PHYSICS-CONSTRAINED FINE-TUNING OF FLOW-MATCHING MODELS FOR GENERATION AND INVERSE PROBLEMS

Anonymous authorsPaper under double-blind review

ABSTRACT

We present a framework for fine-tuning flow-matching generative models to enforce physical constraints and solve inverse problems in scientific systems. Starting from a model trained on low-fidelity or observational data, we apply a differentiable post-training procedure that minimizes weak-form residuals of governing partial differential equations (PDEs), promoting physical consistency and adherence to boundary conditions without distorting the underlying learned distribution. To infer unknown physical inputs, such as source terms, material parameters, or boundary data, we augment the generative process with a learnable latent parameter predictor and propose a joint optimization strategy. The resulting model produces physically valid field solutions alongside plausible estimates of hidden parameters, effectively addressing ill-posed inverse problems in a datadriven yet physics-aware manner. We validate our method on canonical PDE problems, demonstrating improved satisfaction of physical constraints and accurate recovery of latent coefficients. Further, we confirm cross-domain utility through fine-tuning of natural-image models. Our approach bridges generative modelling and scientific inference, opening new avenues for simulation-augmented discovery and data-efficient modelling of physical systems.

1 Introduction

Physical systems with rich spatio-temporal structure can be effectively represented by deep generative models, including diffusion and flow-matching methods (Kerrigan et al., 2024; Erichson et al., 2025; Baldan et al., 2025; Price et al., 2023). Although their dynamics can be highly complex, these systems are often governed by fundamental principles, such as conservation laws, symmetries, and boundary conditions, that constrain the space of admissible solutions. Incorporating such physical structure into generative modelling can improve both sample fidelity and out-of-distribution generalization.

In many scientific domains, including atmospheric and oceanographic modelling, seismic inversion, and medical imaging, we often observe system states without access to the underlying physical parameters that govern them. Crucially, PDE-based constraints are typically parameter-dependent, with residuals that vary according to material properties, source terms, or other latent variables. Prior work has largely focused on simple or global constraints—such as fixed boundaries or symmetries, that apply uniformly across the data distribution. Handling parameter-dependent constraints naively would require training over the joint distribution of solutions and parameters, which is often infeasible because parametric labels are missing, expensive to obtain, or high-dimensional. Addressing this limitation is critical for scientific discovery. Many inverse problems in the natural sciences and engineering require reasoning about unobserved parameters or exploring hypothetical scenarios inaccessible to direct experimentation. A generative model that can enforce parameter-dependent PDE constraints using only observational data would provide a powerful tool for data-efficient simulation, hypothesis testing, and the discovery of new physical phenomena, helping to bridge the gap between raw observations and mechanistic understanding.

This work proposes a framework for fine-tuning flow-matching generative models to enforce parameter-dependent PDE constraints without requiring joint parameter-solution training data. This

work aligns with a growing trend of simulation-augmented machine learning (Karniadakis et al., 2021), where generative models accelerate scientific discovery by efficiently exploring physically plausible solution spaces. Our approach reformulates fine-tuning as a stochastic optimal control problem via Adjoint Matching (Domingo-Enrich et al., 2025), guided by weak-form PDE residuals. By augmenting the model with a latent parameter evolution, we enable joint generation of physically consistent solution–parameter pairs, addressing ill-posed inverse problems. We evaluate our proposed fine-tuning framework on two canonical physical systems and show an application to natural images. We demonstrate denoising and conditional generation capabilities, including robustness to noisy data and the ability to infer latent parameters from sparse observations. Visual and quantitative results, including strong reductions in residuals across tasks and robustness to model misspecification, highlight the flexibility of our method for integrating physical constraints into generative modelling.

To sum up, our contributions are as follows:

- POST-TRAINING ENFORCEMENT OF PHYSICAL CONSTRAINTS: We introduce a finetuning strategy that tilts the generative distribution toward PDE-consistent samples using weak-form residuals, improving physical validity while preserving diversity.
- ADJOINT-MATCHING FINE-TUNING WITH THEORETICAL GROUNDING: Leveraging the
 adjoint-matching framework, we recast reward-based fine-tuning as a stochastic control
 problem, extending flow-matching models to generate latent parameters alongside states,
 enabling inverse problem inference without paired training data.
- BRIDGING GENERATIVE MODELING AND PHYSICS-INFORMED LEARNING: Our approach connects preference-aligned generation with physics-based inference, enabling simulation-augmented models to generate solutions that respect complex physical laws.

2 RELATED WORK

Physics-Constrained Generative Models Integrating physical constraints—such as boundary conditions, symmetry invariances, and partial differential equation (PDE) constraints—into machine learning models improves both accuracy and out-of-distribution generalization. Classical approaches, such as Physics-Informed Neural Networks (PINNs, Raissi et al. (2019)), directly regress solutions that satisfy governing equations. While effective for forward or inverse problems, PINNs do not capture distributions over solutions, making them unsuitable for generative tasks that require sampling diverse plausible outcomes.

In the generative setting, the main challenge is ensuring that the physically constrained samples retain the variability of the underlying generative model, avoiding pathalogical issues such as mode collapse. Bastek et al. (2024) proposes a unified framework for introducing physical constraints into Denoising Diffusion Probabilistic Models (DDPMs, Ho et al. (2020)) at pre-training time, by adding a first-principles physics-residual loss to the diffusion training objective. This loss penalises violations of governing PDEs (e.g. fluid dynamics equations) so that generated samples inherently satisfy physical laws. The method was empirically shown to reduce residual errors for individual samples significantly, while simultaneously acting as a regulariser against overfitting, thereby improving generalisation. To evaluate the physics-residual loss, one needs to compute the expected PDE residual of the final denoised sample conditioned on the current noisy state in the DDPM process. Accurately estimating this expectation requires generating multiple reverse-diffusion trajectories from the same noisy sample, which makes pre-training significantly more expensive. A common alternative is to use Tweedie's formula to approximate the conditional expectation in a single pass, but this shortcut introduces bias, particularly in the final denoising steps.

Zhang & Zou (2025) proposes enforcing constraints through a post-hoc distillation stage, where a deterministic student model is trained from a vanilla diffusion model to generate samples in one-step, regularised by a PDE residual loss. In Wang et al. (2025) the authors introduce PhyDA, diffusion-based data assimilation framework that ensures reconstructions obey PDE-based dynamics, specifically for atmospheric science. An autoencoder is used to encode sparse observations into a structured latent prior for the diffusion model, which is trained with an additional physical residual loss.

Inference- and Post-Training Constraint Enforcement Various works have proposed approaches to enforce PDE constraints at inference time, often in combination with observational constraints, drawing connections to conditional diffusion models (Dhariwal & Nichol, 2021; Ho & Salimans, 2021). Huang et al. (2024) introduce guidance terms within the denoising update of a score-based diffusion model to steer the denoising process towards solutions which are both consistent with data and underlying PDEs. A related approach was considered by Xu et al. (2025), further introducing an adaptive constraint to mitigate instabilities in early diffusion steps. In Christopher et al. (2024), the authors recast the inference-time sampling of a diffusion process as a constrained optimization problem, each diffusion step is projected to satisfy user-defined constraints or physical principles. This allows strict enforcement of hard constraints (including convex and non-convex constraints, as well as ODE-based physical laws) on the generated data. Lu & Xu (2024) consider the setting where the base diffusion model is trained on cheap, low-fidelity simulations, leveraging a similar approach to generate down-scaled samples via projection.

Flow-Matching Models for Simulation and Inverse Problems Flow-matching (FM, Lipman et al. (2023)) has emerged as a flexible generative modelling paradigm for complex physical systems across science, including molecular systems (Hassan et al., 2024), weather (Price et al., 2023) and geology (Zhang et al., 2025). In the context of physics-constrained generative models Utkarsh et al. (2025) introduces a zero-shot inference framework to enforce hard physical constraints in pretrained flow models, by repeatedly projecting the generative flow at sampling time. Similarly, Cheng et al. (2024) proposed the ECI algorithm, to adapt a pretrained flow-matching model so that it exactly satisfies constraints without using analytical gradients. In each iteration of flow sampling, ECI performs: an Extrapolation step (advancing along the learned flow), a Correction step (applying a constraint-enforcement operation), and an Interpolation step (adjusting back towards the model's trajectory). While projection approaches are a compelling strategy for hard constraints, they can be challenging particularly for local constraints such as boundary conditions, as direct enforcement can introduce discontinuities. The above approach mitigates this by interleaving projections with flow steps, however this relies on the flow's ability to rapidly correct such non-physical artifacts.

Baldan et al. (2025) propose Physics-Based Flow Matching (PBFM), which embeds constraints (PDE or symmetries) directly into the FM loss during training. The approach leverages temporal unrolling to refine noise-free final state predictions and jointly minimizes generative and physics-based losses without manual hyperparameter tuning of their tradeoff. To mitigate conflicts between physical constraints and the data loss, they employ the ConFIG (Liu et al., 2024), which combines the gradients of both losses in a way that ensures that gradient updates always minimise both losses simultaneously.

Related to our approach are the works on generative models for Bayesian inverse problems (Stuart, 2010), where the goal is to infer distributions over latent PDE parameters given partial or noisy observations. Conditional diffusion and flow-matching models can be used to generate samples from conditional distributions and posterior distributions, supporting amortized inference and uncertainty quantification (Song et al., 2021; Utkarsh et al., 2025; Zhang et al., 2023). Conditioning is typically achieved either through explicit parameter inputs or guidance mechanisms during sampling, as in classifier-guided diffusion. While effective when large volumes of paired training data is available, these approaches are less relevant to observational settings where parameters are unobserved. In contrast, our approach connects the observed data to the latent parameters only during post-training, requiring substantially smaller volumes of data.

3 Method

FM models are trained to learn and sample from a given distribution of data p_{data} . They approximate this distribution by constructing a Markovian transformation from noise to data, such that the time marginals of this transformation match those of a reference flow $X_t = \beta_t X_1 + \gamma_t X_0$. Specifically FM models learn a vector field $v_t(x)$ that transports noise to data, via the ODE $dX_t = v_t(X_t)$. We can optionally inject a noise schedule $\sigma(t)$ along the trajectory to define an equivalent SDE that preserves the same time marginals (Maoutsa et al., 2020),

$$dX_t = \left(v_t(X_t) + \frac{\sigma(t)^2}{2\eta_t} \left(v_t(X_t) - \frac{\dot{\beta}_t}{\beta_t} X_t\right)\right) + \sigma(t) dB_t =: b_t(X_t) + \sigma(t) dB_t, \quad (1)$$

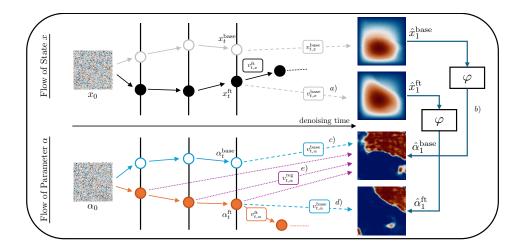


Figure 1: Visual depiction of proposed method. Starting at state x_t^{base} or x_t^{ft} , we use the base vector field $v_{t,x}^{\text{base}}$ to predict the final sample [a]]. Through the inverse predictor φ , we recover the corresponding predicted parameters $\hat{\alpha}_1^{\text{base}}$ and $\hat{\alpha}_1^{\text{ft}}[b]$. These estimates can be used as a target for evolving $\alpha_t^{\text{base}}[c]$ or as a baseline for the fine-tuned evolution of $\alpha_t^{\text{ft}}[d]$. For purposes of regularization, we further consider $v_{t,\alpha}^{\text{reg}}$, pointing from the current α_t^{ft} to the predicted final parameter of the base evolution $\hat{\alpha}_1^{\text{base}}[e]$.

where we combine coefficients β_t and γ_t into $\eta_t = \gamma_t \left(\frac{\dot{\beta}_t}{\beta} \gamma_t - \dot{\gamma}_t \right)$.

Assuming we have access to a FM model which generates samples according to distribution p(x), we seek to adjust this model so as to generate samples from the tilted distribution $p_r(x) \propto e^{\lambda \, r(x)} p(x)$, where r is a reward function and λ characterises the degree of distribution shift induced by finetuning.

To achieve this, we leverage the adjoint-matching framework of Domingo-Enrich et al. (2025). This work reformulates reward fine-tuning for flow-based generative models as a control problem in which the base generative process given by $v_t^{\rm base}$ is steered toward high-reward samples via modifying the learned vector field, which we denote as $v_t^{\rm ft}$ with corresponding drift term $b_t^{\rm ft}$. Our approach is conceptually related to reward- or preference-based fine-tuning of generative models (Christiano et al., 2017; Sun et al., 2024), where a learned or computed reward steers generation toward desired properties. Here, the reward is defined via PDE residuals, encoding knowledge about underlying dynamics and physical constraints to the solutions space as deviations to differential operators or boundary conditions.

Notably, we assume that the distribution generated by the base model p(x) only captures an observed quantity, but does not provide us with corresponding parameters or coefficient fields often needed to evaluate the respective differential operator. In the following, we will present a strategy of jointly recovering unknown parameters and fine-tuning the generation process.

3.1 REWARD

A generative model can reproduce the visual characteristics of empirical data while ignoring the physics that governs it, thereby rendering the samples unusable for downstream scientific tasks. To bridge this gap we impose the known governing equations as *soft constraints*, expressed through differential operators $\mathcal{L}_{\alpha}x=0$ with parameters α . Throughout, a generated sample x is interpreted as the discretisation of a continuous field $x(\xi)$ on a domain Ω . The *strong* PDE residual is defined as

$$\mathcal{R}_{\text{strong}}(x,\alpha) = \left\| \mathcal{L}_{\alpha} x \right\|_{L^{2}(\Omega)}^{2}.$$

In practice, strong residuals involve high-order derivatives that make the optimization landscape unstable. We therefore adopt *weak-form residuals* of the form $\langle \mathcal{L}_{\alpha} x, \psi \rangle_{L^2(\Omega)}$ for suitably chosen test functions $\psi \in \Psi$, which are numerically more stable under noisy or misspecified data. Repeated

applications of integration-by-parts can transfer derivatives from x to ψ . The set Ψ is composed of compactly supported local polynomial kernels. For each evaluation we draw N_{test} such functions; their centers and length-scales are sampled at random. A mollifier enforces $\psi|_{\partial\Omega}=0$, justifying the integration by parts. The resulting residual is

$$\mathcal{R}_{\text{weak}}(x,\alpha) = \frac{1}{N_{\text{test}}} \sum_{i=1}^{N_{\text{test}}} \left| \langle \mathcal{L}_{\alpha} x, \psi^{(i)} \rangle_{L^{2}(\Omega)} \right|^{2}.$$

These randomly sampled local test functions act as stochastic probes of PDE violations, providing a low-variance, data-efficient learning signal. A more detailed description of the test functions used can be found in Appendix D.3. Note that the residual might be augmented by adding soft constraints for boundary conditions.

3.2 Joint Evolution

Fine-tuning is nontrivial in our setting because we must infer latent physical parameters jointly with the generated solutions. On fully denoised samples, we can train an inverse predictor, i.e., $\varphi(x_1) = \alpha_1$, such that the weak PDE residual is minimised. As a naïve approach, this already induces a joint distribution over (x_1, α_1) via the push-forward through φ .

However, we advocate a more principled formulation that evolves both x and α along vector fields, enabling joint sampling of parameters and solutions, as well as a controlled regularisation of fine-tuning through the Adjoint Matching framework as outlined below. In the fine-tuning model, this can be achieved by directly learning the vector field $v_{t,\alpha}^{\rm ft}$ jointly with $v_{t,x}^{\rm ft}$ by augmenting the neural architecture. Since no ground-truth flow of α for the base model is available, at each state (x_t, α_t) we define a surrogate base flow using the inverse predictor φ . Specifically, we consider the one-step estimates

$$\hat{x}_1 = x_t + (1-t) v_t^{\text{base}}(x_t), \qquad \hat{\alpha}_1 = \varphi(\hat{x}_1).$$

The direction from the current state α_t to the predicted final parameter $\hat{\alpha}_1$ serves as a base vector field which we use to evolve alpha, i.e. $v_{t,\alpha}^{\text{base}}(\alpha_t) = (\hat{\alpha}_1 - \alpha_t)/(1-t)$ inducing corresponding drift $b_{t,\alpha}^{\text{base}}$. This surrogate base flow, starting at a noise sample $\alpha_0^{\text{base}} \sim \mathcal{N}(0,I)$, emulates a denoising process of the recovered parameter. We denote by α^{base} the parameter aligned with the base trajectory x^{base} . While the evolution of α^{base} does not influence the trajectory of x^{base} , the inferred vector field can be used to effectively regularize the generation of the fine-tuned model. Similarly, to regularise towards the parameter recovered under the base model, we introduce an additional field $v_{t,\alpha}^{\text{reg}}(\alpha_t^{\text{ft}}) = (\hat{\alpha}_1^{\text{base}} - \alpha_t^{\text{ft}})/(1-t)$. This vector field points from the current parameter estimate of the fine-tuned trajectory α_t^{ft} to the recovered parameter under the base model $\hat{\alpha}_1^{\text{base}}$. The field is used to pull the fine-tuned dynamics towards final samples associated with parameters similar to those of the base trajectory. The introduced vector fields are visualized in Fig. 1.

3.3 ADJOINT MATCHING

Considering an augmented state variable of the joint evolution $\tilde{X}_t = (X_t^T, \alpha_t^T)^T$, we cast fine-tuning as a stochastic optimal control problem:

$$\min_{\tilde{u}} \mathbb{E} \left[\int_{0}^{1} \left(\frac{1}{2} \left\| \tilde{u}_{t}(\tilde{X}_{t}) \right\|^{2} + f(\tilde{X}_{t}) \right) dt + g(\tilde{X}_{1}) \right]
\text{s.t.} \quad d\tilde{X}_{t} = \left(\tilde{b}_{t}^{\text{base}}(\tilde{X}_{t}) + \sigma(t) \, \tilde{u}_{t}(\tilde{X}_{t}) \right) dt + \sigma(t) \, d\tilde{B}_{t}$$
(2)

with control $\tilde{u}_t(\tilde{X}_t)$, running state cost $f(\tilde{X}_t)$, and terminal cost $g(\tilde{X}_1)$. In this formulation, fine-tuning amounts to a point-wise modification of the base drift through application of control \tilde{u} , i.e.

$$\tilde{b}_t^{\mathrm{fit}}(\tilde{X}_t) = \tilde{b}_t^{\mathrm{base}}(\tilde{X}_t) + \sigma(t)\,\tilde{u}_t(\tilde{X}_t).$$

In Domingo-Enrich et al. (2025), Adjoint Matching is introduced as a technique with lower variance and computational cost than standard adjoint methods. The method is based on a *Lean Adjoint* state, which is initialized as

$$\tilde{a}_1^T = \tilde{\lambda} \nabla_{\tilde{x}} g(\tilde{X}_1) = \left(\lambda_x \nabla_x g(X_1, \alpha_1), \, \lambda_\alpha \nabla_\alpha g(X_1, \alpha_1) \right)$$

and evolves backward in time according to

$$\frac{d}{dt}\tilde{a}_t = -\left(\nabla_{\tilde{x}}\tilde{b}_t^{\text{base}}(\tilde{X}_t)^T\tilde{a}_t + \nabla_{\tilde{x}}f(\tilde{X}_t)^T\right) = -\left(\begin{matrix}J_{xx}^T & J_{\alpha x}^T\\J_{x\alpha}^T & J_{\alpha \alpha}^T\end{matrix}\right)\begin{pmatrix}a_{t,x}\\a_{t,\alpha}\end{pmatrix} - \begin{pmatrix}\nabla_x f(X_t,\alpha_t)^T\\\nabla_{\alpha}f(X_t,\alpha_t)^T\end{matrix}\right) (3)$$

where the block-Jacobian is evaluated along the base drift for X and α , which means that $J_{ij} = \nabla_j \, b_{t,i}^{\mathrm{base}}(X_t, \alpha_t)$ for $i, j \in \{x, \alpha\}$. The hyperparameters λ_x and λ_α can be used to regulate the extent to which the fine-tuned distribution departs from the base distribution. The *Adjoint Matching* objective can then be formulated as a consistency loss:

$$\mathcal{L}(\tilde{u}; \tilde{X}) = \frac{1}{2} \int_{0}^{1} \|\tilde{u}_{t}(\tilde{X}_{t}) + \sigma(t) \tilde{a}_{t}\|^{2} dt
= \frac{1}{2} \int_{0}^{1} (\|u_{t,x}(X_{t}, \alpha_{t}) + \sigma(t) a_{t,x}\|^{2} + \|u_{t,\alpha}(X_{t}, \alpha_{t}) + \sigma(t) a_{t,\alpha}\|^{2}) dt.$$
(4)

It can be shown (Domingo-Enrich et al., 2025) that with f=0, this objective is consistent with the tilted target distribution for reward r=-g, if optimized with a *memoryless* noise schedule. This schedule ensures sufficient mixing during generation such that the final sample X_1 is independent of X_0 . To stabilise fine-tuning we introduce a scaled variant of the memoryless noise schedule. Instead of using the canonical choice $\sigma^2(t)=2\eta_t$ identified by Domingo-Enrich et al. (2025), we adopt

$$\sigma^2(t) = (1 - \kappa) \, 2\eta_t, \qquad 0 \le \kappa < 1,$$

which retains the theoretical memoryless property (see Lemma 1 in Appendix D.4) while attenuating the magnitude of the noise variance. The introduction of the scaling factor $0 \le \kappa < 1$ constitutes a simple but novel extension of the adjoint-matching framework. Whereas prior work highlighted a unique schedule, our analysis shows that a family of scaled schedules remains consistent with the memoryless condition. This additional degree of freedom acts as a numerical stabilisation knob, mitigating blow-ups near $t \to 0$ without losing theoretical consistency. Further, it offers a control-fidelity trade-off by regulating the amount of exploration. In practice, this flexibility allows practitioners to adapt fine-tuning to the conditioning of the PDE residuals and the stability of the solver, a feature not available in the original formulation.

Equation 2 is optimized by iteratively sampling trajectories with the fine-tuned model while following a memoryless noise schedule, numerically computing the lean adjoint states by solving the ODE in Equation 3, and taking a gradient descent step to minimize the loss in Equation 4. Note that gradients are only computed through the control \tilde{u}_t and not through the adjoint, reducing the optimization target to a simple regression loss. We state the full training algorithm and implementation details in Appendix D.5.

Adjoint Matching steers the generator toward the reward-tilted distribution, thereby reshaping the entire output distribution rather than correcting individual trajectories. However, when fine-tuning observational data or under system misspecification, we might be interested in retaining sample-specific detail. Empirically we find that this can be effectively encoded by imposing similarity of the inferred coefficients between base and fine-tuned model. Therefore, we add a running state cost

$$f(\alpha) = \lambda_f \left\| v_{t, \alpha}^{\text{ft}}(\alpha) - v_{t, \alpha}^{\text{reg}}(\alpha) \right\|^2$$

which penalises deviations of the fine-tuned α -drift from the direction pointing toward the base estimate $\hat{\alpha}_1^{\text{base}}$. The hyper-parameter λ_f controls a smooth trade-off: $\lambda_f=0$ recovers pure Adjoint Matching, while larger λ_f progressively anchors the final parameters α_1 obtained under the fine-tuned model to their base-model counterparts, thus retaining trajectory-level detail.

4 EXPERIMENTS

We evaluate across three settings: two PDE systems (one with observational noise, one with misspecified boundary data) and a natural-image model. Unlike latent-space fine-tuning for images, our PDE models operate directly in pixel space. High-variance noise during sampling can drive off-manifold trajectories and perturb PDE residuals, motivating $\kappa > 0$ for these models. For base Flow Matching backbones we use U-FNO (Wen et al., 2022) for PDEs and the DiT-based latent FM of

Dao et al. (2023) for images. In all experiments we first sample from the base generator and pre-train the inverse predictor φ to recover α by minimizing the (PDE) residual, then fine-tune. Following Domingo-Enrich et al. (2025), fine-tuning is initialized from the base weights. We augment capacity to condition $v_{t,x}^{\rm ft}$ on α_t and add a separate head for $v_{t,\alpha}^{\rm ft}$. Fine-tuning uses a memoryless noise schedule, while all reported results are generated without injected noise ($\sigma(t)=0$). Implementation details appear in App. D.2.

4.1 DARCY DENOISING

Consider a square domain of $\Omega=[0,1]^2$ where a permeability $\alpha(\xi)$, here serving as the hidden parameter, and a forcing term $f(\xi)$ induce a pressure distribution $x(\xi)$. We observe solutions for samples of permeability $\alpha\sim\mu_{\alpha}$, drawn from a discretized Gaussian Process such that a takes values $\alpha\in\{3,12\}$. Pressure x follows the law

$$-\nabla \cdot (\alpha(\xi) \nabla x(\xi)) - f(\xi) = 0, \quad x \in \Omega.$$

We assume zero Dirichlet boundary conditions and constant forcing. This is a standard dataset used in the Neural Operator literature, established in Li et al. (2020). However, to mimic observational data, we add noise to the generated solutions before training the base Flow Matching model. Details about the dataset are given in Appendix B. Figure 2 compares three Darcy samples generated from the *same* noise seed x_0 : the base draw, fine-tuning with our regularization (here $\lambda_f=1.0$), and fine-tuning without regularization. The base pressure $x^{\rm base}$ is visibly contaminated by high-frequency noise, and the inverse predictor φ correspondingly yields a scattered, artefact-ridden permeability map $\alpha^{\rm base}$. With regularization enabled, fine-tuning attenuates noise in the pressure $x^{\rm ft}$ while remaining close to $\alpha^{\rm base}$. Because $\alpha^{\rm base}$ is itself fragmented, some artefacts persist. In contrast, disabling regularization produces a fully denoised pressure and a markedly more coherent $\alpha^{\rm ft}$, but at the expected expense of erasing sample-specific details present in the base realization. More non-curated samples are given in Appendix F.1.

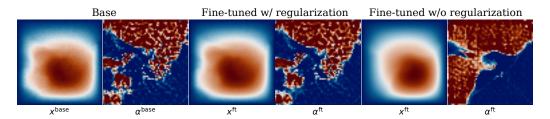


Figure 2: Denoising on Darcy for a fixed noise seed. Left: base pressure x^{base} and inferred permeability α^{base} . Middle: fine-tuned output with parameter regularization ($\lambda_f = 1.0$). Right: fine-tuning with $\lambda_f = 1.0$. Color maps throughout this work taken from Crameri et al. (2020).

Figure 3 quantifies how the hyperparameters $\lambda_x, \lambda_\alpha$ and λ_f mediate the trade-off between staying close to the base model and reducing the PDE residual. Each configuration is evaluated on 128 samples with noise seeds shared across models. In (a), starting from the base setting $\lambda_x = \lambda_\alpha = 0$, we jointly increase $\lambda_x = \lambda_\alpha$ while fixing $\lambda_f = 0$. As expected, the residual decreases with increased λ . However, this also leads to decreasing diversity in the generated samples, here reported as the complement of the mean pairwise SSIM for α . Thus $\lambda_x, \lambda_\alpha$ provide a direct control knob for the tuning the trade-off between fine-tuning and diversity. In (b), we fix $\lambda_x = \lambda_\alpha = 20 \text{K}$ and vary λ_f . Larger λ_f raises the residual, whereas weaker regularization attains lower residuals but departs more from the base, quantified by the mean pointwise relative MSE between fine-tuned and base samples for both x and α . The ablations indicate that the method can be adapted to produce task-appropriate samples. By tuning $(\lambda_x, \lambda_\alpha, \lambda_f)$, practitioners can steer fine-tuning toward specific targets, e.g., stronger residual reduction versus closer adherence to the base sample—informed by domain priors.

Computationally, adaptation is lightweight: fine-tuning on noisy Darcy requires only 20 gradient steps (hyperparameters in App. E.1) and completes in under 15 minutes on a single NVIDIA L40S, after which sampling proceeds at base-model cost with no inference-time adjustments.

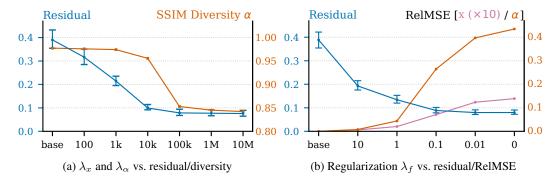


Figure 3: Ablation over $\lambda_x, \lambda_\alpha, \lambda_f$ (128 samples per setting, shared seeds x_0). (a) Increasing $\lambda_x = \lambda_\alpha$ with $\lambda_f = 0$ monotonically lowers the PDE residual while reducing diversity (reported as complement of mean pairwise SSIM). (b) Fixing $\lambda_x = \lambda_\alpha = 20 \mathrm{K}$ and sweeping λ_f trades residual against fidelity to the base via mean pointwise relative MSE for x (scaled up by 10) and α .

4.2 GUIDANCE ON SPARSE OBSERVATIONS

In many realistic settings dense observations of a state variable are available for pre-training a generative model, whereas only a few measurements of the latent parameter can be collected. To sample from the posterior of parameter–state pairs that respect such sparse evidence we steer the generative process through *guidance*. Huang et al. (2024) demonstrate guided sampling towards sparse observations from a model that was pre-trained on the joint parameter-state distribution. Our approach applies a similar guidance mechanism, however, to a model that was pre-trained on noisy state observations alone. We state details on the guiding mechanism in E.2. Figure 4 shows that the guided sampler adheres to sparse measurements while preserving realistic variability in the generated samples. Additional results for different amounts of conditioning observations are given in Appendix F.3.

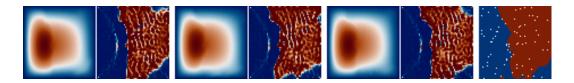


Figure 4: Three samples through guidance towards sparse observations (white markers in right panel) of the permeability, showing a plausible conditional distribution.

4.3 LINEAR ELASTICITY: MISSPECIFICATION

We consider plane–strain linear elasticity on $\Omega=[0,1]^2$ with spatially varying Young's modulus, considered as the hidden parameter $\alpha(\xi)$, and fixed Poisson ratio ν . The displacement $x:\Omega\to\mathbb{R}^2$ solves the static, body-force-free equilibrium. Boundary conditions are fully Dirichlet and identical across samples: the left and right edges are clamped, while the top and bottom edges receive inward sinusoidal normal displacements with zero tangential slip and common amplitude $A_{\rm top}=A_{\rm bot}=0.10$. During fine-tuning, we add a penalty term for deviations from the boundary conditions under modified $A_{\rm bot}=0.075$ to the weak PDE residual. Detailed specifications of this system and data generation can be found in Appendix B. Figure 5 shows resulting displacement fields from the base and fine-tuned model together with heatmaps of their weak residuals. Our fine-tuning method successfully adjusts the lower boundary while retaining the details of the base samples and keeping low PDE residuals. While also ECI (Utkarsh et al., 2025) adjusts to the new boundary condition, the sampled displacements appear physically implausible and show high residual values. Details and non-curated samples are given in the Appendix (E.3, F.2).

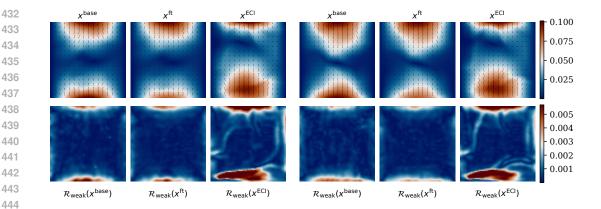


Figure 5: Fine-tuning towards boundary modification, comparing our approach with ECI. Samples for two random seeds shared across models. Top row: displacement fields. Bottom row: corresponding weak residual heatmaps.

NATURAL IMAGES: PARAMETRIC COLOR TRANSFORMATION

To demonstrate cross-domain utility, we apply our method to natural images by introducing a parametric recoloring pathway: analogous to the hidden PDE parameter, α here specifies a polynomial color transform that operates outside the latent space, enabling exploration of image appearances not well supported by the base distribution. We use a class-conditional Latent Flow Matching (LFM) model (Dao et al., 2023) pretrained on ImageNet-1k (Deng et al., 2009) and optimize PickScore (Kirstain et al., 2023) with a globally fixed prompt. As a concrete example, we fine-tune on the class macaw with the prompt "close-up Pop Art of a macaw parrot," yielding the samples in Fig. 6. Joint fine-tuning with recoloring produces markedly more vibrant palettes and, crucially, joint adjustments (e.g., background textures that the recoloring exploits). Details about the recoloring parametrization are given in Appendix E.4 and further non-curated samples are provided in Appendix F.4.



(a) Vanilla Adjoint Matching

(b) Joint model with parametric recoloring

Figure 6: Fine-tuning of LFM model on macaw class using prompt "close-up Pop Art of a macaw parrot", comparing vanilla Adjoint Matching with our joint approach.

CONCLUSION

We have introduced a framework for post-training fine-tuning of flow-matching generative models to enforce physical constraints and jointly infer latent physical parameters informing the constraints. Through a novel architecture, combined with the combination of weak-form PDE residuals with an adjoint-matching scheme our method can produce samples that adhere to complex constraints without significantly affecting the sample diversity. Preliminary experiments on Darcy flow and acoustic wave problems demonstrate the potential of this method to reduce residuals and enable joint solution-parameter generation, supporting its promise for physics-aware generative modelling. Future steps include adaptive approaches to optimising trade-off between constraint enforcement and generative diversity, and extending the framework to more complex and multi-physics systems, including coupled PDEs and stochastic or chaotic dynamics. We would also explore how this methodology can be leveraged for uncertainty quantification and propagation, and downstream tasks such as optimal sensor placement and scientific discovery workflows.

REPRODUCIBILITY STATEMENT

We report datasets, model backbones, training schedules, loss definitions, evaluation metrics, and the key hyperparameters required to reproduce our results in the main text and appendix. Remaining implementation choices are documented in the released configuration files. Upon acceptance, we will open-source the code, including full training scripts. We fixed random seeds where applicable and specify hardware/software versions.

REFERENCES

- Giacomo Baldan, Qiang Liu, Alberto Guardone, and Nils Thuerey. Flow matching meets pdes: A unified framework for physics-constrained generation. *arXiv preprint arXiv:2506.08604*, 2025.
- Jan-Hendrik Bastek, WaiChing Sun, and Dennis Kochmann. Physics-informed diffusion models. In *The Thirteenth International Conference on Learning Representations*, 2024.
- Yoav Chai, Raja Giryes, and Lior Wolf. Supervised and unsupervised learning of parameterized color enhancement. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 992–1000, 2020.
- Chaoran Cheng, Boran Han, Danielle C Maddix, Abdul Fatir Ansari, Andrew Stuart, Michael W Mahoney, and Yuyang Wang. Gradient-free generation for hard-constrained systems. *arXiv* preprint arXiv:2412.01786, 2024.
- Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling laws for contrastive language-image learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 2818–2829, 2023.
- Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.
- Jacob K Christopher, Stephen Baek, and Nando Fioretto. Constrained synthesis with projected diffusion models. *Advances in Neural Information Processing Systems*, 37:89307–89333, 2024.
- Fabio Crameri, Grace Shephard, and Philip Heron. The misuse of colour in science communication. *Nature Communications*, 11, 2020. doi: 10.1038/s41467-020-19160-7.
- Quan Dao, Hao Phung, Binh Nguyen, and Anh Tran. Flow matching in latent space. *arXiv preprint arXiv:2307.08698*, 2023.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.
- Carles Domingo-Enrich, Michal Drozdzal, Brian Karrer, and Ricky T. Q. Chen. Adjoint matching: Fine-tuning flow and diffusion generative models with memoryless stochastic optimal control. In *International Conference on Representation Learning*, 2025.
- N Benjamin Erichson, Vinicius Mikuni, Dongwei Lyu, Yang Gao, Omri Azencot, Soon Hoe Lim, and Michael W Mahoney. Flex: A backbone for diffusion-based modeling of spatio-temporal physical systems. *arXiv preprint arXiv:2505.17351*, 2025.
- Majdi Hassan, Nikhil Shenoy, Jungyoon Lee, Hannes Stärk, Stephan Thaler, and Dominique Beaini. Et-flow: Equivariant flow-matching for molecular conformer generation. *Advances in Neural Information Processing Systems*, 37:128798–128824, 2024.
 - Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021.

- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- Jiahe Huang, Guandao Yang, Zichen Wang, and Jeong Joon Park. Diffusionpde: Generative pdesolving under partial observation. *Advances in Neural Information Processing Systems*, 37: 130291–130323, 2024.
 - George Em Karniadakis, Ioannis G Kevrekidis, Lu Lu, Paris Perdikaris, Sifan Wang, and Liu Yang. Physics-informed machine learning. *Nature Reviews Physics*, 3(6):422–440, 2021.
 - Gavin Kerrigan, Giosue Migliorini, and Padhraic Smyth. Functional flow matching. In *International Conference on Artificial Intelligence and Statistics*, pp. 3934–3942. PMLR, 2024.
 - Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Picka-pic: An open dataset of user preferences for text-to-image generation. *Advances in neural information processing systems*, 36:36652–36663, 2023.
 - Zongyi Li, Nikola Kovachki, Kamyar Azizzadenesheli, Burigede Liu, Kaushik Bhattacharya, Andrew Stuart, and Anima Anandkumar. Fourier neural operator for parametric partial differential equations. *arXiv preprint arXiv:2010.08895*, 2020.
 - Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. In 11th International Conference on Learning Representations, ICLR 2023, 2023.
 - Qiang Liu, Mengyu Chu, and Nils Thuerey. Config: Towards conflict-free training of physics informed neural networks. *arXiv preprint arXiv:2408.11104*, 2024.
 - Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint* arXiv:1711.05101, 2017.
 - Yulong Lu and Wuzhe Xu. Generative downscaling of pde solvers with physics-guided diffusion models. *Journal of scientific computing*, 101(3):71, 2024.
 - Dimitra Maoutsa, Sebastian Reich, and Manfred Opper. Interacting particle solutions of fokker–planck equations through gradient–log–density estimation. *Entropy*, 22(8):802, 2020.
 - Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
 - William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 4195–4205, 2023.
 - Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual reasoning with a general conditioning layer. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
 - Boris T Polyak and Anatoli B Juditsky. Acceleration of stochastic approximation by averaging. *SIAM journal on control and optimization*, 30(4):838–855, 1992.
 - Ilan Price, Alvaro Sanchez-Gonzalez, Ferran Alet, Tom R Andersson, Andrew El-Kadi, Dominic Masters, Timo Ewalds, Jacklynn Stott, Shakir Mohamed, Peter Battaglia, et al. Gencast: Diffusion-based ensemble forecasting for medium-range weather. *arXiv preprint arXiv:2312.15796*, 2023.
 - Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PmLR, 2021.
 - Maziar Raissi, Paris Perdikaris, and George E Karniadakis. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational physics*, 378:686–707, 2019.

- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10684–10695, 2022.
 - Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pp. 234–241. Springer, 2015.
 - Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in neural information processing systems*, 35:25278–25294, 2022.
 - Yang Song, Liyue Shen, Lei Xing, and Stefano Ermon. Solving inverse problems in medical imaging with score-based generative models. In *NeurIPS 2021 Workshop on Deep Learning and Inverse Problems*, 2021.
 - Andrew M Stuart. Inverse problems: a bayesian perspective. Acta numerica, 19:451–559, 2010.
 - Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang Gan, Liang-Yan Gui, Yu-Xiong Wang, Yiming Yang, et al. Aligning large multimodal models with factually augmented rlhf. In *Annual Meeting of the Association for Computational Linguistics*, 2024.
 - Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. *Advances in neural information processing systems*, 30, 2017.
 - Utkarsh Utkarsh, Pengfei Cai, Alan Edelman, Rafael Gomez-Bombarelli, and Christopher Vincent Rackauckas. Physics-constrained flow matching: Sampling generative models with hard constraints. *arXiv preprint arXiv:2506.04171*, 2025.
 - Hao Wang, Jindong Han, Wei Fan, Weijia Zhang, and Hao Liu. Phyda: Physics-guided diffusion models for data assimilation in atmospheric systems. *arXiv preprint arXiv:2505.12882*, 2025.
 - Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.
 - Gege Wen, Zongyi Li, Kamyar Azizzadenesheli, Anima Anandkumar, and Sally M Benson. U-fno—an enhanced fourier neural operator-based deep-learning model for multiphase flow. *Advances in Water Resources*, 163:104180, 2022.
 - Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, et al. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, pp. 38–45, 2020.
 - Ruichen Xu, Haochun Wang, Georgios Kementzidis, Chenhao Si, and Yuefan Deng. Apod: Adaptive pde-observation diffusion for physics-constrained sampling. In *ICML 2025 Workshop on Assessing World Models*, 2025.
 - Hao Zhang, Yuanyuan Li, and Jianping Huang. Diffusionvel: Multi-information integrated velocity inversion using generative diffusion model. In *86th EAGE Annual Conference & Exhibition*, volume 2025, pp. 1–5. European Association of Geoscientists & Engineers, 2025.
 - Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 3836–3847, 2023.
 - Yi Zhang and Difan Zou. Physics-informed distillation of diffusion models for pde-constrained generation. *arXiv preprint arXiv:2505.22391*, 2025.

A USE OF LARGE LANGUAGE MODELS

We employed large language models to polish the manuscript (wording, grammar, and synonyms) and to assist with implementation details, including plotting scripts and code rewriting/refactoring. The research questions, problem formulation, algorithmic design and experimental methodology, however, were conceived by the authors. All LLM-produced text and code were reviewed, adapted, and verified by the authors prior to inclusion.

B DATASET DETAILS

In this section we detail the datasets used throughout our study. Our guiding principle was to select scenarios in which the underlying parameter fields contain sharp discontinuities, thereby inducing rich, non-linear behaviour in the associated state variables and making the inverse problem decidedly non-trivial. Although the Darcy-flow benchmark follows the setup of Li et al. (2020), we regenerate the data so that the sample count, grid resolution, and ground-truth parameters can be controlled precisely. Complete scripts for producing both the Darcy and elasticity datasets will be released to facilitate transparency and reproducibility.

B.1 DARCY FLOW DATASET

The dataset comprises 20,000 independent samples, each a pair (a, u) on the unit square $\Omega = [0, 1]^2$, where $a: \Omega \to \mathbb{R}_{>0}$ is the permeability and $u: \Omega \to \mathbb{R}$ is the steady-state pressure solving

$$-\nabla \cdot (a(\xi) \nabla u(\xi)) = f(\xi), \qquad \xi = (\xi_1, \xi_2) \in \Omega,$$

with homogeneous Dirichlet boundary conditions $u|_{\partial\Omega}=0$ and constant forcing $f\equiv 1$.

Discretization. We use a uniform 64×64 nodal grid with spacing $\Delta x = 1/(64-1)$ and the standard five-point finite-difference scheme. Interface permeabilities are formed by two-point arithmetic averaging. Dirichlet values are imposed strongly, yielding a sparse SPD linear system that is solved with a direct sparse solver.

Permeability sampling. We draw a zero-mean Gaussian random field a_{raw} via a cosine-basis Karhunen-Loève synthesis on Ω associated with the Matérn-type covariance operator

$$\mathcal{C} = (\tau^2 - \Delta)^{-\alpha}, \qquad \alpha = 2, \ \tau = 3,$$

i.e., using the DCT-II basis (Neumann eigenfunctions of $-\Delta$), setting the DC mode to zero to enforce exact zero mean, scaling by the spectrum of C, and applying an orthonormal inverse DCT to obtain a grid realization.

To model sharp contrasts, we threshold the Gaussian field into a piecewise-constant permeability,

$$a(\xi) = \begin{cases} 12, & a_{\text{raw}}(\xi) \ge 0, \\ 3, & a_{\text{raw}}(\xi) < 0. \end{cases}$$

Observational noise. In experiments with noisy observations, we corrupt the pressure with additive Gaussian noise

$$\tilde{u} = u + \sigma \varepsilon, \qquad \varepsilon \sim \mathcal{N}(0, I), \quad \sigma = 10^{-4}.$$

B.2 Linear Elasticity Dataset

The dataset contains $N=10{,}000$ independent samples on the unit square $\Omega=[0,1]^2$. Each sample is a pair (E,u) where $E:\Omega\to\mathbb{R}_{>0}$ denotes the spatially varying Young's modulus and $u:\Omega\to\mathbb{R}^2$ is the plane–strain displacement field under fixed Poisson ratio $\nu=0.30$. The boundary loading is deterministic and identical across samples: the left/right edges are clamped, and sinusoidal normal displacements are prescribed on the top and bottom edges with zero tangential slip (amplitudes $A_{\rm top}=A_{\rm bot}=0.10$). For each E,u solves the static linear elasticity equations (no body force)

$$-\nabla \cdot \sigma(\xi) = 0, \qquad \sigma \ = \ \lambda(E,\nu)\operatorname{tr}(\varepsilon) I + 2\,\mu(E,\nu)\,\varepsilon, \qquad \varepsilon \ = \ \tfrac{1}{2} \big(\nabla u + \nabla u^\top\big), \quad \xi \in \Omega,$$

with Lamé parameters $\lambda(E,\nu)=\frac{\nu\,E}{(1+\nu)(1-2\nu)}$ and $\mu(E,\nu)=\frac{E}{2(1+\nu)}$, and Dirichlet boundary conditions

$$\begin{split} u(\xi_1 = 0, \xi_2) &= 0, \quad u(\xi_1 = 1, \xi_2) = 0 \\ u_x(\xi_1, \xi_2 = 0) &= 0, \quad u_y(\xi_1, \xi_2 = 0) = -A_{\mathrm{bot}} \sin(\pi \, \xi_1) \\ u_x(\xi_1, \xi_2 = 1) &= 0, \quad u_y(\xi_1, \xi_2 = 1) = +A_{\mathrm{top}} \sin(\pi \, \xi_1). \end{split}$$

During fine-tuning, we set $A_{\rm bot}=0.075$ to enforce adaptation of the solutions, simulating a misspecification between observed data and the assumed model.

Discretization. We use a uniform 64×64 grid on Ω and standard second–order finite differences in the small–strain regime. Dirichlet data are imposed strongly. The discrete equilibrium equations are advanced by a stable explicit time–marching (damped gradient) scheme for the static problem until the global residual norm falls below 10^{-6} , or a cap of 2×10^4 iterations is reached.

Coefficient field sampling (piecewise-constant Voronoi medium). Heterogeneous modulus fields E are obtained from a Voronoi tessellation constructed by drawing a fixed number of sites uniformly in Ω , assigning to each Voronoi cell a modulus sampled log-uniformly within [1.0, 10.0], and rasterizing the resulting partition to the computational grid by nearest-site labeling. This produces piecewise-constant E with sharp jumps that emulate multi-phase media. To control interface smoothness, a separable Gaussian blur with standard deviation $\sigma_{\rm blur}=1.0$ (in grid units) is applied to the rasterized field.

C PRE-TRAINING OF FLOW MATCHING MODELS

We adopt the vanilla Flow–Matching (FM) procedure of Lipman et al. (2023). Based on the optimal-transport reference flow

$$X_t = \beta_t X_1 + \gamma_t X_0$$

with $\beta_t = t$, and $\gamma_t = 1 - t$, we can define conditional vector field as training targets for a parametric model $v_{\theta}(x_t, t)$. Given and end-point $x_1 \sim p_{\text{data}}$, the conditional vector field is available as

$$v_t(x|x_1) = \frac{1}{1-t}(x_1-x).$$

This leads to the simplest form of Flow Matching objectives:

$$\mathcal{L}_{\text{FM}}^{\text{OT}} = \mathbb{E} \|v_{\theta}(X_t, t) - (X_1 - X_0)\|^2$$

where $X_0 \sim \mathcal{N}(0, I)$, $X_1 \sim p_{\text{data}}$ and $t \sim U[0, 1]$.

Network Backbone. For PDE data, the mapping v_{θ} is realised with a U-FNO (Wen et al., 2022), which combines Fourier Neural Operator (FNO, Li et al. (2020)) with U-Net (Ronneberger et al., 2015) layers. The FNO layers have an inductive bias towards low-frequency solutions and are therefore particularly suited for modeling PDE data, while the U-Net layers help to refine outputs and produce discontinuities.

Departing from the original design, we use the U–Net skip-connection structure in all layers. Also, we employ a more powerful time-conditioning in the U–Net layers by using FiLM-style Adaptive Group Normalization (AdaGN) time conditioning, i.e., predicting per-channel scale and shift from a sinusoidal time embedding and applying them after GroupNorm in each residual block. This follows the feature-wise linear modulation idea of Perez et al. (2018) and its diffusion U–Net instantiations in Dhariwal & Nichol (2021); Rombach et al. (2022).

We prepend the physical input with fixed sinusoidal encodings for absolute spatial coordinates and the normalised time stamp so the spectral backbone receives explicit space-time context. Padding inside the U-FNO (for the spectral convolution layers) is reflective, which has empirically worked better than replicative. Before training, data is standardised to zero mean and unit variance.

Table 1 compiles the hyperparameters of the U-FNO architecture used for the base FM models in our experiments.

On natural image data, we use a pretrained Latent Flow Matching Model (Dao et al., 2023) based on a DiT (Peebles & Xie, 2023) backbone.

Table 1: U-FNO (2D) backbone hyperparameters used in our experiments.

Hyperparameter	Value	Description
Number of layers	4	Count of spectral and U-Net operator blocks.
Padding mode	reflect	Boundary padding for convolutions/lifts.
Input channels	1 / 2	Channels of input field x .
Output channels	1 / 2	Channels of output field.
Time embedding dim	32	Dimensionality of time conditioning.
Spatial embedding dim	32	Dimensionality of coordinate embedding.
Lifting width	256	Channels in input lifting/projection.
Fourier modes (kx, ky)	[32, 32]	Retained spectral modes per axis.
Spectral block width	32	Channel width within FNO layers.
U-Net base widths	32	Stage-wise channel widths.
Embedding width (U–Net)	64	Channels for conditioning embeddings.
Attention stages	[]	Stages with self-attention (empty \Rightarrow none).
Attention heads	[]	Multi-head count if attention is enabled.
Total number of parameters	≈ 19M	

Optimisation. We train the FM backbone with AdamW (Loshchilov & Hutter, 2017), using a linear warmup of the learning rate from 0 to the base value over the first $p_{\rm wu}$ fraction of training steps, followed by a constant learning rate thereafter. For evaluation stability, we maintain an exponential moving average (EMA) of the network parameters $\theta^{\rm EMA} \leftarrow \beta \, \theta^{\rm EMA} + (1-\beta) \, \theta$, a practice rooted in Polyak averaging (Polyak & Juditsky, 1992) and widely adopted in modern deep models (e.g., Tarvainen & Valpola, 2017).

Table 2: Flow Matching (FM) training hyperparameters and schedule.

Hyperparameter	Description	
Batch size	128	Minibatch size per step.
Base learning rate	1e-4	AdamW step size.
Warmup fraction p_{wu}	0.01	Fraction of total steps used for linear LR warmup.
Epochs	300	Full passes over the dataset.
Optimizer	AdamW	Applied to FM backbone parameters.
Weight decay	0.01	AdamW L2/decoupled decay coefficient.
EMA decay	0.9998	Exponential moving average of weights for evaluation.
Schedule after warmup	constant	LR held constant after warmup.

The FM training configuration is summarized in Table 2. For the Darcy dataset (20k samples), training takes around 12 hours on a single NVIDIA RTXA6000 GPU. For the smaller (10k) elasticity dataset, training takes around 7 hours using the same configuration.

D METHOD: IMPLEMENTATION DETAILS

In this section, we provide further details into the implementation of our fine-tuning method. This includes neural network architectures, the specific design of test functions for the computation of the weak residuals, numerical heuristics, and finally the full training algorithm for fine-tuning. All relevant code is implemented using Python 3.12.3, specifically, neural architectures are implemented in PyTorch Paszke et al. (2019) version 2.7.

D.1 INVERSE PREDICTOR φ

We parametrise the inverse map φ with a two-layer U-FNO, mirroring the spectral–spatial bias of the forward backbone. However, we increase the capacity of the U-Net components and use attention at the two lowest resolutions (stage indices 2 and 3). Since the inverse predictor only operates at t=1, we drop the temporal but keep the spatial conditioning. Exact parameters are stated in Table 3

Table 3: U-FNO Inverse Predictor.		
Hyperparameter	Value	
Number of layers	2	
Padding mode	reflect	
Input channels	1 / 2	
Time embedding dimension	0	
Spatial embedding dimension	32	
Output channels	1	
Lifting width	256	
Fourier modes (k_x, k_y)	[32, 32]	
Spectral block width	32	
U–Net base widths	[64, 64, 96, 128]	
U-Net embedding width	64	
Attention stages	[2, 3]	
Attention heads	Λ	

D.2 ARCHITECTURE MODIFICATIONS

Fine-tuning augments the base U-FNO map $x \mapsto v_{t,x}^{\text{base}}(x,t)$ to a joint, α -conditioned vector field

$$(v_{t,x}^{\mathrm{ft}}, v_{t,\alpha}^{\mathrm{ft}}) = v^{\mathrm{ft}}(x, \alpha, t),$$

implemented as residual corrections around the frozen U–FNO core. Given the padded input x_{pad} (with time/space embeddings) and conditioning fields α and $v_{t,\alpha}^{\mathrm{base}}$, the base feature stack x_* is produced by the original spectral+skip+U–Net blocks. A first correction head (a U–Net) takes $[x_*, v_{t,x}^{\mathrm{base}}, \alpha]$ and outputs an additive refinement, yielding

$$v_{t,x}^{\mathrm{ft}} \; = \; v_{t,x}^{\mathrm{base}} \; + \; \mathcal{U}_{\!x}\!\big(x_*,\, v_{t,x}^{\mathrm{base}},\, \alpha,\, t\big).$$

After unpadding, a *pixel-wise correction* (lightweight channel-wise MLP) further adjusts $v_{t,x}^{\text{ft}}$ using local features and 2D positional channels,

$$v_{t,x}^{\text{ft}} \leftarrow v_{t,x}^{\text{ft}} + \mathcal{M}_{x}(\text{pos}(x_*,t), v_{t,x}^{\text{ft}}),$$

which provides extra capacity for rapid adaptation without altering the global operator. For the parameter dynamics, we adopt a strictly residual strategy that conditions on both α and the baseline field:

$$v_{t,\alpha}^{\text{ft}} = v_{t,\alpha}^{\text{base}} + \mathcal{U}_{\alpha}(x_{\text{pad}}, \alpha, v_{t,\alpha}^{\text{base}}, t),$$

where \mathcal{U}_{α} is a second U–Net. All correction heads are zero-initialized at their final projection layers, so that at the start of fine-tuning $v^{\mathrm{ft}} \equiv (v_{t,x}^{\mathrm{base}}, v_{t,\alpha}^{\mathrm{base}})$ and departures are learned smoothly. Table 4 lists the parameters of the correction U–Nets.

Overall, the modifications to the base architecture add around 6M parameters to the model, resulting in a total \approx 25M parameters.

D.3 WEAK RESIDUALS AND TEST FUNCTIONS

Darcy Flow. A pressure field u that solves our Darcy flow equations fulfills $-\nabla \cdot (a\nabla u) - f = \mathcal{L}_a u = 0$ on $\Omega \subset \mathbb{R}^d$ with homogeneous Dirichlet data $u|_{\partial\Omega} = 0$. For any $\psi \in C^1_0(\Omega)$ we compute the L^2 inner product by multiplying with ψ and integrating:

Table 4: Parameterization of the fine-tuning U–Net heads \mathcal{U}_x and \mathcal{U}_α .

Hyperparameter	Value
U-Net base widths	[64,64,96,128]
Embedding width (time/aux)	64
Hidden lift (internal width)	256
Self-attention stages	[2, 3]
Attention heads	4

$$\langle \mathcal{L}_a u, \psi \rangle_{L^2(\Omega)} = \int_{\Omega} \left(-\nabla \cdot (a \nabla u) - f \right) \psi \ d\xi$$

One application of the divergence theorem yields

$$\langle \mathcal{L}_{a}u, \psi \rangle_{L^{2}(\Omega)} = -\int_{\partial \Omega} (a\nabla u \cdot n) \psi \, d\xi$$
$$+ \int_{\Omega} (a\nabla u \cdot \nabla \psi) - f \psi \, d\xi$$
$$= \int_{\Omega} (a\nabla u \cdot \nabla \psi) - f \psi \, d\xi,$$

where the boundary term vanishes because of $\psi|_{\partial\Omega}=0$. This expression only contains a first-order derivative of u and can be used to compute in the computation of the weak residual. To obtain a dimensionless, coefficient-scaled quantity comparable across locations, we normalize by the local mean permeability over the support of ψ ,

$$\overline{a}_{\psi} := \frac{1}{|\operatorname{supp} \psi|} \int_{\operatorname{supp} \psi} a(\xi) d\xi, \qquad \widetilde{\mathcal{R}}[\psi] := \frac{\mathcal{R}[\psi]}{\overline{a}_{\psi}}.$$

In practice, supp ψ is the compact patch where ψ (or its mollified variant) is nonzero, so that \overline{a}_{ψ} captures the local coefficient scale used to normalize the residual.

Linear Elasticity. For any compactly supported vector test $\psi \in C_0^1(\Omega; \mathbb{R}^2)$, the weak residual is

$$\mathcal{R}[\psi] := \langle \mathcal{L}_E u, \psi \rangle_{\mathrm{L}^2(\Omega)} = \int_{\Omega} \left(-\nabla \cdot \sigma(u; E, \nu) \right) \cdot \psi \, d\xi$$

A single integration by parts yields

$$\mathcal{R}[\psi] = -\int_{\partial\Omega} (\sigma n) \cdot \psi \, dS + \int_{\Omega} \sigma : \nabla \psi \, d\xi = \int_{\Omega} \sigma : \nabla \psi \, d\xi,$$

since $\psi|_{\partial\Omega}=0$. Here " $\sigma:\nabla\psi$ " is the Frobenius product and we denote the stress components by $\sigma_{xx},\sigma_{xy}(=\sigma_{yx}),\sigma_{yy}$. To reuse the same scalar test generator for both momentum equations, we restrict to tests that share a single scalar profile in both components. Concretely, we take

$$\psi^{(x)} = (\psi, 0), \qquad \psi^{(y)} = (0, \psi), \qquad \psi \in C_0^1(\Omega)$$

With this restriction, residuals can be computed component-wise

$$\langle \mathcal{L}_E u, \psi^{(x)} \rangle = \int_{\Omega} \left(\sigma_{xx} \, \partial_{\xi_1} \psi + \sigma_{xy} \, \partial_{\xi_2} \psi \right) d\xi, \qquad \langle \mathcal{L}_E u, \psi^{(y)} \rangle = \int_{\Omega} \left(\sigma_{xy} \, \partial_{\xi_1} \psi + \sigma_{yy} \, \partial_{\xi_2} \psi \right) d\xi.$$

To obtain a dimensionless, coefficient-scaled quantity, we normalize by the local mean modulus over the support of ψ ,

$$\overline{E}_{\psi} := \frac{1}{|\operatorname{supp} \psi|} \int_{\operatorname{supp} \psi} E(\xi) \, d\xi, \qquad \widetilde{\mathcal{R}}^{(x)}[\psi] := \frac{\langle \mathcal{L}_E u, \psi^{(x)} \rangle}{\overline{E}_{\psi}}, \quad \widetilde{\mathcal{R}}^{(y)}[\psi] := \frac{\langle \mathcal{L}_E u, \psi^{(y)} \rangle}{\overline{E}_{\psi}}.$$

For a family of tests $\{\psi_k\}_{k=1}^N$, the scalar residual used in experiments is the ℓ_2 aggregation over components followed by averaging over tests. Therefore, the weak residual for the elasticity experiment is

$$\mathcal{R}_{\text{weak}} := \frac{1}{N} \sum_{k=1}^{N} \left(\left(\widetilde{\mathcal{R}}^{(x)}[\psi_k] \right)^2 + \left(\widetilde{\mathcal{R}}^{(y)}[\psi_k] \right)^2 \right).$$

Wendland—wavelet test family. For estimating the weak residual accurately and to provide a strong learning signal, we need to sample sufficiently many test functions. Evaluating them on the entire computational grid would be prohibitively costly. We therefore consider test functions which are locally supported, such that we can restrict computations to smaller patches. Wendland polynomials are a natural candidate meeting these requirements since they are compactly supported within unit radius and allow for efficient gradient computation. Here, we will describe the test functions in detail.

Each test function is drawn from a radially anisotropic family.

$$\psi_{c,\sigma,b}(x) = \underbrace{\left(1-r(x)\right)_+^4 \left(4r(x)+1\right)}_{\text{Wendland }C^2} \underbrace{\left(1-64\,b\,r(x)^4\right)}_{\text{optional wavelet}},$$

where $r(x) = \sqrt{\sum_j (x_j - c_j)^2/\sigma_j^2}$. Length-scales σ_j are uniformly sampled from the range $[\sigma_{\min} \Delta_j, \sigma_{\max} \Delta_j]$. By multiplying σ_{\min} and σ_{\max} with the grid spacing Δ_j of axis j, we obtain a parametrization that is intuitive to tune since the length-scales of the test functions are given in pixel units. Instead of also sampling center points c uniformly and independently from the full domain Ω , we instead start from the grid coordinates of the data, considering each grid point as one center. We found that this way of ensuring spatial coverage improves training efficiency. To still retain stochasticity, we apply i.i.d jitter to each center point. $b \sim \text{Ber}(p)$ randomly toggles a high-frequency wavelet factor that provides additional variability within the test functions and proved especially effective on noisy data.

To enforce $\psi_i|_{\partial\Omega}=0$, we multiply every test function by a *bridge mollifier*

$$m(\xi) = ((\xi - \xi_{\min})(\xi_{\max} - \xi))/(\xi_{\max} - \xi_{\min})^2,$$

applied per axis. This legitimises the application of integration-by-parts in the derivation of the weak forms. At training time we sample a set of test functions $\{\psi^{(i)}\}_{i=1}^{N_{\text{test}}}$ independently per residual evaluation and define the loss as

$$\mathcal{R}_{\text{weak}}(x,\alpha) = \frac{1}{N_{\text{test}}} \sum_{i} |\langle \mathcal{L}_{\alpha} x, \psi^{(i)} \rangle_{L^{2}}|^{2}.$$

D.4 ADJOINT MATCHING DETAILS

As mentioned in the main paper, we introduce a scaling coefficient κ that allows us to attenuate the magnitude of the noise variance. In this section, we provide a justification for using $\kappa > 0$ and lay out numerical heuristics used in the implementation.

Memoryless Noise. Domingo-Enrich et al. (2025) define a generative process with noise schedule $\sigma^2(t)$ to be memoryless, if and only if $\sigma^2(t) = 2\eta_t + \chi(t)$ with $\chi : [0,1] \to \mathbb{R}$ chosen such that

$$\forall t \in (0,1) \quad \lim_{t' \to 0} \beta_{t'} \, \exp\left(-\int_{t'}^t \frac{\chi(s)}{2\gamma_s^2} \, ds\right) = 0.$$

Specifically, they refer to $\sigma(t) = \sqrt{2\eta_t}$ as the memoryless noise schedule. In our setting of $\beta_t = t$ and $\gamma_t = 1 - t$, the memoryless noise schedule can be simplified to

$$\sigma(t) = \sqrt{\frac{2(1-t)}{t}}.$$

Lemma 1 (Scaling of memoryless noise). Consider a generative process as in 1 with $\beta_t = t$ and $\gamma_t = 1 - t$. For $0 \le \kappa < 1$, the schedule $\sigma^2(t) = (1 - \kappa) 2\eta_t$ is memoryless.

Proof. First, we consider the integral term. The desired $\sigma^2(t) = (1 - \kappa) 2\eta_t$ implies that

$$\chi(t) = -2\kappa \eta_t = -2\gamma_t^2 \frac{\kappa}{t(1-t)}.$$

With this, we can simplify the integral term:

$$-\int_{t'}^{t} \frac{\chi(s)}{2\gamma_s^2} ds = \kappa \int_{t'}^{t} \frac{1}{s(1-s)} ds = \kappa \int_{t'}^{t} \frac{1}{s} + \frac{1}{1-s} ds$$
$$= \kappa \left[\log(s) - \log(1-s) \right]_{t'}^{t}$$
$$= \kappa \left(\log \frac{t}{1-t} - \log \frac{t'}{1-t'} \right).$$

Thus, for an arbitrary fixed $t \in (0, 1)$, it holds that

$$\beta_{t'} \exp\left(-\int_{t'}^{t} \frac{\chi(s)}{2\gamma_{s}^{2}} ds\right) = t' \exp\left(\kappa \left(\log \frac{t}{1-t} - \log \frac{t'}{1-t'}\right)\right)$$

$$= t' \left(\frac{t}{1-t}\right)^{\kappa} \left(\frac{t'}{1-t'}\right)^{-\kappa}$$

$$= \underbrace{\left(\frac{t}{1-t}\right)^{\kappa}}_{\text{const}} \underbrace{t'^{1-\kappa}}_{\to 0} \underbrace{(1-t')^{\kappa}}_{\text{limited}} \xrightarrow[t'\to 0]{} 0.$$

Therefore, scaling down the noise schedule by a factor $1 - \kappa$ is justified and consistent with the theory provided in Domingo-Enrich et al. (2025).

Heuristics. Still, the term η_t causes numerical problems for t=0. Furthermore, it forces the control u to be close to zero for t close to 1. Following the original paper, we instead use

$$\eta_t = \frac{1 - t + h}{t + h},$$

where we choose h as the step size of our numerical ODE/SDE solver. This resolves infinite values and allows for faster fine tuning by letting the fine-tuned model deviate further from the base model close to t=1.

While κ is an effective tool to mitigate residual noise in final samples, increasing slows down training. As another way of improving sample quality without adding computational cost, we consider non-uniform time grids when sampling memoryless rollouts. We observe that the most critical phases of sampling are at the beginning, where noise magnitudes are the highest, and at the end, where final denoising happens. Therefore, we tilt the uniform grid towards the endpoints:

Let $S \in \mathbb{N}$ be the number of steps and define the uniform grid $t_i = i/S$ for $i = 0, \dots, S$. We tilt this grid toward the endpoints by first applying a cosine–ease mapping

$$g(t) = \frac{1}{2} (1 - \cos(\pi t)), \qquad t \in [0, 1],$$

which has $g(0)=0,\,g(1)=1.$ For a mixing parameter $q\in[0,1],$ the tilted times are the convex combination

$$\tilde{t}_i = (1-q) t_i + q g(t_i), \quad i = 0, \dots, S.$$

 $\{\tilde{t}_i\}$ is strictly increasing and distributes grid points more densely near t=0 and t=1 for q>0. For PDE experiments, we use q=0.9. This heuristic was not needed for latent-space models.

Loss Computation. As in the original paper, we do not compute the Adjoint Matching loss (Equation 4) for all simulated time steps, since the gradient signal for successive time steps is similar. Note that for solving the lean adjoint ODE, we do not need to compute gradients through the FM model, therefore we can compute the lean adjoint states efficiently but save computational ressources when computing the Adjoint Matching loss. Again, the last steps in the sampling process are most important for empirical performance. For that reason, we also compute the loss for a fraction of last steps K_{last} . Additionally, we sample K steps from the remaining time steps.

To ensure stable learning, we apply a clipping function to the loss to exclude noisy high-magnitude gradients from training. Empirically, the values provided in Domingo-Enrich et al. (2025) work well in our setting, i.e. we set the loss clipping threshold (LCT) as LCT $_x=1.6\,\lambda_x^2$ and LCT $_\alpha=1.6\,\lambda_\alpha^2$ respectively.

D.5 FULL TRAINING ALGORITHM

1026

1027

1028

1029

1030 1031

1032 1033

1034

1035

1036

1037

1039

1040 1041

1064

1067 1068

1069

1070

1071

1074 1075

1077

1078

1079

Algorithm 1 details the complete optimisation loop used in all experiments. Starting from the pretrained base flow $v^{\rm base}$ we attach two residual heads that (i) condition the state flow $v^{\rm ft}_{t,x}$ on the latent parameter α and (ii) predict the parameter flow $v^{\rm ft}_{t,\alpha}$. Their respective output layers are initialized to zero. The inverse predictor φ is first pre-trained on base samples and then frozen, providing surrogate target flows for α . Each epoch rolls both the base and fine-tuned trajectories on the tilted time grid, solves the lean-adjoint equation, and updates only the fine-tune parameters θ through the Adjoint-Matching loss.

Algorithm 1 Adjoint Matching on Joint Evolution

```
1043
                      Input: initialise core v^{\text{ft}} \leftarrow v^{\text{base}}, add residual-style output heads to get v_x^{\text{ft}}, v_\alpha^{\text{ft}}
1044
                        1: Pretrain \varphi based on x_1 samples generated with v_x^{\text{base}} by minimizing \mathcal{R}_{\text{weak}}(x_1, \varphi(x_1))
2: Freeze weights of v_x^{\text{base}} and \varphi, denote trainable parameters as \theta
1045
1046
                        3: for number of epochs do
1047
                                          x_0 \sim \mathcal{N}(0, I), \quad \alpha_0 \sim \mathcal{N}(0, I)
                        4:
1048
                        5:
                                          \overline{T} \leftarrow \text{GET\_TILTED\_TIME}(T)
                        6:
                                          for i \in |T| do
1049
                                                   t \leftarrow \bar{T}_i, \quad t^+ \leftarrow \bar{T}_{i+1}, \quad h \leftarrow t^+ - t
                        7:
1050
                                                   v_{x,t}^{\text{base}} \leftarrow v^{\text{base}}(x_t, t)
1051
                                                   \begin{array}{l} v_{\alpha,t}^{\text{base}} \leftarrow \frac{\hat{\alpha}_{t}^{\text{ft}} - \alpha_{t}}{1 - t} \\ v_{x,t}^{\text{ft}}, \ v_{\alpha,t}^{\text{ft}} \leftarrow v^{\text{ft}}(x_{t}, \alpha_{t}, v_{\alpha,t}^{\text{base}}, t) \\ x_{t}^{\text{base}} \leftarrow \text{SDE\_STEP}(x_{t}^{\text{base}}, v_{x,t}^{\text{base}}, \sigma(t), h) \end{array}
1052
                        9:
1053
                      10:
1054
                      11:
1055
                                                    \begin{aligned} x_{t^+}^{\text{ft}} \leftarrow & \text{SDE\_STEP}(x_t^{\text{ft}}, v_{x,t}^{\text{ft}}, \sigma(t), h) \\ \alpha_{t^+}^{\text{ft}} \leftarrow & \text{SDE\_STEP}(\alpha_t^{\text{ft}}, v_{\alpha,t}^{\text{ft}}, \sigma(t), h) \end{aligned} 
                      12:
1056
                      13:
1057
                      14:
1058
                                          \tilde{a} \leftarrow \texttt{SOLVE\_LEAN\_ADJOINT}(x^{\texttt{base}}, x^{\texttt{ft}}, \alpha^{\texttt{ft}})
                      15:
                                          \theta \leftarrow \text{GRADIENT\_DESCENT}(\theta, \mathcal{L}(\tilde{u}; X, \tilde{a}))
                      16:
                      17: end for
1061
                      18: return v^{\rm ft}, \varphi
1062
```

E DETAILS: EXPERIMENTAL RESULTS

Here we describe the fine-tuning configuration of the conducted experiments. Specifics about the base model FM training can be found in Appendix C. In all experiments, we use AdamW as the optimizer with a weight decay of 0.01 and all fine-tuning is conducted on a single NVIDIA L40S GPU.

E.1 DARCY

For the Darcy experiments, we use the hyperparameters listed in Table 5. Note that one epoch amounts to exactly one gradient descent step, meaning that we only fine-tune for 20 total steps. The parameters λ_x , λ_α and λ_f are varied in the experiments, see Figure 3 in the main text.

In panel (a) of Figure 3, we report a SSIM-based (Wang et al., 2004) metric for diversity, which is implemented as follows:

1080 1082

Table 5: Darcy fine-tuning hyperparameters.

1	083
1	084
1	085
1	086

1090

1093 1094 1095

1099

1100 1101

1102 1103

1104 1105 1106

1107

1108

1109

1110 1111

1113 1114 1115

1116 1117 1118

1119

1120

1121

1122

1123 1124 1125

1126 1127

1	1	2	8
1	1	2	(
1	1	3	(
1	1	3	
1	1	3	d

1133

Hyperparameter Value Time-tilting factor q 0.9 λ_x varying λ_{α} varying λ_f varying K_{last} 20 K20 Noise scaling κ 0.9 100 Sampling steps (per trajectory) Training epochs 20 Learning rate 0.00002 Batch size 15

Given a batch $\{\alpha_i\}_{i=1}^B$ of single-channel parameter maps scaled to [0,1], we define diversity as the mean complement of the pairwise structural similarity index (SSIM):

$$\mathcal{D}_{\text{SSIM}}(\{\alpha_i\}_{i=1}^B) \ = \ \frac{1}{{B \choose 2}} \sum_{1 \le i < j \le B} \Big(1 - \text{SSIM}(\alpha_i, \alpha_j)\Big), \qquad \alpha_i \in [0, 1]^{H \times W}.$$

With SSIM $\in [0,1]$ (data range = 1), $\mathcal{D}_{SSIM} \in [0,1]$ and larger values indicate greater sample diversity.

E.2 GUIDANCE

For generating the guidance results, we use the same model as in the Darcy experiments to further highlight that we can infer functional joint distributions when starting from noisy observations. Instead of the usual Euler stepping, here we use a Heun sampler following Huang et al. (2024). While this is more expensive, since we need to differentiate through two forward passes of our model to obtain the guidance gradient, it empirically improved faithfulness to the sparse observations significantly. Note, however, that we only guide on sparse observations and not on PDE residuals. The full sampling procedure is show in Algorithm 2. In our experiments, we use 100 steps for sampling and guidance scales of $\gamma_x = \gamma_\alpha = 0.75$.

E.3 ELASTICITY

Fine-tuning in the elasticity experiment is more challenging than the Darcy denoising experiment, which is why we increase the number of fine-tuning epochs to 100. We fix $\lambda_x = \lambda_\alpha = 100k$ and set the regularization $\lambda_f = 1k$. Other parameters are the same as in the Darcy experiments.

We compare our fine-tuning approach with the inference-time correction method ECI presented in Cheng et al. (2024). Our implementation of the sampling correction method is given in Algorithm 3. For ECI, we set M=5 and use 1000 steps when sampling, compared to 100 steps used when sampling from the fine-tuned model. The reported residual heat maps in the main paper show one test function per grid point, where the magnitude indicates local error without aggregating across test functions.

E.4 NATURAL IMAGES: RECOLORING

For natural-image experiments, we adopt the class-conditional Latent Flow Matching (LFM) model of Dao et al. (2023) trained on ImageNet-1k (Deng et al., 2009) with a DiT backbone (Peebles & Xie, 2023). We fix an ImageNet class label y to condition the generator and hold a global text prompt c to define the fine-tuning direction. As a reward we use PickScore v1, a CLIP-H/14-based preference scorer trained on the Pick-a-Pic dataset, that evaluates image-text compatibility via cosine similarity in the CLIP embedding space (Kirstain et al., 2023; Radford et al., 2021; Cherti et al., 2023; Schuhmann et al., 2022), implemented with TRANSFORMERS (Wolf et al., 2020).

```
1134
                 Algorithm 2 Guided Heun sampler with sparse observations
1135
                 Require: initial state x_0; initial parameter \alpha_0; observed targets \alpha_{obs}; index set \mathcal{I}; steps S; guidance
1136
                         scales \gamma_x, \gamma_\alpha; base fields v_t^x, v_t^\alpha; fine-tuned joint field v_t^{\mathrm{II}}
1137
                 Ensure: trajectories \{x^{(i)}\}_{i=0}^{S} and \{\alpha^{(i)}\}_{i=0}^{S}
1138
                   1: x^{(0)} \leftarrow x_0
1139
                   2: \alpha^{(0)} \leftarrow \alpha_0
1140
                   3: t_i \leftarrow i/S for i = 0, \dots, S
1141
                   4: h \leftarrow 1/S
1142
1143
                   5: for i = 0 to S - 1 do
1144
                                Predictor (Euler):
1145
                                v^x \leftarrow v_{t}^x(x^{(i)})
                                v_{\text{base}}^{\alpha} \leftarrow v_{t_i}^{\alpha}(x^{(i)}, \alpha^{(i)}, v^x)
1146
                                (\tilde{v}^x, \tilde{v}^\alpha) \leftarrow v_t^{\text{ft}}(x^{(i)}, \alpha^{(i)}, v_{\text{base}}^\alpha)
1148
                                \hat{x} \leftarrow x^{(i)} + h \, \tilde{v}^x
                 10:
1149
                                \hat{\alpha} \leftarrow \alpha^{(i)} + h \, \tilde{v}^{\alpha}
1150
1151
                 12:
                                if i < S - 1 then
1152
                 13:
                                        Corrector (Heun):
1153
                 14:
                                        v_+^x \leftarrow v_{t_{i+1}}^x(\hat{x})
1154
                                       v_{\text{base},+}^{\alpha} \leftarrow v_{t_{i+1}}^{\alpha}(\hat{x}, \, \hat{\alpha}, \, v_{+}^{x})
                 15:
1155
                                       (\tilde{v}_+^x, \tilde{v}_+^\alpha) \leftarrow v_{t+1}^{\text{ft}}(\hat{x}, \hat{\alpha}, v_{\text{base},+}^\alpha)
                 16:
1156
                                       One-step terminal extrapolation:
                 17:
1157
                                       \hat{\alpha}_1 \leftarrow \hat{\alpha} + (1 - t_{i+1}) \, \tilde{v}_+^{\alpha}
                 18:
                                       Sparse-observation loss: L \leftarrow \frac{1}{|\mathcal{I}|} \sum_{j \in \mathcal{I}} \left\| \alpha_{\text{obs}}[j] - \hat{\alpha}_1[j] \right\|_2^2
1158
                 19:
1159
                 20:
1160
                                       Heun average update: x^{(i+1)} \leftarrow x^{(i)} + \frac{h}{2} \left( \tilde{v}^x + \tilde{v}_+^x \right)
1161
                 21:
1162
1163
                                       \alpha^{(i+1)} \leftarrow \alpha^{(i)} + \frac{\bar{h}}{2} \left( \tilde{v}^{\alpha} + \tilde{v}_{+}^{\alpha} \right)
1164
                 23:
1165
                                       Decaying guidance:
1166
                                       s \leftarrow \sqrt{1 - i/S}
                 25:
1167
                                        \begin{aligned} x^{(i+1)} &\leftarrow x^{(i+1)} - s \, \gamma_x \, \nabla_{x^{(i)}} L \\ \alpha^{(i+1)} &\leftarrow \alpha^{(i+1)} - s \, \gamma_\alpha \, \nabla_{\alpha^{(i)}} L \end{aligned} 
                 26:
1168
                 27:
1169
                 28:
1170
                                       Last step (no correction/guidance):
                 29:
1171
                                       x^{(i+1)} \leftarrow \hat{x}
                 30:
1172
                                       \alpha^{(i+1)} \leftarrow \hat{\alpha}
                 31:
1173
                                end if
                 32:
1174
                 34: return \{x^{(i)}\}_{i=0}^{S}, \{\alpha^{(i)}\}_{i=0}^{S}
1175
1176
1177
1178
```

Concretely, we maximize $R\left(T\left(D(z),\ \alpha\right),\ c\right)$, where D is the latent-space decoder and $T(\cdot,\alpha)$ is a parametric per-pixel recoloring with coefficients α . The recoloring operates in logit space to avoid saturation: with $D(z)=x:\Omega\to(0,1)^3$,

1179

1180

1181

1182 1183

1184

1185 1186

1187

$$x_{\varepsilon}(\xi) = \operatorname{clip}(x(\xi), \varepsilon, 1 - \varepsilon), \qquad \ell(\xi) = \log \frac{x_{\varepsilon}(\xi)}{1 - x_{\varepsilon}(\xi)} \in \mathbb{R}^{3},$$

we apply a residual $r(\xi;\alpha)=\alpha\,\Phi_D^{\mathrm{bias}}\!\big(x(\xi)\big)$ built from RGB monomials up to total degree D and map back via $y(\xi;\alpha)=\sigma(\ell(\xi)+r(\xi;\alpha))$. Equivalently, for channel $c\in\{R,G,B\}$,

$$y_c(\xi;\alpha) = \sigma \Big(\ell_c(\xi) + \sum_{m=1}^{M} \alpha_{cm} \,\phi_m(x(\xi)) + \alpha_{c,0}\Big), \quad \phi_m \in \Phi_D.$$

1205

1206

1207

1208

1209

1210 1211

1212

1213

1214 1215

1216

1217

1218

1219

1221

1222

1223 1224

1225

1226

1227

1228

1229

1230

1231

1232 1233 1234

1235

1236

1237

1238 1239

1240

1241

22: **return** $\{x^{(i)}\}_{i=0}^{S}$

Table 6: Elasticity fine-tuning hyperparameters.

Hyperparameter	Value
Time–tilting factor q	0.9
λ_x	100k
λ_{lpha}	100k
λ_f	1k
$K_{ m last}$	20
K	20
Noise scaling κ	0.9
Sampling steps (per trajectory)	100
Training epochs	100
Learning rate	0.00002
Batch size	15

Algorithm 3 ECI-style sampling with boundary correction

Require: initial state x_0 ; steps S; inner correction iterations M; model drift $v_t(\cdot)$; correction operator $C(\cdot)$; noise sampler $Noise(\cdot)$

```
Ensure: trajectory \{x^{(i)}\}_{i=0}^{S}
 1: x^{(0)} \leftarrow x_0
 2: t_i \leftarrow i/S for i = 0, \dots, S
 3: for i = 0 to S - 1 do
            t \leftarrow t_i
 4:
            t^{\text{next}} \leftarrow t_{i+1}
 5:
            Inner ECI corrections at fixed t:
 6:
            \tilde{x} \leftarrow x^{(i)}
 7:
            for j = 1 to M do
 8:
 9:
                   v \leftarrow v_t(\tilde{x})
10:
                   x_{\text{os}} \leftarrow \tilde{x} + (1-t)v
                                                                                                                              one-step Euler update
                  x_{\text{corr}} \leftarrow \mathcal{C}(x_{\text{os}})
11:
                                                                                                                      apply boundary correction
                   \eta \leftarrow \text{Noise}(\text{shape of } x_0)
12:
                   \hat{x} \leftarrow (1-t) \eta + t x_{\text{corr}}
13:
                                                                                                                      stochastic convex blending
            end for
14:
            Final correction and roll-forward to t^{next}:
15:
16:
            v \leftarrow v_t(\tilde{x})
17:
            x_{os} \leftarrow \tilde{x} + (1-t)v
18:
            x_{\text{corr}} \leftarrow \mathcal{C}(x_{\text{os}})
            \eta \leftarrow \text{Noise}(\text{shape of } x_0)
19:
            x^{(i+1)} \leftarrow (1 - t^{\text{next}}) \eta + t^{\text{next}} x_{\text{corr}}
20:
21: end for
```

This residual parameterization is identity at initialization ($\alpha=0$) and provides a low-dimensional, channel-coupled appearance pathway that can disentangle *content* from *presentation* or reach colorings underrepresented by the base LFM. It is related to CNN-predicted polynomial color transforms (quadratic in Chai et al. (2020)), whereas we use cubic polynomials and learn solely via the reward.

The parameter predictor φ here maps from latent feature tensors z to parameters α . First, a compact scene descriptor is extracted by passing z through convolutional transformations, aggregating information across multiple spatial scales via downsampling and global pooling, and enriching it with low-order channel statistics of z (e.g., moments). The concatenated descriptor is projected into

 a fixed embedding space, refined by a lightweight pre–LayerNorm MLP with a residual connection, and finally mapped by a linear head to the recoloring coefficients α .

Building on the inverse predictor, we augment the base U–Net generator with two lightweight residual heads that couple image dynamics and color–parameter evolution. First, the *image path* predicts the base drift $v_{t,x}^{\text{base}}(x,t)$. Then, we form a compact α -token by flattening the current color parameters α and mapping them through a small MLP. This token is broadcast to a k_{α} -channel feature map and concatenated with $(x, v_{t,x}^{\text{base}})$ into a shallow U–Net "correction" that outputs an additive refinement, yielding

$$v_{t,x}^{\mathrm{ft}} \ = \ v_{t,x}^{\mathrm{base}} \ + \ \mathcal{U}_{x}\!\big(x,\, v_{t,x}^{\mathrm{base}},\, \mathrm{map}(\alpha),\, t\big).$$

Second, the parameter path updates the polynomial recoloring coefficients by a residual on top of the baseline parameter field $v_{t,\alpha}^{\text{base}}$. Here, a dedicated α -projection module mirrors the inverse predictor: multi-scale pooled conv features are fused—at the token level—with tokens from α and $v_{t,\alpha}^{\text{base}}$ via a small fusion MLP, AdaLN/FiLM modulation of a head token, and a light SE rescaling of conv features. The head then predicts a delta added to the baseline,

$$v_{t,\alpha}^{\text{ft}} = v_{t,\alpha}^{\text{base}} + \mathcal{U}_{\alpha}(x, \alpha, v_{t,\alpha}^{\text{base}}).$$

All correction heads are zero-initialized at their final projections, so fine-tuning starts from the base behavior and departs smoothly as training progresses. Table 7 lists the parameters of the correction U–Nets. In total, the modified architecture adds $\approx 9M$ parameters to the $\approx 130M$ parameters of the base DiT backbone.

Table 7: Correction head hyperparameters (image and parameter paths).

Hyperparameter	Value
k_{α}	16
U-Net base widths	[96, 128, 160, 192]
U-Net embedding	96
U-Net hidden lift	256
Attention stages	[2, 3]
Attention heads	8

With the LFM model, we use 40 steps when sampling and use the same training specifications as the original Adjoint Matching paper (Domingo-Enrich et al., 2025). Again, we first pretrain the predictor φ , and then perform joint fine-tuning with our Adjoint Matching formulation. Fine-tuning is performed for 100 epochs with a batch size of 15.

F ADDITIONAL RESULTS

F.1 DENOISING

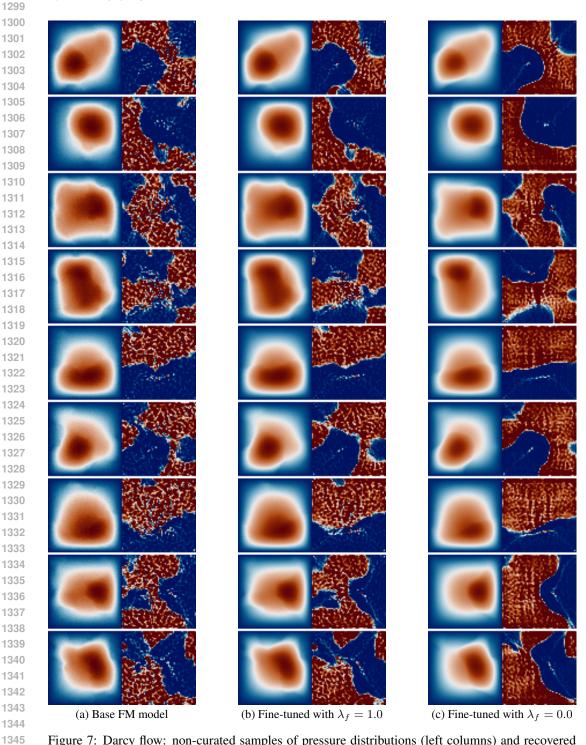


Figure 7: Darcy flow: non-curated samples of pressure distributions (left columns) and recovered permeability fields (right columns). Each row was generated using the same initial noise across the three models. Color scales are normalized per row for the pressure distributions. For the base model, permeabilities are obtained with the pre-trained inverse predictor.

ELASTICITY

Figure 8: Non-curated samples from the elasticity experiment, where fine-tuning has to scale down the lower boundary.

(b) Fine-tuned with model

(a) Base FM model

F.3 GUIDANCE

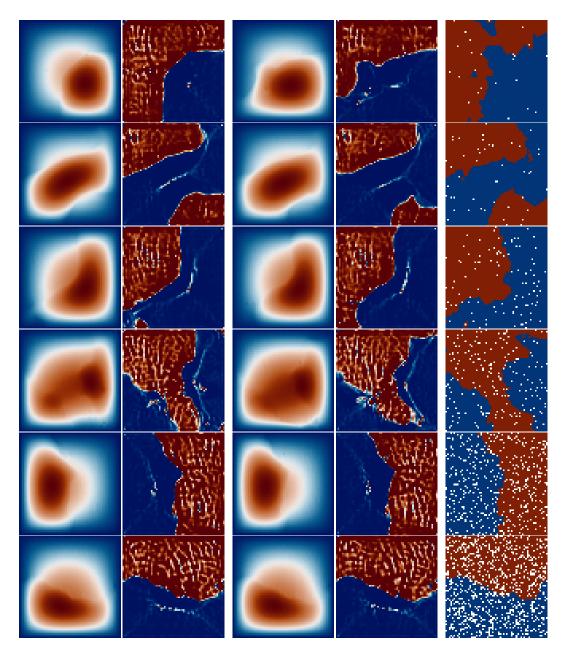


Figure 9: Guided samples with an increasing number of given observations, specifically [25, 50, 100, 200, 300, 400, 500, 750, 1000]. For each number of conditioning points, we generate two samples from independent noise and condition on the same sparse samples, indicated as white markers in the right column. As expected, with more points both x and α become more constraint and less diverse.

F.4 NATURAL IMAGES: RECOLORING

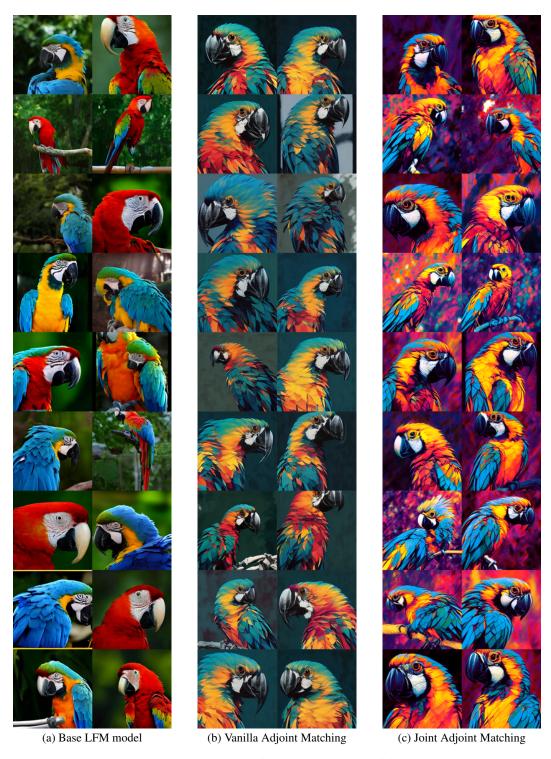


Figure 10: Non-curated independent samples from LFM model conditioned on class "macaw" and using guidance scale 4.0. Models were fine-tuned to maximize PickScore using the prompt "close-up pop art of a macaw parrot".

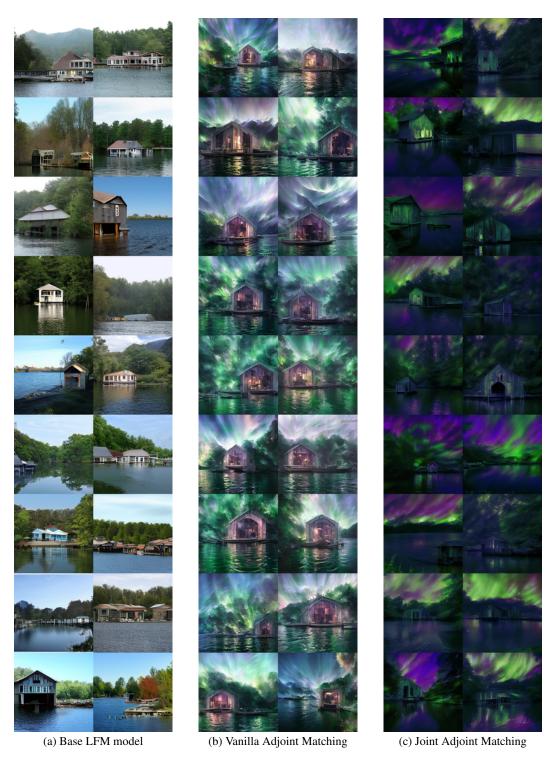


Figure 11: Non-curated independent samples from LFM model conditioned on class "boathouse" and using guidance scale 4.0. Models were fine-tuned to maximize PickScore using the prompt "boathouse with green and purple curtains of northern lights." Our joint model is able to generate the colors demanded in the prompt while retaining diversity in the generated boathouses.