

Accelerating Materials Discovery via Physics-Informed & Data-Driven Constraints

Ju Sun

McKnight Land-Grant Professor
Computer Science & Engineering

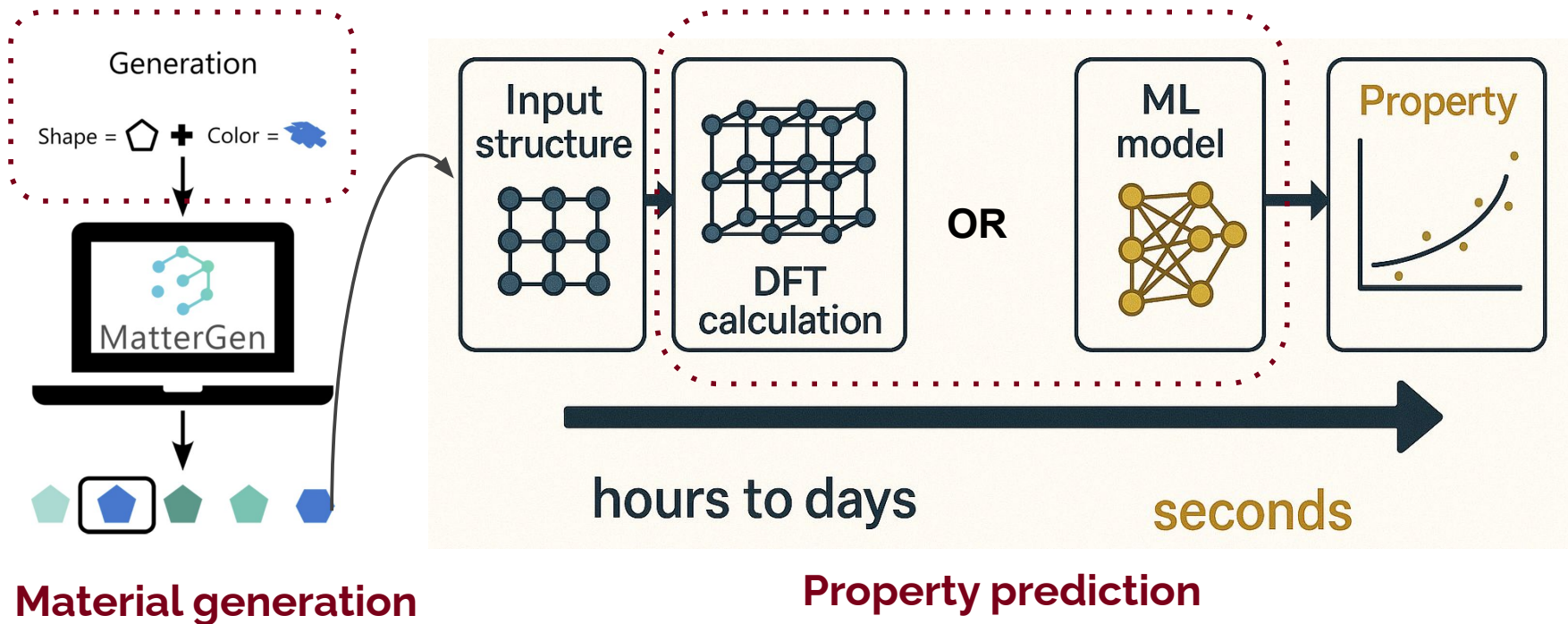
May 15, 2026



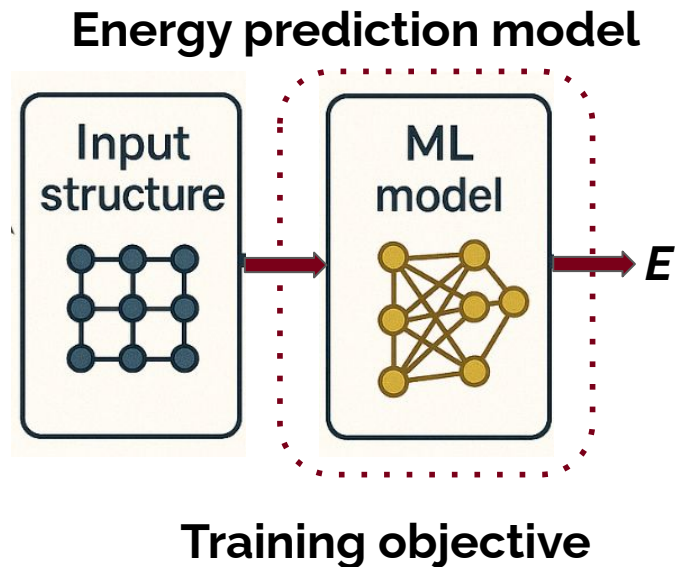
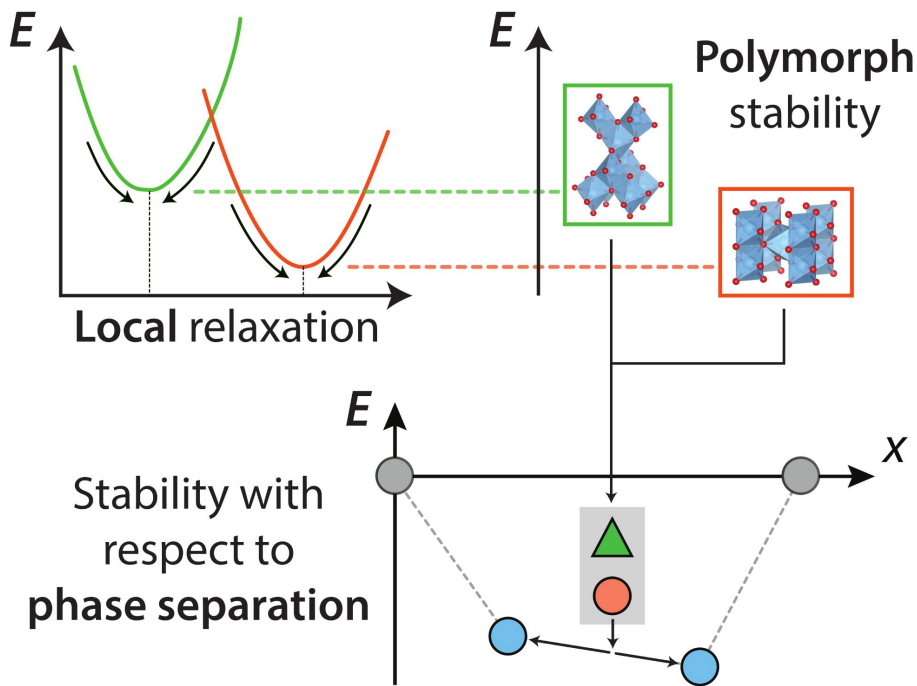
UNIVERSITY OF MINNESOTA

Driven to DiscoverSM

AI-accelerated materials discovery



Property prediction: thermodynamic stability

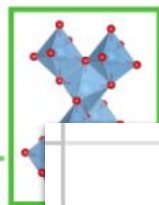


$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N \ell(\mathbf{y}_i, f_{\theta}(\mathbf{x}_i)) + \Omega(\theta)$$

... with physics-informed constraints

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N \ell(\mathbf{y}_i, f_{\theta}(\mathbf{x}_i)) + \Omega(\theta)$$

E



Polymorph
stability

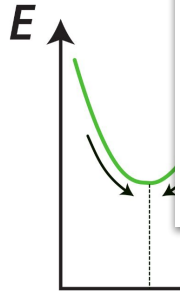


$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N \ell(\mathbf{y}_i, f_{\theta}(\mathbf{x}_i)) + \Omega(\theta)$$

$$\text{s. t. } f_{\theta}(\mathbf{x}_p) \geq f_{\theta}(\mathbf{x}_q) \quad \forall (p, q) \in \mathcal{O}_{\text{W-CSTR-T}},$$

	Test			
Model	F1 Score	L1 Score	Ground State Score	Kendall Tau
MACE	0.7873	0.0317	0.006	0.5966
MACE + CDL	0.82	0.0327	0.005	0.6474
Improvement	0.0327	0.001	-0.001	0.0508

Next: more powerful constraints



- **(W-CSTR-T) Weak explicit constraints informed by phase transition:** For each composition, the energy/atom of any non-ground-state polymorph is greater than that of the ground-state polymorph
- **(W-CSTR-S) Weak explicit constraints informed by phase separation:** For each chemical space, the energy/atom of any TUS material is above the lower convex envelope in the composition-energy space

Local relaxation

E

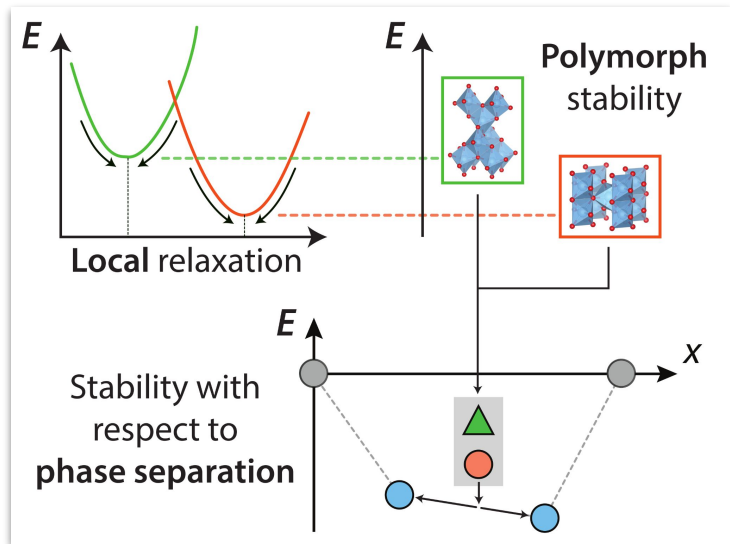
Stability with respect to phase separation

- **(S-CSTR-T) Strong explicit constraints informed by phase transition:** For each composition, the ordering of the energy/atom of all polymorphs is preserved
- **(S-CSTR-S) Strong explicit constraints informed by phase separation:** For each chemical space, the ordering of energy/atom of all materials is preserved, and the energy/atom of any TUS material is above the lower convex envelope in the composition-energy space



ACED: Accelerating Computing-Enabled Scientific Discovery (ACED)

Deep learning with nontrivial “hard” constraints



$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N \ell(\mathbf{y}_i, f_{\theta}(\mathbf{x}_i)) + \Omega(\theta)$$

$$\text{s. t. } f_{\theta}(\mathbf{x}_p) \geq f_{\theta}(\mathbf{x}_q) \quad \forall (p, q) \in \mathcal{O}_{\text{W-CSTR-T}},$$

Robustness evaluation of DL models

$$\begin{aligned} & \max_{\mathbf{x}'} \ell(\mathbf{y}, f_{\theta}(\mathbf{x}')) \\ & \text{s. t. } \mathbf{x}' \in \Delta(\mathbf{x}) = \{\mathbf{x}' \in [0, 1]^n : d(\mathbf{x}, \mathbf{x}') \leq \varepsilon\} \end{aligned}$$

$$\min_{\mathbf{x}' \in [0, 1]^n} d(\mathbf{x}, \mathbf{x}') \quad \text{s. t. } \max_{\ell \neq y} f_{\theta}^{\ell}(\mathbf{x}') \geq f_{\theta}^y(\mathbf{x}')$$

Imbalance Learning

$$\max_{\theta, t} \frac{\sum_{i=1}^N \mathbb{1}\{y_i = +1\} \mathbb{1}\{f_{\theta}(\mathbf{x}_i) > t\}}{\sum_{i=1}^N \mathbb{1}\{f_{\theta}(\mathbf{x}_i) > t\}}$$

$$\text{s. t. } \frac{\sum_{i=1}^N \mathbb{1}\{y_i = +1\} \mathbb{1}\{f_{\theta}(\mathbf{x}_i) > t\}}{\sum_{i=1}^N \mathbb{1}\{y_i = +1\}} \geq \alpha$$

Heuristic handling of hard constraints is problematic

Method	cross-entropy loss		margin loss	
	Viol. (%) ↓	Att. Succ. (%) ↑	Viol. (%) ↓	Att. Succ. (%) ↑
Fast-LPA	73.8	3.54	41.6	56.8
LPA	0.00	80.5	0.00	97.0
PPGD	5.44	25.5	0.00	38.5
PWCF (ours)	0.62	93.6	0.00	100

LPA, Fast-LPA: penalty methods

PPGD: Projected gradient descent

Penalty methods tend to encounter **large constraint violation** (i.e., infeasible solution in theory) or **suboptimal solution**

**Unreliable computation
leads to
unreliable
reliability evaluation**

NCVX: A principled solver for constrained DL



NCVX PyGRANSO
Documentation

Search the docs ...

Introduction

Installation

Settings

Examples

Home

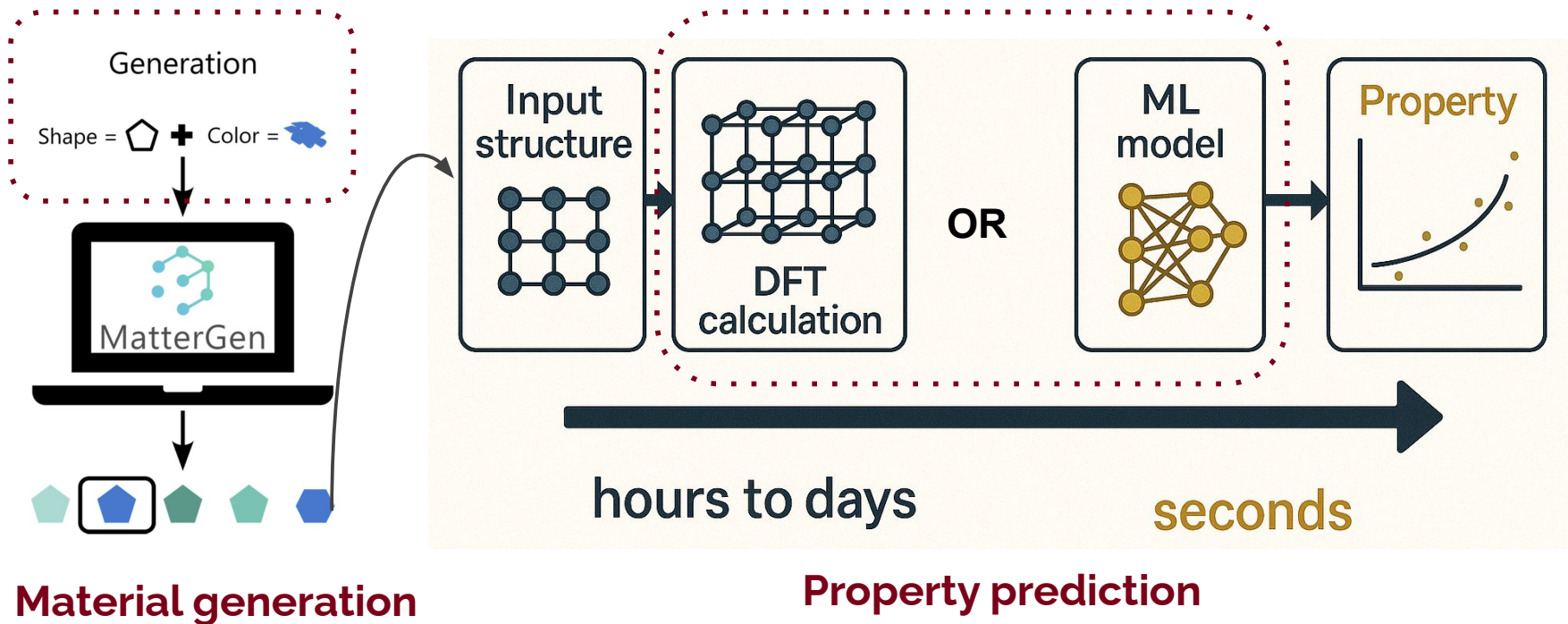
<https://ncvx.org/>



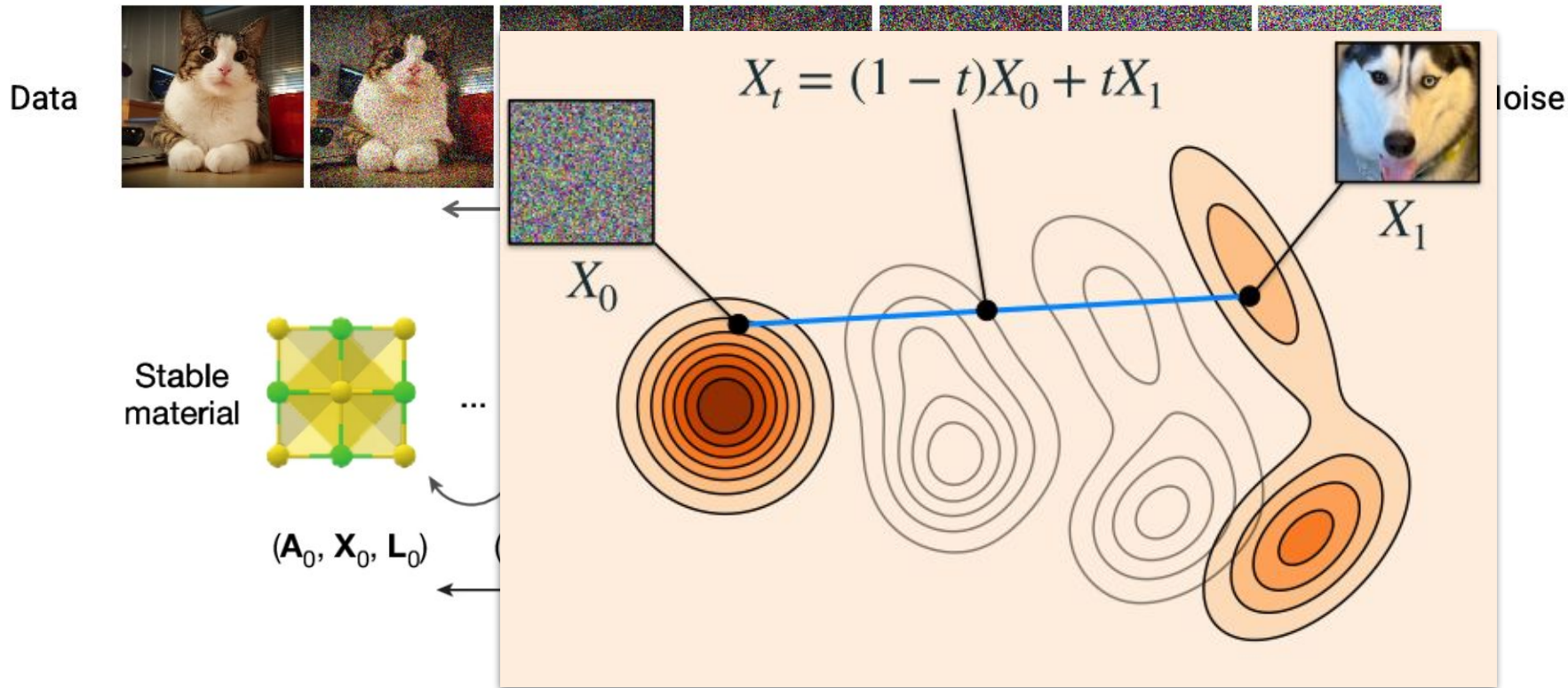
$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}), \text{ s.t. } c_i(\mathbf{x}) \leq 0, \forall i \in \mathcal{I}; c_i(\mathbf{x}) = 0, \forall i \in \mathcal{E}$$

- **First general-purpose solver for hard-constrained deep learning problems**
- **Recently updated to be compatible with PyTorch 2.8**

AI-accelerated materials discovery



Generative models



Data bottlenecks for numerous domains

IMAGES

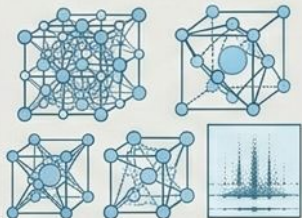


LAION-

Size: 5.8

Key Detail: for SOTA models; English

CRYSTAL STRUCTURES



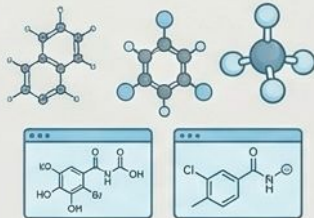
MATERIALS PROJECT / OQMD

Size: >150k unique compounds

Key Detail: Density functional theory (DFT) calculated properties; foundation for new crystal discovery.

VIDEOS

ORGANIC MOLECULES



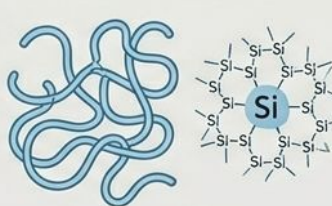
ZINC20 / PubChem

Size: >1.5 Billion distinct conformers

Key Detail: Synthesizable chemical space; pre-training for generative molecular graph and SMILES models.

3D OBJECTS

POLYMERS & AMORPHOUS

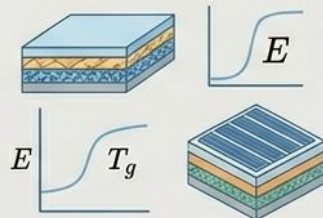


KHASHKHOV / Specialized Sets

Size: >10k annotated entries

Key Detail: Processed from literature & experiments; focuses on complex chain architectures and

POLYMERS & AMORPHOUS

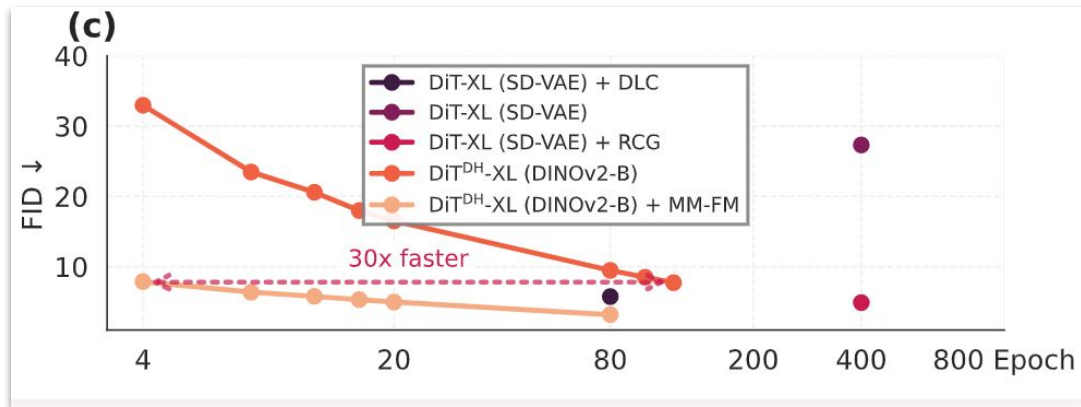
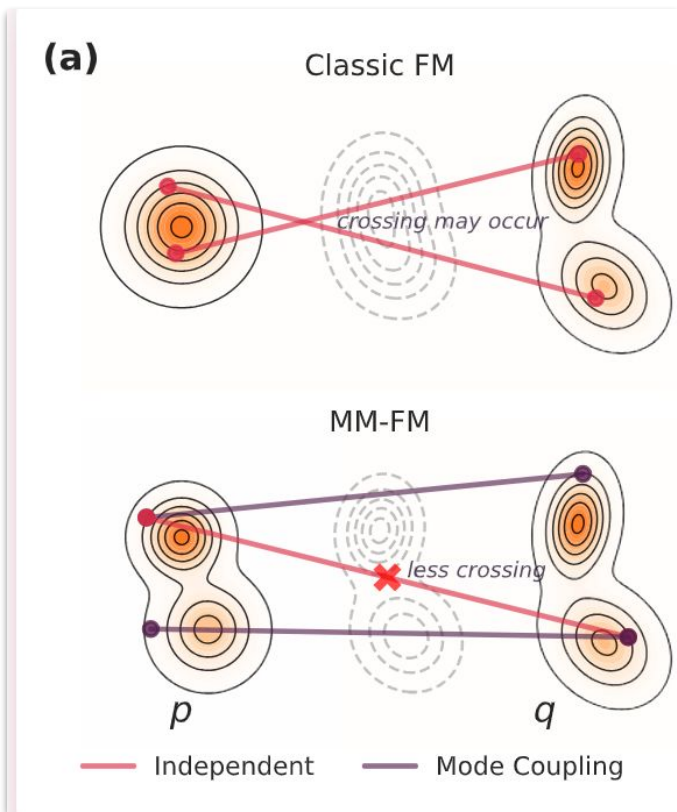


KHASHKHOV / Specialized Sets

Size: >10k annotated entries

Key Detail: Processed from literature & experiments; focuses on chain architectures and composite design.

Maximize data-efficient for generative models



Better
data/training/inference efficiencies

For more details

Flow Matching for Multimodal Distributions

Gaoxiang Luo^{1,*} , Frank Cole^{2,*} , Sihang Zhang³ , Yuxiang Wan¹ , Yulong Lu² , Ju Sun¹

*Equal Contribution ¹Department of Computer Science and Engineering

²School of Mathematics ³School of Statistics

University of Minnesota Twin Cities

CVPR 2026



Paper (Coming soon)



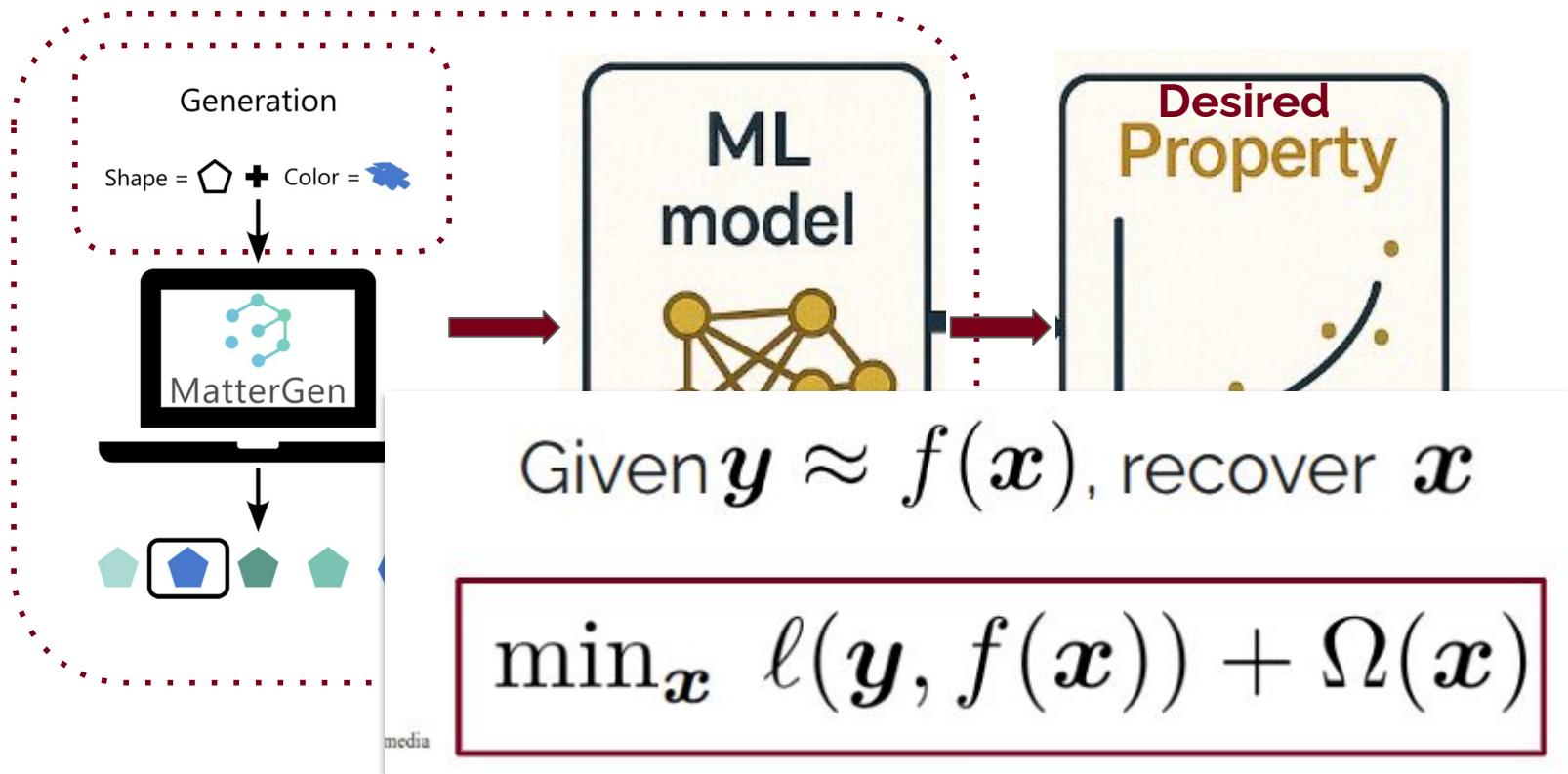
Code (Coming soon)



Models (Coming soon)

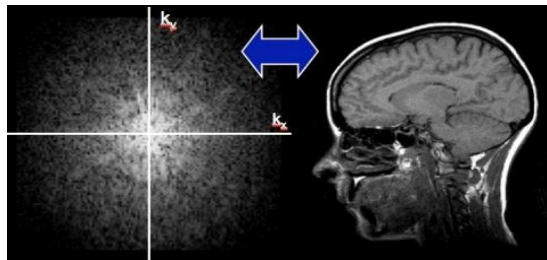
<https://mm-flow.github.io/>

Materials discovery as inverse design

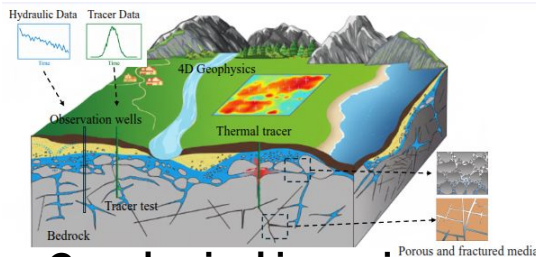


Inverse modeling

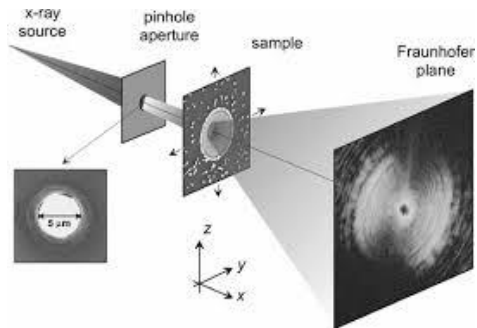
(inverse problems, reconstruction, data assimilation, inverse design/control, conditional generation, ...)



MRI reconstruction



Geophysical inversion



Coherent diffraction imaging



Inverse shape control

Given $\mathbf{y} \approx f(\mathbf{x})$, recover \mathbf{x}

$$\min_{\mathbf{x}} \ell(\mathbf{y}, f(\mathbf{x})) + \Omega(\mathbf{x})$$

Challenges:

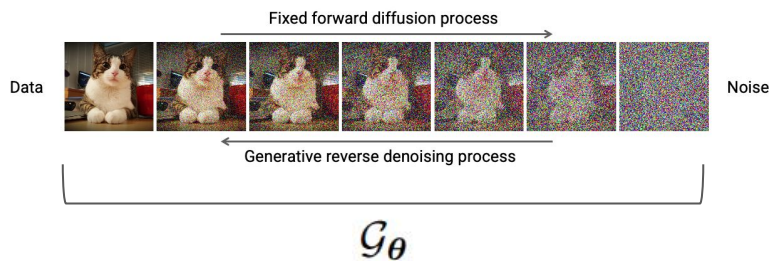
- Linear vs. **nonlinear** f
- Unconstrained vs. **constrained** (e.g., PDE $(x, y) = 0$)
- Explicit vs. **implicit** f

Inverse modeling with generative priors

object-only datasets $\{\mathbf{x}_i\}_{i=1,\dots,N}$



Distribution learning via deep generative models



Plug in pretrained deep generative priors

$$\min_z \mathcal{L}(z) \doteq \ell(\mathbf{y}, \mathcal{A} \circ \mathcal{G}_\theta(z)) + \Omega \circ \mathcal{G}_\theta(z)$$

DMPlug: A Plug-in Method for Solving Inverse Problems with Diffusion Models

Hengkang Wang, Xu Zhang, Taihui Li, Yuxiang Wan, Tiancong Chen, Ju Sun

NeurIPS'24; with **domain-specific** priors

Saving Foundation Flow-Matching Priors for Inverse Problems

Yuxiang Wan, Ryan Devera, Wenjie Zhang, Ju Sun

ICML'26; with **foundation** priors (e.g., stable diffusion models)