

MULTIRESOLUTION CONVOLUTIONAL AUTOENCODERS *

YUYING LIU[†], COLIN PONCE[‡], STEVEN L. BRUNTON[§], AND J. NATHAN KUTZ[¶]

Abstract. We propose a *multi-resolution convolutional autoencoder* (MrCAE) architecture that integrates and leverages three highly successful mathematical architectures: (i) multigrid methods, (ii) convolutional autoencoders and (iii) transfer learning. The method provides an adaptive, hierarchical architecture that capitalizes on a progressive training approach for multiscale spatio-temporal data. This framework allows for inputs across multiple scales: starting from a compact (small number of weights) network architecture and low-resolution data, our network progressively deepens and widens itself in a principled manner to encode new information in the higher resolution data based on its current performance of reconstruction. Basic transfer learning techniques are applied to ensure information learned from previous training steps can be rapidly transferred to the larger network. As a result, the network can dynamically capture different scaled features at different depths of the network. The performance gains of this adaptive multiscale architecture are illustrated through a sequence of numerical experiments on synthetic examples and real-world spatial-temporal data.

Key words. convolutional autoencoder, multiresolution analysis, multigrid, transfer learning, model scaling, multi-scale dynamics

AMS subject classifications. 65T99, 37N10, 37M10

1. Introduction. The multiscale spatio-temporal dynamics observed in many complex systems poses significant challenges for modeling and prediction. Although we are often primarily interested in macroscale phenomena, the microscale dynamics must also be modeled and understood, as it plays an important role in driving the macroscale behavior. Indeed, a given system may have multiple fast and slow time scales as well as a number of micro to macro spatial scales that interact to produce the complex dynamics observed. This makes modeling multiscale systems particularly difficult unless the time scales are disambiguated in a principled way. Coarse graining and mesh-refinement (multigrid methods) are two principled mathematical techniques for addressing multiscale behavior. In the former, the microscale physics are averaged over to produce an effective macroscale model, while in the latter, computational models are refined where the macroscale variables produce large errors. In this paper, we present a *multi-resolution convolutional autoencoder* (MrCAE), a decomposition scheme inspired by multigrid computational methods to progressively refine a multiscale description of spatio-temporal data. It leverages and exploits aspects of multigrid methods and transfer learning to produce an effective multi-scale analysis tool for characterizing large scale spatio-temporal data.

Multigrid methods [27, 38] have been extensively developed for physics-based simulation models where coarse grained models must be progressively refined in order to achieve a required numerical precision while keeping the simulation tractable. Multigrid architectures provide a principled method for targeting the refinement process, constituting a mature field with wide spread applications in the engineering and physical sciences. In contrast, coarse graining methods attempt to construct a macroscale physics model by progressive construction of coarse grained variables and their dynamics. Mathematical algorithms such as the *heterogeneous multiscale*

*Submitted to the editors DATE.

[†]Department of Applied Mathematics, University of Washington, Seattle, WA (yliu814@uw.edu).

[‡]Lawrence Livermore National Lab, Livermore, CA

[§]Mechanical Engineering, University of Washington, Seattle, WA

[¶]Department of Applied Mathematics, University of Washington, Seattle, WA

modeling (HMM) [41, 40] and *equation-free method* [19] provide principled methods for multiscale systems. Additional work has focused on testing for the presence of multiscale dynamics so that analyzing and simulating multiscale systems is more computationally efficient [8, 9].

Data-driven methods, specifically neural networks (NNs), have emerged as an attractive alternative for characterizing multi-scale physics [11, 43, 39, 26, 24, 6, 14, 31, 32]. The structure of *convolutional neural networks* (CNNs) are especially relevant for multiscale data as the convolutional window extracts features of the data at an appropriate, coarse- or fine-grained resolution. As a result, CNNs have demonstrated exceptional performance in image processing tasks [12]. Indeed, in the wake of AlexNet [20], many aspects of CNN design have been thoroughly studied individually, such as the spatial filters [35, 34], nonlinear activation functions [42], width and depth of the network [45], skip connections [15], batch normalization [17], etc [1]. As more complex neural architectures and designs arise, NNs have become increasingly popular, which enables the process of automating architecture engineering by jointly considering all design factors [37, 23, 46]. However, the design choices of NN algorithms are highly automated and often uninterpretable, which makes integrating domain knowledge difficult.

Although CNNs can capitalize on the multiscale features of data, because they begin processing with only small local patches, they often fail to exploit the large-scale, low-dimensional structure typically observed in spatio-temporal data. Such structure is routinely leveraged for reduced order modeling [2, 21, 4, 30, 16, 36] of high-dimensional spatio-temporal systems. This motivates the innovations of the current work; here, we propose a multi-resolution convolutional autoencoder (MrCAE) that begins by processing on downsampled (ie “coarse”) data in order to capture large-scale, low-dimensional structure, and then progressively refines both the data and the neural network while employing transfer learning in building each stage’s neural network. This progress-refinement approach provides a principled framework for leveraging multiresolution and transfer learning ideas, enabling one to build multiscale models for spatio-temporal data that result in more compact networks and more compact encodings than those produced by traditional CAE methods. We envision our proposed architecture to be a critical piece of many multi-scale neural network architectures, especially as encoder-decoder based models have already resulted in many successful applications, such as linearization of dynamics [24, 10], model discovery via sparse regression [6], forecasting [33], etc.

Importantly, our architecture leverages data in an intelligent way: querying data that targets the refinement process. Thus a hierarchy of models is constructed with less data, producing highly compact networks and encodings. The method is well aligned with the arguments in [7]: ‘*Initially, a small model may be preferred, in order to prevent overfitting and to reduce the computational cost of using the model. Later, a large model may be necessary to fully utilize the large dataset.*’

Our paper is organized as follows: the multi-resolution network architecture is proposed in section 2, the progressive training algorithm is presented in section 3, experimental results are in section 4, and the conclusions and discussions follow in section 5. Our code is publicly available at <https://github.com/luckystarufo/MrCAE>.

2. Network Architecture. section 2 details the overall architecture of our MrCAE. It consists of several levels of autoencoders. Each level is associated with the data of a particular resolution. The network is built up recursively: in the base level, a small autoencoder is trained to encode and reconstruct the coarsest data. We then

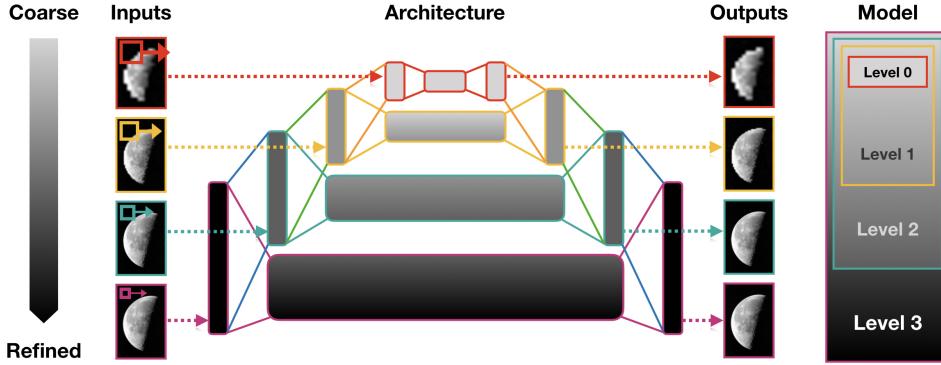


Fig. 1: **A Schematic Overview of the Network Architecture.** In this example, there are 4 levels for the architecture which are colored in red, yellow, cyan and purple respectively. They are recursively built up to process data across different resolutions — architectures built for processing coarser data are later embedded into the next-level architectures to ensure knowledge transfer.

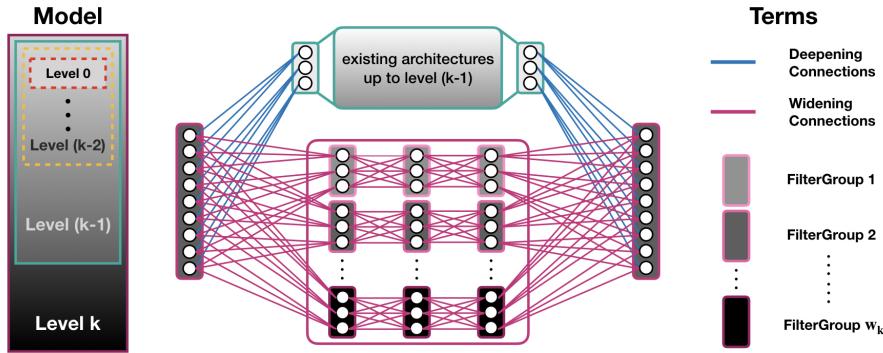


Fig. 2: **Build up Network Architecture at Level k .** Within level k , we perform one deepening operation and a sequence of widening operations. The deepening operation (shown in blue) is the transfer learning step, creating a convolutional filter that connects the new input (fine) to the previous level input (coarse). Widening operations (shown in purple) are performed sequentially by allocating more convolutional filters so that it can capture new, higher-resolution features.

refine the data so that our inputs are of a higher resolution and embed the existing architecture into a new autoencoder (the transfer learning step) to learn the new, refined data. In this way, we construct a hierarchy of trained neural networks. The mathematical details of the construction are given in the subsections that follow.

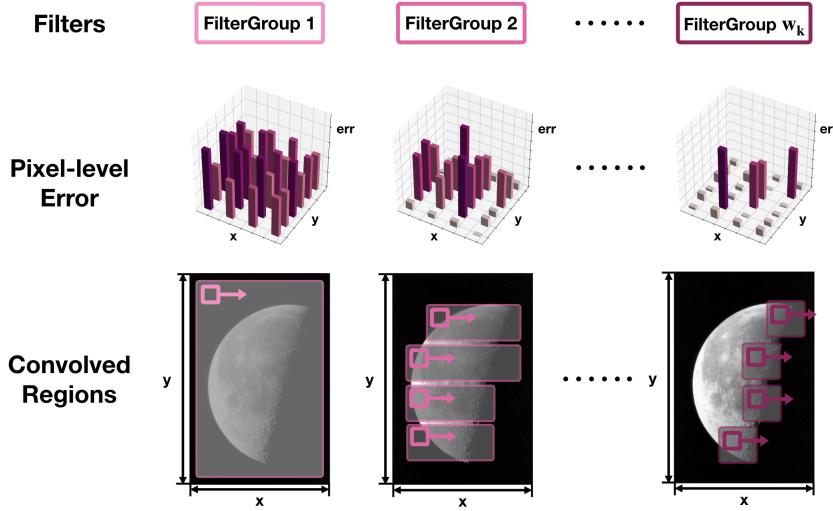


Fig. 3: **Adaptive Filters (I).** For each widening operation shown in section 2, new convolutional filters are only applied to the regions that are still poorly resolved. The progressive refinement ensures a parsimonious use of the parameters.

2.1. Hierarchical Construction. The technical details for constructing each level of the network is demonstrated in section 2: we perform a deepening operation (except for the base level) followed by a sequence of widening operations.

The *deepening operation* inserts a convolutional/deconvolutional layer between the current and previous level inputs/outputs. This is the transfer learning step. We denote the inserted convolutional filter as $c^{(k)}$ and the deconvolutional filter as $d^{(k)}$, and let $f_{\theta_{k-1}}$ be the existing network at level $(k-1)$, $g^{(k)}$ be the network after we apply this deepening operation. Mathematically, we have

$$(2.1) \quad g^{(k)} = d^{(k)} \circ f_{\theta_{k-1}} \circ c^{(k)}$$

In our setup, we make the resolution differ by a factor of 2 along each spatial dimension between two adjacent levels, which is typical of a multi-resolution analysis [21]. As a result, convolutional/deconvolutional filters with a kernel size 3×3 and stride size 2 are used to ensure compatible dimensions. Similar to [7], the two inserted layers are properly initialized to ensure knowledge transfer. More technical details about the initialization are covered in the Appendix A. There is another, perhaps more intuitive, way to understand the design: data at coarser levels can be regarded as obtained by applying a sequence of hard-coded down-sampling (convolutional) operations to the finest data and vice versa. When we proceed to a new level, the network replaces one of these hard-coded operators with trainable convolutional/deconvolutional filters.

The *widening operation* expands the network capacity in order to capture the new, finer-grained features of the higher resolution data. Specifically, it adds a group of new convolutional pathways through the network, and then continues training. See the FilterGroups in Figure ???. Note that each new pathway is separate, not adding any connections to old pathways; this reduces the risks of overfitting. Mathematically speaking, the widening operations at level k create a list of new pathways for the

network, we denote them as $w_1^{(k)}, w_2^{(k)}, \dots, w_{k_n}^{(k)}$. Then we have,

$$(2.2) \quad f_{\theta_k} = g^{(k)} + w_1^{(k)} + \dots + w_{w_k}^{(k)}.$$

2.2. Adaptive Filters. The way we perform a sequence of widening operations is shown in section 2. Suppose we have performed j groups of widening operations and finished the subsequent training. We proceed by computing the mean squared error of a 3×3 region around each pixel. Our $(j+1)^{th}$ group of filters will be only applied to the regions where the reconstructions exhibit high error; details are explained in section 3. Then a subsequent training is performed for the new architecture. We do this recursively until all regions can be reconstructed within some error threshold/bound. This process bears a close resemblance to mesh refinement [3] of multigrid methods [27, 38] and therefore is highly adaptive and ensures a parsimonious use of parameters.

It is worth noting that we do not use nonlinear activation functions, in contrast to the classical designs in which convolutional filters are usually followed by a Rectified Linear Units (ReLU). We find it actually gives us a performance boost in terms of the training progress while also making the filters highly interpretable. In our design, the highly adaptive filters play the role of the nonlinear activation functions: by explicitly specifying where to apply the filters, we are forcing those filters to be activated in the relevant regions while leaving them inhibited in other regions. Additionally, this implies that the resulting, trained network is fully linear, enabling extremely fast matrix-based representations.

3. Training.

3.1. Framework. To accompany our adaptive architecture, we have developed a progressive training framework which is outlined in Algorithm 3.1. The training is performed whenever the architecture changes, thus it can effectively utilize the new pathways for representing new data. Early stopping criteria are implemented as well, however, we reference the full training details in Appendix B in order to maintain clarity and simplicity for the description of the architecture.

Algorithm 3.1 Progressive Training

```

Define  $N$  to be the number of levels
Define  $E$  to be the maximum training epochs
Define training data with increasing resolutions  $\{D^{(0)}, \dots, D^{(N-1)}\}$ 
Initialize a model object  $M$ 
for  $i$  in  $0, 1, \dots, N - 1$ : do
    Perform deepening operation on  $M$ ;
    Perform  $E$  epochs training;
    while Reconstruction at this level is not fully resolved do
        Perform widening operation on  $M$ ;
        Perform  $E$  epochs training;
    end while
end for
return  $M$ 

```

3.2. Loss Function. Our loss function is designed to capture the general, low-rank dynamics as well as identifying outlier/singularities. It is of the same functional form across each level. However, different resolution data are used to calculate this

quantity. Suppose we are training the network at level k ($0 \leq k < N$) with data set $D^{(k)} \in R^{m_k \times n_k \times T_k}$ (m_k, n_k, T_k are the width, height of each snapshot and the number of snapshots), and $\widehat{D}^{(k)} := f_{\theta_k}(D^{(k)})$ is the corresponding reconstruction through the network, where θ_k are the parameters of the current network f_{θ_k} at level k ¹. Let i, j, t to be the indices of the row, column and snapshot. We formulate our loss function $\mathcal{L}(\theta_k, D^{(k)})$ to be the following:

$$(3.1) \quad \mathcal{L}(\theta_k, D^{(k)}) = \omega \mathcal{L}_{mse}(\theta_k, D^{(k)}) + (1 - \omega) \mathcal{L}_{max}(\theta_k, D^{(k)})$$

where

$$(3.2) \quad \mathcal{L}_{mse}(\theta_k, D^{(k)}) = \frac{1}{m_k n_k T_k} \sum_{i=1}^{m_k} \sum_{j=1}^{n_k} \sum_{t=1}^{T_k} (D_{i,j,t}^{(k)} - \widehat{D}_{i,j,t}^{(k)})^2$$

$$(3.3) \quad \mathcal{L}_{max}(\theta_k, D^{(k)}) = \max_{i,j} \frac{1}{T_k} \sum_{t=1}^{T_k} (D_{i,j,t}^{(k)} - \widehat{D}_{i,j,t}^{(k)})^2$$

The first term in (3.1) is the classic mean squared error (MSE) to ensure an overall satisfactory reconstruction. The second term captures the worst case scenario which primarily corresponds to the highly variable dynamics over some regions in the spatio-temporal data. The loss function is a weighted sum of these two terms modulated by a control parameter ω ($0 \leq \omega \leq 1$). Based on our experience, we find a small ω (eg. 0.5) is ideal for dynamics that have a clear separation of spatial scales, so that the network pays more attention to the singularities or highly variable regions. But for those without persistent spatial patterns, we recommend setting ω close to 1.

3.3. Measuring Error. In addition to the loss function in (3.1), we define two other important quantities. The first quantity is developed for tracking the overall training progress. The loss function in (3.1) only gives us the information of how it performs at a specific level. Nothing can be said about the reconstructions at higher levels. In this case, knowledge transfer between adjacent levels cannot be evaluated. What we need is a metric that reflects the overall training progress. Therefore, we define a metric that estimates the reconstruction error with respect to the finest resolution data at each level:

$$(3.4) \quad \mathcal{L}^{global}(\theta, D^{(N-1)}) = \omega \mathcal{L}_{mse}^{global}(\theta, D^{(N-1)}) + (1 - \omega) \mathcal{L}_{max}^{global}(\theta, D^{(N-1)})$$

$$(3.5) \quad \mathcal{L}_{mse}^{global} = \frac{1}{m_{N-1} n_{N-1} T_{N-1}} \sum_{i=1}^{m_{N-1}} \sum_{j=1}^{n_{N-1}} \sum_{t=1}^{T_{N-1}} (D_{i,j,t}^{(N-1)} - \mathcal{U}^{(N-k-1)}(\widehat{D}_{i,j,t}^{(k)}))^2$$

$$(3.6) \quad \mathcal{L}_{max}^{global} = \max_{i,j} \frac{1}{T_{N-1}} \sum_{t=1}^{T_{N-1}} (D_{i,j,t}^{(N-1)} - \mathcal{U}^{(N-k-1)}(\widehat{D}_{i,j,t}^{(k)}))^2$$

¹In fact, at level k , the network keeps growing itself as we perform the deepening operation and the sequence of widening operations, which gives rise to different networks $g_k, f_{\theta_k}^{(1)}, \dots, f_{\theta_k}^{(k_n)}$ respectively. Here, without causing ambiguity, we refer them all as f_{θ_k} .

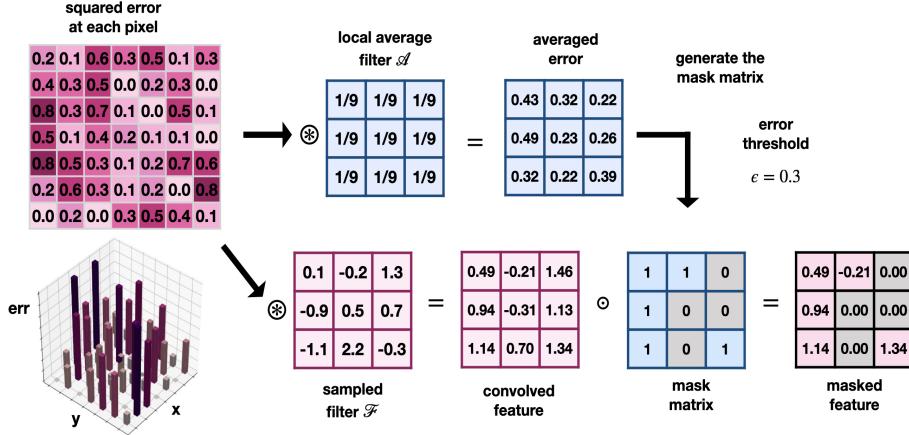


Fig. 4: **Adaptive Filters (II).** This picture illustrates the implementation details of applying adaptive filters. It is done by calculating the mask (shown at the top) and convolved features of all regions (shown at the bottom) followed by a point-wise multiplication.

where

$$(3.7) \quad \mathcal{U}^{(s)} = \underbrace{\mathcal{U} \circ \mathcal{U} \circ \dots \circ \mathcal{U}}_s$$

is a composition of a sequence of up-sampling operation \mathcal{U} which is chosen to be a bi-linear interpolation with a kernel size 3×3 and a stride size 2. This is similar to the prolongation operator from multigrid methods. It should be noted that this metric is only proposed for the purpose of validating the effectiveness of knowledge transfer across different levels of the model, it would be absent in the realistic settings where we have a growing data set so that the finest data $D^{(N-1)}$ may not be accessible initially.

The other metric is used to provide feedback for the adaptive filters. Within a widening operation, the new group of convolutional filters are only applied to the regions that still need to be further refined. Technically, this is achieved by applying the filters to all regions of the inputs first (and therefore obtaining the convolved features), then determining the irrelevant ones. To tell if a convolved feature should be masked out or not, we first calculate the mean squared error of each pixel along the time axis, followed by performing a down-sampling operation \mathcal{A} (local average) since the width and height of the convolved features are both reduced by a factor of 2 compared to the original inputs. Finally, we threshold it by a prescribed tolerance ϵ to obtain the mask \mathcal{M} :

$$(3.8) \quad \mathcal{M} = \mathbb{1} \left\{ \mathcal{A} \left[\frac{1}{T_k} \sum_{t=1}^{T_k} (D_{:, :, t}^{(k)} - \hat{D}_{:, :, t}^{(k)})^2 \right]_{i,j} \geq \epsilon \right\}$$

The process is illustrated in subsection 3.3.

4. Experimental results. In this section, we test the performance of our Mr-CAE on a range of data displaying increasing levels of complexity. Several unsteady

fluid flow fields are analyzed, as fluid dynamics exhibit a range of multiscale behavior and have been the focus of intense modeling efforts, with and without machine learning [36, 5]. In the first part, we show the performance of our network on each individual case. Specifically, we show the reconstructions on some testing snapshots and heat maps of the regions being refined across each level and the error plot with respect to (3.4) that reflects the overall training progress. In the second part, we benchmark our network against a structurally similar network: residual encoder-decoder network (RED-Net[28]). We show our architecture scales better in terms of the number of network parameters and the size of encoding.

All data are resized so that each snapshot has the shape $(2^p - 1) \times (2^q - 1)$ (with some positive integers p and q) before data ingestion into the network, all snapshots are shuffled randomly and split in to training set (70%), validation set (20%) and testing set (10%).

4.1. Individual experiments.

4.1.1. Two oscillatory modes. For the first example (subsection 4.1), we consider two nonlinear spatial modes driven by sinusoidal temporal modes with different frequencies:

$$(4.1) \quad \Phi(x, y, t) = u(x, y) \cos(\omega_0 t) + v(x, y) \cos(\omega_1 t + \frac{\pi}{4})$$

where

$$(4.2) \quad \begin{aligned} u(x, y) &= \cosh\left(\frac{x+1}{\sigma_0}\right) \cosh\left(\frac{y-1}{\sigma_0}\right) \\ v(x, y) &= \frac{1}{(2\pi\sigma_1)^{1/2}} \exp\left(-\frac{(x-1)^2 + (y+1)^2}{2\sigma_1^2}\right) \end{aligned}$$

We generate it in the domain of $[-5, 5] \times [-5, 5]$ using 127×127 grids, and 500 snapshots are uniformly collected from the time interval $[0, 8\pi]$. We set $\omega_0 = 0.5$, $\omega_1 = 4.0$, $\sigma_0 = 10.0$ and $\sigma_1 = 0.25$ in our experiment.

subsection 4.1a shows the synthetic spatio-temporal data. We visualize the two different spatial modes and their associated temporal modes. 6 sampled snapshots are drawn from the test set which are marked in red. This is an example where small spatial (high-frequency) content is modulated by a large, slow-varying background mode.

For this example, we construct a network of 3 levels, with each level having 1, 2 and 3 groups of adaptive filters (a.k.a. widening operations) correspondingly. The output reconstructions across different levels shown in subsection 4.1b become sharper and sharper as the network grows. The regions that adaptive filters being applied are shown in subsection 4.1c. One can clearly see the network faithfully picks up the high-frequency content with more convolutional filters while leaving other regions less parametrized which is consistent with our intuition. We can also see the proposed error metric on the validation set keeps decreasing throughout the training process except for the starting phase of each operation, which is the effects of random initialization, suggesting an effective knowledge transfer.

4.1.2. Two oscillatory modes with one drifting. For the second example, we also consider two nonlinear modes driven by different sinusoidal temporal modes. However, only one of the nonlinear modes is purely spatial, the other one is modulated by time. This synthetic example is meant to replicate the effects of traveling

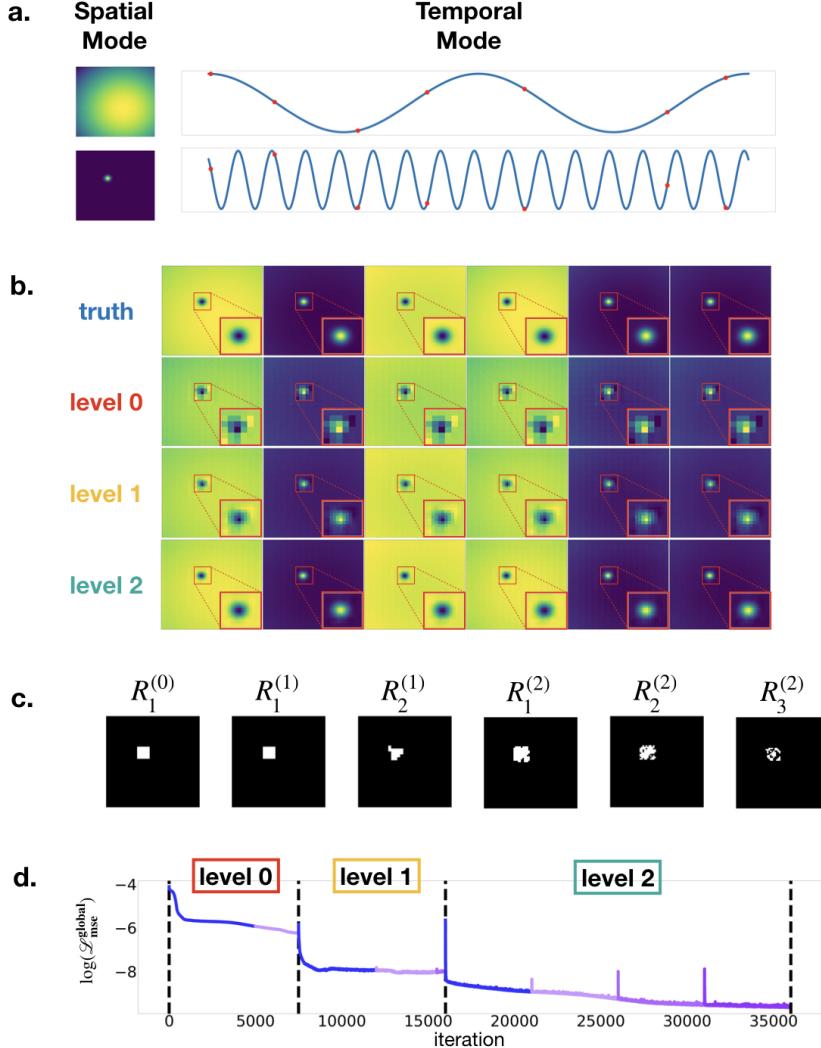


Fig. 5: **a.** Two separated spatial and temporal modes for the example of two oscillatory modes. Red dots correspond to the sampled test snapshots. **b.** Original data and the reconstructions of the sampled test snapshots across each level of the network. **c.** Regions that different groups of adaptive filters apply across different levels of the architecture. Here, $R_j^{(k)}$ represents the regions that the j^{th} group of adaptive filters being applied at level k . **d.** Logarithmic error plot (in terms of the metric presented in subsection 3.3) on the validation set throughout the training.

waves or the drifting dynamics that usually appears in the spatial-temporal data. Mathematically, we have:

$$(4.3) \quad \Phi(x, y, t) = u(x) \cos(\omega_0 t) + v(x, t) \cos(\omega_1 t + \frac{\pi}{4})$$

where

$$(4.4) \quad \begin{aligned} u(x, y) &= \cosh\left(\frac{x+1}{\sigma_0}\right) \cosh\left(\frac{y-1}{\sigma_0}\right) \\ v(x, y, t) &= \frac{1}{(2\pi\sigma_1)^{1/2}} \exp\left(-\frac{(x-3+0.5t)^2 + (y+3-0.5t)^2}{2\sigma_1^2}\right) \end{aligned}$$

We generate the data in the domain of $[-5, 5] \times [-5, 5]$ using 127×127 grids, and 500 snapshots are uniformly collected from the time interval $[0, 4\pi]$. We set $\omega_0 = 0.5$, $\omega_1 = 4.0$, $\sigma_0 = 10.0$, $\sigma_1 = 0.25$ and $v = 1.0$ in our experiment. [subsection 4.1.2](#) shows the two spatial-temporal modes, sampled test snapshots, reconstructions across different levels, refined regions and the logarithmic error plot.

In this example, we construct a network of 3 levels, with 3, 3 and 3 groups of adaptive filters applied on each level. Similarly, we see the network keeps refining itself by growing new pathways, our adaptive filters successfully capture the high-frequency contents drifting along the diagonal and knowledge learned in the previous level can be successfully transferred to the next level.

4.1.3. Channel flow. The channel flow dataset is acquired from the Johns Hopkins Turbulence Database [22, 29, 13]. Data was generated by solving the Navier-Stokes equations in a domain of size $8\pi \times 2 \times 3\pi$ using $2048 \times 512 \times 1536$ nodes. For more details, refer to [18]. We sample the slice of $z = 1.5\pi$ and down-sample the other two spatial dimensions both by a factor of 8, so that the spatial dimensions of our snapshots are 255×63 . 500 snapshots are collected with $\Delta t = 0.052$ between each snapshot.

In this example, we set up a network of 4 levels with each level having 0, 3, 3, 4 groups of adaptive filters. (At level 0, no adaptive filters are used because all pixels can be well-reconstructed with the deepening operation alone.) [subsection 4.1.3](#) shows the increasingly refined outputs at each level and the heat maps of regions that the adaptive filters apply. Notably, we see from the heat maps that the network spends more effort processing the regions near the boundary which intuitively makes sense: there are more detailed small-scale vortical behaviors arising within these regions. The plot of the decreasing error on the validation set also justifies the knowledge transfer of each operation.

4.1.4. Forced isotropic turbulence. The forced isotropic turbulence data set is also acquired from the Johns Hopkins Turbulence Database [22, 29, 44]. The data is generated from a direct numerical simulation of forced isotropic turbulence on a $1024 \times 1024 \times 1024$ periodic grid, using a pseudo-spectral parallel code. The range of spatial dimensions x, y and z are $[0, 2\pi]$. We sample the slice of $z = \pi$ and down-sample the other two spatial dimensions both by a factor of 8 so that each snapshot is of size 127×127 . 503 snapshots are collected with $\Delta t = 0.02$ between each snapshot. This is an example without a clear separation of scales, although it still possesses rich microscopic dynamics across all regions.

For this example, we set up a network of 4 levels with each level having 0, 4, 4, 4 groups of adaptive filters. (At level 0, no adaptive filters are used because all pixels can be well-reconstructed with the deepening operation alone.) As shown in [subsection 4.1.4](#), the network still offers reasonably well and progressively refined reconstructions of sampled test snapshots and the error plot also shows the effectiveness in knowledge transfer. But the heat maps of refined regions do not exhibit clear spatial patterns due to the isotropic nature of the data. In this case, one should be very

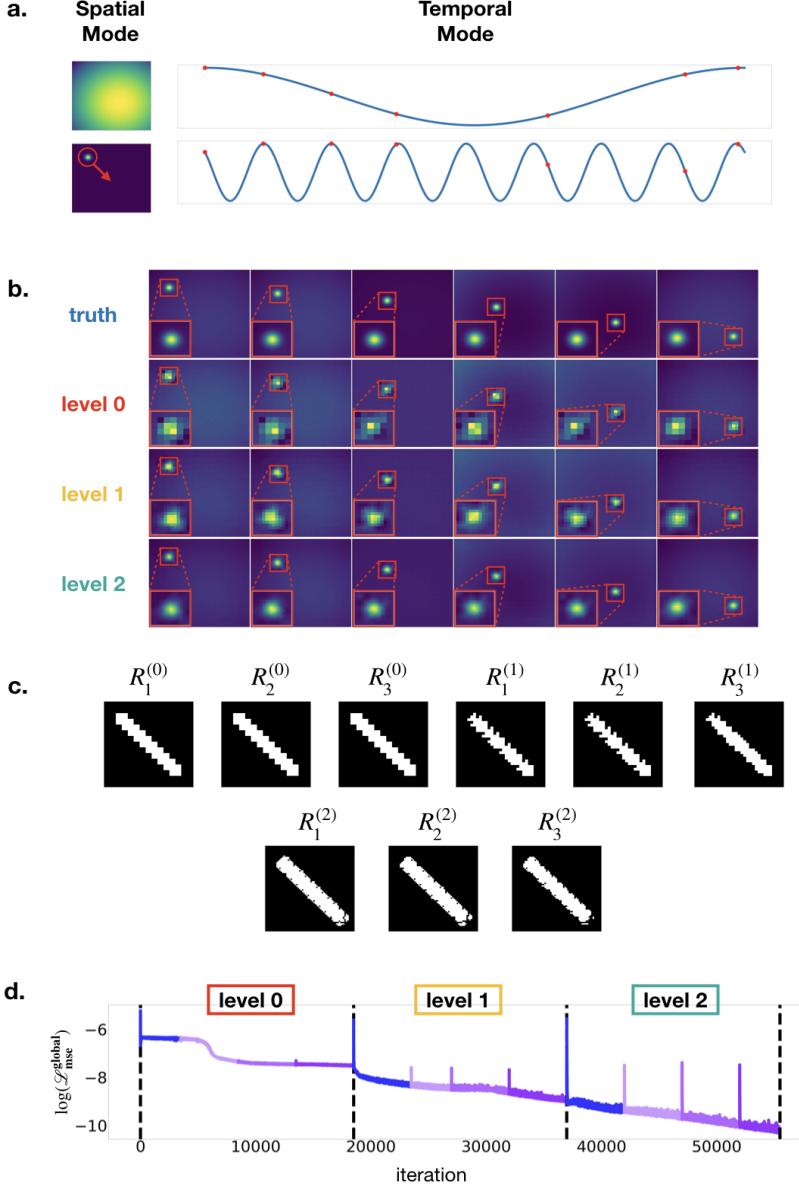


Fig. 6: From top to bottom: **a.** Two separated spatial and temporal modes for the example of two oscillatory modes with one drifting. Red dots correspond to the sampled test snapshots. **b.** Original data and the reconstructions of the sampled test snapshots across each level of the network. **c.** Regions that different groups of adaptive filters apply across different levels of the architecture. Here, $R_j^{(k)}$ represents the regions that the j^{th} group of adaptive filters being applied at level k . **d.** Logarithmic error plot (in terms of the metric presented in subsection 3.3) on the validation set throughout the training.

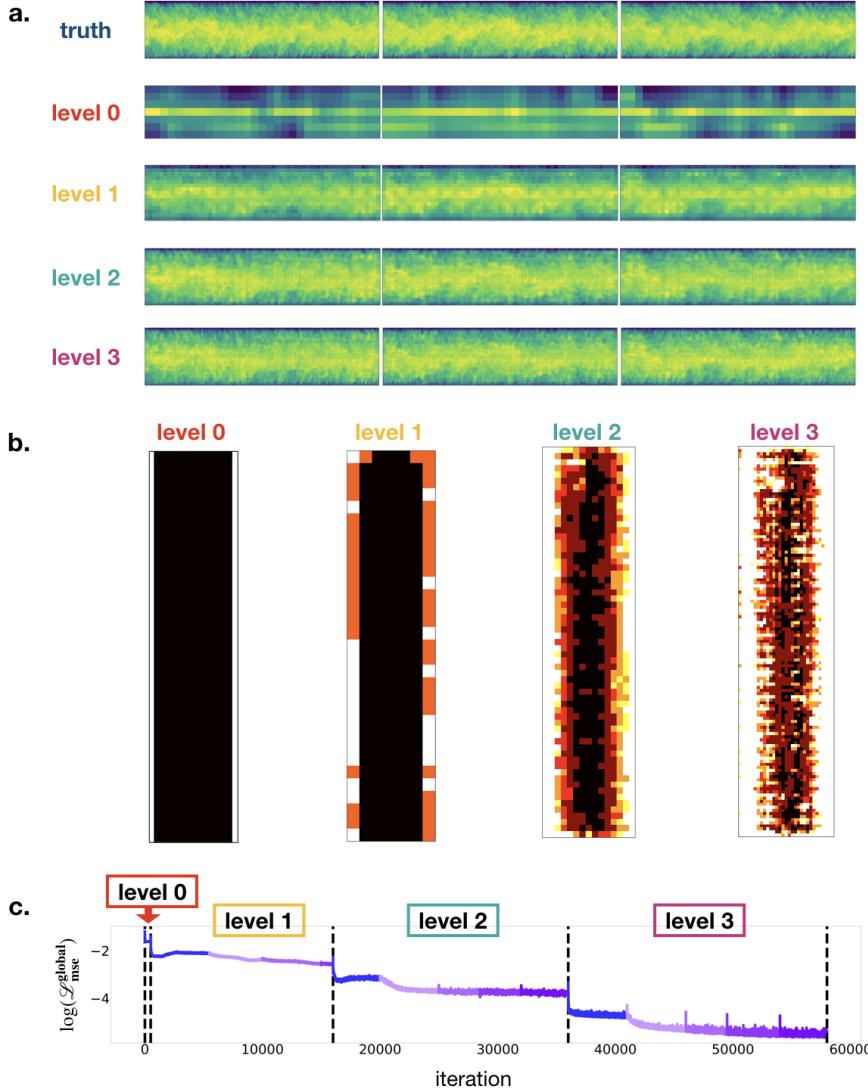


Fig. 7: **a.** Original data and reconstructions of the channel flow example across different levels of the network over the sampled test snapshots. **b.** Heat maps that reflect the regions that adaptive filters apply across different levels of the architecture. The brighter the pixel, the more filters are applied. **c.** Logarithmic error plot (in terms of the metric presented in subsection 3.3) on the validation set throughout the training.

cautious when attempting to use the hierarchical representations encoded in the network since neural networks are fundamentally interpolation methods [25] and lacking of interpretations often suggests the representations are hard to generalize.

4.1.5. Sea surface temperature. In this example, we consider the global sea surface temperature (SST) data. The NOAA OI SST V2 data set is open source and can be downloaded at <http://www.esrl.noaa.gov/psd/>. The data is collected each

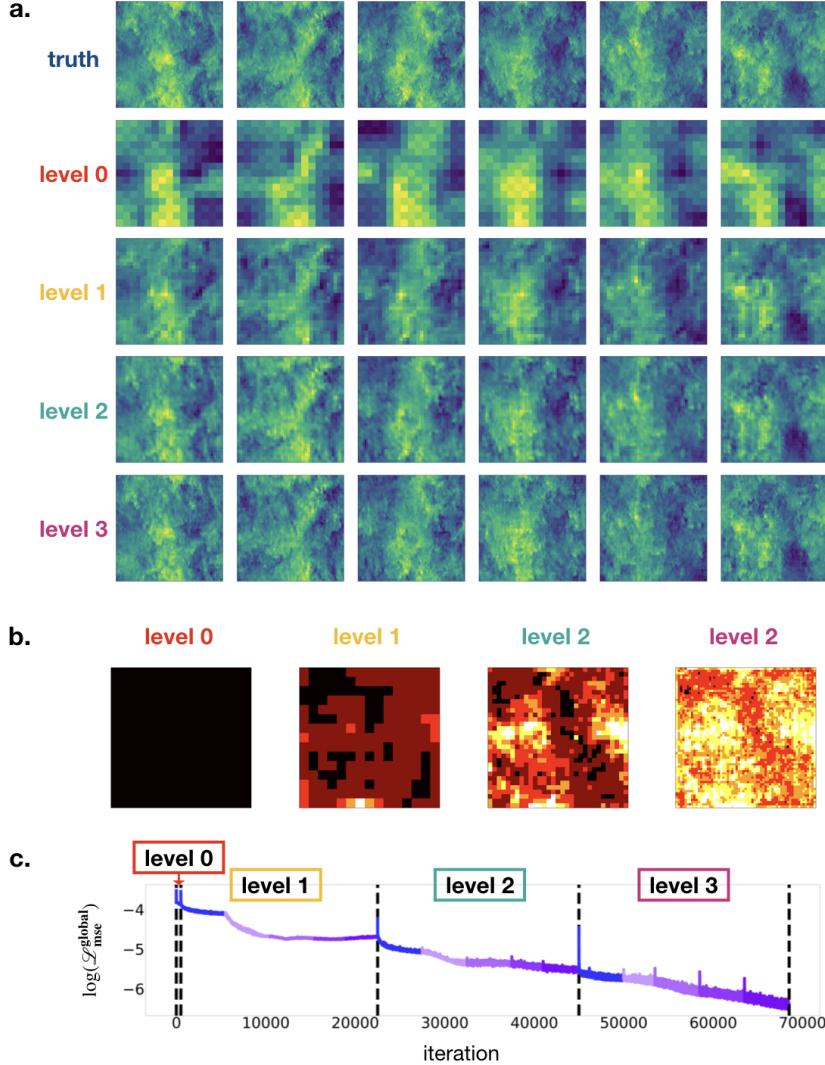


Fig. 8: **a.** Original data and the reconstructions of the forced isotropic turbulence example across different levels of the network over the sampled test snapshots. **b.** Heat maps that reflect the regions that adaptive filters apply across different levels of the architecture. The brighter the pixel, the more filters are used. **c.** Logarithmic error plot (in terms of the metric presented in subsection 3.3) on the validation set throughout the training.

month and spans a 20 year period from 1990 to 2010. In this example, we set up a network of 4 levels with each level having 2, 2, 4, 4 groups of adaptive filters. subsection 4.1.5 shows the increasingly refined reconstructions at higher and higher levels and the heat maps of regions where the adaptive filters apply. One observation is that network explores the middle regions more exhaustively. We conjecture that the underlying reason could be the temperature at the North Pole (top) and the South

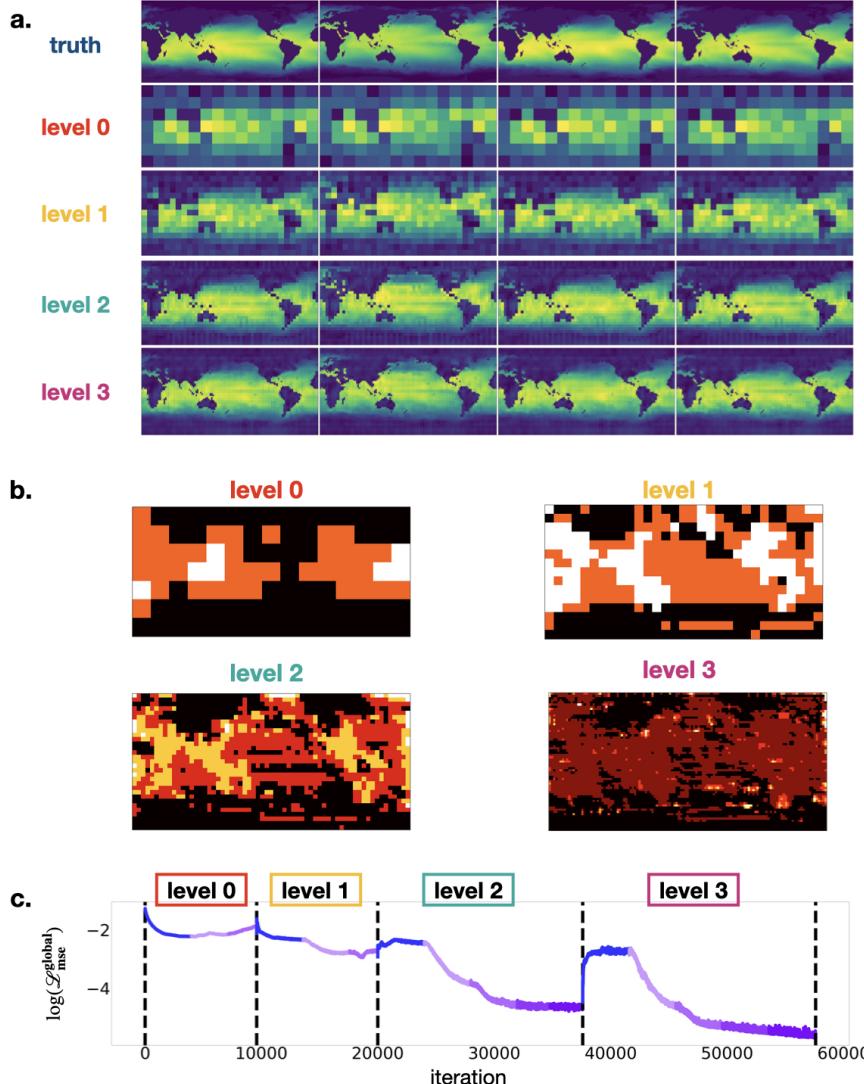


Fig. 9: **a.** Original data and reconstructions of the sea surface temperature example across different levels of the network over the sampled test snapshots. **b.** Heat maps that reflect the regions that adaptive filters apply across different levels of the architecture. The brighter the pixel, the more filters are used. **c.** Logarithmic error plot (in terms of the metric presented in subsection 3.3) on the validation set throughout the training.

Pole (bottom) are less variant. In addition, one can see the boundaries of the highly intensive regions align well with edges of the continents which suggests the transition from the land to the ocean.

4.2. Benchmarks. We benchmark MrCAE against the RED-Net, which is structurally similar to ours. We run experiments of the network with and without rectified nonlinear units (ReLU) which results in two candidate comparisons: RED-Net(Linear)

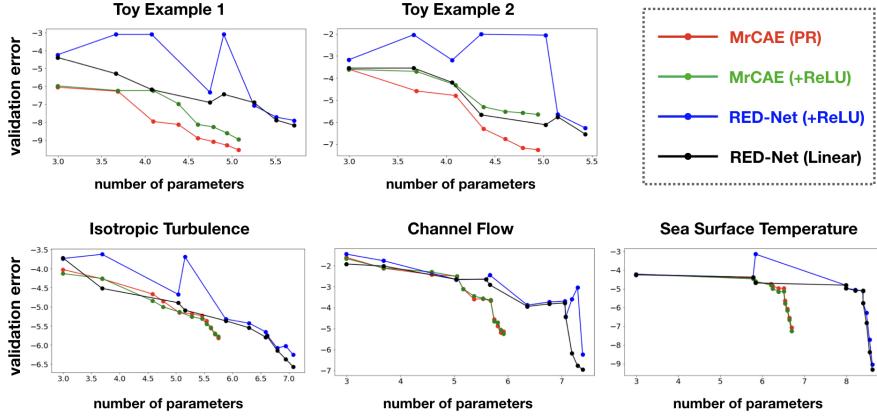


Fig. 10: **MrCAEs require less parameters.** Comparison of the validation error versus number of parameters for our method versus other state-of-the-art methods. Both the number of parameters and the validation error are on logarithmic scales. In this figure, each point represents a particular architecture: we progressively train the networks of all kinds so that their depths and widths increase along the curves. The number of parameters of MrCAEs scale better as the network expands while exhibiting a steady decrease in validation error.

and RED-Net(+ReLU). We also include our MrCAE without adaptive filters, but equipped with the rectified nonlinear units, which is named MrCAE(+ReLU).

Although RED-Nets were not primarily proposed to do progressive training with data of increasing spatio-temporal resolution, for fair comparison, we also use different resolution data to train them and use the metric in subsection 3.3 to evaluate their performance.

The procedure is as follows: throughout the training of our MrCAE network, we obtain a sequence of network architectures at different hierarchical levels. For each specific architecture, we train a corresponding one (with the same network depth and the same number of filters across each layer) for all other networks used for comparison. The key differences among the four candidate networks are as follows:

- **number of parameters:** MrCAEs have sparse connections which requires less parameters whereas connections in RED-Nets are dense. This is due to the fact that different group of filters are independently patched in MrCAEs.
- **size of encoding:** MrCAE(PR) explicitly utilizes a sparse coding scheme whereas the other three types of networks don't have this feature.

Results are shown in subsection 4.2.1 and subsection 4.2.2. Axes in both plots are in logarithmic scales.

4.2.1. Number of parameters. In subsection 4.2.1, we show that MrCAEs use significantly fewer parameters in comparison with similar architectures of RED-Nets. In other words, due to the sparse connections, it scales better when the network grows. Moreover, the error curves for MrCAEs seem to steadily decrease as the network grows which is highly desirable. But for RED-Nets, the curves are quite variable — more parameters does not necessarily offer better results. And despite the parsimonious use of parameters, MrCAEs can achieve reasonable accuracy and sometimes, in the

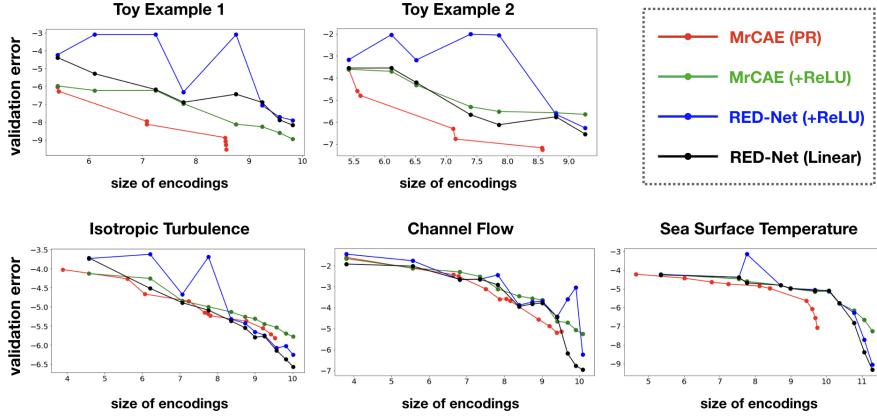


Fig. 11: **MrCAE(PR)** offers better compression ratio. Comparison of the validation error versus size of encoding for our method versus other state-of-the-art methods. Both the size of encoding and the validation error are on logarithmic scales. MrCAE provides smaller sizes of encoding over similar architectures due to the exploitation of spatial patterns in data sets.

toy examples where the dynamics are not that complex, outperform other networks. We conjecture the reasons are of two folds: First, this spatio-temporal data has a clear separation of scales and our adaptive filters are primarily designed to efficiently leverage this feature. And second, the progressive training process may have some guiding effects for the network parameters which lowers the risks of getting stuck in some local minima.

4.2.2. Size of encoding. In subsection 4.2.2, we show the validation error versus the size of encoding. By utilizing the adaptive filters, we are able to hierarchically encode different regions: some filters are only applied to highly variable regions. Therefore in practice, with similar architecture, MrCAE(PR) can always generate a smaller encoding size which leads to a better compression ratio. (It should be noted that unlike other types of networks, in principle, we should count not only the encoded information but also the corresponding positions which is similar to storing sparse matrices when we calculate size of encoding for MrCAE. However, since the position information is shared across all the snapshots and in practice there would be a large number of snapshots, the size of the position information becomes negligible, so we don't include it in our calculation.) In the isotropic turbulence example, this gain is less obvious because of its isotropic nature: there's very little structured spatio-temporal data in this application.

5. Discussion. In summary, we have equipped the highly-successful convolutional autoencoder (CAE) with both adaptive filters and multiscale modeling capabilities, giving rise to a flexible workflow which allows for improved interpretability, control of the modeling framework, and in-depth understanding over modern end-to-end architectures. Our proposed MrCAE architecture integrates and leverages three highly successful mathematical architectures: (i) multigrid methods, (ii) convolutional autoencoders and (iii) transfer learning. Instead of training an extremely large black

box model end-to-end, our model progressively utilizes a refinement strategy to build a hierarchical structure which leverages increased data due to improved spatial resolution. As a consequence, it enables automatic data augmentation across different spatial resolutions, and outputs insightful intermediate models and results for distinct spatial scales. These intermediate models can either guide the next-level architecture design or be put into immediate use which may be favored by online algorithms.

In addition, we develop a masking mechanism for the use of adaptive convolutional filters which resembles the mesh refinement process. It fully exploits the spatial patterns (which often shows up in many interested spatio-temporal physical systems) in the data set and therefore the size of encoding is small, the use of parameters is parsimonious in nature and the interpretation for each convolutional filter is clear. Though it is highly effective in resolving reconstruction errors, cautions must be taken when there is no clear spatio-temporal separation of scales in the data set.

There are many ongoing challenges and promising directions that motivate future works. Our network is currently very shallow. However, many successful applications achieved by neural networks rely heavily on its depth, as these architectures are reported to create more abstractions and therefore finding more efficient representations [12]. Moreover, our original objective for building up this framework is to do multi-scale forecasting. By adopting such an encoder-decoder architecture, now we are able to build different forecasting tools for different parts of the hierarchical representations. Then, through the decoder, we are able to map it back to the original space. It will be very interesting to see how different forecasting algorithms will be built into the network and how they will affect the encoder-decoder architecture.

Appendix A. Remarks on filter initialization. As we perform the deepening operation or widening operation, we are essentially adding properly initialized convolutional and deconvolutional filters to the existing network.

For a deepening operation, one convolutional filter C_0 and one deconvolutional filter D_0 are added to the network to connect the inputs and outputs from two adjacent levels. Both of the two filters are of kernel size 3×3 and applied with stride size 2 without padding, and they are initialized with the following matrices:

$$(A.1) \quad C_0 = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} + \begin{bmatrix} \epsilon_{11}^{(c)} & \epsilon_{12}^{(c)} & \epsilon_{13}^{(c)} \\ \epsilon_{21}^{(c)} & \epsilon_{22}^{(c)} & \epsilon_{23}^{(c)} \\ \epsilon_{31}^{(c)} & \epsilon_{32}^{(c)} & \epsilon_{33}^{(c)} \end{bmatrix}$$

$$(A.2) \quad D_0 = \begin{bmatrix} \frac{1}{4} & \frac{1}{2} & \frac{1}{4} \\ \frac{1}{2} & 1 & \frac{1}{2} \\ \frac{1}{4} & \frac{1}{2} & \frac{1}{4} \end{bmatrix} + \begin{bmatrix} \epsilon_{11}^{(d)} & \epsilon_{12}^{(d)} & \epsilon_{13}^{(d)} \\ \epsilon_{21}^{(d)} & \epsilon_{22}^{(d)} & \epsilon_{23}^{(d)} \\ \epsilon_{31}^{(d)} & \epsilon_{32}^{(d)} & \epsilon_{33}^{(d)} \end{bmatrix}$$

In a nutshell, C_0 is initialized with a down-sampling operator while D_0 is initialized with an interpolation operator so that the transition is smooth between each level. We also add some uniformly distributed random noise ($\epsilon_{ij}^{(c)}, \epsilon_{ij}^{(d)} \in [-\epsilon, \epsilon]$) to break the symmetry.

To perform a widening operation, we create a group of convolutional filters and deconvolutional filters where the channel number is specified by the user. Also for the sake of ensuring a smooth transition, we initialize all entries of these filters with independent, uniformly distributed random variables from $[-\epsilon, \epsilon]$ so that immediately

after the change of the architecture, reconstructions won't be affected by much.

Appendix B. Remarks on training. In addition to specifying the tolerances (threshold for residuals on all pixels) for different training phases as shown in [Algorithm 3.1](#), we set a maximum number of epochs to stop the training. We also implement a stopping criterion as the convergence slows down: if the relative error decrease every ten epochs is less than 0.001, we stop the training and move to the next operation to expand the network capacity.

Acknowledgements. This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344 and was supported by the LLNL-LDRD Program under Project No. 19-ERD-019, LLNL-JRNL-808088-DRAFT. SLB acknowledges support from the Army Research Office (ARO W911NF-17-1 0306). JNK acknowledges support from the Air Force Office of Scientific Research (AFOSR) grant FA9550-17-1-0329.

This document was prepared as an account of work sponsored by an agency of the United States government. Neither the United States government nor Lawrence Livermore National Security, LLC, nor any of their employees makes any warranty, expressed or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States government or Lawrence Livermore National Security, LLC. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States government or Lawrence Livermore National Security, LLC, and shall not be used for advertising or product endorsement purposes.

REFERENCES

- [1] M. Z. ALOM, T. M. TAHA, C. YAKOPCIC, S. WESTBERG, P. SIDIKE, M. S. NASRIN, B. C. VAN ESESN, A. A. S. AWWAL, AND V. K. ASARI, *The history began from alexnet: A comprehensive survey on deep learning approaches*, arXiv preprint arXiv:1803.01164, (2018).
- [2] P. BENNER, S. GUGERCIN, AND K. WILLCOX, *A survey of projection-based model reduction methods for parametric dynamical systems*, SIAM review, 57 (2015), pp. 483–531.
- [3] M. J. BERGER AND P. COLELLA, *Local adaptive mesh refinement for shock hydrodynamics*, Journal of computational Physics, 82 (1989), pp. 64–84.
- [4] S. L. BRUNTON AND J. N. KUTZ, *Data-driven science and engineering: Machine learning, dynamical systems, and control*, Cambridge University Press, 2019.
- [5] S. L. BRUNTON, B. R. NOACK, AND P. KOUMOUTSAKOS, *Machine learning for fluid mechanics*, Annual Review of Fluid Mechanics, 52 (2020), pp. 477–508.
- [6] K. CHAMPION, B. LUSCH, J. N. KUTZ, AND S. L. BRUNTON, *Data-driven discovery of coordinates and governing equations*, arXiv preprint arXiv: 1904.02107, (2019).
- [7] T. CHEN, I. GOODFELLOW, AND J. SHLENS, *Net2net: Accelerating learning via knowledge transfer*, arXiv preprint arXiv:1511.05641, (2015).
- [8] G. FROYLAND, G. A. GOTTWALD, AND A. HAMMERLINDL, *A computational method to extract macroscopic variables and their dynamics in multiscale systems*, SIAM Journal on Applied Dynamical Systems, 13 (2014), pp. 1816–1846.
- [9] G. FROYLAND, G. A. GOTTWALD, AND A. HAMMERLINDL, *A trajectory-free framework for analysing multiscale systems*, Physica D: Nonlinear Phenomena, 328 (2016), pp. 34–43.
- [10] C. GIN, B. LUSCH, S. L. BRUNTON, AND J. N. KUTZ, *Deep learning models for global coordinate transformations that linearize PDEs*, arXiv preprint arXiv:1911.02710, (2019).
- [11] R. GONZALEZ-GARCIA, R. RICO-MARTINEZ, AND I. KEVREKIDIS, *Identification of distributed parameter systems: A neural net based approach*, Computers & chemical engineering, 22 (1998), pp. S965–S968.
- [12] I. GOODFELLOW, Y. BENGIO, AND A. COURVILLE, *Deep learning*, MIT press, 2016.
- [13] J. GRAHAM, K. KANOV, X. YANG, M. LEE, N. MALAYA, C. LALESCU, R. BURNS, G. EYINK, A. SZALAY, R. MOSER, ET AL., *A web services accessible database of turbulent channel flow and its use for testing a new integral wall model for LES*, Journal of Turbulence, 17 (2016), pp. 181–215.
- [14] J. HE AND J. XU, *Mgnet: A unified framework of multigrid and convolutional neural network*, Science china mathematics, 62 (2019), pp. 1331–1354.
- [15] K. HE, X. ZHANG, S. REN, AND J. SUN, *Deep residual learning for image recognition*, in The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), June 2016.
- [16] J. S. HESTHAVEN, G. ROZZA, B. STAMM, ET AL., *Certified reduced basis methods for parametrized partial differential equations*, vol. 590, Springer, 2016.
- [17] S. IOFFE AND C. SZEGEDY, *Batch normalization: Accelerating deep network training by reducing internal covariate shift*, arXiv preprint arXiv:1502.03167, (2015).
- [18] K. KANOV, R. BURNS, C. LALESCU, AND G. EYINK, *The johns hopkins turbulence databases: an open simulation laboratory for turbulence research*, Computing in Science & Engineering, 17 (2015), pp. 10–17.
- [19] I. G. KEVREKIDIS, C. W. GEAR, J. M. HYMAN, P. G. KEVREKIDIS, O. RUNBORG, C. THEODOROPOULOS, AND OTHERS, *Equation-free, coarse-grained multiscale computation: Enabling macroscopic simulators to perform system-level analysis*, Communications in Mathematical Sciences, 1 (2003), pp. 715–762.
- [20] A. KRIZHEVSKY, I. SUTSKEVER, AND G. E. HINTON, *Imagenet classification with deep convolutional neural networks*, in Advances in neural information processing systems, 2012, pp. 1097–1105.
- [21] J. N. KUTZ, *Data-driven modeling & scientific computation: methods for complex systems & big data*, Oxford University Press, 2013.
- [22] Y. LI, E. PERLMAN, M. WAN, Y. YANG, C. MENEVEAU, R. BURNS, S. CHEN, A. SZALAY, AND G. EYINK, *A public turbulence database cluster and applications to study Lagrangian evolution of velocity increments in turbulence*, Journal of Turbulence, (2008), p. N31.
- [23] C. LIU, B. ZOPH, M. NEUMANN, J. SHLENS, W. HUA, L.-J. LI, L. FEI-FEI, A. YUILLE, J. HUANG, AND K. MURPHY, *Progressive neural architecture search*, in Proceedings of the European Conference on Computer Vision (ECCV), 2018, pp. 19–34.
- [24] B. LUSCH, J. N. KUTZ, AND S. L. BRUNTON, *Deep learning for universal linear embeddings of nonlinear dynamics*, Nature communications, 9 (2018), p. 4950.
- [25] S. MALLAT, *Understanding deep convolutional networks*, Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences, 374 (2016),

- p. 20150203.
- [26] A. MARDT, L. PASQUALI, H. WU, AND F. NOÉ, *VAMPnets: Deep learning of molecular kinetics*, Nature Communications, 9 (2018).
 - [27] S. F. MCCORMICK, *Multigrid methods*, SIAM, 1987.
 - [28] X. PENG, R. S. FERIS, X. WANG, AND D. N. METAXAS, *Red-net: A recurrent encoder-decoder network for video-based face alignment*, International Journal of Computer Vision, 126 (2018), pp. 1103–1119.
 - [29] E. PERLMAN, R. BURNS, Y. LI, AND C. MENEVEAU, *Data exploration of turbulence simulations using a database cluster*, in Proceedings of the 2007 ACM/IEEE conference on Supercomputing, ACM, 2007, p. 23.
 - [30] A. QUARTERONI, G. ROZZA, ET AL., *Reduced order methods for modeling and computational reduction*, vol. 9, Springer, 2014.
 - [31] M. RAISSI, P. PERDIKARIS, AND G. KARNIADAKIS, *Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations*, Journal of Computational Physics, 378 (2019), pp. 686–707.
 - [32] M. RAISSI, A. YAZDANI, AND G. E. KARNIADAKIS, *Hidden fluid mechanics: Learning velocity and pressure fields from flow visualizations*, Science, 367 (2020), pp. 1026–1030.
 - [33] S. H. RUDY, J. N. KUTZ, AND S. L. BRUNTON, *Deep learning of dynamics and signal-noise decomposition with time-stepping constraints*, Journal of Computational Physics, 396 (2019), pp. 483–506.
 - [34] K. SIMONYAN AND A. ZISSERMAN, *Very deep convolutional networks for large-scale image recognition*, arXiv preprint arXiv:1409.1556, (2014).
 - [35] C. SZEGEDY, W. LIU, Y. JIA, P. SERMANET, S. REED, D. ANGUELOV, D. ERHAN, V. VANHOUCKE, AND A. RABINOVICH, *Going deeper with convolutions*, in Proceedings of the IEEE conference on computer vision and pattern recognition, 2015, pp. 1–9.
 - [36] K. TAIRA, S. L. BRUNTON, S. DAWSON, C. W. ROWLEY, T. COLONIUS, B. J. McKEON, O. T. SCHMIDT, S. GORDEYEV, V. THEOFILIS, AND L. S. ÜKEILEY, *Modal analysis of fluid flows: An overview*, AIAA Journal, 55 (2017), pp. 4013–4041.
 - [37] M. TAN AND Q. V. LE, *Efficientnet: Rethinking model scaling for convolutional neural networks*, arXiv preprint arXiv:1905.11946, (2019).
 - [38] U. TROTTERBERG, C. W. OOSTERLEE, AND A. SCHULLER, *Multigrid*, Elsevier, 2000.
 - [39] C. WEHMEYER AND F. NOÉ, *Time-lagged autoencoders: Deep learning of slow collective variables for molecular kinetics*, The Journal of Chemical Physics, 148 (2018), pp. 1–9.
 - [40] E. WEINAN, *Principles of multiscale modeling*, Cambridge University Press, 2011.
 - [41] E. WEINAN, B. ENGQUIST, AND OTHERS, *The heterogeneous multiscale methods*, Communications in Mathematical Sciences, 1 (2003), pp. 87–132.
 - [42] B. XU, N. WANG, T. CHEN, AND M. LI, *Empirical evaluation of rectified activations in convolutional network*, arXiv preprint arXiv:1505.00853, (2015).
 - [43] L. YANG, D. ZHANG, AND G. E. KARNIADAKIS, *Physics-informed generative adversarial networks for stochastic differential equations*, arXiv preprint arXiv:1811.02033, (2018).
 - [44] P. YEUNG, D. DONZIS, AND K. SREENIVASAN, *Dissipation, enstrophy and pressure statistics in turbulence simulations at high reynolds numbers*, Journal of Fluid Mechanics, 700 (2012), pp. 5–15.
 - [45] S. ZAGORUYKO AND N. KOMODAKIS, *Wide residual networks*, arXiv preprint arXiv:1605.07146, (2016).
 - [46] B. ZOPH AND Q. V. LE, *Neural architecture search with reinforcement learning*, arXiv preprint arXiv:1611.01578, (2016).