PACKETLSTM: DYNAMIC LSTM FRAMEWORK FOR STREAMING DATA WITH VARYING FEATURE SPACE

Anonymous authors

Paper under double-blind review

ABSTRACT

We study the online learning problem characterized by the varying input feature space of streaming data. Although LSTMs have been employed to effectively capture the temporal nature of streaming data, they cannot handle the dimensionvarying streams in an online learning setting. Therefore, we propose a dynamic LSTM-based novel method, called packetLSTM, to model the dimension-varying streams. The packetLSTM's dynamic framework consists of an evolving packet of LSTMs, each dedicated to processing one input feature. Each LSTM retains the local information of its corresponding feature, while a shared common memory consolidates global information. This configuration facilitates continuous learning and mitigates the issue of forgetting, even when certain features are absent for extended time periods. The idea of utilizing one LSTM per feature coupled with a dimension-invariant operator for information aggregation enhances the dynamic nature of packetLSTM. This dynamic nature is evidenced by the model's ability to activate, deactivate, and add new LSTMs as required, thus seamlessly accommodating varying input dimensions. The packetLSTM achieves state-of-the-art results on five datasets, and its underlying principle is extended to other RNN types, like GRU and vanilla RNN.

026 027 028

029

025

004

010 011

012

013

014

015

016

017

018

019

021

023

1 INTRODUCTION

Online learning, characterized by streaming data, where data instances arrive one by one, has been studied extensively (Gama, 2012; Neu & Olkhovskaya, 2021; Agarwal et al., 2008). Recently, there has been a growing focus on online learning in environments with varying input feature spaces. Examples include movie sentiment classification and crowdedness severity prediction (He et al., 2023; Agarwal et al., 2024). These varying input features, termed haphazard inputs (Agarwal et al., 2023), are denoted as $X^t \in \mathbf{R}^{d^t}$, where d^t indicates the dimensionality of input, varying over time t. The field of haphazard inputs is expanding, prompting the introduction of new methods, applications, and appropriate datasets as elaborated in section A of the Appendix.

The current landscape is focused on developing new methods. Predominantly, haphazard inputs are modeled using classical approaches like naive Bayes (Katakis et al., 2005), decision stumps (Schreckenberger et al., 2022; 2023), and linear classifiers (Beyazit et al., 2019), favored for their dynamic architectures. However, there is a push towards developing dynamic deep learning solutions (Agarwal et al., 2022; 2023), motivated by the capabilities of neural networks. Nevertheless, current methodologies have not adequately leveraged the streaming nature of data. To bridge this gap, we advocate using Recurrent Neural Networks (RNNs) (Hochreiter & Schmidhuber, 1997; Zhang et al., 2021a), which can effectively exploit the temporal dynamics of data.

We introduce a novel architecture, termed packetLSTM, designed to dynamically adapt to varying input feature space. This framework employs a unique ensemble of Long Short-Term Memory (LSTM) units, each dedicated to a specific input feature. The packetLSTM allows for robust interaction among its LSTMs, fostering the integration of global information while preserving feature-specific knowledge within each unit's short-term memory. The packetLSTM facilitates continuous learning without the risk of catastrophic forgetting in online learning environments (Hoi et al., 2021). The feature-based learning, coupled with a dimension-invariant aggregation operator, allows packetLSTM to dynamically activate, deactivate, and add new LSTMs as needed, leading to adeptly managing haphazard inputs. The main contributions of our work are as follows. (1) We intro-

054 duce the first RNN-based framework, packetLSTM, to effectively handle haphazard inputs in online 055 learning. (2) The packetLSTM exhibits learning without forgetting capabilities in an online learning setting. We demonstrate this capability in a challenging scenario where certain features are absent 057 for extended time periods. (3) The principles underlying the packetLSTM framework are adapt-058 able and can be extended to other types of RNNs, like Gated Recurrent Units (GRUs) and vanilla RNNs. We substantiate this adaptability by developing packetGRU and packetRNN models. (4) We achieve state-of-the-art results across five datasets, as shown in Figure 1. (5) We introduce a strong 060 baseline based on the Transformer called HapTransformer and demonstrate that HapTransformer 061 outperforms other baselines; however, it is still inferior to packetLSTM. 062

063 064

065

2 RELATED WORKS

066 Haphazard Inputs The initial approach 067 to address haphazard inputs utilized naive 068 Bayes and χ^2 statistics to dynamically 069 incorporate features, update existing feature statistics, and select a feature subset 070 (Katakis et al., 2005). This method was fur-071 ther expanded by employing an ensemble of 072 naive Bayes classifiers for predictive anal-073 ysis (Wenerstrom & Giraud-Carrier, 2006). 074 Subsequent research has focused on in-075 ferring unobserved features from observed 076 ones using various techniques, including 077 graph methods (He et al., 2019; Sajedi & Razzazi, 2024), and Gaussian copula (He 079 et al., 2021; Zhuo et al., 2024), followed 080 by the application of classifiers across the complete feature space. Concurrently, an-081 other line of research projects data into a shared subspace to learn a linear classifier 083 using empirical risk minimization (Beyazit



Figure 1: Performance comparison of packetLSTM with other models. The legend and outer labels are methods and datasets, respectively. Exact balanced accuracy values are provided in Table 1.

084 et al., 2019) or online gradient descent (Zhou & Matsushima, 2023). Distinctly, You et al. (2024) 085 maintains an informativeness matrix of each feature to update a linear classifier. Another research direction explores the use of decision trees to handle haphazard inputs. Specifically, Schreckenberger et al. (2022) proposed Dynamic Forest, which uses an ensemble of decision stumps, each 880 based on a feature from a selected subset of all seen features. However, this approach can result in 089 numerous decision stumps, prompting Schreckenberger et al. (2023) to introduce a refined approach 090 that constructs only one decision stump for each feature. Lee et al. (2023a) proposed to utilize adaptive random forest on a fixed set of features, created through imputation assuming large buffer 091 storage. Despite the dynamic nature of the above classical methods in modifying their architectures, 092 the era of big data necessitates adopting deep learning approaches to effectively model haphazard 093 inputs. To date, seminal contributions in this domain include works by Agarwal et al. (2022; 2023), 094 which are based on neural networks. The Auxiliary Network (Agarwal et al., 2022) incorporates 095 parallel hidden layers to process each input feature. Conversely, Aux-Drop (Agarwal et al., 2023) 096 implements selective and random dropouts within its layers to handle haphazard inputs. Although these models operate under specific assumptions, recent advancements by the same authors (Agar-098 wal et al., 2024) provide simple solutions to mitigate these assumptions in haphazard inputs. While 099 all the methods discussed above handle haphazard inputs, none effectively exploits the temporal dy-100 namics of streaming data, which we achieve using RNNs, specifically LSTMs. The comparison of 101 haphazard inputs with other varying feature space fields is presented in section V of Appendix.

102

RNNs RNNs are among the most popular models in sequential data analysis (Salehinejad et al., 2017; Allen-Zhu & Li, 2019). Various techniques are developed to capture the temporal dynamics of data, including stacking LSTMs to establish a hierarchical framework (Hermans & Schrauwen, 2013; Wang et al., 2017), as well as designing specialized time-modeling LSTM units (Che et al., 2018; Kazemi et al., 2019). In this article, we utilize Time-LSTM (Zhu et al., 2017) as our LSTM units because of its demonstrated capability in modeling both short-term and long-term interest.

108 LSTMs have also been extensively used in the field of multi-modality (Xie & Wen, 2019; Liu et al., 2018; Xu et al., 2020; Lee et al., 2023b). However, to the best of our knowledge, the RNN-based 110 models proposed in the multi-modal domain or any other domain are not capable of modeling data 111 with varying input dimensions in an online learning setting as discussed in section B of the Ap-112 pendix. Therefore, we propose a new dynamic RNN framework in this article, filling a significant gap in the research landscape of RNNs and haphazard inputs. 113

- 3
- 115 116

114

PRELIMINARIES

117 **Notations** In this article, we represent time in superscript and feature id in subscript. For example, 118 x_2^3 represents the value of feature 2 (F₂) at time t₃. For ease of readability, we slightly adjust the 119 notation of time here. Instead of denoting by $v_i^{t_1}$, we use v_i^1 . When time is referenced individually, 120 it is denoted as t_1 . Moreover, t indicates a random time, and t-1 denotes the time preceding t. A 121 complete list of notations is provided in section C of the Appendix.

122

123 **Characteristics** Haphazard inputs exhibit six characteristics which are illustrated in Figure 2(b). 124 These characteristics are: (1) Streaming data, which are received sequentially and processed with-125 out storage. (2) Missing data, which are features present in some instances but may be absent in subsequent ones like F_1 at time t_2 . (3) Missing features, that are not received from the onset; how-126 ever, their availability is known like F_4 . (4) Sudden features, that arrives unexpectedly without prior 127 indication of their existence like F_3 at time t_2 . (5) Obsolete features, which can cease to exist after 128 any instance, such as F_2 . (6) Unknown number of total features, results from the combined effect of 129 missing data, missing features, sudden features, and obsolete features. 130

131 **Feature Space** The characteristics of haphazard inputs result in a varying feature space. We define 132 a universal feature space (\mathbb{F}^t) as the set of features encountered till time t. The specific set of features 133 present at time t is represented by \mathbb{F}^t and is termed current feature space as shown in Figure 2(b). 134 The universal feature space will grow with time due to the emergence of sudden features and missing 135 features. For example, $\mathbb{F}^1 = \{F_1, F_2\}$ at time t_1 , and $\mathbb{F}^2 = \{F_1, F_2, F_3\}$ at time t_2 in Figure 2(b). 136 In an ideal condition, \mathbb{F}^t can contract with the removal of obsolete features; however, since the 137 cessation of obsolete features is unknown, $\overline{\mathbb{F}}^t$ may not decrease in practice.

138

139 **Mathematical Formulation** The haphazard input received at time t can be represented by X^t , 140 where $X^t \in \mathbf{R}^{|\mathbb{F}^t|}$. Here, $|\cdot|$ represents the cardinality of a set. The corresponding ground truth is 141 denoted by y^t , where $y^t \in [0,1]^C$ and C is the number of classes. This paper deals with the binary 142 classification problem but can be easily extended to multi-class scenarios. The model, denoted by f, operates in an online learning setting with $f^{t-1}: X^t \to y^t$, where f^0 represents the initialized 143 state of the model. After processing $\{X^t, y^t\}$, the model's state is represented by f^t . At each time 144 t, the model receives X^t , and f^{t-1} processes X^t to yield a prediction \hat{y}^t . Upon the revelation of y^t , 145 the loss $l^t = H(y^t, \hat{y}^t)$ is computed, where H is a loss function. The model then updates from f^{t-1} 146 to f^t for the subsequent instances based on l^t . This iterative process continues for each instance. 147

4 METHOD

149 150

148

151 The packetLSTM consists of a pack of LSTMs, each dedicated to a distinct feature, as illustrated 152 in the gray box in Figure 2(a). We utilize LSTMs due to their proven effectiveness in capturing 153 temporal dynamics of data (Kazemi et al., 2019; Wang et al., 2017; Zhang et al., 2021a).

154 Due to the varying dimensions of input feature space, a single LSTM cannot process all features, 155 necessitating one LSTM per feature. Each LSTM (L_j) receives inputs comprising the feature value 156 (x_i^t) , time delay (Δ_i^t) , its previous short-term memory $(h_i^{t_-})$, and common long-term memory 157 (c^{t-1}) , as labeled *Input* in Figure 2(a). The Δ_i^t measures the time difference between the current 158 time t and the last observed time t_{-} for feature j. This temporal information is critical to the model 159 because feature availability varies (Zhu et al., 2017; Che et al., 2018). 160

At each instance, only LSTMs corresponding to available features are activated. For example, at time 161 t_2 , the absence of feature F_1 results in the deactivation of LSTM L_1 , depicted by a sketched gray



Figure 2: (a) The packetLSTM architecture based on (b) the data snapshot.

box in Figure 2(a). Thus, packetLSTM dynamically handles all the characteristics of the haphazard inputs by activating, deactivating, or adding new LSTMs as needed.

The common long-term memory (c^t) , which aggregates information from all long-term memories of active LSTMs, facilitates the interaction among features, crucial for enriched knowledge acquisition (Zhang et al., 2021b). This c^t contains global information, while the short-term memory of each LSTM retains the local information about individual features, alleviating the issue of catastrophic forgetting in an online learning setting. The aggregation operator, being dimension-invariant, accommodates the variable number of features. All active short-term memories are aggregated to generate a predictive short-term memory (h^t) . These memories are subsequently concatenated and processed through a fully connected classifier to produce the final output \hat{y}^t (see Figure 2(a)).

The model parameters are updated based on the loss $l^t = H(y^t, \hat{y}^t)$ in an online learning setting. New LSTMs are introduced with initialized memories (h^0, c^{t-1}) and a time delay (Δ) set to zero for sudden (e.g., F_3 at t_2) and missing features (e.g., F_4 at t_3), as shown in Figure 2.

208

192 193

Working Principle Based on the current feature space \mathbb{F}^t at time t, the output of each LSTM is given by $h_j^t, c_j^t = L_j(x_j^t, \Delta_j^t, h_j^{t-}, c^{t-1}) \forall F_j \in \mathbb{F}^t$. Numerous LSTM variants exist in the literature that model time delay (Δ_j^t) , such as decay in GRU-D (Che et al., 2018), time gate in Phased LSTM (Neil et al., 2016), and Time2Vec (Kazemi et al., 2019). In this article, we employ Time-LSTM (Zhu et al., 2017) because of its highly demonstrated capability in modeling time information. Time-LSTM utilizes time delay to capture the short-term interest for current instances and preserves these delays to address long-term interest for future instances. We specifically utilize Time-LSTM 3 (T-3) among its three versions – Time-LSTM 1, Time-LSTM 2, and T-3 – because it integrates the input and forget gates (Greff et al., 2016), offering a model that is both less computationally intensive and concise. The formulation of T-3 (and other time modeling variants) for feature j at time t, within the packetLSTM framework, represented by L_j , is discussed in the section D of the Appendix.

219 Our packetLSTM framework differs from T-3 in its approach to handling input feature spaces. Un-220 like T-3, which utilizes a single LSTM unit for a fixed feature space, our approach employs a distinct 221 LSTM unit for each feature to manage a varying feature space. This necessitates a different input 222 configuration for each LSTM compared to the T-3. Specifically, rather than using the previous long-223 term memory of each LSTM as an input, the packetLSTM framework utilizes a common long-term 224 memory, which facilitates interaction among features. Moreover, while T-3 accepts the short-term 225 memory from time t-1 as input at time t, the packetLSTM inputs the short-term memory from 226 time t_{-} at time t, where t_{-} is not necessarily equal to t - 1 in haphazard inputs.

Next, a dimension-invariant aggregation operator (*AGG*) determines the common long-term memory and the predictive short-term memory as

$$c^{t} = AGG(\bigcup_{j,\forall F_{j}\in\mathbb{F}^{t}} \{c_{j}^{t}\}), \qquad h^{t} = AGG(\bigcup_{j,\forall F_{j}\in\mathbb{F}^{t}} \{h_{j}^{t}\}).$$
(1)

The AGG operator, defined as $AGG : [\mathbf{R}]^{|F^t|} \to \mathbf{R}$, accepts a variable number of inputs, specifically, $|F^t|$ inputs. Given that c_j^t and h_j^t are vectors, the AGG operator performs aggregation elementwise. Common examples of dimension-invariant aggregation operators include mean, maximum, minimum, and summation. Finally, the prediction \hat{y}^t is generated by a fully connected neural network (FCN), applied on the concatenation of c^t and h^t as $\hat{y}^t = \text{FCN}(concat(c^t, h^t))$.

237 238

239

230 231

5 EXPERIMENTS

240 Datasets We consider 5 datasets – magic04 (Bock et al., 2004), imdb (Maas et al., 2011), a8a 241 (Kohavi et al., 1996), SUSY (Baldi et al., 2014), and HIGGS (Baldi et al., 2014) - with details 242 provided in section E of the Appendix. The motivation to choose these datasets is three-fold. First, they include both real (imdb) and synthetic datasets. Second, the number of instances varies from 243 19020 in magic04 to 1M in HIGGS. Third, the number of features ranges from 8 in SUSY to 7500 244 in imdb. Therefore, the diversity in the number of features and instances allows us to determine the 245 efficacy of packetLSTM effectively. The imdb dataset is haphazard in nature. However, the synthetic 246 dataset needs to be transformed into haphazard inputs. Following the baseline papers (Beyazit et al., 247 2019; Agarwal et al., 2023), we transform synthetic datasets based on the probability values p, where 248 p = 0.25 means only 25% of features are available at each time instance. We consider p = 0.25, 249 0.5, and 0.75. The synthetic dataset preparation is further discussed in section E of the Appendix. 250

Metrics We compare all models using five metrics: number of errors, accuracy, Area Under the Receiver Operating Characteristic curve (AUROC), Area Under the Precision-Recall Curve (AUPRC), and balanced accuracy. Each metric is discussed in section F of the Appendix. The balanced accuracy is the primary comparison metric in the main manuscript, while detailed comparisons using the other metrics are presented in section K of the Appendix. We adhere to the standard evaluation protocol for haphazard inputs, which is detailed in section G of the Appendix.

257

Baseline We consider 10 baseline models, including NB3(Katakis et al., 2005), FAE (Wenerstrom & Giraud-Carrier, 2006), DynFo (Schreckenberger et al., 2022), ORF³V (Schreckenberger et al., 2023), OLVF (Beyazit et al., 2019), OCDS (He et al., 2019), OVFM (He et al., 2021), Aux-Net (Agarwal et al., 2022), Aux-Drop (Agarwal et al., 2023), and OLIFL(You et al., 2024). Additionally, there are a few other models – see section H of the Appendix – applicable to the field of haphazard inputs. However, we could not include these models because of the lack of open-source code and the challenges associated with implementing them.

Implementation Details We ran all the models five times, except for the deterministic models –
 NB3, FAE, OLVF, and OLIFL – which consistently produce identical outcomes across runs. The results are reported as mean ± standard deviation (std). For packetLSTM, max and mean aggregation operators were used in synthetic and real datasets, respectively. Additionally, Z-score normalization was implemented in an online manner, as discussed in section 6. The implementation details and hyperparameter search are discussed in the section I and J of the Appendix, respectively.

270Table 1: Comparison of models on all datasets based on balanced accuracy. The deterministic mod-271els — NB3, FAE, OLVF, and OLIFL — underwent a single execution, and the non-deterministic272models were executed 5 times, with the mean \pm standard deviation reported. A [‡] symbol indicates273non-deterministic models that were run only once on specific datasets due to substantial time con-274straints, and [†] denotes the real datasets. % Im. denotes the % improvement in the performance of275packetLSTM compared to the previous best-performing baseline (denoted by *italics*) in each dataset.

Dataset	p	NB3	FAE	OLVF	OLIFL	OCDS	OVFM	DynFo	ORF ³ V	Aux-Net	Aux-Drop	packetLSTM
:-04	0.25	50.01	50.01	53.18	53.06	51.89±0.10	51.94±0.00	52.75±0.30	47.94±0.22	50.09±0.07	56.04±0.53	61.33±0.07
magic04	0.30	49.99	50.00	56.19	60.75	53.40 ± 0.43 53.76 ± 1.07	54.13 ± 0.08 58.79 ± 0.04	55.12 ± 0.06 56.75 ± 0.02	48.36 ± 0.11 49.32 ± 0.04	50.09 ± 0.03 50.05 ± 0.07	63.18 ± 0.61	$73.64{\pm}0.11$
imdb†		81.56	82.18	80.08	54.36	$50.14{\pm}0.02$	77.43 [‡]	$57.98{\pm}0.29$	76.47±0.11	67.41 [‡]	$73.10{\pm}0.19$	$85.06{\pm}0.04$
a8a	0.25 0.50 0.75	50.01 50.01 50.01	50.00 50.00 50.00	60.67 66.46 70.60	53.57 54.63 56.56	$\begin{array}{c} 54.75{\pm}0.87\\ 64.04{\pm}1.01\\ 68.81{\pm}1.10\end{array}$	$\begin{array}{c} 58.66{\pm}0.00\\ 66.02{\pm}0.00\\ 70.95{\pm}0.00\end{array}$	$\begin{array}{c} 50.01{\pm}0.03\\ 50.11{\pm}0.01\\ 50.13{\pm}0.01\end{array}$	$\substack{49.99\pm0.00\\50.01\pm0.00\\49.99\pm0.00}$	$\begin{array}{c} 50.00{\pm}0.00\\ 50.00{\pm}0.00\\ 50.00{\pm}0.00\end{array}$	$\begin{array}{c} 50.00{\pm}0.01\\ 55.33{\pm}1.99\\ 62.87{\pm}0.93\end{array}$	60.53±0.16 67.73±0.21 71.11±0.16
SUSY	0.25 0.50 0.75	50.00 50.00 50.00	49.90 50.01 50.12	51.12 53.21 55.98	51.23 53.59 56.39	$\begin{array}{c} 52.11{\pm}0.19\\ 54.03{\pm}0.28\\ 54.84{\pm}0.48\end{array}$	$\begin{array}{c} 58.00{\pm}0.00\\ 62.85{\pm}0.00\\ 68.51{\pm}0.00\end{array}$	$\begin{array}{c} 54.69{\pm}0.01\\ 58.27{\pm}0.00\\ 60.94{\pm}0.01\end{array}$	$\substack{49.37 \pm 0.01 \\ 48.33 \pm 0.02 \\ 47.53 \pm 0.03 }$	$\begin{array}{c} 50.53{\pm}1.17\\ 57.89{\pm}7.19\\ 53.67{\pm}8.13\end{array}$	61.98±0.10 68.79±0.14 73.55±0.11	$62.77{\pm}0.01$ $69.28{\pm}0.01$ $73.85{\pm}0.02$
HIGGS	0.25 0.50 0.75	50.00 50.00 50.00	50.16 50.01 50.55	50.57 51.21 51.98	50.56 51.50 52.48	49.97±0.07 50.06±0.06 49.97±0.05	50.61±0.01 51.40±0.01 52.66±0.00	50.18 [‡] 50.21 [‡] 50.16 [‡]	49.86±0.03 49.82±0.02 49.75±0.03	49.99 ± 0.00 49.99 ± 0.01 49.98^{\ddagger}	51.17±0.05 53.09±0.05 55.55±0.11	52.22±0.04 55.47±0.04 58.71±0.05

Results The balanced accuracy (b) of all the models is presented in Table 1. The packetLSTM has a statistically significant better performance than all the models across each dataset except in a8a with p = 0.25 and 0.75 (hypothesis supported by paired t-tests at 99% significance level). To quantify the performance enhancement of packetLSTM over the previously best-performing baseline (*pbb*), we calculate the % improvement as $(b_{packetLSTM} - b_{pbb}) * 100/b_{pbb}$, and these results are shown in the last column of Table 1. Among synthetic datasets, packetLSTM surpasses the *pbb* by at least 9.44% and 2.05% in magic04 and HIGGS, respectively. Moreover, packetLSTM consistently exhibits superior performance across all p within the SUSY dataset. In a8a, packetLSTM achieves the best result in 2 out of 3 cases, except at p = 0.25, where its performance is comparable. In the real dataset (imdb), packetLSTM outperforms the *pbb* by a margin of 3.5%.

6 ABLATION STUDIES

We conduct ablation studies within packetLSTM to assess the impact of each component and iden-tify optimal variants, adhering to the hyperparameters described in Implementation Details of section 5, unless specified otherwise. We report the mean of the balanced accuracy in the main manuscript (Table 2), with std detailed in section L of the Appendix. We calculate the number of wins for each component across all dataset combinations, defining a win as achieving the highest balanced accu-racy or being within a 0.05 margin of the top value, accounting for variability indicated by the std. For example, in the magic04 dataset at p = 0.25, the packetLSTM with the Min aggregator reports a slightly higher balanced accuracy (61.36) than the Max (61.33) but exhibits greater std (0.14 vs. 0.07). This variability suggests that some runs using the Min aggregator achieved lower balanced accuracy than the Max. Therefore, a 0.05 margin mitigates such discrepancies.

Which aggregation operator performs the best? The Max operator, in general, performed best with 7 wins, as shown in the AGG component of Table 2. However, to determine the best operators, we compare them one by one (see \rightarrow in Table 2). (1) The performance difference between Min and Max is negligible, indicating both are suitable choices. (2) Between Sum and Mean, there is no evident superiority. However, the Sum is sensitive to the number of features, while the Mean remains relatively stable unless the feature distribution changes substantially. We tested this by al-ternating p (from 0.25 to 0.75) every 100 instances in the a8a dataset, which has the most features among the synthetic datasets. The Mean operator (balanced accuracy = 62.77) substantially out-performed the Sum operator (54.16), establishing the Mean as the superior operator. (3) Table 2 indicates that the efficacy of the Mean increases as data availability decreases from p = 0.75 to 0.25. Further experiments at p = 0.05 and 0.95, shown in section M of the Appendix, validate this trend. Consequently, Max is recommended for scenarios with high data availability, while the Mean is preferable otherwise. This is further supported by the imdb dataset (p = 0.0165), where the Mean (85.06) significantly surpassed the Max (77.90). Consequently, we utilized the Max aggregator for all the synthetic datasets and the Mean operator for the real dataset.

. di	Variants		magic04		imdb		a8a			SUSY			HIGGS		w
Con		25	50	75		25	50	75	25	50	75	25	50	75	
	Mean	60.98	67.06	72.30	85.06	58.52	65.15	67.50	63.16	69.59	74.07	52.01	55.27	58.04	3
5	Sum	61.27	67.93	73.21	83.63	58.26	65.15	68.80	63.19	69.29	74.39	52.24	55.64	59.06	5
AC	Min	61.36	68.28	73.58	77.64	60.53	67.78	71.13	62.78	69.29	73.85	52.21	55.50	58.70	6
	Max	61.33	68.31	73.64	77.90	60.53	67.73	71.11	62.77	69.28	73.85	52.22	55.47	58.71	7
\rightarrow	Min - Max	0.03	-0.03	-0.06	-0.26	0.00	0.05	0.02	0.01	0.01	0.00	-0.01	0.03	-0.01	
\rightarrow	Sum - Mean	0.29	0.87	0.91	-1.43	0.30	0.00	1.30	0.03	-0.30	0.32	0.03	0.37	1.02	
\rightarrow	Mean-Max	-0.35	-1.25	-1.34	7.16	-2.01	-2.58	-3.61	0.39	0.31	0.22	-0.21	-0.2	-0.67	
еЪ	None	58.83	65.12	69.87	85.18	58.77	66.52	70.30	62.30	68.70	73.25	51.37	54.05	57.06	0
oq	Decay	58.73	64.99	69.81	85.30	58.82	67.14	70.76	62.37	68.89	73.40	51.41	54.09	57.01	1
Σ	TimeLSTM-1	61.13	68.14	73.63	85.06	60.45	67.63	71.09	62.77	69.33	73.91	52.23	55.53	58.70	8
me	TimeLSTM-2	61.32	68.34	73.65	85.15	60.33	67.62	71.09	62.79	69.31	73.88	52.26	55.53	58.65	9
	TimeLSTM-3	61.33	68.31	73.64	85.06	60.53	67.73	71.11	62.77	69.28	73.85	52.22	55.47	58.71	10
at	Universal	61.37	68.33	73.69	85.06	60.32	67.58	71.15	62.79	69.27	73.84	52.26	55.52	58.69	11
0 Eu	Current	61.33	68.31	73.64	85.06	60.53	67.73	71.11	62.77	69.28	73.85	52.22	55.47	58.71	13
H	Only LTM	60.94	67.66	73.09	85.15	59.36	66.86	70.34	62.78	69.32	73.85	52.20	52.62	57.69	5
ũ	Only STM	61.15	68.26	73.38	85.17	60.02	67.57	70.77	62.82	69.33	73.89	52.18	52.19	57.63	5
ĉ	Both	61.33	68.31	73.64	85.06	60.53	67.73	71.11	62.77	69.28	73.85	52.22	55.47	58.71	12
	None	59.14	65.44	70.53	85.61	59.97	67.89	71.22	63.12	69.65	74.19	51.46	54.33	57.39	4
0 N	Min Max	57.69	63.33	68.42	85.23	59.94	67.76	71.15	62.37	69.06	73.61	51.21	53.10	55.52	0
Ĺ j	Decimal Scaling	55.03	60.36	65.72	85.23	49.99	50.00	50.00	50.96	58.53	63.31	50.00	50.00	50.00	0
Smc	Mean Norm	60.23	66.71	71.86	84.89	58.03	65.94	70.07	63.19	69.68	74.31	51.86	54.88	58.11	3
lor	Unit Vector	50.01	54.23	60.02	85.32	58.12	65.67	68.97	56.52	65.43	71.44	50.75	52.86	55.53	0
4	Z-score	61.33	68.31	73.64	85.06	60.53	67.73	71.11	62.77	69.28	73.85	52.22	55.47	58.71	7
	Vanilla RNN	60.70	66.76	72.14	83.55	51.16	65.24	64.50	62.84	68.06	73.98	51.96	54.26	56.77	7
NN	GRU	60.34	66.24	71.23	81.49	50.96	60.86	73.38	62.73	69.16	73.72	51.80	54.25	56.79	3
щ	LSTM	58.83	65.12	69.87	85.18	58.77	66.52	70.30	62.30	68.70	73.25	51.37	54.05	57.06	4

Table 2: Ablation study on the components (Comp.) of packetLSTM. Here, w stands for the number of wins a variant registered out of all the 13 dataset combinations. The *variants* marked in bold and italics are employed in packetLSTM. \rightarrow denotes the difference between two variants.

Is time modeling beneficial? We investigated 4 variants of time modeling: T-3, Time-LSTM 2, Time-LSTM 1, and Decay (Che et al., 2018). As observed in the Time Model component, T-3 performed the best. Moreover, LSTM, without any time modeling (*None*), performed the worst, affirming the beneficial impact of time modeling for haphazard inputs. Moreover, all the 4 variants outperform the *pbb* in 11 out of 13 scenarios except in a8a (p = 0.25) and SUSY (p = 0.75).

Are current feature space sufficient for the predictive capability of packetLSTM? The packetLSTM utilizes current feature space for the predictive short-term memory. However, it is also possible to consider the universal feature space by aggregating the short-term memories from all LSTMs linked to the universal features for the predictive short-term memory. The current feature space has 13 wins compared to 11 wins of universal feature space (see the Feat component). Therefore, the current feature space is sufficient for enhancing the predictive capability of packetLSTM.

360 What is the importance of common long-term and predictive short-term memory? The final 361 prediction in packetLSTM concatenates common long-term and predictive short-term memory. We 362 assessed the importance of each memory type by comparing the original packetLSTM, termed *Both*, 363 with its variants that either exclude predictive short-term memory (*Only LTM*) or common long-364 term memory (*Only STM*), as shown in the Concat component. *Both* significantly outperformed 365 others, securing 12 wins, confirming that both memory types are crucial. While common long-term 366 memory holds global information, predictive short-term memory provides local information, and 367 individually, they surpass *pbb* in most datasets except in a8a and HIGSS (p = 0.5).

Does the data need to be normalized? We observed significant variation in the range of feature values across datasets. For example, in magic04, the first feature ranges from 4.28 to 334.18. Therefore, we performed streaming normalization by using five methods: Min-Max, Decimal Scaling, Mean-Norm, Unit Vector, and Z-score, as detailed in section N of the Appendix. Table 2 shows that Z-score performed the best, underscoring the need for streaming normalization in handling haphazard inputs. Remarkably, the packetLSTM, even without any normalization (labeled as *None*), surpassed *pbb* in all datasets except for a8a (p = 0.25), showcasing the superiority of packetLSTM.

375

324

325

326

327 328

348

349

350

351

352

359

Is the concept of packet architecture adaptable? The packetLSTM framework's principles are adaptable, leading us to develop packetRNN and packetGRU architectures. We outline these architectures' formulations and performance in section O of the Appendix and the RNN component

378 of Table 2, respectively. Both packetRNN and packetGRU deliver competitive performances com-379 pared to *pbb*. Specifically, packetRNN outperforms *pbb* in the magic04, imdb, SUSY, and HIGGS 380 datasets, while packetGRU excels in magic04, a8a (p = 0.75), SUSY, and HIGGS. These results 381 confirm the effectiveness of the packet architecture. We evaluate packetLSTM without any time 382 modeling component for a fair comparison among packet architectures. The packetRNN emerges as the best performer, although it significantly underperforms in the a8a dataset. Consequently, we use 383 LSTM for the packet architecture. Furthermore, LSTM's ability to incorporate long-term memory 384 allows it to encode global information, a characteristic absent in both GRU and vanilla RNN. 385

386

Complexity Analysis The time complexity of packetLSTM is $\sum_{t=1}^{T} O(g(|\mathbb{F}^t|) * s^2 + P)$, where T is the number of instances, $O(g(|\mathbb{F}^t|) * s^2)$ denotes the complexity for processing all activated 387 388 LSTMs at time t, and O(P) encompasses constant-time operations such as aggregations and final 389 predictions. Here, s is the hidden size of the LSTM, and the function $q(|\mathbb{F}^t|)$ varies between 1 and 390 $|\mathbb{F}^t|$, often close to 1. Thus, the packetLSTM scales well in terms of time complexity with feature 391 size. The space complexity at time t is $O(|\overline{\mathbb{F}^t}| * L + K)$, where O(L) is the space complexity 392 of one LSTM and O(K) accounts for the final prediction network and aggregated memories. The 393 introduction of sudden and missing features increases the space complexity, showing potential scal-394 ability limits. However, packetLSTM easily handles the 7500 features in the imdb dataset, affirming 395 its capability to handle high-dimensional data. Detailed calculations of time and space complexities 396 and a strategy to mitigate space complexity are provided in section P and Q of the Appendix. 397

398 **packetLSTM vs Single LSTM** The effectiveness of packetLSTM is further evidenced through its 399 comparison with a single LSTM model, which necessitates the availability of all features. This re-400 quirement can be met through imputation, although it requires contradicting the sixth characteristic 401 of haphazard inputs where the total number of features is unknown. Nonetheless, for the purpose of 402 evaluating packetLSTM, we compare it against a single LSTM, employing three imputation techniques: forward fill, mean of the last five observed values of a feature, and Gaussian copula (Zhao 403 et al., 2022). Table 3 demonstrates that packetLSTM significantly outperforms single LSTM with all 404 imputation-based techniques across all datasets, except in a8a at p = 0.75. Further details on the sin-405 gle LSTM model, hyperparameter search, the Gaussian copula model, and its specific performance 406 issues with the a8a and imdb datasets are discussed in section R of the Appendix. 407

408 409

410

7 CHALLENGING SCENARIOS

411 We design three challenging experiments to explicitly elucidate the efficacy of packetLSTM to (1) 412 handle sudden features, (2) handle obsolete features, and (3) demonstrate the learning without forgetting capability. We considered the HIGGS and SUSY datasets for these experiments due to their 413 substantial number of instances. In the main manuscript, we discuss the results corresponding to 414 HIGGS, while similar conclusions for SUSY are discussed in section S.2 of the Appendix. For 415 comparative analysis, we opted for OLVF, OLIFL, OVFM, and Aux-Drop baselines based on their 416 superior performance in the HIGGS dataset, as indicated in Table 1. We divided the dataset into 5 417 intervals, each consisting of 20% of the total instances, with successive intervals containing the next 418 20% of instances. For the exact values of Figure 3 and 4, refer to section S.1 of the Appendix.

419 420

Sudden Features Here, each subsequent interval includes an additional 20% of features, starting
with 20% in the first interval and increasing to 100% by the fifth interval. This progression, termed
the trapezoidal data stream (Zhang et al., 2016; Liu et al., 2022), is illustrated in Figure 3 by a pinkshaded region. Model performance is expected to enhance with increased data volume. All models,
except OLVF, exhibit this increasing performance trend. Notably, packetLSTM outperforms other
models in each interval, demonstrating its superior capability in handling sudden features.

426

427 Obsolete features Here, all features are initially present in the first interval, then only the first 80%
 428 remain in the second interval, decreasing sequentially as shown in Figure 3. Intuitively, the model's
 429 performance should deteriorate as data volume decreases, a pattern observed in all models except
 430 Aux-Drop. The Aux-Drop's performance drops significantly in the second interval, suggesting a
 431 misleading impression of improvement in the third. The packetLSTM outperforms all other methods in each interval, demonstrating its superior ability to handle obsolete features.



Figure 3: Model performance in the scenario of sudden and obsolete features. The mean balanced accuracy in each data interval (e.g., 0 - 0.2 M) is shown in its middle (0.1 M). The pink-colored y-axis represents the Features. For e.g., '1-4' means features 1 to 4 are available in instances shaded with pink color. Both y-axes are shared by the two graphs.



Figure 4: Mitigating catastrophic forgetting: Models performance in reappearing features scenario and its capability of learning without forgetting. The mean balanced accuracy is depicted in the middle of each data interval. Please refer to section S.1 of the Appendix for more details.

463 **Reappearing** We assess the model's performance in a scenario where feature sets are disjoint across consecutive intervals and reappear after a gap, as shown in Figure 4. The first 50% of fea-464 tures are present in the first, third, and fifth intervals, while the remaining 50% appear in the second 465 and fourth intervals. Ideally, effective learning without forgetting would result in improved perfor-466 mance upon the reappearance of the same features. However, apart from packetLSTM and OLIFL, 467 all models show declining performance as they progress from the first to the third and fifth intervals 468 and from the second to the fourth interval. Even OLIFL fails in the SUSY dataset as its performance 469 declines from the second to the fourth interval (see section S.2 of the Appendix). The packetLSTM's 470 ability to retain knowledge is attributed to each feature's local information stored in its correspond-471 ing LSTM's short-term memory. To quantify knowledge retention, we compare packetLSTM with 472 a version retrained (packetLSTM-retraining) at each interval using the available features (see the 473 rightmost graph of Figure 4). The packetLSTM's balanced accuracy improves by 1.61 and 1.89 at 474 the third and fifth intervals, respectively, and by 0.93 at the fourth interval, demonstrating effective learning without forgetting. Notably, there is also a performance increase in the second interval 475 despite those features not being previously observed. This increase is due to the global information 476 learned during the first interval, aiding subsequent predictions. Overall, packetLSTM outperforms 477 all other baselines in each interval. Further discussion on mitigating catastrophic forgetting is pro-478 vided in section S.1 of the Appendix. It is important to distinguish that catastrophic forgetting within 479 an online learning setting differs from online continual learning (see section T of the Appendix). 480

481

442

443

444

445

446 447

448

449

450

451

452 453

454

455

456

457 458

459

460

461 462

8 TRANSFORMER ON HAPHAZARD INPUTS

482 483

Despite the lack of application of the Transformer (Vaswani et al., 2017) in the field of haphazard 484 inputs, its inherent ability to manage variable-size inputs makes it a natural choice. Therefore, we 485 also investigate Transformer-based methodologies for modeling haphazard inputs.

Dataset	n	Single	e LSTM + Impu	utation	Transforme	r + Padding	Set Transformer	HapTransformer	nacketLSTM
Dutabet	Р	Ffill	KNN Mean	Gaussian	Only Values	Pairs	bet Hunsternier	Thep Transformer	patients
	0.25	52.56±0.33	$57.42 {\pm} 0.07$	$57.82 {\pm} 0.10$	56.41±0.14	$53.95 {\pm} 0.07$	57.27±0.11	60.25±0.25	61.33±0.07
magic04	0.50	$59.86 {\pm} 0.07$	64.01 ± 0.12	64.52 ± 0.11	58.66±0.19	56.61±0.12	59.50 ± 0.07	66.62 ± 0.14	68.31±0.15
imdb	0.75	$68.97{\pm}0.10$	$70.91 {\pm} 0.10$	$71.32{\pm}0.05$	$63.93 {\pm} 0.08$	$60.09 {\pm} 0.14$	$62.04{\pm}0.79$	72.11±0.25	$73.64{\pm}0.11$
imdb		$50.42{\pm}0.06$	$50.37{\pm}0.07$	-	$55.58{\pm}0.17$	$51.09{\pm}0.26$	$51.06 {\pm} 0.12$	$69.40 {\pm} 0.28$	$85.06{\pm}0.04$
	0.25	$50.12{\pm}0.08$	$54.12{\pm}0.18$	-	49.98±0.01	$50.00 {\pm} 0.00$	50.00 ± 0.00	57.76±0.23	60.53±0.16
a8a	0.50	59.54 ± 0.15	65.16 ± 0.26	-	50.02 ± 0.02	50.00 ± 0.00	50.00 ± 0.00	65.95±0.21	67.73±0.21
	0.75	$68.74{\pm}0.20$	$71.31 {\pm} 0.31$	-	$50.05 {\pm} 0.02$	$50.00 {\pm} 0.00$	50.00 ± 0.00	$69.37 {\pm} 0.10$	71.11 ± 0.16
	0.25	$57.20 {\pm} 0.01$	59.31±0.03	59.59±0.01	58.03±0.03	58.78±0.11	58.17±0.03	62.31±0.02	62.77±0.01
SUSY	0.50	64.16 ± 0.02	66.26±0.03	66.87±0.02	62.22 ± 0.10	63.97±0.06	62.03±0.01	69.13±0.01	$69.28 {\pm} 0.01$
	0.75	$71.01 {\pm} 0.02$	$72.18{\pm}0.01$	$72.85 {\pm} 0.02$	$67.60 {\pm} 0.13$	$69.59 {\pm} 0.07$	$65.16 {\pm} 0.03$	$73.74{\pm}0.03$	$73.85{\pm}0.02$
	0.25	$50.54{\pm}0.01$	51.42 ± 0.02	51.77±0.22	50.43±0.01	$50.29 {\pm} 0.01$	50.12±0.01	51.82±0.05	52.22±0.04
HIGGS	0.50	$52.92 {\pm} 0.02$	$54.54 {\pm} 0.03$	$54.84 {\pm} 0.02$	51.07 ± 0.02	$50.58 {\pm} 0.01$	50.28 ± 0.01	55.18 ± 0.04	55.47±0.04
a8a SUSY HIGGS	0.75	$56.76 {\pm} 0.02$	$58.27 {\pm} 0.03$	$58.39 {\pm} 0.02$	$52.58 {\pm} 0.03$	$51.13 {\pm} 0.02$	$50.38 {\pm} 0.01$	$58.44 {\pm} 0.04$	$58.71 {\pm} 0.05$

Table 3: Comparison of packetLSTM vs single LSTM with imputation, transformer with padding, Set Transformer, and HapTransformer on all datasets based on balanced accuracy.

> **Padding** We consider padding inputs with zeros or truncating them, which necessitates specifying a fixed input length (f_l) . If the number of features in an instance $(|\mathbb{F}^t|)$ is less than f_l , the Transformer pads the input with zeros, and if $|\mathbb{F}^t|$ exceeds f_l , it truncates the excess features, potentially leading to information loss. A potential, albeit inefficient, solution is to set an excessively high f_l . Nonetheless, to assess how packetLSTM compares to Transformer, we conducted two experiments: *Only Values*, padding available feature values, and *Pairs*, pairing each feature value with its feature ID and padding the sequence. The packetLSTM outperforms Transformer with padding across all dataset scenarios (see Table 3). Further details can be found in section U.1 of the Appendix.

Natural Language Given the variable size of inputs, natural language can be seen as an application where features arrive one by one, and most are missing. We compared packetLSTM with DistilBERT (Sanh et al., 2019) and BERT (Devlin et al., 2019), detailed in section U.2 of the Appendix. We observe that both DistilBERT and BERT are unable to perform classification on haphazard inputs with a balanced accuracy of around 50 in all cases.

Set Transformer The Set Transformer (Lee et al., 2019), designed for variable-length inputs, has
been previously used in offline learning. We employ it for online learning to handle haphazard
inputs. Results in Table 3 indicate that packetLSTM significantly outperforms Set Transformer
across all datasets, possibly due to Set Transformer's assumption of permutation invariance, which
does not hold for haphazard inputs. More details are provided in section U.3 of the Appendix.

HapTransformer To tackle permutation invariance, we introduce HapTransformer, which trans-forms each feature into a distinct learnable embedding, a technique similarly utilized in prior re-search (Huang et al., 2020; Gorishniy et al., 2021; Somepalli et al., 2021). However, these models do not accommodate variable-sized inputs. The learnable embeddings are subsequently processed by the Set Transformer's decoder, allowing implicit communication of feature identities. Details on the architecture and hyperparameters are provided in section U.4 of the Appendix. Despite its strengths, packetLSTM surpasses HapTransformer in all dataset scenarios (see Table 3). HapTrans-former struggles particularly with datasets having high feature counts like imdb and a8a, and requires significant computational time. However, it remains a strong baseline, outperforming other baseline models in the SUSY and HIGGS datasets, as detailed in Tables 1 and 3.

9 CONCLUSION

In conclusion, our work introduces packetLSTM, a dynamic framework for handling streaming data
 with varying feature dimensions in real-time learning environments. Significantly, the underlying
 principles of packetLSTM are extendable to other architectures, such as vanilla RNN and GRUs,
 demonstrating the model's adaptability to different neural architectures. The packetLSTM not only
 outperforms existing methods across various datasets but also showcases flexibility in adapting to
 changing data conditions without forgetting previously learned information. Our research opens
 new avenues for further exploration of dynamic neural architectures and sets a new benchmark for
 online learning with haphazard inputs.

540 10 ETHICS STATEMENT

541 542

The packetLSTM framework introduces a novel method for handling streaming data with variable 543 features, potentially impacting numerous fields such as healthcare, finance, autonomous systems, en-544 vironmental monitoring, and personalized recommendation systems. In healthcare, it can improve real-time patient monitoring by adapting to new or absent data types, enhancing patient care. Finan-546 cial sectors can benefit from more stable predictive models for trading and risk management due to 547 their ability to process erratic market data. Autonomous systems can gain reliability by effectively 548 managing inconsistencies in sensory data, while environmental monitoring may achieve greater accuracy in tracking ecological changes, aiding policy decisions. In digital platforms, it may refine 549 personalized recommendations by adjusting to user behavior changes, and in educational technol-550 ogy, it may personalize content to student needs, improving engagement and learning outcomes. 551

552 Due to the potential application of packetLSTM in the above-discussed crucial and sensitive fields, 553 it is necessary to consider ethical concerns. The packetLSTM framework must navigate several 554 ethical issues to align with the European Union's Artificial Intelligence Act (EU AI Act). Deploy-555 ing packetLSTM technology requires careful consideration, including data privacy, transparency in 556 decision-making, and addressing data biases to prevent discrimination and ensure fairness across 557 diverse user groups. By addressing these challenges, packetLSTM can significantly enhance effi-558 ciency and personalization across multiple domains while upholding high ethical standards.

559 560

561

11 Reproducibility

All the information to reproduce the result is available in sections 4, 5, 6, 7, and 8 of the main manuscript. Additional information is also provided in sections I, J, R, and U of the Appendix.
The code is attached to this submission in the supplementary material with sufficient instructions to faithfully reproduce all the experimental results. The link to the datasets is provided in section I of the Appendix. We report the resources used to conduct the experiments in section I of the Appendix. Moreover, for the benchmark results, we also report the time taken by each individual model in section K of the Appendix.

569 570

571

579

580

581

582

591

References

- Rohit Agarwal, Krishna Agarwal, Alexander Horsch, and Dilip K Prasad. Auxiliary network: Scalable and agile online learning for dynamic system with inconsistently available inputs. In *Inter- national Conference on Neural Information Processing*, pp. 549–561. Springer, 2022.
- Rohit Agarwal, Deepak Gupta, Alexander Horsch, and Dilip K. Prasad. Aux-drop: Handling haphazard inputs in online learning using auxiliary dropouts. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL https://openreview.net/forum?id= R9CgBkeZ6Z. Reproducibility Certification.
 - Rohit Agarwal, Arijit Das, Alexander Horsch, Krishna Agarwal, and Dilip K Prasad. Online learning under haphazard input conditions: A comprehensive review and analysis. *arXiv preprint arXiv:2404.04903*, 2024.
- Sumeet Agarwal, V Vijaya Saradhi, and Harish Karnick. Kernel-based online machine learning and support vector reduction. *Neurocomputing*, 71(7-9):1230–1237, 2008.
- Zeyuan Allen-Zhu and Yuanzhi Li. Can sgd learn recurrent neural networks with provable general ization? Advances in Neural Information Processing Systems, 32, 2019.
- Pierre Baldi, Peter Sadowski, and Daniel Whiteson. Searching for exotic particles in high-energy physics with deep learning. *Nature Communications*, 5(1):4308, 2014.
- Ege Beyazit, Jeevithan Alagurajah, and Xindong Wu. Online learning from data streams with varying feature spaces. In *AAAI Conference on Artificial Intelligence*, volume 33, pp. 3232–3239, 2019.

- ⁵⁹⁴ RK Bock, A Chilingarian, M Gaug, F Hakl, Th Hengstebeck, M Jiřina, J Klaschka, E Kotrč, P Savický, S Towers, et al. Methods for multidimensional event classification: a case study using images from a cherenkov gamma-ray telescope. *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment*, 516(2-3): 511–528, 2004.
- Zhengping Che, Sanjay Purushotham, Kyunghyun Cho, David Sontag, and Yan Liu. Recurrent neural networks for multivariate time series with missing values. *Scientific Reports*, 8(1):6085, 2018.
- Zhong Chen, Yi He, Di Wu, Huixin Zhan, Victor Sheng, and Kun Zhang. Robust sparse online
 learning for data streams with streaming features. In *Proceedings of the 2024 SIAM International Conference on Data Mining (SDM)*, pp. 181–189. SIAM, 2024.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL https://aclanthology.org/ N19-1423.
- Jiahua Dong, Yang Cong, Gan Sun, Tao Zhang, Xu Tang, and Xiaowei Xu. Evolving metric learning for incremental and decremental features. *IEEE Transactions on Circuits and Systems for Video Technology*, 32(4):2290–2302, 2021.
- Joao Gama. A survey on learning from data streams: current and future trends. *Progress in Artificial Intelligence*, 1:45–55, 2012.
- Felix A Gers and Jürgen Schmidhuber. Recurrent nets that time and count. In *Proceedings of the IEEE-INNS-ENNS International Joint Conference on Neural Networks. IJCNN 2000. Neural Computing: New Challenges and Perspectives for the New Millennium*, volume 3, pp. 189–194. IEEE, 2000.
- Yury Gorishniy, Ivan Rubachev, Valentin Khrulkov, and Artem Babenko. Revisiting deep learning
 models for tabular data. Advances in Neural Information Processing Systems, 34:18932–18943,
 2021.
- Klaus Greff, Rupesh K Srivastava, Jan Koutník, Bas R Steunebrink, and Jürgen Schmidhuber. Lstm:
 A search space odyssey. *IEEE Transactions on Neural Networks and Learning Systems*, 28(10):
 2222–2232, 2016.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pp. 1026–1034, 2015.

630

634

635

636

- Yi He, Baijun Wu, Di Wu, Ege Beyazit, Sheng Chen, and Xindong Wu. Online learning from capricious data streams: a generative approach. In *International Joint Conference on Artificial Intelligence, Main track*, 2019.
- Yi He, Jiaxian Dong, Bo-Jian Hou, Yu Wang, and Fei Wang. Online learning in variable feature spaces with mixed data. In *IEEE International Conference on Data Mining*, pp. 181–190. IEEE, 2021.
- Yi He, Christian Schreckenberger, Heiner Stuckenschmidt, and Xindong Wu. Towards utilitarian
 online learning-a review of online algorithms in open feature space. In *International Joint Con- ferences on Artificial Intelligence Organization*, 2023.
- Michiel Hermans and Benjamin Schrauwen. Training and analysing deep recurrent neural networks.
 Advances in Neural Information Processing Systems, 26, 2013.
- 647 Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8): 1735–1780, 1997.

648 Steven CH Hoi, Doyen Sahoo, Jing Lu, and Peilin Zhao. Online learning: A comprehensive survey. 649 Neurocomputing, 459:249–289, 2021. 650 Bo-Jian Hou, Lijun Zhang, and Zhi-Hua Zhou. Learning with feature evolvable streams. Advances 651 in Neural Information Processing Systems, 30, 2017. 652 653 Chenping Hou and Zhi-Hua Zhou. One-pass learning with incremental and decremental features. 654 IEEE Transactions on Pattern Analysis and Machine Intelligence, 40(11):2776–2792, 2018. doi: 655 10.1109/TPAMI.2017.2769047. 656 Xin Huang, Ashish Khetan, Milan Cvitkovic, and Zohar Karnin. Tabtransformer: Tabular data 657 modeling using contextual embeddings. arXiv preprint arXiv:2012.06678, 2020. 658 659 Ioannis Katakis, Grigorios Tsoumakas, and Ioannis Vlahavas. On the utility of incremental feature 660 selection for the classification of textual data streams. In Advances in Informatics: 10th Panhel-661 lenic Conference on Informatics, pp. 338-348. Springer, 2005. 662 Seyed Mehran Kazemi, Rishab Goel, Sepehr Eghbali, Janahan Ramanan, Jaspreet Sahota, Sanjay 663 Thakur, Stella Wu, Cathal Smyth, Pascal Poupart, and Marcus Brubaker. Time2vec: Learning a 664 vector representation of time. arXiv preprint arXiv:1907.05321, 2019. 665 666 Cheol Ho Kim, Jung-Hoon Lee, Hawon Shin, and Ock Kee Baek. Robust and adaptive incremental 667 learning for varying feature space. IEEE Access, 12:64177–64192, 2024. doi: 10.1109/ACCESS. 2024.3395996. 668 669 Donald Knuth. Seminumerical algorithms. The Art of Computer Programming, 2, 1981. 670 671 Ron Kohavi et al. Scaling up the accuracy of naive-bayes classifiers: A decision-tree hybrid. In Second International Conference on Knowledge Discovery and Data Mining, volume 96, pp. 672 202-207, 1996. 673 674 Juho Lee, Yoonho Lee, Jungtaek Kim, Adam Kosiorek, Seungjin Choi, and Yee Whye Teh. Set 675 transformer: A framework for attention-based permutation-invariant neural networks. In Interna-676 tional conference on machine learning, pp. 3744–3753. PMLR, 2019. 677 Jung-Hoon Lee, Sungyup Lee, Cheol Ho Kim, and OK Baek. A study on imputation-based online 678 learning in varying feature spaces. In 14th IEEE International Conference on Information and 679 Communication Technology Convergence, pp. 1759–1764, 2023a. 680 681 Yi-Lun Lee, Yi-Hsuan Tsai, Wei-Chen Chiu, and Chen-Yu Lee. Multimodal prompting with missing 682 modalities for visual recognition. In Proceedings of the IEEE/CVF Conference on Computer 683 Vision and Pattern Recognition, pp. 14943–14952, 2023b. 684 Diyang Li and Bin Gu. When online learning meets ode: learning without forgetting on variable 685 feature space. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 37, pp. 686 8545-8553, 2023. 687 688 Zhizhong Li and Derek Hoiem. Learning without forgetting. IEEE Transactions on Pattern Analysis 689 and Machine Intelligence, 40(12):2935–2947, 2017. 690 Jiaying Liu, Yanghao Li, Sijie Song, Junliang Xing, Cuiling Lan, and Wenjun Zeng. Multi-modality 691 multi-task recurrent neural network for online action detection. IEEE Transactions on Circuits 692 and Systems for Video Technology, 29(9):2667-2682, 2018. 693 Yanfang Liu, Xiaocong Fan, Wenbin Li, and Yang Gao. Online passive-aggressive active learning 694 for trapezoidal data streams. IEEE Transactions on Neural Networks and Learning Systems, 2022. 696 Andrew Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts. 697 Learning word vectors for sentiment analysis. In 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies, pp. 142–150, 2011. 699 Daniel Neil, Michael Pfeiffer, and Shih-Chii Liu. Phased lstm: Accelerating recurrent network 700 training for long or event-based sequences. Advances in Neural Information Processing Systems, 29, 2016.

702 703 704	Gergely Neu and Julia Olkhovskaya. Online learning in mdps with linear function approximation and bandit feedback. <i>Advances in Neural Information Processing Systems</i> , 34:10407–10417, 2021.
705 706 707	Peijia Qin and Liyan Song. Online learning in varying feature spaces with informative variation. In <i>International Conference on Intelligent Information Processing</i> , pp. 19–33. Springer, 2024.
708 709	Reza Sajedi and Mohammadreza Razzazi. Data stream classification in dynamic feature space using feature mapping. <i>The Journal of Supercomputing</i> , pp. 1–19, 2024.
710 711 712	Hojjat Salehinejad, Sharan Sankar, Joseph Barfett, Errol Colak, and Shahrokh Valaee. Recent advances in recurrent neural networks. <i>arXiv preprint arXiv:1801.01078</i> , 2017.
713 714	Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. <i>arXiv preprint arXiv:1910.01108</i> , 2019.
715 716 717 718	Christian Schreckenberger, Christian Bartelt, and Heiner Stuckenschmidt. Dynamic forest for learn- ing from data streams with varying feature spaces. In 28th International Conference on Cooper- ative Information Systems, pp. 95–111. Springer, 2022.
719 720 721	Christian Schreckenberger, Yi He, Stefan Lüdtke, Christian Bartelt, and Heiner Stuckenschmidt. Online random feature forests for learning in varying feature spaces. In <i>AAAI Conference on Artificial Intelligence</i> , volume 37, pp. 4587–4595, 2023.
722 723 724	Gowthami Somepalli, Micah Goldblum, Avi Schwarzschild, C Bayan Bruss, and Tom Goldstein. Saint: Improved neural networks for tabular data via row attention and contrastive pre-training. <i>arXiv preprint arXiv:2106.01342</i> , 2021.
725 726 727 728 729 730 731	 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), Ad- vances in Neural Information Processing Systems, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/ file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf.
732 733 734	Liyuan Wang, Xingxing Zhang, Hang Su, and Jun Zhu. A comprehensive survey of continual learning: Theory, method and application. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 2024.
735 736 737	Yunbo Wang, Mingsheng Long, Jianmin Wang, Zhifeng Gao, and Philip S Yu. Predrnn: Recurrent neural networks for predictive learning using spatiotemporal lstms. <i>Advances in Neural Information Processing Systems</i> , 30, 2017.
738 739 740	Brent Wenerstrom and Christophe Giraud-Carrier. Temporal data mining in dynamic feature spaces. In Sixth International Conference on Data Mining, pp. 1141–1145. IEEE, 2006.
741 742 743	Kai Xie and Ying Wen. Lstm-ma: A lstm method with multi-modality and adjacency constraint for brain image segmentation. In 2019 IEEE International Conference on Image Processing (ICIP), pp. 240–244. IEEE, 2019.
744 745 746 747	Bo Xu, Cheng Lu, Yandong Guo, and Jacob Wang. Discriminative multi-modality speech recognition. In <i>Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition</i> , pp. 14433–14442, 2020.
748 749 750	Huigui Yan, Jiale Liu, Da Han, Dianlong You, Hongtao Wu, Zhen Chen, and Xindong Wu. Online learning from incomplete data streams for multi-classification. <i>Information Sciences</i> , pp. 121411, 2024a.
751 752 753	Huigui Yan, Jiale Liu, Jiawei Xiao, Shina Niu, Siqi Dong, Dianlong You, and Limin Shen. Online learning for data streams with bi-dynamic distributions. <i>Information Sciences</i> , pp. 120796, 2024b.
754 755	Dianlong You, Huigui Yan, Jiawei Xiao, Zhen Chen, Di Wu, Limin Shen, and Xindong Wu. Online learning for data streams with incomplete features and labels. <i>IEEE Transactions on Knowledge and Data Engineering</i> , 2024.

756 757 758 750	Chaoyun Zhang, Marco Fiore, Iain Murray, and Paul Patras. Cloudlstm: A recurrent neural model for spatiotemporal point-cloud stream forecasting. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 35, pp. 10851–10858, 2021a.
760 761 762	Qin Zhang, Peng Zhang, Guodong Long, Wei Ding, Chengqi Zhang, and Xindong Wu. Online learning from trapezoidal data streams. <i>IEEE Transactions on Knowledge and Data Engineering</i> , 28(10):2709–2723, 2016.
763 764 765	Xiang Zhang, Marko Zeman, Theodoros Tsiligkaridis, and Marinka Zitnik. Graph-guided network for irregularly sampled multivariate time series. In <i>International Conference on Learning Representations</i> , 2021b.
766 767 768 769	Zhen-Yu Zhang, Peng Zhao, Yuan Jiang, and Zhi-Hua Zhou. Learning with feature and distribution evolvable streams. In <i>International Conference on Machine Learning</i> , pp. 11317–11327. PMLR, 2020.
770 771 772	Yuxuan Zhao, Eric Landgrebe, Eliot Shekhtman, and Madeleine Udell. Online missing value impu- tation and change point detection with the gaussian copula. In <i>Proceedings of the AAAI Confer-</i> <i>ence on Artificial Intelligence</i> , volume 36, pp. 9199–9207, 2022.
773 774 775	Han Zhou and Shin Matsushima. Online learning under capricious feature data streams. In 37th Annual Conference of the Japanese Society for Artificial Intelligence, 2023.
776 777	Peng Zhou, Shuai Zhang, Lin Mu, and Yuanting Yan. Online learning from capricious data streams via shared and new feature spaces. <i>Applied Intelligence</i> , pp. 1–17, 2024.
778 779 780	Yu Zhu, Hao Li, Yikang Liao, Beidou Wang, Ziyu Guan, Haifeng Liu, and Deng Cai. What to do next: Modeling user behaviors by time-lstm. In <i>International Joint Conference on Artificial Intelligence</i> , volume 17, pp. 3602–3608, 2017.
781 782 783 784	Sheng-Da Zhuo, Jin-Jie Qiu, Chang-Dong Wang, and Shu-Qiang Huang. Online feature selection with varying feature spaces. <i>IEEE Transactions on Knowledge and Data Engineering</i> , pp. 1–14, 2024. doi: 10.1109/TKDE.2024.3377243.
786 787	
788	
789	
790	
791	
792	
793	
794	
795	
796	
797	
790	
800	
801	
802	
803	
804	
805	
806	
807	
808	
809	

810 PRACTICAL APPLICATIONS А

811 812

With the progress of the field of haphazard inputs, more datasets are expected to become avail-813 able. For instance, Schreckenberger et al. (2023) introduced the crowdsense dataset in 2023, col-814 lected from environmental sensors (measuring sound pressure, eCO2 level, and eTVOC level) across 815 Spain's 56 largest cities over 790 days from January 1st (2020) to February 28th (2022). The feature 816 space varies as the new sensing data continuously emerges while many old sensors stop to provide 817 data. The aim is to predict government restriction severity based on sensed regional crowdedness. 818 We could not include this dataset in our study because it has only 790 instances, which is very 819 small for a deep-learning model. Nevertheless, we anticipate the emergence of larger datasets in this 820 domain. Note that we demonstrate the superior performance of packetLSTM compared to article (Schreckenberger et al., 2023) in our study in Table 1 in the main manuscript. 821

822 Another potential application is in the study of sub-cellular organisms, as discussed in (Agarwal 823 et al., 2024). The mitochondria undergo morphological changes during the drug discovery process, 824 leading to fusion, fission, and kiss-and-run events. Quantifying these events is crucial for biologists, 825 yet direct observation is impractical due to hundreds of mitochondria within a single cell. Fusion and fission introduce obsolete and sudden features, while kiss-and-run events can lead to sudden, 826 obsolete, and missing features. The combination of haphazard input and segmentation models may 827 effectively model the mitochondrial dynamics. 828

829 830

831 832 833

834 835 836

837

838

В MULTI-MODAL LEARNING VS PACKETLSTM

The multi-modality requires handling data from different data domains, like images, audio, etc. Some studies in the multi-model domain utilize LSTMs (more generally RNNs), as discussed below.

Brain Image Segmentation LSTM-MA (Xie & Wen, 2019) converts the multi-modal image slices to feature sequences using pixel and superpixel constraints. These are then fed to the LSTM, followed by a fully connected neural network to predict the final class label. Finally, the segmentation image is obtained by combining all the classified nodes of a slice.

839 840 841

843

Online Action Detection Liu et al. (2018) proposed an RNN framework for online action de-842 tection and forecasting on the fly from the untrimmed stream data. The data have two modalities, i.e., RGB video and skeleton sequence data. The paper uses deep ConvNet and motion networks 844 to extract embedding from RGB and skeleton data, respectively, and the embeddings are processed 845 individually by stacked LSTM.

846 847

Speech Recognition AE-MSR (Xu et al., 2020) works on two modalities, namely, audio and 848 video, to perform speech recognition. Here, both audio and video data are passed to an AE sub-849 network to generate visual features and enhanced audio magnitude. The generated features are fur-850 ther passed to different element-wise attention-gated recurrent units (EleAtt-GRU) encoders. This 851 generates visual and audio context, respectively. Both contexts are further passed to a decoder con-852 sisting of EleAtt-GRU to generate outputs. 853

Within the multi-modal domain, to the best of our knowledge, the total number of features/modalities 854 always remains the same and is known at the outset (t = 0). Some recent work also considers missing 855 modality in the multi-modal domain (Lee et al., 2023b), but still, the total number of modalities is 856 known at the outset (t = 0). Therefore, the above method cannot accommodate haphazard inputs. 857 Whereas the packetLSTM can dynamically change its architecture to adapt to varying feature spaces 858 in real-time environments. 859

860 861

С NOTATIONS

862 863

All the notations used in this article are provided in Table 4 with their corresponding meanings.

865		
866	Notations	Meaning
867	+	A random time or the t^{th} instance
868	$\overset{\iota}{X^t}$	The set of input values arrived at time t
869	u^t	Ground truth at time t
870	$\overset{g}{C}$	Number of output classes
871	Ř	Real numbers
872	d^t	The number of input features arrived at time t
873	t-1	The instance preceding time t
874	F_{i}	Feature <i>j</i>
875	$ec{\mathbb{F}^t}$	Universal feature space
876	\mathbb{F}^{t}	Current feature space
877	f^t	The trained model after time t
878	\hat{y}^t	Prediction at time t by model f
879	l^t	Loss of the model f at time t
880	H	Loss function
881	L_j	LSTM corresponding to feature j
882	t_k	The k^{th} time or the instance
883	x_j^k	The value of feature j at time t_k
884	Δ_j^t	Time delay of feature j at time t
885	h_j^t	Short-term memory of L_j at time t
886	c_j^t	Long-term memory of L_j at time t
887	h_i^{t}	Last observed short-term memory of L_i before time t
888	\check{h}^t	Predictive short-term memory at time t
220	c^t	Common long-term memory at time t
800	h^0	Initialized values of short-term memory
090	c^0	Initialized values of long-term memory
091	i_{j}^{t}	Input gate at time t of L_j
892	$T \tilde{1}_{i}^{t}$	Time gate 1 at time t of L_i
893	$T2_{i}^{t}$	Time gate 2 at time t of L_i
894	$o_i^{t'}$	Output gate at time t of L_i
895	\tilde{c}_{i}^{t}	Cell state at time t of L_i
896	W, w, b	Weight parameters of LSTM
897	s	Size of the short-term and long-term memory
898	$\sigma_{v \in \{c,h\}}$	Tanh function
899	$\sigma_{v \in \{t, T1, T2, \Delta, o\}}$	Sigmoid functions
900	f_l	fixed input length
901	$ \mathbb{F}^t $	Number of available features at time t

Table 4: Notations and their meaning in this article.

D DIFFERENT TIME-MODELING VARIANTS WITHIN THE CONTEXT OF PACKETLSTM

In this article, we present the packetLSTM framework with four different time modeling variants. Primarily, we employ Time-LSTM 3 (Zhu et al., 2017) within the packetLSTM framework. Additionally, we explore the use of Time-LSTM 1, Time-LSTM 2, and decay (Che et al., 2018) into the packetLSTM architecture. The following paragraphs discuss all these time-modeling variants within the context of packetLSTM.

Time-LSTM 3 The formulation of Time-LSTM 3 (and other time modeling variants) for feature j at time t, within the packetLSTM framework, represented by L_j , is given by $h_j^t, c_j^t =$

 $L_i(x_i^t, \Delta_i^t, h_i^{t_-}, c^{t-1}) \forall F_i \in \mathbb{F}^t$. Mathematically, the computation of L_i is expressed as

$$i_{j}^{t} = \sigma_{i}(x_{j}^{t}W_{j,xi} + h_{j}^{t-}W_{j,hi} + c^{t-1} \odot w_{j,ci} + b_{j,i}),$$
(2)

$$T1_{j}^{t} = \sigma_{T1}(x_{j}^{t}W_{j,xT1} + \sigma_{\Delta}(\Delta_{j}^{t}W_{j,T1}) + b_{j,T1}),$$

$$T2_{j}^{t} = \sigma_{T2}(x_{j}^{t}W_{j,xT2} + \sigma_{\Delta}(\Delta_{j}^{t}W_{j,T2}) + b_{j,T2}),$$
(3)

$$\tilde{c}_{j}^{t} = (1 - i_{j}^{t} \odot T1_{j}^{t}) \odot c^{t-1} + i_{j}^{t} \odot T1_{j}^{t} \odot \sigma_{c}(x_{j}^{t}W_{j,x\tilde{c}} + h_{j}^{t-}W_{j,h\tilde{c}} + b_{j,\tilde{c}}),$$
(5)

$$c_j^t = (1 - i_j^t) \odot c^{t-1} + i_j^t \odot T2_j^t \odot \sigma_c(x_j^t W_{j,xc} + h_j^{t-1} W_{j,hc} + b_{j,c}),$$
(6)

$$o_j^t = \sigma_o(x_j^t W_{j,xo} + \Delta_j^t W_{j,o} + h_j^{t-} W_{j,ho} + \tilde{c}_j^t \odot w_{j,\tilde{c}o} + b_{j,o}),$$

$$h_j^t = o_j^t \odot \sigma_h(\tilde{c}_j^t).$$

$$(7)$$

(8)

> where $i_i^t, T1_i^t, T2_i^t$, and o_j^t represents the input gate, time gate 1, time gate 2, and output gate of the LSTM L_i at time t, respectively. The cell state c_i^t influences the current prediction through the output gate and the short-term memory h_i^t . The functions $\sigma_i, \sigma_{T1}, \sigma_{T2}, \sigma_{\Delta}$, and σ_o are sigmoid functions, while σ_c and σ_h are tanh functions. The $W_{j,hi}$ denotes the weight associated with the short-term memory within the input gate of L_i . This definition extends analogously to all W parameters. The symbol \odot indicates the Hadamard product, and all the w parameters are peephole connection weights (Gers & Schmidhuber, 2000).

> **Time-LSTM 2** Similar to Time-LSTM 3, Time-LSTM 2 has two gates. Therefore, equations 2, 3, 4, 7, and 8 hold true for Time-LSTM 2, along with an addition of forget gate (f_i^t) , an update of the equation of cell state (\tilde{c}_i^t) , and long-term memory (c_i^t) .

 $f_{i}^{t} = \sigma_{f}(x_{i}^{t}W_{j,xf} + h_{i}^{t-}W_{j,hf} + c^{t-1} \odot w_{j,cf} + b_{j,f}),$ (9)

$$\tilde{c}_j^t = f_j^t \odot c^{t-1} + i_j^t \odot T \mathbf{1}_j^t \odot \sigma_c(x_j^t W_{j,x\tilde{c}} + h_j^{t-} W_{j,h\tilde{c}} + b_{j,\tilde{c}}), \tag{10}$$

$$c_{j}^{t} = f_{j}^{t} \odot c^{t-1} + i_{j}^{t} \odot T2_{j}^{t} \odot \sigma_{c}(x_{j}^{t}W_{j,xc} + h_{j}^{t-}W_{j,hc} + b_{j,c}).$$
(11)

Time-LSTM 1 Time-LSTM 1 only contains a single time gate. Therefore, the formulation of Time-LSTM 1 within the packetLSTM framework is given by equations 2, 9, 3, 10, 7, and 8.

Decay The decay mechanism introduced by Che et al. (2018) attenuates the short-term memory by a learnable factor γ_i^t . Subsequently, the process adheres to the conventional operations of a vanilla LSTM, which is given by

$$\begin{aligned} \gamma_{j}^{t} &= exp\{-max(0, W_{j,\gamma}\Delta_{j}^{t} + b_{j,\gamma})\},\\ \tilde{h}_{j}^{t_{-}} &= \gamma_{j}^{t} \odot h_{j}^{t_{-}},\\ i_{j}^{t} &= \sigma_{i}(x_{j}^{t}W_{j,xi} + \tilde{h}_{j}^{t_{-}}W_{j,hi} + b_{j,i}),\\ f_{j}^{t} &= \sigma_{f}(x_{j}^{t}W_{j,xf} + \tilde{h}_{j}^{t_{-}}W_{j,hf} + b_{j,f}),\\ \tilde{c}_{j}^{t} &= \sigma_{c}(x_{j}^{t}W_{j,x\tilde{c}} + \tilde{h}_{j}^{t_{-}}W_{j,h\tilde{c}} + b_{j,\tilde{c}}),\\ c_{j}^{t} &= f_{j}^{t}c^{t-1} + i_{j}^{t}\tilde{c}_{j}^{t},\\ o_{j}^{t} &= \sigma_{o}(x_{j}^{t}W_{j,xo} + \tilde{h}_{j}^{t_{-}}W_{j,ho} + b_{j,o}),\\ h_{j}^{t} &= o_{j}^{t}\sigma_{h}(c_{j}^{t}). \end{aligned}$$
(12)

E DATASETS

The detailed descriptions of each dataset are provided in Table 5 and further elaborated below:

• magic04 (Bock et al., 2004): It is a Monte Carlo simulated dataset for the registration of high-energy particles in an atmospheric Cherenkov telescope. The binary classification

Dataset	# Instances	# Features	Imbalance	Missing Values %	Туре
magic04	19020	10	64.84%	25%, 50%, 75%	Synthetic
imdb	25000	7500	50%	98.35%	Real
a8a	32561	123	75.92%	25%, 50%, 75%	Synthetic
SUSY	1M	8	45.79%	25%, 50%, 75%	Synthetic
HIGGS	1M	21	52.97%	25%, 50%, 75%	Synthetic

Table 5: Dataset Description: # {Instances, Features} represents the number of {Instances, Features}.
974

task associated with magic04 is to distinguish between a shower image caused by primary gammas (1) and cosmic rays in the upper atmosphere (0). The magic04 dataset can be accessed via this link¹.

- imdb (Maas et al., 2011): It is a movie sentiment classification dataset with labels as positive (1) or negative (0). The original imdb dataset consists of training and test subsets, each containing 25000 instances. Following the previous literature in the field of haphazard inputs (Beyazit et al., 2019; Agarwal et al., 2024), we consider the training subset of the data provided in this link² and the 7500 most prevalent words within this subset, which represents the 7500 features. This approach adheres to the standard practices of previous works, ensuring fair comparison.
- **a8a** (Kohavi et al., 1996): It is the income data from a census conducted in 1994. The task is to classify where the income exceeds \$50k per year. The dataset is pre-processed, resulting in 123 features. The a8a dataset can be accessed via this link³.
 - SUSY (Baldi et al., 2014): It is a Monte Carlo simulated data of kinematic properties of particles with a binary classification task of predicting between a signal process (1) where SUSY particles are produced and a background process (0) with the same detectable particles. Following article (Agarwal et al., 2023), the first 8 features are considered. The SUSY data can be accessed via this link⁴.
- HIGGS (Baldi et al., 2014): Similar to SUSY, it is a Monte Carlo simulation data associated with a binary classification task to differentiate between a signal process (1) where new theoretical Higgs bosons are produced and a background process (0) with identical decay products but distinct kinematic features. We consider the first 21 features of the HIGGS dataset (Agarwal et al., 2023). The HIGGS data can be accessed via this link⁵.

Synthetic Dataset Preparation We created haphazard input datasets from synthetic datasets based on probability values p, as defined in the baseline papers (Beyazit et al., 2019; Agarwal et al., 2023). Specifically, $100 \times p\%$ of features at each time instance is simulated as available independently of each other following a uniform distribution. For example, if p = 0.25, 25% of features at each time instance is only available. From each synthetic dataset, three subsets are generated corresponding to p values of 0.25, 0.5, and 0.75. This approach creates a spectrum of datasets, from highly unavailable data to those with extensive data availability, thereby facilitating the test-ing of models under varying conditions of data accessibility. Notably, 98.35% of data values are unavailable in the imdb dataset (see Table 5).

F METRICS

Number of Errors This is defined as the number of instances incorrectly classified by the model.
 A model is deemed more effective when it exhibits a lower number of errors. However, note that the number of errors may not be a suitable metric for imbalanced data.

¹https://archive.ics.uci.edu/dataset/159/magic+gamma+telescope

^{1023 &}lt;sup>2</sup>https://ai.stanford.edu/~amaas/data/sentiment/

^{1024 &}lt;sup>3</sup>https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/binary.html

^{1025 &}lt;sup>4</sup>https://archive.ics.uci.edu/dataset/279/susy

⁵https://archive.ics.uci.edu/dataset/280/higgs

Accuracy The Model's accuracy can be determined by calculating the total number of correct predictions divided by the total number of instances. Similar to the number of errors, accuracy may not serve as a reliable metric in the context of imbalanced datasets, as it can be misleadingly high when the majority class is predicted correctly while neglecting the minority class.

AUROC The Area Under the Receiver Operating Characteristic Curve (AUROC) measures a model's capability to differentiate between positive and negative classes. The ROC curve plots the true positive rates against the false positive rates. The area of this curve, the AUROC, is a value for which higher numbers indicate superior model performance.

1035

AUPRC AUPRC stands for Area Under the Precision-Recall Curve. It is similar to AUROC; how ever, it utilizes precision and recall rather than true and false positive rates. AUPRC is particularly
 valuable for assessing the performance of models on imbalanced datasets, providing insight into the
 model's ability to identify positive instances.

Balanced Accuracy This is calculated as the average of specificity and sensitivity. This metric offers a more nuanced assessment of model performance across imbalanced datasets compared to traditional accuracy, as it equally weighs the correct identification rates of both positive and negative classes.

1044 1045

1046 G EVALUATION PROTOCOL

1047

We adhere to the standard evaluation protocol for haphazard inputs and, more generally, online learn-1048 ing. As described in the *Mathematical Formulation* paragraph of Section 3 in the main manuscript, 1049 the model is evaluated and then trained iteratively. There is no separate evaluation set. The model 1050 first receives input features at time t, makes a prediction, and then the ground truth is revealed. The 1051 model calculates the cross-entropy loss using the prediction and ground truth. Each instance is pro-1052 cessed once and not revisited, reflecting the online learning characteristics of batch size 1 and epoch 1053 1. The prediction logits and the labels for each instance are stored, and upon processing all the 1054 inputs, the balanced accuracy and other metrics are determined based on these data. This evaluation 1055 method is consistently applied across all models in this study.

1056 1057

1058

H BASELINES NOT INCLUDED FOR COMPARISON

The following models apply to the field of haphazard inputs: OLCF (Zhou & Matsushima, 2023), OIL (Lee et al., 2023a), DCDF2M (Sajedi & Razzazi, 2024), OFSVF (Zhuo et al., 2024), DFLS (Li & Gu, 2023), RSOL (Chen et al., 2024), OVFIV (Qin & Song, 2024), RAIL (Kim et al., 2024), OLBD (Yan et al., 2024b), OLCDS (Zhou et al., 2024), and OLIDS_{PLM} (Yan et al., 2024a). However, the inclusion of these additional models in our article was not feasible due to the lack of open-source code and the challenges associated with implementing them, which stemmed from either insufficient detail or the complexity of the models.

1066

1067 I IMPLEMENTATION DETAILS

All the models were implemented using the PyTorch framework, and the experiments were conducted on an NVIDIA DGX A100 machine. Since the model operates in an online learning setting, CPUs were exclusively used to sequentially process the data.

1072

packetLSTM The size for both the long-term and short-term memory components was set at 64.
The final fully connected network is a two-layer neural network with a ReLU non-linear activation function. Stochastic gradient descent was employed for back-propagation. Cross-entropy loss and AdamW optimizer were used. The learning rate for magic04, imdb, a8a, SUSY, and HIGGS is set as 0.0006, 0.0008, 0.0009, 0.0008, and 0.0002, respectively.

- 1078
- **Baselines** In the case of Aux-Drop baseline (Agarwal et al., 2023), the online deep learning framework was adopted as the foundational model. The implementation strategies for the majority of

1080	Table 6: Description of all the hyperparameters used in each model and their search values. This
1081	table is adapted from article (Agarwal et al., 2024). {} represents individual values, and [] represents
1082	a range.

Model	Hyperparameter(s)	Description	Search
NB3	n	proportion of total features to consider as the number of top fea- tures	{.2, .4, .6, .8, 1}
	m	number of instances after which learner's output is considered	5
	f	threshold on the difference of youngest learner's feature set and current top M features	0.15
FAE	n	number of consecutive instances a learner is under the threshold	3
	r	before being removed minimum number of instances between the 2 consecutive learn-	10
		ers	
	N M	number of instances to calculate accuracy over proportion of total features to train new feature forest on	$\{.2, .4, .6, .8, 1\}$
	С	loss bounding parameter for instance classifier	{.0001, .01, 1}
OLVF	\bar{C}	loss bounding parameter for feature classifier	$\{.0001, .01, 1\}$
	λ	regularization parameter	$\{.01, .1, .5, .5, .7, .9, 1\}$ $\{.0001, .01, 1\}$
	Т	number of instances after which 'p' (weighing factor) is updated	{8, 16}
	α	absorption scale parameter used in equation 10 of (He et al., 2010)	{.0001, .001, .01, .1, 1}
OCDS	β_0	absorption scale introduced by us for the 1st term in eq. 9 of He	{.0001, .001, .01, .1, 1}
	ß.	et al. (2019) tradeoff parameter used in equation 0 of (He at al. 2010)	£0001_001_01_1_1
	$\beta_1 \\ \beta_2$	tradeoff parameter used in equation 9 of (He et al., 2019)	$\{.0001, .001, .01, .01, .1, 1\}$
	decay_choice (dc)	decay update rules choices (see original code (He et al., 2021))	{0, 1, 2, 3, 4}
01.179. <i>1</i>	contribute_error_rate (ce)	used in the original code implementation of classifiers in (He et al., 2021)	{.001, .005, .01, .02, .05}
OVFM	window_size (ws)	a (buffer-like) window to store data instances	20
	decay_coef_change (dc) batch_size_denominator (bs)	used in update step in case of learning rate decay	{True, False} {8, 10, 20}
	batch_c (bc)	added to the denominator for stability in learning rate decay	(0, 10, 20)
	α_{ρ}	impact on weight update	$\{.1, .5\}$
	δ	fraction of features to consider (bagging parameter)	$\{.001, .01\}$
DynEo	ϵ	penalty if the split of decision stump is not in the current instance threshold for error rate	$\{.001, .01\}$
Dymo	M	number of learners in the ensemble	$\{500, 1000\}$
	N Al	buffer size of instances lower bounds for the undate strategy	20 05
	$\theta 2$	upper bounds for the update strategy	{.6, .75}
	forestSize (fS)	number of stumps for every feature	$\{3, 5, 10\}$
	updateStrategy (uS)	strategy to replace stumps	{0ldest, random}
ORF3V	replacementChance (rC)	probability to not replace stump for "random" update strategy	.7
	α	weight update parameter	$\{.01, .1, .3, .5, .9\}$
	δ	pruning parameter	.001
	n_base_layer (nb) n_end_layers (ne)	number of base layers number of end layers	5
Aux-Net	n_nodes_layers (n)	number of nodes in each layer	50
	lr b	learning rate discount rate	{.001, .005, .01, .05, .1, .3, .5}
	S	smoothing parameter	.2
	max_num_hidden_layers (nl)	number of hidden layers	6
	n_neuron_aux_layer (na)	total number of neurons in the AuxLayer	$\sim 5 \times \text{num}_{\text{feat}}$
Aux-Drop	b	discount rate	.99
	s lr	learning rate	{.001, .005, .01, .05, .1, .3, .5}
	dropout_p (d)	dropout rate in the AuxLayer	{.3, .5}
OLIFL	option C	The options to determine τ , where loss* = loss/inner_product(X) Tradeoff parameter of loss	$\{loss^*, min(C, 2loss^*)\}$
	- hidden_size	size of short-term and long-term memory	{32. 64}
acketLSTM	lr	learning rate	[0.05, 0.0001]
	aggregate_by (agg)	dimension-invariant aggregation function	{Mean, Max, Min, Sum}

baseline models were in accordance with the guidelines provided by Agarwal et al. (2024), except
for the OLIFL model (You et al., 2024), for which the implementation from the OLIFL article was
utilized. The OVFM model (He et al., 2021), which typically requires buffer storage to store the inputs and thus contradicts the principles of online learning, was adapted by limiting the buffer storage to two instances.

1134 Table 7: Best hyperparameters of the models based on balanced accuracy. In synthetic datasets, we 1135 consider p = 0.5 for hyperparameter search. Experiments marked with § employed heuristically set 1136 hyperparameters without a search process due to the high computational demands. The best value 1137 of λ for OLVF and hidden_size for packetLSTM is always found to be .0001 and 64, respectively, 1138 for all the datasets.

HIGGS	SUSY	a8a	imdb	magic04	Models
0.2	1.0	0.2	0.4	0.6	NB3 (n)
0.4	0.6	0.2	0.8	1.0	FAE (M)
.01, .0001, 1	.01, .01, 1	1, .0001, 1	.01, .0001, 1	.0001, .0001, 1	OLVF (C, \overline{C} , B)
1, .08	1,.3	1, 2	1, .007	0, -	OLIFL (option, C)
8, .0001, .0001, .01, .0001	16, .0001, .01, .0001, .0001	16, 1, 1, .0001, .0001	16, 1, .0001, .01, .01 [§]	16, .01, .01, .0001, .0001	OCDS (T, α , β_0 , β_1 , β_2)
4, .001, F, 20 [§]	3, .005, F, 20	OVFM (dc, ce, dc, bs)			
.5, .5, .2, .001, .7, 1000, .6§	.5, .5, .4, .001, .7, 1000, .6§	.5, .5, .03, .001, .7, 1000, .6	.5, .8, .001, .001, .7, 1000, .6§	.5, .5, .1, .001, .7, 1000, .6	DynFo (α , β , δ , ϵ , γ , M, θ 2)
5, 10, random, .1 [§]	5, 10, random, .18	10, 10, oldest, .1	10, 5, oldest, .01§	10, 5, random, .01	$ORF^{3}V$ (fS, rI, uS, α)
0.05§	0.05§	.01§	.01§	.5	Aux-Net (lr)
100, .05, .3§	100, .05, .3§	500, .01, .3	30000, .01, .3§	100, .01, .3	Aux-Drop (na, lr, d)
.0002, max	.0008, max	.0009, max	.0008, mean	.0006, max	packetLSTM (lr, agg)

1147 1148

1149

1150

1139

J HYPERPARMETER SEARCHING

1151 1152 We followed the strategy employed in article (Agarwal et al., 2024), where the best hyperparameters 1153 for the synthetic dataset are found at p = 0.5 and subsequently used for p = 0.25 and 0.75. The best 1154 hyperparameters for all baseline models, except OLIFL, are derived from article (Agarwal et al., 1154 2024), and the same hyperparameter search protocol is applied to both OLIFL and packetLSTM. 1155 Details of the hyperparameter search values and the selected best values are presented in Table 6 1156 and Table 7. The hyperparameter search process for packetLSTM and OLIFL is discussed next.

1157 1158

packetLSTM We determine the best values of hyperparameters sequentially. First, we fixed the 1159 learning rate and aggregation operator to 0.0001 and mean, respectively, and varied the hidden sizes 1160 of the short-term and long-term memory, being set at 32 and 64. We found the best hidden size to 1161 be 64. Next, we experimented with four aggregation operators, namely, mean, max, min, and sum, 1162 maintaining a fixed learning rate of 0.0001 and a hidden size of 64. We found max and mean as 1163 the best operators for synthetic and real datasets, respectively. Finally, with the best hidden size 1164 and aggregation operator established, we tested a range of learning rates (0.05, 0.01, 0.005, 0.001, 1165 0.0005, 0.0001). We searched in the vicinity of the optimal learning rate found in the previous step. 1166 For example, 0.0005 was found to yield good results on the magic04 dataset, prompting further 1167 exploration of nearby values - 0.0002, 0.0003, 0.0004, 0.0006, 0.007, 0.0008, and 0.0009. Among 1168 these, 0.0006 emerged as the most effective learning rate. The same process was followed for all the datasets. 1169

1170

1171 1172 1173 1174 OLIFL We fixed the option second (min(C, 2loss/inner_product(X))) to determine τ and tested a range of C values (10, 1, 0.1, 0.01, 0.001, 0.00001, 0.000001). Similar to packetLSTM, we searched in the vicinity of the optimal C value found in the previous step. In addition to this, we also experimented with the first option (loss/inner_product(X)) to determine τ .

1175 1176

K BENCHMARKING RESULTS ON OTHER METRICS

In the main manuscript, we compared the models in terms of their balanced accuracy. We also provide the comparison of models on other metrics and time in Table 8.

- 1181 1182
- 1183
- 1184
- 1185

L ABLATIONS STUDIES RESULTS WITH STANDARD DEVIATION

1186 The ablation study conducted in section 6 of the main manuscript reports only the mean balanced 1187 accuracy across five runs, as shown in Table 2. The comprehensive results, including both the mean and the standard deviation, are presented in Table 9. 1188Table 8: Comparison of models on all the metrics and their execution time. The deterministic models1189- NB3, FAE, OLVF, and OLIFL - underwent a single execution. In contrast, the non-deterministic1190models were executed 5 times, with the mean \pm standard deviation reported. The bAcc, Acc, ROC,1191PRC, Err, and Time stand for balanced accuracy, accuracy, AUROC, AUPRC, number of errors, and1192execution time, respectively. A \ddagger symbol indicates non-deterministic models that were run only once1193on specific datasets due to substantial time constraints.

1195	Dataset	p	Metric	NB3	FAE	OLVF	OLIFL	OCDS	OVFM	DynFo	ORF ³ V	Aux-Net	Aux-Drop	packetLSTM
1196			bAcc Acc	50.01 64.81	50.01 64.84	53.18 59.59	53.06 57.22	51.89±0.10 62.99±0.03	51.94±0.00 65.04±0.00	52.75±0.30 59.04±0.37	47.94±0.22 37.96±0.19	50.09±0.07 63.31±0.42	56.04±0.53 67.63±0.30	61.33±0.07 70.95±0.06
1107		0.25	ROC	49.91	48.66	44.02	46.93	48.79±0.26 62.96±0.22	57.22±0.00 70.14±0.00	50.29±0.15 65.56±0.11	48.00±0.09 63.04±0.05	50.13±0.25 64.80±0.19	57.99±0.76 68.46±0.59	67.79±0.08 77.48±0.09
1197			Err Time	6694 1.46	6687 72.7	7686 2.09	8137 2.16	7039.60±5.13 4.30±0.01	6650.00±0.00 70.19±0.09	7789.40±69.64 46.59±5.63	11800.00±36.69 952.28±23.98	6978.40±79.97 376.94±3.42	6156.60±57.79 111.51±7.59	5524.60±10.63 69.58±5.50
1198			bAcc	50.02 64.81	50.00 64.83	54.6 61.06	57.28 58.83	53.40±0.45 54.26+2.00	54.13±0.08 65.02±0.04	55.12±0.06 64.12±0.06	48.56±0.11 35.06±0.12	50.09±0.03 63.38±0.60	59.29±0.48 68.91±0.37	68.31±0.15 75.25±0.13
1199	magic04	0.50	ROC	50.51	47.97	45.38	42.71	55.43±1.54	61.15±0.06	53.24±0.04	47.27±0.07	49.90±0.34 64.61±0.22	62.39±0.34	76.65±0.12
1200			Err Time	6693 1.6	6690 72.16	7407 2.69	7831	8699.00±380.40 5.44±0.65	6653.80±8.14 76.92±11.83	6824.40±10.95 1110.78±38.39	12351.40±22.33 1109.44±17.12	6965.60±113.62 478.10±30.42	5912.80±70.86 116.15±1.19	4707.00±23.94 71.77±1.93
1201			bAcc Acc	49.99 64.82	50 64.83	56.19 62.31	60.75 62.67	53.76±1.07 54.83±2.14	58.79±0.04 61.83±0.02	56.75±0.02 66.58±0.01	49.32±0.04 34.82±0.03	50.05±0.07 63.32±0.72	63.18±0.61 70.66±0.38	73.64±0.11 78.78±0.10
1000		0.75	ROC PRC	49.34 64.15	46.5 59.67	46.89 61.4	39.24 64.53	55.90±1.74 67.60±1.39	60.95±0.02 77.80±0.05	53.79±0.04 67.60±0.01	44.26±0.13 59.68±0.14	50.16±0.24 64.88±0.33	68.48±0.92 75.93±0.43	82.64±0.14 87.38±0.16
1202			Err Time	6691 1.78	6689 51.8	7168 2.16	7101 2.09	8591.00±406.63 3.20±0.00	7259.00±4.36 62.21±7.17	6356.60±1.52 1134.59±13.96	12397.40±6.58 1259.35±21.46	6977.20±136.68 693.89±19.41	5580.80±72.82 122.99±2.65	4035.2±18.94 70.85±2.45
1203			bAcc	81.56	82.18	80.08	54.36	50.14±0.02	77.43 [‡]	57.98±0.29	76.47±0.11	67.41 [‡]	73.10±0.19	85.06±0.04
1204	imdb		ROC	88.08	46.17	48.57	45.64	50.04±0.01	83.92‡	43.11±0.14	55.53±0.02	71.34‡	79.67±0.29	92.66±0.01
1205			PRC Err	84.93 4609	47.95 4455	48.29 4980	59.91 11410	50.13±0.02 12464.20±6.14	82.50 ⁺ 5643 [‡]	45.68±0.15 10504.00±73.27	55.28±0.01 5882.80±27.38	71.27 ⁺ 8148 [‡]	77.74±0.40 6725.60±48.12	92.49±0.02 3734.8±11.12
1205			Time	1257.01 50.01	4117.33	22.17	1570.06 53.57	47818.07±2783.21 54 75±0.87	109199.99 [‡] 58.66±0.00	4072.86±138.29 50.01±0.03	2768.24±23.26 49.99±0.00	223699.04 [‡] 50.00±0.00	46437.42±18114.58 50.00±0.01	4734.65±153.79 60.53±0.16
1206			Acc	75.9	75.91	72.19	75.78	73.07±0.34	78.00±0.00 77.70±0.00	75.73±0.02 50.86±0.09	75.90±0.00 47.88±0.05	75.91±0.00 54.45±0.18	75.85±0.07	77.93±0.03 75.78±0.19
1207		0.25	PRC	77.2	75.82	85.07	66.67	82.31±0.88 8768 60±100 75	91.36±0.00 7165.00±0.00	76.35±0.06	74.49±0.03	77.78±0.23	81.27±2.08	90.13±0.09
1208			Time	22.94	25.7	4.96	37.16	19.17±0.02	2139.21±88.33	4394.51±72.96	207.23±4.25	5484.63±54.75	421.91±119.74	240.96±4.76
1000			Acc	75.92	75.91	76.93	54.65 76.47	74.23±0.64	80.42±0.00	50.11±0.01 75.84±0.01	50.01±0.00 75.91±0.00	50.00±0.00 75.91±0.00	55.33±1.99 77.34±0.55	80.64±0.12
1209	a8a	0.50	PRC	76.55	75.84	88.6	60.82	89.40±0.51	94.12±0.00	76.88±0.06	74.27±0.02	79.40±0.05	89.28±0.93	92.97±0.18 92.97±0.09
1210			Time	23.52	26.18	6.59	42.82	51.67±3.29	2989.01±178.88	4858.85±133.63	7844.00±0.00 395.44±0.80	17639.23±2014.52	235.54±8.39	6302.2±37.75 243.05±4.79
1211			bAcc Acc	50.01 75.92	50 75.91	70.63 79.66	56.56 77.70	68.81±1.10 74.57±0.73	70.95±0.00 82.48±0.00	50.13±0.01 75.85±0.01	49.99±0.00 75.88±0.00	50.00±0.00 75.91±0.00	62.87±0.93 79.94±0.34	71.11±0.16 82.20±0.12
1919		0.75	ROC PRC	50.88 76.16	49.74 75.84	72.43 90	43.44 55.82	77.80±0.57 91.64±0.31	86.79±0.00 95.35±0.00	52.04±0.08 77.11±0.06	46.94±0.05 74.07±0.03	62.44±0.51 81.26±0.30	79.66±1.14 92.35±0.43	84.96±0.09 94.20±0.05
1010			Err Time	7842 28.16	7843 37.84	6622 6.37	7271 48.37	8279.80±238.28 32.28±3.32	5705.00±0.00 3712.20±459.57	7863.00±4.06 5263.24±92.75	7852.00±0.00 569.20±11.54	7843.40±0.55 28776.93±2045.19	6533.00±111.55 577.31±768.95	5796.20±37.90 246.70±6.71
1213			bAcc	50 54.2	49.9 47.09	51.12 49.59	51.23 50.09	52.11±0.19 50.26±0.30	58.00±0.00 60.06±0.00	54.68±0.01 54.85±0.01	49.37±0.01 50.49±0.01	50.53±1.17 54.39±0.93	61.98±0.10 63.79±0.09	62.77±0.01 64.29±0.01
1214		0.25	ROC	49.88	50.41	56.26	48.77	55.93±0.32	57.52±0.01	50.18±0.00	49.42±0.01	51.03±1.80	69.26±0.12	69.68±0.01
1215			Err	457957	529137 12100 93	504062	499125	497382.00±2960.08 163.96±1.52	399356.00±16.84 3771 36±499 81	451496.67±76.94 654852.68±8956.39	495106.60±134.17 24252.71±37.64	456143.60±9290.25 16455.42±167.84	362096.60±871.73 5773.05±129.01	357084.2±63.48 3538.95+32.03
1016			bAcc	50	50.01	53.21	53.59	54.03±0.28	62.85±0.00	58.27±0.00	48.33±0.02	57.89±7.19	68.79±0.14	69.28±0.01
1210	SUSY	0.50	ROC	54.2 50.03	40.87	61.86	46.41	51.77±0.36 57.95±0.32	64.34±0.00	58.76±0.00 50.04±0.00	48.42±0.03 49.06±0.01	60.14±8.30	76.48±0.15	76.78±0.01
1217			Err	45.79 457957	45.51 531338	62.65 483988	53.28 472973	54.01±0.65 482305.60±3558.98	66.74±0.00 356403.2±18.85	45.88±0.00 412369.00±36.87	44.85±0.01 515823.80±285.70	57.49±10.17 392908.20±60770.38	75.44±0.10 300952.40±1155.92	76.19±0.02 297545.2±145.48
1218			bAcc	50	50.12	55.98	56.39	255.58±0.50 54.84±0.48	68.51±0.00	60.94±0.01	23928.85±169.87 47.53±0.03	53.67±8.13	73.55±0.11	73.85±0.02
1219		0.75	ROC	54.2 49.96	47.04 49.51	54.72 65.15	56.05 43.61	52.24±0.54 58.41±0.32	69.71±0.00 72.28±0.00 74.04±0.00	61.61±0.01 49.91±0.00	47.50±0.03 48.36±0.02	57.06±6.93 55.41±10.29	74.40±0.08 81.07±0.07	74.63±0.02 81.42±0.02
1220			Err	457952	529567 12324-28	452821	439520 96.24	477624.00±5394.23 166.07±0.20	302922.6±10.21 2578 27±161 13	383919.4±129 114179.9±2433.89	524984.20±345.24 25690 16±79 67	429402.00±69311.90 23987.22±1326.00	255994.60±801.85 5787.78±123.95	253676.6±186.74 4188.41±86.62
1001			bAcc	50	50.16	50.57	50.56	49.97±0.07	50.61±0.01	50.18 [‡]	49.86±0.03	49.99±0.00	51.17±0.05	52.22±0.04
1221		0.25	ROC	50.03	50.09	49.99	49.44	49.96±0.05 49.96±0.07	51.13±0.00	49.99 [‡]	49.08±0.03 50.00±0.01	52.72±0.02 50.04±0.02	52.36±0.17	54.27±0.03
1222			Err	470403	480197	482630	476766	52.82±0.06 500440.40±475.98	53.80±0.00 474238.00±73.15	492724 [‡]	503204.00±312.46	472809.80±233.06	464302.00±186.18	458531.4±279.81
1223			bAcc	50	50.01	51.21	51.50	50.06±0.06	51.40±0.01	50.21‡	42814.04±94.55 49.82±0.02	49.99±0.01	53.09±0.05	4390.83±09.00 55.47±0.04
1224	HICCE	0.50	Acc ROC	52.96 49.97	52.56 50.12	52.22 50.29	53.06 48.50	50.06±0.05 50.04±0.08	53.09±0.01 52.42±0.00	50.78 [‡] 50.01 [‡]	48.77±0.02 49.98±0.01	52.73±0.04 50.05±0.01	54.81±0.04 55.60±0.06	56.46±0.02 58.59±0.02
1005	nious	0.50	PRC	52.9 470388	53 474407	53.29 477760	56.59 469385	52.86±0.05 499377 20+487 22	54.86±0.00 469145.00±96.03	52.99 [‡] 492219 [‡]	52.96±0.01 512329.40+189.51	53.01±0.02 472704 40+358 46	57.81±0.12 451891 40+365 32	60.62±0.03 435406 2+242 82
1223			Time	142.97	7885.78	145.83	193.74	267.35±0.54	8241.69±231.79	1788308.41‡	42623.97±97.07	44123.61±283.06	6039.45±565.24	4500.17±165.67
1226			bAcc Acc	50 52.97	50.55 49.24	51.98 52.75	52.48 53.77	49.97±0.05 49.97±0.05	52.66±0.00 53.73±0.00	50.16 [†] 50.78 [†]	49.75±0.03 47.84±0.03	49.98 [‡] 52.72 [‡]	55.55±0.11 56.77±0.08	58.71±0.05 59.25±0.05
1227		0.75	ROC PRC	49.95 52.94	50.1 53.48	50.66 53.66	47.52 56.83	49.98±0.02 52.80±0.02	54.01±0.00 56.07±0.00	49.94 [†] 52.91 [†]	49.88±0.01 52.89±0.00	50.06 [‡] 53.04 [‡]	58.96±0.18 60.78±0.22	62.78±0.05 64.53±0.04
1228			Err Time	470336 153.05	507591 549606.27	472523 125.45	462305 191.47	500300.00±475.25 200.58±2.95	462742.00±14.27 8095.05±218.48	492196 [†] 801655.55 [†]	521619.00±301.39 44079.06±134.05	472790 [‡] 65500.41 [‡]	432298.00±791.48 5762.40±170.11	407498.6±538.70 5186.53±50.48
			-											

M MEAN VS MAX AGGREGATION OPERATOR

Table 10 presents the comparison between the Mean and Max aggregation operators in the packetL-STM framework. It is evident from Table 10 that the performance of the Mean operator relative to the Max increases as data availability decreases from p = 0.95 to 0.05. The best hyperparameters found at p = 0.5 is used for p = 0.05 and 0.95.

N STREAMING NORMALIZATION

The details of each streaming normalization technique are discussed below.

1242Table 9: Ablation study on the components (Comp.) of packetLSTM architecture. This is the copy1243of Table 2 with standard deviation (std). The results are reported as mean \pm std of 5 runs. This result1244corresponds to the HIGGS dataset.

240														
040	Variants		magic04		imdb		a8a			SUSY			HIGGS	
240	Cott	25	50	75		25	50	75	25	50	75	25	50	75
247	Mean	60.98±0.08	67.06±0.10	72.30±0.07	85.06±0.04	58.52±0.21	65.15±0.14	67.50±0.19	63.16±0.02	69.59±0.01	74.07±0.05	52.01±0.01	55.27±0.04	58.04±0.02
	Min	61.36±0.14	68.28±0.08	73.58±0.10	77.64±0.21	60.53±0.27	67.78±0.31	71.13±0.12	62.78±0.01	69.29±0.01	73.85±0.00	52.24±0.02	55.50±0.03	58.70±0.07
248	Max	61.33±0.07	68.31±0.15	73.64±0.10	77.90±0.23	60.53±0.16	67.73±0.20	71.11±0.15	62.//±0.01	69.28±0.01	73.85±0.02	52.22±0.04	55.4/±0.04	58.71±0.05
0.40	None Decay	58.83±0.12	65.12 ± 0.14 64.99 ± 0.27	69.87±0.12	85.18±0.08	58.77±0.33	66.52±0.22	70.30±0.25	62.30±0.01	68.70±0.00	73.25±0.02 73.40±0.18	51.37±0.02	54.05±0.07	57.06±0.09
249	E TimeLSTM-1	61.13±0.09	68.14±0.07	73.63±0.14	85.06±0.06	60.45±0.14	67.63±0.13	71.09±0.17	62.77±0.05	69.33±0.01	73.91±0.01	52.23±0.04	55.53±0.06	58.70±0.06
250	≣ TimeLSTM-2 ∏ TimeLSTM-3	61.32 ± 0.07 61.33 ± 0.07	68.34 ± 0.19 68.31 ± 0.15	73.65 ± 0.19 73.64 ± 0.10	85.15±0.06 85.06±0.04	60.33±0.28 60.53±0.16	67.62±0.22 67.73±0.20	71.09 ± 0.26 71.11 ± 0.15	62.79 ± 0.12 62.77 ± 0.01	69.31 ± 0.01 69.28 ± 0.01	73.88±0.01 73.85±0.02	52.26 ± 0.02 52.22 ± 0.04	55.53±0.03 55.47±0.04	58.65±0.06 58.71±0.05
251	Universal	61.37±0.07	68.33±0.17	73.69±0.12	85.06±0.11	60.32±0.33	67.58±0.23	71.15±0.26	62.79±0.01	69.27±0.01	73.84±0.01	52.26±0.02	55.52±0.03	58.69±0.05
231	E Current	61.33±0.07	68.31±0.15	/3.64±0.10	85.06±0.04	60.53±0.16	67.73±0.20	/1.11±0.15	62.//±0.01	69.28±0.01	73.85±0.02	52.22±0.04	55.47±0.04	58./1±0.05
252	Only LTM	60.94±0.07	67.66±0.07	73.09±0.21 73.38±0.11	85.15±0.06 85.17±0.10	59.36±0.34	66.86±0.29	70.34±0.31	62.78±0.02	69.32±0.01	73.85±0.01 73.89±0.02	52.20±0.04	52.62±0.86	57.69±0.10
232	Both	61.33±0.07	68.31±0.15	73.64±0.10	85.06±0.04	60.53±0.16	67.73±0.20	71.11±0.15	62.77±0.01	69.28±0.01	73.85±0.02	52.22±0.04	55.47±0.04	58.71±0.05
253	None	$59.14{\pm}0.08$	$65.44 {\pm} 0.16$	$70.53 {\pm} 0.10$	$85.61 {\pm} 0.04$	59.97±0.11	$67.89{\pm}0.25$	$71.22{\pm}0.21$	$63.12 {\pm} 0.02$	$69.65{\pm}0.01$	$74.19 {\pm} 0.02$	$51.46 {\pm} 0.01$	$54.33 {\pm} 0.08$	$57.39 {\pm} 0.08$
	Min Max	57.69±0.26	63.33±0.20	68.42±0.16	85.23±0.05	59.94±0.25	67.76±0.15	71.15±0.19	62.37±0.05	69.06±0.03	73.61±0.06	51.21±0.02	53.10±0.04	55.52±0.10
254	Decimal Scaling	55.03±0.19	60.36±0.62	65.72±0.35	85.23±0.05	49.99±0.01	50.00±0.01	50.00±0.00	50.96±1.20	58.53±0.87	63.31±0.49	50.00±0.00	50.00±0.00	50.00±0.00
204	Unit Vector	50.01±0.04	54 23+0 24	60.02±0.08	85 32±0.12	58.12±0.25	65.67±0.08	68 97±0.14	56 52±0.01	65 43±0.02	71.44+0.03	51.80 ± 0.01 50.75 ± 0.02	52 86±0.03	55.53 ± 0.11
255	Z-score	61.33±0.07	68.31±0.15	73.64±0.10	85.06±0.04	60.53±0.16	67.73±0.20	71.11±0.15	62.77±0.01	69.28±0.01	73.85±0.02	52.22±0.04	55.47±0.04	58.71±0.05
200	Vanilla RNN	60.70±0.13	66.76±0.25	72.14±0.21	83.55±0.13	51.16±2.32	65.24 ± 0.40	64.50±6.30	62.84±0.02	68.06±2.61	73.98±0.01	51.96±0.05	54.26±0.13	56.77±0.13
256	RU GRU	60.34 ± 0.19	66.24 ± 0.13	71.23 ± 0.21	81.49 ± 0.17	50.96 ± 1.44	$60.86 {\pm} 0.95$	73.38 ± 0.11	62.73 ± 0.01	$69.16{\pm}0.01$	73.72 ± 0.01	$51.80 {\pm} 0.04$	$54.25 {\pm} 0.03$	56.79 ± 0.06
200	LSTM	$58.83 {\pm} 0.12$	$65.12{\pm}0.14$	$69.87 {\pm} 0.12$	$85.18{\pm}0.08$	58.77±0.33	$_{66.52\pm0.22}$	$70.30 {\pm} 0.25$	$62.30 {\pm} 0.01$	$68.70{\pm}0.00$	73.25 ± 0.02	$51.37 {\pm} 0.02$	$54.05 {\pm} 0.07$	$57.06{\pm}0.09$

Table 10: Comparison between Mean vs. Max aggregation operator in the packetLSTM. The mean balanced accuracy of 5 runs is reported here. We consider the synthetic datasets here and the models are run for five p values, namely, 0.05, 0.25, 0.5, 0.75, and 0.95. All the results corresponding to p = 0.25, 0.5, and 0.75 are taken from Table 2 of the main manuscript. All the values of other parameters of packetLSTM are exactly the same, and only the aggregation operator is varied here.

Characteristics	magic04					a8a			SUSY				HIGGS							
	0.05	25	50	75	0.95	0.05	25	50	75	0.95	0.05	25	50	75	0.95	0.05	25	50	75	0.95
Mean Max	56.78 56.74	60.98 61.33	67.06 68.31	72.30 73.64	76.85 78.13	50.14 50.01	58.52 60.53	65.15 67.73	67.50 71.11	66.48 73.33	58.01 57.57	63.16 62.77	69.59 69.28	74.07 73.85	76.60 76.57	50.39 50.43	52.01 52.22	55.27 55.47	58.04 58.71	76.24 77.07
Mean-Max	0.04	-0.35	-1.25	-1.34	-1.28	0.13	-2.01	-2.58	-3.61	-6.85	0.44	0.39	0.31	0.22	0.03	-0.04	-0.21	-0.2	-0.67	-0.83

- Min-Max Normalization: Here, we utilize two placeholders, namely, Mx_j^t and Mn_j^t , which represents the maximum and minimum value of feature *j* till time *t*. The min-max normalization is then defined as $\frac{x_j^t Mn_j^t}{Mx_j^t Mn_j^t}$.
- Decimal Scaling: Here, all the values are scaled down by a predefined threshold as $\frac{x_j^2}{10^m}$. For all the datasets, we set m = 3.
- Z-score Normalization: Here, the values of each feature are normalized based on their running means (μ) and standard deviation (σ) as ^{x_j^t μ_j^t}/_{σ_j^t}, where μ_j^t = μ_j^{t_-} + ^{x_j^t μ_j^{t_-}}/_{k_j^t} and (σ_j^t)^2 = ^{v_j^t}/_{k_j^{t_-1}}. The notation k_j^t denotes the count of the feature j till time t and v_j^t = v_j^{t_-} + (x_j^t μ_j^t)(x_j^t μ_j^t). We utilize the above way of computing running mean and variance because of its superior numerical stability (Knuth, 1981).
- Mean Normalization: Here, the values of each feature are subtracted by their running means as $x_j^t \mu_j^t$.
- Unit Vector Normalization: Here, we consider the whole input feature as a vector and normalize the vector as $\frac{x_j^t}{||X^t||_2}$, where $|| \cdot ||_2$ represents the L_2 norm.
- 1286 1287

1290

19/15

1257

1270

1272

1274

1276

1278 1279

1280

1281

1283

1284

1285

O PACKETRNN AND PACKETGRU

packetRNN The packetRNN framework is illustrated in Figure 5(a). Unlike the LSTM, which incorporates both short-term and long-term memory, the RNN possesses only a single memory element, known as the hidden state (h_j^t) . Consequently, the packetRNN lacks a mechanism for integrating global information and instead maintains local information within its hidden state. These hidden states are combined using a dimension-invariant aggregation operator to generate a common hidden state for final predictions. The mathematical working of a Vanilla RNN unit within the packetRNN



Figure 5: (a) The packetRNN architecture based on (b) the data snapshot. In this framework, the Vanilla RNN may be substituted with a GRU to establish the packetGRU framework.

1329 framework is given by

1324

1327 1328

1330

1331

1332

1334 1335

1336 1337 1338

1349

$$h_{j}^{t} = tanh(x_{j}^{t}W_{j,xh} + h_{j}^{t-}W_{j,hh} + b_{j,h}),$$
(13)

1333 and the final output is given by

$$\hat{y}^t = FCN(AGG(\bigcup_{j,\forall F_j \in \mathbb{F}^t} \{h_j^t\})).$$
(14)

packetGRU The vanilla RNN in Figure 5(a) can be substituted with GRU to establish the packet-GRU framework. Similar to the vanilla RNN, the GRU contains only one memory component but is enhanced with two gates: the update gate (u_j^t) and the reset gate (r_j^t) . The operations of GRU within the packetGRU framework are defined by the following equations:

$$\begin{array}{ll} 1343 \\ 1344 \\ 1344 \\ 1345 \\ 1345 \\ 1346 \\ 1347 \\ 1348 \end{array} \qquad \begin{array}{ll} u_{j}^{t} = \sigma(x_{j}^{t}W_{j,xu} + h_{j}^{t-}W_{j,hu} + b_{j,u}), \\ r_{j}^{t} = \sigma(x_{j}^{t}W_{j,xr} + h_{j}^{t-}W_{j,hr} + b_{j,r}), \\ \tilde{h}_{j}^{t} = tanh(x_{j}^{t}W_{j,xh} + r_{j}^{t} \odot h_{j}^{t-}W_{j,hh} + b_{j,h}), \\ h_{j}^{t} = u_{i}^{t} \odot h^{t-} + (1 - u_{j}^{t}) \odot \tilde{h}_{j}^{t}, \end{array}$$

$$\begin{array}{l} (15) \\ h_{j}^{t} = u_{j}^{t} \odot h^{t-} + (1 - u_{j}^{t}) \odot \tilde{h}_{j}^{t}, \end{array}$$

and the final output of packetGRU is given by the equation 14.

1350 Р **COMPLEXITY ANALYSIS** 1351

1352 **Time Complexity** The time complexity of packetLSTM is $\sum_{t=1}^{T} O(g(|\mathbb{F}^t|) * s^2 + P)$, where T is 1353 the number of instances, $O(g(|\mathbb{F}^t|) * s^2)$ is the time complexity to process all the activated LSTMs 1354 at time t corresponding to \mathbb{F}^t features, and O(P) broadly denotes the constant time required to 1355 perform other fixed operations like aggregations and final prediction. We utilize the torch.matmul() 1356 function for matrix multiplication of LSTMs. The time required for each LSTM is dependent on its hidden size (s) and performs 3 matrix multiplication of size (3s, s+1) and (s+1, 1) requiring 1357 a time complexity of $O(s^2)$. However, we vectorize the operation of $|\mathbb{F}^t|$ LSTMs by single matrix 1358 multiplication of $(|\mathbb{F}^t|, 3s, s+1)$ and $(|\mathbb{F}^t|, s+1, 1)$. Note that the time required by torch.matmul() 1359 does not scale linearly with $|\mathbb{F}^t|$; rather, it just takes a small overhead depending on the type of 1360 hardware and other dependencies. Here, we denote this overhead as a function of $|\mathbb{F}^t|$ as $g(|\mathbb{F}^t|)$ 1361 where $1 \leq q(|\mathbb{F}^t|) \leq |\mathbb{F}^t|$. This is further corroborated by the time required by packetLSTMs on 1362 each synthetic dataset with different p values (see Table 8). For example, the time required to process 1363 the whole HIGGS dataset takes 4396.83 and 4500.17 seconds for p = 0.25 and 0.5, respectively. 1364 Even though the value of $|\mathbb{F}^t|$ doubles from p = 0.25 to 0.5, the time required doesn't increase by 1365 the same factor. Therefore, the time required by the $|\mathbb{F}^t|$ LSTMs at time t would be $O(g(|\mathbb{F}^t|) * s^2)$. Finally, since it is an online learning task and each instance is processed sequentially, the total 1367 time complexity of T instances would be $\sum_{t=1}^{T} O(g(|\mathbb{F}^t|) * s^2 + P)$. It is difficult to find the exact 1368 form of the function g. However, based on the time required by packetLSTM on each dataset with 1369 increasing p, it can be safely assumed that the value of $q(|\mathbb{F}^t|)$ is closer to 1 than $|\mathbb{F}^t|$. Therefore, 1370 the packetLSTM model demonstrates scalability in terms of time complexity corresponding to the number of features. 1371

1372

1373 **Space Complexity** The space complexity of packetLSTM is directly dependent on the space complexity of an LSTM. Let us denote the space complexity of an LSTM by O(L). At each time t, we 1374 have a universal feature space of $|\mathbb{F}^t|$ cardinality, therefore, corresponding $|\mathbb{F}^t|$ LSTMs are present 1375 in the packetLSTMs. The space complexity of the final prediction network and aggregation function 1376 are fixed and denoted by O(K). Therefore, the total space complexity of packetLSTM at time t 1377 can be given by $O(|\mathbb{F}^t| * L + K)$. Note that the space complexity of packetLSTM increases with 1378 the arrival of sudden and missing features. Therefore, the limitation of packetLSTM is that it is not 1379 scalable in terms of space complexity corresponding to the feature size. However, it is noteworthy 1380 that packetLSTM effectively manages up to 7500 features, as demonstrated with the imdb dataset. 1381 Additionally, we present a strategy to further mitigate space complexity limitation in section Q of 1382 the Appendix.

1383

1384 **Number of Parameters** Here, we provide the number of learnable parameters in the packetLSTM 1385 framework with Time-LSTM 3 as the time modeling unit. The mathematical formulation is given 1386 by equations 2-8. The input gate (i), time gate 1 (T1), time gate 2 (T2), cell state (\tilde{c}), long-term 1387 memory (c), and output gate (o) requires $s^2 + 3s$, 3s, 3s, $s^2 + 2s$, $s^2 + 2s$, and $s^2 + 4s$ learnable 1388 parameters, respectively. Therefore, the total number of learnable parameters in an LSTM unit is $4s^2 + 17s$. The fully connected layer accounts for $2s^2 + 4s + 2$ parameters. Therefore, the total 1389 number of learnable parameters at each instance is $|\mathbb{F}^t| * (4s^2 + 17s) + 2s^2 + 4s + 2$. For all 1390 the experiments, the value of s is 64. In this article, the maximum number of parameters for each 1391 dataset would correspond to the worst-case scenario where all the features are present, resulting in a 1392 maximum of \sim 183K, \sim 131M, \sim 2M, \sim 148K, and \sim 375K learnable parameters for magic04, imdb, 1393 a8a, SUSY, and HIGGS, respectively. Consequently, packetLSTM requires 1.40 MB, 999.45 MB, 1394 15.26 MB, 1.13 MB, and 2.85 MB memory with 64-bit precision for magic04, imdb, a8a, SUSY, 1395 and HIGGS, respectively. Considering that a small large language model of 1 billion parameters is 1396 common nowadays, packetLSTM can handle \sim 57K features with 1 billion parameters. We believe that a practical application would have a number of features way less than \sim 57K.

1398



- 1400
- 1401

DROPPING FEATURES TO RESOLVE SPACE COMPLEXITY Q

The space complexity increases with the number of features, as discussed above. However, pack-1402 etLSTM easily handles even the 7500 features in the imdb dataset. The total number of learnable 1403 parameters of packetLSTM for the imdb dataset is ~131M. Therefore, packetLSTM with 1B param1404 eters and a hidden size of 64 can handle around \sim 57K features. Hence, we argue that packetLSTM 1405 can deal with high-dimensional data. However, we also propose a solution to curb the space com-1406 plexity by defining a maximum limit (say l_f) on the number of LSTMs. When the number of 1407 features in the universal feature space $|\overline{\mathbb{F}}^t| > l_f$, we drop $|\overline{\mathbb{F}}^t| - l_f$ features. Here, dropping the 1408 feature means that the corresponding LSTM is reinitialized and assigned to some new features that 1409 arrived at time t. The dropped feature can come in future instances (> t). However, this feature will 1410 then be considered as a sudden feature. We employ KL-Divergence to determine the dropped features. After processing instance t-1, packetLSTM has the short-term memory of each LSTM and 1411 the common long-term memory. The KL divergence between each short-term and common long-1412 term memory is determined. The feature corresponding to the short-term memory, which has the 1413 lowest KL-Divergence, is dropped. This is because the common long-term memory already holds 1414 the dropped feature's short-term memory information and is used for final prediction. We also put 1415 a limit on the number of times a feature is seen (i_l) by packetLSTM before it can be considered for 1416 dropping. We experimented with the imdb dataset since it has the highest number of features (7500). 1417 The best hyperparameters found for packetLSTM on the imdb dataset in Table 7 are used here. The 1418 l_f and i_l are set to 100 and achieved a balanced accuracy of 78.92, which still performs better than 1419 7 out of 10 baseline models (see Table 1 in main manuscript). The 78.92 balanced accuracy is lower 1420 than the packetLSTM without any limits (85.06). So, there is a tradeoff between performance and 1421 the space complexity.

1422

1423 1424 R SINGLE LSTM

1425 1426

1427

1428 1429

1430

1431

1440

The three employed imputation techniques are:

- Forward fill: The missing value of a feature is imputed with its last observed value.
- Mean: The missing value of a feature is imputed with its forward mean determined in a rolling manner. That is, the last K observed values of a feature are considered to calculate its mean. Here K = 5.
- 1432 Gaussian Copula: (Zhao et al., 2022) proposed using Gaussian copula for online streaming 1433 data with a known N, where N is the total number of features. This method argues that data 1434 points are generated from a latent Gaussian vector, which is then transformed to match the 1435 marginal distributions of observed features. However, Gaussian Copula requires storing a 1436 matrix of K instances to perform imputation. The method fails if the matrix's determinant 1437 is 0 or if a feature's values within the matrix are identical. The magic 04 and SUSY require 1438 K = 5, and HIGGS needs K = 30. The method does not work for a8a and imdb for even 1439 K = 300, excluding its application to these datasets.

The hyperparameter search of a single LSTM is performed similarly to packetLSTM. The hyper-1441 parameters of the single LSTM model are hidden size, number of layers, and learning rate. We 1442 searched for the best hidden size among 32, 64, 128, and 256. We determined the optimal num-1443 ber of layers between 1, 2, 3, and 4. Similar to packetLSTM, we tested a range of learning rates 1444 (0.001, 0.0005, 0.0001, 0.00005). We further searched in the vicinity of the optimal learning rate 1445 found in the previous step. The optimal hidden size and number of layers for each dataset are found 1446 to be 32 and 1, respectively. The best learning rate is 0.001, 0.0006, 0.0001, 0.0002, and 0.0008 1447 for magic04, a8a, SUSY, HIGGS, and imdb, respectively. Similar to packetLSTM, we employed 1448 Z-score streaming normalization.

- 1449 1450
- 1451 1452

1454

S CHALLENGING SCENARIOS ON HIGGS AND SUSY

1453 S.1 HIGGS

Here, we provide the exact value of the balanced accuracy used in Figure 3 and 4 from section 7 of
the main manuscript. Table 11 provides the results associated with the experiments on sudden, obsolete, and reappearing features on the HIGGS dataset. Additionally, Table 12 details the comparative results of packetLSTM versus packetLSTM retrained across five data intervals.

Data Interval	OLVF	OLIFL	OVFM	Aux-Drop	packetLSTM
			Sudden		
First	51.00	51.07	52.06 0.00	52 51 10.05	54 20 1 0 04
FIISt	51.09	52.10	52.90 ± 0.00	55.51 ± 0.05	54.50±0.04
Second	54.45	55.10	54.49 ± 0.00	55.77 ± 0.04	50.98±0.05
Inira	54.14	54.68	54.53 ± 0.00	57.98 ± 0.10	60.20±0.13
Fourth	54.08	54.93	54.82 ± 0.00	60.09 ± 0.16	62.37±0.12
Fifth	53.86	55.15	55.31 ± 0.00	61.24 ± 0.20	63.28±0.06
			Obsolete		
First	52 56	50.20	54.62 ± 0.01	5764 0 19	50 00 1 0 08
FIISt	51.00	52.52 52.27	54.02 ± 0.01	57.04 ± 0.18	59.90±0.08
Second	51.80	52.27	52.21 ± 0.01	55.48 ± 0.10	58.15±0.09
Third	51.70	51.94	52.16 ± 0.01	55.69 ± 0.11	57.01±0.03
Fourth	51.48	51.80	51.41 ± 0.01	52.38 ± 0.09	53.77±0.05
Fifth	50.59	50.40	50.73 ± 0.00	50.93 ± 0.07	$51.32{\pm}0.02$
		R	leappearing		
First	51.42	53.34	53.03±0.00	$53.52 {\pm} 0.04$	58.74±0.08
Second	52.07	51.39	$52.58 {\pm} 0.01$	53.20 ± 0.02	53.27±0.07
Third	50.60	54.78	50.79 ± 0.00	52.98 ± 0.19	60.42±0.11
Fourth	51.30	51.88	51.28 ± 0.01	51.65 ± 0.05	53.94±0.08
Fifth	50.63	55.27	$50.78 {\pm} 0.01$	50.79 ± 0.09	60.54±0.13
	Data Interval First Second Third Fourth Fifth First Second Third Fourth Fifth Fifth Fifth Fifth Fifth Fifth Fifth Fifth Fifth	Data Interval OLVF First 51.09 Second 54.45 Third 54.14 Fourth 54.08 Fifth 53.86 First 53.56 Second 51.80 Third 51.70 Fourth 51.48 Fifth 50.59 First 51.42 Second 52.07 Third 50.60 Fourth 51.30 Fifth 50.63	Data IntervalOLVFOLIFLFirst 51.09 51.07 Second 54.45 53.10 Third 54.45 53.10 Third 54.14 54.68 Fourth 54.08 54.93 Fifth 53.86 55.15 First 53.56 52.32 Second 51.80 52.27 Third 51.70 51.94 Fourth 51.48 51.80 Fifth 50.59 50.40 FirstSecond 52.07 Third 51.42 53.34 Second 52.07 51.39 Third 50.60 54.78 Fourth 51.30 51.88 Fifth 50.63 55.27	Data IntervalOLVFOLIFLOVFMFirst 51.09 51.07 52.96 ± 0.00 Second 54.45 53.10 54.49 ± 0.00 Third 54.45 53.10 54.49 ± 0.00 Fourth 54.08 54.53 ± 0.00 Fourth 54.08 54.93 Fifth 53.86 55.15 First 53.56 52.32 Second 51.80 52.27 Second 51.80 52.27 Second 51.48 51.401 Fourth 51.48 51.80 First 50.59 50.40 First 51.42 53.34 Socond 52.07 51.39 ± 0.00 First 51.42 53.34 Socond 52.07 51.39 ± 0.00 Second 52.07 51.39 Second 52.07 51.39 First 51.42 53.34 Socond 52.07 51.39 Second 52.07 51.39 Second 52.07 51.39 First 51.42 53.34 Socond 52.07 50.79 ± 0.00 Fourth 51.30 51.88 Socond 52.27 50.78 ± 0.01 Third 50.63 55.27 Socond 50.78 ± 0.01	Data IntervalOLVFOLIFLOVFMAux-DropFirst 51.09 51.07 52.96 ± 0.00 53.51 ± 0.05 Second 54.45 53.10 54.49 ± 0.00 55.77 ± 0.04 Third 54.14 54.68 54.53 ± 0.00 57.98 ± 0.10 Fourth 54.08 54.93 54.82 ± 0.00 60.09 ± 0.16 Fifth 53.86 55.15 55.31 ± 0.00 61.24 ± 0.20 ObsoleteFirstSacond 51.80 52.27 52.21 ± 0.01 57.64 ± 0.18 Second 51.70 51.94 52.16 ± 0.01 55.69 ± 0.11 Fourth 51.48 51.80 51.41 ± 0.01 52.38 ± 0.09 Fifth 50.59 50.40 50.73 ± 0.00 50.93 ± 0.07 First 51.42 53.34 53.03 ± 0.00 53.52 ± 0.04 Second 52.07 51.39 52.58 ± 0.01 53.20 ± 0.02 Third 50.60 54.78 50.79 ± 0.00 52.98 ± 0.19 Fourth 51.30 51.88 51.28 ± 0.01 51.65 ± 0.05 Fifth 50.63 55.27 50.78 ± 0.01 50.79 ± 0.09

1458 Table 11: The mean±standard deviation of the balanced accuracy of sudden, obsolete, and reap-1459 pearing experiments conducted in section 7 of the main manuscript. This table reflects the values of 1460 the results presented in Figure 3 and 4 and corresponds to the HIGGS dataset.

Table 12: The exact values represented in the experiment 'Learning Without Forgetting' in Figure 4 1484 from section 7 of the main manuscript. This result corresponds to the HIGGS dataset. 1485

1486						
1487	Model	First	Second	Third	Fourth	Fifth
1488 1489	packetLSTM	58.74 ± 0.08 58.74 \pm 0.08	53.27±0.07 53.22±0.06	60.42±0.11 58.81±0.05	53.94±0.08	60.54±0.13 58.65±0.07
1405	packetLSTW-Kettalling	J0.74±0.00	33.22 ± 0.00	30.01±0.05	JJ.01±0.00	J0.0J±0.07

¹⁴⁹⁰ 1491

1482 1483

4.4.0.4

1492

1493

The reappearing experiment in the main manuscript and Figure 4 shows that packetLSTM mitigates 1494 catastrophic forgetting. We further justify this by comparing packetLSTM with Aux-Drop (the best 1495 baseline method as shown in Table 1).

1496 Aux-Drop: The performance of Aux-Drop decreases from 53.52 in the first interval to 52.98 in the 1497 third interval. The performance further decreases to 50.93 in the fifth interval. Aux-Drop also shows 1498 a performance decline from 55.48 in the second interval to 52.38 in the fourth interval. Therefore, 1499 Aux-Drop suffers from catastrophic forgetting.

1500 packetLSTM: The performance of packetLSTM increases from 58.74 in the first interval to 60.42 1501 in the third interval. The performance further increases to 60.54 in the fifth interval. The packetL-1502 STM also shows a performance increase from 53.27 in the second interval to 53.94 in the fourth 1503 interval. Thus, packetLSTM mitigates catastrophic forgetting. 1504

It can be argued that the performance increase of packetLSTM in the third interval is not due to 1505 the mitigation of catastrophic forgetting; it is instead due to a likely situation of better predictive 1506 capability of the features present in the third interval compared to the first interval. 1507

1508 To refute the above argument, we retrained the packetLSTM in each interval – a new initialized 1509 packetLSTM was considered for each interval – and referred to it as packetLSTM-Retraining (see Figure 4 and Table 12). The performance of the packetLSTM (60.42) outperforms packetLSTM-1510 Retraining (58.81) in the third interval, which shows that packetLSTM can retain its learning from 1511 the first interval. Similar performance increase is observed in the fourth and fifth intervals.



Figure 6: Model performance in the scenario of sudden and obsolete features. The mean balanced accuracy in each data interval (e.g., 0 - 0.2 M) is shown in its middle (0.1 M). The pink-colored y-axis represents the Feature(s). For e.g., '1-2' means features 1 to 2 are available in instances shaded with pink color. Both y-axes are shared by the two graphs. This graph corresponds to the SUSY dataset.



Figure 7: Model performance in the scenario of reappearing features. The mean balanced accuracy in each data interval (e.g., 0 - 0.2 M) is shown in its middle (0.1 M). The pink-colored y-axis represents the Features. For e.g., '1-4' means features 1 to 4 are available in instances shaded with pink color. Both y-axes are shared by the two graphs. This graph corresponds to the SUSY dataset.

S.2 SUSY

We perform three challenging experiments – sudden features, obsolete features, and reappearing features – on the SUSY dataset, replicating the experiments conducted on the HIGGS dataset. The data creation settings for the SUSY dataset are identical to those employed for the HIGGS dataset. We considered the OLVF, OLIFL, OVFM, and Aux-Drop baselines for comparative analysis. The performance of each method is illustrated in Figures 6 and 7, with precise values detailed in Tables 13 and 14.

Sudden Features Similar to the approach used for the HIGGS dataset, 20% of the features are sequentially added at each interval in the SUSY dataset, with values rounded to the nearest integer. Given that SUSY comprises 8 features, the initial data interval includes features 1 and 2, the subsequent interval contains features 1 through 3, and this incremental addition continues as depicted in the leftmost graph of Figure 6. All models exhibit a pattern of performance improvement as sudden features are introduced at each interval, as evident in Figure 6. Notably, the packetLSTM consistently outperforms all baselines in each data interval, followed by Aux-Drop, as detailed in Table 13. This underscores the efficacy of packetLSTM in handling sudden features.

1564

1547

1548 1549

1523

1524

1525

1526

1527

Obsolete Features The data arrival pattern in this experiment is the inverse of that observed in the sudden feature experiment, as illustrated in the middle graph of Figure 6. Generally, the performance

29

1568						
1569	Data Interval	OLVF	OLIFL	OVFM	Aux-Drop	packetLSTM
1570	First	51.85	55.06	68 86+0.00	7053 ± 0.04	70.60+0.03
1571	Second	52 44	55.00	68.88 ± 0.00	70.33 ± 0.04 70.72+0.07	70.00±0.03 70.83+0.02
1572	Third	61.42	59.99	69.18 ± 0.00	71.59 ± 0.03	70.03±0.02 71.73±0.02
1573	Fourth	61.46	60.20	69.12 ± 0.00	$71.64{\pm}0.05$	71.83±0.03
1574	Fifth	63.30	63.51	$75.45 {\pm} 0.00$	$77.37{\pm}0.08$	77.37±0.05
1575				Obsolete		
1576	First	62 87	63 03	7575 ± 0.00	76.62 ± 0.06	76 99+0 03
1577	Second	61 55	61.09	69.17 ± 0.00	70.02 ± 0.00 71.61+0.11	70.99 ± 0.05 71.80 ±0.06
1578	Third	61.67	61.20	69.27 ± 0.00	71.79 ± 0.06	71.78 ± 0.01
1579	Fourth	52.37	54.89	68.56 ± 0.00	70.64 ± 0.06	70.80±0.02
1580	Fifth	51.80	54.30	$69.00{\pm}0.00$	$71.08{\pm}0.04$	$71.09{\pm}0.02$
1581			R	eannearing		
1582	First	61.91	59.99	69.43 ± 0.00	70.67 ± 0.02	70.78+0.03
1583	Second	53.20	55.59	70.68 ± 0.00	70.64 ± 0.78	71.68±0.03
1584	Third	62.08	62.23	64.29 ± 0.03	71.01 ± 0.04	71.09±0.02
1585	Fourth	53.39	54.93	$56.87 {\pm} 0.62$	71.98±0.04	$71.76 {\pm} 0.03$
1586	Fifth	62.28	62.42	$54.96{\pm}1.36$	$71.26{\pm}0.05$	71.26±0.01
1587						

1566 Table 13: The mean±standard deviation of the balanced accuracy of sudden, obsolete, and reap-1567 pearing experiments presented in Figure 6 and 7 and corresponds to the SUSY dataset.

Table 14: The exact values represented in the experiment 'Learning Without Forgetting' in the rightmost graph of Figure 7. This result corresponds to the SUSY dataset.

Model	First	Second	Third	Fourth	Fifth
packetLSTM	$70.78 {\pm} 0.03$	$71.68 {\pm} 0.03$	$71.09{\pm}0.02$	71.76±0.03	71.26±0.01
packetLSTM-Retraining	$70.78 {\pm} 0.03$	$71.68 {\pm} 0.03$	$71.00 {\pm} 0.02$	71.71 ± 0.02	$71.18 {\pm} 0.03$

1594 1595

1603

1588

1589

1596 of all models declines as features become obsolete in each subsequent interval. Interestingly, an 1597 increase in performance from the fourth to the fifth interval is observed in models such as OVFM, 1598 Aux-Drop, and packetLSTM. A plausible explanation for this phenomenon may relate to the nature of the SUSY data, where the exclusion of feature 6 notably impacts performance more severely in the fourth interval. The packetLSTM consistently outperforms other models throughout each data interval, followed by Aux-Drop, demonstrating its robust capability in handling obsolete features effectively.

Reappearing The performance of packetLSTM, Aux-Drop, and OLVF improves when previously 1604 encountered sets of features reappear, as depicted in Figure 7. Furthermore, packetLSTM exhibits performance gains of 0.09, 0.05, and 0.08 in the third, fourth, and fifth intervals, respectively, com-1606 pared to its retrained counterpart, as shown in the rightmost graph of Figure 7 and Table 14. This enhancement reaffirms the 'learning without forgetting' capability of packetLSTM. Notably, pack-1608 etLSTM is the sole method that consistently demonstrates the 'learning without forgetting' capabil-1609 ity across both the HIGGS and SUSY datasets. Moreover, packetLSTM consistently gives the best 1610 result in four out of five data intervals, followed by Aux-Drop, as detailed in Table 13, underscoring 1611 its superior performance.

1612 1613 1614

Т CATASTROPHIC FORGETTING IN ONLINE LEARNING VERSUS ONLINE CONTINUAL LEARNING

1615 1616

1617 In this article, we present the 'learning without forgetting' capability of packetLSTM. The 'learning without forgetting' capability is also referred to as mitigating catastrophic forgetting in the existing 1618 literature (Hoi et al., 2021). We explore this capability in an online learning setting, where the model 1619 receives and processes instances sequentially without any buffer storage.

1620 The problem of catastrophic forgetting is extensively studied in a parallel and widely recognized 1621 field of online continual learning (Li & Hoiem, 2017; Wang et al., 2024). Online continual learning 1622 is the field of machine learning where a model continually learns new tasks. In this scenario, tasks 1623 are delivered sequentially, yet all the instances for each task are presented at once. Consequently, the learning for each task is conducted in an offline setting, utilizing all data related to a task before 1624 transitioning to the subsequent one. For example, a model may initially learn from data related to 1625 task 1 in an offline batch mode. Subsequently, it proceeds to learn from data pertaining to task 2 in 1626 a similar batch mode, and this pattern continues. The introduction of new tasks may result in the 1627 addition of new classes, a shift in data distribution, or the arrival of an entirely new task domain. In 1628 the context of online continual learning, catastrophic forgetting refers to the challenge a model faces 1629 in retaining its ability to perform previously learned tasks. 1630

Although online continual learning shares similarities with online learning regarding sequential task learning, it differs significantly in its execution within individual tasks. In online continual learning, the model adheres to a batch training paradigm, which diverges from the principles of online learning algorithms. Consequently, the phenomenon of catastrophic forgetting, as observed in online learning, distinctly differs from that in online continual learning. It is worth noting that the packetLSTM framework, with certain modifications, may potentially be adapted for use in online continual learning to address catastrophic forgetting. However, this application extends beyond the scope of our article.

1639

1652

- 1640 U TRANSFORMER
- 1641 1642 U.1 PADDING

1643 The two input padding methods are:

- Only Values: For each dataset, zeros are added to the sequence following the available feature values until f_l , where $f_l = N$, where N is the total number of features. Note that this contradicts the sixth characteristic of haphazard inputs. However, to compare packetLSTM with Transformer, we assume that N is known. The specific N for each dataset is defined in the column labeled '#Features' in Table 5 of the Appendix. An example of an input instance at time t where feature 1 and j is available is $[x_1^t, x_j^t, 0, ..., 0]$.
 - *Pairs*: Each available feature value is paired with its corresponding feature ID, and the sequence of these pairs is padded with zeros till f_l . The feature ID ranges from 1 to N, and $f_l = 2N$. An Example is $[x_1^t, 1, x_j^t, j, 0, ..., 0]$

1654Padded inputs are initially processed through a linear embedding layer, which outputs embedding of1655dimension d. The input dimension for this layer is N and 2N for Only Values and Pairs, respectively.1656The embedding is then passed to an encoder. The encoder consists of n_l encoder layers. Each1657encoder layer comprises n_h heads and maintains dimensions of size d for both input and output.1658The encoder's output feeds into the same fully connected neural network utilized in packetLSTM,1659consisting of a linear layer of dimensions (d, d), a ReLU activation, and another linear layer leading1660to the output classes, culminating in a softmax layer for predictions.

The hyperparameters search is performed sequentially, as in the packetLSTM. We searched the optimal value of d among 32, 64, 128, 256, and 512, n_h among 1, 2, 4, 8, and 16, n_l among 1, 2, 3, and 4. Similar to packetLSTM, we tested a range of learning rates (0.001, 0.0005, 0.0001, 0.00005). We further searched in the vicinity of the optimal learning rate found in the previous step. Optimal hyperparameters include a d of 128 for magic04 and 32 for other datasets, n_h of 4 for SUSY and imdb, 8 for magic04, and 16 for a8a and HIGGS. The n_l is 1 for all the datasets, with the learning rate of 0.0002 for magic04 and a8a and 0.0001 for SUSY, HIGGS, and imdb. The Z-score is employed for streaming normalization.

- 1668
- 1669 U.2 NATURAL LANGUAGE

We experimented with two models, DistilBERT (Sanh et al., 2019) and BERT (Devlin et al., 2019),
 to process haphazard inputs. The inputs to the model are created in two ways:

- 1673
- Values: Sequences of available feature values, formatted as a string (example " $x_1^t x_j^t$ ").

Table 15: Performance of DistilBERT and BERT on haphazard inputs. The synthetic dataset is considered for p = 0.5, and the performance of 3 runs is reported here. Since the mean balanced accuracy is around 50 in all cases, we did not conduct further experiments.

Model	Data Type	Metric	magic04	imdb	a8a	SUSY	HIGGS
		bAcc	50.89±1.27	$50.48 {\pm} 0.06$	$50.0 {\pm} 0.00$	$50.02 {\pm} 0.04$	$50.0 {\pm} 0.00$
		Acc	65.21 ± 0.57	50.48 ± 0.06	75.92 ± 0.0	54.21 ± 0.02	52.95 ± 0.0
	Values	ROC	50.69 ± 2.02	50.08 ± 0.03	49.64 ± 0.03	50.13 ± 0.08	50.05 ± 0.01
		PRC	65.17±1.26	50.09±0.08	83.01±0.19	45.94±0.14	52.98±0.01
		Err	6616.67±107.95	$123/9.0\pm15.75$	/841.6/±0.4/	457861.67±236.68	470490.67 ± 22.9
DistilBERT		Time	487.88±0.05	/03.84±3.01	800.3±3.32	19/0/.84±048.30	20383.08±2032.13
		bAcc	50.01 ± 0.02	50.35 ± 0.06	50.0 ± 0.00	-	-
		Acc	55.41 ± 13.31	50.35 ± 0.06	75.92 ± 0.0		
	Input Pairs	ROC	50.01 ± 0.07	50.12 ± 0.03	49.62 ± 0.01		
		Fre	73.83 ± 7.09 8481 67 ± 2530.07	30.00 ± 0.02 12412 33 ± 14.38	85.01 ± 0.00 7841 67 ± 0.47		
		Time	496.47+3.26	846.39+8.61	938.85+2.88		
		bAcc	50.0+0.01	50.26+0.02	50.0+0.00		
		Acc	64.83±0.0	50.26 ± 0.02 50.26 ± 0.22	75.92 ± 0.00		
	17.1	ROC	50.0 ± 0.05	50.04 ± 0.02	49.67±0.06		
	values	PRC	80.71±1.03	$50.02 {\pm} 0.01$	$82.87 {\pm} 0.02$		
		Err	$6689.0 {\pm} 0.82$	12434.0 ± 55.37	$7842.0 {\pm} 0.0$		
BERT		Time	619.53 ± 10.21	1137.5 ± 3.14	1470.69 ± 6.11		
		bAcc	$50.0 {\pm} 0.00$	$50.44 {\pm} 0.02$	$50.0 {\pm} 0.00$	-	-
		Acc	$64.82 {\pm} 0.01$	50.44 ± 0.02	75.92 ± 0.0		
	Input Pairs	ROC	49.93 ± 0.03	50.02 ± 0.06	49.62 ± 0.01		
	1	PRC	80.91±1.15	50.07±0.04	82.95±0.11		
		Err	6691.33 ± 1.25	12390.33 ± 4.64	/842.0±0.0		
		Time	031.01±2.44	1000.33±3.00	1203.76±4.31		

1696

1677

169

1698

• Input Pairs: Sequences where each feature value is paired with its feature ID (example " $[x_1^t, 1] [x_i^t, j]$ ").

We utilize the default hidden size of both DistilBERT and BERT models, and learning rates are determined by hyperparameter search among 0.001, 0.0005, 0.0001, and 0.00005. The balanced accuracy on the imdb and synthetic datasets (with p = 0.5) is around 50 for both models (see Table 15), with the highest being 50.89 for DistilBERT using *Values* on the magic04 dataset. Given that the models did not learn to classify haphazard inputs and considering the extensive computation time required for the BERT models, we did not conduct further experiments.

1705 1706

1707

U.3 SET TRANSFORMER

We utilize the encoder and decoder provided in the Set Transformer article (Lee et al., 2019), which is permutation invariant and handles variable length inputs. The hyperparameters are hidden size, number of heads, number of induction (inducing points), and learning rate. Please refer to Set Transformer (Lee et al., 2019) for more details about the architecture.

We conduct a hyperparameter search similar to packetLSTM. We searched the optimal value of 1712 hidden size and the number of induction points among 32, 64, 128, 256, and 512, and the number 1713 of heads among 1, 2, 4, and 8. Similar to packetLSTM, we tested a range of learning rates (0.001, 1714 0.0005, 0.0001, 0.00005). We further searched in the vicinity of the optimal learning rate found in 1715 the previous step. The optimal hidden sizes are 64 for magic04 and 32 for the rest of the dataset. 1716 The best value of the number of induction points is 64 for magic04 and HIGGS, 32 for imdb and 1717 a8a, and 256 for SUSY. The optimal number of layers is 1 for magic04 and HIGGS, 2 for imdb and 1718 a8a, and 8 for SUSY. The best learning rate is 0.0001, 0.00007, 0.0002, 0.00008, and 0.00006 for 1719 magic04, imdb, a8a, SUSY, and HIGGS, respectively.

1720

1721 U.4 HAPTRANSFORMER

1723 We create embeddings of dimension (d) for each feature, initialized with He initialization (He et al., 1724 2015), and designed to be learnable to capture feature representations effectively. At each time 1725 instance, only embeddings for available features are utilized and are collectively denoted as E_t . 1726 Given that the shape of E_t would be $(1, |\mathbb{F}^t|, d)$, it varies in length due to the changing number of 1727 features at time t ($|\mathbb{F}^t|$). Therefore, the decoder of the Set Transformer (Lee et al., 2019), which 1728 handles variable length inputs, is utilized to process E_t . The decoder accepts an input of size d for 1729 Table 16: Performance of packetLSTM with complete features 1730 (p = 1) on synthetic datasets. All the metrics are reported as mean \pm standard deviation of five runs. The bAcc, Acc, 1731 ROC, PRC, Err, and Time stand for balanced accuracy, accu-1732 racy, AUROC, AUPRC, number of errors, and execution time, 1733 respectively. 1734

Μ	letric	magic04	a8a	SUSY	HIGGS
b	Acc	79.65±0.13	73.57±0.15	$77.42 {\pm} 0.05$	62.66±0.18
Α	cc	83.68±0.12	$83.37 {\pm} 0.08$	78.11±0.05	62.98±0.19
R	OC	87.79±0.06	87.54 ± 0.11	$84.92 {\pm} 0.05$	$67.95 {\pm} 0.28$
Pl	RC	90.49±0.11	$95.49 {\pm} 0.06$	85.17±0.06	69.60 ± 0.24
E	rr	3103.2±21.93	5415.4±26.19	218943.6±492.38	370198±1863.46
Ti	ime	$79.02{\pm}8.97$	$261.74{\pm}14.66$	$4146.19 {\pm} 98.9$	5589.51±94.65



Figure 8: Performance of packetL-STM on various probability of features availability.

1743 each feature and includes n_h number of heads. The output from the decoder is then passed through 1744 a softmax layer to generate predictions. 1745

We conduct a hyperparameter search for d, n_h , and learning rate following the methodology used 1746 for packetLSTM. We determined the optimal d among 32, 64, 128, 256, and 512, and n_h among 1747 1, 2, 4, and 8. The best learning rate was determined among 0.001, 0.0005, 0.0001, and 0.00005. 1748 We further searched in the vicinity of the optimal learning rate found in the previous step. Optimal 1749 values identified include a d of 64, 512, 256, 32, and 64, n_b of 2, 1, 2, 2, and 4 for magic04, imdb, 1750 a8a, SUSY, and HIGGS, respectively. The learning rate is found to be 0.0005, 0.00007, 0.00009, 1751 0.00009, and 0.00008 for magic04, imdb, a8a, SUSY, and HIGGS, respectively.

1752

1728

1

1741 1742

1753 1754

V HAPHAZARD INPUTS VS OTHER FIELDS OF VARYING FEATURE SPACE

1755 The field of varying feature space is also studied as feature evolvable streams (Hou et al., 2017; 1756 Zhang et al., 2020) and incremental and decremental features (Hou & Zhou, 2018; Dong et al., 1757 2021). However, both these fields assume some form of structure in their data stream, which con-1758 tradicts the characteristics of haphazard inputs. Specifically, the feature evolvable streams work in 1759 batches and assume that there is an overlap period between the transition where old features vanish 1760 and new features occur. Similarly, incremental and decremental features assume that the data arrives 1761 in batches, where the initial batch consists of both vanishing and surviving features and subsequent 1762 batches encompass surviving and newly augmented features. The assumption of the batch and the structure of the data within the batch limits the applicability of both fields in haphazard inputs. 1763

1764 1765

1766

W PACKETLSTM WITH COMPLETE FEATURES

1767 We set p = 1 in the synthetic datasets to determine the upper bound performance of packetLSTM 1768 when all the features are present at each time instance. The corresponding results are provided in 1769 Table 16. We present a comparison of packetLSTM's performance across varying p values (0.25, 0.5, 1770 0.75, and 1) in Figure 8. Unsurprisingly, the performance of packetLSTM increases with decreasing 1771 haphazardness (i.e., increasing p value) in the data. The most significant decline in performance occurs at p = 0.25, which is expected due to the increased haphazardness. 1772

- 1773
- 1774
- 1775 1776
- 1777
- 1778
- 1779
- 1780
- 1781