

Learning in PINNs: Phase transition, total diffusion, and generalization

Sokratis J. Anagnostopoulos^{a,1}, Juan Diego Toscano^{b,1},
Nikolaos Stergiopoulos^a, George Em Karniadakis^{b,c}

^a*Laboratory of Hemodynamics and Cardiovascular Technology,
EPFL, Lausanne, 1015, VD, Switzerland*

^b*Division of Applied Mathematics, Brown University, Providence, 02912, RI, USA*

^c*School of Engineering, Brown University, Providence, 02912, RI, USA*

Abstract

We investigate the learning dynamics of fully-connected neural networks through the lens of gradient signal-to-noise ratio (SNR), examining the behavior of first-order optimizers like Adam in non-convex objectives. By interpreting the drift/diffusion phases in the information bottleneck theory, focusing on gradient homogeneity, we identify a third phase termed “total diffusion”, characterized by equilibrium in the learning rates and homogeneous gradients. This phase is marked by an abrupt SNR increase, uniform residuals across the sample space and the most rapid training convergence. We propose a residual-based re-weighting scheme to accelerate this diffusion in quadratic loss functions, enhancing generalization. We also explore the information compression phenomenon, pinpointing a significant saturation-induced compression of activations at the total diffusion phase, with deeper layers experiencing negligible information loss. Supported by experimental data on physics-informed neural networks (PINNs), which underscore the importance of gradient homogeneity due to their PDE-based sample inter-dependence, our findings suggest that recognizing phase transitions could refine ML optimization strategies for improved generalization.

Keywords: Information bottleneck, PINNs convergence/generalization, gradient SNR transition, total diffusion phase, residual homogeneity

¹The first two authors contributed equally to this work

1. Introduction

1.1. Phase transitions in deep learning

The optimization process in deep learning can vary significantly in terms of smoothness and convergence rate, depending on various factors such as the complexity of the model, the quality/quantity of the data or the loss landscape characteristics. However, for non-convex problems this process has often been observed to be far from smooth and steady; instead it is rather dominated by discrete, successive phases. Recent studies have shed light on several key aspects influencing these phases and the overall optimization dynamics [1–10].

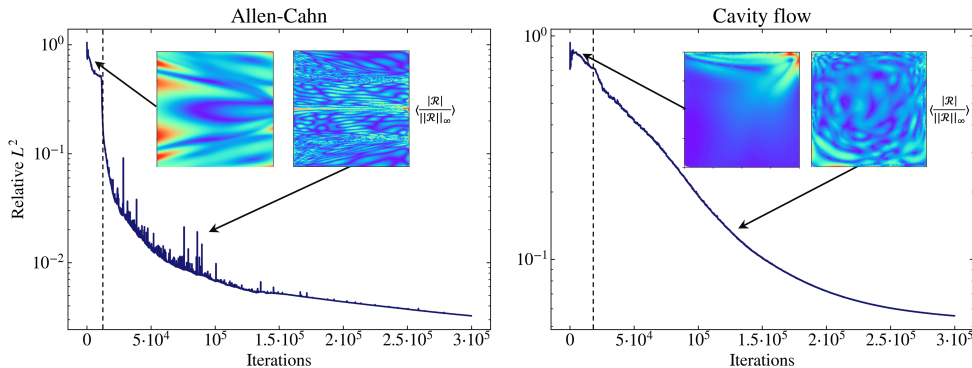


Figure 1: **Phase transition in PINNs:** The test error between the prediction and the exact solution converges faster after total diffusion (dashed lines), which occurs with an abrupt phase transition defined by homogeneous residuals. Although the convergence starts during the onset of the diffusion phase, the optimal training performance is met when the gradients of different batches become equivalent, indicating a general agreement on the direction of the optimizer steps (total diffusion).

The importance of gradient noise in escaping local optima of non-convex optimization has been explored, demonstrating its role in guaranteeing polynomial time convergence to a global optimum [1]. The authors of the same work suggest the existence of a phase transition for a perturbed gradient descent GD algorithm, from escaping local optima to converging to a global solution as the artificial noise decreases. In a later work, a phenomenon called “super-convergence” has been highlighted, where models trained with a two-phase cyclical learning rate may lead to improved regularization balance and generalization [2]. Furthermore, recent investigations have discovered a two-phase learning regime for full-batch gradient descent (GD), characterized by distinct behaviors [3]. During the “lazy” phase ($\eta < 2/\lambda_0$), the model behaves linearly about its initial parameters [4]. However, optimal performance is usually achieved during the “catapult” phase ($2/\lambda_0 < \eta < \eta_{max}$), which

exhibits instabilities due to the increased curvature. The findings of this study are also supported using the NTK method [11], which has been proven to be a very useful tool for studying deep learning dynamics in the infinite width limit. Other researchers have demonstrated that neural networks tend towards a self-organized critical state, characterized by scale invariance, where both trainable and non-trainable parameters are driven to a stochastic equilibrium [5], which is also present in many biological systems.

An intriguing study showed empirically that contrary to the theoretical clues on stability, GD very often trains at the “edge of stability”, which defines a regime of hypercritical sharpness (max Hessian eigenvalue), that hovers just above $2/\eta$, during which there is non-monotonic convergence [6], present also for adaptive optimizers [7]. The authors also posed questions regarding the mystery of the progressive sharpening preceding the “edge of stability” phase. A follow-up paper explains this equilibrium state using the cubic Taylor expansion, showing that GD has an inherent self-stabilizing mechanism, which decreases the sharpness whenever extreme oscillations are present [8]. Analyzing the dynamics of popular optimizers, such as Adam [12], has also revealed insights into the non-convex learning process. For instance, the concept of rotational equilibrium studies a steady state of angular network parameter updates, offering a deeper understanding of the phase transitions experienced by the weight vector during optimization [9]. Additionally, the “double-descent” phenomenon has emerged as a notable discovery, revealing a critical range of network width that negatively impacts generalization performance, with implications for the success of large language models [10], while this transition can also be met epoch-wise. Collectively, these findings emphasize the complexity of deep learning optimization and the existence of discrete learning phases, offering valuable insights for improving optimization algorithms and understanding deep neural network behavior. A possible speculation could be that most of these phase transitions largely coincide, but this investigation is beyond the scope of our study.

1.2. *Physics-informed learning*

Physics-informed neural networks (PINNs) have emerged as an alternative to classical numerical methods for solving general non-linear partial differential equations (PDEs) [13]. Unlike traditional methods which use numerical integration to propagate the solution from initial or boundary conditions, PINNs approximate the PDE solution through a gradient-based optimization process. This optimization process involves minimizing a constrained loss function designed to concurrently satisfy the governing physical laws and initial/boundary conditions, while fitting any available experimental data (inverse solution) [14–18].

PINNs have gained increased popularity due to their superiority in inverse problems, motivating researchers to develop multiple methods to ease the optimization process. Some of these studies focus on the multi-layer perceptron structure via domain decompositions [19], architecture reformulations [20–24], input dimension expansions [25, 26], sequential learning [27, 28], or adaptive activation functions [29]. Similarly, other studies attempt to simplify the problem by reducing the number of loss terms via exact boundary imposition [30–32] or PDE reformulations [33, 34].

On the other hand, other approaches have focused on the imbalances of the constrained optimization. Maintaining a global balance (i.e., between each loss term) and local balance (i.e., between training samples) is crucial to ensure that all constraints are adequately minimized across the domain. Some researchers have attempted to deal with these imbalances by adaptively resampling points in critical areas [35–38] while other studies involve re-weighting methods, which assign adaptive global or local multipliers. In general, global multipliers modify the average contribution of each loss term [20, 39–41] while local multipliers adjust the local contribution of each training sample [34, 42–47].

1.3. Information bottleneck method

Information bottleneck (IB) theory provides an interpretation of neural network training and performance from the information theory perspective. It formulates a methodology for estimating the optimal trade-off between compression and prediction in supervised learning, and constructs an idealized principle for creating a compressed or “bottlenecked” representation T (layer activations) of an input variable \mathcal{X} , that retains as much information as possible about an output variable \mathcal{Y} [48, 49]. The theory utilizes the concept of mutual information I , which measures the amount of information obtained about one random variable by observing another, suggesting that the optimal model representations retain all the relevant information about the output while discarding irrelevant details about the input, hence creating an “information bottleneck”. A remarkable consequence of the IB is that deep learning is governed by two discrete phases: the fitting and the diffusion [50–52], which are separated by a phase transition, captured by the signal-to-noise ratio (SNR) of the gradients. The theory suggests that most of the learning occurs during the slow diffusive phase, where the model is able to generalize well.

Although the IB theory has remained a subject of increasing interest over the years, there have been ongoing debates regarding some of its claims, mainly the compression phenomenon and its contribution to generalization [53–55]. In addition, since providing an accurate estimate of the mutual information is non-trivial, there has been ongoing research on different methodologies and how they affect the

interpretation of the process [53, 56]. A recent study provides a clearer view on quantifying I and supports that the compression mechanism is driven by progressive clustering of T , which is also present in deterministic DNNs [55]. Finally, in the same study, the authors justify the binning approximation as a measure of clustering and suggest that compression is not necessarily related to generalization.

1.4. Main contributions

The aim of this work is to provide a better understanding of the deep learning process, with a focus on PINNs applications. First, we analyze the training dynamics of full-batch GD using the gradient SNR and demonstrate how the equilibrium of the iteration-wise SNR used in Adam does not guarantee uniform learning across the sample space, which can lead to instabilities, overfitting and bad generalization. This weakness caused by the quadratic loss averaging can have detrimental effects on PINNs, where the samples (collocation points) should ideally converge simultaneously to achieve a low relative error with the exact solution.

We provide novel evidence for the existence of phase transitions proposed by the IB theory (fitting/diffusion) for PINNs trained with Adam, and reveal a third phase which we call “total diffusion”, where optimal convergence occurs (Fig. 1). We show that this phase is tightly linked with gradient homogeneity, suggesting that PINNs learn best when the residuals are diffused in a highly homogeneous state, where there is an equivalence of batches and a stable equilibrium of the optimizer can be reached. We propose a re-weighting technique based on the moving average of the residuals named residual-based attention (RBA) [42] to further support the hypothesis that homogeneous residuals lead to better generalization. We compare the performance of the vanilla and the RBA models, showing that the latter can achieve total diffusion faster, leading to better generalization for most cases.

Finally, we demonstrate how the SNR and the residual diffusion coincide with information compression, which we interpret as the “binarization” of the activations. During the abrupt transition into “total diffusion,” the weights rapidly increase, causing saturation of the neuron activations. We also deploy a relative-binning approach and observe that although there is no major information loss for deeper layers, there is a hierarchy for the information carried by each layer, which is consistent with the IB theory. We further demonstrate that it is the middle layers that saturate the most during the “binarization” process, a clue reminiscent of an encoder-decoder process within the fully connected architecture.

2. Analysis

2.1. Gradient signal and noise for SGD

Let \mathcal{X} be a dataset of n training samples $x_i|_{i=1}^n$ and $\mathcal{R}(f(x_i, \theta), y_i)$ be the residuals between the prediction of a model f with parameters θ and the correct labels y_i . The minimization problem we aim to solve using the quadratic loss is defined as:

$$\min_{\theta \in \Theta} \mathcal{L}(\theta) = \frac{1}{n} \sum_{i \in \mathcal{X}} \mathcal{R}(f(x_i; \theta), y_i)^2. \quad (1)$$

A common gradient-based optimization method for problems of this form is Stochastic Gradient Descent (SGD), which updates the model parameters at every iteration t , using mini-batches \mathcal{B} of size m as:

$$\theta^{t+1} = \theta^t - \eta \cdot \nabla_{\theta} \mathcal{L}_{\mathcal{B}}, \quad (2)$$

where

$$\nabla_{\theta} \mathcal{L}_{\mathcal{B}} = \frac{1}{m} \sum_{i \in \mathcal{B}} \nabla_{\theta} \mathcal{R}(f(x_i, \theta), y_i)^2, \quad (3)$$

where $\mathcal{B} \subset \mathcal{X}$ and η is the learning rate. Since each batch gradient $\nabla_{\theta} \mathcal{L}_{\mathcal{B}}$ is an unbiased estimator of $\nabla_{\theta} \mathcal{L}$, it can be expressed as a composition of the expected gradient and some noise ϵ :

$$\nabla_{\theta} \mathcal{L}_{\mathcal{B}} = \nabla_{\theta} \mathcal{L} + \epsilon_{\mathcal{B}} = \mathbb{E}[\nabla_{\theta} \mathcal{L}_{\mathcal{B}}] + \epsilon_{\mathcal{B}}, \quad (4)$$

where $\nabla_{\theta} \mathcal{L}$ is the true gradient and $\epsilon_{\mathcal{B}}$ is a zero-mean ($\mathbb{E}[\epsilon_{\mathcal{B}}] = 0$) variable, measuring the deviation of the batch gradient from the true gradient.

Thus, Eq. 2 can be expanded as:

$$\theta^{t+1} = \theta^t - \eta \cdot (\mathbb{E}[\nabla_{\theta} \mathcal{L}_{\mathcal{B}}] + \epsilon_{\mathcal{B}}), \quad (5)$$

where $\eta \cdot \mathbb{E}[\nabla_{\theta} \mathcal{L}_{\mathcal{B}}]$ is the deterministic direction and $\eta \cdot \epsilon_{\mathcal{B}}$ is the stochastic direction that the optimizer will take when it performs a step on each batch \mathcal{B} . The variance of $\nabla_{\theta} \mathcal{L}_{\mathcal{B}}$ is also calculated from Eq. 4 as follows:

$$\begin{aligned} \nabla_{\theta} \mathcal{L}_{\mathcal{B}} - \mathbb{E}[\nabla_{\theta} \mathcal{L}_{\mathcal{B}}] &= \epsilon_{\mathcal{B}} \\ \mathbb{E}[\nabla_{\theta} \mathcal{L}_{\mathcal{B}} - \mathbb{E}[(\nabla_{\theta} \mathcal{L}_{\mathcal{B}})]^2] &= \mathbb{E}[\epsilon_{\mathcal{B}}^2] \\ \text{Var}(\nabla_{\theta} \mathcal{L}_{\mathcal{B}}) &= \mathbb{E}[\epsilon_{\mathcal{B}}^2]. \end{aligned} \quad (6)$$

In this study, we measure the generalization performance of the model with the relative L^2 error (test error), which is calculated as:

$$\text{Relative } L^2 = \frac{\|f(x; \theta) - u_e\|_2}{\|u_e\|_2}, \quad (7)$$

where $f(x; \theta)$ is the prediction of the model and u_e is the exact solution, unseen by the model. Note that for PINNs, the term to be minimized in Eq. 1 is $\mathcal{R}(f(x_i; \theta))^2$, where \mathcal{R} are the PDE residuals.

In IB theory, a usual metric to quantify the training dynamics is the batch-wise Signal-to-Noise Ratio (SNR), which is defined as:

$$\text{SNR} = \frac{\|\mu\|_2}{\|\sigma\|_2} = \frac{\|\mathbb{E}[\nabla_{\theta} \mathcal{L}_{\mathcal{B}}]\|_2}{\|\text{std}[\nabla_{\theta} \mathcal{L}_{\mathcal{B}}]\|_2} \quad (8)$$

where $\|\mu\|_2$ and $\|\sigma\|_2$ are the L^2 norms of the batch-wise mean and std of $\nabla_{\theta} \mathcal{L}_{\mathcal{B}}$ for all network parameters. The numerator represents the gradient signal, which drives the optimizer to a consistent direction, while the denominator is the gradient noise or variation from the learning signal. Based on IB theory, the task of efficient neural network training lies in the optimal extraction of relevant information (signal) from irrelevant or misleading information (noise). In this work, we study the training dynamics of Adam trained with the entire dataset \mathcal{X} , and calculate $\nabla_{\theta} \mathcal{L}_{\mathcal{B}}$ without performing an update of θ for each \mathcal{B} so that we can investigate the batch-wise behavior on the same iteration t . The term which dominates in Eq. 5 defines the optimization regime, which is stochastic for $\text{SNR} < 1$ or deterministic for $\text{SNR} > 1$ (Fig. 2). The IB theory states that neural gradients undergo a phase transition from high to low SNR, which corresponds to a fitting phase of high gradient agreement and a diffusion phase when the gradient fluctuations dominate the signal and achieve a compressed representation. However, in all our case studies, a third phase also emerges for prolonged training (more than 10^4 iterations), which we call “total diffusion”. In this phase, we observe a sudden jump in the SNR due to gradient homogeneity, which implies a high degree of uniformity in the samples and their corresponding sensitivity to the optimizer steps. While the neural network starts converging from the diffusion phase, as shown in previous studies, we suggest that it is during total diffusion when the neural network can converge optimally. Therefore, diffusion may be considered an exploratory stochastic phase, while total diffusion exploits a consistent direction. Note that a sudden increase in the SNR is also met in the literature [50, 52, 53] although it is not explicitly addressed, probably due to the shorter training periods.

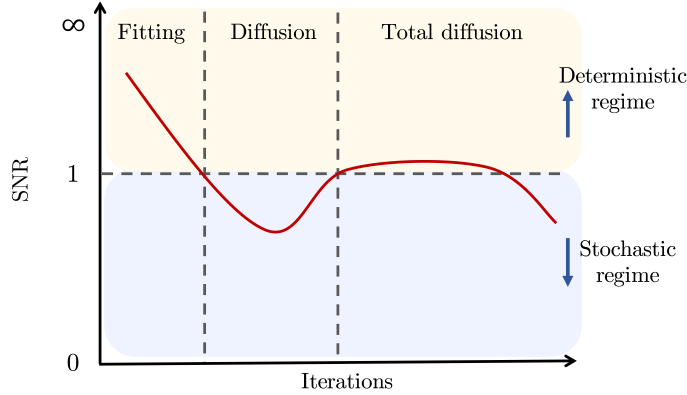


Figure 2: **Gradient-based optimization regimes:** Indicative SNR training curve at each full-batch iteration. For $\text{SNR} \gg 1$, the deterministic term dominates, while for $\text{SNR} \ll 1$, each step becomes more stochastic. The first two stages of learning are defined as “fitting” ($\text{SNR} \gg 1$) and “diffusion” ($\text{SNR} < \mathcal{O}(1)$). The “total diffusion” starts when the batch gradients are approximately equivalent, met with an abrupt increase of the SNR, which typically stabilizes above $\mathcal{O}(1)$. During the final stages, SNR decreases as the signal (numerator) tends to zero and some noise (denominator) persists.

Figure 3 shows three indicative cases where $\text{SNR} > 1$, $\text{SNR} = 0$ and $\text{SNR} = 1$, for a network of three parameters θ_j and three samples x_i . When $\text{SNR} > 1$, the directions of samples x_i agree, resulting in a deterministic $\nabla_{\theta}\mathcal{L}$ of large magnitude. When $\text{SNR} = 0$, the vectors cancel out, indicating convergence. Finally, for $\text{SNR} = 1$, there is an equilibrium between determinism and stochasticity in the direction of $\nabla_{\theta}\mathcal{L}$, for which the magnitude is non-zero. Thus, SNR measures directional and magnitude agreement.

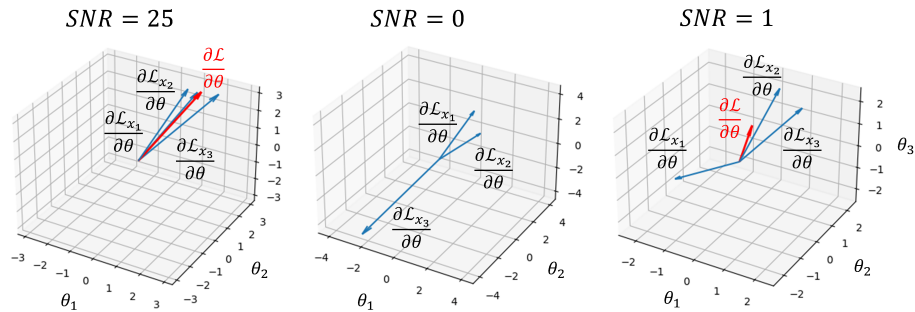


Figure 3: **Batch-wise SNR directions:** Indicative directions in the parameter space for each SNR case. For $\text{SNR} > 1$, there is an agreement of directions among samples x_i , resulting in a deterministic $\frac{\partial\mathcal{L}}{\partial\theta}$ of large magnitude. When $\text{SNR} = 0$, the directions cancel out, indicating convergence to a local minimum. Finally, for $\text{SNR} = 1$, there is an equilibrium between determinism and stochasticity in the direction of $\frac{\partial\mathcal{L}}{\partial\theta}$, for which the magnitude is non-zero.

To provide a more intuitive understanding of the SNR, we can expand Eq. 8 as:

$$\text{SNR} = \frac{\|\mathbb{E}[\nabla_{\theta}\mathcal{L}_{\mathcal{B}}]\|_2}{\|\sqrt{\text{Var}(\nabla_{\theta}\mathcal{L}_{\mathcal{B}})}\|_2} = \frac{\|\mathbb{E}[\nabla_{\theta}\mathcal{L}_{\mathcal{B}}]\|_2}{\|\epsilon_{\mathcal{B},rms}\|_2}, \quad (9)$$

where $\epsilon_{\mathcal{B},rms}$ is the root mean square of the batch deviations $\epsilon_{\mathcal{B}}$. Therefore, from Eq. 9 when $\text{SNR} < 1$, the magnitude of the typical stochastic component $\|\epsilon_{rms}\|_2$ is greater than the magnitude of the deterministic component $\|\mathbb{E}[\nabla_{\theta}\mathcal{L}_{\mathcal{B}}]\|_2$ towards the update of θ . This implies that the gradient noise will have a more pronounced effect than the signal in the direction of the parameter updates. Especially for cases where $\text{SNR} \ll 1$, the erratic behavior of the gradient due to the dominating noise among samples can be a characteristic of highly non-convex regions, caused by sample-disagreement, where the optimizer is unable to make meaningful progress. Similarly, when $\text{SNR} \gg 1$, there is a high certainty among samples for the next optimizer step. While this might be desirable in convex problems, this behavior can trap the optimization process in poor local minima from the early training stages for a non-convex loss landscape.

2.2. Adam and step-wise SNR

Adaptive-learning optimizers use SNR-related metrics to correct the next step for each parameter. Specifically for Adam [12], at each iteration t , the parameters are updated with the following rule:

$$\theta^{t+1} = \theta^t - \eta \cdot \mathcal{C}, \quad (10)$$

where $\Delta\theta = \eta \cdot \mathcal{C}$ is the parameter step and \mathcal{C} the learning rate correction given by:

$$\mathcal{C} = \frac{\hat{m}}{\sqrt{\hat{v}} + \epsilon} = \frac{\mathbb{E}_{\beta_1}[\nabla_{\theta}\mathcal{L}]\sqrt{1 - \beta_2^t}}{\sqrt{\mathbb{E}_{\beta_2}[(\nabla_{\theta}\mathcal{L})^2](1 - \beta_1^t) + \epsilon}}, \quad \forall \theta \in \Theta \quad (11)$$

where \hat{m} is the first moment and \hat{v} is the second moment of the full-batch gradients, and \mathbb{E}_{β} is the exponential moving average (EMA) of past steps. For all iterations past the initial training stage, the effect of the corrective terms is negligible with $\sqrt{1 - \beta_2^t} \approx (1 - \beta_1^t) \approx 1$. Thus, for $\epsilon \ll \sqrt{\mathbb{E}_{\beta_2}[(\nabla_{\theta}\mathcal{L})^2]}$, Eq. 11 can be simplified as:

$$\mathcal{C} = \frac{\mathbb{E}_{\beta_1}[\nabla_{\theta}\mathcal{L}]}{\sqrt{\mathbb{E}_{\beta_2}[(\nabla_{\theta}\mathcal{L})^2]}}, \quad \forall \theta \in \Theta. \quad (12)$$

Based on Eq. 12, we can define the Signal-to-Rms ratio (SRR) as:

$$\text{SRR}_{\mathcal{C}} = \frac{\|\mathbb{E}_{\beta_1}[\nabla_{\theta}\mathcal{L}]\|_2}{\|\sqrt{\mathbb{E}_{\beta_2}[(\nabla_{\theta}\mathcal{L})^2]}\|_2}, \quad (13)$$

where the term $\|\mathbb{E}_{\beta_1}[\nabla_{\theta}\mathcal{L}]\|_2$ is the expected gradient (signal) and $\|\sqrt{\mathbb{E}_{\beta_2}[(\nabla_{\theta}\mathcal{L})^2]}\|_2$ is the Rms of the gradient (noise). Hence, $\text{SRR}_{\mathcal{C}}$ represents the un-centered SNR for each parameter θ , with respect to the optimizer steps t (resembling Eq. 8).

Suppose that for an interval of T iterations, the optimal control of the learning rate correction for each θ is achieved, then \mathcal{C} will be at an equilibrium \mathcal{C}^* for which $\hat{m} \propto \sqrt{\hat{v}}$ and:

$$\mathcal{C}^* \approx \text{const} \Rightarrow \text{SRR}_{\mathcal{C}^*} \approx \text{const}, \quad \forall t \in T. \quad (14)$$

Therefore, if the optimal step is $\Delta\theta^* = \eta \cdot \mathcal{C}^*$, then for $\mathcal{C} < \mathcal{C}^*$ the current step $\Delta\theta$ decreases relative to the optimal step and vice versa. In simpler terms, in highly deterministic regions, the optimizer takes larger steps, while in highly stochastic regions, the optimizer is encouraged to take smaller steps, with a theoretical training equilibrium at $\text{SRR}_{\mathcal{C}} = \text{SRR}_{\mathcal{C}^*}$.

2.3. Training dynamics and parameter sensitivity

Each step of the optimizer makes a parameter update that aims to better fit the outputs based on the inputs. However, the parameters often have a different sensitivity to the loss, corresponding to scale variance of the landscape (Fig. 4). To find a good minimum and exploit the full network architecture, the optimizer must respond to each parameter θ with similar sensitivity to fully capture their inter-dependencies during each backward pass. For example, if a parameter θ_1 has a very sharp minimum which constrains another parameter θ_2 from progressing in that dimension, it requires a higher sensitivity of the optimizer, such that the different scales are comparable and meaningful progress can be made in all possible dimensions. Adam aims to do that by adjusting the learning rate correction \mathcal{C} (sensitivity) based on each parameter’s gradient response as discussed in Section 2.2, leading to an adaptive scaling effect of the landscape for each dimension, as shown in Figure 4.

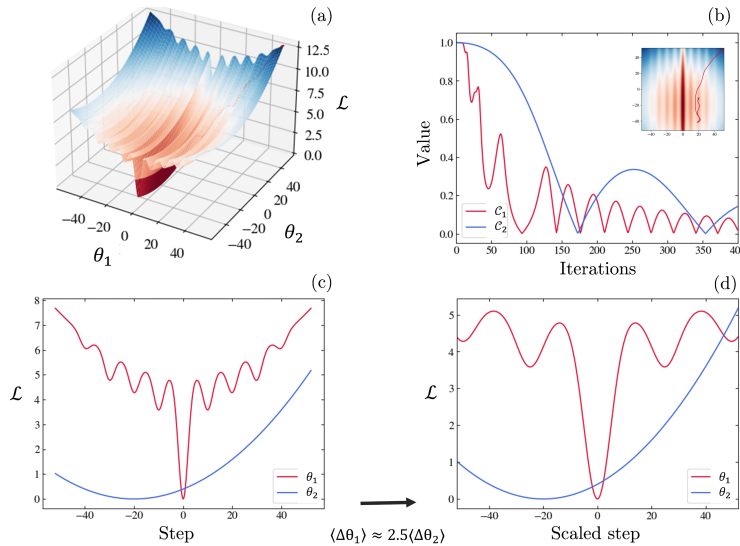


Figure 4: **Landscape scaling with Adam:** Indicative loss landscape for two parameters (a), convergence of learning rate corrections (b), initial slices at the global minimum ($\theta_1 = 0, \theta_2 = -10$) (c) and scaled slice for θ_1 with the average learning rate correction (d). Adam assigns a larger learning rate for θ_1 on average since it is more sensitive to the loss than θ_2 . This has an adaptive scaling effect as the steps on θ_1 gradually decrease faster than for θ_2 .

However, learning rate adjustments take into account the full-batch gradient (average gradient of each sample $\nabla_{\theta} \mathcal{L}_{x_i}$). This can result in prolonged training times, instabilities, or getting stuck in a sub-optimal space where not all the loss landscape trajectories have been explored. To illustrate this weakness, in Figure 5, we draw an indicative training path of a single parameter θ and three samples x , similar to [57]. At t_1 , there is general agreement in the sample’s direction. However, this path later splits into three possible branches: in branch A, the optimizer reaches a state where two (or more) samples cancel out. In branch B, the residual of sample x_3 remains unchanged during the recent steps, so for a quadratic loss, this means that $\nabla_{\theta} \mathcal{L}_{x_3} = 0$. Finally, branch C shows the path of a training equilibrium where the model is learning non-uniformly from different parts of \mathcal{X} . This process eventually converges at a stationary solution with $\nabla_{\theta} \mathcal{L} = 0$, where the sample-wise gradients cancel out while being disproportional (as shown for step t_{∞}), indicating that the model overfits to certain x_i while underfitting to others. So while Adam aims to maintain a consistent direction for the full-batch trajectory ($\text{SRR}_{\mathcal{L}} \approx \text{const}$), the individual sample-wise gradients can have a variable agreement between consecutive steps ($\text{SNR} \ll 1$).

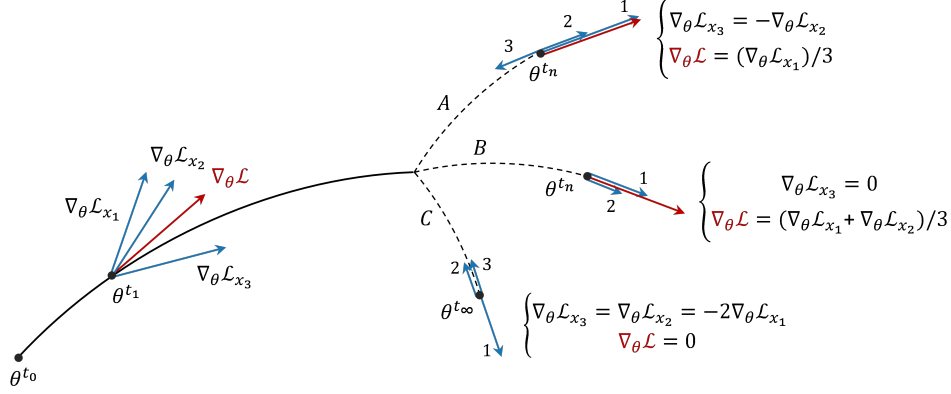


Figure 5: **Training trajectory:** At each step t , Adam aims to maintain a consistent direction based on the gradient signal $\mathbb{E}[\nabla_{\theta}\mathcal{L}]$. However, this does not guarantee agreement in the sample-wise gradients, which can lead to converging in a local minimum where the sample-wise gradients cancel out while being disproportional. Such a scenario could indicate that the model overfits to certain x_i while underfitting to others.

This can also be interpreted by the concept of sample typicality [58], where the gradients of certain samples of a subset \mathcal{H} provide the most informative search direction, while the remaining samples are grouped in subset \mathcal{J} . It has been shown that controlling the antagonism of \mathcal{H} and \mathcal{J} by assigning higher sampling probability to the first, accelerates the training convergence of SGD. Therefore, the subset \mathcal{H} consists of gradients that cancel out or that have become zero, as shown in Figure 5. Sample typicality is formally defined below.

Definition 1: A set of m typical training samples \mathcal{H} is called highly representative subset of \mathcal{X} if it satisfies:

$$\frac{1}{m} \sum_{i \in \mathcal{H}} \nabla_{\theta} \mathcal{L}_i = \nabla_{\theta} \mathcal{L}, \quad \text{with } \mathcal{H} \subset \mathcal{X} \quad (15)$$

Along this line, we can express the concept of gradient homogeneity for a given sample set \mathcal{X} when a size m exists for which all equally sized subsets (or batches) are highly representative of the true gradient. This scenario implies an ideal gradient uniformity among samples, which is theoretically possible in the absence of gradient outliers and heteroscedastic noise in the labels y_i .

Definition 2: A set \mathcal{X} satisfies gradient homogeneity when the gradients of all its subsets \mathcal{H}_i with fixed size m are approximately equal to the true gradient:

$$\frac{1}{m} \sum_{i \in \mathcal{H}_i} \nabla_{\theta} \mathcal{L}_i = \nabla_{\theta} \mathcal{L}, \quad \forall \mathcal{H}_i \in \mathcal{X} \quad (16)$$

Proposition: Starting from the expression $rms(x)^2 = \langle x \rangle^2 + std(x)^2$ it can be shown that the SNR is tightly linked with gradient homogeneity among samples through the following expression:

$$SNR = \frac{|SRR_{\mathcal{B}}|}{\sqrt{1 - SRR_{\mathcal{B}}^2}}, \quad (17)$$

where $SRR_{\mathcal{B}} = \|\mathbb{E}[\nabla_{\theta}\mathcal{L}_{\mathcal{B}}]\|_2 / \|\sqrt{\mathbb{E}[(\nabla_{\theta}\mathcal{L}_{\mathcal{B}})^2]}\|_2$ quantifies the degree of gradient homogeneity between batches, bounded between $[0, 1]$. When Eq. 16 holds, then $SRR_{\mathcal{B}} = 1$ and $SNR \rightarrow \infty$. However, given that the batches are not identical, SNR is bounded by the channel capacity at maximum compression [50], so in practice $SRR_{\mathcal{B}} < 1$. Moreover, from Eq. 17 it can be shown that:

$$SNR > 1, \text{ for } SRR_{\mathcal{B}} \in \left(\frac{\sqrt{2}}{2}, 1\right), \quad (18)$$

$$SNR \approx SRR_{\mathcal{B}}, \text{ for } SRR_{\mathcal{B}} \ll 1. \quad (19)$$

This means that when $SNR > 1$ the batch gradients become highly homogeneous (Eq. 18), while for low gradient homogeneity the SNR becomes approximately equal to $SRR_{\mathcal{B}}$ (Eq. 19).

2.4. Residual homogeneity

As previously discussed, although Adam moderates each parameter step to leverage a consistent gradient signal, this does not necessarily guarantee that each sample x_i is equally satisfied at the stationary solution (or at any point during training). In other words, the converged solution can have heterogeneous gradients for different regions of the sample space, which can imply heterogeneous residuals \mathcal{R} .

To analyse this statement, first we assume that a domain Ω has homogeneous residuals across the sample space if there exists a fixed volume A for sub-domains Ω_i such that:

$$\int_{\Omega_a} |\mathcal{R}(x)| dx \approx \int_{\Omega_b} |\mathcal{R}(x)| dx, \forall \Omega_a, \Omega_b \in \Omega \quad (20)$$

$$\text{with } A(\mathcal{R}(x_i)) = const, \forall x_i \in d$$

Then, starting from the first-order Taylor expansion of the loss function \mathcal{L} , for two consecutive steps t and $t + 1$, we can make the following approximation:

$$\mathcal{L}(\theta^{t+1}) \approx \mathcal{L}(\theta^t) + (\theta^{t+1} - \theta^t) \nabla_{\theta} \mathcal{L}(\theta^t). \quad (21)$$

Substituting the GD update rule in Eq. 21:

$$\mathcal{L}(\theta^{t+1}) - \mathcal{L}(\theta^t) \approx -\eta \cdot \|\nabla_{\theta} \mathcal{L}(\theta^t)\|_2^2 \quad (22)$$

At a stationary solution, $\mathcal{L}(\theta^{t+1}) - \mathcal{L}(\theta^t) = 0$, so given $\eta > 0$:

$$\|\nabla_{\theta} \mathcal{L}(\theta^t)\|_2^2 = 0. \quad (23)$$

Expanding the norm for p parameters and substituting \mathcal{L} gives:

$$\sum_{j=1}^p \left(\frac{\partial \mathcal{L}}{\partial \theta_j} \right)^2 = 0 \Rightarrow \sum_{j=1}^p \left[\frac{\partial}{\partial \theta_j} \left(\frac{1}{n} \sum_{i=1}^n \mathcal{R}_i^2 \right) \right]^2 = 0 \quad (24)$$

By applying the linearity of differentiation and simplifying, we get the following:

$$\sum_{j=1}^p \left(\sum_{i=1}^n \mathcal{R}_i \frac{\partial \mathcal{R}_i}{\partial \theta_j} \right)^2 = 0. \quad (25)$$

So for each parameter j :

$$\sum_{i=1}^n \mathcal{R}_i \frac{\partial \mathcal{R}_i}{\partial \theta_j} = 0. \quad (26)$$

Eq. 26 states that although the terms $\mathcal{R}_i \frac{\partial \mathcal{R}_i}{\partial \theta_j}$ cancel out for all individual samples x_i , this condition is not enough to guarantee that the residual magnitudes $|\mathcal{R}_i|$ are homogeneous (or 0) $\forall x_i \in d$. This can lead to a sub-optimal solution where Eq. 20 is not satisfied, implying that different parts of the domain have converged while others have not. However, to reach a low relative L^2 error (Eq. 7), residual homogeneity must also hold (Eq. 20).

2.5. Adaptive sample-wise weighting

A state where the sample gradients have equal sensitivity on the parameter updates can be achieved by enforcing uniform residuals for a certain interval of steps. If during that interval, \mathcal{R} remains uniform (Eq. 20), then $\nabla_{\theta} \mathcal{R}$ will also tend to be uniform. The residuals can still vary but at a similar rate. Therefore, the gradients $\nabla_{\theta} \mathcal{L}$ will also be uniform for any parameter θ , since they are proportional to the product $\mathcal{R} \nabla_{\theta} \mathcal{R}$. Moreover, homogeneous residuals imply that information is propagating across the sample space evenly, which is a crucial aspect of training PINNs [39]. In this work, the constraint of Eq. 20 is softly enforced during training with a residual-based attention (RBA) re-weighting scheme. This mechanism serves

as a form of regularization by dynamically adjusting the contribution of individual samples, based on the residual history from preceding steps. To analyze how RBA works, we consider the following modified quadratic loss:

$$\mathcal{L}(\theta, x) = \frac{1}{n} \sum_{i=1}^n \mathcal{R}_i(\theta, x_i)^2 \rightarrow \mathcal{L}^* = \frac{1}{n} \sum_{i=1}^n (\lambda_i \cdot \mathcal{R}_i(\theta, x_i))^2, \quad (27)$$

where \mathcal{R}_i is the PDE residual at sample i and λ_i is its corresponding weight. Note that for simplicity, we consider a single-constraint optimization. For PINNs, this means that we focus on the residual loss while the boundary/initial condition losses are simply treated as added bias terms.

The RBA method calculates the weights λ as a moving average of the relative residuals:

$$\lambda_i^{t+1} = \gamma \lambda_i^t + \eta^* \frac{|\mathcal{R}_i|}{\|\mathcal{R}_i\|_\infty} \quad (28)$$

Normalizing with $\|\mathcal{R}_i\|_\infty$ ensures that the maximum multiplier will be bounded between $[0, 1]$ when $\gamma = 1 - \eta^*$. Thus, λ_i is a weighting coefficient which re-scales each term $\mathcal{R}_i \frac{\partial \mathcal{R}_i}{\partial \theta_j}$ with respect to their relative $|\mathcal{R}|$ component. At convergence, for each parameter j , substituting \mathcal{L}^* in Eq. 26 gives:

$$\sum_{i=1}^n \left(\lambda_i^2 \mathcal{R}_i \frac{\partial \mathcal{R}_i}{\partial \theta_j} + \lambda_i \mathcal{R}_i^2 \frac{\partial \lambda_i}{\partial \theta_j} \right) = 0. \quad (29)$$

We can simplify this expression by dropping the second term in the sum since $\frac{\partial \lambda_i}{\partial \theta_j} \approx 0$, as it does not vary rapidly during training. This happens because η^* is chosen to be small enough such that λ represents the averaged residual magnitudes of many past iterations. Therefore, Eq. 29 becomes:

$$\sum_{i=1}^n \lambda_i^2 \mathcal{R}_i \frac{\partial \mathcal{R}_i}{\partial \theta_j} = \sum_{i=1}^n \left\langle \frac{|\mathcal{R}_i|}{\|\mathcal{R}_i\|_\infty} \right\rangle^2 \mathcal{R}_i \frac{\partial \mathcal{R}_i}{\partial \theta_j} = 0. \quad (30)$$

Hence, samples with consistently larger residuals have a higher influence on the next optimizer step. Given that the larger residuals can be reduced during training, this process has an equilibrium where $\lambda_i \approx 1 \Rightarrow |\mathcal{R}_i| \approx \|\mathcal{R}_i\|_\infty, \forall x_i \in d$. If this happens, each scaling term will tend to one, and the residuals will be approximately homogeneous across the sample space \mathcal{X} (Eq. 20). It should be noted, however, that in cases where the training method or the available sample space are not capable of reducing the residuals beyond a certain threshold \mathcal{R}^* (e.g. shallow network, heteroscedastic or sparse data etc.), RBA will also converge to the equilibrium

where $\lambda_i \approx \frac{|\mathcal{R}_i^*|}{\|\mathcal{R}_i^*\|_\infty}$. Another example would be exceptionally stiff problems (e.g., Kuramoto–Sivashinsky equation) requiring auxiliary domain decomposition methods.

3. Results

3.1. Total diffusion and convergence

To demonstrate the theoretical section in practice, we train four standard PINN benchmarks: Allen-Cahn, Helmholtz, Burgers, and lid-driver cavity flow (Figs. 6, 7, 8, 9 respectively). All figures show the relative L^2 convergence, the SNR curves, and three snapshots (a-c) of the prediction $f(x, \theta)$. We always compare the RBA with the vanilla model and plot the RBA weights for both. The three learning phases (Fitting, Diffusion and Total diffusion) are marked by the SNR evolution of the RBA model. We note that we only split the collocation points into batches and treat the initial and boundary conditions as constant terms added to each batch. The SNR plots have been smoothed using a simple moving average (SMA) of 10 previous iterations. Further implementation details are listed in the Appendix.

As evident from these results, it is during the total diffusion that the models experience the steepest convergence of the relative L^2 error. Moreover, the RBA models induce higher homogeneity of the residuals, resulting in better generalization. During this interval, the SNR increases towards the deterministic regime where it stabilizes, as the gradients agree on the final direction of the optimizer. Note that in the Burgers case of Fig. 8, the discontinuity limits the gradient homogeneity due to the steep gradient, which is always present, so although both the SNR and the residuals undergo a phase transition, the SNR does not become larger than 1. Interestingly, the first three benchmarks converge non-monotonically after total diffusion, which could be related to the “edge of stability” phase [6]. As expected, the more pronounced fluctuations are present in the Burgers case due to the discontinuity. On the other hand, the convergence of the cavity flow appears monotonic (Fig. 9), which could be attributed to the SNR being larger than 1. These results align with the intuition that training in the diffusion phase with residual homogeneity is essential, especially in the context of PINNs, where individual samples are explicitly interdependent based on the governing PDEs. The code for all cases will be made available on github.com/soanagno/total-diffusion.

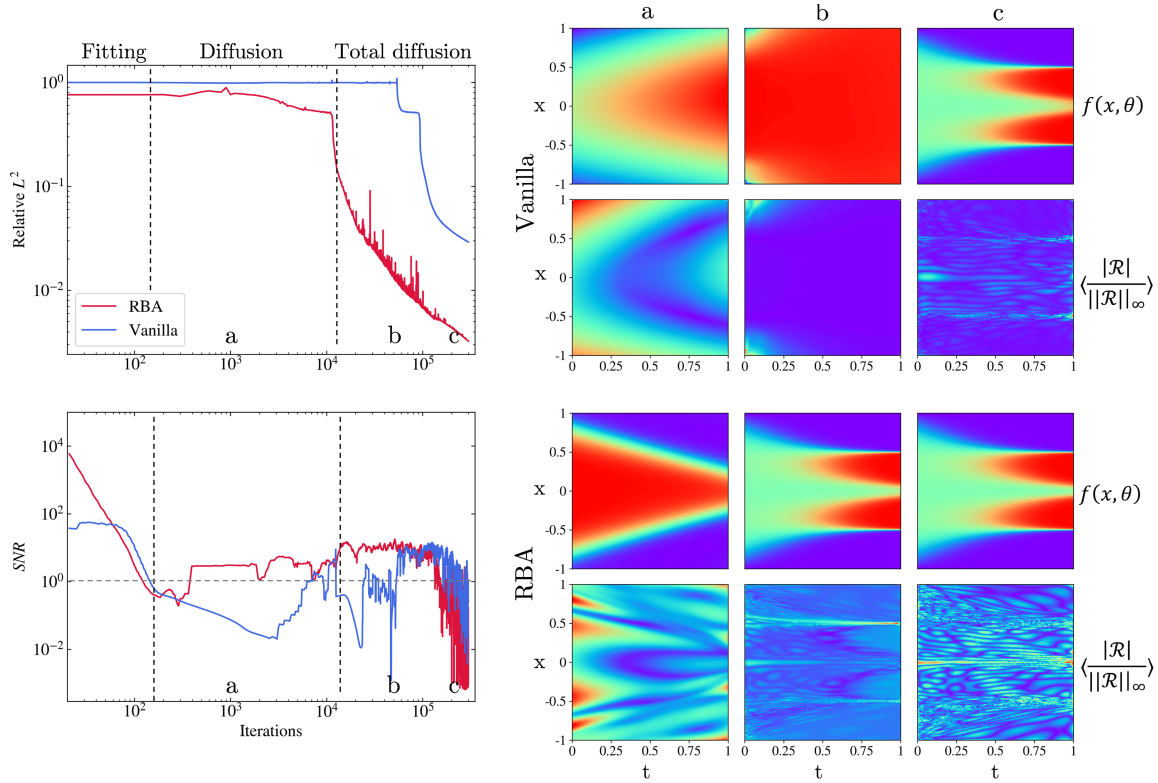


Figure 6: **1D Allen-Cahn training dynamics: Fitting:** Both vanilla and RBA experience a rapid decrease in their SNR prior to diffusion, while the relative L^2 is still stagnant. **Diffusion:** Both models enter diffusion as the SNR is in the stochastic regime. The vanilla model does not converge, and the residuals \mathcal{R} are not homogeneous, while RBA slowly converges with a characteristic double-jet shape for the residuals (snapshot (a)). **Total diffusion:** The RBA model enters the total diffusion phase first, while \mathcal{R} appears highly homogeneous, accompanied by a step L^2 convergence. The vanilla SNR shows a small step increase towards total diffusion, and its L^2 decreases, but the residuals are still heterogeneous (snapshot (b)). On the other hand, the RBA \mathcal{R} appears highly diffused, accompanied by a steep L^2 convergence. At snapshot (c), the vanilla finally has converged to a sub-optimal solution where \mathcal{R} appear more diffused than before but still have a few points which have not converged, causing a less-pronounced color spectrum of \mathcal{R} . The RBA method keeps decreasing the L^2 further, while \mathcal{R} are highly homogeneous, implying that different regions of \mathcal{X} are equally satisfied. **Accuracy:** For a 10% L^2 error, vanilla takes 100k iterations while RBA takes 10k, thus RBA is 10 times faster than vanilla.

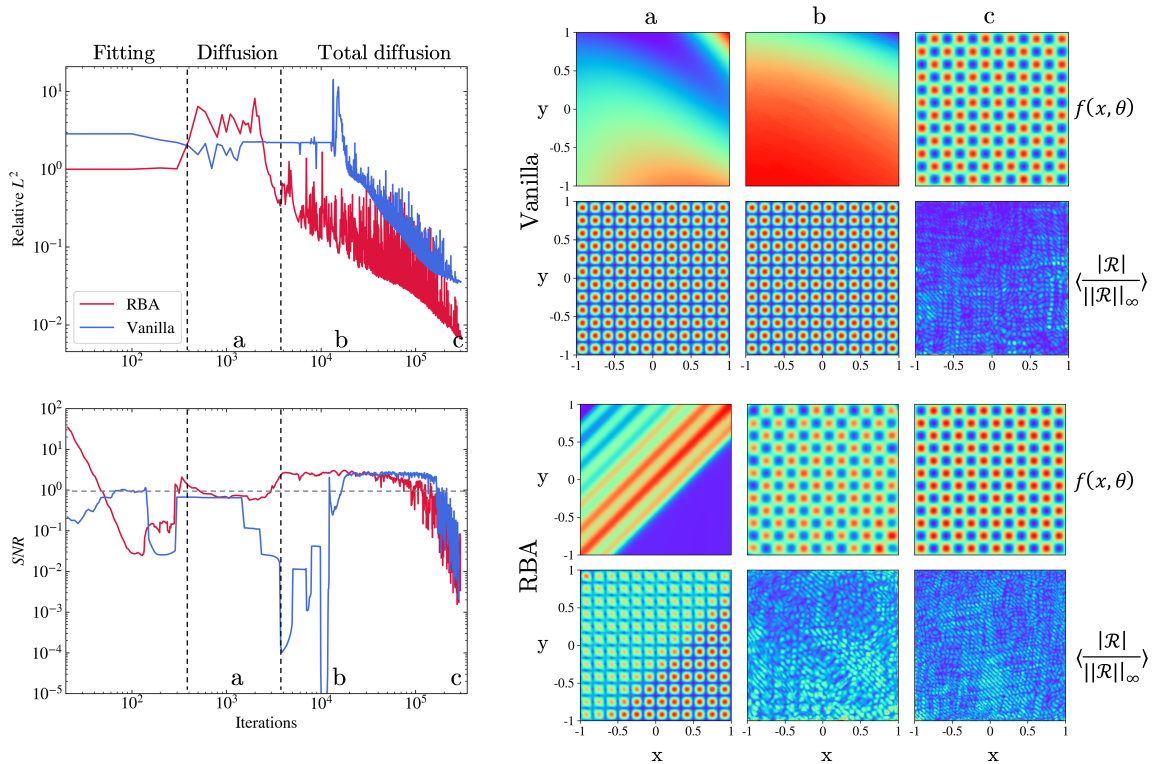


Figure 7: **Helmholtz training dynamics: Fitting:** The RBA model undergoes the usual rapid transition from high to low SNR. However, the vanilla model starts from the stochastic regime. This implies that the optimizer is not able to find a meaningful direction. **Diffusion:** The SNR of both models is in the stochastic regime but vanilla SNR further decreases while the SNR of RBA is constant and close to total diffusion. While the residuals of the vanilla are not homogeneous, they start diffusing for RBA (snapshot (a)). **Total diffusion:** The vanilla SNR dives deeper into the stochastic regime where the L^2 cannot be decreased, while the RBA goes into total diffusion and initiates the main convergence trajectory of the relative L^2 error. The vanilla \mathcal{R} are still heterogeneous, but for RBA, they appear more disordered and homogeneous (snapshot (b)). At snapshot (c), both models are now in total diffusion with corresponding relative L^2 convergence. However, the residuals of the vanilla model have not yet diffused as they appear for the RBA, which results in a suboptimal solution compared to the RBA. The RBA method keeps decreasing the L^2 further, while \mathcal{R} are highly homogeneous, implying that different regions of \mathcal{X} are equally satisfied. Note that the increased noise in the L^2 convergence for the RBA may be associated with higher exploration of the non-convex loss landscape. **Accuracy:** Again, 10% of L^2 error is reached by the vanilla model within 100k iterations while RBA takes 10k. Thus, RBA is 10 times faster than vanilla.

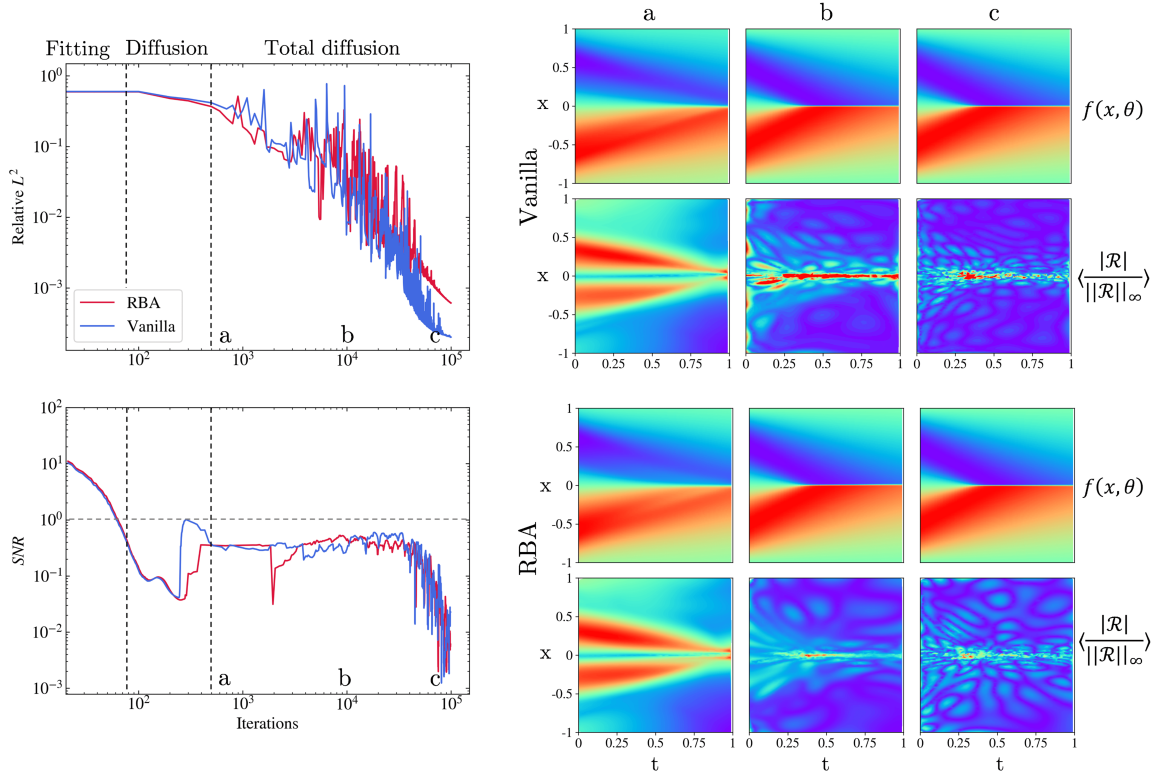


Figure 8: **Burgers training dynamics:** Both models follow similar SNR paths, which undergo a high-to-low transition. Here, the SNR remains below 1 during total diffusion, probably due to the persistent steep gradient that causes a discontinuity at $x = 0$, hence increased stochasticity in the gradients. To make the residual homogeneity more visible for the vanilla case, we have limited the maximum range of their values (only for Burgers), since the high residuals at the discontinuity dominated the rest of the domain. Note that the Burgers equation is not considered a stiff problem so the training quickly enters total diffusion, as opposed to the other cases (around 500 iterations). Here, the RBA model performs a bit worse than vanilla since the weighting multipliers focus on a discontinuous region, which cannot be resolved. Furthermore, the relative L^2 convergence is characterized by large fluctuations due to the discontinuity, which is also reflected in the SNR order of magnitude.

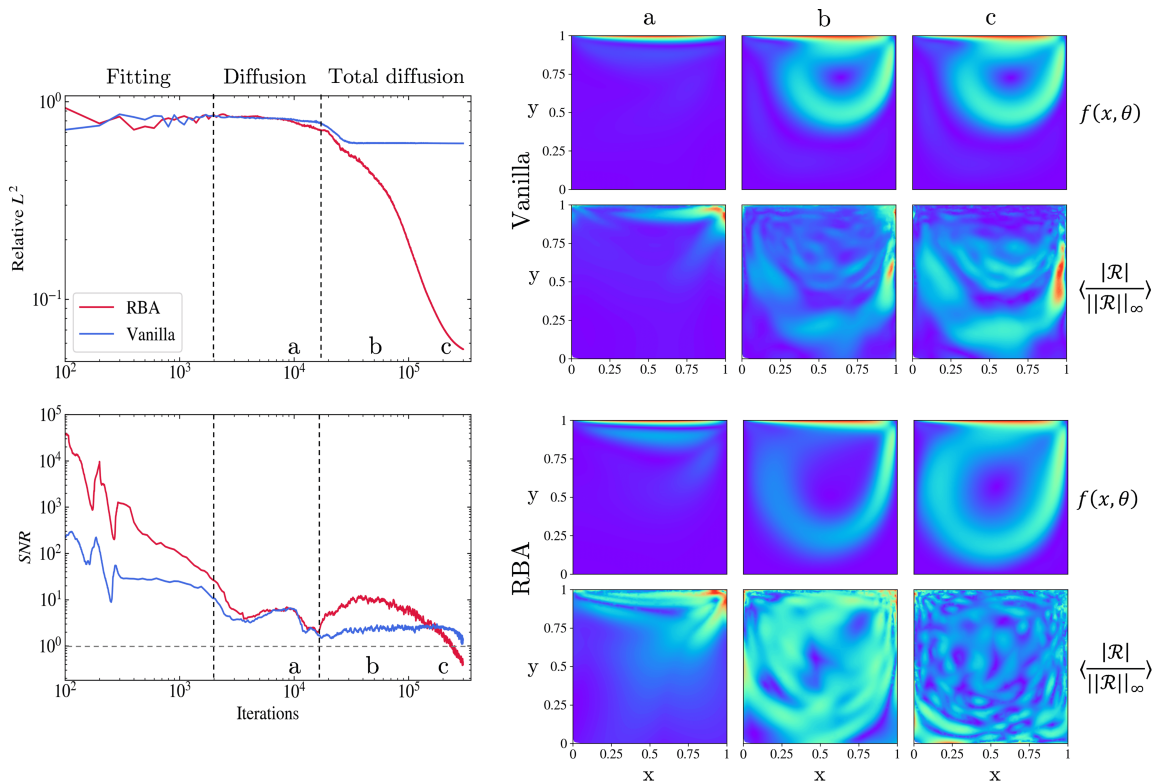


Figure 9: **Lid-driven cavity flow (Re 1000) training dynamics:** Both models undergo a phase transition from high-to-low SNR which does not reduce the L^2 error. During diffusion, the error starts converging slowly, and the SNR further decreases. The abrupt transition into total diffusion is met by a sudden increase of the SNR and a simultaneous drop in the relative L^2 . However, the vanilla model converges to a stationary solution with high test error due to the heterogeneous residuals on the right boundary, which hinder the information flow to the rest of the domain. Since the region with accumulated high residuals cannot be resolved, the vanilla cannot generalize. The RBA model, however, distributes the residuals evenly across the cavity and generalizes with a final L^2 error of about 6%. Note how the SNR remains above 1, which is reflected by higher agreement/determinism of the steps and smoother convergence relative to the previous benchmarks.

3.2. Correlation of Adam and batch SNR

As previously shown, the SNR is related to gradient homogeneity through $\text{SRR}_{\mathcal{B}}$: high SNR implies $\text{SRR}_{\mathcal{B}} \approx 1$, while low SNR implies $\text{SNR} \approx \text{SRR}_{\mathcal{B}}$ (consistent with Eqs. 18, 19). Furthermore, we also observe a correlation with $\text{SRR}_{\mathcal{C}}$ (the step-wise SNR of Adam) for all benchmark cases, as shown in Figure 10. More specifically, for all RBA models $\text{SRR}_{\mathcal{C}}$ closely follows the curve of SNR, while at total diffusion, they both reach a stable equilibrium which coincides with the steepest L^2 convergence (vertical dashed lines Fig. 10). However, this is not always the case for the vanilla models, as they spend more time in the stochastic regime, where Adam cannot reach equilibrium due to the high stochasticity/disagreement in the batch gradients. Therefore, it appears that Adam is able to converge more optimally when the SNR is at an equilibrium with homogeneous gradients. Notably, although SNR undergoes a high-to-low transition, the optimal convergence is met during the total diffusion.

During that final phase, the abrupt increase in the SNR remains several orders of magnitude lower than the fitting phase SNR.

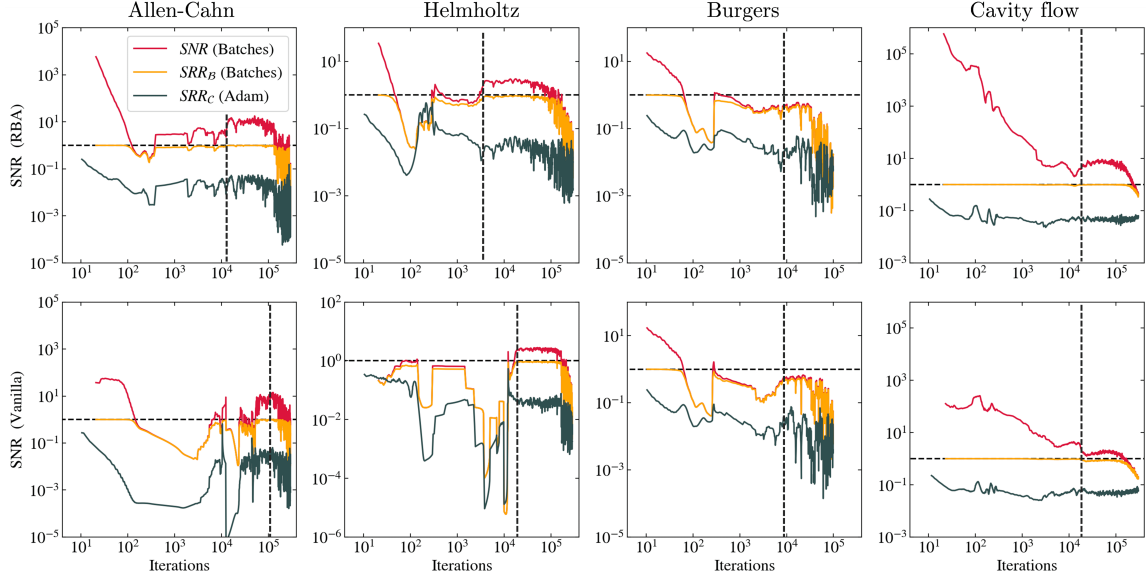


Figure 10: **Correlation of SNR, SRR_B and SRR_C :** For all models the SNR is related to gradient homogeneity through SRR_B , since high SNR implies $SRR_B \approx 1$, while low SNR implies $SNR \approx SRR_B$ (consistent with Eqs. 18, 19). Furthermore, there is a high correlation of SNR with SRR_C , (the step-wise SNR of Adam) for all benchmark cases. For all RBA models, the SRR_C closely follows the curve of SNR, while at total diffusion, they both reach a stable equilibrium which coincides with the steepest L^2 convergence (vertical dashed lines). However, for the vanilla models, this is not always the case as they spend more time in the stochastic regime, where Adam cannot reach an equilibrium due to the high stochasticity in the batch gradients. Therefore, it appears that Adam is able to converge more optimally when the SNR is at an equilibrium above $\mathcal{O}(1)$. Notably, although SNR undergoes a high-to-low transition, the optimal convergence is met at the total diffusion when the gradients become homogeneous. This causes an abrupt increase in the SNR, which usually remains several orders of magnitude lower than the fitting phase SNR.

3.3. Information compression

In IB theory, information compression serves as the primary explanation for the phase transition of the gradients, which occurs when the layers' activations saturate at the bounds of the activation function, given that such bounds exist. During this phase, the network maximizes the information of the layer activations T with respect to the output $I(T; \mathcal{Y})$, while minimizing (compressing) their information based on the input $I(\mathcal{X}; T)$. However, this claim has been debated through several studies. For instance, the phase transition has been replicated even with ReLU activation, which is not bounded like tanh [53]. Another important observation is that for a continuous range the information is infinite, so the measured information can be quite sensitive to the discretization method which is usually adopted [56]. Thus, to interpret how the different layers carry the information from the input to the output, the range of the activations must be carefully discretized prior to applying the appropriate strategy of quantifying the information [55]. However, a universal method that yields the exact quantification of I remains an open question.

This section aims to investigate the compression phenomenon by focusing on $I(\mathcal{X}; T)$ and its behavior during the phase transition. The mutual information, under the assumption that T is a deterministic function of \mathcal{X} , is given by:

$$I(\mathcal{X}; T) = H(T) - H(T|\mathcal{X}) = H(T) = \sum_{i=1}^{i=N} p_i \log_2(p_i). \quad (31)$$

To discretize the activations for each layer, we calculate the mean range of each layer’s neurons and create 30 equidistant bins within that range. This relative-binning approach has the advantage of considering the relative range of each layer, which might vary depending on the network properties and stage of training. Therefore, each neuron takes 30 possible values, and each layer l represents a string of p_l digits (representation), where p are the parameters of the layer. If each unique input x_i is correlated to a unique representation of the layer, we assume that the layer carries the maximum information possible to distinguish the inputs, equal to $\log_2(n)$ bits. The entropy of Eq. 31 is finally calculated using the probability of occurrence of each unique representation T for each input \mathcal{X} .

Although there is some lossy compression after total diffusion (decrease in I), most of the information is available across all layers from the early iterations, implying that the network can be sensitive enough to pass the information to the deeper layers, given that I is calculated with the relative-binning (Fig. 11(a, b)). We note that measuring all layers with the same fixed range (as some of the earlier IB studies) produced mixed results. During the fitting phase, the range of the deep layers was smaller than the layers close to the input, but after diffusion, the range of the deep layers was larger. Thus, since the first layers don’t saturate, this lead to higher I in the deep layers during diffusion, which, based on the Markovian chain assumption is not consistent. However, we suggest that the information compression does not necessarily imply dramatic information loss in the layers but mainly saturation of the activations. As shown in Figure 11, the deeper layers become more saturated (binary) as the training progresses, with a sudden increase during the total diffusion. Interestingly, the middle layers appear to saturate more than the last layers, which can be interpreted by an encoding-decoding process. In essence, it appears that the information required to distinguish most inputs is present from early on. Still, the network becomes more binary and achieves a compressed representation during total diffusion, mainly for the deeper layers.

To quantify how efficient each layer is at conveying the information of the inputs, we measure the percentage of binary activations (saturated at the boundaries of tanh) (Fig. 11(c, d)). At total diffusion, most layers become highly saturated, indicating a compressed/efficient representation. This is also captured by the network parameter norm $\|\theta\|$, which show that the middle layers compress more (Fig. 11(e, f)). This phenomenon is reminiscent of an encoding-decoding process and was also present for Burgers and cavity flow (not shown here). Only for Helmholtz, the layers are compressed with a hierarchical order (ascending $\|\theta\|$ from the first to the last layer). Moreover, although the vanilla model achieves a compressed representation, the layers undergo a less smooth transition (Fig. 11(b, d)).

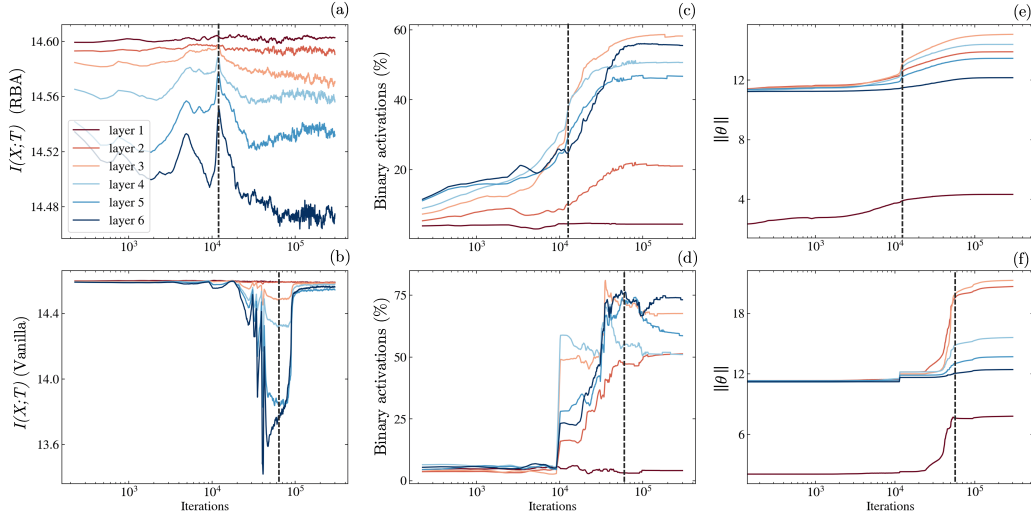


Figure 11: **Compressed representation and phase transition:** $I(\mathcal{X}; T)$ curves for each layer in bits (a, b), percentage of binary activations per layer (c, d) and parameter norm (e, f). Due to the discretization method, there is only a negligible I decrease after the total diffusion for RBA. This implies that even the deeper layers respond to the inputs during training, although with a smaller range for the first iterations (a). For vanilla, there is a more considerable information loss at the onset of diffusion, which bounces back to total diffusion (b). The parameters θ rapidly increase at total diffusion due to high gradient agreement among batches (dashed lines). The tanh activation saturates at its boundaries, compressing the representation of the input information (c-f). Note that the middle layers show the highest compression, reminiscent of an encoding-decoding architecture.

The “binarization” of the layers is further demonstrated in Figure 12, where we show the activations for a random input prior to and after total diffusion. The deeper layers achieve a binary representation of the input, which is more compressed than the first layers. However, a more compressed representation does not necessarily lead to higher information loss.

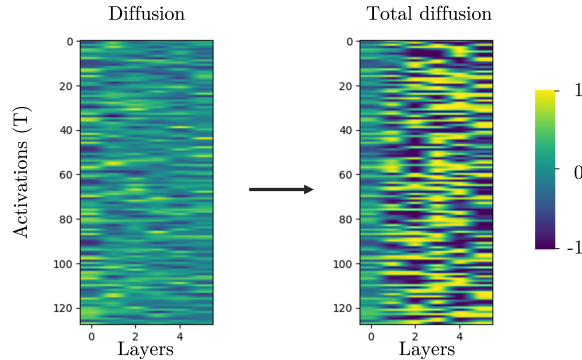


Figure 12: **Saturation of deep layers:** Indicative activation values for a random input during the diffusion (left) and total diffusion (right) phase. As the weights increase, the tanh activation saturates at its boundaries, compressing the representation of the input information.

4. Summary

In this study, we empirically demonstrated the existence of phase transitions proposed by the IB theory (fitting/diffusion) for PINNs trained with the Adam optimizer. Additionally, we identified a previously undisclosed third phase, which we named “total diffusion,” characterized by batch gradient agreement and fast convergence. Our findings highlight a strong correlation between this phase and gradient homogeneity, indicating that PINNs achieve optimal learning when residuals are diffused uniformly, leading to increased stability of the optimizer learning rate correction. To further support this hypothesis, we introduce a re-weighting technique (RBA) to improve the vanilla PINN performance, which induces homogeneous residuals faster and improves generalization.

Moreover, we investigate the relationship between signal-to-noise ratio (SNR) phase transition and information compression, observing a compression phenomenon due to the saturation of activations. As models transition abruptly into “total diffusion”, the network parameters rapidly increase, resulting in the saturation of neuron activations. Employing a relative-binning strategy, we find that while there is negligible information loss in deeper layers, there exists a hierarchy in the information content across layers, aligning with the principles of the IB theory. Our analysis reveals that middle layers exhibit the most pronounced saturation during the “binarization” process, reminiscent of an encoder-decoder mechanism within the fully-connected architecture. We believe this novel connection between IB theory and PINNs could open new possibilities for quantifying the performance of the many versions of PINNs and even neural operators.

Acknowledgements

SJA and NS thank the Swiss National Science Foundation grant “Hemodynamics of Physiological Aging” (Grant nr 205321_197234). JDT and GEK acknowledge support by the NIH grant R01AT012312, the DOE SEA-CROGS project (DESC0023191), the MURI-AFOSR FA9550-20-1-0358 project, and the ONR Vannevar Bush Faculty Fellowship (N00014-22-1-2795).

Appendix A. SNR details

In this section, we present some additional results regarding the SNR curves. In Figure A.13, we show the layer-wise SNR where each layer seems to follow a similar pattern with all the phase transitions occurring at the same iteration, while decreasing SNR hierarchies are forming (shallow to deep), as observed in IB theory. For Allen-Cahn, the layers almost coincide, since the information carried is almost identical. In Figure A.14, we plot the mean and std curves for Allen-Cahn, which provides a clearer picture on why the SNR fraction follows this specific trend.

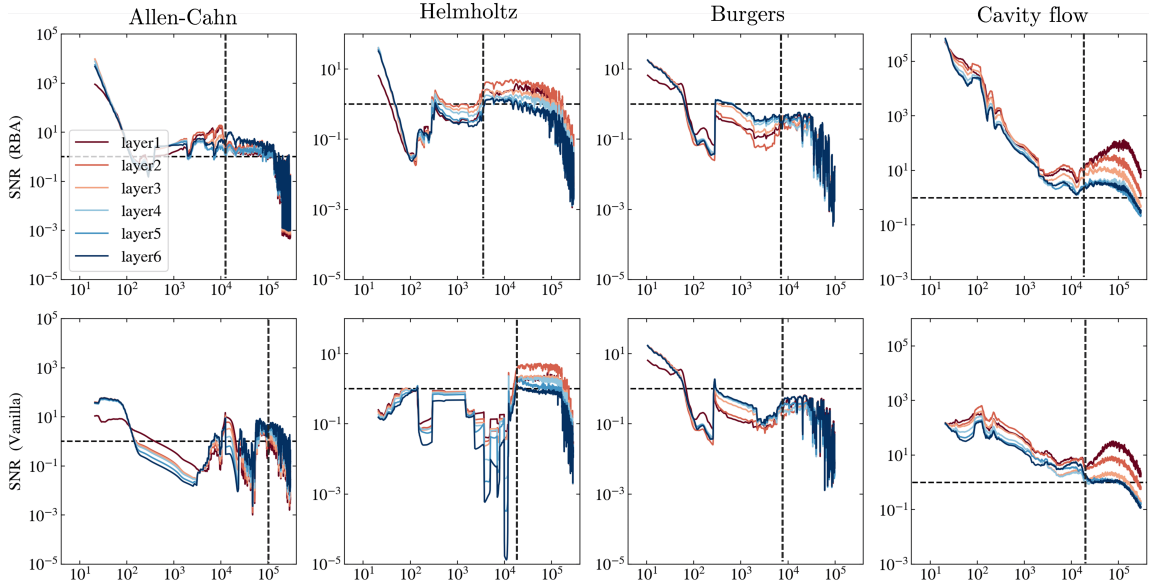


Figure A.13: **Layer-wise SNR:** We plot the SNR progression for each layer of the DNN and observe that each layer follows a similar phase transition. Moreover, for Helmholtz and cavity there is an decreasing SNR hierarchy from shallow to deep layers, which aligns with the IB theory. For Allen-Cahn and Burgers the layers largely coincide, which could be attributed to similar information transfer (Fig. 11).

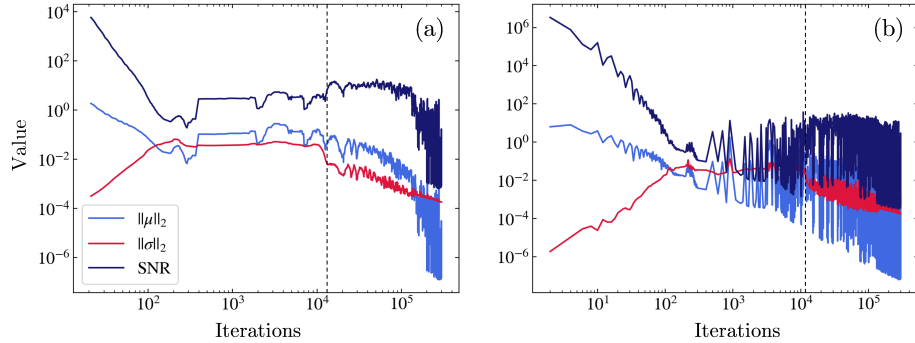


Figure A.14: **Layer-wise SNR:** Both fitting and total diffusion (dashed lines) are characterized by an increased SNR. However, during fitting, the SNR starts high due to low std of the batches since in the beginning of training the direction distinguishing each batch has not been determined (a). The increase of the SNR during total diffusion is also accompanied by stability, as both the mean and std decrease with similar rate. In the late training stages, the mean tends to 0 (as the optimizer converges) while some gradient noise persists (a). It is noteworthy that after total diffusion there are large oscillations present in the non-smoothed curves (b), which implies a certain self-stabilization regime, which is able to overcome local minima while maintaining a stable SNR. This could also be attributed to the hypercritical sharpness phenomenon, which implies a highly non-convex but stable equilibrium.

Appendix B. Implementation Details

Appendix B.1. Allen Cahn Equation

The general form of the 1D Allen-Cahn PDE is given by:

$$\frac{\partial u}{\partial t} = \epsilon \frac{\partial^2 u}{\partial x^2} + f(u), \quad (\text{B.1})$$

where $\epsilon > 0$ represents the diffusion coefficient, and $f(u)$ denotes a non-linear function of u . For this study, we have set $\epsilon = 10^{-4}$ and defined $f(u) = -5u^3 + 5u$. The initial condition is given by:

$$u(0, x) = x^2 \cos(\pi x), \quad \forall x \in [-1, 1], \quad (\text{B.2})$$

To solve the Allen-Cahn PDE, we trained a physics-informed neural network (PINN) for $3 \cdot 10^5$ iterations using the standard Adam optimizer [12] and an exponential learning rate scheduler, evaluated on 25600 collocation points. In general, the loss function $\mathcal{L}(\mathbf{x}, \theta)$ is comprised of weighted terms to ensure adherence to the initial condition, boundary condition, and the PDE residual:

$$\mathcal{L}(\mathbf{x}, \theta) = \mathcal{L}_{ic}(\mathbf{x}, \theta) + \mathcal{L}_{bc}(\mathbf{x}, \theta) + \mathcal{L}_{\mathcal{R}}(\mathbf{x}, \theta),$$

where $\mathcal{L}_{ic}(\mathbf{x}, \theta)$ and $\mathcal{L}_{bc}(\mathbf{x}, \theta)$ are the loss terms for the initial and periodic boundary conditions, respectively, and $\mathcal{L}_{\mathcal{R}}(\mathbf{x}, \theta)$ is the term enforcing the PDE residuals from Equation B.1. RBA weights are applied solely for the PDE residual term $\mathcal{L}_{\mathcal{R}}$ [42] as discussed in Section 2.5.

Appendix B.2. Helmholtz Equation

The 2D Helmholtz PDE is expressed as follows:

$$\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} + k^2 u - q(x, y) = 0, \quad (\text{B.3})$$

where $q(x, y)$ is the forcing term defined as:

$$\begin{aligned} q(x, y) = & - (a_1 \pi)^2 \sin(a_1 \pi x) \sin(a_2 \pi y) \\ & - (a_2 \pi)^2 \sin(a_1 \pi x) \sin(a_2 \pi y) \\ & + k \sin(a_1 \pi x) \sin(a_2 \pi y), \end{aligned} \quad (\text{B.4})$$

leading to the analytical solution $u(x, y) = \sin(a_1 \pi x) \sin(a_2 \pi y)$ [43]. The boundary conditions are defined as:

$$u(-1, y) = u(1, y) = u(x, -1) = u(x, 1) = 0, \quad (\text{B.5})$$

For this study, we set $a_1 = 6$, $a_2 = 6$, and $k = 1$. The loss function for this problem is given by:

$$\mathcal{L}(\mathbf{x}, \theta) = \mathcal{L}_{bc}(\mathbf{x}, \theta) + \mathcal{L}_{\mathcal{R}}(\mathbf{x}, \theta),$$

where $\mathcal{L}_{bc}(\mathbf{x}, \theta)$ ensures the model satisfies the boundary conditions (B.5), and $\mathcal{L}_{\mathcal{R}}(\mathbf{x}, \theta)$ enforces the PDE residuals from Equation B.3. Similar to the previous case, RBA only for the PDE residual term $\mathcal{L}_{\mathcal{R}}$ [42].

We trained the PINN for $3 \cdot 10^5$ iterations and utilized the Adam optimizer [12] with an exponential learning rate scheduler, on a dataset of 25,600 collocation points.

Appendix B.3. Burgers' Equation

The Burgers' equation is defined as:

$$u_t + uu_x = -\nu u_{xx}, \quad (\text{B.6})$$

where u represents the velocity field, subject to the viscosity $\nu = 1/(100\pi)$. The initial condition and boundary conditions are described as follows:

$$u(0, x) = -\sin(\pi x), \quad \forall x \in \Omega, \quad (\text{B.7})$$

$$u(t, -1) = u(t, 1) = 0, \quad \forall t \geq 0, \quad (\text{B.8})$$

defined over the domain $\Omega = (-1, 1) \times (0, 1)$, where $\mathbf{x} = (x, y)$ signifies the spatial coordinates. The loss function is:

$$\mathcal{L} = \mathcal{L}_{bc} + \mathcal{L}_{\mathcal{R}}$$

where \mathcal{L}_{bc} and $\mathcal{L}_{\mathcal{R}}$ denotes the loss functions from the boundary conditions and residual points to satisfy equations B.8 and B.6 respectively. We minimize the total loss using Adam optimizer with exponential learning rate scheduler for 10^5 iterations. For this problem, we use 10000 collocation points.

Appendix B.4. Lid-driven flow

The governing physical laws for the lid-driven cavity problem are encapsulated by the 2D steady Navier-Stokes (NS) equations, represented by the momentum equations:

$$(\mathbf{v} \cdot \nabla)\mathbf{v} = -\nabla p + \frac{1}{Re}\nabla^2\mathbf{v}, \quad (\text{B.9})$$

and the continuity equation:

$$\nabla \cdot \mathbf{v} = 0, \quad (\text{B.10})$$

where $\mathbf{v}(\mathbf{x}) = (u(\mathbf{x}), v(\mathbf{x}))$ represents the velocity field within the cavity, p denotes the pressure, $Re = 1000$ is the Reynolds number, and $\mathbf{x} = (x, y) \in \Omega = (0, 1) \times (0, 1)$, with Ω denoting the two-dimensional cavity domain [20]. The boundary conditions are specified as:

$$\mathbf{v}(\mathbf{x}) = \begin{cases} (1, 0) & \text{if } \mathbf{x} \in \Gamma_1, \\ (0, 0) & \text{if } \mathbf{x} \in \Gamma_2, \end{cases} \quad (\text{B.11})$$

where Γ_1 is the top boundary and Γ_2 encompasses the sides and bottom of the cavity. For the top boundary condition, the edges were smoothed to avoid singularities forming at the points where Γ_1 coincides with Γ_2 , hence creating a well-posed problem. Following [20], we simplify the problem by utilizing the streamfunction ψ , where $\mathbf{v} = \nabla \times \psi$, and approximate u, v , and p as follows:

$$u(x, y) = \frac{\partial \psi_{NN}}{\partial y}(x, y), \quad (\text{B.12})$$

$$v(x, y) = -\frac{\partial \psi_{NN}}{\partial x}(x, y), \quad (\text{B.13})$$

$$p(x, y) = p_{NN}(x, y), \quad (\text{B.14})$$

where ψ_{NN} and p_{NN} are neural network outputs for the stream function and pressure, respectively. This approximation inherently satisfies the continuity equation, $\nabla \cdot \mathbf{v} = \nabla \cdot (\nabla \times \psi) = 0$, allowing us to focus on the momentum equations (B.9). The total loss function is defined as:

$$\mathcal{L} = \mathcal{L}_{bc} + \mathcal{L}_{ns}, \quad (\text{B.15})$$

where the combined loss terms for boundary conditions (\mathcal{L}_{bc}) and Navier-Stokes equations (\mathcal{L}_{ns}) are given by:

$$\mathcal{L}_{bc} = \sum_{b=1}^2 \langle (e_{h,b})^2 \rangle_h, \quad (\text{B.16})$$

$$\mathcal{L}_{ns} = \sum_{n=1}^2 \langle (\lambda_j \cdot e_{j,n})^2 \rangle_j, \quad (\text{B.17})$$

here, $e_{h,b}$ and $e_{j,n}$ denote the residuals for boundary conditions and Navier-Stokes equations at points h and j , respectively. Finally, λ_j represents the RBAs applied to PDE residual points.

To address this system of PDEs, a physics-informed neural network (PINN) was trained for 3×10^5 with Adam, coupled with an exponential learning rate scheduler. This training process was conducted on a set of 25,600 collocation points.

References

- [1] M. Zhou, T. Liu, Y. Li, D. Lin, E. Zhou, T. Zhao, Toward understanding the importance of noise in training neural networks, in: International Conference on Machine Learning, PMLR, 2019, pp. 7594–7602.
- [2] L. N. Smith, N. Topin, Super-convergence: Very fast training of neural networks using large learning rates, in: Artificial intelligence and machine learning for multi-domain operations applications, Vol. 11006, SPIE, 2019, pp. 369–386.
- [3] A. Lewkowycz, Y. Bahri, E. Dyer, J. Sohl-Dickstein, G. Gur-Ari, The large learning rate phase of deep learning: the catapult mechanism, arXiv preprint arXiv:2003.02218 (2020).
- [4] J. Lee, L. Xiao, S. Schoenholz, Y. Bahri, R. Novak, J. Sohl-Dickstein, J. Pennington, Wide neural networks of any depth evolve as linear models under gradient descent, Advances in neural information processing systems 32 (2019).
- [5] M. I. Katsnelson, V. Vanchurin, T. Westerhout, Self-organized criticality in neural networks, arXiv preprint arXiv:2107.03402 (2021).
- [6] J. Cohen, S. Kaur, Y. Li, J. Z. Kolter, A. Talwalkar, Gradient descent on neural networks typically occurs at the edge of stability, in: International Conference on Learning Representations, 2020.
- [7] J. M. Cohen, B. Ghorbani, S. Krishnan, N. Agarwal, S. Medapati, M. Badura, D. Suo, D. Cardoze, Z. Nado, G. E. Dahl, et al., Adaptive gradient methods at the edge of stability, arXiv preprint arXiv:2207.14484 (2022).
- [8] A. Damian, E. Nichani, J. D. Lee, Self-stabilization: The implicit bias of gradient descent at the edge of stability, arXiv preprint arXiv:2209.15594 (2022).
- [9] A. Kosson, B. Messmer, M. Jaggi, Rotational equilibrium: How weight decay balances learning across neural networks (2023).
- [10] P. Nakkiran, G. Kaplun, Y. Bansal, T. Yang, B. Barak, I. Sutskever, Deep double descent: Where bigger models and more data hurt, Journal of Statistical Mechanics: Theory and Experiment 2021 (12) (2021) 124003.
- [11] A. Jacot, F. Gabriel, C. Hongler, Neural tangent kernel: Convergence and generalization in neural networks, Advances in neural information processing systems 31 (2018).
- [12] D. P. Kingma, J. Ba, Adam: A method for stochastic optimization, arXiv preprint arXiv:1412.6980 (2014).
- [13] M. Raissi, P. Perdikaris, G. 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 (2019) 686–707.

- [14] M. Raissi, A. Yazdani, G. E. Karniadakis, Hidden fluid mechanics: Learning velocity and pressure fields from flow visualizations, *Science* 367 (6481) (2020) 1026–1030.
- [15] K. A. Boster, S. Cai, A. Ladrón-de Guevara, J. Sun, X. Zheng, T. Du, J. H. Thomas, M. Nedergaard, G. E. Karniadakis, D. H. Kelley, Artificial intelligence velocimetry reveals in vivo flow rates, pressure gradients, and shear stresses in murine perivascular flows, *Proceedings of the National Academy of Sciences* 120 (14) (2023) e2217744120.
- [16] V. Oommen, B. Srinivasan, Solving inverse heat transfer problems without surrogate models: a fast, data-sparse, physics informed neural network approach, *Journal of Computing and Information Science in Engineering* 22 (4) (2022) 041012.
- [17] S. Cai, H. Li, F. Zheng, F. Kong, M. Dao, G. E. Karniadakis, S. Suresh, Artificial intelligence velocimetry and microaneurysm-on-a-chip for three-dimensional analysis of blood flow in physiology and disease, *Proceedings of the National Academy of Sciences* 118 (13) (2021) e2100697118.
- [18] Q. Cao, S. Goswami, T. Tripura, S. Chakraborty, G. E. Karniadakis, Deep neural operators can predict the real-time response of floating offshore structures under irregular waves, *Computers & Structures* 291 (2024) 107228.
- [19] A. Kopaničáková, H. Kothari, G. E. Karniadakis, R. Krause, Enhancing training of physics-informed neural networks using domain-decomposition based preconditioning strategies, *arXiv preprint arXiv:2306.17648* (2023).
- [20] S. Wang, Y. Teng, P. Perdikaris, Understanding and mitigating gradient flow pathologies in physics-informed neural networks, *SIAM Journal on Scientific Computing* 43 (5) (2021) A3055–A3081.
- [21] T. Salimans, D. P. Kingma, Weight normalization: A simple reparameterization to accelerate training of deep neural networks, *Advances in neural information processing systems* 29 (2016).
- [22] S. Wang, B. Li, Y. Chen, P. Perdikaris, Piratenets: Physics-informed deep learning with residual adaptive networks, *arXiv preprint arXiv:2402.00326* (2024).
- [23] F. Jiang, X. Hou, M. Xia, Densely multiplied physics informed neural network, *arXiv preprint arXiv:2402.04390* (2024).
- [24] Q. Zhang, C. Wu, A. Kahana, Y. Kim, Y. Li, G. E. Karniadakis, P. Panda, Artificial to spiking neural networks conversion for scientific machine learning, *arXiv preprint arXiv:2308.16372* (2023).
- [25] W. Guan, K. Yang, Y. Chen, Z. Guan, A dimension-augmented physics-informed neural network (dapinn) with high level accuracy and efficiency, *arXiv preprint arXiv:2210.13212* (2022).

- [26] S. Wang, H. Wang, P. Perdikaris, On the eigenvector bias of fourier feature networks: From regression to solving multi-scale pdes with physics-informed neural networks, *Computer Methods in Applied Mechanics and Engineering* 384 (2021) 113938.
- [27] A. Krishnapriyan, A. Gholami, S. Zhe, R. Kirby, M. W. Mahoney, Characterizing possible failure modes in physics-informed neural networks, *Advances in Neural Information Processing Systems* 34 (2021) 26548–26560.
- [28] C. L. Wight, J. Zhao, Solving allen-cahn and cahn-hilliard equations using the adaptive physics informed neural networks, *arXiv preprint arXiv:2007.04542* (2020).
- [29] A. D. Jagtap, K. Kawaguchi, G. E. Karniadakis, Adaptive activation functions accelerate convergence in deep and physics-informed neural networks, *Journal of Computational Physics* 404 (2020) 109136.
- [30] N. Sukumar, A. Srivastava, Exact imposition of boundary conditions with distance functions in physics-informed deep neural networks, *Computer Methods in Applied Mechanics and Engineering* 389 (2022) 114333.
- [31] C. Leake, D. Mortari, Deep theory of functional connections: A new method for estimating the solutions of partial differential equations, *Machine learning and knowledge extraction* 2 (1) (2020) 37–55.
- [32] L. Lu, R. Pestourie, W. Yao, Z. Wang, F. Verdugo, S. G. Johnson, Physics-informed neural networks with hard constraints for inverse design, *SIAM Journal on Scientific Computing* 43 (6) (2021) B1105–B1132.
- [33] X. Jin, S. Cai, H. Li, G. E. Karniadakis, Nsfnets (navier-stokes flow nets): Physics-informed neural networks for the incompressible navier-stokes equations, *Journal of Computational Physics* 426 (2021) 109951.
- [34] S. Basir, Investigating and mitigating failure modes in physics-informed neural networks (pinns), *arXiv preprint arXiv:2209.09988* (2022).
- [35] L. Lu, X. Meng, Z. Mao, G. E. Karniadakis, Deepxde: A deep learning library for solving differential equations, *SIAM review* 63 (1) (2021) 208–228.
- [36] C. Wu, M. Zhu, Q. Tan, Y. Kartha, L. Lu, A comprehensive study of non-adaptive and residual-based adaptive sampling for physics-informed neural networks, *Computer Methods in Applied Mechanics and Engineering* 403 (2023) 115671.
- [37] B. Zapf, J. Haubner, M. Kuchta, G. Ringstad, P. K. Eide, K.-A. Mardal, Investigating molecular transport in the human brain from mri with physics-informed neural networks, *Scientific Reports* 12 (1) (2022) 15475.

- [38] S. Deguchi, M. Asai, Dynamic & norm-based weights to normalize imbalance in back-propagated gradients of physics-informed neural networks, *Journal of Physics Communications* 7 (7) (2023) 075005.
- [39] S. Wang, S. Sankaran, P. Perdikaris, Respecting causality for training physics-informed neural networks, *Computer Methods in Applied Mechanics and Engineering* 421 (2024) 116813.
- [40] Z. Xiang, W. Peng, X. Liu, W. Yao, Self-adaptive loss balanced physics-informed neural networks, *Neurocomputing* 496 (2022) 11–34.
- [41] D. Liu, Y. Wang, A dual-dimer method for training physics-constrained neural networks with minimax architecture, *Neural Networks* 136 (2021) 112–125.
- [42] S. J. Anagnostopoulos, J. D. Toscano, N. Stergiopoulos, G. E. Karniadakis, Residual-based attention in physics-informed neural networks, *Computer Methods in Applied Mechanics and Engineering* 421 (2024) 116805.
- [43] L. McClenny, U. Braga-Neto, Self-adaptive physics-informed neural networks using a soft attention mechanism, *arXiv preprint arXiv:2009.04544* (2020).
- [44] S. Basir, I. Senocak, Physics and equality constrained artificial neural networks: Application to forward and inverse problems with multi-fidelity data fusion, *Journal of Computational Physics* 463 (2022) 111301.
- [45] S. Basir, I. Senocak, An adaptive augmented lagrangian method for training physics and equality constrained artificial neural networks, *arXiv preprint arXiv:2306.04904* (2023).
- [46] G. Zhang, H. Yang, F. Zhu, Y. Chen, et al., Dasa-pinns: Differentiable adversarial self-adaptive pointwise weighting scheme for physics-informed neural networks, *SSRN* (2023).
- [47] H. Son, S. W. Cho, H. J. Hwang, Enhanced physics-informed neural networks with augmented lagrangian relaxation method (al-pinns), *Neurocomputing* (2023) 126424.
- [48] N. Tishby, F. C. Pereira, W. Bialek, The information bottleneck method, *arXiv preprint physics/0004057* (2000).
- [49] N. Tishby, N. Zaslavsky, Deep learning and the information bottleneck principle, in: *2015 IEEE Information Theory Workshop (ITW)*, IEEE, 2015, pp. 1–5.
- [50] R. Shwartz-Ziv, N. Tishby, Opening the black box of deep neural networks via information, *arXiv preprint arXiv:1703.00810* (2017).
- [51] Z. Goldfeld, Y. Polyanskiy, The information bottleneck problem and its applications in machine learning, *IEEE Journal on Selected Areas in Information Theory* 1 (1) (2020) 19–38.

- [52] R. Shwartz-Ziv, Information flow in deep neural networks, arXiv preprint arXiv:2202.06749 (2022).
- [53] A. M. Saxe, Y. Bansal, J. Dapello, M. Advani, A. Kolchinsky, B. D. Tracey, D. D. Cox, On the information bottleneck theory of deep learning, *Journal of Statistical Mechanics: Theory and Experiment* 2019 (12) (2019) 124020.
- [54] M. A. Alomrani, A critical review of information bottleneck theory and its applications to deep learning, arXiv preprint arXiv:2105.04405 (2021).
- [55] Z. Goldfeld, E. v. d. Berg, K. Greenewald, I. Melnyk, N. Nguyen, B. Kingsbury, Y. Polyanskiy, Estimating information flow in deep neural networks, arXiv preprint arXiv:1810.05728 (2018).
- [56] S. S. Lorenzen, C. Igel, M. Nielsen, Information bottleneck: Exact analysis of (quantized) neural networks, arXiv preprint arXiv:2106.12912 (2021).
- [57] P. Domingos, Every model learned by gradient descent is approximately a kernel machine, arXiv preprint arXiv:2012.00152 (2020).
- [58] X. Peng, L. Li, F.-Y. Wang, Accelerating minibatch stochastic gradient descent using typicality sampling, *IEEE transactions on neural networks and learning systems* 31 (11) (2019) 4649–4659.