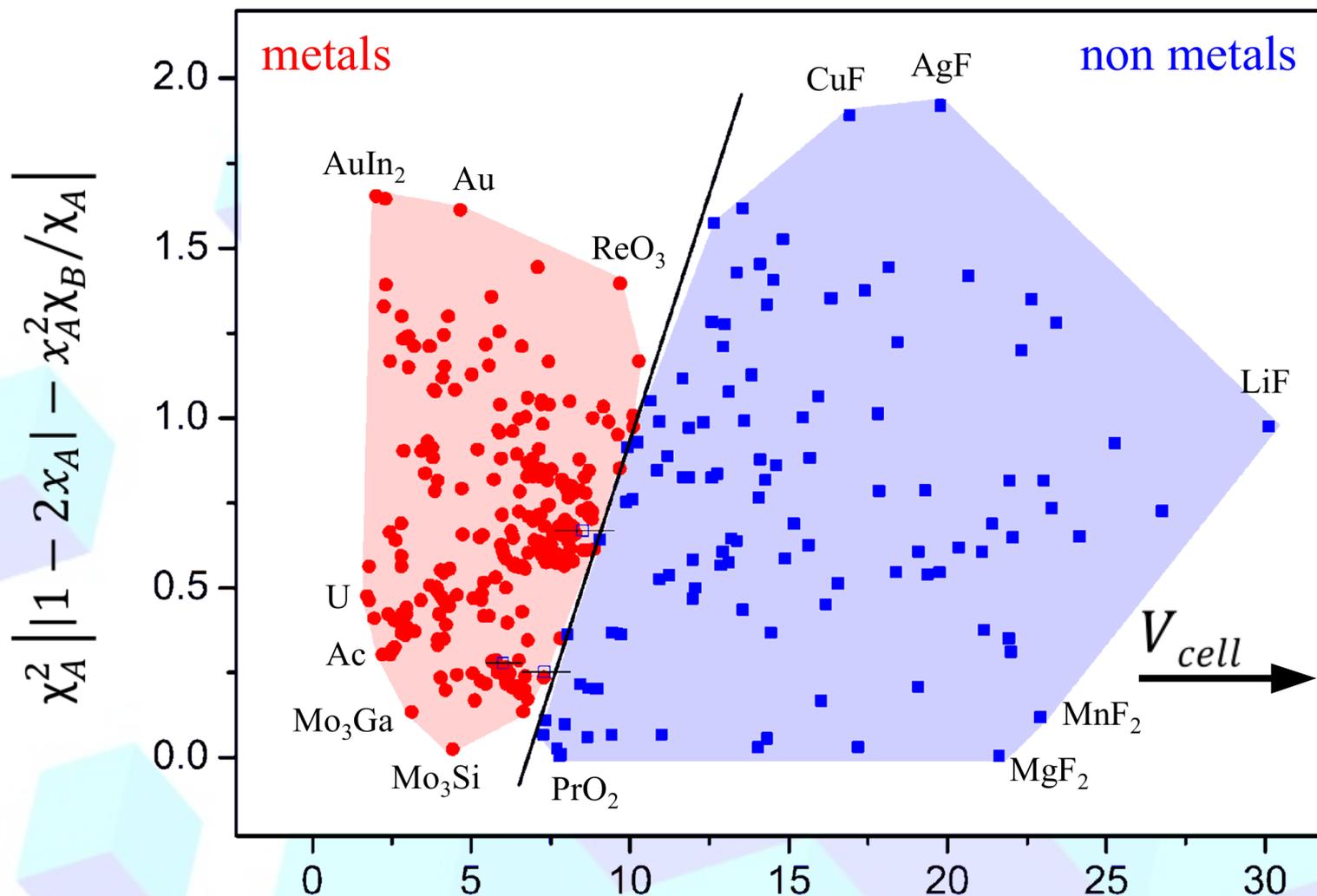
Iterative generation of feature subspaces



Topological
insulator

d''_1

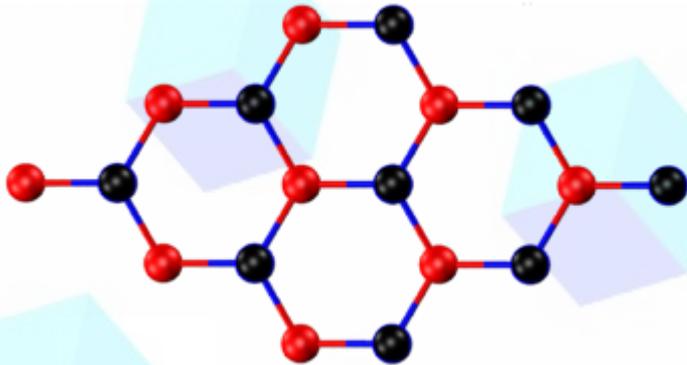
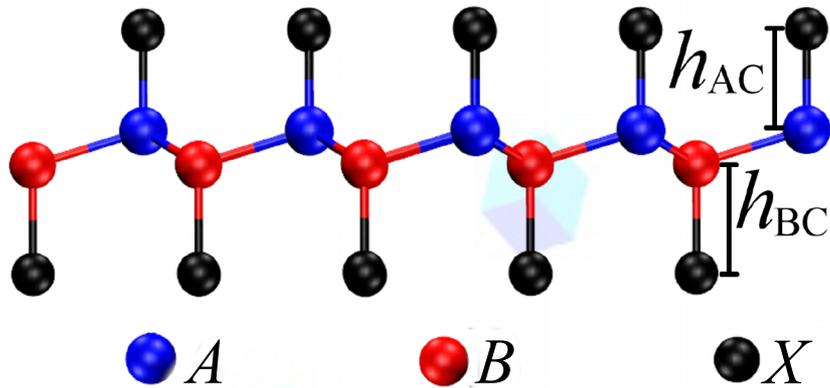
SISSO: metal/nonmetal classification of binary materials



x Atomic fraction
 IE Ionization energy
 χ Electronegativity

$$\frac{\sum V_{atom} / V_{cell} \cdot x_A}{\chi_A} \cdot \frac{IE_B \sqrt{\chi_B}}{\chi_A} \text{ (eV)}$$

SISSO: predicting novel honeycomb (~2D) topological insulators

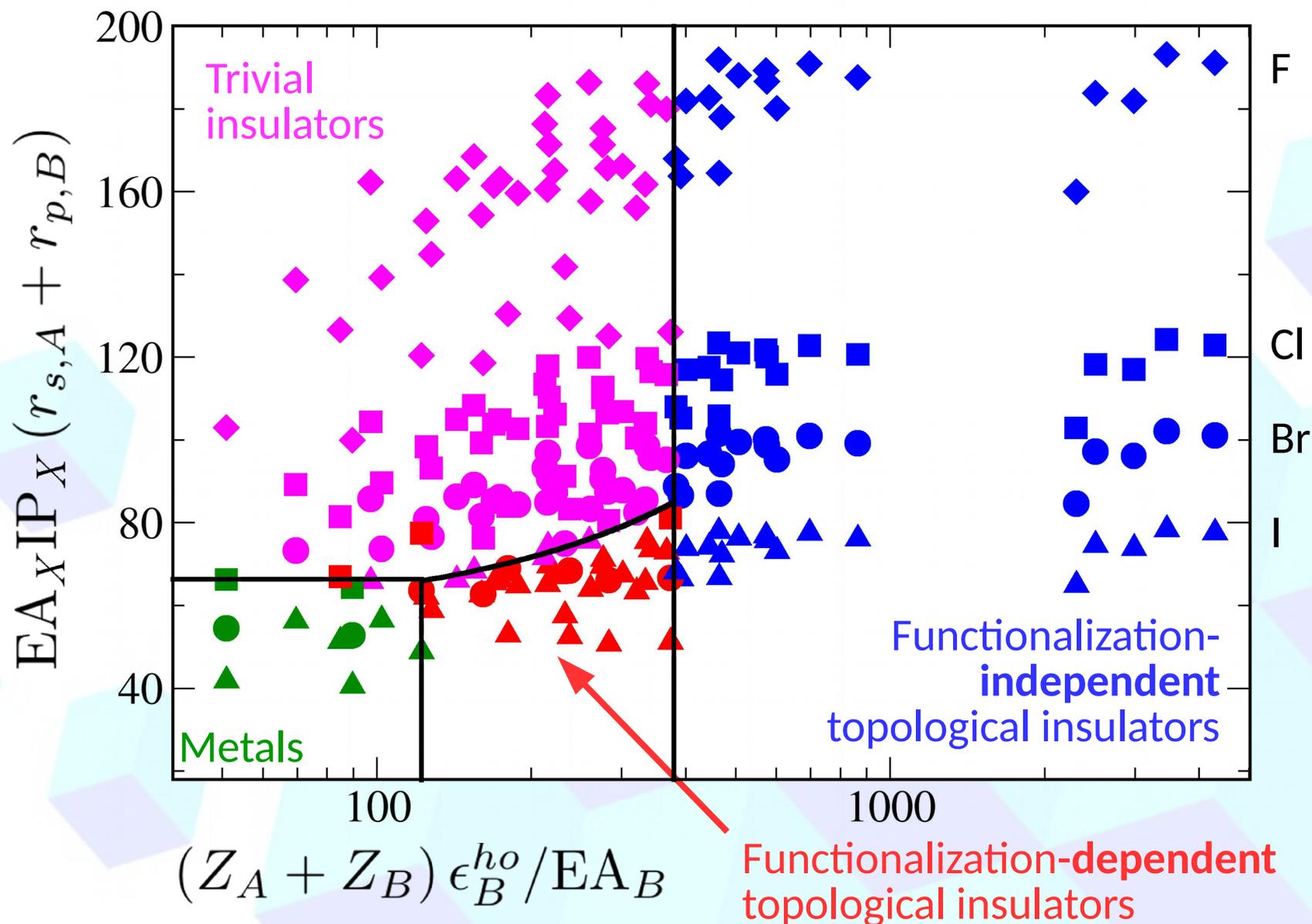


A II III IV
 B VI V IV
 X VII VII VII

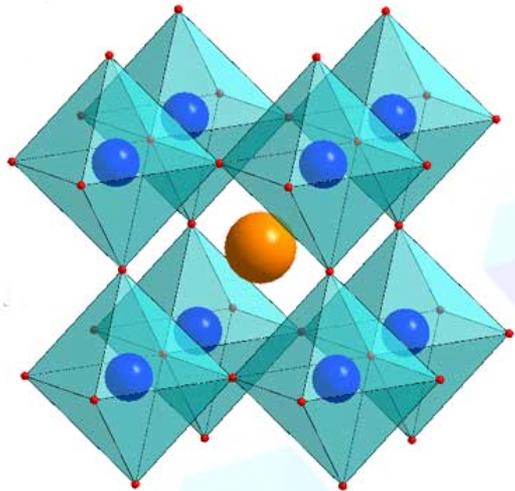
		II-VI						18 VIIA
		III-V				VII	2 4.0026	
		IIIA		IV-IV		VIA	VIIA	
		5 10.81	6 12.011	7 14.007	8 15.999	9 18.998	10 20.180	
		B	C	N	O	F	Ne	
		BORON	CARBON	NITROGEN	OXYGEN	FLUORINE	NEON	
		13 26.982	14 28.085	15 30.974	16 32.06	17 35.45	18 39.948	
		Al	Si	P	S	Cl	Ar	
		ALUMINIUM	SILICON	PHOSPHORUS	SULPHUR	CHLORINE	ARGON	
3.546	IB	30 65.38	31 69.723	32 72.64	33 74.922	34 78.971	35 79.904	36 83.798
u		Zn	Ga	Ge	As	Se	Br	Kr
		ZINC	GALLIUM	GERMANIUM	ARSENIC	SELENIUM	BROMINE	KRYPTON
07.87		48 112.41	49 114.82	50 118.71	51 121.76	52 127.60	53 126.90	54 131.29
g		Cd	In	Sn	Sb	Te	I	Xe
		CADMIUM	INDIUM	TIN	ANTIMONY	TELLURIUM	IODINE	XENON
96.97		80 200.59	81 204.38	82 207.2	83 208.98	84 (209)	85 (210)	86 (222)
u		Hg	Tl	Pb	Bi	Po	At	Rn
		MERCURY	THALLIUM	LEAD	BISMUTH	POLONIUM	ASTATINE	RADON
(280)		112 (285)	113 (...)	114 (287)	115 (...)	116 (291)	117 (...)	118 (...)
g		Cn	Uut	Fl	Uup	Lv	Uus	Uuo
		COPERNICIUM	UNUNTRIUM	FLEROVIUM	UNUNPENTIUM	LIVERMORIUM	UNUNSEPTIUM	UNUNOCTIUM

Data source: high throughput DFT (FHI-aims, Carlos Mera Acosta)

SISSO: predicting novel honeycomb (~2D) topological insulators



Perovskites' stability: Improving on Goldschmidt Tolerance Factor



ABX_3

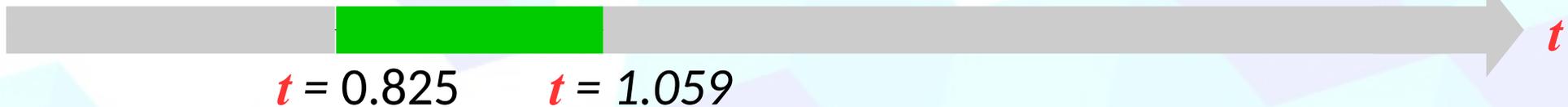
$$t = \frac{r_A + r_X}{\sqrt{2}(r_B + r_X)}$$

→ Ionic radius

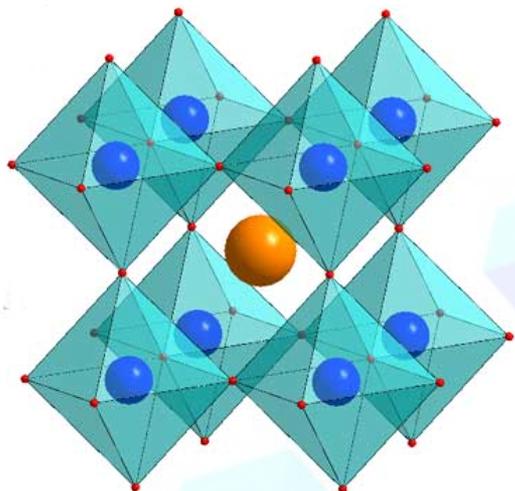
H																	He	
Li	Be											B	C	N	O	F	Ne	
Na	Mg											Al	Si	P	S	Cl	Ar	
K	Ca	Sc	Ti	V	Cr	Mn	Fe	Co	Ni	Cu	Zn	Ga	Ge	As	Se	Br	Kr	
Rb	Sr	Y	Zr	Nb	Mo	Tc	Ru	Rh	Pd	Ag	Cd	In	Sn	Sb	Te	I	Xe	
Cs	Ba			Hf	Ta	W	Re	Os	Ir	Pt	Au	Hg	Tl	Pb	Bi	Po	At	Rn
Fr	Ra																	
		La	Ce	Pr	Nd	Pm	Sm	Eu	Gd	Tb	Dy	Ho	Er	Tm	Yb	Lu		
		Ac	Th	Pa	U	Np	Pu	Am	Cm	Bk	Cf	Es	Fm	Md	No	Lr		

Goldschmidt* **stable** perovskites:

accuracy 74%



Perovskites' stability: Improving on Goldschmidt Tolerance Factor



$$t = \frac{r_A + r_X}{\sqrt{2}(r_B + r_X)}$$

Ionic radius

$$\tau = \frac{r_X}{r_B} - n_A \left(n_A - \frac{r_A/r_B}{\ln(r_A/r_B)} \right)$$

Oxidation state

$1/\mu = \text{Octahedral factor}$

Goldschmidt* **stable** perovskites:

accuracy 74%

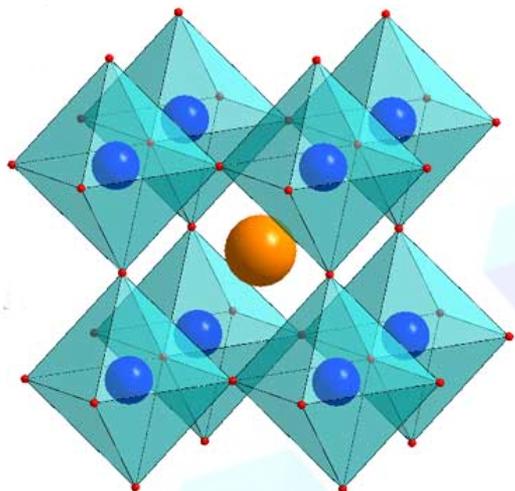


Our **stable** perovskites:

accuracy 92%



Perovskites' stability: Improving on Goldschmidt Tolerance Factor



$$t = \frac{r_A + r_X}{\sqrt{2}(r_B + r_X)}$$

Ionic radius

$$\tau = \frac{r_X}{r_B} - n_A \left(n_A - \frac{r_A/r_B}{\ln(r_A/r_B)} \right)$$

Oxidation state

$1/\mu = \text{Octahedral factor}$

Goldschmidt* **stable** perovskites:

accuracy 74%

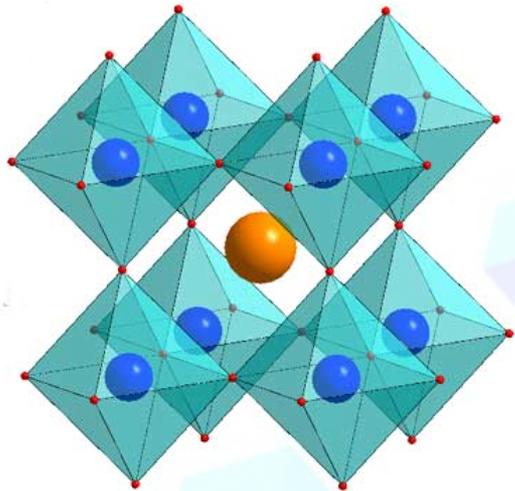


Our **stable** perovskites:

accuracy 99%



Perovskites' stability: Improving on Goldschmidt Tolerance Factor



ABX_3

$$t = \frac{r_A + r_X}{\sqrt{2}(r_B + r_X)}$$

Ionic radius

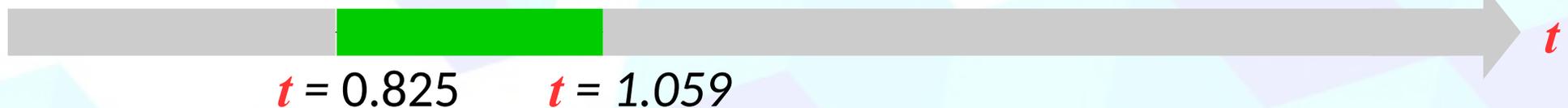
$$\tau = \frac{r_X}{r_B} - n_A \left(n_A - \frac{r_A/r_B}{\ln(r_A/r_B)} \right)$$

Oxidation state

$1/\mu = \text{Octahedral factor}$

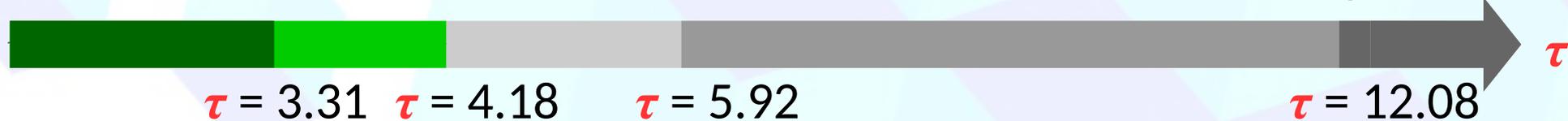
Goldschmidt* **stable** perovskites:

accuracy 74%



Our **stable** perovskites:

accuracy 100%



... and more

Continuous property

- Adsorption energy of O on metal-oxide surfaces
- Adsorption energy and OCO angle of adsorbed CO₂ on metal-oxide surfaces
- Adsorption energy of metal atoms on metal-alloys surfaces

Features: atoms (of the surface) and pristine surface

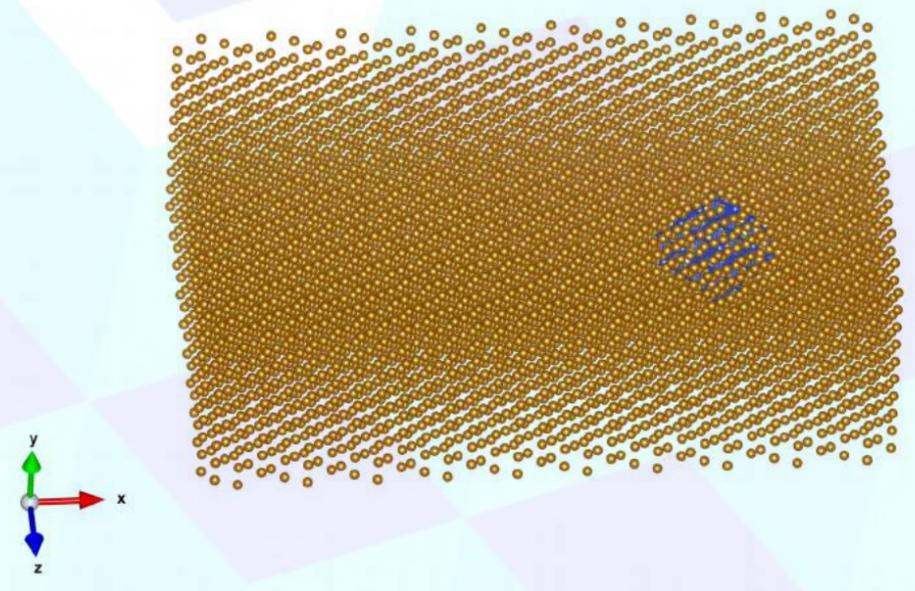
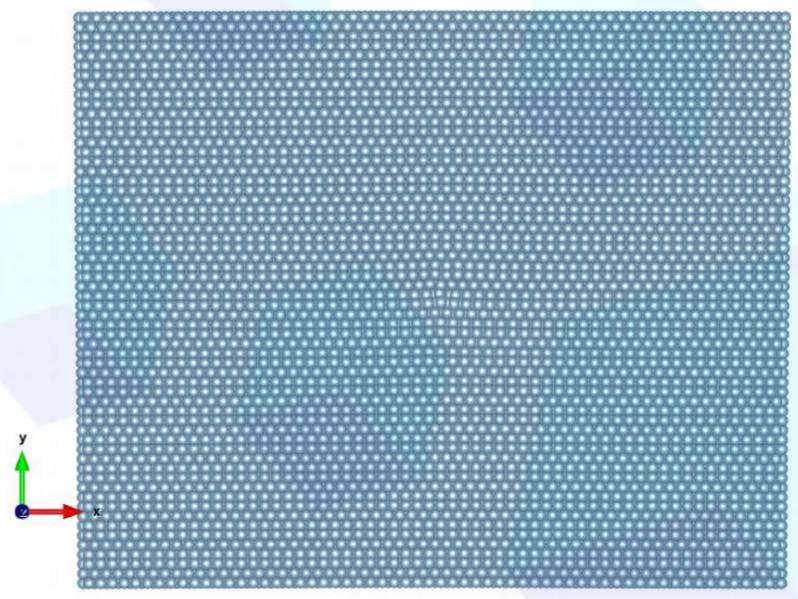
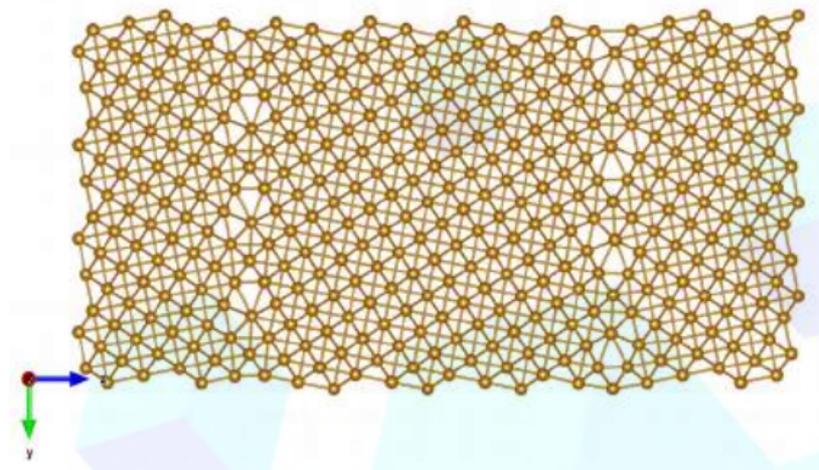
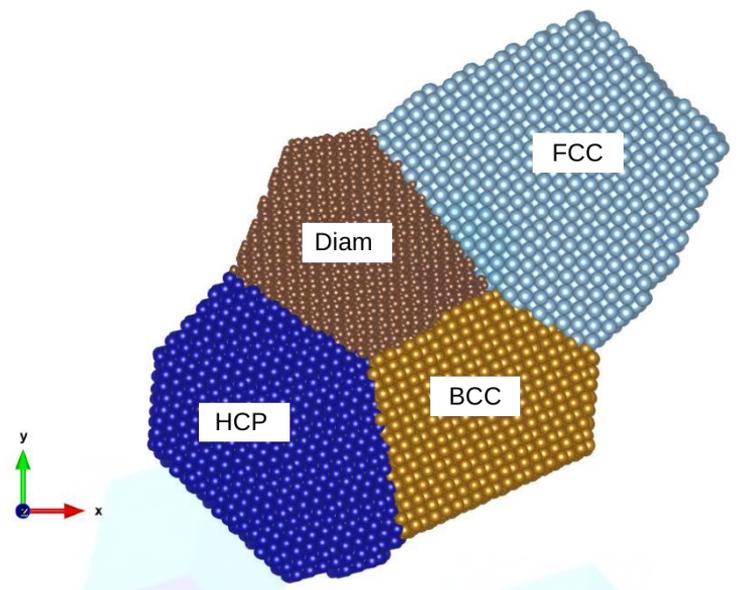
Classification

- Tetradymite 5-component 3d topological insulators (vs trivial insulators) arXiv:1808.04733

Features: atoms

... and further more

Convolutional neural networks for (local) crystal-structure recognition



SISSO and metal/insulator proof of concepts

Runhai Ouyang, Emre Ahmetcik, Stefano Curtarolo

Topological Insulators

Carlos Mera Acosta, Adalberto Fazzio, Runhai Ouyang, Christian Carbogno

Perovskites

Christopher J. Bartel, Christopher Sutton, Bryan R. Goldsmith, Runhai Ouyang, Charles B. Musgrave

Convolutional neural networks for structure recognition

Angelo Ziletti, Andreas Leitherer, Devinder Kumar

NOMAD PIs

Angel Rubio (MPSD Hamburg), Risto Nieminen (Aalto Univ. Helsinki), Francesc Illas (Univ. Barcelona), Daan Frenkel (Univ. Cambridge), Claudia Draxl (HU Berlin), Alessandro De Vita (Kings College London), Kristian Thygesen (DTU Lyngby); Kimmo Koski (CSC Helsinki), Stefan Heinzl (MPSCD Garching), Jose Maria Cela (BSC Barcelona), Dieter Kranzlmüller (LRZ Munich); Ciaran Clissman (Pintail Dublin)

All the above

Matthias Scheffler