# Materials Informatics: The Marriage of Materials and Data Sciences

Omer Kaspi,<sup>a</sup> Hadar Binyamin,<sup>a</sup> Olga Girshevitz,<sup>b</sup> Abraham Yosipof,<sup>c</sup> Hanoch Senderowitz<sup>a</sup>

<sup>a</sup>Department of Chemistry, Bar Ilan University, Israel <sup>b</sup>Bar Ilan Institute of Nanotechnology and Advanced Materials, Bar-Ilan University, Israel <sup>c</sup>Department of Information Systems, College of Law & Business, Israel

AIDD School, October 21<sup>st</sup>, Leuven, Belgium

## What Can You Do With Simple Tools?

- No neural networks
- No complex ML algorithms
- But a few cool applications

#### **Materials Informatics is Rapidly Growing**



Senderowitz and Tropsha, JCIM 2018, 58, 1313-1314

#### **The Rise and Rise of Material Informatics**

FIZ Karlsruhe – Leibniz Institute for Information Infrastructure







> 160000 compounds

#### **Material Genome Initiative (NIST)**

The Materials Genome Initiative is a federal multiagency initiative for discovering, manufacturing, and deploying advanced materials twice as fast and at a fraction of the cost compared to traditional methods. The initiative creates policy, resources, and infrastructure to support U.S. institutions in the adoption of methods for accelerating materials development.



## AFLOW (https://aflowlib.org/)

Welcome to AFLOW, a globally available database of **3,528,653** material compounds with over **733,959,824** calculated properties, and growing.



#### **Too Much Information or Too Little Information?**

#### Materials space is HUGE ~10<sup>30-50</sup> candidates



## Drowning in Information but Starving for Knowledge (John Naisbitt)



Tremendous opportunities for materials-informatics Important lessons to be taught from chemoinformatics

#### **Structure-Property Predictions**

- Accurate experimental data (e.g., activities, dependent variables)
- Descriptors (independent variables)
  - Structure-derived , process-derived (measured; calculated)
- A mathematical model
  - *e.g.*, quantitative, qualitative, linear, non-linear
- Model validation
  - Models developed on a training set and tested on an independent test set
  - Models should be simple and interpretable





#### **The Compounds**

#### • Chemoinformatics

- Structures typically well-defined
- Potential exceptions: Polymers and mixtures (but constituting building blocks/components are known)
- True even for combinatorial chemistry
- Materials informatics
  - Structures sometimes welldefined
  - Not true for combinatorial material synthesis





## **The Data**

#### Chemoinformatics

- Primarily concerned with pharmacokinetics / pharmacodynamic related activities
- Diversity comes from the targets / ligands
- Medium / large / very large data sets

#### • Materials informatics

- Diversity comes from activities and the nature of the materials
- Solubility of materials
- \* Biological activities
- Young's modulus
- Thermal conductivity
- Atomization energies

- Glass transition temperatures
- \* Half decomposition temperature
- \* Melting point of ionic liquids
- Viscosity
- Photovoltaic properties
- Tiny / small / medium / large / very large data sets

#### **The Descriptors**

#### Chemoinformatics

- Typically nD (n = 1,5) "classical" descriptors
- Limited usage of QM-derived descriptors
- Materials informatics
  - \* Typically nD (n = 1,5) "classical" descriptors
  - \* Heavy reliance on QM descriptors
  - \* Usage of experimental conditions as descriptors
  - Heavy reliance on measured descriptors (for undefined structures)

#### Focus on Spectra (Raw data, Images)



## **The Algorithms**

- Chemoinformatics and materials informatics
  - Unsupervised methods
    - Data reduction techniques (e.g., PCA)
    - Clustering
  - Supervised methods
    - Classification models (e.g., Random Forests)
    - > Quantitative models (e.g., MLR, SVM, *kNN*, neural networks, DL)

#### Validation

- Chemoinformatics
  - \* "OECD" principles available and frequently followed
- Materials informatics
  - Insufficient external validation
  - Inappropriate control for chance correlation

#### **Prediction of Impact Sensitivity for Energetic Materials**

- Prediction of impact sensitivity of nitro compounds
- 161 compounds, specific and global models MLR, "OECD" validation



calculated 1.5

1.0

0.5

1.0

1.5

experimental

2.0

2.5

3.0

 Good models for nitramine and nitroaliphatic but not for nitroaromatic compounds

Fayet et al., Process Safety Progress (Vol.31, No.3)

#### **Prediction of Impact Sensitivity for Energetic Materials**

- Prediction of impact sensitivity of nitro compounds from "physical principles"
- Sensitivity Index (SI)
  - Number of atoms
  - Dissociation energy of the weakest X-NO<sub>2</sub> bond
  - \* Energy released upon the decomposition of 1 mole of compound



Mathieu, J. Phys. Chem. A 2013, 117, 2253-2259

## **Material Cartography**

#### • Purpose

- Displaying material space (AFLOWLIB)
- Similarity-based Identification of specific materials
- QMSPR models
- Descriptors
  - Band structure fingerprints
  - SiRMS (fragment-like)
  - QM
- Methods
  - Clustering, RF, PLS



Isayev et al., Chem. Mater. 2015, 27, 735–743

## **60 Seconds on Solar Cells**

- 1. Generation of the charge carriers (electrons and holes) due to the absorption of photons
- Separation of the photo-generated charge carriers in a junction via *n*-type (high electron conductivity) and *p*-type (high hole conductivity) semi-conductors
- 3. collection of the photo-generated charge carriers at the termini of the junction Inside a photovoltaic cell
- Key Parameters
  - \* Open circuit voltage ( $V_{OC}$ )
  - \* Short circuit current  $(J_{SC})$
  - Internal quantum efficiency (IQE)
  - ✤ Fill factor (FF)
  - Power Conversion Efficiency (PCE)

$$=\frac{FF \times J_{SC} \times V_{OC}}{P_{in}}$$



#### **Efficiencies of Solar Cells**



#### **Statistical Modeling of Solar Cells**

- Goals
  - \* Identify factors responsible for PV properties
  - \* Build predictive model for PV properties
  - Experimental design



## **Dye Sensitized Solar Cells (Grätzel Cells)**

- Illumination  $\rightarrow$  photon absorbed by dye
- Photo-excited electrons from sensitized dye transferred to TiO<sub>2</sub>
- Dye regenerated by electrolyte
- Electrolyte reduced by counter electrode







## The Importance of the Dye

- Photon Harvesting
- Electron Injection
- Overall cell performances
- Metal-based dyes (Ruthenium)
  - Rare and expensive metals
  - Complex synthesis
  - Low molar extinction coefficients (How strongly a compound absorbs light at a particular wave-length)



Binyamin and Senderowitz, NPJ computational Materials, 8, 142 (2022)

#### The Importance of the Dye

- Metal-free dyes
  - \* Inexpensive
  - ✤ Easy synthesis
  - High molar extinction coefficients
  - \* Tunable photovoltaic and electrochemical properties
  - \* The main disadvantage: lower PCE (power conversion efficiency)



#### The Dye Sensitized Solar Cells Database

- Over 4000 entries of experimentally tested DSSCs
- Preprocessing:
  - Only metal-free dyes
  - No duplicates
  - Exclude cases with co-sensitizers
    & co-adsorbents
  - Uniform testing conditions
- <u>Resulting database: ~1400</u>
  <u>entries</u>

Scaffold	Number of entries
Triphenylamine	621 (42.44%)
Phenothiazine	270 (18.46%)
Carbazole	182 (12.44%)
Indoline	115 (7.86%)
Coumarin	53 (3.62%)
Diphenylamine	33 (2.26%)
Bodipy	16 (1.09%)
Imidazole	14 (0.96%)
Cyanine	14 (0.96%)

Venkatraman et al., J Cheminform. 2018, 10(1):18

#### The Dye Sensitized Solar Cells Database

Scaffold	Structure	'Actives' set	'Decoys' Set
Carbzole		Carbazole dyes subset (165 compounds)	ZINC carbzole library (5916 compounds)
Indoline	NH	Indoline dyes subset (115 compounds)	ZINC indoline library (19620 compounds)
Phenothiazine	S S S S S S S S S S S S S S S S S S S	Phenothiazine dyes subset (254 compounds)	ZINC phenothiazine library (392 compounds)
Triphenylamine		Triphenylamine dyes subset (520 compounds)	ZINC triphenyamine library (488 compounds)

## **Photovoltaphores: Pharmacophores for Dyes**

A pharmacophore can be considered as the highest common denominator of a group of molecules exhibiting a specific "activity"



#### **Photovoltaphores from Active Dyes**



#### **Photovoltaphores from Inactive Dyes**



#### **Model Validation**

- Model tested using small-scale VS campaign
  - Active compounds from DSSCDB
  - Inactive compounds from ZINC (similar scaffolds)
- Models validated using ROC curve and numerical metrics





#### **Virtual Hits**

- HOMO/LUMO
  - \* DFT
  - \* B3LYP
  - & 6-31G(d,p)
- UV-vis absorption spectrum
  - \* TD-DFT
  - ✤ CAM-B3LYP
  - \* 6-31G(d,p)



(ZINC000014358980)

## **Virtual Hits**

- Luminescent, absorbs in UVvis and NIR regions, high molar extinction coefficient
- HOMO remote from semiconductor's conduction band and below the energy level of the redox electrolyte
- LUMO close to semiconductor's surface, and higher than the semiconductor conduction band potential
- Hydrophobic periphery in order to enhance the cell's stability





## Solar Cells Based on Metal Oxides (MO)

- Material
  - Abundant
  - Environmentally safe
  - \* Optimizeable
  - Low cost
- Fabrication
  - Cheap fabrication methods
- Operation
  - Long term operation (stability)

## But Cells Not Efficient Enough New Metal Oxides (MO) Required



#### **Combinatorial Material Science**

- ~60 "useful" elements leading to
  - ~30K inorganic compounds
  - ✤ 3600 binary compounds (ABO<sub>X</sub>); mostly known
  - ✤ 216K ternary compounds (ABCO<sub>x</sub>) many of which unknown



#### **Combinatorial Material Synthesis**



#### **Analysis**



#### **Visualizing the Solar Space**

#### ~1350 cells from 5 libraries

	TiO <sub>2</sub>   Cu-O	TiO <sub>2</sub>   Cu <sub>2</sub> O	TiO <sub>2</sub>   Co <sub>3</sub> O <sub>4</sub>   MoO <sub>3</sub>	TiO <sub>2</sub>   Co <sub>3</sub> O <sub>4</sub>	TiO <sub>2</sub>   CuO-NiO-In <sub>2</sub> O <sub>3</sub>
# Cells	169	338	169	338	338
$V_{OC} [mV]$	31-380	6.6-354	24-620	172-443	111-509
$J_{SC}$ [µA cm <sup>-2</sup> ]	73-290	13.9-406	25-Oct	5.5-11	Jul-54
$P_{max} \ [\mu W \ cm^{-2}]$	0.02-1.02	0-1.26	0.0018-0.11	0.016-0.042	0.015-0.11
FF [%]	23-63	0.16-40	23-41	32-53	26-41
$R_s$ [Ohm cm <sup>2</sup> ]	$18-23 \times 10^3$	$9-0.5 \times 10^6$	$5x10^{3}-0.6x10^{6}$	$8x10^{3}-0.5x10^{6}$	$5x10^{3}-0.6x10^{6}$
$R_{sh}$ [Ohm cm <sup>2</sup> ]	$10^{3}$ -1.9x10 <sup>6</sup>	$5x10^3$ -0.6x10 <sup>6</sup>	$0.7 \times 10^6 - 3.3 \times 10^6$	$2.8 \times 10^{6} - 12 \times 10^{6}$	$0.15 \times 10^{6} - 2.6 \times 10^{6}$
IQE [%]	0.57-1.16	0.1-2.67	0.05-0.3	0.06-0.32	0.02-0.47

#### **Visualizing the Solar Space**



#### **Visualizing the Solar Space**



	PCA	Kernel PCA	Isomap	Diffusion map	Original space
Trust	0.82	0.78	0.70	0.78	
1NN Classification	0.92	0.91	0.92	0.92	0.95
3NN Classification	0.92	0.90	0.93	0.91	0.95
5NN Classification	0.90	0.90	0.93	0.90	0.94

## **Principle Component Analysis (PCA)**

 $TiO_2/Cu-O$ 

 $TiO_2/Fe_2O_3$  (treated at different temperatures)





**PC2** 

TEMP

BGP\_Fe2O3





#### Machine Learning: kNN and GFA

End point	$Q_{LOO}^2$	No applicabi	lity domain With applicability domain			Descriptors	
		$Q_{\rm ext}^2$ ( $R^2$ )	MAE	$Q_{\rm ext}^2$ (R <sup>2</sup> )	MAE	%coverage	
J <sub>sc</sub> (Ag) V <sub>oc</sub> (Ag) IQE(Ag)	0.87 0.86 0.77	0.86 (0.88) 0.73 (0.74) 0.80 (0.84)	0.01 0.02 0.05	0.86 (0.89) 0.75 (0.77) 0.83 (0.86)	0.01 0.02 0.04	83 % 75 % 87 %	Ratio, BGP, D <sub>center</sub> T <sub>7702</sub> , J <sub>max</sub> T <sub>7702</sub> , Ratio, R <sub>a</sub>



Model	$R_{\rm CV}^2$	$Q_{\rm ext}^2$ ( $R^2$ )	MAE
$J_{SC} = 0.062 + 0.0004 \times T_{Cu-O} - 430384.1022/R_{a}$	0.88	0.86 (0.87)	0.01
$V_{OC} = 0.011 \times J_{max} + 1.201 \times 10^{-5} \times T_{TiO_{2}} \times D_{center} - 0.04 - 6.62 \times 10^{-13} \times T_{Cu-O} \times R_{a}$	0.62	0.54 (0.55)	0.03
$IQE = 1.784 \times Ratio + 0.072/Ratio - 2642279.244/(2356681.705 + R_{a})$	0.65	0.74 (0.74)	0.06

#### Yosipof et al., Molecular Informatics, 2015, 34, 367--79

#### Virtual Cell

10



#### **Forensic Informatics**

#### shoeprints

#### faked coins

ammunition





#### gunshot residues

#### glass fragments

DNA



#### **Glass Fragments**

Every year, the Israeli Police Force reports on dozens of crimes that involved glass fragment as evidence which did not realize their forensic potential



#### **Case from Israeli Police**







# PIXE as a Tool for Evidence Characterization



#### **Case from the Israeli Police: Assocoation**





#### Perfect Match!!!!

However, association requires a set of reference structures



#### **Increasing Complexity: Classification**





1000

1 2 3 4 5

6 7 8 9 10 X-ray energy, keV

9 10 11 12 13 14 15 16



Kaspi et al., Talanta 2021, 234, 122608

#### **Machine Learning**



#### **Results**



	Random Train   Test (%)		Train – Surface		Train – Bulk	
			Test – Bulk (%)		Test – Surface (%)	
	RF	<i>k</i> NN	RF	kNN	RF	kNN
Recall	$0.76\pm0.07$	$0.74\pm0.07$	0.86	0.83	0.87	0.77
Precision	$0.80\pm0.09$	$0.77\pm0.10$	0.88	0.86	0.82	0.75
F1 - Score	$0.73\pm0.07$	$0.71\pm0.08$	0.84	0.82	0.84	0.75

#### **Increasing Diversity**

#### Uniting results from three different laboratories



Kaspi et al., Forensic Science International 333(1-3):111216

#### **Further Increasing Diversity**

## Uniting results obtained with different analytical techniques from different laboratories

	PIXE (RBI+BINA)	PIGE +INAA	LA-ICP-MS	EDS
	P: 0.87 ± 0.04	P: 0.78 ± 0.05	P: 0.92 ± 0.02	P: 0.81 ± 0.04
PIXE (RBI+BINA)	R: 0.90 ± 0.03	R: 0.81 ± 0.04	R: 0.93 ± 0.02	R: 0.84 ± 0.03
	F1: 0.87 ± 0.05	F1: 0.77 ± 0.05	F1: 0.92 ± 0.03	F1: 0.81 ± 0.05
		P: 0.92 ± 0.06	P: 0.86 ± 0.06	P: 0.49 ± 0.06
PIGE + INAA		R: 0.87 ± 0.10	R: 0.83 ± 0.09	R: 0.50 ± 0.09
		F1: 0.89 ± 0.08	F1: 0.83 ± 0.07	F1: 0.46 ± 0.07
			P: 0.79 ± 0.09	P: 0.88 ± 0.06
LA-ICP-MS			R: 0.76 ± 0.12	R: 0.89 ± 0.06
			F1: 0.76 ± 0.10	F1: 0.86 ± 0.07
EDS				P: 0.51 ± 0.09
				R: 0.52 ± 0.11
				F1: 0.48 ± 0.10

Kaspi et al., JCIM submitted for publications

#### Conclusions

- Statistical modeling is useful in the field of materials sciences
  - Insight
  - ✤ Experimental design
- Examples discussed
  - Solar cells
  - Forensics
  - Many others
- Challenges
  - Well curated large datasets
  - New descriptors for non well-defined compositions

# Collaborating with experimentalists is indispensable

#### **Acknowledgments**









Omer Kaspi

Avi Yosipof

Hadar Binyamin

Olga Girshevitz

- Arie Zaban (Bar-Ilan University)
- Iva Bogdanović Radović (Ruđer Bošković Institute)
- Jyrki Räisänen (Helsinki University)
- Israel Police Force (Division of Identification and Forensic)



