NVIDIA's GPU Technology Conference (GTC) 2019 San Jose, CA Mar 18–21, 2019

# Materials Discovery with Artificial Intelligence

2019 03. 21.

#### Hyo Sug Lee & Youn-Suk Choi



SAMSUNG ADVANCED INSTITUTE OF TECHNOLOGY [S9967]

#### Contents

- Trend in Computer-Aided Materials Discovery
- High-Throughput Computational Screening & Exhaustive Enumeration
- Deep-Learning-based Evolutionary Design
- Deep-Learning-based Inverse Design
- Efficacy of Computer-Aided Materials Discovery

For accelerated materials discovery



SAMSUNG ADVANCED INSTITUTE OF TECHNOLOGY

Prediction of materials property based on machine learning

Build-up of Materials vs. Property DB → Materials Informatics





Materials design based on machine learning

– Inverse QSAR → Inverse Design



\*GAN: Generative Adversarial Network



In-silico technologies for materials discovery



# High-Throughput Computational Screening & Exhaustive Enumeration

"Landscape of phosphorescent light-emitting energies of homoleptic Ir(III)complexes predicted by a graph-based enumeration and deep learning", GI01.02.02, *2018 MRS fall meeting* 

 Property prediction with high-performance computing for largescale exploration of materials candidates



ML (Machine Learning)-assisted HTCS for higher efficiency



- Exhaustive enumeration based on graph-theory
  - "Graphs"
    - Mathematical structures used to model pairwise relations between objects.
    - Made up of nodes and edges.
    - In chemistry, graph is used to model molecules, where nodes represent atoms and edges represent bonds.



**\* Exhaustive enumeration**: Systematical enumeration of all possible molecules for optimal solution search



Complete list of non-isomorphic graphs



http://www.cadaeic.net/graphpics.htm

- Landscape of phosphorescent light-emitting energies of homoleptic Ir(III)-complex core structures
  - Ir(III)-complexes
    - Widely used as phosphorescent OLED dopants.
    - Figuring out the full landscape of emission color is important for discovering high-performing molecules in target color regions.



- Approach
  - Consider the nodes in graph as rings and edges as ring-connections.
  - Limited the total number rings between 3 and 5.
  - Exclude non-planar type (5-21) and invalid structures as dopant.
    - $\rightarrow$  Only 11 graphs are valid among the total 29 graphs.



SUNG ADVANCED INSTITUTE OF TECHNOLOGY

#### Enumeration

- For 5- and 6-membered rings.
- Substitute some carbons of each molecule with nitrogen atoms (max. five).
  - → Total 9,919,469 (~10M) core structures



#### ⇒ 4. Substitute some carbon atoms with nitrogen atoms

- Property prediction
  - Trained a deep-neural-network model with simulated T<sub>1</sub> data
    - Input: ECFP(Extended Connectivity FingerPrints) of molecular structures
    - Outputs: T<sub>1</sub> energy (phosphorescent light-emitting wavelength)



By simulating the properties of only 0.8% molecules, we can fully scan the chemical space of 10M!

- Results
  - Distribution of T<sub>1</sub> values
  - Blue-color emitting materials are rare compared with red and green



### Conclusions

- In materials discovery, deep-learning-based HTCS is a good alternative to conventional trial-and-error type approach.
- Moreover, exhaustive enumeration makes it possible to systematically explore the whole chemical space.
- With the proposed exhaustive enumeration method based on graph theory and deep learning, the whole landscape of 10M phosphorescent Ir-dopants could be scanned with just 0.8% computational cost compared with the pure simulation-based approach.

"Evolutionary design of organic molecules based on deep learning and genetic algorithm", COMP, ACS fall 2018 National Meeting

# **Evolutionary Design**

- A generic population-based metaheuristic optimization technique
- Uses bio-inspired operators to reach near-optimal solutions
  - ; mutation, crossover, and selection in case of genetic algorithm



#### Proposed approach



MSUNG ADVANCED INSTITUTE OF TECHNOLOGY

- Deep learning models
  - [DNN] 3 hidden layers, 500 hidden units in each layer
  - [RNN] 3 hidden layers, 500 long short-term memory units



- Validation test
  - Design target: change the S<sub>1</sub> (light-absorbing wavelength) of seed molecules
  - Training data: M.W. 200~600 g/mol from PubChem (10,000~50,000 molecules)

No. of	Prediction accuracy of DNN <sup>*1</sup> (R, MAE)			Success rate of
training data	S <sub>1</sub>	номо	LUMO	decoding <sup>*2</sup> (RNN)
1 50,000	0.973, 0.198	0.945, 0.172	0.955, 0.209	86.7%
2 30,000	0.930, 0.228	0.934, 0.191	0.945, 0.224	85.3%
3 10,000	0.913, 0.278	0.885, 0.244	0.917, 0.287	83.2%

%1. No. of test data=No. of training data/10%2. Chemical validity is evaluated with RDKit, No. of test data=5,000



- Evolution toward the increase and decrease of S<sub>1</sub> (eV)
  - Seed: randomly selected 50 molecules (3.8<S<sub>1</sub><4.2)
  - Number of training data = 10k, 30k, 50k



- Evolution under the constraint of HOMO and LUMO (eV)
  - Seed: randomly selected 50 molecules (3.8<S<sub>1</sub><4.2)



NG ADVANCED INSTITUTE OF TECHNOLOGY

Examples of evolved molecules (No. of training data = 50k)



- Constraint (eV)
  - -7.0<HOMO<-5.0
  - LUMO<0.0

#### Conclusions

- A fully data-driven evolutionary molecular design based on deep-learning models (DNN & RNN) was proposed and automatically evolved seed molecules toward target without any pre-defined chemical rules.
- Unlike HTCS, the closed-loop evolutionary workflow guided by deep-learning automatically derived target molecules and found rational design paths by elucidating the relationship between structural features and their effect on the molecular properties.

npj Comput. Mater., 4, 67, 2018

Paradigm shift of ML in computer-aided materials discovery



Implementation of inverse-design model



SAMSUNG ADVANCED INSTITUTE OF TECHNOLOGY 28

- Inverse design of light-absorbing organic molecules (1/2)
  - Training DB
    - 50k molecules sampled from PubChem (M.W. 200~600)
    - DFT calculations for S<sub>1</sub>

#### Distribution of $\lambda_{max}$ of the inverse-designed molecules



**※ About 10% of the designed molecules were found in PubChem even though those were not included in the randomly selected training library.** 

ADVANCED INSTITUTE OF TECHNOLOGY

Inverse design of light-absorbing organic molecules (2/2)



Examples of inverse-designed molecules which share the moieties with well-known dye materials

- Inverse design of hosts for blue phosphorescent OLED (1/3)
  - Target:  $T_1 \ge 3.00 \text{ eV}$

INSTITUTE OF TECHNOLOGY

- Training DB
  - In-house library of 6,000 molecules by combinatorial enumeration (with nine linker (L) and fifty-seven terminal fragments (R) which are frequently employed in OLED hosts; symmetric R-L-R & R-R type enumeration).
  - Property labeling with DFT calculations.



MSI

Inverse design of hosts for blue phosphorescent OLED (2/3)

**Examples of inverse-designed host materials** 



Inverse design of hosts for blue phosphorescent OLED (3/3)

The connection rules of the inverse-designed molecules



#### Conclusions

- A fully data-driven inverse design method successfully extracted the latent materials design rules and proposed target molecular structures without any external intervention.
- The inverse design model successfully proposed new candidates by modifying the assemble rules and creating new fragments.

# **Efficacy of Computer-Aided Materials Discovery**

#### Simulation-based Screening

#### HTCS for pre-defined chemical space



1<sup>st</sup> trial: 1M Candidates QC simulations take 1.5 years Fail to find the target structure

2<sup>nd</sup> trial: 1M Candidates QC simulations take 1.5 years Fail to find the target structure

**3<sup>rd</sup> trial:** 1M Candidates QC simulations take 1.5 years Succeed to find the right structure

Total TAT took 4.5 years



[Step1] Building the training dataset Needs only QC sim. for 50k molecules (27 days)

#### **Inverse Design**

#### Inverse Design [Step2] Deep learning model training with GPU (3 days)

#### **Full search**



[Step3] QC simulations for the proposed molecules (1 day)

#### Total TAT takes 1 month

more than 50X speed up (4.5 years vs. 1 month)

"The inverse design learns by itself the molecular design rules inherent in the libraries and can reduce the effort of researchers and total time to reach the goal"

\* QC simulation tool : turbomole Total computational resources=10,000 CPU In case of 10 CPU computing per molecule, the simulation requires about 13 hrs.



#### **Prospects for AI-based Materials Development**



ADVANCED INSTITUTE OF TECHNOLOGY

SAMSUNG

# Thank you (Q&A)

SAMSUNG

ysuk.choi@samsung.com