# Patterns, Correlations, and Discovery of Descriptors in Big Data of Materials

Bryan R. Goldsmith
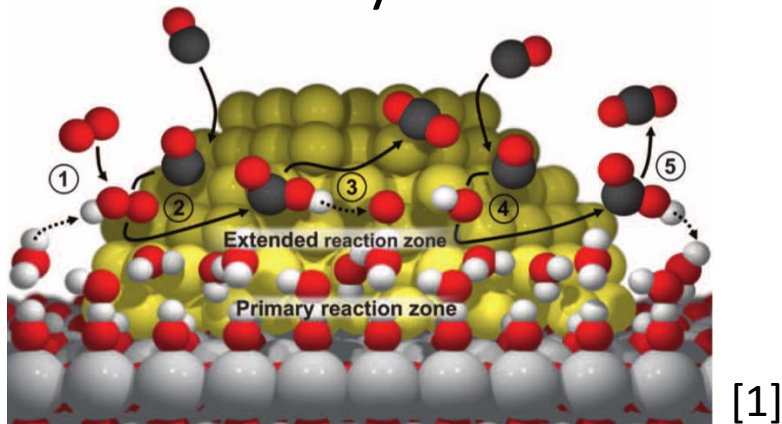
Fritz-Haber-Institut der Max-Planck-Gesellschaft
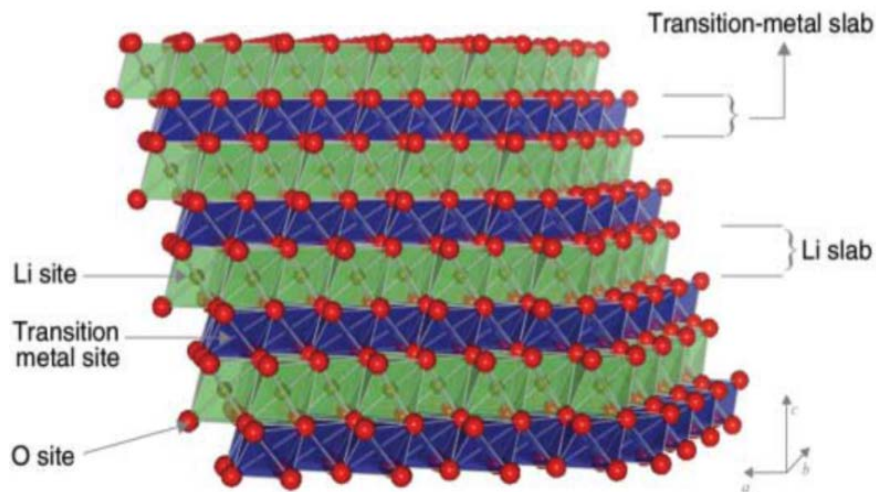Theory Department

# Designing materials requires an understanding
# of the mechanisms underlying a materials function

## Catalysts



① ② ③ ④ ⑤

**Extended reaction zone**

**Primary reaction zone**

[1]

## Thermoelectrics



Conducting plate

Electron

$T_{hot}$

p-type

n-type

E

Hole

E

Conducting plate

$T_{cold}$

Resistor

Electrical current

## Batteries



Transition-metal slab

Li slab

Li site

Transition metal site
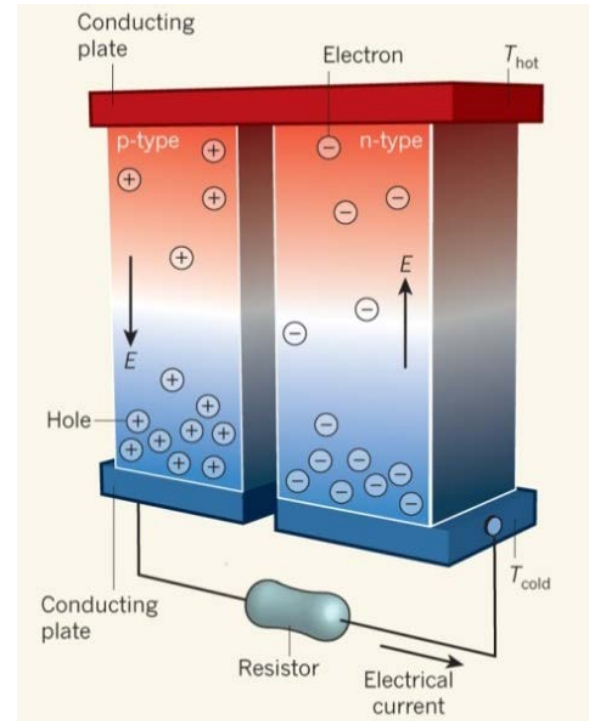
O site

[2]

[1] J. Saavedra *et al. Science* 345, 1599 (2014)

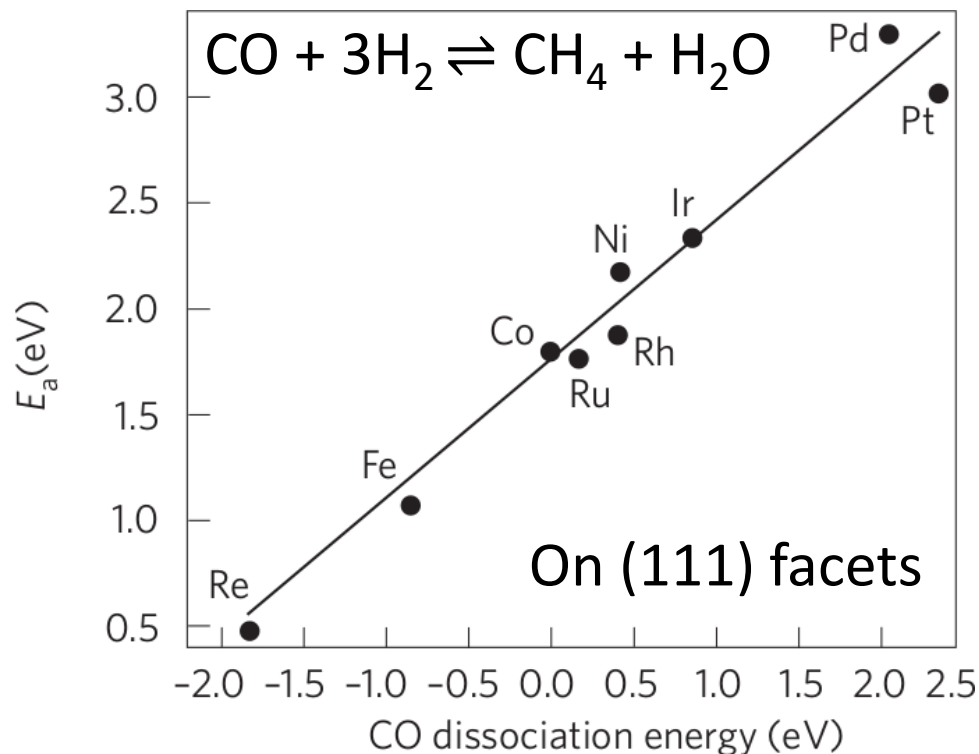[2] K. Kang *et al. Science* 311, 977 (2006)

# Identifying physically meaningful *descriptors* that describe materials properties is critical

Descriptor → Property

Descriptor = function(atomic or material features)

Predict new materials

Increase material understanding



$CO + 3H_2 \rightleftharpoons CH_4 + H_2O$

On (111) facets

J. K. Nørskov *et al. Nat. Chem.* 1, 37 (2009)

The development of data-analytic tools can facilitate the discovery of descriptors

# This talk will focus on two different data-analytics tools to find descriptors of materials

**1. Compressed sensing to find interpretable descriptors**

Application 1: octet binary semiconductors

**2. Subgroup discovery to find local patterns and their descriptions**

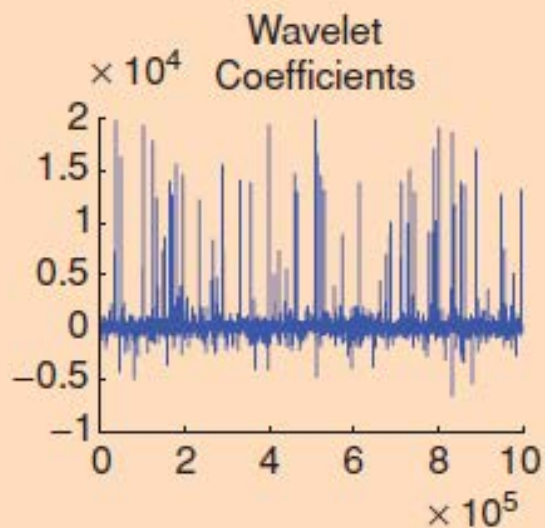Application 1: gold clusters in the gas phase (sizes 5-14 atoms)

Application 2: octet binary semiconductors

# Compressed sensing allows the construction of sparse models with high accuracy

Original image

Sparse
in the basis set

Recovered with 10%
measurements



Emmanuel Candès, Terence Tao, and David Donoho

# Compressed sensing allows the construction of sparse models with high accuracy

### $l_0$-norm minimization

$$\min \|\beta\|_0 \text{ subject to } y = D\beta$$

$l_0$-**norm:** total # of non-zero coefficients

Target material property

Material feature matrix

coefficients

Emmanuel Candès, Terence Tao, and David Donoho

$l_0$-norm minimization is too expensive
to perform for large feature matrix **D**

Instead minimize $l_1$-norm (LASSO)
as approximation of $l_0$-norm

**$l_1$-norm:**
Sum of absolute value
of coefficients

$$\hat{\beta}_{LASSO}(\lambda) = \underset{\beta}{\mathrm{argmin}}\left(\frac{1}{2}\|y - \boldsymbol{D}\beta\|_2^2 + \lambda\|\beta\|_1\right)$$
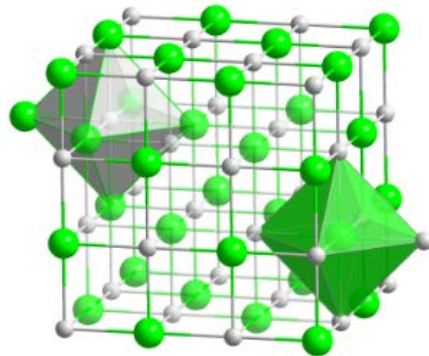
Root mean
squared error

Regularization
parameter

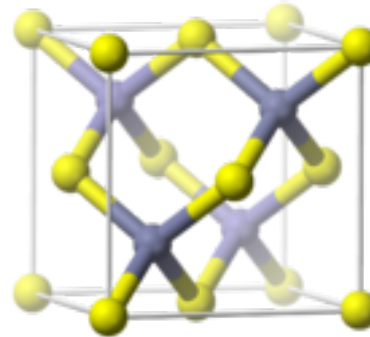# Find descriptors that predict the crystal structure energy differences between 82 octet binary compounds

Target property = $E_{rocksalt} - E_{zincblende}$

### Rocksalt (RS)          Zincblende (ZB)



vs.

Predict energy differences between 0.01 – 1.5 eV

**LASSO+$l_0$:** L. M. Ghiringhelli, J. Vybiral, S. V. Levchenko, C. Draxl, M. Scheffler, *PRL 114,* 105503 (2015)

# 1. Build a large feature space of candidate descriptors

## Start from atomic properties (primary features)

$$IP(A) \quad EA(A) \quad IP(B) \quad EA(B)$$
$$H(A) \quad L(A) \quad H(B) \quad L(B)$$

$$r_s(A) \quad r_p(A) \quad r_d(A)$$
$$r_s(B) \quad r_p(B) \quad r_d(B)$$

## Build up feature matrix $D$ of 4500 compound features

operator set: $\{+,-,\exp,^2,\div,\times\}$

$$r_s(A)^2, (r_p(A) + r_s(A))^2$$
$$\exp(r_s(A)), \exp(r_p(A) \pm r_s(A))$$

$$|IP(A) \pm IP(B)|$$
$$|L(B) \pm H(A)|$$
$$|r_p(A) \pm r_s(A)|$$

# 2. Feature selection using two step scheme: **LASSO+$l_0$**

## Step 1: **LASSO**

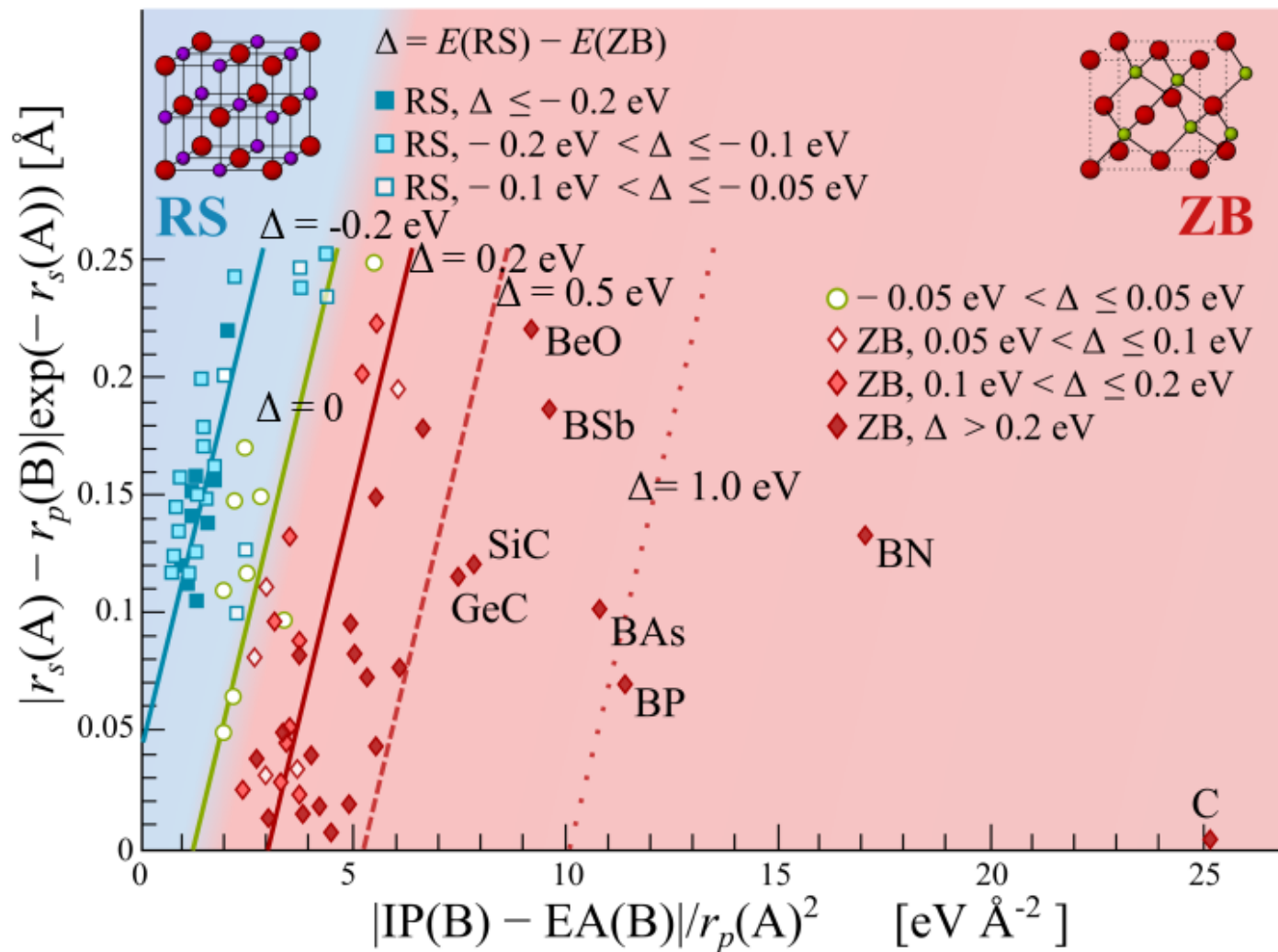Reduce the feature space to find candidate set of possible descriptors (e.g., 4500 → 50 features)

$$\hat{\beta}_{LASSO}(\lambda) = \underset{\beta}{\operatorname{argmin}} \left( \frac{1}{2} \|y - \boldsymbol{D}\beta\|_2^2 + \lambda\|\beta\|_1 \right)$$

## Step 2: $l_0$ **minimization**

Find the 'best' $n$-dimensional model from the candidate set $\boldsymbol{D}`$

$$\hat{\beta}_{l_0}(\lambda) = \underset{\beta}{\operatorname{argmin}} \left( \frac{1}{2} \|y - \boldsymbol{D}`\beta\|_2^2 + \lambda\|\beta\|_0 \right)$$

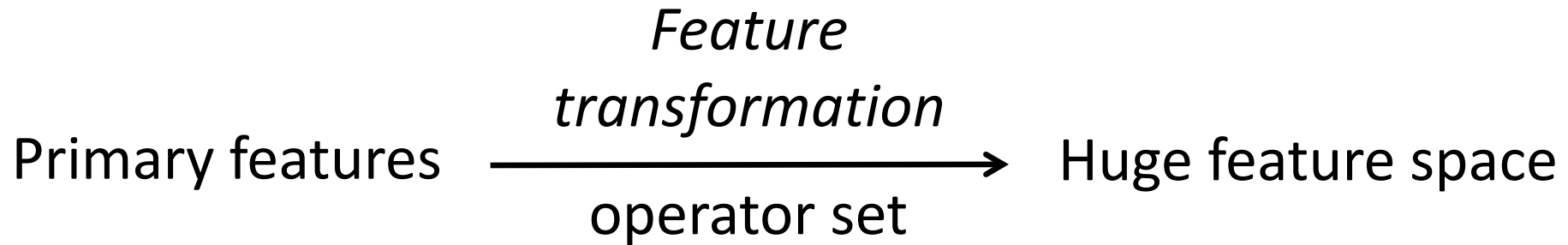# A highly predictive two-dimensional descriptor is found for the 82 octet binary semiconductors

# LASSO has stability issues for huge feature space of correlated features

## Managing high dimensional and correlated feature spaces by combining screening and compressed sensing



Runhai Ouyang, Luca M. Ghiringhelli,
Emhre Ahmetcik, Matthias Scheffler

# 1. Systematically construct a huge feature space (**one trillion**)

Primary features $\xrightarrow[\text{operator set}]{\textit{Feature transformation}}$ Huge feature space

$$\hat{R} = \{+, -, \times, \div, \exp, log, ^{-1}, ^2, ^3, sqrt, |-|\}$$

## 2. Select top ranked features using Sure Independent Screening (SIS)[1]

Sure independent screening

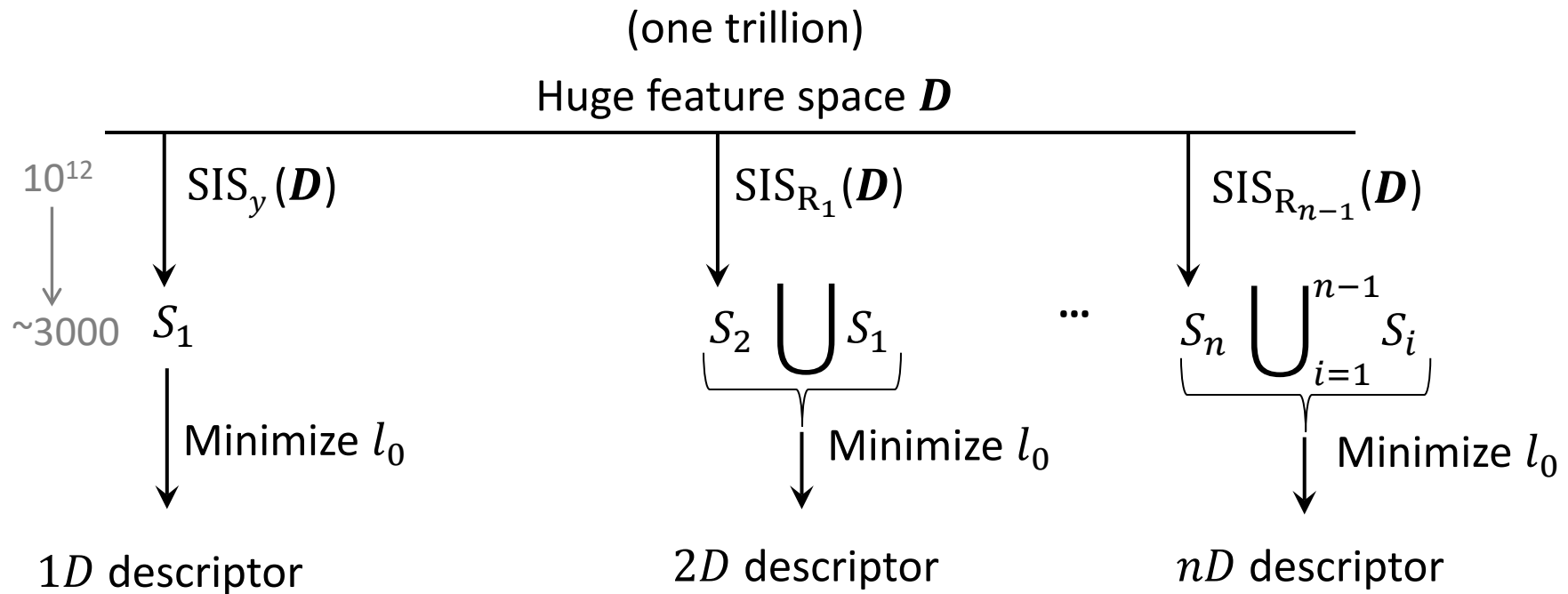➢ Select the $N$ largest components of $\boldsymbol{D}^T \mathrm{y}$

Results in a subspace of $N$ features that are most correlated with the property of interest

$\boldsymbol{D}$: matrix of the feature space

$y$: target property $(E_{\text{rocksalt}} - E_{\text{zincblende}})$

[1] J. Fan and J. Lv, J. R. Statist. Soc. B 70, 849 (2008)

# 3. Iteratively apply sure independent screening with a sparse approximation algorithm

(one trillion)

Huge feature space $\boldsymbol{D}$

$10^{12}$

$\downarrow$

~3000

| $\text{SIS}_y(\boldsymbol{D})$ | $\text{SIS}_{R_1}(\boldsymbol{D})$ | ... | $\text{SIS}_{R_{n-1}}(\boldsymbol{D})$ |
|---|---|---|---|
| $S_1$ | $S_2 \bigcup S_1$ | | $S_n \bigcup_{i=1}^{n-1} S_i$ |
| Minimize $l_0$ | Minimize $l_0$ | | Minimize $l_0$ |
| $1D$ descriptor | $2D$ descriptor | | $nD$ descriptor |

SIS = sure independent screening

$S_i$ = feature subspace        $R_i$ = Residual of target property using
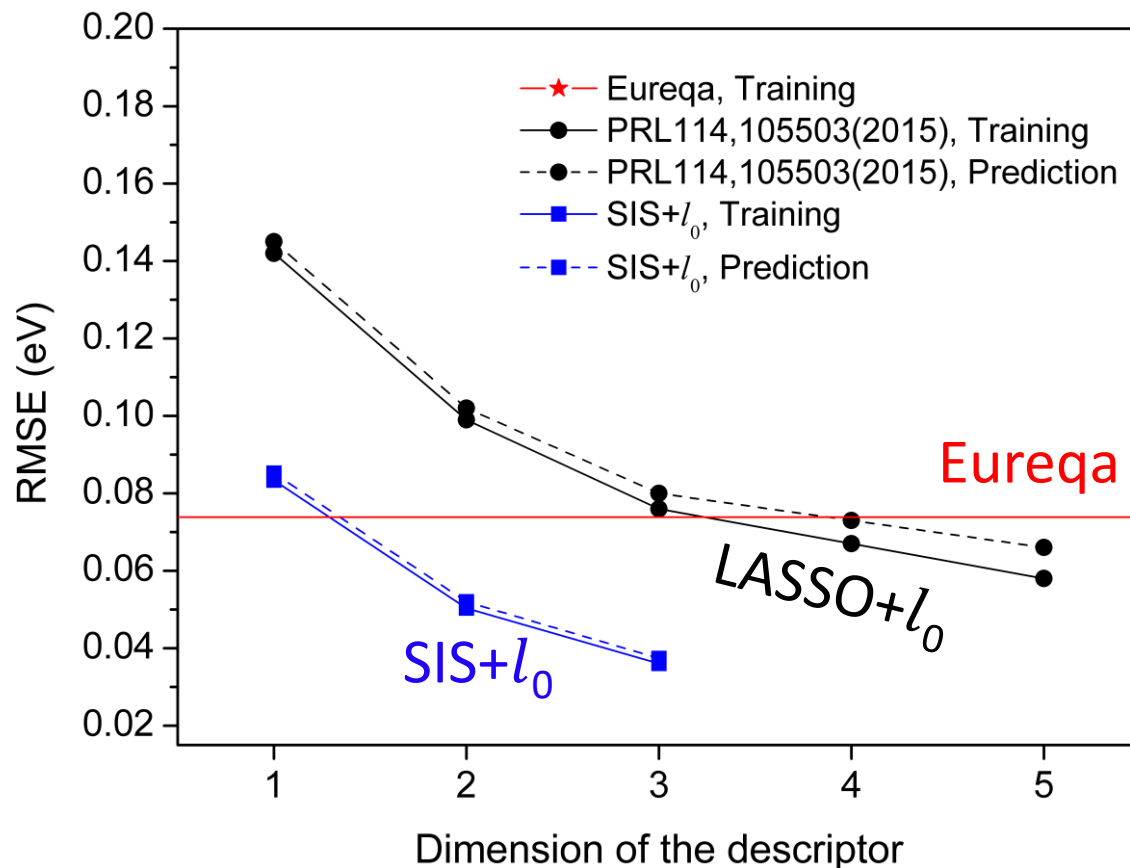
$y$ = target property            the previous iterations least squares prediction

# Accurate and interpretable descriptors are found from a **one trillion** dimension feature space

Target property = $E_{\text{rocksalt}} - E_{\text{zincblende}}$

**1D descriptor:**
$$\frac{|r_s(B) - r_d(A)| * r_s(B)}{r_d(A) * (r_p(A)^3 + r_p(B)^3)}$$



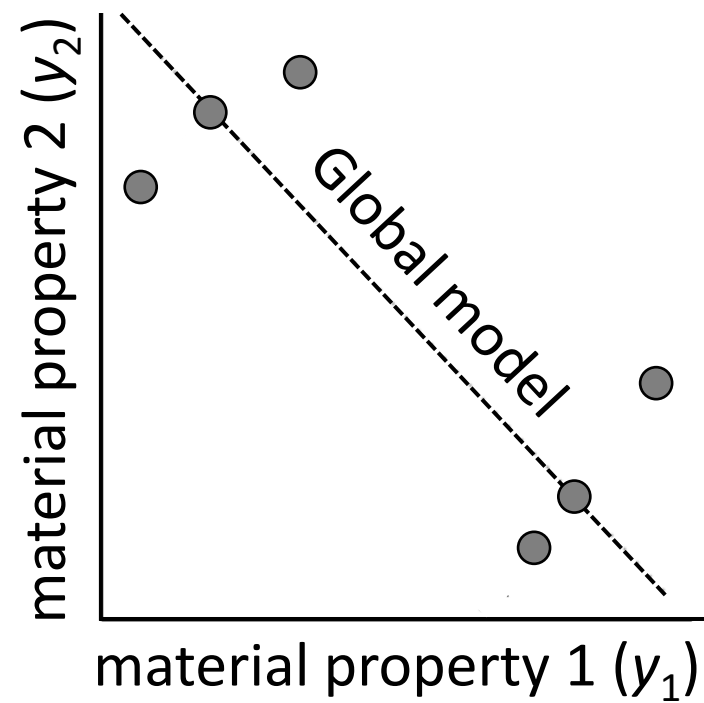R. Ouyang, E. Ahmetcik,  L.  M. Ghiringhelli, M. Scheffler, *unpublished*

# What if data is heterogeneous and the global model is too complex to be interpretable?

Underlying mechanisms can change across materials

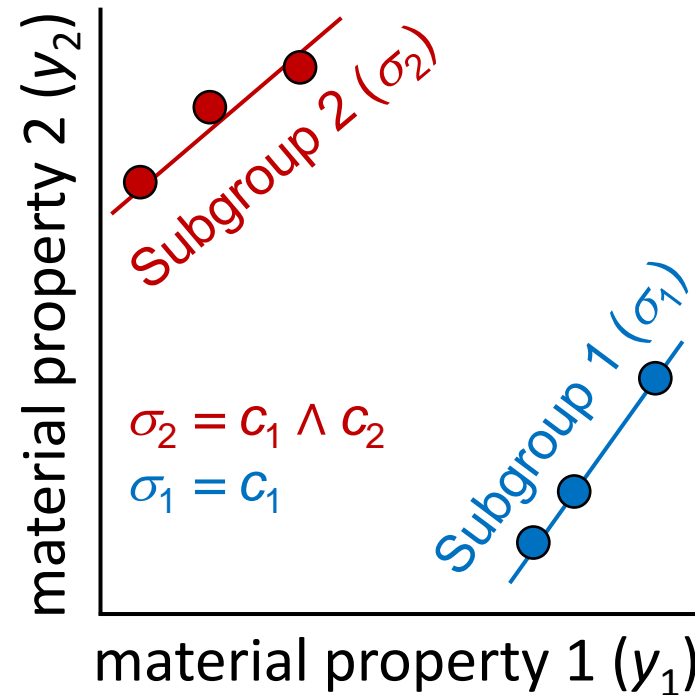Relations between subsets of data may be important

Goal: find *physically interpretable local models* of a target property in materials data

# Subgroup discovery: find physically interpretable local models of a target property in materials data
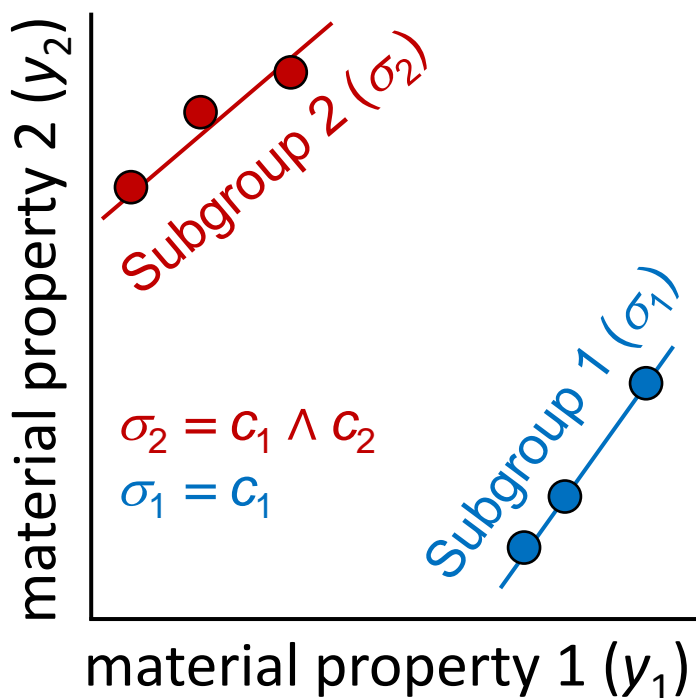


Mario Boley, me,
Luca Ghiringhelli



$\sigma_2 = c_1 \wedge c_2$
$\sigma_1 = c_1$

M. Atzmueller, *WIREs Data Min. Knowl. Discov. 5* (2015)
W. Duivesteijn, A. J. Feelders, A. Knobbe, *Data Min. Knowl. Discov. 30*, 47 (2016)

# Subgroup discovery: how it works

*Descriptive features,* $a_1, \ldots, a_m \in A$

*e.g.,* energy, bonding topology, number of atoms

*Target features* $y_1, \ldots, y_n \in Y$

*e.g.,* HOMO-LUMO energy gap

*Basic selectors,* $c_1, \ldots, c_k \in C \rightarrow \{$false, true$\}$

*e.g.,* Is there an even number of atoms?

Find *selector* $\sigma = c_1(\cdot) \wedge \cdots \wedge c_l(\cdot)$

that maximizes *quality* $q = \left( \frac{|\text{ext}(\sigma)|}{|P|} \right)^{\alpha} u(Y_{\sigma})^{1-\alpha}$

$\frac{|\text{ext}(\sigma)|}{|P|}$ is the coverage of points where $\sigma$ is true

$u(Y_{\sigma})$ is the utility function (optimization criteria)



material property 2 ($y_2$)

Subgroup 2 ($\sigma_2$)

Subgroup 1 ($\sigma_1$)

$\sigma_2 = c_1 \wedge c_2$
$\sigma_1 = c_1$

material property 1 ($y_1$)

Subgroup discovery (SGD) implemented
by Mario Boley using Creedo and realKD library

# Two applications of subgroup discovery

## 1. Gas Phase Gold Clusters (size 5-14)

Display interesting optical, chemical, and electronic properties



*ca.* 25,000 gold cluster configurations in total

## 2. Classification of 82 Octet Binary Semiconductors

Rocksalt    vs.    Zincblende

# Rediscover simple insight about HOMO-LUMO gap

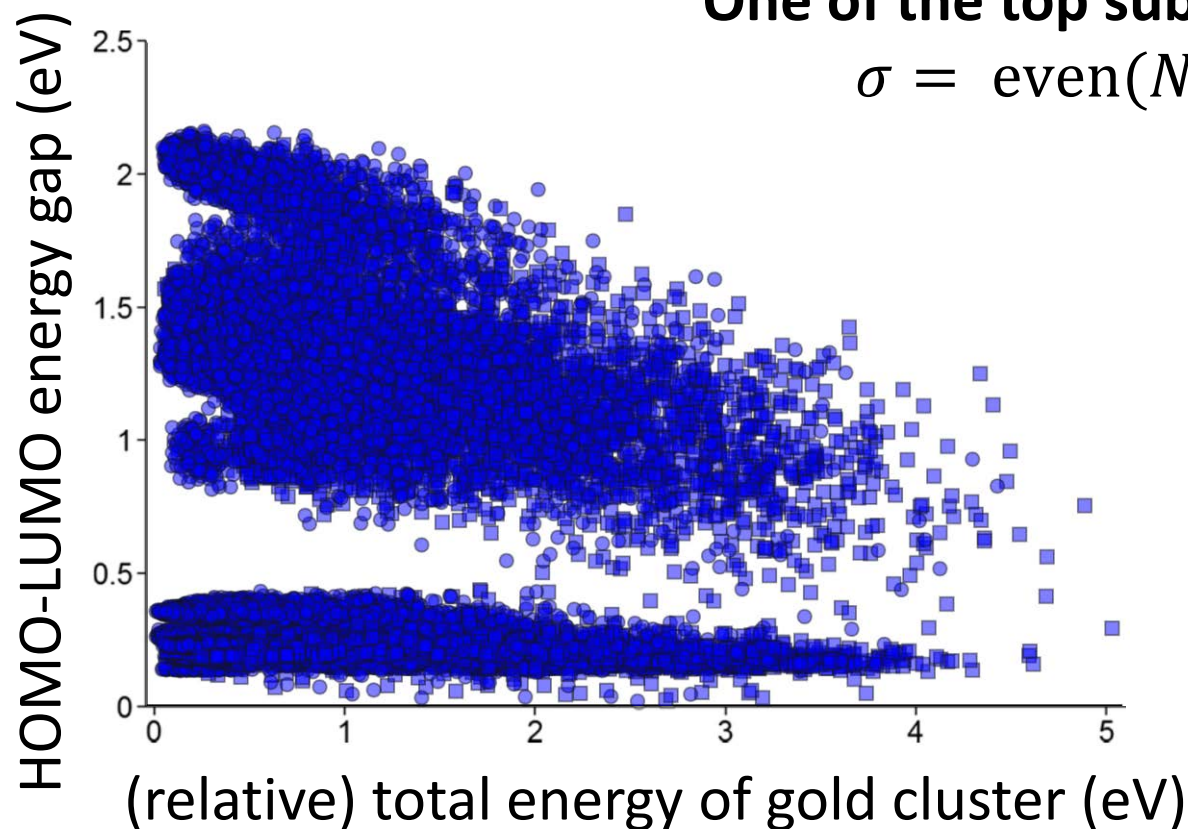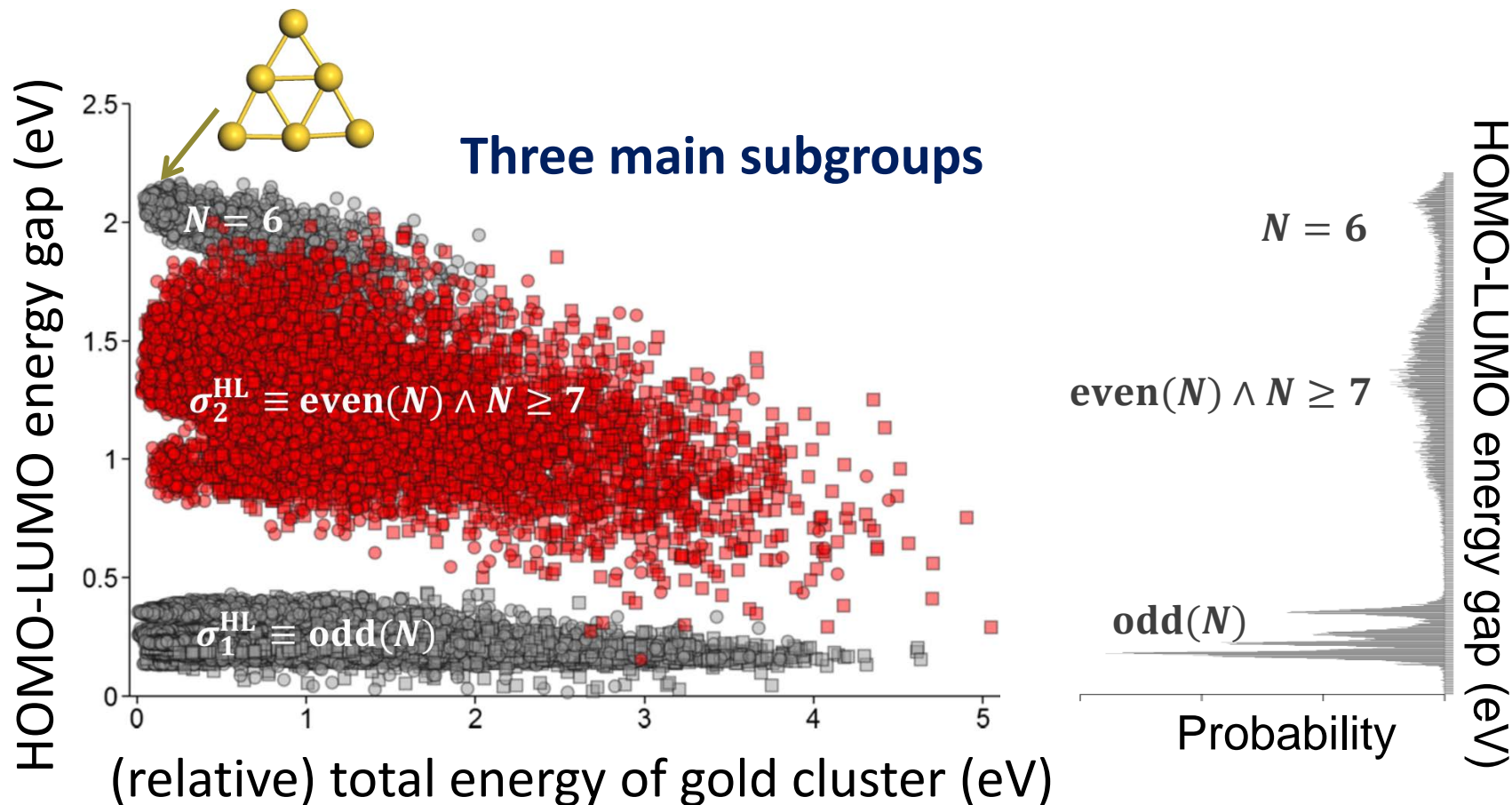25,000 gold cluster configurations (sizes 5-14) in the gas phase

*Choose* **target property**
HOMO-LUMO energy gap

*Choose* **variation reduction utility function**
$$u(Y') = (\mathrm{std}(Y) - \mathrm{std}(Y'))/\mathrm{std}(Y)$$

**One of the top subgroup selectors found**
$$\sigma = \mathrm{even}(N) \wedge N \geq 7$$



HOMO-LUMO energy gap (eV)

(relative) total energy of gold cluster (eV)

# Rediscover simple insight about HOMO-LUMO gap

25,000 gold cluster configurations (sizes 5-14) in the gas phase

***Choose*** **target property**
HOMO-LUMO energy gap

***Choose*** **variation reduction utility function**
$$u(Y') = (\text{std}(Y) - \text{std}(Y'))/\text{std}(Y)$$

**Three main subgroups**



HOMO-LUMO energy gap (eV)

$N = 6$

$\sigma_2^{\text{HL}} \equiv \textbf{even}(N) \wedge N \geq 7$

$\sigma_1^{\text{HL}} \equiv \textbf{odd}(N)$

(relative) total energy of gold cluster (eV)

$N = 6$

$\textbf{even}(N) \wedge N \geq 7$

$\textbf{odd}(N)$

Probability

HOMO-LUMO energy gap (eV)

# Equilibrium state has maximum electronic hardness at 0 K

**Electronic hardness = resistance to electron density deformation**

$$\frac{\partial \mu}{\partial N}\bigg|_v \approx \frac{1}{2}(E_{\text{LUMO}} - E_{\text{HOMO}})$$

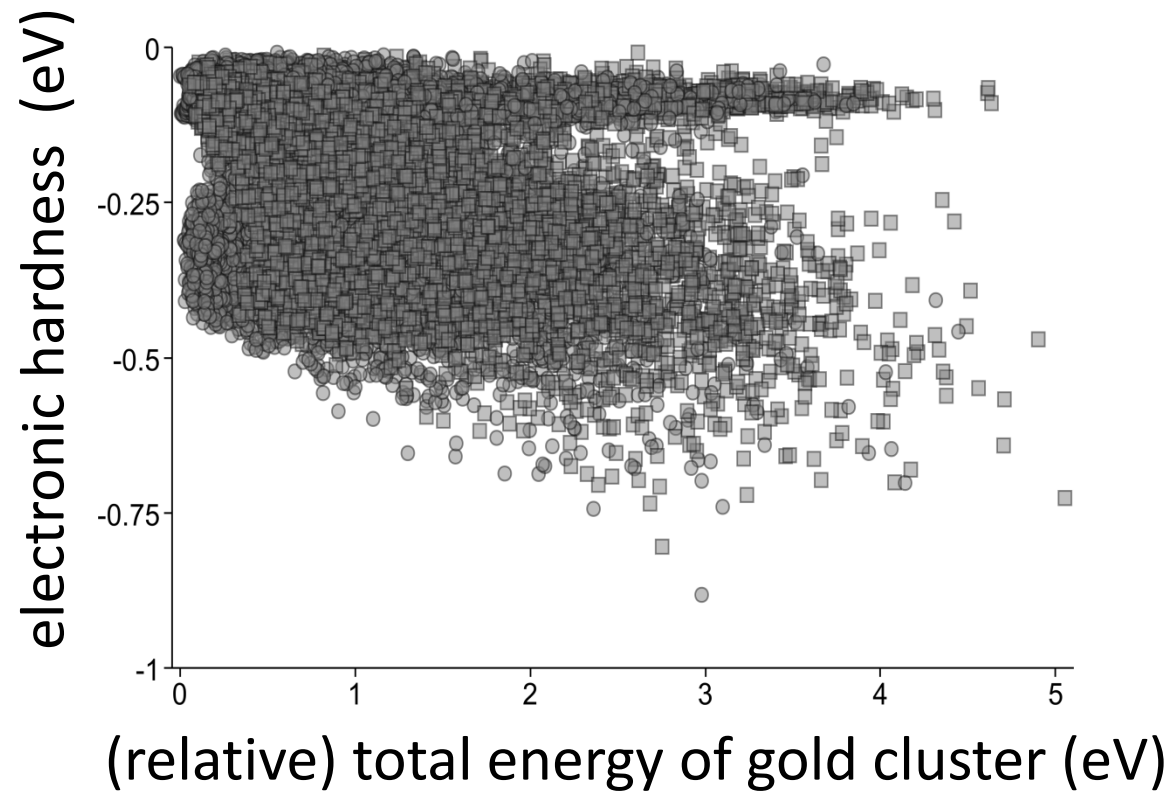Can electronic hardness be a descriptor for cluster isomer stability?

*Choose target property*
electronic hardness and formation energy

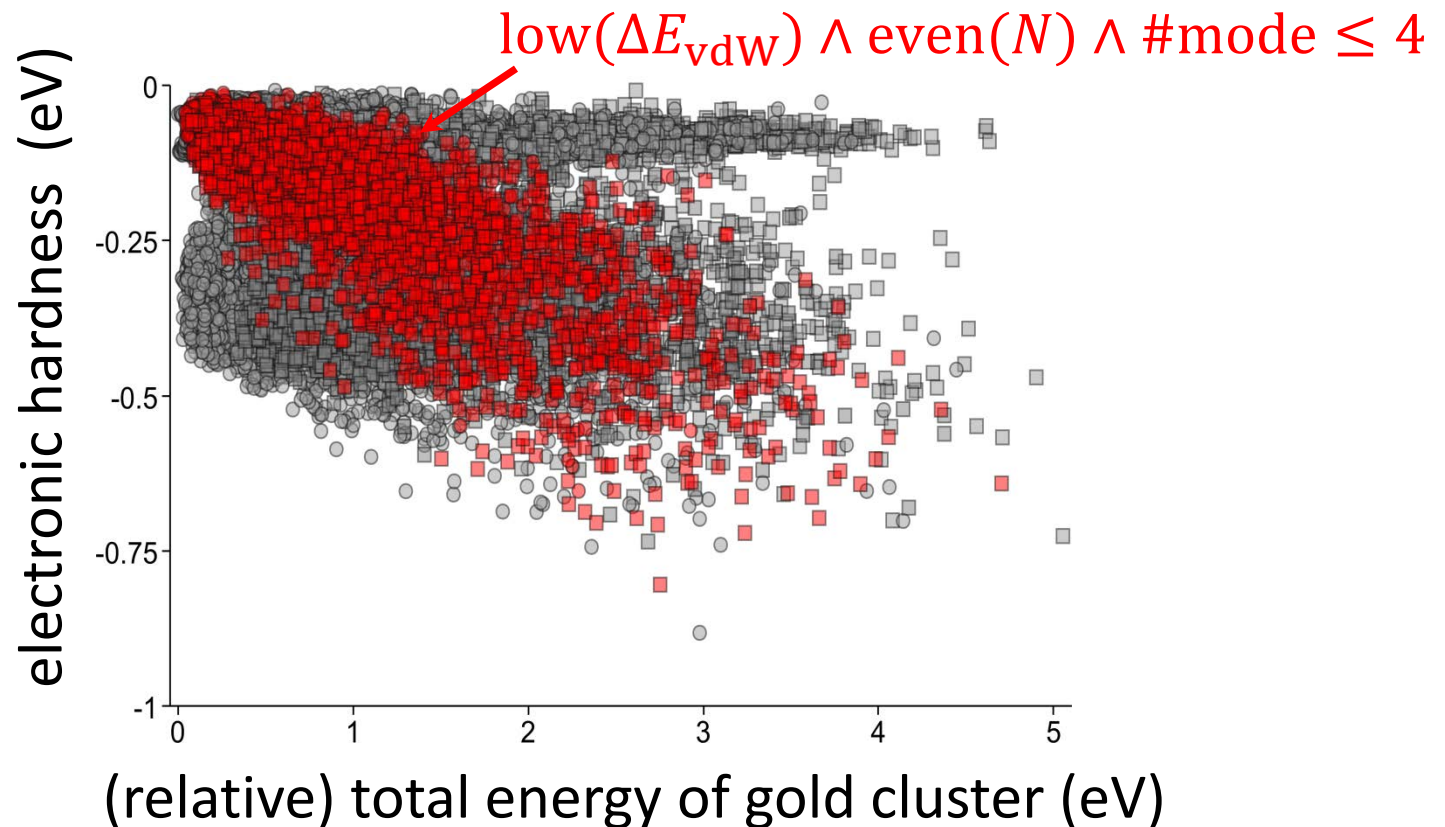**Find a linear model:** $u(Y'_1, Y'_2)$

**Top subgroup selector found**

$$\text{low}(\Delta E_{\text{vdW}}) \wedge \text{even}(N) \wedge \#\text{mode} \leq 4$$

# Electronic hardness can be a qualitative descriptor for stability beyond the ground state (among the subgroup)
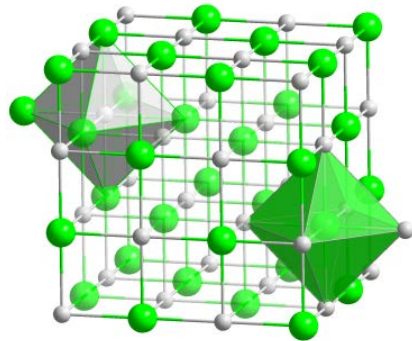
Subgroup population = 30% of even sized clusters



$\text{low}(\Delta E_{\text{vdW}}) \wedge \text{even}(N) \wedge \#\text{mode} \leq 4$

electronic hardness (eV)

(relative) total energy of gold cluster (eV)

# Can subgroup discovery find local models for the 82 octet binary materials that describe ZB and RS?

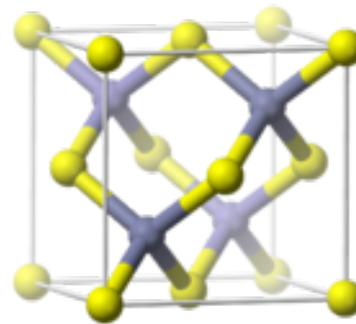$Target = \text{sign}(E_{rocksalt} - E_{zincblende})$

**Subgroup 1**

**Subgroup 2**

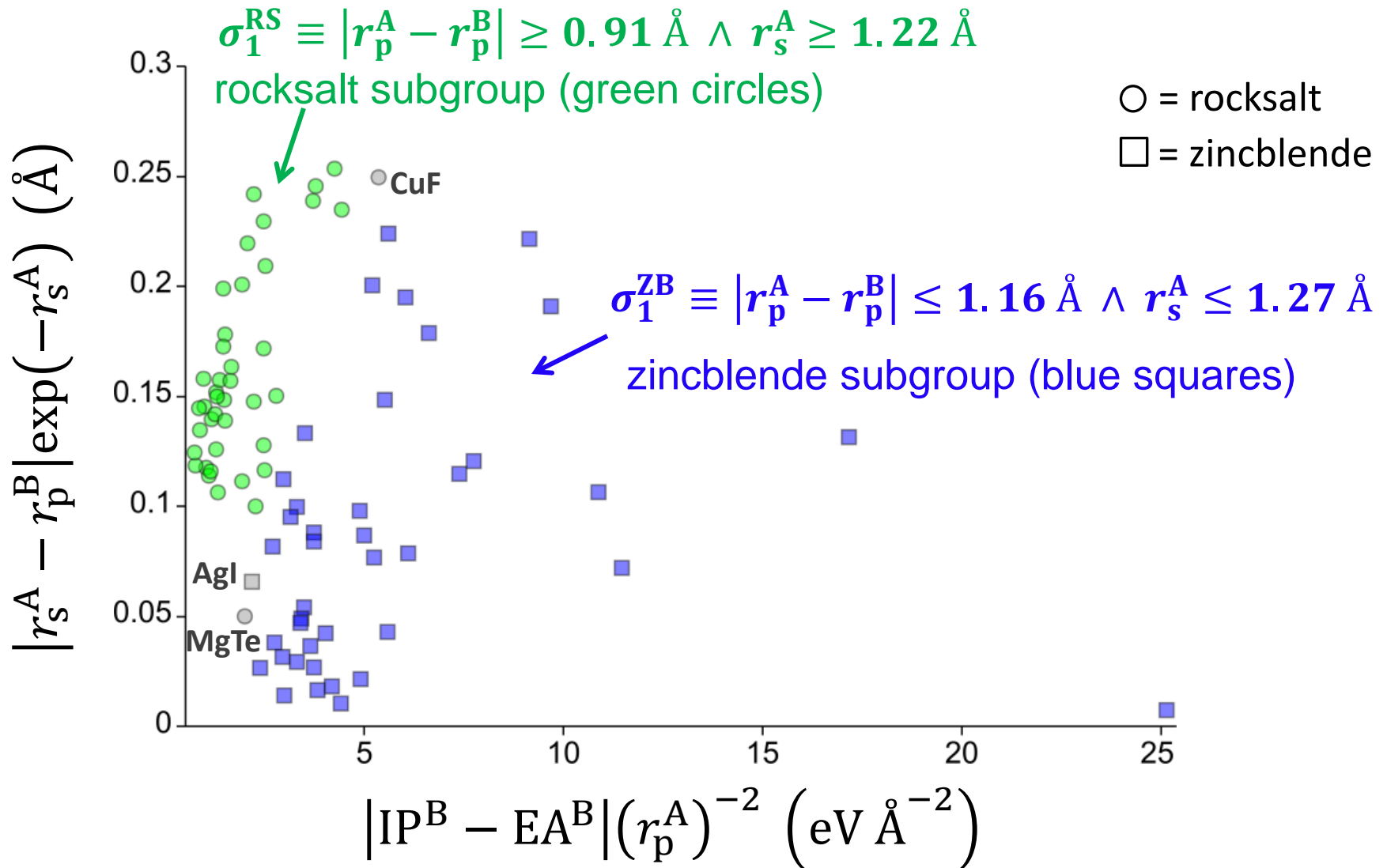Rocksalt (RS)            Zincblende (ZB)



vs.

# Subgroup discovery classifies 79 of 82 octet binary semiconductors using a 2D descriptor
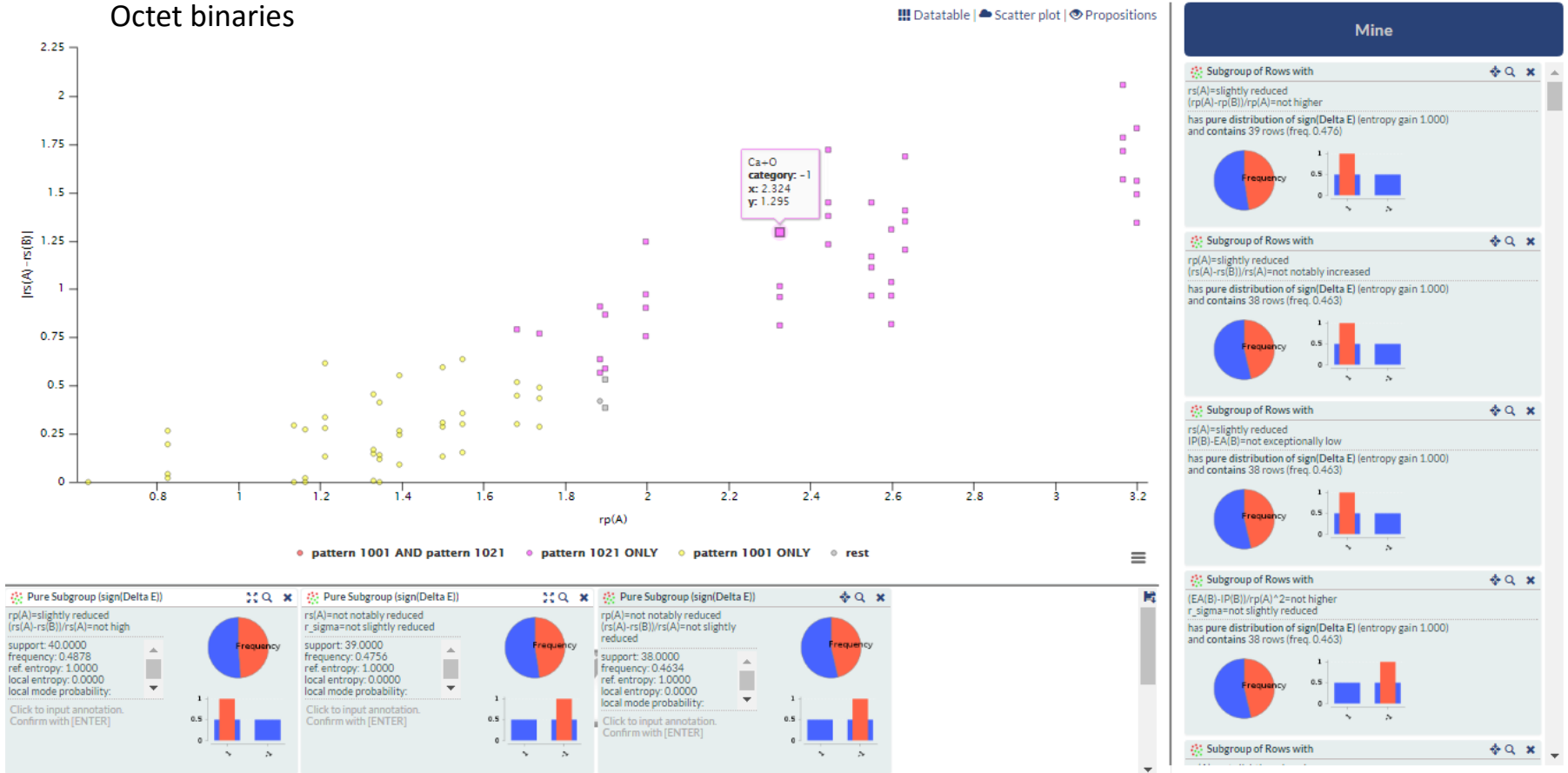


$$\sigma_1^{RS} \equiv \left| r_p^A - r_p^B \right| \geq 0.91 \,\text{Å} \,\wedge\, r_s^A \geq 1.22 \,\text{Å}$$

rocksalt subgroup (green circles)

$$\sigma_1^{ZB} \equiv \left| r_p^A - r_p^B \right| \leq 1.16 \,\text{Å} \,\wedge\, r_s^A \leq 1.27 \,\text{Å}$$

zincblende subgroup (blue squares)

○ = rocksalt
□ = zincblende

CuF

AgI

MgTe

y-axis: $\left| r_s^A - r_p^B \right| \exp\left(-r_s^A\right)$ (Å)

x-axis: $\left| IP^B - EA^B \right| \left(r_p^A\right)^{-2}$ $\left(\text{eV Å}^{-2}\right)$

# We have made subgroup discovery, among other data analytic tools, available online as interactive tutorials



Subgroup discovery (SGD) was implemented by Mario Boley using Creedo and realKD library, www.realkd.org

# Big-data analytics tools for materials science are being developed by the NOMAD team

Compressed sensing (LASSO, $l_0$, LASSO+$l_0$)

Sure independent screening+$l_0$

Subgroup discovery

The question of causal inference is still open (advice is welcome!)

Bigger data and harder problems

Crystal structure maps

Transparent conducting oxides

…. and much more!

Subgroup discovery: B. Goldsmith, M. Boley, J. Vreeken, L. Ghiringhelli, M. Scheffler, *unpublished*

SIS+L0: R. Ouyang, E. Ahmetcik,  L. Ghiringhelli, M. Scheffler, *unpublished*

LASSO+L0: L. Ghiringhelli, J. Vybiral, S. V. Levchenko, C. Draxl, M. Scheffler, *PRL 114*, 105503 (2015)

# Acknowledgements

## The entire NOMAD consortium

*Special thanks:*

| | |
|---|---|
| Matthias Scheffler | Runhai Ouyang |
| Luca Ghiringhelli | Christopher Sutton |
| Mario Boley | Matthias Rupp |
| Jilles Vreeken | Claudia Draxl |

Fritz Haber Institute of the Max Planck Society
Theory Department



**Alexander von Humboldt**
Stiftung / Foundation

**mpii** max planck institut informatik