

Learning Materials Science from Text

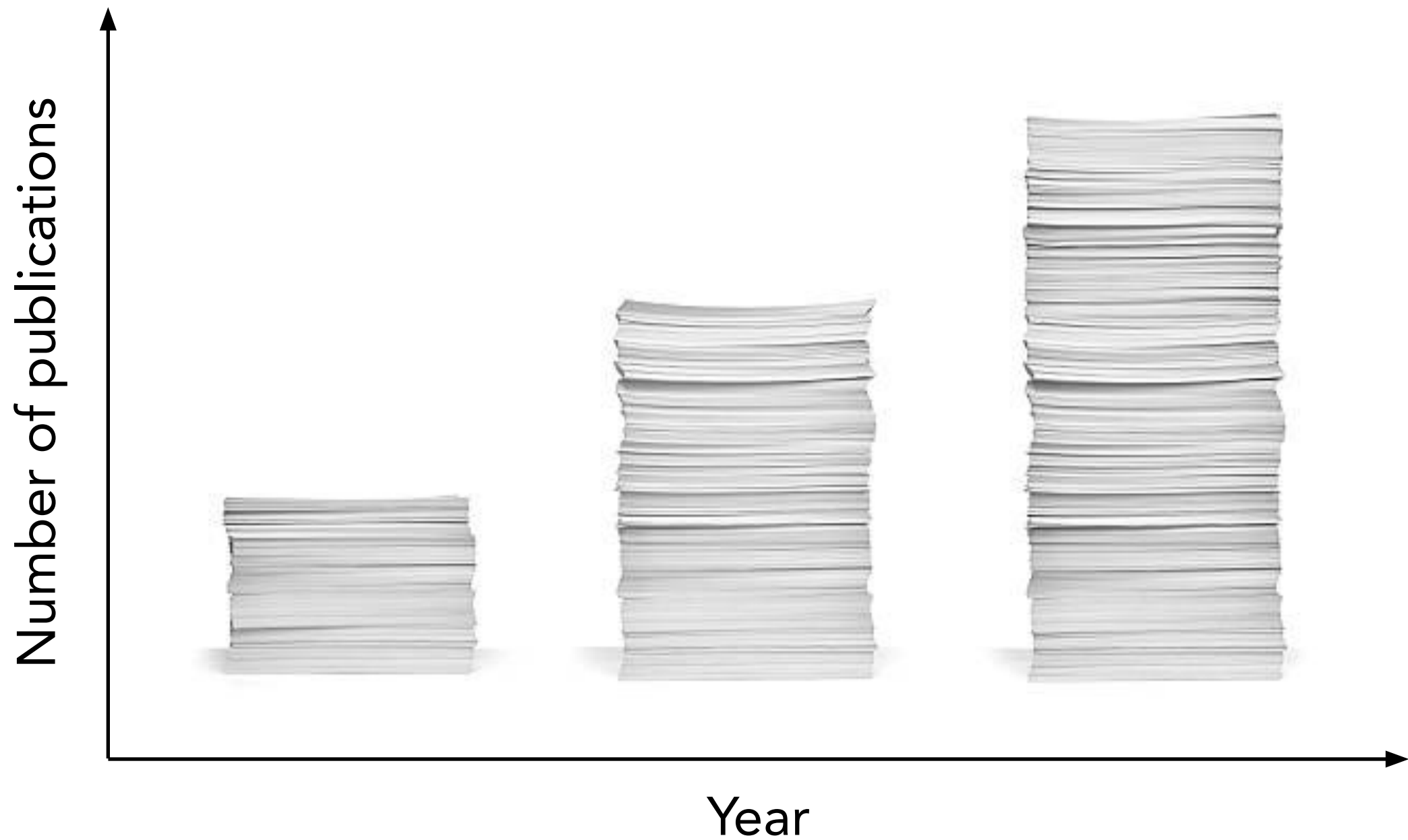
Vahe Tshitoyan

Machine Learning @ Google ATAP*

December 14, 2019

* The presented work was carried out while I was a Postdoc @ Berkeley Lab

Science is in text



Problem



Solution



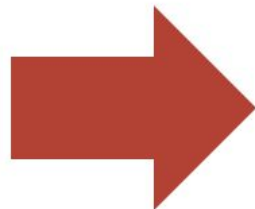
Computers speak in numbers

Word vectors (embeddings)

Map each word to a vector in space (sequence of numbers)

Vocabulary:

Man, woman, boy,
girl, prince,
princess, queen,
king, monarch



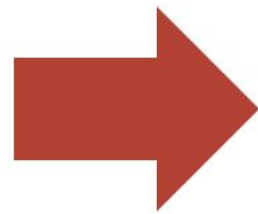
	Femininity	Youth	Royalty
Man	0	0	0
Woman	1	0	0
Boy	0	1	0
Girl	1	1	0
Prince	0	1	1
Princess	1	1	1
Queen	1	0	1
King	0	0	1
Monarch	0.5	0.5	1

Each word gets a
1x3 vector

Similar words...
similar vectors

Vector algebra with embeddings

Vocabulary:
Man, woman, boy,
girl, prince,
princess, queen,
king, monarch



	Femininity	Youth	Royalty
Man	0	0	0
Woman	1	0	0
Boy	0	1	0
Girl	1	1	0
Prince	0	1	1
Princess	1	1	1
Queen	1	0	1
King	0	0	1
Monarch	0.5	0.5	1

Each word gets a
1x3 vector

Similar words...
similar vectors

- Words with similar meanings have similar vectors

Also...

- $\text{queen} - \text{king} = \text{woman} - \text{man} = \text{girl} - \text{boy}$
- $(\text{prince} + \text{princess} + \text{queen} + \text{king})/4 = \text{monarch}$

How do we "machine learn" word
vectors?

Word2Vec: skip-gram

Sentence: Electrochemical properties of LiCoO₂ thin film.

Task: train a neural network to predict the context (target) word given the centre word

Training example 1:

of →? Electrochemical

of →? properties

of →? LiCoO₂

of →? thin

Word2Vec: skip-gram

Sentence: Electrochemical properties of LiCoO₂ thin film.

Task: train a neural network to predict the context (target) word given the centre word

Training example 2:

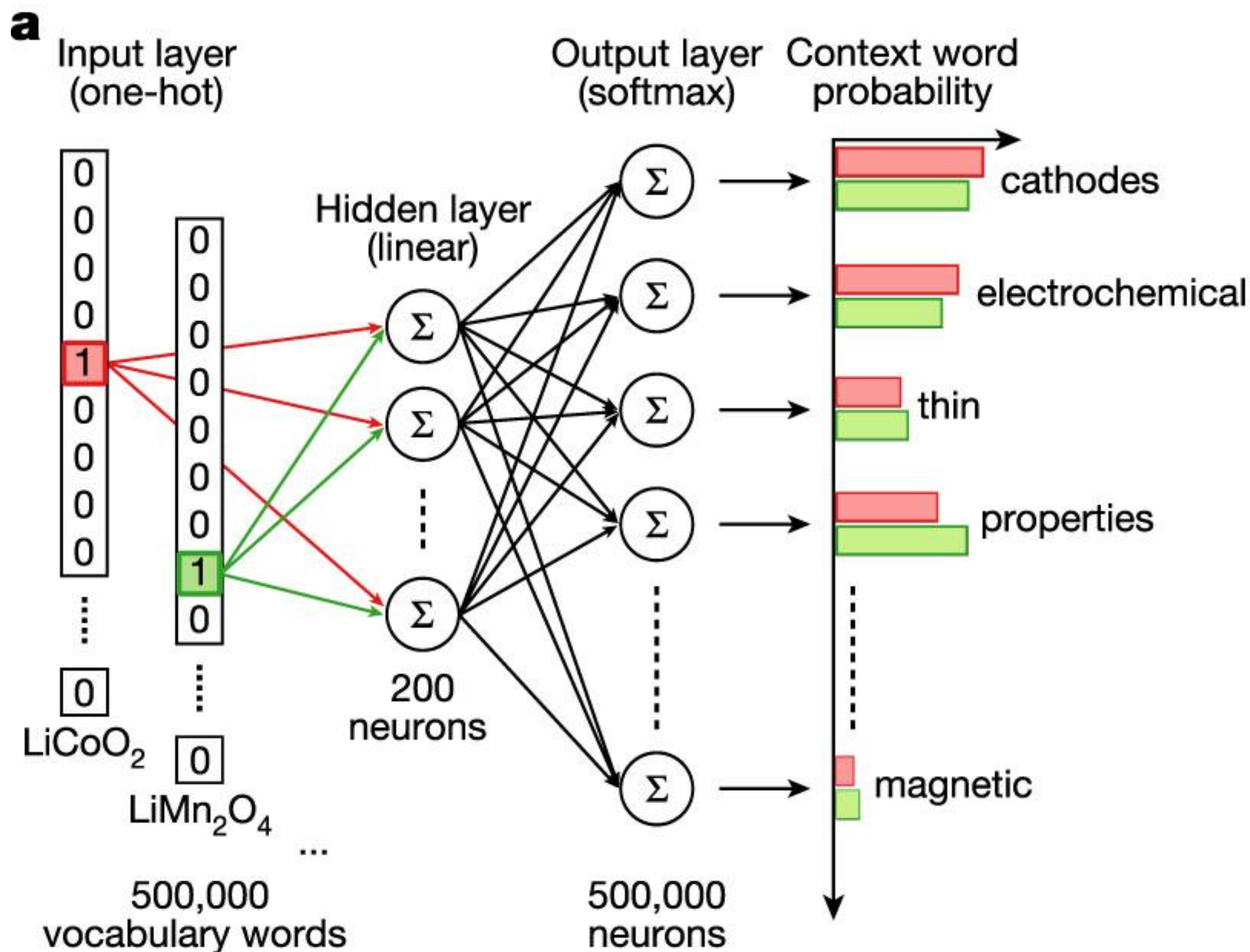
LiCoO₂ →? properties

LiCoO₂ →? of

LiCoO₂ →? thin

LiCoO₂ →? film

The Neural Network



Word2Vec: skip-gram

Sentence: Electrochemical properties of LiCoO₂ thin film.

Task: train a neural network to predict the context (target) word given the centre word

Training example 2:

LiCoO₂ →? properties

LiCoO₂ →? of

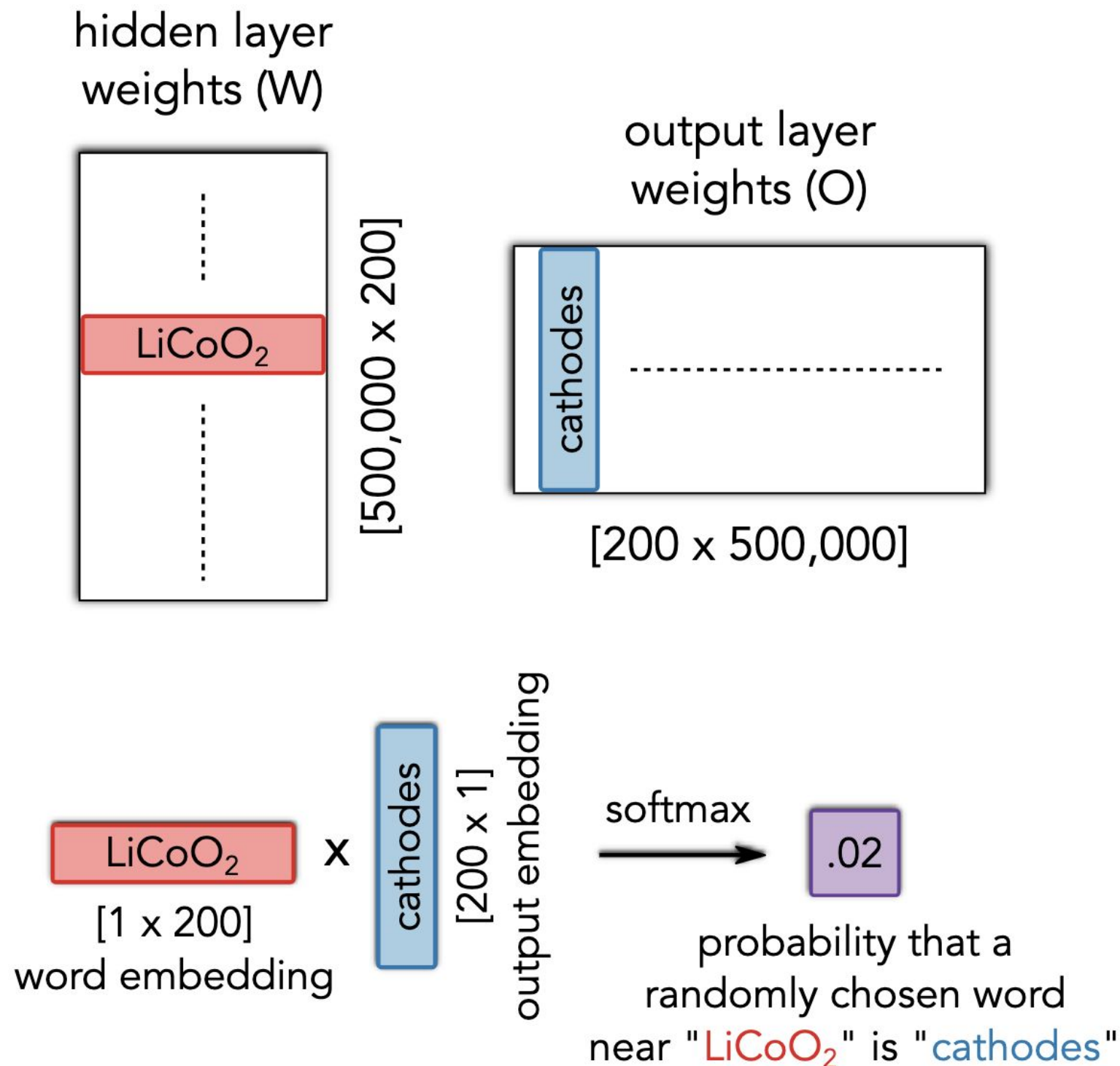
LiCoO₂ →? thin

LiCoO₂ →? film

Objective: learn vector representations for words

If two different words have similar contexts, then the neural network needs to output similar results for these two words. The network will output similar context predictions for these two words if the word vectors are similar.

Outcomes of the training



Training Data

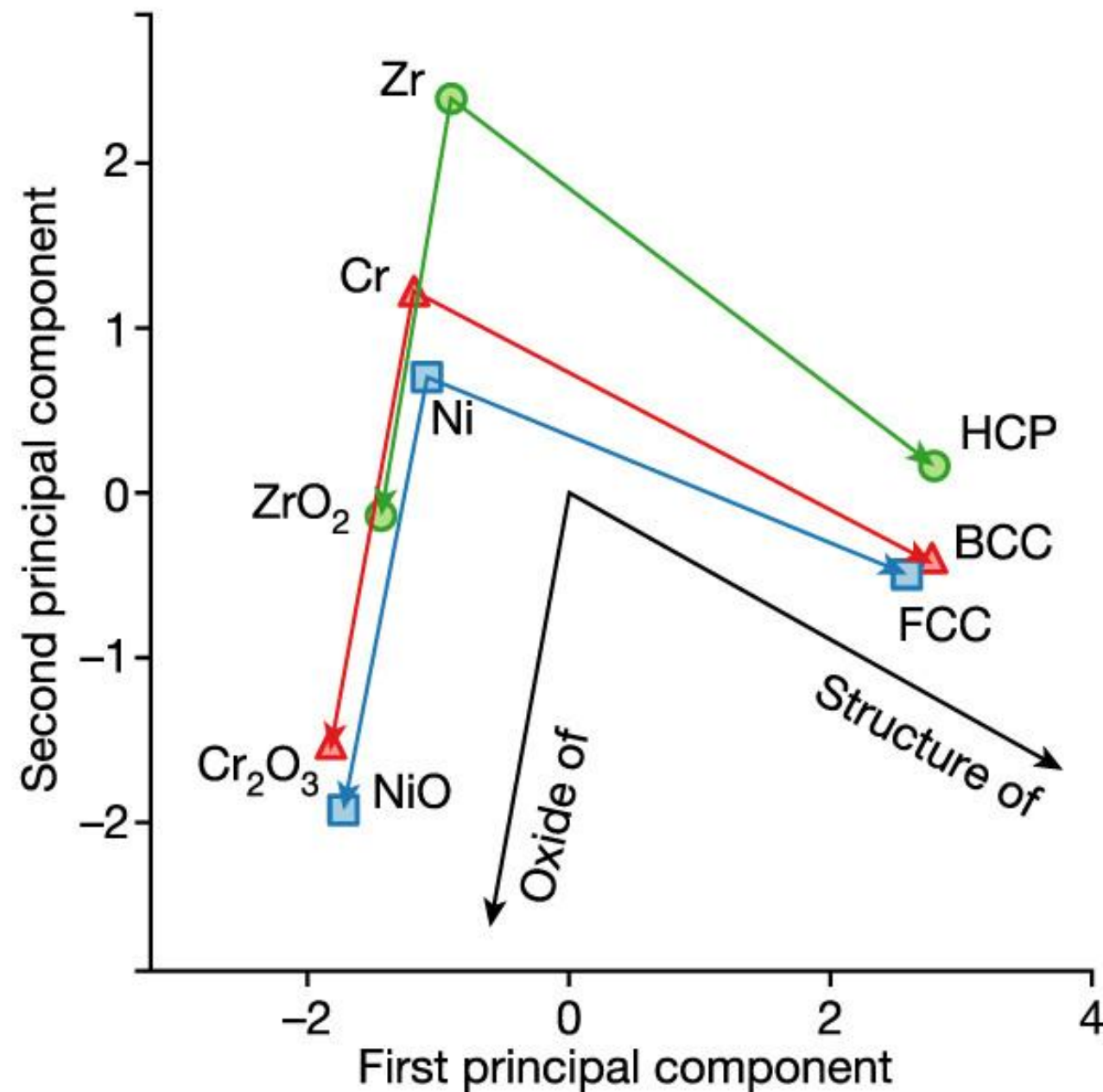
- A total of 3.3 million abstracts downloaded
- 1.5 million were classified as inorganic materials science
 - ~300 million words across the abstracts
 - ~500k unique vocabulary words and phrases
 - ~100k unique material formulae

Word similarity

Words most similar to ____ are?

- **LiCoO₂**: LiMn₂O₄, LiNi_{0.5}Mn_{1.5}O₄, LiNi_{0.8}Co_{0.2}O₂,
LiNi_{0.8}Co_{0.15}Al_{0.05}O₂
- **ferromagnetic**: ferrimagnetic, antiferromagnetic,
anti-ferromagnetic, paramagnetic

Analogies



CoO - Co + Al = Al₂O₃

Li - lithium + helium = He

temperature - K + Pa = pressure

LiCoO₂ - cathode + anode = graphite

A famous vector representation in chemistry?

The periodic table of elements

(0, 0)

2D representation : H - (1, 1)
Li - (1, 2)
...

H

He

Li

Be

B

C

N

O

F

Ne

Na

Mg

Al

Si

P

S

Cl

Ar

K

Ca

Sc

Ti

V

Cr

Mn

Fe

Co

Ni

Cu

Zn

Ga

Ge

As

Se

Br

Kr

Rb

Sr

Y

Zr

Nb

Mo

Tc

Ru

Rh

Pd

Ag

Cd

In

Sn

Sb

Te

I

Xe

Cs

Ba

Hf

Ta

W

Re

Os

Ir

Pt

Au

Hg

Tl

Pb

Bi

Po

At

Rn

Fr

Ra

Rf

Db

Sg

Bh

Hs

Mt

Ds

Rg

Cn

Nh

Fl

Mc

Lv

Ts

Og

2D representation :

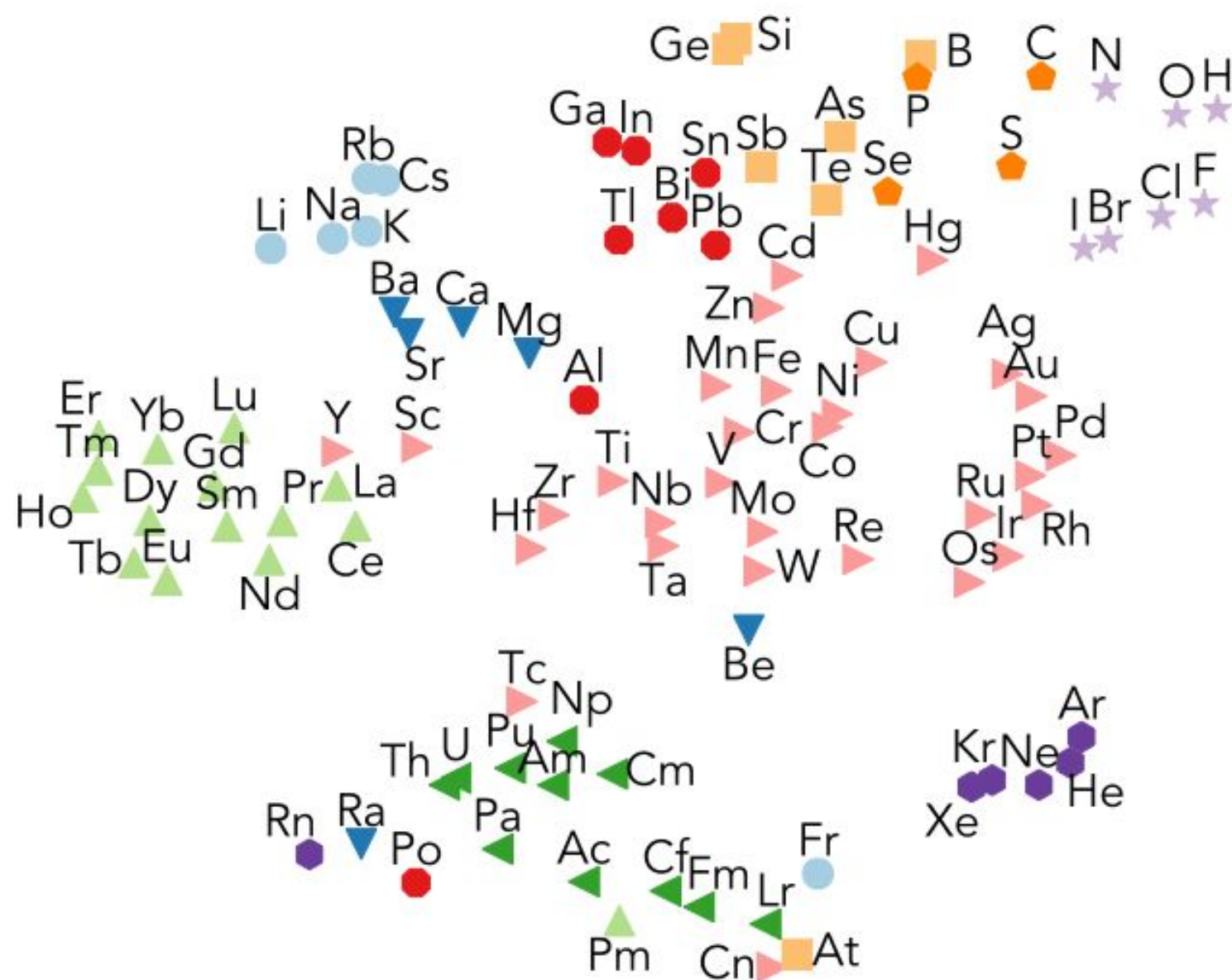
H - (1, 1)

Li - (1, 2)

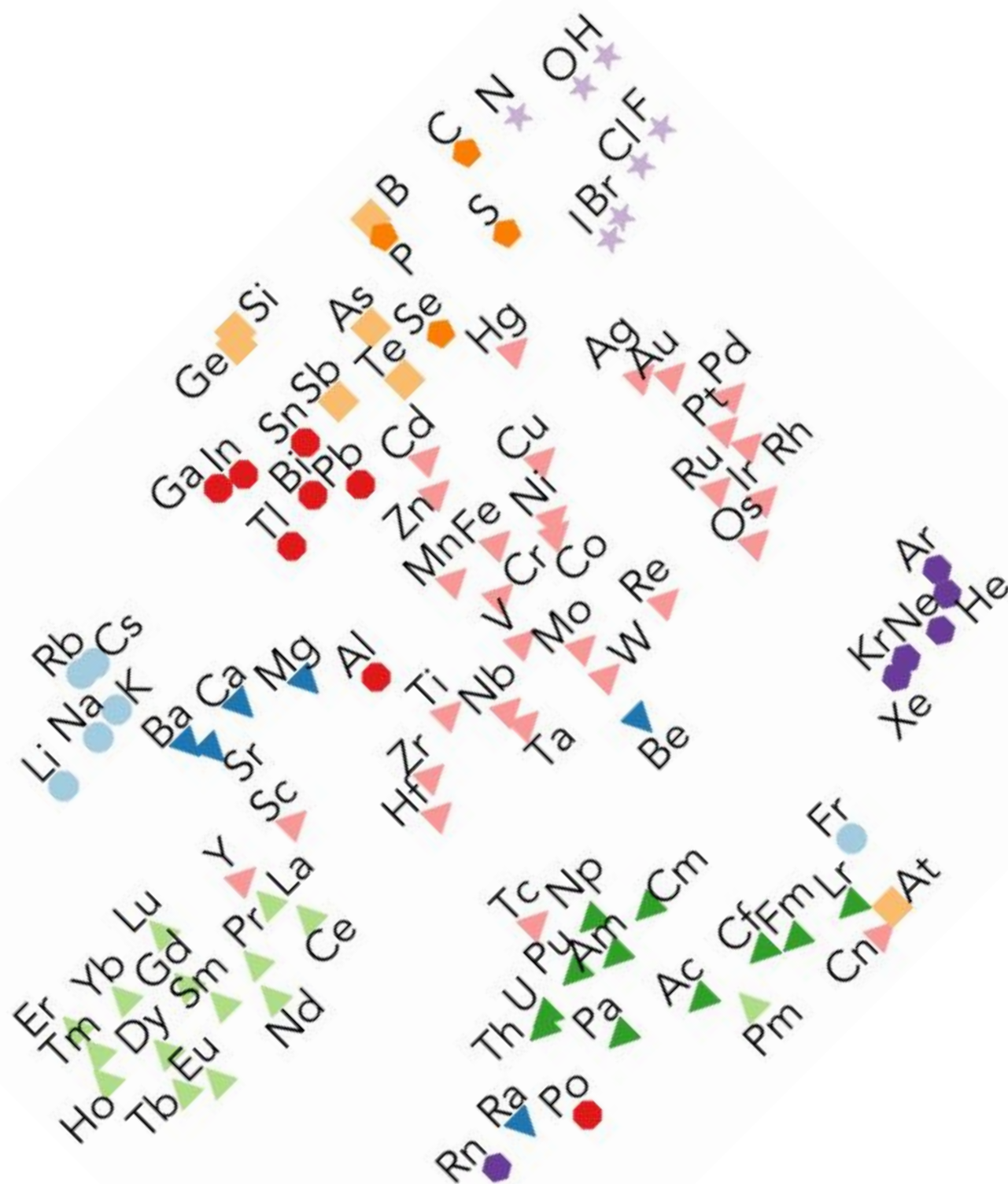
...

La	Ce	Pr	Nd	Pm	Sm	Eu	Gd	Tb	Dy	Ho	Er	Tm	Yb	Lu
Ac	Th	Pa	U	Np	Pu	Am	Cm	Bk	Cf	Es	Fm	Md	No	Lr

Word vectors of elements in 2D

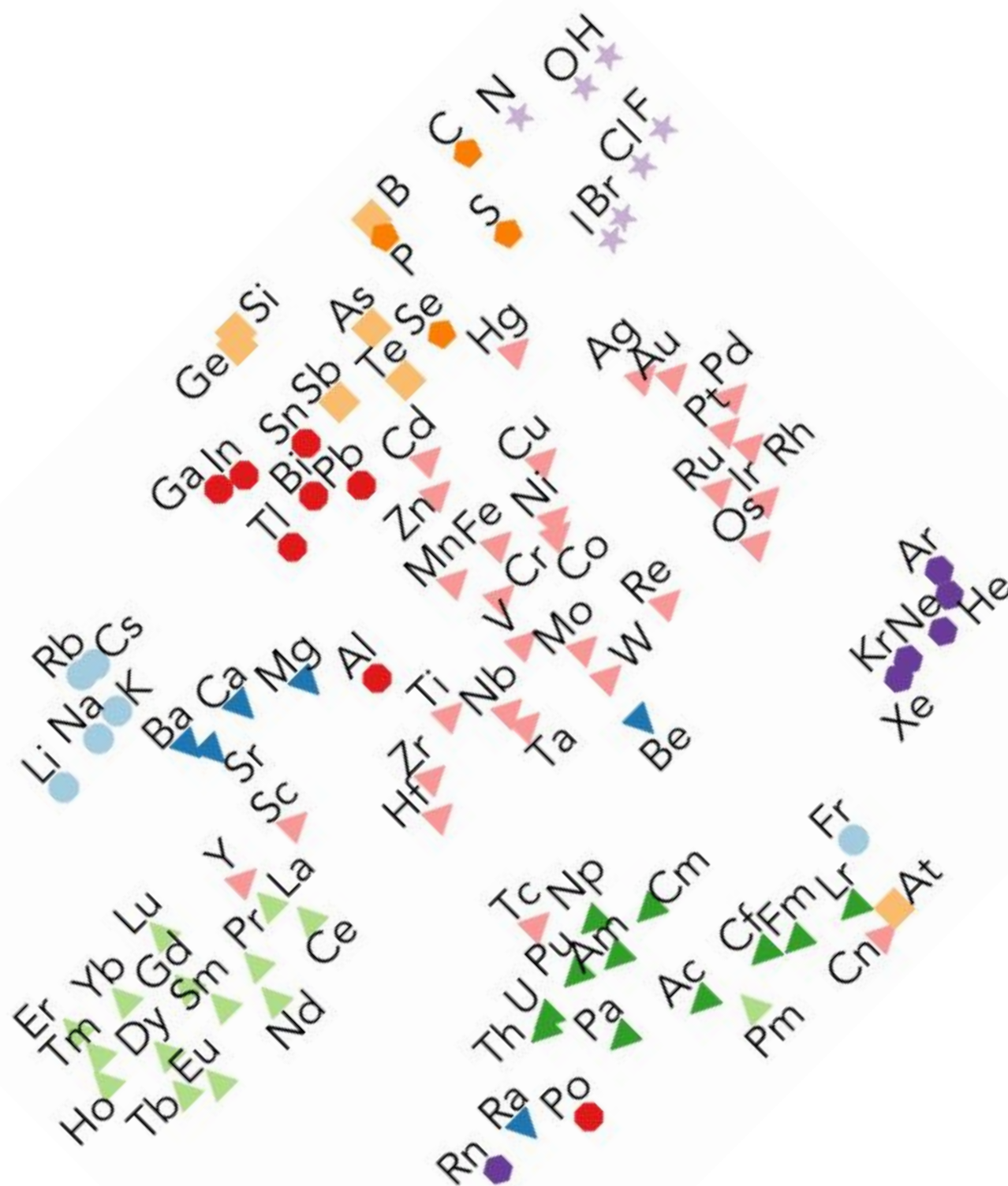


Word vectors of elements in 2D



- alkali metal
- alkaline earth metal
- lanthanide
- actinide
- transition metal
- post-transition metal
- metalloid
- polyatomic nonmetal
- diatomic nonmetal
- noble gas

Word vectors of elements in 2D



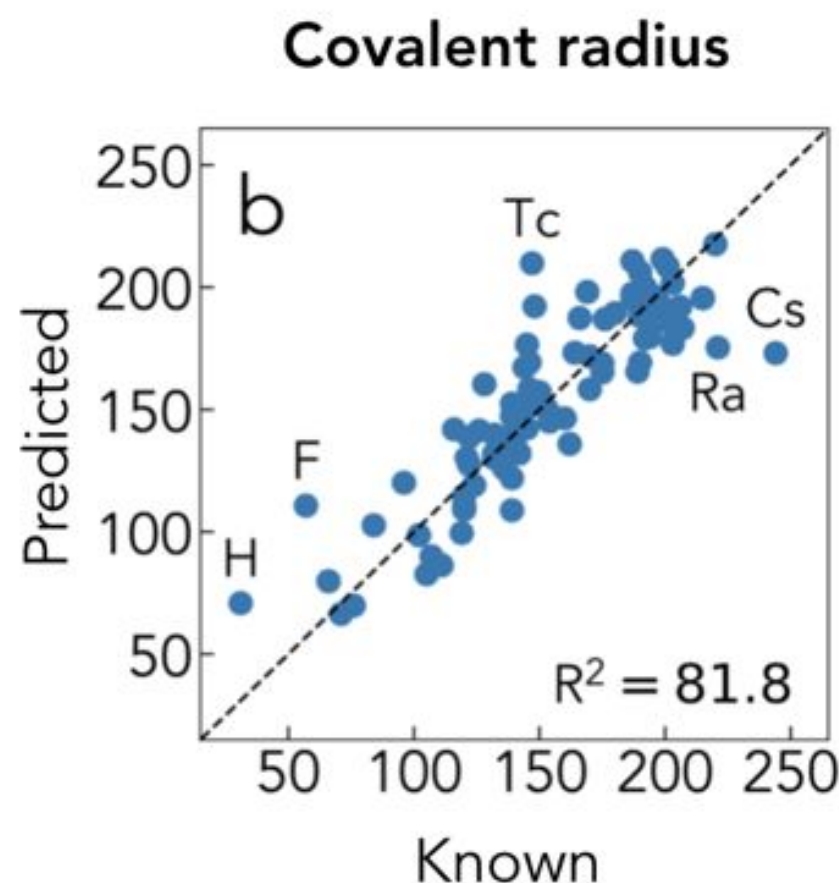
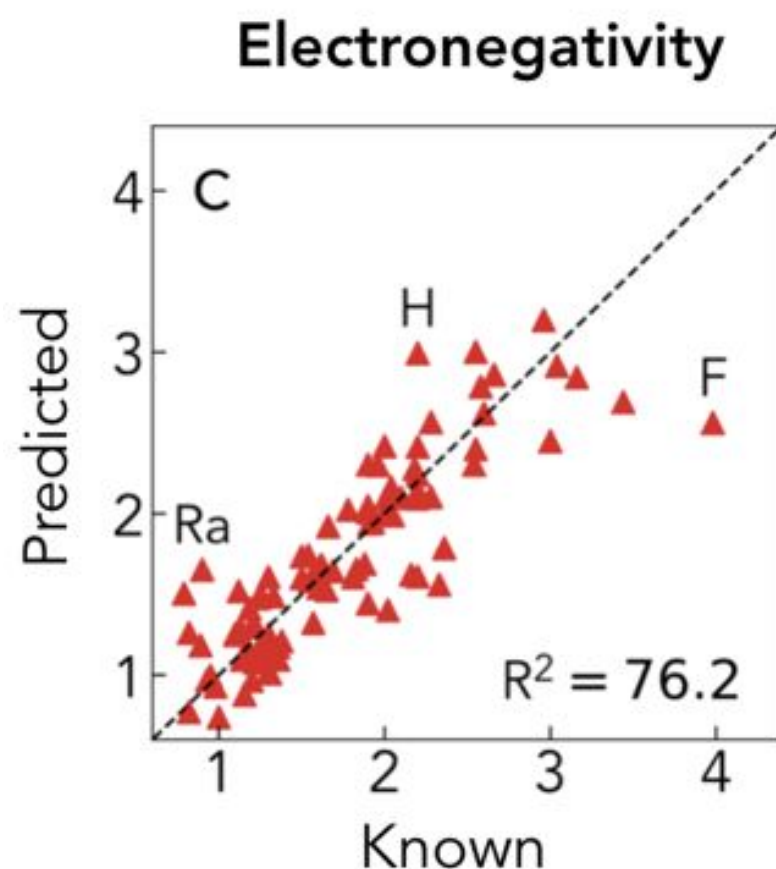
- alkali metal
- ▼ alkaline earth metal
- ▲ lanthanide
- ▼ actinide
- ▶ transition metal

H																	He
Li	Be											B	C	N	O	F	Ne
Na	Mg											Al	Si	P	S	Cl	Ar
K	Ca	Sc	Ti	V	Cr	Mn	Fe	Co	Ni	Cu	Zn	Ga	Ge	As	Se	Br	Kr
Rb	Sr	Y	Zr	Nb	Mo	Tc	Ru	Rh	Pd	Ag	Cd	In	Sn	Sb	Te	I	Xe
Cs	Ba		Hf	Ta	W	Re	Os	Ir	Pt	Au	Hg	Tl	Pb	Bi	Po	At	Rn
Fr	Ra		Rf	Db	Sg	Bh	Hs	Mt	Ds	Rg	Cn	Nh	Fl	Mc	Lv	Ts	Og

La	Ce	Pr	Nd	Pm	Sm	Eu	Gd	Tb	Dy	Ho	Er	Tm	Yb	Lu
Ac	Th	Pa	U	Np	Pu	Am	Cm	Bk	Cf	Es	Fm	Md	No	Lr

- post-transition metal
- metalloid
- ◆ polyatomic nonmetal
- ★ diatomic nonmetal
- noble gas

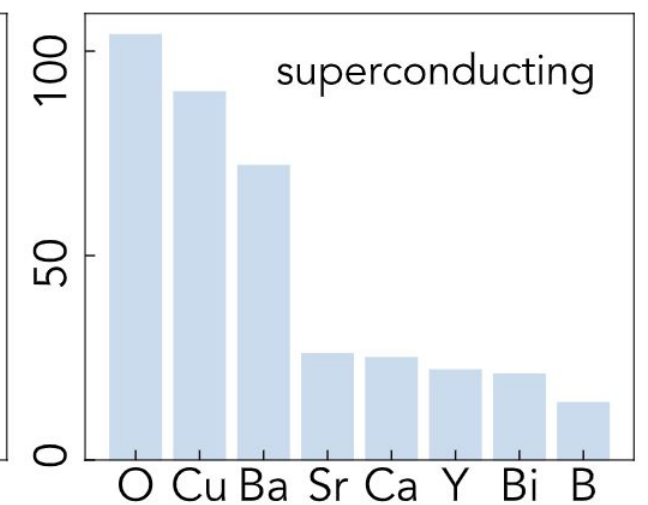
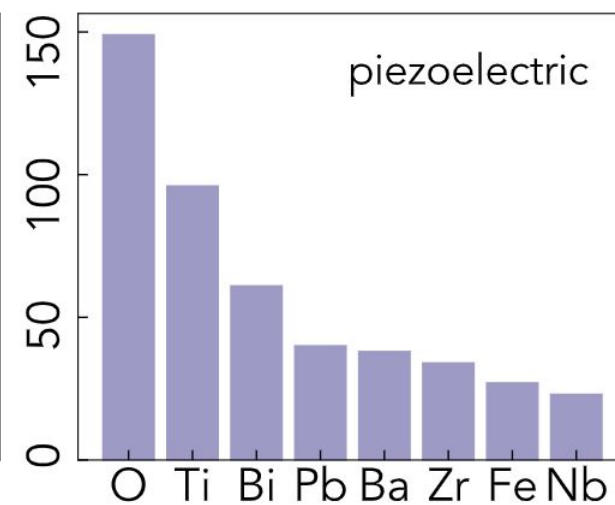
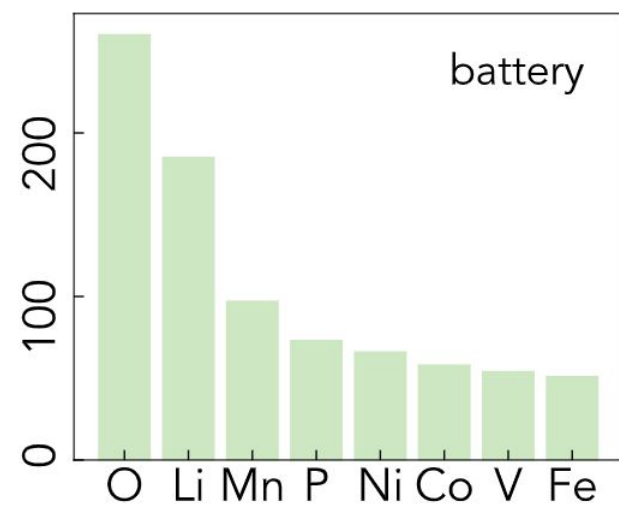
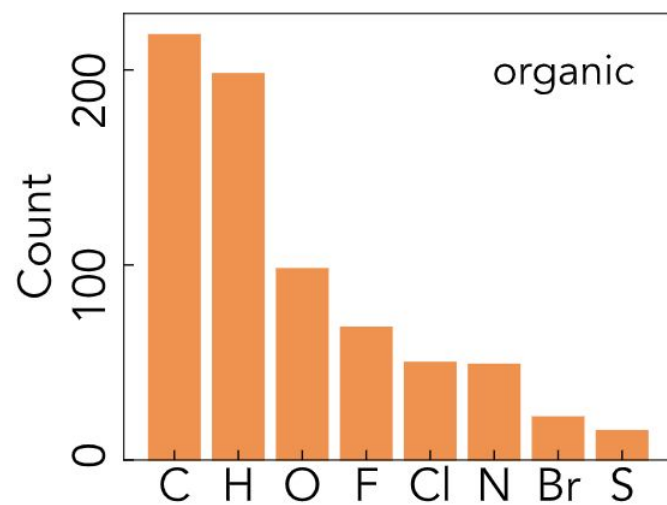
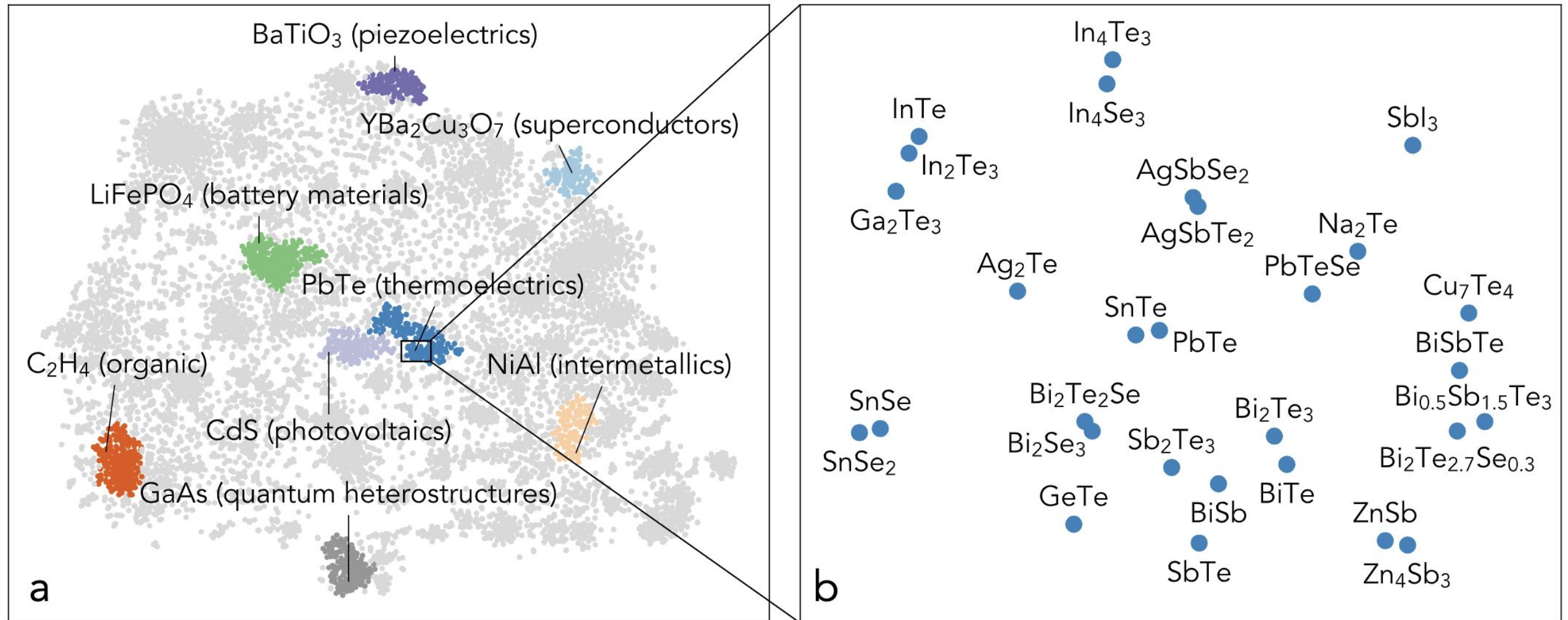
Directions in latent latent space correspond to elemental properties



Property	R^2
Mendelev number	74.4
Atomic weight	72.4
Melting T	73.9
Covalent radius	81.8
Electronegativity	76.2
Row	80.7
Column	69.2

Basic linear regression models fit using only the unsupervised, text-derived embeddings as feature vectors.

Similar materials cluster together



...what is the most similar material to
the word "thermoelectric"?

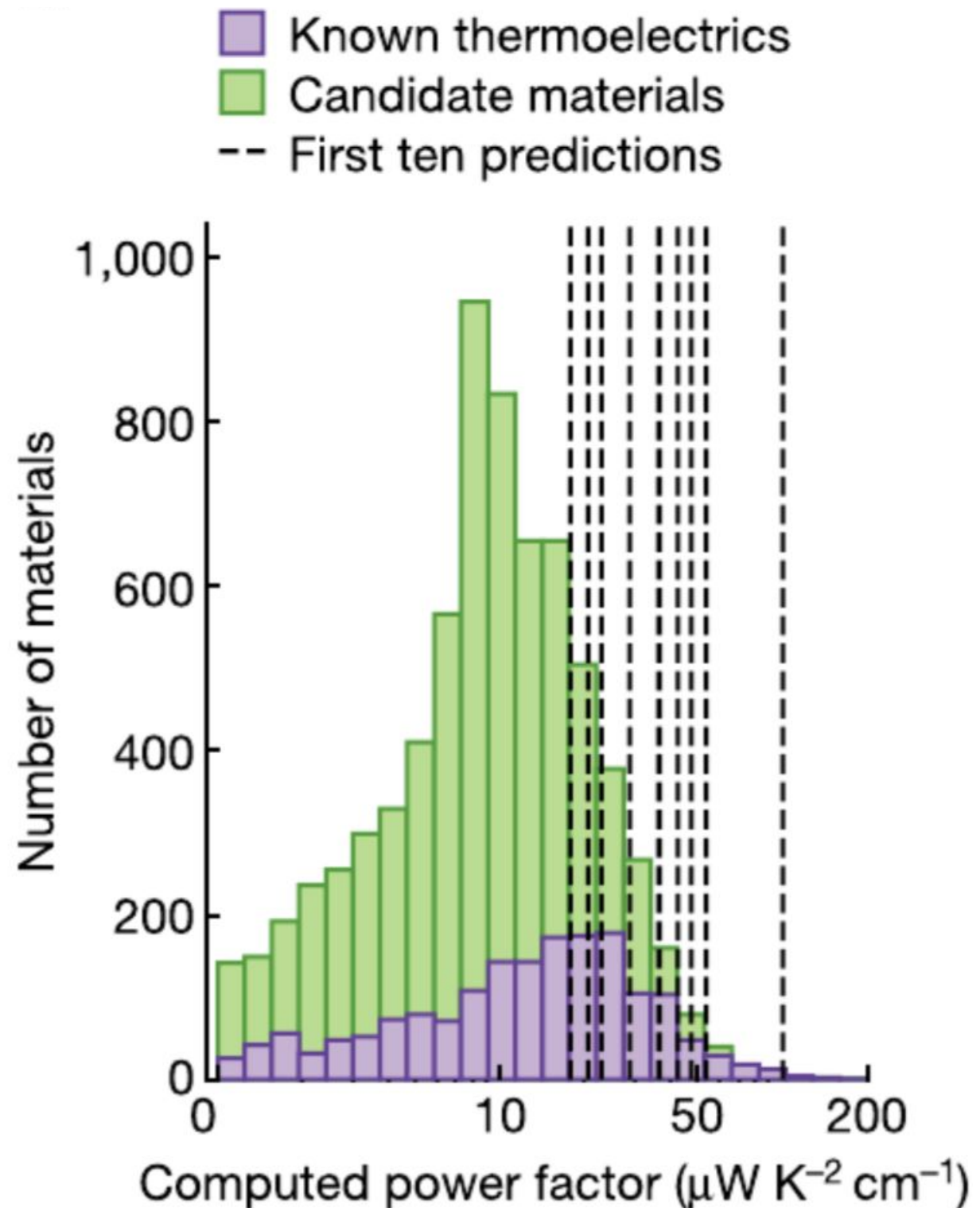
Predictions of new Thermoelectrics

Cosine similarity to
'thermoelectric'

1.	Bi_2Te_3	✓
2.	MgAgSb	✓
3.	PbTe	✓
...		✓
326.	Li_2CuSb	?
...		✓
328.	In_4Te_3	✓
...		✓
345.	$\text{Cu}_3\text{Nb}_2\text{O}_8$?
...		✓

✓ Known thermoelectrics

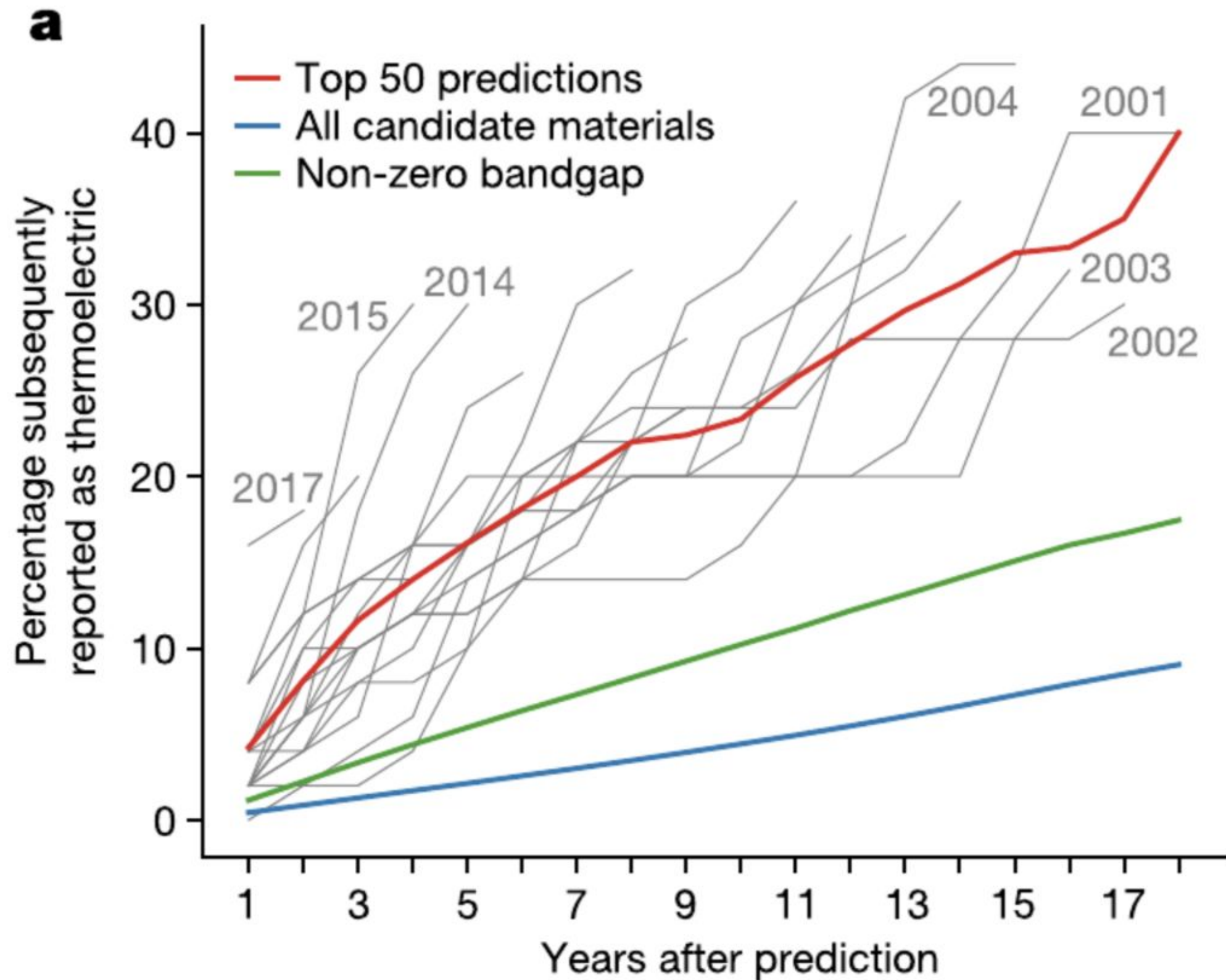
? Predictions



We can examine the plausibility of these predictions by performing experiments “in the past”.

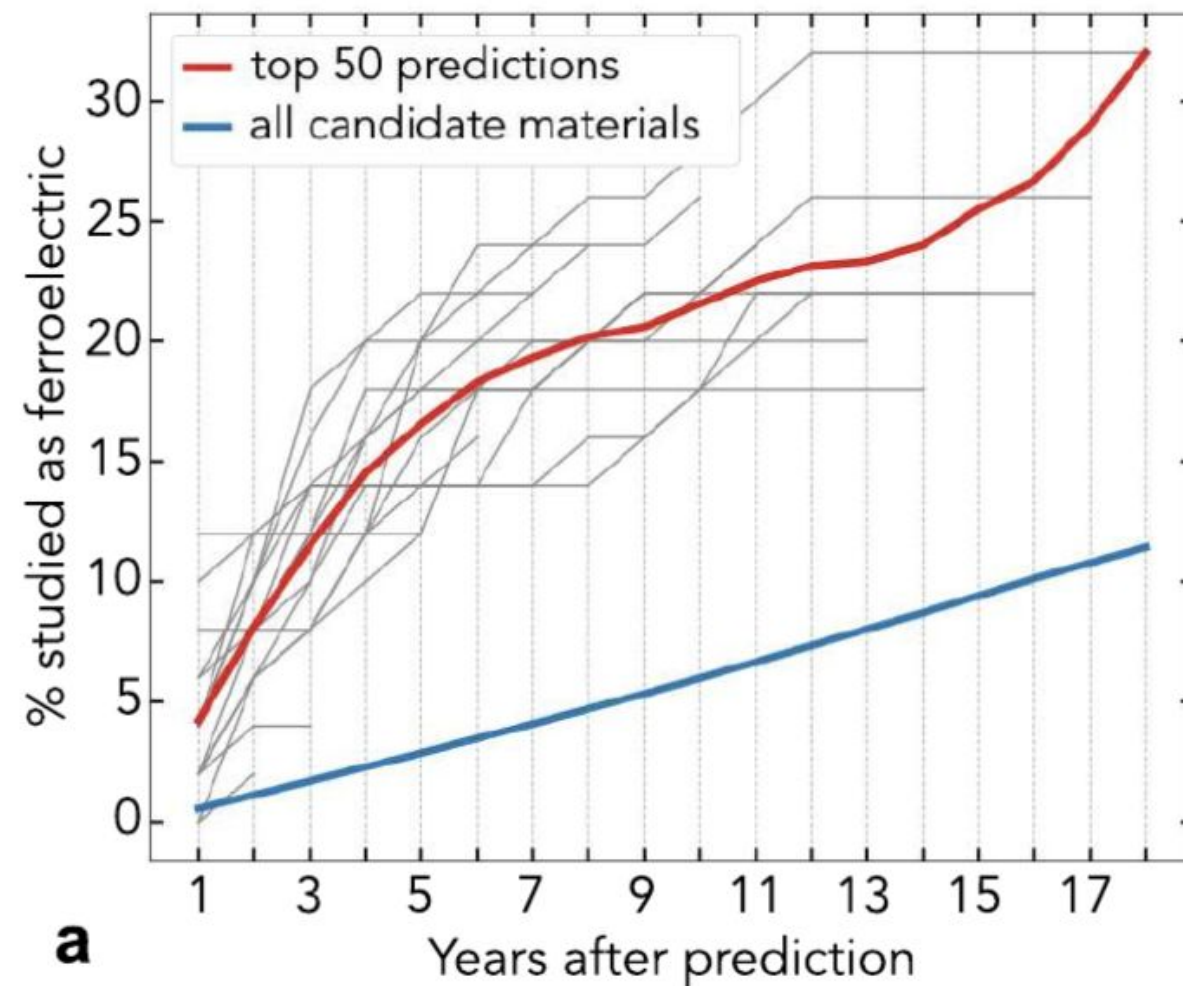
(i.e. remove all abstracts published after a given year, re-train the model, and rank candidates.)

Historical validation

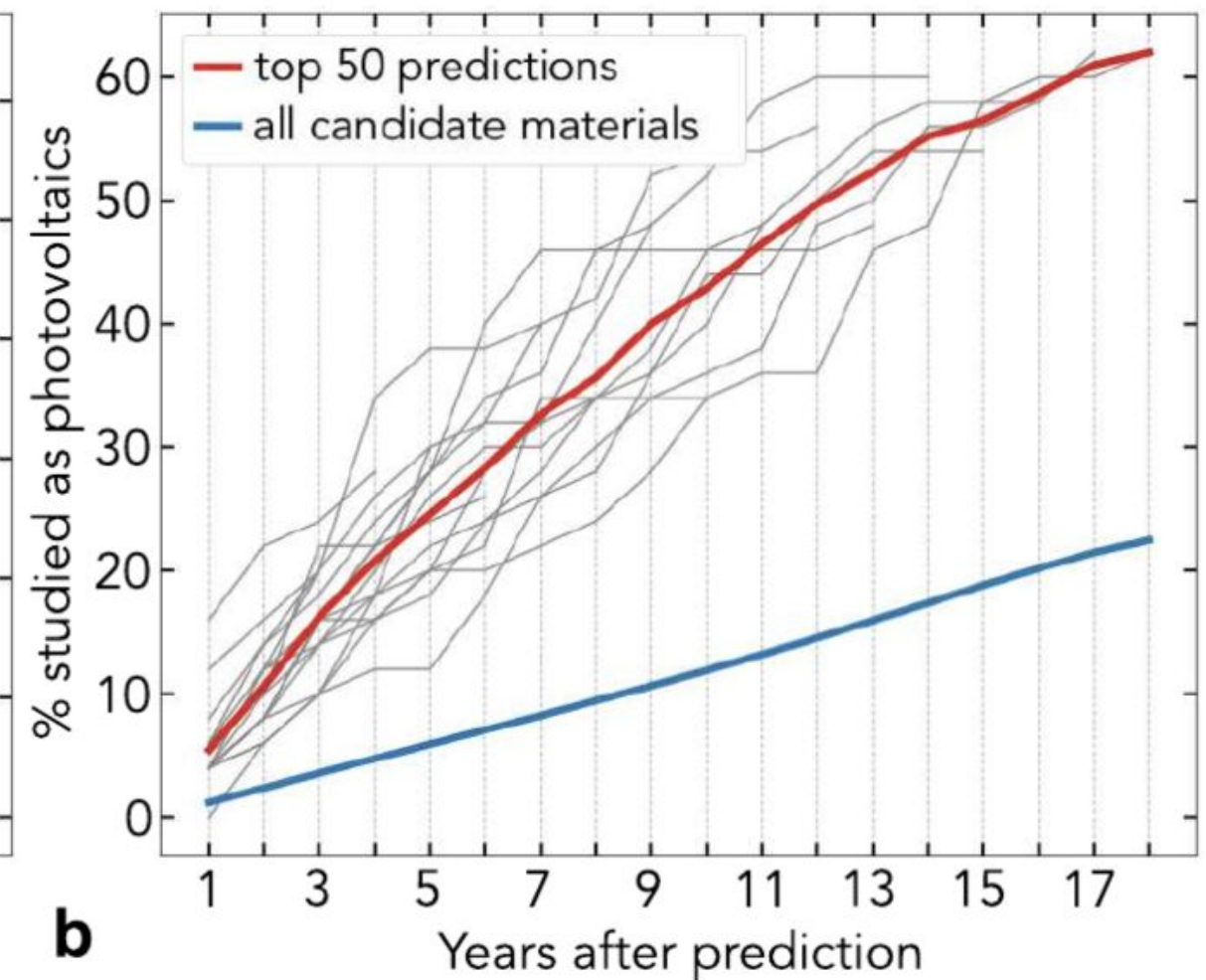


Changing the application is just a matter of changing the query words

"ferroelectric"



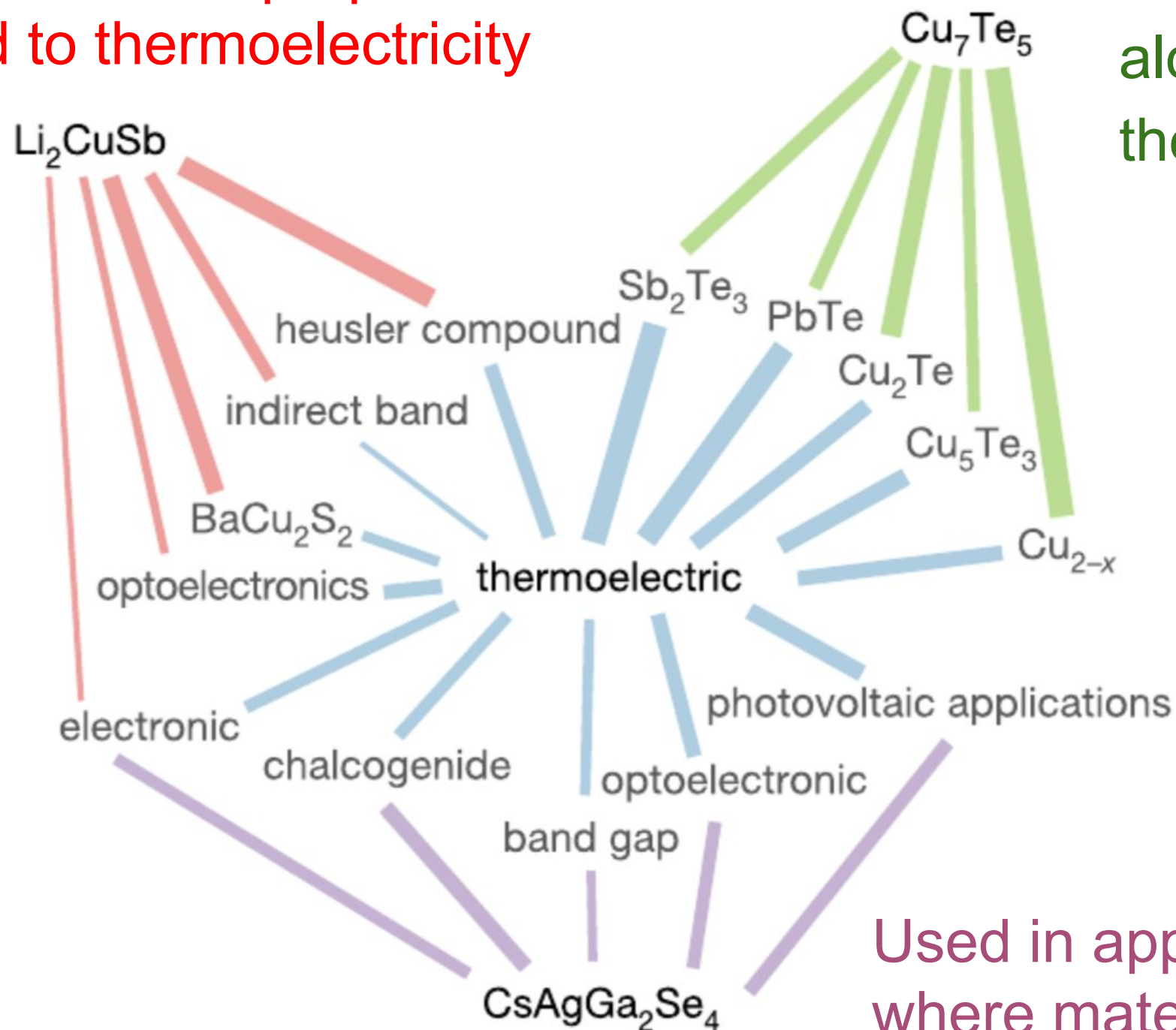
"photovoltaics"

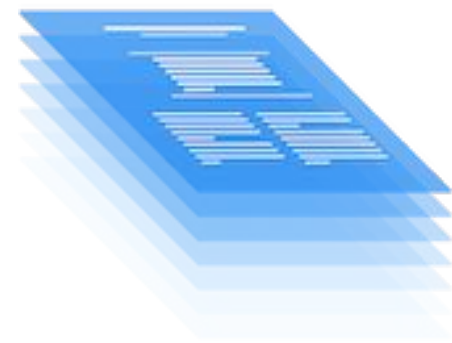


Explaining the predictions

Associated with properties related to thermoelectricity

Commonly referenced alongside known thermoelectrics





MATSCHOLAR

BETA

<https://www.matscholar.com>

Database of all published results

Information retrieval from literature:

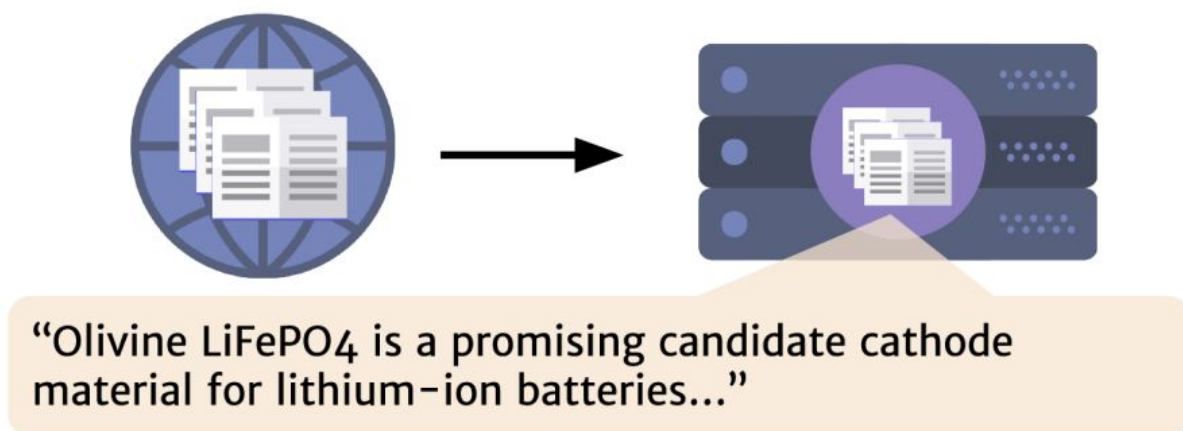
- Material (MAT)
- Material descriptor (DSC)
- Symmetry / Phase (SPL)
- Property (PRO)
- Application (APL)
- Synthesis method (SMT)
- Characterization method (CMT)

GaN thin films are widely used in laser diodes . We have deposited wurtzite GaN films by metalorganic chemical vapor deposition , and the band gap was measured by photoluminescence spectroscopy .

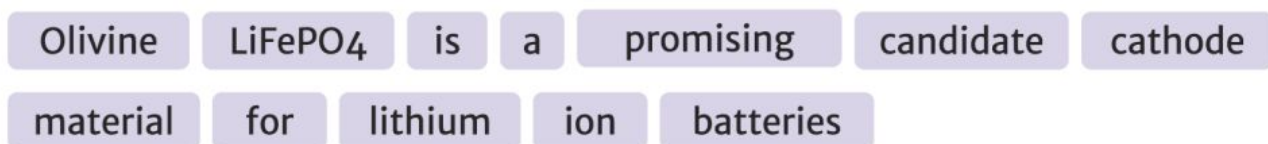
Extracted Entity Tags:

PRO CMT SPL DSC SMT MAT APL

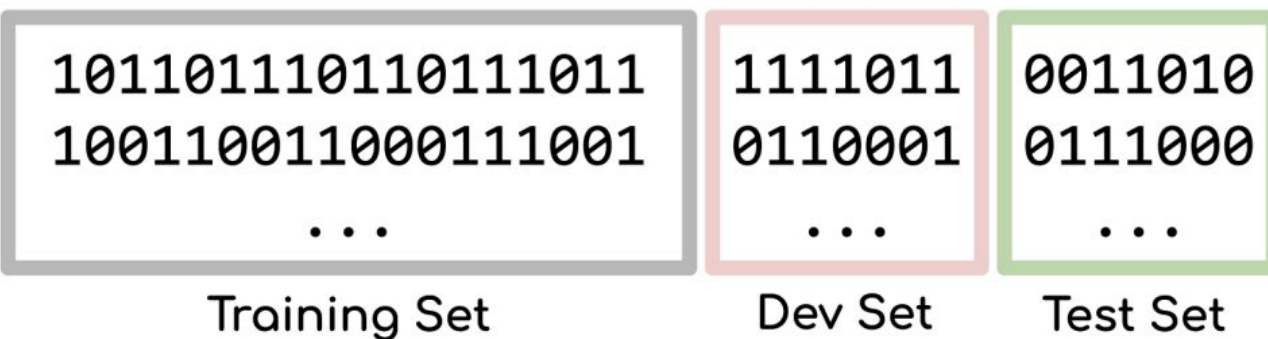
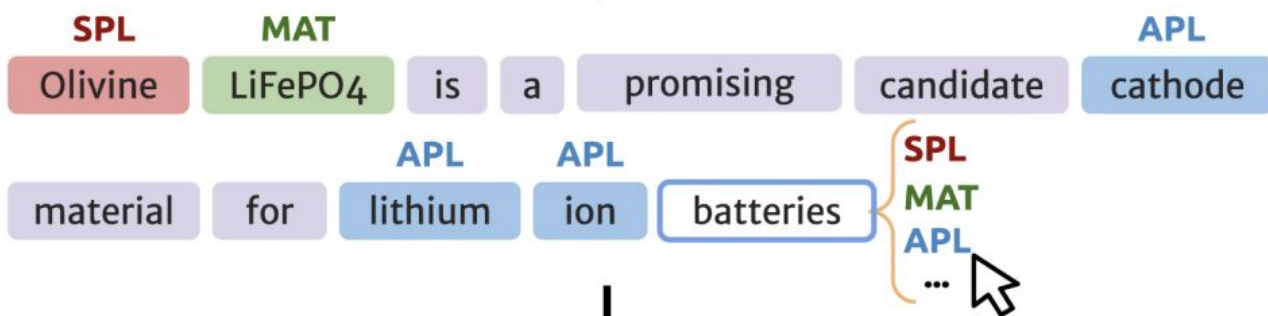
1. Abstracts collected and stored in Matscholar corpus



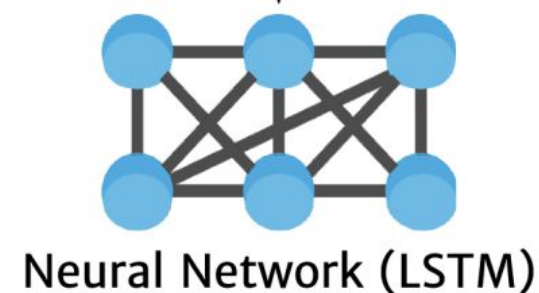
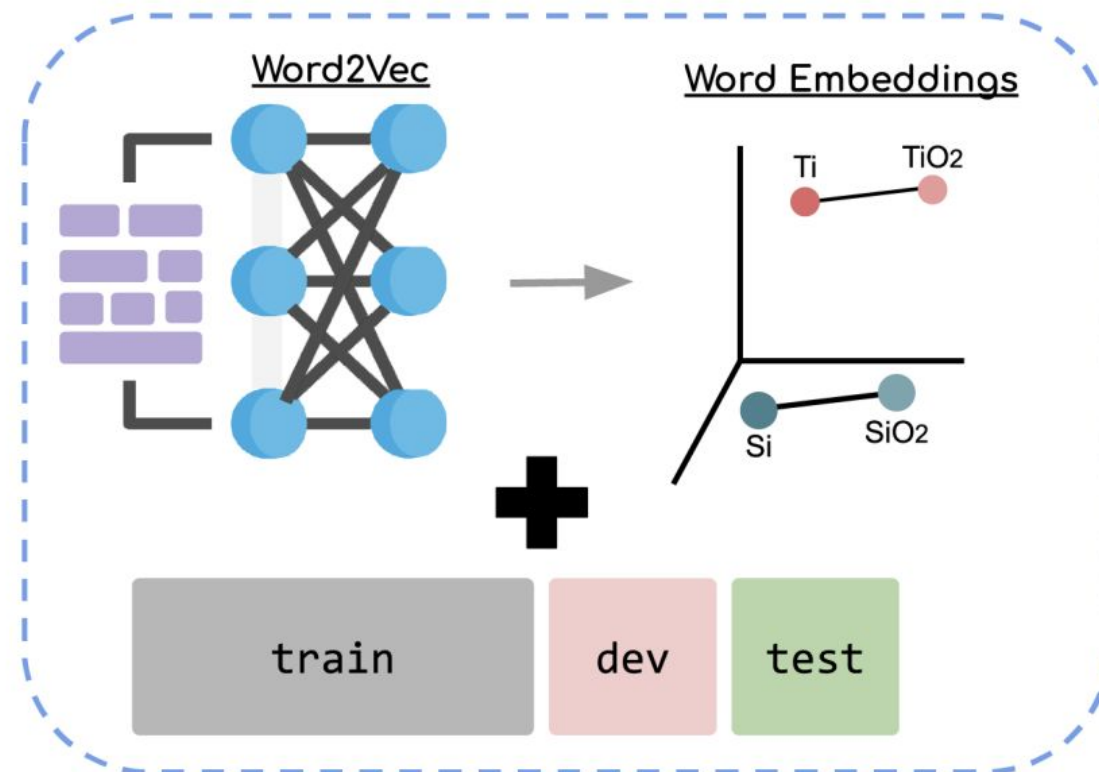
2. Tokenization



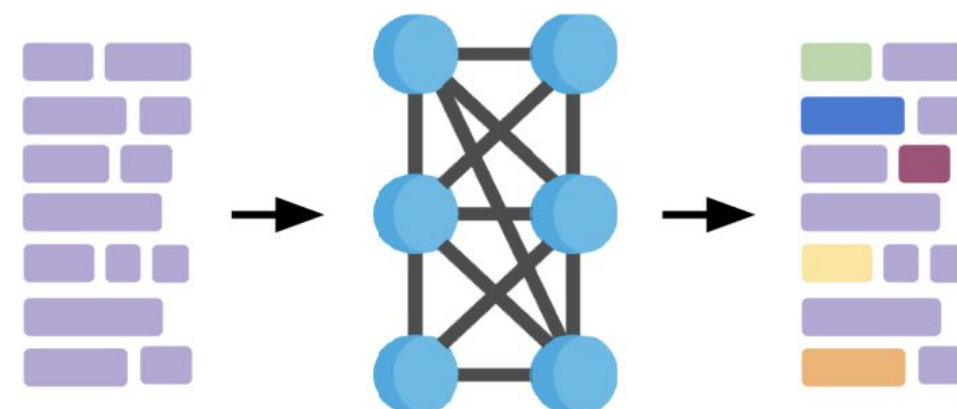
3. Labeling



4. Train model



5. Extract entities with model



Named Entities in Numbers

> 3.3 million
Abstracts Collected

31 million
Properties

19 million
Materials Mentions

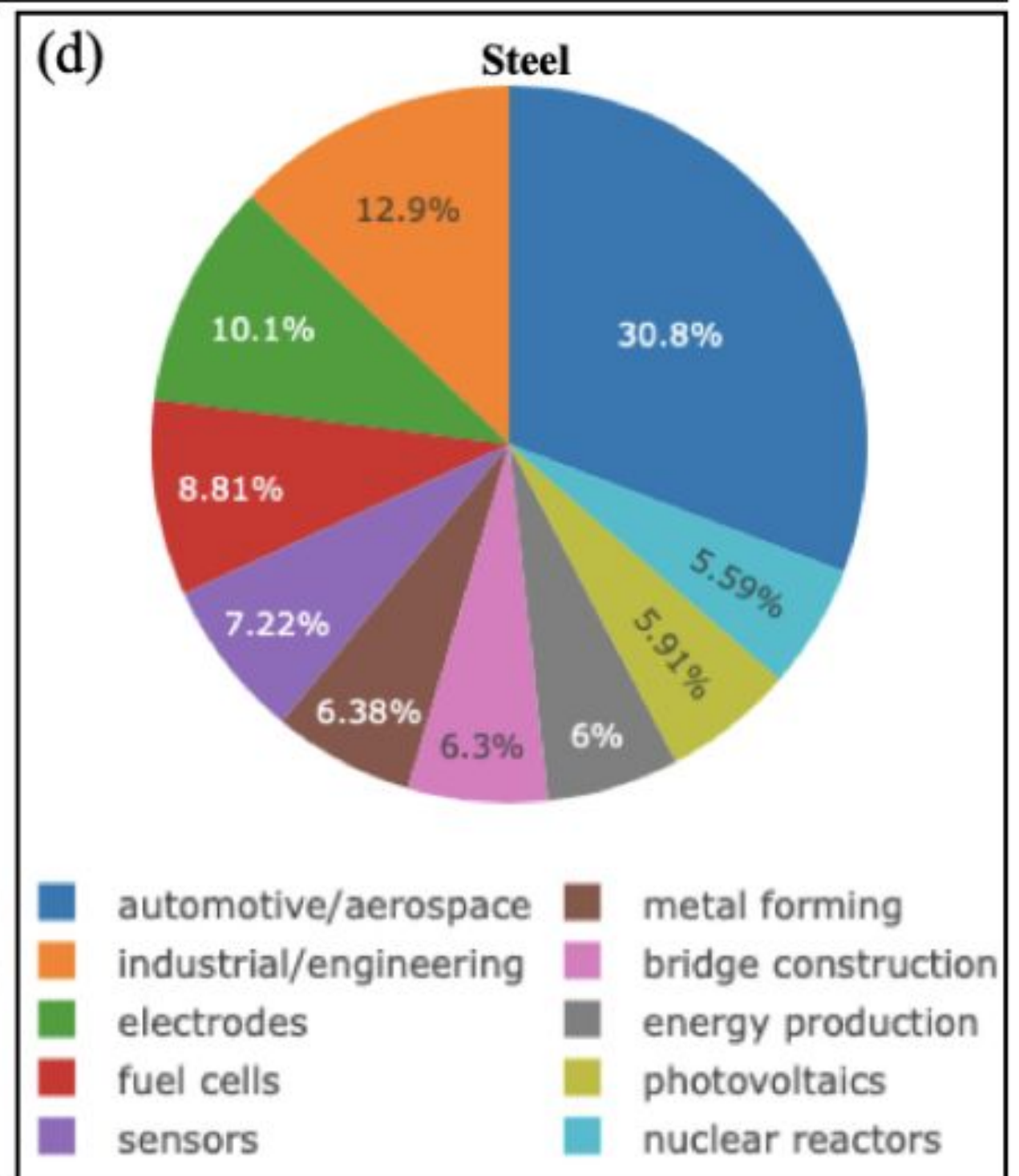
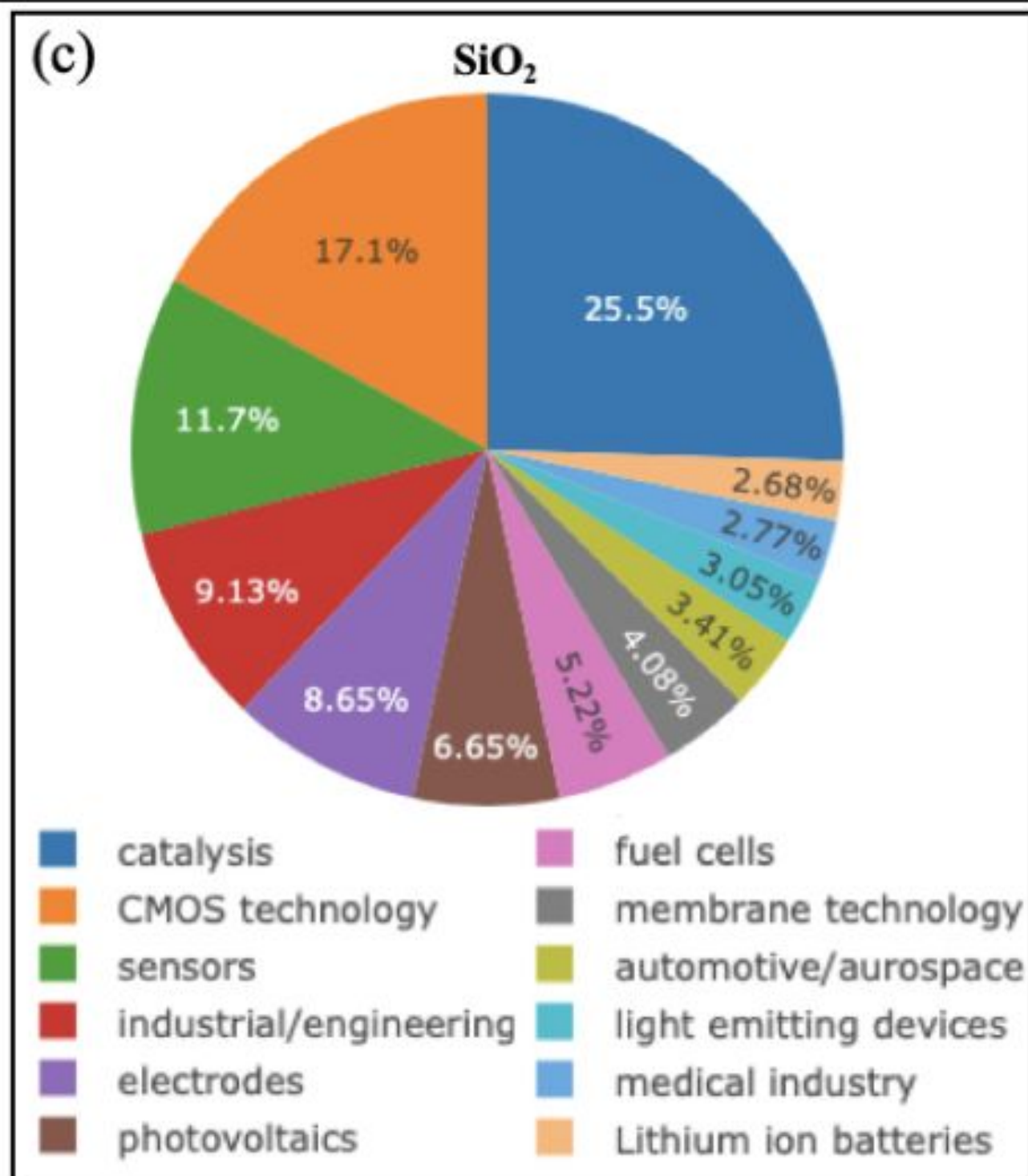
8.8 million
Characterization Methods

7.5 million
Applications

5 million
Synthesis Methods

Working on full text extraction

What are the most common applications of a material?



References

- Tshitoyan et al. [Nature 571, 95–98 \(2019\)](#).
 - <https://github.com/materialsintelligence/mat2vec>
- Weston et al. [J. Chem. Inf. Model. \(2019\)](#).
- REST API: <https://github.com/materialsintelligence/matscholar>
- Website: <https://www.matscholar.com>

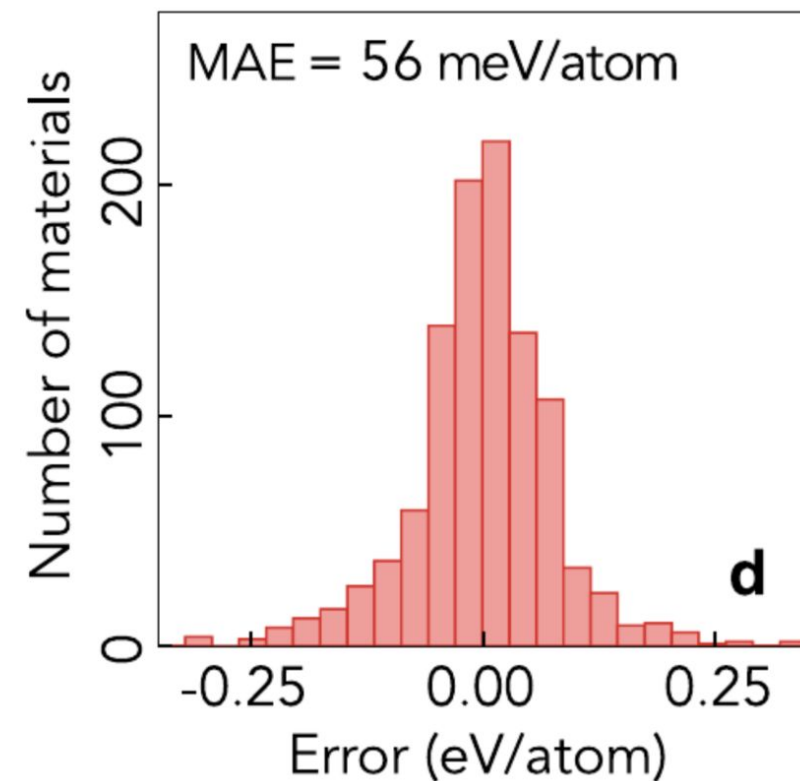
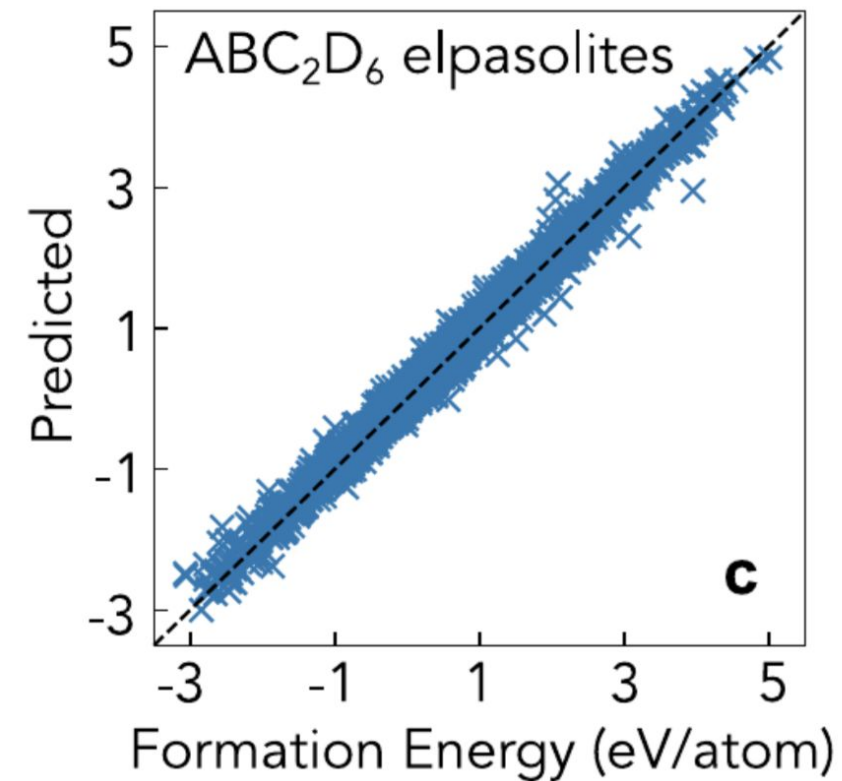
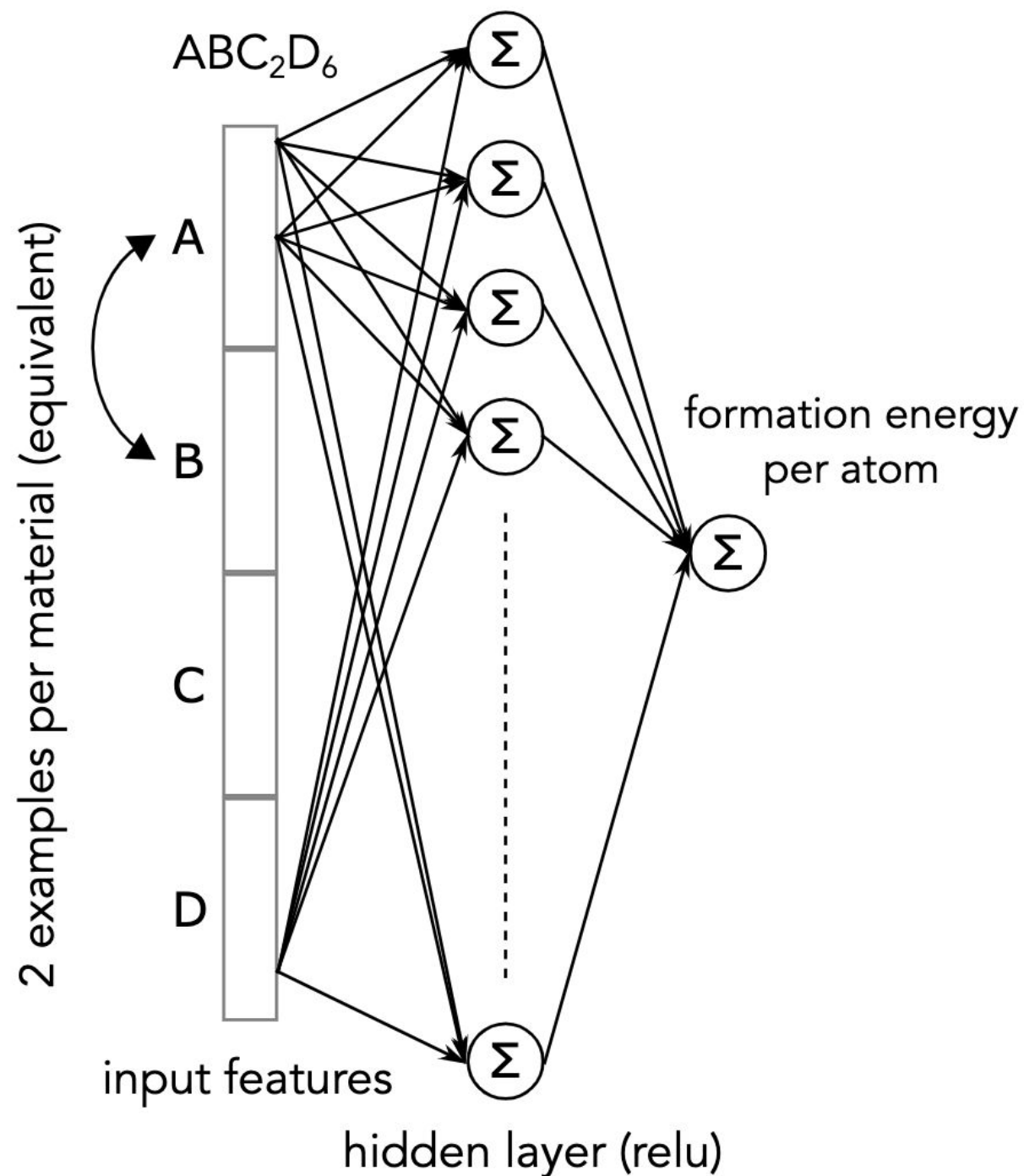


Contact

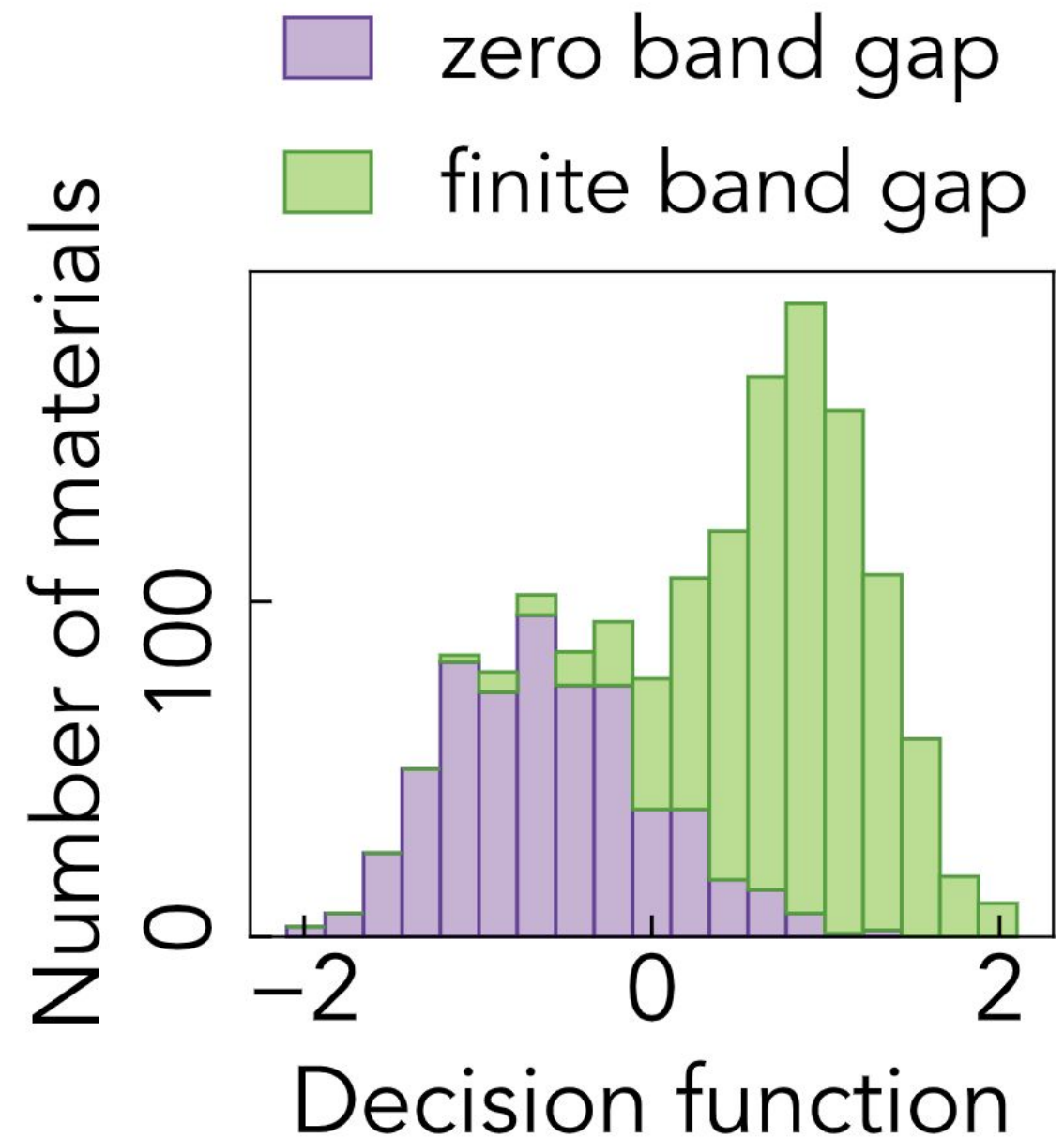
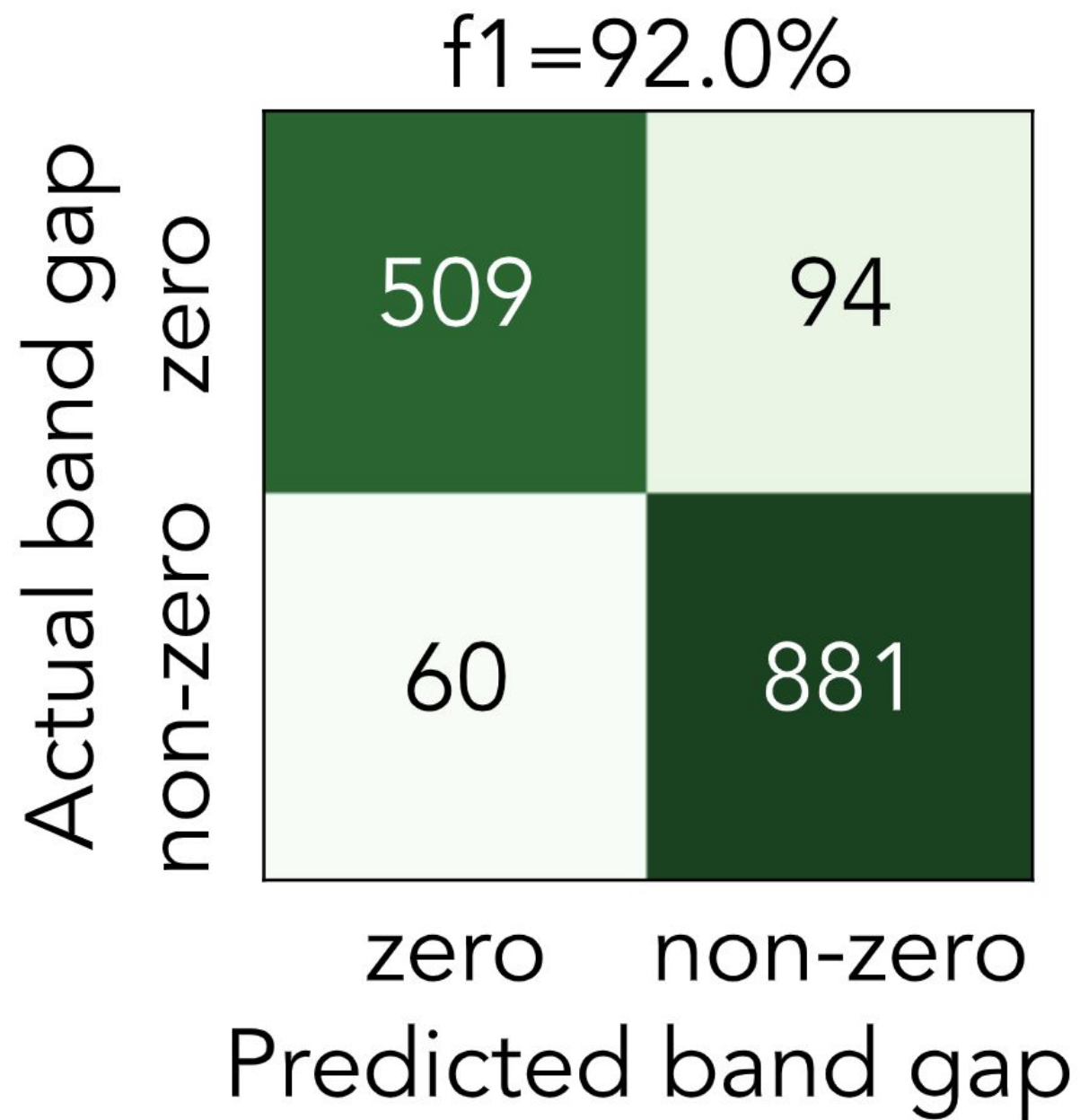
- Email: vahe.tshitoyan@gmail.com
- LinkedIn: </in/vahe-tshitoyan/>
- **Project contact:** Anubhav Jain - ajain@lbl.gov

Extra slides

Formation Energy Prediction



Band gap vs no band gap



The team at Berkeley



Gerbrand Ceder



Anubhav Jain



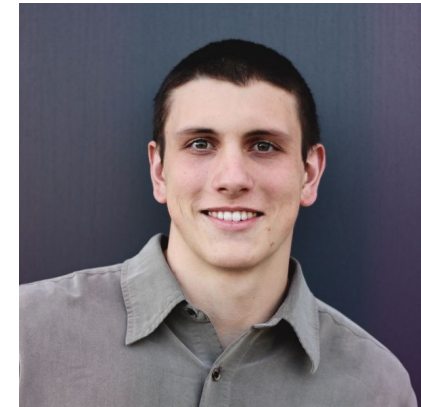
Kristin Persson



Leigh Weston



John Dagdelen



Alex Dunn



Olga Kononova



Ziqin Rong



TOYOTA
RESEARCH INSTITUTE