

Adaptive Neural Network-Based Approximation to Accelerate Eulerian Fluid Simulation

Wenqian Dong
University of California, Merced
wdong5@ucmerced.edu

Zhen Xie
University of California, Merced
zxie10@ucmerced.edu

Jie Liu
University of California, Merced
jliu279@ucmerced.edu

Dong Li
University of California, Merced
dli35@ucmerced.edu

ABSTRACT

The Eulerian fluid simulation is an important HPC application. The neural network has been applied to accelerate it. The current methods that accelerate the fluid simulation with neural networks lack flexibility and generalization. In this paper, we tackle the above limitation and aim to enhance the applicability of neural networks in the Eulerian fluid simulation. We introduce Smart-fluidnet, a framework that automates model generation and application. Given an existing neural network as input, Smart-fluidnet generates multiple neural networks before the simulation to meet the execution time and simulation quality requirement. During the simulation, Smart-fluidnet dynamically switches the neural networks to make best efforts to reach the user’s requirement on simulation quality. Evaluating with 20,480 input problems, we show that Smart-fluidnet achieves 1.46x and 590x speedup comparing with a state-of-the-art neural network model and the original fluid simulation respectively on an NVIDIA Titan X Pascal GPU, while providing better simulation quality than the state-of-the-art model.

ACM Reference Format:

Wenqian Dong, Jie Liu, Zhen Xie, and Dong Li. 2019. Adaptive Neural Network-Based Approximation to Accelerate Eulerian Fluid Simulation. In *The International Conference for High Performance Computing, Networking, Storage, and Analysis (SC ’19)*, November 17–22, 2019, Denver, CO, USA. ACM, New York, NY, USA, 14 pages. <https://doi.org/10.1145/3295500.3356147>

1 INTRODUCTION

The fluid simulation aims to study the flow of fluid materials and has been widely applied to multiple disciplines such as chemical physics and material science [1–3]. However, the simulation of fluid dynamics usually requires prohibitively high computational resources [4, 5] and thus limits its application in the related fields.

The neural network-based machine learning model, as a tool to learn and model complicated (non)linear relationships between input and output datasets, has shown preliminary success in HPC problems (e.g., detecting neutrinos [6], developing Bose-Einstein Condensates state [7], and recognizing extreme weather events [8]). Using neural networks, scientists are able to augment existing simulations by improving accuracy and significantly reducing latency.

The neural network has been applied to accelerate fluid simulation as well [9–11]. By replacing some execution phases with neural networks, the most recent work reports 14.6× to 716× performance improvement [10].

However, the current methods to accelerate the fluid simulation using neural networks have fundamental limitations. First,

the current method to apply the neural work to the fluid simulation lacks flexibility. In particular, given the simulation code, the current method usually generates just one neural network model. At runtime, this model is used throughout the whole execution to approximate some target computation. This method ignores the fact that replacing the target computation at different execution phases of the fluid simulation can have different implications on the simulation quality. At some execution phase, using another neural network model may be able to generate higher simulation quality without losing performance. Hence, using multiple neural network models instead of one is a better strategy. However, with just one single neural network model, the current method lacks such flexibility.

Second, the current method to apply the neural network to the fluid simulation lacks a generalization ability. In particular, given a fluid simulation code with the neural network applied, there is no guarantee that the application can run the simulation to completion and meet the requirement on simulation quality for every input problem. Using a single neural network model for all input problems either leaves the performance opportunities on the table (discussed in the last paragraph), have large simulation quality violations, or both.

Third, there is no *systematic* approach to construct and apply neural network models to the fluid simulation. How to construct neural networks to meet the user requirement on performance (execution time) and simulation quality is challenging. Currently, domain scientists build neural networks intuitively. There is no systematic approach to help them build and apply neural networks for HPC applications. The recent work on Auto-Keras [12] and AutoML [13] aims to automatically generate a neural network model with high accuracy. However, they lack concerns on high performance (execution time), and focus on image processing or natural language processing. Hence, they are not directly usable by HPC.

In this paper, we tackle the above limitations and aim to enhance the applicability of neural networks in HPC applications (particularly the Eulerian fluid simulation). Given an existing neural network model as input, our system uses a systematic approach to construct multiple neural network models and dynamically switches them at runtime during the execution of the fluid simulation to meet the user requirements on performance and simulation quality.

In order to tackle the above limitations, we must address three challenges. First, we must automatically generate multiple neural network models to enable high flexibility and better generality when applying neural networks. Given the user requirements on

performance and simulation quality, we aim to generate multiple neural network models, each of which has different topologies and different implications on performance and simulation quality. We should not expect the domain scientists to manually construct models.

Second, how to select neural network models at runtime to enable the best performance without violating the quality requirement is a challenge. We must have a method to predict the impact of applying a neural network model at a certain execution phase on the final simulation quality.

Third, a neural network model can approximate the fluid simulation with high accuracy for some input problems but not for all. How to construct neural network models to provide a high-quality approximation for a large number of input problems and ensure overall performance benefit is another challenge.

In this paper, we focus on a Eulerian fluid dynamic simulation code (mantaflow [14]) and introduce a framework (named “Smart-fluidnet”) to address the above three challenges. Smart-fluidnet has three major components. (1) It includes a model construction plugin for Auto-Keras to extend its functionality to enable automatic construction of multiple neural network models. (2) Smart-fluidnet also includes a multilayer perceptrons model (MLP) to guide the model selection process to meet the requirement on performance and simulation quality. The neural network models selected by MLP is able to cover more input problems to ensure overall performance benefit. (3) Smart-fluidnet includes a runtime component integrated into mantaflow and dynamically switches the neural network models to make best efforts to meet the user requirement on simulation quality. The runtime component is based on a metric and a lightweight runtime algorithm that can predict the final simulation quality in the middle of the fluid simulation.

We summarize the major contributions of this paper as follows:

- A systematic approach and a framework (Smart-fluidnet) to accelerate the Eulerian fluid simulation; Our evaluation shows that using 20,480 input problems for the simulation, Smart-fluidnet achieves 46% performance improvement over the Tompson’s model [10] (a state-of-the-art neural network model) on average and 590× speedup over the original fluid simulation on average, while providing better simulation quality than the Tompson’s model.
- A new methodology that constructs, selects, and applies multiple neural network models (instead of one) to address the fundamental limitation of model flexibility and generalization in the existing neural network-based approximation. We demonstrate great potential of using this methodology to meet user requirements on the simulation quality and execution time.

2 BACKGROUND

In this section, we provide background on the Eulerian fluid simulation and neural network-based approximation. In the rest of the paper, the term *performance* means execution time, not prediction accuracy in the machine learning field. We also use the terms *approximation model* and *neural network model* interchangeably.

2.1 Eulerian Fluid Simulation

The Eulerian fluid simulation, in essence, solves the Navier-Stokes equations. The Navier-Stokes equations describe the fluid movement under a continuous velocity field \vec{u} and a pressure field p . When the fluid has zero viscosity, one uses the incompressible Navier-Stokes equations, which can be expressed as the Euler equations as follows:

$$\frac{\partial \vec{u}}{\partial t} = -\vec{u} \cdot \nabla \vec{u} - \frac{1}{\rho} \nabla p + \vec{\tau}, \quad (1)$$

$$\nabla \cdot \vec{u} = 0 \quad (2)$$

Equation 1 is a vector equation called “momentum equation”. This equation can make the velocity field stay divergence-free. Equation 2 is the incompressibility condition, which enforces fluid volume to remain constant throughout the simulation. In the above two equations, t is time, ρ represents fluid density, and $\vec{\tau}$ represents gravity.

Mantaflow performs fluid simulation by iteratively solving Equations 1 and 2. In this paper, we use MAC (marker-and-cell) grids [15] to discretize fluid flows, and use the finite difference (FD) method to calculate partial derivative on each grid [16, 17].

For each velocity component that borders a grid cell, the FD method iteratively applies updates on velocity and pressure. In a grid cell, the pressure is sampled at the grid cell center and the velocity is sampled at the centers of the vertical faces of the grid cell. The above method is common and can simplify the handling of solid-cell boundaries conditions.

To solve Equations 1 and 2, mantaflow uses a standard operator splitting method [18–20] to split up Equation 1 (Equation 2 is used as a constraint) into three parts. The three parts are advection, adding external force, and pressure projection. Algorithm 1 depicts the implementation of the Eulerian simulation in mantaflow, which includes the above three parts.

The Eulerian simulation in mantaflow includes N time steps. The first part (Lines 4-5) of each time step is to solve the momentum equation to get an auxiliary velocity field \vec{u}^B . \vec{u}^B is a velocity approximation which is not divergence-free, and the pressure-gradient term (∇p) used during the solving process is computed in the previous time step. At Line 7, the divergence-free pressure \vec{p}_η is computed by solving a Poisson’s equation which includes the divergence of \vec{u}^B and a scaled gradient of the pressure. At Line 18, a divergence-free velocity field, $\vec{u}^{\eta+1}$ is calculated, by subtracting off the pressure gradient from the approximate velocity field \vec{u}^B . In the above process, solving the Poisson’s equation is the most crucial and time-consuming step to preserve the divergence-free constraint on the velocity and maintain simulation accuracy.

The process of solving the Poisson’s equation in mantaflow (Line 7 in Algorithm 1) is based on the Preconditioned Conjugate Gradient (PCG) method, which involves large computation that iteratively converges to meet a convergence criteria (Lines 8-17). Mantaflow uses a multi-grid approach [21] as a preprocessing step of the PCG method. The pre-conditioner (Line 10) applied in mantaflow is the Modified Incomplete Cholesky L0 preconditioner, called “MICCG(0)”. In this paper and the existing work [10], neural networks are used to approximate this PCG method.

Algorithm 1 Velocity Update in the Euler Equation

Require: Simulation time step N ;

- 1: Start with an initial divergence-free velocity field \vec{u}^0
- 2: Determine a good time step Δt to go from time t_n to time t_{n+1} .
- 3: **for** $n \leftarrow 1$ to N **do**
- 4: **Advection.** Set $\vec{u}^A = \text{advect}(\vec{u}^n, \Delta t, q)$;
- 5: **Add body force.** $\vec{u}^B = \vec{u}^A + \Delta t \vec{f}$;
- 6: **Pressure projection.** set $\vec{u}^{n+1} = \text{Project}(\Delta t, \vec{u}^B)$;
- 7: 1) Solve the Poisson eq. $\nabla \cdot \vec{p}_n = \frac{1}{\Delta t} \nabla \cdot \vec{u}^B$
- 8: //Use a PCG solver to to update \vec{p}_n .
- 9: Set initial guess $\vec{p}_n=0$ and residual vector $r=d$ (if $r=0$, then return \vec{p}_n)
- 10: Set search direction $\vec{s} = \text{Appl Preconditioner}(r)$;
- 11: **while** residual doesn't reach the convergence criteria
- 12: **do**
- 13: Set $\vec{s} = \frac{r^T r}{s^T A s}$;
- 14: Calculate the residual $r = r - A \vec{p}_n$;
- 15: Update the solution $\vec{p}_n = \vec{p}_n + \vec{s}$;
- 16: Update the conjugated direction $\vec{s} = r + \vec{s}$;
- 17: **end while**
- 18: 2) Apply velocity update: $\vec{u}^{n+1} = \vec{u}^B - \Delta t \nabla \vec{p}_n$;
- 19: **end for**
- 20: **return** 0

In this paper, we simulate a 2D smoke plume [22, 23]. The simulation output in mantaflow is a smoke *dense matrix* of a rendered smoke frame. The smoke dense matrix represents density blurring of the plume, which reflects fluid movement. After using neural networks to approximate the computation in the fluid simulation, the output dense matrix can be different from that in the original simulation (using the mantaflow's PCG-based solver), which means we have quality loss. The simulation quality loss (Q_{loss}) is formally defined by:

$$Q_{loss} = \frac{1}{N \times M} \sum_{i=1}^N \sum_{j=1}^M |ij - ij^*|, \quad (3)$$

where ij refers to the smoke density matrix generated in the original simulation, and ij^* represents the smoke density matrix generated after applying neural network-based approximation. ij and ij^* are matrix elements with the coordinate (i, j) . In essence, Equation 3 calculates the average relative error of all matrix elements. After applying the neural network-based approximation, we want to avoid quality loss (i.e., we do not want to lose simulation accuracy).

2.2 Neural Network-Based Approximation

The neural network is a general-purpose method that can be used to learn and model complicated linear and non-linear relationships between input and output datasets. Hence, the neural network has been used to approximate some conventional algorithms in an application to improve performance [9–11]. The neural networks are expected to generate similar outputs as the conventional algorithms when fed with the same inputs as for the conventional algorithms.

A neural network can be represented as a directed acyclic graph where nodes of the graph are connected neurons. Embedded in the graph, there are a number of parameters (or “weights”). Those

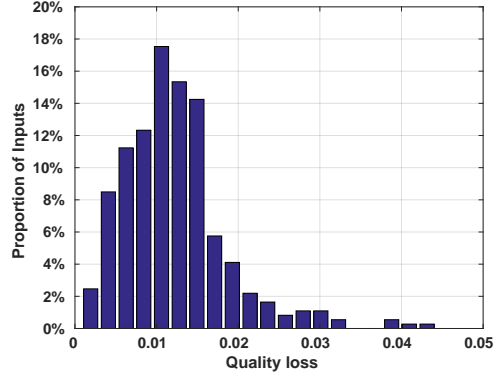


Figure 1: Distribution of quality loss for the Tompson’s model with different input problems [10].

neurons and weights are organized as layers. The process of obtaining the values of those weights are called *training*. Once the neural network is trained offline using training datasets, it can be used online within the fluid dynamic simulation to improve performance. During training, an *objective function* is used to calculate the model accuracy loss, so that the weights can be adjusted to reach better model accuracy.

In this paper, we use Convolutional Neural Networks (CNN) to approximate computation (i.e., using the PCG solver to solve the Poisson’s equation) in the Eulerian fluid simulation. This computation is the most time-consuming part of the Eulerian fluid simulation. Our profiling results reveal that this computation takes 70-80% of total simulation time.

Recent work [10] introduces an unsupervised learning framework to generate a CNN model with five stages of convolution and Rectified Linear Unit (ReLU) layers to approximate computation (the PCG solver) in the Eulerian fluid simulation. The inputs of the CNN are the divergence of the velocity field, denoted as $\nabla \cdot \vec{u}_t^*$, and the geometry field, denoted as t_{-1} . The output of the CNN is the pressure field, denoted as \hat{p}_t . The mapping f_{con} from input to output by this CNN model can be represented as follows:

$$\hat{p}_t = f_{con}(\nabla \cdot \vec{u}_t^*, t_{-1}; W), \quad (4)$$

where W is the CNN model parameters. The predicted \hat{p}_t is used to update velocity \vec{u}_t in Equation 1. In our study, this CNN (named as the Tompson’s model) is used as input in Smart-fluidnet to generate new neural networks for online fluid simulation.

The objective function of the Tompson’s model, i.e., *Di Norm*, is the sum of weighted L-2 norm of the divergence of the predicted velocity \vec{u}_t over all fluid cells (mesh volumes) in the rendered smoke frame. *Di Norm* is defined as follows:

$$Di\ Norm = \sum_i w_i \{ \nabla \cdot \vec{u}_t \}_i^2, \quad (5)$$

where w_i is a weighting term for each fluid cell to emphasize the divergence of grids on geometry boundaries, i.e., $w_i = \max(1, k-d_i)$. d_i is the distance field. It takes the value 0 for solid cells or the minimum Euclidean distance to the nearest solid cell for fluid-cells.

Table 1: Execution time and simulation quality loss of three models for solving the Poisson’s equation.

Method	Execution Time (ms)	Avg. Quality Loss
PCG	2.34×10^8	--
Tompson [10]	7.19×10^4	1.3×10^{-2}
Yang [11]	3.20×10^4	4.9×10^{-2}

2.3 Motivation of Our Work

The existing work to solve the Poisson’s equation includes the PCG solver in mantaflow and two neural network models (Tompson [10] and Yang [11]). We study the implications of the three methods on simulation quality and execution time. To study the impact of each model, we evaluate 20,480 different input problems of the fluid simulation and report the average simulation quality loss and execution time. The simulation quality loss Q_{loss} is calculated by comparing the simulation output and the real physical measurements on fluid flow. Table 1 and Figure 1 show the results.

Table 1 shows that PCG achieves the highest simulation quality as an exact solution but the worst performance (the longest execution time). On the other hand, the Yang’s model performs $10^4 \times$ faster than PCG but causes about $10^2 \times$ loss in the simulation quality. There is a clear trade-off between simulation quality and performance.

Figure 1 shows the distribution of quality loss for various input problems when we use the Tompson’s model. The figure reveals that given various input problems, the model generates different simulation quality. For most input problems, the simulation quality loss is between 0.01 and 0.02. Given a user-defined quality requirement (e.g., the quality loss should be less than 0.01), the simulation may not meet the quality requirement for most input problems (e.g., around 65.42% input problems can not meet user requirement when the requirement is 0.01). We have the same observation for the Yang’s model.

The above results reveal that it is imperative to use multiple models to explore the trade-off between performance and simulation quality and maximize the possibility of reaching the user requirement on the simulation quality for various input problems.

3 OVERVIEW

Figure 2 gives the workflow of Smart-fluidnet. The workflow includes offline and online phases. During the offline phase, Smart-fluidnet takes an existing neural network model as an input and constructs a set of neural network models by model transformation. We introduce four operations (deleting, narrowing, pooling and dropout) to transform the input neural model into multiple neural network models. Then we choose neural network models that are promising for high performance improvement and high quality based on the Pareto optimality analysis.

After model construction, Smart-fluidnet further chooses models based on the user requirement about simulation quality. Given various input problems of the fluid simulation, we introduce an MLP-based model to predict the probability of each model to reach the user requirement on the simulation quality. Considering the

possible cost of restarting the simulation when the simulation quality does not meet the user requirement, we choose those models that have a sufficiently high probability to benefit the performance improvement. After the above offline phase, we have a handful of neural network models ready for online approximation. At runtime, given an input problem of the Eulerian fluid simulation, Smart-fluidnet dynamically switches neural network models with the most promising neural network to meet the user requirement on the simulation quality. The model switching is based on a metric and a linear regression model to predict whether the current model can reach the requirement on the simulation quality at the end of the simulation.

The Smart-fluidnet framework consists of three main modules: approximation model construction (Section 4), offline output-quality control (Section 5), and quality-aware runtime design (Section 6). Their relationships are depicted in Figure 2. We explain them in detail as follows.

4 APPROXIMATE MODEL CONSTRUCTION

Given an input neural network, we transform it to construct multiple neural network models with different network architectures.

The new neural network models can be more accurate or more efficient (i.e., using less execution time) than the input neural network; Having such a mixture of different models provides flexible computation approximation during the online fluid simulation.

To generate more accurate neural network models, the user can use an existing framework such as Auto-Keras to generate and train models. Auto-Keras uses Bayesian optimization to generate a single model with the best accuracy. We change Auto-keras to generate and train five models with better accuracy. We generate five models, because generating more than five highly accurate models often causes more than five models to be selected after applying MLP (Section 5), which causes large runtime overhead when making the decision on switching models; While generating less than five models will lead to insufficient candidates after the model selection (Section 5).

Besides the above, we also aim to generate less accurate but faster models. We introduce new transformation operations into Auto-Keras to simplify the neural network architecture, which will shorten the execution time. We describe our transformation operations as follows.

Operation 1: deleting a layer of the neural network. This operation is denoted as *shallow*(G, L), where G is the neural network graph of the input neural network and L is the layer to be deleted. This operation not only shortens the depth of neural network but also reduces memory consumption, thus makes the fluid simulation time shorter.

Operation 2: narrowing a layer of the neural network. This operation reduces neurons in an intermediate layer. This operation is denoted as *narrow*(G, L, r), where r is the number of neurons to be reduced at the layer L , and L can be either a fully-connected layer or a convolutional layer.

Operation 3: pooling. This operation, denoted as *poolin* (G, L, m), downsamples a layer L using a pooling matrix m . We can apply two pooling strategies, i.e., *maxpoolin* and *averagepoolin*. A special

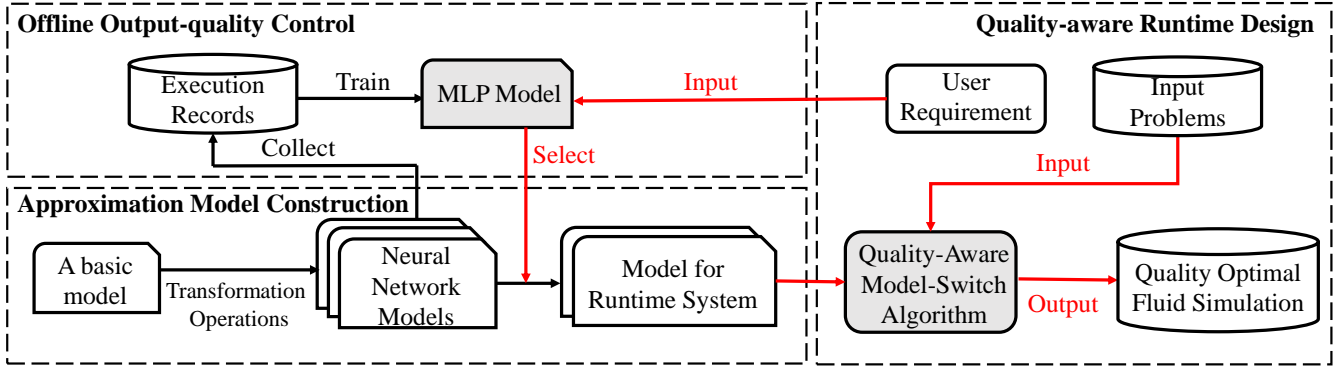


Figure 2: Workflow of the proposed Smart-fluidnet.

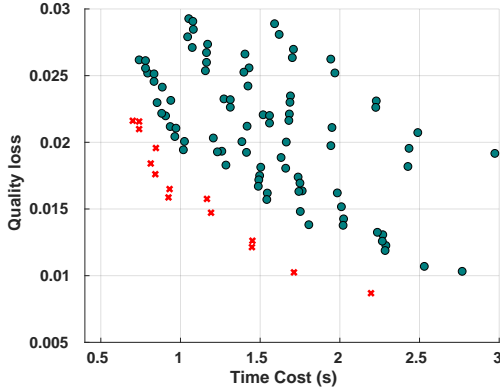


Figure 3: Scatter plot of quality loss and time cost of different neural network models.

case of m is a 2×2 matrix which can discard 75% neurons in the intermediate layers.

Operation 4: dropout. This operation, denoted as $dropout(G, L, p)$, drops neurons at a layer L with a given probability p . It offers a more flexible way to reduce the number of neurons, compared with the second operation by controlling the value of p . This operation is useful to increase the generalization capability of the model.

Based on the above four operations, we transform the input neural network into new ones. Given an input neural network, we first apply a $shallow(G, L)$ operation on each of the intermediate layers. We do not apply the operation more than twice in the input model, and hence do not delete more than one intermediate layers. After applying $shallow(G, L)$, we generate five new models.

Second, we apply $narrow(G, L, r)$ operation on the five new models. A big value of r means a large number of neurons will be reduced. Based on our experimental results, the simulation quality loss can be large (more than 20%) for most of the input problems, if $r > \frac{|L|}{2}$, where $|L|$ is the total number of neurons. Therefore, we empirically use $r = \frac{|L|}{10}$. For each new model, we randomly choose r neurons to apply the $narrow(G, L, r)$ operation; Furthermore, for each new model, we apply $narrow(G, L, r)$ ten times, each of which generates a new model. In our case, in total, we have 55 new models

(five new models after applying $shallow(G, L)$ and 50 more after applying $narrow(G, L, r)$).

Third, for the 55 new models, we apply $poolin(G, L, m)$ operations. In particular, for each new model, we randomly replace any of two neighbor-neurons with a new neuron using max pooling. The total number of neurons to be replaced is constrained to be 10% of the total neurons. After applying $poolin(G, L, m)$, we have 55 more new models (110 models in total).

At last, to enrich our neuron network models, we randomly select 18 out of the 110 models to apply the $dropout(G, L, p)$ operation. In particular, in each of the 18 models, we randomly drop out neurons. The total number of neurons to drop is limited to 10% of the total neurons. After applying $dropout(G, L, p)$, we have 18 more new models. In total, we have 128 models.

We apply the four operations in the above order, because the operation that tends to reduce more neurons than other operations will be performed earlier. This method allows us to efficiently generate new models. Using a different order can take longer time to generate models or be prone to generate less accurate models.

After the above model generation and in combination with the accurate models generated by Auto-Keras (five models), we have 133 models in total. Afterward, we use *Pareto optimality* to reduce the number of models for online approximation. We select models that have the lowest time cost, the lowest quality loss, or both.

To explain the idea of Pareto optimality, we use Figure 3 to illustrate it. Figure 3 shows the result of our model selection approach. In this figure, each point represents a model; The red points are those selected model for further analysis (Section 5); The green points are discarded models. In Figure 3, the quality loss and execution time for each model are collected during the model construction. This is a common method to collect the model information in AutoML work [24, 25]. Section 7 gives details on the hardware platform and input datasets to collect the results in Figure 3. We can observe that those models located in the leftmost part of Figure 3 either have the lowest time cost, the lowest quality loss, or both. Those models (14 models) are selected based on the Pareto optimality method (we name them “model candidates” in the later discussion).

Sensitivity Study. The above process of constructing neural networks involves a couple of parameters. We summarize them as follows and change those parameters to study their impact on the

simulation quality. In our study, we use 100 input problems. Using the simulation quality of PCG as the baseline, we calculate the average quality loss of all input problems when using the Tompson’s model. This average quality loss is used as the user requirement for quality loss in our sensitivity study.

(1) The number of layers to prune (Operation 1). Our current method prunes one layer at most. For sensitivity study, We prune more than one layer, but find that it leads to a large quality violation (20% quality loss on average), which is not good.

(2) The percentage of neurons to apply pooling (Operation 3). Our current method applies pooling to 10% of total neurons. We set the percentage of neurons to apply pooling as 5%, 20% and 30%, and evaluate the impact of this parameter on the simulation quality. Our experiments show that using 20% and 30%, the fluid simulation has a large quality violation (35% and 50% quality loss on average, respectively), which is not good. However, using 5%, the fluid simulation has the similar quality loss on average as using 10%. In addition, since using 10% can lead to better performance than 5%, we choose 10% in our study.

(3) The dropout rate (Operation 4). Our current method drops out 10% of total neurons. We also try 5% and 15% as the dropout rate. Our experiments show that 5% and 10% outperform 15% in terms of the simulation quality loss. In particular, the quality losses with 5% and 10% as the dropout rate are 0.156 and 0.164 respectively, while the quality loss with 15% as the dropout rate is higher (0.239). Since 5% and 10% have the similar quality loss and using 10% leads to less execution time, we adopt 10% as our dropout rate.

(4) The number of neural network models to apply the dropout operation. Our current method chooses 18 models. Our study reveals that choosing 15-20 models are enough for the Pareto optimality and MLP (see Section 5) to generate 2-5 models. Using less than 15 models to apply the dropout operation, however, we could generate no qualified model after applying MLP. Using more than 20 models, we could have more than 5 qualified models after applying MLP. Having more than 5 models means that the runtime system may suffer from large runtime overhead for selecting a model to use. As a result, we choose 18 (in between 15 and 20) as our parameter.

5 OFFLINE OUTPUT-QUALITY CONTROL

In many fluid simulation cases, users can have specific requirement on the simulation quality and execution time. We use the notation, $U(q, t)$, to represent the user requirement, where q and t are the user requirement on the quality loss and execution time respectively. The final quality loss and execution time of the fluid simulation should be less than q and t , respectively. The quality control should be aware of the *success rate* of neural network models, namely the ratio of those input problems with which the fluid simulation can reach the simulation quality and time requirement to the total number of input problems.

The awareness of the success rate can be developed from statistical knowledge by executing neural network models on various input problems. In this section, we design a non-linear MLP model to develop such awareness. In the following subsections, we first introduce the construction of training samples for training the MLP model and then give details on how to construct and apply the MLP.

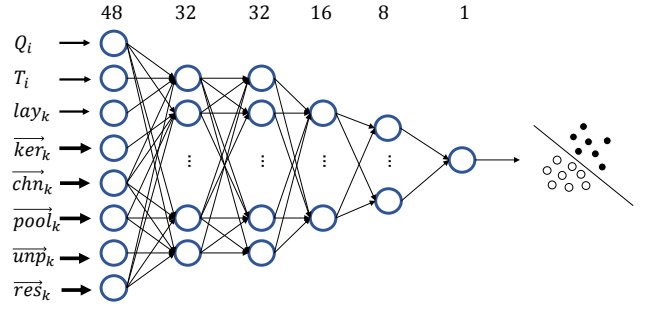


Figure 4: The network architecture of our MLP model.

5.1 Construction of Training Samples

We collect execution records (i.e., simulation quality and execution time) for the 14 neural network models after the construction of those models (Section 4). Based on the execution records, we generate training samples to train the MLP model.

Collection of Execution Records. For each of the 14 neural network models, we get N execution records by running N input problems ($N = 20, 480$ input problems in this paper). Each of the N execution records includes the simulation quality q_n^k and execution time t_n^k , where $n \in [1, N]$ and k refers to a neural network NN_k from the 14 neural network models. Each execution record is represented as ER_n^k . Such execution record is collected during the model construction (Section 4). Note that our model construction process, similar to other AutoML work [24, 25], includes not only changing model architecture but also training models. Hence we can collect execution records during the model construction.

Sample Generation. Given N execution records, we can generate samples to train the MLP. Each sample is represented with a feature vector using Equation 6.

$$F = q, t, l_k, \overrightarrow{ker}_k, \overrightarrow{chn}_k, \overrightarrow{pool}_k, \overrightarrow{unp}_k, \overrightarrow{res}_k \quad (6)$$

In Equation 6, the feature vector F includes quality and execution time requirements (q and t), and architecture information for the neural network NN_k (i.e., l_k , \overrightarrow{ker}_k , \overrightarrow{chn}_k , \overrightarrow{pool}_k , \overrightarrow{unp}_k , and \overrightarrow{res}_k , representing the number of layers, kernel sizes, channel number, pooling size, unpooling size, and residual connection of each layer respectively). Each of the last five architecture information is a vector composed of nine components that indicate properties of each layer of the neural network NN_k . Therefore, each input feature vector has $3 + 5 * 9 = 48$ components in total.

Given a neural network NN_k , we generate a sample by randomly picking up a user requirement (q and t), and then use Equation 6 to build the feature vector of the sample. Based on the N execution records and the user requirement, we calculate that how many of N execution records meet the user requirement. The ratio of those execution records to N , denoted as $r_{k; q, t}$, is the label of the sample. By choosing different combinations of q and t , we can generate as many samples as possible.

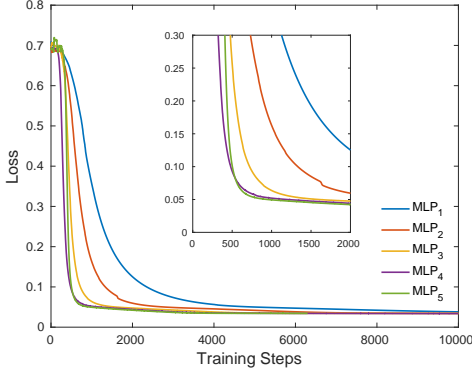


Figure 5: Training losses of five MLPs.

5.2 MLP Model Construction and Loss Function

Given a neural network NN_k and user requirement q and t , we can build a feature vector F using Equation 6. Our MLP model (see Equation 7) takes such a feature vector as input and generates an output $\hat{r}_{k;q;t}$, which is a floating point number indicating the probability that NN_k meets the user requirement for any input problem.

$$\hat{r}_{k;q;t} = \hat{r}_{MLP} F_{k;q;t} \quad (7)$$

Figure 4 shows the network topology of our MLP. It includes six hidden layers and a 48-neuron input layer. The numbers of neurons in the six hidden layers are 32, 32, 16, 16, 8 and 8 respectively. All the neurons in the hidden layers use ReLU as activation to increase the non-linearity of the model. The neurons in the last hidden layer use a sigmoid function as activation. Using the samples constructed in Section 5.1 to train the MLP, we aim to minimize the loss between the model output $\hat{r}_{k;q;t}$ and the ground truth label $r_{k;q;t}$.

Alternative MLP Topologies. Besides the above MLP topology, we try four alternative MLP topologies, in order to find one for best accuracy. Among the four MLPs, two of them are deeper than our current MLP, while the other two are shallower. When building an alternative MLP, we follow the rule that in the topology of a neural network, a deeper layer generally has less number of neurons than shallower ones (i.e., the number of neurons gradually decreases across layers). Many discriminative neural networks, such as Alex-Net[26], VGG-Net[27], are constructed, following this rule. The architectures of the four models plus our current MLP model are briefly described as follows:

- MLP_1 has 4 layers with 48, 32, 16 and 1 neurons;
- MLP_2 has 5 layers with 48, 32, 16, 8 and 1 neurons;
- MLP_3 (our current MLP model) has 6 layers with 48, 32, 32, 16, 8 and 1 neurons;
- MLP_4 has 7 layers with 48, 64, 32, 32, 16, 8 and 1 neurons;
- MLP_5 has 8 layers with 48, 64, 64, 32, 32, 16, 8 and 1 neurons.

Figure 5 presents the training loss curves of the above five MLPs (MLP_3 is our current MLP model). We find that the convergence speed of MLP_3 is faster than those of MLP_1 and MLP_2 , and offers lower training loss (i.e., higher prediction accuracy). Compared with MLP_3 model, MLP_4 and MLP_5 do not have significant advantages in terms of convergence speed and loss, although they have deeper

topologies. Hence, MLP_3 exhibits a balanced trade-off between prediction accuracy and model size, and is thus chosen as our MLP model in this paper.

5.3 Usage of MLP

Given a user-specified simulation quality (q) and time cost (t), we use MLP to calculate $\hat{r}_{k;q;t}$ for a given neural network model NN_k . A larger value of $\hat{r}_{k;q;t}$ represents a higher success rate. In other words, it is highly possible that NN_k can meet the simulation quality and time requirement on an input problem.

Given the user-specified simulation quality (q) and time requirement (t), the neural network NN_k and MLP prediction result ($\hat{r}_{k;q;t}$), we use the following method to decide if NN_k should be selected for the runtime system for the fluid simulation. In particular, considering the probability that the user requirement on the simulation quality is violated and the user has to re-run the simulation without using any neural network, the simulation time is calculated based on Equation 8.

$$T_{total} = \hat{r}_{k;q;t} \times T_{M_k} + (1 - \hat{r}_{k;q;t}) \times T', \quad (8)$$

where T' is the execution time without using any neural network and T_{NN_k} is the execution time using the neural network NN_k . We compare T_{total} with the user requirement on the execution time t . Only those neural networks that have T_{total} less than t is selected.

The above selection method considers the impact of violating the simulation quality requirement on the simulation time, and ensure that if NN_k is repeatedly employed for many input problems, there is a performance benefit.

6 QUALITY-AWARE RUNTIME DESIGN

After applying MLP, multiple neural networks are selected. We use a runtime technique to schedule those neural networks to optimize performance and meet the simulation quality requirement. In order to determine which neural network should be used at runtime, we need to evaluate the model being used in terms of the final quality loss, and switch to a suitable one if necessary. However, without running the simulation to completion, we cannot know the final quality loss. Thus we construct a metric called *CumDi Norm* (defined in Equation 9), which is used to set up a bridge between *Di Norm* (see Equation 5) measurable at runtime and the final simulation quality loss Q_{loss} (Section 6.1). Based on *CumDi Norm* and the predicted final quality loss, we introduce a quality-aware model-switch algorithm (Section 6.2) to select the best neural network to accelerate the fluid simulation.

6.1 Prediction of Simulation Quality Loss

The objective function *Di Norm* provides a goal that our neural network aims to achieve. Using *Di Norm*, we can know how the neural work performs in terms of prediction accuracy. However, there is a missing link between the prediction accuracy of the neural network and the final simulation quality loss Q_{loss} .

***CumDi Norm*: A Metric for Runtime Quality Control.** To explore the relationship between *Di Norm* and final simulation quality Q_{loss} , we calculate *CumDi Norm* (i.e., the accumulation of *Di Norm*) and Q_{loss} at each simulation time step (denoted as Q_{loss}^t).

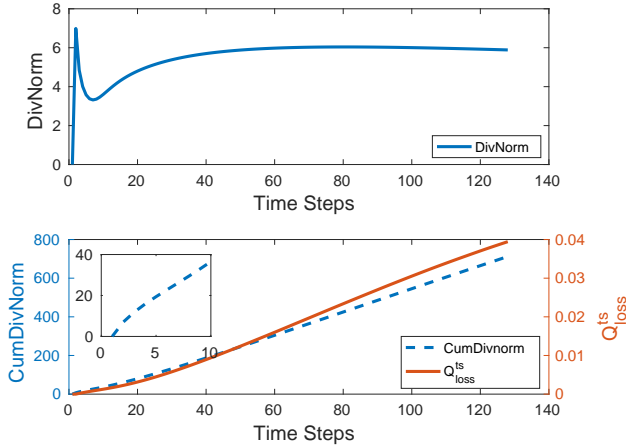


Figure 6: Relationship between $CumDi\ Norm$ and Q_{loss}^{ts} .

The accumulation of $Di\ Norm$ over n time steps is defined in Equation 9.

$$CumDi\ Norm = \sum_{i=1}^n Di\ Norm_i. \quad (9)$$

In order to observe how these variables are correlated, Figure 6 depicts how $Di\ Norm$, $CumDi\ Norm$, and Q_{loss} vary across all time steps of the fluid simulation using an input problem with the grid size 1028×1028 . We have the following observations. These observations are valid for other input problems as well.

- **Observation 1:** $Di\ Norm$ dramatically increases at the first few time steps and then gradually converges to a stable value;
- **Observation 2:** Q_{loss}^{ts} and $CumDi\ Norm$ have similar increasing tendency (except the first few time steps).

The above observations indicate that $CumDi\ Norm$ and Q_{loss}^{ts} calculated at each time step are correlated. In order to quantify the relationship between $CumDi\ Norm$ and Q_{loss}^{ts} at each time step, we use the Pearson's product moment correlation coefficient (rp) [28] and the Spearman's rank correlation coefficient (rs) [29] to statistically reveal the correlation between the two variables.

The two coefficients are defined as follows. Given two input vectors X and Y (the vector length is n), the calculation of rp and rs to quantify the correlation between X and Y is formulated as follows:

$$rp = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}, \quad (10)$$

$$rs = 1 - \frac{6 \sum_{i=1}^n d_i^2}{n(n^2 - 1)}. \quad (11)$$

where x_i and y_i are the values of X and Y for the i -th component, and d counts the pairwise disagreement (i.e., x_i is not equal to y_i) between the two vectors (see [30] for a more detailed description). In general, the coefficients that belong to $[0.10-0.29]$ represent weak association, $[0.29-0.49]$ represent medium association, and above 0.49 represent strong association [31].

In our study, we use 20,480 input problems, each of which will have 128 simulation time steps. The $CumDi\ Norm$ and Q_{loss}^{ts} of

Algorithm 2 The quality-aware model-switch runtime algorithm

Require: The user requirement $U(q; t)$.

- 1: Choose a neural network model M_k with the highest success rate according to MLP.
 - 2: **while** t does not reach the final time step **do**
 - 3: Send M_k to predict the final simulation quality.
 - 4: **Prediction of quality loss:**
 - 5: 1) Build a linear regression model with $Di\ Norm$ values measured in the last five time steps;
 - 6: 2) Predict $CumDi\ Norm^{final}$ by the linear regression model;
 - 7: 3) Predict Q'_{loss} of the current neural network model by the KNN algorithm;
 - 8: **Model Switch:**
 - 9: **if** Q'_{loss} is close to q **then**
 - 10: Continue using the current neural network model for L steps;
 - 11: **else if** Q'_{loss} less than q **then**
 - 12: Switch to a faster (less accurate) neural network model;
 - 13: **else if** Q'_{loss} is larger than q **then**
 - 14: Switch to a slower (more accurate) neural network model;
 - 15: **else if** Cannot find any model **then**
 - 16: Restart by the PCG method;
 - 17: **end if**
 - 18: $t \leftarrow t + L$; // L is the check interval.
 - 19: **end while**
 - 20: **return** 0
-

each time step will be calculated to build the two input vectors. Using Equations 10 and 11, we have $rp = 0.61$ and $rs = 0.79$, which indicates a strong correlation.

Based on the above discussion, it is possible to use $CumDi\ Norm$ to predict Q_{loss}^{ts} in the final time step (i.e., the final quality loss Q_{loss}). Note that we cannot calculate Q_{loss}^{ts} at runtime, because it involves PCG, which is too expensive. In the following discussion, we first discuss how to predict $CumDi\ Norm$ in the final time step ($CumDi\ Norm^{final}$), and then we discuss how to predict Q_{loss} based on the predicted $CumDi\ Norm^{final}$.

Predicting $CumDivNorm^{final}$. We introduce a lightweight approach to predict $CumDi\ Norm^{final}$. Our prediction approach is based on Figure 6. In this figure, $CumDi\ Norm$ quickly grows at first, and then the growth rate remains stable. This trend is general across all 20,480 input problems we test.

The stable growth rate of $CumDi\ Norm$ makes it possible to predict $CumDi\ Norm^{final}$ in the middle of the fluid simulation. In particular, we use five time steps to build a linear regression model ($f_k(x) = ax + b$), by the least square method, where x is the time step and $f_k(x)$ is the predicted $CumDi\ Norm$. Note that the data used to build the linear regression model must be collected *after* the growth rate becomes stable. Hence, we skip the first five time steps and build the regression model after each five steps. Also, in each five time steps (a check interval) to build the model, we skip the first two to make sure the trend is stable and only use the remaining three to build the model. The above model-building process happens every five time steps, and the model is used to predict and check Q_{loss} (see the following discussion). Hence, we fix the check interval to be five time steps in the rest of the paper. In Section 7.4, we study the impact of the check interval on the simulation quality.

Predicting Q_{loss} based on $CumDivNorm^{nal}$. We use a method based on the k-nearest neighbor (KNN) algorithm to predict Q_{loss} . Our method includes offline and online phases. During the offline phase, we test the neural network models selected by MLP with 128 small input problems. For each test, we collect a pair of data ($CumDi Norm^{final}$, Q_{loss}) and put them into a historical database. The offline phase is fast because we use small input problems. During the online phase, at a time interval, to predict Q_{loss} , we check $CumDi Norm^{final}$ in the database and find k pairs whose $CumDi Norm^{final}$ are the closest to the predicted $CumDi Norm^{final}$ in the current time interval. We use the average of Q_{loss} in the k pairs as the predicted Q_{loss} in the current time interval. In our evaluation, we use different values for k , but find that $k \in [4, 6]$ is usually sufficient to give an accurate prediction, hence we choose $k = 4$ to reduce runtime overhead.

For example, assuming that the predicted $CumDi Norm^{final}$ in a time interval is 108. To predict Q_{loss} for this time interval, we select four pairs from the database, i.e., (101, 0.09), (112, 0.11), (105, 0.10), and (109, 0.11), whose $CumDi Norm^{final}$ is closest to the given $CumDi Norm$ (108). Then the predicted Q_{loss} for this time interval is $0.1025 = (0.09 + 0.11 + 0.10 + 0.11)/4$. We organize all data pairs as a binary search tree, such that finding the four pairs is cheap.

6.2 Quality-Aware Model-Switch Algorithm

With the ability to predict the simulation quality loss, we introduce a quality-aware model-switch algorithm. Algorithm 2 depicts this runtime algorithm. After applying MLP, we have several promising neural network models and their probabilities to meet the user-specified requirement. We also know the execution time (i.e., the inference time) of each network model. During the simulation, the neural network with the highest probability to meet the user-specified requirement is selected as the first model to approximate computation in the fluid simulation. Then we calculate $CumDi Norms$ in the first check interval, build a linear regression model, calculate $CumDi Norm^{final}$ using the regression model, and predict Q_{loss} by the KNN algorithm. After that, we compare the predicted Q_{loss} (annotated with Q'_{loss} in the rest of the paper) with the user requirement q . If Q'_{loss} is close to q , we predict that the current neural network model can meet the user requirement. The runtime algorithm continues to use the current neural network model. But if Q'_{loss} is larger (or smaller) than q , then the runtime algorithm chooses an accurate (or fast) model with better (or worse) accuracy. If all the neural network models cannot meet q , we restart the simulation and use the traditional simulation method (i.e., the PCG method). The above model switch process happens periodically (the period is the check interval). We calculate $CumDi Norms$ at the end of every check interval to determine if the model switch is necessary.

An Example. Figure 7 gives an example to further explain our runtime algorithm. In this example, we have five neural network models and the user requirements on the simulation quality loss and execution time are 0.013 and 6.64s respectively.

During the offline phase, five neural network models are constructed by the model transformation (Section 4) and selected by MLP (Section 5); We record the possibility and execution time

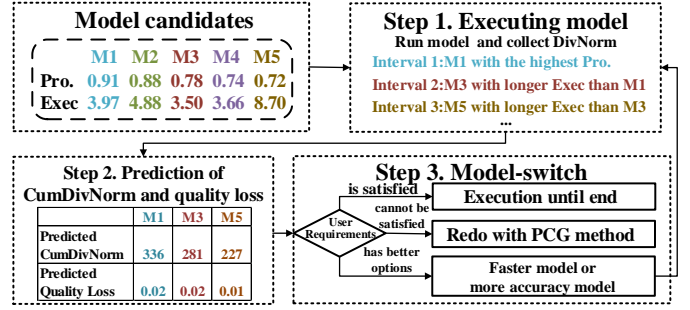


Figure 7: An example to explain our runtime algorithm.

of the five neural network models (shown as Step 1 in Figure 7). At runtime, we use the first neural network ($M1$) which has the highest probability (91%) to meet the user requirement on the simulation quality. Then we skip the first five steps and fit the values of $Di Norms$ in the first check interval into a linear regression model, and predict that the $CumDi Norm$ closes to 395 at the final time step. Using the KNN method, we predict Q_{loss} as 0.019, which is much larger than the user requirement (0.013). So we switch to a more accurate neural network model, i.e., $M3$ (a model with higher accuracy than $M1$ and the highest probability among remaining neural network models). After using $M3$ for another five time steps (one check interval), we predict Q'_{loss} of $M3$ as 0.015, which is still larger than the user requirement. So we switch to another more accurate neural network model, i.e., $M5$. We predict Q'_{loss} of $M5$ as 0.013, which can meet the user requirement, we then use $M5$ during the next check interval.

Compared with using a single neural network model, our runtime algorithm introduces additional computation (i.e., predicting $CumDi Norm$ and applying the KNN method). However, the computation of the simple linear regression algorithm to predict $CumDi Norm$ and the traversal of the binary tree to apply the KNN method are lightweight. Such overhead can be easily outweighed by the performance benefit introduced by adaptive neural network-based approximation. We evaluate performance (including the runtime overhead) in Section 7.2.

7 EVALUATION

We evaluate our framework to examine its impact on performance and simulation quality of the Eulerian fluid simulation.

Platform. We conduct all experiments on a high-end server with 24 Intel Xeon E6-2760 v3 CPU cores running at 2.30GHz. The server is equipped with an NVIDIA Titan X (Pascal) GPU. We use cuDNN 5.0 on this GPU to run neural networks.

Fluid Simulation. During the fluid simulation, we run 128 time steps (the default number of time steps in mantaflow) for each input problem. To comprehensively evaluate the performance, we use multiple grid sizes for each input problem, including $128*128$, $256*256$, $512*512$, $768*768$, and $1024*1024$.

Input Datasets. We generate training and evaluation datasets by mantaflow, which is an open-source framework for the fluid simulation. The training dataset is used to optimize the parameters of the neural network models and MLP model, while the evaluation

Table 2: Percentage of input problems with which the simulation reaches the requirement on quality.

Grid size	128*128	256*256	512*512	768*768	1024*1024
Tompson	68.22%	67.16%	85.27%	71.06%	63.38%
Smart- uidnet	88.27%	87.14%	91.36%	86.47%	70.5%

Figure 8: Performance (execution time) of the Tompson's model and Smart- uidnet.

Figure 9: Variation of quality loss with various grid sizes for input problems.

dataset is used to evaluate the performance of the online algorithm during the uid simulation. Each of the two datasets contains 20,480 input problems. There is no overlapping between the training and test datasets. To generate the total 40,960 problems, we initialize velocity by a pseudo-random turbulent field [2], and generate occupancy grids with the border wall by introducing some objects in the simulation domain. Those objects are from the NTU 3D Model Dataset [33].

Neural Networks . We use Auto-Keras [2] with extension (Section 4) to search qualified convolutional neural network architectures. Auto-Keras constructs neural networks using Python, but the uid simulation is implemented in Lua/Torch7. Hence we re-implement and optimize the neural networks with the Torch7 package.

7.1 Model Speedup and Accuracy

We conduct experiments to measure performance. We take the PCG solver as the baseline method. PCG is the traditional method used in the Eulerian uid simulation and does not include any neural network. Compared with neural network-based approximation, PCG has the highest simulation quality but the performance is very bad. All performance (execution time) reported in this section is shown as speedup with respect to the performance of PCG.

Figure 8 shows the results for Smart- uidnet and the Tompson's model with different grid sizes. The Tompson's model represents the state-of-the-art neural network to accelerate the Eulerian uid simulation. In all test cases, Smart- uidnet performs better than the Tompson's model and is 46% better on average. The largest improvement over the Tompson's model is 25%. Besides the execution time, we also study the simulation quality of the Tompson's

Figure 10: Performance (execution time) for the Tompson's model and Smart- uidnet.

model and Smart- uidnet. We use 20,480 input problems to evaluate each method. We use the simulation quality of PCG as the ground truth and study the quality loss of the Tompson's model and Smart- uidnet. Since Smart- uidnet requires the user to specify a requirement on the quality loss, we use the average quality loss of all input problems when using the Tompson's model, as the user requirement (the target).

Figure 9 presents boxplots to show the results of the quality loss for all input problems with ve selected grid sizes. We draw two observations from Figure 9: (1) The outputs of Smart- uidnet are closer to the target value than those of the Tompson's model; (2) The variances of Smart- uidnet are smaller than those of Tompson's model. These two observations reveal that Smart- uidnet can give more consistent simulation quality than the Tompson's model, which is crucial for dealing with largely diversified input problems.

To further study the consistency of simulation quality with various input problems, we study how many input problems can lead to the simulation with satisfactory simulation quality. Table 2 shows the results. Table 2 reveals that Smart- uidnet leads to a larger percentage of a high-quality simulation than the Tompson's model in all cases. The difference between Smart- uidnet and Tompson's model is as large as 44.67% (when the grid size is 1024*1024).

7.2 Analysis on Runtime System

In this section, we analyze the effectiveness of our runtime system. Taking the grid size of 1024*1024 as an example, the average quality loss and execution time of the Tompson's model are 0.013 and 6.64 seconds respectively. We take this quality loss and execution time as the target (i.e., the user requirement) of Smart- uidnet. Speedup We show the performance (the speedup of execution time) of running the 14 neural network models alone without model switching. We use the performance of PCG as the baseline to calculate speedup. We also show the performance of Smart- uidnet and compare it with the 14 individual models. Figure 10 shows that the performances of the 14 neural network models are quite

¹In the boxplots, the boxes are bounded by 25th and 75th percentiles of the quality loss; The central marks of the boxes indicate the median; The '+' markers outside the boxes indicate the extreme outliers [34].

Figure 11: Variation of simulation quality in different model candidates.

Table 3: Execution time distribution for the five neural network models used by Smart-uidnet at runtime.

Grid size	M7	M5	M10	M2	M13
Prob.(MLP)	86.12%	82.16%	79.43%	74.60%	70.38%
Time Distr.	50.56%	18.10%	11.12%	4.07%	16.15%

different, with the speedup ranging from 541.25 to 141.17. The performance of Smart-uidnet is close to the median performance (440.1) of the 14 neural network models. This is the result of dynamically using different neural network models at runtime. Quality. We compare the 14 neural network models, the Tompson's model, and Smart-uidnet, in terms of quality loss. We calculate the quality loss using the method in Section 7.1. Figure 11 shows the results. Similar to Figure 9 in Section 7.1, the figure shows the distribution and variation using the boxplots. Figure 11 reveals that the variation of the quality loss in Smart-uidnet with various input problems is much smaller than any of the 14 neural network models alone. With Smart-uidnet, 91.05% of the input problems' simulation quality meet the user requirement. With the shortest and longest models (among the 14 neural network models), 12.52% and 92.71% of the input problems' simulation quality meet the user requirement. Figures 10 and 11 include the results for using only the fastest model M1 or the most accurate model M14 throughout the simulation. M1 is 1:18 faster than Smart-uidnet, but achieves the user requirement on the simulation quality in only 12.52% of the input problems (for Smart-uidnet, it is 91.05); M14 achieves the user requirement on the simulation quality in 92.71% of the input problems, which is close to Smart-uidnet, but the performance of M14 is 3:12 worse than Smart-uidnet.

Table 3 shows the time distribution of five neural network models used by Smart-uidnet for all input problems. The second row of the table shows the probability of reaching the target when using each neural network model alone, which is predicted by MLP; The third row shows the percentage of the execution time of the uid simulation for the five neural network models (i.e., the time distribution). The table shows that the model with the highest probability, M7, takes 50.56% of the total execution time, which is the longest execution time among the five models. We also notice that M5, which is the fastest model among the five neural network models, takes 18.1% of the total execution time (the second longest execution time among the five models). These indicate that Smart-uidnet makes best efforts to reach the user requirement on the simulation quality and execution time.

Figure 12: Success rate of reaching target quality with or without using MLP.

7.3 Evaluation of MLP Effectiveness

In this section, we evaluate the effectiveness of MLP. We compare the success rate of our runtime system with and without MLP. The success rate means using all input problems for tests (20,480), how many of them reach the simulation quality requirement with our runtime system. Similar to Section 7.1, we use the average quality loss when using all input problems for the uid simulation with the Tompson's model, as the user requirement. Without MLP, we have 14 neural network models to be used by the runtime before the uid simulation, while with MLP, we have five.

Without MLP, we use the fastest neural network model (but less accurate) in the beginning and then switch to more accurate models until we find a model that can reach the user requirement on the simulation quality. We use this model in the remaining of the uid simulation. Figure 12 shows the results. The figure reveals that Smart-uidnet with MLP causes higher success rates than without MLP. The success rate with MLP is 88.86% on average and can be up to 91.36%. This result shows that without MLP, the runtime system can use those neural network models that have a lower possibility to reach the quality target, in order to have better performance. With MLP, we avoid applying those models, hence improving the success rate. We also compare performance with and without MLP. For the grid sizes of 128, 256, 512, 768, and 1024, the corresponding performance when we use MLP, which is normalized by the performance without MLP, is 97%, 84%, 92%, 79%, and 83% respectively. With MLP, we perform better in all cases.

7.4 Sensitivity Study: Check Interval

We study the impact of the check interval on execution time and success rate. We change the check interval, and Figure 13 shows the results. Figure 13 shows that the success rate decreases when the interval increases. Such a decrease is because the model switching is too slow to achieve high simulation quality. We also observe an unusual increase in the success rate when the interval changes from 14 to 16. We attribute such an increase to the statistical variance of using the linear regression method to make the prediction. Although using 16 seems to be useful to increase prediction accuracy for our input problems, using 5 achieves the highest success rate.

Therefore, we use 5 as the check interval throughout our evaluation. We do not use a check interval smaller than 5, because we have to skip the first two time steps to ensure that the growth rate of

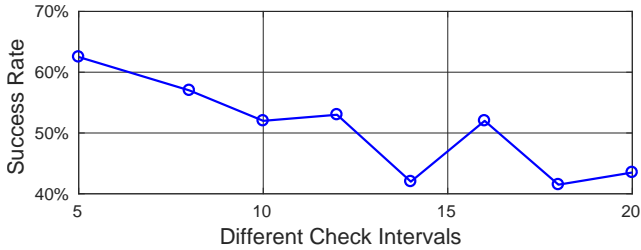


Figure 13: Impact of the check interval on the success rate.

Table 4: Resource usage of different methods.

Methods	FLOP (single step)	GPU Memory
PCG	~1,250 M	332 MB
Tompson	243.79 M	299 MB
Smart-fluidnet	110.97 M	1,069 MB

CumDi Norm is stable and using less than three time steps to build the linear regression model cannot give high prediction accuracy.

7.5 Evaluation of Resource Usage

We evaluate resource usage (FLOP and GPU memory consumption) by PCG, the Thompson’s and Smart-fluidnet with the grid size of 512×512 . Since all input problems for such a grid size has the same resource usage, we randomly choose an input problem and present the evaluation results in Table 4.

We find that Smart-fluidnet requires much less FLOP than PCG and the Thompson’s, which explains why Smart-fluidnet has better performance. On the other hand, Smart-fluidnet consumes more memory than PCG and the Thompson’s, because Smart-fluidnet uses five neural network models on GPU (but not running them simultaneously). However, the memory consumption of Smart-fluidnet is still smaller than the GPU memory capacity (12 GB). If GPU memory is not sufficient, then we may either use fewer neural network models or run them on the CPU.

8 RELATED WORK

Applying neural networks to HPC applications. Neural networks (especially deep neural networks (DNNs)) have been employed in HPC applications recently. In [35], a neural network is used to generate input data for modulating the simulation process. Wigley et al. [7] propose a neural network-based online optimization process for the Bose-Einstein condensates (BEC). In [36], neural networks are used to accelerate thermal-hydraulic modeling and enables faster assessment for the dynamic thermal-hydraulic system. Richard et al. [35] apply a neural network to ITER magnets to predict the occurrence of the interruption, which can be used to adjust the reaction to continue generating power and avoid ITER damage. Mathuriya et al. [37] build a CNN model to determine the physical model that describes our universe.

In our work, we use neural network models to approximate computation in an HPC application. Different from the existing efforts, we aim to address the limitation of model flexibility and generality in the existing work [10].

Acceleration of the fluid simulation. Since the fluid simulation is an important HPC application, many research efforts have been focusing on improving its performance. Popovic et al. use a multi-grid approach to pre-process data for the PCG method in the traditional fluid simulation [38].

Molemaker et al. use iterated orthogonal projections and Michael Lentine et al. apply a coarse-grid correction method [39] to approximate the Poisson’s equation [40]. These two approaches are effective, but both of them are inexact and only competitive in low-resolution settings. Some recent efforts [41] address the above limitation by using a data-driven approach or leveraging some useful statistical characteristics in the data distribution. Our work is different from them because we use neural network models (not traditional solvers) to improve the performance of the fluid simulation.

Some recent efforts attempt to build a neural network model that makes prediction for stability, collisions, forces, and velocities of data objects in images or videos [42]. Given pressure data from previous frames, voxel occupancy, and velocity divergence, Yang et al. use a patch-based neural network to predict the next pressure [43]. Jonathan et al. use a convolutional neural network to predict the pressure value in the Eulerian fluid simulation [10]. The existing NN-based approximation approaches lack flexibility and generalization, discussed in Section 1.

Model compression. During the process of model construction, we use some operations to generate simpler neural network models. The recent work using model compression aims to generate simpler models [44]. There are multiple techniques for model compression, including quantization parameters [45, 46], layer pruning [47, 48], binarized networks [49, 50], low rank approximation [51], and knowledge distillation [52, 53]. Different from the existing work that focuses on resource-constrained execution environment (e.g., mobile devices), we focus on simplifying model without resource constrains and study how to build the models to meet the simulation quality for an HPC application.

9 CONCLUSIONS

Using machine learning (especially neural networks) to approximate computation in HPC applications and improve performance has shown preliminary success recently. However, using this approach faces fundamental limitations due to the lack of model flexibility and generality. In this paper, we focus on a specific HPC application and introduce a systematic approach to address the above limitation. In particular, we introduce a framework (Smart-fluidnet) that automatically uses multiple neural network models at runtime to approximate computation and make best efforts to meet the user requirements on simulation quality and execution time. The framework includes a series of techniques to construct and select neural network models. Based on the comprehensive evaluation, we show that Smart-fluidnet is $1.46\times$ and $590\times$ faster than a state-of-the-art neural network model and the original fluid simulation respectively on an NVIDIA Pascal GPU, while providing better simulation quality than the state-of-the-art model.

Acknowledgement. This work was partially supported by U.S. National Science Foundation (CNS-1617967, CCF-1553645 and CCF-1718194) and Chameleon Cloud. We thank reviewers and our shepherd for their constructive comments.

REFERENCES

- [1] Anne H de Boer, Paul Hagedoorn, Robert Woolhouse, and Ed Wynn. Computational fluid dynamics (cfD) assisted performance evaluation of the twincer™ disposable high-dose dry powder inhaler. *Journal of Pharmacy and Pharmacology*, 64(9):1316–1325, 2012.
- [2] Sanghun Choi, Shinjiro Miyawaki, and Ching-Long Lin. A feasible computational fluid dynamics study for relationships of structural and functional alterations with particle depositions in severe asthmatic lungs. *Computational and mathematical methods in medicine*, 2018.
- [3] David R Rutkowski, Scott B Reeder, Luis A Fernandez, and Alejandro Roldán-Alzate. Surgical planning for living donor liver transplant using 4d flow mri, computational fluid dynamics and in vitro experiments. *Computer Methods in Biomechanics and Biomedical Engineering: Imaging & Visualization*, 6(5):545–555, 2018.
- [4] Yidong Xia, Ansel Blumers, Zhen Li, Lixiang Luo, Yu-Hang Tang, Joshua Kane, Hai Huang, Matthew Andrew, Milind Deo, and Jan Goral. A gpu-accelerated package for simulation of flow in nanoporous source rocks with many-body dissipative particle dynamics. *arXiv preprint arXiv:1903.10134*, 2019.
- [5] Troy Snyder and Minel Braun. A cfd-based frequency response method applied in the determination of dynamic coefficients of hydrodynamic bearings. part 1: Theory. *Lubricants*, 7(3):23, 2019.
- [6] Alexander Radovic. Neutrino Identification with a Convolutional Neural Network in the NOvA Detectors. In *International Conference on High Energy Physics*, 2016.
- [7] P B. Wigley, P J. Everitt, Anton Hengel, John Bastian, M A. Sooriyabandara, Gordon McDonald, Kyle Hardman, C D. Quinlivan, Manju Perumbil, Carlos claiton Noschang kuhn, I R. Petersen, Andre Luiten, J Hope, N Robins, and Michael Hush. Fast machine-learning online optimization of ultra-cold-atom experiments. *Scientific Reports*, 6:25890, 2015.
- [8] Evan Racah, Christopher Beckham, Tegan Maharaj, Samira Kahou, Mr. Prabhath, and Chris Pal. ExtremeWeather: A Large-scale Climate Dataset for Semi-supervised Detection, Localization, and Understanding of Extreme Weather Events. In *NIPS*, 2017.
- [9] Byungsoo Kim, Vinicius C Azevedo, Nils Thuerey, Theodore Kim, Markus Gross, and Barbara Solenthaler. Deep fluids: A generative network for parameterized fluid simulations. In *Computer Graphics Forum*, volume 38, pages 59–70. Wiley Online Library, 2019.
- [10] Jonathan Tompson, Kristofer Schlachter, Pablo Sprechmann, and Ken Perlin. Accelerating eulerian fluid simulation with convolutional networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 3424–3433. JMLR. org, 2017.
- [11] Cheng Yang, Xubo Yang, and Xiangyun Xiao. Data-driven projection method in fluid simulation. *Computer Animation and Virtual Worlds*, 27(3-4):415–424, 2016.
- [12] Haifeng Jin, Qingquan Song, and Xia Hu. Efficient neural architecture search with network morphism. *arXiv preprint arXiv:1806.10282*, 2018.
- [13] Hugo Jair Escalante, Wei-Wei Tu, Isabelle Guyon, Daniel L Silver, Evelyne Viegas, Yuqiang Chen, Wenyuan Dai, and Qiang Yang. Automl@ neurips 2018 challenge: Design and results. *arXiv preprint arXiv:1903.05263*, 2019.
- [14] Nils Thuerey and Tobias Pfaff. Mantaflow. <http://mantaflow.com>, 2016.
- [15] Francis H Harlow and J Eddie Welch. Numerical calculation of time-dependent viscous incompressible flow of fluid with free surface. *The physics of fluids*, 8(12):2182–2189, 1965.
- [16] Alan Jeffrey. *Applied partial differential equations: an introduction*. Academic Press, 2003.
- [17] Zhaosheng Yu. A DLM/FD method for fluid/flexible-body interactions. *Journal of computational physics*, 207(1):1–27, 2005.
- [18] Joseph L Steger and RF Warming. Flux vector splitting of the inviscid gasdynamic equations with application to finite-difference methods. *Journal of computational physics*, 40(2):263–293, 1981.
- [19] NN Yanenko. Simple schemes in fractional steps for the integration of parabolic equations. In *The Method of Fractional Steps*, pages 17–41. Springer, 1971.
- [20] Alexandre Joel Chorin. Numerical solution of the navier-stokes equations. *Mathematics of computation*, 22(104):745–762, 1968.
- [21] Aleka McAdams, Eftychios Sifakis, and Joseph Teran. A parallel multigrid poisson solver for fluids simulation on large grids. In *Proceedings of the 2010 ACM SIGGRAPH/Eurographics Symposium on Computer Animation*, pages 65–74. Eurographics Association, 2010.
- [22] Olivier Gènevaux, Arash Habibi, and Jean-Michel Dischler. Simulating fluid-solid interaction. In *Graphics Interface*, volume 2003, pages 31–38, 2003.
- [23] Ronald Fedkiw, Jos Stam, and Henrik Wann Jensen. Visual simulation of smoke. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, pages 15–22. ACM, 2001.
- [24] Evan R. Sparks, Ameet Talwalkar, Daniel Haas, Michael J. Franklin, Michael I. Jordan, and Tim Kraska. Automating model search for large scale machine learning. In *Proceedings of the Sixth ACM Symposium on Cloud Computing*, 2015.
- [25] Tian Li, Jie Zhong, Ji Liu, Wentao Wu, and Ce Zhang. Ease.ml: Towards multi-tenant resource sharing for machine learning workloads. *Proc. VLDB Endow.*, 11(5), 2018.
- [26] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [27] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [28] MD Nefzger and James Drasgow. The needless assumption of normality in pearson’s r. *American Psychologist*, 12(10):623, 1957.
- [29] Sten Henrysson. Gathering, analyzing, and using data on test items. *Educational measurement*, 2, 1971.
- [30] DJ Best and DE Roberts. Algorithm as 89: the upper tail probabilities of spearman’s rho. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 24(3):377–379, 1975.
- [31] Nian Shong Chok. *Pearson’s versus Spearman’s and Kendall’s correlation coefficients for continuous data*. PhD thesis, University of Pittsburgh, 2010.
- [32] Theodore Kim, Nils Thuerey, Doug James, and Markus Gross. Wavelet turbulence for fluid simulation. In *ACM Transactions on Graphics (TOG)*, volume 27, page 50. ACM, 2008.
- [33] Jiantao Pu and Karthik Ramani. On visual similarity based 2d drawing retrieval. *Computer-Aided Design*, 38(3):249–259, 2006.
- [34] Robert Dawson. How significant is a boxplot outlier? *Journal of Statistics Education*, 19(2), 2011.
- [35] L Savoldi Richard, R Bonifetto, Stefano Carli, A Froio, A Foussat, and R Zanino. Artificial neural network (ann) modeling of the pulsed heat load during iter cs magnet operation. *Cryogenics*, 63:231–240, 2014.
- [36] Stefano Carli, R Bonifetto, L Savoldi, and R Zanino. Incorporating artificial neural networks in the dynamic thermal-hydraulic model of a controlled cryogenic circuit. *Cryogenics*, 70:9–20, 2015.
- [37] Amrita Mathuriya, Deborah Bard, Peter Mendygral, Lawrence Meadows, James Armemann, Lei Shao, Siyu He, Tuomas Kärrä, Diana Moise, Simon J Pennycook, et al. Cosmoflow: using deep learning to learn the universe at scale. In *SC18: International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 819–829. IEEE, 2018.
- [38] Aleka McAdams, Eftychios Sifakis, and Joseph Teran. A Parallel Multigrid Poisson Solver for Fluids Simulation on Large Grids. In MZoran Popovic and Miguel Otaduy, editors, *Eurographics/ ACM SIGGRAPH Symposium on Computer Animation*. The Eurographics Association, 2010.
- [39] Michael Lentine, Wen Zheng, and Ronald Fedkiw. A novel algorithm for incompressible flow using only a coarse grid projection. In *ACM Transactions on Graphics (TOG)*, volume 29, page 114. ACM, 2010.
- [40] Jeroen Molemaker, Jonathan M. Cohen, Sanjit Patel, and Jonyong Noh. Low viscosity flow simulations for animation. pages 9–18, 01 2008.
- [41] Tyler De Witt, Christian Lessig, and Eugene Fiume. Fluid simulation using laplacian eigenfunctions. *ACM Trans. Graph.*, 31(1):10:1–10:11, February 2012.
- [42] Adam Lerer, Sam Gross, and Rob Fergus. Learning physical intuition of block towers by example. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning-Volume 48*, pages 430–438. JMLR. org, 2016.
- [43] Yang Hong, Kuo-Lin Hsu, Soroosh Sorooshian, and Xiaogang Gao. Precipitation estimation from remotely sensed imagery using an artificial neural network cloud classification system. *Journal of Applied Meteorology*, 43(12):1834–1853, 2004.
- [44] Yu Cheng, Duo Wang, Pan Zhou, and Tao Zhang. A survey of model compression and acceleration for deep neural networks. *arXiv preprint arXiv:1710.09282*, 2017.
- [45] 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 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2704–2713, 2018.
- [46] Shuang Wu, Guoqi Li, Feng Chen, and Luping Shi. Training and inference with integers in deep neural networks. *arXiv preprint arXiv:1802.04680*, 2018.
- [47] Tien-Ju Yang, Yu-Hsin Chen, and Vivienne Sze. Designing energy-efficient convolutional neural networks using energy-aware pruning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5687–5695, 2017.
- [48] Pavlo Molchanov, Stephen Tyree, Tero Karras, Timo Aila, and Jan Kautz. Pruning convolutional neural networks for resource efficient inference. *arXiv preprint arXiv:1611.06440*, 2016.
- [49] Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. Xnor-net: Imagenet classification using binary convolutional neural networks. In *European Conference on Computer Vision*, pages 525–542. Springer, 2016.
- [50] Matthieu Courbariaux, Itay Hubara, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Binarized neural networks: Training deep neural networks with weights and activations constrained to +1 or -1. *arXiv preprint arXiv:1602.02830*, 2016.
- [51] Cheng Tai, Tong Xiao, Yi Zhang, Xiaogang Wang, et al. Convolutional neural networks with low-rank regularization. *arXiv preprint arXiv:1511.06067*, 2015.
- [52] Bharat Bhuvan Sau and Vineeth N Balasubramanian. Deep model compression: Distilling knowledge from noisy teachers. *arXiv preprint arXiv:1610.09650*, 2016.

- [53] Raphael Gontijo Lopes, Stefano Fenu, and Thad Starner. Data-free knowledge distillation for deep neural networks. *arXiv preprint arXiv:1710.07535*, 2017.