

A Machine Learning Approach for Rate Prediction in Multicast File-stream Distribution Networks

Yujia Mu, Yuanlong Tan, Malathi Veeraraghavan*, Cong Shen

The Charles L. Brown Department of Electrical and Computer Engineering, University of Virginia, USA

Abstract—Large-volume scientific data is one of the prominent driving forces behind next generation networking. In particular, Software Defined Network (SDN) makes leveraging path-based network multicast services practically feasible. In our prior work, we have developed a cross-layer architecture for supporting reliable file-streams multicasting over SDN-enabled Layer-2 network, and implemented the architecture for a meteorology data distribution application in atmospheric science. However, it is challenging to determine an optimal rate for this application with the varying type, volume, and quality of meteorological data. In this paper, we propose a Quality of Service (QoS)-driven rate management pipeline to determine the optimal rate based on the input traffic characteristics and performance constraints. Specifically, the pipeline employs a feedtype classifier using Multi-Layer Perception (MLP) to recognize the type of meteorological data and a delay prediction regressor using stacked Long Short-Term Memory (LSTM) to predict per-file delay for the file-streams. Finally, we determine the optimal rate for the given file-streams using the trained regressor. We implement this pipeline to test the real-world file-stream data collected from a trial deployment, and the results show that our regressor outperforms all baselines by selecting the optimal rate in the presence of varying file set sizes.

I. INTRODUCTION

Due to the significant increase of scientific activities, data transfer among the geographically distributed research communities has become a major bottleneck and attracted much interest. For example, in the Unidata Internet Data Distribution (IDD) project, the University Corporation of Atmospheric Research (UCAR) uses Local Data Manager (LDM) to distribute 30 different types of meteorological data (e.g., satellite imagery, lightning data, computer-model output, and radar data) to almost 240 institutions on a near-real-time basis [1].

The current LDM, *LDM6*, is an Application-Layer Multicast (ALM) solution that uses multiple unicast transmissions to construct a multicast. With the rapid growth of data volume, ALM requires much more sender computation power and network bandwidth than the network multicast. Compared with the complexity of IP multicast, the use of SDN greatly simplifies the mechanisms for configuring forwarding table information in network switches to build a multicast tree.

In our prior work [2], [3], we have developed and implemented a cross-layer architecture for supporting reliable file-

stream multicasting over Layer-2 (L2) path-based services. The solution adopts a transport layer protocol, called File Transfer Multicast Protocol (FMTP) [4], and leverages L2 multipoint VLAN/MPLS paths. Compared with the IP-routed service, a path-based service offers in-sequence delivery and rate guarantees. With rate guarantees, flow control and congestion control are handled by receivers that receive packets at a fixed rate of the multipoint VLAN in the path setup phase.

The problem we study in this paper is how to determine the optimal VLAN/sending rates for this multicast file-stream distribution system. The silence periods between files within the file-streams, and the varying rate at which the meteorological data is generated, make the problem challenging. Different feedtypes have drastically different characteristics and it is complicated to optimize the rate when all the feedtypes are considered together. In particular, the relationship between delay and rate is not monotonic but depending on various system parameters include file-stream characteristics (size of files, inter-arrival time between files within a file-stream), the number of receivers in the multicast tree, and round-trip times (RTTs) on the path from the sender to each of the receivers.

In this paper, we propose a QoS-driven rate management pipeline that centers on a machine learning (ML) engine, which aims at finding the optimal VLAN/sending rates in a path-based multicast networking based on the predefined requirement of QoS. We first design a feedtype classifier using MLP to identify the feedtype of the meteorological data based on the statistical features. Then, we train a delay prediction regressor using stacked LSTM to predict the per-file delay of the file-streams for the corresponding feedtype. The ML engine then combines the input traffic characteristics with the QoS requirement to determine the sending rate for the next time interval. We validate our method on the real-world file-stream data collected from a trial deployment, and the results show that the performance of our LSTM regressor outperforms all the baselines with respect to various metrics. Although our approach is proposed based on the meteorology data, the approach can be extended to other applications.

The rest of this paper is organized as follows. Related works are reviewed in Section II. The file-stream distribution system is presented in Section III. In Section IV, we describe the proposed QoS-driven rate management pipeline, followed by an experimental evaluation in Section V. Finally Section VI concludes the paper.

The work was partially supported by the US National Science Foundation (NSF) under Grant OAC-1659174, and in part by a Virginia Commonwealth Cyber Initiative (CCI) cybersecurity research collaboration grant.

* Prof. Malathi Veeraraghavan contributed significantly to this work before her untimely passing earlier this year. She did not have a chance to read the submitted version of this paper.

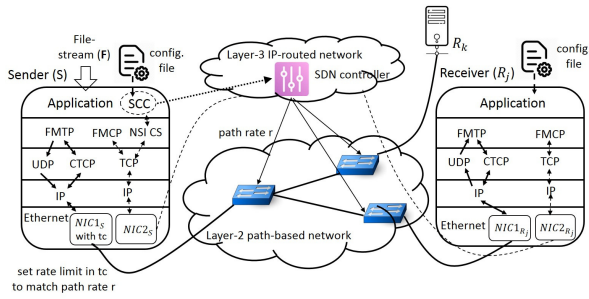


Fig. 1. Dynamic Reliable File-Stream Multicast (DRFSM) Architecture

II. RELATED WORK

One of the conventional multicast services applied to file-stream distribution is NACK-Oriented Reliable Multicast (NORM) [5]. Since it is designed for IP-multicast that does not offer rate guarantees, NORM adopts a rate-based congestion control mechanism. A multicast-based replication for Hadoop HDFS is proposed to save network bandwidth [6], [7], by developing a congestion-controlled reliable multicast socket (CCRMSocket) that uses IP-multicast to provide reliable multicast service. In comparison, our Dynamic Reliable File-Stream Multicast (DRFSM) service avoids data-plane congestion control by using SDN controllers to configure a rate-guaranteed multipoint VLAN, whose rate depends on QoS.

Inspired by its significant development, ML-based approaches for solving complex network problems have attracted much interest [8]–[10]. Mirza *et al.* applied support vector regression to predict network throughput based on file size and measurements of path properties, e.g., queuing delay [8]. Kong *et al.* proposed a packet loss predictor using random forest, which could predict the probability of packet loss caused by congestion, lower the frequency of sending rate reduction, and achieve higher throughput [9]. Both works do not directly adjust the sending rate, which thus does not solve the multicast rate selection problem of this paper. In addition, Mei *et al.* used LSTM based on a set of collected bandwidth traces for accurate prediction of mobile bandwidth [10].

III. PROBLEM FORMULATION

We proposed a DRFSM cross-layer architecture, illustrated in Fig. 1, for supporting scientific file-streams distribution. It relies on i) Layer-2 (L2) multipoint Virtual LAN (VLAN) service, and ii) a reliable transport-layer protocol of FMTP.

This architecture requires two types of network services: L2 path-based service for data-plane and L3 IP-routed service for control-plane. In a multi-receiver context, sequenced delivery and rate guarantees would simplify the key transport-layer functions of error control, flow control, and congestion control. With sequenced delivery, error control is simpler, because a receiver can assume that a packet is dropped when it receives an out-of-sequence block, and then sends a Negative ACKnowledgment (NACK) requesting a block retransmission. With rate guarantees on the paths from the sender to receivers, flow control and congestion control are handled by receivers

agreeing to handle packets at the fixed rate of the multipoint network path in the control-plane path setup phase.

The rate is selected by the sender based on the characteristics of the file-stream and the application delay requirements. The key point is that active sending rate adjustments in the data-plane, i.e., during file-stream transmission, are avoided so that one receiver with insufficient compute or network resources does not slow down data delivery to other receivers.

A. Latency Requirements

The latency of each file consists of six parts, and we usually focus on the following four: (i) transmission delay, (ii) sender-buffer queuing delay, (iii) retransmission delay, and (iv) propagation delay. The switch/router queuing delay and processing delay are out of the scope of this work.

Except for propagation delay that is determined by the link length between the sender and receiver, the other three types of delay are related to the sending rate. The relation of delay and rate is not monotonic but depending on various parameters, which is empirically explained in our previous work [3].

We use file latency as a metric to evaluate the network performance. It is meaningful to describe the file-stream of latencies at percentiles. To that end, we first define holding interval τ as the period during which the assigned sending rate is unchanged. The k th holding interval is denoted as $(t + (k-1)\tau, t + k\tau)$. The file index i is reset to 1 for the first file arriving in each holding interval. The arrival time of file i is denoted as a_i . Thus, the delay requirement $R = (W, \alpha, \tau)$ is set such that only a fraction α of files would experience delay greater than the threshold W in holding interval τ at a receiver. Finally, the delay requirement is defined as $\frac{|V_k|}{N_k} \leq \alpha$, where V_k is the set of files that experienced latency greater than W in the k th holding interval and N_k is the number of files that arrived in this interval. Note that file i belongs to set V_k if $w_i > W$ and $t + (k-1)\tau \leq a_i < t + k\tau$, where w_i denotes the latency experienced by file i .

B. Feedtype Classification and Delay Prediction Regression

The UCAR Unidata IDD project distributes 30 types of meteorological data (feedtypes), which are very different in both distribution and quality. Therefore, it is very challenging to predict precise latency when considering mixed feedtypes. Fig. 2 illustrates this challenge by plotting the arrival rates of three feedtypes: NGRID, CONDUIT, and NEXRAD2, on a one-hour basis over 166 hours. The traffic characteristics among them are very distinguished, which means it is feasible to accurately identify the incoming feedtype. We thus adopt a cascading pipeline, where we first design a classifier to recognize the feedtype and then design a regressor to predict the file latency and find an optimal rate using the output of the feedtype classifier as input.

Within the holding interval, we use $Y = (Y_1, Y_2, \dots, Y_m)$ to denote the latencies of the m files, where Y_i describes the time interval that file i consumes from the sender to the receiver. We use $X = (X_1, X_2, \dots, X_m)$ to represent input samples, where $X_i = (x_i^1, x_i^2, \dots, x_i^n)$ are n measured features for sample i .

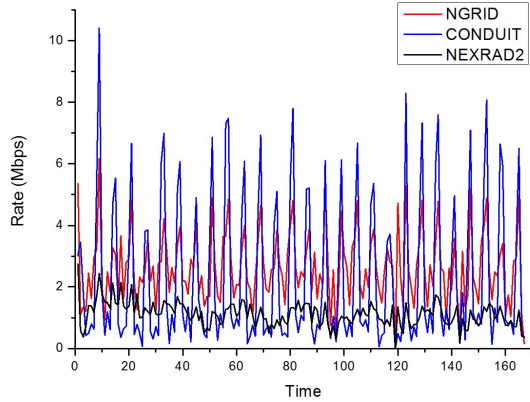


Fig. 2. Average arrival rates of files versus time (hours) of NGRID, CONDUIT and NEXRAD2 feedtype.

Here we choose $n = 5$ features. Description of the features is given in Table I. Our objective is to find the minimum rate where the predicted latency satisfies the delay requirements.

TABLE I
INPUT AND OUTPUT DESCRIPTION

Feature	Description
<i>Input</i>	
X_i^1	Inter-arrival time of file i at sender
X_i^2	File size of file i at sender
X_i^3	Transmission rate
X_i^4	RTT
X_i^5	Number of receivers
<i>Output</i>	
Y_i	Latency of received file i

IV. METHODOLOGY

This section describes how we develop the ML engine to determine the optimal rate. We propose a QoS-driven rate management pipeline, which is composed of two steps in the training stage as shown in Fig. 3. First, based on the different characteristics of file-streams, we train a classifier using MLP to identify its feedtype. Then, combining the original time-series data, network state, and sending rate with the corresponding feedtype, a stacked LSTM based regressor is trained, which predicts the per-file delay. Accordingly, in the testing stage, the incoming file-stream feedtype is first decided by the classifier, and then we apply the trained regressor to predict the delay under various possible rates, and finally we select the sending rate as one that is the minimum rate whose predicted delay of the file-stream satisfies the predefined delay requirement.

A. MLP Classifier

The MLP classifier consists of three parts: feature extraction, training, and testing. The first step is feature extraction.

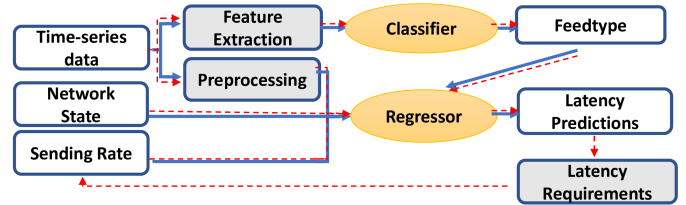


Fig. 3. Overview of the proposed QoS-driven rate management. Blue arrows indicate the process in the training stage, while red dashed arrows indicate the testing process.

We extract a set of seven statistical features of both inter-arrival time and file sizes when we set the holding interval $\tau = 1$ hour. The specific features we use are given in Table II. All the data is normalized with 0 mean and unit variance.

The training data was collected on Nov. 14-19, 2018. Overall, the file streams are divided into 7 classes. We train an MLP model ($128*64*16*7$) based on the Keras deep learning framework. The hyperparameters we choose are decided by the highest validation accuracy, where the validation data is 20% of the total training data. In order to avoid overfitting, we add a 0.2 dropout rate to each layer (except the output layer). Categorical cross-entropy loss is used as the target loss function for training with the Adam optimizer. In addition, we compare the results of other ML baselines with the MLP model, including Linear Discriminant Analysis (LDA), Support Vector Machine (SVM), K-Nearest Neighbors (KNN), Logistic Regression (LogReg), Decision Tree (DT), and Random Forest (RF). In the testing step, we use the data collected on Nov. 20-21, 2018. Then we put the trained model to transform the testing data and display the results.

TABLE II
STATISTICAL FEATURE DESCRIPTION

Feature	Description
max_inter_arr	The maximum value of inter-arrival time
mean_inter_arr	The average value of inter-arrival time
median_inter_arr	The median value of inter-arrival time
min_size	The minimum value of file size
max_size	The maximum value of file size
mean_size	The average value of file size
median_size	The median value of file size

B. Stacked LSTM Regressor

The reason we have chosen LSTM-based regressor is to exploit the potential time correlation among the adjacent files. Intuitively, if the delay of a file-stream is relatively high, it means that the files at the sender-side are either queued or to be retransmitted, and it has a high probability that the delay of these files will also be high. On the other hand, if the delay of a file-stream is relatively low, it means that the network is in low contention, and the file delay is likely to be small in the next time slot. The LSTM method [11] has a strong

ability to handle the entire sequence of data and keep track of the long-term dependencies. Thus, in this work, we train an LSTM network to predict the delay of a file-stream for each feetype. The model then combines different rates with the delay requirements to determine the minimal rate for the next time interval.

Data Preprocessing. The original features are of different scales, which necessitates data preprocessing. However, for the meteorology data, we note that most of the files have small delays and only a small portion of files have large delays, as can be seen from Fig. 4(a). If a common preprocessing method (e.g., normalization and standardization) is applied to all of the data, most will become indistinguishable after preprocessing. In order to avoid this problem, we choose the quantile method that transforms the original data into a uniform distribution. The results are shown in Fig. 4.b. In this case, files with similar delay can be converted to different values to enable a more effective prediction.

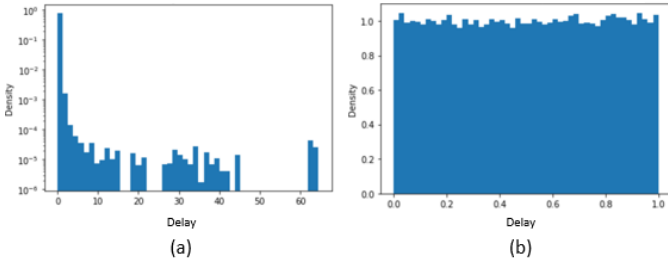


Fig. 4. Delay distribution of training data before and after quantile preprocessing, as shown in (a) and (b), respectively.

LSTM. LSTM is a powerful variant of RNN [12] that can solve the gradient disappearance problem. The LSTM controls the addition or deletion of state information of cells through a module that includes the gating mechanism. In other words, the gate selectively controls the information whether to pass or remove. Furthermore, a *stacked* LSTM architecture can be defined as an LSTM model comprised of multiple LSTM layers, which provides a sequence output rather than a single value output. The stacked LSTM has a more powerful generalization capability and can solve complex problems.

Taking advantage of the stacked LSTM architecture, we build our ML model as shown in Figure 5. Specifically, the first LSTM layer takes time-series features as the input, and the subsequent LSTM layer takes the output from the earlier LSTM layer as new input. This variation of LSTM enables the later layers to capture longer-term dependencies of the input sequence. In order to provide a mapping of the prediction to $[0, 1]$, a dense layer with the sigmoid function is placed on the last LSTM layer to estimate the delay. The optimization objective is to minimize the MAE loss between the output of the regressor and the corresponding file-stream delay.

V. EXPERIMENTS AND RESULTS

A. Experimental Setting

Local Data Manager (LDM) has been the *de facto* software program for real-time meteorological data distribution since

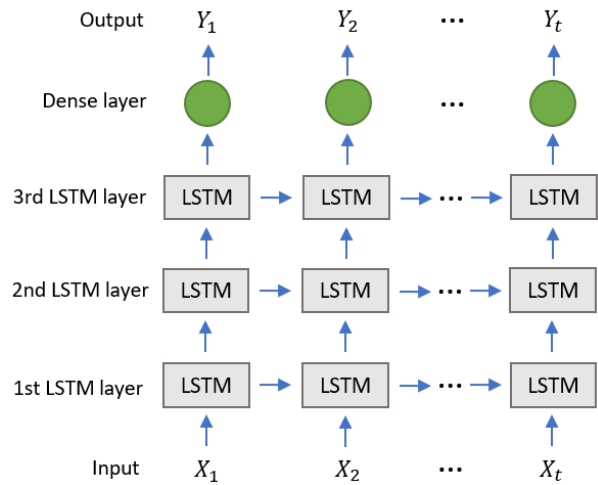


Fig. 5. Architecture of multilayer LSTM model. The gray rectangle LSTM represents an LSTM unit.

1994. We use the latest version LDM7 in the experiment, which has been described in our prior papers [2], [13]. The data has been collected on *Chameleon*, a highly reconfigurable, NSF-supported network testbed [14]. In the experiment, we set the holding interval $\tau = 1$ hour and deliver 7 types of meteorological data. Then we collect the timestamps and file sizes at the sender-side and receiver-side from the log files measured on *Chameleon*, respectively. Based on the ground-truth knowledge about the data, we label feetypes for all the samples in the dataset and calculate the delay for each file.

As the data collection results in successful delivery for all of the files, the retransmission component, which happens frequently in practice, is missing from the dataset. We address this issue by running the experiments on trial deployments. In our previous work, we have deployed LDM7 at eight university campuses throughout the United States that are connected via the regional REN (Research-and-Education Network) and Internet2 [3], as depicted in Fig. 6. Most of the regional RENs serving the corresponding university are connected to the nearest Internet2. We only experiment with the NGRID feetype in the experiments, which is the most representative among all the feetypes. Based on our previous analysis of the network performance with multiple receivers, we notice that when the sending rates are not significant, these receivers tend to have similar throughput as the case of a single receiver. Therefore, the NGRID distribution is started on the UCAR LDM7 publisher. Then the UVA LDM7 subscribers send NGRID subscription requests to the UCAR to receive the files with $\text{RTT} = 40$ ms. We deliver the files with six different sending rates and collect ten-hours NGRID data (approximately three million files) on September 6, 2020.

B. Discrimination Power of Different Feetypes

We apply the MLP method to the feetype classification problem and compare it with the state of the art baselines. Then we run the trained model to transform the test data and

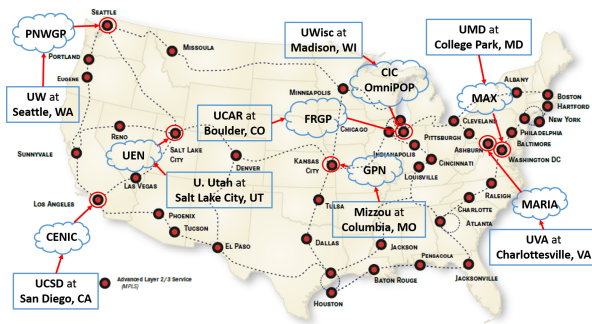


Fig. 6. Trial deployment of LDM7 across Internet2.

display the results, which indicate that the accuracy of our MLP model is 97.7%. The other metrics for different types are shown in Table III, which shows that the MLP model can recognize the feedtype accurately with respect to these metrics.

TABLE III
CLASSIFICATION RESULTS FOR DIFFERENT FEEDTYPES

Feedtypes	Precision	Recall	F1-score
NGRID	0.971	0.892	0.93
CONDUIT	0.925	1.	0.961
NEXRAD2	1.	1.	1.
NEXRAD3	1.	0.973	0.986
DIFAX	1.	1.	1.
HDS	0.947	0.973	0.96
IDS DDPLUS	1.	1.	1.
Average	0.978	0.977	0.977

One practical difficulty with the MLP method is that it is hard to decide how much data to use for training, as the accuracy of this model depends on the amounts of data. Figure 7 compares the accuracy with different training dataset sizes. For a fair comparison, the testing set is fixed in each experiment. From the results, we see that the accuracy increases with the training dataset size, which indicates that more training samples lead to better test accuracy. In comparison, few examples produce a lower testing accuracy, perhaps because the chosen model overfits the training set or the training set is not sufficiently representative of the problem.

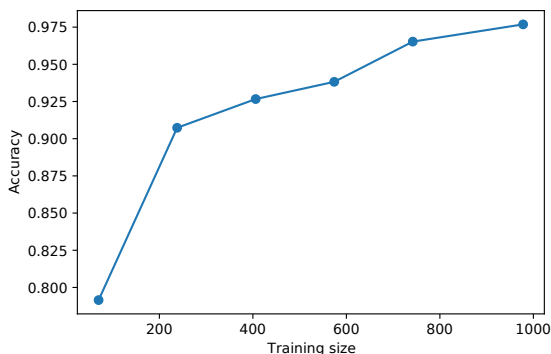


Fig. 7. MLP accuracy with different training dataset sizes.

In addition, we compare the results of other ML baselines with the MLP model. Figure 8 reveals that RF outperforms other methods while LDA has the lowest accuracy. Although MLP does not have the best testing accuracy, we see that it is competitive with other methods. We also note that the performance of MLP can continue to improve as more data becomes available to the model, but the hyperparameters of the model must be adjusted to support the increases in data.

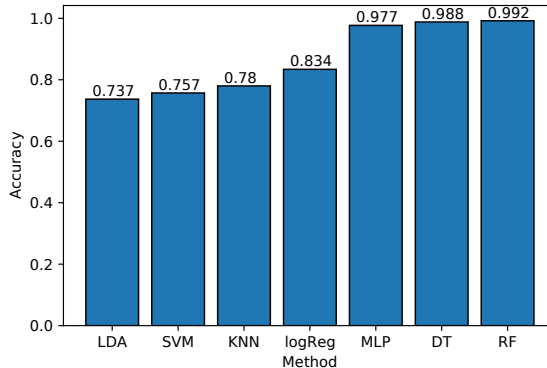


Fig. 8. Accuracy compared with baselines.

C. Performance of Delay Prediction Regressor

After recognizing the feedtype, we train an LSTM regressor with the original time-series data, network state, and sending rate, which predicts the delay of each file in a file-stream. Here we train the regressor with respect to the NGRID feedtype, which is representative among all the feedtypes.

To investigate the impact of the LSTM method, we conduct experiments on different dataset sizes for delay prediction with both LSTM and MLP algorithms. Table IV reports the MAE performance with different datasets. We can see that LSTM outperforms MLP in terms of all the dataset sizes, which indicates that LSTM can keep track of the long-term dependencies across the delay of files. To further visualize the effect of prediction, Figure 9 shows the normalized delay distribution of prediction when the dataset size equals 200. It can be observed that the prediction has almost the same distribution as the original delay.

For our final set of results of finding the optimal rate, we set the QoS requirement as $W = 0.2$ seconds and $\alpha = 8\%$, which means the delay of no more than 8% of files can exceed the threshold of 0.2 seconds. Figure 10 compares the original and predicted α with six different rates, respectively. Blue bars indicate the α of the original delay, where the optimal rate is 100 Mbps. Since the file-stream with rates greater than 100 Mbps (including 100 Mbps) meet the requirements, the smallest value (100 Mbps) should be chosen as the optimal rate in order to save bandwidth. In addition, it is not true that the higher the rate, the smaller α is (i.e., most of the delay of files are smaller than the threshold W), which indicates the relationship of delay and rate is not monotonic. For example, when the rate increases to 400 Mbps, α is greater than 200 Mbps according to the blue bar in Figure 10. The

plausible reason is that packet loss occurs, which increases the delay of most files. We conduct the experiments to verify the performance of our LSTM model in terms of correctly predicting the optimal rate. For the same file-stream, we adjust the rate and have different delay predictions through the LSTM regressor. Based on our selection rules, we end up choosing 100 Mbps as the optimal rate, which is consistent with the true result. To have a deeper understanding, we further notice that the predicted values of α for each rate are similar to the true α . The results suggest that the proposed delay prediction regressor using stacked LSTM has a strong ability for rate selection based on the predefined requirement of QoS. We have comprehensive results considering other settings, e.g. RTT=100ms. Due to the space limit, we will report them in our future work.

TABLE IV
MAE RESULTS FOR LSTM MODEL COMPARED WITH MLP

Set Size	MLP	LSTM
50	0.0534	0.0469
100	0.074	0.0679
150	0.0678	0.0591
200	0.0691	0.0596

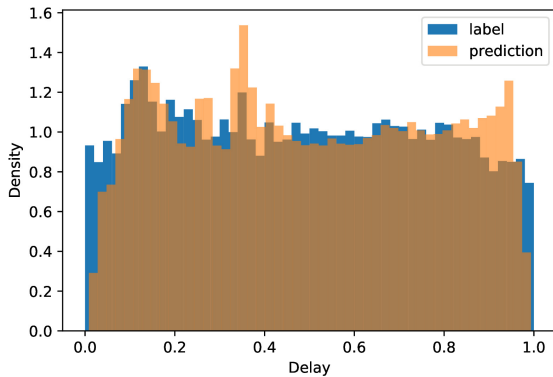


Fig. 9. Normalized delay distribution of prediction.

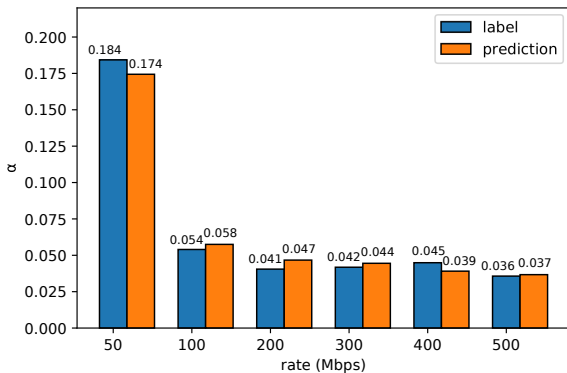


Fig. 10. The α value when threshold $W = 0.2s$. Blue bars indicate the α of original delay, while orange bars represent the predicted α .

VI. CONCLUSION AND FUTURE WORK

In this paper, we have proposed a QoS-driven rate management pipeline to determine the optimal sending rate for a reliable network-multicast application for scientific data transfer. Based on a set of seven features (for the specific meteorological data) such as inter-arrival time and file size, we successfully reached 97.7% accuracy of feetype classification among the six most popular feetypes. We then trained the delay prediction regressor with the meteorological data collected from the trial deployment, which achieved MAE with 0.0469 when the dataset size is 50. More importantly, the LSTM model outperformed the MLP method in terms of the MAE, and demonstrated strong potential to select an optimal rate in real-world networks. In the future, we plan to investigate incorporating more features of the architecture, e.g., more various RTTs, the different sets of receivers, and more feetypes, to train and evaluate the proposed method for a wide range of scientific data use cases.

REFERENCES

- [1] "Unidata collaborative projects." [Online]. Available: <https://www.unidata.ucar.edu/projects/index.html>
- [2] S. Chen, X. Ji, M. Veeraraghavan, S. Emmerson, J. Slezak, and S. G. Decker, "A cross-layer multicast-push unicast-pull (MPUP) architecture for reliable file-stream distribution," in *IEEE 40th Annual Computer Software and Applications Conference (COMPSAC)*, vol. 1, 2016, pp. 535–544.
- [3] Y. Tan, S. Chen, S. Emmerson, Y. Zhang, and M. Veeraraghavan, "Advances in reliable file-stream multicasting over multi-domain software defined networks (SDN)," in *IEEE 28th International Conference on Computer Communication and Networks (ICCCN)*, 2019, pp. 1–11.
- [4] J. Li, M. Veeraraghavan, S. Emmerson, and R. D. Russell, "File multicast transport protocol (FMTP)," in *IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing*, 2015, pp. 1037–1046.
- [5] B. Adamson, C. Bormann, M. Handley, and J. Macker, "Nack-oriented reliable multicast (NORM) transport protocol," *Internet Engineering Task Force (IETF) RFC*, vol. 5740, 2009.
- [6] K. Shvachko, H. Kuang, S. Radia, and R. Chansler, "The hadoop distributed file system," in *IEEE 26th symposium on mass storage systems and technologies (MSST)*, 2010, pp. 1–10.
- [7] J. Wu and B. Hong, "Multicast-based replication for hadoop HDFS," in *IEEE/ACIS 16th International Conference on Software Engineering, Artificial Intelligence, Networking and Parallel/Distributed Computing (SNPD)*, 2015, pp. 1–6.
- [8] M. Mirza, J. Sommers, P. Barford, and X. Zhu, "A machine learning approach to tcp throughput prediction," *ACM SIGMETRICS Performance Evaluation Review*, vol. 35, no. 1, pp. 97–108, 2007.
- [9] Y. Kong, H. Zang, and X. Ma, "Improving TCP congestion control with machine intelligence," in *Proceedings of the 2018 Workshop on Network Meets AI & ML*, 2018, pp. 60–66.
- [10] L. Mei, R. Hu, H. Cao, Y. Liu, Z. Han, F. Li, and J. Li, "Realtime mobile bandwidth prediction using lstm neural network and bayesian fusion," *Computer Networks*, vol. 182, p. 107515, 2020.
- [11] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [12] J. T. Connor, R. D. Martin, and L. E. Atlas, "Recurrent neural networks and robust time series prediction," *IEEE Trans. Neural Netw.*, vol. 5, no. 2, pp. 240–254, 1994.
- [13] X. Ji, Y. Liang, M. Veeraraghavan, and S. Emmerson, "File-stream distribution application on software-defined networks (SDN)," in *IEEE 39th Annual Computer Software and Applications Conference*, vol. 2, 2015, pp. 377–386.
- [14] K. Keahey, P. Riteau, D. Stanzione, T. Cockerill, J. Mambretti, P. Rad, and P. Ruth, "Chameleon: a scalable production testbed for computer science research," in *Contemporary High Performance Computing*. CRC Press, 2019, pp. 123–148.