

Synthetic Geology: Structural Geology Meets Deep Learning

Paper + Code



Simon Ghyselincks^{1,2}
George Turkiyyah³

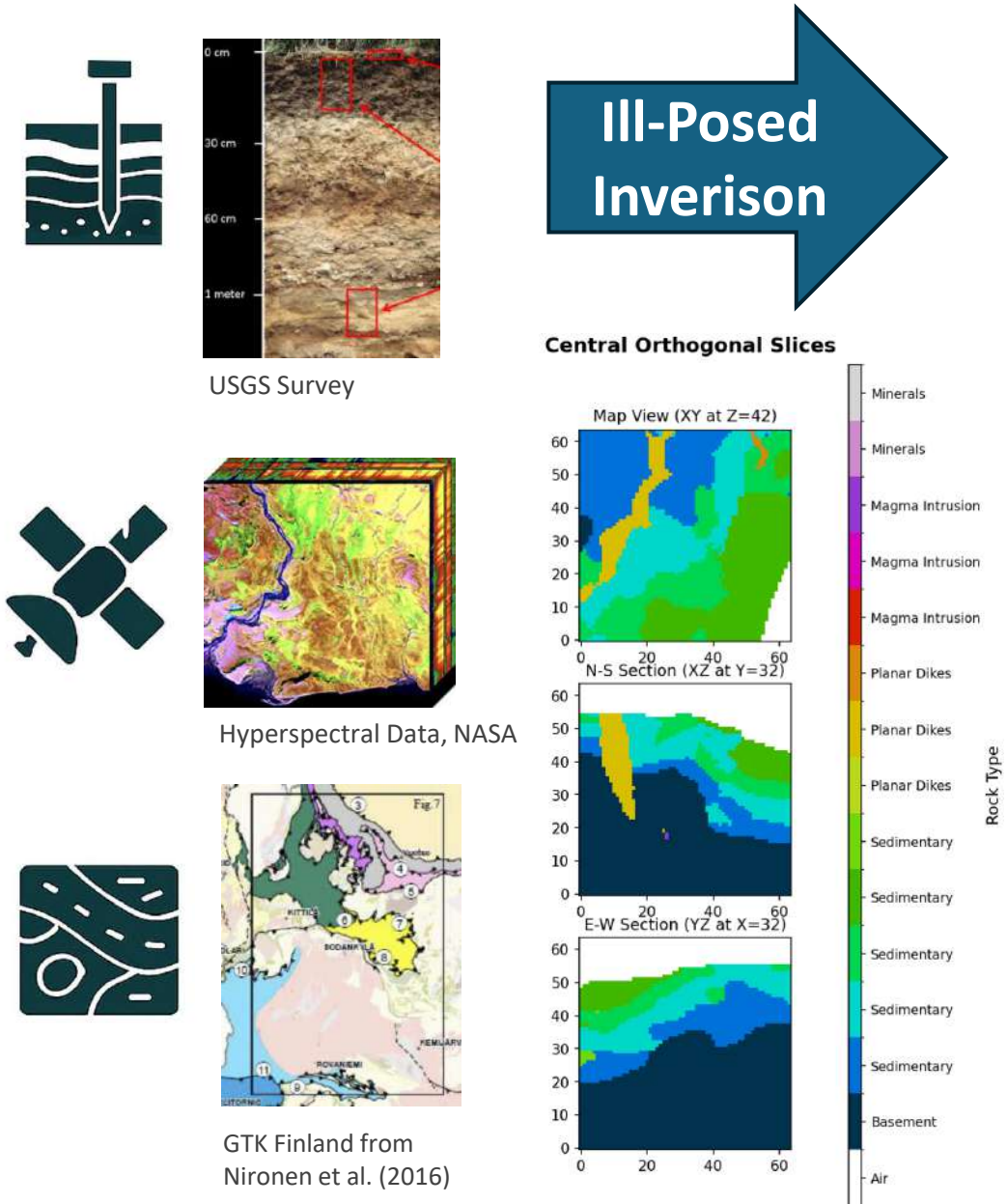
Valeriia Okhmak³
David Keyes³

Stefano Zampini³
Eldad Haber¹

¹University of British Columbia ²Vector Institute ³KAUST

A Hidden Earth

Geological observed data is highly heterogenous, multi-modal, and incomplete. There is no known **ground truth** for **structural geology** models.



This framework deals with **discrete categorical rock types** to make probabilistic reconstructions of the earth from observed data.

Inverse Problems

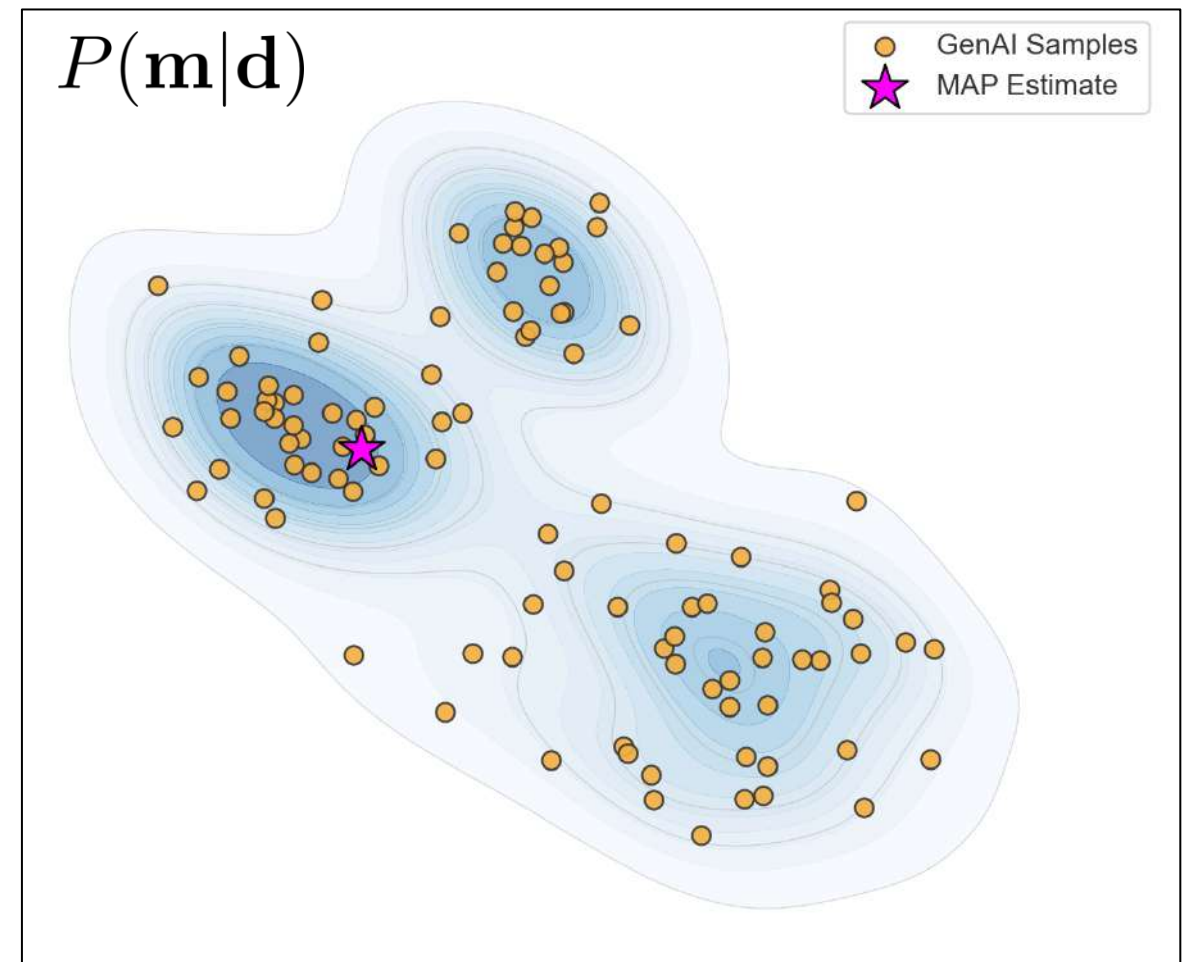
Classical Regularized Inversion:
Seek a single maximum likelihood solution using Maximum a Posteriori
 $\mathbf{A}\mathbf{m} + \epsilon = \mathbf{d}$

$$P(\mathbf{m}|\mathbf{d}) = \frac{P(\mathbf{d}|\mathbf{m})P(\mathbf{m})}{P(\mathbf{d})}$$

$$\hat{\mathbf{m}}_{\text{MAP}} = \arg \max_{\mathbf{m}} [P(\mathbf{d}|\mathbf{m})P(\mathbf{m})]$$

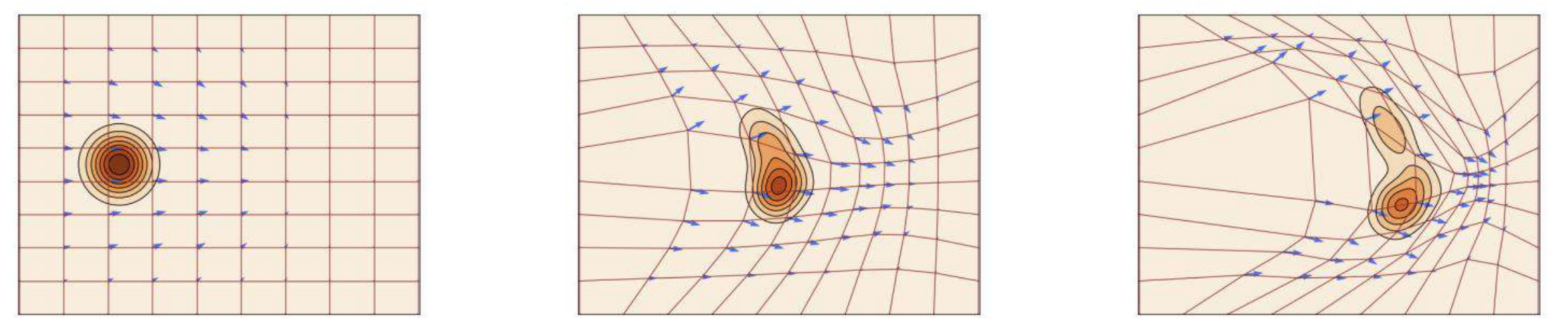
$$\phi(\mathbf{m}) = \frac{1}{2} \|\mathbf{A}\mathbf{m} - \mathbf{d}\|^2 + \alpha R(\mathbf{m})$$

Misfit Prior



Posterior Sampling with Gen AI

Direct learning is intractable due to the partition function $P(\mathbf{m}|\mathbf{d}) = \frac{1}{Z(\mathbf{d})} e^{-E(\mathbf{m}|\mathbf{d})}$, where $Z(\mathbf{d}) = \int_{\Omega} e^{-E(\mathbf{m}|\mathbf{d})} d\mathbf{m}$

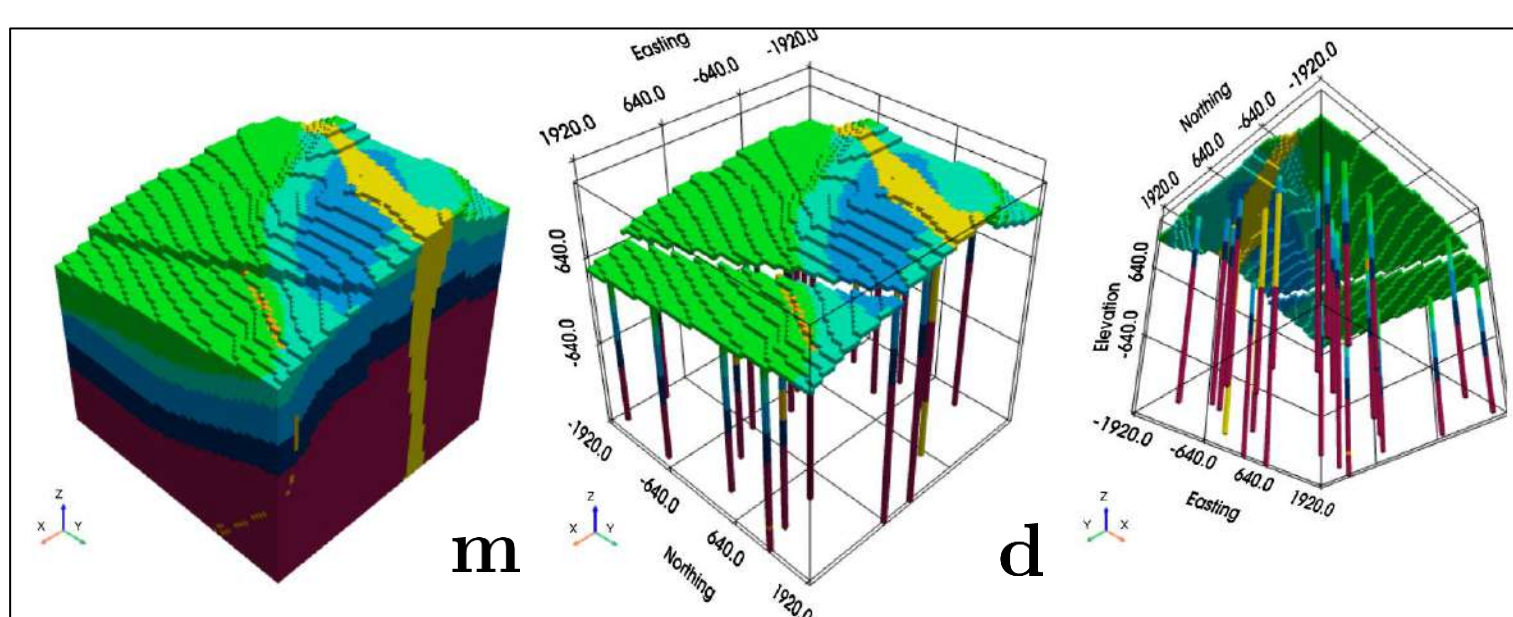
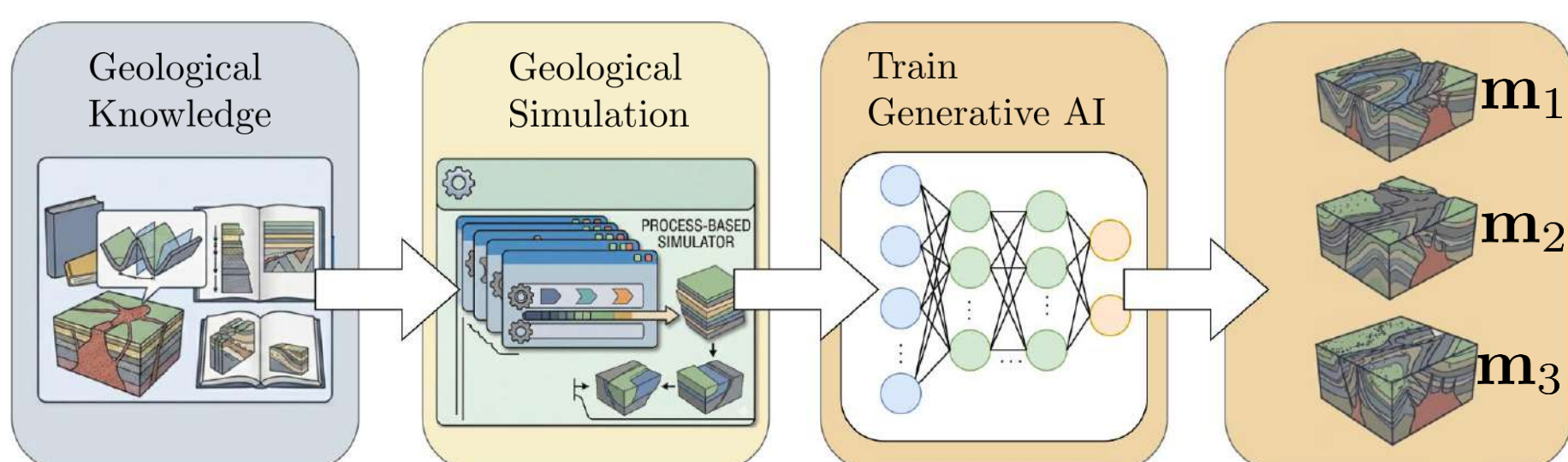


In **conditional flow matching**, the conditional probability distribution over all models is indirectly learned via a time dependent velocity field $\hat{\mathbf{v}}(\mathbf{m}, t; \xi)$ with model params ξ . Discrete data must first be embedded into a continuous space.

$$\frac{d\mathbf{m}}{dt} = \hat{\mathbf{v}}(\mathbf{m}_t^{(e)}, t; \xi) \quad \mathbf{m}_0^{(e)} \sim N(0, \mathbf{I}) \quad t \in [0, 1]$$

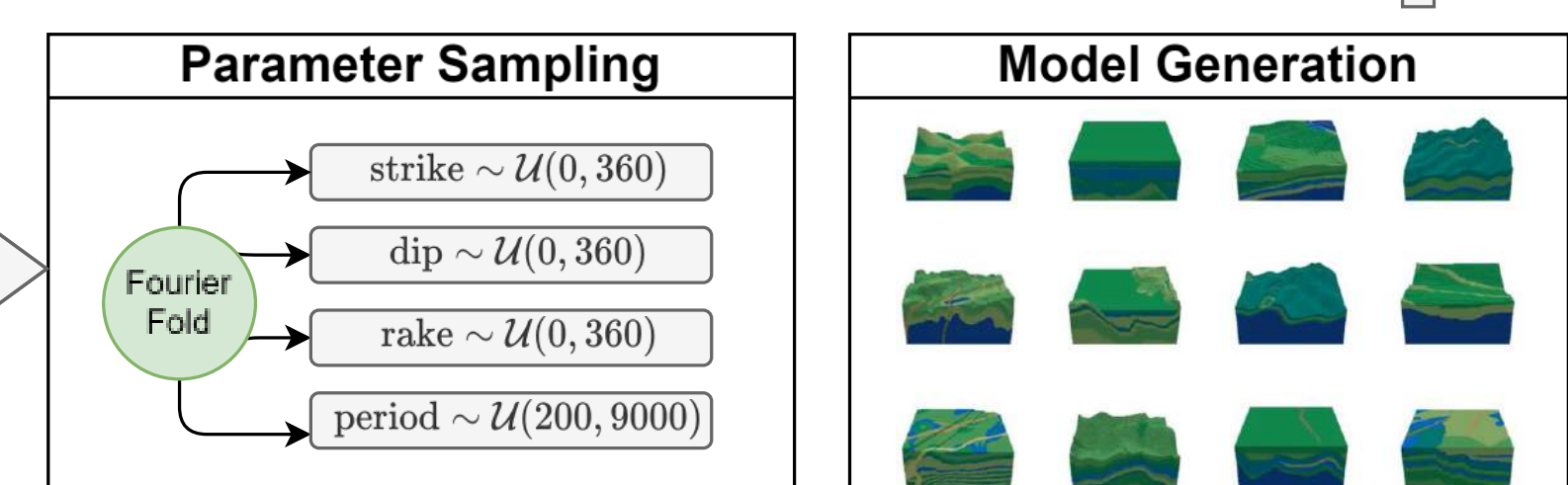
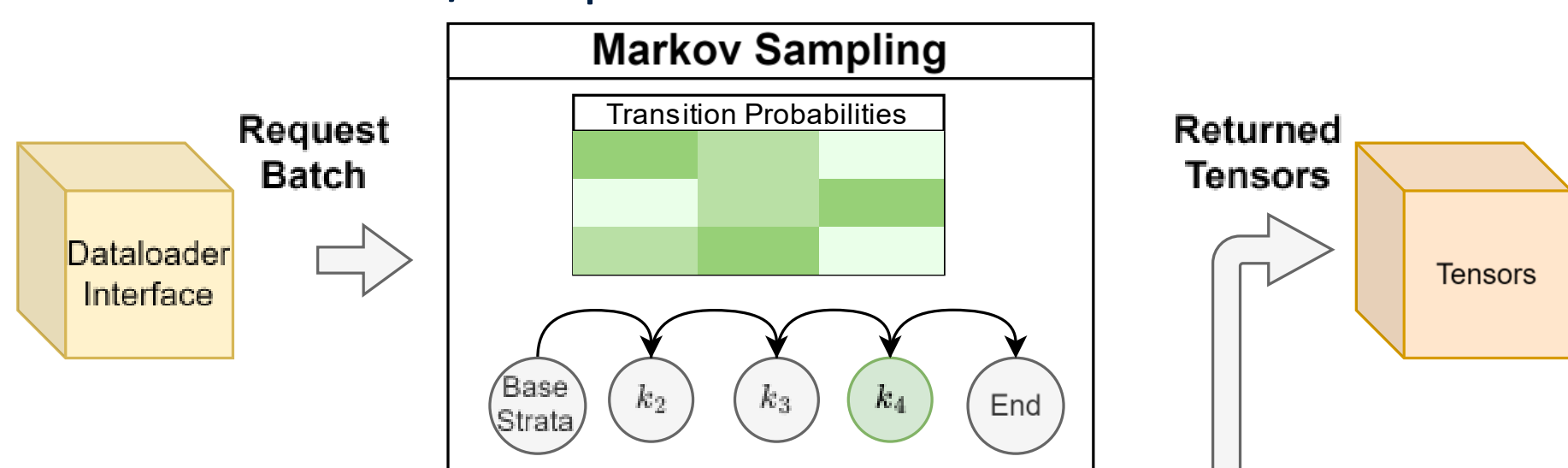
Data Driven Geology

A **simulator** encodes **geological knowledge** into complete ground truth samples and observables to **train a generative model** that outputs plausible geology

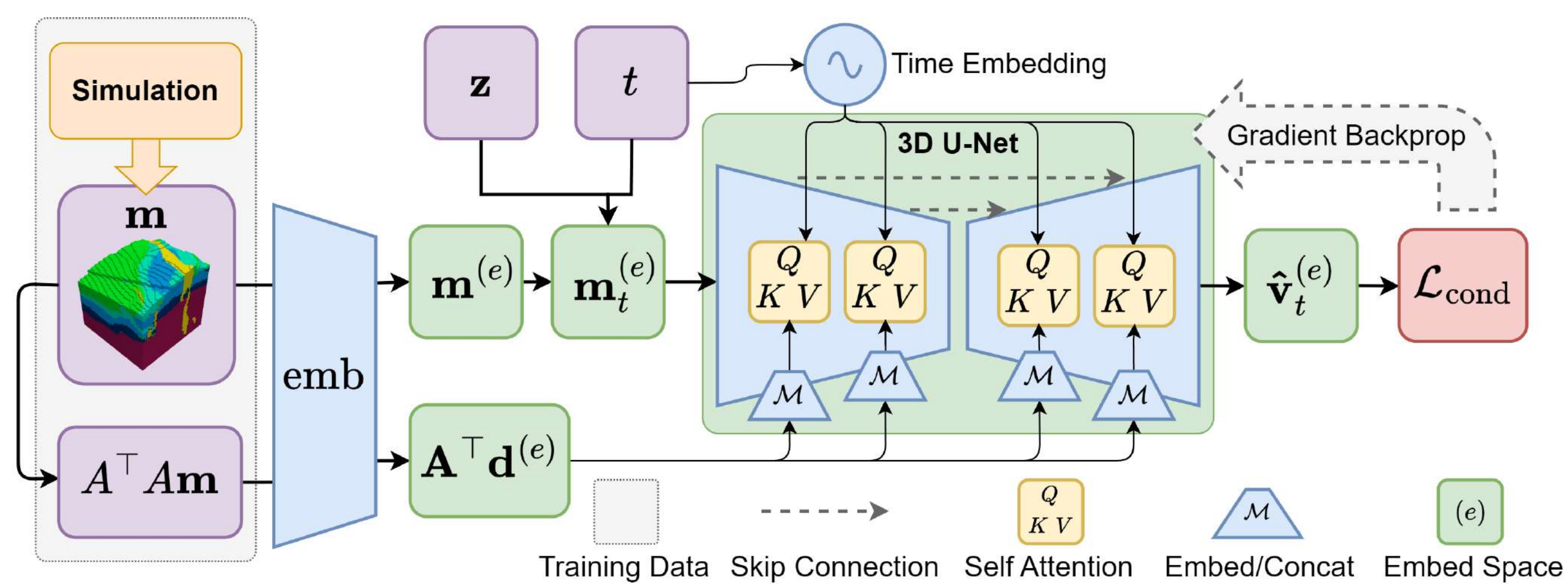


Synthetic Data Pipeline

- Unlimited sampling
- Zero disk storage
- Customizable / adaptable



Training Architecture



$$\mathcal{L}_{\text{cond}} = \|\hat{\mathbf{v}}(\mathbf{m}_t, \mathbf{d}, t; \xi) - \mathbf{v}_t\|_2^2 + \lambda \cdot t \cdot \|\hat{\mathbf{d}} - \mathbf{d}\|_2^2$$

Results and Applications

