WISJ Machine Learning Summer School 2025
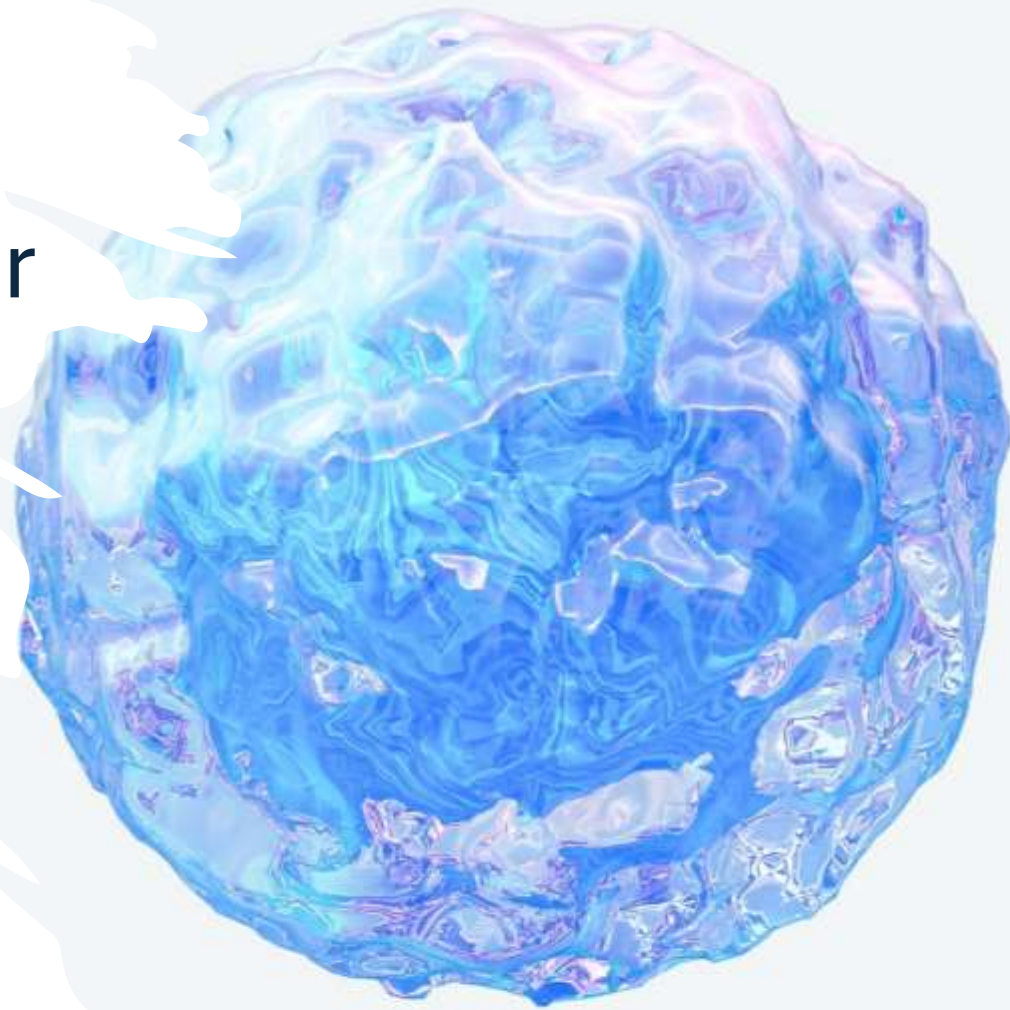
# LLMs and Foundation Models for Material Science

Indra Priyadarsini, Ph.D.
Lisa Hamada, Ph.D.

Research Scientist
IBM Research - Tokyo

Need

↓

Challenges

↓

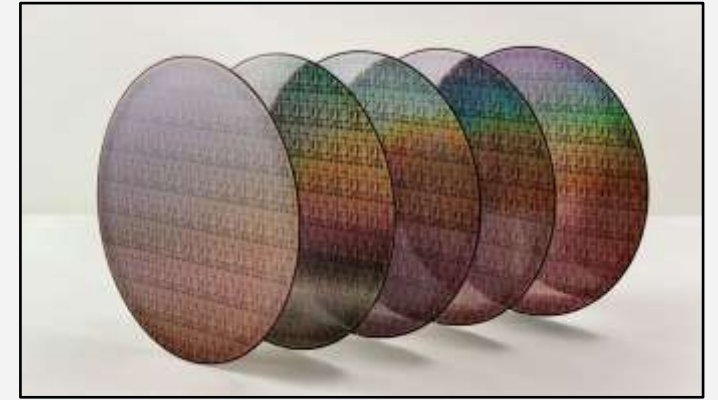Opportunities

# New materials are potentially high impact

but discovery is  very difficult …

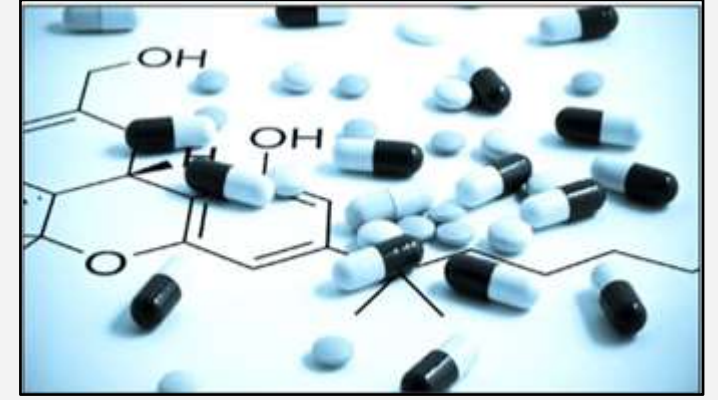Energy and climate solutions

Greener feedstocks

Sustainable semiconductor processing

Clean water & food supplies
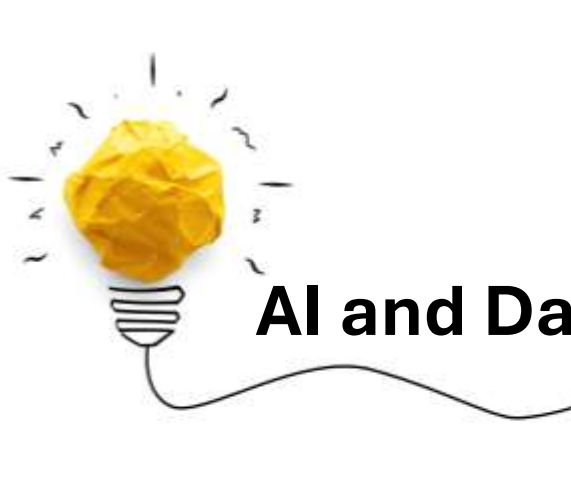
Plastics recycling

Novel therapies

Finding a good suitable material is like finding a needle in a haystack
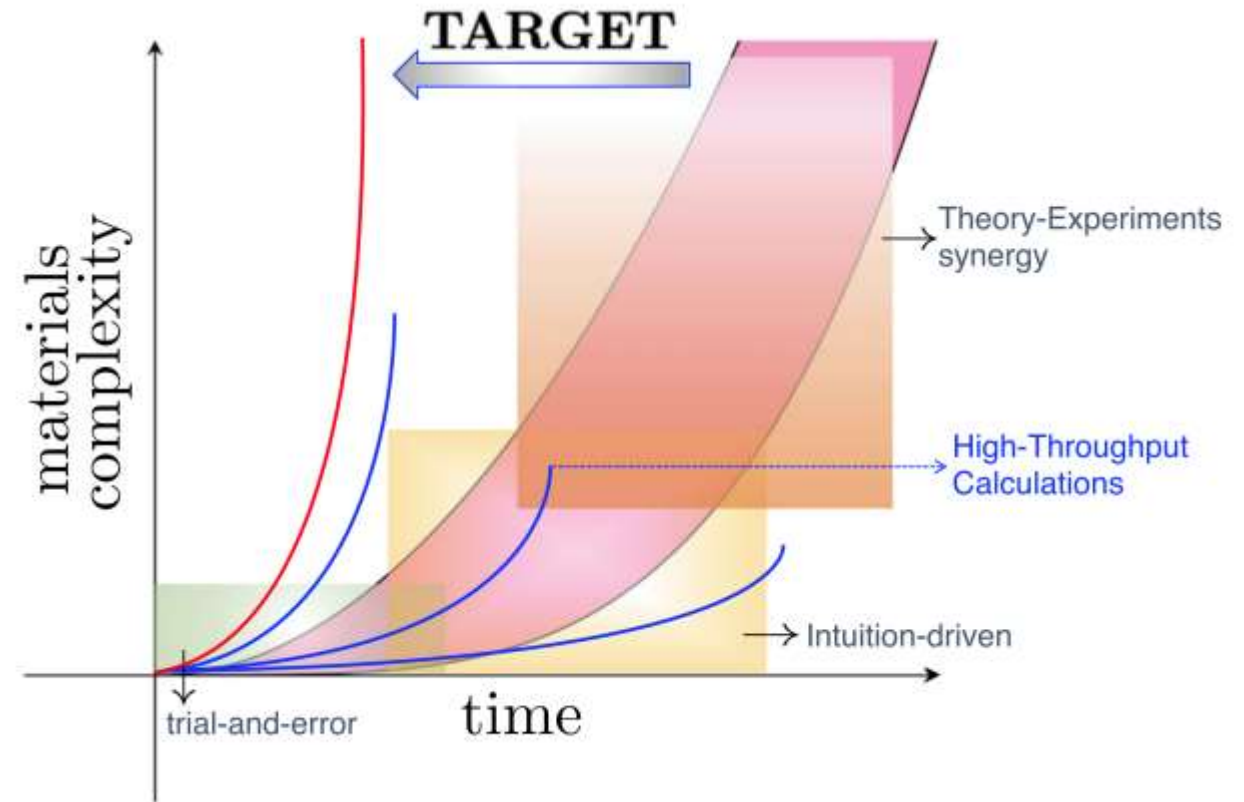
# How do we guide experiments towards materials with desired properties?

Given desired material properties,

$10^{60}$ possible candidates!

**AI and Data-driven Modelling**



T. Lookman, et al. *npj Computational Materials* 5.1 (2019): 21.

# Accelerated Discovery – Discovery Technology Foundations

Accelerating scientific workflows with AI foundation models, toolkits and discovery platform services

| *Deep Search* | *Simulation (ST4SD)* | *Generative (GT4SD)* | *Lab Automation (RXN)* |
|---|---|---|---|
| Deep Search Toolkit for Scientific Discovery | Simulation Toolkit for Scientific Discovery | Generative Toolkit for Scientific Discovery | AI for Chemistry |

**Collection and integration of knowledge**

ds4sd.github.io

**Acceleration of simulation**

github.com/st4sd

**Generation of new chemicals**

github.com/gt4sd

**Automated synthesis**

rxn.app.accelerate.science

## *Foundation Model for Materials*

Foundation Model for materials, chemistry, biology

# Global Presence

## 3000
Researchers

## 79
Years

## 18
Sites



Albany
Yorktown
Almaden
Cambridge

Dublin
Daresbury
Hursley
Zurich

Haifa

Delhi

Tokyo
Shin-Kawasaki

Bangalore

Nairobi

Singapore

Rio de Janeiro
Sao Paulo

Johannesburg

**Members working on Foundation Model for Materials and Chemistry**

# AI is gaining focus in materials, but efforts are fragmented and limited

Uni-modal models for discrete chemical classes focus on limited tasks such as prediction of individual properties, de novo generation of molecules, or prediction of synthesis pathways, etc.
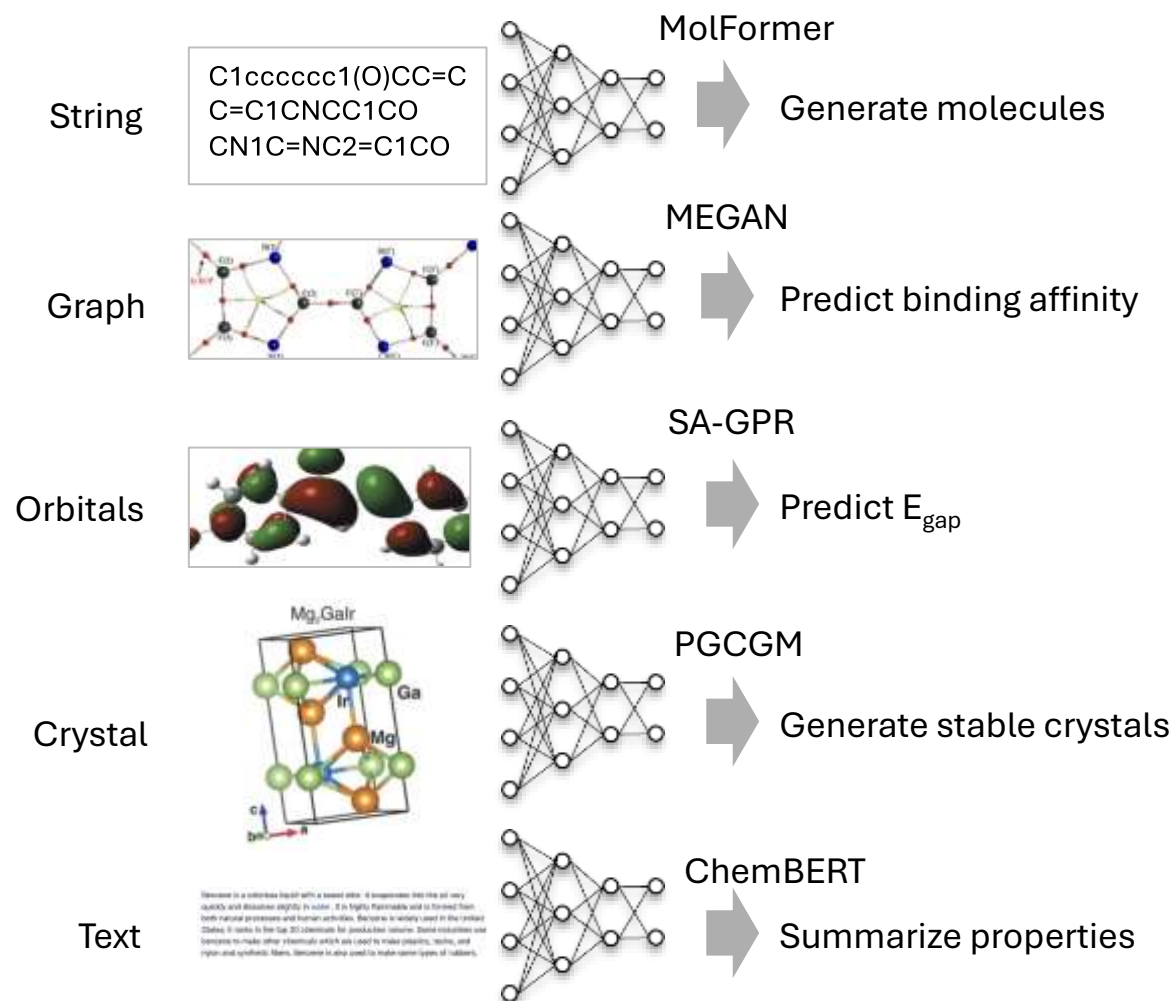
Limitations:
- Scale & data – **today's models are small & limited**
  - Small parameter size : ~100M
  - Small ground-truth data: ~1M
  - Single modality
  - Few chemical classes
  - Limited tasks

- Performance – **today's models perform poorly**
  - Insufficient accuracies
  - Heavy emphasis still on human-workflow
  - No synergies between models

- Many redundant efforts reinventing similar models

String

C1cccccc1(O)CC=C
C=C1CNCC1CO
CN1C=NC2=C1CO

MolFormer

Generate molecules

Graph

MEGAN

Predict binding affinity

Orbitals

SA-GPR

Predict $E_{gap}$

Crystal

$Mg_yGaIr$

PGCGM

Generate stable crystals
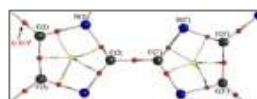
Text

ChemBERT

Summarize properties

# Greater impact can be achieved by centralizing global efforts to build a foundation model

A foundation model trained with multi-modal data sets can be applied to diverse classes of downstream application tasks.

- **Increase scale and grow data:**
  - Larger parameter models: 10B+
  - More diverse data
  - Multi-modality & multi-representation

- **Improve performance and impact:**
  - Richer representation by fusing data
  - Higher accuracy in predictions
  - Higher fidelity generation
  - Integrated knowledge

**More modalities & representations**

**Fused multi-modal FMs and framework**

**Downstream tasks**

Prediction of comprehensive properties

Generation of molecules, reaction pathways, materials characterization

Prediction of experimental outcomes

**Conversational user interface**

# Multi-modal representations enhance modeling capabilities



Spectrum

Basic properties

Atom positions

Electron density

Energy diagram

Fingerprint

Graph

CN1C=NC2=C1C(=O)N(C(=O)N2C)C

SMILES

Synthesis

Text description

# Multi-modal representations enhance modeling capabilities



Spectrum

Basic properties

Atom positions

Electron density

Energy diagram

Fingerprint

Graph

CN1C=NC2=C1C(=O)N(C(=O)N2C)C
SMILES

Synthesis

Text description

Fused representation  (0.1, 2.5, 8.3, ...,3.7)

**Downstream Tasks #1**
*High-accuracy predictions*

Downstream models

# Multi-modal representations enhance modeling capabilities



Spectrum

Basic properties

LUMO (eV)    -0.8.

Dipole moment    3.80
(debye)

Atom positions

Electron density

Energy diagram

Fingerprint

Graph

CN1C=NC2=C1C(=
O)N(C(=O)N2C)C
SMILES

Synthesis

...e is a methylxanthine
...structurally related to adenosine
...ingestion, caffeine binds to adenosi
...mediated downregulation of CNS a

Text description

# Multi-modal representations enhance modeling capabilities



Spectrum

Basic properties

LUMO (eV)    -0.8.
Dipole moment    3.80
(debye)

Atom positions

Electron density

Energy diagram

Fingerprint

Graph

SMILES

CN1C=NC2=C1C(=
O)N(C(=O)N2C)C

Synthesis

Text description

*Predict properties
from a molecular structure*

# Multi-modal representations enhance modeling capabilities



*Spectrum*

Basic properties

*Atom positions*

*Electron density*

*Energy diagram*

*Fingerprint*

*Text description*

*Graph*

CN1C=NC2=C1C(= O)N(C(=O)N2C)C

SMILES

*Synthesis*

*Generate molecular structures from properties*

LUMO (eV)          -0.8.-

Dipole moment      3.80
(debye)

# Multi-modal representations enhance modeling capabilities



Basic properties

Spectrum

*Atom positions*

*Electron density*

*Energy diagram*

*Fingerprint*

*Graph*

**Generate molecular structures from properties and spectrum**

CN1C=NC2=C1C(= O)N(C(=O)N2C)C

SMILES

*Synthesis*

*Text description*

**(Query)** Generate a brand new molecular structure having the following properties and emission spectrum with it's synthetic path.

- $\Delta E$ = 2.3 eV
- $T_{melt}$ = 150 ℃
...

**(Answer)** Here is the 1st candidate.

# Overview of Model Architecture

**Modality-specific Foundation Models**

**Fused Foundation Model(s)**

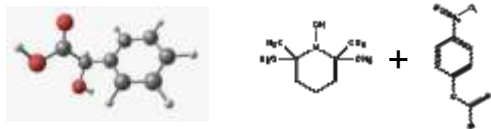**Downstream tasks**

C1CCC...

[C][C=][C]...

SMILES model
*Transformer*

SELFIES model
*BART*

Graph model
*GNN + MHG*

3D pos. model
*Transformer*

Elec. density model
*Transformer*

Fingerprint encoder
*Graph algorithm*

Spectrum model
*GRU, LSTM*

The molecule is a jak inhibitor, immunomodulator, protein tyrosine kinase inhibitor, protein kinase inhibitor and belongs to the autoimmune disease treatment class of molecules.

Text model
*MolT5, Mistral, etc.*

...

Late fusion algorithms
• Mixture of Experts
• Dynamic fusion
• Contrastive learning

Encoder          Decoder

Fused representation vectors

Powerful representation for predictions
• Feature selection
• Domain constraints

$T_g$

$\lambda_{abs}$

Cross-modal inferences
• GFlowNet
• Contrastive learning

SMILES → → SMILES

Props → → Props

Text → → Text

# Overview of Model Architecture

# (1) SMILES Model

Transformers-based encoder-decoder model pre-trained on 91 M samples curated from PubChem (4 billion molecular tokens). We also include a **Mamba-based** variant for faster inference, and **Polymer SMILES** version.

Eduardo Almeida Soares



Eduardo Soares, et al. "A Large Encoder-Decoder Family of Foundation Models For Chemical Language." arXiv preprint arXiv:2407.20267 (2024). (Under Review NeurIPS)

Eduardo Soares, et al. "MoLMamba: A Large State-Space-based Foundation Model for Chemistry." American Chemical Society (ACS) Fall Meeting. 2024.

**Open-released in July '24**

## Variety of models

| Model | Num Parameters | Architecture | Input | Open- Source |
|---|---|---|---|---|
| **Base** | 289M | Transformers | SMILES | Yes |
| **Large** | 738M | Transformers | SMILES | No |
| **XL** | 2.5B | Transformers | SMILES | No |
| **MoE** | 8 X 289M | Transformers | SMILES | Yes |
| **SSM** | 336M | SSM - Mamba | SMILES | In process |
| **PSMILES** | 289M | Transformers | Polymers | No |

## Key Performance example (MoleculeNet and QM9)

| Methods | BBBP | HIV | BACE | ClinTox | SIDER | TOX21 | QM9 |
|---|---|---|---|---|---|---|---|
| MolFM | 72.9 | 78.8 | 78.8 | 79.7 | 64.2 | 77.2 | - |
| MoLFormer | 90.9 | **80.5** | 86.3 | 91.2 | 65.5 | 80.4 | 1.59 |
| **IBM.materials.smi-TED289M (Frozen)** | 91.4 | **80.5** | 85.6 | 93.5 | **66.01** | 81.5 | 7.49 |
| **IBM.materials.smi-TED289M (Fine-Tuned)** | **92.2** | 76.8 | **88.2** | **94.3** | 65.6 | **81.9** | **1.32** |

# (2) SELFIES Model

BART model is trained with SELFIES strings in a self-supervised manner with/without masking tokens.
**1 billion** training samples extracted from ZINC22 and PubChem are used.

Indra Priyadarsini



**Input**

[C][=C][C][=C][C][Ring1][Branch1]

SELFIES strings

**Bidirectional Encoder**

**Latent representation**

**Autoregressive Decoder**

**BART** Transformer Model

**Output**

[C][=C][C][=C][C][Ring1][Branch1]

SELFIES strings

**To be released in August !!**

## Variety of models

| Model | Num Parameters | Dataset | Tokens | Samples Trained |
|-------|----------------|---------|--------|-----------------|
| **Mini** | 2.2 M | ZINC 22 | 173 | 8 B |
| **Base** | 354M | ZINC 22 | 173 | 1B |
| **Base-mix** | 354M | ZINC + PubChem | 3160 | 1B |
| **Large** | 1 B | ZINC + PubChem | 3160 | 100 M |

## Key Performance example

| Model | BBBP | HIV | BACE | ClinTox | Sider | Tox21 |
|-------|------|-----|------|---------|-------|-------|
| ChemBERTa | 64.3 | 62.2 | 79.9 | 73.3 | - | 72.8 |
| MolFormer-ZINC | 89.9 | 78.4 | 87.7 | 82.2 | 66.8 | 83.2 |
| MolFormer-XL | 93.7 | 82.2 | 88.2 | 94.8 | 69 | **84.7** |
| SELFormer | 90.2 | 68.1 | 83.2 | - | **74.5** | 65.3 |
| SELF-BART (mini) | 92.6 | 74.2 | 87 | 88.3 | 62.4 | 75.1 |
| SELF-BART (base) | **95.2** | **83** | **88.8** | **96.9** | 65 | 76.5 |

Indra Priyadarsini et al., "A transformer based large-scale molecular representation model." Materials Research Society (MRS) Fall Meeting. 2023.

# (3) Molecular Graph



Akihiro Kishimoto

Input

G (V,E)

**Can understand Molecules as graphs**

Encoder
**GNN**

Latent representation

Decoder
**MHG+RNN**

Training Production rule sequence

"0, 3, 7, 1, 1, 2, 4, -1"

**Can always generate *valid* molecules**

Downstream task Property prediction

- $\Delta E$ = 2.3 eV
- $T_{melt}$ = 150 ℃
- …

Kishimoto, Akihiro, et al. "Autoencoder based on Graph and Recurrent Neural Networks and Application to Property Prediction." *Materials Research Society (MRS) Fall Meeting*. 2023.

Kishimoto, Akihiro, et al. "MHG-GNN: Combination of Molecular Hypergraph Grammar with Graph Neural Network." *AI4Mat Workshop @ NeurIPS* (2023).

**Example of MHG-GNN performance evaluation (R2 score ↑)**

| | Polymer | | | Photoresist | | Chromophore |
|---|---|---|---|---|---|---|
| Method | Density | Bulk Modulus | Refractive Index | Homo | Lumo | $\Lambda_{max}$ on NIR |
| MHG-GNN | **0.578** | **0.516** | **0.865** | **0.896** | **0.845** | **0.845** |
| ECFP6 | 0.523 | 0.482 | 0.823 | 0.791 | 0.782 | 0.708 |
| Modred | 0.567 | 0.505 | 0.859 | 0.894 | 0.830 | 0.842 |

JSR

NeurIPS '23 AI4Mat

# (4) Topological distance descriptor

Algorithm based model for extracting topological distance between substructures pairs considering intramolecular interaction.

Lisa Hamada



- Can target any substructure without size limitations
- Can handle distances of more than five bonds (a limitation in GNN)
- Consider multiple exception handling scenarios

L. Hamada, *et al*. "Molecular Descriptors Accounting for Intramolecular Interactions and Application to Chemical Property Prediction." *American Chemical Society (ACS) Fall Meeting*. 2022.

## Key Performance example (R² score)

| Dataset | Cmp. ($CHCl_3$) | | Cmp. ($CH_3OH$) | |
|---|---|---|---|---|
| Structure | Large | | Large - Medium | |
| property | $\lambda_{abs.}^{max}$ | $\Delta\lambda$ | $\lambda_{abs.}^{max}$ | $\Delta\lambda$ |
| MolFormer | 0.83 | 0.56 | 0.87 | 0.45 |
| MolCLR | 0.73 | 0.43 | 0.78 | 0.40 |
| Mordred | 0.84 | 0.51 | 0.86 | 0.46 |
| Atom Pair | 0.84 | 0.58 | 0.84 | 0.50 |
| **Our method** | **0.94** | **0.77** | **0.93** | **0.72** |

## High explainability of Chemical Insight

**Cmp. ($CH_3OH$) $\lambda_{abs.}^{max}$**

| Rank | Sub1 | Sub2 | Rep. Mol. |
|---|---|---|---|
| 1st | | | Squarium |
| 3rd | | | DPP |

**Cmp. ($CH_3OH$) $\Delta\lambda$**

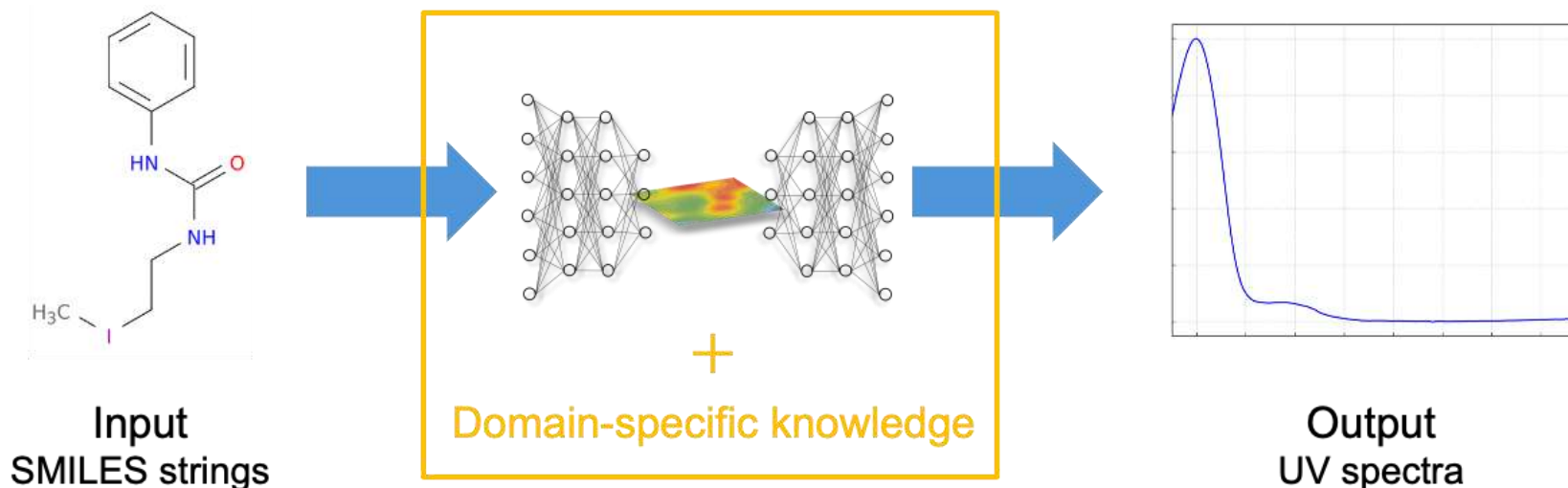| Rank | Sub1 | Sub2 | Rep. Mol. |
|---|---|---|---|
| 2nd | | | Donor / Acceptor |
| 4th | | | |

# (5) Spectrum

Hajime Shinohara

UV spectrum is limited data availability due to the experimental setting.
By implementing domain-specific knowledge into the model, the prediction performance has been increased.



**Input**
SMILES strings

**+ Domain-specific knowledge**

**Output**
UV spectra

Shinohara, Hajime, et al. "Pre-Treatment Methods for Machine Learning in Finer UV Spectrum Inference." *Materials Research Society (MRS) Fall Meeting*. 2023.

## Domain-specific knowledge implementation of UV spectrum from organic molecules

Ex)
Peak position addition
Curvature limitation
Curriculum learning method

## Key results example



By implementing domain-specific knowledge of UV spectrum from organic molecules, the prediction performance of the spectrum has increased in various models even with small dataset (~3k)
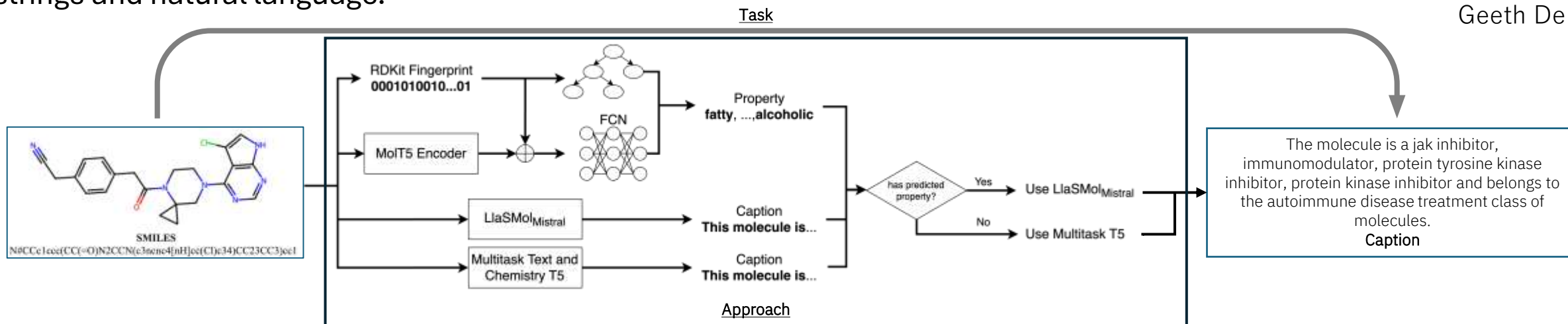
| Model | FCN | LSTM | GRU | RNN |
|---|---|---|---|---|
| Baseline | 0.1513 | 0.1712 | 0.1597 | 0.1733 |
| Our model | **0.1377** | **0.1451** | **0.1445** | **0.1716** |

# (6) Text description - Molecule Captioning from SMILES

**"Translating"** between molecules encoded in SMILES strings and natural language.

Geeth De Mel



## Data Sources

- PubChem
- Chemical Function (CheF)
- ChemFOnt

## Models

- MolT5 Small/Base/Large - 60M ~ 770M
- Multitask text and Chemistry Small/Base augm - 60M ~ 223M
- Meditron - 7B
- Mistral - 7B
- XGBoost-based molecule property predictor

## Results

| Molecule Type | Model | BLEU-2 | ROUGE-L | METEOR |
|---|---|---|---|---|
| Has Predicted Props. | Multitask T5 | 82.15 | 60.20 | 87.05 |
| | LlaSMol_Mistral | **82.66** | **60.54** | **87.70** |
| No Props. Predicted | Multitask T5 | **43.12** | **50.67** | **51.87** |
| | LlaSMol_Mistral | 35.24 | 47.95 | 45.50 |

Translation metrics by molecular type on dev. set

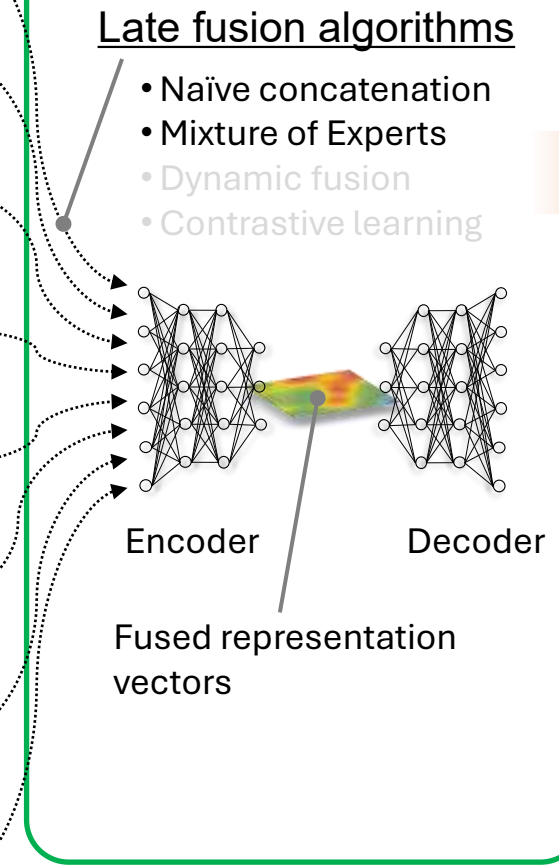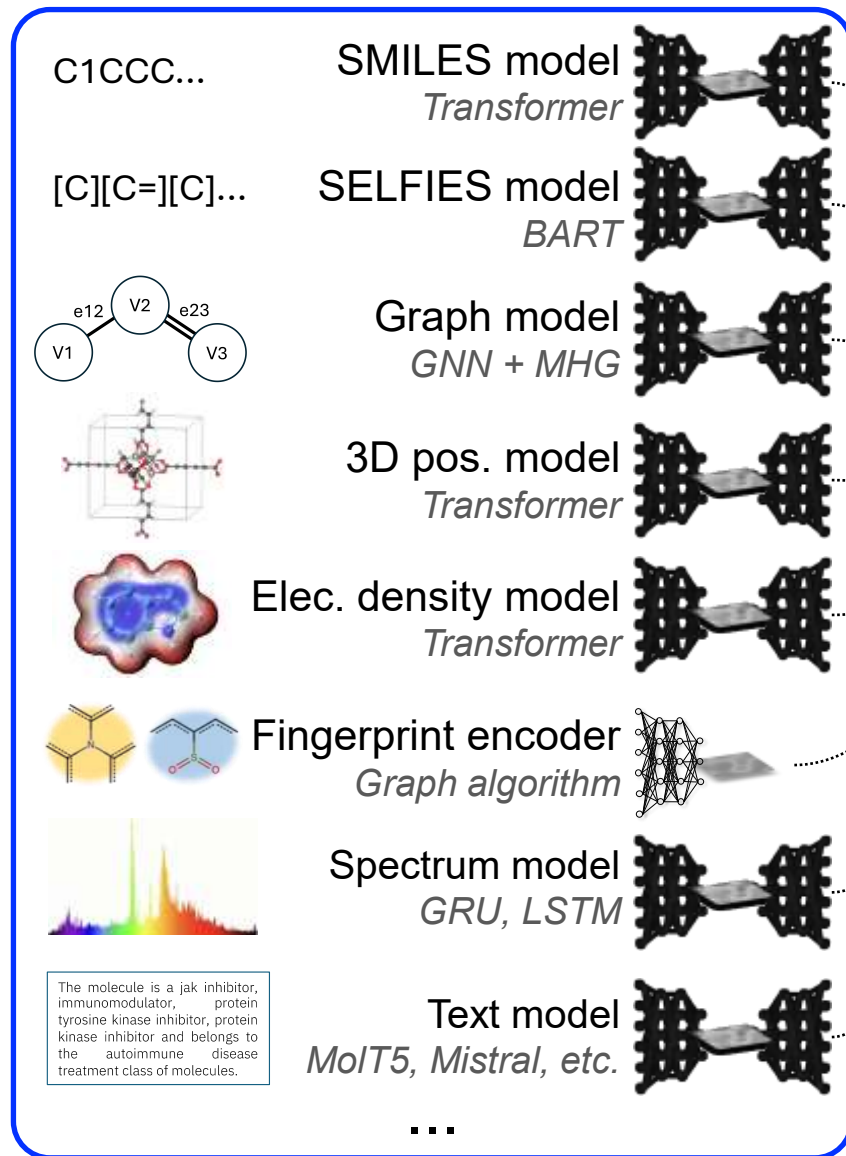| Model | Overall Increase | Translation Metric Increase | Prop. Metric Increase | BLEU-2 | BLEU-4 | ROUGE-L | METOR | Overall Prop. F1 |
|---|---|---|---|---|---|---|---|---|
| *baselines* | | | | | | | | |
| MolT5-Small | 0.00 | 0.00 | 0.00 | 70.90 | 51.20 | 54.40 | 70.10 | 7.88 |
| Meditron-7b | 13.15 | 5.50 | 15.70 | **79.20** | **57.60** | 57.50 | 75.70 | 8.93 |
| *ours* | | | | | | | | |
| Multitask T5 | **15.31** | 5.23 | **18.67** | 78.22 | 56.73 | 57.28 | 76.27 | **19.10** |
| LlaSMol_Mistral | 10.59 | 4.68 | 12.56 | 78.84 | 57.17 | 56.50 | 74.87 | 15.35 |
| Ensembled | 15.21 | **5.52** | 18.44 | 78.70 | 57.04 | **57.51** | **76.72** | 19.09 |

Overall increase from MolT5-Small baseline and translation metrics results on dev. set

# Overview of Model Architecture
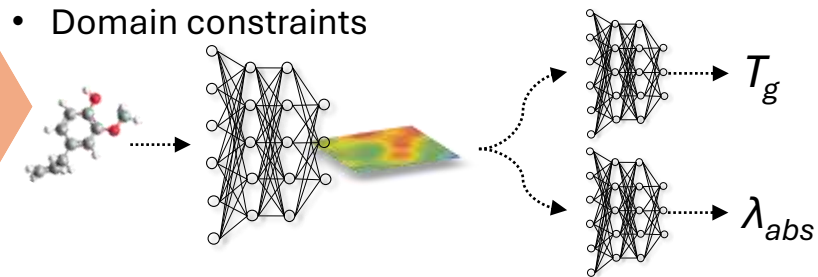


**Modality-specific Foundation Models**

C1CCC...

SMILES model
*Transformer*

[C][C=][C]...

SELFIES model
*BART*

Graph model
*GNN + MHG*

3D pos. model
*Transformer*

Elec. density model
*Transformer*

Fingerprint encoder
*Graph algorithm*

Spectrum model
*GRU, LSTM*

The molecule is a jak inhibitor, immunomodulator, protein tyrosine kinase inhibitor, protein kinase inhibitor and belongs to the autoimmune disease treatment class of molecules.

Text model
*MolT5, Mistral, etc.*

...

**Fused Foundation Model(s)**

Late fusion algorithms
• Naïve concatenation
• Mixture of Experts
• Dynamic fusion
• Contrastive learning

Encoder          Decoder

Fused representation vectors

**Downstream tasks**

Powerful representation for predictions

• Feature selection
• Domain constraints

$T_g$

$\lambda_{abs}$

Cross-modal inferences
• GFlowNet
• Contrastive learning

SMILES → → SMILES
Props → → Props
Text → → Text

# Multi-modal feature representations improve downstream prediction accuracy



Soares, Eduardo, et al. "A Multi-View approach based on Graphs and Chemical Language Foundation Model for Molecular Properties Prediction." AAAI 2024.

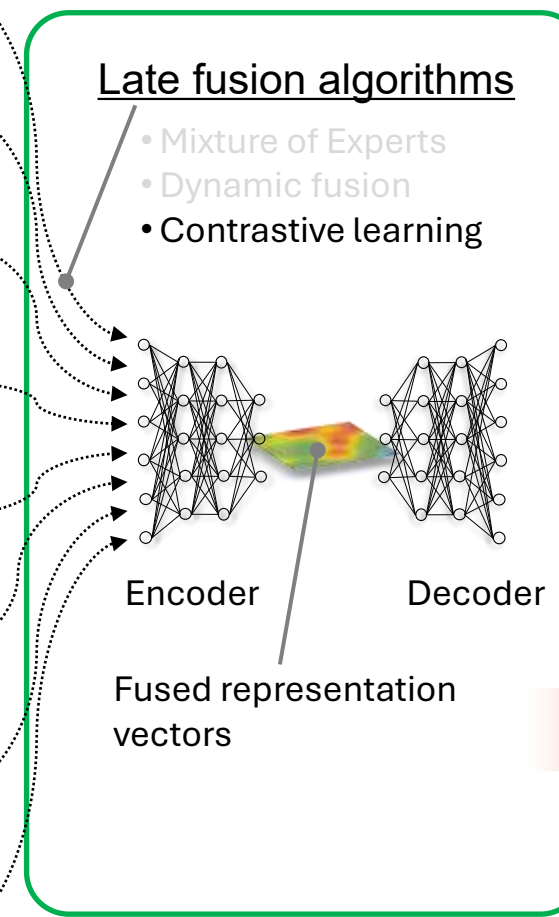| | | BBBP | HIV | BACE | ClinTox | SIDER | Tox21 |
|---|---|---|---|---|---|---|---|
| **Single** modality | SELFIES (BART) | 95.22 | 83.04 | 88.77 | 96.86 | 64.95 | 76.53 |
| | Graph (MHG) | 93.54 | 82.89 | 89.53 | 87.46 | 66.93 | 79.32 |
| | SMILES (MolFormer-XL) | 93.70 | 82.20 | 88.21 | 94.80 | **69.00** | **84.70** |
| **Multi**-modality | SELFIES + Graph | 95.35 | 83.93 | 88.66 | 92.31 | 66.02 | 78.81 |
| | SELFIES + SMILES | 96.42 | 83.99 | 89.80 | 99.82 | 64.41 | 80.48 |
| | SMILES + Graph | **96.60** | 84.93 | **90.48** | 99.59 | 65.20 | 78.15 |
| | SELFIES + SMILES + Graph | 96.06 | **85.28** | 90.00 | **99.94** | 65.65 | 79.19 |

# Overview of Model Architecture

**Modality-specific Foundation Models**

**Fused Foundation Model(s)**

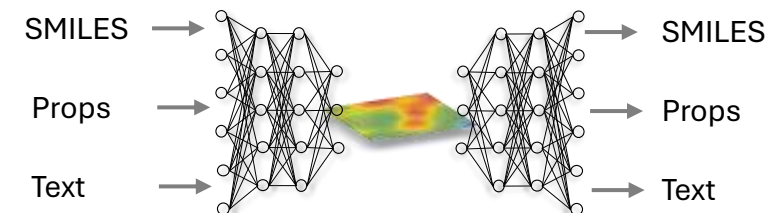**Downstream tasks**



C1CCC...

SMILES model
*Transformer*

[C][C=][C]...

SELFIES model
*BART*

Graph model
*GNN + MHG*

3D pos. model
*Transformer*

Elec. density model
*Transformer*

Fingerprint encoder
*Graph algorithm*

Spectrum model
*GRU, LSTM*

The molecule is a jak inhibitor, immunomodulator, protein tyrosine kinase inhibitor, protein kinase inhibitor and belongs to the autoimmune disease treatment class of molecules.

Text model
*MolT5, Mistral, etc.*

...

Late fusion algorithms
- Mixture of Experts
- Dynamic fusion
- **Contrastive learning**

Encoder          Decoder

Fused representation
vectors

Powerful representation for predictions
- Feature selection
- Domain constraints

$T_g$

$\lambda_{abs}$

Cross-modal inferences
- GFlowNet
- Contrastive learning

SMILES → → SMILES
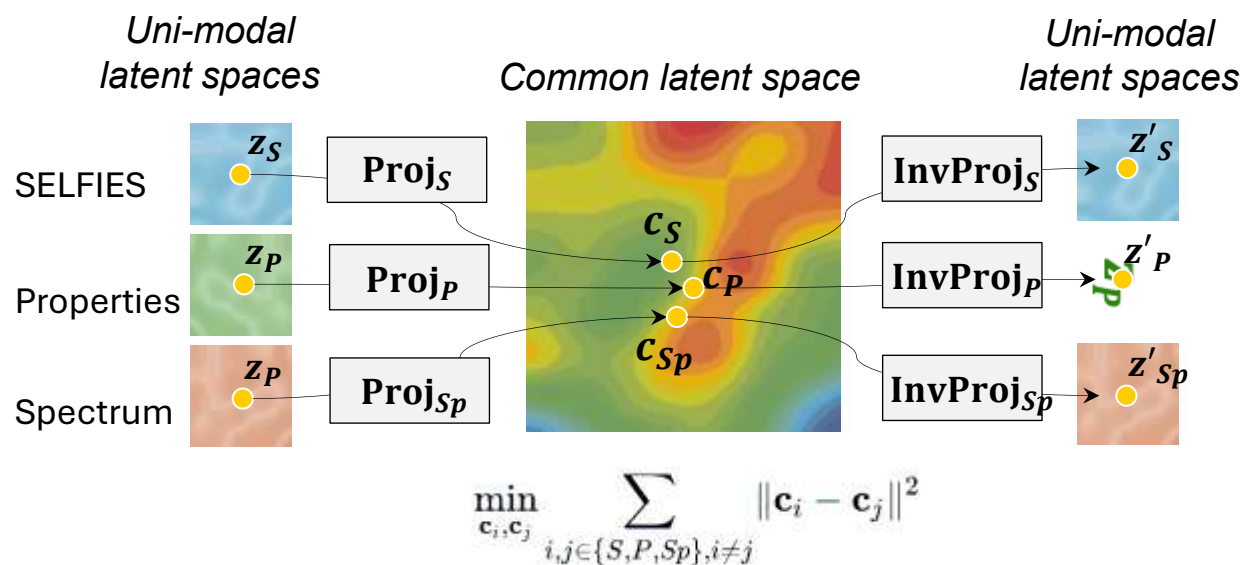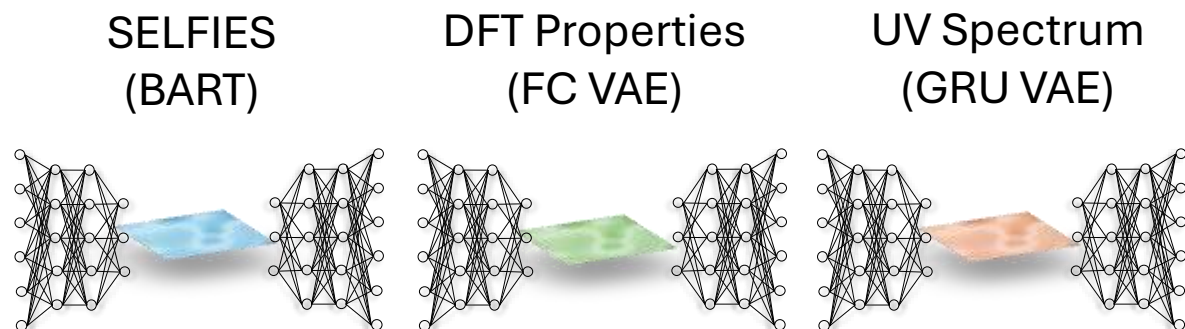
Props → → Props

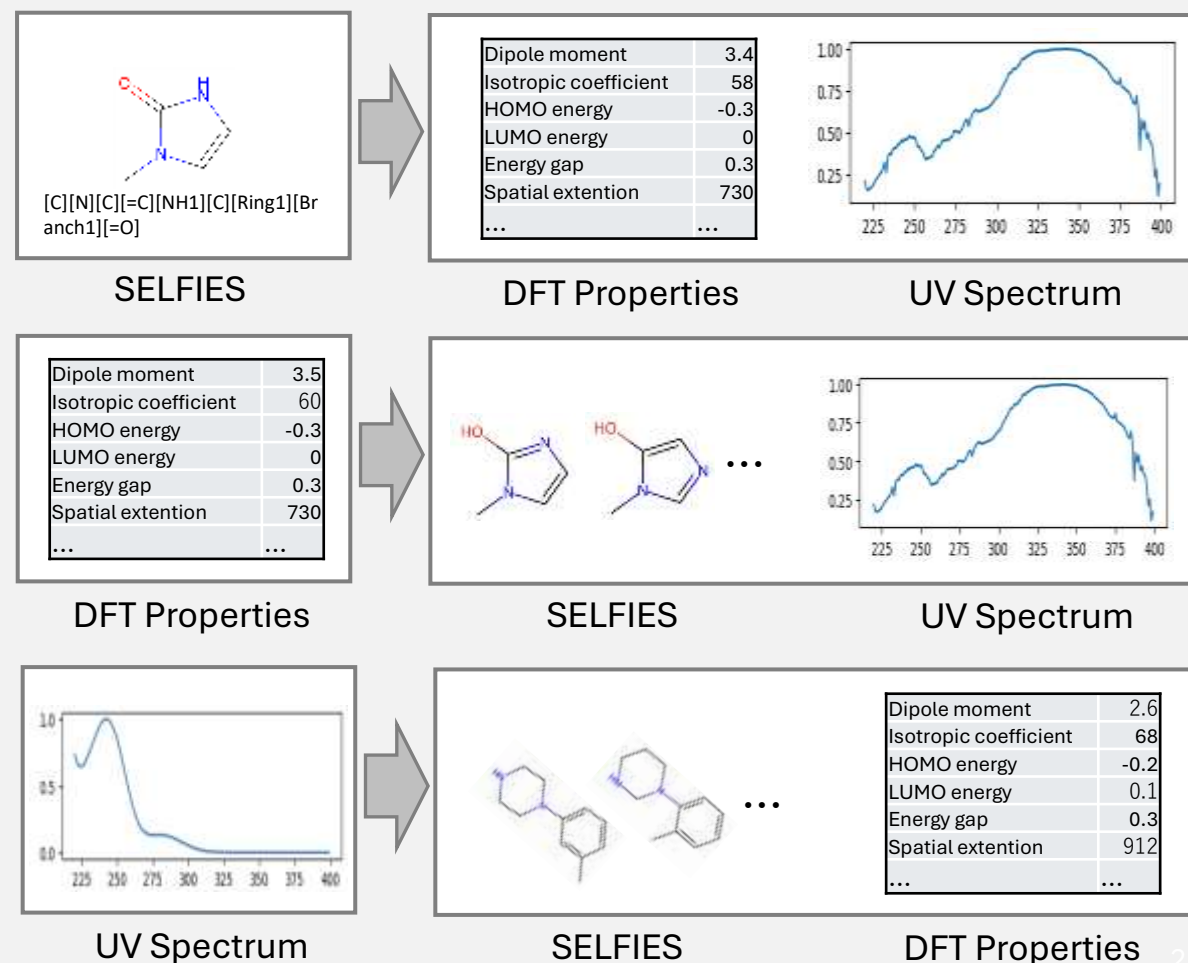Text → → Text

# Cross-Modal Inferences

Projectors/inverse-projectors are trained to align different modal representations on a common latent space so that modality-to-modality (cross-modal) inferences are achieved.

Seiji Takeda

SELFIES
(BART)

DFT Properties
(FC VAE)

UV Spectrum
(GRU VAE)

*Uni-modal latent spaces*

*Common latent space*

*Uni-modal latent spaces*

SELFIES  $z_S$  $\mathbf{Proj}_S$  $c_S$  $\mathbf{InvProj}_S$  $z'_S$

Properties  $z_P$  $\mathbf{Proj}_P$  $c_P$  $\mathbf{InvProj}_P$  $z'_P$

Spectrum  $z_P$  $\mathbf{Proj}_{Sp}$  $c_{Sp}$  $\mathbf{InvProj}_{Sp}$  $z'_{Sp}$

$$\min_{\mathbf{c}_i, \mathbf{c}_j} \sum_{i,j \in \{S,P,Sp\}, i \neq j} \|\mathbf{c}_i - \mathbf{c}_j\|^2$$

## Example of cross-modal inferences

[C][N][C][=C][NH1][C][Ring1][Branch1][=O]

**SELFIES**

| | |
|---|---|
| Dipole moment | 3.4 |
| Isotropic coefficient | 58 |
| HOMO energy | -0.3 |
| LUMO energy | 0 |
| Energy gap | 0.3 |
| Spatial extention | 730 |
| ... | ... |

**DFT Properties**

**UV Spectrum**

| | |
|---|---|
| Dipole moment | 3.5 |
| Isotropic coefficient | 60 |
| HOMO energy | -0.3 |
| LUMO energy | 0 |
| Energy gap | 0.3 |
| Spatial extention | 730 |
| ... | ... |

**DFT Properties**

**SELFIES**

**UV Spectrum**

**UV Spectrum**

**SELFIES**

| | |
|---|---|
| Dipole moment | 2.6 |
| Isotropic coefficient | 68 |
| HOMO energy | -0.2 |
| LUMO energy | 0.1 |
| Energy gap | 0.3 |
| Spatial extention | 912 |
| ... | ... |

**DFT Properties**

# Battery Materials Discovery Empowered by AI and Foundation Models
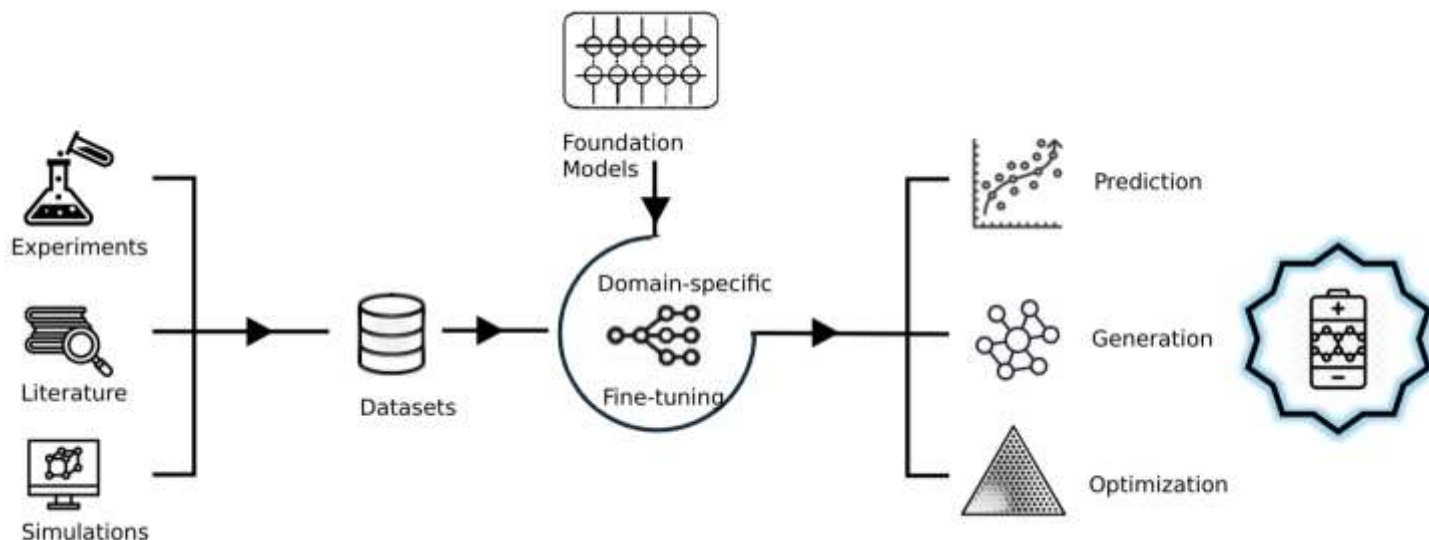
Young-Hye Na



**Battery Materials**
- Large design space and multivariable relationships
- Complexities at multiple-length scale

✓ **Foundation Models (FM) are fine-tuned for domain specific tasks** using labeled datasets derived from literature, simulations, or lab-experimentation

✓ **Our fine-tuned FM models effectively map material structures, compositions, and device performance,** predicting and optimizing the properties and performance of complex mixed materials (e.g., electrolyte formulations)

✓ **These models, along with our customized AI-workflows,** are currently being used to discover new electrolytes for our industrial partners
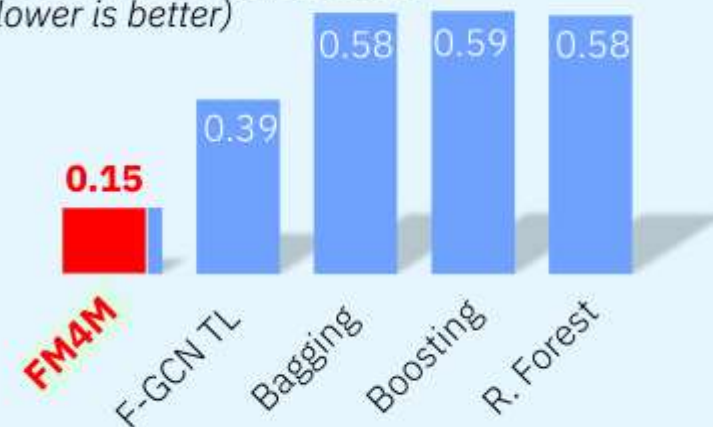
## AI-assisted Electrolyte Discovery Workflow



## Prediction of Battery Performance (LCE)

Indrapriyadarsini et.al, "Improving Performance Prediction of Electrolyte Formulations with Transformer-based Molecular Representation Model." ML4LMS @ *ICML*. 2024.



IBM **Research/ Sustainable Battery Materials Discovery/ Young-Hye Na: yna@us.ibm.com**
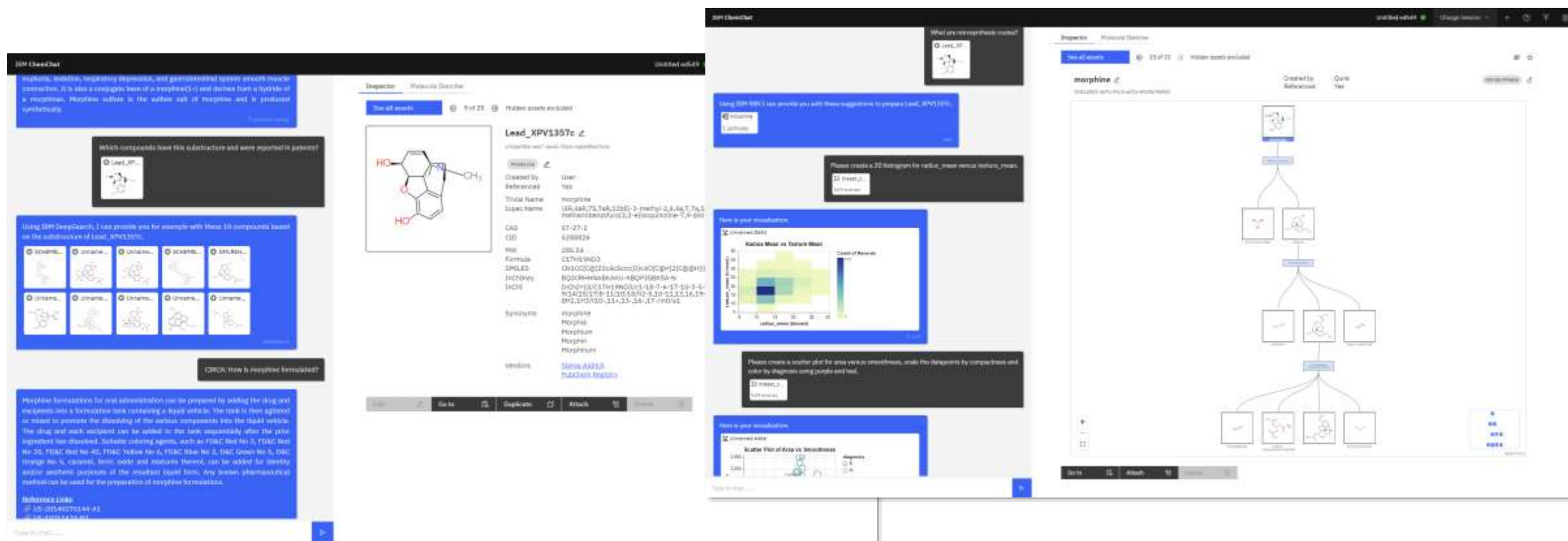
# ChemChat Material Science and Data Visualization Assistant

**Tim Erdmann**

Molecules:
Identification & description, Vendor check, Property calculations & predictions by cheminformatics, conventional AI and FM models, PFAS determination after EPA & ECHA, Retrosynthesis prediction, Generation by target properties, Locating molecules and substructure matches in patents

Documents:
Querying 30M material and chemistry-focused patents, Querying user-defined collection of pdfs

Visualizations:
11 chart types incl. request for aggregations and interactive elements

# The AI Alliance & Open Innovation

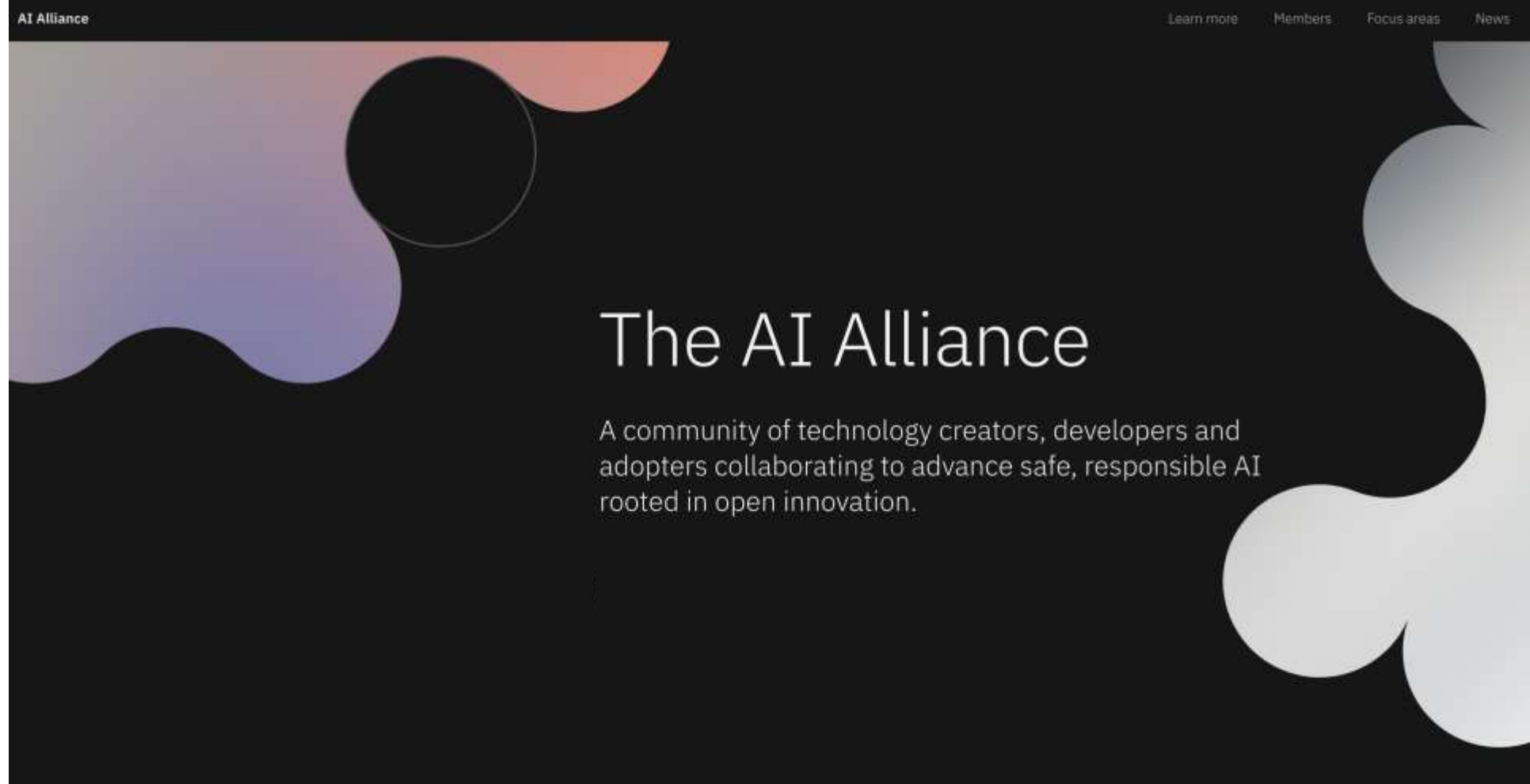**Build** and support open technologies across software, models and tools.

**Enable** developers and scientists to understand, experiment, and adopt open technologies.

**Advocate** for open innovation with organizational and societal leaders, policy and regulatory bodies, and the public.



https://thealliance.ai/

AI Alliance     Learn more   Members   Focus areas   News

# The AI Alliance

A community of technology creators, developers and adopters collaborating to advance safe, responsible AI rooted in open innovation.

## WG4M – Working group for Materials and Chemistry

- Community development of open-source code and models
  - Defining target domains/use-cases (PFAS, polymer etc)
  - Rich model family of modalities and architectures
- Open-source of base models at https://github.com/IBM/materials (one model released, more planned)
- Sharing: model architecture and contact points for members, and workstream-based engagement approach
- 20+ materials company and worldwide researchers involved

# AI Alliance Working Group for Materials & Chemistry (WG4M)

**(1) Kick-off Technical Workshop**
(16th May 2024, Tokyo)
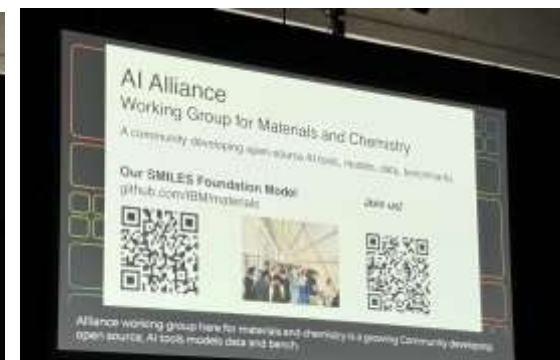


Keynote Talks & Lightning Pitches

Group discussion & Readout

**(2) Social Gathering @ ICML AI4Science**
(25th July 2024, Vienna)

# AI Alliance Working Group for Materials & Chemistry (WG4M)

**Community Development**

- 20+ materials companies
- Bi-weekly calls with material companies
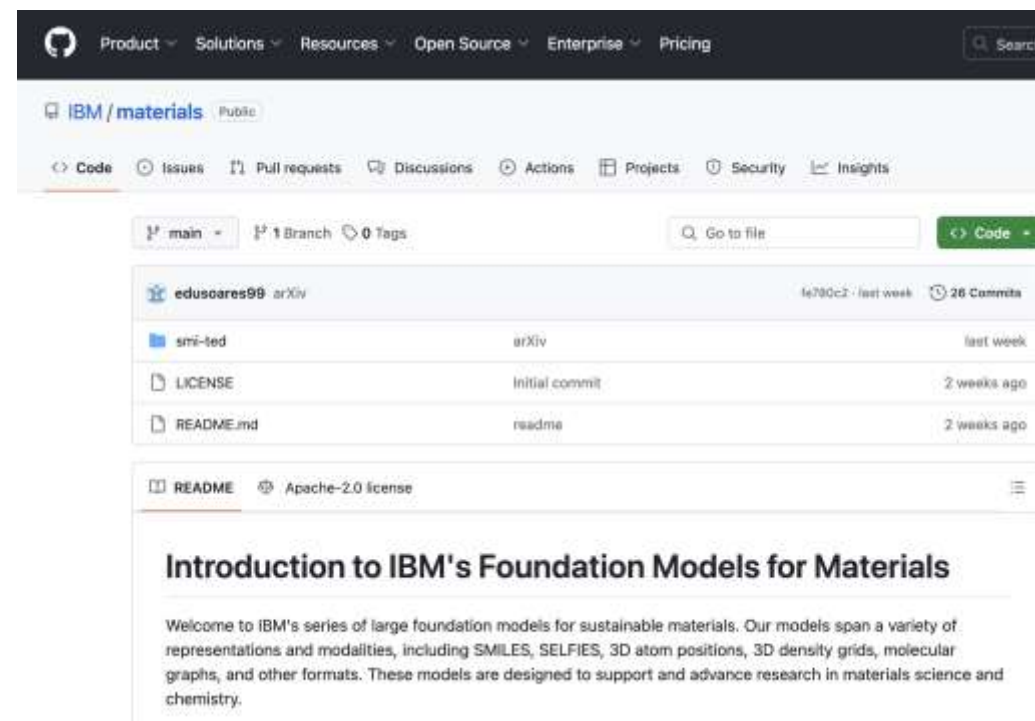- Seeking co-chairs, topic leads



**Model development**

Open-source of base models at github.com/IBM/materials

- SMILES Transformer model was released on 25 July
- Other modalities models (SELFIES, Graph, etc.) will soon be released!



https://thealliance.ai/

# Thank you

https://www.linkedin.com/in/indra-ipd/

# Backup

# Working Group for Materials (WG4M)

**Vision**

To build an open development community for Materials Informatics (MI) that transcends national and organizational boundaries, and to accelerate the speed of global materials development by creating a suite of MI tools and models that are accessible to everyone.

1. Open development and consumption of MI AI tools and models

2. Data sharing (to the extent possible)

3. Creation of benchmarks

4. Provision of platforms for information exchange and discussion

5. Development of talent in Materials x AI