# Learning Descriptors for Materials Properties with Symbolic Regression and Compressed Sensing

Luca M. Ghiringhelli
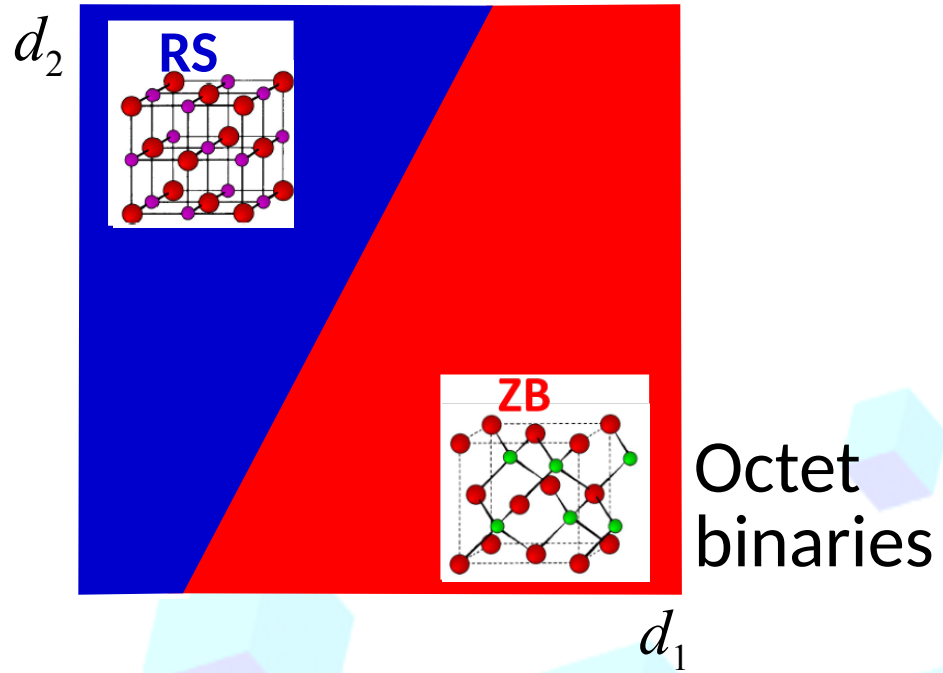FRITZ-HABER-INSTITUT
MAX-PLANCK-GESELLSCHAFT

Big Data Summer
A summer school of the BiGmax Network
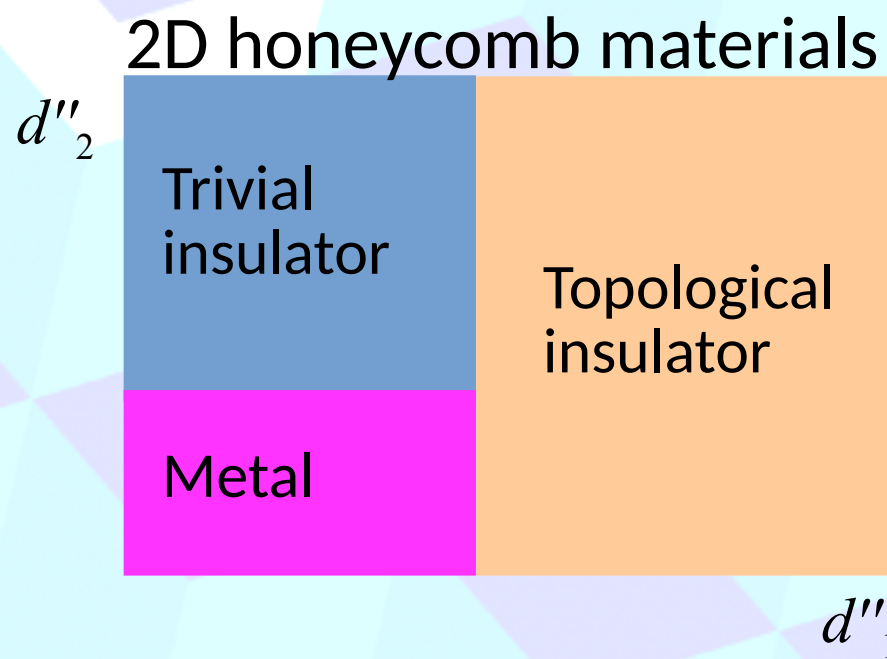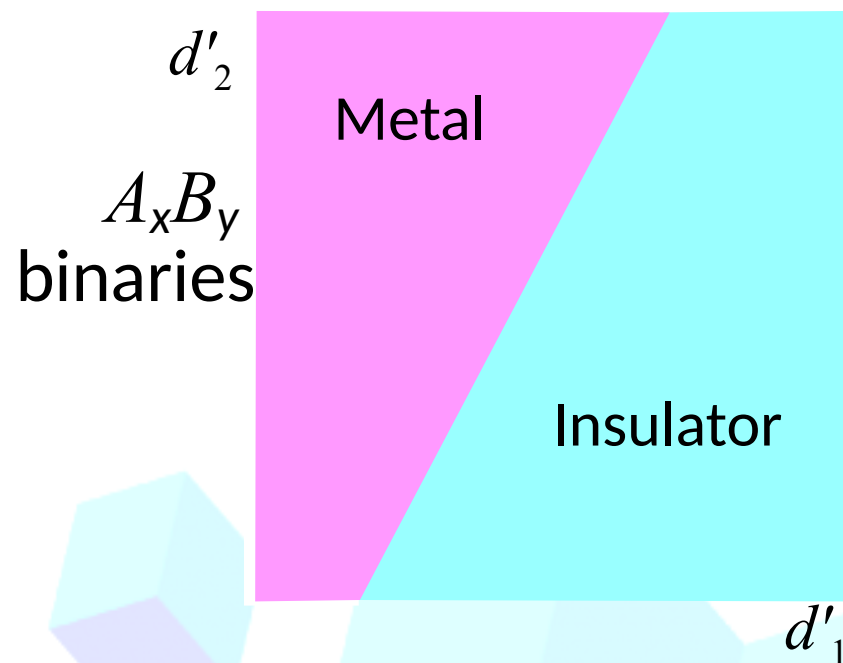Platja d'Aro, Spain, September 9 - 13, 2019

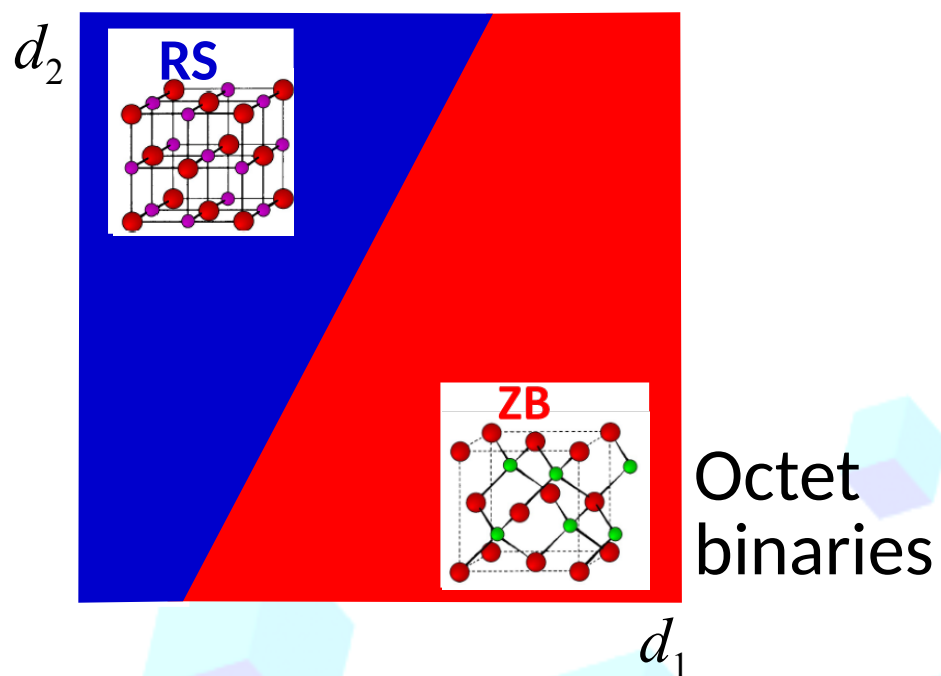# Charts/maps of materials

$d_2$

**RS**

**ZB**

Octet binaries

$d_1$

# Charts/maps of materials

# Charts/maps of materials



$d_2$

**RS**

**ZB**

Octet
binaries

$d_1$

$d'_2$

$A_xB_y$
binaries

Metal

Insulator

$d'_1$

$E_{ads}(CO_2)$ on oxides

$d*_2$

$|E_{ads}|$

$d*_1$

2D honeycomb materials

$d''_2$

Trivial
insulator

Topological
insulator

Metal

$d''_1$

# 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 $ABX_3$):
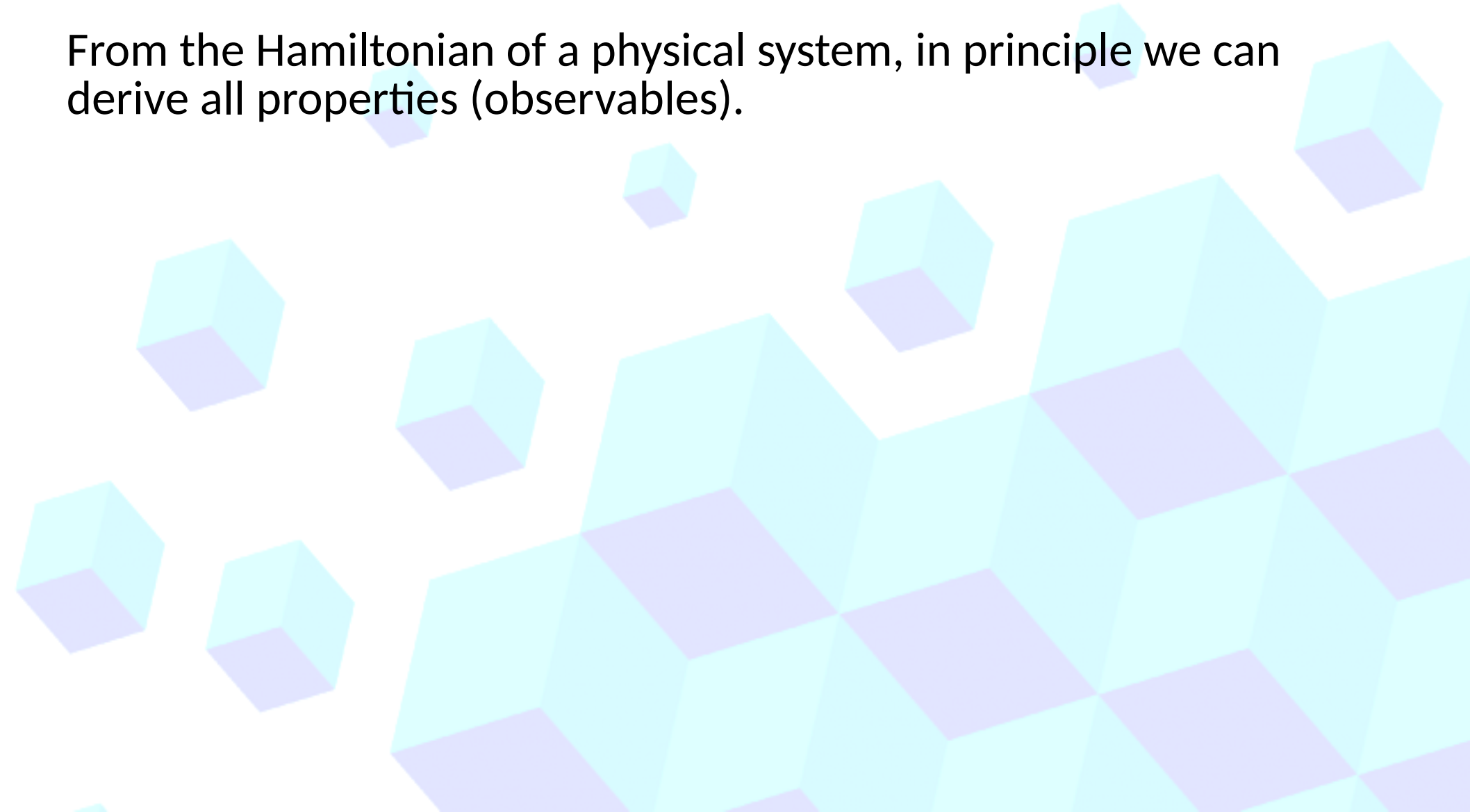
# Building maps of materials properties
# A quantum many-body problem

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 240 000 *inorganic* materials are known to exist (Springer Materials)

# 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)
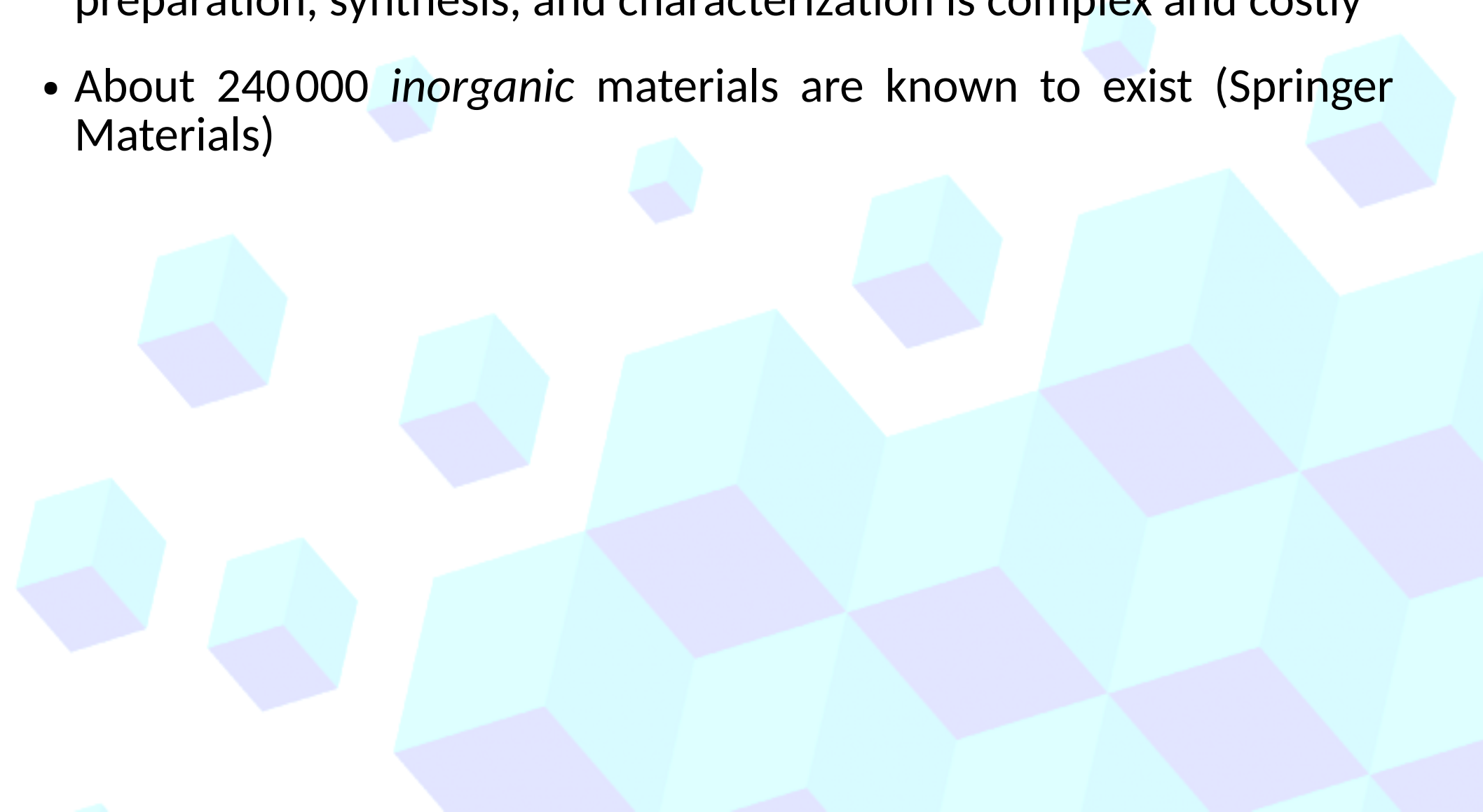
- Basic properties determined for very few of them

# 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

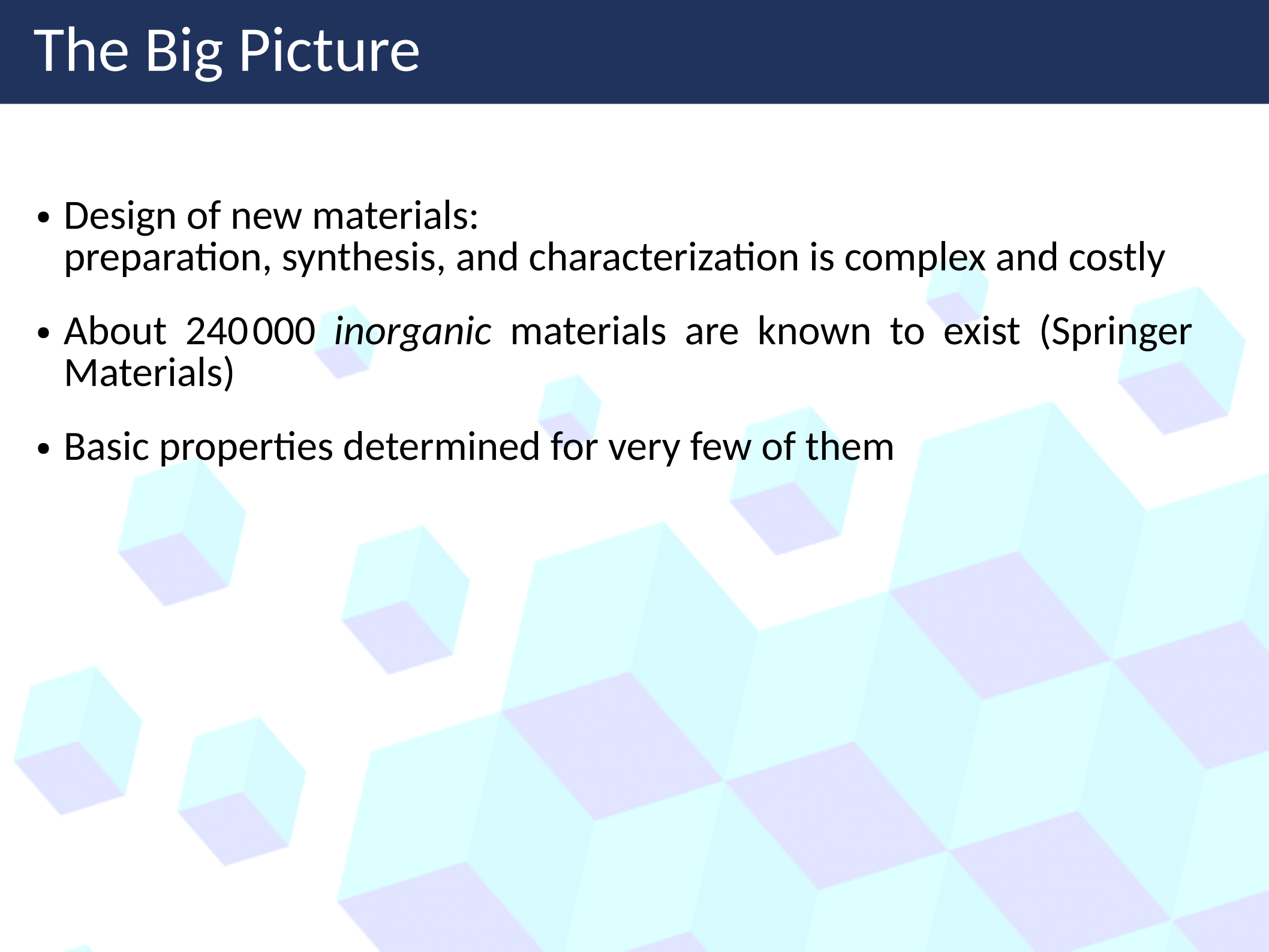- Number of possible materials: practically infinite
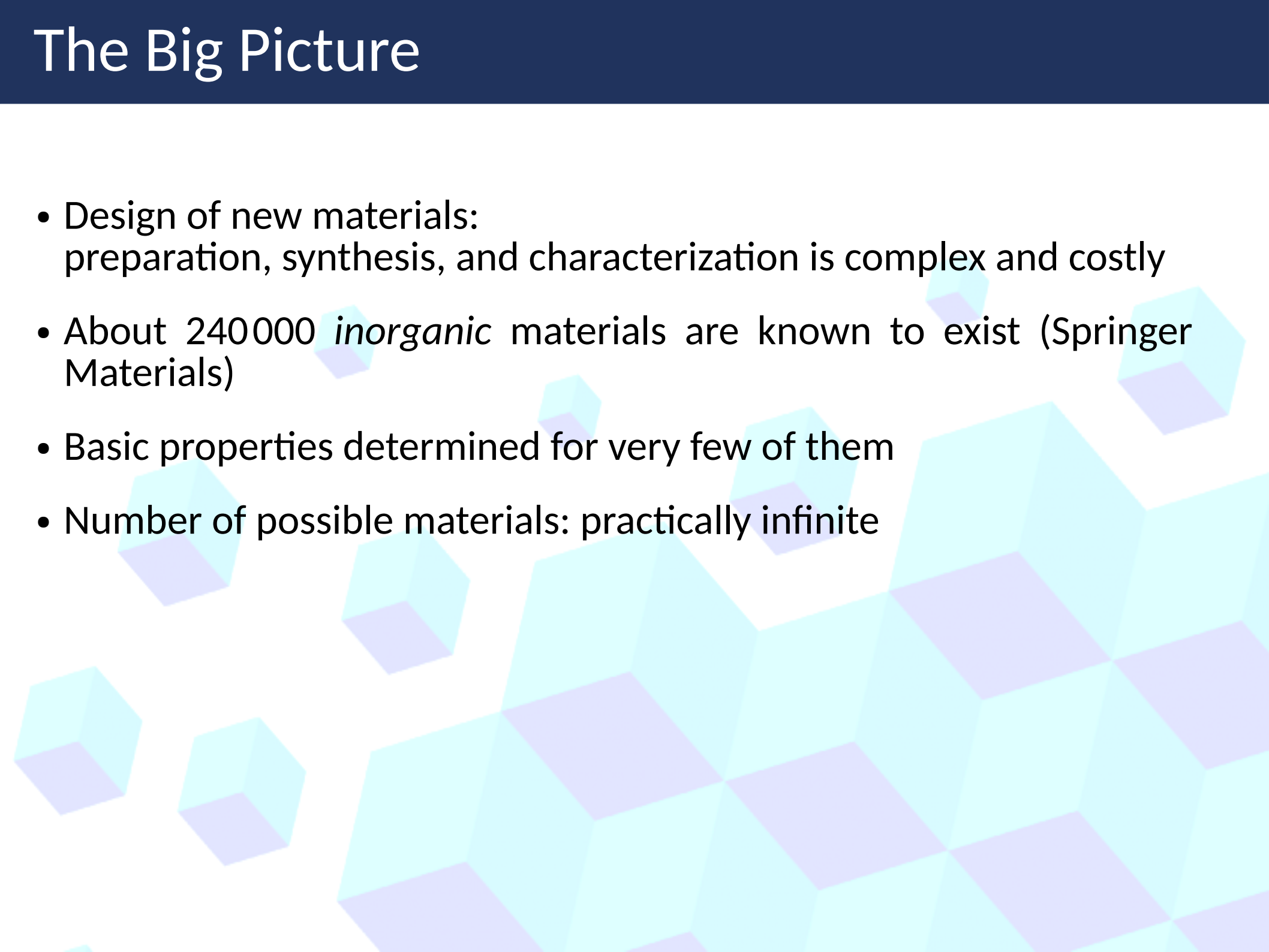
# The Big Picture

- 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

# 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

# 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. — $R^2O$ | Gruppe II. — $RO$ | Gruppe III. — $R^2O^3$ | Gruppe IV. $RH^4$ $RO^2$ | Gruppe V. $RH^3$ $R^2O^5$ | Gruppe VI. $RH^2$ $RO^3$ | Gruppe VII. $RH$ $R^2O^7$ | Gruppe VIII. — $RO^4$ |
|---|---|---|---|---|---|---|---|---|
| 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. — R²O | Gruppe II. — RO | Gruppe III. — R²O³ | Gruppe IV. RH⁴ RO² | Gruppe V. RH³ R²O⁵ | Gruppe VI. RH² RO³ | Gruppe VII. RH R²O⁷ | Gruppe VIII. — RO⁴ |
|---|---|---|---|---|---|---|---|---|
| 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**
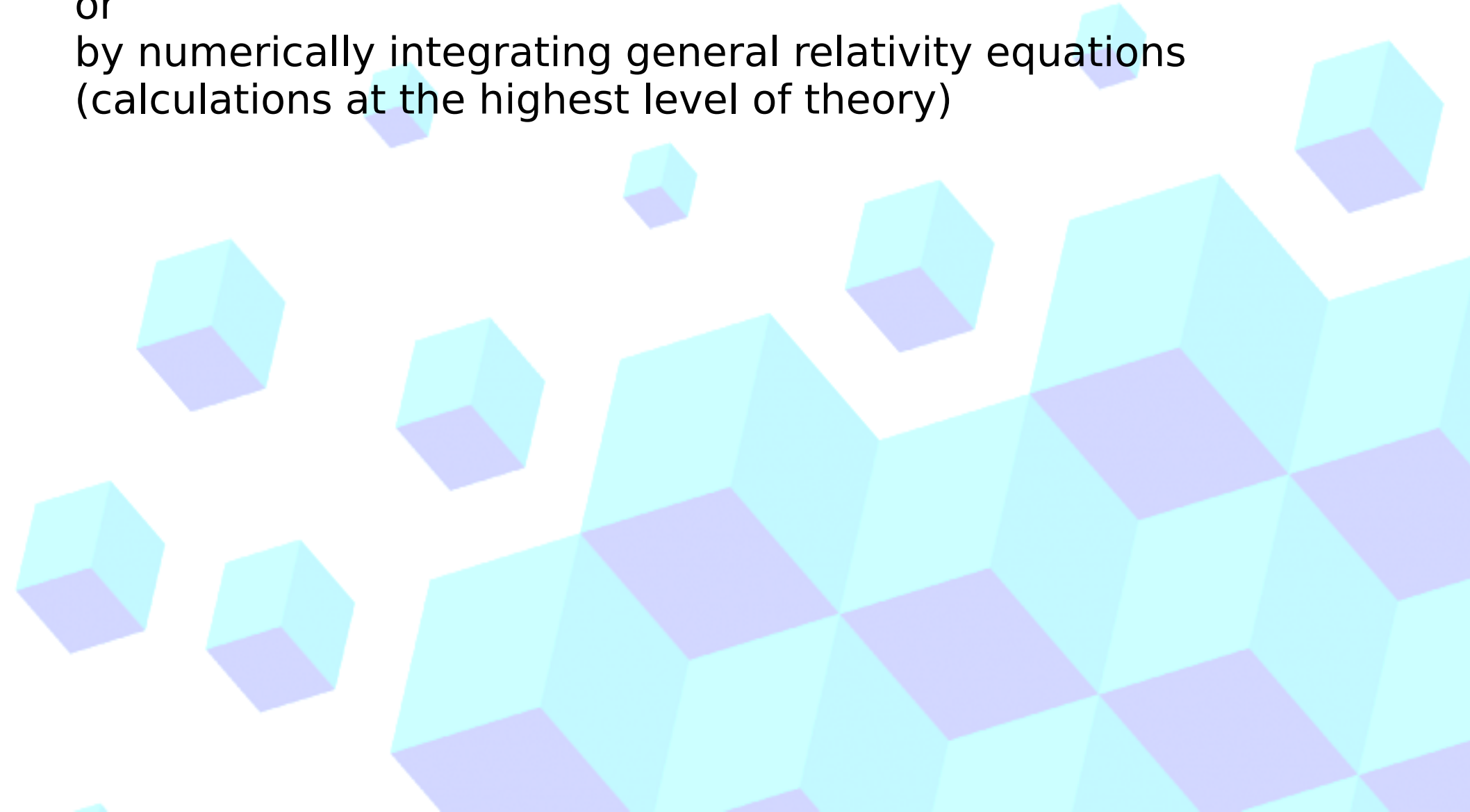


Mendeleev's **1871** periodic table

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)

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)



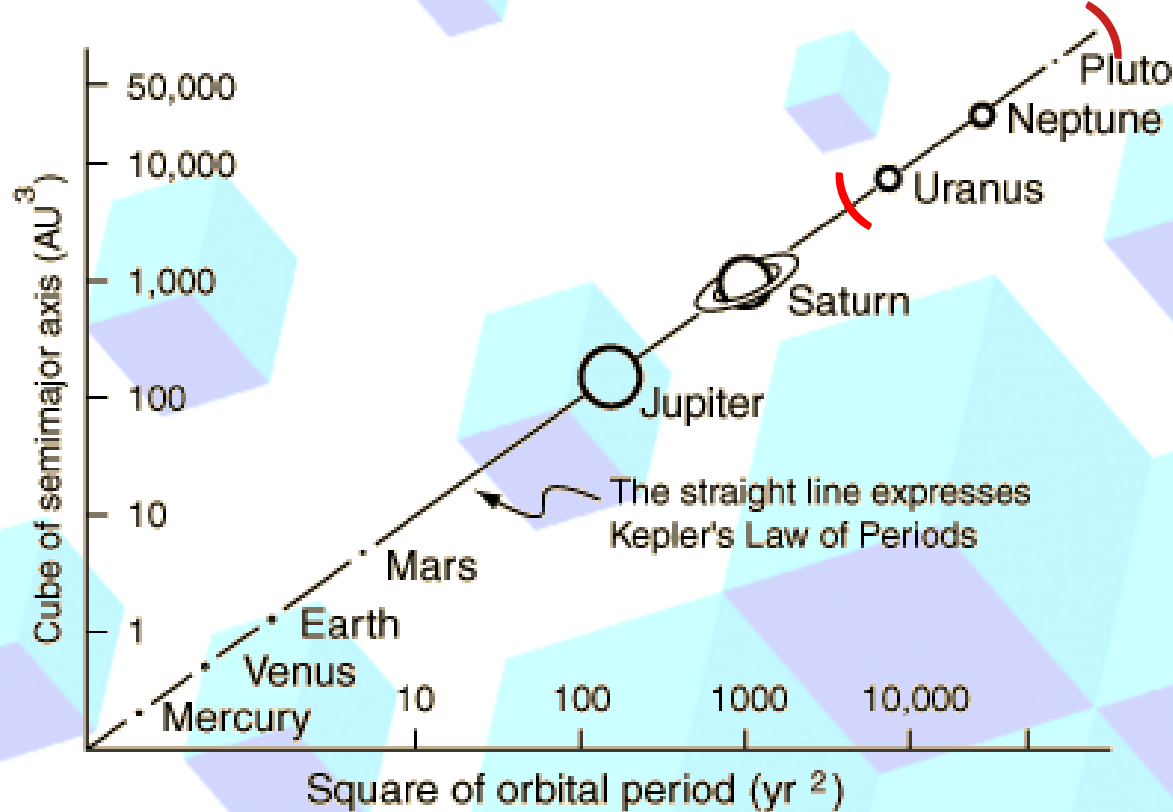(Orbital period)² = C (orbit's major
axis)³

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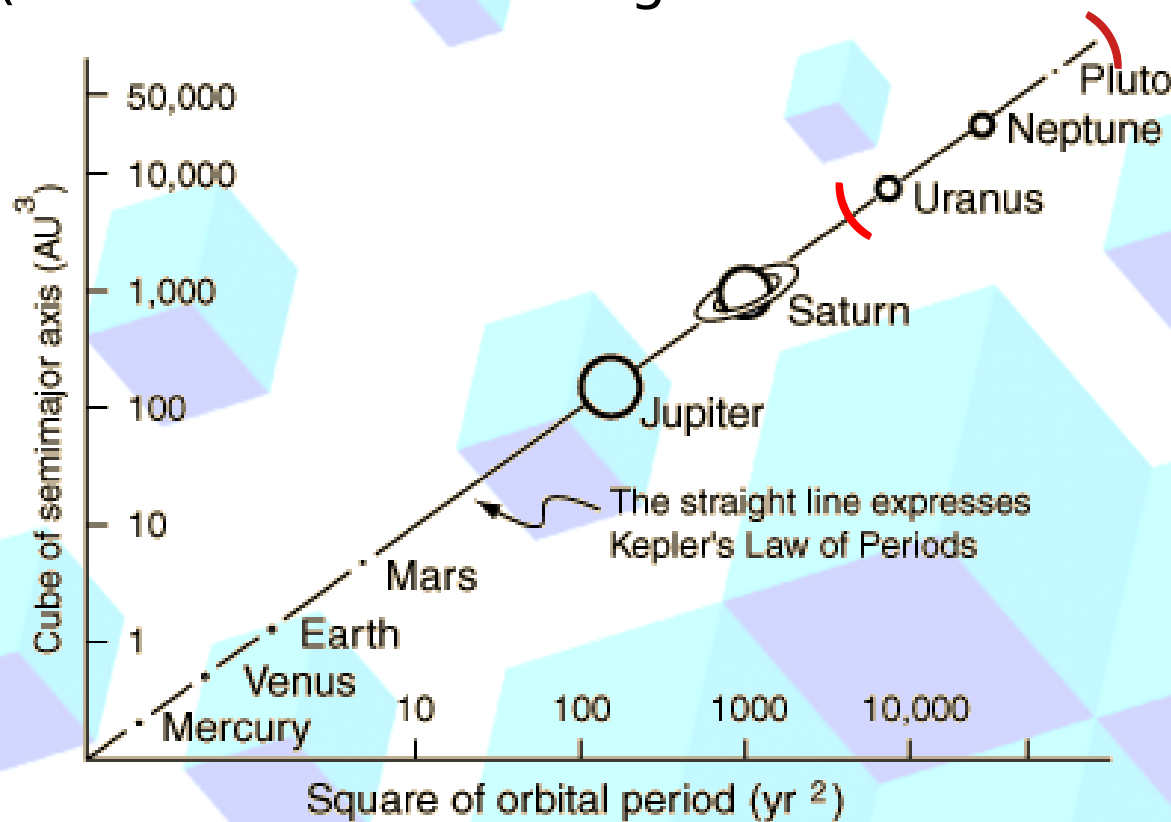(calculations at the highest level of theory)

Data
(collected by
Tycho Brahe)

↓

Statistical learning
(performed by
Johannes Kepler)

↓

Physical law
(assessed by
Isaac Newton)

(Orbital period)² = C (orbit's major axis)³



50,000
10,000
1,000
100
10
1

Cube of semimajor axis (AU³)

Pluto
Neptune
Uranus
Saturn
Jupiter

The straight line expresses
Kepler's Law of Periods

Mars
Earth
Venus
Mercury

10    100    1000    10,000

Square of orbital period (yr²)

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

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

**Descriptor**
Find the *appropriate*
descriptor $\boldsymbol{d}_i$,
build a table:
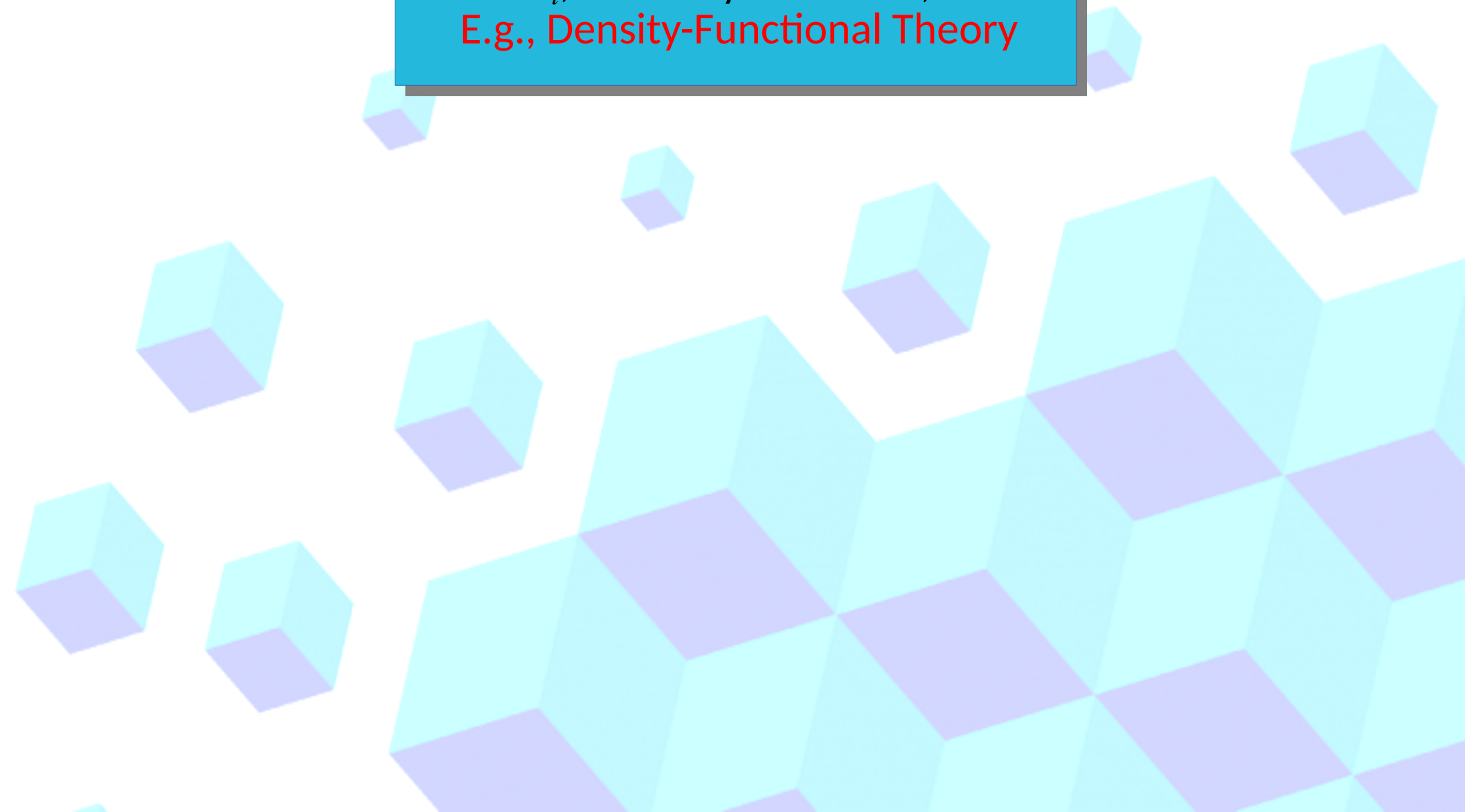$\left| \, i \, \middle| \, \boldsymbol{d}_i \, \middle| \, P_i \, \right|$

# 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 $\boldsymbol{d}_i$,
build a table:
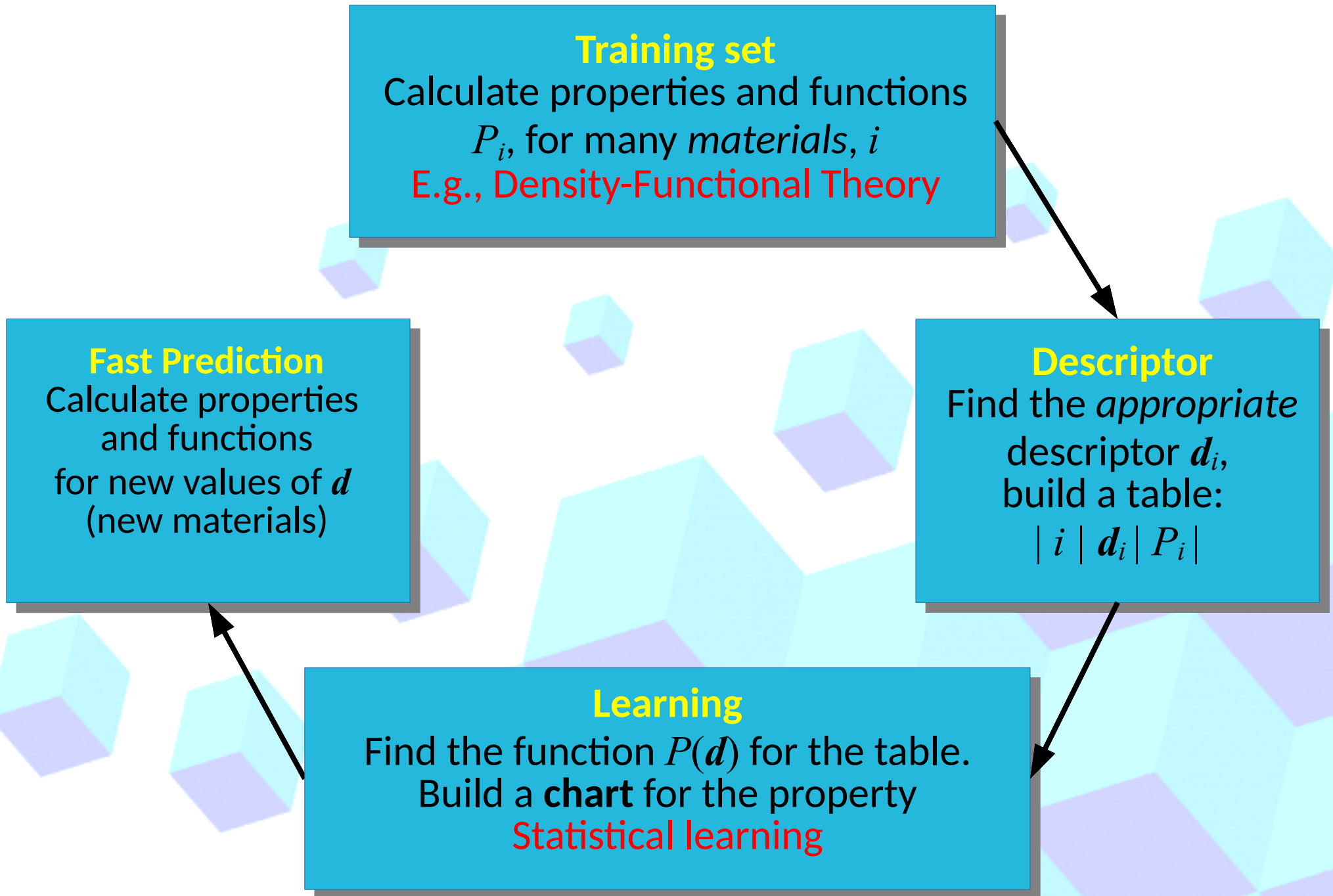$$|\,i\,|\,\boldsymbol{d}_i\,|\,P_i\,|$$

**Learning**
Find the function $P(\boldsymbol{d})$ for the table.
Build a **chart** for the property
Statistical learning

# 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:
$$| \, i \, | \, \boldsymbol{d}_i \, | \, P_i \, |$$

**Learning**
Find the function $P(\boldsymbol{d})$ for the table.
Build a **chart** for the property
Statistical learning

**Fast Prediction**
Calculate properties
and functions
for new values of $\boldsymbol{d}$
(new materials)

# 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.
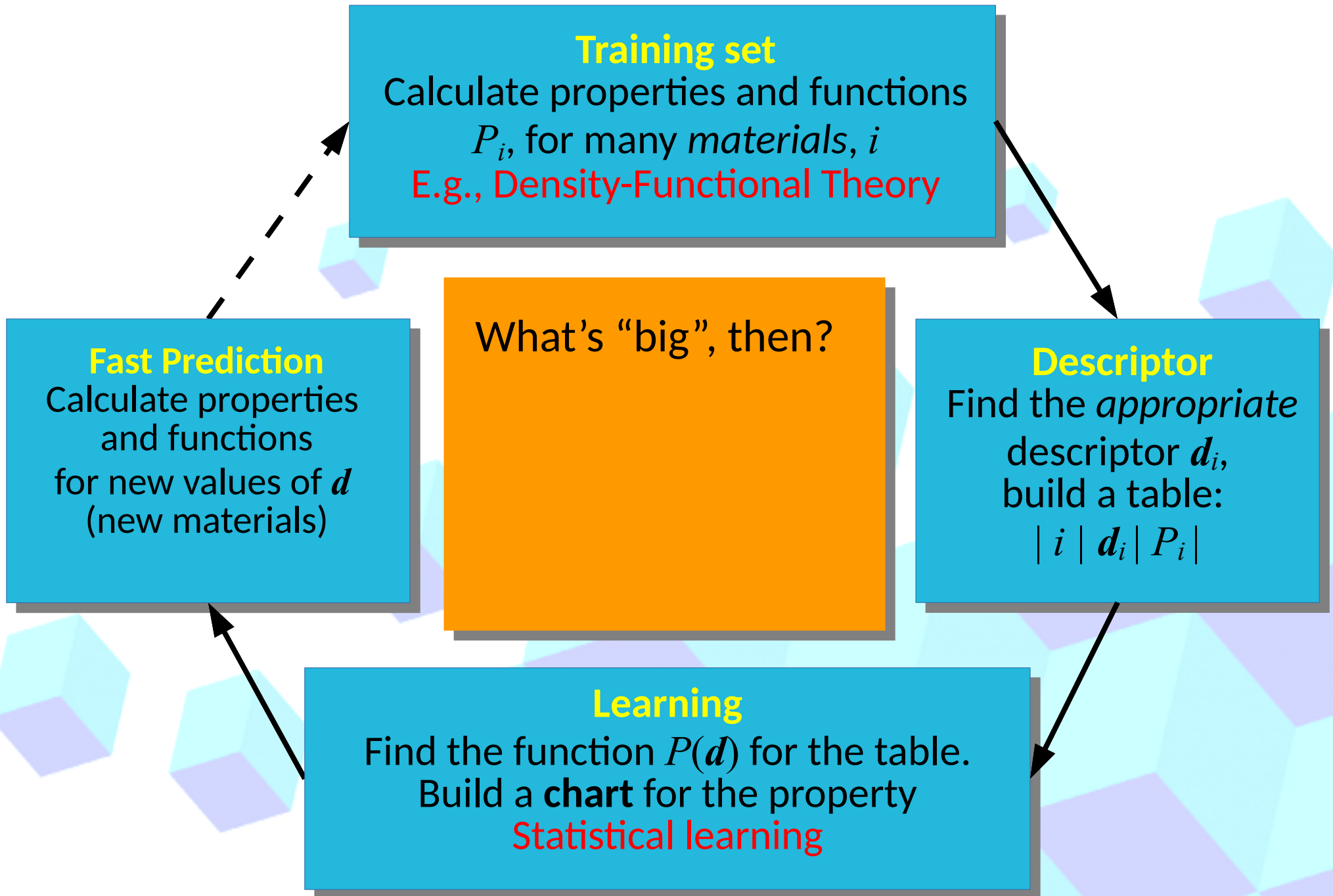Build a **chart** for the property
Statistical learning

**Fast Prediction**
Calculate properties
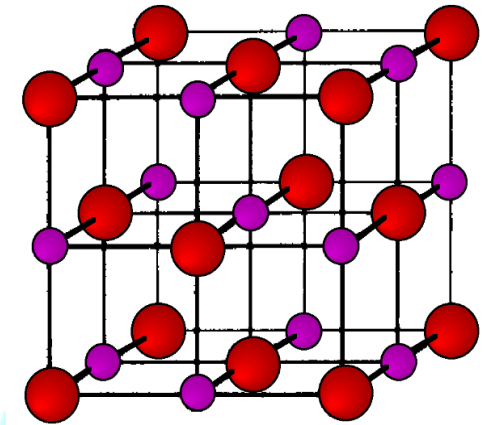and functions
for new values of $d$
(new materials)

# Supervised big-data analysis: a flow chart

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

What's "big", then?

**Descriptor**
Find the *appropriate*
descriptor $d_i$,
build a table:
$| i | d_i | P_i |$

**Fast Prediction**
Calculate properties
and functions
for new values of $d$
(new materials)

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

# Supervised big-data analysis: a 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)

**What's "big", then?**

- Volume
- Velocity
- Variety
- Veracity issue

**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

# 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 $\boldsymbol{d}_i$,
build a table:
$$| \, i \, | \, \boldsymbol{d}_i \, | \, P_i \, |$$

**Learning**
Find the function $P(\boldsymbol{d})$ for the table.
Build a **chart** for the property
Statistical learning

# An example: predicting crystal structures from the composition

## 82 octet AB binary compounds



**Rock salt**



**Zinc blende**

# An example: predicting crystal structures from the composition



82 octet AB binary compounds

Rock salt

Rock salt/Zinc blende

Zinc blende
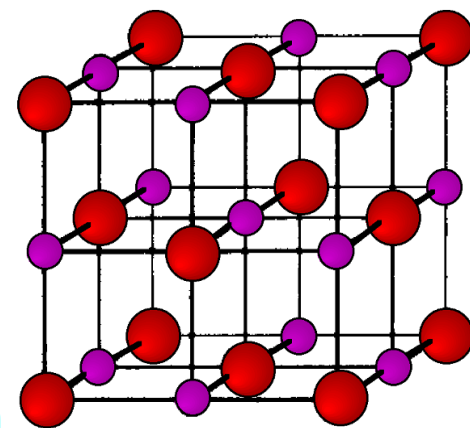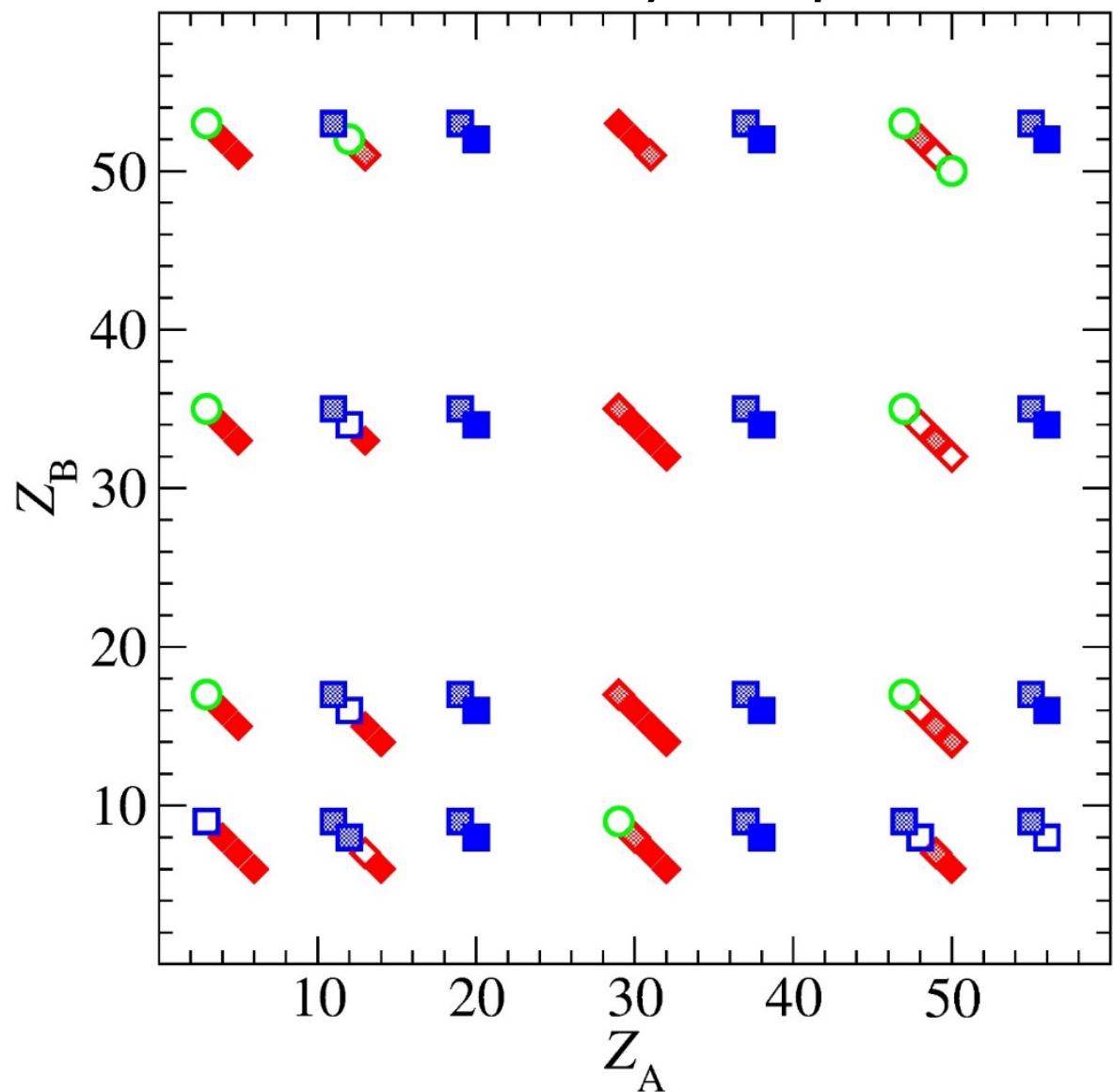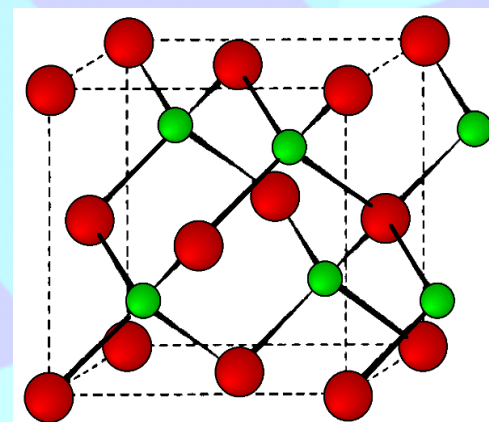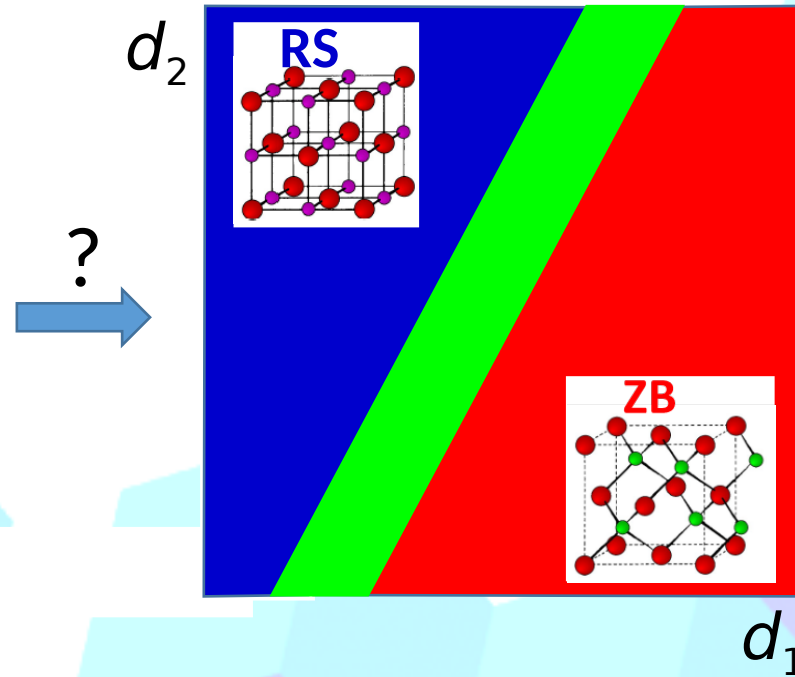
Rock salt

Zinc blende

# An example: predicting crystal structures from the composition

## 82 octet AB binary compounds



$d_2$    RS

$?$

ZB

$d_1$

■ 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).
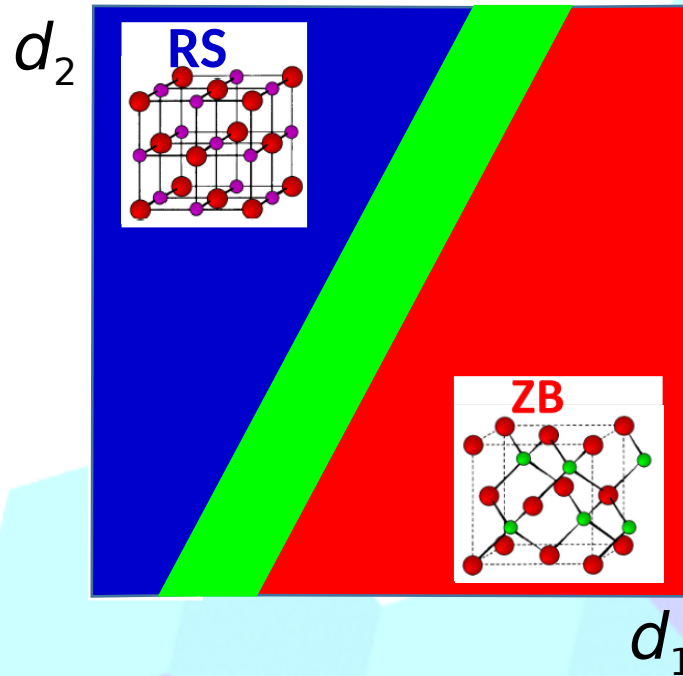**Y. Saad, D. Gao, T. Ngo, S. Bobbitt, J. R. Chelikowsky, and W. Andreoni**, Phys. Rev. B 85, 104104 (2012).

# 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
- ... ?

$\downarrow$

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



$d_2$

**RS**

**ZB**

$d_1$

■ 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).

example: Sn (Tin)

example: Sn (Tin)

# Primary (atomic) features



example: Sn (Tin)

KS levels [eV]

LUMO

Valence *p* (HOMO)

Valence *s*

KS levels [eV]

Radial probability densities

Valence *s*

Valence *p*

Average radius

Radius @ max

Turning point

# 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
- ... ?

$E$(Rock salt) – $E$(Zinc blende)



$d_2$

**RS**

**ZB**

$d_1$

- ■ 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)

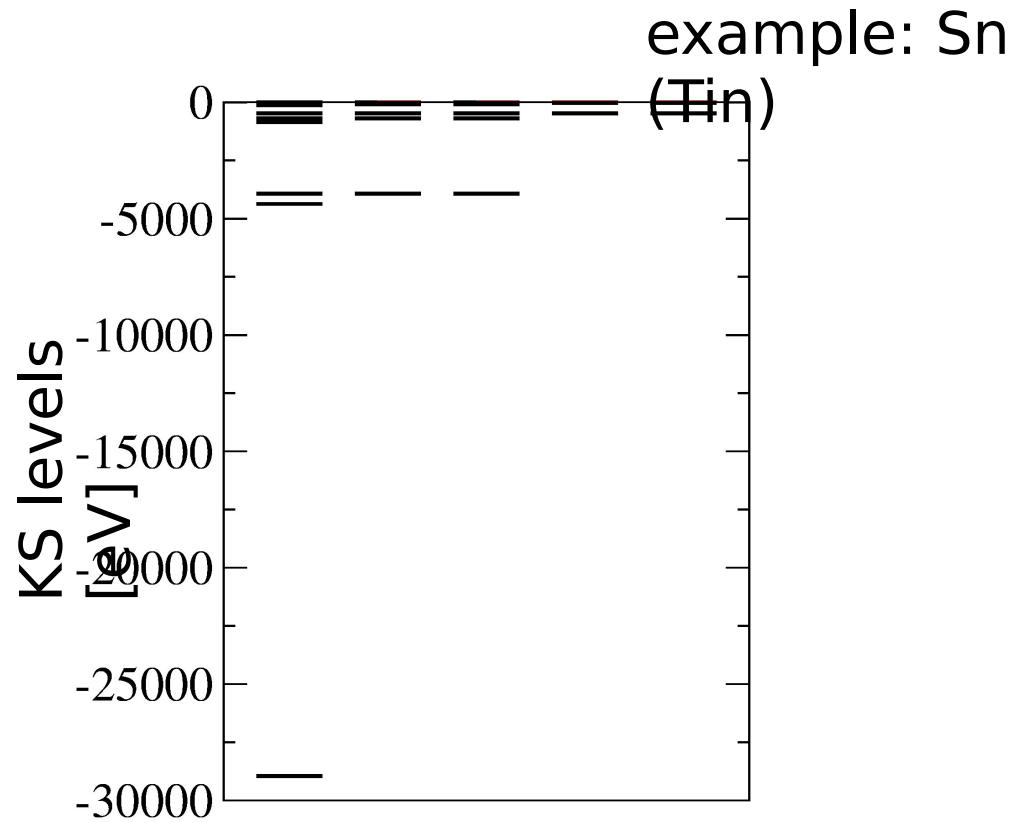J. R. Chelikowsky and J. C. Phillips, Phys. Rev. B 33, 2453 (1978)

A. Zunger, Phys. Rev. B 22, 5839 (1980).
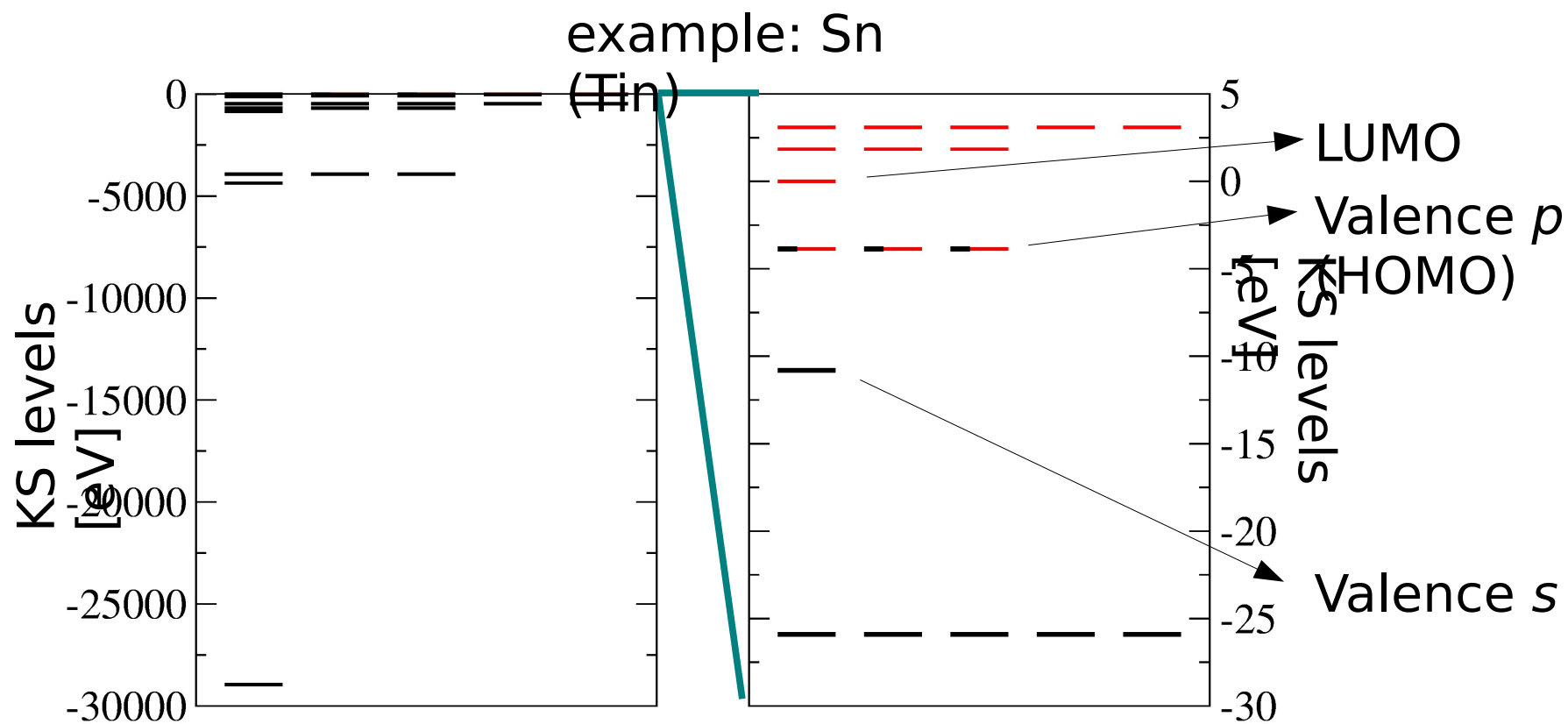
D. G. Pettifor, Solid State Commun. 51, 31 (1984).

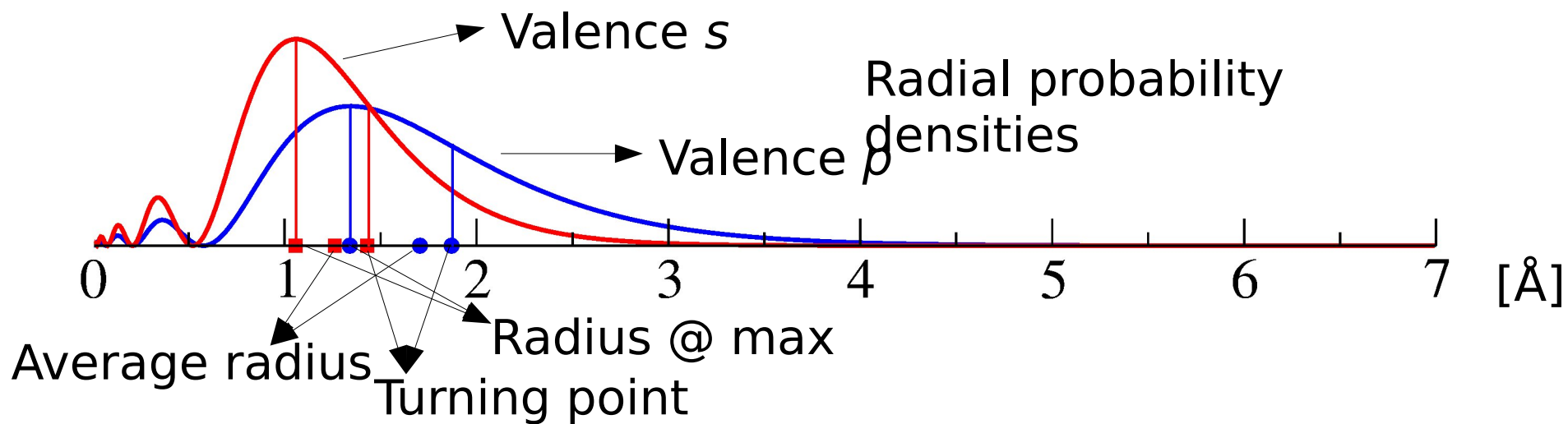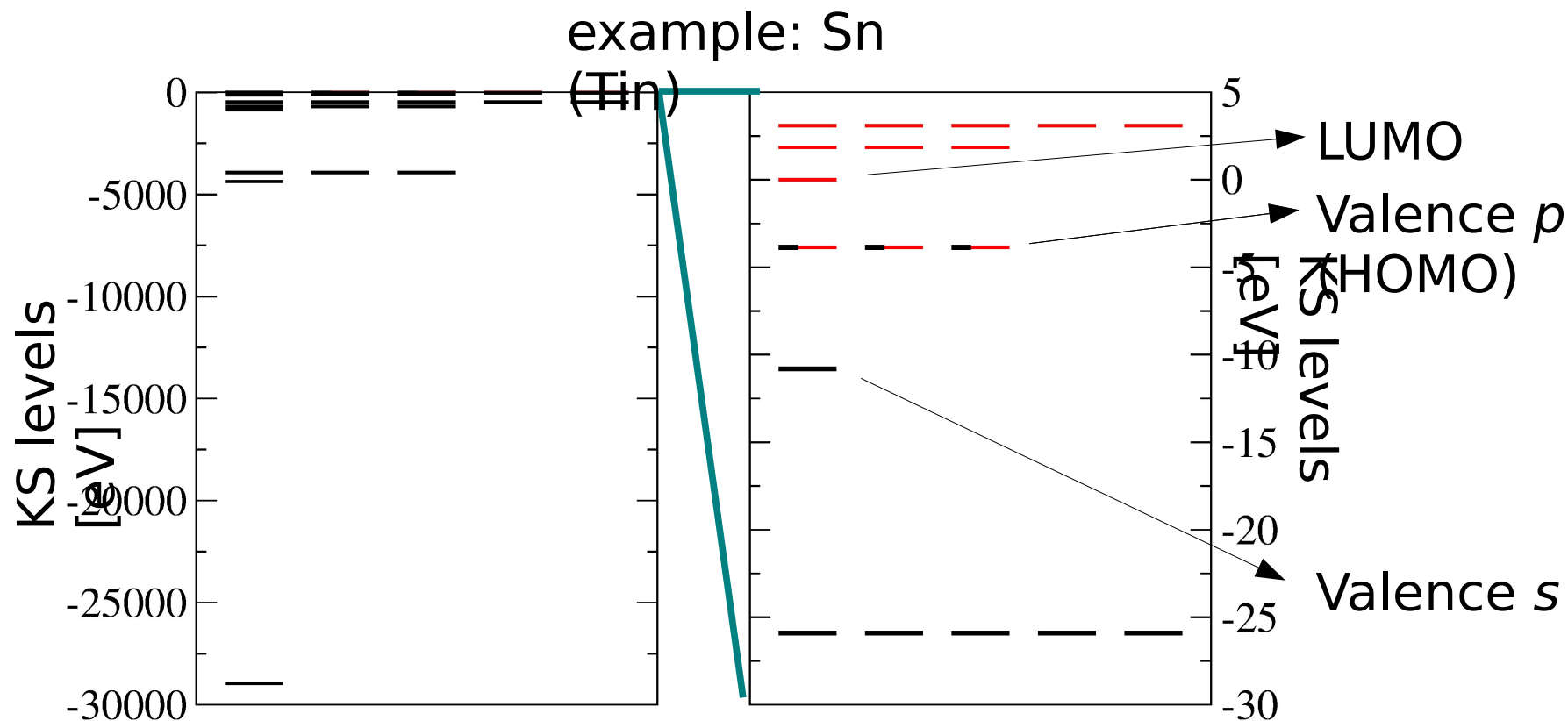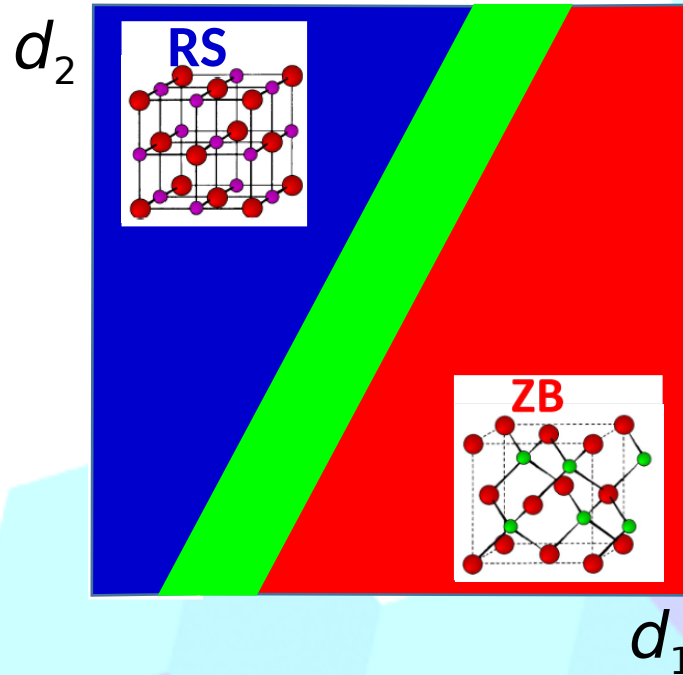Y. Saad, D. Gao, T. Ngo, S. Bobbitt, J. R. Chelikowsky, and W. Andreoni, Phys. Rev. B 85, 104104 (2012).

# Compressed sensing

Aim: finding descriptors and learning predictive models
Ansatz:
$$P = c_1 d_1 + c_2 d_2 + \dots c_n d_n$$
Where
$P$ is the property of interest
$d_1, \dots d_n$ 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.

Aim: finding descriptors and learning pre~~dictive~~ models
Ansatz:
$$P = c_1 d_1 + c_2 d_2 + \dots c_n d_n$$
Where
$P$ is the property of interest
$d_1, \dots d_n$ 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**

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

Ansatz: atomic features

- Valence number $Z_v$ est
- 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
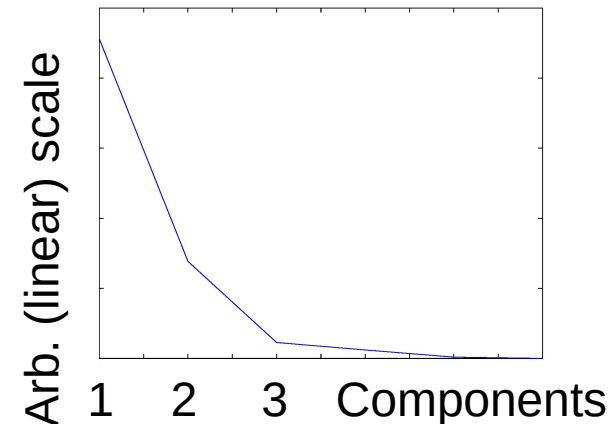constrained to being orthogonal to all the
preceding components

Arb. (linear) scale

1    2    3    Components

Saad, …, Chelikowsky, and Andreoni, PRB  85, 104104 (2012)

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

lineariy uncorreiated (caiied principai 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
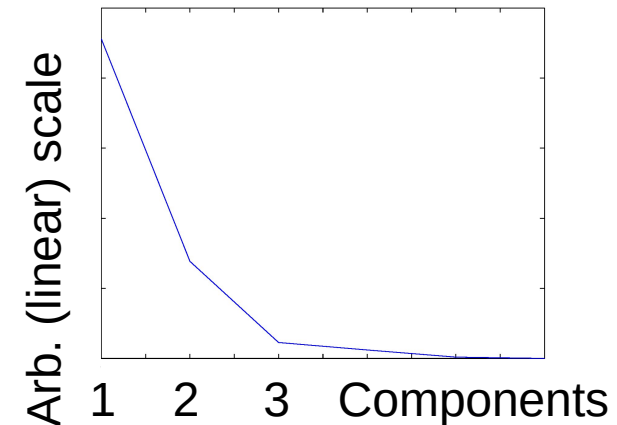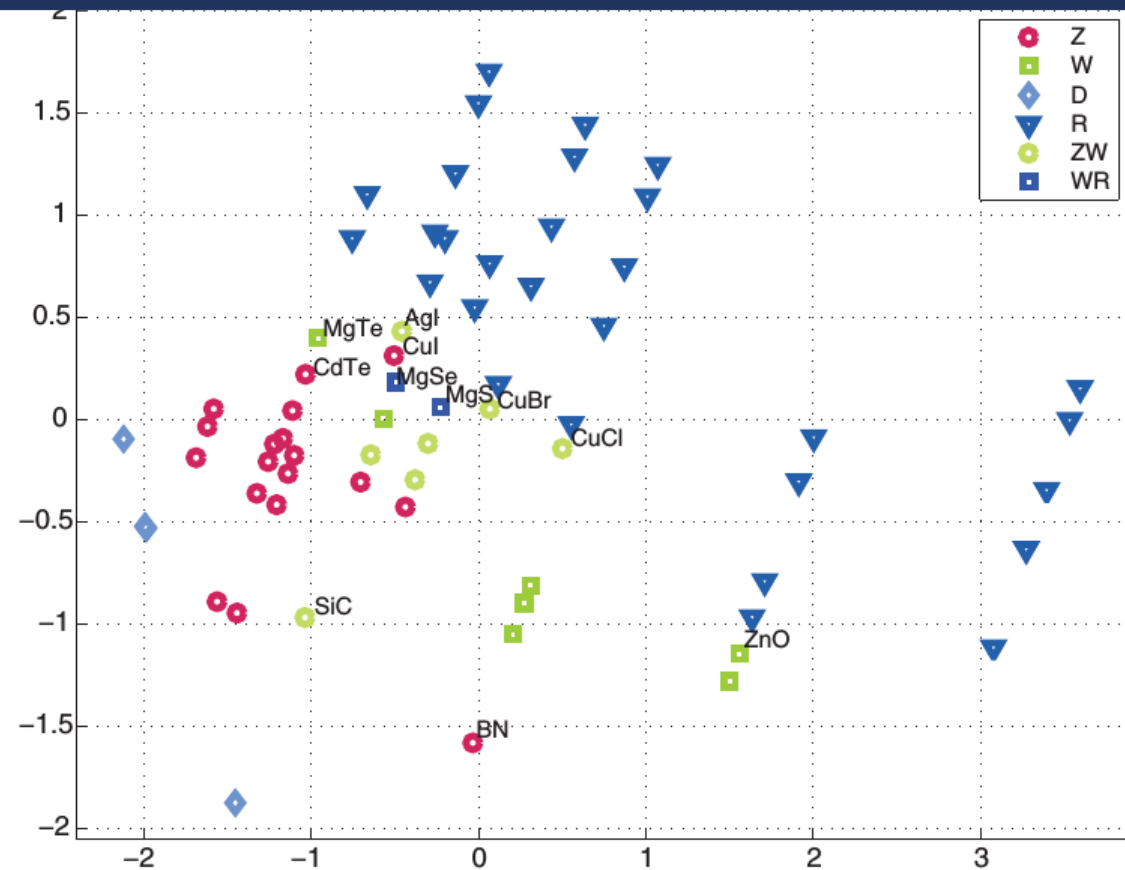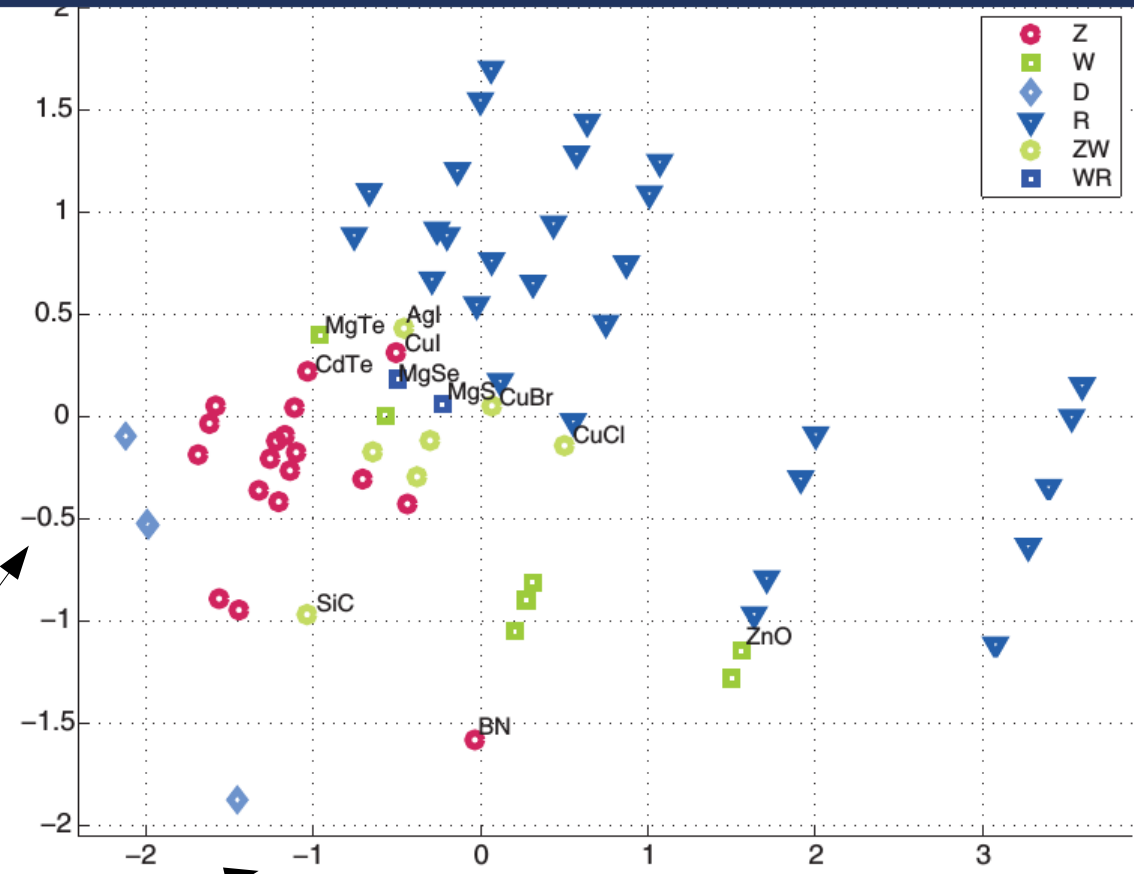


Saad, …, Chelikowsky, and Andreoni, PRB 85, 104104 (2012)

Ansatz: atomic features

- Valence number $\qquad$ $Z_v$
- Energy of valence $s$ orbital $\quad$ $E_s$
- Energy of valence $p$ orbital $\quad$ $E_p$
- Radius of valence $s$ orbital $\quad$ $r_s$
- Radius of valence $p$ orbital $\quad$ $r_p$
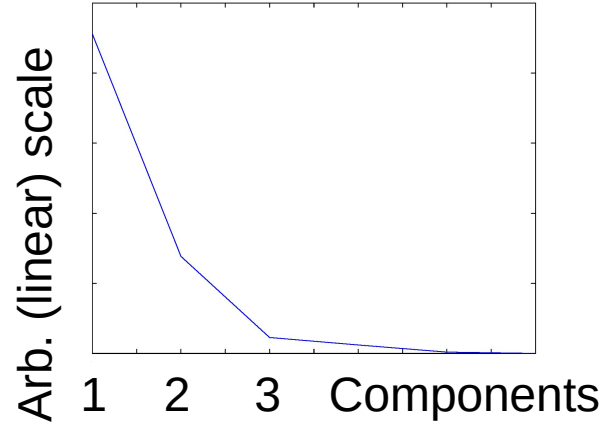
$r_s$, $r_p$, $E_s/\sqrt{Z_v}$, $E_p/\sqrt{z_v}$, for A and B atoms



What's on the axes?

Linear combination of (possibly all) the initial dimensions



Saad, …, Chelikowsky, and Andreoni, PRB 85, 104104 (2012)

# 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

$E$(Rock salt) − $E$(Zinc blende)

$d_2$

RS

ZB

- Rock salt
- RS / ZB
- Zinc blende

$d_1$

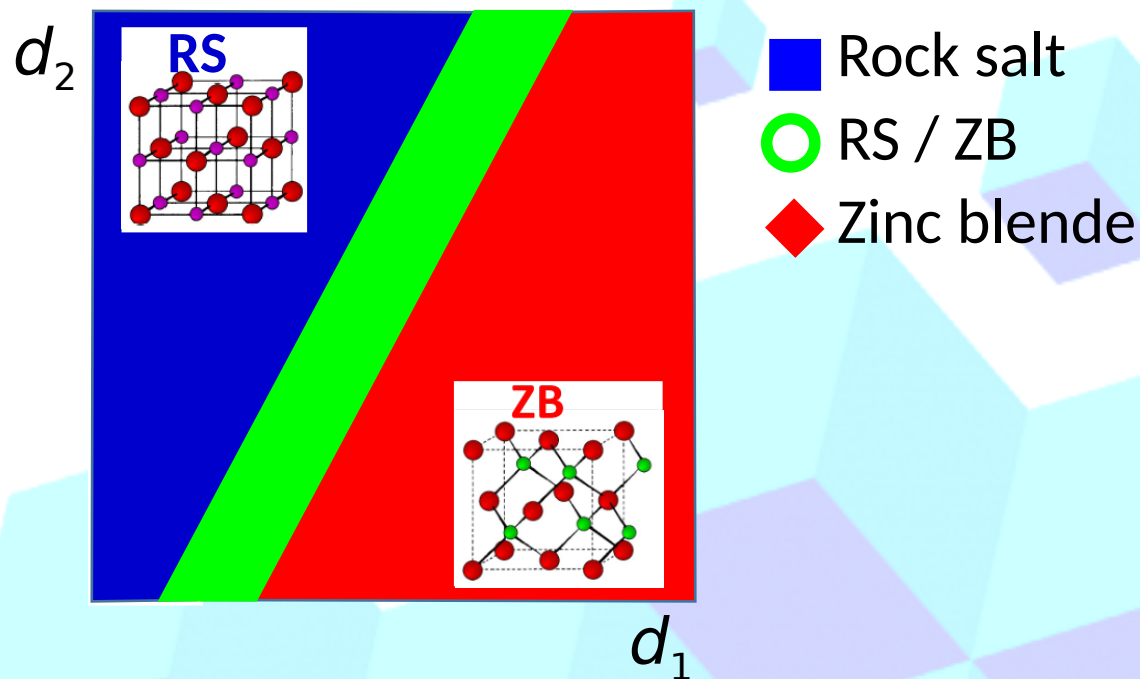$$P = c_1 d_1 + c_2 d_2 + \ldots c_n d_n$$

# 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**

$E$(Rock salt) − $E$(Zinc blende)



$d_2$

**RS**

**ZB**

- Rock salt
- RS / ZB
- Zinc blende

$d_1$

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

# Systematic construction of the feature space

# Systematic construction of the feature space

# Systematic construction of the feature space

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 | |
|---|---|
| | Exponential |
| Constant value | Natural logarithm |
| Input variable | Power |
| Addition | Square root |
| Subtraction | Logistic function |
| Multiplication | Minimum |
| Division | Maximum |
| Negation | 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



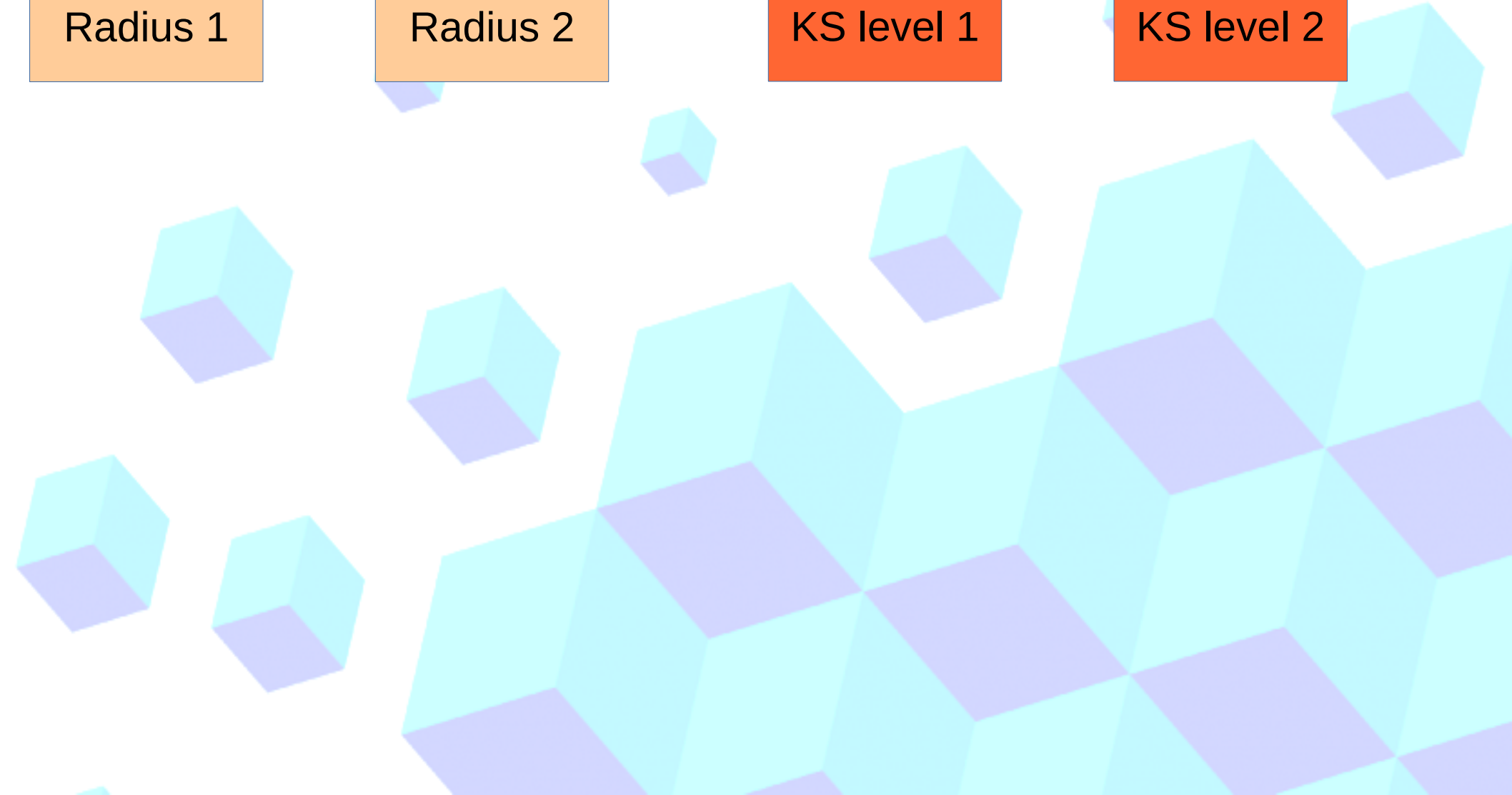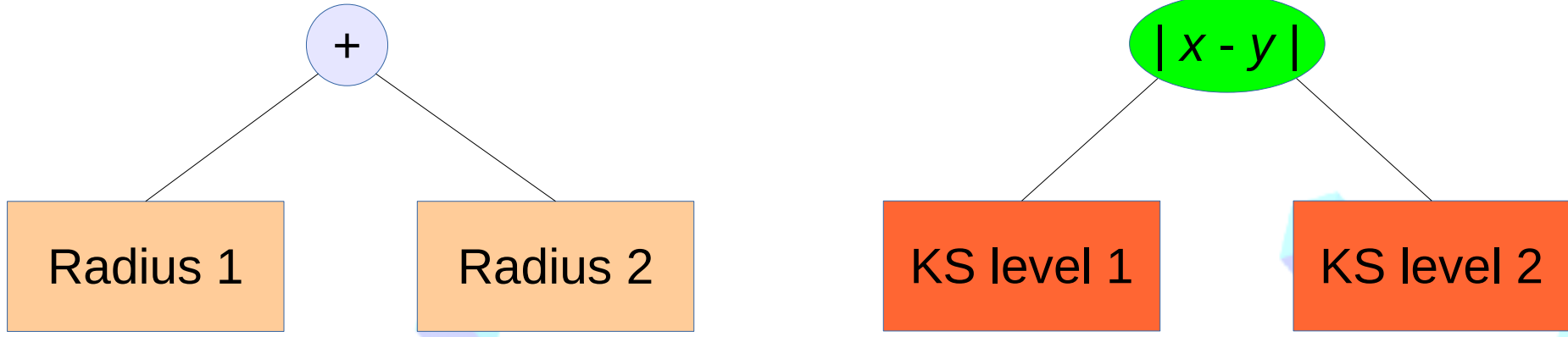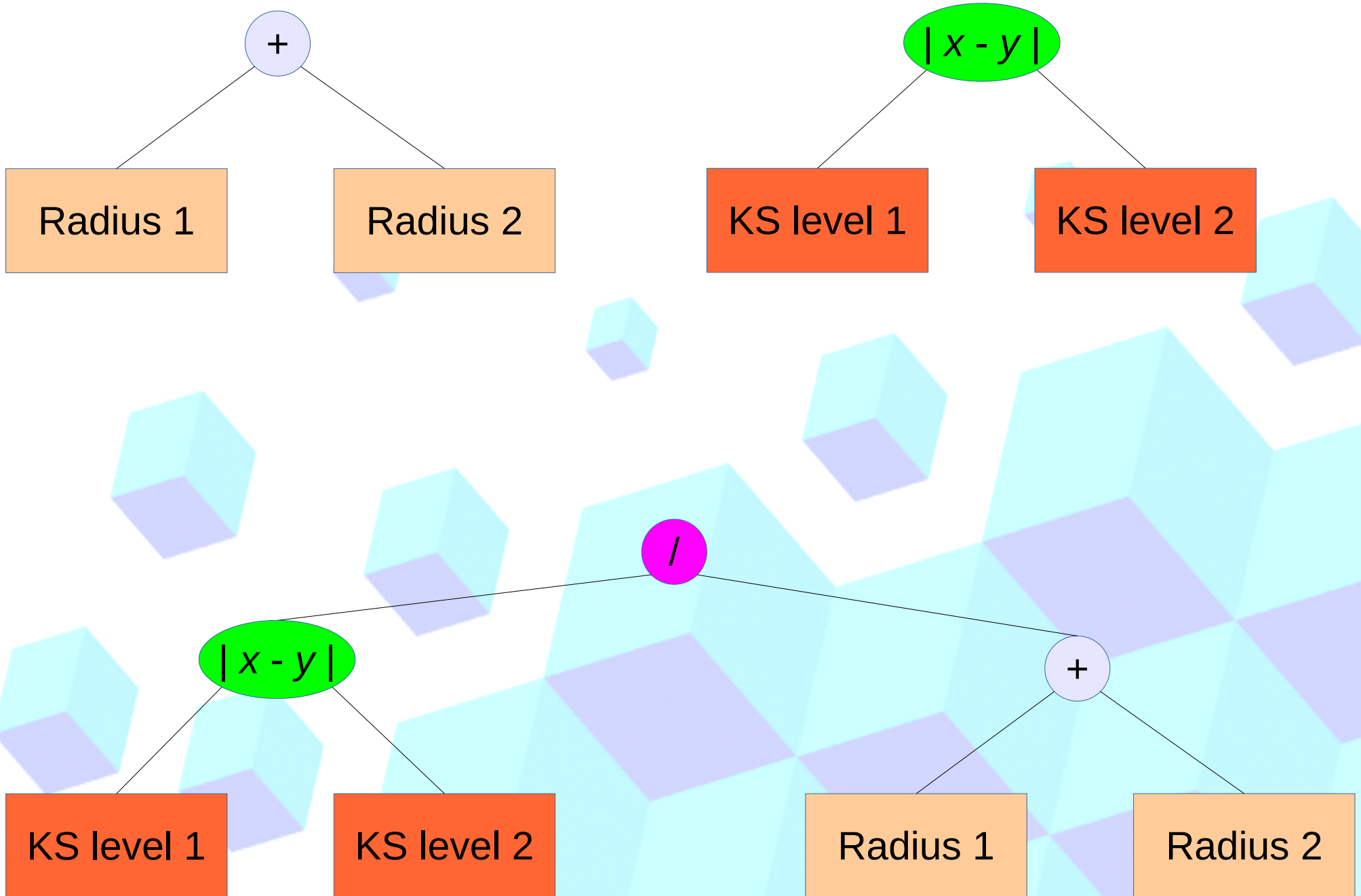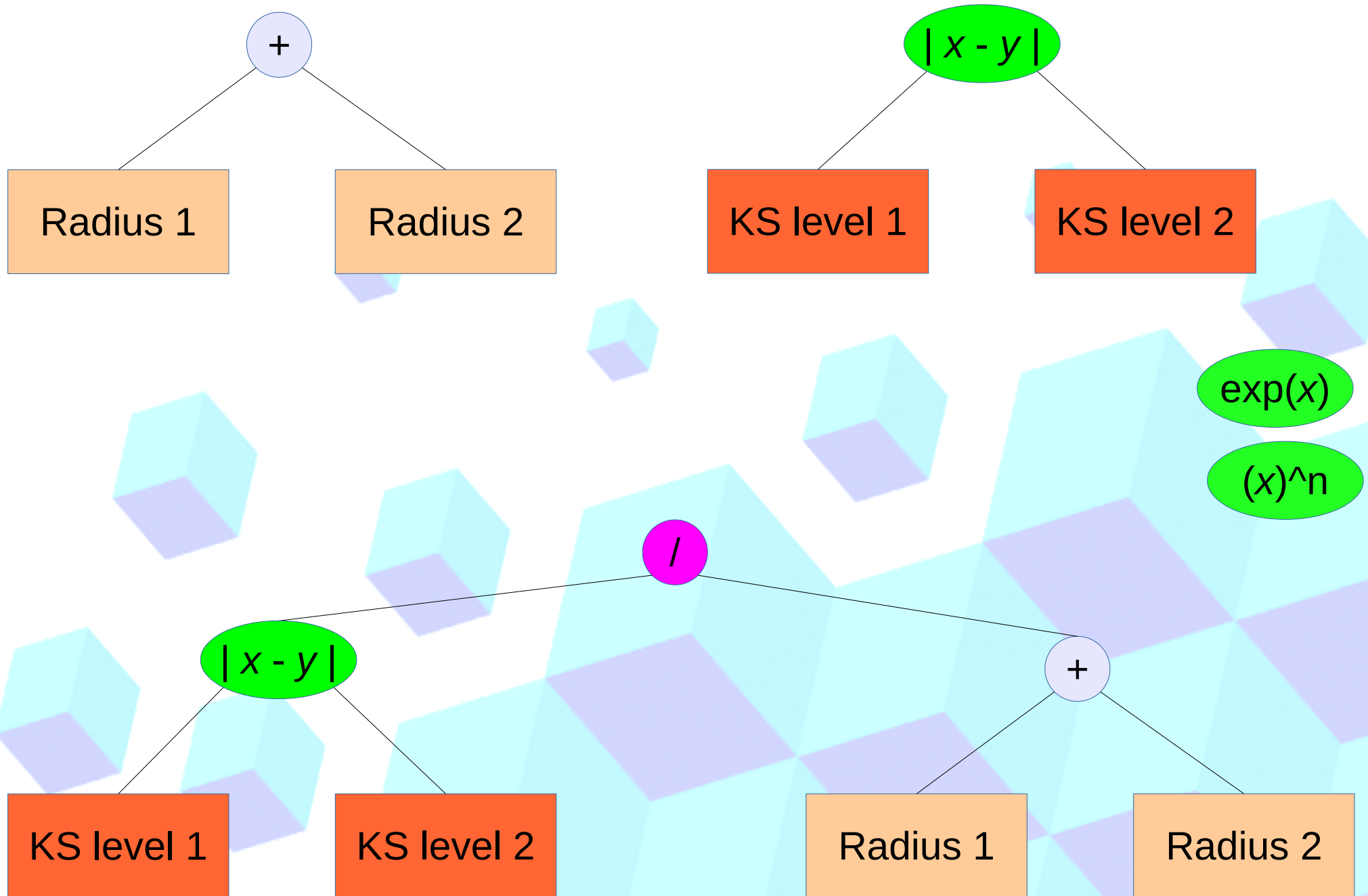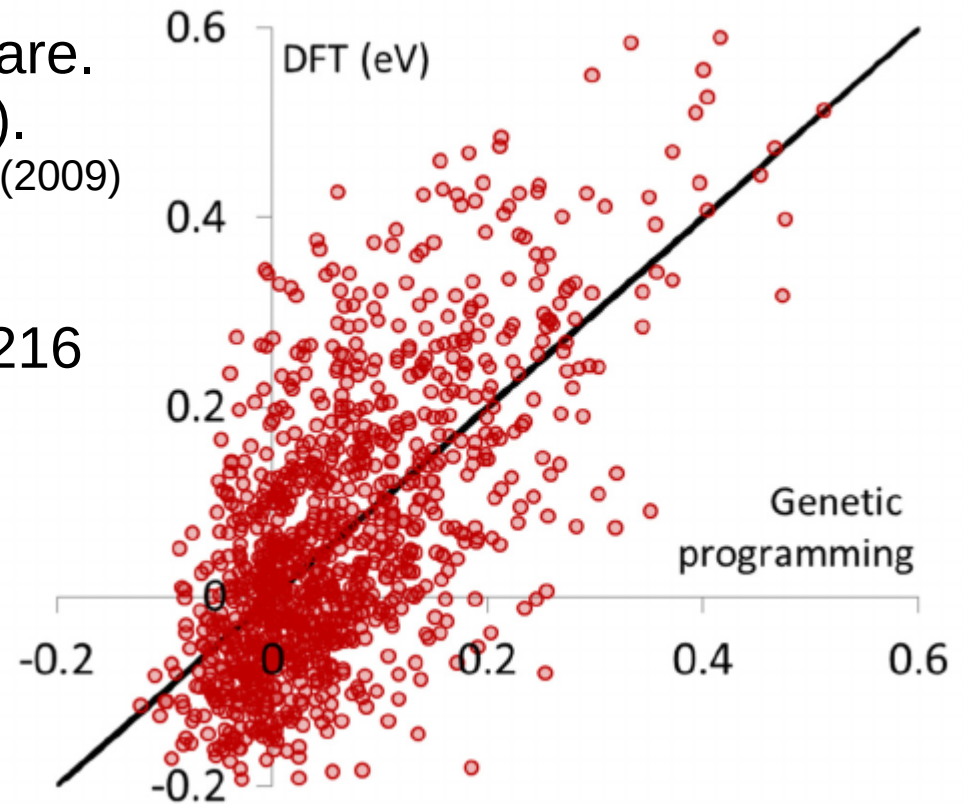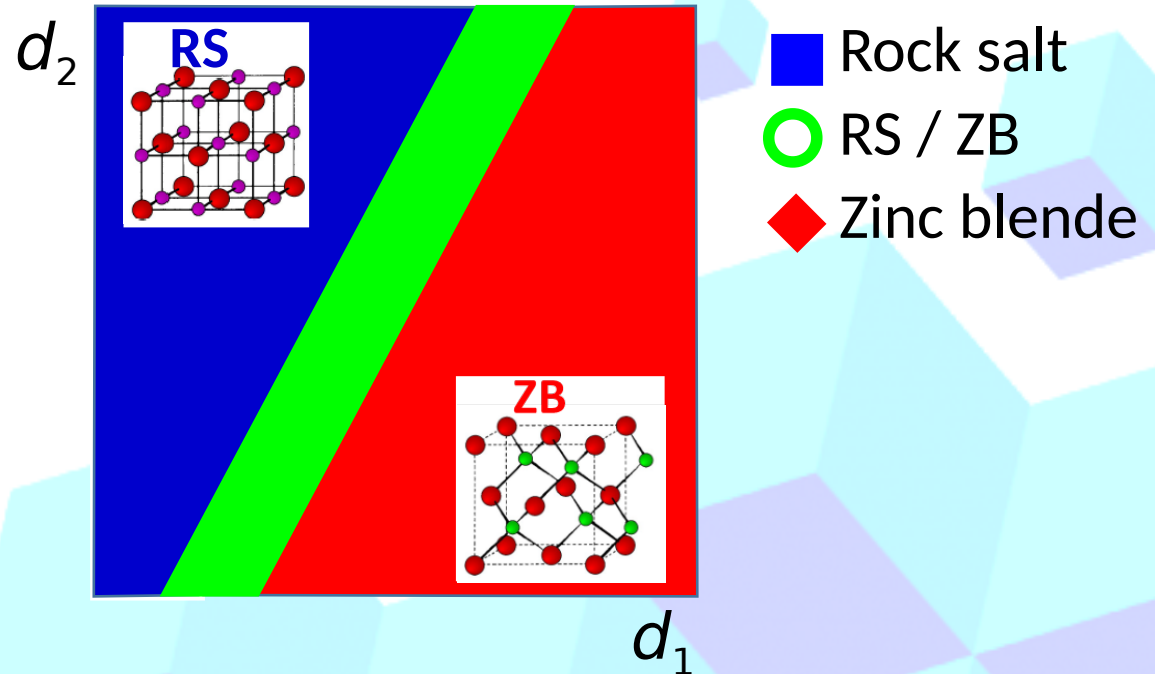$d_2$

**RS**

**ZB**

$d_1$

- Rock salt
- RS / ZB
- Zinc blende

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

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

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

# Compressed sensing: the quest for descriptors and predictive models

Ideal method: regression with $\ell_0$ regularization

$$\underset{\boldsymbol{c} \in \mathbb{R}^M}{\mathrm{argmin}} \left( \| \boldsymbol{P} - \boldsymbol{D}\boldsymbol{c} \|_2^2 + \lambda \| \boldsymbol{c} \|_0 \right)$$

Optimal solution
Non-polinomial complexity
Small # columns in $\boldsymbol{D}$

# Compressed sensing: the quest for descriptors and predictive models

Ideal method: regression with $\ell_0$ regularization

$$\underset{\boldsymbol{c}\in\mathbb{R}^M}{\operatorname{argmin}}\left(\left\|\boldsymbol{P}-\boldsymbol{D}\boldsymbol{c}\right\|_2^2 + \lambda\left\|\boldsymbol{c}\right\|_0\right)$$

Optimal solution
Non-polinomial complexity
Small # columns in $\boldsymbol{D}$

| | |
|---|---|
| $\left\|\boldsymbol{c}\right\|_0$ | # of nonzero elements of $\boldsymbol{c}$ |
| $\left\|\boldsymbol{c}\right\|_2$ | Euclidean. Square root of sum of squares of the elements of $\boldsymbol{c}$) |

# Compressed sensing: the quest for descriptors and predictive models

Ideal method: regression with $\ell_0$ regularization

$$\underset{\boldsymbol{c} \in \mathbb{R}^M}{\operatorname{argmin}} \left( \|\boldsymbol{P} - \boldsymbol{D}\boldsymbol{c}\|_2^2 + \lambda\|\boldsymbol{c}\|_0 \right)$$

Optimal solution
Non-polinomial complexity
Small # columns in $\boldsymbol{D}$

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

For matrices D with uncorrelated columns: LASSO

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

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

| | |
|---|---|
| $\|\boldsymbol{c}\|_1$ | "Manhattan". Sum of absolute values of the elements of $\boldsymbol{c}$ |

# Compressed sensing: the quest for descriptors and predictive models



min $\ell_1$ norm

min $\ell_2$ norm

$$\operatorname*{argmin}_{\boldsymbol{c} \in \mathbb{R}^M} \lVert \boldsymbol{P} - \boldsymbol{Dc} \rVert_2^2 + \lambda \lVert \boldsymbol{c} \rVert_1$$

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

$\lVert \boldsymbol{c} \rVert_1$    "Manhattan". Sum of absolute values of the elements of $\boldsymbol{c}$

## 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**[a,1], **Rongjie Lai**[b,1], **Russel Caflisch**[c,1], and **Stanley Osher**[c,1,2]

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

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

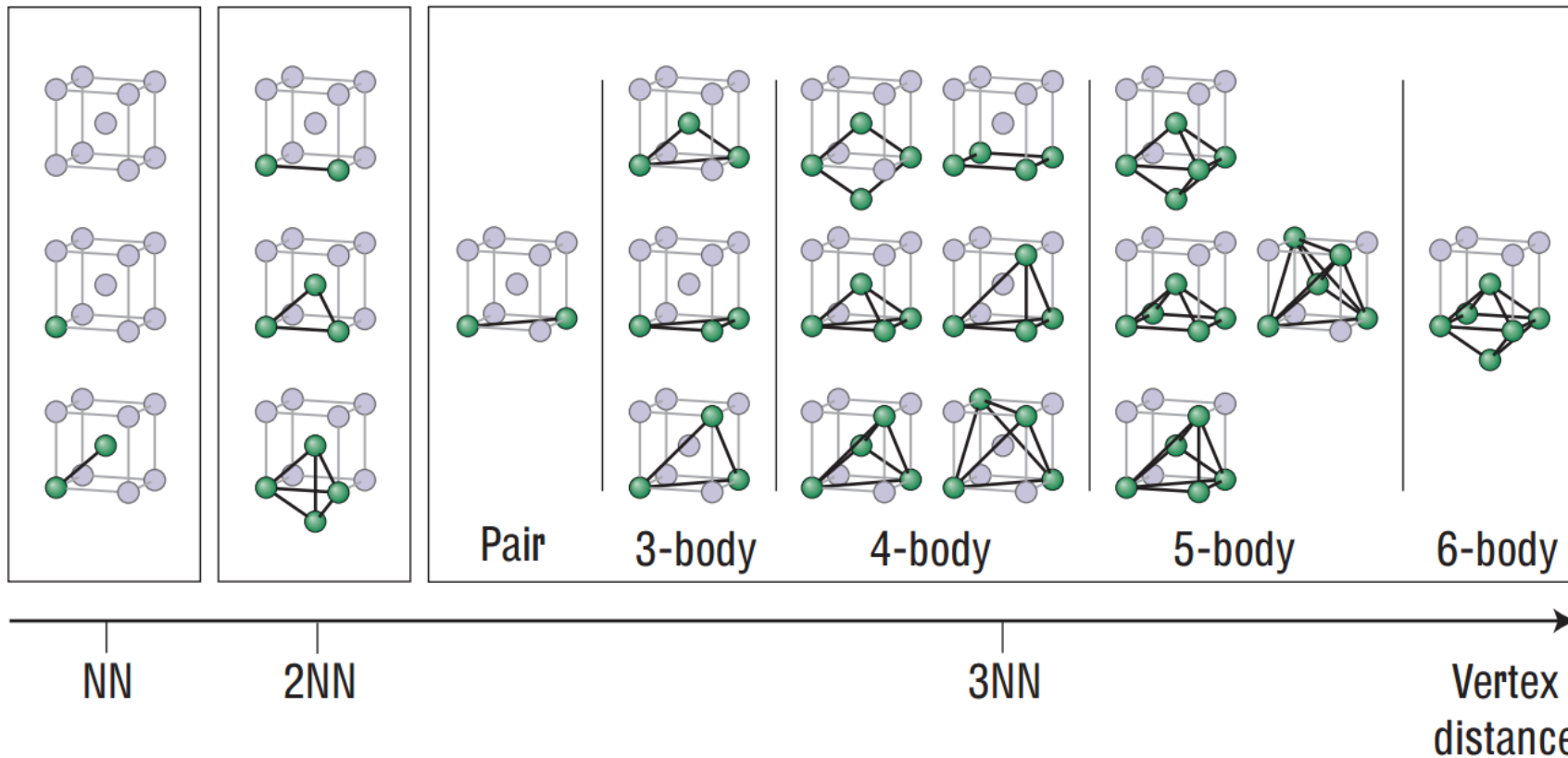Lance J. Nelson and Gus L. W. Hart

*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*

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

# 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
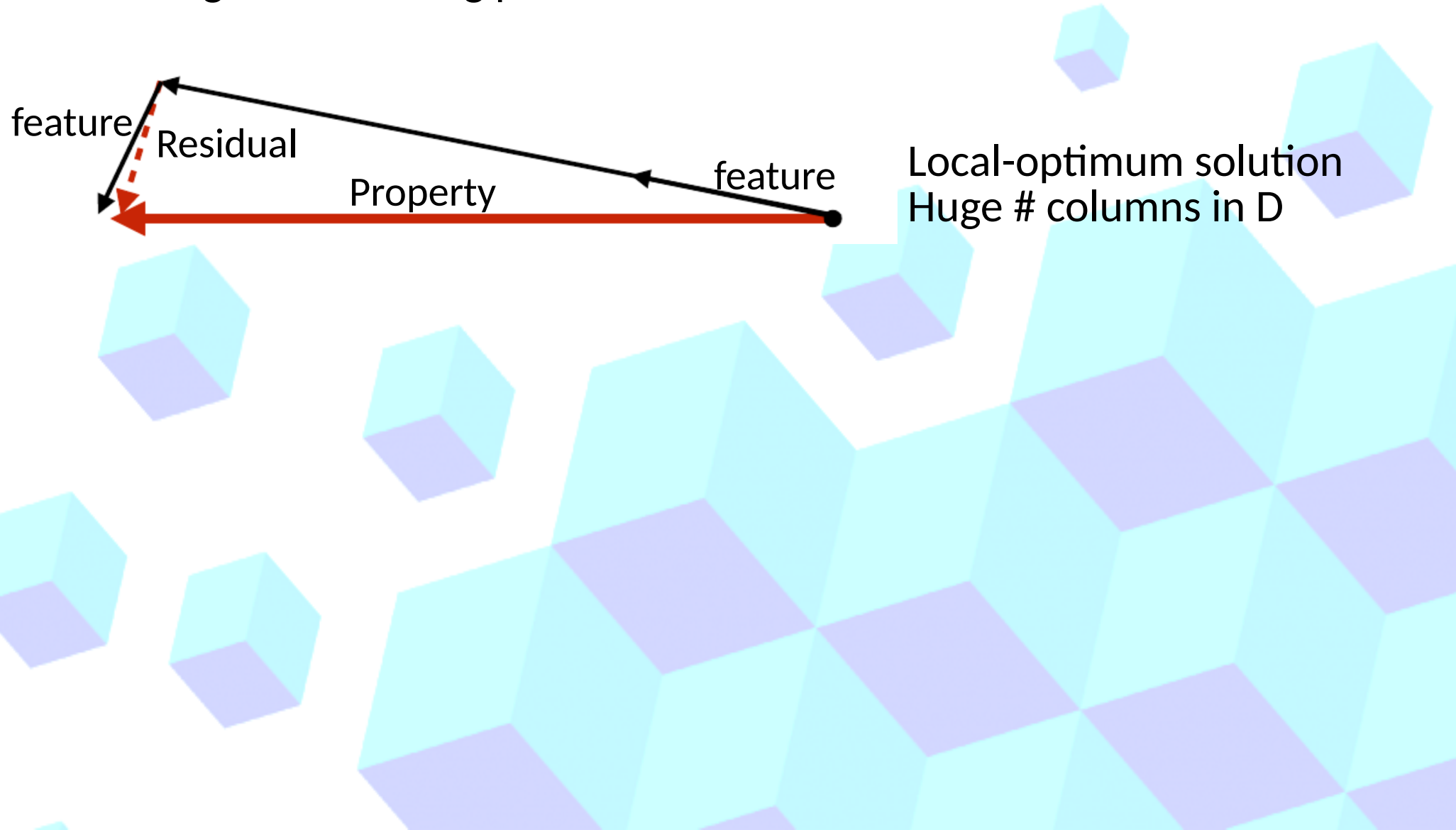- Billions of non-linear functions of the above



$d_2$

RS

ZB

$d_1$

- Rock salt
- RS / ZB
- Zinc blende

$E$(Rock salt) − $E$(Zinc blende)

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

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

# Compressed sensing: the quest for descriptors and predictive models

From orthogonal matching pursuit ….

feature

Residual

Property

feature

Local-optimum solution
Huge # columns in D

# Compressed sensing: the quest for descriptors and predictive models

From orthogonal matching pursuit ....



feature

Residual

Property

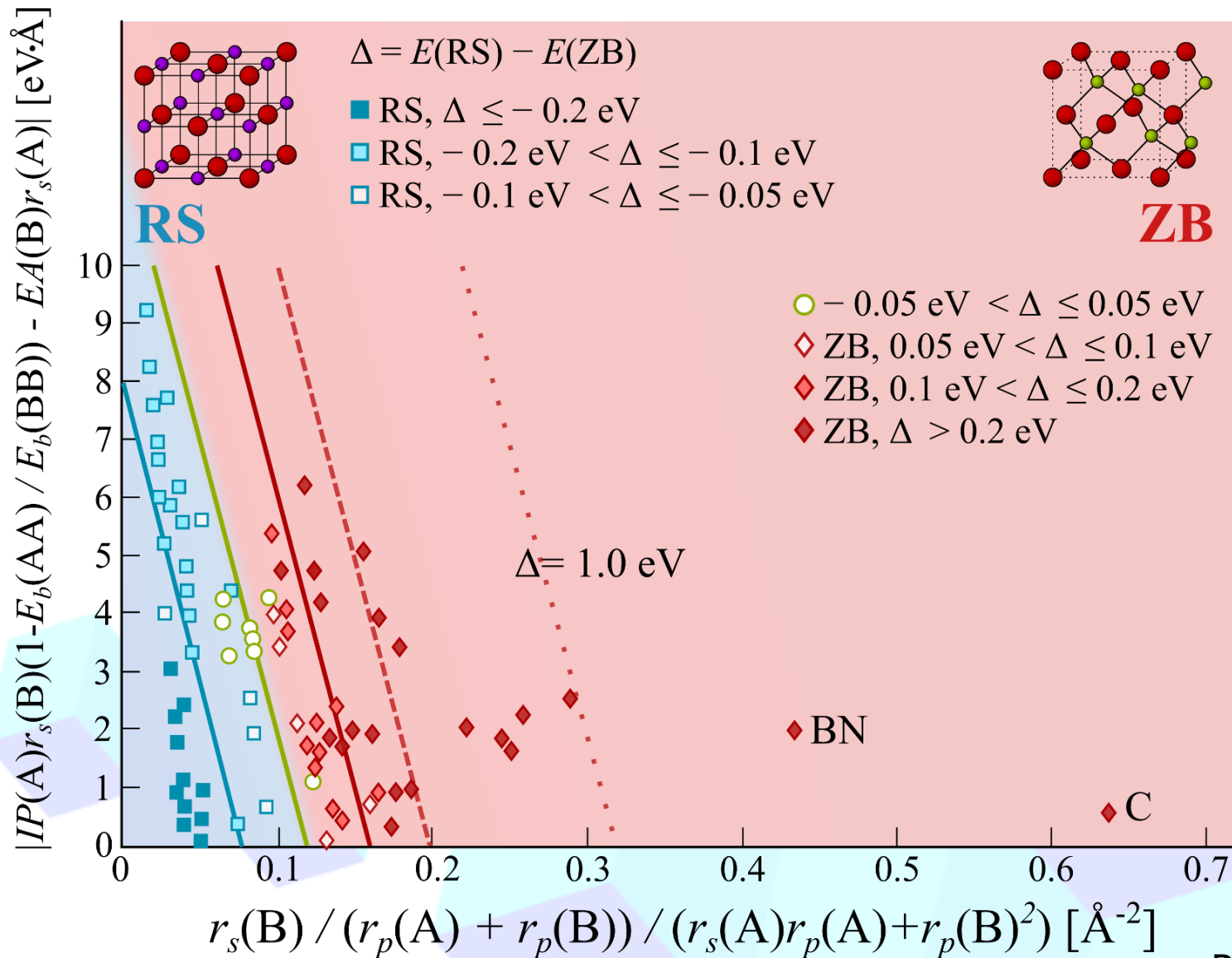feature

Local-optimum solution
Huge # columns in D

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



features $S_1$

Residual$_1$

Property

features $S_0$

Proxy of
global-optimum solution
Huge # columns in D

R. Ouyang *et al*. PRM **2**, 083802 (2018), published 7 August 2018)

# 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

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

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.

# Charts/maps of materials

$$\underset{\boldsymbol{c} \in \mathbb{R}^M}{\mathrm{argmin}}(\|\boldsymbol{P} - \boldsymbol{D}\boldsymbol{c}\|_2^2 + \lambda\|\boldsymbol{c}\|_0)$$

New cost function to be minimized: overlap of *convex* domains
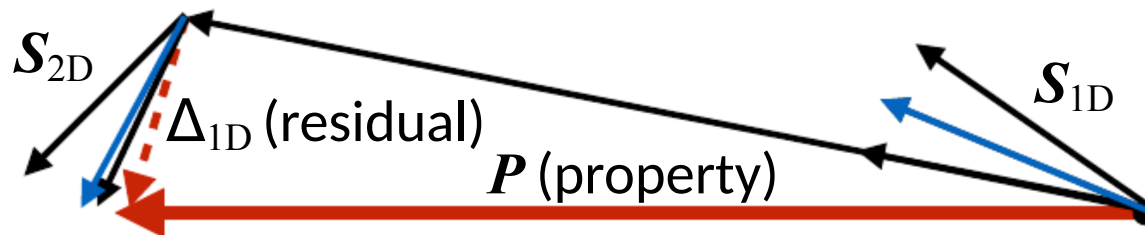
$d'_2$

Metal

$A_xB_y$
binaries

Insulator

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))

$d'_1$

$E_a$ aterials

$d*_2$

opological
nsulator

Iterative generation of feature subspaces

$\boldsymbol{S}_{2D}$

$\Delta_{1D}$ (residual)

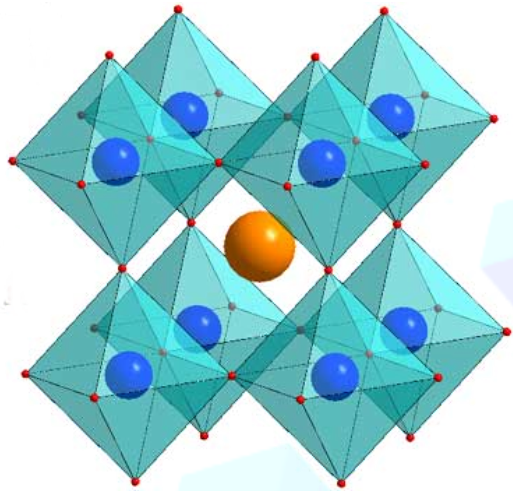$\boldsymbol{S}_{1D}$

$\boldsymbol{P}$ (property)

$d''_1$

# Perovskites' stability: an improved Goldschmidt Tolerance Factor



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

$ABX_3$

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

# Perovskites' stability: an improved Goldschmidt Tolerance Factor

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

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

→ Oxidation state

↘ 1 / μ = Octahedral factor

$ABX_3$

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

Our stable perovskites: **τ** < 4.18, accuracy 92%

Bartel, Sutton, Goldsmith, Ouyang, Musgrave, LMG &Scheffler, Sci. Adv. *5, eaav0693 (2019)*

# Perovskites' stability: an improved Goldschmidt Tolerance Factor

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

Oxidation state

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

$1 / \mu$ = Octahedral factor

$ABX_3$

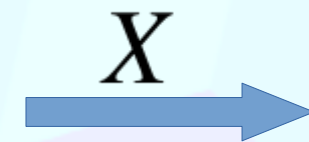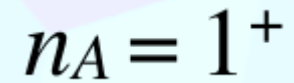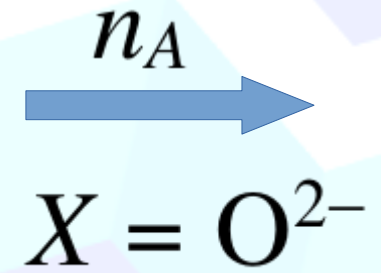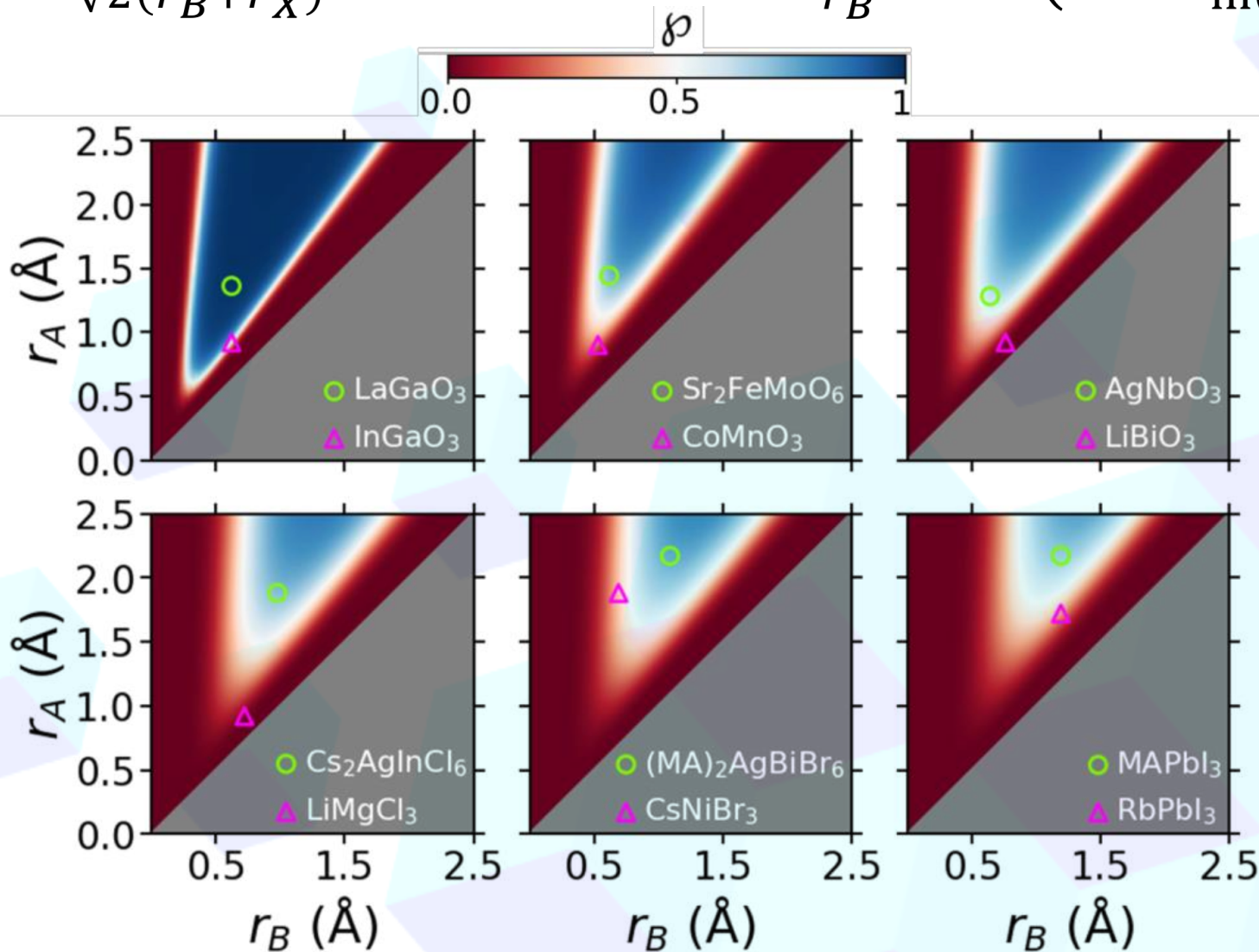Goldschmidt* stable perovskites: $0.825 < t < 1.059$, accuracy 79%

Our stable perovskites:         $\tau < 4.18$,         accuracy 92%
$\tau < 3.31$ or $\tau > 5.92$, 99% accuracy (1/3 of the training data)
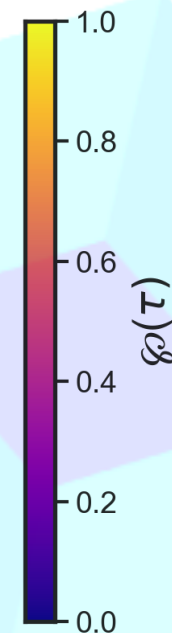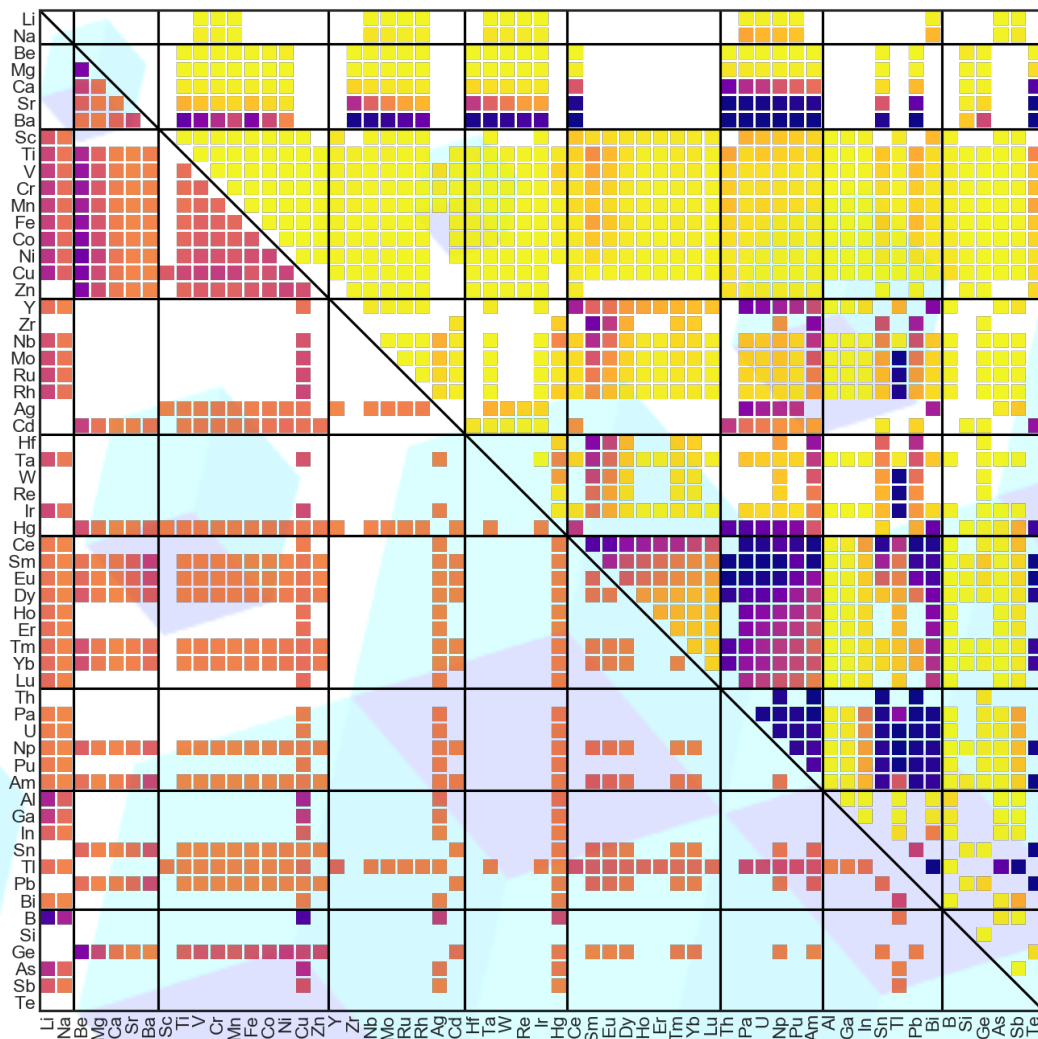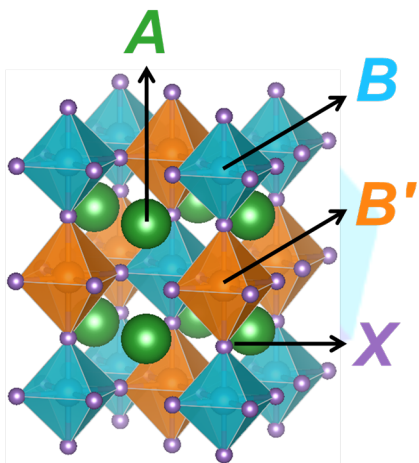$\tau < 3.31$ or $\tau > 12.08$, 100% accuracy (1/4 of the training data)

Bartel, Sutton, Goldsmith, Ouyang, Musgrave, LMG &Scheffler, Sci. Adv. *5, eaav0693 (2019)*

# Improved Goldschmidt Tolerance Factor: Materials design

$$t = \frac{r_A + r_X}{\sqrt{2}(r_B + r_X)} \qquad \longrightarrow \qquad \tau = \frac{r_X}{r_B} - n_A\left(n_A - \frac{r_A/r_B}{\ln(r_A/r_B)}\right)$$
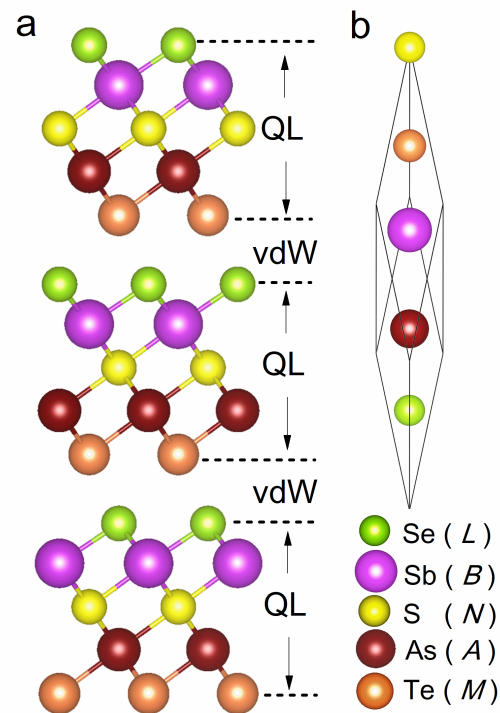
# Improved Goldschmidt Tolerance Factor: Extension of the materials space

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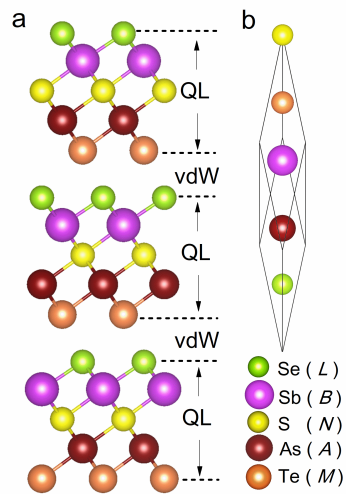


La$_2$BB'O$_6$

Cs$_2$BB'Cl$_6$

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



a b QL vdW QL vdW QL

Se ( *L* )
Sb ( *B* )
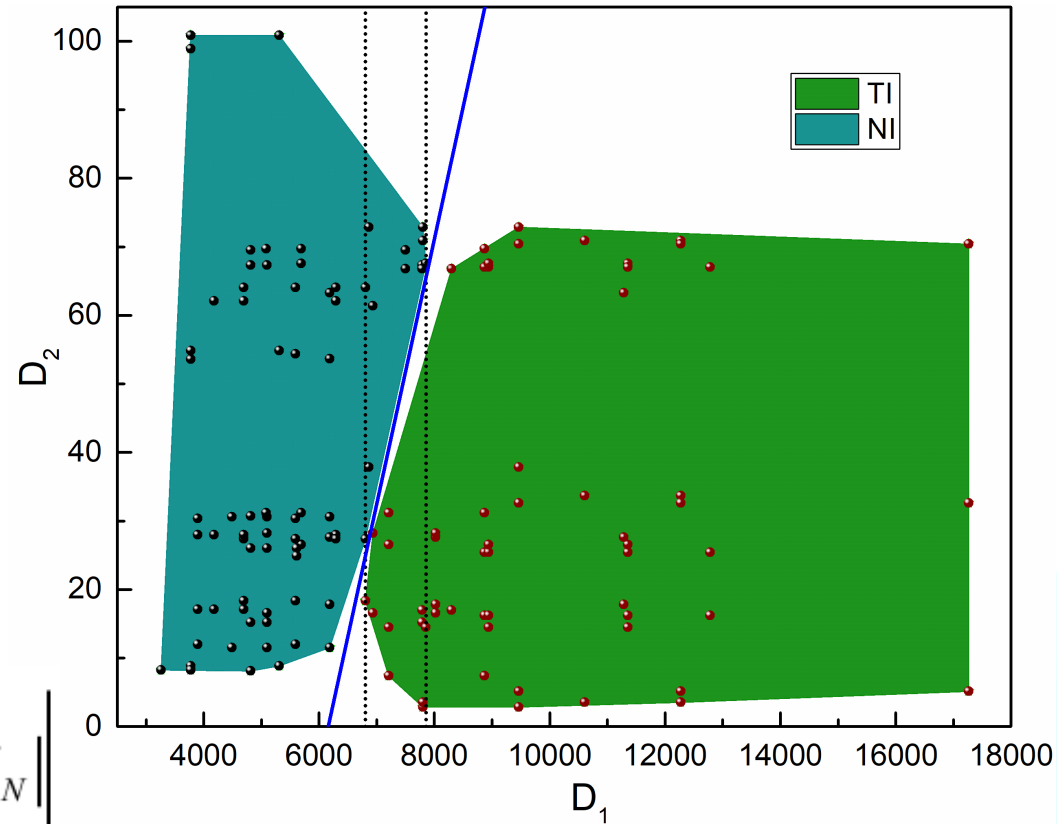S ( *N* )
As ( *A* )
Te ( *M* )

# SISSO: predicting new tetradymite topological insulators

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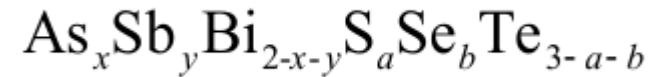*AB-LNM*
AB = {As,Sb,Bi}
LNM = {S, Se,Te}



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

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

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



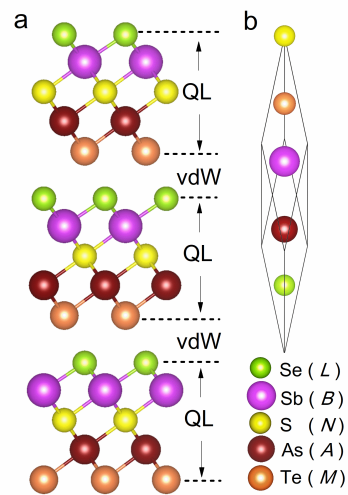$$As_x Sb_y Bi_{2-x-y} S_a Se_b Te_{3-a-b}$$

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

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

Cao, Liu, Ouyang, LMG, Zhou, Scheffler, Zhang, Carbogno, submitted (2019)

# SISSO: predicting new tetradymite topological insulators

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

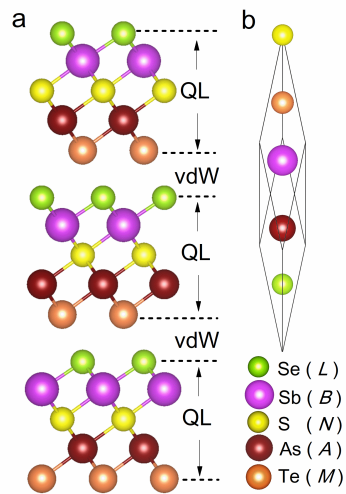Se ( $L$ )
Sb ( $B$ )
S  ( $N$ )
As ( $A$ )
Te ( $M$ )

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

$$D_2 = \left| \frac{(\chi_M + \chi_N) \cdot Z_E}{\chi_A} - (Z_M + Z_N) - |Z_M - Z_N| \right|$$
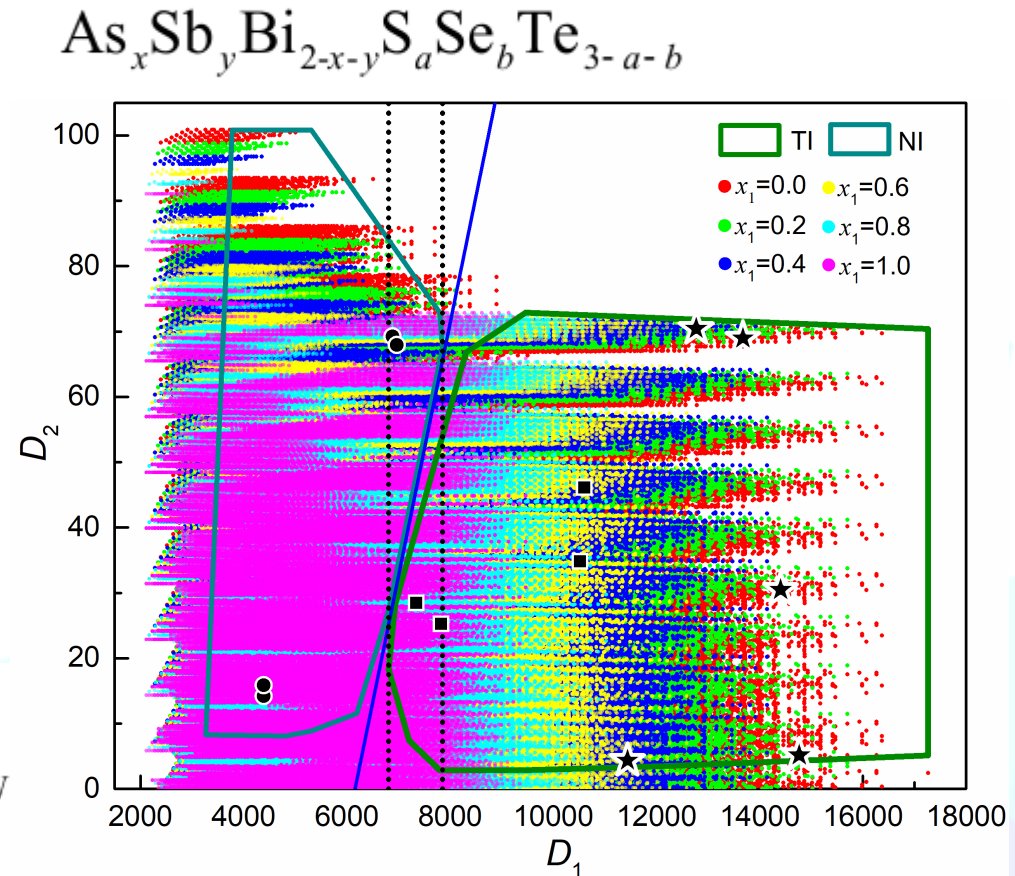
$As_x Sb_y Bi_{2-x-y} S_a Se_b Te_{3-a-b}$

TI    NI
$x_1$=0.0    $x_1$=0.6
$x_1$=0.2    $x_1$=0.8
$x_1$=0.4    $x_1$=1.0

$D_2$

$D_1$

Cao, Liu, Ouyang, LMG, Zhou, Scheffler, Zhang, Carbogno, submitted (2019)

# Acknowledgements

**Compressed sensing and SISSO**

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

**Application of SISSO to perovskites**

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

**Application of SISSO to topological insulators**

Guohua Cao, Runhai Ouyang, Huijun Liu, Zizhen Zhou, Zhenyu Zhang, Christian Carbogno

**And**

Matthias Scheffler