# SparseST: Exploiting Data Sparsity in Spatiotemporal Modeling and Prediction

Junfeng Wu, Hadjer Benmeziane, Kaoutar El Maghraoui, Liu Liu, Yinan Wang

arXiv:2511.14753v1 [cs.LG] 18 Nov 2025

*Abstract*—Spatiotemporal data mining (STDM) has a wide range of applications in various complex physical systems (CPS), i.e., transportation, manufacturing, healthcare, etc. Among all the proposed methods, the Convolutional Long Short-Term Memory (ConvLSTM) has proved to be generalizable and extendable in different applications and has multiple variants achieving state-of-the-art performance in various STDM applications. However, ConvLSTM and its variants are computationally expensive, which makes them inapplicable in edge devices with limited computational resources. With the emerging need for edge computing in CPS, efficient AI is essential to reduce the computational cost while preserving the model performance. Common methods of efficient AI are developed to reduce redundancy in model capacity (i.e., model pruning, compression, etc.). However, spatiotemporal data mining naturally requires extensive model capacity, as the embedded dependencies in spatiotemporal data are complex and hard to capture, which limits the model redundancy. Instead, there is a fairly high level of data and feature redundancy that introduces an unnecessary computational burden, which has been largely overlooked in existing research. Using the sequence of images as an example, (1) the informative pixels at each frame are usually spatially sparse (i.e., only foreground pixels contain meaningful information in action recognition) and (2) the informative features in the first-order difference of two adjacent images are possibly sparse with a fixed or slowly evolving background. Therefore, we developed a novel framework SparseST, that pioneered in exploiting data sparsity to develop an efficient spatiotemporal model. In addition, we explore and approximate the Pareto front between model performance and computational efficiency by designing a multi-objective composite loss function, which provides a practical guide for practitioners to adjust the model according to computational resource constraints and the performance requirements of downstream tasks.

*Note to Practitioners*—This paper is motivated by the need to reduce the computational cost on edge devices with limited computational resources. In various complex physical systems (CPS), spatiotemporal data mining is an important task and naturally requires extensive model capacity, as the embedded dependencies in spatiotemporal data are complex and hard to capture. Instead of model compression, we aim to trade the redundant information in the input data or features for computational efficiency. In addition, our framework provides practitioners with direct control over their preference for efficiency and model accuracy. And the approximated Pareto front of the trade-off between efficiency and accuracy gives the guideline to adjust the model according to computational resource constraints and the performance requirements of downstream tasks.

*Index Terms*—Spatiotemporal Data Mining, Efficient AI, Data and Feature Sparsity, Convolutional Long Short-term Memory, Multi-objective Pareto Front Learning

## I. INTRODUCTION

Spatiotemporal data mining (STDM) focuses on data with both spatial and temporal dimensions, capturing how features evolve over time and space, with a wide prevalence in real-world applications such as video analytics, traffic flow forecasting [1], and industrial process monitoring [2]. However, deep neural networks (DNN) for the application of STDM are computationally expensive in advanced manufacturing settings, including but not limited to process monitoring, quality control, and anomaly detection. With the increasing amount of data provided by different sensors in smart manufacturing systems, the size and complexity of DNN models are also increased, which means deploying DNN models is more and more computationally expensive. However, most edge computing or embedded devices in manufacturing systems only have limited memory, storage, and computational capacity, leading to the pressing need to improve the computational efficiency of DNN models while preserving their performance. To handle the issue mentioned above, efficient AI methods are developed to reduce the computational cost and redundancy in DNN models [3], [4].

We aim to develop efficient AI methods to accelerate the classical Convolutional Long Short-Term Memory (ConvLSTM) model for the downstream task of spatiotemporal prediction and anomaly detection. ConvLSTM was proposed by Shi et al. [5], which integrates convolutional operations into the LSTM architecture. Unlike traditional LSTM models that use fully connected layers in gate computation, ConvLSTM employs convolution to capture the spatial patterns of input data and hidden state features, effectively combining the strengths of both Convolutional Neural Network (CNN) and LSTM to better handle spatiotemporal data.

ConvLSTM models are computationally expensive, as they integrate convolutional operations into recurrent architectures. Applying existing model compression techniques to reduce model redundancy of ConvLSTM is nontrivial for several reasons: (1) Recurrent Neural Network (RNN) based model is susceptible to vanishing and exploding gradient [6], and aggressive pruning or quantization can further deteriorate this issue, making it more difficult to train the model; (2) low-rank factorization methods introduce high computational

Junfeng Wu and Yinan Wang are with the Department of Industrial and Systems Engineering, Rensselaer Polytechnic Institute, Troy, NY, 12180. E-mail: {wuj26, wangy88}@rpi.edu

Liu Liu is with the Department of Electrical, Computer, and Systems Engineering and the Department of Computer Science, Rensselaer Polytechnic Institute, Troy, NY, 12180. E-mail: liu.liu@rpi.edu

Hadjer Benmeziane and Kaoutar El Maghraoui are with the IBM Research. E-mail: Hadjer.Benmeziane@ibm.com, kelmaghr@us.ibm.com

Manuscript received XXX; revised YYY. (Corresponding Author: Yinan Wang)

overhead for decomposition and impose strict structural constraints on the factorized components, which can limit the representational capacity of the compressed model [7]. Rather than directly compressing the model, which can compromise training stability and limit representation capacity, we adopt a complementary perspective: reducing redundancy by exploiting the inherent sparsity of input data and intermediate features. This strategy improves efficiency while preserving the integrity of the ConvLSTM architecture.

We consider data with high spatial sparsity, such as images that contain only a small portion of informative pixels within a mostly black background. Such spatially sparse data can be efficiently processed by the 2D Sparse Convolution [8], which discards Multiply–Accumulate (MAC) operations involving zero elements. However, this method has two main drawbacks: (1) 2D Sparse Convolution can not exploit temporal sparsity in spatiotemporal data. For example, the first-order difference of pixel values between two adjacent image frames is often sparse due to slow temporal change; (2) 2D Sparse Convolution can cause a feature dilation effect, which means feature sparsity decreases as model depth increases [8], thereby reducing computational efficiency in deeper layers. To address the above limitations, we adopt the Delta Network (DN) Algorithm [9] to leverage temporal sparsity by applying a threshold to zero out feature values with a small first-order temporal difference. Therefore, it effectively leverages the temporal redundancy in spatiotemporal data and mitigate the problem of feature dilation in deeper layers.

In addition to the importance of efficiency in spatiotemporal modeling, there is a pressing need for practitioners to understand how to control the trade-off between efficiency and model performance, enabling them to adjust this balance according to computational resource constraints and performance requirements of downstream tasks. To address this, we formulate our problem as a multi-objective Pareto front learning task. Specifically, we design a multi-objective composite loss for two conflicting objectives: performance and efficiency, and train the neural network using the Smooth Tchebycheff Scalarization (STCH) method [10], which can effectively explore Pareto optimal points in both convex and non-convex regions of the Pareto front. Finally, we train a multi-task Gaussian Process (GP) to jointly model the dependency between model performance and efficiency. By approximating the Pareto front, we provide a guide for practitioners to fine-tune the preference weight and achieve a desired trade-off.

In this work, we propose SparseST, a framework designed as an efficient variant of ConvLSTM, which leverages the spatial and temporal sparsity present in raw input data and intermediate feature representations. In the meantime, we investigate the trade-off between model performance and efficiency by approximating the Pareto front. Our experimental results demonstrate notable computational savings with respect to (w.r.t.) Floating Point Operations (FLOPs) compared with baseline models on spatiotemporal prediction and anomaly detection tasks. The approximated Pareto front further shows the trade-off trend as expected. We summarize the contribution of this work as follows:

1) We propose SparseST, a framework that integrates 2D Sparse Convolution and the DN algorithm into the ConvLSTM architecture. By exploiting both data and feature sparsity, SparseST accelerates spatiotemporal modeling while preserving model performance.
2) We formulate the trade-off between model performance and efficiency as a multi-objective optimization problem. By approximating the Pareto front, our framework provides practical guidelines for practitioners to adjust preference weights based on computational resource constraints and downstream performance requirements.

The remainder of this paper is organized as follows. Section II reviews related literature on efficient AI methods and multi-objective Pareto front learning. Section III presents the preliminary background on the ConvLSTM model and the spatially sparse tensor. Section IV demonstrates the proposed SparseST from the concepts of Spatially Sparse Convolution and the DN Algorithm to the full model architecture, computational cost analysis, and the design of the multi-objective composite loss. Section V shows the experimental details, including dataset descriptions, hyperparameter settings, evaluation metrics, and results analysis for the two downstream tasks. Finally, Section VI concludes this work and discusses the limitations.

## II. LITERATURE REVIEW

This section reviews the literature on efficient AI from the perspectives of model compression and data or feature redundancy. Particularly, the related work on Multi-objective Pareto Front Learning is summarized, which is crucial to developing the SparseST framework. We also compare our method with the existing work to show the new methodological contribution.

### A. Efficient AI

Training a deep ConvLSTM model requires a massive number of MAC operations in each layer, both in forward and backward propagation. Model compression is a widely used method to reduce MAC operations for better efficiency, including pruning, quantization, low-rank factorization, etc. Another approach is to reduce MAC operations by exploiting the intrinsic spatial and temporal sparsity of sequential data.

*1) Model Compression:*

*a) Pruning:* Pruning eliminates those redundant or less important model parameters that do not contribute much to model performance, therefore reducing storage size, computational cost, and energy consumption without a substantial loss in model accuracy.

There are various pruning methods, such as weight pruning, neuron pruning, layer pruning, etc. Han et al. [11] proposed an iterative process of pruning and retraining. In each iteration, the model is first trained to learn informative weights. Then, unimportant weights with small magnitudes are pruned. Finally, the pruned model is retrained to fine-tune the pruned weights and regain model performance. However, the majority of parameters removed by this method are from the fully connected layers, where the computational cost is much lower than that of the convolutional layers [12]. To effectively compress CNN layers, Li et al. follow the same iterative

pruning and retraining process as [11] and proposed a new method to prune unimportant convolutional filters evaluated by the summation of their absolute weights [12]. Similarly, Liu et al. [13] proposed a network slimming method for pruning the entire channel of CNN by applying a sparsity-inducing $l_1$ regularizer to the scaling factors in batch normalization (**BN**) layers, effectively identifying and pruning less important channels.

Pruning is also explored in RNN models. Narang et al. [14] pruned the weights of RNNs during the model training and achieved a weight sparsity of $90\%$ with a small loss in accuracy by multiplying each weight matrix with a binary mask at each update step. At regular intervals, the masks are updated by a threshold function controlled by a set of hyperparameters. This technique does not require additional retraining steps used in [11]. However, it introduced an additional step of hyper-parameter tuning to control the threshold, which is difficult to implement in practice [14]. Wen et al. [15] explored the structured sparsity of LSTMs and proposed the concept of Intrinsic Sparse Structures (ISS). They novelly identified the basic structure inside RNNs, which can be considered as a group, thus making it possible to remove weights independently across all the basic structures within ISS.

*b) Quantization:* Quantization is a model compression technique to reduce the number of bits required for the representation and storage of model weights. For example, weights in 32-bit floating-point full precision representation can be quantized to a lower precision representation, such as 16-bit or 8-bit. Ott [16] designed multiple experiments to test different methods to quantize the weights of three major RNN types: vanilla RNN, GRU, and LSTM, showing for the first time how low-precision quantization of weights can be performed during RNN training. Besides weight quantization, gradient, and neuron activation can also be quantized to reduce the network size and memory requirement during model training and inference. Wei et al. [17] introduce TernGrad, a method for gradient quantization that reduces the communication cost in distributed deep learning by quantizing gradients to ternary values $(-1, 0, 1)$. This approach significantly reduces the amount of data that needs to be communicated between workers and the parameter server during training while maintaining high accuracy. Hubara et al. [18] quantized both weights and neuron activations of vanilla LSTM and CNN models, achieving low-precision training and inference without significant loss in accuracy.

*c) Low-rank factorization:* The low-rank factorization method aims to use matrix or tensor decomposition to estimate the informative parameters of the DNNs. For example, A matrix $W$ of size $m \times n$ can be approximated by the product of two smaller matrices $A$ and $B$, where $A$ is of size $m \times k$ and $B$ is of size $k \times n$, with $k < \min(m, n)$. This reduces the number of parameters from $mn$ to $mk + kn$. Common factorization techniques for 2D matrices include Singular Value Decomposition (SVD), Principal Component Analysis (PCA), and their variant [3] [19] [20]. For high-dimensional tensors in the case of convolutional filters, factorization methods such as Canonical Polyadic (CP) decomposition, Batch Normalization

(BN) low-rank decomposition, etc. [21] [22]. Wang et al. [23] proposed a CPAC-Conv layer by adopting CP-decomposition to compress convolutional kernels, with formulas for both forward and backward propagations derived. The value of the decomposed kernels indicates significant feature maps, which are informative for feature selection. Lu et al. [24] first undertook a systematic study and investigated redundancies in recurrent architectures such as RNN and LSTM. They compared the compression performance for different low-rank factorization methods [25] [26] and found that a hybrid strategy of using structured matrices in the bottom layers and shared low-rank factors on the top layers is particularly effective [24]. Later advancement shows various tensor decomposition-based compression methods, including tensor train [27], tensor ring [28], and block-term [29], can bring several orders of magnitude fewer parameters for large-scale RNNs while still maintaining high classification/prediction performance. However, methods mentioned above suffer from high computational overhead and limited representation ability of compressed RNN models due to extra flatten or permutation operations and strict constraints on either the shapes or the combination manners of the low-rank tensor components [7].

Although the model compression methods mentioned above have achieved prominent compression ratios, they only reduced the computational cost from the perspective of model redundancy, but didn't exploit the data or feature redundancy. Besides, applying existing model compression techniques to ConvLSTM is problematic due to issues such as unstable gradients, dimension mismatch, extensive retraining, degraded model capacity, etc. The following methods, which leverage the spatiotemporal sparsity of input data and intermediate features, will not suffer from the above issues, yet are still efficient in reducing the computational costs without sacrificing the model performance.

*2) Data and Feature Sparsity:*

*a) Spatially Sparse Convolution:* Spatially Sparse Convolution is a convolution operation optimized for processing spatially sparse data, for example, 3D point cloud data obtained from a LiDAR scanner, images of handwritten numbers in a sparse background with mostly zero pixel values, etc. Note that Spatially Sparse Convolution is different from the conventional dense CNN with weight pruning or compression. The word "sparse" only means sparsely distributed informative entries in data tensors.

Graham et al. [30] first introduced 2D Sparse Convolution for handwriting recognition and image classification tasks. A handwritten character can be represented as a sparse matrix, in which only pixels indicating the character are informative. Instead of conducting convolution operations over the entire sparse matrix (i.e., conventional dense CNN), 2D Sparse Convolution only conducts convolution operations over nonzero features (i.e., pixels indicating characters) and skips those pixels that only represent backgrounds (i.e., usually are zeros). Therefore, 2D Sparse Convolution significantly reduces computational cost by exploiting the sparsity in matrix-like data. Afterward, a sparse 3D CNN is explored to efficiently process 3D point cloud data [31] [32]. However, the current Sparse Convolution still dilates the number of nonzero values

over layers of convolution operations, which reduces the level of sparsity of feature matrices over the layers of the neural network.

In [33] [8], Graham et al. introduced Submanifold Sparse Convolution on a 3D point cloud data segmentation task. This method only conducts the convolution with the corresponding receptive field centered at a nonzero element, thereby keeping the same sparsity pattern throughout the layers of the network without dilating the feature maps. Later, Yan [34] adopted Submanifold Sparse Convolution [8] for the 3D object detection task and optimized the data flow to make it more GPU-friendly compared with [8]. In the more recent work by [35] and [36], Tang et al. propose the TorchSparse library for processing 3D point cloud data in the application of autonomous driving, with optimized computational regularity and data movement compared to previous work.

In summary, although Spatially Sparse Convolution is well studied, especially in 3D object detection and segmentation tasks, to the best of our knowledge, there are still significant gaps to directly apply it to STDM, which includes (1) Submanifold Sparse Convolution is inherently inapplicable to STDM since it keeps the same sparsity pattern across layers, which hinders the model to learn about temporal dynamics of spatiotemporal data; (2) The dilation issue of Sparse Convolution results in degraded computational efficiency in deeper layers; (3) Sparse Convolution can only efficiently process spatially sparse data (with a sparse ground), not applicable to spatiotemporal data with a fixed but not sparse background. To mitigate the issues mentioned above, we resort to the DN Algorithm.

*b) Delta Network Algorithm:* The DN Algorithm, first proposed by Neil et al. [9], is a computational reuse method to induce temporal sparsity of sequential data and intermediate features in the RNN model. It sets a delta threshold to zero out elements below that threshold, resulting in sparse delta vectors. Replacing the dense state vectors with sparse delta vectors, which contain the temporal difference of the states between two adjacent time steps, can reduce both computational costs and required memories [37]. Gao et al. [38] first adopted the DN Algorithm to accelerate the Gated Recurrent Unit (GRU) on a specifically designed FPGA hardware. The following work by Gao et al. [39] further reduced the memory requirement of edge computing devices by storing RNN parameters in off-chip memory and explored the trade-off of model performance and sparsity level under different delta threshold settings. However, this work only exploited temporal data sparsity. Later, Gao et al. [37] proposed Spartus, a hardware accelerator for LSTM models using a structured pruning method with the DN algorithm, making it possible to implement the LSTM model on the edge device with small memory for real-time online speech recognition tasks.

The existing work of applying the DN Algorithm to RNN models is about the co-design of hardware and algorithms with applications on speech recognition and audio digit recognition tasks. However, it has never been used to explore the temporal correlation of spatiotemporal data. More importantly, the DN algorithm is a perfect match for Sparse Convolution to mitigate the dilation issue and relax the constraint of spatially

sparse data (with a sparse background) to a fixed and dense background.

### B. Multi-objective Pareto Front Learning

Learning the Pareto front is a fundamental task in multi-objective optimization, which aims to recover the full set of all Pareto optimal solutions illustrating the trade-off between conflicting objectives. Optimizing the linear scalarization of multiple objectives (also referred to as the weighted sum method) is a widely used technique to find Pareto optimal solutions. The expression of linear scalarization is shown as follows.

$$\min_{\boldsymbol{x} \in \mathcal{X}} \sum_{i=1}^{m} w_i f_i(\boldsymbol{x}), \tag{1}$$

where $w_i$ is the preference weight for the $i^{\text{th}}$ objective $f_i$.

It converts the multi-objective optimization problem into a single-objective one by weighting the objectives using the preference vector. However, linear scalarization fails to reach Pareto optimal points on the non-convex part of the Pareto front [40]. Therefore, Tchebycheff Scalarization (TCH) [41], [42] is developed for scalarization to find Pareto optimal points for both convex and non-convex parts of the Pareto front, i.e.

$$\min_{\boldsymbol{x} \in \mathcal{X}} \max_{1 \le i \le m} \left\{ w_i \left( f_i(\boldsymbol{x}) - z_i^* \right) \right\} \tag{2}$$

where $m$ is the number of objectives, $z_i^*$ is the optimal value for the $i$ th objective. There is also a promising necessary and sufficient condition for TCH scalarization to find all Pareto optimal solutions [43]. Despite its long-standing recognition in the multi-objective optimization community, TCH is rarely used for gradient-based optimization due to the nonsmoothness of the $max$ function in Equation 2. Consequently, even when all objective functions $f_i$ are differentiable, the TCH formulation remains non-differentiable, making it infeasible for gradient-based optimization [10]. For example, Parallel Efficient Global Optimization (ParEGO) [44] is one of the multi-objective Bayesian optimization (MOBO) methods that employs Tchebycheff scalarization to approximate the Pareto front, which cannot be directly optimized by gradient-based methods. Therefore, Lin et al. [10], [45] proposed a smoothed variant of Tchebycheff scalarization (STCH) for gradient-based multi-objective optimization, making it well-suited for integration with the training of modern neural networks using gradient-based optimizers such as Adam [46].

Besides scalarization, there are also Pareto-based methods, such as evolutionary algorithms (EA) [47], [48], hypervolume-maximization approaches [49], [50], and Pareto-based MOBO [51], which produce a set of Pareto-optimal points in each run rather than a single point as in scalarization methods. These approaches also avoid the challenging task of selecting an appropriate preference vector for each objective to reflect the true user preference. However, since preference vectors are absent in the formulation of Pareto-based methods, there is no direct control over the user preference, potentially leading to wasted computation on regions of little interest on the Pareto front.

TABLE I: Comparison with Existing Work

| Existing work | Architecture | Type of Sparsity | Limitations | SparseST Improvement (ours) |
|---|---|---|---|---|
| Sparse Convolution [8] | CNN | Spatial only | – Active sites dilation issue <br> – No temporal sparsity | – DN algorithm to mitigate dilation <br> – Exploits temporal correlation |
| DN Algorithm [9] | RNN/GRU | Temporal only | – Fixed delta threshold $\Theta$ <br> – No spatial sparsity | – Learnable delta threshold $\Theta$ <br> – Leverages sparse/fixed background |
| Spartus [37] | LSTM on FPGA | Weight pruning + temporal | – Speech/audio only <br> – Compresses model weights <br> – No spatial data sparsity | – Applied to video prediction/anomaly detection <br> – No capacity loss from compression <br> – Tunable performance–efficiency trade-off |

Considering the benefits of STCH, which allows gradient-based optimization and provides direct control over the user preference, we adopt STCH for our multi-objective optimization framework. To approximate the full Pareto front, we aim to find the optimal model parameter $\boldsymbol{\theta}^*$ such that $\boldsymbol{x}^*(\boldsymbol{w}) = h_{\boldsymbol{\theta}^*}(\boldsymbol{w})$ is the corresponding Pareto optimal solution for any given preference vector $\boldsymbol{w}$, where $h$ is the surrogate model that maps any valid preference vector $\boldsymbol{w}$ to its corresponding Pareto optimal solution. We describe in detail the surrogate model used in Section IV-E.

### C. Comparison with Existing Work

In addition to the research gaps mentioned at the end of each subsection above, we want to further articulate our work as a new methodological contribution rather than an incremental improvement. See table I for a summary of the comparison with the existing work.

## III. PRELIMINARY BACKGROUND

In this section, we review the basics of the ConvLSTM network and introduce the concepts of the spatially sparse tensor.

### A. Convolutional LSTM network

A ConvLSTM unit is composed of an input gate $i_t$, a forget gate $f_t$ and an output gate $o_t$, indexed by time step $t$. The spatial feature of the input tensor at the current time step $\mathcal{X}_t$ and hidden state $\mathcal{H}_{t-1}$ from the last time step are firstly extracted by the convolution filters with their corresponding weights. Each of the three gates $i_t$, $f_t$, $o_t$ at each time step, with a range of values between 0 and 1, acts as a valve to control how much information should be passed through. The update equations for a single LSTM unit are as follows [5]:

$$
\begin{aligned}
i_t &= \sigma \left( \mathcal{W}_{xi} * \mathcal{X}_t + \mathcal{W}_{hi} * \mathcal{H}_{t-1} \right) \\
f_t &= \sigma \left( \mathcal{W}_{xf} * \mathcal{X}_t + \mathcal{W}_{hf} * \mathcal{H}_{t-1} \right) \\
o_t &= \sigma \left( \mathcal{W}_{xo} * \mathcal{X}_t + \mathcal{W}_{ho} * \mathcal{H}_{t-1} \right) \\
\tilde{\mathcal{C}}_t &= \tanh \left( \mathcal{W}_{xc} * \mathcal{X}_t + \mathcal{W}_{hc} * \mathcal{H}_{t-1} \right) \\
\mathcal{C}_t &= f_t \odot \mathcal{C}_{t-1} + i_t \odot \tilde{\mathcal{C}}_t \\
\mathcal{H}_t &= o_t \odot \tanh \left( \mathcal{C}_t \right)
\end{aligned}
\tag{3}
$$

where $\mathcal{W}$ denotes the weight matrices of convolutional filters, and $\sigma$, $\tanh$, $\odot$, $*$ represent the Sigmoid activation function, Hyperbolic Tangent activation function, element-wise multiplication, and convolution operator, respectively.
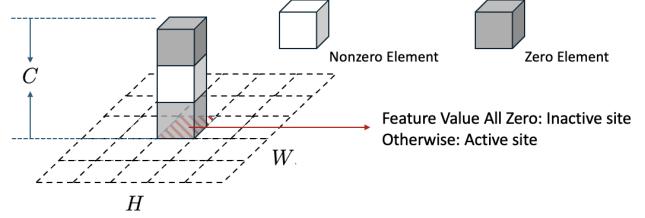


Fig. 1: Example of an active site. Each spatial location on the grid is an inactive site if all the feature values of the column vector located at this site are zero. Otherwise, it is defined as an active site [8].

### B. Spatially Sparse Tensor

Following the notation in [8], an input tensor has $(d+1)$ dimensions, where $d$ is the spatial dimension with an additional feature dimension. Figure 1 illustrates an example of an active site on a grid with spatial dimension $H$ and $W$, in which case $d = 2$.

A spatially sparse tensor denotes a tensor with spatial locations that are mostly inactive sites, which can be represented by two matrices: a coordinate matrix to store the spatial locations of all active sites, and a feature matrix to store the corresponding feature values. For input tensor $\mathcal{T} \in \mathbb{R}^{C_{in} \times H \times W}$ with the spatial dimension $H \times W$ and feature dimension $C_{in}$, we can represent it as follows:

$$
\mathcal{L} = \begin{bmatrix} h_1 & w_1 \\ \vdots & \vdots \\ h_N & w_N \end{bmatrix}, \mathcal{F} = \begin{bmatrix} F_1^T \\ \vdots \\ F_N^T \end{bmatrix}
\tag{4}
$$

where $N$ is the total number of active sites, and $F_i \in \mathbb{R}^{C_{in}}$ is the column feature vector for $i$-th active site.

## IV. METHODOLOGY: SPARSEST FRAMEWORK

In this section, we first introduce the preliminary knowledge on 2D sparse convolution and the delta network algorithm. Then, the unit structure and full model architecture of our proposed SparseST are introduced to demonstrate how 2D Sparse Convolution can be integrated with the delta algorithm to leverage both the spatial and temporal sparsities from the input data. In addition, we give the computational cost analysis of the proposed model. Finally, to explicitly control the trade-off between model performance and efficiency, we explain
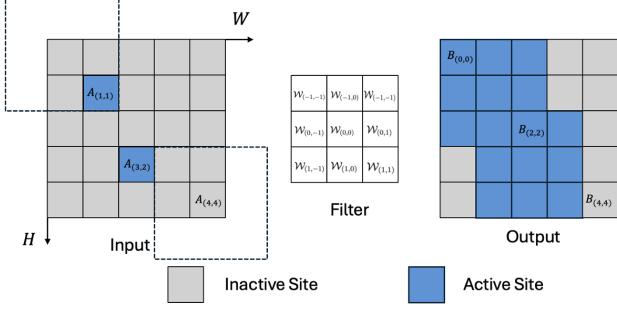
Fig. 2: An example of 2D Sparse Convolution. Active and inactive sites in the input and output feature maps are represented as blue and gray blocks, respectively. Input feature map is zero-padded to keep the spatial size unchanged. When the receptive field of the convolutional filter (dotted frame) is at the top left corner, there is only one active site $A_{(1,1)}$, yielding the corresponding output $B_{(0,0)} = \mathcal{W}_{(1,1)} \times A_{(1,1)}$, discarding MAC computations of all other inactive sites. We can finally get the output feature map following the same principle as described above.

how to formulate the multi-objective optimization problem and design the composite loss function.

### A. 2D Sparse Convolution

In [33] [8], two types of Spatially Sparse Convolution are introduced and compared: Sparse Convolution and Sub-manifold Sparse Convolution. The latter preserves the same level of sparsity by preventing active sites in feature maps from expanding during forward propagation across network layers. However, this constraint limits our framework to extrapolate new active sites. Therefore, this type of convolution is inherently inapplicable to spatiotemporal data mining. To avoid this limitation, we adopt the 2D Sparse Convolution in our framework, which discards the convolution calculation of inactive sites and reduces the computational cost compared with the dense convolution operation [33] [8].

For input tensor $\mathcal{T} \in \mathbb{R}^{C_{in} \times H \times W}$ in the form of Equation (4), the output tensor $\mathcal{Y}$ at spatial location $(u, v)$ can be represented as:

$$
\mathcal{Y}_{u,v} = \begin{cases} \sum_{(u+i,v+j)\in R} \mathcal{W}_\delta^\top F_{u+i,v+j} & \text{if } (u+i, v+j) \in \mathcal{L}. \\ 0 & \text{otherwise} \end{cases}
$$
(5)

where $F_{u+i,v+j} \in \mathbb{R}^{C_{in}}$ is the column feature vector at input spatial location $(u + i, v + j)$, $\mathcal{W}_\delta^\top \in \mathbb{R}^{C_{out} \times C_{in}}$ is the transpose of corresponding filter weight partition and $R$ is a receptive field coordinate set.

Figure 2 shows an example of 2D Sparse Convolution. It is noteworthy that the number of active sites colored in blue in the output is getting larger compared with that of the input, which means 2D Sparse Convolution will dilate the number of active sites, therefore reducing the computational efficiency in the following layers. To mitigate the dilation issue of 2D

Sparse Convolution, as well as to leverage the data or feature sparsity introduced by temporal correlation, we introduce the general formulation of the DN Algorithm in the next part, followed by our proposed framework, SparseST.

### B. Delta Network Algorithm

Consider a multiplication of a weight matrix $\mathbf{W}$ with an input vector sequence $\mathbf{X} = \{\mathbf{x_1}, \mathbf{x_2}, ..., \mathbf{x_T}\}$ where the component vector $\mathbf{x_t}$ is indexed by time step $\mathbf{t}$, the output vector sequence $\mathbf{Y} = \{\mathbf{y_1}, \mathbf{y_2}, ..., \mathbf{y_T}\}$ can be calculated in a recursive way:

$$
\begin{aligned}
\mathbf{y_0} &= \mathbf{0}, \mathbf{x_0} = \mathbf{0} \\
\Delta \mathbf{x_t} &= \mathbf{x_t} - \mathbf{x_{t-1}} \\
\mathbf{y_t} &= \mathbf{W}\mathbf{x_t} = \mathbf{W}\Delta\mathbf{x_t} + \mathbf{y_{t-1}}
\end{aligned}
$$
(6)

where the delta vector $\Delta \mathbf{x_t}$ is the difference of two input sequence elements from adjacent time steps and $\mathbf{y_{t-1}}$ is the multiplication result from the last time step, which is stored in the delta memory.

In the delta vector update Equations (7) [9], we set a delta threshold $\Theta$ to zero out the below-threshold elements in the delta vector $\Delta \mathbf{x_t}$. Therefore, the multiplication of entire columns in a weight matrix with zero elements in the delta vector can be skipped.

$$
\begin{aligned}
\hat{x}_{i,t-1} &= \begin{cases} x_{i,t-1} & \text{if } |x_{i,t} - \hat{x}_{i,t-1}| > \Theta \\ \hat{x}_{i,t-2} & \text{otherwise} \end{cases} \\
\Delta x_{i,t} &= \begin{cases} x_{i,t} - \hat{x}_{i,t-1} & \text{if } |x_{i,t} - \hat{x}_{i,t-1}| > \Theta \\ 0 & \text{otherwise} \end{cases}
\end{aligned}
$$
(7)

where $\Delta x_{i,t}$ is the $i$-th element in the delta vector $\Delta \mathbf{x_t}$ at time step $t$, $\hat{x}_{i,t-1}$ is used to prevent the accumulation of delta approximation error. For example, when the absolute difference between $x_{i,t}$ and $\hat{x}_{i,t-1}$ is less than the threshold $\Theta$, the $\Delta x_{i,t}$ will be set to zero, and therefore introducing an approximation error less than $\Theta$. To prevent the accumulation of this approximation error, we need to update $\hat{x}_{i,t-1}$ to reset the approximation error when there is an update above the threshold.

### C. SparseST Unit Structure

The proposed unit structure of our SparseST framework is built by innovatively integrating sparse convolution and the delta algorithm into the ConvLSTM unit, which is shown in Figure 3. It aims to fully exploit the spatial and temporal sparsity to reduce the computational complexity. In our proposed model, the delta tensors $\Delta \mathcal{X}_t$ and $\Delta \mathcal{H}_{t-1}$ for both the input and hidden states are generated as the difference between two adjacent time steps. And then the 2D Sparse Convolution is leveraged to skip computations on inactive sites following the example in Figure 2 .

When updating the three gates $i_t$, $f_t$, $o_t$ and memory cell $\mathcal{C}_t$ at time step $t$, we can recursively get the result by keeping a record of previous sum-product results stored in the delta
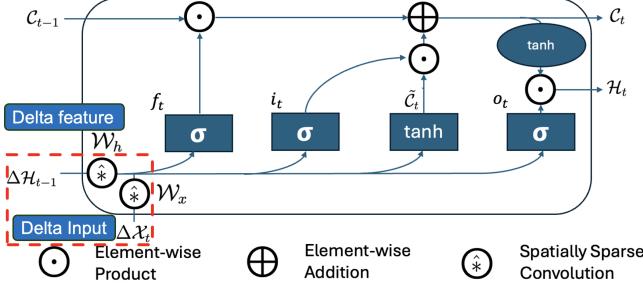
Fig. 3: Unit structure for our proposed SparseST model.

memory $M_{t-1}$ for each update. The update equation for a single SparseST unit is shown in Equation (8).

$$
\begin{aligned}
\{i, f, o\}_t &= \sigma(M_t^{(g)}) \quad g \in \{i, f, o\} \\
&= \sigma(\mathcal{W}_x^{(g)} \hat{\ast} \Delta\mathcal{X}_t + \mathcal{W}_h^{(g)} \hat{\ast} \Delta\mathcal{H}_{t-1} + M_{t-1}^{(g)}) \\
\tilde{\mathcal{C}}_t &= \tanh(M_t^{(c)}) \\
&= \tanh(\mathcal{W}_x^{(c)} \hat{\ast} \Delta\mathcal{X}_t + \mathcal{W}_h^{(c)} \hat{\ast} \Delta\mathcal{H}_{t-1} + M_{t-1}^{(c)}) \\
\mathcal{C}_t &= f_t \odot \mathcal{C}_{t-1} + i_t \odot \tilde{\mathcal{C}}_t \\
\mathcal{H}_t &= o_t \odot \tanh(\mathcal{C}_t)
\end{aligned}
\tag{8}
$$

The delta tensor update equation for $\Delta\mathcal{X}_t$ and $\Delta\mathcal{H}_t$ as shown in Equation (9), where $\delta$ represents the delta thresholding function defined in Equation (7). Note that all of the delta thresholds for each unit are learnable and optimized along with model weights during training.

$$
\begin{aligned}
\Delta\mathcal{X}_t &= \delta(\mathcal{X}_t - \mathcal{X}_{t-1}, \Theta_\mathcal{X}) \\
\Delta\mathcal{H}_t &= \delta(\mathcal{H}_t - \mathcal{H}_{t-1}, \Theta_\mathcal{H})
\end{aligned}
\tag{9}
$$

See Algorithm 1 for the detailed update flow in a single SparseST unit.

---

**Algorithm 1:** SparseST Unit Update Flow

**Input** : Training data sequence $\mathcal{X} = \{\mathcal{X}_1, \mathcal{X}_2, ..., \mathcal{X}_T\}$
**Output:** Trained unit parameter
$\quad\quad \boldsymbol{\theta} = \{\mathcal{W}_x, \mathcal{W}_h, \Theta_\mathcal{X}, \Theta_\mathcal{H}\}$

1 **for** t in iterations **do**
2 $\quad$ Calculate delta tensors $\Delta\mathcal{X}_t$ and $\Delta\mathcal{H}_{t-1}$ by Equation (6)
3 $\quad$ Process $\Delta\mathcal{X}_t$ and $\Delta\mathcal{H}_{t-1}$ using 2D Spatially Sparse Convolution by Equation (5)
4 $\quad$ Update the gates, cell state and delta memories in the unit by Equation (8)
5 $\quad$ Update delta tensors using delta thresholding by Equation (9)
6 $\quad$ Update model parameter $\boldsymbol{\theta}$ by Backpropagation through Time (BPTT) [52]
7 **return** Trained unit parameter $\boldsymbol{\theta}$

---

In our design, Sparse Convolution reduces computation by skipping inactive sites but dilates active spatial sites across layers. The DN algorithm can suppress redundant temporal information by creating delta tensors and mitigate the dilation issue by thresholding. Their hybridization ensures that both spatial and temporal redundancies are minimized: spatially through zero skipping and temporally through delta thresholding.

### D. Computational Cost Analysis

In this section, we analyze the computational cost of a SparseST unit and compare it with the baseline ConvLSTM unit.

Table II shows the computational cost comparison between 2D dense convolution and Sparse Convolution. The acceleration ratio (AR), which measures the reduction in computational cost, is defined as the spatial sparsity of the input tensor, which is the ratio of inactive sites.

| Dense Convolution | Sparse Convolution | Acceleration Ratio |
|---|---|---|
| $HWK^2 C_{in} C_{out}$ | $DK^2 C_{in} C_{out}$ | $1 - D/HW$ |

TABLE II: Computational cost comparison between 2D dense convolution and Sparse Convolution, where $H$ and $W$ are the spatial size of input tensor, $K$ is the size of convolutional filter, $D$ is the number of active sites in the input tensor, $C_{in}$ and $C_{out}$ are the number of channels of input and output tensors, respectively.

In Equation (3) of a single ConvLSTM unit, the total computational cost for gate updates is as follows:

$$
FLOP_{gate} = 4 \times HWK^2 C_{out}(C_{in} + C_{out}) \tag{10}
$$

The update of cell state $\mathcal{C}_t$ involves $2 * HWC_{out}$ number of FLOPs for element-wise multiplication and $HWC_{out}$ number of FLOPs for element-wise addition. The hidden state $\mathcal{H}_t$ involves $HWC_{out}$ number of FLOPs for element-wise multiplication. Therefore, the total computational cost for state updates is as follows:

$$
FLOP_{states} = 4 \times HWC_{out} \tag{11}
$$

The total computational cost for a single ConvLSTM cell is as follows:

$$
\begin{aligned}
FLOP_{dense} &= FLOP_{gate} + FLOP_{states} \\
&= 4 \times HWC_{out}[K^2(C_{out} + C_{in}) + 1]
\end{aligned}
\tag{12}
$$

Note that we didn't include the computational cost for activation functions (sigmoid and tanh function) for simplicity.

For the computational cost of a single SparseST unit, we follow the same analysis as the ConvLSTM unit and combine the result in Table II. Therefore, we have the following:

$$
\begin{aligned}
FLOP_{sparse} &= 4 \times DK^2 C_{out}(C_{in} + C_{out}) + 4 \times HWC_{out} \\
&= 4 * C_{out}[DK^2(C_{out} + C_{in}) + HW]
\end{aligned}
\tag{13}
$$

The acceleration ratio of a single SparseST unit compared with the baseline ConvLSTM unit can be calculated as follows:

$$AR = (FLOP_{dense} - FLOPs_{sparse})/FLOP_{dense}$$
$$= (HW - D)/(HW + \epsilon) \qquad (14)$$
$$\approx 1 - D/HW$$

where $\epsilon = 1/K^2(C_{in} + C_{out})$ is negligible compared with $HW$ in the above denominator. Note that the result in Equation (14) is exactly the ratio of inactive sites (spatial sparsity) of the inputs for a single unit. For a tiny numerical example, if the delta thresholds $\Theta_{\mathcal{X}}$ and $\Theta_{\mathcal{H}}$ introduce the spatial sparsity of 60% and 40% in $\Delta\mathcal{X}_t$ and $\Delta\mathcal{H}_{t-1}$, respectively, we report the AR defined in Equation (14) as 50%, which is the average ratio of spatial sparsity for both delta tensors in one SparseST unit.

*E. Composite Loss for Multi-objective Optimization*

From Section IV-C, we have two sets of parameters to optimize during training. One is the weight matrices $\mathcal{W}$ in Equation 8. Another is the delta threshold $\Theta$ in Equation 9, which directly controls the number of active sites $D$ in Equation 14 for computing AR. A larger $\Theta$ corresponds to a smaller $D$ for a higher AR, since we introduce more sparsity in the delta tensor. However, we achieve this higher efficiency at the cost of model performance.

Therefore, we consider model performance and efficiency to be two conflicting optimization objectives. We aim to enable practitioners to set their preferences under different settings of computational resources, thereby balancing the trade-off between these two objectives.

First, the formulations of these two objectives are defined as follows: the mean square error (MSE) computed as the average squared difference between the predicted and ground truth pixel values across all channels and spatial dimensions for each data sample $\mathcal{L}_{\mathrm{MSE}}(\boldsymbol{\theta})$, and the average unit occupancy computed as the average of the ratio of active sites across all units in the model $\mathcal{L}_{\mathrm{Occupancy}}(\boldsymbol{\theta})$.

$$\mathcal{L}_{\mathrm{MSE}}(\boldsymbol{\theta}) = \frac{1}{n}\sum_{i=1}^{n}(y_i - f(x_i; \boldsymbol{\theta}))^2 \qquad (15)$$

$$\mathcal{L}_{\mathrm{Occupancy}}(\boldsymbol{\theta}) = \frac{1}{NHW}\sum^{N}\sum^{H}\sum^{W}\mathbb{1}(\mathcal{T}_{\boldsymbol{\theta}}) \qquad (16)$$

The indicator function for counting the active sites can be defined as follows:

$$\mathbb{1}(\mathcal{T}_{\boldsymbol{\theta}}) = \begin{cases} 0 & \text{if } \mathcal{T}_{:,i,j}(\boldsymbol{\theta}) = 0 \\ 1 & \text{otherwise} \end{cases} \qquad (17)$$

where $\boldsymbol{\theta}$ is the network parameter, including model weights and learnable delta thresholds, $N$ is the total number of units, $H$ and $W$ are the spatial dimensions for tensor $\mathcal{T} \in \mathbb{R}^{C \times H \times W}$. In our case, the tensor $\mathcal{T}$ can be either delta input $\Delta\mathcal{X}_t$ or delta feature $\Delta\mathcal{H}_{t-1}$.

Given the two objectives defined above, we formulate the problem into a multi-objective optimization (MOO) problem and solve it with the Smooth Tchebycheff Scalarization

method [10] considering the benefits that (1) it allows gradient-based optimization method with modern deep learning optimizer; (2) it provides direct control over the user preference towards each objective; (3) it has a complete theoretical guarantee to explore the Pareto optimal points, including those on the non-convex part of Pareto front.

**Proposition 1.** *(Smooth Tchebycheff Scalarization Method)*
*For a general MOO problem:*

$$\min_{\mathbf{x} \in \mathbf{X}} \mathbf{F}(\mathbf{x}) = [F_1(\mathbf{x}), F_2(\mathbf{x}), \cdots, F_k(\mathbf{x})]^T \qquad (18)$$

*where $\mathbf{x} \in \mathbb{R}^n$ is a vector of design variables, $\mathbf{X}$ is the feasible design space, and $\mathbf{F}(\mathbf{x}) \in \mathbb{R}^k$ is a vector of objective functions $F_i(\mathbf{x}) : \mathbb{R}^n \to \mathbb{R}$.*

*It can be transformed as the following optimization problem by the smooth Tchebycheff scalarization method:*

$$g_\mu^{(\mathrm{STCH})}(\mathbf{x} \mid \boldsymbol{w}) = \mu\log\left(\sum_{i=1}^{k}e^{\frac{w_i(f_i(\mathbf{x})-z_i^*)}{\mu}}\right) \qquad (19)$$

*where $\boldsymbol{w} \in \mathbb{R}^k$ is the preference vector, and $w_i \geq 0, \sum_{i=1}^{k}w_i = 1$ for all $i$, $z_i^*$ is the ideal value for each objective, $\mu > 0$ is the smoothing factor.*

During model training, we minimize the following composite loss function to minimize the MSE loss $\mathcal{L}_{\mathrm{MSE}}$ and minimize the average unit occupancy $\mathcal{L}_{\mathrm{Occupancy}}$, thereby maximizing the average unit sparsity for efficiency.

$$\begin{aligned} g_\mu^{(\mathrm{STCH})}(\boldsymbol{\theta} \mid \boldsymbol{w}) = \mu\log\Big( &\exp\big(\frac{w_{\mathrm{MSE}}\big(\mathcal{L}_{\mathrm{MSE}}(\boldsymbol{\theta})-z_1^*\big)}{\mu}\big) \\ &+ \exp\big(\frac{(1-w_{\mathrm{MSE}})\big(\mathcal{L}_{\mathrm{Occupancy}}(\boldsymbol{\theta})-z_2^*\big)}{\mu}\big)\Big). \end{aligned} \qquad (20)$$

where $w_{MSE}$ is the preference weight on the MSE loss, which provides a direct control for users to balance the trade-off between model performance and efficiency.

**Proposition 2.** *(Ability to Find All Pareto Solutions, quoted from [10])*
*Under mild conditions, there exists a $\mu^*$ such that, for any $0 < \mu < \mu^*$, every Pareto solution of the original multi-objective optimization problem (18) is an optimal solution of the Smooth Tchebycheff scalarization problem (19) with some valid preference $\boldsymbol{\lambda}$.*

See the mild conditions [10] in Appendix B. These two propositions indicate that optimizing the loss function in Equation (20) can recover Pareto optimal points in both convex and non-convex regions of the Pareto front. Consequently, the remaining challenge is to find a surrogate model $h$, described in Section II-B, to explore and approximate the entire Pareto front efficiently.

Lin et al. [10] employed a two-layer multi-layer perceptron (MLP) as the surrogate model. However, in our setting, acquiring objective values is computationally expensive since it requires fully training a deep neural network until convergence for each preference vector to obtain a single Pareto optimal point. This high cost makes it impractical to generate the large number of Pareto optimal points needed to train a neural network surrogate model effectively. To efficiently explore the

Pareto front with minimal rounds of neural network training, we adopt a multi-task GP regression framework. Specifically, we train a GP as the surrogate model with a vector output to approximate the mapping from preference weights to Pareto-optimal outcomes. We begin by initializing the GP with a set of Pareto optimal points obtained by training the neural network within a predefined range of preference weights. At each iteration, an acquisition function based on the sum of marginal variance selects the next preference weight to explore. The neural network is then trained with this new preference, yielding a new Pareto optimal solution that is used to update the GP. At the end of the training loop, the trained GP can be queried with arbitrary preference weights to generate an approximated Pareto front, which is presented in Section V.

## V. EXPERIMENTS

In this section, we design two experiments on two different downstream tasks to present the effectiveness of our framework. One is to perform a spatiotemporal prediction task on the Moving MNIST dataset. The other is anomaly detection task on real industrial dataset.

### A. Case Study: Spatiotemporal Prediction



Fig. 4: A sample of the Moving MNIST dataset.

*1) Dataset:* The Moving MNIST dataset [53] is a synthetic spatiotemporal benchmark dataset, which contains sequences of moving handwritten digits. It is widely used for evaluating sequence modeling and video prediction models. Each data sample in this dataset has 20 grayscale images of size $64 \times 64$. Each image contains two digits of size $28 \times 28$, which is randomly chosen from the original MNIST dataset [54], bouncing inside the $64 \times 64$ frame with a trajectory following a trigonometric pattern. The initial starting position of each digit is randomly picked inside the frame, and the velocity direction is randomly chosen on a unit circle with a magnitude over a fixed range [53]. If the moving digit touches the frame boundary, the direction of its velocity will be inverted, but the magnitude will remain unchanged. Figure 4 shows a sample of the Moving MNIST dataset. The first row shows the moving trajectory of the two digits "1" and "2" for the first ten time steps, and the second row shows the trajectory for the next ten time steps.

*2) Experimental Settings:* Our task is next-frame spatiotemporal prediction with an autoencoder structure supervisedly. And our training, validation, and testing datasets have 8000, 2000, and 10000 samples, respectively. During model training, for each data sample with a length of 20 frames, the encoder takes in the 1st to 10th frames, and the decoder outputs the 2nd to 11th frames, which means that the model is trained to reconstruct the 2nd to 10th frames and make a prediction

on the 11th frame. During model testing, we load the trained model with the lowest validation loss and predict the 11th to 20th images for each sample recursively.

Two baseline models are included for this case study. One is the original ConvLSTM, and another is SparseConvLSTM (ConvLSTM only with 2D Sparse Convolution). SparseConvLSTM is considered one of the baseline models due to the high spatial sparsity of the moving MNIST dataset. We train the baseline models by minimizing the MSE loss, and the proposed SparseST model by minimizing the proposed composite loss in Equation 20. All models are trained by BPTT [52] for 200 epochs with early stopping applied using a patience of 10. Optimization is conducted using the Adam optimizer implemented in PyTorch [55], with an initial learning rate of $10^{-4}$ and a decay rate of 0.5.

*3) Evaluation Metrics:* We compare our proposed model with the baseline methods by MSE loss for the model performance, and the analytical computational savings using the AR shown in Equation (14). We report the average AR over all SparseST units during testing.

*4) Results:* Figure 5 shows a visualization of the predicted results in two randomly selected data samples. For each sample, the first and second rows consist of the first ten frames as input for model training and the last ten frames as the ground truth, respectively. The following rows show the prediction results for different models. As we can see, ConvLSTM and SparseConvLSTM have comparable model performance, indicating that 2D Sparse Convolution achieves similar predictive accuracy to the standard dense Conv2D in the ConvLSTM structure. For SparseST, the model performance degrades when the preference weight $w_{MSE}$ decreases. This trend is expected, since setting a lower weight to the performance objective will shift the focus of optimization towards reducing the computational complexity, therefore achieving a higher acceleration rate at the expense of reduced predictive accuracy.

The numerical result in Table III validates the above analysis. It is worth noting that SparseConvLSTM slightly outperforms ConvLSTM while achieving a 32% reduction in computational cost. This suggests that on the Moving MNIST dataset, 2D Sparse Convolution effectively removes redundant computation caused by spatial sparsity. It preserves performance and even yields a slight improvement. In our proposed SparseST framework, we further exploit the temporal sparsity, which arises from subtle temporal changes of the informative pixels (moving digits) between adjacent time steps. By integrating the Delta Network algorithm, SparseST can leverage the temporal sparsity by optimizing the delta thresholds for each unit. As a result, our framework can achieve higher acceleration ratios, with the trade-off between model performance and acceleration controlled by the preference weight.

### B. Case Study: Anomaly Detection

*1) Dataset:* The Industrial Process Anomaly Detection (IPAD) dataset [56] is specifically curated for the video anomaly detection task in industrial manufacturing scenarios.

(a) Moving digits "7" and "8".



(b) Moving digits "3" and "7".

Fig. 5: Visualizations of model prediction results on Moving MNIST dataset.

It features a comprehensive collection of video sequences containing both real-world data acquired by live video recordings and simulated data generated by SolidWorks, a professional 3D modeling software. The dataset consists of 16 distinct types of industrial machinery, including but not limited to conveyor belts, lift tables, cutting machines, drilling machines, machine grippers, and cranes, providing a rich diversity of normal and abnormal operational patterns.

*2) Experimental Settings:* In this case study, we select the crane scenario from the IPAD dataset, where a crane delivers a cargo from the bottom-left corner into the center of the frame. Each training sequence contains only normal frames, while testing sequences include both normal and abnormal segments with annotated frame-level labels. Each sequence consists of approximately 650 to 700 frames. To enhance the visibility of temporal dynamics between adjacent frames and considering the computational resource constraints, we sampled each sequence every 5 frames and downsampled the size of each frame to $3 \times 64 \times 64$. Figure 6 shows several samples of frames in both training and testing sequences.

| Model Type | MSE loss | GFLOPs | AR |
|---|---|---|---|
| ConvLSTM | 0.001715 | 1322.88 | N/A |
| SparseConvLSTM | 0.001607 | 899.57 | 32.00% |
| SparseST ($w_{MSE} = 1.00$) | 0.002750 | 771.78 | 41.66% |
| SparseST ($w_{MSE} = 0.90$) | 0.003859 | 430.89 | 67.43% |
| SparseST ($w_{MSE} = 0.75$) | 0.004067 | 318.32 | 75.94% |
| SparseST ($w_{MSE} = 0.50$) | 0.005739 | 232.73 | 82.41% |
| SparseST ($w_{MSE} = 0.25$) | 0.009493 | 133.38 | 89.92% |
| SparseST ($w_{MSE} = 0.10$) | 0.017967 | 119.62 | 90.26% |

TABLE III: Comparison of model performance, computational cost (GFLOPs) and acceleration ratio in different settings. GFLOPs is reported as the average number of FLOPs over the testing dataset during inference. For the parameter count, all models have the parameter size close to 7.71 MB. Our proposed SparseST model only introduces delta thresholds for each unit as additional parameters, which is negligible.
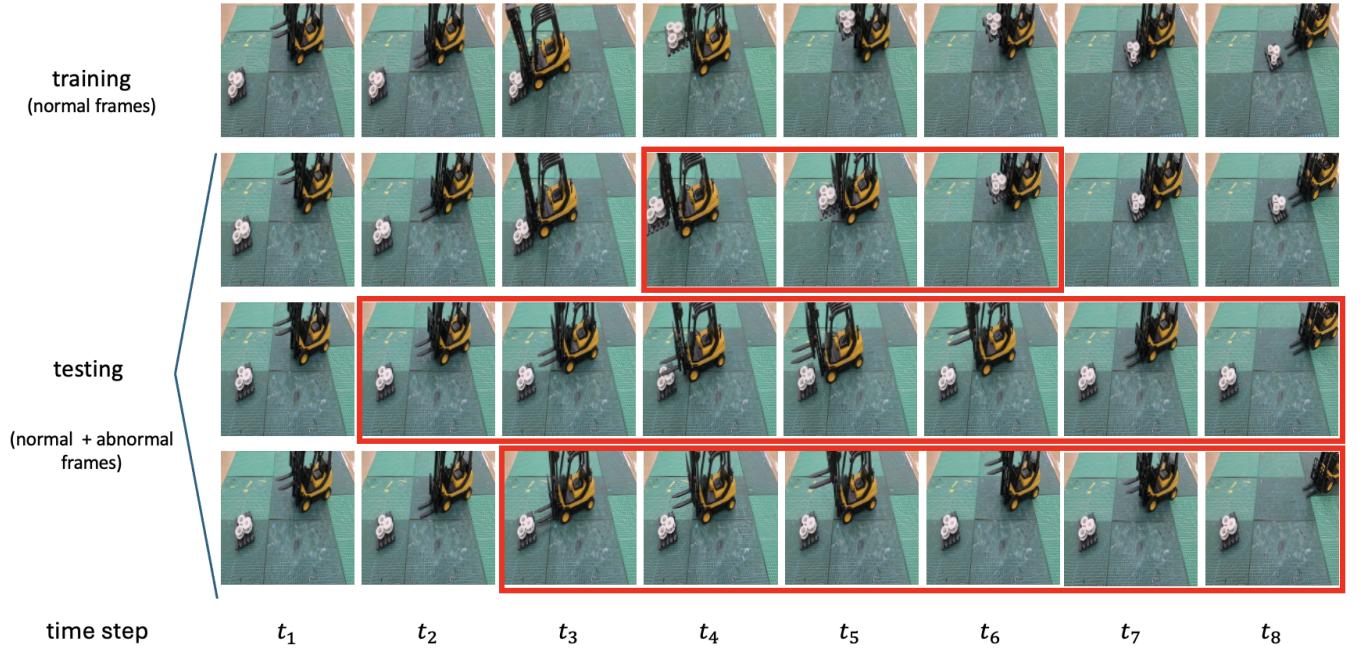


Fig. 6: Sampled frames of a normal sequence. The first row shows a sample of frames of one training sequence; The last three rows show samples of frames of testing sequences with different types of anomaly. Frames highlighted in red are labeled as abnormal.

Figure 7 shows the flow chart of reconstruction-based anomaly detection. During training and validation, we adopt a sliding window approach with a window size of 21 frames and a stride of 1 to capture the long-range spatiotemporal dependencies within the sequence data. The model is trained to reconstruct all the input frames within each window. During testing, we only reconstruct the central frame of each sliding window and record the MSE between the reconstructed and ground truth frames as the anomaly score. As only normal frames are included in the training phase, in the testing phase, abnormal frames with features deviated from the normal patterns tend to yield higher reconstruction errors and thus, enable effective frame-level anomaly detection. All hyperparameters for training are kept the same as those in the previous experiments.

*3) Results:* Figure 8 compares the receiver operating characteristic (ROC) curves of all models, along with their corresponding preference weight, AUC (area under the ROC curve), and acceleration ratio (AR) in the legend.

We can observe that SparseConvLSTM model has a similar AUC score to the ConvLSTM model, and the AR is only 7.74%. This is because the IPAD dataset consists of RGB images with little to no spatial sparsity compared with the Moving MNIST dataset. In this case, 2D Sparse Convolution itself cannot effectively leverage sparsity to improve efficiency without the DN algorithm.

For the SparseST models, we integrate the DN algorithm to leverage the temporal sparsity of data, and the AR can be largely boosted. A lower preference weight of MSE introduces lower values of AUC (worse anomaly detection performance) and higher values of AR (lower computational costs), indicating greater efficiency at the expense of model performance. In particular, the SparseST models perform better (higher AUC scores) than the baseline model when the preference weight

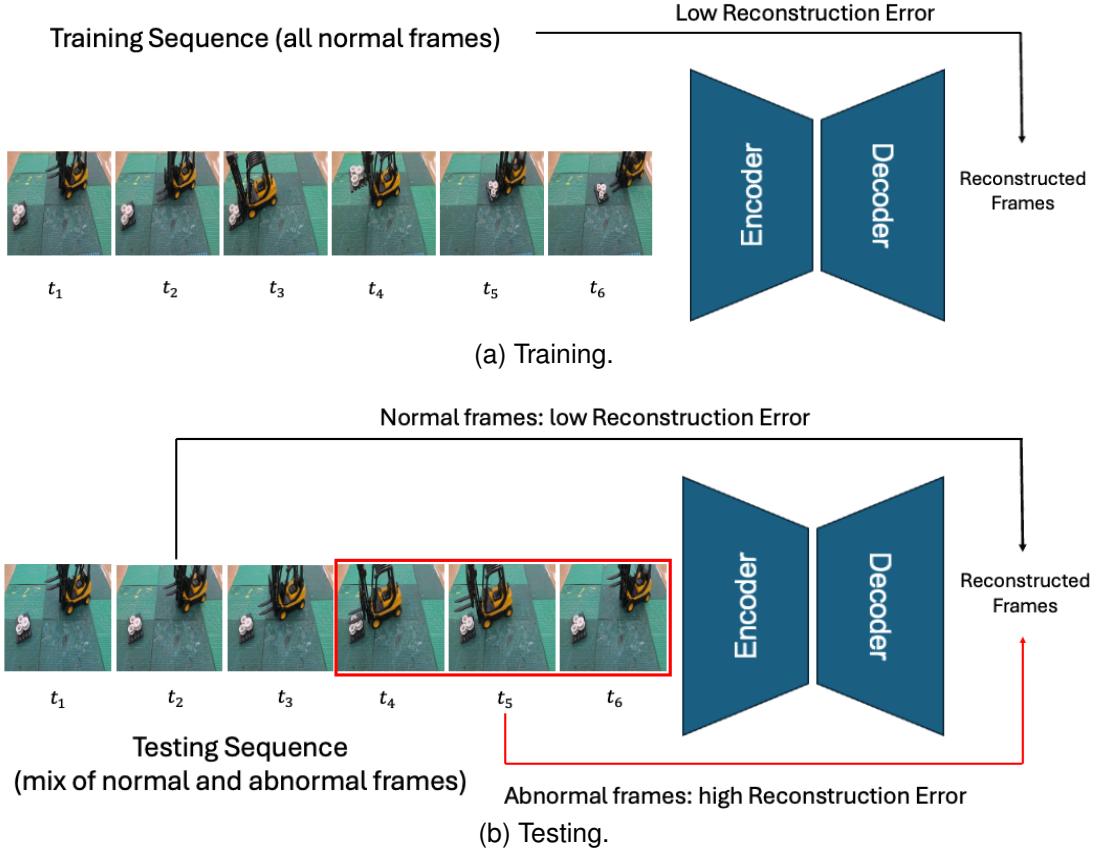(a) Training.



(b) Testing.

Fig. 7: Workflow of reconstruction-based anomaly detection.

of MSE is greater than 0.25. This means, for reconstruction-based anomaly detection tasks, there is excess information within the data that can be traded for efficiency. Despite the spatiotemporal sparsity introduced by the framework, the information remained in the delta input or features is sufficient to enable the model to distinguish the normal pattern from the abnormal. However, when the preference weight of MSE is decreased to 0.10, the performance of SparseST turns worse than the baseline ConvLSTM, implying that the loss of information introduced by spatial and temporal sparsity in the data stream is too significant to provide sufficient information for anomaly detection.

Figure 9 shows the anomaly score comparison between normal and abnormal frames using different models on a single testing sequence. The MSE for each frame in the sequence is plotted as the anomaly score. As highlighted in the figure, two distinct peaks appear consistently across all models, indicating that each model successfully detects two types of anomalies in this sequence: (1) when the crane approaches the cargo but fails to load it due to the arm being positioned too high, and (2) when the crane moves backwards to the destination, it fails to move the cargo with it. Although there is an increasing trend of anomaly scores across SparseST models when the preference weight for MSE decreases, the difference between anomaly scores on normal and abnormal frames is still significant enough for anomaly detection. This suggests that in this case study, our proposed framework can effectively trade redundant

information in the data for efficiency.

### C. Pareto Front Analysis

Following the methodology mentioned in Section IV-E, we obtain Figure 10 to illustrate the approximated Pareto fronts for the two competing objectives: model performance on the downstream tasks (measured by MSE) and computational efficiency (measured by occupancy) in the above two case studies. They include both the initial points and the new candidate points with the largest sum of marginal variances of the two objectives in each iteration. Each point is annotated by the value of its associated preference weight $w_{MSE}$. Ideally, all of these points are supposed to be Pareto optimal (non-dominated). However, we obtain these points by training a neural network, which means the convergence and global optimality of the training process cannot be fully guaranteed. It is worth noting that the approximated Pareto fronts are non-convex, which empirically demonstrates the effectiveness of the smoothed Tchebycheff scalarization to recover Pareto-optimal solutions in non-convex regions of the objective space. And this is what standard linear scalarization fails to guarantee. A clear trade-off trend between the two objectives shows up along with the Pareto fronts. As the preference weight $w_{MSE}$ increases, the model achieves lower MSE (better accuracy) at the cost of higher occupancy (lower acceleration). Conversely, decreasing $w_{MSE}$ will favor sparsity and computational efficiency, but reduce predictive accuracy. This visualization can
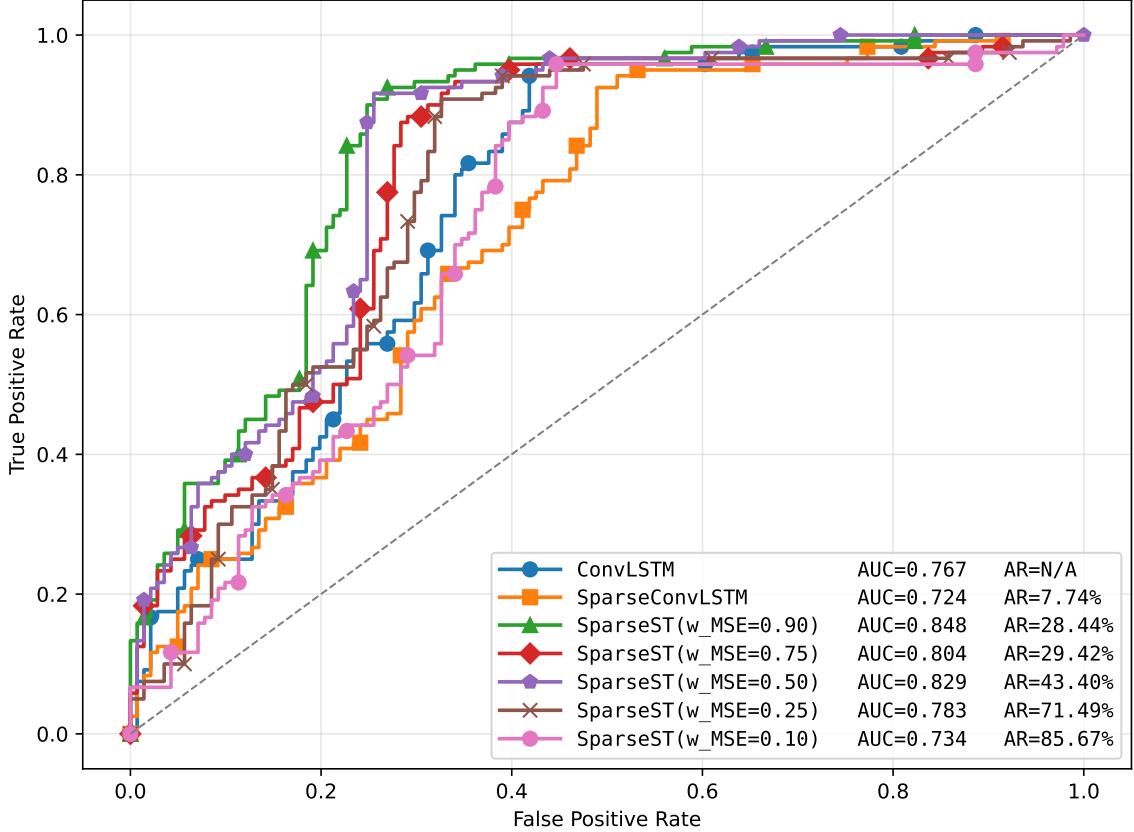
Fig. 8: ROC curve comparison between different models.

provide practical guidance for practitioners to select appropriate preference weights tailored to the computational resource constraints or accuracy requirements of a downstream task. For example, one can choose a point on the Pareto front that satisfies a specific range of MSE or occupancy, then finetune the preference weight according to the explored points.

## VI. SUMMARY

In this work, we investigate efficient AI for STDM with the goal of exploiting spatial and temporal sparsity in raw data and intermediate features. Our motivation is to trade redundant information for efficiency, particularly in spatiotemporal data with fixed or slowly changing backgrounds, where modeling the full complex dependency is unnecessary.

To address this, we integrate 2D Sparse Convolution and the DN algorithm within a ConvLSTM backbone and propose the SparseST framework. This design not only improves computational efficiency but also overcomes key limitations: the DN algorithm mitigates the active site dilation problem of 2D Sparse Convolution and effectively leverages temporal correlations to exploit sparsity between consecutive frames.

Experimental results demonstrate that SparseST achieves substantial computational savings while maintaining accuracy comparable to baseline models. Moreover, for downstream tasks that do not rely on full spatiotemporal dependency modeling (e.g., anomaly detection), our framework provides improved efficiency without sacrificing or even improving the model performance.

Furthermore, we formulate the performance–efficiency trade-off as a multi-objective optimization problem and design a composite loss function to balance the two objectives. By applying STCH scalarization within a multi-task GP regression, we approximate the Pareto front and provide practical guidance for adjusting preferences between computational cost and predictive performance, depending on downstream task requirements.

## ACKNOWLEDGMENTS

## CONFLICT OF INTEREST

All authors declare that they have no known conflicts of interest in terms of competing financial interests or personal relationships that could have an influence or are relevant to the work reported in this paper

## APPENDIX

**Definition 1.** [57] Pareto dominance: A point $\mathbf{x}^* \in \mathbf{X}$ dominates another point $\mathbf{x}'$ when both: (1) $\mathbf{x}^*$ is not worse than $\mathbf{x}'$ on any objective, i.e. $F_i(\mathbf{x}^*) \leq F_i(\mathbf{x}')$, $\forall i \in \{1, ..., k\}$; (2) $\mathbf{x}^*$ is better than $\mathbf{x}'$ on at least one objective, i.e. $\exists j \in \{1, ..., k\}$ s.t. $F_j(\mathbf{x}^*) < F_j(\mathbf{x}')$.
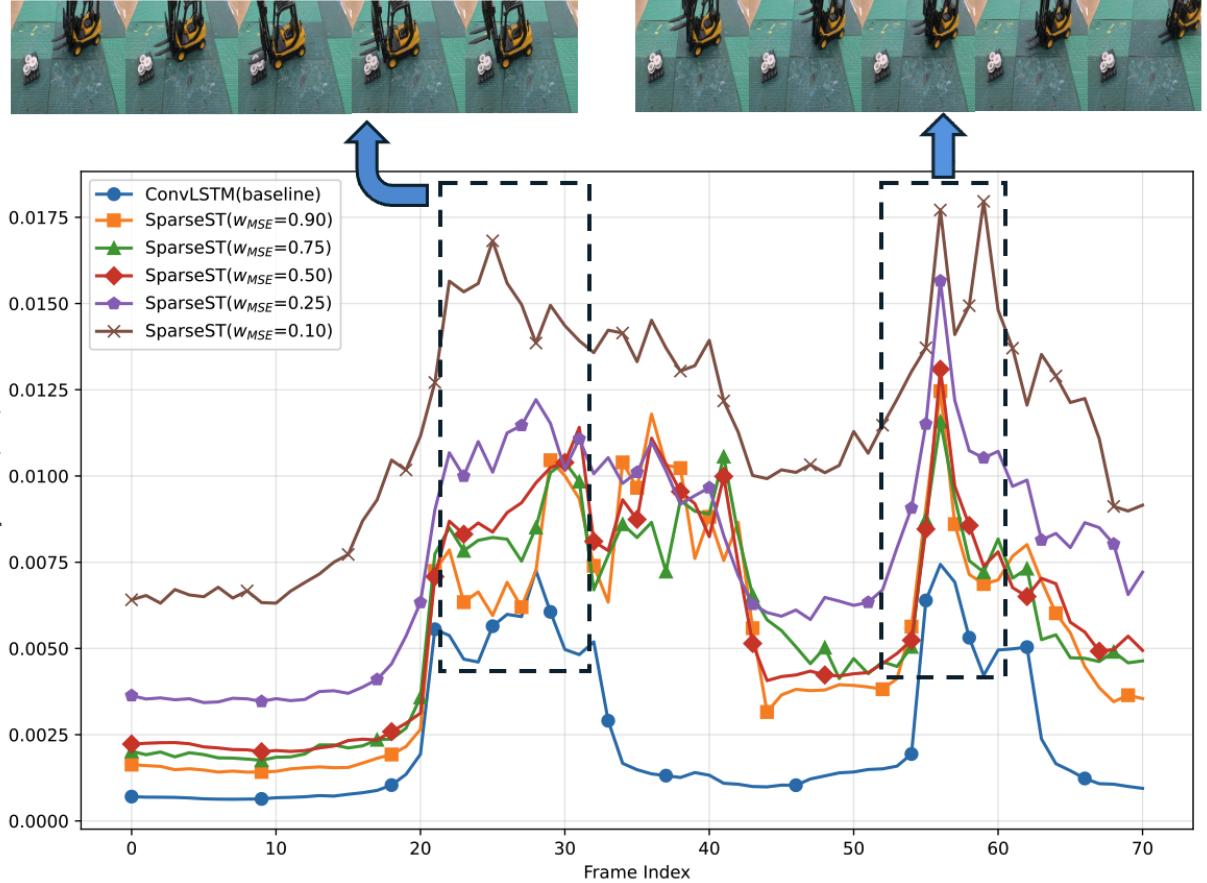
Fig. 9: Difference of anomaly score between normal and abnormal frames.

**Definition 2.** [57] Pareto optimality: A point $\mathbf{x}^* \in \mathbf{X}$ is Pareto optimal if it is not dominated by any other points. The Pareto frontier is the k-dimensional manifold of the objective values of all Pareto optimal solutions.

**Definition 3.** [58] Weak Pareto optimality: A point $\mathbf{x}^* \in \mathbf{X}$ is weakly Pareto optimal iff there does not exist another point $\mathbf{x} \in \mathbf{X}$, such that $F_i(x) < F_i(x^*)$, $\forall i \in \{1, ..., k\}$.

Define the Pareto front of a multi-objective optimization problem as a surface $\phi(x) = 0$, where $x \in \mathbb{R}^m$ is a vector in the objective space.

**Assumption 1.** *[59] For each point $x$ on the Pareto front $\phi$, assume:*
- $\nabla \phi(\mathbf{y}) \succ 0$, *which means* $\frac{\partial \phi(\mathbf{y})}{\partial y_i} > 0$ *for all* $i = 1, \ldots, m$;
- *all second-order derivatives of $\phi(\mathbf{y})$ are bounded (i.e., the Hessian $\nabla^2 \phi(y)$ has finite eigenvalues).*

The first statement assumes $\phi$ is strictly increasing in each objective, indicating that any weakly Pareto optimal solution must also be Pareto optimal. The second statement is a standard regularity assumption to ensure the scalarization is smooth enough.

### APPENDIX

The GP training curves for both experiments are shown in Figure 11. In Figure, we show the predictions of the trained

multi-task GP on each of the two objectives.

### REFERENCES

[1] Y. Zhang, P. Wang, B. Wang, X. Wang, Z. Zhao, Z. Zhou, L. Bai, and Y. Wang, "Adaptive and interactive multi-level spatio-temporal network for traffic forecasting," *IEEE Transactions on Intelligent Transportation Systems*, 2024.

[2] J. Lee, B. Bagheri, and H.-A. Kao, "Recent advances and trends of cyber-physical systems and big data analytics in industrial informatics," in *International proceeding of int conference on industrial informatics (INDIN)*, pp. 1–6, Citeseer, 2014.

[3] Y. Cheng, D. Wang, P. Zhou, and T. Zhang, "A survey of model compression and acceleration for deep neural networks," *arXiv preprint arXiv:1710.09282*, 2017.

[4] T. Choudhary, V. Mishra, A. Goswami, and J. Sarangapani, "A comprehensive survey on model compression and acceleration," *Artificial Intelligence Review*, vol. 53, pp. 5113–5155, 2020.

[5] X. Shi, Z. Chen, H. Wang, D.-Y. Yeung, W.-K. Wong, and W.-c. Woo, "Convolutional lstm network: A machine learning approach for precipitation nowcasting," *Advances in neural information processing systems*, vol. 28, 2015.

[6] N. Jain, M. Wyawahare, V. Mankar, and T. Paratkar, "A survey on efficient neural network compression techniques," in *International Conference on Computational Intelligence*, pp. 285–298, Springer, 2022.

[7] M. Yin, S. Liao, X.-Y. Liu, X. Wang, and B. Yuan, "Compressing recurrent neural networks using hierarchical tucker tensor decomposition," 2020.

[8] B. Graham and L. van der Maaten, "Submanifold sparse convolutional networks," 2017.

[9] D. Neil, J. H. Lee, T. Delbruck, and S.-C. Liu, "Delta networks for optimized recurrent network computation," in *International conference on machine learning*, pp. 2584–2593, PMLR, 2017.

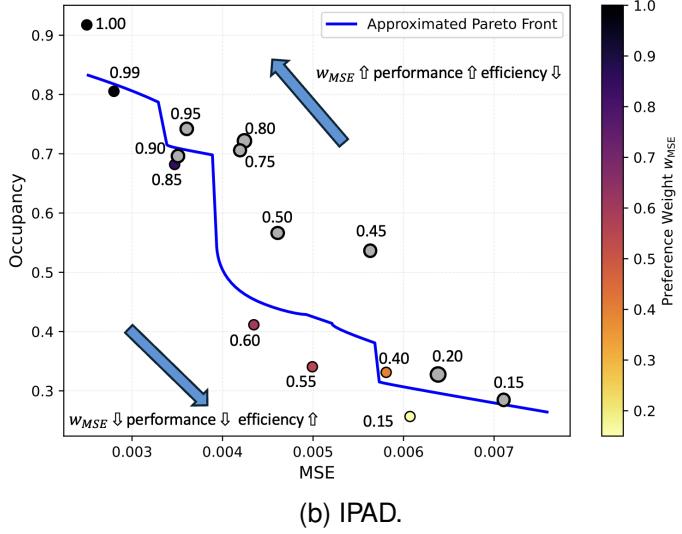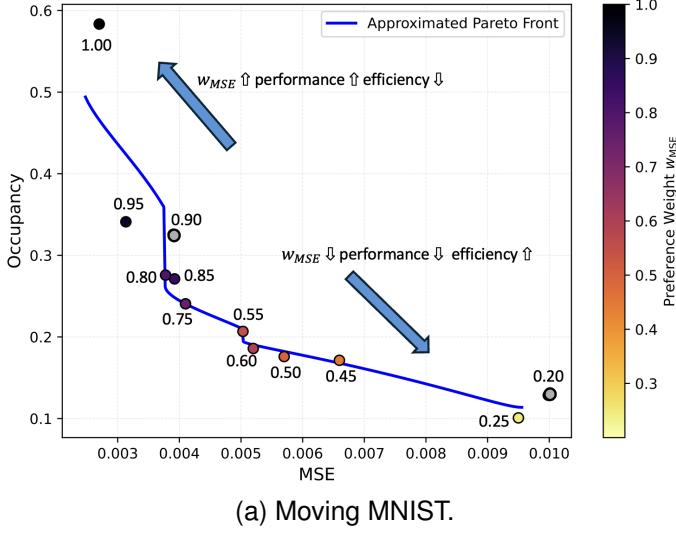(a) Moving MNIST.



(b) IPAD.

Fig. 10: Visualization of Approximated Pareto Front. The points colored in gray are dominated, while non-dominated points are colored w.r.t the preference weight colorbar.



Fig. 11: Change of negative marginal log likelihood (MLL) loss over GP training iterations.



(a) Moving MNIST.



(b) IPAD.

Fig. 12: Multi-task GP predictions on each objective for both experiments. The blue region is the 95% confidence interval for each preference weight on the x-axis.

[10] X. Lin, X. Zhang, Z. Yang, F. Liu, Z. Wang, and Q. Zhang, "Smooth tchebycheff scalarization for multi-objective optimization," *arXiv preprint arXiv:2402.19078*, 2024.

[11] S. Han, J. Pool, J. Tran, and W. Dally, "Learning both weights and connections for efficient neural network," *Advances in neural information processing systems*, vol. 28, 2015.

[12] H. Li, A. Kadav, I. Durdanovic, H. Samet, and H. P. Graf, "Pruning filters for efficient convnets," 2017.

[13] Z. Liu, J. Li, Z. Shen, G. Huang, S. Yan, and C. Zhang, "Learning efficient convolutional networks through network slimming," in *Proceedings of the IEEE international conference on computer vision*, pp. 2736–2744, 2017.

[14] S. Narang, E. Elsen, G. Diamos, and S. Sengupta, "Exploring sparsity in recurrent neural networks," 2017.

[15] W. Wen, Y. He, S. Rajbhandari, M. Zhang, W. Wang, F. Liu, B. Hu, Y. Chen, and H. Li, "Learning intrinsic sparse structures within long short-term memory," 2018.

[16] J. Ott, Z. Lin, Y. Zhang, S.-C. Liu, and Y. Bengio, "Recurrent neural networks with limited numerical precision," *arXiv preprint arXiv:1608.06902*, 2016.

[17] W. Wen, C. Xu, F. Yan, C. Wu, Y. Wang, Y. Chen, and H. Li, "Terngrad: Ternary gradients to reduce communication in distributed deep learning," *Advances in neural information processing systems*, vol. 30, 2017.
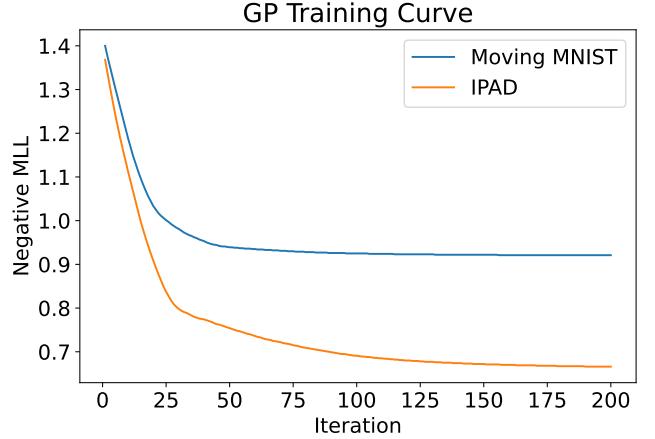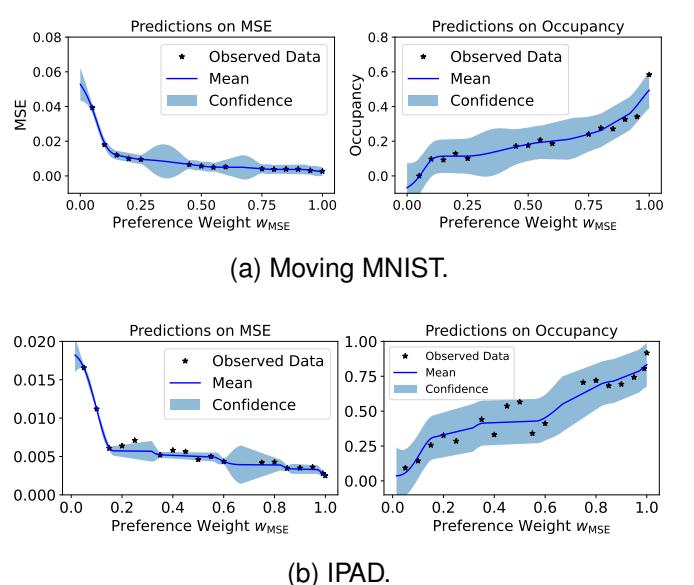
[18] I. Hubara, M. Courbariaux, D. Soudry, R. El-Yaniv, and Y. Bengio, "Quantized neural networks: Training neural networks with low precision weights and activations," *Journal of Machine Learning Research*, vol. 18, no. 187, pp. 1–30, 2018.

[19] M. Denil, B. Shakibi, L. Dinh, M. Ranzato, and N. De Freitas, "Predicting parameters in deep learning," *Advances in neural information processing systems*, vol. 26, 2013.

[20] Y. Lu, A. Kumar, S. Zhai, Y. Cheng, T. Javidi, and R. Feris, "Fully-adaptive feature sharing in multi-task networks with applications in person attribute classification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5334–5343, 2017.

[21] V. Lebedev, Y. Ganin, M. Rakhuba, I. Oseledets, and V. Lempitsky, "Speeding-up convolutional neural networks using fine-tuned cp-decomposition," *arXiv preprint arXiv:1412.6553*, 2014.

[22] C. Tai, T. Xiao, Y. Zhang, X. Wang, *et al.*, "Convolutional neural networks with low-rank regularization," *arXiv preprint arXiv:1511.06067*, 2015.

[23] Y. Wang, W. Guo, and X. Yue, "Tensor decomposition to compress convolutional layers in deep learning," *IISE Transactions*, vol. 54, no. 5, pp. 481–495, 2022.

[24] Z. Lu, V. Sindhwani, and T. N. Sainath, "Learning compact recurrent

neural networks," in *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 5960–5964, IEEE, 2016.

[25] H. Sak, A. W. Senior, and F. Beaufays, "Long short-term memory recurrent neural network architectures for large scale acoustic modeling," 2014.

[26] B. E. Kingsbury, T. N. Sainath, and V. Sindhwani, "Low-rank matrix factorization for deep belief network training with high-dimensional output targets," Feb. 16 2016. US Patent 9,262,724.

[27] Y. Yang, D. Krompass, and V. Tresp, "Tensor-train recurrent neural networks for video classification," 2017.

[28] Y. Pan, J. Xu, M. Wang, J. Ye, F. Wang, K. Bai, and Z. Xu, "Compressing recurrent neural networks with tensor ring for action recognition," 2018.

[29] J. Ye, L. Wang, G. Li, D. Chen, S. Zhe, X. Chu, and Z. Xu, "Learning compact recurrent neural networks with block-term tensor decomposition," 2018.

[30] B. Graham, "Spatially-sparse convolutional neural networks," 2014.

[31] M. Engelcke, D. Rao, D. Z. Wang, C. H. Tong, and I. Posner, "Vote3deep: Fast object detection in 3d point clouds using efficient convolutional neural networks," 2017.

[32] B. Graham, "Sparse 3d convolutional neural networks," 2015.

[33] B. Graham, M. Engelcke, and L. Van Der Maaten, "3d semantic segmentation with submanifold sparse convolutional networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 9224–9232, 2018.

[34] Y. Yan, Y. Mao, and B. Li, "Second: Sparsely embedded convolutional detection," *Sensors*, vol. 18, no. 10, p. 3337, 2018.

[35] H. Tang, Z. Liu, X. Li, Y. Lin, and S. Han, "Torchsparse: Efficient point cloud inference engine," *Proceedings of Machine Learning and Systems*, vol. 4, pp. 302–315, 2022.

[36] H. Tang, S. Yang, Z. Liu, K. Hong, Z. Yu, X. Li, G. Dai, Y. Wang, and S. Han, "Torchsparse++: Efficient training and inference framework for sparse convolution on gpus," in *Proceedings of the 56th Annual IEEE/ACM International Symposium on Microarchitecture*, pp. 225–239, 2023.

[37] C. Gao, T. Delbruck, and S.-C. Liu, "Spartus: A 9.4 top/s fpga-based lstm accelerator exploiting spatio-temporal sparsity," *IEEE Transactions on Neural Networks and Learning Systems*, 2022.

[38] C. Gao, D. Neil, E. Ceolini, S.-C. Liu, and T. Delbruck, "Deltarnn: A power-efficient recurrent neural network accelerator," in *Proceedings of the 2018 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, pp. 21–30, 2018.

[39] C. Gao, A. Rios-Navarro, X. Chen, S.-C. Liu, and T. Delbruck, "Edge-drnn: Recurrent neural network accelerator for edge inference," *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, vol. 10, no. 4, pp. 419–432, 2020.

[40] S. P. Boyd and L. Vandenberghe, *Convex optimization*. Cambridge university press, 2004.

[41] V. J. Bowman Jr, "On the relationship of the tchebycheff norm and the efficient frontier of multiple-criteria objectives," in *Multiple Criteria Decision Making: Proceedings of a Conference Jouy-en-Josas, France May 21–23, 1975*, pp. 76–86, Springer, 1976.

[42] R. E. Steuer and E.-U. Choo, "An interactive weighted tchebycheff procedure for multiple objective programming," *Mathematical programming*, vol. 26, no. 3, pp. 326–344, 1983.

[43] E. U. Choo and D. R. Atkins, "Proper efficiency in nonconvex multicriteria programming," *Mathematics of Operations Research*, vol. 8, no. 3, pp. 467–470, 1983.

[44] J. Knowles, "Parego: A hybrid algorithm with on-line landscape approximation for expensive multiobjective optimization problems," *IEEE transactions on evolutionary computation*, vol. 10, no. 1, pp. 50–66, 2006.

[45] X. Lin, Y. Liu, X. Zhang, F. Liu, Z. Wang, and Q. Zhang, "Few for many: Tchebycheff set scalarization for many-objective optimization," *arXiv preprint arXiv:2405.19650*, 2024.

[46] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[47] Q. Zhang and H. Li, "Moea/d: A multiobjective evolutionary algorithm based on decomposition," *IEEE Transactions on evolutionary computation*, vol. 11, no. 6, pp. 712–731, 2007.

[48] K. Deb, A. Pratap, S. Agarwal, and T. Meyarivan, "A fast and elitist multiobjective genetic algorithm: Nsga-ii," *IEEE transactions on evolutionary computation*, vol. 6, no. 2, pp. 182–197, 2002.

[49] N. Beume, B. Naujoks, and M. Emmerich, "Sms-emoa: Multiobjective selection based on dominated hypervolume," *European journal of operational research*, vol. 181, no. 3, pp. 1653–1669, 2007.

[50] X. Zhang, X. Lin, B. Xue, Y. Chen, and Q. Zhang, "Hypervolume maximization: A geometric view of pareto set learning," *Advances in Neural Information Processing Systems*, vol. 36, pp. 38902–38929, 2023.

[51] X. Lin, Z. Yang, X. Zhang, and Q. Zhang, "Pareto set learning for expensive multi-objective optimization," *Advances in neural information processing systems*, vol. 35, pp. 19231–19247, 2022.

[52] P. J. Werbos, "Backpropagation through time: what it does and how to do it," *Proceedings of the IEEE*, vol. 78, no. 10, pp. 1550–1560, 1990.

[53] N. Srivastava, E. Mansimov, and R. Salakhutdinov, "Unsupervised learning of video representations using lstms," 2016.

[54] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.

[55] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Köpf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, "Pytorch: An imperative style, high-performance deep learning library," pp. 8024–8035, 2019.

[56] J. Liu, Y. Yan, J. Li, W. Zhao, P. Chu, X. Sheng, Y. Liu, and X. Yang, "Ipad: Industrial process anomaly detection dataset," *IEEE Transactions on Circuits and Systems for Video Technology*, 2024.

[57] M. Ruchte and J. Grabocka, "Scalable pareto front approximation for deep multi-objective learning," in *2021 IEEE international conference on data mining (ICDM)*, pp. 1306–1311, IEEE, 2021.

[58] R. T. Marler and J. S. Arora, "Survey of multi-objective optimization methods for engineering," *Structural and multidisciplinary optimization*, vol. 26, no. 6, pp. 369–395, 2004.

[59] D. Li, "Convexification of a noninferior frontier," *Journal of Optimization Theory and Applications*, vol. 88, no. 1, pp. 177–196, 1996.