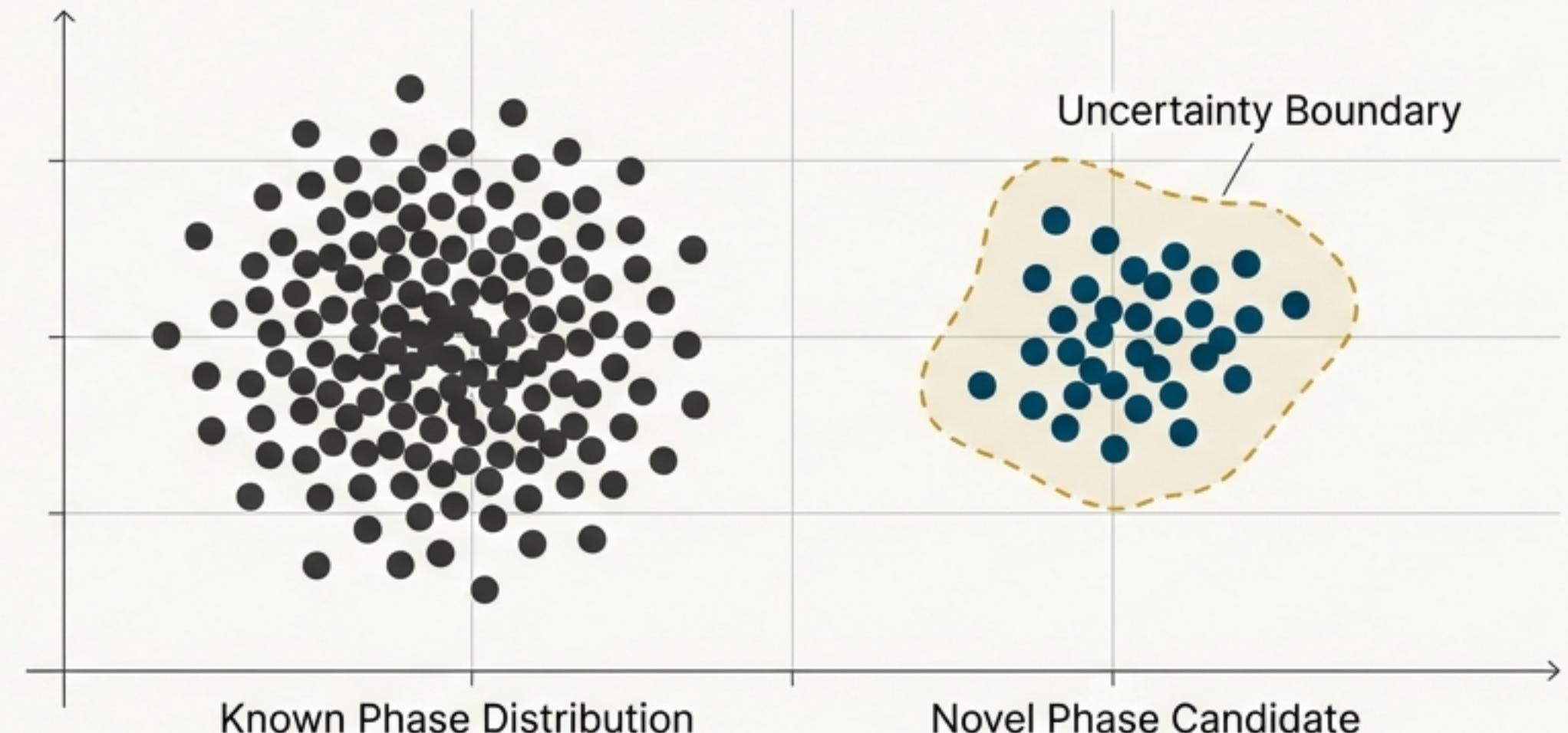


Physics-Informed Generative Models for the Autonomous Discovery of Novel Phases of Matter

A research proposal for the BSc Computing with AI programme.



Proposed by: Alex J. Sutton

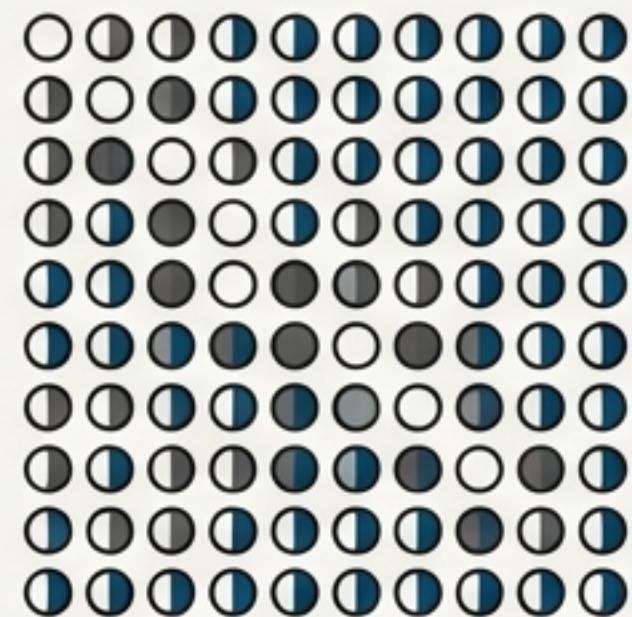
Supervised by: Dr. Eleanor Vance

This project introduces a framework to automatically identify both known and unknown phases of matter by combining generative machine learning, fundamental physics principles, and robust uncertainty estimation.

The Study of Phase Transitions is Central to Modern Physics

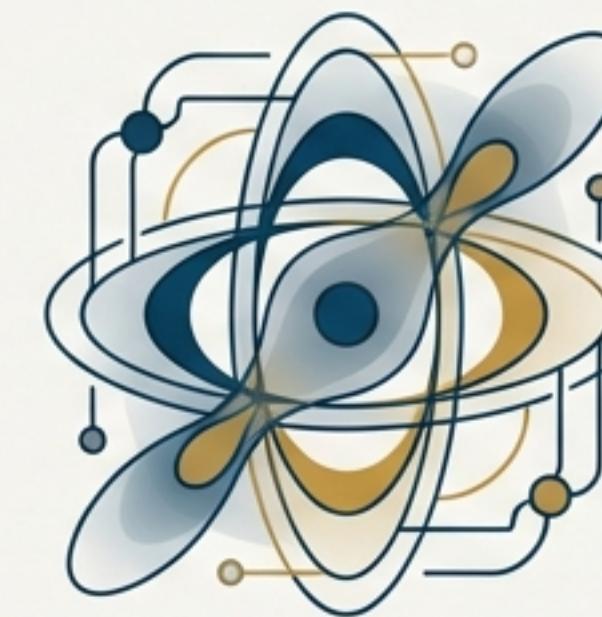
Understanding phase transitions is fundamental to describing the behaviour of materials, from simple magnets to complex quantum systems.

Classical Foundation



The 2D Ising model serves as a canonical example, with its well-defined critical temperature ($T_c \approx 2.269$) separating ordered and disordered phases.

Modern Frontier

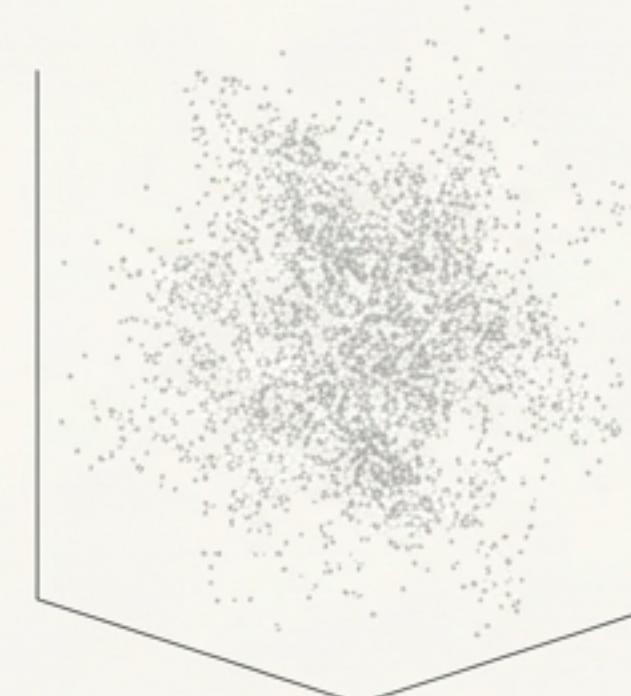


However, many contemporary materials, such as quantum magnets and topological insulators, exhibit non-trivial transitions that are far more difficult to characterise and discover.

Generative Models Are Becoming a Powerful Tool for Analysing Physical Systems

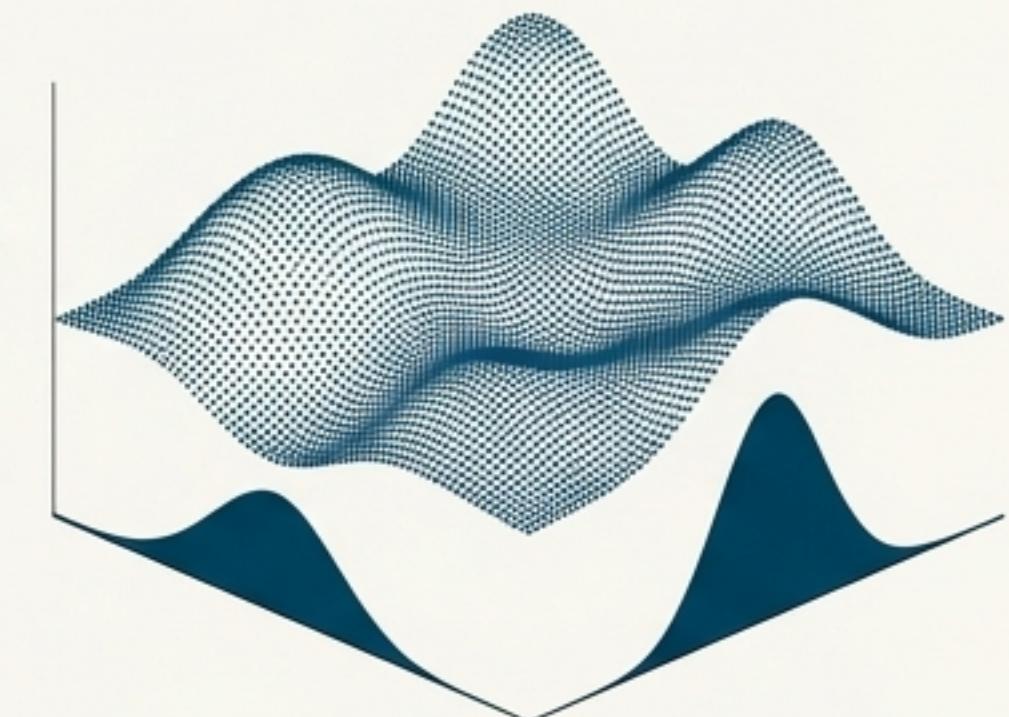
The Challenge of Dimensionality

High-dimensional data from simulations makes manual analysis of complex systems intractable.



The Power of Generative Models

Unlike classifiers that only learn labels, generative models learn the full data distribution, offering a richer understanding of the underlying physics.



MIT (2024)

Successfully used diffusion models to identify phase boundaries in an unsupervised manner (without labels).

Arnold et al. (2024)

Demonstrated that generative models can map known classical phases with over 95% accuracy.

Key Takeaway: These recent breakthroughs prove the principle, but a crucial capability remains elusive: the *discovery* of entirely unknown phases.

A Critical Gap Exists Between Characterising Known Phases and Discovering New Ones



1. Reliance on a priori Knowledge:

Most models are trained via supervised learning to predict known phases, or use simple clustering (PCA, t-SNE) that only detects clear boundaries.



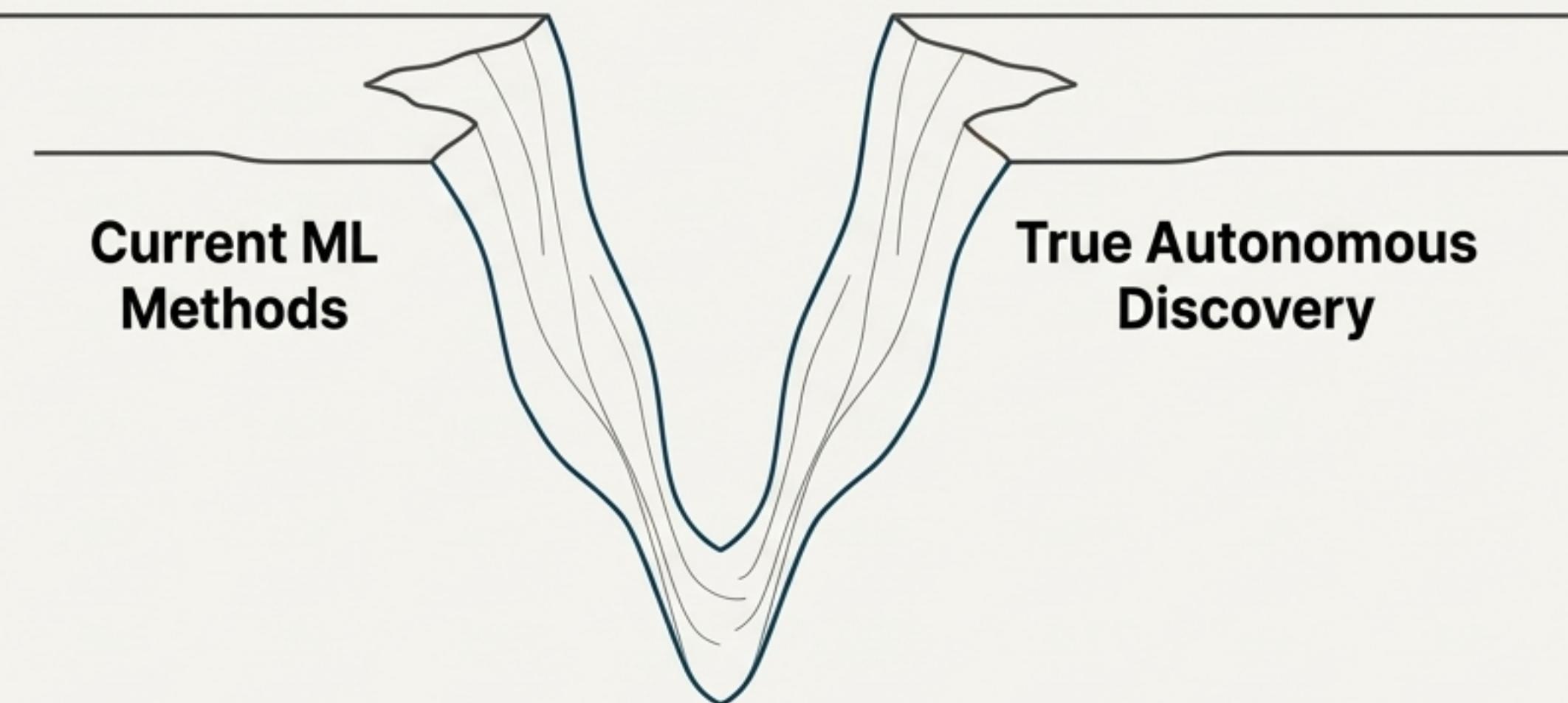
2. Lack of Uncertainty Quantification:

Models provide predictions but no reliable measure of confidence, making it impossible to know when they are encountering data outside their training distribution.



3. Absence of Physical Constraints:

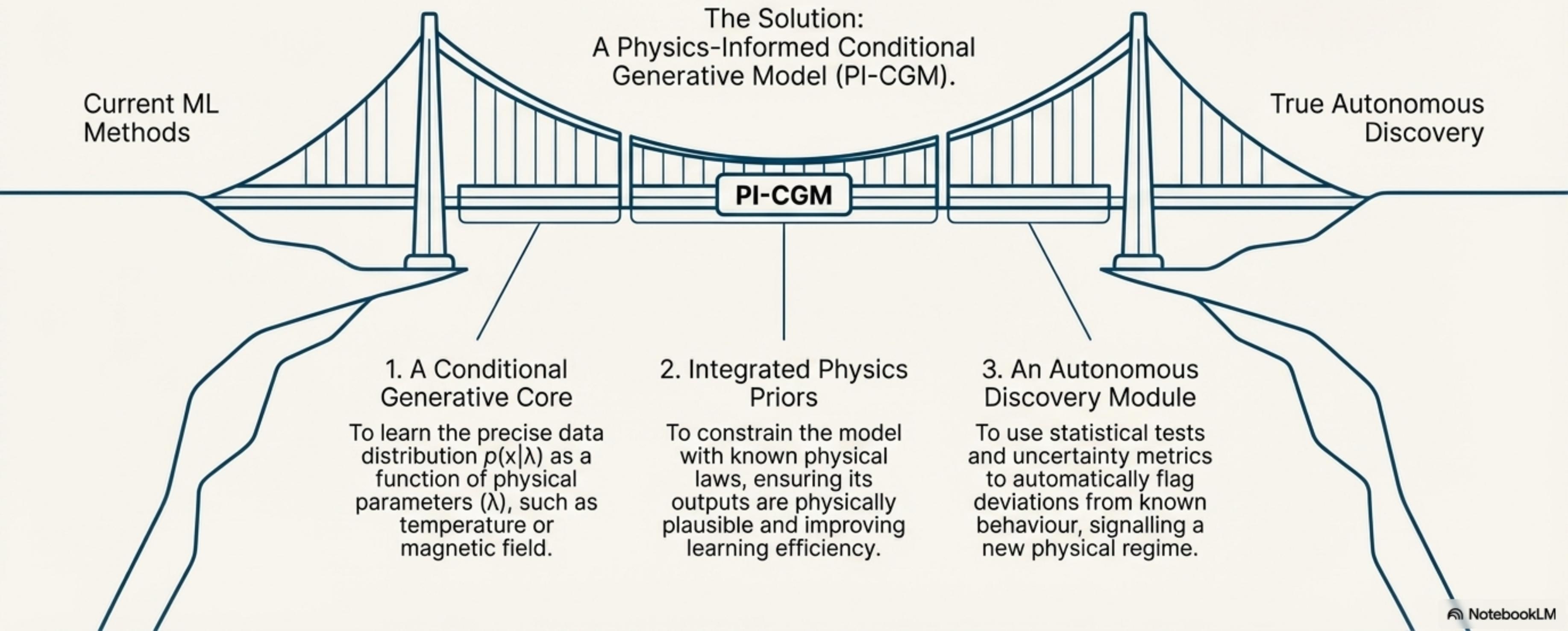
They often operate as 'black boxes,' failing to integrate fundamental physical laws (e.g., energy conservation, symmetries) into the learning process.



The Missing Piece

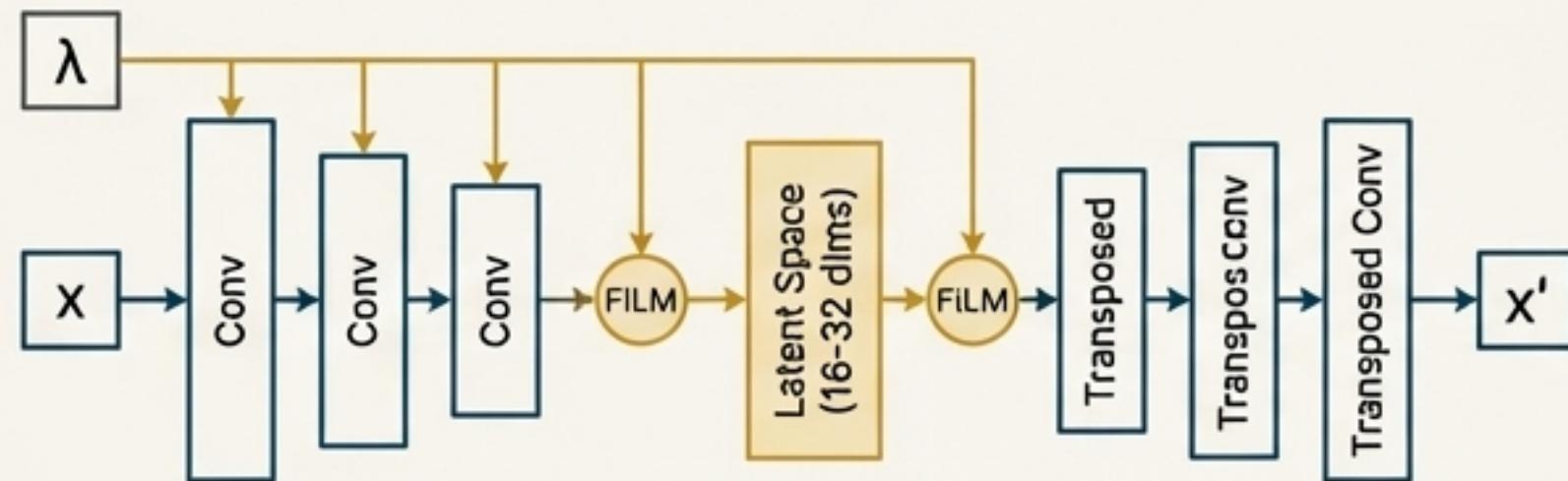
A system that models the probability density of physical data and can statistically identify when a new sample significantly deviates, signalling a potential new phase.

We Propose a Physics-Informed Generative Model to Bridge This Gap



The Blueprint: A Choice of Two Powerful Generative Architectures

Option 1: Conditional Variational Autoencoder (VAE)

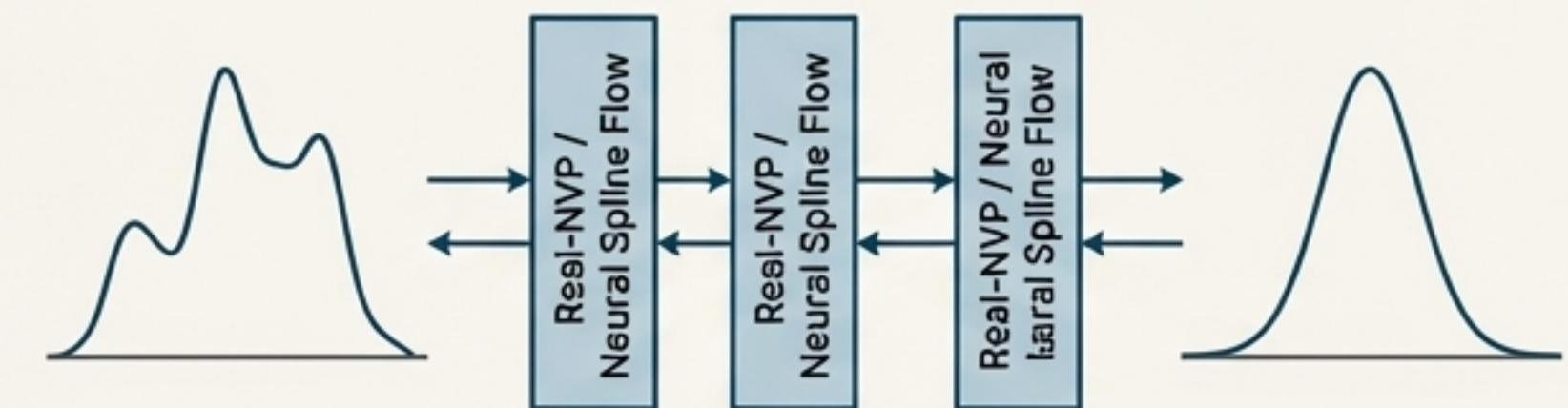


Architecture: Encoder (3 Conv layers + MLP) \rightarrow Latent Space (16-32 dims) \rightarrow Decoder (Transposed Conv).

Conditioning: Physical parameters (λ) are injected via Feature-wise Linear Modulation (FiLM) layers.

- ✓ Computationally efficient and excellent for learning compressed latent representations.

Option 2: Normalizing Flow



Architecture: Real-NVP or Neural Spline Flow.

- ✓ Provides exact log-likelihood estimation (not an approximation like the VAE's ELBO), which is ideal for the statistical tests in our discovery module.

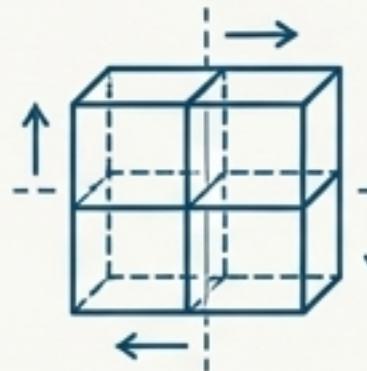
Proven Performance: Glow-like models have achieved Negative Log-Likelihood (NLL) scores of <0.02 on Ising configurations, demonstrating their power in this domain.

The Blueprint: Integrating Physical Laws and an Autonomous Discovery Engine

Integrating Physics Priors

Symmetry Constraints

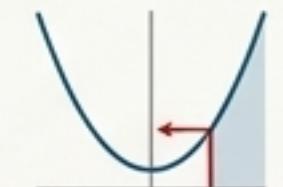
Model architecture will be designed to respect known symmetries (e.g., translational, reflection) of the physical system.



Energy-Consistency Loss

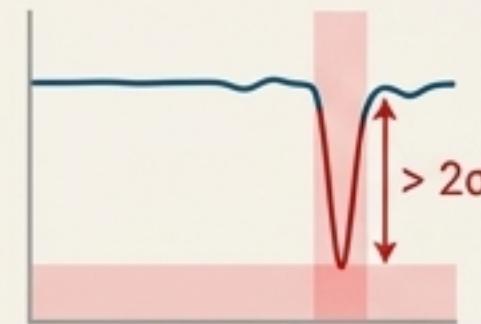
A mini-loss term is added to the main objective function to penalise physically inconsistent predictions:

$$L_{\text{phys}} = \|E_{\text{pred}} - E_{\text{true}}\|^2 \text{ with } \lambda_{\text{phys}} = 0.2$$



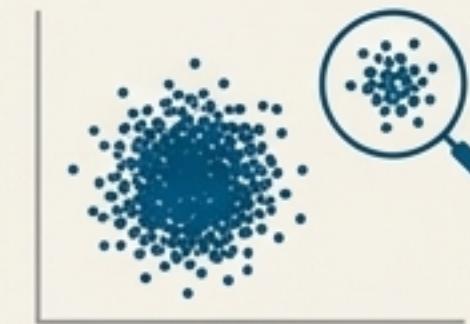
The Autonomous Discovery Module (A Three-Signal Approach)

1. Likelihood Anomaly



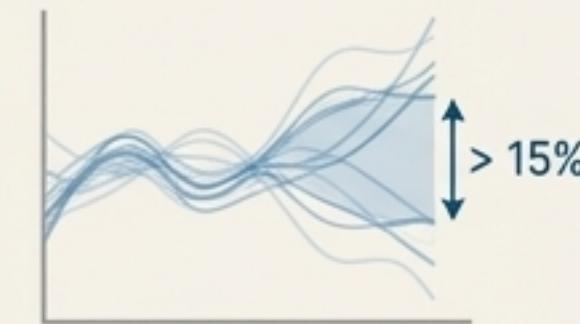
A sharp drop in model likelihood for a new data point (Z -score > 2.5 or $> 2\sigma$) indicates it's statistically different from the learned phases.

2. Latent Space Clustering



The HDBSCAN algorithm is applied to the model's latent space to find new, dense clusters of data points that correspond to novel phases.

3. Model Disagreement



High prediction variance from an ensemble of models ($> 15\%$ disagreement) serves as a strong uncertainty signal.

The Proving Ground: First, Validating the Framework on Well-Understood Classical Systems

Strategy

To first benchmark and validate the PI-CGM framework on large-scale, simulated datasets where the ground truth (phase transitions) is known with high precision.

Data Generation Details

Algorithms: Metropolis-Hastings and Wolff cluster algorithms will be used for efficient and accurate sampling.

Model	Lattice Size	# Samples	Parameters
2D Ising	64×64	50,000	$T = 1.0 \rightarrow 4.0$
3-State Potts	32×32	30,000	$T = 0.5 \rightarrow 2.5$
XY Model	32×32	20,000	Coupling $J=1$

The Proving Ground: Next, Tackling the Complexity of Quantum Systems

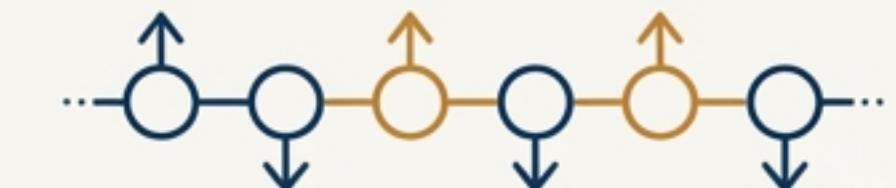
Strategy

To apply the validated framework to more challenging quantum mechanical systems, where phase transitions are subtle and manual discovery is extremely difficult.

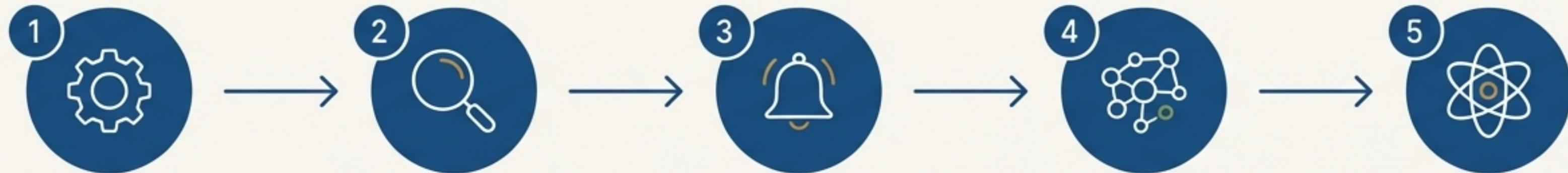
Data Generation Details

Methods: Data will be generated using established, high-fidelity numerical tools like Exact Diagonalization and Variational Monte Carlo (VMC).

Model	System Size	Data Size	Tool
Heisenberg chain	$N = 12$	5,000 states	NetKet
Transverse-field Ising	$N = 10$	8,000 states	QuSpin



The Autonomous Phase Discovery Pipeline in Action



Train

The PI-CGM is trained on data across a range of physical parameters (λ).

Scan & Compute

The model scans new data points, computing their likelihood under the learned distributions.

Detect Anomalies

The discovery module flags statistical anomalies (e.g., likelihood Z-score > 2.5) as potential new physical regimes.

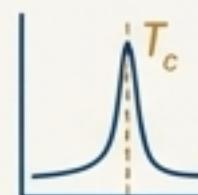
Cluster in Latent Space

HDBSCAN is used to confirm if these anomalies form a coherent new cluster.

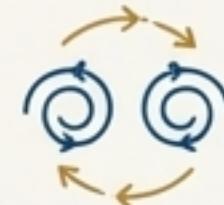
Confirm with Physics

The system cross-references flagged regions with known physical indicators of a phase transition (e.g., peaks in susceptibility, jumps in entanglement entropy) for final validation.

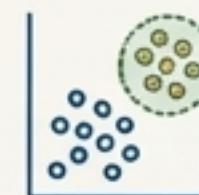
Example Expected Outputs



Automatically identify the critical temperature of the 2D Ising model to within ± 0.03 accuracy.



Successfully detect the Kosterlitz-Thouless transition in the XY model, a notoriously difficult task for standard CNNs.



Flag a potential novel cluster in the quantum Heisenberg dataset for further investigation.

Success is Defined by Clear, Quantitative Performance Targets

The framework's performance will be rigorously evaluated against established metrics from both machine learning and physics.

Physics Accuracy

Transition Point Estimation Error

$\leq 5\%$

of the known value for benchmark systems.

Model Confidence

Uncertainty Calibration (ECE)

< 0.1

Ensuring the model's confidence is reliable.

Likelihood Separation

$> 1.5\text{--}2.0 \sigma$

Standard deviations between different phases.

Discovery Power

Latent Clustering Purity (ARI)

> 0.8

Adjusted Rand Index for known clusters.

Reconstruction Error

< 0.05

Target for VAE-based models.

The Project Will Deliver a Reusable Framework and Publishable Scientific Insights

Primary Scientific Contributions



- **A Validated Framework:** A novel, uncertainty-aware generative methodology for the autonomous discovery of physical phases.



- **Empirical Benchmarks:** Rigorous testing and recovery of known classical and quantum phase transitions, providing a baseline for future work.



- **Potential for New Discovery:** A tool capable of flagging previously unknown or uncharacterised physical regimes in complex datasets.

Concrete Deliverables



- **Open-Source Code:** A complete, well-documented physics-ML pipeline implemented in PyTorch/JAX.



- **Generated Datasets:** Access to the curated classical and quantum datasets for community use.



- **A Publishable Study:** A manuscript detailing the methodology, results, and findings, suitable for a peer-reviewed journal.

The Project is Feasible and Built on a Realistic 12-Month Plan

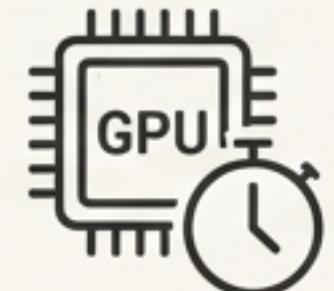
Leveraging a Mature Software Ecosystem

- ML Frameworks: PyTorch / JAX
- Physics Libraries: NetKet, QuSpin
- Simulation Code: Open-source MCMC implementations

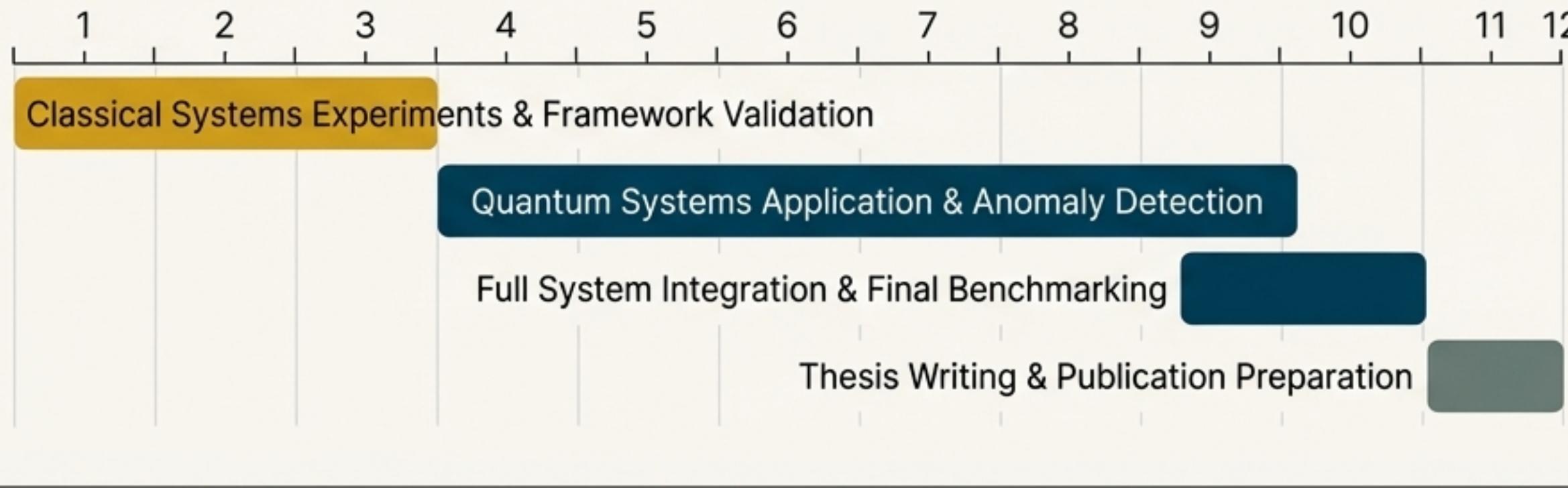


Validated Compute Requirements

- Hardware: A single GPU with 8–12 GB of VRAM is sufficient.
- Training Time: Each model experiment requires only 3–6 hours of training.



Project Timeline (12 Months)



Estimated Budget

£3,100–£3,600 (Primarily for cloud compute credits and documentation services).

A Feasible and Impactful Contribution to the Future of Computational Discovery

This project introduces a physics-informed generative system that:

- ✓ Learns the full distribution of complex physical systems.
- ✓ Autonomously detects novel phases using statistical and uncertainty-based signals.
- ✓ Provides physics-based validation for its discoveries.

The proposed research is a direct and novel approach to one of the key challenges in computational science. It is methodologically sound, technically feasible, and promises to yield significant contributions to machine learning for physics, computational materials science, and quantum simulation.