

Convolutional Neural Network and Long Short-Term Memory Models for Ice-Jam Predictions

Fatemehalsadat Madaeni¹, Karem Chokmani¹, Rachid Lhissou¹, Saied Homayouni¹, Yves Gauthier¹, and Simon Tolszczuk-Leclerc²

¹INRS-ETE, Université du Québec, Québec City, G1K 9A9, Canada

²EMGeo Operations, Natural Resources Canada, Ottawa, K1S 5K2, Canada

Correspondence to: Fatemehalsadat Madaeni (Fatemehalsadat.Madaeni@ete.inrs.ca)

Abstract. In cold regions, ice jams frequently result in severe flooding due to a rapid rise in water levels upstream of the jam. Sudden floods resulting from ice jams threaten human safety, and cause damage to properties and infrastructure. Hence, ice-jam prediction tools can give an early warning to increase response time and minimize the possible damages. However, ice-jam prediction has always been a challenge as there is no analytical method available for this purpose. Nonetheless, ice jams form when some hydro-meteorological conditions happen, a few hours to a few days before the event. Ice-jam prediction can be addressed as a binary multivariate time-series classification. Deep learning techniques have been widely used for time-series classification in many fields such as finance, engineering, weather forecasting, and medicine. In this research, we successfully applied convolutional neural networks (CNN), long short-term memory (LSTM), and combined convolutional-long short-term memory (CNN-LSTM) networks to predict the formation of ice jams in 150 rivers in the Province of Quebec (Canada). We also employed machine learning methods including support vector machine (SVM), k-nearest neighbors classifier (KNN), decision tree, and multilayer perceptron (MLP) for this purpose. The hydro-meteorological variables (e.g., temperature, precipitation, and snow depth) along with the corresponding jam or no-jam events are used as model inputs. Ten percent of the data was excluded from the model and set aside for testing, and 100 reshuffling and splitting iterations were applied to 80% of the remaining data for training and 20% for validation. The developed deep learning models achieved improvements in performance in comparison to the developed machine learning models. The results show that the CNN-LSTM model yields the best results in the validation and testing with F1 scores of 0.82 and 0.92, respectively. This demonstrates that CNN and LSTM models are complementary, and a combination of both further improves classification.

1 Introduction

Predicting ice-jams gives an early warning of possible flooding events, but there is no analytical solution to predict these events due to the complex interactions between the hydro-meteorological variables (e.g., temperature, precipitation, snow depth, and solar radiation) involved. To date, a small number of empirical and statistical prediction methods such as threshold methods, multi-regression models, logistic regression models, and discriminant function analysis have been developed for ice jams (Barnes-Svarney and Montz, 1985; Mahabir et al., 2006; Massie et al., 2002; White, 2003; White and Daly, 2002; Zhao et al., 2012). However, these methods are site-specific and have high rates of false-positive errors (White, 2003). The numerical models developed for ice-jam prediction (e.g., ICEJAM (Flato and Gerard, 1986, cf.; Carson et al., 2011), RIVJAM (Beltaos, 1993), HEC-RAS (Brunner, 2002), ICESIM

38 (Carson et al., 2001 and 2003), and RIVICE (Lindenschmidt, 2017)) have several limitations. More particularly, the
39 mathematical formulations used in these models are complex and need many parameters, which are often unavailable
40 as they are challenging to measure in ice conditions. The subsequent simplifications necessary to model application
41 decrease model accuracy (Shouyu & Honglan, 2005). A detailed overview of the previous models for ice-jam
42 prediction based on hydro-meteorological data is presented in Madaeni et al. (2020).

43 Prediction of ice-jam occurrence can be considered as a binary multivariate time-series classification (TSC) problem
44 when the time series of various hydrometeorological variables can be used to classify jam or no jam events. Time-
45 series classification (particularly multivariate) has been widely used in various fields, including biomedical
46 engineering, clinical prediction, human activity recognition, weather forecasting, and finance. Multivariate time series
47 provide more patterns and improve classification performance compared to univariate time series (Zheng et al., 2016).
48 Time-series classification is one of the most challenging problems in data mining and machine learning.

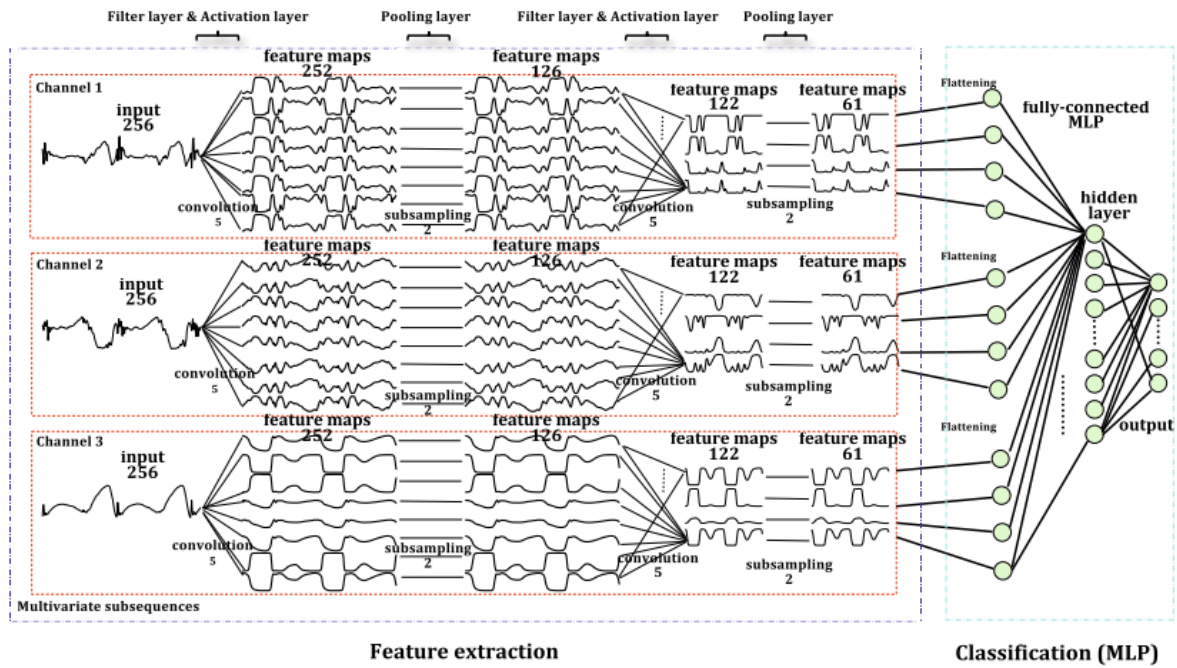
49 Most existing TSC methods are feature-based, distance-based, or ensemble methods (Cui et al., 2016). Feature
50 extraction is challenging due to the difficulty of handcrafting useful features to capture intrinsic characteristics from
51 time-series data (Karim et al., 2019a; Zheng et al., 2014). Hence, distance-based methods work better in TSC (Zheng
52 et al., 2014). Among the hundreds of methods developed for TSC, the leading classifier with the best performance
53 was an ensemble nearest neighbor approach with dynamic time warping (DTW) (Fawaz et al., 2019a; Karim et al.,
54 2019a).

55 In the k-nearest neighbors (KNN) classifier, the test instance is classified by the majority vote of its k-nearest neighbors
56 in the training dataset. The entire dataset is necessary to make a prediction based of KNN, which requires a lot of
57 processing memory. Hence, it is computationally expensive and time-consuming when the database is large. It is also
58 sensitive to irrelevant features and data scale. Furthermore, the number of neighbors included in the algorithm should
59 be carefully selected. The KNN classifier is very challenging to be used for multivariate TSC. The dynamic time
60 warping approach is a robust alternative for Euclidean distance (the most widely used time-series distance measure)
61 to measure the similarity between two time series by searching for an optimal alignment (minimum distance) between
62 them (Zheng et al., 2016). However, the combined KNN with DTW is time-consuming and inefficient for long
63 multivariate time series (Lin et al., 2012; Zheng et al., 2014). Traditional classification and data mining algorithms
64 developed for TSC have high computational complexity or low prediction accuracy. This is due to the size and inherent
65 complexity of time series, seasonality, noise, and feature correlation (Lin et al., 2012).

66 There are some machine learning methods available for TSC such as KNN and support vector machine (SVM).
67 However, the focus of this research is on the deep learning models that have greatly improved sequence classification
68 and that perform well with multivariate TSC. Deep learning methods work with 2-D multivariate time series and their
69 deeper architecture could further improve classification especially for complex problems. This explains why deep
70 learning methods generally have more accurate and robust results than other currently used methods (Wu et al., 2018).
71 However, their training is more time consuming and their interpretation is more difficult.

72 Deep learning involves neural networks that use multiple layers where nonlinear transformation is used to extract
73 higher-level features from the input data. Although deep learning has recently shown promising performance in
74 various fields such as image and speech recognition, document classification, and natural language processing, only a

75 few studies were dedicated to using deep learning for TSC (Gu et al., 2018; Fawaz et al., 2019a). Various studies show
 76 that deep neural networks significantly outperform the ensemble nearest neighbor with DTW (Fawaz et al., 2019a).
 77 The main benefit of deep learning networks is automatic feature extraction, which reduces the need for expert
 78 knowledge and removes engineering bias during classification as the probabilistic decision (e.g., classification) is
 79 taken by the network (Fawaz et al., 2019b).
 80 The most widely used deep neural networks for TSC are multi-layer perceptron (MLP; i.e., fully connected deep
 81 neural networks), convolutional neural networks (CNNs), and long short-term memory networks (LSTM). The
 82 application of CNNs for TSC has recently become increasingly popular, and different types of CNN are being
 83 developed with superior accuracy for this purpose (Cui et al., 2016). Zheng et al. (2014) and Zheng et al. (2016)
 84 introduce a multi-channel deep convolutional neural network (MC-DCNN) for multivariate TSC, where each variable
 85 (i.e., univariate time series) is trained individually to extract features and finally concatenated using an MLP to perform
 86 classification (Fig. 1). The authors showed that their model achieves a state-of-the-art performance in terms of
 87 efficiency and accuracy on a challenging dataset. The drawback of their model and similar architectures (e.g.,
 88 Devineau et al., 2018a) is that they do not capture the correlation between variables as the feature extraction is carried
 89 out separately for each variable.

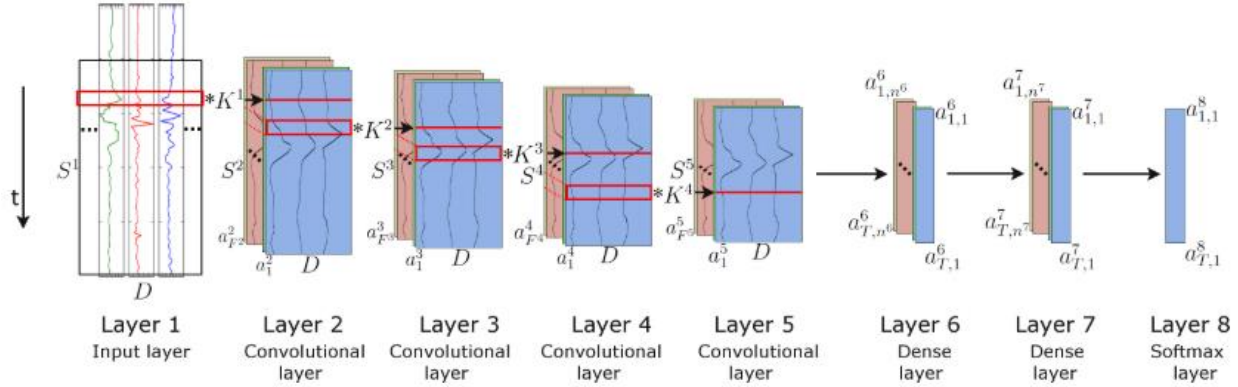


90 **Feature extraction** **Classification (MLP)**

91 **Figure 1. A two-stage MC-DCNN architecture for activity classification. This architecture consists of a three-channel input,**
 92 **two filter layers, two pooling layers, and two fully connected layers (after Zheng et al., 2014).**

93 Brunel et al. (2019) present CNNs adapted for TSC in cosmology using 1-D filters to extract features from each
 94 channel over time and a convolution in depth to capture the correlation between the channels. They compared the

95 results from LSTMs with those from CNNs, and demonstrated that CNNs had better results. Nevertheless, both deep
 96 learning approaches are very promising.
 97 The combination of CNNs and LSTM units has already yielded promising results in problems requiring temporal
 98 information classification, such as human activity recognition (Li et al., 2017; Mutegeki and Han, 2020), text
 99 classification (Luan and Lin, 2019; She and Zhang, 2018; Umer et al., 2020), video classification (Lu et al., 2018 and
 100 Wu et al., 2015), sentiment analysis (Ombabi et al., 2020; Sosa, 2017; Wang et al., 2016; Wang et al., 2019), typhoon
 101 formation forecasting (Chen et al., 2019), and arrhythmia diagnosis (Oh et al., 2018). In this architecture, convolutional
 102 operations capture features and LSTMs capture time dependencies on the extracted features. Ordóñez and Roggen
 103 (2016) propose a deep convolutional LSTM model (DeepConvLSTM) for activity recognition (Fig. 2). Their results
 104 are compared to the results from standard feedforward units showing that DeepConvLSTM reaches a higher F1 score
 105 and better decision boundaries for classification. Furthermore, they noticed that the LSTM model is also promising
 106 with relatively small datasets. Furthermore, LSTMs perform better with longer temporal dynamics, whereas the
 107 convolution filters can only capture the temporal dependencies dynamics within the length of the filter.



108

109 **Figure 2. Architecture of the DeepConvLSTM framework used for activity recognition (after Ordóñez and Roggen, 2016).**

110 The project presented in this paper is part of a greater project called DAVE, which aims at developing a tool for
 111 regional ice jam watches and warnings, based on the integration of three aspects: current ice cover conditions; hydro-
 112 meteorological patterns associated with breakup ice jams; and channel predisposition to ice-jam formation. The
 113 outputs will be used to develop an ice jam monitoring and warning module that will transfer the knowledge to the end
 114 users managing ice jam consequences.

115 The objective of this research is to develop deep learning models to predict breakup ice-jam events to be used as an
 116 early warning system of possible flooding. While most TSC research in deep learning is performed on 1-D channels
 117 (Hatami et al., 2018), our approach consists of using deep learning frameworks for multivariate TSC, applied to ice-
 118 jam prediction. Through our comprehensive literature review, we noticed that CNN (e.g., Brunel et al., 2019; Cui et
 119 al., 2016; Devineau et al., 2018b; Kashiparekh, 2019; Nosratabadi et al., 2020; Yan et al., 2020; Yang et al., 2015; Yi
 120 et al., 2017; Zheng et al., 2016), LSTM (e.g., Fischer and Krauss, 2018; Lipton et al., 2015; Nosratabadi et al., 2020;
 121 Torres et al., 2021), and a combined CNN-LSTM (e.g., Karim et al., 2017; Livieris et al., 2020; Ordóñez and Roggen,
 122 2016; Sainath et al., 2015; Xingjian et al., 2015) have been widely used for TSC. Numerous applications of CNN,
 123 LSTM, and their hybrid versions are currently used in the field of hydrology (Althoff et al., 2021; Apaydin et al.,

124 2020; Barzegar et al., 2021, 2020; Kratzert et al., 2018; Wunsch et al., 2020; Zhang et al., 2018). Although deep
125 learning methods seem promising to address the requirements of ice-jam predictions, none of these methods yet have
126 been explored for ice-jam prediction.

127 Although machine learning methods have been widely used in time series forecasting of hydro-meteorological data,
128 they have been used less frequently in the prediction of ice jams (Graf et al., 2022). Semenova et al. (2020) used KNN
129 to predict ice jams using hydro-meteorological variables such as precipitation, snow depth, water level, water
130 discharge, and temperature. They developed their model with data collected from the confluence of Sukhona River
131 and Yug River in Russia between 1960 and 2016 and achieved accuracy of 82%. Sarafanov et al. (2021) presented an
132 ensemble-based model of machine learning methods and a physical snowmelt-runoff model to account for the
133 advantages of physical models (interpretability) and machine learning models (low forecasting error). Their hybrid
134 models proposed an automated approach for short-term flood forecasting in Lena River, Poland, using hydro-
135 meteorological variables (e.g., maximum water level, mean daily water and air temperature, mean daily water
136 discharge, relative humidity, snow depth, and ice thickness). They applied an automated machine learning approach
137 based on the evolutionary algorithm to automatically identify machine learning models, tune hyperparameters, and
138 combine stand-alone models into ensembles. Their model was validated on ten hydro gauges for two years, showing
139 that the hybrid model is much more efficient than stand-alone models with a Nash–Sutcliffe efficiency coefficient of
140 0.8. Graf et al. (2022) developed an MLP and extreme gradient boosting model to predict ice jams with data from
141 1983 to 2013, in Warta River, Poland. They employed water and air temperatures, river flow, and water level as inputs
142 to their models, showing that both machine learning methods provide promising results. In Canada, De Coste et al.
143 (2021) developed a hybrid model including a number of machine learning models (e.g., KNN, SVM, random forest,
144 and gradient boosting) for St. John River (New Brunswick). The most successful ensemble model combining 6
145 different member models was produced with a prediction accuracy of 86% over 11 years of record.

146 We developed three deep learning models; a CNN, an LSTM, and a combined CNN-LSTM for ice-jam predictions,
147 and compared the results. The previous studies show that these models successfully capture features, the correlation
148 between features (through convolution units) and time dependencies (through memory units) which are subsequently
149 used for TSC. The combined CNN-LSTM can reduce errors by compensating for the internal weaknesses of each
150 model. In the CNN-LSTM model, CNNs capture features, then the LSTMs identifies time dependencies on the
151 captured features.

152 Furthermore, we also developed some machine learning methods as simpler methods for ice-jam prediction. And their
153 results are compared with those obtained from the deep learning models.

154 **2 Materials and Methods**

155 **2.1 Data and study area**

156 It is known that specific hydro-meteorological conditions lead to ice-jam occurrence (Turcotte and Morse, 2015; and
157 White, 2003). For instance, breakup ice jams occur when a period of intense cold is followed by a rapid peak discharge
158 resulting from spring rainfall and snowmelt runoff (Massie et al., 2002). Accumulated freezing degree days (AFDD)
159 can be used as a proxy for intense cold periods. Sudden spring runoff increase, however, is not often available at the

160 jam location and can be represented by liquid precipitation and snow depth a few days prior to ice-jam occurrence
 161 (Zhao et al., 2012). Prowse and Bonsal (2004) and Prowse et al. (2007) assessed various hydroclimatic explanations
 162 for river ice freeze-up and breakup, concluding that shortwave radiation is the most critical factor influencing the
 163 mechanical strength of ice and consequently the possibility of breakup ice jams to occur. Turcotte and Morse (2015)
 164 explain that accumulated thawing degree day (ATDD), an indicator of warming periods, partially covers the effect of
 165 shortwave radiation. In previous studies addressing ice-jam and breakup predictions, discharge and changes in
 166 discharge, water level and changes in water level, AFDD, ATDD, precipitation, solar radiation, heat budget, and
 167 snowmelt or snowpack are the most frequently used variables (Madaeni et al., 2020).
 168 The inputs used in this study are historical ice-jam or no ice-jam occurrence (Fig. 3) as well as hydro-meteorological
 169 variables from 150 rivers in Quebec, namely liquid precipitation (mm), minimum and maximum temperature (°C),
 170 AFDD (from August 1 of each year; °C), ATDD (from January 1 of each year; °C), snow depth (cm) and net radiation
 171 ($W m^{-2}$). The net solar radiation, which represents the total energy available to influence the climate, is calculated as
 172 the difference between incoming and outgoing energy. If the median temperature is greater than 1, precipitation is
 173 considered to be liquid. The statistics of hydro-meteorological data used in the models are presented in Table 1. The
 174 source, time period, and spatial resolution of the input variables are shown in Table 2.

175 **Table 1. Statistics of hydro-meteorological variables used in the models.**

Statistic	Liquid Precipitation (mm)	Minimum temperature (°C)	Maximum temperature (°C)	Net radiation ($W m^{-2}$)	ATDD (°C)	AFDD (°C)	Snowdepth (cm)
Minimum	0.00	-40.00	-25.97	-67.77	0.00	-	0.00
Maximum	50.87	12.05	27.48	222.69	280.82	-35.41	121.86
Mean	1.04	-9.41	0.98	59.75	8.83	-898.48	15.99
Median	0.00	-7.73	1.68	59.41	1.27	-890.74	11.50

176

177

178 **Table 2. Source, duration, and spatial resolution of hydro-meteorological data used in the models.**

Data	Source	Duration	Spatial resolution
Min and Max temperature*	Daily Surface Weather Data (Daymet; Thornton et al., 2020)	1979–2019	1 km
Liquid precipitation	Canadian Precipitation Analysis (CaPA; Mahfouf et al., 2007)	2002–2019	10–15 km
Liquid precipitation	North American Regional Reanalysis (NARR; Mesinger et al., 2006)	1979–2001	30 km
Infrared radiation emitted by the atmosphere	North American Regional Reanalysis (NARR)	1979–2019	30 km
Infrared radiation emitted from the surface	North American Regional Reanalysis (NARR)	1979–2019	30 km
Snow depth	North American Regional Reanalysis (NARR)	1979–2019	30 km

179 * The average was used to derive the AFDD and the ATDD.

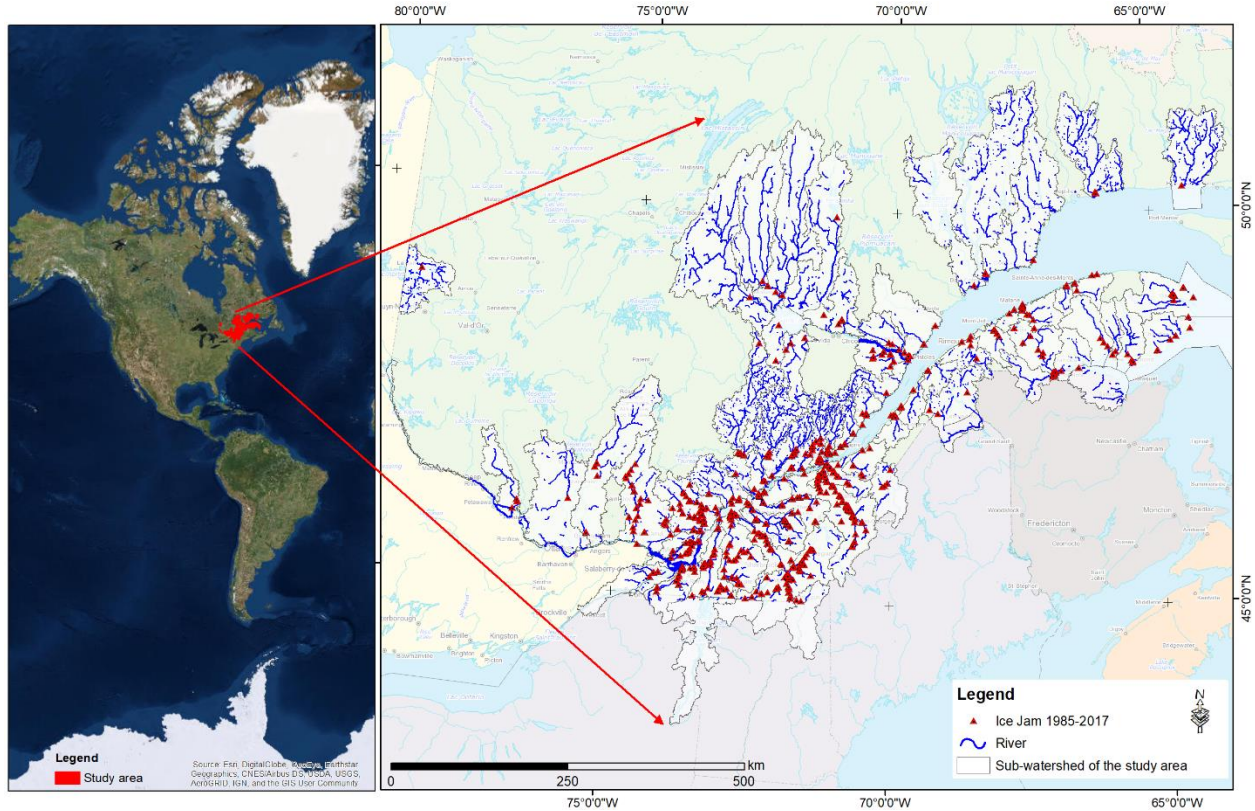
180

181 Ice-jam database was provided by the provincial public safety department (Ministère de la sécurité publique du
 182 Québec; MSPQ; Données Québec, 2021) for 150 rivers, mainly in the St. Lawrence River Basin. The database comes
 183 from digital or paper events reported by local authorities under the jurisdiction of the MSPQ from 1985 to 2014.
 184 Moreover, other data used to build this database were provided by field observations collected by the Vigilance/Flood

185 application from 2013 to 2019. It contains 995 recorded jam events that are not validated and contain many
186 inaccuracies, mainly in the toponymy of the rivers, location, dating, and jam event redundancy.

187 The names of the watercourse of several recorded jams are not given, wrong or misspelled. The toponymy of the rivers
188 was corrected using the National Hydrographic Network (NHN; National Hydrographic Network - Natural Resources
189 Canada (NRCan)), the GeoBase of the Quebec Hydrographic Network (National Hydro Network - NHN - GeoBase
190 Series - Natural Resources Canada), and the Toporama Web map service (The Atlas of Canada - Toporama - Natural
191 Resources Canada) of the Sector of Earth Sciences. Other manual corrections had to be carried out on ice jam data.
192 For example, ice jams location is sometimes placed on the riverbanks at a small distance (less than 20 m) from the
193 river polygon. In this case, the location of the ice jam is manually moved inside the river polygon. In other cases, the
194 ice-jam point is placed further on the flooded shore at a distance between 20 m and 200 m from the true ice jam. This
195 was corrected based on images with very high spatial resolution, based on the sinuosity and the narrowing of the river,
196 the history of ice jams at the site, and press archives. In addition, some ice jams were placed too far from the river due
197 to wrong coordinates in the database. A single-digit correction in longitude or latitude returned the jam to its exact
198 location. There are certain cases where the date of jam formation is verified by searching press archives, notably when
199 the date of formation is missing or several jams with the same dates and close locations in a section of a river are
200 present.

201 The ice jam database contains many duplicates. This redundancy can be explained by the merging of two databases,
202 a double entry during ice jam monitoring, or numerous recordings for an ice jam that lasted for several days. To
203 remediate this, the duplicates were removed from the database. The corrected ice-jam database contains 850 jams for
204 150 rivers, mainly in southern Quebec (Fig. 3). Ice jams formed in November and December (freeze-up jams) are
205 removed from the model since the processes involved are different from breakup ice jams (included from January 15).
206 The final breakup ice-jam database used in this study includes 504 jam events.



207

208

Figure 3. Study area and historic ice-jam locations recorded in Quebec from 1985 to 2017.

209

210 **2.2 Machine learning models for TSC**

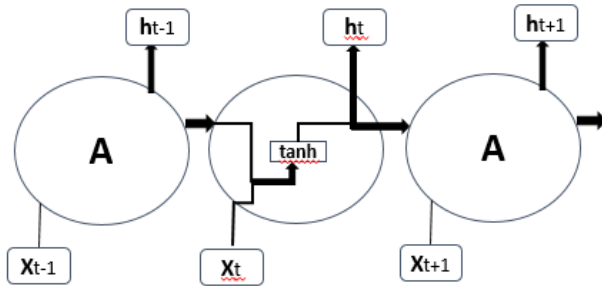
211 The most common machine learning techniques used for TSC are SVM (Rodríguez and Alonso, 2004; Xing and
 212 Keogh, 2010), KNN (Li et al., 2013; Xing and Keogh, 2010), decision tree (DT; Brunello et al., 2019; Jović et al.,
 213 2012), and multilayer perceptron (MLP; del Campo et al., 2021; Nanopoulos et al., 2001). These models and their
 214 applications in TSC are beyond the scope of this study and will not be further addressed.

215 We developed the mentioned machine learning methods and compared their results with those of deep learning
 216 models. After some trials and errors, the parameters that are changed from the default values for each machine learning
 217 model are as follows. We developed an SVM with a polynomial kernel with a degree of 5 that can distinguish curved
 218 or nonlinear input space. The KNN is used with 3 neighbors used for classification. The decision tree model is applied
 219 with all the default values. The shallow MLP is used with “lbfgs” solver (which can converge faster and perform better
 220 for small datasets), alpha of 1 e-5, and 3 layers with 7 neurons in each layer.

221 **2.3 Deep learning models for TSC**

222 The most common and popular deep neural networks for TSC are MLPs, CNNs, and LSTMs (Brownlee, 2018b; and
 223 Torres et al., 2021). Although it is very powerful approach, MLP networks are limited by the fact that each input (i.e.,
 224 time-series element) and output are treated independently, which means that the temporal or space information is lost
 225 (Lipton et al., 2015). Hence, an MLP needs some temporal information in the input data to model sequential data,

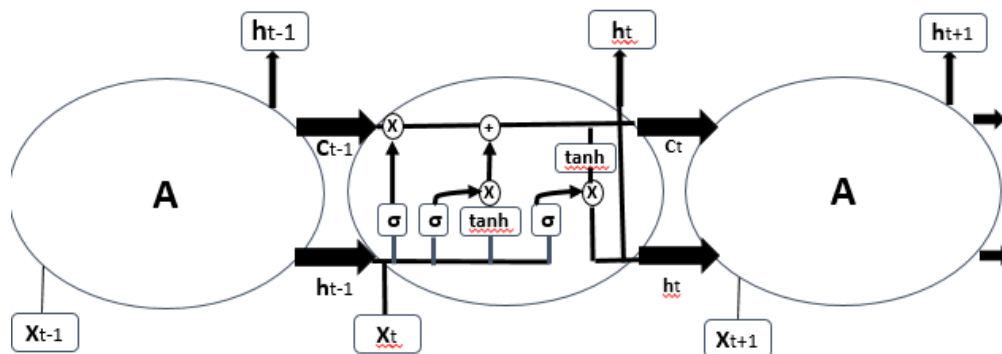
226 such as time series (Ordóñez and Roggen, 2016). In this regard, recurrent neural networks (RNNs) are specifically
 227 adapted to sequence data through the direct connections between individual layers (Jozefowicz et al., 2015). Recurrent
 228 neural networks perform the same repeating function with a straightforward structure, e.g., a single tanh (hyperbolic
 229 tangent) layer, for every input of data (x_t), and the inputs are related to each other within their hidden internal state,
 230 which allows it to learn the temporal dynamics of sequential data (Fig. 4).



231
 232 **Figure 4. An RNN with a single tanh layer, where A is a chunk of the neural network, x is input data, and h is output data.**

233 Recurrent neural networks are rarely used in TSC due to their significant limitations: RNNs mostly predict outputs
 234 for each time-series element; they are sensitive to the first examples seen, and it is challenging to capture long-term
 235 dependencies due to vanishing gradients, exploding gradients, and their complex dynamics (Devineau et al., 2018b;
 236 Fawaz et al., 2019b).

237 Long short-term memory RNNs are developed to improve the performance of RNNs by integrating a memory
 238 component to model long-term dependencies in time-series problems (Brunel et al., 2019; Karim et al., 2019a). Long
 239 short-term memory networks do not have the problem of exploding gradients. The LSTMs have four interacting neural
 240 network layers in a very special way (Fig. 5). An LSTM has three sigmoid (σ) layers to control how much of each
 241 component should be let through by outputting numbers between zero and one. The input to an LSTM goes through
 242 three gates (“forget”, “input”, and “output gates”) that control the operation performed on each LSTM block (Ordóñez
 243 and Roggen, 2016). The first step is the “forget gate” layer that gets the output of the previous block (h_{t-1}), the input
 244 for the current block (x_t), and the memory of the previous block (C_{t-1}) and gives a number between 0 and 1 for each
 245 number in the cell state (C_{t-1} ; Olah, 2015). The second step is called the “input gate” with two parts, a sigmoid layer
 246 that decides which values to be updated and a tanh layer that creates new candidate values for the cell state. The new
 247 and old memories are then be combined and control how much the new memory should influence the old memory.
 248 The last step (output gate) gives the output by applying a sigmoid layer deciding how much new cell memory goes to
 249 output, and multiplies it by tanh applied to the cell state (resulting in values between -1 and 1).



250
251 **Figure 5. Structure of LSTM block with four interacting layers.**

252 Recently, convolutional neural networks challenged the assumption that RNNs (e.g., LSTMs) have the best
253 performance when working with sequences. Convolutional NNs perform well when processing sequential data such
254 as speech recognition and sentence classification, similar to TSC (Fawaz et al., 2019b).

255 Convolutional NNs are the most widely used deep learning methods in TSC problems (Fawaz et al., 2019b). They
256 learn spatial features from raw input time series using filters (Fawaz et al., 2019b). Convolutional NNs are robust and
257 need a relatively small amount of training time compared to RNNs or MLPs. They work best for extracting local
258 information and reducing the complexity of the model.

259 A CNN is a kind of neural network with at least one convolutional (or filter) layer. A CNN usually involves several
260 convolutional layers, activation functions, and pooling layers for feature extraction, followed by dense layers used as
261 classifiers (Devineau et al., 2018b). The reason to use a sequence of filters is to learn various features from time series
262 for TSC. A convolutional layer consists of a set of learnable filters that compute dot products between local regions
263 in the input and corresponding weights. With high-dimensional inputs, it is impractical to connect subsequent neurons
264 to all the neurons of the previous layer. Therefore, each neuron in CNNs is connected to only a local region of the
265 input—receptive field—whose size is equivalent to that of the filter (Fig. 6). This feature reduces the number of
266 parameters by limiting the number of connections between neurons in different layers. The input is first convolved
267 with a learned filter, and then an element-wise nonlinear activation function is applied to the convolved results (Gu et
268 al., 2018). The pooling layer performs a downsampling operation such as maximum or average, reducing the spatial
269 dimension. One of the most powerful features of CNNs is called weight or parameter sharing, where all neurons share
270 filters (weights) in a particular feature map (Fawaz et al., 2019b). This allows to reduce the number of parameters.



271
 272 **Figure 6. Structure of a convolution layer including two sets of filters.**

273
 274 **2.4 Model libraries**

275 In an Anaconda (Analytics, C., 2016) environment, Python is implemented to develop CNN, LSTM, and CNN-LSTM
 276 networks for TSC. To build and train networks, the networks are implemented in Theano (Bergstra et al., 2010) using
 277 the Lasagne library (Dieleman et al., 2015). Other core libraries used for importing, preprocessing, training data, and
 278 visualization of results include Pandas (Reback et al., 2020), NumPy (Harris et al., 2020), Scikit-Learn (Pedregosa et
 279 al., 2011), and Matplotlib. PyLab (Hunter, J. D., 2007). The Spyder (Raybaut, 2009) package of Anaconda can be
 280 used as an interface otherwise, the command window can be used without any interface.

281 To develop machine learning models, Scikit-Learn machine learning libraries are used except for NumPy, Pandas,
 282 and Scikit-Learn preprocessing libraries.

283 **2.5 Preprocessing**

284 The data is comprised of variables with varying scales, and the machine-learning algorithms can benefit from rescaling
 285 the variables to all have one single scale. Scikit-learn (Pedregosa et al., 2011) is a free library for machine learning in
 286 Python that can be used to preprocess data. We examined Scikit-learn MinMaxScaler (scaling each variable between
 287 0 and 1), Normalizer (scaling individual samples to the unit norm), and StandardScaler (transforming to zero mean
 288 and unit variance separately for each feature). The results show that MinMaxScaler (Eq. (1)) leads to the most accurate
 289 results. Validation data rescaling is carried out based on minimum and maximum values of the train data.

290
$$X_{\text{scaled}} = \frac{X - X_{\text{min}}}{X_{\text{max}} - X_{\text{min}}} \quad (1)$$

291 For each jam or no jam event, the data from 15 days preceding the event was used to predict the event on the 16th
 292 day. A balanced dataset with the same number of jam and no-jam events (1008 small sequences) was generated,
 293 preventing the model from becoming biased to jam or no-jam events. The hydro-meteorological data related to no-
 294 jam events were constructed by extracting data from the reaches of no-jam records. Model generalizations were
 295 assessed by extracting 10% of data for testing. With the remaining data, 80% was used for training and 20% for
 296 validation. We used ShuffleSplit subroutine from the Scikit-learn library, where the database was randomly sampled
 297 during each re-shuffling and splitting iteration to generate training and validation sets. We applied 100 re-shuffling

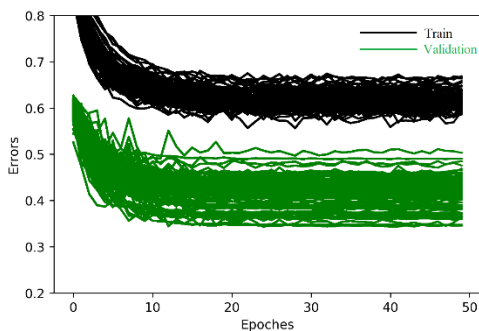
298 and splitting iterations for training and validation. There are 726, 181, and 101 small sequences with the size of (16,
299 7), 16 days of data for the seven variables; for training, validation, and test, respectively.

300 2.6 Training

301 Training a deep neural network with an excellent generalization to new unseen inputs is challenging. As a benchmark,
302 a CNN model with parameters and layers similar to previous studies (e.g., Ordóñez and Roggen, 2016) is developed.
303 The model shows underfitting or overfitting with various architectures and parameters. To overcome underfitting,
304 deeper models and more nodes can be added to each layer; however, overfitting is more challenging to resolve. The
305 ice-jam dataset for Quebec contains 1008 balanced sequence instances (with a length of 16), which is considered to
306 be a small amount of data in the context of deep learning. Deep learning models tend to overfit small datasets by
307 memorizing inputs rather than training. This is due to the fact that small datasets may not appropriately describe the
308 relationship between input and output spaces.

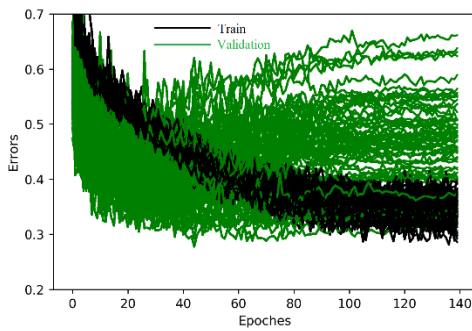
309 2.6.1 Overcome overfitting

310 There are various ways to resolve the problem of overfitting, including acquiring more data, data augmentation (e.g.,
311 cropping, rotating, and noise injection), dropout (Srivastava et al., 2014), early stopping, batch normalization (Ioffe
312 and Szegedy, 2015), and regularization. Acquiring more data is not possible with ice-jam records. We added the
313 Gaussian noise layer (from the Lasagne library), where the noise values are Gaussian-distributed with zero mean and
314 a standard deviation of 0.1 to the input. The noise layers applied to the CNN and LSTM models significantly overcome
315 the overfitting problem through data augmentation. However, the performance of the CNN-LSTM model dramatically
316 deteriorates when a noise layer is added (Fig. 7). Adding a noise layer to other layers does not improve any of the
317 developed models for ice-jam prediction.



318
319 **Figure 7. Train and validation errors over epochs for CNN-LSTM model with a noise layer.**

320 Early stopping is another efficient method that halts the training procedure at a point where further training would
321 decrease training loss, but validation loss starts to increase. Neural networks need a loss function to guide optimization
322 problem resolution. The loss function is similar to an objective function for process-based hydrological models.
323 Among the developed models, only LSTM needs early stopping at 40 epochs (Fig. 8). More detailed explanations
324 about the methods that are used in this study to overcome overfitting (e.g., batch normalization and L2 regularization)
325 can be found in the Appendix.



326

327

Figure 8. Train and validation errors over epochs for an LSTM model showing overfitting after 40 epochs.

328

2.6.2 Model hyperparameters

329

Finding hyperparameter values has been challenging due to the complex architecture of deep learning models and the large number of parameters (Garbin et al., 2020). The best model architecture was identified by assessing model performance with different layer and parameters such as the number of layers (noise, batch normalization, convolutional, pooling, LSTM, dropout, and dense) as well as different pooling sizes and strides, different batch sizes, various scaling of data (standardization and normalization), various filter sizes, number of units in LSTM and dense layers, the type of the activation functions, regularization and learning rates, weight decay and number of filters in convolutional layers. Hyperparameters are optimized through manual trial and error searches as grid search experiments suffer from poor dimensional coverage in dimensions (Bergstra and Bengio, 2012). Another reason is that manual experiments are much easier to conduct and interpret when investigating the effect of one hyperparameter of interest. The optimized hyperparameters are presented in Table 3. The most important parameters of the models are explained below and additional information is available in the Appendix.

340

Table 3. Common and selected values for different parameters of the models.

Parameter	Common values	Selected value	Source
Mini-batch size	16, 32, 64	16	Bengio (2012); Devineau et al. (2018b); Masters and Luschi (2018)
Number of convolution filters	32, 64, 128	128	Brownlee (2017); Maggiori et al. (2017)
Filter size	3, 5, 7	(5,1) and (5,3)	Devineau et al. (2018b); Maggiori et al. (2017)
Number of LSTM units	32, 64, 128	128	Brownlee (2017); Karim et al. (2019b); Ordóñez and Roggen (2016)
Number of dense layer units	16, 32, 128, 256	32	Karim et al. (2019a); Livieris et al. (2020); Fawaz et al. (2019b)
Momentum in SGD	0.5, 0.99, 0.9	0.9	Brownlee (2018a)

341

342

2.6.2.1 Number of layers

343

The depth of models is related to the sequence length (Devineau et al., 2018a), as deeper networks need more data to provide better generalization (Fawaz et al., 2019a). In previous studies focused on CNNs, there were usually one, two, or three convolution stages (Zheng et al., 2014). We tried different numbers of CNN, LSTM, and dense layers and the

345

346 best combination was obtained with three CNN layers, two LSTIM layers, and two dense layers. The sequence length
 347 in this study is small (16), and model performance was not improved by simply increasing depth.

348 **2.6.2.2 Number and size of convolution filters**

349 Data with more classes need more filters, and longer time series need longer filters to capture longer patterns and
 350 consequently to produce accurate results (Fawaz et al., 2019a). However, longer filters significantly increase the
 351 number of parameters and potential for overfitting small datasets, while a small filter size risks poor performance. In
 352 this study, two convolutional layers with 1-D filters of size (5, 1) and stride of (1, 1) to capture temporal variation for
 353 each variable separately. Furthermore, one convolutional layer with 2-D filters of size (5, 3) and stride of (1, 1) was
 354 used to capture the correlation between variables via depth-wise convolution of input time series. A big stride might
 355 cause the model to miss valuable data used in predicting and smoothing out the noise in the time series. The layers in
 356 CNNs have a bias for each channel, sharing across all positions in each channel.

357 **2.6.2.4 Adaptive learning rates**

358 The adaptive learning rate decreases the learning rate and consequently weights over each epoch. We tried different
 359 base learning and decay rates for each model and found that the learning rate significantly impacts the model’s
 360 performance. Finally, we chose a base learning rate of 0.1, 0.01, and 0.001 for LSTM, CNN, and CNN-LSTM,
 361 respectively. A decay rate of 0.8 was used for CNN and CNN-LSTM, and a rate of 0.95 for the LSTM model. Table 4
 362 shows the adaptive learning rates for CNN, LSTM, and CNN-LSTM calculated using Equation 2 for each epoch.

363
$$\text{adaptive learning rate} = \text{base learning rate} \times \text{decay}^{\text{epoch}} \quad (2)$$

364 The experiments show that the learning rate is the most critical parameter influencing the model performance. A small
 365 learning rate can cause the loss function to get stuck in local minima, and a large learning rate can result in oscillations
 366 around global minima without reaching it.

367 Our CNN-LSTM model is deeper than the other two models, and deeper models are more prone to a vanishing gradient
 368 problem. To overcome the vanishing gradients, it is generally recommended to use lower learning rates, e.g., lower
 369 than 1 e-4. Interestingly, we found that our CNN-LSTM model works better with lower learning rates than the other
 370 two models.

371

372 **Table 4. The adaptive learning rate for 50 epochs.**

Epochs	Learning rate		
	CNN	CNN-LSTM	LSTM
1	0.008	8.00E-04	0.095
2	0.006	6.40E-04	0.09
3	0.005	5.12E-04	0.086
4	0.004	4.10E-04	0.081
.	.	.	.
.	.	.	.

40	1.30E-06	1.33E-07	0.013
.	.	.	-
50	1.40E-07	1.43E-08	-

373

374

375 2.6.5 Model evaluation

376 The network on the validation set is evaluated after each epoch during training to monitor the training progress. During
 377 validation, all non-deterministic layers are switched to deterministic. For instance, noise layers are disabled, and the
 378 update step of the parameters is not performed.

379 The classification accuracy cannot appropriately represent model performance for unbalanced datasets, as the model
 380 can show a high accuracy by biasing towards the majority class in the dataset (Ordóñez and Roggen, 2016). While we
 381 built a balanced dataset (with the same number of jam and no jam events), randomly selecting test data, shuffling the
 382 inputs, and splitting data into train and validation sets can result in a slightly unbalanced dataset. In our case, the
 383 number of jams and no jams for train and validation and test sets is presented in Table 5. Therefore, the F1 score (Eq.
 384 (3)), which considers each class equally important, is used to measure the accuracy of binary classification. The F1
 385 score, as a weighted average of the precision (Eq. (4)) and recall (Eq. (5)), ranges between 0 and 1, where 0 is the
 386 worse score and 1 is the best. In equations 4 and 5, TP, FP, and FN are true positive, false positive, and false negative,
 387 respectively.

388 **Table 5. The number of jam and no jam events used in the train and validation and test datasets.**

	Train and validation	Test
Jam	456	48
No jam	451	53

389
$$F1 = 2 \times \frac{precision \times recall}{precision + recall} \quad (3)$$

390
$$Precision = \frac{TP}{TP + FP} \quad (4)$$

391
$$Recall = \frac{TP}{TP + FN} \quad (5)$$

392 Although the model accuracy is usually used to examine the performance of deep learning models, the model size
 393 (i.e., number of parameters) provides a second metric, which represents required memory and calculations, to be
 394 compared among models with the same accuracy (Garbin et al., 2020).

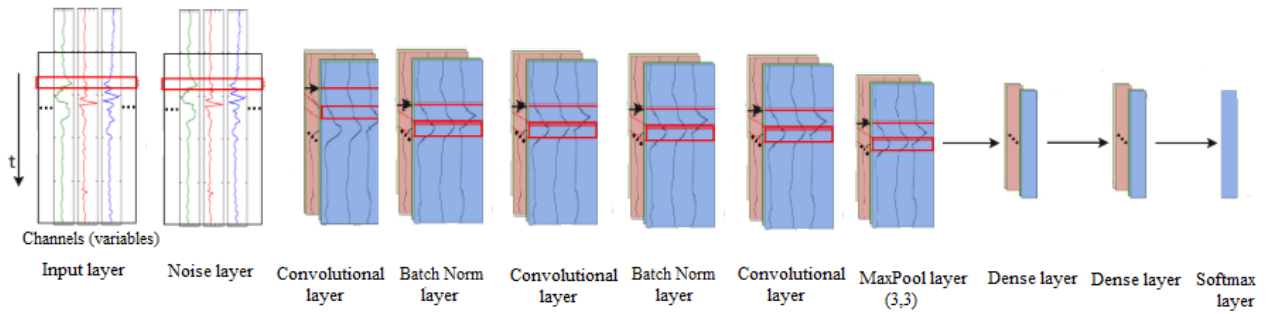
395 After training the model, the well-trained network parameters are saved to a file and are later used to test network
 396 generalizations using the test dataset, composed of data that is not used during training and validation.

397 **2.7 Architecture of models**

398 The final architecture of CNN, LSTM, and CNN-LSTM models are presented in Figures. 9, 10, and 11, respectively.
 399 The layers, their output shapes, and their number of parameters are respectively presented in Tables 6, 7, and 8 for
 400 CNN, LSTM, and CNN-LSTM models.

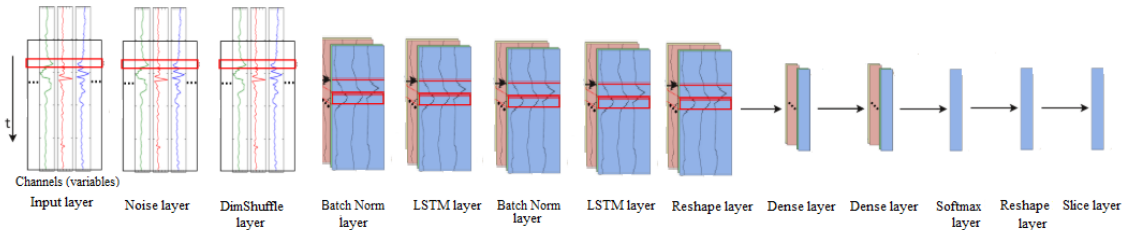
401 Convolutional NN models often include pooling layers to reduce data complexity and dimensionality. However, it is
 402 not always necessary for every convolutional layer to be followed by a pooling layer in the time-series domain
 403 (Ordóñez and Roggen, 2016). For instance, Fawaz et al. (2019a) do not apply any pooling layers in their TSC models.
 404 We tried max-pooling layers after different convolutional layers in CNN and CNN-LSTM networks and found that a
 405 pooling layer following only the last convolutional layer improves the performance of both models. This can be due
 406 to subsampling the time series and using time series with a length of 16 that eliminates the need for decreased
 407 dimensionality.

408



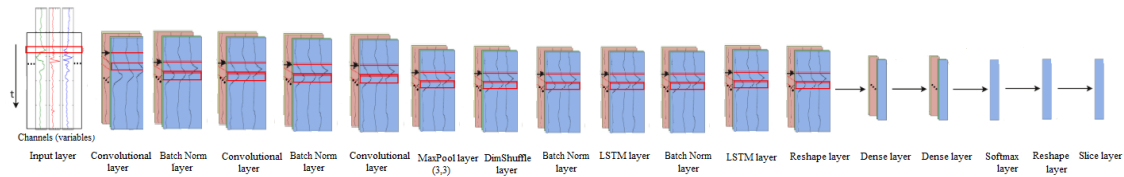
409

410 **Figure 9. Architecture of the CNN model for ice-jam prediction (adapted after Ordóñez and Roggen, 2016).**



411

412 **Figure 10. Architecture of the LSTM model for ice-jam prediction (adapted after Ordóñez and Roggen, 2016).**



413

414 **Figure 11. Architecture of the CNN-LSTM model for ice-jam prediction (adapted after Ordóñez and Roggen, 2016).**

415 **Table 6. Layers, output shapes, and number of parameters for the CNN model.**

Layers	Output shape	Number of parameters
Input	(16, 1, 16, 7)	0

GaussianNoise	(16, 1, 16, 7)	0
Conv2D	(16, 128, 16, 7)	640
BatchNorm	(16, 128, 16, 7)	512
Nonlinearity	(16, 128, 16, 7)	0
Conv2D	(16, 128, 16, 7)	81920
BatchNorm	(16, 128, 16, 7)	512
Nonlinearity	(16, 128, 16, 7)	0
Conv2D	(16, 128, 16, 7)	245888
MaxPool2D	(16, 128, 5, 2)	0
Dense	(16, 32)	40992
Dense	(16, 32)	1056
Softmax	(16, 2)	66

416

417

Table 7. Layers, output shapes, and number of parameters for the LSTM model.

Layers	Output shape	Number of parameters
Input	(16, 1, 16, 7)	0
GaussianNoise	(16, 1, 16, 7)	0
Dimshuffle	(16, 16, 1, 7)	0
BatchNorm	(16, 16, 1, 7)	64
LSTM	(16, 16, 128)	70272
BatchNorm	(16, 16, 128)	64
Nonlinearity	(16, 16, 128)	0
LSTM	(16, 16, 128)	132224
Reshape	(256, 128)	0
Dense	(256, 32)	4128
Dense	(256, 32)	1056
Softmax	(256, 2)	66
Reshape	(16, 16, 2)	0
Slice	(16, 2)	0

418

419

Table 8. Layers, output shapes, and number of parameters for the CNN-LSTM model.

Layers	Output shape	Number of parameters
Input	(16, 1, 16, 7)	0
Conv2D	(16, 128, 16, 7)	640
BatchNorm	(16, 128, 16, 7)	512
Nonlinearity	(16, 128, 16, 7)	0
Conv2D	(16, 128, 16, 7)	81920
BatchNorm	(16, 128, 16, 7)	512
Nonlinearity	(16, 128, 16, 7)	0
Conv2D	(16, 128, 16, 7)	245888
MaxPool2D	(16, 128, 5, 2)	0

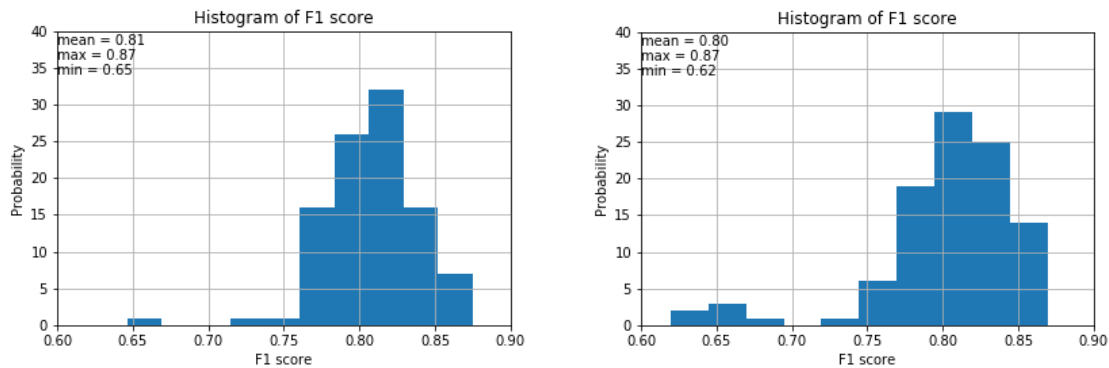
Dimshuffle	(16, 5, 128, 2)	0
BatchNorm	(16, 5, 128, 2)	20
LSTM	(16, 5, 128)	197760
BatchNorm	(16, 5, 128)	20
Nonlinearity	(16, 5, 128)	0
LSTM	(16, 5, 128)	132224
Reshape	(80, 128)	0
Dense	(80, 32)	4128
Dense	(80, 32)	1056
Softmax	(80, 2)	66
Reshape	(16, 5, 2)	0
Slice	(16, 2)	0

420

421 3 Results and Discussion

422 3.1 Weight initialization

423 Among all the methods available for weight initialization in the Lasagne library, GLOT uniform (i.e., Xavier,
 424 Glorot and Bengio, 2010) and He initializations (He et al., 2015) are the most popular initialization techniques to set
 425 the initial random weights in convolutional layers. The results reveal that in our case, these methods yield comparable
 426 F1 scores. However, the histograms of F1 scores reveal that GLOT uniform yields slightly better results (Fig. 12).



427

428 **Figure 12. Histograms of F1 score for CNN using He (left) and GLOT uniform (right) weight initialization with 100**
 429 **random train-validation splits.**

430 3.2 Model evaluation

431 3.2.1 Learning curves and F1 scores

432 Line plots of the loss (i.e., learning curves), which are loss over each epoch, are widely used to assess model
 433 performance in machine learning. Furthermore, line plots clearly indicate common learning problems, such as
 434 underfitting and overfitting. The learning curves for CNN, LSTM, and CNN-LSTM models are presented in Fig. 13.
 435 The LSTM model starts to overfit at epoch 40, so an early stopping is conducted. CNN-LSTM performs better than

436 the other two models, as its training loss is the lowest and is lower than its validation loss. Histograms of F1 scores
437 (Fig. 14 and Table 9) show that CNN-LSTM outperforms the other two models since it results in the highest average
438 and the highest minimum F1-scores for validation (0.82 and 0.75, respectively). Figure 13 shows that the training
439 error of the CNN model is lower than that of the LSTM model, indicating that it trained more efficiently. However, it
440 is not true for the validation error. The validation error is less than the training error in the LSTM model because of
441 the regularization methods used. As LSTM models are often harder to regularize, agreeing with previous studies
442 (e.g., Devineau et al., 2018b).

443 The LSTM network is validated better than the CNN model since its average and minimum F1 scores for validation
444 are better than the CNN model (by 1% and 32%, respectively), and also LSTM yielded no F1 scores below 0.74 (Fig.
445 14 and Table 9).

446 As shown in Figure 13, training loss is higher than validation loss in some of the results. There are some reasons
447 explaining that. Regularization reduces the validation loss at the expense of increasing training loss. Regularization
448 techniques such as the application of noise layers are only used during training, but not during validation resulting in
449 smoother and usually better functions in validation. There is no noise layer in CNN-LSTM model that could result in
450 a lower training error than the validation error. However, other regularization methods such as L2 regularization are
451 used in all the models, including the CNN-LSTM model.

452 Furthermore, batch normalization uses the mean and variance of each batch during the training phase, whereas it uses
453 the mean and variance of the whole training dataset in the validation phase. Additionally, training loss is averaged
454 over each epoch, while validation losses are calculated upon completion of each training epoch. Hence, the training
455 loss includes error calculations with fewer updates.

456 Among the developed machine learning models, SVM shows the best validation performance (Figure 15 and
457 Table 10). However, F1 scores of deep learning models are much higher than those of machine learning models with
458 an average of 6% higher F1 score resulted from CNN-LSTM model compared to the SVM model (Tables 9 and 10).

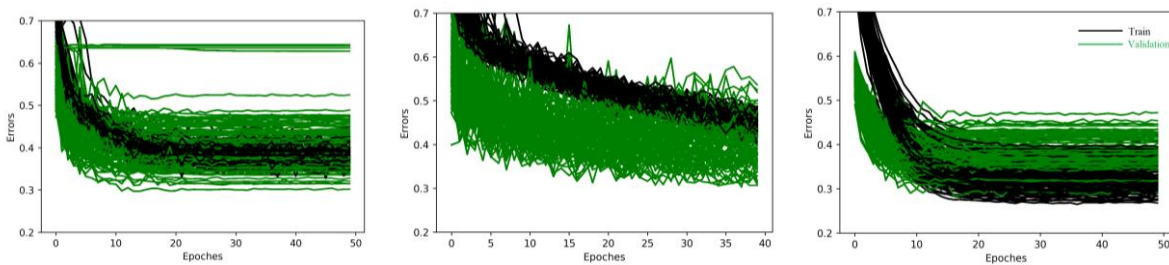


Figure 13. Training and validation errors over epochs for CNN (left), LSTM (middle), and CNN-LSTM (right) models with 100 random train-validation splits.

459

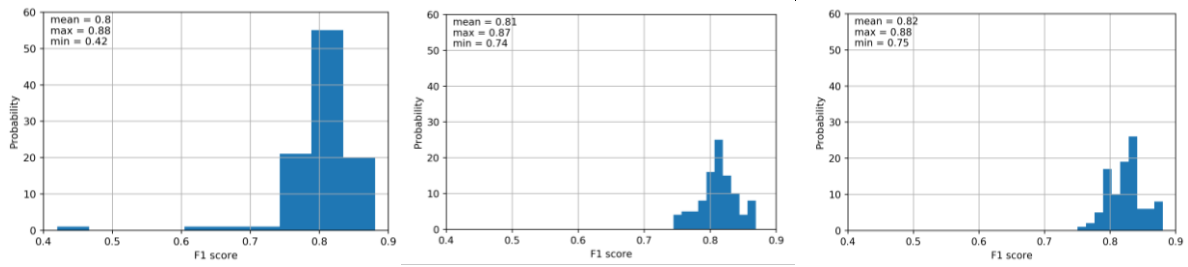


Figure 14. F1 scores of the validation phase for CNN (left), LSTM (middle), and CNN-LSTM (right) models with 100 random train-validation splits.

460

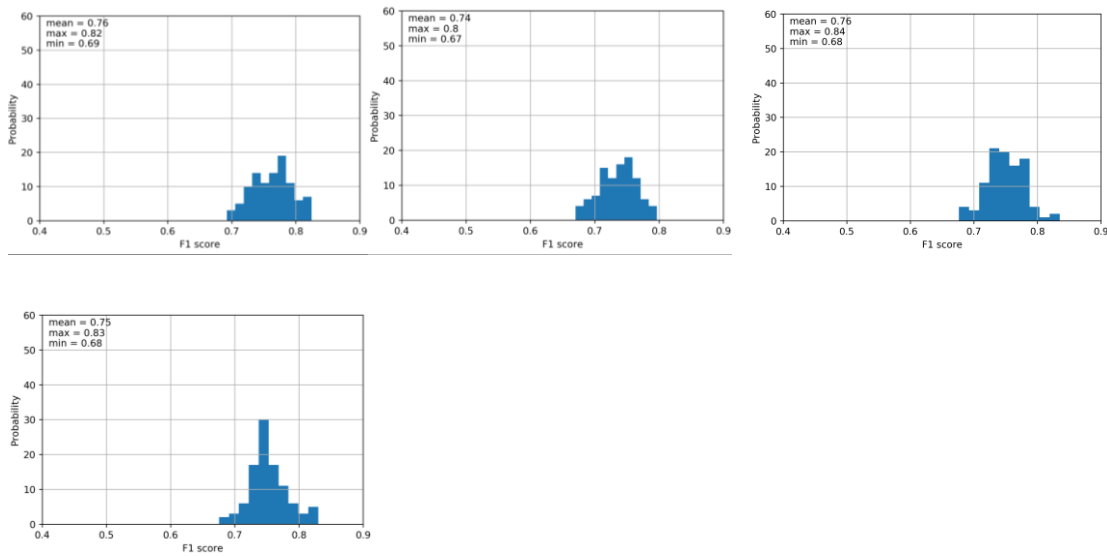


Figure 15. Histograms of F1 scores of the validation phase for SVM (top left), DT (top middle), KNN (top right), and MLP (bottom left) models with 100 random train-validation splits.

461

Table 9. F1 scores of the validation phase for CNN, LSTM, and CNN-LSTM models with 100 random train-validation splits.

Model	F1 score		
	mean	max	min
CNN	0.80	0.88	0.42
LSTM	0.81	0.87	0.74
CNN-LSTM	0.82	0.88	0.75

464

465

466

Table 10. F1 scores of the validation phase for SVM, DT, and KNN and MLP models with 100 random train-validation splits.

Model	F1 score		
	mean	max	min

SVM	0.76	0.82	0.69
DT	0.74	0.80	0.67
KNN	0.75	0.84	0.68
MLP	0.75	0.83	0.68

467 **3.2.2 Number of parameters and run time**

468 The total number of parameters in CNN, LSTM, and CNN-LSTM networks are 371586, 207874, and 664746,
469 respectively. The CNN-LSTM model had the best performance with the highest number of parameters. Even though
470 the number of parameters for the LSTM model is less than that of the CNN model, the LSTM model shows better
471 validation performance. Furthermore, the number of parameters in the CNN-LSTM model is much higher than the
472 other two models, without a large increase in computation time. All three models were trained within 24 hours using
473 100 shuffle splits for training and validation. The models were run on a CPU with four cores, 3.4 GHz clock speed,
474 and 12 GB RAM. On the other hand, a few minutes were enough to train the machine learning models with 100 shuffle
475 splits for training and validation. Although the training time for deep learning models is much higher than that of
476 machine learning models, their superior performance, in this case, justifies their application.

477 **3.3 Order of input variables**

478 It is not clear whether the order of input variables in the input file influences multivariate TSC or not when using 2-D
479 filters and 2-D max-pooling layers. In the benchmark model, the variables were entered from left to right in the
480 following random order: precipitation, minimum temperature, maximum temperature, net radiation, ATDD, AFDD,
481 and snow depth. Another run was conducted by changing the order of the variables for the following random order:
482 snow depth, maximum temperature, precipitation, AFDD, net radiation, minimum temperature, and ATDD. Both
483 models yielded the same mean and minimum F1 scores, whereas the maximum F1 score of the benchmark model is
484 higher (0.88) than that of the comparative run (0.86). Therefore, it can be concluded that the order does not
485 significantly impact the results.

486 **3.4 Testing**

487 To examine the ability of the models to generalize to new unseen data, we randomly set aside 10% of the data from
488 the training and validation phase of all the developed deep learning and machine learning models. A CNN, an LSTM,
489 and a CNN-LSTM model were trained and the well-trained parameters were saved and used to assess the model's
490 ability to generalize. An SVM, a DT, a KNN, and an MLP model were also trained. The trained models were saved
491 and used for testing. The test dataset is nearly balanced with 101 samples with the size of (16, 7), including 48 jam
492 events and 53 no-jam events.

493 The results of the test models show that the CNN-LSTM model had the best F1 score of 0.92 (Table 11). Tables 9 and
494 11 show that although LSTM had a slightly superior validation performance, CNN and LSTM models performed the
495 same in testing.

496 Testing results of machine learning models are presented in Table 11. Among the machine learning models, KNN
497 yields the best results with F1 scores of 78%. By comparing the best deep learning model (CNN-LSTM) with the best

498 machine learning model (KNN), it can be calculated that the deep learning model outperforms the machine learning
499 model by a difference of 14% (F1 score of 92% and 78%, respectively).

500

501 **Table 11. Test F1 scores for the developed deep learning and machine learning models.**

Models	F1 score
CNN-LSTM	0.92
CNN	0.80
LSTM	0.80
KNN	0.78
SVM	0.75
DT	0.71
MLP	0.70

502 **3.5 Model comparison**

503 Classifiers can be combined and used in pattern recognition problems to reduce errors by covering for one another's
504 internal weaknesses (Parvin et al., 2011). Combined classifiers may be less accurate than the most accurate classifier,
505 however, the accuracy of the combined model is always superior to the average accuracy of individual models.
506 Combining two models improved our results compared to convolution-only or LSTM-only networks in both training
507 and testing, supporting the previous studies (e.g., Sainath et al., 2015). It can be because the CNN-LSTM model
508 incorporates both the temporal dependency of each variable by using LSTM networks and the correlation between
509 variables through CNN models. The combined CNN-LSTM model efficiently benefit from automatic feature learning
510 by CNN plus the native support for time series by LSTM.

511 Although LSTM slightly outperformed CNN in the validation phase, these models had comparable performance in
512 the testing phase. The CNN is able to partially include both temporal dependency and the correlation between variables
513 by using 1-D and 2-D filters, respectively. Although the LSTM is unable to incorporate the correlations between
514 variables, it gives promising results with relatively small dataset. Another difference is that LSTM captures longer
515 temporal dynamics, while the CNN only captures temporal dynamics comprised within the length of its filters.

516 Even though our training data in supervised ice-jam prediction is small, the results reveal that deep learning techniques
517 can give accurate results, which agrees with a previous study conducted by Ordóñez and Roggen (2016) in activity
518 recognition. The excellent performance of CNN and CNN-LSTM models may be partially due to the characteristic of
519 CNN that decreases the total number of parameters which does training with limited training data easier (Gao et al.,
520 2016). However, we expect our models to be improved in the future by a larger dataset.

521 Among the developed machine learning models, SVM showed the best performance in validation, whereas KNN
522 worked the best in testing. However, the performance of deep learning models is much better than machine learning
523 models in both validation and testing. The machine learning models do not consider correlations between variables.
524 However, it is not the only reason that deep learning models worked better than machine learning models. As the
525 LSTM also does not consider correlations between variables but worked better than machine learning models. This
526 indicates that there are other aspects of deep learning models that contribute to their high performance level. For

527 instance, deep learning models perform well for the problems with complex-nonlinear dependencies, time
528 dependencies, and multivariate inputs.

529 The developed CNN-LSTM model can be used for future predictions of ice jams in Quebec to provide early warning
530 of possible floods in the area by using historic hydro-meteorological variables and their predictions for some days in
531 advance.

532 **3.6 Discussion on the interpretability of deep learning models**

533 Even though the developed deep learning models performed well in predicting ice jams in Quebec, the interpretability
534 of the results with respect to the physical processes involved in ice jams is still essential. It is because although deep
535 learning models have achieved superior performance in various tasks, these complicated models use a large number
536 of parameters and might sometimes exhibit unexpected behaviours (Samek et al., 2017; Zhang et al., 2021). This is
537 because the real-world environment is still much more complex than any model. Furthermore, the models may learn
538 some spurious correlations in the data and make correct predictions for the “wrong” reason (Samek and Müller, 2019).
539 Hence, interpretability is especially important in some real-world applications like flood and ice-jam predictions where
540 an error could have catastrophic consequences. Nonetheless, interpretability can be used to extract novel domain
541 knowledge and hidden laws of nature in the research fields with limited domain knowledge (Alipanahi et al., 2015)
542 like ice-jam prediction.

543 However, the nested nonlinear structure and the “black box” nature of deep neural networks make the interpretation
544 of their underlying mechanisms and decisions a significant challenge (Montavon et al., 2018; Zhang et al., 2021;
545 Wojtas and Chen, 2020). That is why, deep neural network interpretability still remains a young and emerging field
546 of research. Nevertheless, there are various methods available to facilitate the understanding of decisions made by a
547 deep learning model such as feature importance ranking, sensitivity analysis, layer-wise relevance propagation, and
548 the global surrogate model. However, the interpretability of developed deep learning models for ice-jam prediction is
549 beyond the scope of this study and it will be investigated in our future works.

550

551 **3.7 Model transferability**

552 The transferability of a model between river basins is highly desirable but has not yet been achieved because most
553 river ice-jam models are site specific (Mahabir et al., 2007). The developed models in this study can be used to predict
554 future ice jams some days before the event not only for Quebec but can also be transferred to eastern parts of Ontario
555 and western New Brunswick, since these areas have the similar hydro-meteorological conditions. For other locations,
556 the developed models could be retrained with a small amount of fine-tuning using labelled instances rather than
557 building from scratch. This interesting feature is due to the logic behind the model, which could be transferred to other
558 sites with small modifications. To transfer a model from one river basin to another, historical records of ice jams and
559 equivalent hydro-meteorological variables (e.g., precipitation, temperature, and snow depth) must be available as
560 model inputs for each site.

561

562 **4 Conclusion**

563 The main finding from this project is that the developed deep models successfully predicted ice jams in Quebec, and
564 performed much better than the developed machine learning models. The results show that the CNN-LSTM model is
565 superior to the CNN-only and LSTM-only networks in both validation and testing phases, although the LSTM and
566 CNN performed well.

567 To our best knowledge, this study is the first to apply deep learning models to the ice-jam prediction. The developed
568 models are promising tools for the prediction of ice jams in Quebec and other similar river basins in Canada with re-
569 training and a small amount of fine-tuning.

570 The developed models do not apply to freeze-up jams that occur in early winter and are based on different processes
571 than breakup jams. We studied only breakup ice jams as usually they result in flooding and are more dangerous than
572 freeze-up jams. Furthermore, there is a lack of data availability for freeze-up ice jams in Quebec and only 89 records
573 of freeze-up jams are available which is too small.

574 The main limitation of this study is the availability of ice jam records. Indeed, small datasets may lead deep learning
575 models to overfit the data. Another limitation of the presented work is the lack of interpretability of the results with
576 respect to the physical characteristics of the ice jam. This is a topic of future research and our next step is to explore
577 that.

578 It should also be noted that hydro-meteorological variables are not the only drivers of ice-jam formation.
579 Geomorphological features such as river slope, sinuosity, physical barriers (such as an island or a bridge), channel
580 narrowing, and river confluence also govern the formation of ice jams. In the future, a geospatial model using deep
581 learning will be developed to examine the impacts of these geospatial parameters on ice-jam formation.

582 **Author contribution**

583 Fatemehalsadat Madaeni designed and carried out the experiments under Karem Chokmani and Saeid Homayouni
584 supervision. Fatemehalsadat Madaeni developed the model code and performed the simulations using hydro-
585 meteorological and ice-jam data provided and validated by Rachid Lhissou. Fatemehalsadat Madaeni wrote the bulk
586 of the paper with conceptual edits from Karem Chokmani and Saeid Homayouni. Yves Gauthier and Simon
587 Tolszczuk-Leclerc helped in the refinement of the objectives and the revision of the methodological developments.

588 **Acknowledgment**

589 This study is part of the DAVE project, funded by the Defence Research and Development Canada (DRDC), Canadian
590 Safety and Security Program (CSSP), with partners from Natural Resources Canada (NRCan), and Environment and
591 Climate Change Canada.

592 **References**

593 Alipanahi, B., DeLong, A., Weirauch, M. T., & Frey, B. J. (2015). Predicting the sequence specificities of DNA-and
594 RNA-binding proteins by deep learning. *Nature biotechnology*, 33(8), 831-838.

595 Althoff, D., Rodrigues, L. N., & Bazame, H. C. (2021). Uncertainty quantification for hydrological models based on
596 neural networks: the dropout ensemble. *Stochastic Environmental Research and Risk Assessment*, 35(5), 1051-1067.

597 Analytics, C. (2016). *Anaconda Software Distribution: Version 2-2.4. 0*.

598 Apaydin, H., Feizi, H., Sattari, M. T., Colak, M. S., Shamshirband, S., & Chau, K. W. (2020). Comparative analysis
599 of recurrent neural network architectures for reservoir inflow forecasting. *Water*, 12(5), 1500.

600 Barnes-Svarney, P. L., & Montz, B. E. (1985). An ice jam prediction model as a tool in floodplain management. *Water
601 Resources Research*, 21(2), 256-260.

602 Barzegar, R., Aalami, M. T., & Adamowski, J. (2020). Short-term water quality variable prediction using a hybrid
603 CNN-LSTM deep learning model. *Stochastic Environmental Research and Risk Assessment*, 1-19.

604 Barzegar, R., Aalami, M. T., & Adamowski, J. (2021). Coupling a hybrid CNN-LSTM deep learning model with a
605 Boundary Corrected Maximal Overlap Discrete Wavelet Transform for multiscale Lake water level
606 forecasting. *Journal of Hydrology*, 598, 126196.

607 Beltaos, S. (1993). Numerical computation of river ice jams. *Canadian Journal of Civil Engineering*, 20(1), 88-99.

608 Bengio, Y. (2012). Practical recommendations for gradient-based training of deep architectures. *Neural networks:
609 Tricks of the trade*, 437-478.

610 Bergstra, J., & Bengio, Y. (2012). Random search for hyper-parameter optimization. *Journal of machine learning
611 research*, 13(2).

612 Bergstra, J., Breuleux, O., Bastien, F., Lamblin, P., Pascanu, R., Desjardins, G., ... & Bengio, Y. (2010). Theano: A
613 CPU and GPU math compiler in Python. *Proc. 9th Python in Science Conf*, 3-10.

614 Brownlee, J. (2017). *Long short-term memory networks with python: develop sequence prediction models with deep
615 learning. Machine Learning Mastery*.

616 Brownlee, J. (2018a). *Better deep learning: train faster, reduce overfitting, and make better predictions. Machine
617 Learning Mastery*.

618 Brownlee, J. (2018b). *Deep learning for time series forecasting: predict the future with MLPs, CNNs and LSTMs in
619 Python. Machine Learning Mastery*.

620 Brunel, A., Pasquet, J., PASQUET, J., Rodriguez, N., Comby, F., Fouchez, D., & Chaumont, M. (2019). A CNN
621 adapted to time series for the classification of Supernovae. *Electronic Imaging*, 2019(14), 90-1.

622 Brunello, A., Marzano, E., Montanari, A., & Sciavicco, G. (2019). J48SS: A novel decision tree approach for the
623 handling of sequential and time series data. *Computers*, 8(1), 21.

624 Brunner, G. W. (2002). Hec-ras (river analysis system). North American Water and Environment Congress &
625 Destructive Water, 3782-3787.

626 Carson, R. W., Beltaos, S., Healy, D., & Groeneveld, J. (2003). Tests of river ice jam models—phase 2. Proceedings
627 of the 12th Workshop on the Hydraulics of Ice Covered Rivers, Edmonton, Alta, 19-20.

628 Carson, R., Beltaos, S., Groeneveld, J., Healy, D., She, Y., Malenchak, J., ... & Shen, H. T. (2011). Comparative
629 testing of numerical models of river ice jams. Canadian Journal of Civil Engineering, 38(6), 669-678.

630 Chen, R., Wang, X., Zhang, W., Zhu, X., Li, A., & Yang, C. (2019). A hybrid CNN-LSTM model for typhoon
631 formation forecasting. GeoInformatica, 23(3), 375-396.

632 Cui, Z., Chen, W., & Chen, Y. (2016). Multi-scale convolutional neural networks for time series classification. arXiv
633 preprint arXiv:1603.06995.

634 De Coste, M., Li, Z., Pupek, D., & Sun, W. (2021). A hybrid ensemble modelling framework for the prediction of
635 breakup ice jams on Northern Canadian Rivers. Cold Regions Science and Technology, 189, 103302.

636 del Campo, F. A., Neri, M. C. G., Villegas, O. O. V., Sánchez, V. G. C., Domínguez, H. D. J. O., & Jiménez, V. G.
637 (2021). Auto-adaptive multilayer perceptron for univariate time series classification. Expert Systems with
638 Applications, 181, 115147.

639 Devineau, G., Moutarde, F., Xi, W., & Yang, J. (2018a). Deep learning for hand gesture recognition on skeletal data.
640 13th IEEE International Conference on Automatic Face & Gesture Recognition, 106-113.

641 Devineau, G., Xi, W., Moutarde, F., & Yang, J. (2018b). Convolutional neural networks for multivariate time series
642 classification using both inter-and intra-channel parallel convolutions. Reconnaissance des Formes, Image,
643 Apprentissage et Perception.

644 Dieleman, S., Schlüter, J., Raffel, C., Olson, E., Sønderby, S.K., Nouri, D., ... & Degraeve, J. (2015). Lasagne: First
645 release. (Version v0.1). Zenodo. Retrieved from <http://doi.org/10.5281/zenodo.27878>.

646 Données Québec: Historique (publique) d'embâcles répertoriés au MSP - Données Québec. Retrieved from
647 <https://www.donneesquebec.ca/recherche/dataset/historique-publique-d-embacles-repertories-au-msp>. (Last access:
648 15 June 2021).

649 Fawaz, H. I., Forestier, G., Weber, J., Idoumghar, L., & Muller, P. A. (2019a). Deep neural network ensembles for
650 time series classification. International Joint Conference on Neural Networks, 1-6.

651 Fawaz, H. I., Forestier, G., Weber, J., Idoumghar, L., & Muller, P. A. (2019b). Deep learning for time series
652 classification: a review. Data Mining and Knowledge Discovery, 33(4), 917-963.

653 Fischer, T., & Krauss, C. (2018). Deep learning with long short-term memory networks for financial market
654 predictions. European Journal of Operational Research, 270(2), 654-669.

655 Gao, Y., Hendricks, L. A., Kuchenbecker, K. J., & Darrell, T. (2016). Deep learning for tactile understanding from
656 visual and haptic data. *International Conference on Robotics and Automation*, 536-543.

657 Garbin, C., Zhu, X., & Marques, O. (2020). Dropout vs. batch normalization: an empirical study of their impact to
658 deep learning. *Multimedia Tools and Applications*, 1-39.

659 Glorot, X., & Bengio, Y. (2010). Understanding the difficulty of training deep feedforward neural networks.
660 *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, 249-256.

661 Graf, R., Kolerski, T., & Zhu, S. (2022). Predicting Ice Phenomena in a River Using the Artificial Neural Network
662 and Extreme Gradient Boosting. *Resources*, 11(2), 12.

663 Gu, J., Wang, Z., Kuen, J., Ma, L., Shahroudy, A., Shuai, B., ... & Chen, T. (2018). Recent advances in convolutional
664 neural networks. *Pattern Recognition*, 77, 354-377.

665 Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., ... & Oliphant, T. E.
666 (2020). Array programming with NumPy. *Nature*, 585(7825), 357-362.

667 Hatami, N., Gavet, Y., & Debayle, J. (2018). Classification of time-series images using deep convolutional neural
668 networks. *Tenth International Conference on Machine Vision*.

669 He, K., Zhang, X., Ren, S., & Sun, J. (2015). Delving deep into rectifiers: Surpassing human-level performance on
670 imagenet classification. *Proceedings of the IEEE International Conference on Computer Vision*, 1026-1034.

671 Hunter, J. D. (2007). Matplotlib: A 2D graphics environment. *IEEE Annals of the History of Computing*, 9(03), 90-
672 95.

673 Ioffe, S., & Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal
674 covariate shift. *International Conference on Machine Learning*, 448-456.

675 Jović, A., Brkić, K., & Bogunović, N. (2012). Decision tree ensembles in biomedical time-series classification. *Joint*
676 *DAGM (German Association for Pattern Recognition) and OAGM Symposium*, 408-417.

677 Jozefowicz, R., Zaremba, W., & Sutskever, I. (2015). An empirical exploration of recurrent network architectures.
678 *International Conference on Machine Learning*, 2342-2350.

679 Karim, F., Majumdar, S., & Darabi, H. (2019a). Insights into LSTM fully convolutional networks for time series
680 classification. *IEEE Access*, 7, 67718-67725.

681 Karim, F., Majumdar, S., Darabi, H., & Chen, S. (2017). LSTM fully convolutional networks for time series
682 classification. *IEEE access*, 6, 1662-1669.

683 Karim, F., Majumdar, S., Darabi, H., & Harford, S. (2019b). Multivariate lstm-fcns for time series
684 classification. *Neural Networks*, 116, 237-245.

685 Kashiparekh, K., Narwariya, J., Malhotra, P., Vig, L., & Shroff, G. (2019). ConvTimeNet: A pre-trained deep
686 convolutional neural network for time series classification. *International Joint Conference on Neural Networks*, 1-8.

687 Kratzert, F., Klotz, D., Brenner, C., Schulz, K., & Herrnegger, M. (2018). Rainfall–runoff modelling using long short-
688 term memory (LSTM) networks. *Hydrology and Earth System Sciences*, 22(11), 6005-6022.

689 Li, D., Djulovic, A., & Xu, J. F. (2013). A Study of kNN using ICU multivariate time series data. *Proc. Int. Conf.*
690 *Data Mining*, eds. R. Stahlbock and GM Weiss, 211-217.

691 Li, X., Zhang, Y., Zhang, J., Chen, S., Marsic, I., Farneth, R. A., & Burd, R. S. (2017). Concurrent activity recognition
692 with multimodal CNN-LSTM structure. *arXiv preprint arXiv:1702.01638*.

693 Lin, J., Williamson, S., Borne, K., & DeBarr, D. (2012). Pattern recognition in time series. *Advances in Machine*
694 *Learning and Data Mining for Astronomy*, 1, 617-645.

695 Lindenschmidt, K. E. (2017). RIVICE—a non-proprietary, open-source, one-dimensional river-ice
696 model. *Water*, 9(5), 314.

697 Lipton, Z. C., Berkowitz, J., & Elkan, C. (2015). A critical review of recurrent neural networks for sequence
698 learning. *arXiv preprint arXiv:1506.00019*.

699 Livieris, I. E., Pintelas, E., & Pintelas, P. (2020). A CNN–LSTM model for gold price time-series forecasting. *Neural*
700 *computing and applications*, 32(23), 17351-17360.

701 Lu, N., Wu, Y., Feng, L., & Song, J. (2018). Deep learning for fall detection: Three-dimensional CNN combined with
702 LSTM on video kinematic data. *IEEE journal of biomedical and health informatics*, 23(1), 314-323.

703 Luan, Y., & Lin, S. (2019). Research on text classification based on CNN and LSTM. *International Conference on*
704 *Artificial Intelligence and Computer Applications*, 352-355.

705 Madaeni, F., Lhissou, R., Chokmani, K., Raymond, S., & Gauthier, Y. (2020). Ice jam formation, breakup and
706 prediction methods based on hydroclimatic data using artificial intelligence: A review. *Cold Regions Science and*
707 *Technology*, 103032.

708 Maggiori, E., Tarabalka, Y., Charpiat, G., & Alliez, P. (2017). High-resolution aerial image labeling with
709 convolutional neural networks. *IEEE Transactions on Geoscience and Remote Sensing*, 55(12), 7092-7103.

710 Mahabir, C., Hicks, F. E., & Fayek, A. R. (2007). Transferability of a neuro-fuzzy river ice jam flood forecasting
711 model. *Cold Regions Science and Technology*, 48(3), 188-201.

712 Mahabir, C., Hicks, F., & Fayek, A. R. (2006). Neuro-fuzzy river ice breakup forecasting system. *Cold regions science*
713 *and technology*, 46(2), 100-112.

714 Mahfouf, J. F., Brasnett, B., & Gagnon, S. (2007). A Canadian precipitation analysis (CaPA) project: Description and
715 preliminary results. *Atmosphere-ocean*, 45(1), 1-17.

716 Massie, D.D., White, K.D., Daly, S.F. (2002). Application of neural networks to predict ice jam occurrence. *Cold*
717 *Regions Science and Technology*, 35 (2), 115–122.

718 Masters, D., & Luschi, C. (2018). Revisiting small batch training for deep neural networks. *arXiv preprint*
719 *arXiv:1804.07612*.

720 Mesinger, F., DiMego, G., Kalnay, E., Mitchell, K., Shafran, P. C., Ebisuzaki, W., ... & Shi, W. (2006). North
721 American regional reanalysis. *Bulletin of the American Meteorological Society*, 87(3), 343-360.

722 Montavon, G., Samek, W., & Müller, K. R. (2018). Methods for interpreting and understanding deep neural
723 networks. *Digital Signal Processing*, 73, 1-15.

724 Mutegeki, R., & Han, D. S. (2020). A CNN-LSTM approach to human activity recognition. *International Conference*
725 *on Artificial Intelligence in Information and Communication*, 362-366.

726 Nanopoulos, A., Alcock, R., & Manolopoulos, Y. (2001). Feature-based classification of time-series
727 data. *International Journal of Computer Research*, 10(3), 49-61.

728 National Hydro Network - NHN - GeoBase Series - Natural Resources Canada. Retrieved from
729 <https://open.canada.ca/data/en/dataset/a4b190fe-e090-4e6d-881e-b87956c07977>.

730 National Hydrographic Network - Natural Resources Canada. Retrieved from [https://www.nrcan.gc.ca/science-and-](https://www.nrcan.gc.ca/science-and-data/science-and-research/earth-sciences/geography/topographic-information/geobase-surface-water-program-geeau/national-hydrographic-network/21361)
731 [data/science-and-research/earth-sciences/geography/topographic-information/geobase-surface-water-program-](https://www.nrcan.gc.ca/science-and-research/earth-sciences/geography/topographic-information/geobase-surface-water-program-geeau/national-hydrographic-network/21361)
732 [geeau/national-hydrographic-network/21361](https://www.nrcan.gc.ca/science-and-research/earth-sciences/geography/topographic-information/geobase-surface-water-program-geeau/national-hydrographic-network/21361).

733 Nosratabadi, S., Mosavi, A., Duan, P., Ghamisi, P., Filip, F., Band, S. S., ... & Gandomi, A. H. (2020). Data science
734 in economics: comprehensive review of advanced machine learning and deep learning methods. *Mathematics*, 8(10),
735 1799.

736 Oh, S. L., Ng, E. Y., San Tan, R., & Acharya, U. R. (2018). Automated diagnosis of arrhythmia using combination of
737 CNN and LSTM techniques with variable length heart beats. *Computers in biology and medicine*, 102, 278-287.

738 Olah, C. (2015). Understanding LSTM Networks. Retrieved from [https://colah.github.io/posts/2015-08-](https://colah.github.io/posts/2015-08-Understanding-LSTMs/)
739 [Understanding-LSTMs/](https://colah.github.io/posts/2015-08-Understanding-LSTMs/).

740 Ombabi, A. H., Ouarda, W., & Alimi, A. M. (2020). Deep learning CNN–LSTM framework for Arabic sentiment
741 analysis using textual information shared in social networks. *Social Network Analysis and Mining*, 10(1), 1-13.

742 Ordóñez, F. J., & Roggen, D. (2016). Deep convolutional and lstm recurrent neural networks for multimodal wearable
743 activity recognition. *Sensors*, 16(1), 115.

744 Parvin, H., Minaei, B., Beigi, A., & Helmi, H. (2011). Classification ensemble by genetic algorithms. *International*
745 *Conference on Adaptive and Natural Computing Algorithms*, 391-399.

746 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Duchesnay, E. (2011). Scikit-
747 learn: Machine learning in Python. *The Journal of Machine Learning Research*, 12.

748 Prowse, T. D., & Bonsal, B. R. (2004). Historical trends in river-ice break-up: a review. *Hydrology Research*, 35 (4-
749 5), 281-293.

750 Prowse, T. D., Bonsal, B. R., Duguay, C. R., & Lacroix, M. P. (2007). River-ice break-up/freeze-up: a review of
751 climatic drivers, historical trends and future predictions. *Annals of Glaciology*, 46, 443-451.

752 Raybaut, P. (2009). Spyder-documentation. Retrieved from pythonhosted. org.

753 Reback, J., McKinney, W., Den Van Bossche, J., Augspurger, T., Cloud, P., Klein, A., ... & Seabold, S. (2020).
754 pandas-dev/pandas: Pandas 1.0. 3. Zenodo.

755 Rodríguez, J. J., & Alonso, C. J. (2004). Support vector machines of interval-based features for time series
756 classification. *International Conference on Innovative Techniques and Applications of Artificial Intelligence*, 244-
757 257.

758 Sainath, T. N., Vinyals, O., Senior, A., & Sak, H. (2015). Convolutional, long short-term memory, fully connected
759 deep neural networks. *International Conference on Acoustics, Speech and Signal Processing*, 4580-4584.

760 Samek, W., & Müller, K. R. (2019). Towards explainable artificial intelligence. *Explainable AI: Interpreting,*
761 *Explaining and Visualizing Deep Learning*, 5-22.

762 Samek, W., Wiegand, T., & Müller, K. R. (2017). Explainable artificial intelligence: Understanding, visualizing and
763 interpreting deep learning models. *arXiv preprint arXiv:1708.08296*.

764 Sarafanov, M., Borisova, Y., Maslyayev, M., Revin, I., Maximov, G., & Nikitin, N. O. (2021). Short-Term River Flood
765 Forecasting Using Composite Models and Automated Machine Learning: The Case Study of Lena
766 River. *Water*, 13(24), 3482.

767 Semenova, N., Sazonov, A., Krylenko, I., & Frolova, N. (2020). Use of classification algorithms for the ice jams
768 forecasting problem. *E3S Web of Conferences*.

769 She, X., & Zhang, D. (2018). Text classification based on hybrid CNN-LSTM hybrid model. *11th International*
770 *Symposium on Computational Intelligence and Design*, 185-189.

771 Shouyu, C., & Honglan, J. (2005). Fuzzy Optimization Neural Network Approach for Ice Forecast in the Inner
772 Mongolia Reach of the Yellow River/Approche d'Optimisation Floue de Réseau de Neurones pour la Prévision de la
773 Glace Dans le Tronçon de Mongolie Intérieure du Fleuve Jaune. *Hydrological sciences journal*, 50(2).

774 Sosa, P. M. (2017). Twitter sentiment analysis using combined LSTM-CNN models. *Eprint Arxiv*, 1-9.

775 Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014). Dropout: a simple way to
776 prevent neural networks from overfitting. *The journal of machine learning research*, 15(1), 1929-1958.

777 The Atlas of Canada - Toporama - Natural Resources Canada. Retrieved from
778 <https://atlas.gc.ca/toporama/en/index.html>.

779 Thornton, M.M., Shrestha, R., Wei, Y., Thornton, P.E., Kao, S. & Wilson, B.E. (2020). Daymet: Daily Surface
780 Weather Data on a 1-km Grid for North America, Version 4. ORNL DAAC, Oak Ridge, Tennessee, USA.

781 Torres, J. F., Hadjout, D., Sebaa, A., Martínez-Álvarez, F., & Troncoso, A. (2021). Deep Learning for Time Series
782 Forecasting: A Survey. *Big Data*, 9(1), 3-21.

783 Turcotte, B., & Morse, B. (2015). River ice breakup forecast and annual risk distribution in a climate change
784 perspective. 18th Workshop on the Hydraulics of Ice Covered Rivers, CGU HS Committee on River Ice Processes
785 and the Environment, Quebec.

786 Umer, M., Imtiaz, Z., Ullah, S., Mehmood, A., Choi, G. S., & On, B. W. (2020). Fake news stance detection using
787 deep learning architecture (cnn-lstm). *IEEE Access*, 8, 156695-156706.

788 Wang, J., Yu, L. C., Lai, K. R., & Zhang, X. (2016). Dimensional sentiment analysis using a regional CNN-LSTM
789 model. *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*, 225-230.

790 Wang, J., Yu, L. C., Lai, K. R., & Zhang, X. (2019). Tree-structured regional CNN-LSTM model for dimensional
791 sentiment analysis. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 28, 581-591.

792 White, K. D. (2003). Review of prediction methods for breakup ice jams. *Canadian Journal of Civil Engineering*,
793 30(1), 89-100.

794 White, K. D., & Daly, S. F. (2002). Predicting ice jams with discriminant function analysis. 21st International
795 Conference on Offshore Mechanics and Arctic Engineering, 683-690.

796 Wojtas, M., & Chen, K. (2020). Feature importance ranking for deep learning. *arXiv preprint arXiv:2010.08973*.

797 Wu, J., Yao, L., & Liu, B. (2018). An overview on feature-based classification algorithms for multivariate time series.
798 3rd International Conference on Cloud Computing and Big Data Analysis, 32-38.

799 Wu, Z., Wang, X., Jiang, Y. G., Ye, H., & Xue, X. (2015). Modeling spatial-temporal clues in a hybrid deep learning
800 framework for video classification. 23rd ACM International Conference on Multimedia, 461-470.

801 Wunsch, A., Liesch, T., & Broda, S. (2020). Groundwater Level Forecasting with Artificial Neural Networks: A
802 Comparison of LSTM, CNN and NARX. *Hydrology and Earth System Sciences Discussions*, 2020, 1-23.

803 Xing, Z., Pei, J., & Keogh, E. (2010). A brief survey on sequence classification. *ACM Sigkdd Explorations*
804 *Newsletter*, 12(1), 40-48.

805 Xingjian, S. H. I., Chen, Z., Wang, H., Yeung, D. Y., Wong, W. K., & Woo, W. C. (2015). Convolutional LSTM
806 network: A machine learning approach for precipitation nowcasting. *Advances in neural information processing*
807 *systems*, 802-810.

808 Yan, J., Mu, L., Wang, L., Ranjan, R., & Zomaya, A. Y. (2020). Temporal convolutional networks for the advance
809 prediction of ENSO. *Scientific reports*, 10(1), 1-15.

810 Yang, J., Nguyen, M. N., San, P. P., Li, X. L., & Krishnaswamy, S. (2015). Deep convolutional neural networks on
811 multichannel time series for human activity recognition. *Twenty-fourth International Joint Conference on Artificial
812 Intelligence*.

813 Yi, S., Ju, J., Yoon, M. K., & Choi, J. (2017). Grouped convolutional neural networks for multivariate time
814 series. *arXiv preprint arXiv:1703.09938*.

815 Zhang, D., Lin, J., Peng, Q., Wang, D., Yang, T., Sorooshian, S., ... & Zhuang, J. (2018). Modeling and simulating of
816 reservoir operation using the artificial neural network, support vector regression, deep learning algorithm. *Journal of
817 Hydrology*, 565, 720-736.

818 Zhang, Y., Tiño, P., Leonardis, A., & Tang, K. (2021). A survey on neural network interpretability. *IEEE Transactions
819 on Emerging Topics in Computational Intelligence*.

820 Zhao, L., Hicks, F. E., & Fayek, A. R. (2012). Applicability of multilayer feed-forward neural networks to model the
821 onset of river breakup. *Cold Regions Science and Technology*, 70, 32-42.

822 Zheng, Y., Liu, Q., Chen, E., Ge, Y., & Zhao, J. L. (2014). Time series classification using multi-channels deep
823 convolutional neural networks. *International Conference on Web-Age Information Management*, 298-310.

824 Zheng, Y., Liu, Q., Chen, E., Ge, Y., & Zhao, J. L. (2016). Exploiting multi-channels deep convolutional neural
825 networks for multivariate time series classification. *Frontiers of Computer Science*, 10(1), 96-112.

826

827