

Efficient Neural PDE-Solvers using Quantization Aware Training

Winfried van den Dool¹, Tijmen Blankevoort², Max Welling¹, Yuki M. Asano¹
¹QUVA Lab at the University of Amsterdam, ²Qualcomm AI Research*

w.v.s.o.vandendool@uva.nl

Abstract

In the past years, the application of neural networks as an alternative to classical numerical methods to solve Partial Differential Equations has emerged as a potential paradigm shift in this century-old mathematical field. However, in terms of practical applicability, computational cost remains a substantial bottleneck. Classical approaches try to mitigate this challenge by limiting the spatial resolution on which the PDEs are defined. For neural PDE solvers, we can do better: Here, we investigate the potential of state-of-the-art quantization methods on reducing computational costs. We show that quantizing the network weights and activations can successfully lower the computational cost of inference while maintaining performance. Our results on four standard PDE datasets and three network architectures show that quantization-aware training works across settings and three orders of FLOPs magnitudes. Finally, we empirically demonstrate that Pareto-optimality of computational cost vs performance is almost always achieved only by incorporating quantization.

1. Introduction

In many scientific fields, mathematical models of observed phenomena are expressed in terms of Partial Differential Equations, with the solution commonly represented by a function of a single time variable and one or more spatial variables. While many PDEs can be described compactly (*e.g.* the famous Schrödinger’s equation from quantum physics, or Navier-Stokes’ equations from fluid dynamics) it is usually impossible to write down explicitly the formulas for their exact solutions. Rather, there is a vast amount of scientific research on methods to numerically approximate solutions.

When it comes to practical applications, computational cost and available resources play a significant role in this field of research. For most methods a natural trade-off occurs where one may use available resources to either speed

up the computation, or increase the resolution on which the PDE is defined, leading to a slower but more accurate result. This latter strategy is of particular interest for PDEs that are highly non-linear and potentially exhibit chaotic behavior. Small eddies in a turbulent fluid may, for example, be overlooked when the spatial resolution is too low, leading to increasing errors in the solution function over time. One of the most important fields where computational resources form a bottleneck is climate prediction. The challenge of determining how many degrees the planet will warm over the next decades is a computational one, with supercomputers running massive ensemble-type methods. Insufficient resolution is quite often the main cause for sub-optimal modeling accuracy [21, 1, 16]. Indeed, over time climate models have become better for a large part simply because the resolution of the grid that they are defined on was allowed to increase, due to more efficient computation and more available resources [24, 26, 11].

It is precisely in this context of the computational cost of the traditional solvers that neural PDE solvers are becoming an interesting alternative. A single forward pass in a neural network is extremely fast compared to the iterative solving procedure of classical methods, and most of the operations are large matrix multiplications that can be easily handled by GPUs. While a neural network does need training first, possibly requiring data generated by classical solvers, its efficiency benefit comes from “recycling”: After training is complete, the neural network has learned to generalize to different initial conditions, whereas conventional solvers only solve the PDE for one specific configuration of initial conditions at a time. In practice, a deep learning weather model would have to be trained only once, but could then be applied every day to predict the next day’s weather. Furthermore, when training data is sufficiently abundant in the real world, using that data directly comes with the added advantage that complex phenomena do not need to be first accurately modeled. Indeed, neural PDE solvers trained on real weather data alone have shown promising results [3, 7, 16].

Our goal is not to improve on state-of-the-art models in the conventional sense. Rather, our starting point is to borrow from existing networks and architectures and focus on

*an initiative of Qualcomm Technologies, Inc.

their reducing their inference cost. With neural PDE solvers becoming mainstream for practical use, this measure becomes more important than test loss alone. Our contribution is in finding the optimal way to deal with the loss and computational cost trade-off. There are several orthogonal approaches to achieve lower computational costs. On the one hand, as in classical PDE solvers, the foremost way to make computation manageable is to reduce the resolution on which the PDE is defined. On the other hand, deep learning researchers have developed a variety of different techniques specifically for neural networks, most notably quantization, the focus of this work. Overall, we make the following contributions:

- We provide an evaluation of 3 neural network-based PDE solvers [9, 17, 6] under 5 weight quantization scenarios. We use state-of-the-art quantization methods to provide exemplary benchmarks on 4 of the most common datasets.
- We are the first to investigate both spatial resolution and model (weights and activations) resolution simultaneously as hyper-parameters in designing neural PDE-solvers. We develop a hybrid approach of applying quantization and modifying spatial resolution depending on the problem to achieve compute- and accuracy-optimal results.
- We provide an extensive analysis of the trade-off between computational cost and errors for different quantized models and demonstrate that a certain level of quantization is almost always necessary to be Pareto-optimal on the accuracy-cost curve.

2. Related Works

Neural PDE solvers. There are already many deep learning methods that aim to solve PDEs [6, 23, 9, 4, 5, 20, 10]. We base our research on the most common neural PDE surrogate architectures, namely the Fourier Neural Operator [17], UNet [23, 9] and a type of Transformer [6]. In terms of datasets, there are two recent larger-scale attempts at proper benchmarking of the field, namely PDEArena [9] and PDEBench [27]; from both of which we borrow our training and evaluation datasets.

Quantization. With regards to quantization methods we rely on AIMET [25]¹, which provides state-of-the-art quantization techniques. In particular, we utilize Quantization Aware Training [12, 13], utilizing the straight-through estimator [2] to approximate gradients of rounding operators. We allow quantization ranges to be trainable parameters as well, as in [8].

¹AIMET is a product of Qualcomm Innovation Center, Inc. (BSD-3)

Low precision PDE solving. Recently, there has been increasing interest in running classical methods for solving climate models with lower precision as well. In [14] reduced precision of Float32 and Float16 has been used, as opposed to the standard precision of Float64. However, these solvers are based on traditional approaches, and not on neural networks.

In contrast to previous PDE surrogates, i.e. (deep learning) functions that aim to approximate the original PDE solution, we especially care for computational cost. While it is generally investigated that quantization decreases the accuracy of deep learning models as well as that standard PDE solvers’ accuracy depends on their grid size, we explicitly compare both techniques in terms of how much they reduce computational costs. To make this possible we take a particular interest in the cost-loss trade-off of both techniques separately as well as their combined effects. Also, while a classification algorithm may exhibit a clear breaking point at which it can no longer accurately predict the majority of correct labels, the prediction of a PDE solution is measured in terms of (mean squared) errors between functions. As such there is no clear definition of what a “correct” prediction is, and consequentially there is no clear measure of when a PDE surrogate is successful. This means that a wide regime of models of varying qualities becomes interesting to investigate, both in the low-cost as well as the low-loss regions of the cost-loss trade-off.

3. Method

3.1. Theoretic background and notations

As we work with synthetic data, we are given numeric solutions to a given known PDE. Let $X \subset \mathbb{R}^d$ be a spatial domain and $[0, T]$ be a time window on which $\mathbf{u} : X \times [0, T] \rightarrow \mathbb{R}^n$ is the solution to the partial differential equation

$$\frac{\partial \mathbf{u}}{\partial t} = F \left(\mathbf{x}, \mathbf{u}, \frac{\partial \mathbf{u}}{\partial \mathbf{x}}, \frac{\partial^2 \mathbf{u}}{\partial \mathbf{x}^2}, \dots \right), \quad (1)$$

with initial condition $\mathbf{u}^0(\mathbf{x})$ at time $t = 0$ and boundary conditions defined by the operator $\mathcal{B}[\mathbf{u}](t, \mathbf{x}) = 0$ when \mathbf{u} is on the boundary ∂X of the domain X . Here F can be any function, though we specify the precise form for different datasets in the Supplementary Information (SI). Assume that at some time t , $\mathbf{u}(t, \mathbf{x})$, as well as prior values of \mathbf{u} , are known. Let $\tau > 0$ be some step-size with $t + \tau < T$. We can then try to predict one time step into the future by

$$\mathbf{u}(t + \tau, \mathbf{x}) = \mathbf{u}(t, \mathbf{x}) + \int_t^{t+\tau} \frac{\partial \mathbf{u}(t, \mathbf{x})}{\partial t} dt. \quad (2)$$

$$\forall t \in [0, T - \tau], \mathbf{x} \in X$$

Given that theoretically the value of $\mathbf{u}(t, \mathbf{x})$ determines, through its spatial derivatives and the PDE, all future values, a

first attempt at defining a useful operator may be to replace the right-hand side of the above equation by an operator \mathcal{F}_τ such that

$$\mathbf{u}(t + \tau, \mathbf{x}) = \mathcal{F}_\tau[\mathbf{u}(t, \mathbf{x})] \quad \forall t \in [0, T - \tau], \mathbf{x} \in X, \quad (3)$$

indeed ignoring values of $\mathbf{u}(t', \mathbf{x})$ for $t' < t$. We are then interested in approximating this operator \mathcal{F}_τ by a neural network \mathcal{G} . However, in practice the spatial and time domains are discretized on regular grids $\mathcal{X} \subset X$ and $\mathcal{T} \subset [0, T]$ with $|\mathcal{X}| = N_x$ and $|\mathcal{T}| = N_t$. Crucially, when $\mathbf{u}(t, \mathbf{x})$ is defined on \mathcal{X} rather than the full domain X , its spatial derivatives are no longer determined, breaking the connection from $\mathbf{u}(t, \mathbf{x})$ to $\frac{\partial \mathbf{u}}{\partial t}$ through the PDE F . Consequentially, the full future cannot be exactly determined from $\mathbf{u}(\mathbf{x}, t)$ defined on \mathcal{X} alone. Note that even knowing the spatial derivatives exactly on the grid \mathcal{X} at time t is not enough, as they are required at all later times as well if we are to solve the integral in equation 2. Now the values of $\mathbf{u}(\mathbf{x}, t')$ for $t' < t$ may still provide helpful information, and we may want to include those as inputs in our neural network approximation of \mathcal{F} . One can think for example that such values may allow the backward difference approximation of $\partial \mathbf{u}(t, \mathbf{x}) / \partial t$ as an alternative to using the PDE F with the approximated spatial derivatives.

Incorporating multiple past time steps, the neural network approximation is given by the operator \mathcal{G} , defined by

$$\mathbf{u}(t_n + \tau, \mathbf{x}) = \mathcal{G}_\tau[(\mathbf{u}(t_i, \mathbf{x}))_{i=1}^n] \quad (4)$$

for $t_i, t_n + \tau \in \mathcal{T}, \mathbf{x} \in \mathcal{X}$.

In what follows we assume \mathbf{u} to be defined on $\mathcal{T} \times \mathcal{X}$.

3.2. Training setup

Let the true $\mathbf{u}(\mathbf{x}, t)$ be defined everywhere on our grid $\mathcal{T} \times \mathcal{X}$. For each such trajectory $\mathbf{u}(\mathbf{x}, t)$ in a mini-batch we first randomly select a starting time, after which we take a fixed set of subsequent input indices $\mathcal{T}_{\text{input}} \subset \mathcal{T}$, and target indices, $\mathcal{T}_{\text{target}} \subset \mathcal{T}$. We have $N_t - |\mathcal{T}_{\text{input}}|$ input steps available, but we are only training on $|\mathcal{T}_{\text{target}}|$ time steps. Therefore, to make full use of the available data we repeat each epoch $\lceil (N_t - |\mathcal{T}_{\text{input}}|) / |\mathcal{T}_{\text{target}}| \rceil$ times, selecting new random indices every time so that in expectation the full dataset is used every epoch.

For $|\mathcal{T}_{\text{target}}| > 1$ outputs we may use either temporal bundling (letting \mathcal{G} output multiple subsequent time steps in one forward pass, see for instance [5]), or a recurrent approach, where the last outputs are fed to the network as new inputs, or a combination of both. Similar to [5] we may apply *pushforward*, meaning that we only backpropagate the loss through the last part of the target indices.

Let $\mathbf{u}(t, \mathbf{x})$ be the true targets and $\hat{\mathbf{u}}(t, \mathbf{x})$ the corresponding neural network predictions for $(t, \mathbf{x}) \in \mathcal{T}_{\text{target}} \times \mathcal{X}$.

The loss, assuming no pushforward is applied, is then defined by

$$\frac{\sum_{\mathbf{x} \in \mathcal{X}} \sum_{t \in \mathcal{T}_{\text{target}}} \sum_{i=1}^{N_{\text{fields}}} \|\mathbf{u}_i(\mathbf{x}, t) - \hat{\mathbf{u}}_i(\mathbf{x}, t)\|_2^2}{|\mathcal{X}| |\mathcal{T}_{\text{target}}| N_{\text{fields}}}, \quad (5)$$

where N_{fields} is the dimensionality of $\mathbf{u}(t, \mathbf{x})$.

3.3. Quantization

By default, the weights and activations of the neural network \mathcal{G} are defined as floating point numbers using 32-bit precision. We can, however, store the weights and activations as integer values, to make (matrix) multiplications much more efficient. A floating-point weight matrix \mathbf{W} can be expressed as a single scalar multiplied by a matrix of integer values, and a remainder term ϵ .

$$\mathbf{W} = s_W \cdot \mathbf{W}_{\text{int}} + \epsilon_W$$

Similarly a vector of activations \mathbf{v} may be expressed as

$$\mathbf{v} = s_v \cdot (\mathbf{v}_{\text{int}} - \mathbf{z}_v) + \epsilon_v,$$

where we have added a zero-point \mathbf{z}_v to allow asymmetric quantization of the activations. Matrix-vector multiplication now becomes:

$$\begin{aligned} \mathbf{W}\mathbf{v} &= (s_W \cdot \mathbf{W}_{\text{int}} + \epsilon_W)(s_v \cdot (\mathbf{v}_{\text{int}} - \mathbf{z}_v) + \epsilon_v) \quad (6) \\ &= s_W s_v \mathbf{W}_{\text{int}} \mathbf{v}_{\text{int}} - s_W s_v \mathbf{W}_{\text{int}} \mathbf{z}_v + \text{error terms} \end{aligned}$$

As the second term depends only on the weight matrix \mathbf{W} and quantization parameters s_W, s_v and \mathbf{z}_v , it can be computed beforehand. If we then ignore the error terms, only the first term is remaining during inference: a matrix-vector multiplication of integer values. For a given bitwidth b , there are 2^b possible integer values to choose from. Given a zero-point z and scale factor s , the quantization grid limits q_{min} and q_{max} are then determined by $-sz$ and $s(2^b - 1 - z)$. Any values lying outside this range will be clipped to its limits, incurring a *clipping error*. The full quantization function $q(\cdot)$ is given by

$$\mathbf{x}_{\text{int}} = q(\mathbf{x}; s_x, \mathbf{z}_x, b) = \mathbf{C} \left(\left\lfloor \frac{\mathbf{x}}{s_x} \right\rfloor + \mathbf{z}_x; 0, 2^b - 1 \right), \quad (7)$$

where $\lfloor \cdot \rfloor$ is the round-to-nearest operator and \mathbf{C} is a clamping function defined as:

$$\mathbf{C}(x; a, c) = \begin{cases} a, & x < a \\ x, & a \leq x \leq c \\ c, & x > c \end{cases} \quad (8)$$

To get a quantized network we start with a pre-trained floating point network and then optimize for the values of

s and z that minimize the error terms that result from both clipping and rounding errors. This happens independently per layer, using arbitrary dummy data or training data in the forward pass. Next, we also apply Quantization Aware Training (QAT). This is a fine-tuning step, further training the quantized network using stochastic gradient descent on the original loss function. Training with quantized weights and activations is possible using the straight-through estimator [2].

Noting that the majority of the computational cost comes from (matrix) multiplications, we only quantize the inputs and weights of neural network layers that are based on matrix multiplications, and let the outputs (together with the biases) be accumulated in floating point format.

3.4. Changing resolutions

We change resolutions using bilinear interpolation for the 2D datasets and linear interpolation for the 1D datasets. We leave experimentation with learnable resolution operators for future research. However, we have chosen (bi)linear interpolation deliberately for its speed, which is the prime reason to change resolutions in the first place. If the network is defined on a different resolution than the data, we apply resolution-altering operators before and after it during both validation and training, thus also backpropagating losses through the resolution-altering operators.

3.5. Model cost computation

The proxy we use for model inference efficiency is based on counting multiplication and addition operations per layer. For a given quantized network module with M multiplication and A addition operations, and integer bitwidths of b_w , b_a for the weights and activations respectively, the cost is defined as

$$M \cdot b_w \cdot b_a + A \cdot b_a,$$

where multiplications are considered between weights and inputs, and addition operations only apply to outputs. (We do not quantize bias vectors.) The full network cost is the sum of the cost over all layers. If a network works on a different resolution, we also take the costs of altering the resolution before and after the forward pass into account. Details on the number of multiplications and additions per layer are provided in the code. We have observed no increased losses from quantizing any parameter to bitwidth 16 under any circumstance, and therefore assume that all floating point operations can be harmlessly replaced by their Int16 counterpart. This enables us to measure the computational cost of non-quantized operations as if they were Int16 operations and to generally use fixed-point integer operations as a measure throughout all our comparisons of different model inference costs. To find the number of multiplication and addition operations in a given module we use [22],

slightly adapted for our needs. The deepspeed model profiler lets you choose either FLOPs or MACs as a measure for module cost. We use the fact that one MAC operation consists of a single addition and multiplication and assume that the remaining operations are all addition-type in terms of complexity. Hence, when deepspeed outputs X FLOPs and Y MACs for a certain network module, we know that there are Y multiplications and $X - Y$ remaining FLOPs that we consider additions. A few exceptions and special cases, as well as deliberate deviations from deepspeed, can be found in Appendix A.3.

4. Experiments

We use three popular neural-surrogate architectures: FNO, UNet and Transformer, applied to datasets based on 4 PDEs: Diffusion-Sorption (1D), Burgers’ (1D), Navier-Stokes (2D) and Darcy (2D). The datasets are described in more detail in Appendix A.1. We show benchmark results of the original networks, as well as the results of the networks after quantization and scaling. For each result, we present the test loss as well as the computational cost of inference. We first compare quantization and reduction of spatial resolution as two orthogonal approaches to decrease computational cost, showing their impacts separately. Finally, we apply both methods simultaneously on higher-resolution versions of the datasets.

4.1. Implementation

For a given spatial resolution, we first train the floating point network in a regular fashion, i.e., aiming for the lowest possible loss. We use Adam [15] with weight decay and a cosine annealing learning rate scheduler [18] with a linear warmup. Next, we quantize the network, first finding the best quantization parameters using AIMET’s built-in optimization tool on 20% of the training data. Then we fine-tune the quantized network, training again using Adam and cosine annealing with linear warmup, but without weight decay, with a smaller learning rate, and for fewer epochs. Specific parameters for training and fine-tuning are provided in Appendix A.2, as well as the length of the input and output trajectories, i.e., the number of time steps used. If a model that operates on a different resolution than the original data resolution is unrolled several times, i.e., when training and testing on a longer output trajectory, the resolution is only altered after the first input, and after the output just before the loss is computed: intermediate values that are fed back into the network are kept in the network resolution. The quantization regimes used throughout the experiments are: [w4a4, w4a8, w8a8, w8a16], with wXaY referring to quantizing weights to bitwidth X and activations to bitwidth Y.² We do not quantize the

²We found it better to use a higher activation bitwidth than weight bitwidth when using different bitwidths for weights and activations

Val. MSE loss	FNO	UNet	Transformer
Diff-Sorp (1D)	4.12e-8	4.26e-8	4.33e-8
Burgers' (1D)	5.62e-4	2.19e-3	1.40e-3
Navier-Stokes (2D)	2.85e-3	1.10e-3	5.36e-3
Darcy (2D)	1.52e-2	6.92e-3	8.47e-3

Table 1: **Overview of unquantized model performances.** We report Mean Squared Error (MSE) loss for the unquantized neural networks that we analyze in this paper.

FLOPS (in millions)	FNO	UNet	Transformer
Diff-Sorp (1D)	15	7	4
Burgers' (1D)	21	14	9
Darcy (2D)	150	34	69
Navier-Stokes (2D)	210	181	118

Table 2: **Model FLOPS across different datasets.** In this paper, we evaluate how we can increase the efficiency of models that span almost three orders of magnitude in FLOP inference cost.

first and last layer of each model, having observed that this improves performance with negligible impact on inference cost. Unless mentioned otherwise, we use scaling regimes, i.e., factors by which we scale input data resolution, of [0.7, 0.5, 0.3, 0.2]. For 2D data, we scale both dimensions by this factor, thus the actual number of spatial points depends quadratically on the scaling factor.

4.2. Results

To briefly introduce the models and datasets³ used we present an overview of the unquantized, standard-resolution models' loss and inference cost on all datasets in Tables 1 and 2. Here we use the number of floating point operations (flops), as taken from deepspeed [22], directly as a measure of model inference cost. We note that, as there is no clear measure of when a PDE surrogate is successful (compared to *e.g.* classification algorithms), a wide regime of varying MSE losses and costs is interesting to investigate. In Tables 1 and 2 we can see that the PDEs have vastly different possible solutions in terms of MSE and cost.

4.2.1 Quantization compared to reducing resolution

In what follows we refer to models being applied on reduced resolution as "scaled" models, although it is actually their input data that is scaled. Implicitly such models are smaller as a result of them being applied on lower-resolution data, but we do not change hyperparameters like layer width or number of layers.

³Datasets were solely downloaded and evaluated by QUVA.

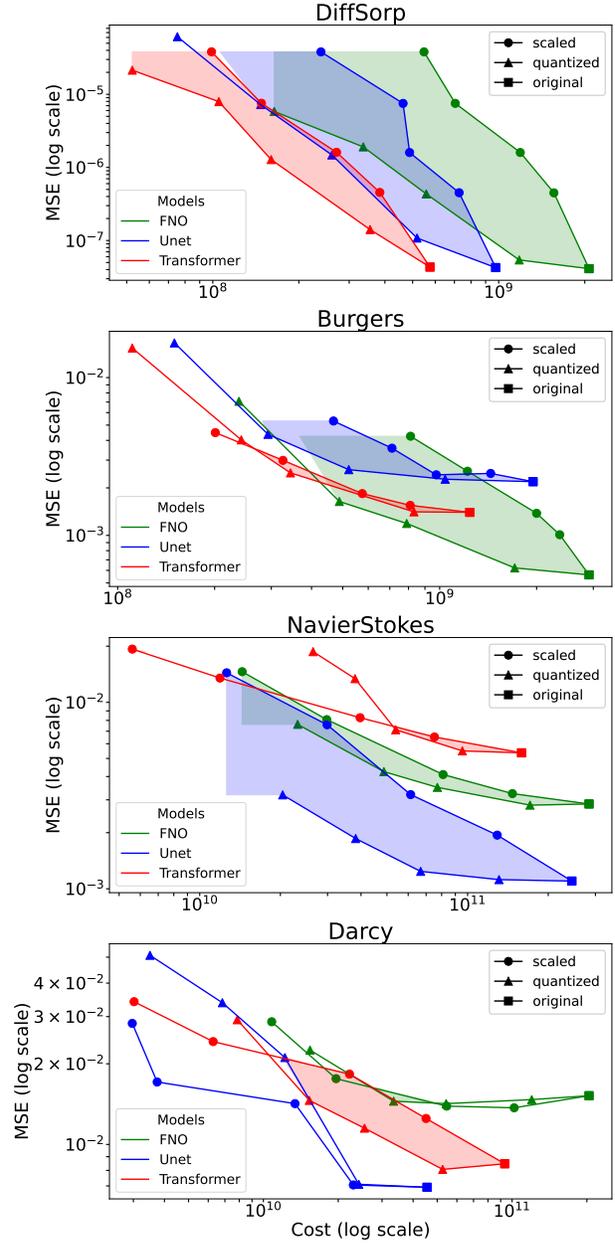


Figure 1: **Quantization for more robust cost-reduction.** We compare model quantization against scaling spatial resolution as ways to reduce compute costs. We show the cost in FLOPS versus the Mean Squared Error (MSE) across 4 PDE datasets and three model types.

The results of our first set of experiments are presented in Figure 1, where we compare quantization against spatial scaling across models and architectures. The plots can be interpreted as follows: If the triangles (quantized models) lie below the circles (scaled models), quantization is a more robust way to reduce inference costs. We observe this is

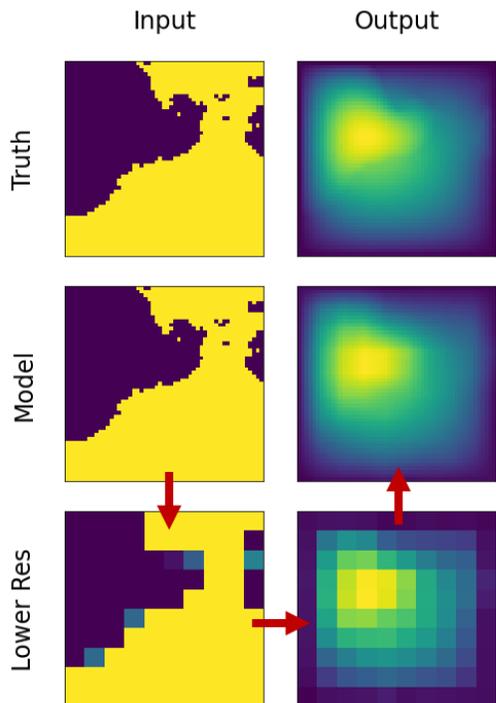


Figure 2: **Spatial resolution reduction on Darcy PDE.** Example of a scaled FNO model on Darcy data. For clarity, we represent the full forward pass of what we describe as a “scaled model” by the red arrows. Note that for this PDE, the true output, i.e. the target, is actually quite blurry to begin with, making the substantial lowering of resolution that is applied appear less problematic.

the case for all model and dataset combinations, except for the U-Net on the Darcy dataset. The reason why scaling is relatively more successful on the Darcy dataset can be understood by looking at some examples. Comparing Figure 3 to Figure 2, the Darcy targets are already quite blurry. Generally, the dataset appears less sensitive to details in the input, making it potentially less likely that important information is lost in the process of changing spatial resolution in a forward pass.

4.2.2 Combining quantization and scaling

The datasets we use have been generated in a higher resolution than is actually used in the experiments. Note that this is common practice: the numerical solvers are applied to a much finer grid to make sure that the data is a relatively accurate representation of the underlying PDE [9, 27, 17]. However, we now consider the hypothetical situation where one is indeed interested in minimizing loss on the higher resolution. Results for these experiments are found in Figures 4.

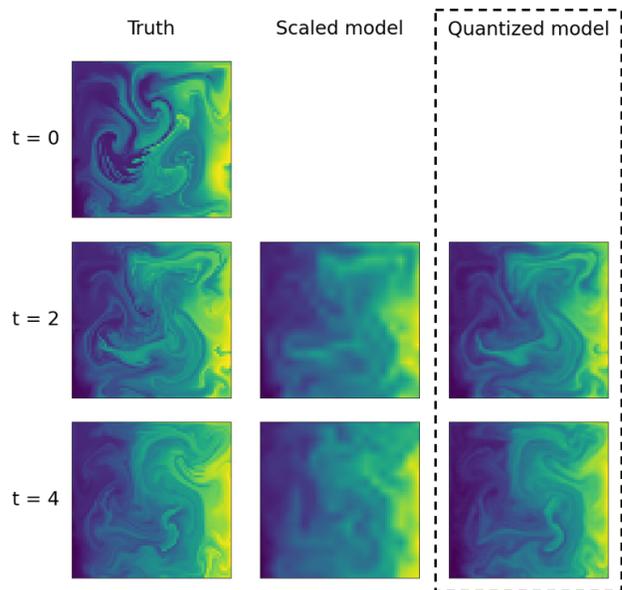
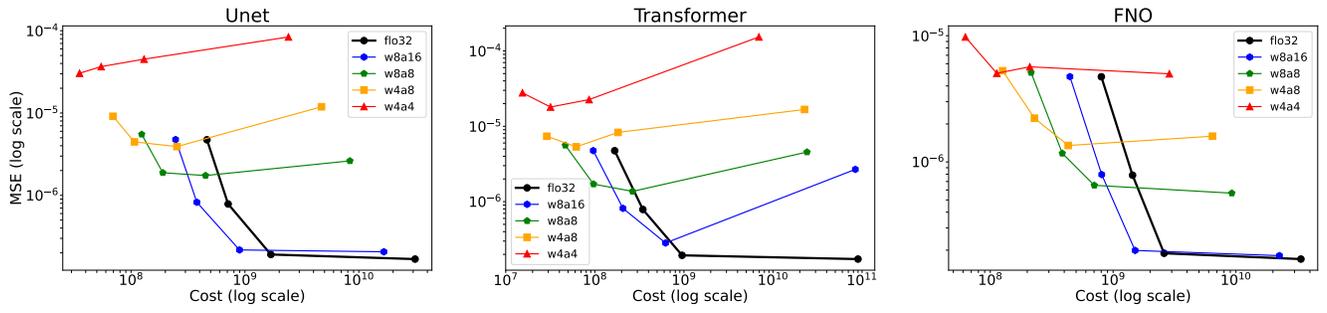


Figure 3: **Quantization retains details across time.** Example of two different UNet prediction unrollings on Navier-Stokes data. The scaled version uses a scaling factor of 0.3, the quantized version is based on w4a8 quantization. The left column ($t = 0$) represents the last input time step used by the models, the other columns are predictions. Both models have similar inference costs, yet we see the quantized model showing more detailed predictions.

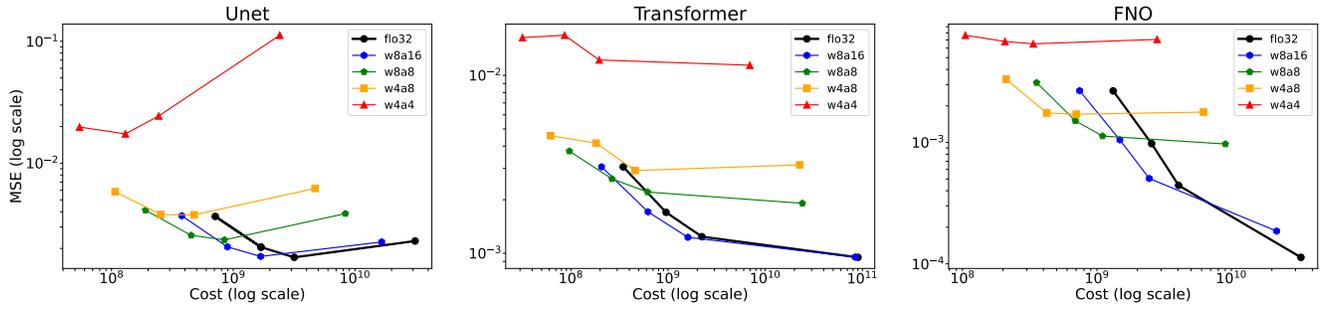
First, we find that when the original resolution is relatively high, reducing it may not affect performance very much for some of these datasets (e.g., Figure 4a), and yields a good reducing in compute cost. However, as more scaling is applied, we note that switching to a lower-bitwidth quantization regime quickly becomes better than scaling alone. In particular, we find that across almost all PDEs and neural PDE solvers, the cost vs error Pareto-optimal curve is achieved by quantization. In other words, we find that quantization is another “knob” that one can turn to reduce costs, besides only changing resolution. Not only is this “knob” available, but the results show that is essential for optimal efficiency. Note that the particular choice of how much quantization and how much scaling ought to be applied depends on the operating point, that can, for example, be specified by a lower bound on the MSE.

5. Discussion and conclusion

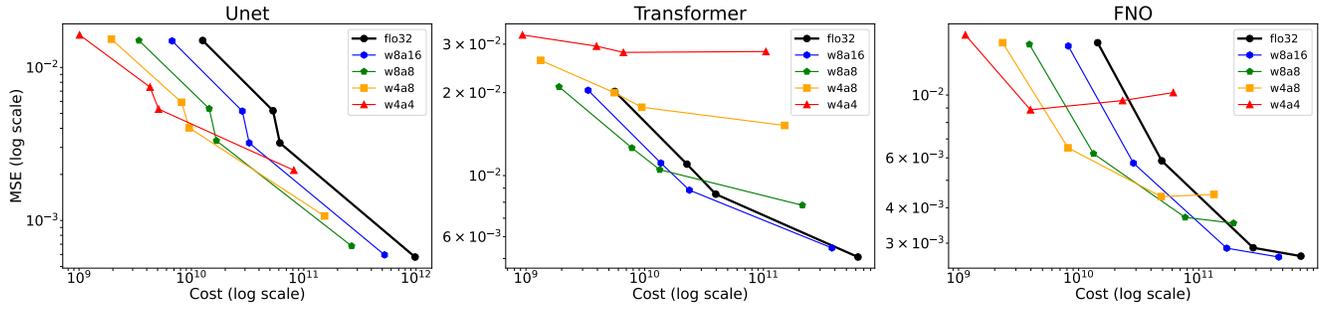
We have shown that quantization can be an effective solution to make neural PDE-solvers more cost-efficient. In particular, it can achieve better results in terms of prediction loss vs inference cost trade-off than the conventional method of making computational cost manageable, which is reducing spatial resolution. To some extent the results



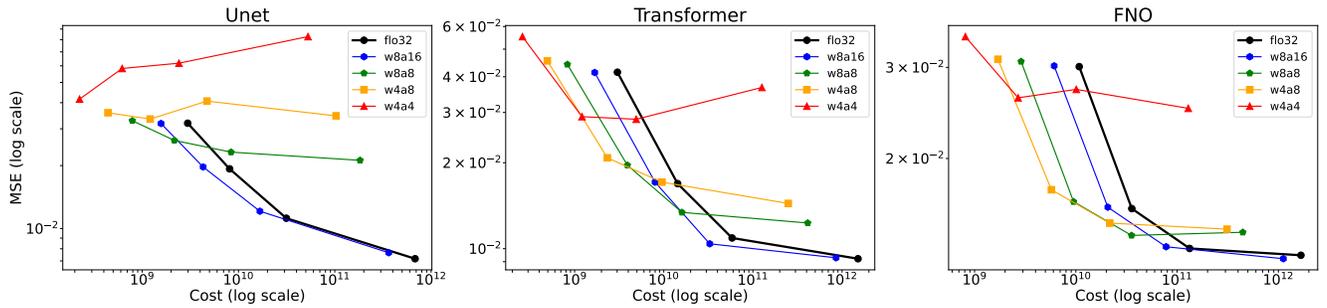
(a) 1D Diffusion-Sorption Equations.



(b) 1D Burgers' Equations.



(c) 2D Navier-Stokes Equations.



(d) 2D Darcy Equations.

Figure 4: **Optimal weight and activation quantization levels yield Pareto-optimal performances.** We vary quantization levels and spatial scaling factors and observe Pareto-optimal performances compared to non-quantized (gray) models. The four points on each line correspond to different scaling levels; the right-most point is the original resolution.

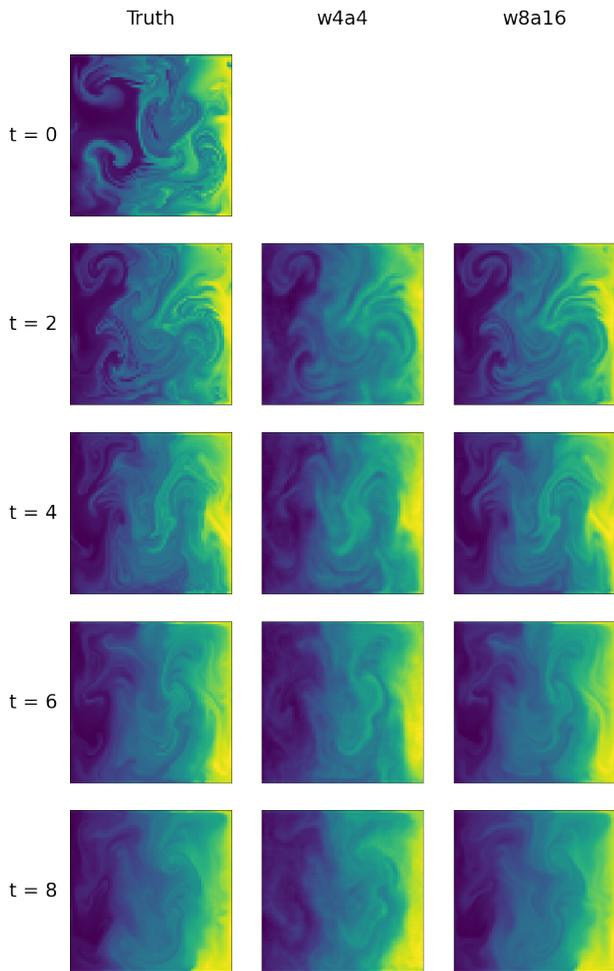


Figure 5: **Two differently quantized models unrolled for more timesteps.** Example of two different UNet prediction unrollings on Navier-Stokes data.

depend on the type of data at hand: If the original data is of extremely, and perhaps unnecessarily, high resolution, there may be at first a clear benefit to be obtained by simply reducing spatial resolution. However, we similarly find that there is no increase in loss by quantizing network weights and activations to `Int16`, with weights even quantized to `Int8` in some cases. We note that both cost-reducing methods, quantization and reducing spatial resolution, have a certain breaking point, depending on the model and dataset at hand, at which their effect on prediction loss increases rapidly. Our results indicate that when one of these methods reaches its breaking point, one can further improve on the loss vs inference trade-off by continuing with the other method. To achieve Pareto-optimal solutions it is necessary to apply both methods simultaneously.

In this paper, we have manually selected levels of quantization and resolution scaling, showing empirically the im-

portance of quantization as a new knob that can be turned to achieve better results for similar inference cost in existing models. However, in future research, we will investigate how a neural PDE-solver network can automatically detect the appropriate scaling or quantization levels to optimize the prediction loss vs inference cost trade-off. Another direction for future research is to experiment with real-world weather and climate prediction data. Although our current research is based on toy data, we actually expect the results to be even more pronounced when working with real-world data: The recurring theme (see for instance [21, 19, 1]) is that insufficient computational power inhibits the use of higher data resolutions which would solve many modeling problems and prediction inaccuracies. We propose quantization as a strategy to free up computational resources, allowing models to be defined on these desired higher resolutions.

Acknowledgements

We thank SURF for the support in using the National Supercomputer Snellius. This work is financially supported by Qualcomm Technologies Inc., the University of Amsterdam and the allowance Top consortia for Knowledge and Innovation (TKIs) from the Netherlands Ministry of Economic Affairs and Climate Policy.

References

- [1] Peter Bauer, Alan Thorpe, and Gilbert Brunet. The quiet revolution of numerical weather prediction. *Nature*, 525(7567):47–55, 2015.
- [2] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv e-prints*, pages arXiv–1308, 2013.
- [3] Kaifeng Bi, Lingxi Xie, Hengheng Zhang, Xin Chen, Xiaotao Gu, and Qi Tian. Pangu-weather: A 3d high-resolution model for fast and accurate global weather forecast. *arXiv e-prints*, pages arXiv–2211, 2022.
- [4] Johannes Brandstetter, Rianne van den Berg, Max Welling, and Jayesh K Gupta. Clifford neural layers for pde modeling. *arXiv e-prints*, pages arXiv–2209, 2022.
- [5] Johannes Brandstetter, Daniel E Worrall, and Max Welling. Message passing neural pde solvers. In *International Conference on Learning Representations*, 2021.
- [6] Shuhao Cao. Choose a transformer: Fourier or galerkin. *Advances in neural information processing systems*, 34:24924–24940, 2021.
- [7] L. Espoholt, S. Agrawal, and C. et al. Sønderby. Deep learning for twelve hour precipitation forecasts. *Nat Commun*, 13(1):5145, 2022.
- [8] Steven K Esser, Jeffrey L McKinstry, Deepika Bablani, Rathinakumar Appuswamy, and Dharmendra S Modha. Learned step size quantization. In *International Conference on Learning Representations*, 2020.
- [9] Jayesh K. Gupta and Johannes Brandstetter. Towards multi-spatiotemporal-scale generalized pde modeling, 2022.
- [10] Shudong Huang, Wentao Feng, Chenwei Tang, and Jiancheng Lv. Partial differential equations meet deep neural networks: A survey. *arXiv e-prints*, pages arXiv–2211, 2022.
- [11] C. E. Iles, R. Vautard, J. Strachan, S. Joussaume, B. R. Eggen, and C. D. Hewitt. The benefits of increasing resolution in global and regional climate simulations for european climate extremes. *Geoscientific Model Development*, 13(11):5583–5607, 2020.
- [12] Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for efficient integer-arithmetic-only inference. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2704–2713, 2018.
- [13] Sambhav Jain, Albert Gural, Michael Wu, and Chris Dick. Trained quantization thresholds for accurate and efficient fixed-point inference of deep neural networks. *Proceedings of Machine Learning and Systems*, 2:112–128, 2020.
- [14] Tom Kimpson, E. Adam Paxton, Matthew Chantry, and Tim Palmer. Climate-change modelling at reduced floating-point precision with stochastic rounding. *Quarterly Journal of the Royal Meteorological Society*, 149(752):843–855, mar 2023.
- [15] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv e-prints*, pages arXiv–1412, 2014.
- [16] Thorsten Kurth, Shashank Subramanian, Peter Harrington, Jaideep Pathak, Morteza Mardani, David Hall, Andrea Miele, Karthik Kashinath, and Anima Anandkumar. Four-castnet: Accelerating global high-resolution weather forecasting using adaptive fourier neural operators. In *Proceedings of the Platform for Advanced Scientific Computing Conference*, pages 1–11, 2023.
- [17] Zongyi Li, Nikola Borislavov Kovachki, Kamyar Aziz-zadenesheli, Kaushik Bhattacharya, Andrew Stuart, Anima Anandkumar, et al. Fourier neural operator for parametric partial differential equations. In *International Conference on Learning Representations*, 2020.
- [18] Ilya Loshchilov and Frank Hutter. SGDR: Stochastic gradient descent with warm restarts. In *International Conference on Learning Representations*, 2017.
- [19] Philipp Neumann, Peter Düben, Panagiotis Adamidis, Peter Bauer, Matthias Brück, Luis Kornbluh, Daniel Klocke, Bjorn Stevens, Nils Wedi, and Joachim Biercamp. Assessing the scales in numerical weather and climate predictions: will exascale be the rescue? *Philosophical Transactions of the Royal Society A*, 377(2142):20180148, 2019.
- [20] 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.
- [21] David A Randall, Richard A Wood, Sandrine Bony, Robert Colman, Thierry Fichefet, John Fyfe, Vladimir Kattsov, Andrew Pitman, Jagadish Shukla, Jayaraman Srinivasan, et al. Climate models and their evaluation. In *Climate change 2007: The physical science basis. Contribution of Working Group I to the Fourth Assessment Report of the IPCC (FAR)*, pages 589–662. Cambridge University Press, 2007.
- [22] Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD ’20*, page 3505–3506, New York, NY, USA, 2020. Association for Computing Machinery.
- [23] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention—MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18*, pages 234–241. Springer, 2015.
- [24] Tapio Schneider, João Teixeira, Christopher S Bretherton, Florent Brient, Kyle G Pressel, Christoph Schär, and A Pier Siebesma. Climate goals and computing the future of clouds. *Nature Climate Change*, 7(1):3–5, 2017.
- [25] Sangeetha Siddegowda, Marios Fournarakis, Markus Nagel, Tijmen Blankevoort, Chirag Patel, and Abhijit Khobare. Neural network quantization with ai model efficiency toolkit (aimet). *arXiv e-prints*, pages arXiv–2201, 2022.
- [26] J. Slingo, P. Bates, and P. et al. Bauer. Ambitious partnership needed for reliable climate prediction. *Nat. Clim. Chang.*, 12:499–503, 2022.

- [27] Makoto Takamoto, Timothy Praditia, Raphael Leiteritz, Daniel MacKinlay, Francesco Alesiani, Dirk Pflüger, and Mathias Niepert. Pdebench: An extensive benchmark for scientific machine learning. *Advances in Neural Information Processing Systems*, 35:1596–1611, 2022.

A. Appendix

A.1. PDEs and Data

We train our neural networks on datasets that contain solutions \mathbf{u} for different boundary and initial conditions so that it is able to generalize across these conditions, without the need for retraining. Datasets are obtained as follows:

- Generate a pair of initial conditions $\mathbf{u}^0(\mathbf{x})$ and boundary conditions $\mathcal{B}[\mathbf{u}](t, \mathbf{x}) = 0$ and evaluate these values on the relevant subsets of our grid $\mathcal{T} \times \mathcal{X}$.
- Use a conventional high-accuracy numerical solver to obtain $\mathbf{u}(t, \mathbf{x})$ for all $(\mathbf{x}, t) \in \mathcal{T} \times \mathcal{X}$.
- Pick a series of input indices, and subsequent target indices, from the time interval \mathcal{T} , starting from a possibly randomly chosen location. The values of $\mathbf{u}(\cdot, \mathbf{x})$ at these indices will form the inputs and targets in our training scheme.

Burger’s equation The Burger’s equation is a common PDE that arises in fluid dynamics and nonlinear wave phenomena. In 1D the PDE, given the domain that we use, is given by

$$\frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} = \frac{\nu}{\pi} \frac{\partial^2 u}{\partial x^2}, \quad x \in (0, 1) \quad t \in (0, 2] \quad (9)$$

where u represents the speed of the fluid at a certain place and time, and ν is the viscosity coefficient. The Burger’s equation describes the conservation of mass and momentum in a one-dimensional fluid flow, taking into account both convection effects ($u \frac{\partial u}{\partial x}$) and diffusion effects ($\nu \frac{\partial^2 u}{\partial x^2}$).

We use the 1D Burger’s equation dataset from [27]. It is defined with a spatial resolution of 1024, with periodic boundary conditions, and temporal resolution of 200. The dataset consists of 9000 train and 1000 test trajectories started from samples of different initial conditions that are formed using a superposition of randomly chosen sinusoidal waves. A viscosity coefficient of $\nu = 0.001$ is used.

Darcy’s Law The steady state 2D Darcy flow equation is a partial differential equation (PDE) that describes the flow of fluid through a porous medium. We use the PDE and domain expressed as

$$-\nabla(a(x)\nabla u(x)) = f(x), \quad x \in (0, 1)^2, \quad (10)$$

$$u(x) = 0, \quad x \in \partial(0, 1)^2,$$

where $a(x)$ is a diffusion coefficient based on the permeability of the porous medium and the dynamic viscosity of the fluid, $u(x)$ represents the pressure of the fluid, and f

represents any external sources or sinks of fluid within the domain. We set f to constant 1 and train an operator that maps $a(x)$ to the solution $u(x)$.

We use the Darcy flow dataset from [17]. It is defined on a spatial grid of 421×421 . We use 1024 train elements ($(a(x), u(x))$ pairs) and 100 validation elements. Details for how $a(x)$ is randomly generated for each data element can be found in [6].

Navier-Stokes equation The 2D Navier-Stokes Equation and domain that we use for our experiments is given by

$$\frac{\partial \mathbf{v}(\mathbf{x}, t)}{\partial t} = -\mathbf{v}(\mathbf{x}, t) \cdot \nabla \mathbf{v}(\mathbf{x}, t) + \nu \nabla^2 \mathbf{v}(\mathbf{x}, t) \quad (11)$$

$$- \nabla p(\mathbf{x}, t) + \mathbf{f}(\mathbf{x}), \quad x \in (0, 32)^2, \quad t \in (0, 21].$$

It describes the flow of a fluid in terms of its velocity components \mathbf{v} , the viscosity ν , and a buoyancy term \mathbf{f} . We assume incompressibility, so $\nabla \cdot \mathbf{v} = 0$, and Dirichlet boundary conditions ($\mathbf{v} = 0$).

The dataset is taken from [9]. A viscosity of $\nu = 0.01$ is used, and a buoyancy factor of $\mathbf{f} = (0, 0.5)^T$. While generating the data, the pressure field p is solved first, before subtracting its spatial gradients. In addition to the two velocity field components a scalar field $s(\mathbf{x})$ is introduced that is being transported through the velocity field. Its evolution is determined by

$$\frac{\partial s}{\partial t} = -\mathbf{v}(\mathbf{x}, t) \nabla s, \quad (12)$$

with Neumann boundaries $\frac{\partial s}{\partial x} = 0$ on the edge of the domain. For more details, see [9, 4]. The full dataset consists of 2080 train samples and 1088 test samples.

Diffusion-Sorption Equation The diffusion-sorption equation models a diffusion process that is retarded by a sorption process. The 1D PDE is given by:

$$\frac{\partial u}{\partial t} = D/R(u) \frac{\partial^2 u}{\partial x^2} \quad x \in (0, 1) \quad t \in (0, 500], \quad (13)$$

where $D = 0.0005$ is the effective diffusion coefficient, and $R(u) = 1 + 2.16u^{-0.126}$ is the retardation factor hindering the diffusion process. This equation is applicable to, for example, groundwater contaminant transport.

The boundary conditions are $u(t, 0) = 1$ and $u(t, 1) = D \frac{\partial u}{\partial x}(t, 1)$. The dataset, taken from [27], is discretized into 1024 spatial steps and 501 time steps. There are 9000 train trajectories and 1000 test trajectories, each based on different randomly generated initial conditions using $u(0, x) \sim U(0, 0.2)$ for $x \in (0, 1)$.

Hyperparams	DiffSorp	Burgers'	N.S.	Darcy
Epochs	200	200	100	400
QAT Epochs	100	50	50	100
Batch size	50	50	16	4
Learning rate	1e-3	1e-3	1e-3	5e-4
Weight decay	1e-6	1e-6	1e-6	1e-6
QAT learn. rate	1e-4	1e-4	1e-4	1e-4
Input steps	5	5	4	1
Output steps	5	5	1	1
Train steps	10	20	1	1
Test steps	10	20	1	1
Subsample t	2	5	1	1
Subsample x	32	16	2	8

Table 3: Dataset-related hyperparameters for all models per experiment. The steps refer to consecutive time steps for the time-dependent PDEs, while the Darcy PDE can optionally be interpreted as having inputs at $t = 0$ and outputs at $t = 1$.

A.2. Hyperparameter specifications

We summarize the hyperparameters used per dataset in Table 3.

In the second session of experiments we did not subsample the spatial grid, except for the Darcy dataset for which we subsampled every 2 grid points. The first three scaling levels applied in figure 4 correspond (from left to right) to 0.01, 0.02, 0.05 for the DiffSorp data, 0.02, 0.05, 0.1 for the Burgers data, 0.1, 0.2, 0.5 for the Navier-Stokes data and 0.05, 0.1, 0.2 for the Darcy data. The loss measure used in all datasets is the MSE as described in equation 5. However, for the Darcy dataset, we also normalize each element in the sum by dividing by the squared targets, and take the squared root of the resulting sum.

The UNet is taken from [9], but in order to make its size comparable to the other models we use 16 hidden channels for the Navier-Stokes dataset and 8 hidden channels for the other datasets. The FNO model is taken from [17], using 4 layers, a width of 128, 32 modes for the 2D datasets, and 16 modes for the 1D datasets. The Transformer is taken from [6]. It uses 6 encoder layers, 128 hidden channels and a Galerkin attention type for the 2D datasets, and 4 encoder layers, 32 hidden channels and Fourier attention type for the 1D datasets.

A.3. Inference Cost Calculation Details

We describe a few differences compared to the regular deepspeed library [22]. Most standard deep learning operations rely on big matrix multiplications and as such deepspeed outputs the number of MACs used in their corresponding modules. On the other hand, there are some operations that have no MACs, and deepspeed simply outputs

the number of FLOPs. However, because we care to differentiate addition and multiplication operations for our proxy measure of inference cost, we add some manual changes to deepspeed so that multiplications are properly accounted for.

- We assume bilinear interpolation to be three times the cost of linear interpolation, and for linear interpolation we assume 2 multiplications and 4 additions per output point.
- The FNO model uses Fast Fourier Transforms, which are not encountered for in deepspeed. To be able to take these into account in our proxy for model inference cost we assume a complexity of $N \lceil \log_2 N \rceil$ (additions and multiplications), which we divide by 2 when the real-valued FFT is used.
- In deepspeed no MACs are assigned to the einsum operator. Although it can in theory represent various different types of computations, in our code we only use it for basic matrix multiplications (in the FNO model). We thus change the deepspeed output accordingly.