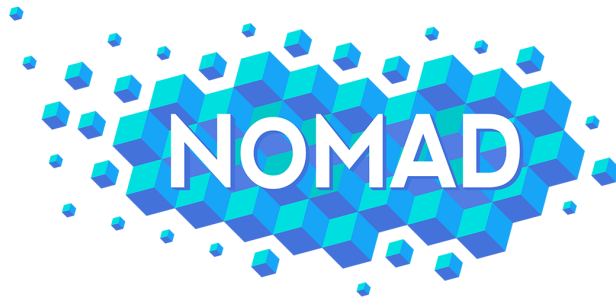


BARCELONA  
2019

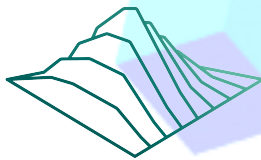


NOVEL MATERIALS DISCOVERY



MAX-PLANCK-GESELLSCHAFT

# (Big-)data analytics for materials science



Luca M. Ghiringhelli

FRITZ-HABER-INSTITUT  
MAX-PLANCK-GESELLSCHAFT

Hands-On DFT and Beyond:

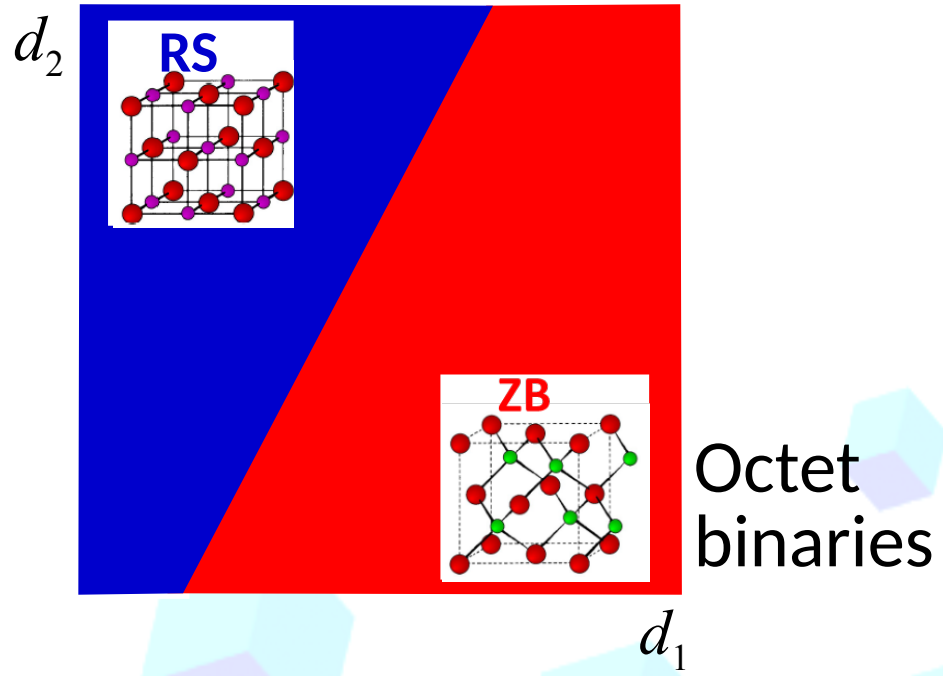
High-throughput screening and big-data analytics,  
towards exascale computational materials science

University of Barcelona, Spain, August 26th to September 6th, 2019

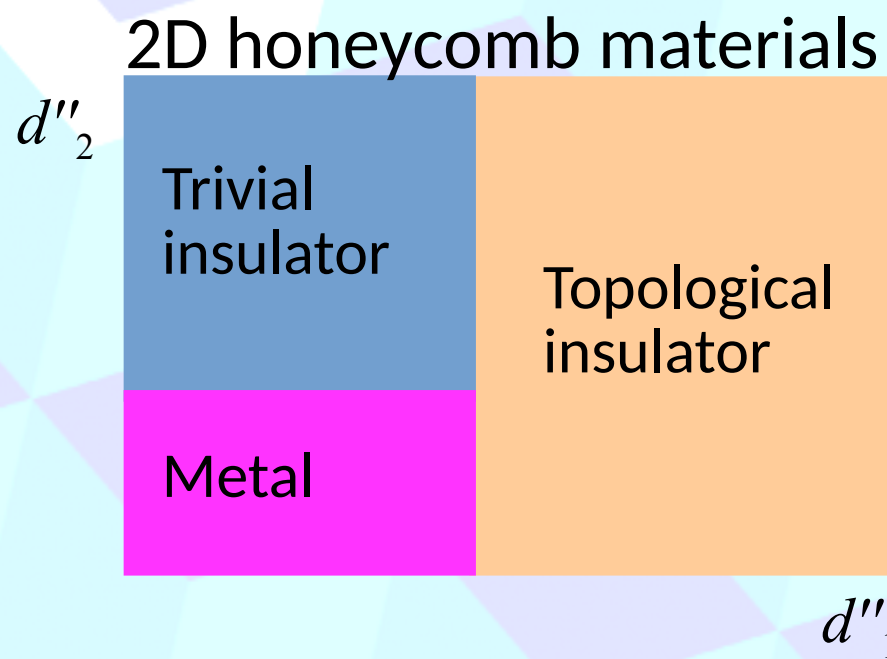
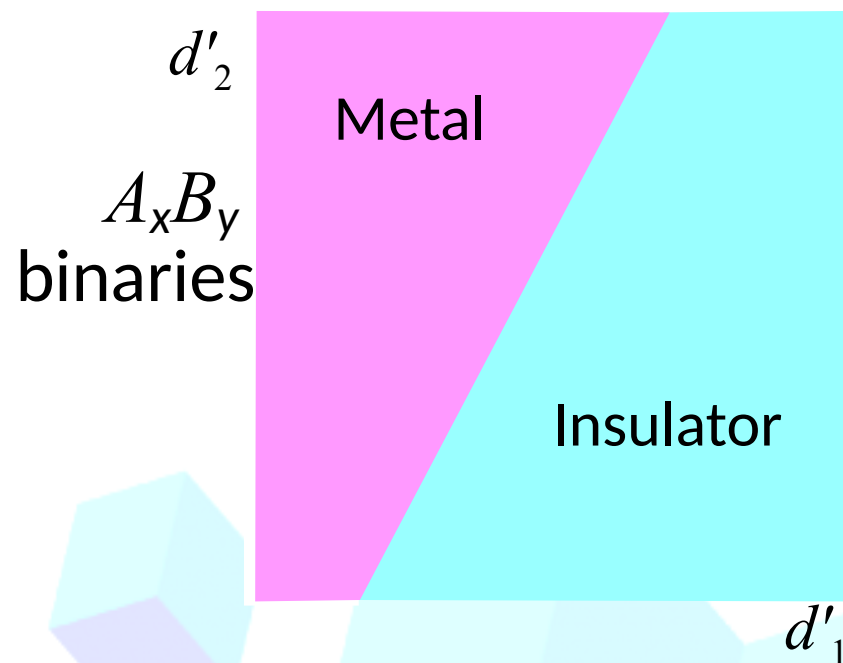
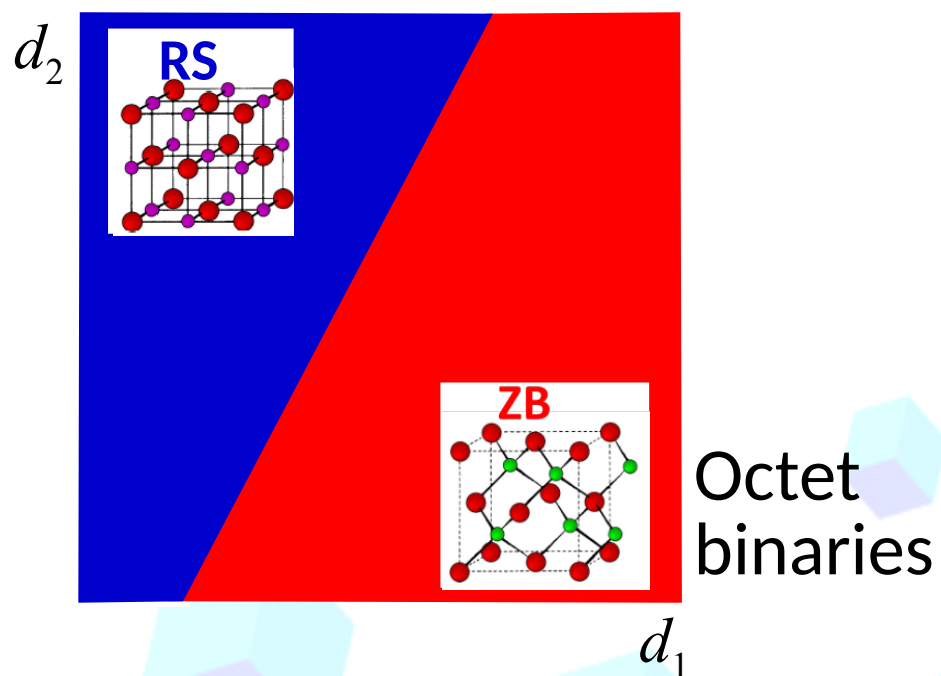
# Charts/maps of materials



# Charts/maps of materials

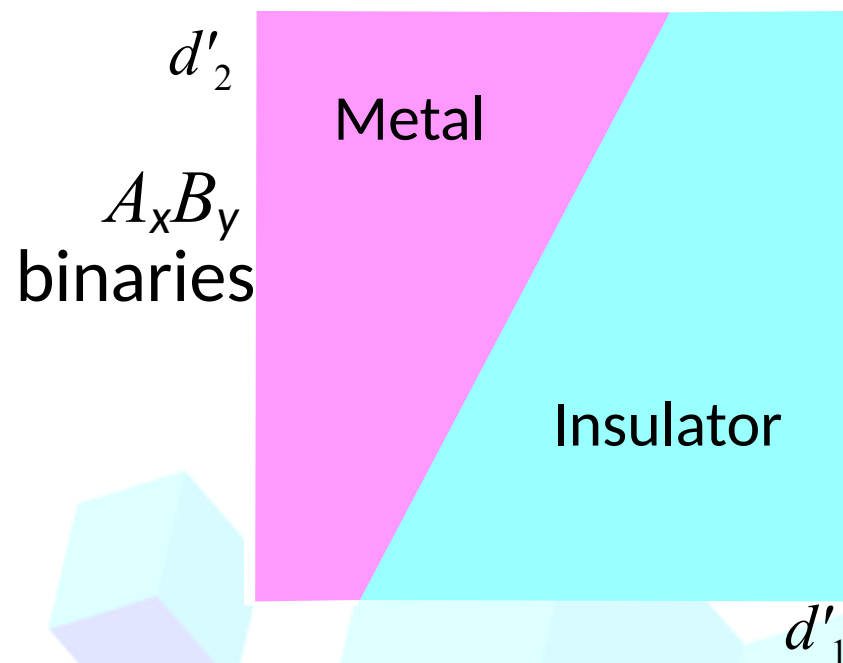
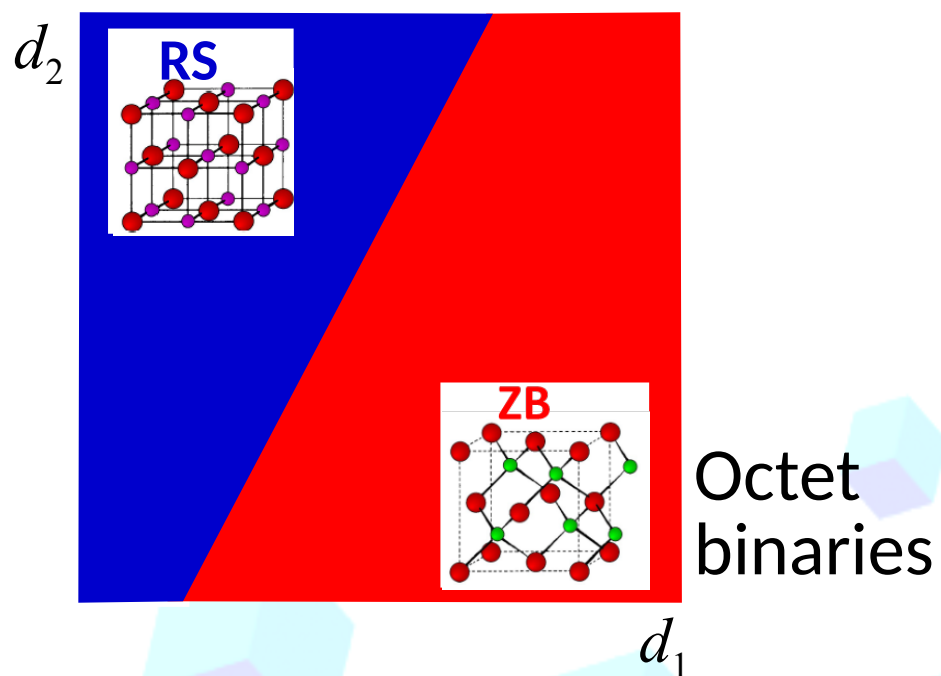


# Charts/maps of materials

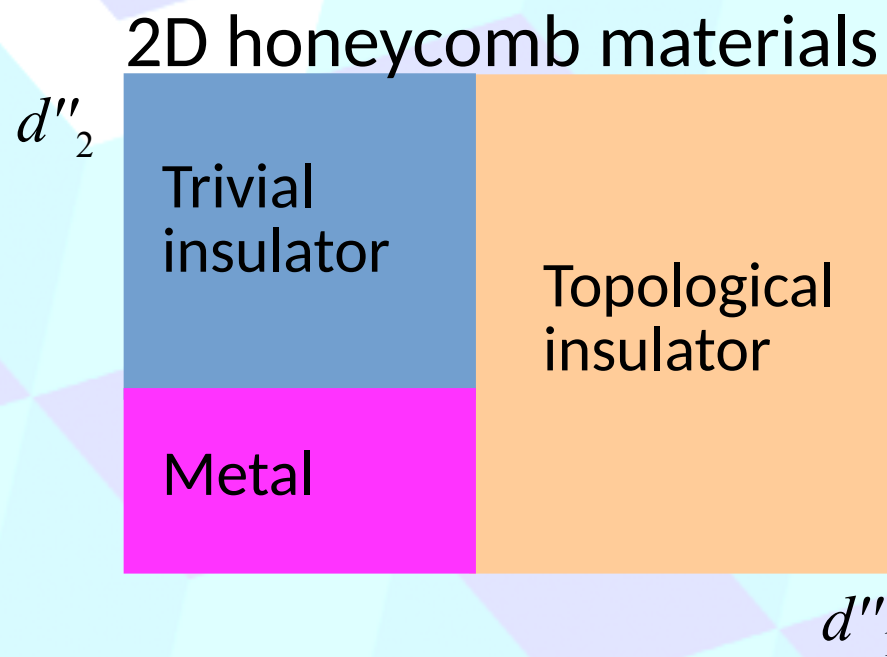
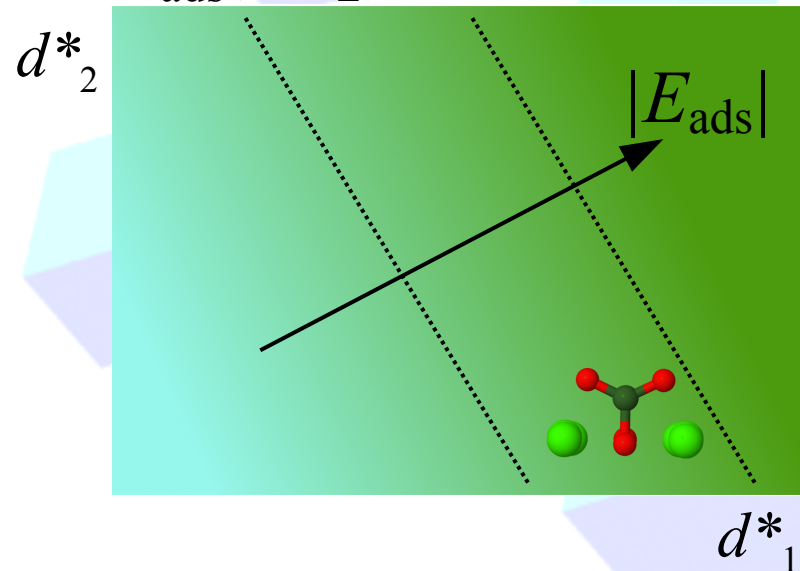




# Charts/maps of materials



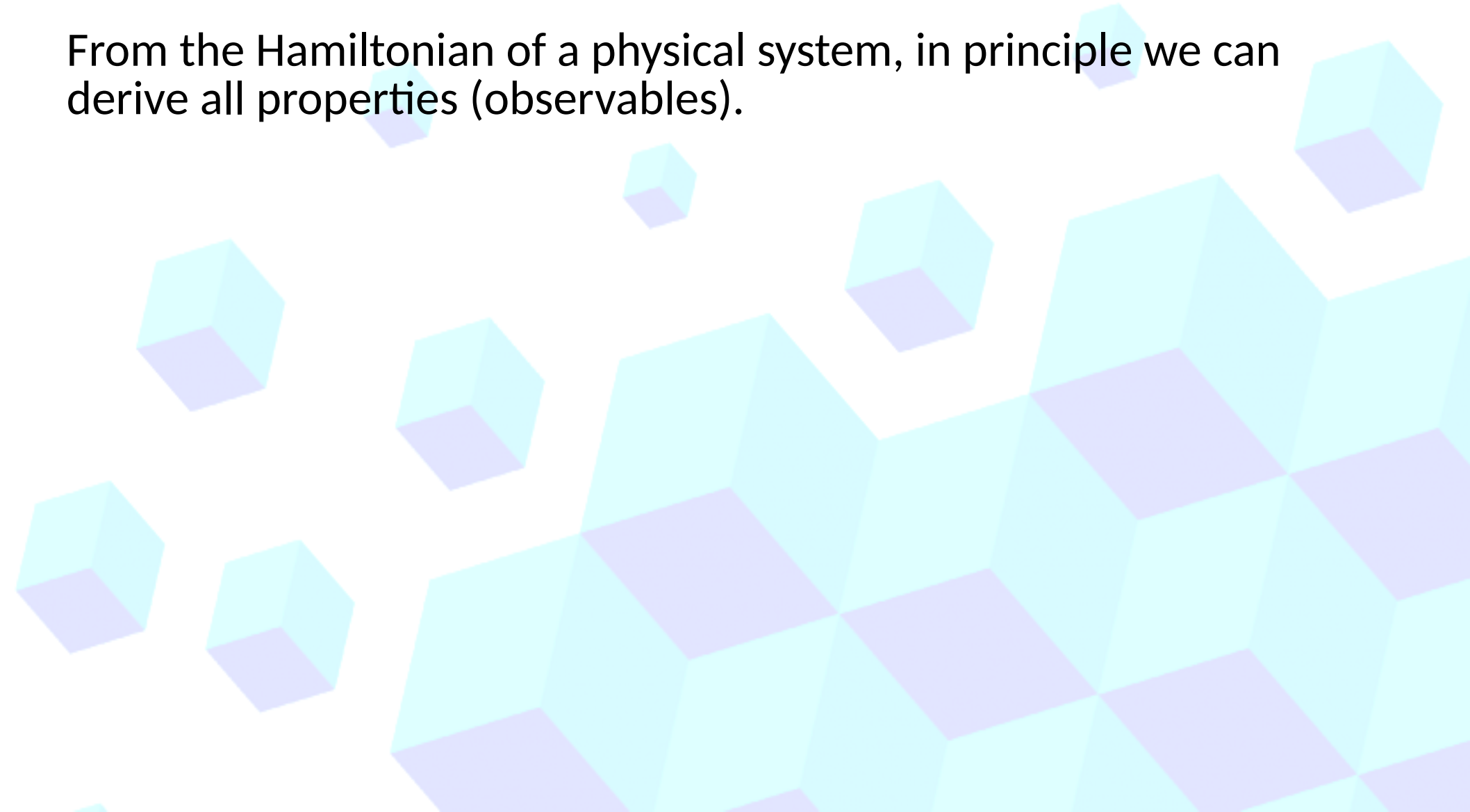
$E_{\text{ads}}(\text{CO}_2)$  on oxides



# Building maps of materials properties

## A quantum many-body problem

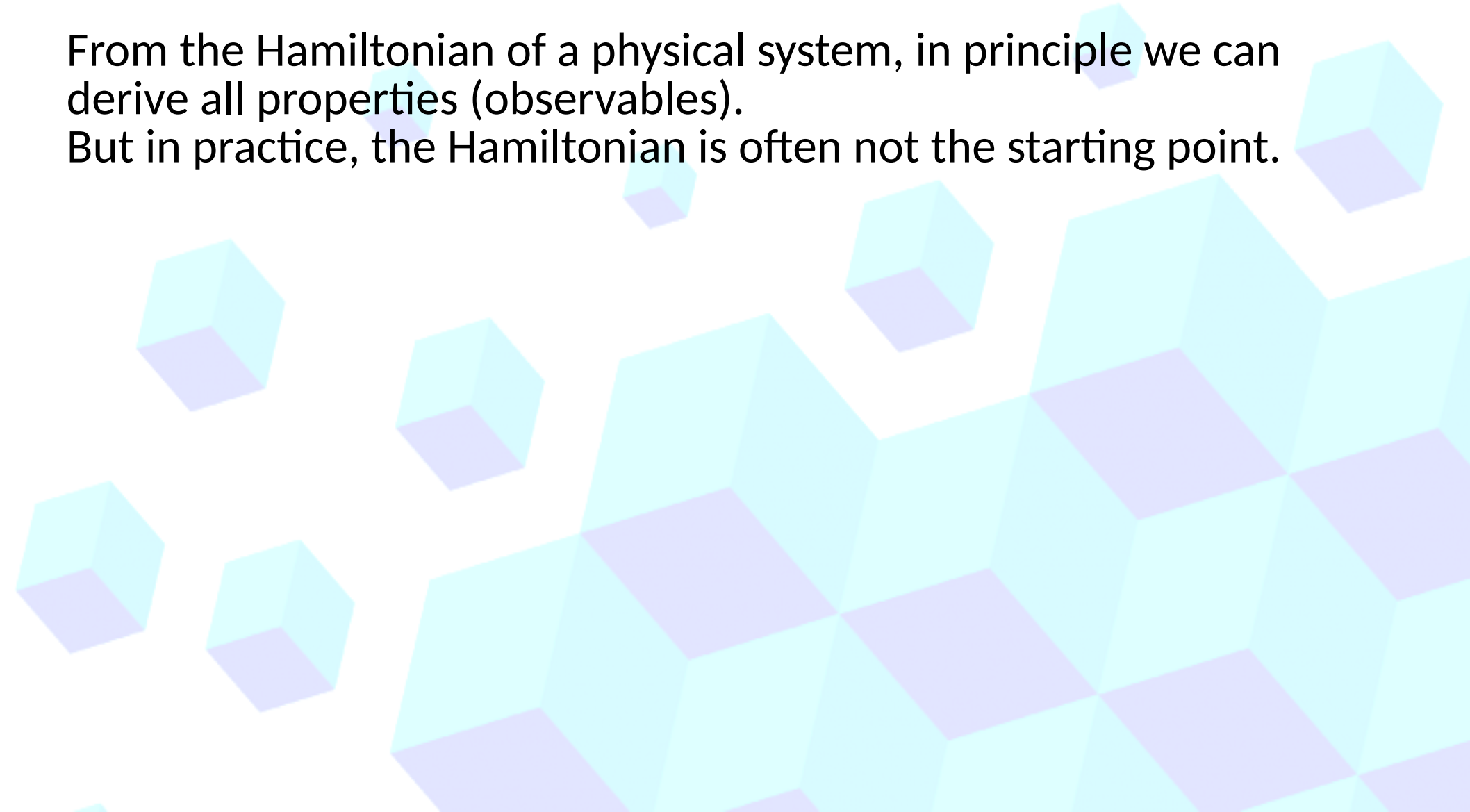
From the Hamiltonian of a physical system, in principle we can derive all properties (observables).



# Building maps of materials properties

## A quantum many-body problem

From the Hamiltonian of a physical system, in principle we can derive all properties (observables).  
But in practice, the Hamiltonian is often not the starting point.



# Building maps of materials properties

## A quantum many-body problem

From the Hamiltonian of a physical system, in principle we can derive all properties (observables).

But in practice, the Hamiltonian is often not the starting point.

For instance, given a class of chemical compositions (e.g., via prototype formula, such as  $ABX_3$ ):

# Building maps of materials properties

## A quantum many-body problem

From the Hamiltonian of a physical system, in principle we can derive all properties (observables).

But in practice, the Hamiltonian is often not the starting point.

For instance, given a class of chemical compositions (e.g., via prototype formula, such as  $ABX_3$ ):

- what is the most stable crystal structure of each material in the class?

# Building maps of materials properties

## A quantum many-body problem

From the Hamiltonian of a physical system, in principle we can derive all properties (observables).

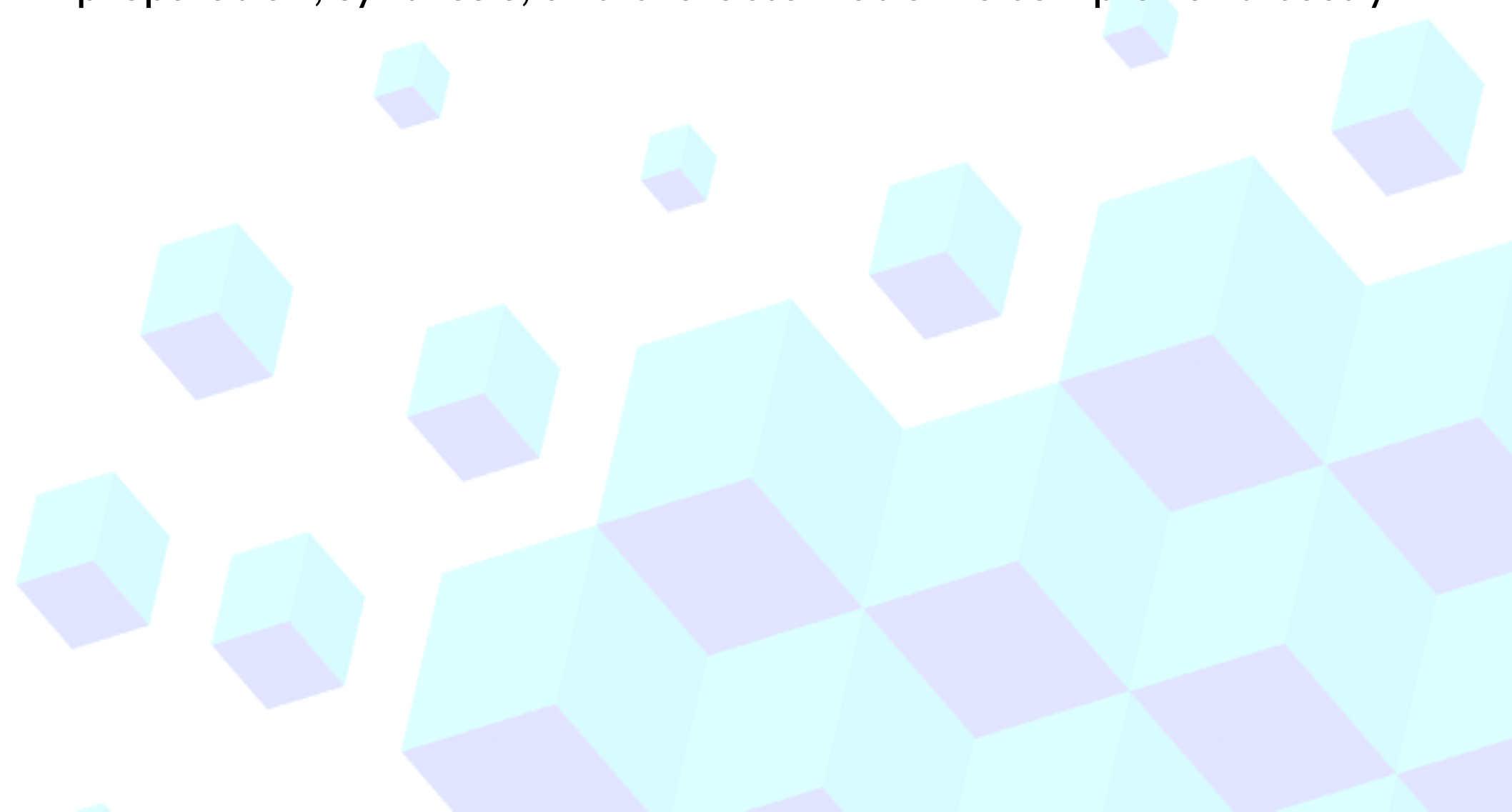
But in practice, the Hamiltonian is often not the starting point.

For instance, given a class of chemical compositions (e.g., via prototype formula, such as  $ABX_3$ ):

- what is the most stable crystal structure of each material in the class?
- which materials are metals / topological insulators / superconductors ?
- which material has the highest melting point?
- which materials has a surface optimal for catalysing some chemical reaction?

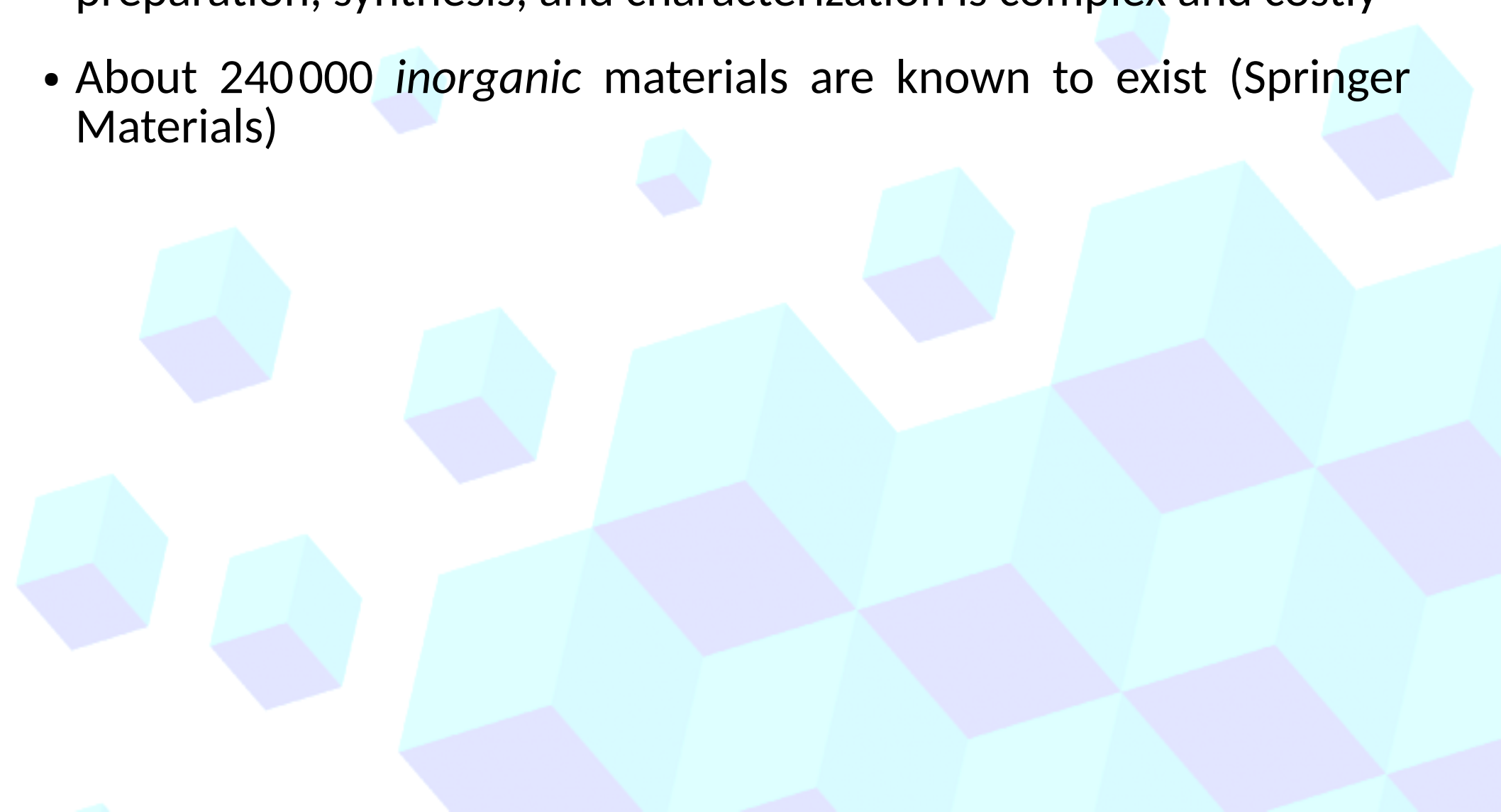
# The Big Picture

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



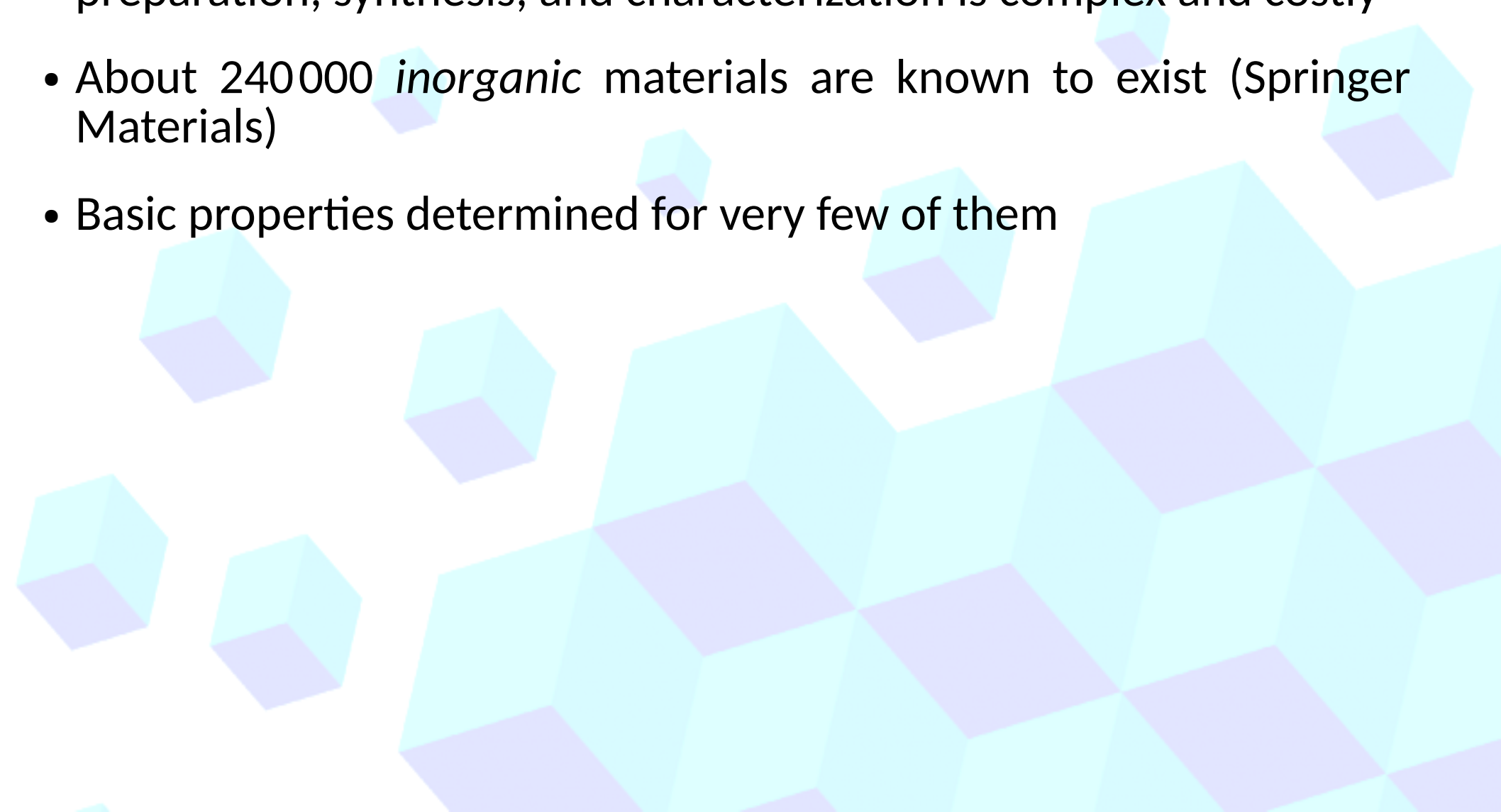
# The Big Picture

- Design of new materials: preparation, synthesis, and characterization is complex and costly
- About 240000 *inorganic* materials are known to exist (Springer Materials)

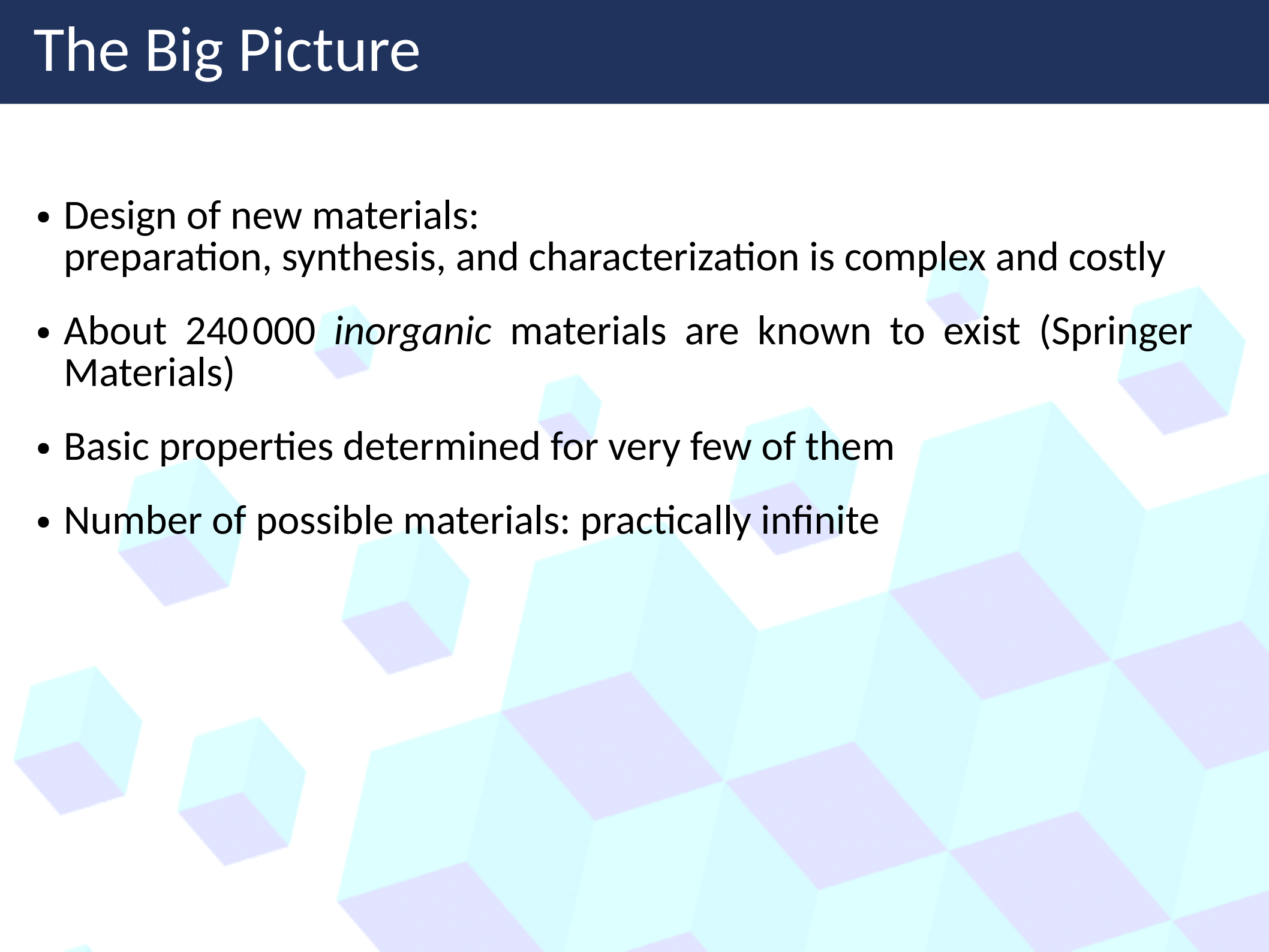




# The Big Picture

- Design of new materials: preparation, synthesis, and characterization is complex and costly
  - About 240000 *inorganic* materials are known to exist (Springer Materials)
  - Basic properties determined for very few of them
- 

# The Big Picture

- Design of new materials: preparation, synthesis, and characterization is complex and costly
  - About 240000 *inorganic* materials are known to exist (Springer Materials)
  - Basic properties determined for very few of them
  - Number of possible materials: practically infinite
- 

# The Big Picture

- Design of new materials: preparation, synthesis, and characterization is complex and costly
  - About 240000 *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

# The Big Picture

- Design of new materials: preparation, synthesis, and characterization is complex and costly
  - About 240000 *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

# We have a dream

From the **periodic table of the elements**  
to **charts of materials**



# We have a dream

From the **periodic table of the elements**  
to **charts of materials**

Reihen	Gruppe I. — R <sup>2</sup> O	Gruppe II. — RO	Gruppe III. — R <sup>2</sup> O <sup>3</sup>	Gruppe IV. RH <sup>4</sup> RO <sup>2</sup>	Gruppe V. RH <sup>3</sup> R <sup>2</sup> O <sup>5</sup>	Gruppe VI. RH <sup>2</sup> RO <sup>3</sup>	Gruppe VII. RH R <sup>2</sup> O <sup>7</sup>	Gruppe VIII. — RO <sup>4</sup>
1	H=1							
2	Li=7	Be=9.4	B=11	C=12	N=14	O=16	F=19	
3	Na=23	Mg=24	Al=27.3	Si=28	P=31	S=32	Cl=35.5	
4	K=39	Ca=40	—=44	Ti=48	V=51	Cr=52	Mn=55	Fe=56, Co=59, Ni=59, Cu=63.
5	(Cu=63)	Zn=65	—=68	—=72	As=75	Se=78	Br=80	
6	Rb=85	Sr=87	?Yt=88	Zr=90	Nb=94	Mo=96	—=100	Ru=104, Rh=104, Pd=106, Ag=108.
7	(Ag=108)	Cd=112	In=113	Sn=118	Sb=122	Te=125	J=127	
8	Cs=133	Ba=137	?Di=138	?Ce=140	—	—	—	— — — —
9	(—)	—	—	—	—	—	—	
10	—	—	?Er=178	?La=180	Ta=182	W=184	—	Os=195, Ir=197, Pt=198, Au=199.
11	(Au=199)	Hg=200	Tl=204	Pb=207	Bi=208	—	—	
12	—	—	—	Th=231	—	U=240	—	— — — —

Mendeleev's 1871 periodic table

# We have a dream

From the **periodic table of the elements**  
to **charts of materials**

Reihen	Gruppe I. — R <sup>2</sup> O	Gruppe II. — RO	Gruppe III. — R <sup>2</sup> O <sup>3</sup>	Gruppe IV. RH <sup>4</sup> RO <sup>2</sup>	Gruppe V. RH <sup>3</sup> R <sup>2</sup> O <sup>5</sup>	Gruppe VI. RH <sup>2</sup> RO <sup>3</sup>	Gruppe VII. RH R <sup>2</sup> O <sup>7</sup>	Gruppe VIII. — RO <sup>4</sup>
1	H=1							
2	Li=7	Be=9.4	B=11	C=12	N=14	O=16	F=19	
3	Na=23	Mg=24	Al=27.3	Si=28	P=31	S=32	Cl=35.5	
4	K=39	Ca=40	—=44	Ti=48	V=51	Cr=52	Mn=55	Fe=56, Co=59, Ni=59, Cu=63.
5	(Cu=63)	Zn=65	—=68	—=72	As=75	Se=78	Br=80	
6	Rb=85	Sr=87	?Yt=88	Zr=90	Nb=94	Mo=96	—=100	Ru=104, Rh=104, Pd=106, Ag=108.
7	(Ag=108)	Cd=112	In=113	Sn=118	Sb=122	Te=125	J=127	
8	Cs=133	Ba=137	?Di=138	?Ce=140	—	—	—	— — — —
9	(—)	—	—	—	—	—	—	
10	—	—	?Er=178	?La=180	Ta=182	W=184	—	Os=195, Ir=197, Pt=198, Au=199.
11	(Au=199)	Hg=200	Tl=204	Pb=207	Bi=208	—	—	
12	—	—	—	Th=231	—	U=240	—	— — — —

Mendeleev's 1871 periodic table

# We have a dream

From the **periodic table of the elements**  
to **charts of materials**

Reihen	Gruppe I. — R <sup>2</sup> O	Gruppe II. — RO	Gruppe III. — R <sup>2</sup> O <sup>3</sup>	Gruppe IV. RH <sup>4</sup> RO <sup>2</sup>	Gruppe V. RH <sup>3</sup> R <sup>2</sup> O <sup>5</sup>	Gruppe VI. RH <sup>2</sup> RO <sup>3</sup>	Gruppe VII. RH R <sup>2</sup> O <sup>7</sup>	Gruppe VIII. — RO <sup>4</sup>
1	H=1							
2	Li=7	Be=9.4	B=11	C=12	N=14	O=16	F=19	
3	Na=23	Mg=24	Al=27.3	Si=28	P=31	S=32	Cl=35.5	
4	K=39	Ca=40	—=44	Ti=48	V=51	Cr=52	Mn=55	Fe=56, Co=59, Ni=59, Cu=63.
5	(Cu=63)	Zn=65	—=68	—=72	As=75	Se=78	Br=80	
6	Rb=85	Sr=87	?Yt=88	Zr=90	Nb=94	Mo=96	—=100	Ru=104, Rh=104, Pd=106, Ag=108.
7	(Ag=108)	Cd=112	In=113	Sn=118	Sb=122	Te=125	J=127	
8	Cs=133	Ba=137	?Di=138	?Ce=140	—	—	—	— — — —
9	(—)	—	—	—	—	—	—	
10	—	—	?Er=178	?La=180	Ta=182	W=184	—	Os=195, Ir=197, Pt=198, Au=199.
11	(Au=199)	Hg=200	Tl=204	Pb=207	Bi=208	—	—	
12	—	—	—	Th=231	—	U=240	—	— — — —

Mendeleev's 1871 periodic table

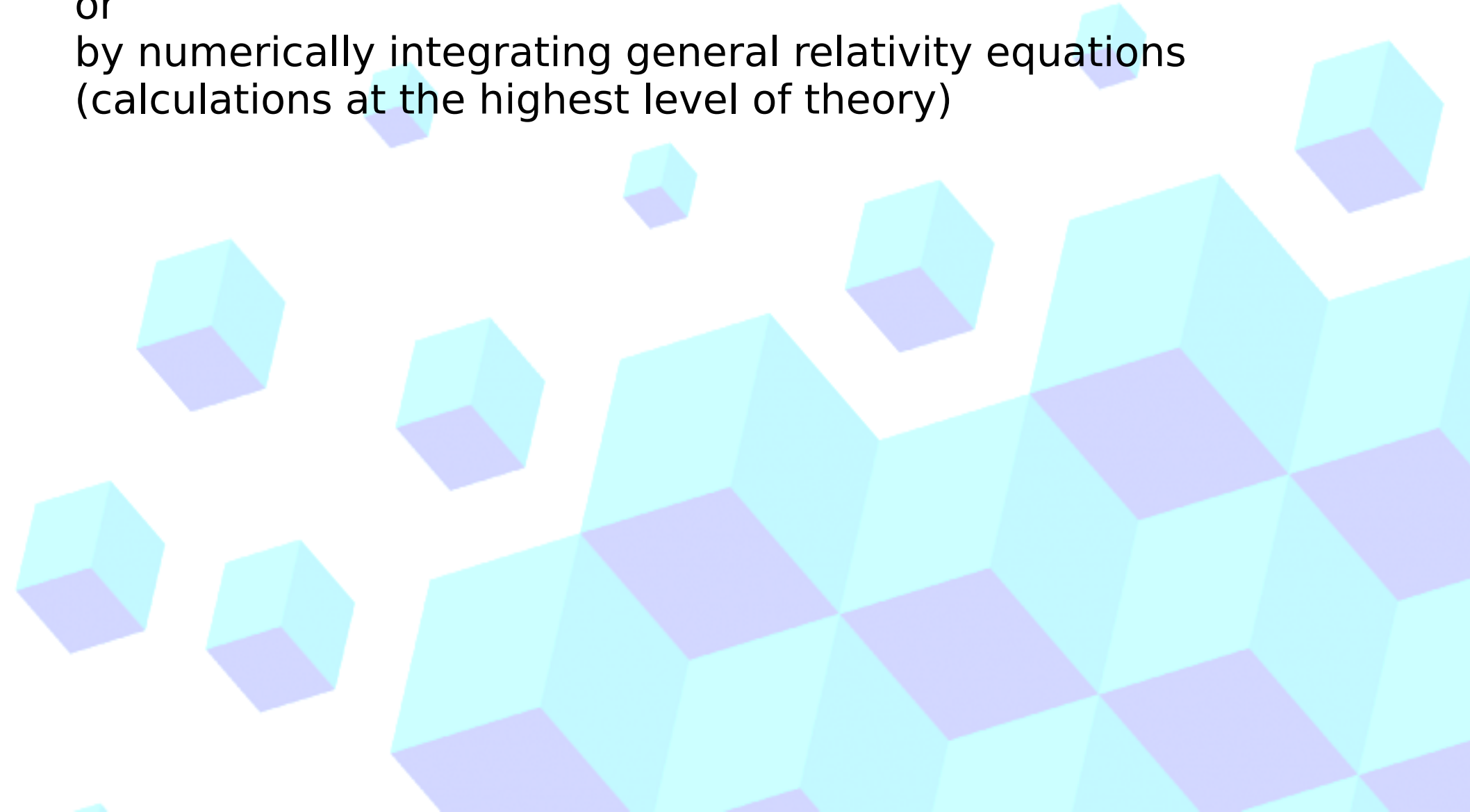


# Learning → Discovery

Suppose

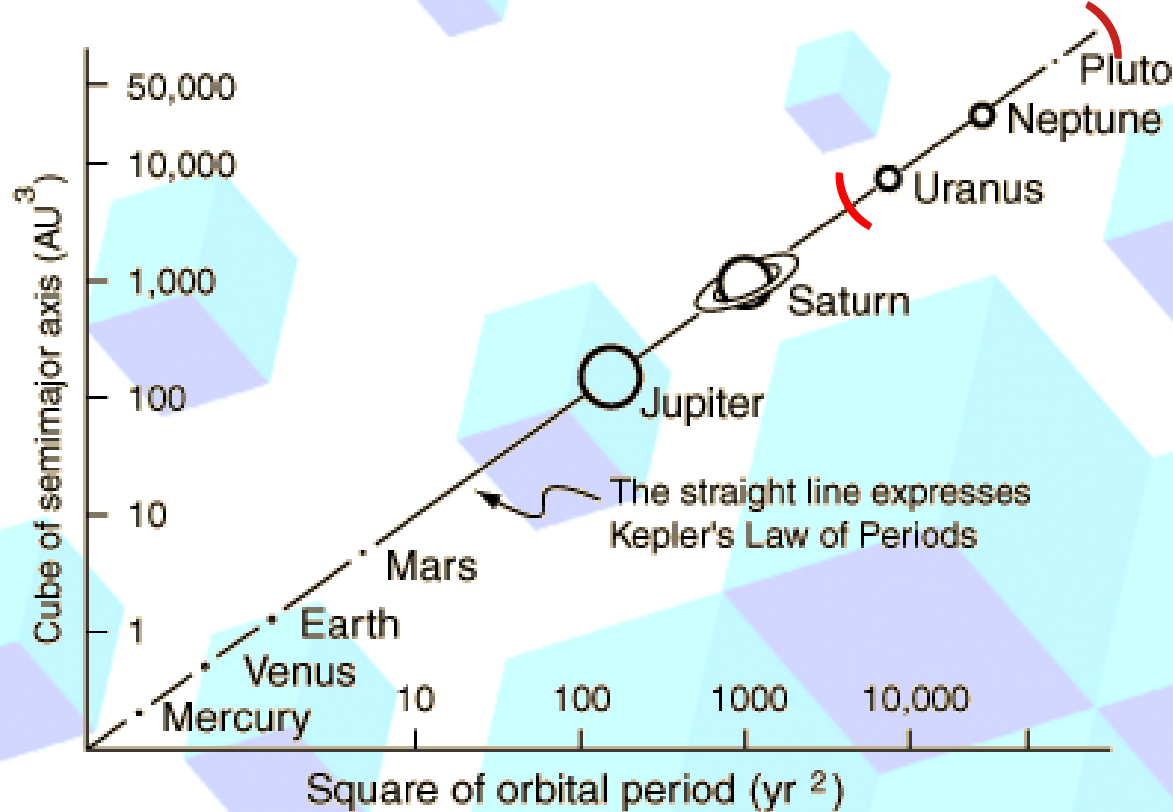
to know the trajectories of all planets in the solar system,  
from accurate observations (experiment)

or  
by numerically integrating general relativity equations  
(calculations at the highest level of theory)



# Learning → Discovery

Suppose  
to know the trajectories of all planets in the solar system,  
from accurate observations (experiment)  
or  
by numerically integrating general relativity equations  
(calculations at the highest level of theory)



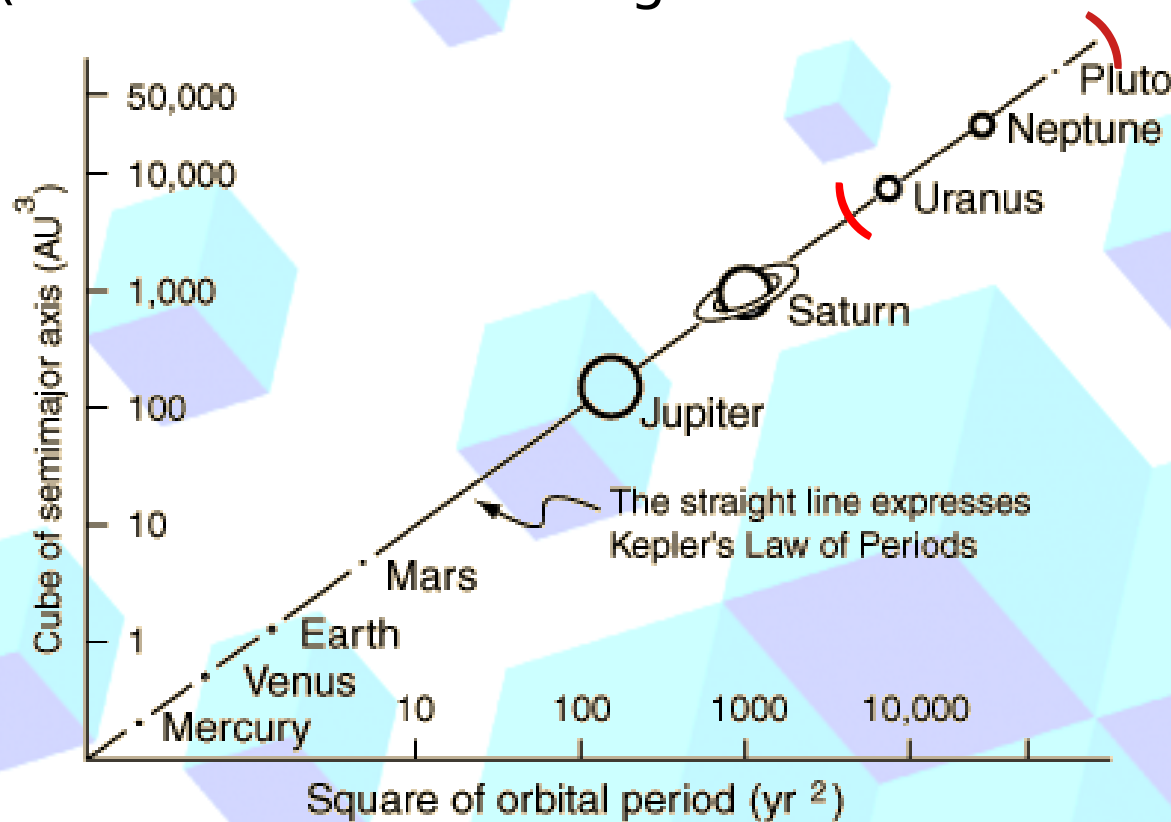
$$(\text{Orbital period})^2 = C (\text{orbit's major axis})^3$$

# Learning → Discovery

Suppose

to know the trajectories of all planets in the solar system,  
from accurate observations (experiment)

or  
by numerically integrating general relativity equations  
(calculations at the highest level of theory)



$$(\text{Orbital period})^2 = C (\text{orbit's major axis})^3$$

Data  
(collected by  
Tycho Brahe)

Statistical learning  
(performed by  
Johannes Kepler)

Physical law  
(assessed by  
Isaac Newton)

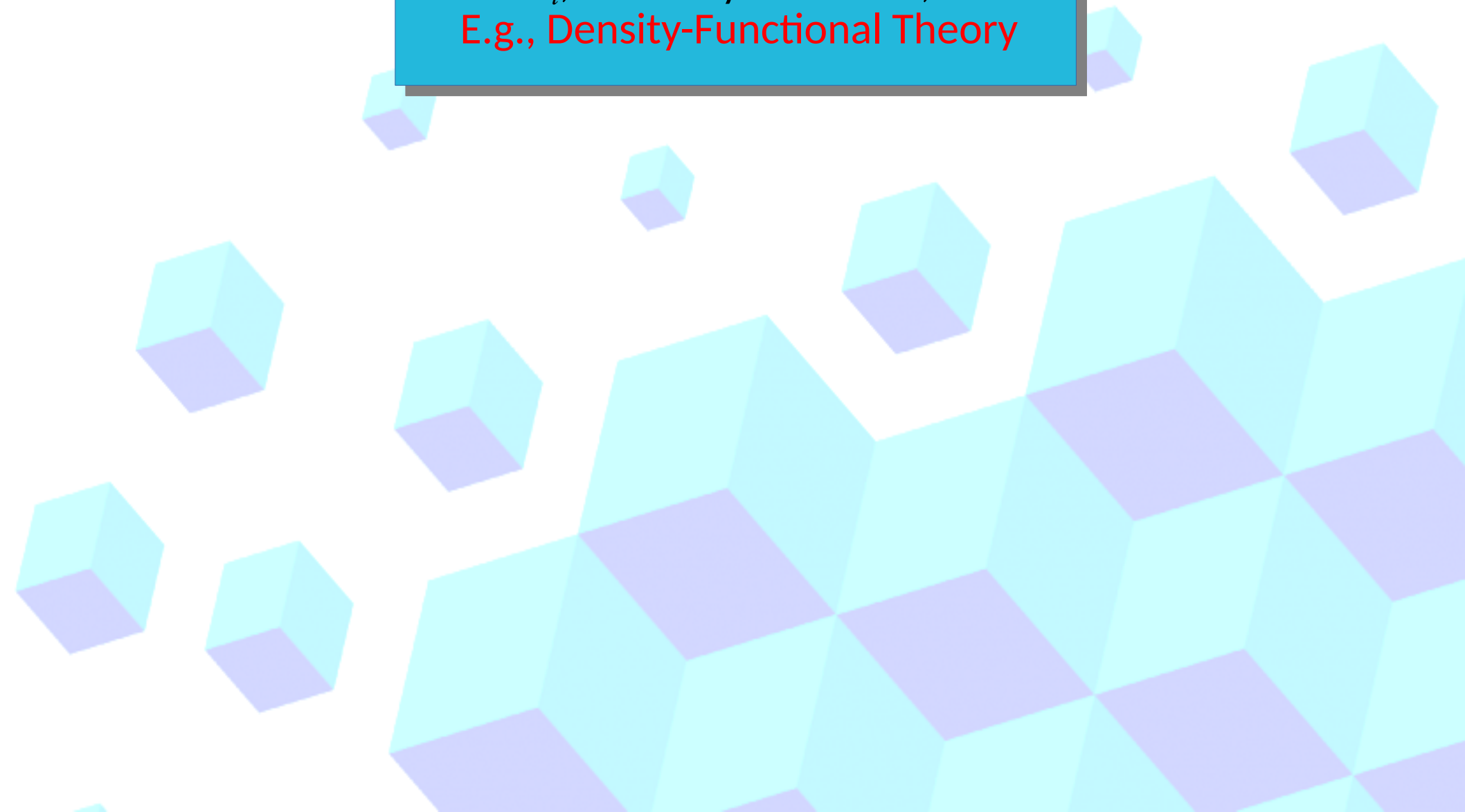
# Supervised (big-)data analysis: a flow chart

## Training set

Calculate properties and functions

$P_i$ , for many *materials*,  $i$

E.g., Density-Functional Theory



# Supervised (big-)data analysis: a 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:

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

# Supervised (big-)data analysis: a 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:

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

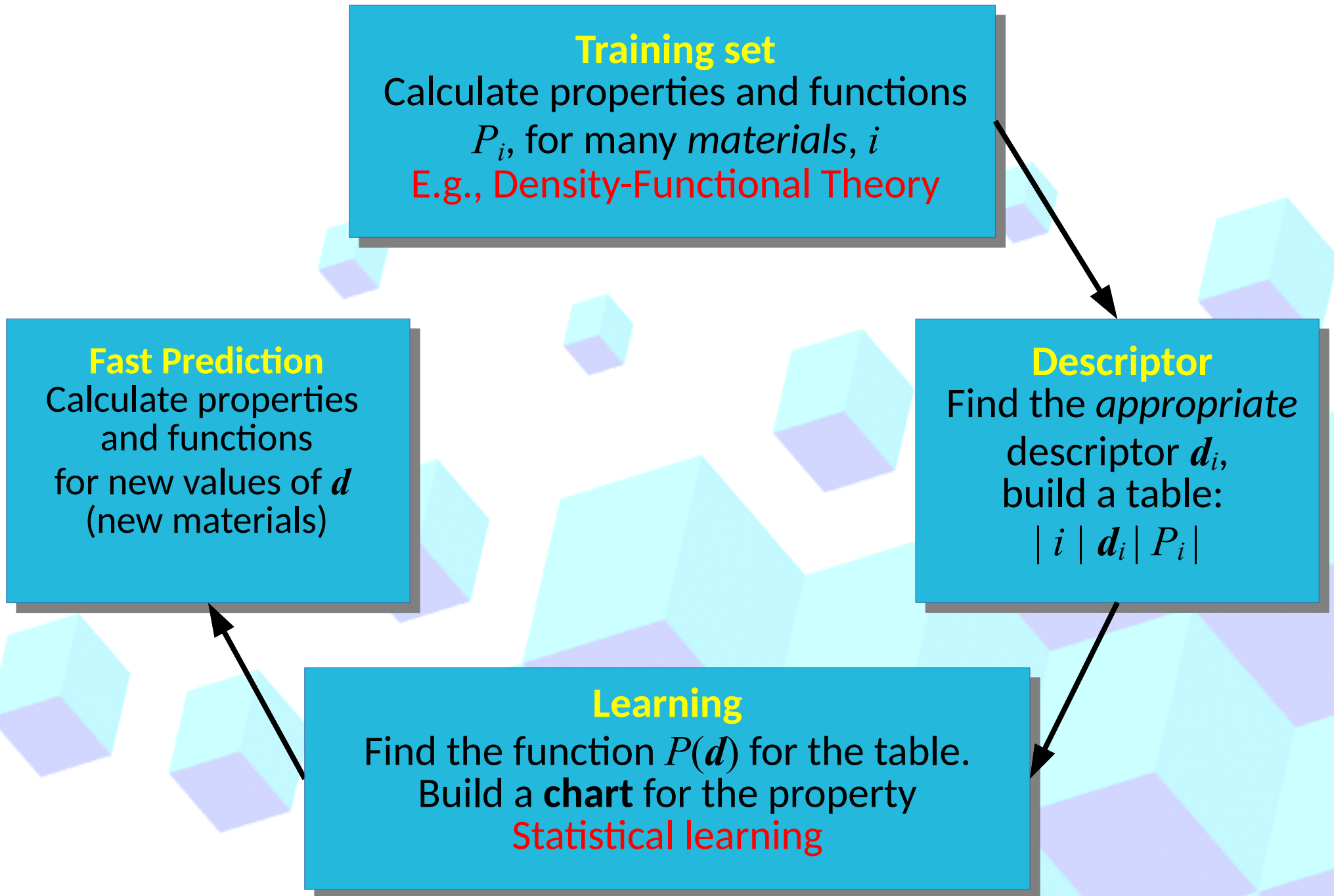
## Learning

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

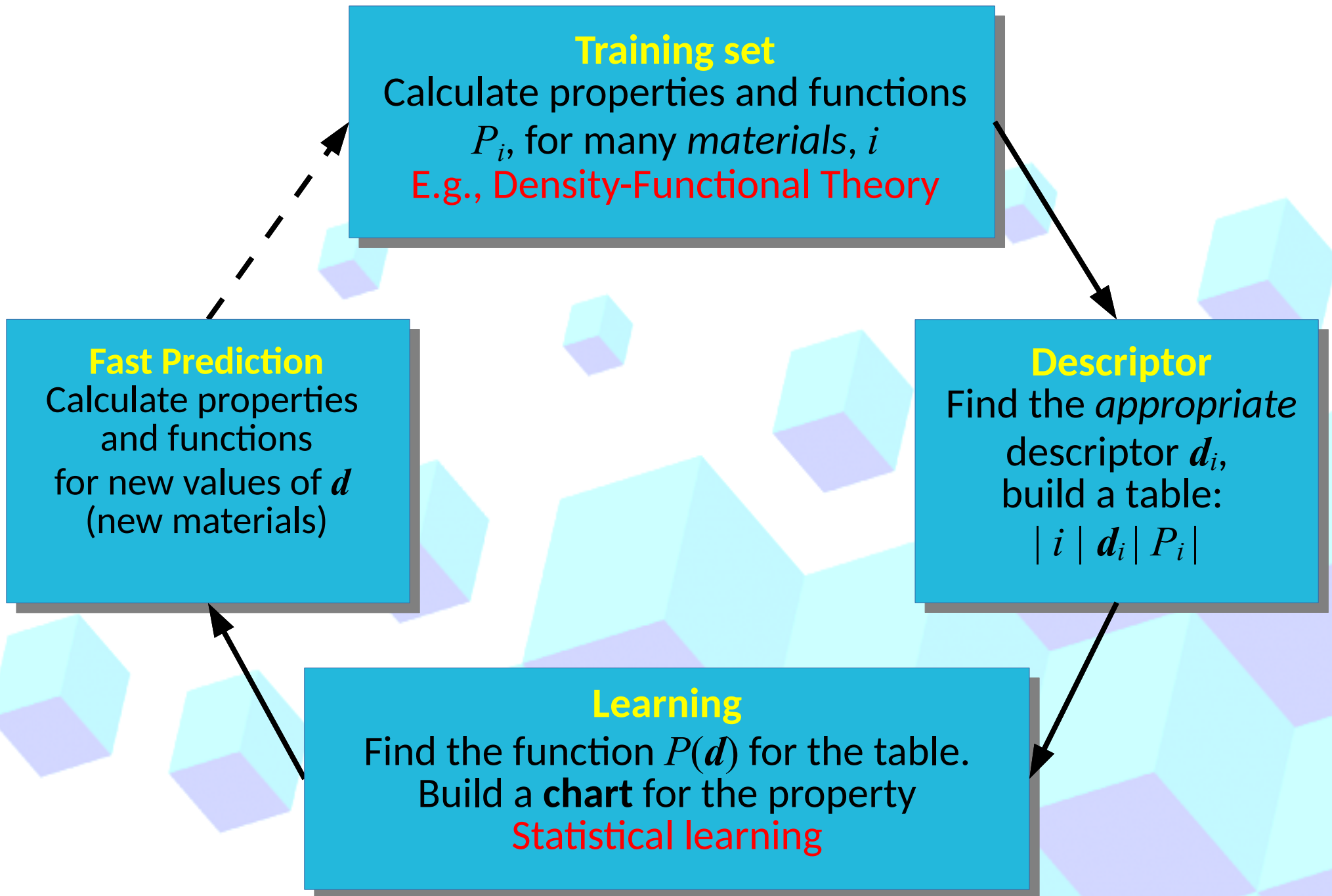
Build a **chart** for the property

Statistical learning

# Supervised (big-)data analysis: a flow chart

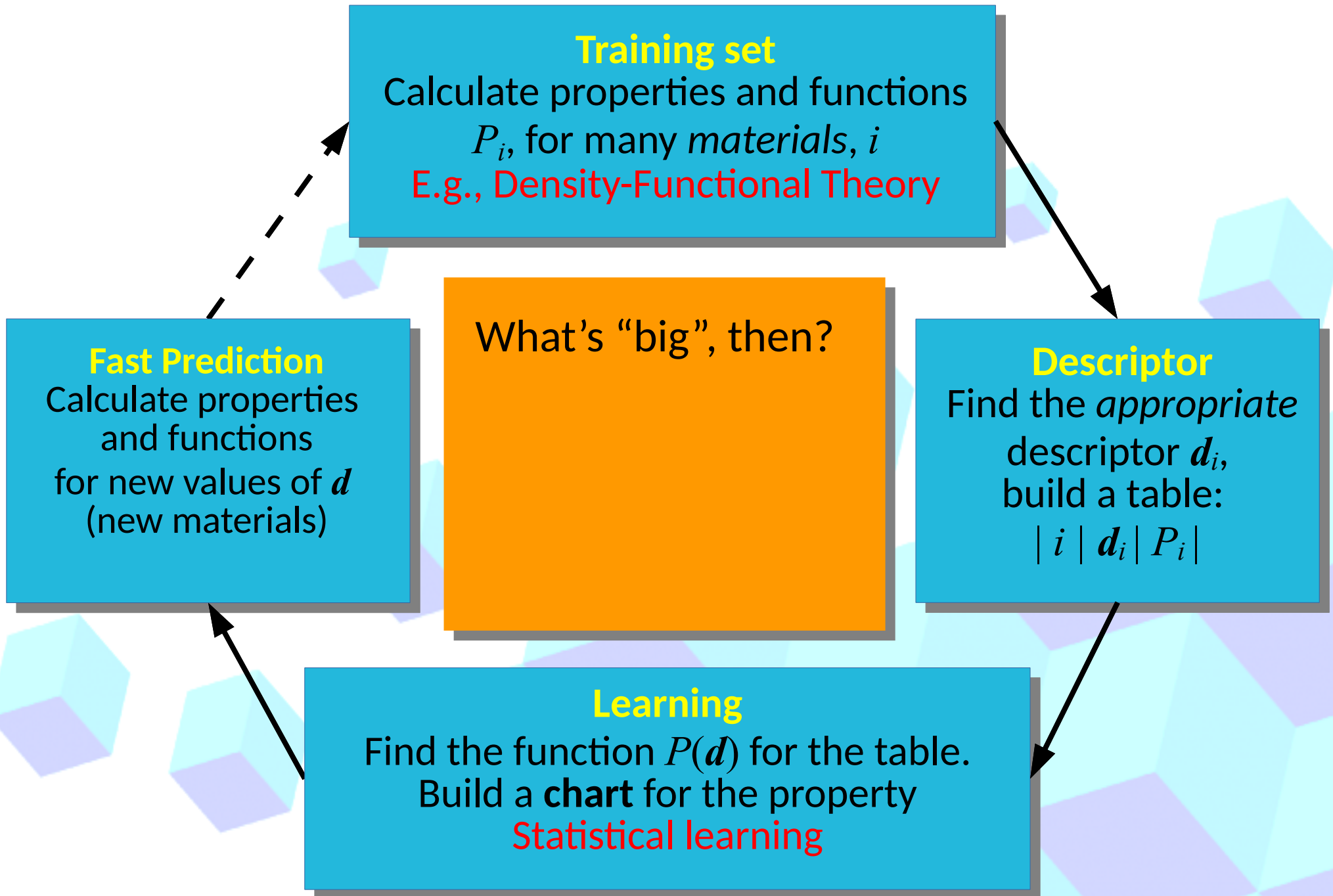


# Supervised (big-)data analysis: a flow chart

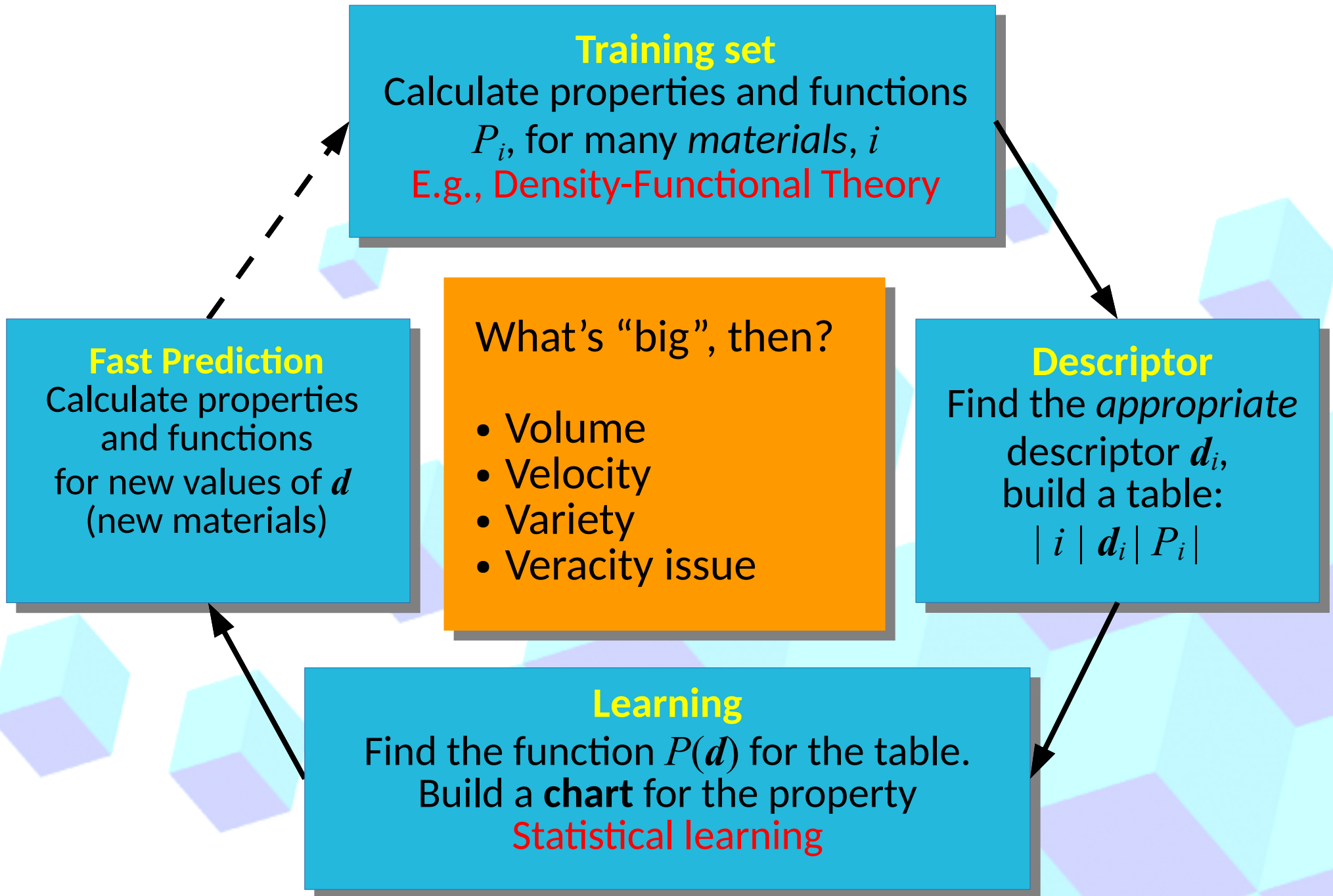




# Supervised big-data analysis: a flow chart



# Supervised big-data analysis: a flow chart



# Descriptor? Don't we know it from the start?

## Training set

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

$\{R_I, Z_I\} \rightarrow$  Hamiltonian

$\{R_I\} \rightarrow$  Geometry

- translational, rotational, permutational invariant

- coarse graining  $\{R_I\}$ ?

$\{Z_I\} \rightarrow$  Chemistry

## 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

# Regression: Mathematical formulation

Figure of merit to be optimized:

$$\operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^M} \sum_{j=1}^N \left( P_j - \sum_{l=1}^M d_{j,l} c_l \right)^2 = \operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^M} \| \mathbf{P} - \mathbf{D}\mathbf{c} \|_2^2 \quad \swarrow \ell_2 \text{ norm}$$

# Ridge Regression: Mathematical formulation

Figure of merit to be optimized:

$$\operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^M} \sum_{j=1}^N \left( P_j - \sum_{l=1}^M d_{j,l} c_l \right)^2 = \operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^M} \| \mathbf{P} - \mathbf{D}\mathbf{c} \|_2^2 \quad \swarrow \ell_2 \text{ norm}$$

Regularization (prefer “lower complexity” in the solution)

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

# Ridge Regression: Mathematical formulation

Figure of merit to be optimized:

$$\operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^M} \sum_{j=1}^N \left( P_j - \sum_{l=1}^M d_{j,l} c_l \right)^2 = \operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^M} \|\mathbf{P} - \mathbf{D}\mathbf{c}\|_2^2 \quad \swarrow \ell_2 \text{ norm}$$

Regularization (prefer “lower complexity” in the solution)

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

# Ridge Regression: Mathematical formulation

Figure of merit to be optimized:

$$\operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^M} \sum_{j=1}^N \left( P_j - \sum_{l=1}^M d_{j,l} c_l \right)^2 = \operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^M} \|\mathbf{P} - \mathbf{D}\mathbf{c}\|_2^2 \quad \swarrow \ell_2 \text{ norm}$$

Regularization (prefer “lower complexity” in the solution)

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

Explicit solver:

$$\mathbf{c} = \left( \mathbf{D}^\top \mathbf{D} + \lambda \mathbf{I} \right)^{-1} \mathbf{D}^\top \mathbf{P}$$

Alternative view, via Hilbert space representation theorem:

$$\mathbf{c} = \sum_j \alpha_j \mathbf{d}_j \quad \text{Sum over data points!}$$

# Kernel Ridge Regression: Mathematical formulation

$$\operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^M} \|\mathbf{P} - \mathbf{D}\mathbf{c}\|_2^2 + \lambda \|\mathbf{c}\|_2^2 \quad \Rightarrow \quad \mathbf{c} = \left( \mathbf{D}^\top \mathbf{D} + \lambda \mathbf{I} \right)^{-1} \mathbf{D}^\top \mathbf{P}$$

↓

$$\mathbf{c} = \sum_j \alpha_j \mathbf{d}_j$$



# Kernel Ridge Regression: Mathematical formulation

$$\operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^M} \|\mathbf{P} - \mathbf{D}\mathbf{c}\|_2^2 + \lambda \|\mathbf{c}\|_2^2 \quad \Rightarrow \quad \mathbf{c} = \left( \mathbf{D}^\top \mathbf{D} + \lambda \mathbf{I} \right)^{-1} \mathbf{D}^\top \mathbf{P}$$

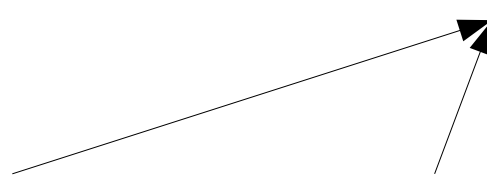
⇓

$$\mathbf{c} = \sum_j \alpha_j \mathbf{d}_j$$

⇓

$$\operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^M} \|\mathbf{P} - \mathbf{K}\boldsymbol{\alpha}\|_2^2 + \lambda \boldsymbol{\alpha}^\top \mathbf{K}\boldsymbol{\alpha}$$

$$K_{ij} = \langle \mathbf{d}_i, \mathbf{d}_j \rangle \quad \text{Linear kernel}$$



# Kernel Ridge Regression: Mathematical formulation

$$\operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^M} \|\mathbf{P} - \mathbf{D}\mathbf{c}\|_2^2 + \lambda \|\mathbf{c}\|_2^2 \quad \Rightarrow \quad \mathbf{c} = \left( \mathbf{D}^\top \mathbf{D} + \lambda \mathbf{I} \right)^{-1} \mathbf{D}^\top \mathbf{P}$$

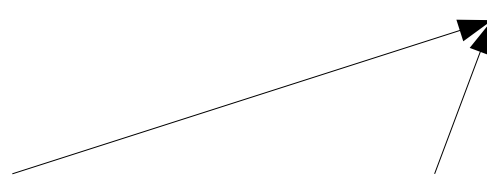
⇓

$$\mathbf{c} = \sum_j \alpha_j \mathbf{d}_j$$

⇓

$$\operatorname{argmin}_{\boldsymbol{\alpha} \in \mathbb{R}^M} \|\mathbf{P} - \mathbf{K}\boldsymbol{\alpha}\|_2^2 + \lambda \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha} \quad \Rightarrow \quad \boldsymbol{\alpha} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{P}$$

$K_{ij} = \langle \mathbf{d}_i, \mathbf{d}_j \rangle$  **Linear kernel**



# Kernel Ridge Regression: Mathematical formulation

$$\operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^M} \|\mathbf{P} - \mathbf{D}\mathbf{c}\|_2^2 + \lambda \|\mathbf{c}\|_2^2 \quad \Rightarrow \quad \mathbf{c} = \left(\mathbf{D}^\top \mathbf{D} + \lambda \mathbf{I}\right)^{-1} \mathbf{D}^\top \mathbf{P}$$

⇓

$$\mathbf{c} = \sum_j \alpha_j \mathbf{d}_j$$

⇓

$$\operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^M} \|\mathbf{P} - \mathbf{K}\boldsymbol{\alpha}\|_2^2 + \lambda \boldsymbol{\alpha}^\top \mathbf{K}\boldsymbol{\alpha} \quad \Rightarrow \quad \boldsymbol{\alpha} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{P}$$

$$K_{ij} = \langle \mathbf{d}_i, \mathbf{d}_j \rangle \quad \text{Linear kernel}$$

## Non-linear kernel

$$\mathbf{c} = \sum_j \alpha_j \Phi(\mathbf{d}_j)$$

$$K_{ij} = \langle \Phi(\mathbf{d}_i), \Phi(\mathbf{d}_j) \rangle$$

$$\operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^M} \|\mathbf{P} - \mathbf{K}\boldsymbol{\alpha}\|_2^2 + \lambda \boldsymbol{\alpha}^\top \mathbf{K}\boldsymbol{\alpha} \quad \Rightarrow \quad \boldsymbol{\alpha} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{P}$$

# Kernel Ridge Regression: Mathematical formulation

## Non-linear kernel

$$\mathbf{c} = \sum_j \alpha_j \Phi(\mathbf{d}_j)$$

$$K_{ij} = \langle \Phi(\mathbf{d}_i), \Phi(\mathbf{d}_j) \rangle$$

$$\operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^M} \|\mathbf{P} - \mathbf{K}\boldsymbol{\alpha}\|_2^2 + \lambda \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha} \quad \Rightarrow \quad \boldsymbol{\alpha} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{P}$$

$$K_{ij} = \langle \mathbf{d}_i, \mathbf{d}_j \rangle$$

**Linear kernel**

$$K_{ij} = (\langle \mathbf{d}_i, \mathbf{d}_j \rangle + b)^n$$

**Polynomial kernel**

$$K_{ij} = \exp\left(\frac{\|\mathbf{d}_i - \mathbf{d}_j\|^2}{2\sigma^2}\right)$$

**Gaussian (radial basis function) kernel**

$$K_{ij} = \exp\left(\frac{\|\mathbf{d}_i - \mathbf{d}_j\|}{\sigma}\right)$$

**Laplacian kernel**

# Kernel Ridge Regression: Mathematical formulation

## Non-linear kernel

$$\mathbf{c} = \sum_j \alpha_j \Phi(\mathbf{d}_j)$$

$$K_{ij} = \langle \Phi(\mathbf{d}_i), \Phi(\mathbf{d}_j) \rangle$$

$$\operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^M} \|\mathbf{P} - \mathbf{K}\boldsymbol{\alpha}\|_2^2 + \lambda \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha} \quad \Rightarrow \quad \boldsymbol{\alpha} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{P}$$

$$K_{ij} = \langle \mathbf{d}_i, \mathbf{d}_j \rangle$$

$$K_{ij} = (\langle \mathbf{d}_i, \mathbf{d}_j \rangle + b)^n$$

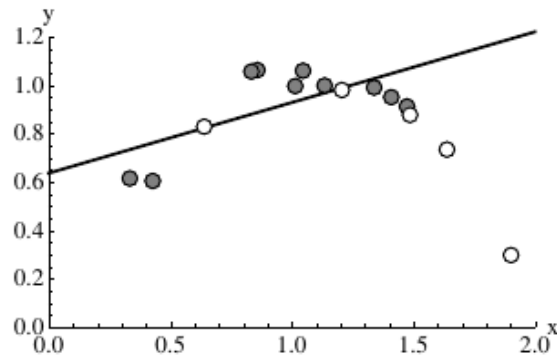
$$K_{ij} = \exp\left(\frac{\|\mathbf{d}_i - \mathbf{d}_j\|^2}{2\sigma^2}\right)$$

$$K_{ij} = \exp\left(\frac{\|\mathbf{d}_i - \mathbf{d}_j\|}{\sigma}\right)$$

In all cases,  
a kernel introduces a  
**similarity measure**

# Regularized regression in practice: beware of overfitting

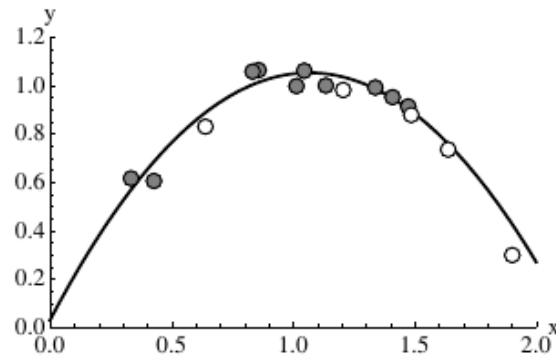
## Underfitting



Training/  
validation  
error

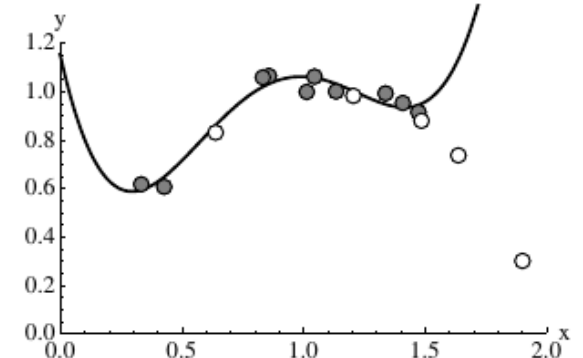
0.123 / 0.443

## Fitting



0.044 / 0.068

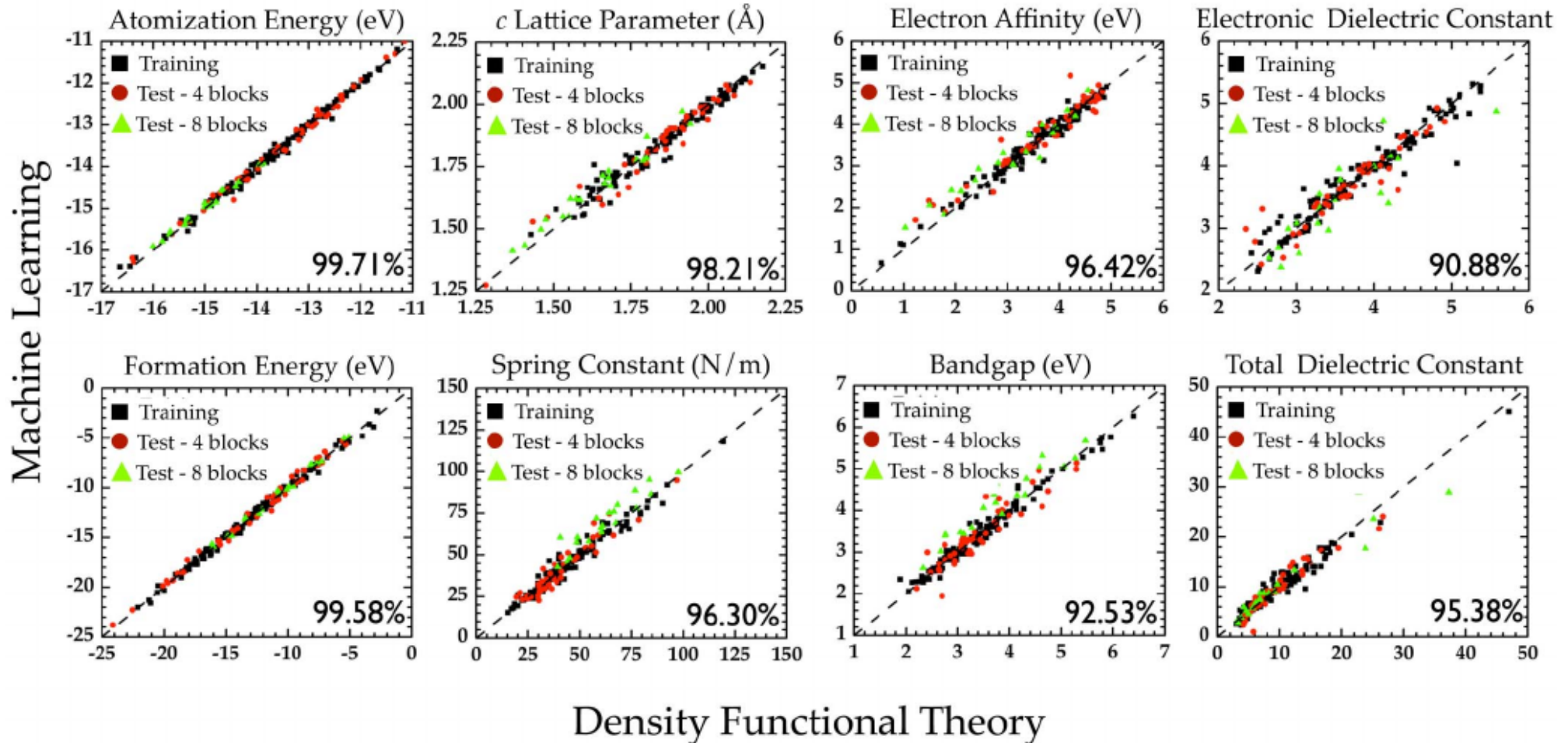
## Overfitting



0.036 / 0.939

# KRR success stories: 1D polymers “eugenetics”

Data: 175 linear 4-blocks periodic polymers. 7 blocks:  $\text{CH}_2$ ,  $\text{SiF}_2$ ,  $\text{SiCl}_2$ ,  $\text{GeF}_2$ ,  $\text{GeCl}_2$ ,  $\text{SnF}_2$ ,  $\text{SnCl}_2$ ,  
Descriptor: 20 dimensions [# building blocks of type  $i$ , of  $ii$  pairs, of  $iii$  triplets]





# KRR success stories: $n$ -grams for kaggle

Research Prediction Competition

## Nomad2018 Predicting Transparent Conductors

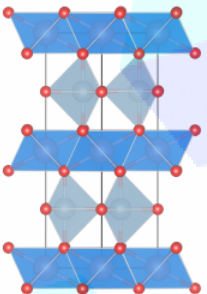
Predict the key properties of novel transparent semiconductors

€5,000 Prize Money

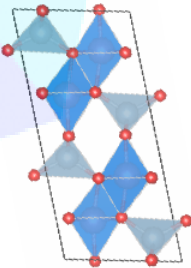
883 teams · 12 days ago

Overview Data Kernels Discussion Leaderboard Rules Team Host My Submissions **Late Submission**

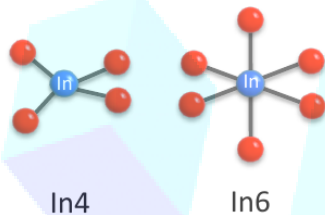
p63/mmc



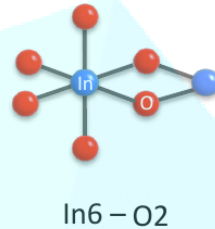
C2/m



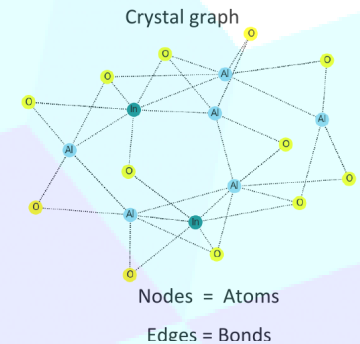
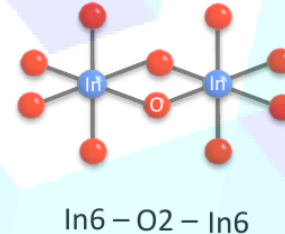
Uni-gram



Bi-gram



Tri-gram



Input features: count number of sequences of various lengths



Unigrams: 2 O3, 2 Ga4, 1 O2, 1 In5

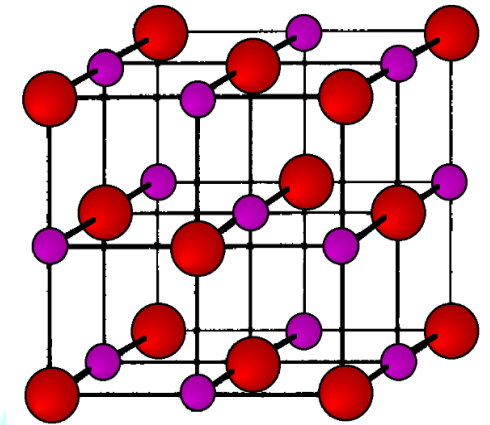
Bigrams: 2 O3-Ga4, 2 Ga4-O2, 1 O3-In5



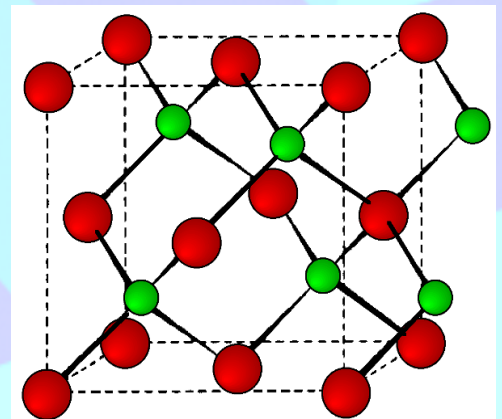
# An example: predicting crystal structures from the composition

82 octet AB binary compounds

hydrogen 1 H 1.008	beryllium 4 Be 9.0122											boron 5 B 10.811	carbon 6 C 12.011	nitrogen 7 N 14.007	oxygen 8 O 15.999	fluorine 9 F 18.998	helium 2 He 4.0026		
lithium 3 Li 6.941	beryllium 4 Be 9.0122											aluminum 13 Al 26.982	silicon 14 Si 28.086	phosphorus 15 P 30.974	sulfur 16 S 32.065	chlorine 17 Cl 35.453	neon 10 Ne 20.180		
sodium 11 Na 22.990	magnesium 12 Mg 24.305	scandium 21 Sc 44.956	titanium 22 Ti 47.867	vanadium 23 V 50.942	chromium 24 Cr 51.996	manganese 25 Mn 54.938	iron 26 Fe 55.845	cobalt 27 Co 58.933	nickel 28 Ni 58.693	argon 18 Ar 39.948									
potassium 19 K 39.098	calcium 20 Ca 40.078	yttrium 39 Y 88.906	zirconium 40 Zr 91.224	niobium 41 Nb 92.906	molybdenum 42 Mo 95.94	technetium 43 Tc [98]	ruthenium 44 Ru 101.07	rhodium 45 Rh 102.91	palladium 46 Pd 106.42	argon 18 Ar 39.948									
rubidium 37 Rb 85.468	strontium 38 Sr 87.62	lutetium 71 Lu 174.97	hafnium 72 Hf 178.49	tantalum 73 Ta 180.95	tungsten 74 W 183.84	rhenium 75 Re 186.21	osmium 76 Os 190.23	iridium 77 Ir 192.22	platinum 78 Pt 195.08	argon 18 Ar 39.948									
caesium 55 Cs 132.91	barium 56 Ba 137.33	* 57-70 Lr	* 71 Lu	* 72 Hf	* 73 Ta	* 74 W	* 75 Re	* 76 Os	* 77 Ir	* 78 Pt	gold 79 Au 196.97	mercury 80 Hg 200.59	indium 49 In 114.82	tin 50 Sn 118.71	antimony 51 Sb 121.76	tellurium 52 Te 127.60	bromine 35 Br 79.904	iodine 53 I 126.90	radon 86 Rn [222]
francium 87 Fr [223]	radium 88 Ra [226]	* 89-102 Ac	* 103 Lr	* 104 Rf	* 105 Db	* 106 Sg	* 107 Bh	* 108 Hs	* 109 Mt	* 110 Uun	* 111 Uuu	* 112 Uub	thallium 81 Tl 204.38	lead 82 Pb 207.2	bismuth 83 Bi 208.98	polonium 84 Po [209]	astatine 85 At [210]	radon 86 Rn [222]	radon 86 Rn [222]
copper 29 Cu 63.546	zinc 30 Zn 65.39											gallium 31 Ga 69.723	germanium 32 Ge 72.61	arsenic 33 As 74.922	selenium 34 Se 78.96	iodine 53 I 126.90	radon 86 Rn [222]		
silver 47 Ag [107.8682]	cadmium 48 Cd [112.4114]											indium 49 In 114.82	tin 50 Sn 118.71	antimony 51 Sb 121.76	tellurium 52 Te 127.60	iodine 53 I 126.90	radon 86 Rn [222]		
Lanthanide series		lanthanum 57 La 138.91	cerium 58 Ce 140.12	praseodymium 59 Pr 140.91	neodymium 60 Nd 144.24	promethium 61 Pm [145]	samarium 62 Sm 150.36	europium 63 Eu 151.96	gadolinium 64 Gd 157.25	terbium 65 Tb 158.93	dysprosium 66 Dy 162.50	holmium 67 Ho 164.93	erbium 68 Er 167.26	thulium 69 Tm 168.93	ytterbium 70 Yb 173.04				
Actinide series		actinium 89 Ac [227]	thorium 90 Th 232.04	protactinium 91 Pa 231.04	uranium 92 U 238.03	neptunium 93 Np [237]	plutonium 94 Pu [244]	americium 95 Am [243]	curium 96 Cm [247]	berkelium 97 Bk [247]	californium 98 Cf [251]	einsteinium 99 Es [252]	fermium 100 Fm [257]	mendelevium 101 Md [258]	nobelium 102 No [259]				



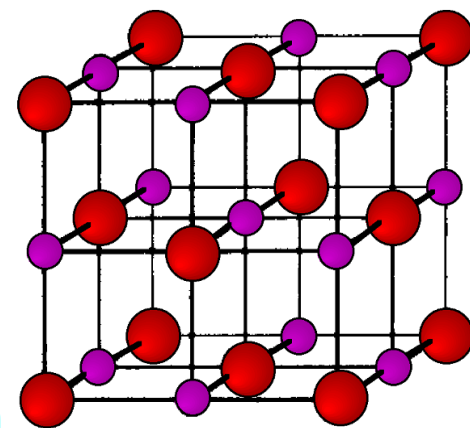
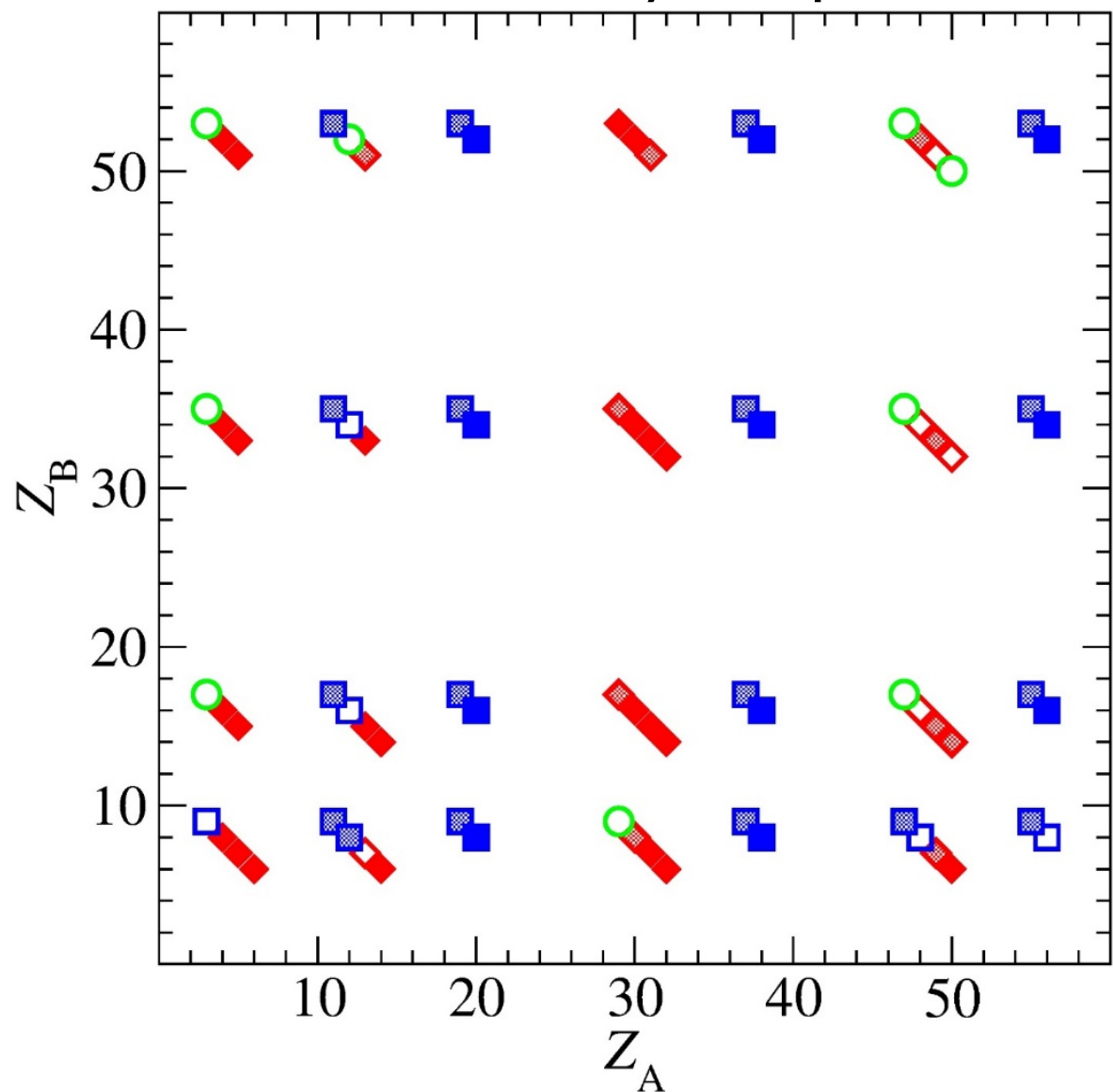
Rock salt



Zinc blende

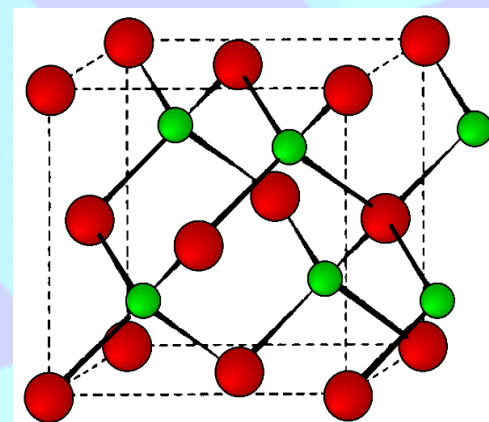
# An example: predicting crystal structures from the composition

82 octet AB binary compounds



Rock salt

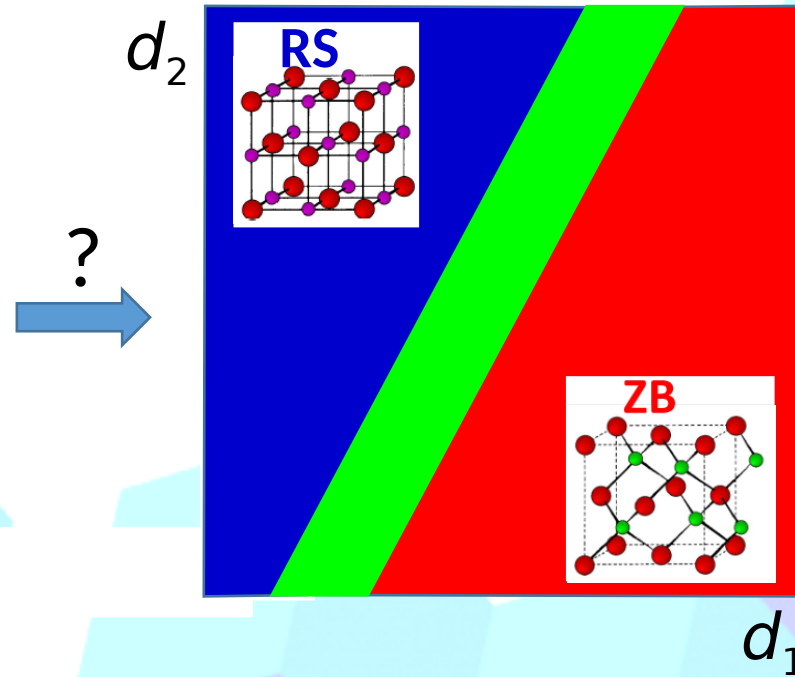
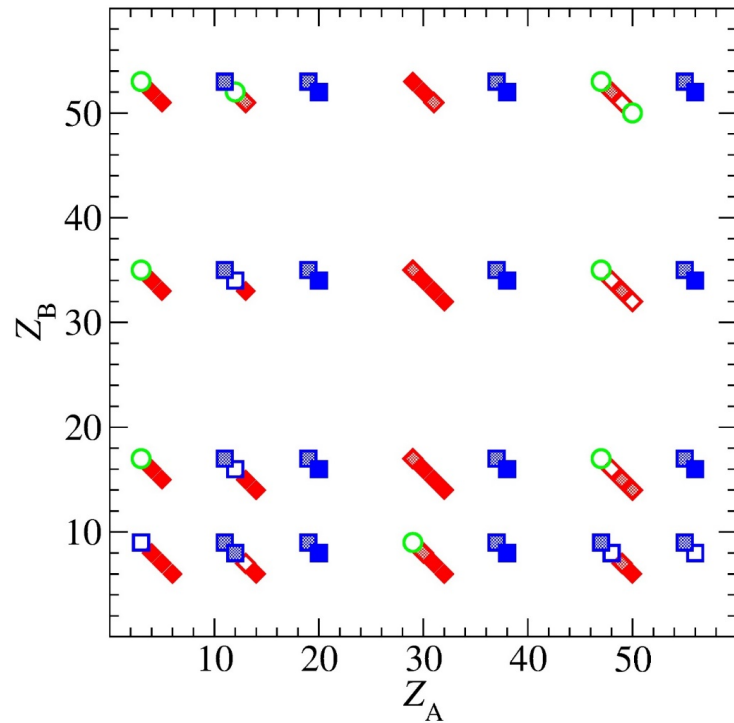
- Rock salt
- Rock salt/Zinc blende
- ◆ Zinc blende



Zinc blende

# An example: predicting crystal structures from the composition

## 82 octet AB binary compounds

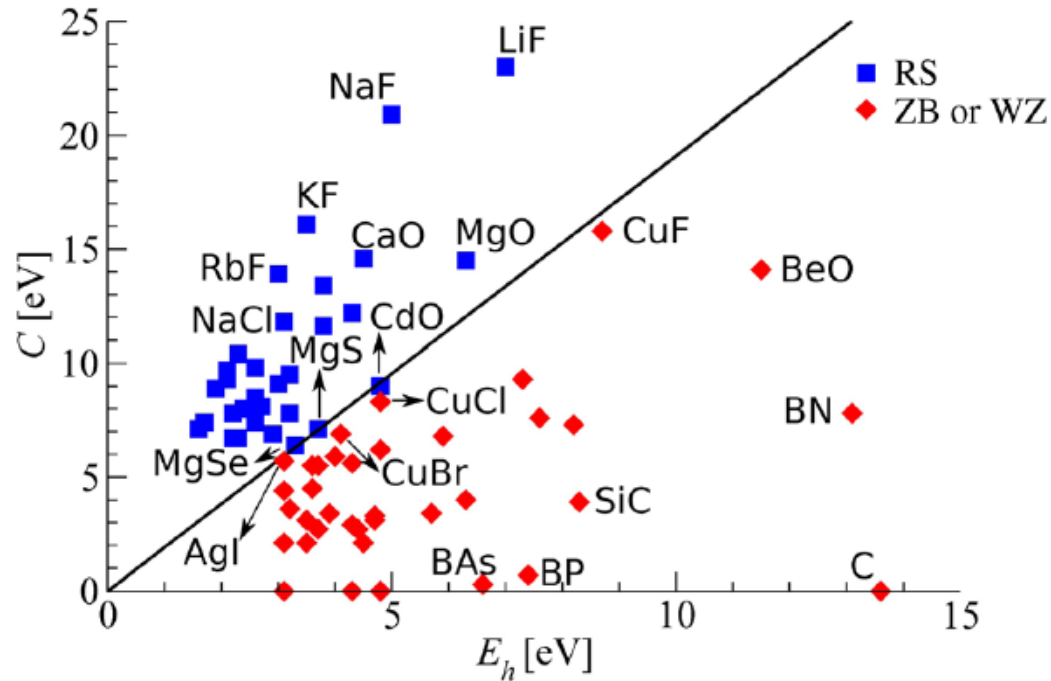


- Rock salt
- Rock salt/Zinc blende
- ◆ Zinc blende

J. A. van Vechten, *Phys. Rev.* 182, 891 (1969).  
J. C. Phillips, *Rev. Mod. Phys.* 42, 317 (1970).  
J. John and A.N. Bloch, *Phys. Rev. Lett.* 33, 1095 (1974)  
J. R. Chelikowsky and J. C. Phillips, *Phys. Rev. B* 33, 2453 (1978)  
A. Zunger, *Phys. Rev. B* 22, 5839 (1980).  
D. G. Pettifor, *Solid State Commun.* 51, 31 (1984).  
Y. Saad, D. Gao, T. Ngo, S. Bobbitt, J. R. Chelikowsky, and W. Andreoni, *Phys. Rev. B* 85, 104104 (2012).

# An example: predicting crystal structures from the composition

## 82 octet AB binary compounds



The descriptor proposed by Phillips and van Vechten in 1969-70 depends on:

- lattice parameter
- electrical conductivity

J. A. van Vechten, *Phys. Rev.* 182, 891 (1969).

J. C. Phillips, *Rev. Mod. Phys.* 42, 317 (1970).

J. John and A.N. Bloch, *Phys. Rev. Lett.* 33, 1095 (1974)

J. R. Chelikowsky and J. C. Phillips, *Phys. Rev. B* 33, 2453 (1978)

A. Zunger, *Phys. Rev. B* 22, 5839 (1980).

D. G. Pettifor, *Solid State Commun.* 51, 31 (1984).

Y. Saad, D. Gao, T. Ngo, S. Bobbitt, J. R.

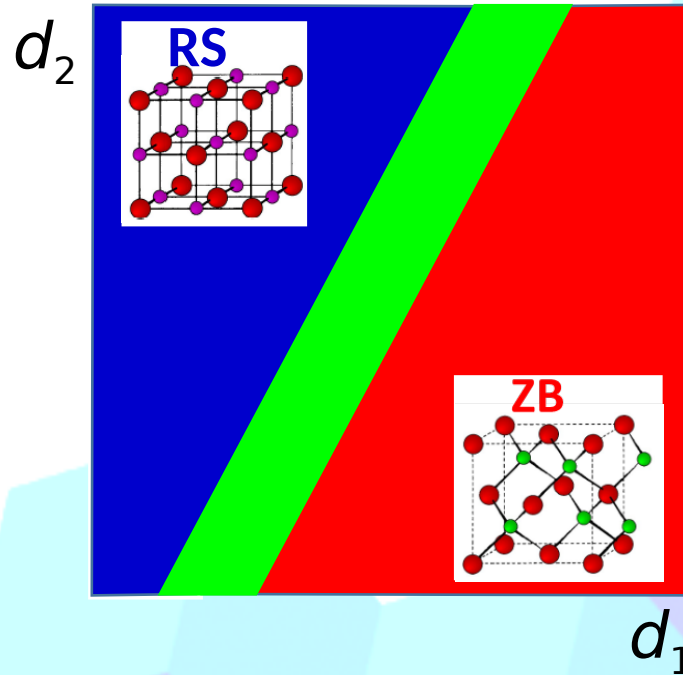
Chelikowsky, and W. Andreoni, *Phys. Rev. B* 85, 104104 (2012).

# An example: predicting crystal structures from the composition

## 82 octet AB binary compounds

Ansatz: atomic features

- HOMO
- LUMO
- Ionization Potential
- Electron Affinity
- Radius of valence  $s$  orbital
- Radius of valence  $p$  orbital
- Radius of valence  $d$  orbital
- ... ?



- Rock salt
- Rock salt/Zinc blende
- ◆ Zinc blende

$E(\text{Rock salt}) - E(\text{Zinc blende})$

J. A. van Vechten, *Phys. Rev.* 182, 891 (1969).

J. C. Phillips, *Rev. Mod. Phys.* 42, 317 (1970).

J. John and A.N. Bloch, *Phys. Rev. Lett.* 33, 1095 (1974)

J. R. Chelikowsky and J. C. Phillips, *Phys. Rev. B* 33, 2453 (1978)

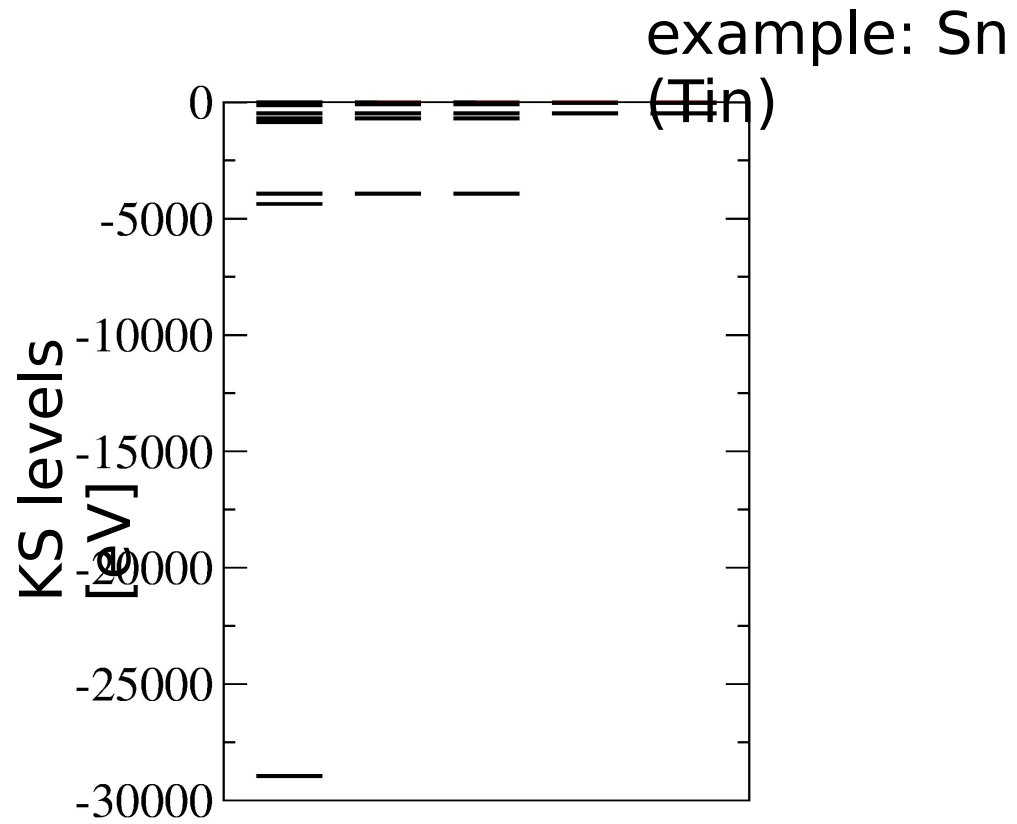
A. Zunger, *Phys. Rev. B* 22, 5839 (1980).

D. G. Pettifor, *Solid State Commun.* 51, 31 (1984).

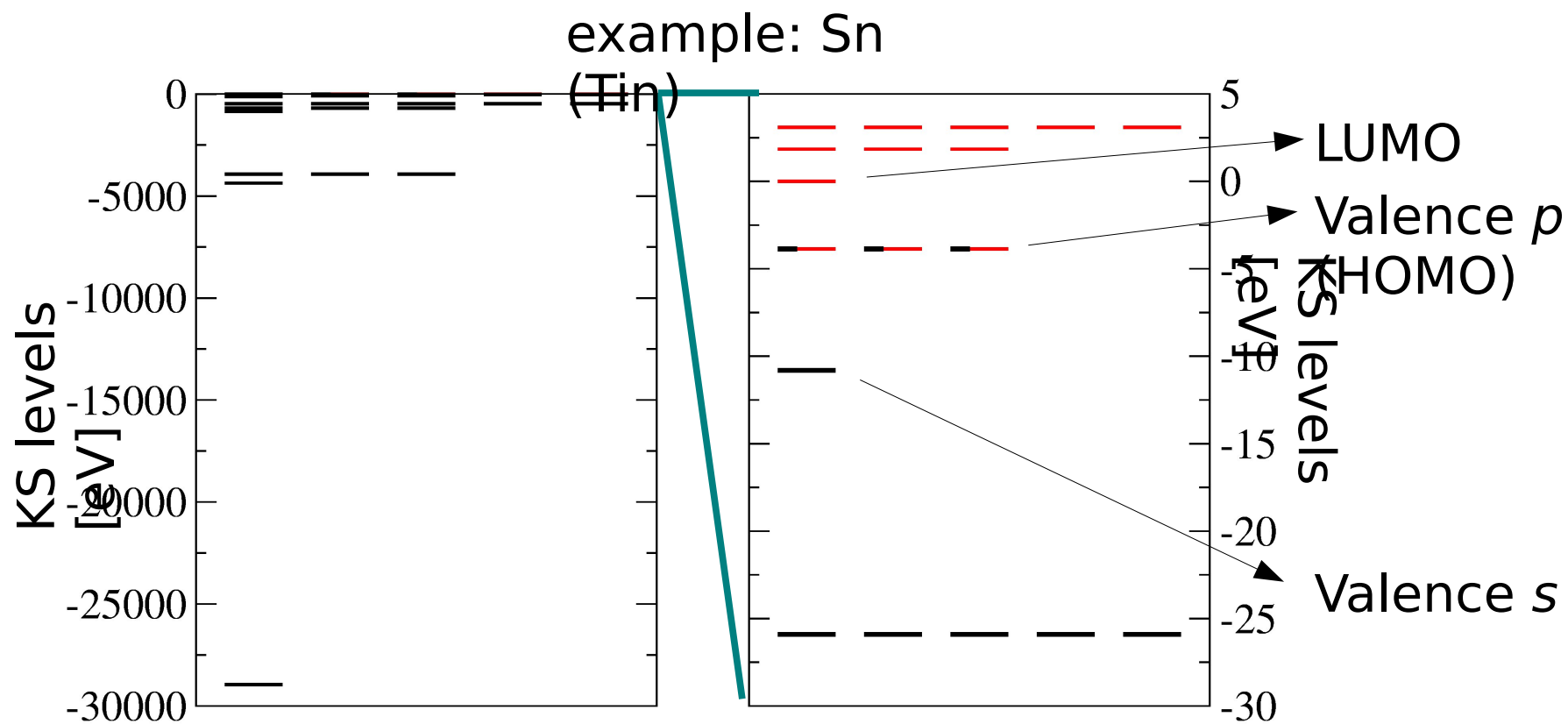
Y. Saad, D. Gao, T. Ngo, S. Bobbitt, J. R.

Chelikowsky, and W. Andreoni, *Phys. Rev. B* 85, 104104 (2012).

# Primary (atomic) features

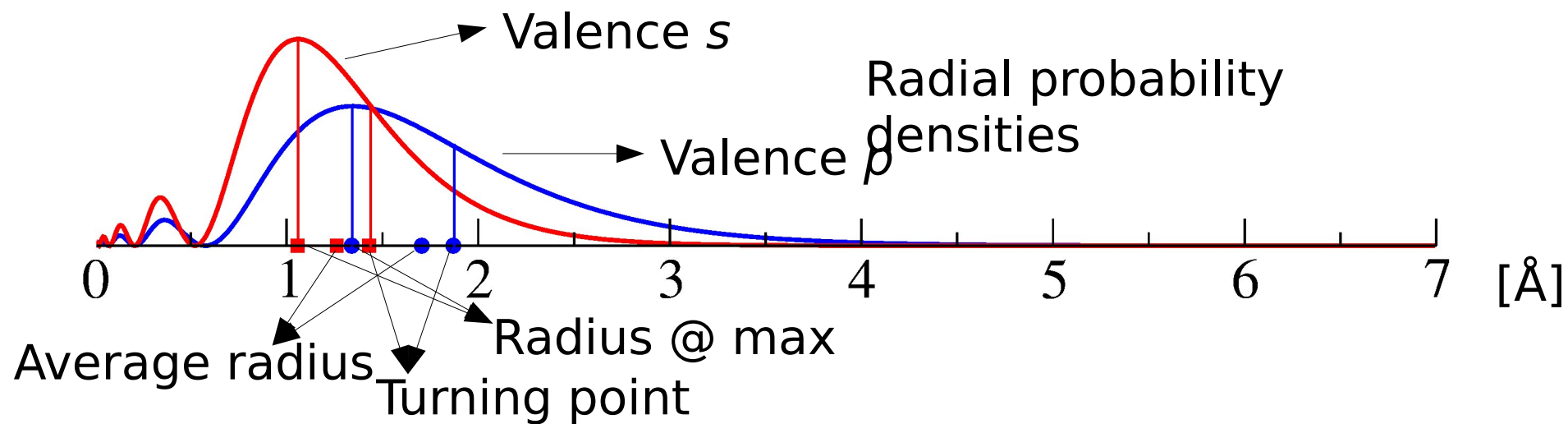
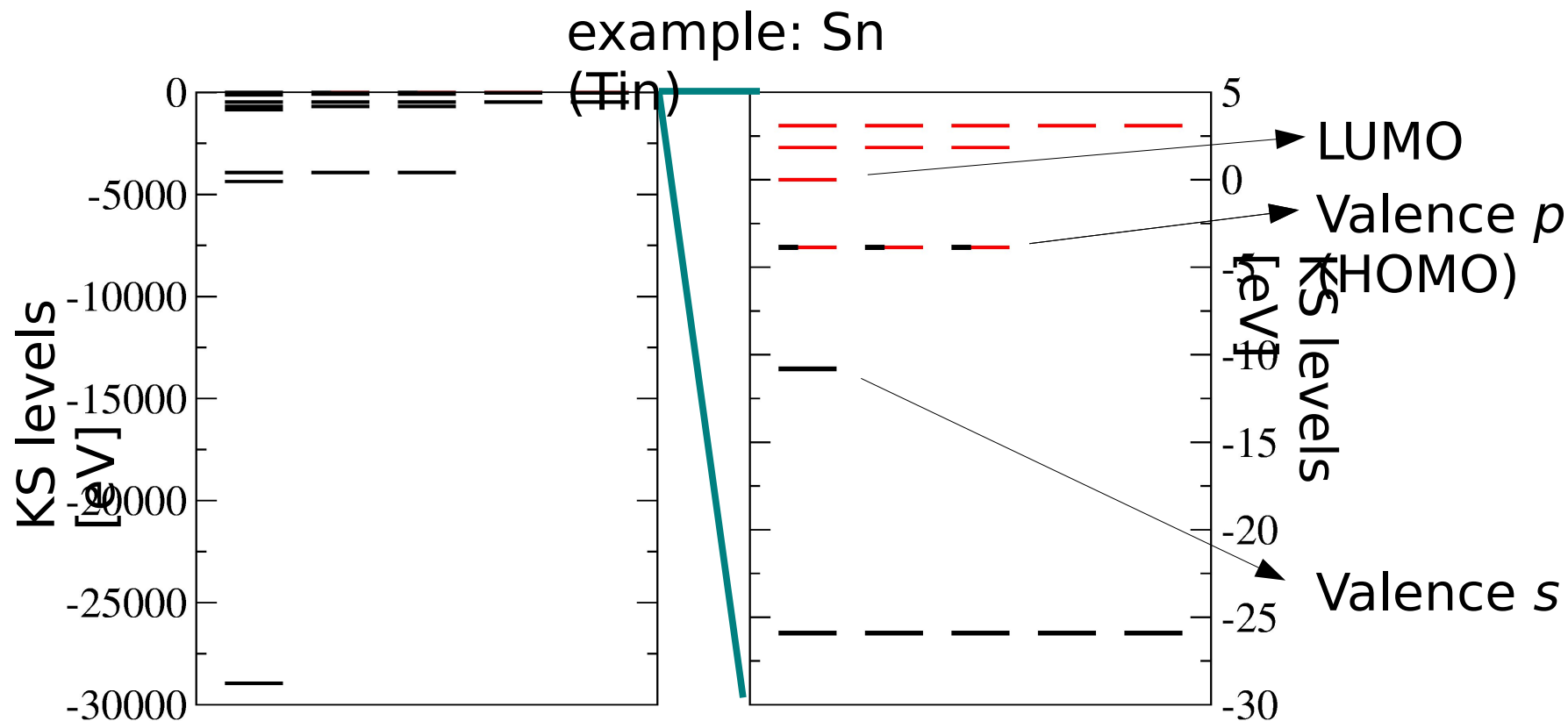


# Primary (atomic) features





# Primary (atomic) features



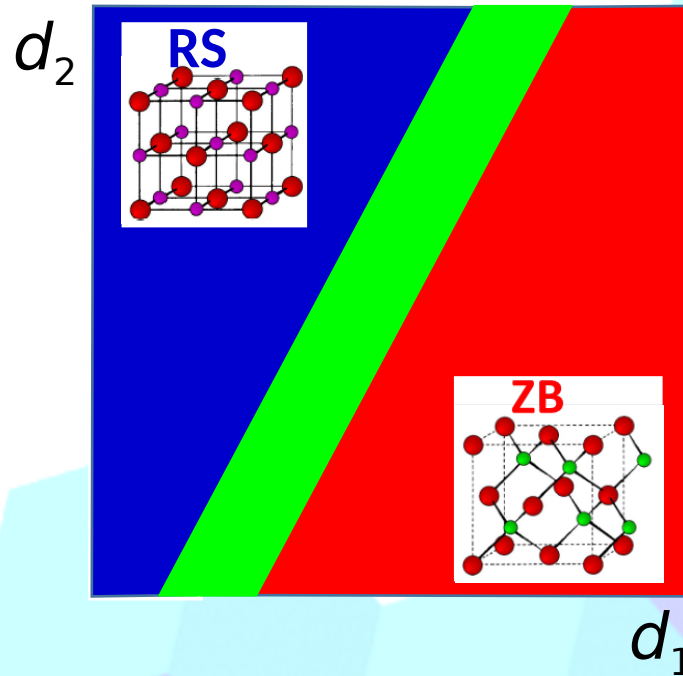


# An example: predicting crystal structures from the composition

## 82 octet AB binary compounds

Ansatz: atomic features

- HOMO
- LUMO
- Ionization Potential
- Electron Affinity
- Radius of valence  $s$  orbital
- Radius of valence  $p$  orbital
- Radius of valence  $d$  orbital
- ... ?



- Rock salt
- Rock salt/Zinc blende
- ◆ Zinc blende

$E(\text{Rock salt}) - E(\text{Zinc blende})$

- J. A. van Vechten, *Phys. Rev.* 182, 891 (1969).
- J. C. Phillips, *Rev. Mod. Phys.* 42, 317 (1970).
- J. John and A.N. Bloch, *Phys. Rev. Lett.* 33, 1095 (1974)
- J. R. Chelikowsky and J. C. Phillips, *Phys. Rev. B* 33, 2453 (1978)
- A. Zunger, *Phys. Rev. B* 22, 5839 (1980).
- D. G. Pettifor, *Solid State Commun.* 51, 31 (1984).
- Y. Saad, D. Gao, T. Ngo, S. Bobbitt, J. R. Chelikowsky, and W. Andreoni, *Phys. Rev. B* 85, 104104 (2012).

# Compressed sensing

Aim: finding descriptors and learning predictive models

Ansatz:

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

Where

$\mathbf{P}$  is the property of interest

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

$c_1, \dots, c_n$  are unknown coefficients, with the extra constraint that these (nonzero) coefficients should be as few as possible.

# Compressed sensing

Aim: finding descriptors and learning predictive models

Ansatz:

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

Where

$\mathbf{P}$  is the property of interest

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

$c_1, \dots, c_n$  are unknown coefficients, with the extra constraint that these (nonzero) coefficients should be as few as possible.

With a foreword on  
**dimensionality reduction**

# Linear dimensionality reduction: Principal components

Pearson, K. "On Lines and Planes of Closest Fit to Systems of Points in Space".  
Philosophical Magazine 2, 559 (1901)

# Linear dimensionality reduction: Principal components

Pearson, K. "On Lines and Planes of Closest Fit to Systems of Points in Space".  
Philosophical Magazine 2, 559 (1901)

Orthonormal transformation of coordinates, converting a set of (possibly) linearly correlated coordinates into a new set of linearly uncorrelated (called principal or normal) components, such that the first component has the largest variance and each subsequent has the largest variance constrained to being orthogonal to all the preceding components

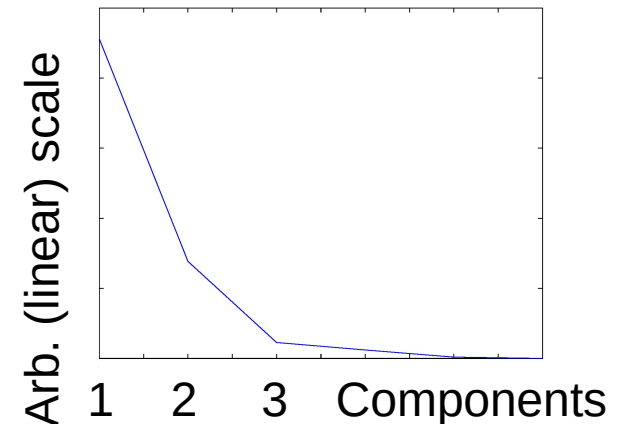
# Linear dimensionality reduction: Principal components

## Ansatz: atomic features

- Valence number  $Z_v$
- Energy of valence  $s$  orbital  $E_s$
- Energy of valence  $p$  orbital  $E_p$
- Radius of valence  $s$  orbital  $r_s$
- Radius of valence  $p$  orbital  $r_p$

$r_s, r_p, E_s/\sqrt{Z_v}, E_p/\sqrt{Z_v}$ ,  
for A and B atoms

linearly uncorrelated (called principal or normal) components, such that the first component has the largest variance and each subsequent has the largest variance constrained to being orthogonal to all the preceding components



# Linear dimensionality reduction: Principal components

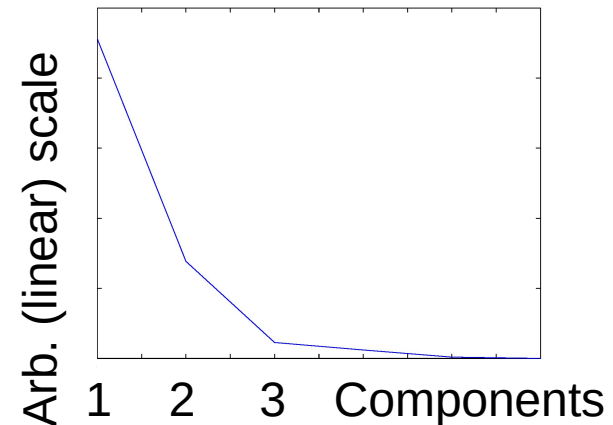
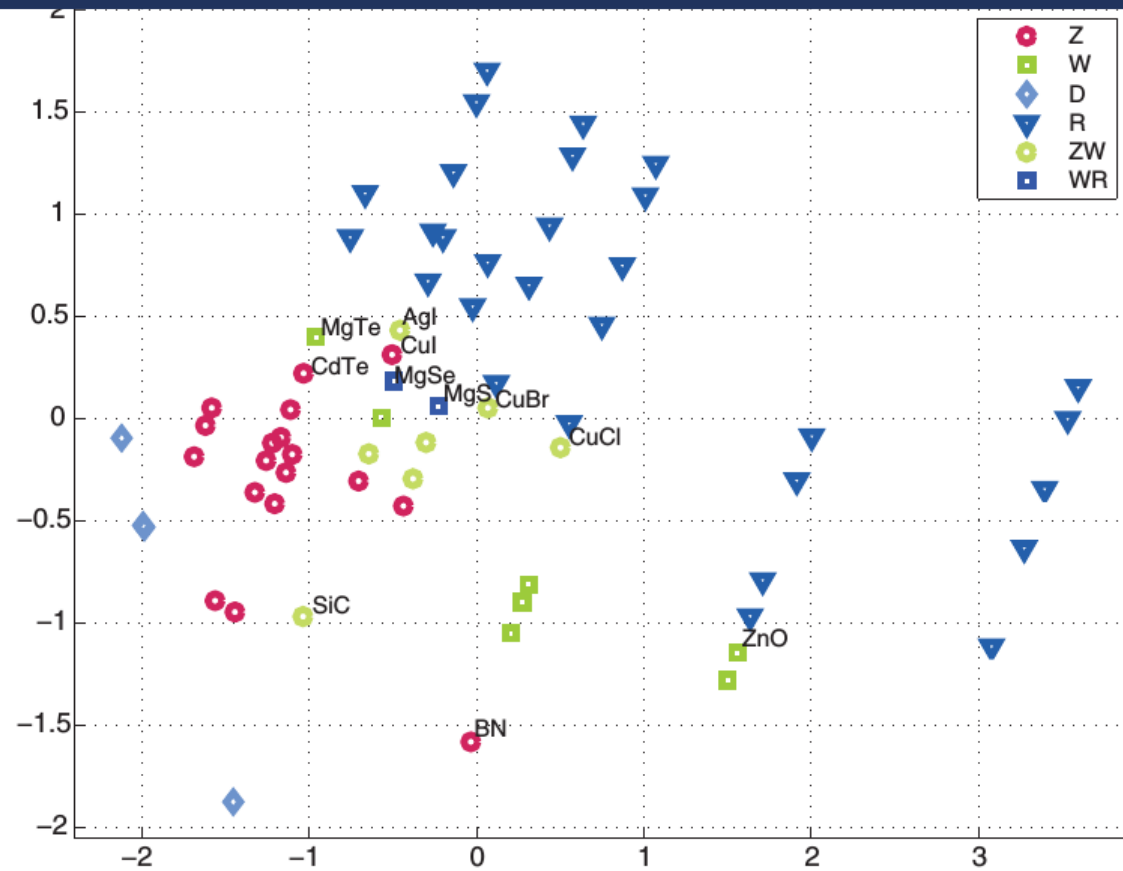
## Ansatz: atomic features

- Valence number  $Z_v$
- Energy of valence  $s$  orbital  $E_s$
- Energy of valence  $p$  orbital  $E_p$
- Radius of valence  $s$  orbital  $r_s$
- Radius of valence  $p$  orbital  $r_p$

$$r_s, r_p, E_s/\sqrt{Z_v}, E_p/\sqrt{Z_v},$$

for A and B atoms

linearly uncorrelated (called principal or normal) components, such that the first component has the largest variance and each subsequent has the largest variance constrained to being orthogonal to all the preceding components



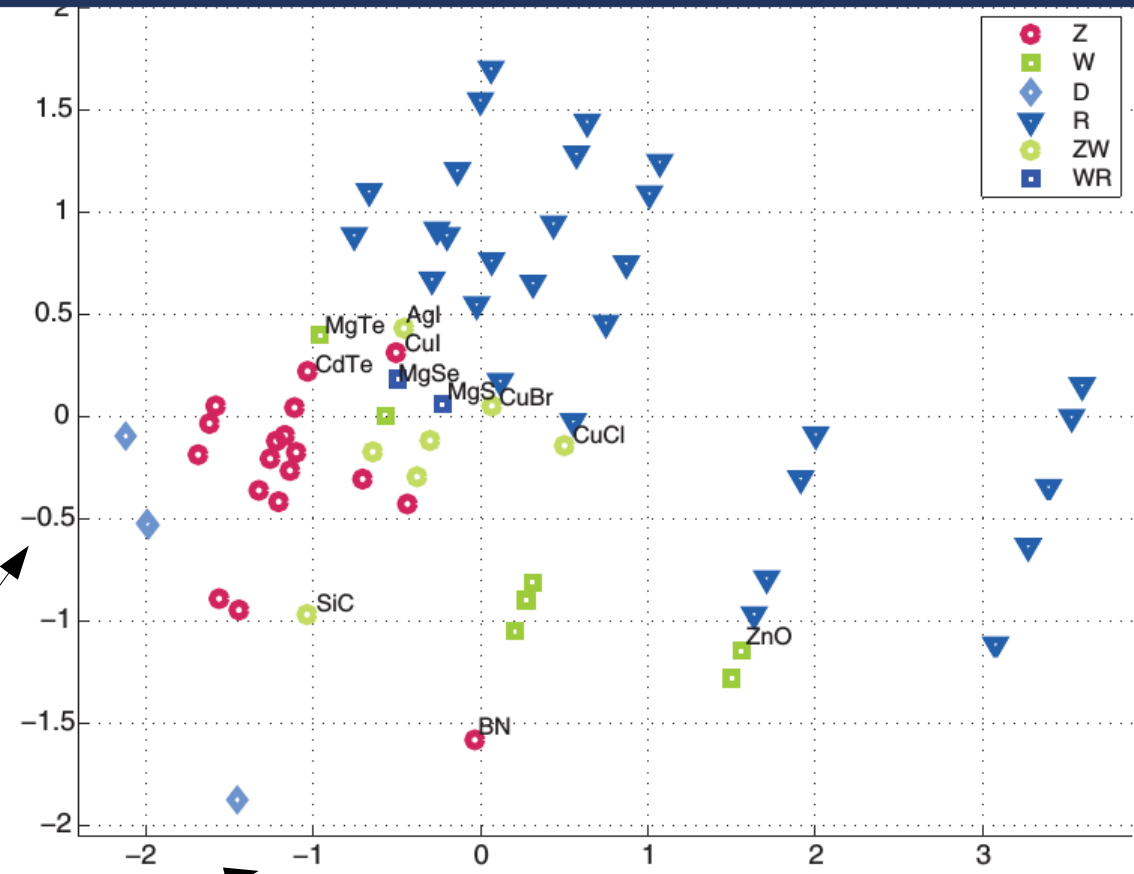
# Linear dimensionality reduction: Principal components

Ansatz: atomic features

- Valence number
- Energy of valence  $s$  orbital
- Energy of valence  $p$  orbital
- Radius of valence  $s$  orbital
- Radius of valence  $p$  orbital

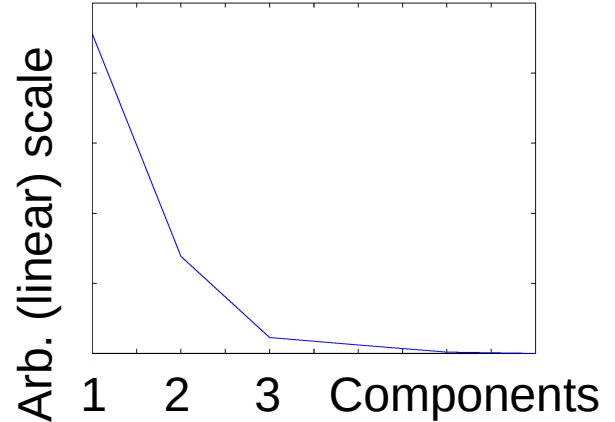
$Z_v$   
 $E_s$   
 $E_p$   
 $r_s$   
 $r_p$

$r_s, r_p, E_s/\sqrt{Z_v}, E_p/\sqrt{Z_v}$   
for A and B atoms



What's on the axes?

Linear combination of (possibly all) the initial dimensions



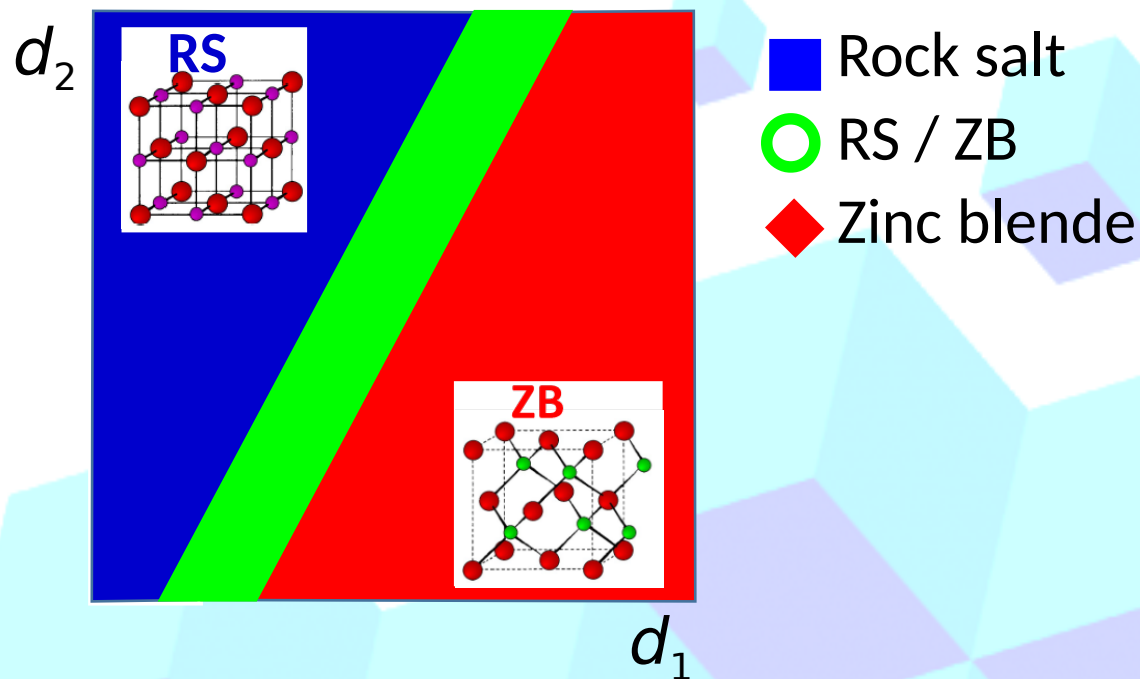


# Compressed sensing: the quest for descriptors and predictive models

## 82 octet AB binary compounds

Ansatz: atomic features

- HOMO
- LUMO
- Ionization Potential
- Electron Affinity
- Radius of valence  $s$  orbital
- Radius of valence  $p$  orbital
- Radius of valence  $d$  orbital
- Thousands to billions of non-linear functions of the above



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

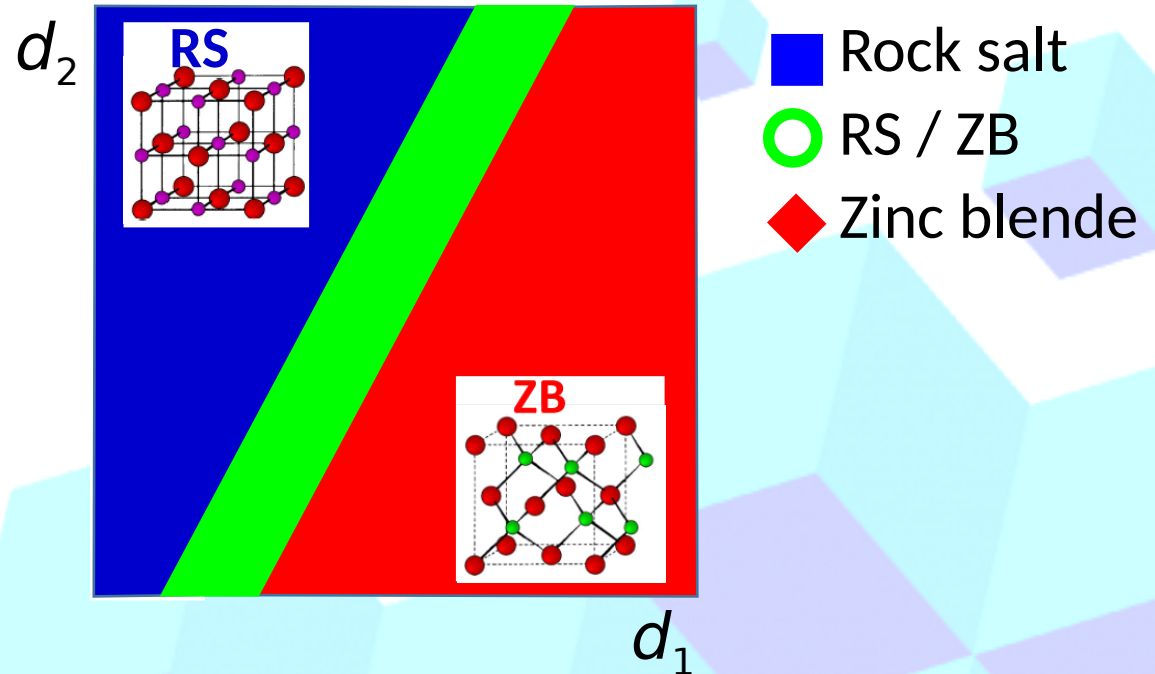
$E(\text{Rock salt}) - E(\text{Zinc blende})$

# Compressed sensing: the quest for descriptors and predictive models

## 82 octet AB binary compounds

Ansatz: atomic features

- HOMO
- LUMO
- Ionization Potential
- Electron Affinity
- Radius of valence  $s$  orbital
- Radius of valence  $p$  orbital
- Radius of valence  $d$  orbital
- Thousands to billions of non-linear functions of the above

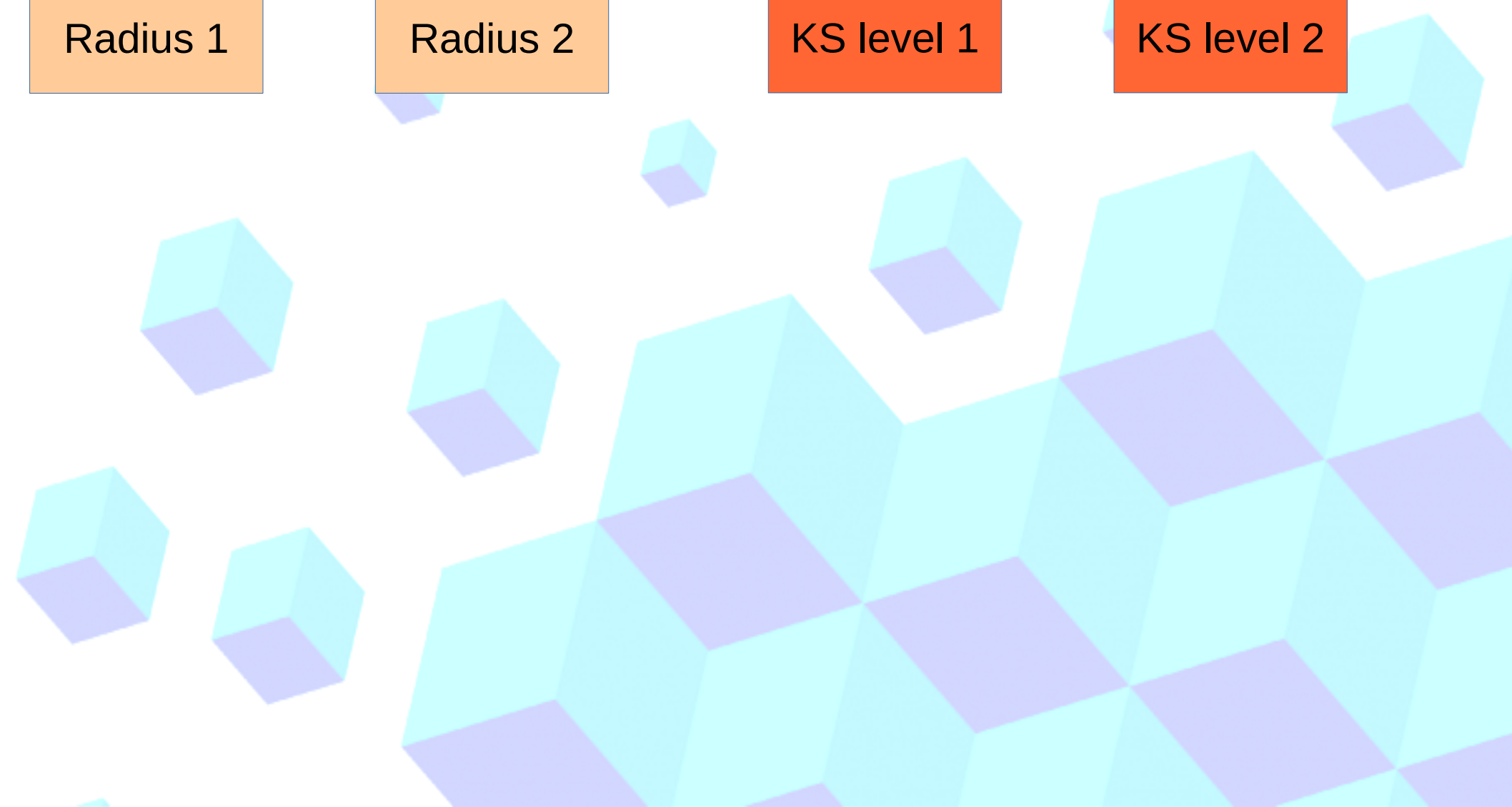
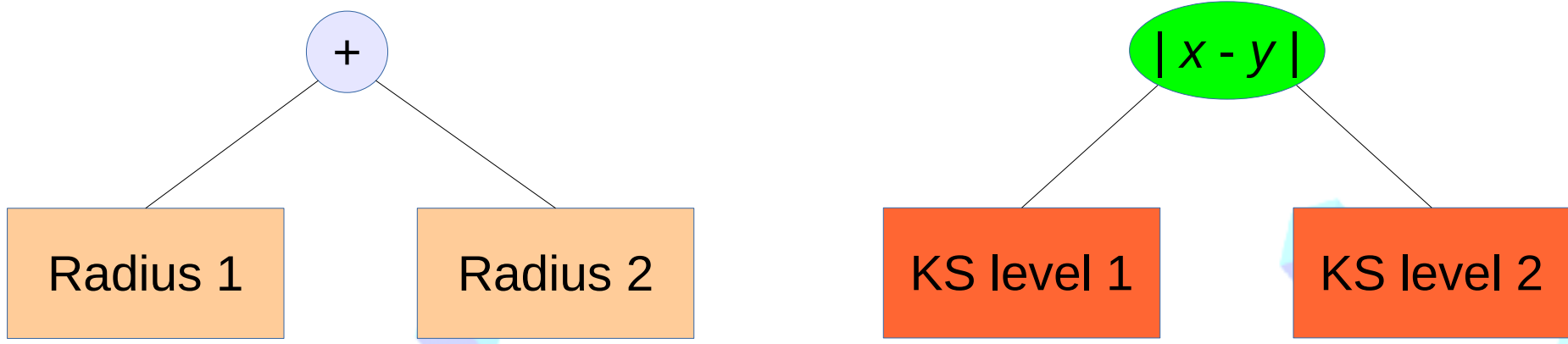


**Symbolic Regression**

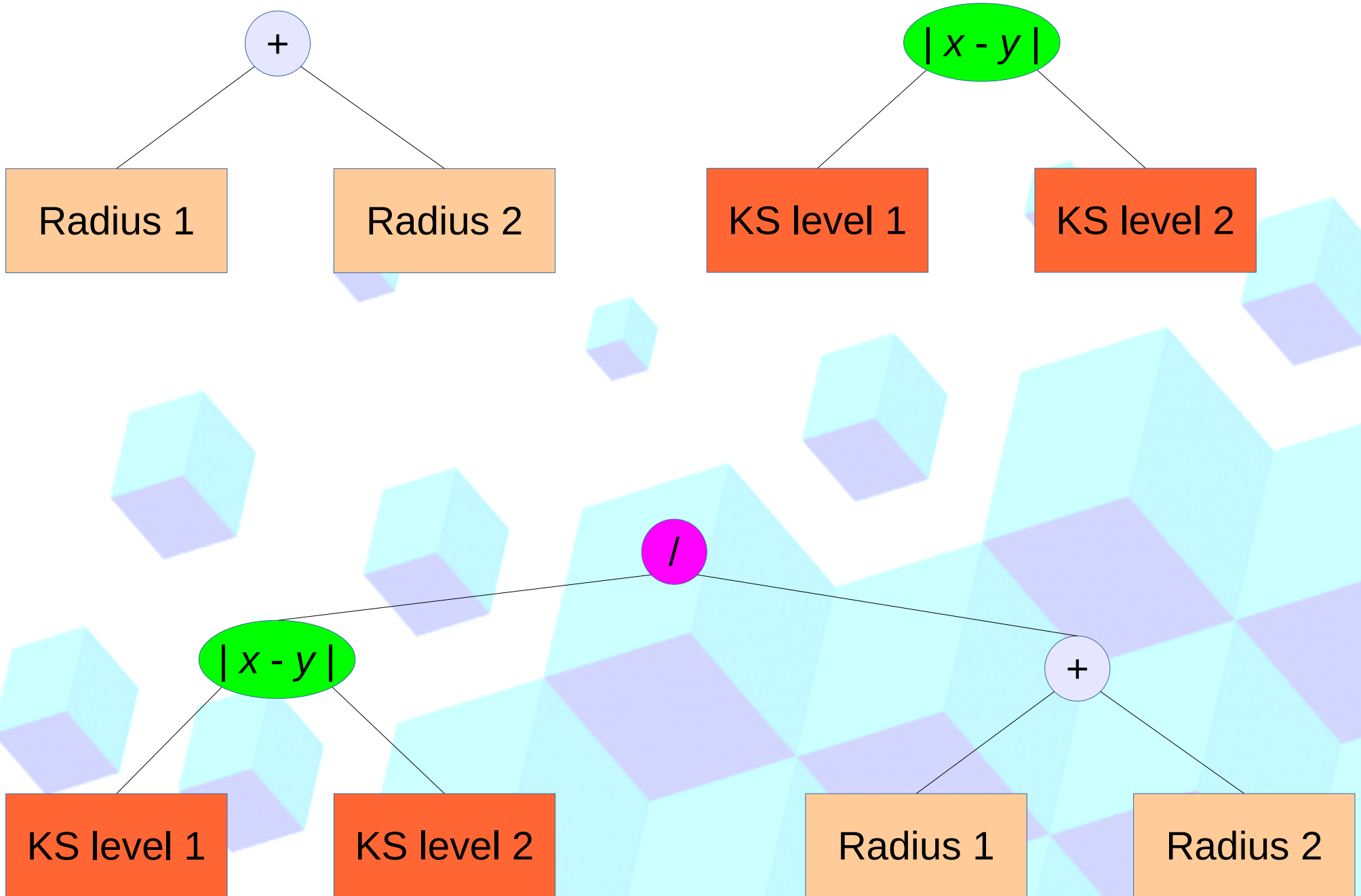
$$P = c_1 d_1 + c_2 d_2 + \dots c_n d_n$$

$E(\text{Rock salt}) - E(\text{Zinc blende})$

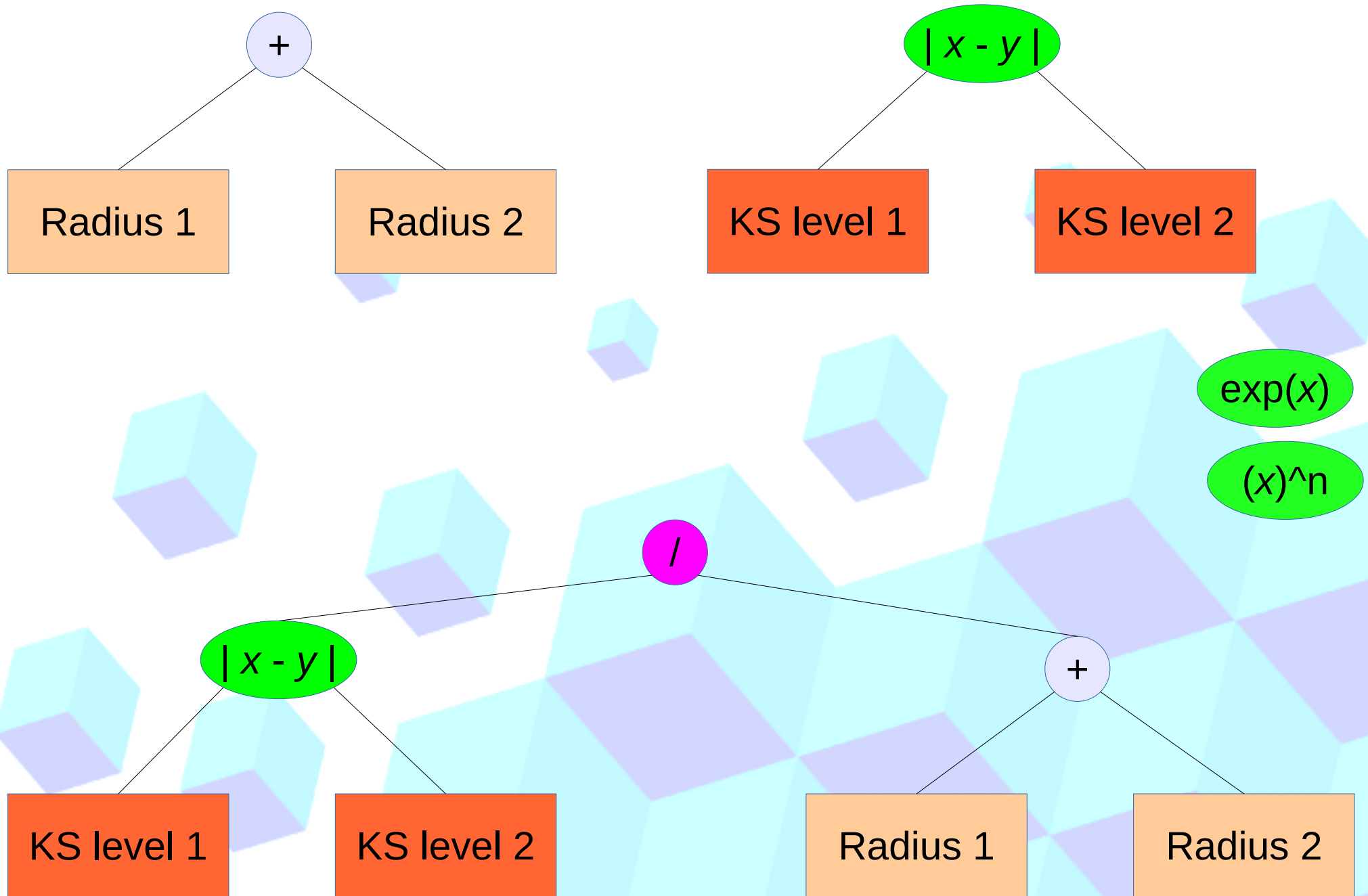
# Systematic construction of the feature space



# Systematic construction of the feature space



# Systematic construction of the feature space



# Systematic construction of the feature space: EUREQA

EUREQA: genetic programming software.  
Global optimization (genetic algorithm).  
Schmidt M., Lipson H., Science, Vol. 324, No. 5923, (2009)

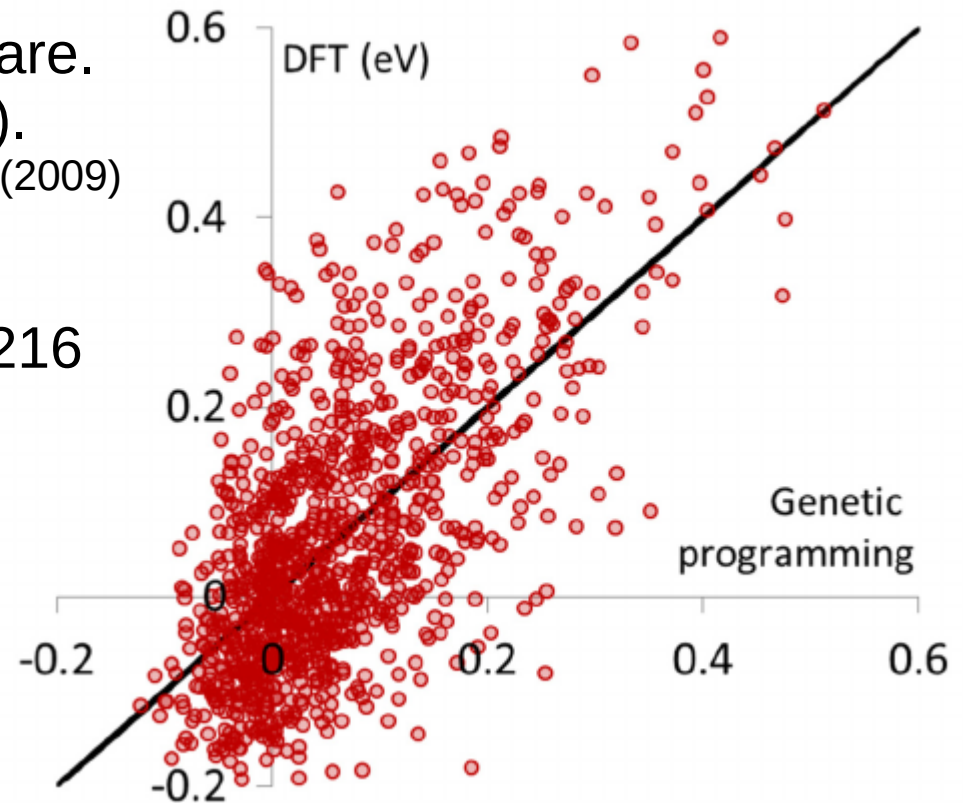
T. Müller et al. PRB **89** 115202 (2014):  
Data: ~1000 amorphous structures of 216  
Si atoms (saturated)

Property: hole trap depth

$$\frac{\min(1.66355, a) \max(5.37551, c) - f - bd}{g} - h \max(3.42929, e),$$

Descriptor (candidates: 242)

- a The largest distance between a H atom and its nearest Si neighbor
- b The shortest distance between a Si atom and its sixth-nearest Si neighbor
- c The maximum bond valence sum on a Si atom
- d The smallest value for the fifth-smallest relative bond length around a Si atom
- e The fourth-shortest distance between a Si atom and its eighth-nearest neighbor
- f The second-shortest distance between a Si atom and its fifth-nearest neighbor
- g The third-shortest distance between a Si atom and its sixth-nearest neighbor
- h The H-Si nearest-neighbor distance for the hydrogen atom with the fourth-smallest difference between the distances to the two Si atoms nearest to a H atom



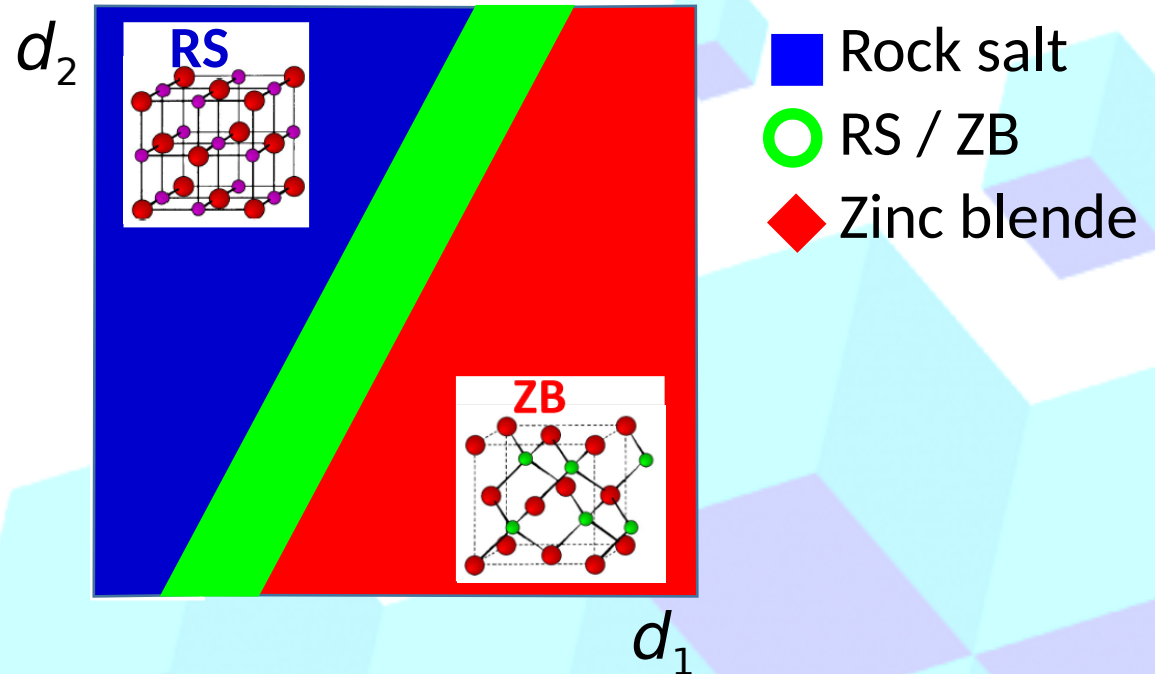
Building block	
Constant value	Exponential
Input variable	Natural logarithm
Addition	Power
Subtraction	Square root
Multiplication	Logistic function
Division	Minimum
Negation	Maximum
	Absolute value

# Compressed sensing: the quest for descriptors and predictive models

## 82 octet AB binary compounds

Ansatz: atomic features

- HOMO
- LUMO
- Ionization Potential
- Electron Affinity
- Radius of valence  $s$  orbital
- Radius of valence  $p$  orbital
- Radius of valence  $d$  orbital
- Thousands of non-linear functions of the above



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

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

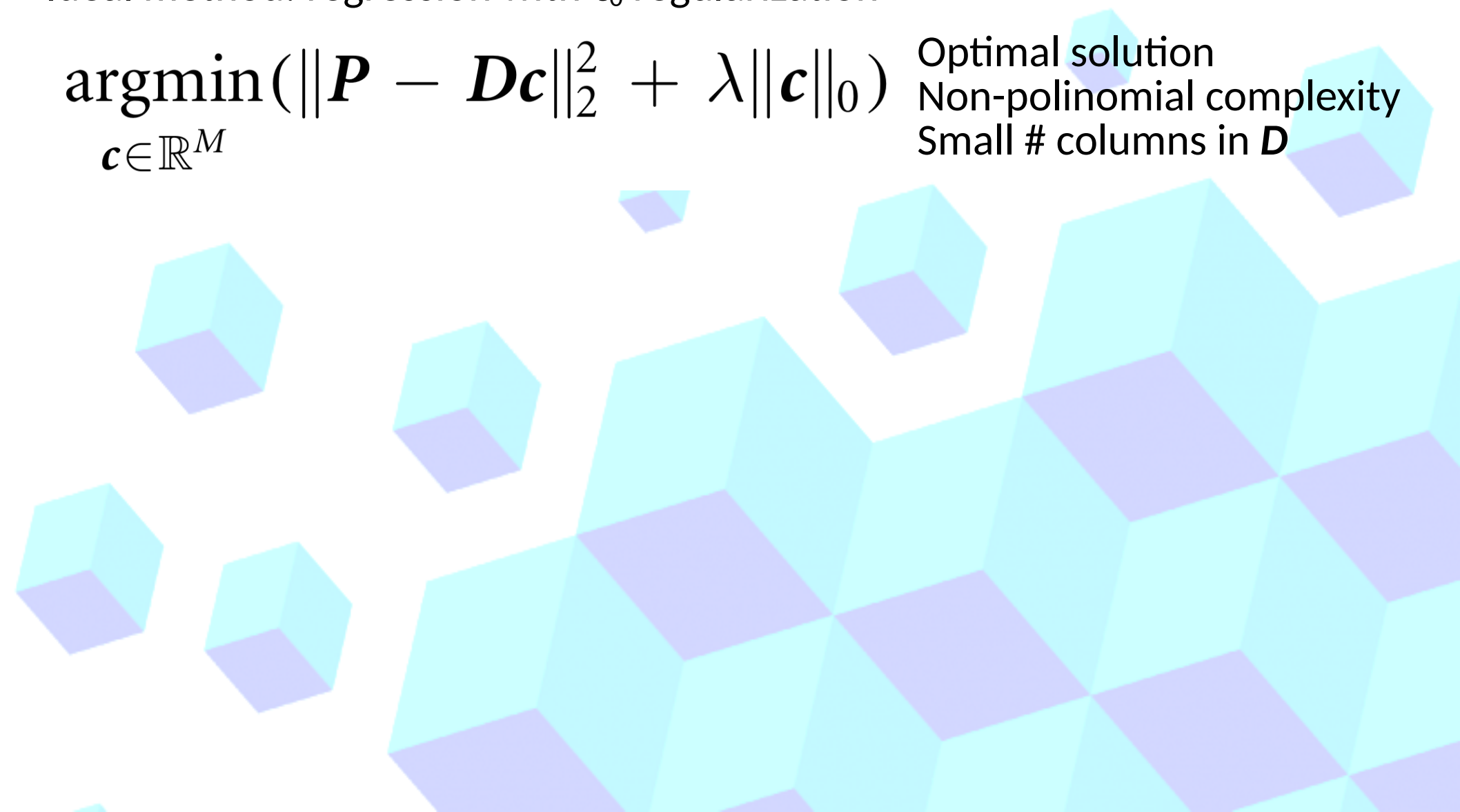
$E(\text{Rock salt}) - E(\text{Zinc blende})$

# Compressed sensing: the quest for descriptors and predictive models

Ideal method: regression with  $\ell_0$  regularization

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

Optimal solution  
Non-polynomial complexity  
Small # columns in  $\mathbf{D}$





# Compressed sensing: the quest for descriptors and predictive models

Ideal method: regression with  $\ell_0$  regularization

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

Optimal solution  
Non-polynomial complexity  
Small # columns in  $\mathbf{D}$

$\ \mathbf{c}\ _0$	# of nonzero elements of $\mathbf{c}$
$\ \mathbf{c}\ _2$	Euclidean. Square root of sum of squares of the elements of $\mathbf{c}$

# Compressed sensing: the quest for descriptors and predictive models

Ideal method: regression with  $\ell_0$  regularization

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

Optimal solution  
Non-polynomial complexity  
Small # columns in  $\mathbf{D}$

$\|\mathbf{c}\|_0$  # of nonzero elements of  $\mathbf{c}$   
 $\|\mathbf{c}\|_2$  Euclidean. Square root of sum of squares of the elements of  $\mathbf{c}$

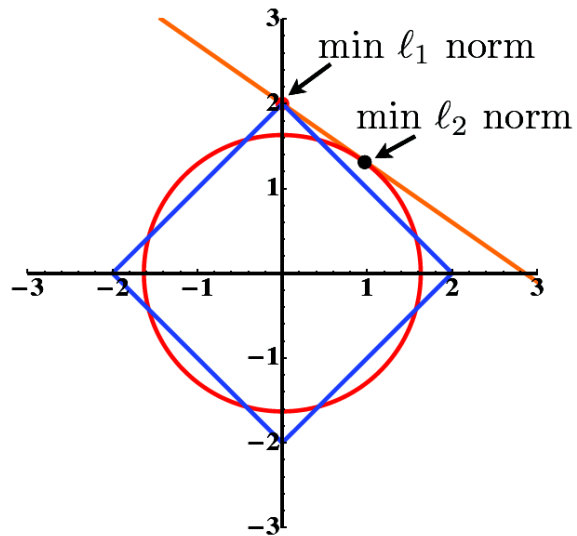
For matrices  $\mathbf{D}$  with uncorrelated columns: LASSO

$$\operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^M} \|\mathbf{P} - \mathbf{D}\mathbf{c}\|_2^2 + \lambda \|\mathbf{c}\|_1$$

(Possibly) optimal solution  
Convex optimization  
Moderate # columns in  $\mathbf{D}$

$\|\mathbf{c}\|_1$  “Manhattan”. Sum of absolute values of the elements of  $\mathbf{c}$

# Compressed sensing: the quest for descriptors and predictive models



$$\operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^M} \|\mathbf{P} - \mathbf{D}\mathbf{c}\|_2^2 + \lambda \|\mathbf{c}\|_1$$

(Possibly) optimal solution  
Convex optimization  
Moderate # columns in  $\mathbf{D}$

$\|\mathbf{c}\|_1$  “Manhattan”. Sum of absolute values of the elements of  $\mathbf{c}$

## Lattice Anharmonicity and Thermal Conductivity from Compressive Sensing of First-Principles Calculations

Fei Zhou (周非)

*Physical and Life Sciences Directorate, Lawrence Livermore National Laboratory, Livermore, California 94550, USA*

Weston Nielson, Yi Xia, and Vidvuds Ozoliņš

*Department of Materials Science and Engineering, University of California, Los Angeles, California 90095-1595, USA*

(Received 22 April 2014; published 27 October 2014)



## Compressed modes for variational problems in mathematics and physics

Vidvuds Ozoliņš<sup>a,\*</sup>, Rongjie Lai<sup>b,1</sup>, Russel Caflisch<sup>c,1</sup>, and Stanley Osher<sup>c,1,2</sup>

Departments of <sup>a</sup>Materials Science and Engineering, and <sup>c</sup>Mathematics, University of California, Los Angeles, CA 90095-1555; and <sup>b</sup>Department of Mathematics, University of California, Irvine, CA 92697-3875

Contributed by Stanley Osher, October 8, 2013 (sent for review September 3, 2013)

PHYSICAL REVIEW B 87, 035125 (2013)

## Compressive sensing as a paradigm for building physics models

Lance J. Nelson and Gus L. W. Hart

*Department of Physics and Astronomy, Brigham Young University, Provo, Utah 84602, USA*

Fei Zhou (周非) and Vidvuds Ozoliņš\*

*Department of Materials Science and Engineering, University of California, Los Angeles, California 90095, USA*

(Received 26 June 2012; revised manuscript received 26 September 2012; published 18 January 2013)

## Compressive sensing as a paradigm for building physics models

Lance J. Nelson and Gus L. W. Hart

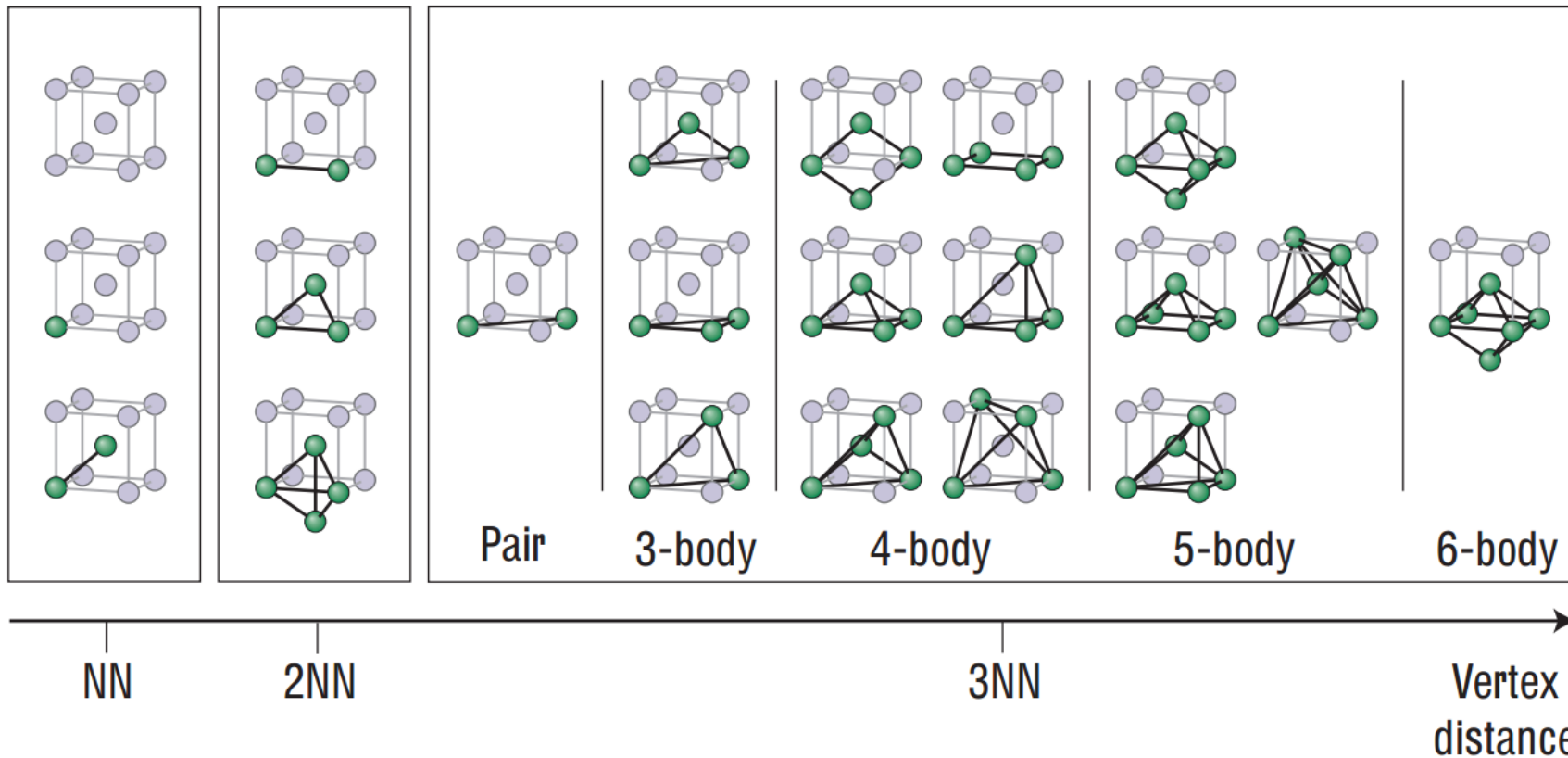
*Department of Physics and Astronomy, Brigham Young University, Provo, Utah 84602, USA*

Fei Zhou (周非) and Vidvuds Ozoliņš\*

*Department of Materials Science and Engineering, University of California, Los Angeles, California 90095, USA*

(Received 26 June 2012; revised manuscript received 26 September 2012; published 18 January 2013)

$$E(\sigma) = E_0 + \sum_f \bar{\Pi}_f(\sigma) J_f$$

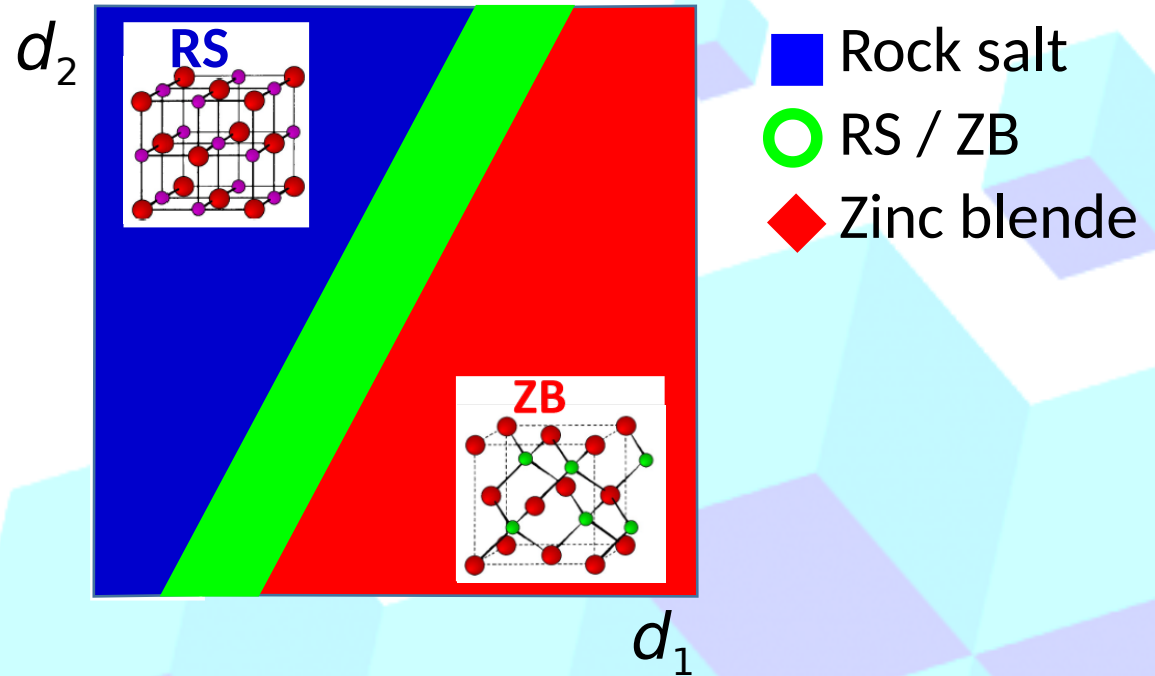


# Compressed sensing: the quest for descriptors and predictive models

## 82 octet AB binary compounds

Ansatz: atomic features

- HOMO
- LUMO
- Ionization Potential
- Electron Affinity
- Radius of valence  $s$  orbital
- Radius of valence  $p$  orbital
- Radius of valence  $d$  orbital
- Billions of non-linear functions of the above



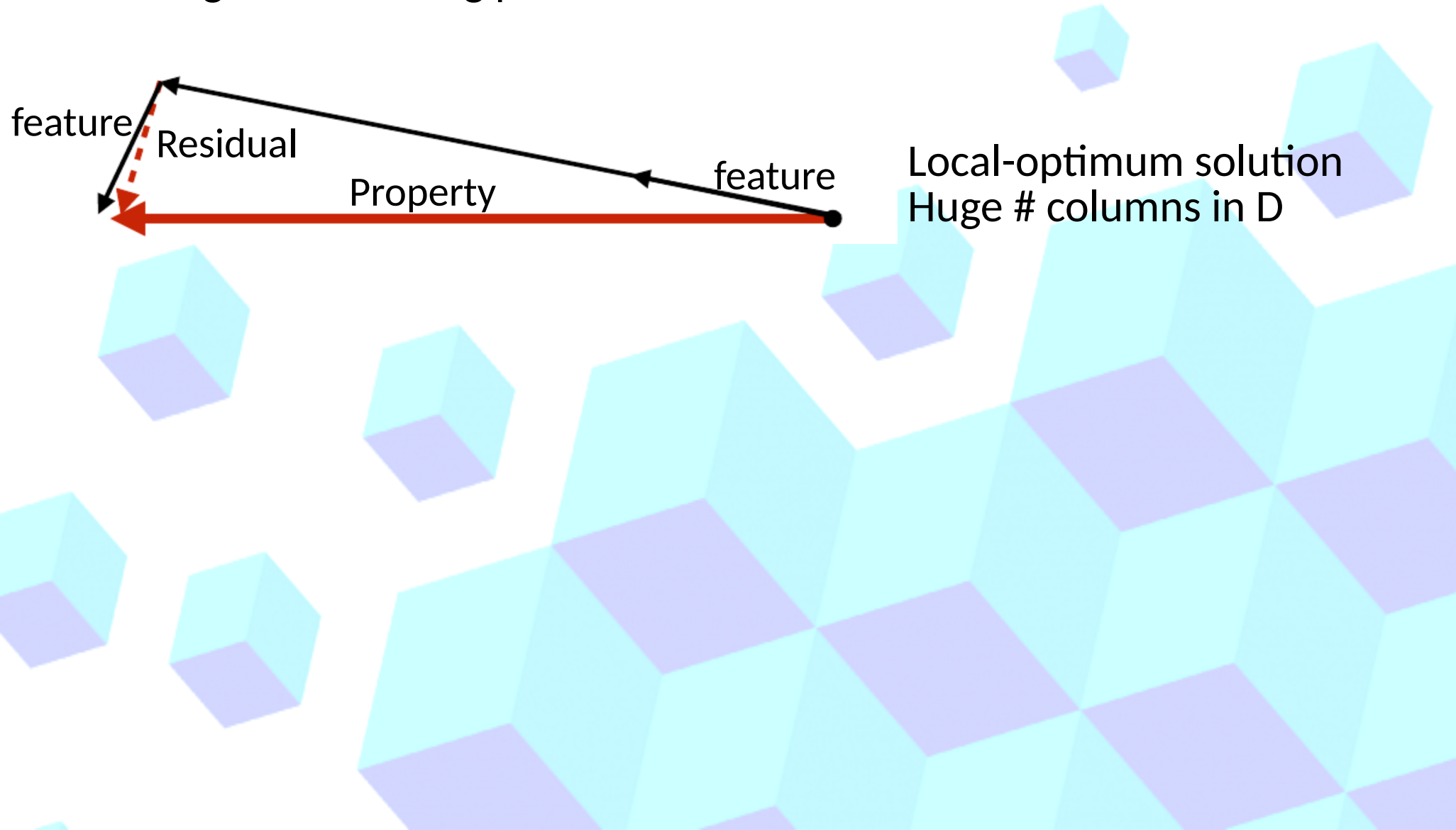
$$P = c_1 d_1 + c_2 d_2 + \dots c_n d_n$$

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

$E(\text{Rock salt}) - E(\text{Zinc blende})$

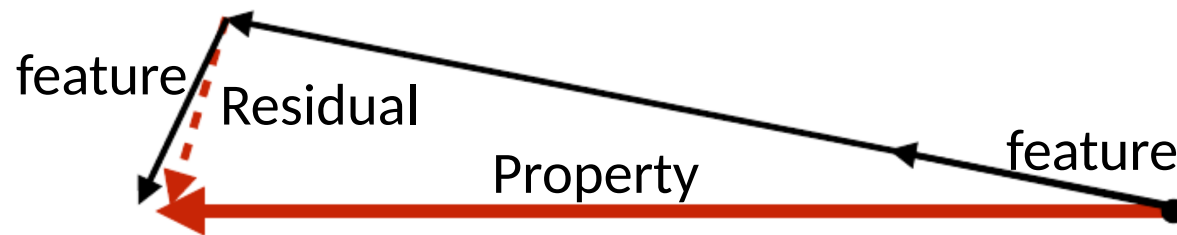
# Compressed sensing: the quest for descriptors and predictive models

From orthogonal matching pursuit ....



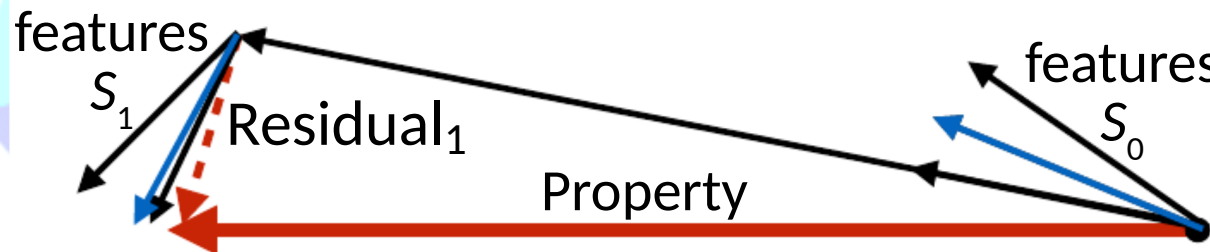
# Compressed sensing: the quest for descriptors and predictive models

From orthogonal matching pursuit ....



Local-optimum solution  
Huge # columns in D

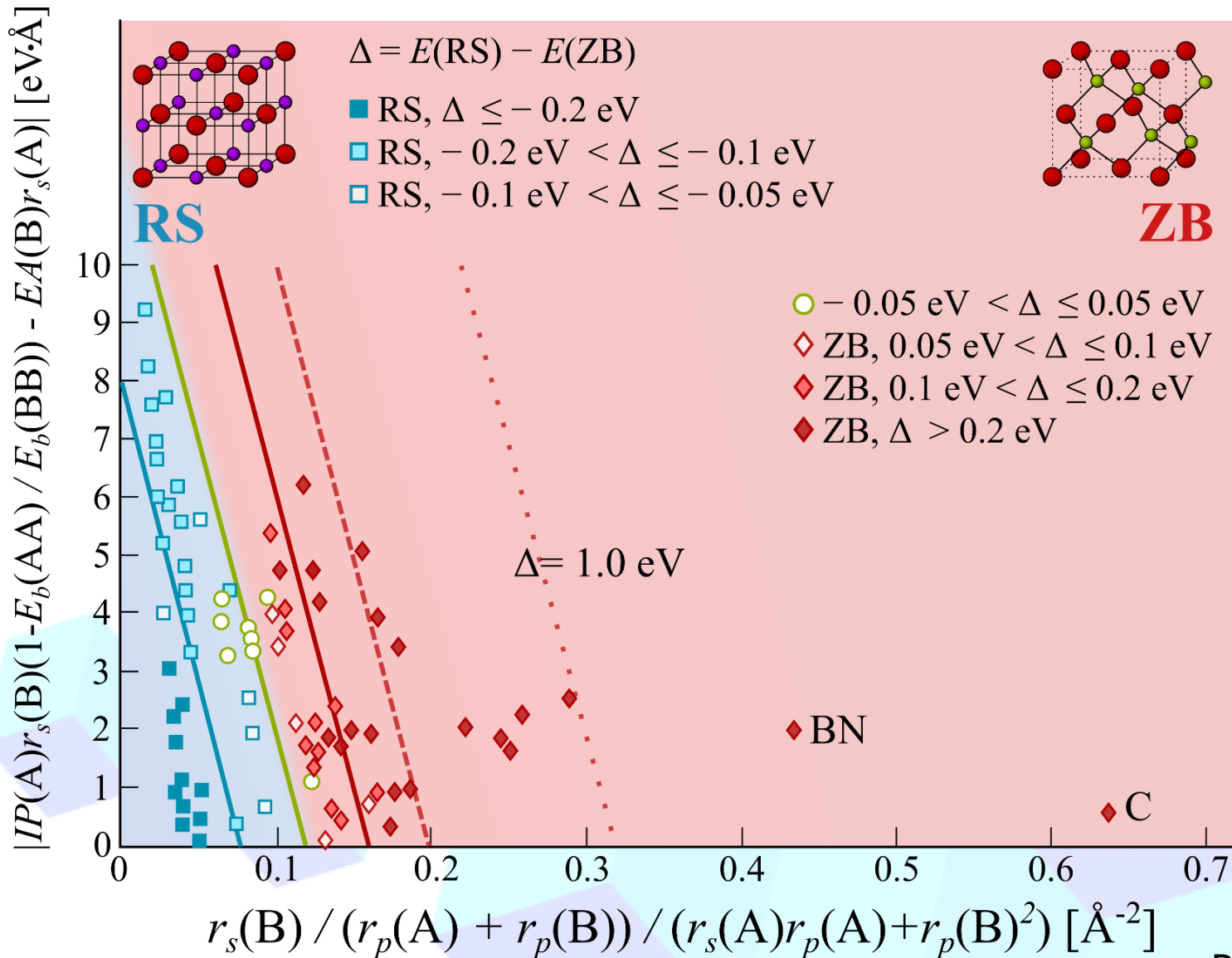
... to Sure Independence Screening + Sparsifying Operator (SISSO)



Proxy of  
global-optimum solution  
Huge # columns in D



# Compressed sensing: the quest for descriptors and predictive models



Structure map with SISSO, starting from 7 atomic + 6 dimer features  
 Feature space:  $10^{11}$  features

# Compressed sensing: the quest for descriptors and predictive models

$$\arg \min_{\mathbf{c}} (\|\mathbf{P} - \mathbf{D}\mathbf{c}\|_2^2 + \lambda\|\mathbf{c}\|_0)$$

Compressed-sensing-based model identification:  
Shares concepts with

Regularized regression. But: Massive sparsification.

Dimensionality reduction. But supervised, and yielding sparse, “inspectable” descriptors

Feature/Basis-set selection/extraction. But: non-greedy solver.

Symbolic regression. But: deterministic solver.

# One descriptor to rule them all: Multi-task learning

$$\{P^{(1)}, P^{(2)}, \dots, P^{N^T}\}$$



# One descriptor to rule them all: Multi-task learning

$$\{P^{(1)}, P^{(2)}, \dots, P^{N^T}\} \longrightarrow P^k = \mathbf{d} \cdot \mathbf{c}^k$$



# One descriptor to rule them all: Multi-task learning

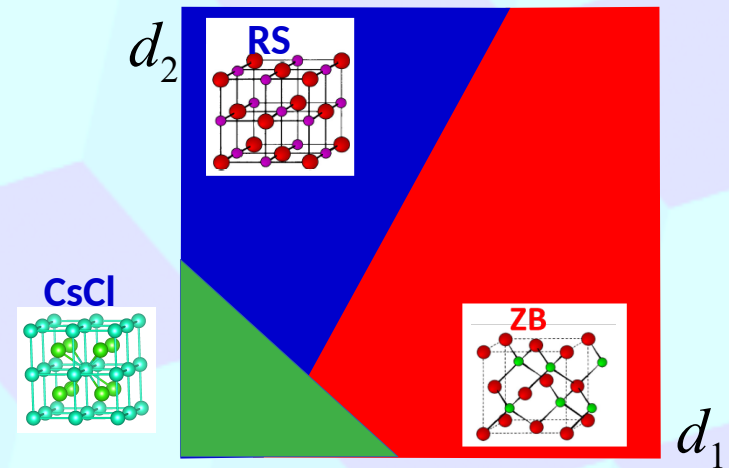
$$\{P^{(1)}, P^{(2)}, \dots, P^{N^T}\} \longrightarrow P^k = \mathbf{d} \cdot \mathbf{c}^k$$



# One descriptor to rule them all: Multi-task learning

$$\{P^{(1)}, P^{(2)}, \dots, P^{N^T}\} \longrightarrow P^k = \mathbf{d} \cdot \mathbf{c}^k$$

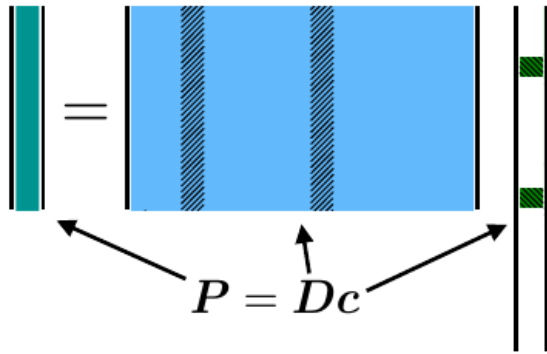
Application: multi-phase stability diagram  
Properties: crystal-structure formation energies



# One descriptor to rule them all: Multi-task learning

$$\{P^{(1)}, P^{(2)}, \dots, P^{N^T}\} \longrightarrow P^k = \mathbf{d} \cdot \mathbf{c}^k$$

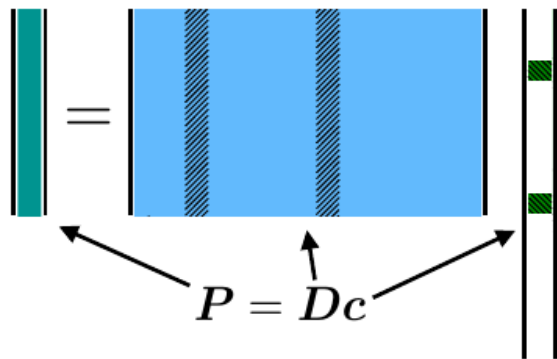
ST-SISSO



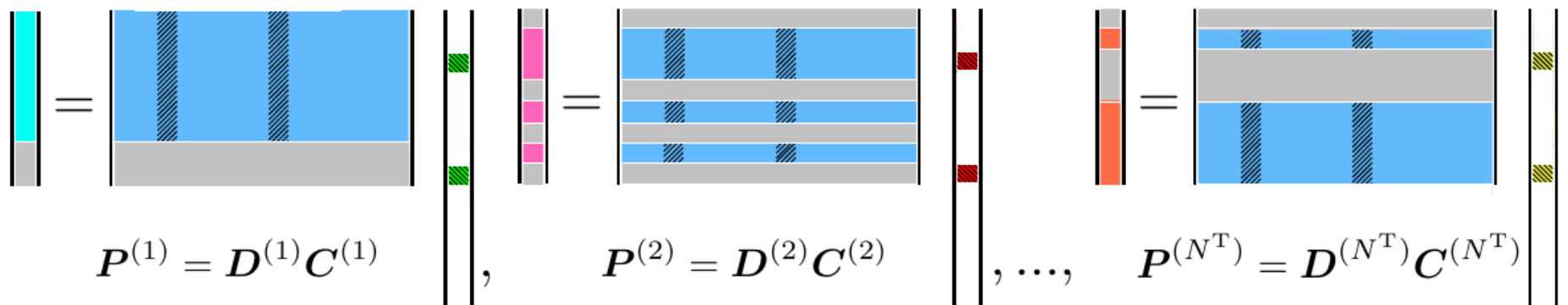
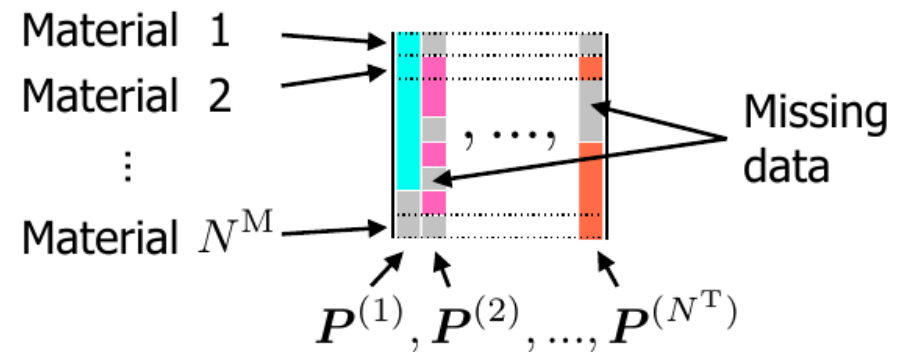
# One descriptor to rule them all: Multi-task learning

$$\{P^{(1)}, P^{(2)}, \dots, P^{N^T}\} \longrightarrow P^k = \mathbf{d} \cdot \mathbf{c}^k$$

ST-SISSO



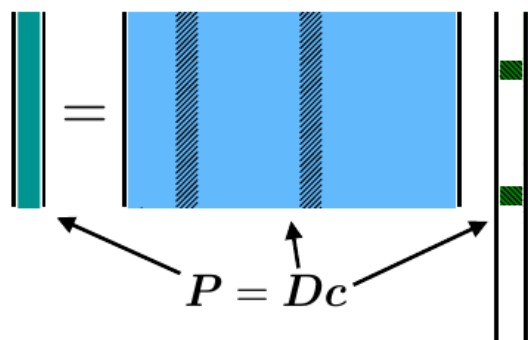
MT-SISSO



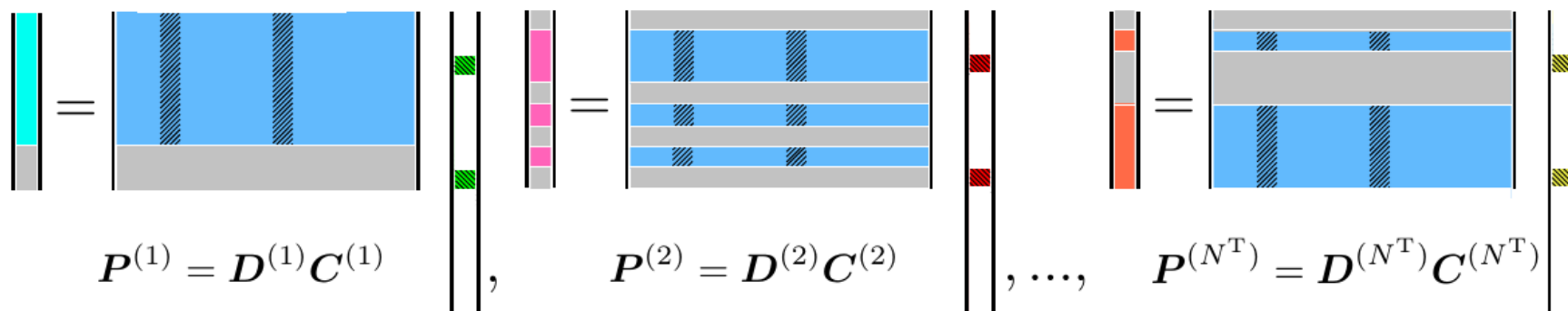
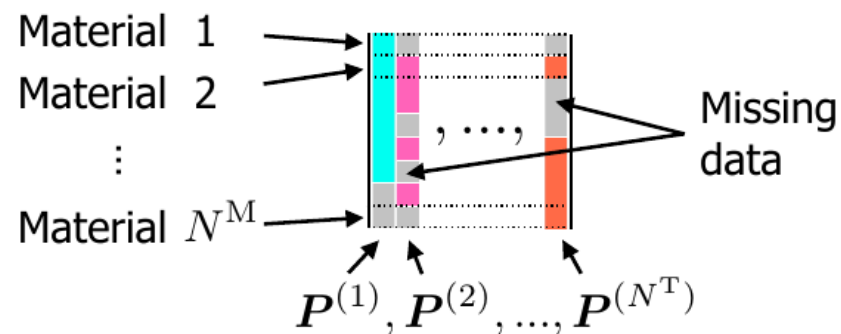


# One descriptor to rule them all: Multi-task learning

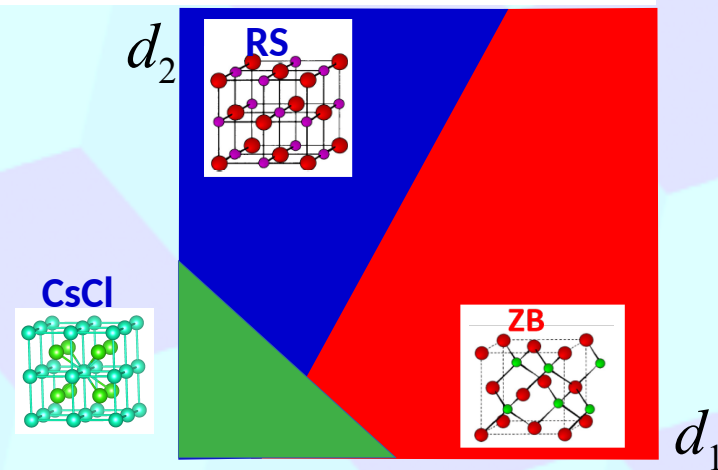
ST-SISSO



MT-SISSO



Application: multi-phase stability diagram  
Properties: crystal-structure formation energies



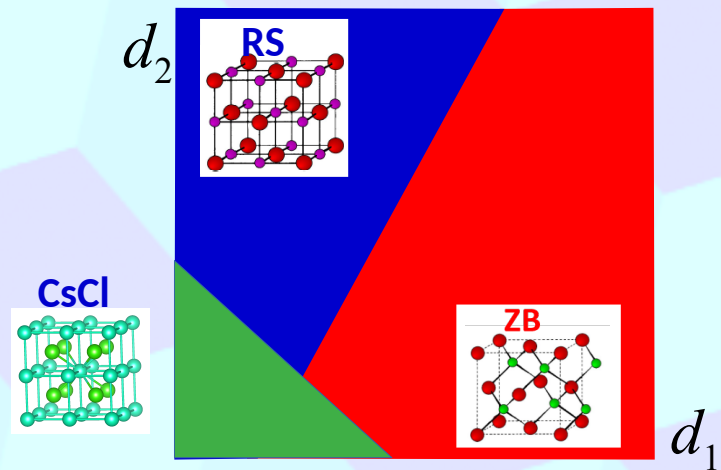
# One descriptor to rule them all: Multi-task learning

$$\{P^{(1)}, P^{(2)}, \dots, P^{N^T}\} \longrightarrow P^k = \mathbf{d} \cdot \mathbf{c}^k$$

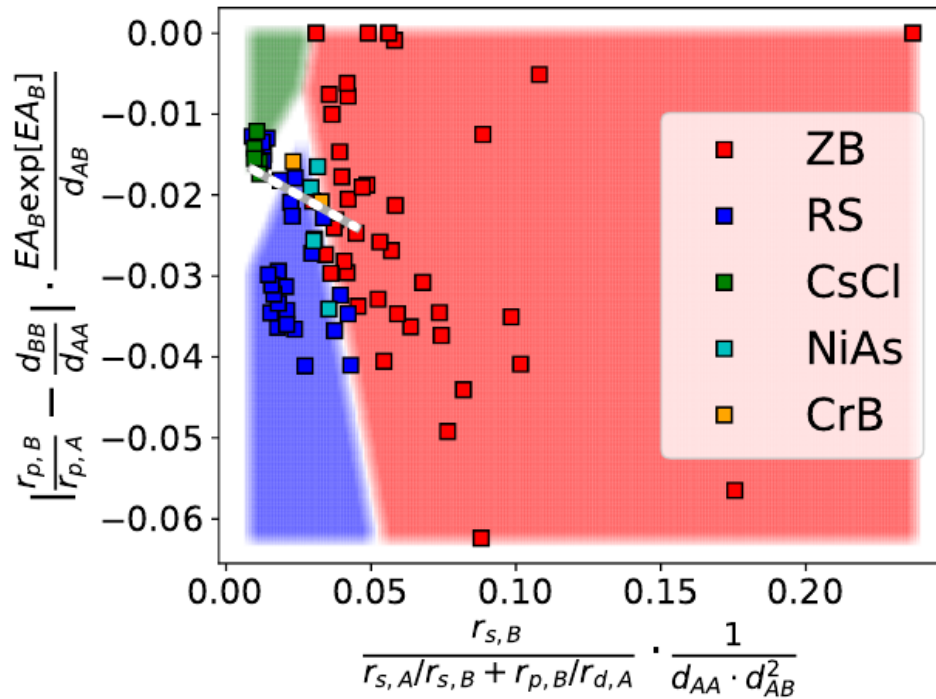
$$\arg \min_{\mathbf{c}} (\|\mathbf{P} - \mathbf{D}\mathbf{c}\|_2^2 + \lambda \|\mathbf{c}\|_0)$$

$$\arg \min_{\mathbf{C}} \sum_{k=1}^{N^T} \frac{1}{N_k^M} \|\mathbf{P}^k - \mathbf{D}^k \mathbf{C}^k\|_2^2 + \lambda \|\mathbf{C}\|_0$$

Application: multi-phase stability diagram  
Properties: crystal-structure formation energies



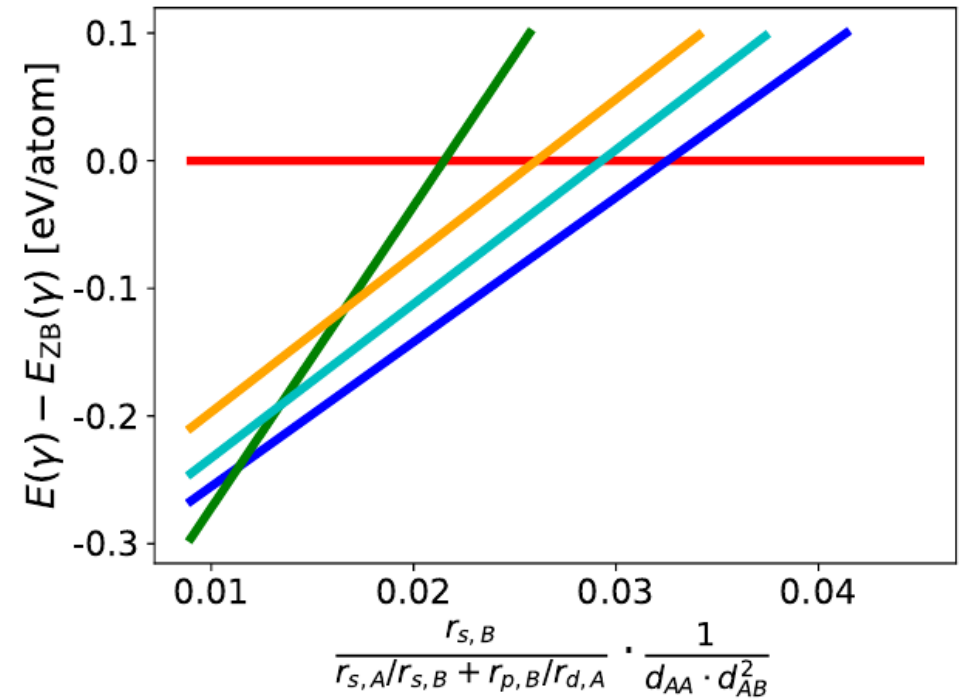
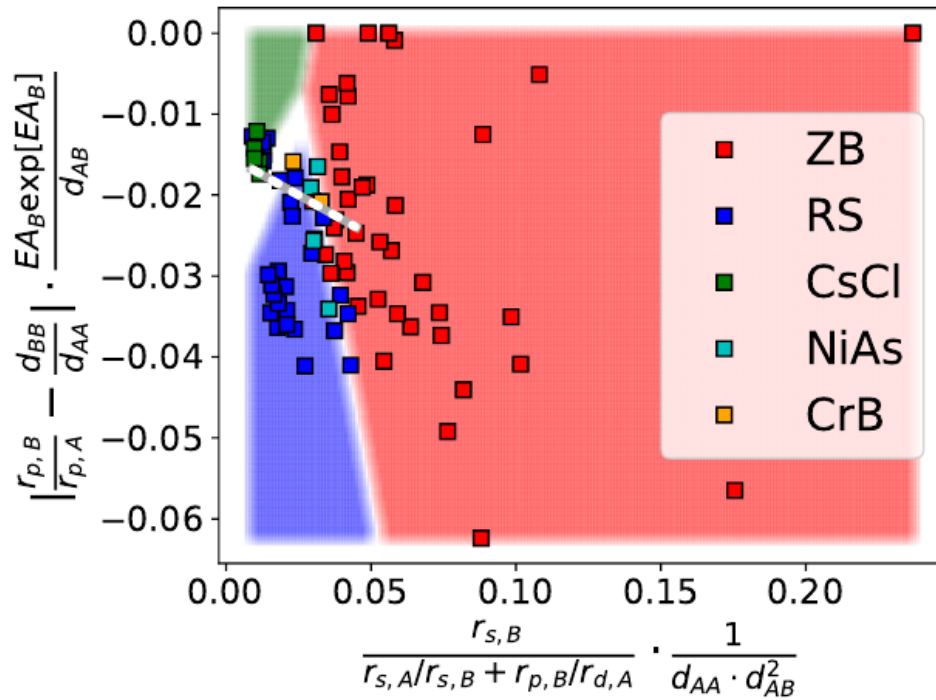
# One descriptor to rule them all: Multi-task SISSO Energy differences among 5 crystal structures.



Ouyang, Ahmetcik, Carbogno, Scheffler & LMG, J. Phys. Mater. 2, 024002 (2019)

# One descriptor to rule them all: Multi-task SISSO

## Energy differences among 5 crystal structures.



# A general scheme about training, cross-validation, test

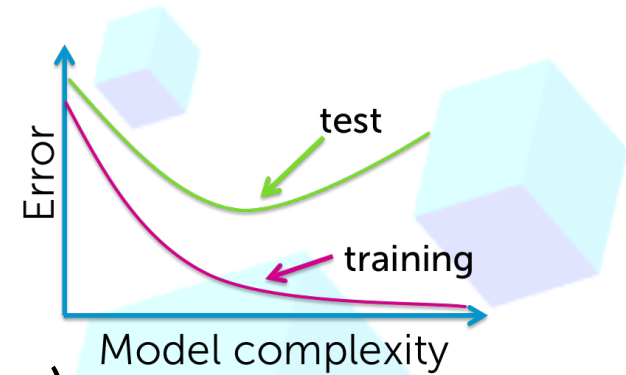
**Training:** input (features, descriptor) + labels (values target property)  
→ yields one model which minimizes a cost function (incl. regularization)



# A general scheme about training, cross-validation, test

**Training:** input (features, descriptor) + labels (values target property)  
→ yields one model which minimizes a cost function (incl. regularization)

**Cross-validation:** used to tune model-complexity  
- perform training  $n$  times on different split of data.  
Training + test/validation sets  
→ yields one model that minimizes the test (validation) error



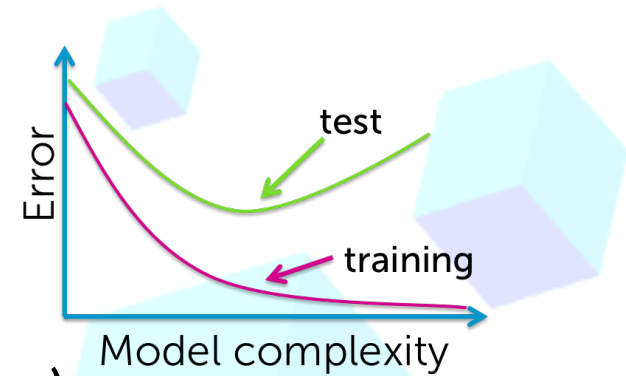
# A general scheme about training, cross-validation, test

**Training:** input (features, descriptor) + labels (values target property)  
→ yields one model which minimizes a cost function (incl. regularization)

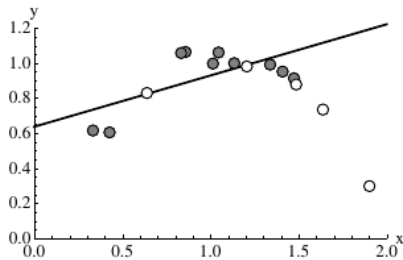
**Cross-validation:** used to tune model-complexity  
- perform training  $n$  times on different split of data.

Training + test/validation sets

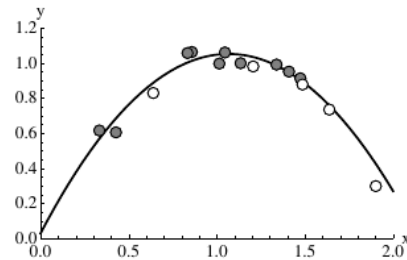
→ yields one model that minimizes the test (validation) error



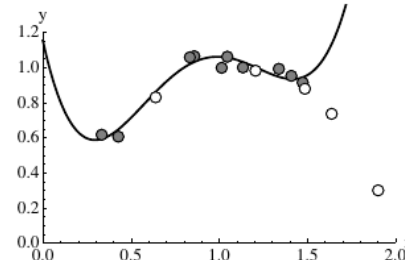
Underfitting



Fitting



Overfitting



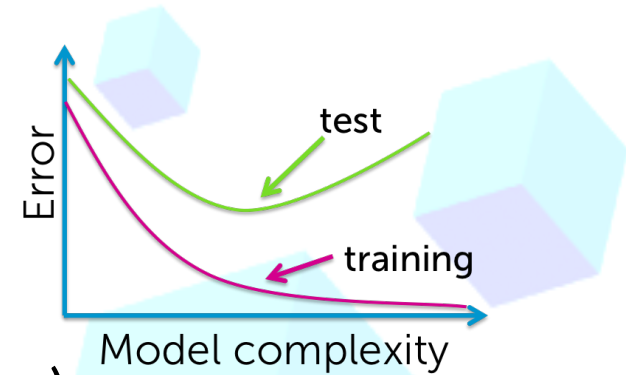
# A general scheme about training, cross-validation, test

**Training:** input (features, descriptor) + labels (values target property)  
→ yields one model which minimizes a cost function (incl. regularization)

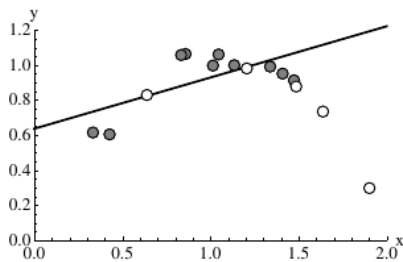
**Cross-validation:** used to tune model-complexity  
- perform training  $n$  times on different split of data.

Training + test/validation sets

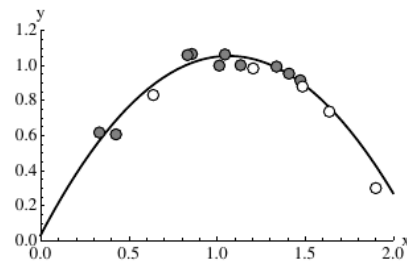
→ yields one model that minimizes the test (validation) error



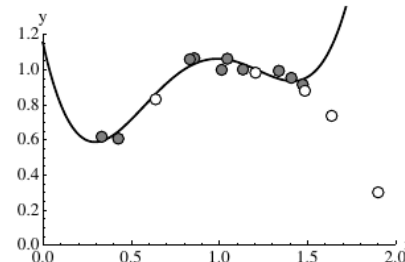
Underfitting



Fitting



Overfitting

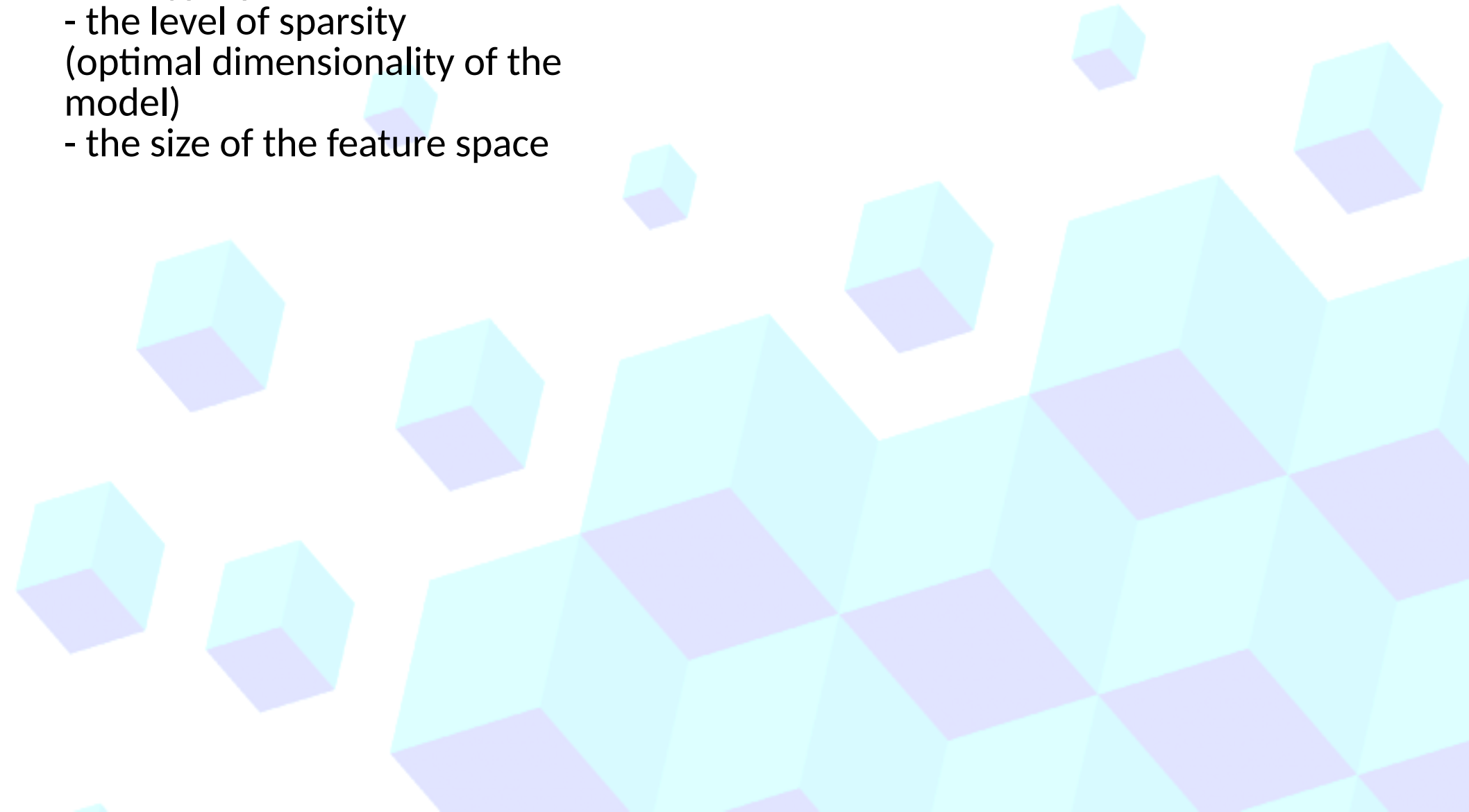


**Test:** evaluation of the performance of the model on data never used for training (i.e., the whole cross-validation procedure), aka left-out set



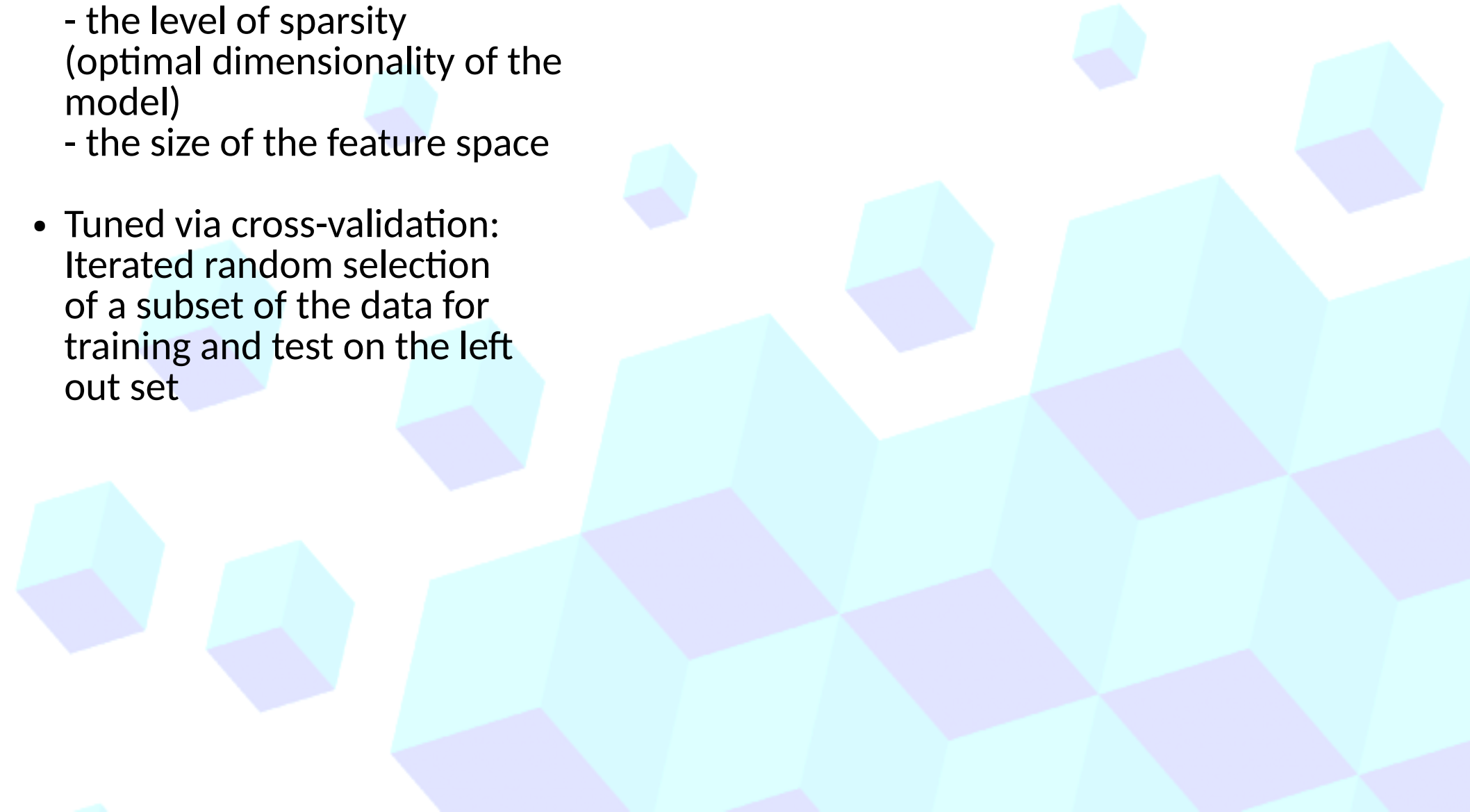
# Data-driven model complexity

- In compressed sensing the “hyperparameters” are
  - the level of sparsity (optimal dimensionality of the model)
  - the size of the feature space



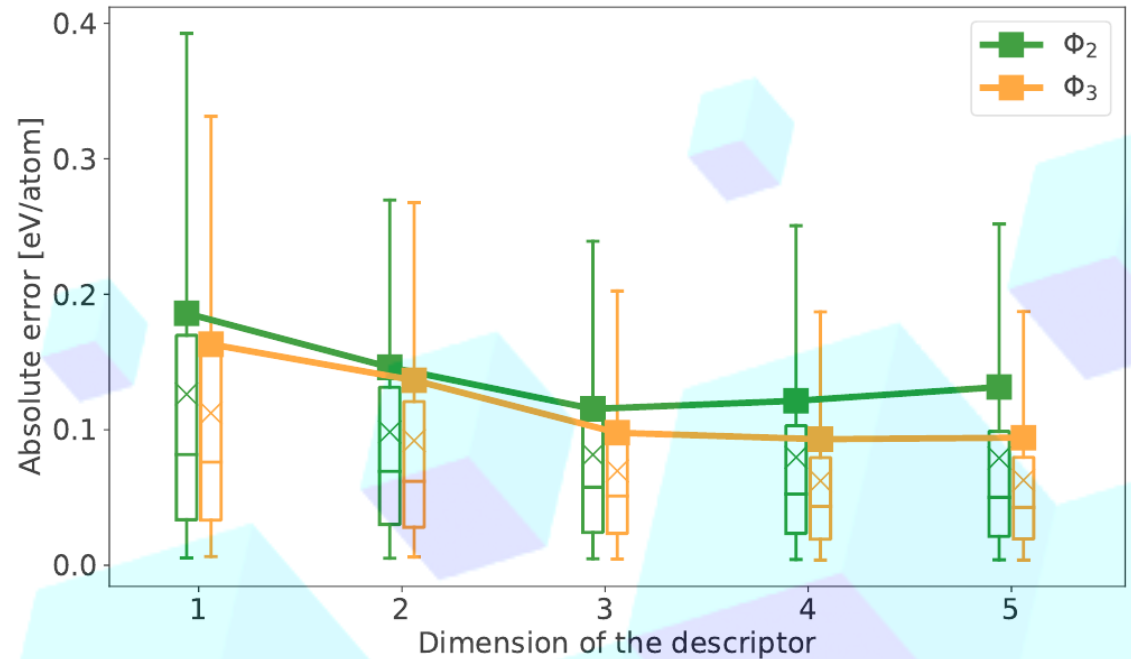
# Data-driven model complexity

- In compressed sensing the “hyperparameters” are
  - the level of sparsity (optimal dimensionality of the model)
  - the size of the feature space
- Tuned via cross-validation: Iterated random selection of a subset of the data for training and test on the left out set



# Data-driven model complexity

- In compressed sensing the “hyperparameters” are
  - the level of sparsity (optimal dimensionality of the model)
  - the size of the feature space
- Tuned via cross-validation: Iterated random selection of a subset of the data for training and test on the left out set



# Charts/maps of materials

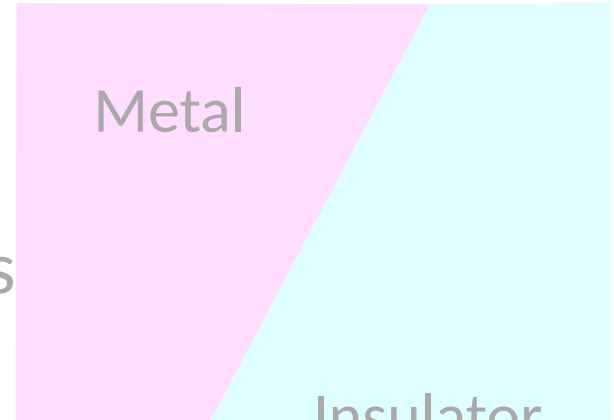
$$\operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^M} (\underbrace{\|\mathbf{P} - \mathbf{D}\mathbf{c}\|_2^2}_{\text{overlap of convex domains}} + \lambda \|\mathbf{c}\|_0)$$

New cost function to be minimized:  
overlap of *convex* domains

1. # points in the *convex* overlap domain
2. Area of the domain overlap
3. Distance between domains

Good also for multi-categorical problems  
(see A. F. Bialon *et al.*, Chem. Mater. **28**, 2550 (2016))

$d'_2$   
 $A_x B_y$   
binaries

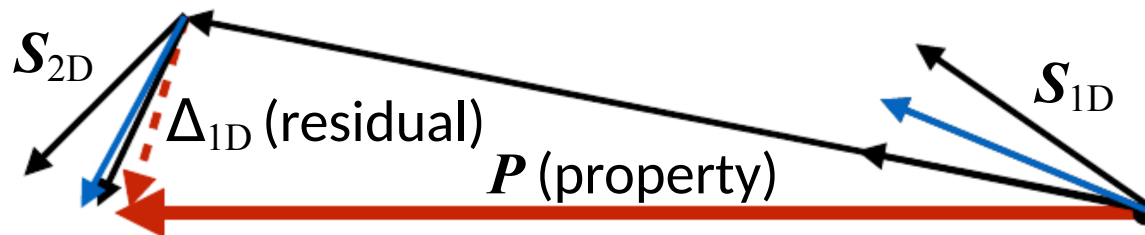


Insulator

$d'_1$

materials

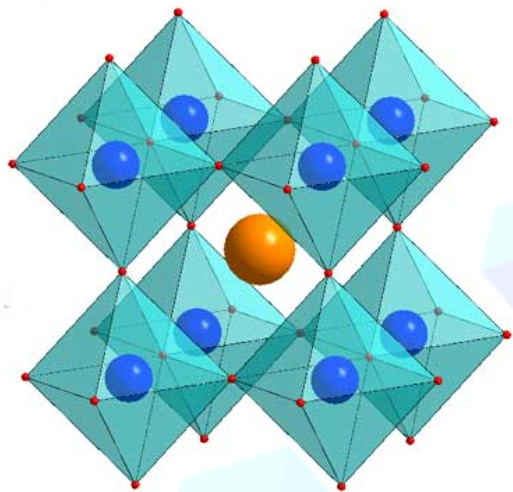
Iterative generation of feature subspaces



topological  
insulator

$d''_1$

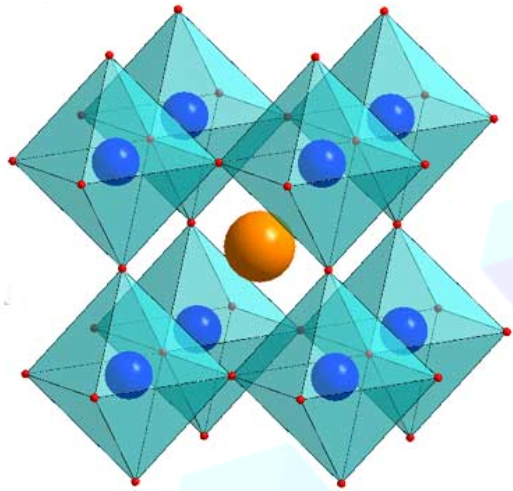
# Perovskites' stability: an improved Goldschmidt Tolerance Factor



$$t = \frac{r_A + r_X}{\sqrt{2}(r_B + r_X)} \longrightarrow \text{Ionic radius}$$

Goldschmidt\* stable perovskites:  $0.825 < t < 1.059$ , accuracy 79%

# Perovskites' stability: an improved Goldschmidt Tolerance Factor



$ABX_3$

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

Ionic radius

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

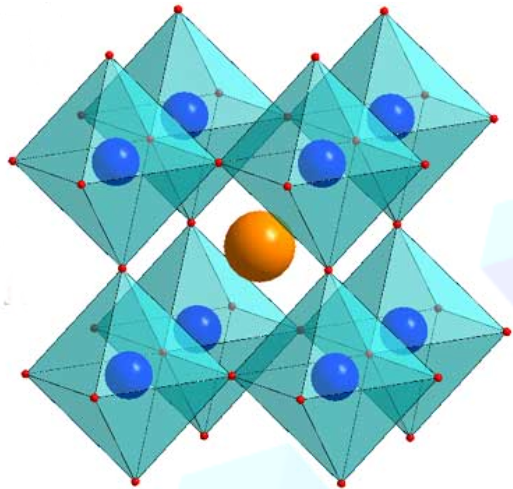
Oxidation state

$1/\mu =$  Octahedral factor

Goldschmidt\* stable perovskites:  $0.825 < t < 1.059$ , accuracy 79%

Our stable perovskites:  $\tau < 4.18$ , accuracy 92%

# Perovskites' stability: an improved 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

$1/\mu =$  Octahedral factor

Goldschmidt\* stable perovskites:  $0.825 < t < 1.059$ , accuracy 79%

Our stable perovskites:  $\tau < 4.18$ , accuracy 92%

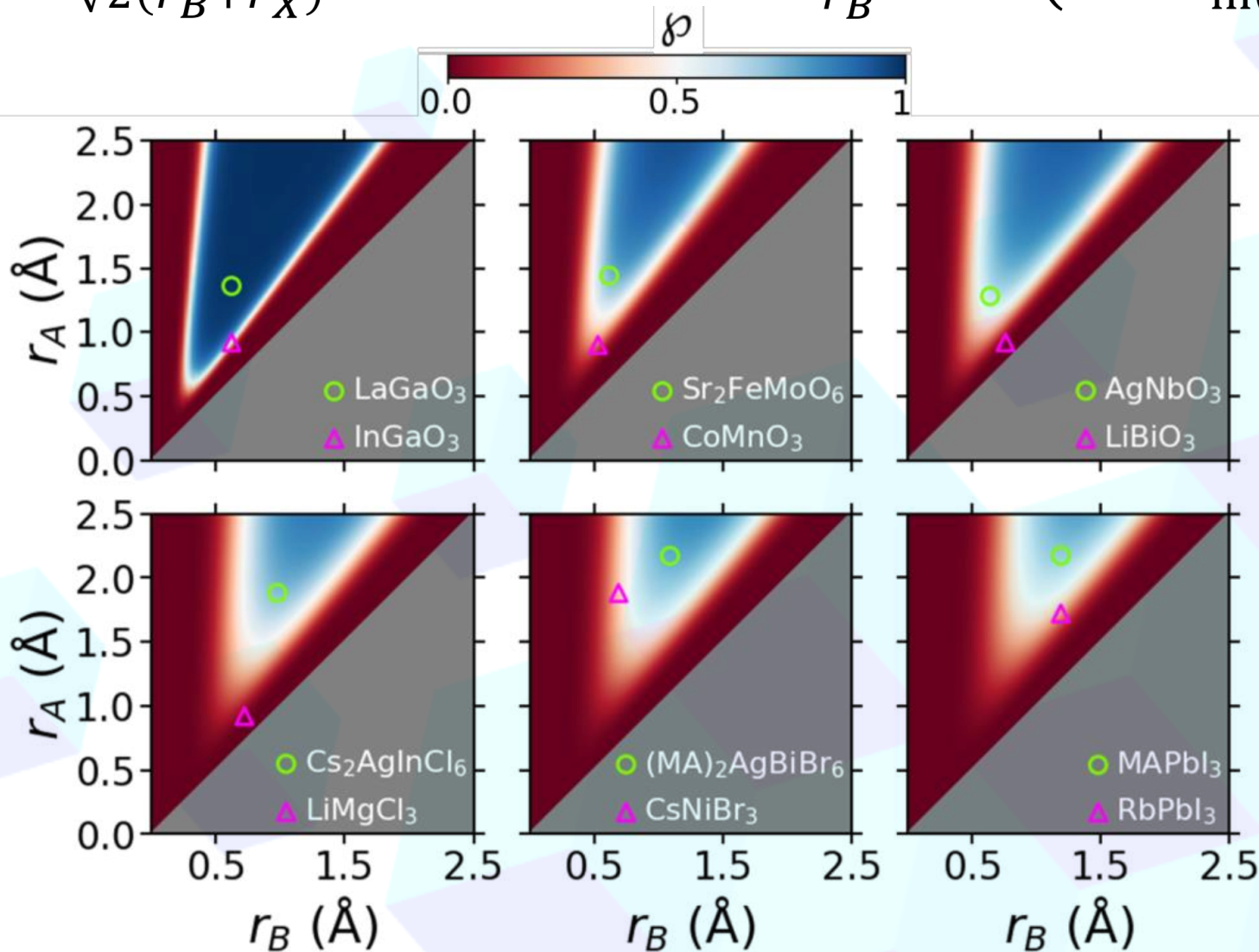
$\tau < 3.31$  or  $\tau > 5.92$ , 99% accuracy (1/3 of the training data)

$\tau < 3.31$  or  $\tau > 12.08$ , 100% accuracy (1/4 of the training data)



# Improved Goldschmidt Tolerance Factor: Materials design

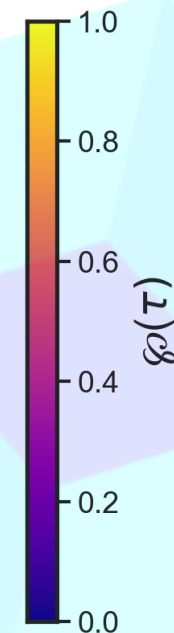
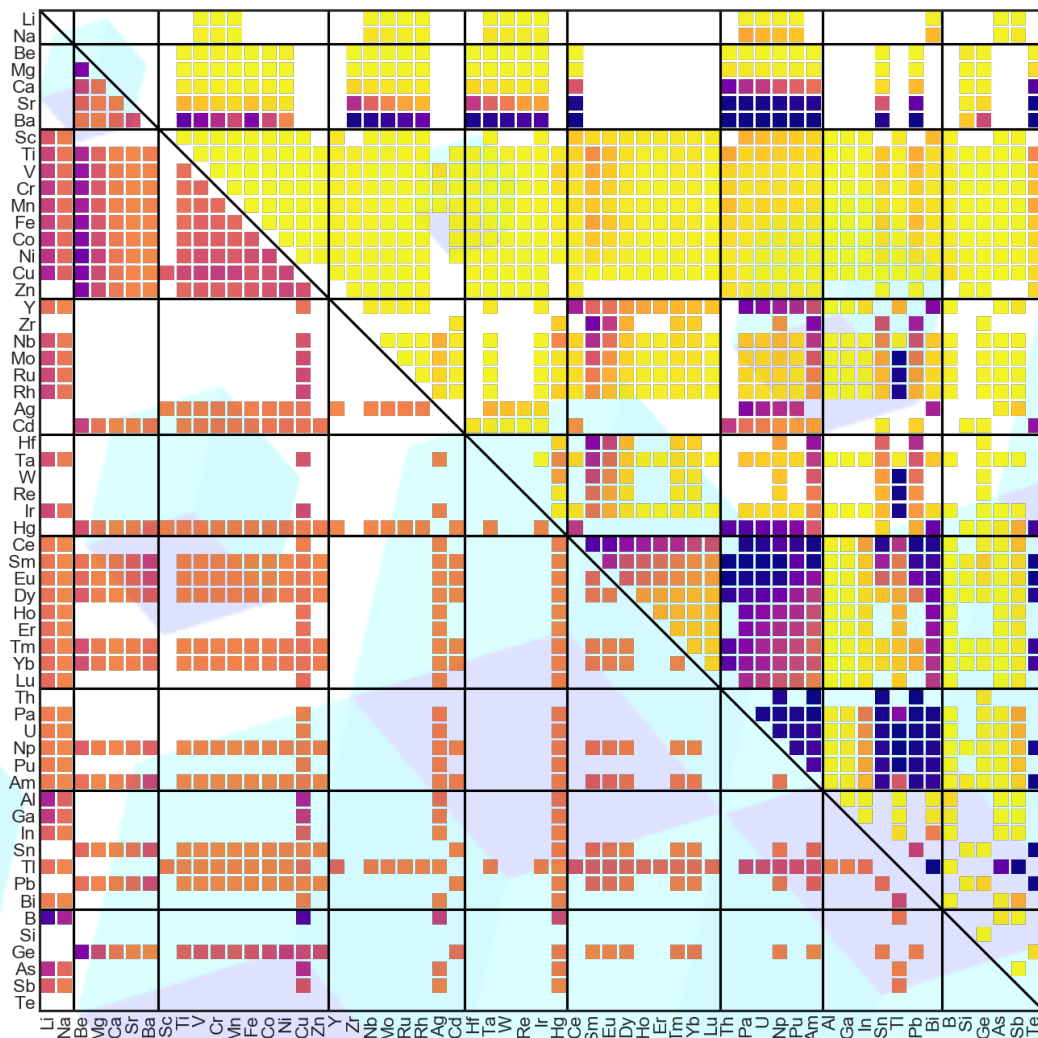
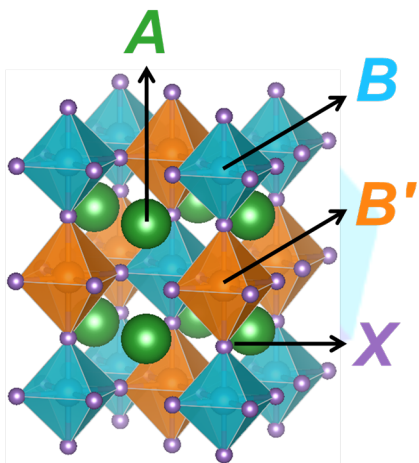
$$t = \frac{r_A + r_X}{\sqrt{2}(r_B + r_X)} \longrightarrow \tau = \frac{r_X}{r_B} - n_A \left( n_A - \frac{r_A/r_B}{\ln(r_A/r_B)} \right)$$





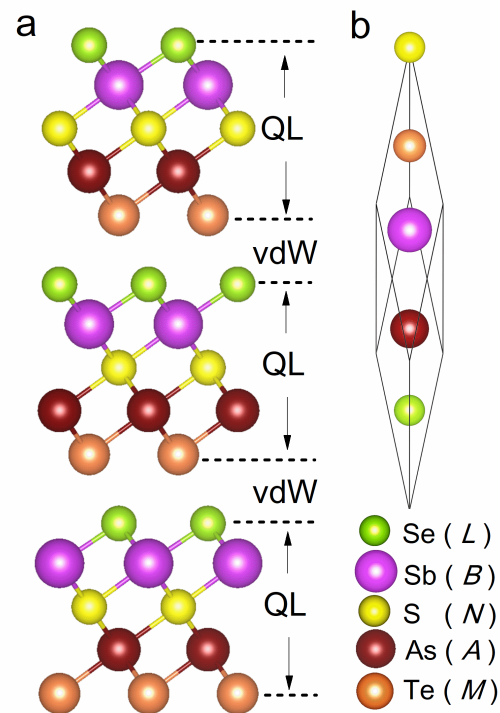
# Improved Goldschmidt Tolerance Factor: Extension of the materials space

$$t = \frac{r_A + r_X}{\sqrt{2}(r_B + r_X)} \longrightarrow \tau = \frac{r_X}{r_B} - n_A \left( n_A - \frac{r_A/r_B}{\ln(r_A/r_B)} \right)$$



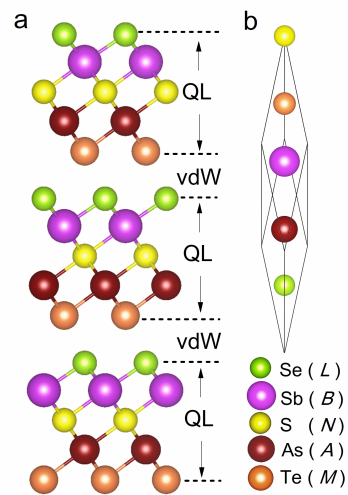
# SISSO: predicting new tetradymite topological insulators

Prototype formula:  
 $AB-LNM$   
 $AB = \{As, Sb, Bi\}$   
 $LNM = \{S, Se, Te\}$



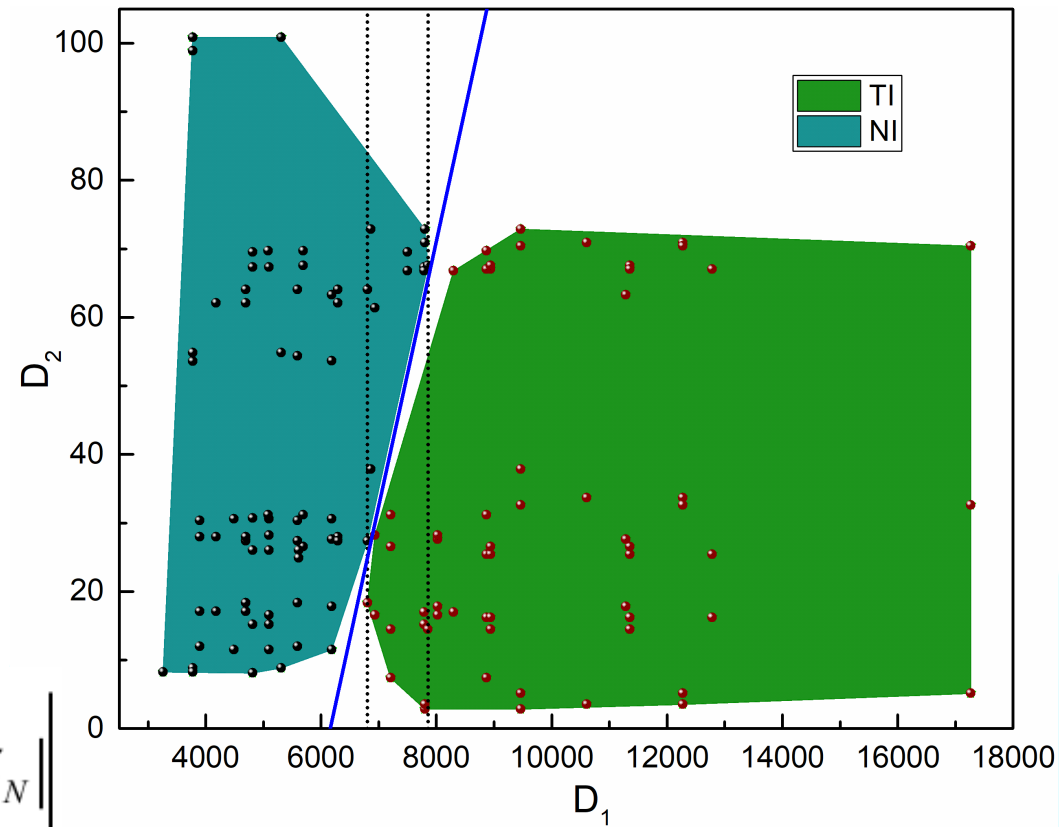
# SISSO: predicting new tetradymite topological insulators

Prototype formula:  
 $AB-LNM$   
 $AB = \{As, Sb, Bi\}$   
 $LNM = \{S, Se, Te\}$



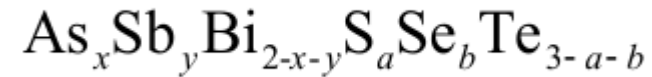
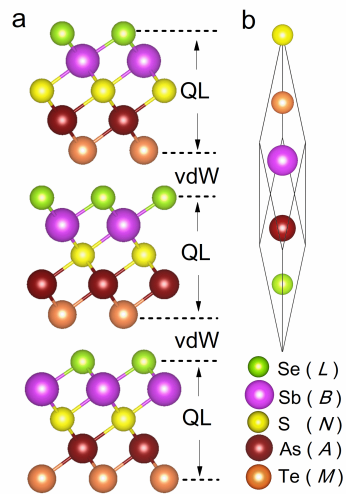
$$D_1 = (Z_A + Z_B) \cdot (Z_L + Z_M) - |Z_A Z_M - Z_B Z_L|$$

$$D_2 = \left| \frac{(\chi_M + \chi_N) \cdot Z_E}{\chi_A} - (Z_M + Z_N) - |Z_M - Z_N| \right|$$



# SISSO: predicting new tetradymite topological insulators

Prototype formula:  
 $AB-LNM$   
 $AB = \{As, Sb, Bi\}$   
 $LNM = \{S, Se, Te\}$

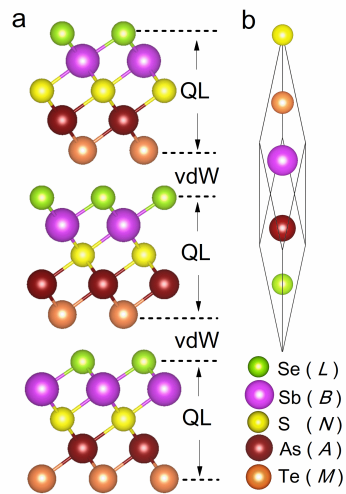


$$D_1 = (Z_A + Z_B) \cdot (Z_L + Z_M) - |Z_A Z_M - Z_B Z_L|$$

$$D_2 = \left| \frac{(\chi_M + \chi_N) \cdot Z_E}{\chi_A} - (Z_M + Z_N) - |Z_M - Z_N| \right|$$

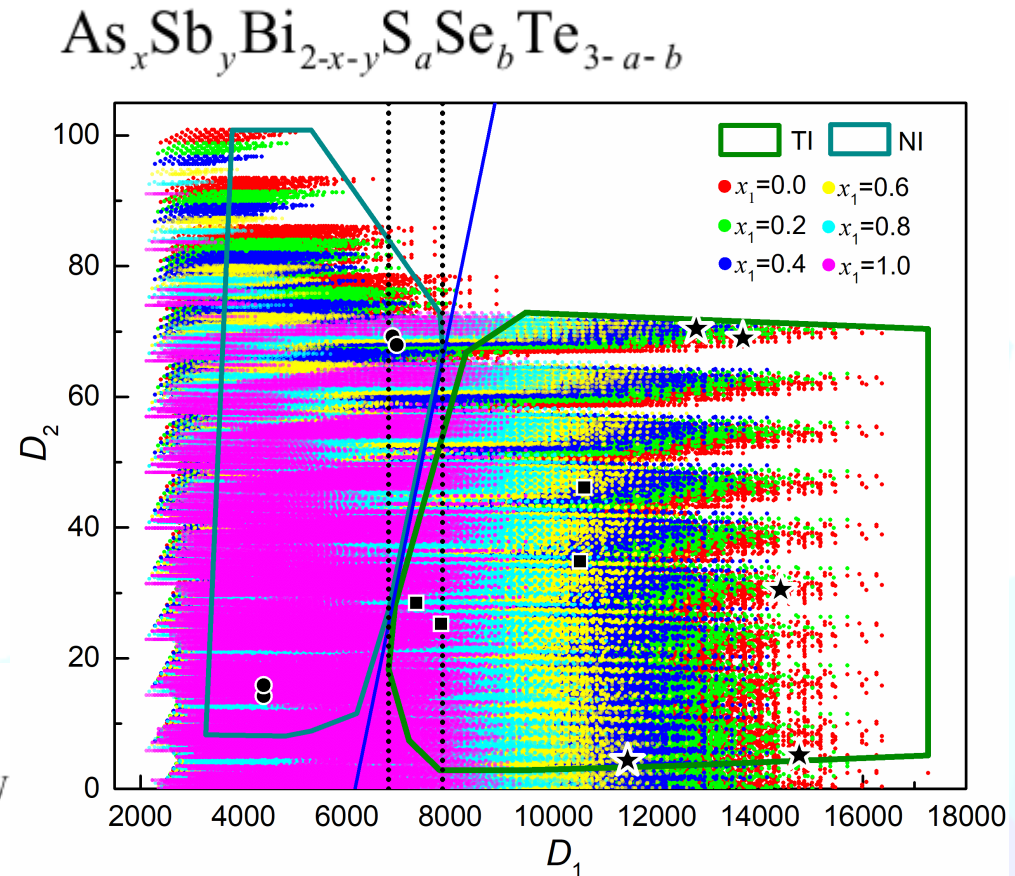
# SISSO: predicting new tetradymite topological insulators

Prototype formula:  
 $AB-LNM$   
 $AB = \{As, Sb, Bi\}$   
 $LNM = \{S, Se, Te\}$



$$D_1 = (Z_A + Z_B) \cdot (Z_L + Z_M) - |Z_A Z_M - Z_B Z_L|$$

$$D_2 = \left| \frac{(\chi_M + \chi_N) \cdot Z_E}{\chi_A} - (Z_M + Z_N) - |Z_M - Z_N| \right|$$



# Acknowledgements

## **Compressed sensing, SISO, multi-task SISO, and metal/insulator proof of concepts**

Jan Vybiral, Runhai Ouyang, Emre Ahmetcik, Stefano Curtarolo, Christian Carbogno, Sergey Levchenko, Claudia Draxl

## **Application of SISO to perovskites**

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

## **Transparent conducting oxide: NOMAD-kaggle competition**

Christopher Sutton, Angelo Ziletti, Takenori Yamamoto, Claudia Draxl, Daan Frenkel, Kristian Thygesen, Samuel Kaski, Bernhard Schölkopf

## **All the above**

Matthias Scheffler



*NOMAD has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No. 676580.*