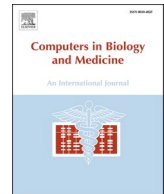




Since January 2020 Elsevier has created a COVID-19 resource centre with free information in English and Mandarin on the novel coronavirus COVID-19. The COVID-19 resource centre is hosted on Elsevier Connect, the company's public news and information website.

Elsevier hereby grants permission to make all its COVID-19-related research that is available on the COVID-19 resource centre - including this research content - immediately available in PubMed Central and other publicly funded repositories, such as the WHO COVID database with rights for unrestricted research re-use and analyses in any form or by any means with acknowledgement of the original source. These permissions are granted for free by Elsevier for as long as the COVID-19 resource centre remains active.



# A novel approach for COVID-19 Infection forecasting based on multi-source deep transfer learning

Sonakshi Garg<sup>1</sup>, Sandeep Kumar, Pranab K. Muhuri<sup>\*</sup>

Department of Computer Science, South Asian University, Akbar Bhawan, Chanakyapuri, New Delhi, 110021, India

## ARTICLE INFO

### Keywords:

COVID-19  
Province-specific data  
Multi-source domain dataset  
LSTM  
Deep transfer learning  
Contagious disease  
Coronavirus infection forecasting

## ABSTRACT

COVID-19 is a contagious disease; so, predicting its future infections in a provincial region requires the consideration of the related data (i.e., rates of infection, mortality and recovery, etc.) over a period of time. Clearly, the COVID-19 data of a particular provincial region can be easily modelled as a time-series. However, predicting the future COVID-19 infections in a particular region is quite challenging when the availability of COVID-19 dataset of the province is of little quantity. Accordingly, ML models when deployed for such tasks usually results in low infection prediction accuracy. To overcome such issues of low variance and high bias in a model due to data scarcity, multi-source transfer learning (MSTL) along with deep learning may be quite useful and effective. Therefore, this paper proposes a novel technique based on multi-source deep transfer learning (MSDTL) to efficiently forecast the future COVID-19 infections in the provinces with insufficient COVID-19 data. The proposed approach is a novel contribution as it considers the fact that future COVID-19 transmission in a region also depends on its population density and economic conditions (GDP) for accurate forecasting of the infections to tackle the pandemic efficiently. The importance of this feature selection is experimentally proved in this paper. Our proposed approach employs the well-known recurrent neural network architecture, the Long-short term memory (LSTM), a popular deep-learning model for history-dependent tasks. A comparative analysis has been performed with existing state-of-art algorithms to portray the efficiency of LSTM. Thus, formation of MSDTL approach enhances the predictive precision capability of the LSTM. We evaluate the proposed methodology over the COVID-19 dataset from sixty-two provinces belonging to different nations. We then empirically evaluate the performance of the proposed approach using two different evaluation metrics, viz. The mean absolute percentage error and the coefficient of determination. We show that our proposed MSDTL based approach is better in terms of the accuracy of the future infection prediction, and produces improvements up to 96% over its without-TL counterpart.

## 1. Introduction

COVID-19 is a contagious disease effected from a novel form of coronavirus related with the severe acute respiratory syndrome (SARS) [1]. Later in the year 2019 in Wuhan, a city in Hubei, China, the foremost international outbreak of COVID-19, a coronavirus pathogen, was reported. In March 2020, the World Health Organization (WHO) declared the Coronavirus epidemic as a global pandemic. On March 12, 2020 the outbreak of COVID-19 has spread to 118 countries around the globe with a total of 125,260 confirmed cases and 4613 deaths as reported by WHO [2]. However, the infection count of COVID-19 has

reached 200 million with deaths of 4 million worldwide until August 12, 2021. The human-to-human transmission is through inhalation of micro-droplets exhaled by infected persons. Hand sanitization, maintaining social distance, and similar other regulatory health guidelines are the only preventions. Therefore, finding suitable approaches for the accurate prediction of future transmissions of COVID-19 has become a fertile research ground. This is because, any such authentic predictions of future COVID-19 infections can help the governments, healthcare administrators, and related policymakers to prepare for better handling of the COVID-19-related challenges [3]. This may help governments to visualize the overwhelming scenario to look around the future

<sup>\*</sup> Corresponding author.

E-mail addresses: [sonakshi.garg12@gmail.com](mailto:sonakshi.garg12@gmail.com) (S. Garg), [2431sandeep@gmail.com](mailto:2431sandeep@gmail.com) (S. Kumar), [pranabmuhuri@cs.sau.ac.in](mailto:pranabmuhuri@cs.sau.ac.in), [pranabmuhuri@gmail.com](mailto:pranabmuhuri@gmail.com) (P.K. Muhuri).

<sup>1</sup> Present address: Department of Computing Science, Umeå University 90187, Umea, Sweden.

happenings which will result in efficient planning for vaccination drives. Thus, it may be beneficial for administrators to prepare and deployment of resources proficiently. Hence, future awareness among the community needs to be generated.

To forecast future transmission of the infection, data for the current scenario is needed along with past observations for precise predictions. Hence, there is a need to consider time-series data of a province for suitable prediction modelling to obtain future transmission of expected cases in its upcoming days. It is well-known that bigger is the dataset, better is the performance of a prediction model [4]. However, there are many provinces for which available datasets are insufficient for training to generate authentic forecasts. To overcome this limitation of data scarcity, the concept of Transfer Learning (TL) may be quite useful. The name transfer learning implies improvement of the traditional machine learning problem by learning a task and gaining knowledge in a certain domain, which is then transferred to improve learning in some related tasks. This method of transferring knowledge signifies improvement in developing machine learning as effective as human learning [5]. Each domain in TL architecture is composed of a feature space  $F$  along with  $P(X)$ , marginal probability distribution. Mathematically, it can be expressed as follows:

$$\{F, P(X)\} \text{ where } X = \{x_1, \dots, x_n\} \in F \quad (1)$$

Generally, if some domains are dissimilar, then they may have diverse feature spaces ( $F$ ) or diverse marginal probability distributions  $P(X)$ . For a definite domain, a task (expressed by  $T = \{Y, f(\cdot)\}$ ) comprises two parameters: a label space ( $Y$ ) along with an objective prediction function  $f(\cdot)$  that can be learned from training dataset comprising of  $\{x_i, y_i\}$  pairs i.e.,  $x_i \in X$  and  $y_i \in Y$ . Therefore, the function  $f(\cdot)$  is useful for prediction of the equivalent label  $f(z)$ , for a new instance.  $z$ .

For TL architecture, we have source domain  $D_s$ , source learning task  $T_s$ , along with target domain  $D_t$  and target learning task  $T_t$ . A function  $f_t(\cdot)$  is calculated to progress the learning of target prediction in  $D_t$  using information from source domain and task, with  $D_s \neq D_t$ , or  $T_s \neq T_t$ . Thus, the source domain data is useful for enhancing the performance by transferral of knowledge during training for target task. However, the case where source and target domains are identical (i.e.,  $= D_t$ ) and their learned tasks are similar (i.e.,  $T_s = T_t$ ), then it becomes a classical machine learning problem.

Rosenstein et al. anticipated that when dissimilarity is observed between two tasks, the performance of the approach destroys even though knowledge transfer has been utilized with the help of brute force techniques [51]. If two tasks are highly dissimilar, then transfer learning results in inefficient performance. Thus, the proposed approach utilizes multi-source domain dataset from wide range of countries possessing similar data-distributions in terms of demographic and economic conditions. Apart from the COVID-19 data, population, geographical area and GDP of numerous provinces from such countries are considered in the proposed approach.

Several studies have been proposed to show the correlation between community and infectious spread of a disease (like coronavirus) in a locality [6,7]. It has been observed that the greater the population, more is the chance of infectious spread of a contagious disease. Similar relation is also observed between geographical area of a densely populated locality and infectious transmission [8]. GDP is also a deciding factor in prediction of the disease transmission [9]. GDP is composed of several sectors such as agriculture, industries, education, healthcare, available resources etc. A country with deprived GDP is an indication of poor facilities. Hence, we consider a collection of COVID data along with population, area and GDP belonging to a set of related provinces which follow different data distributions.

Therefore, this paper proposes a novel LSTM-based approach in the architecture of multi-source deep-TL to forecast COVID-19 transmission. LSTM model is a well-known deep-learning model used in time-series problems that require an abundant amount of data for training. When

LSTM is trained initially, it is initialized with random weights, which may cause low variance and high bias i.e., it may underfit the data, and overfitting can be contained by splitting the data into training and validation sets. Similarly, in the COVID-19 scenario, the availability of data of a specific region (province) in the world is very little in quantity. However, forecasting by LSTM using random weights with a little amount of data may result in very poor accuracy. So, to tide over these limitations, this paper proposes multi-source domain transfer learning. This results in a more generalized model which maximizes the accuracy of prediction.

The main contributions of this paper are summarized as follows:

1. Any contagious disease transmission is dependent on economic conditions, geographical parameters, population density, and several other demographical factors of an area. So, for the first time, this paper considers the population, area, and GDP of a province along with coronavirus infection data of confirmed cases, mortality, and recovery rate for forecasting future infectious spread. We gather a collection of datasets belonging to a set of related provinces to achieve better performance in forecasting.
2. Initially, we train each province-specific dataset on a simple LSTM model for future infection forecasting in that province. However, due to the unavailability of huge province-specific data, it faces a data scarcity problem which results in poor predictions. Therefore, this paper proposes a multi-source deep-transfer learning model to precisely forecast the infectious transmission of the disease.
3. Thus, multi-source deep-TL provides a warmth start for the sub-optimal weights which are obtained in initial training for enhancing the performance. Therefore, we perform fine-tuning of the initially trained model to obtain accurate forecasting.
4. From the empirical results, we show that the proposed methodology accurately forecasts the infectious spread of the disease in the upcoming days, and hence builds a robust predictive model.
5. This is the first such effort that exploits multi-source TL along with deep learning to forecast COVID-19 transmission.

The remaining paper is organized in the following manner. Thorough discussions on the background and the research related to our work are provided in Section II. After that, the dataset and comprehensive step-by-step description of our proposed approach are given in Section III. Section IV elaborately explains the experimental results in a comparative fashion. In Section 5, we provide a detailed discussion on the outcome of the work. Finally, Section VI discourses the conclusion of the work and mentions future directions.

## 2. Background and related works

This section presents the recent studies in the field of our proposed approach. It discusses works on novel coronavirus and significant approaches to TL. Much research was not done for the detection and prediction of COVID-19 infection and recovery. Recently, certain studies were published for the prediction of COVID-19 infection using statistical and mathematical aspects. One of the widely used models for observing the dynamics of various diseases was the Susceptible Infectious Recovered (SIR) model. It provides epidemic growth with the help of time-series differential equations [54]. Berger et al. [18] studied Susceptible Exposed Infectious Recovered (SEIR) epidemiology model to analyze the importance of testing and a conditional quarantine period. In the case of asymptomatic infection, if testing is not performed and strict quarantine policies are not observed then it may lead to the contagious spread of infection. Therefore, to control the epidemic, testing at a higher rate and strict quarantine are crucial measures. Godio et al. [19] studied SEIR epidemiological model for forecasting infection in Italian regions. They used a stochastic technique for training the model parameters using a Particle Swarm Optimization (PSO). They compared their results with the official data and forecasted infections for other

countries such as South Korea and Spain. Ghanbari et al. [20] discussed the Shannon entropy-based thermodynamics model for the prediction and propagation of COVID-19. Li et al. [21] anticipated Gaussian distribution theory using forward and backward propagation analysis for forecasting the COVID-19 infection. They used the current situation report of the Hubei for the predictions of the epidemic trends in Italy, Iran, and South Korea. Zhao et al. [22] used linear, a non-temporal mathematical model for infection transmission of coronavirus. Abir et al. [23] propose PCovNet, an LSTM, and variational autoencoder-based anomaly detection model to diagnose COVID-19 infections in the pre-symptomatic stage from resting heart rate with the help of wearable devices such as fitness trackers or smartwatches. Their findings suggest the usage of wearable devices along with a deep-learning framework as a secondary tool for the diagnosis of COVID-19 infection.

Shen et al. [55] applied the Logistic Growth model to predict the evolution of infected patients of COVID-19 in China, South Korea, and Iran using the rate of growth and population as the parameters. This model is time-independent, thus resulting in poor performance since time-dependent characteristics of COVID-19 are very crucial and shouldn't be ignored. Zhang et al. [56] studied the Poisson model which is based on power-law and exponential law to estimate the COVID-19 spread in Canada, France, Germany, Italy, the UK, and the USA. However, this model can only be used for short-term predictions; thus, an enhanced model is required for predicting infectious rates. Wang [57] discussed the applications, limitations, and potentials of mathematical models for COVID-19. He assumed and constrained the transmission rates as a constant to simplify the mathematical analysis and data fitting of the models. Such limitations may not allow a model to capture real-life fluctuations in transmission rates thereby affecting model generalization. Above discussed models also didn't consider the economic and population impact on infection rates of the pandemic. Thus, he insisted to use machine learning along with mathematical modelling to achieve better performance. Kumar et al. [58] performed a comparative analysis between statistical and deep learning models to predict COVID-19 infections, and the comparison was conducted based on mse and rmse values. Upon analysis, they found that, for most of the time-series data of the countries, deep learning-based models LSTM and GRU outperformed statistical ARIMA and SARIMA models. Sah et al. [59] conducted a comparative analysis using prophet, statistical ARIMA, and stacked LSTM-GRU model for efficient forecasting of COVID-19 infections. They found that stacked LSTM-GRU models outperformed the other existing models in terms of R square and RMSE values, thus they are proposed for future use. Jakka et al. [60] employed different learning models such as sigmoid modelling, SEIR, ARIMA, and LSTM for estimating the number of confirmed COVID-19 cases so that preventive measures can be performed, and the epidemic can be handled efficiently. They found that LSTM models resulted in more promising outcomes, and it can help decision-makers to take necessary actions for the effect of interventions. All the mathematical models that were studied above show limitations in capturing temporal components of the data, which is very crucial for forecasting COVID-19 infections, as upcoming infections in the next few days depend on past observations. Later, the techniques of artificial intelligence were studied.

Artificial intelligence and data mining are vital components for successful technology in the medical sector [24]. Dianbo et al. [25] used a clustering approach for forecasting of COVID-19 outbreak with the help of news alerts from GLEAM, an agent-based mechanistic model [25]. Yao et al. [26] used lung ultrasound images instead of chest CT for the detection of infection. They found that this method would be used to reflect both the infection duration and disease severity in a person. Esther et al. [27] found that oral rinses and posterior oropharyngeal saliva could be used as an alternative swab collection method for the detection of SARS-CoV-2 RNA by RT-PCR. Narin et al. [28] used chest X-ray images and applied several convolutional neural network models such as InceptionV3 and ResNet 50 for the diagnosis of infection. Qin

et al. [29] incorporated social media platforms for daily case prediction. Social media search indexes (SMSI) were collected for dry cough, fever, pneumonia, coronavirus, and chest distress from December 31, 2019 to February 9, 2020. The lagged series of SMSI were used to forecast new suspects using several techniques such as subset selection, lasso regression elastic net, ridge regression, and forward selection. The subset selection was the optimal method obtained during validation. Huang et al. [30] propose the LightEfficientNetV2 CNN model for COVID-19 classification with the help of CT scans and X-ray images. Their findings resulted in an accuracy of up to 98.3% using the proposed approach. Ahmad et al. [31] make use of deep-CNN models for extracting significant features and classifying infected patients using X-ray images. They employed ImageDataGenerator to overcome the small dataset size problem, and their findings resulted in an accuracy of up to 97.68% as compared to the baseline models. Li et al. [61] propose MultiR-Net, a 3D deep learning model for combined COVID-19 classification and lesion segmentation, with the help of U-Net to achieve real-time and interpretable COVID-19 chest CT diagnosis. Their proposed approach resulted in an accuracy of up to 93.23% making it suitable for future uses. Jin et al. [62] propose an ensemble hybrid model based on CNN, GRU, deep belief networks, Q-learning, and SVM to achieve forecasting of COVID-19 infections. Their analysis ensures the accuracy, robustness, and generalization of their approach. However, the efficient forecasting of infections is interdependent on several factors such as population density and economic conditions of a region which were not considered by them. Altan et al. [46] proposed a hybrid model consisting of 2D curvelet transformation, chaotic salp swarm algorithm, and deep learning for the detection of COVID-19 infections using X-ray images. Thus, deep-learning techniques were efficiently used for the detection and transmission of the disease. However, these models were not suitable for real-time based dataset and specifically for coronavirus, because it is time series-based dataset.

Therefore, it becomes crucial to introduce a sequential network model for real-time forecasting of infection. Karasu et al. [47] anticipated a technique for forecasting crude oil prices as it possesses non-linear dynamics in a time-series environment. They studied efficient techniques for extracting relevant features and then applied the time-series model. Altan et al. [48] proposed a hybrid model consisting of LSTM and grey wolf optimizer for forecasting wind speed required for the generation of wind power. They provided an efficient forecasting algorithm using LSTM. Wang et al. [50] proposed an LSTM model along with a rolling update mechanism to accurately forecast COVID-19 infections for long-term projections. But their study focused only on the data of confirmed cases for training, which is not alone sufficient to provide the optimal results.

Certain mathematical and statistical models [10,11] were proposed to predict the evolution of COVID-19. Statistical algorithms, for example, ARIMA (Autoregressive Integrated Moving Average) depends immensely on assumptions. Many of the data-driven techniques work linearly [12] neglecting the temporal attributes of the data. In many cases, these models resulted in low accuracy and were not able to fit the data properly. Thus, to overcome the constraints of statistical models, this paper employs a deep-learning-based Recurrent Neural Network (RNN) to forecast real-time infection cases. RNN is an integral deep learning method useful for temporal components in the time-series analysis [13]. It is a sequential neural network that remembers things learned from prior inputs while generating outputs [14,15] but they were not capable of learning historical dependencies [16]. So, to overcome this shortcoming, Hochreiter et al. have designed a long-term dependency LSTM model to regulate information flow in memory cells using hidden layer units and multiplicative gates [17].

Kirbas et al. [49] performed a comparative analysis using ARIMA, NARNN, and LSTM models to forecast COVID-19 infections in European countries. Upon analysis, LSTM performed better which could be used further in the future [49]. Kumar et al. [32] studied the SARIMA technique for forecasting cumulative COVID-19 cases in the top 16 affected

countries. However, they assumed the parameters using the auto-ARIMA methodology, making it difficult to fit the data in every situation. Matamoros et al. [33] used the ARIMA model for forecasting per region future infections. They proposed the correlation between countries of similar geographical areas and COVID-19 infections. They inferred that, if the geographical area of the two regions was similar, then the future infections advancements in such regions would also be similar. But future infectious spread not only depends on the geographical area of the region but also on population density, education, health care facilities, and the economy of a region. These parameters were not considered by them. So, further works were necessary to precisely predict future infections. Huang et al. [34] proposed a CNN-LSTM model for COVID-19 forecasting in China. But their analysis proposed forecasting's for a specific area, which might not perform well in some other regions. Arora et al. [35] used daily confirmed cases time-series data for the forecast of disease in all states and union territories of India. They analyzed certain variants of the LSTM model such as Stacked LSTM, Convolutional LSTM, and Bi-directional LSTM. The optimal method was used further for one-day and then one week of forecasting. Their study covered only the data of confirmed cases as a parameter, and it was specific to a particular region. However, a more generalized model with a wide range of features was required for accurate forecasting of COVID-19 infections. Chimmula and Zhang [36] anticipated LSTM networks for forecasting in Canada. All the studies proposed above were either utilized for forecasting in a specific geographical region or considered very few parameters. Consequently, these models have restricted the scope of forecast in a certain country and were not able to preserve spacio-temporal components of the data. Thus, to forecast infection in province-specific region, the above proposed models face data shortage issue as historical data taken from a specific -province is insufficient for precise forecast.

On the other hand, TL approaches performed better in many such scenarios. Kumar et al. [37] studied transfer learning approach for GDP prediction using carbon emission dataset. Loey et al. [38] studied hybrid technique of transfer learning and machine learning algorithms for detection of face mask in this pandemic. Gautam [39] predicted coronavirus cases and deaths of several countries using TL technique. The model was trained on countries like Italy and US and tested on different set of countries. But the epidemic spread of COVID-19 depends also on other parameters such as population, geographical area, and economic status of a region, which are not studied by the authors. Su et al. [40] presented a multilevel thresholding image segmentation method based on an enhanced multiverse optimizer, to improve the efficient processing of COVID-19 chest films, with the addition of vertical and horizontal search mechanisms. Their findings provided an effective approach that can be used in a medical organization with the diagnosis of infection.

However, the evolution of COVID-19 infection not only depends on daily confirmed cases and fatality rate but also on other parameters of geographical area, population, and GDP [6–9]. Armando et al. [41] discussed the relationship between positive coronavirus cases and transport accessibility in an area. The higher the accessible any certain geographical area, the easier the virus reaches its population. Wu et al. [42] focused on meteorological parameters and experimented that relative humidity and temperature had a positive effect on the spread of COVID-19.

The related studies reviewed until now indicate that mathematical modeling and some artificial intelligence approaches played a crucial role in forecasting, detecting, and propagation of the COVID-19 pandemic, which could help governments to manage the situation efficiently. To the best of our knowledge, no work has been reported so far using population, geographical parameters, and economic conditions along with mortality and recovery rate for the prediction of transmission of COVID-19 infection using the multi-source deep-TL model. Therefore, the current study proposes to bridge this gap.

### 3. Proposed approach

This section presents our proposed multi-source deep transfer learning-based novel approach to robustly forecast the COVID-19 transmission in a province within a country. Initially, in the first subsection, the dataset used in this paper is described. In the next subsection, the proposed multi-source deep-TL approach is explained.

#### a) Datasets

This paper proposes a multi-source deep-TL-based approach for multi-step forecasting of future infectious transmission of COVID-19 in a province. The multivariate dataset includes confirmed cases of COVID-19, number of fatalities, and number of patients recovered along with geographical parameters i.e., Population, Area, and GDP per capita of a province. The dataset utilized in this paper for experimentation is extracted from the John Hopkins University repository, from the first case of COVID-19 recorded till December 23, 2020 [43]. The multivariate dataset generated is in the time-series format with month, date, and year; therefore, preserving temporal components of data. The input parameters useful for training purposes are given in Table 1, and the output parameter is given in Table 2.

Fig. 1 depicts a visual representation of total COVID-19 tests per 1000 individuals vs GDP per capita for all countries around the globe [44]. It represents the correlation between appropriate testing of COVID-19 being done and the GDP of a nation. Thus, if the GDP of any country is higher, which is an indication of a sustainable economy with enough resource availability, then higher is the testing rate for the detection of infections nationwide. Thus, this would impact the authentication of the statistical report of infected cases. Countries with deprived GDP provide inappropriate testing reports due to a lack of resources and economic-social hardships. Thus, it can be visualized from Fig. 1 that countries with higher GDP are providing a high amount of testing. Therefore, these developed nations are providing accurate statistics that can be utilized for forecasting future infections. Therefore, we have utilized several provinces of seven developed nations such as: Australia, Canada, Denmark, France, Netherlands, the United Kingdom, and the United States for efficient forecasting of COVID-19 infectious spread.

#### b) Proposed multi-step deep TL methodology

A recurrent Neural Network (RNN) is certainly the most promising deep learning approach to handle temporal components as it provides the capability to learn sequentially [27–29]. However, due to the problems of vanishing and exploding gradients, an enhanced model of LSTM is preferred. It provides the solution by presenting a new memory state in RNN which learns parameters for a long duration [16]. Thus, it is suitable for time-series modeling.

In a deep neural network framework, the initial layer is focused on capturing basic patterns and features of the data. Then deeper layers are aimed to extract more complex patterns and advanced features. The success of LSTM neural networks lies with the usage of multiplicative gates, named: input gate, forget gate, and output gate. The architecture of LSTM is depicted diagrammatically in Fig. 2. We employ this LSTM

**Table 1**  
Input parameters description.

| INPUTS | DESCRIPTION            |
|--------|------------------------|
| 1      | Observation Date       |
| 2      | No. of confirmed cases |
| 3      | No. of Deaths          |
| 4      | Recovery rate          |
| 5      | Population             |
| 6      | Area (sq. km)          |
| 7      | GDP per capita (US \$) |

**Table 2**  
Output parameter description.

| OUTPUT | DESCRIPTION                  |
|--------|------------------------------|
| 1      | Future Coronavirus Infection |

neural network regression technique for deep-transfer learning for multistep forecasting of COVID-19 infectious transmission in a province.

**Algorithm 1.** Multi-source deep-TL for COVID-19 infection forecasting using LSTM.

*Algorithm 1:* Multi-source deep-TL for COVID-19 infection forecasting using LSTM.

*Output:* Future forecasting of coronavirus infection.

*Input:* Multi-source domain dataset (Confirmed cases, Deaths, No. of Recoveries, Population, Area, GDP per capita)

1. Multi-source domain datasets creation:
  - a. Gather multi-source domain datasets (Table 1) from several province of various countries of the world and pre-process them by removing the outliers and performing normalization.
2. Train a deep transfer learning model over this dataset.
3. Validate the loss of the model over this multi-source domain dataset by calculating the mean absolute percentage error (MAPE) and coefficient of determination.
4. Gain knowledge by learning the sub-optimal weights.
5. For each province do:
  - a. Retrain the above-trained model (Step 4) on an individual province using deep TL. Compute MAPE and coefficient of determination on each province to maximize the performance.
  - b. The fine-tuned model is used to precisely forecast the future infectious spread of that province.
6. Return future forecast of COVID-19 infection spread for every province.

The flow diagram in Fig. 3 (a) describes COVID-19 infection forecasting using individual province specific dataset. Single province multivariate dataset is taken at a time which is trained on LSTM model and multistep forecast of future coronavirus infection in that province for next few days are observed. This resulted in imprecise prediction. Since, LSTM is a deep-learning model, it requires enormous amount of trainable data. The availability of COVID-19 dataset for a specific province is very little in quantity. When LSTM is trained using random weights initialization, it exhibits low variance and high bias. That is, it may underfit<sup>2</sup> the data due to insufficient data availability. Therefore, to overcome this limitation of under fitting due to data scarcity, we have utilized the multi-source deep-TL technique as visualized in Algorithm 1. The flow diagram of this Algorithm 1 is also demonstrated diagrammatically in Fig. 3 (b).

The stepwise description of the Algorithm 1 as follows:

#### Step 1 Multi-source domain dataset creation:

Multi-source domain COVID-19 datasets are collected for several provinces of the world with features represented in Table 1. An individual province data represents a single source domain, as it is exclusive. We have gathered source domain data from several provinces, which possess huge data distribution differences from each other, with respect

<sup>2</sup> High variance and low bias (i.e., overfitting) maybe checked by splitting the data into train data and validation data.

to variances in their geographical locations and their economic conditions. Thus, it is termed as a multi-source dataset.

#### Step 2 Deep-Transfer Learning:

At this step, the supervised deep-learning regression model LSTM is utilized. It is a type of RNN model, that has the characteristic of learning long-term dependencies in the historical-dependent sequential dataset. The collective time-series dataset representing sixty-two provinces is gathered providing multivariate parameters (generated in Table 1) for training the model to forecast future regression of infectious spread. The

LSTM network consists of three layers with 128, 64, and 32 neurons respectively, and an output layer with 6 neurons making 6-days ahead predictions. The rectified linear (ReLU) activation function is used for the LSTM blocks. ReLU is a piecewise linear function that will output the input directly if it is positive, otherwise, it will output zero. The network is trained for 256 epochs and a batch size of 64 is deployed. Furthermore, during the compilation period, the 'Adam' optimizer is considered which captures the desired properties of AdaGrad and RMSprop techniques to generate an optimization algorithm that is best suited to handle gradients.

To enhance the performance of the model, several callback functions are also utilized in the proposed network which is: 'Early Stopping', 'Model Checkpoint', and 'Reduce LR On plateau'. The Early Stopping callback monitors the performance of the training data in terms of the 'validation loss' metric and, it gets triggered when the performance stops improving thus, the last training epoch gets recorded i.e. when no improvement is observed for '10' continuous epochs, then the training of the model is stopped. While ReduceLRonPlateau is employed to reduce the learning rate by a factor of 0.1 if no improvement is observed for the '3' epochs. Further, these callbacks will terminate the training process once triggered, but the resulted model at the end of the training process may not be the best-performing model on the validation dataset. Thus, to achieve this an additional callback known as 'Model Checkpoint' is required. The best performing model while training gets saved and later used in validation process by this technique.

#### Step 3 Prediction precision evaluation:

The multivariate dataset comprises historical dependent data with

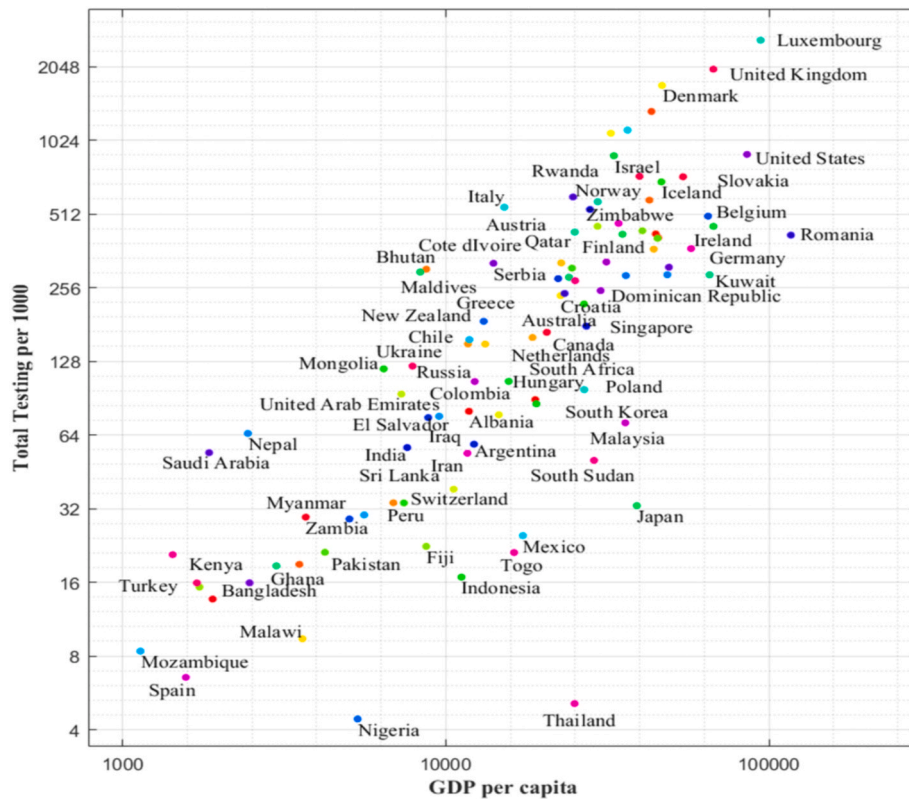


Fig. 1. Representation of COVID-19 Testing vs GDP.

COVID-19 cases starting from January 26, 2020 to December 23, 2020 for all the provinces. Therefore, 80% of the dataset ranging from 26/1/2020–18/10/2020 for every province is used for training the model. Then it is validated on the remaining 20% dataset i.e.- from October 19, 2020 to December 23, 2020 for their prediction preciseness in estimating future infection transmission. The statistical performance has been evaluated using two different statistical measures, which are: mean-absolute percentage error (MAPE) and coefficient of determination. They are well-known evaluation metrics for regression approaches. MAPE is computed using the Eq. (2) as follows:

$$M = \frac{1}{N} \sum_{t=1}^n \left| \frac{O_t - F_t}{O_t} \right| \tag{2}$$

Here,  $F_t$  is the forecasted output value,  $O_t$  is the observed value, and  $N$  is the total number of observations. It measures accuracy as a percentage. However, an inverse proportion relation is observed between MAPE and the forecasted accuracy. Thus, the lesser the MAPE outcome, the better

the designed model.

Step 4 Knowledge Extraction:

The best model obtained in initial training is saved and loaded consisting of sub-optimal weights. These sub-optimal weights are thus feasible, which can further be optimised during retraining. Therefore, instead of using random weights, sub-optimal weights are being used with lower learning rate for training the prediction model. Hence, transfer learning has been successfully utilized.

Step 5 Fine-tuning of the trained model for each province:

Now, the sub-optimal weights (trainable parameters weights) obtained from initial training are considered for parameter-tuning. Consequently, one province dataset containing input parameter (Table 1) is considered and retrained on LSTM network model using sub-optimal weights which are already attained from initial training. However, this process is repeated for all sixty-two-provinces' dataset and infectious spread is forecasted for each province. The approach of using multi-source deep-TL have shown robustness and resulted in enhanced forecasting, and thus, representing minor forecast error.

Step 6 COVID-19 Infection Forecasting

Finally, the trained multi-source deep-TL model is used to forecast the future coronavirus cases. The source domain dataset consists of sixty-two provinces which are:

- { 'Alabama', 'Arizona', 'French Guiana', 'Ohio', 'Aruba', 'Alaska', 'Bermuda', 'Bonaire', 'Alberta', 'Sint Eustatius and Saba', 'British Columbia', 'Arkansas', 'California', 'French Polynesia', 'Gibraltar', 'Cayman Islands', 'Channel Islands', 'Australian Capital Territory', 'Kentucky', 'Greenland', 'Idaho', 'Colorado', 'Victoria', 'Ontario', 'South Carolina', 'Western Australia', 'Curacao', 'Wyoming', 'Virginia', 'Pennsy

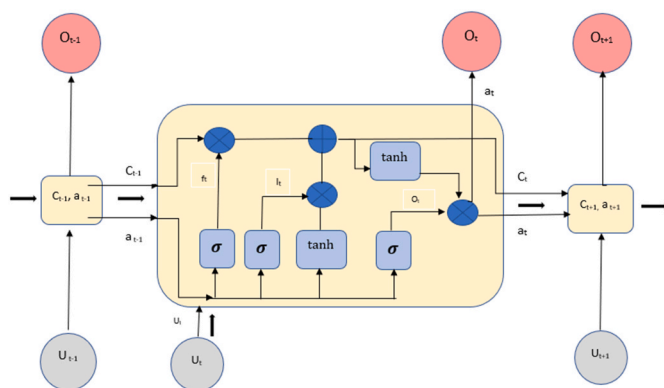


Fig. 2. LSTM architecture [45].

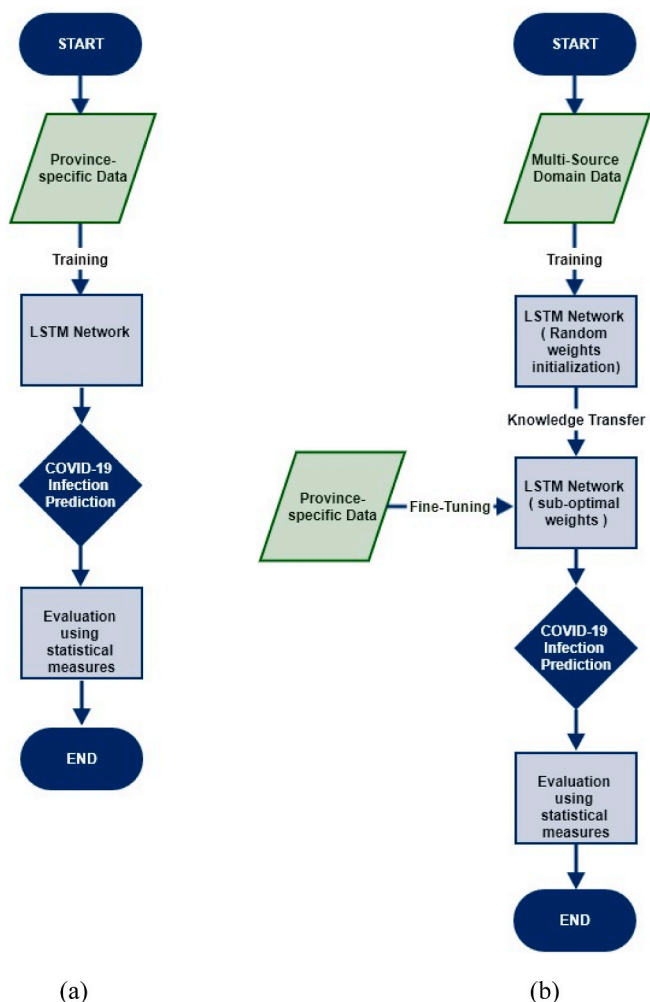


Fig. 3. Training deep learning model: (a) without TL (b) with TL.

Ivania', 'Indiana', 'South Dakota', 'Faroe Islands', 'Florida', 'Georgia', 'Guadeloupe', 'Hawaii', 'Washington', 'Isle of Man', 'Kansas', 'Manitoba', 'Illinois', 'Martinique', 'Reunion', 'Mayotte', 'New Brunswick', 'New found land and Labrador', 'New South Wales', 'New York', 'Nova Scotia', 'Oregon', 'Puerto Rico', 'Quebec', 'Queensland', 'Iowa', 'Saskatchewan', 'Sint Martin', 'South Australia', 'St. Martin', 'Texas', 'Turks and Caicos Islands', 'West Virginia', 'Wisconsin' }.

The selection of these provinces is totally arbitrary. One can choose any set of provinces' datasets. There are no specific inclusion-exclusion criteria that are considered for the selection of these provinces, they are just randomly chosen. Now, it comprises seven developed countries. However, the major emphasis is on the selection of developed countries datasets for the experimentation. Developed countries have certainly higher GDP i.e., significantly available resources for COVID-19 testing to generate definite statistics report, which is beneficial for precise evaluation of the model. Their statistics are more reliable as compared to developing nations statistics, since developing countries lack with enough resources to precisely perform COVID-19 testing. Thus, the motivation is to perform experimentation with considering datasets of developed nations, and in future works this trained model can be utilized for efficient forecasting of COVID-19 infections in developing nations.

#### 4. Experimental results and analysis

In this section, the human transmission of COVID-19 infection in the

future is experimentally forecasted using the multi-source deep-TL model. It consists of two sections: the first one is on experimental results evaluation, and another sub-section graphically depicts the COVID-19 Infection estimation of the next few days using the proposed methodology.

#### A. Experimental results evaluation

This section describes the outcomes obtained using the multi-source deep-TL model for the prediction of COVID-19 infectious spread. Here, we have employed two different types of techniques for forecasting: the first one is termed as a without-TL approach. In this approach, a province-specific dataset is considered which is trained on the LSTM model with randomly initialized weights, and thus, the forecasted accuracy is observed. This procedure is followed for all the provinces individually. Since data available in each province is not in large quantity, it may result in poor performance. The second type of technique is termed as with-TL which can overcome the limitations of data shortage. In the second approach, a collective dataset comprising of all sixty-two provinces is gathered. Considering it as a borne start, the sub-optimal weights are initialized using LSTM model and collective training is performed. Subsequently, the individual province-specific data is considered at a time, which is retrained on the initially trained model using the already obtained sub-optimal weights and lower learning rate. Thus, fine-tuning using multi-source deep TL is performed. This approach is termed as TL approach. One of the common performance indicators used to measure forecast accuracy is the root mean square error (RMSE). But they are sensitive to outliers. Therefore, the performance of the proposed model is measured using two different statistical measures, which are mean absolute percentage error (MAPE) and coefficient of determination (R-squared test). Both are commonly used to measure prediction accuracy. This sub-section is now further divided into two parts, where the initial part describes the performance of the proposed approach of COVID-19 infections forecasting using the mean absolute percentage error evaluation metric and the second part presents the performance of the proposed methodology using the coefficient of determination evaluation metric.

##### i. Performance evaluation using MAPE

Table 3 Mean absolute percentage error.

| Province Name            | Mean absolute percentage error |       | % Improvement | % Deterioration |
|--------------------------|--------------------------------|-------|---------------|-----------------|
|                          | Without-TL                     | TL    |               |                 |
| Alabama                  | 5.9                            | 1.4   | 76.2          | -               |
| Alaska                   | 32.7                           | 7.3   | 77.6          | -               |
| Arizona                  | 35.8                           | 2.5   | 93.01         | -               |
| Bermuda                  | 5.9                            | 1.1   | 81.35         | -               |
| California               | 3.66                           | 1.0   | 72.67         | -               |
| Faroe Islands            | 2.25                           | 0.57  | 74.66         | -               |
| Florida                  | 99.31                          | 3.58  | 96.39         | -               |
| Georgia                  | 36.25                          | 1.65  | 95.44         | -               |
| Illinois                 | 10.5                           | 8.37  | 20.28         | -               |
| Indiana                  | 66.1                           | 1.50  | 97.73         | -               |
| Manitoba                 | 5.37                           | 4.21  | 21.60         | -               |
| New South Wales          | 1.53                           | 1.20  | 21.56         | -               |
| Pennsylvania             | 3.7                            | 1.72  | 53.51         | -               |
| Quebec                   | 4.11                           | 0.62  | 84.91         | -               |
| South Carolina           | 22.49                          | 1.38  | 93.86         | -               |
| Texas                    | 13.8                           | 1.07  | 92.24         | -               |
| Virginia                 | 22.9                           | 3.68  | 83.93         | -               |
| Wyoming                  | 13.77                          | 4.17  | 69.71         | -               |
| Greenland                | 4.86                           | 8.37  | -             | -72.22          |
| Martinique               | 48.58                          | 50.86 | -             | -4.69           |
| Queensland               | 4.03                           | 10.16 | -             | -152            |
| Turks and Caicos Islands | 10.22                          | 22.06 | -             | -115            |



The performance of the proposed methodology in forecasting COVID-19 transmission has been initially evaluated using mean absolute percentage error. We have shown the results of some provinces after performing training using the proposed approach. We have represented provinces' outcomes in terms of MAPE using without-TL as well as TL technique as depicted in Table 3. Further, it describes the percentage of improvement in terms of the prediction error from the classical (without-TL) approach to a multi-source deep-TL technique for the provinces, where our proposed algorithm performs better as visualized in column 3 of Table 3. However, a progressive decline can be observed for a few provinces where the classical approach is more suitable. Thus, they are depicted as deterioration percentages in column 4 of Table 3.

The provinces which are used for result visualizations are: 'Alabama', 'Alaska', 'Arizona', 'Bermuda', 'California', 'Faroe Islands', 'Florida', 'Georgia', 'Illinois', 'Indiana', 'Manitoba', 'New South Wales', 'Pennsylvania', 'Quebec', 'South Carolina', 'Texas', 'Virginia', 'Wyoming', 'Greenland', 'Martinique', 'Queensland', 'Turks and Caicos Islands'.

In Table 3, the best results of forecasting accuracy depicting minimum error in terms of MAPE between the two approaches are marked in bold. There are certain provinces in which the proposed technique of multi-source deep-TL model has performed outstandingly as compared to without-TL approach. For the provinces such as Arizona, Bermuda, Florida, Georgia, Indiana, Quebec, South Carolina, Texas, and Virginia, the percentage of improvement in terms of prediction accuracy, from simple model using without-TL to deep-TL model is significantly higher than 80% as visualized in Table 3. Arizona has resulted 93.1% improvement, thus provided better forecast accuracy when performed using multi-source deep-TL approach. Similar behavior is observed for other provinces such as Florida with 96.39% improvement, and Quebec with 84.91% improvement.

Consequently, for certain province such as Alabama, Alaska, California, Faroe Islands and Wyoming, the improvement percentage varies from 65 to 80%. For Alabama, classical model results 5.9% error whereas our proposed multi-source deep-TL-based model obtained 1.4% of error. Thus, improving the performance by 76%. Provinces like Illinois, Manitoba, New South Wales, and Pennsylvania have attained improvement up to 65%. Illinois has obtained 20.28% improvement whereas New South Wales has resulted in 21.6% of improvement. These results have proved that the forecasting precision of the multi-source deep-TL approach is better than without-TL.

ii. Performance evaluation using the coefficient of determination

Now, the performance of the proposed approach in forecasting COVID-19 infectious spread is again evaluated using the coefficient of determination, which is another statistical evaluation metric. The coefficient of determination is a statistical measure to examine the relationship between two variables. It describes the amount of variability of one parameter caused due to its relationship with another parameter. This coefficient is commonly termed R-squared (or R<sup>2</sup>), and is usually referred to as the "goodness of fit." A value close to 1 indicates a perfect fit model, and thus, a highly reliable model for future forecasts. It can be computed using the mathematical formula given as Eq. (3).

$$R^2 = 1 - \frac{\sum_{i=1}^N (\hat{Y}_i - Y_i)^2}{\sum_{i=1}^N \left( Y_i - \frac{\sum_{k=1}^N Y_k}{N} \right)^2} \tag{3}$$

where Y<sub>i</sub> presents the actual value of a variable for i = 1, 2 .... N and represent the forecasted value of a variable.

We have presented the performance of the proposed approach using an R-squared statistical measure in Table 4. It comprises four columns: the first column gives the province name, the second column represents the R-squared value using the classical training model (without-TL approach), and the third column presents the performance using the

Table 4  
Coefficient of determination.

| Province Name            | Coefficient of determination |              | Improvement in R <sup>2</sup> |
|--------------------------|------------------------------|--------------|-------------------------------|
|                          | Without-TL                   | With-TL      |                               |
| Alabama                  | -18.94                       | <b>-0.64</b> | 18.3                          |
| Alaska                   | -81.21                       | <b>-4.09</b> | 77.12                         |
| Arizona                  | -1378                        | <b>0.42</b>  | 1378.42                       |
| Bermuda                  | -11.66                       | <b>0.34</b>  | 12                            |
| California               | -19.78                       | <b>0.48</b>  | 20.26                         |
| Faroe Islands            | -49.2                        | <b>-29.3</b> | 19.9                          |
| Florida                  | -12026                       | <b>-2.64</b> | 12023.36                      |
| Georgia                  | -2177                        | <b>0.84</b>  | 2177.84                       |
| Illinois                 | -17.44                       | <b>0.5</b>   | 17.94                         |
| Indiana                  | -474                         | <b>-5.22</b> | 468.78                        |
| Manitoba                 | -27.41                       | <b>0.67</b>  | 28.08                         |
| New South Wales          | -93.1                        | <b>-33.1</b> | 60                            |
| Pennsylvania             | -21.08                       | <b>-0.03</b> | 21.05                         |
| Quebec                   | -12.89                       | <b>0.71</b>  | 13.6                          |
| South Carolina           | -1328                        | <b>-1.81</b> | 1326.19                       |
| Texas                    | -42.9                        | <b>0.41</b>  | 43.31                         |
| Virginia                 | -316.4                       | <b>-7.19</b> | 309.21                        |
| Wyoming                  | -74.01                       | <b>-8.47</b> | 65.54                         |
| Greenland                | <b>0.2</b>                   | -0.66        | -0.86                         |
| Martinique               | <b>-7.53</b>                 | -9.31        | -1.78                         |
| Queensland               | <b>-125.8</b>                | -234.5       | -108.7                        |
| Turks and Caicos Islands | <b>-19.16</b>                | -53.17       | -34.01                        |

proposed with-TL approach (for each province, the best R<sup>2</sup> values are marked with bold in Table 4), and the last column presents the improvements in the coefficient of determination using the with-TL approach over the without-TL approach. It is computed by subtracting the R<sup>2</sup> values given by the without-TL from the same proposed approach. The higher the magnitude of the values in the fourth column of Table 4, the better the performance of the with-TL model over the without-TL model. That is, the goodness of fit (R<sup>2</sup>) is best in the former (the with-TL model).

There are numerous provinces in which the proposed approach of forecasting COVID-19 infectious spread using the multi-source deep-TL model has performed outstandingly when compared to classical training (without-TL model). A few such provinces are: 'Alaska', 'Arizona', 'California', 'Florida', 'Georgia', 'Illinois', 'Indiana', 'Manitoba', 'Quebec', 'South Carolina', 'Virginia', 'and Wyoming'. A maximum of 0.84 similarity has resulted in actual and forecasted infections for the province of Georgia. The value is very close to 1, which means precise forecasting is obtained using the proposed approach. For the province of Arizona, an improvement of 1378.42 in R<sup>2</sup> has been observed using the TL model over the without-TL approach. The province of Florida is also able to provide a precise prediction model using TL model by resulting in an improvement of 12023.36 in terms of R<sup>2</sup> value. More generalized 'goodness of fit' (R<sup>2</sup> values closer to 1) model's is observed in 'Manitoba', 'Alaska', 'Florida', 'Georgia', 'Indiana', 'South Carolina', 'Virginia', 'Wyoming' with an improvement of 28.08, 77.12, 12023.36, 2177.84, 468.78, 1326.19, 309.21, 65.54, respectively (as shown in 4th column of Table 4). Similar results are obtained for other provinces where the proposed approach has increased the efficacy of the model over classical training model.

However, there are 4 provinces out of the total 62 provinces in total, where the proposed approach resulted in poor performance in R<sup>2</sup> metric. Similar trends of performance degradation of the proposed approach for these 4 provinces were also observed during their MAPE result analysis. These provinces are: 'Greenland', 'Martinique', 'Queensland' and 'Turks and Caicos Island'. The province of Greenland has resulted in -0.86 degradation (negative improvement) with respect to the proposed with-TL approach. Similarly, for the province of Queensland, a decline of -108.7 in R<sup>2</sup> is resulted.

B. COVID-19 Infection estimation of the next few days

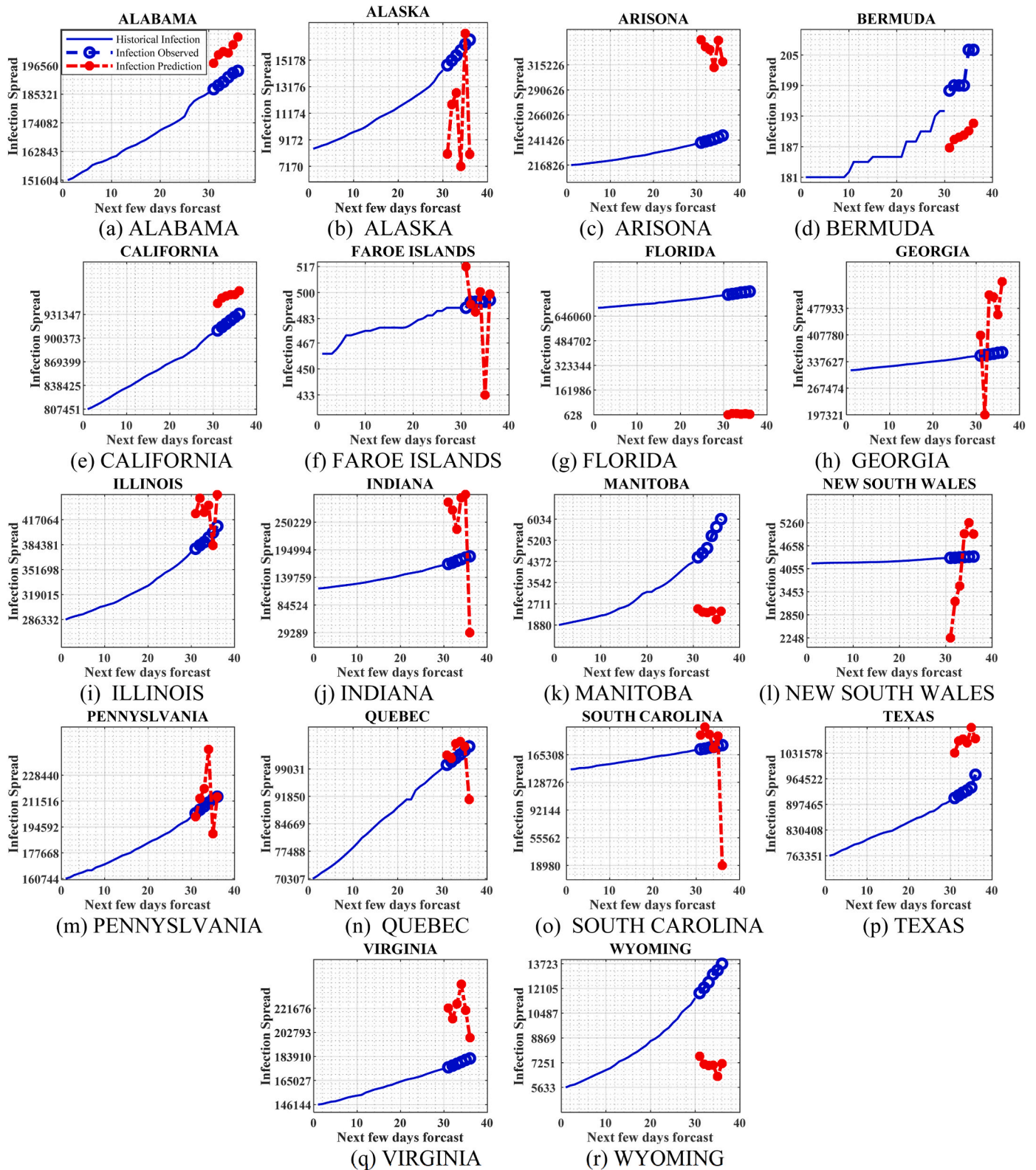


Fig. 4. COVID-19 Infection forecast using without-TL approach.

The COVID-19 Infections forecasting using the proposed approach of multi-source deep transfer learning is depicted graphically in Figs. 4–6. Fig. 4 depicts COVID-19 infection using the without-TL approach. These figures have 3 lines for visualization. A solid blue line depicts historical Infection in the past 30 days, red-dotted lines demonstrate infection forecasting using without-TL approach in the upcoming 6 days, whereas

actual observed infection in the next 6 days is represented by blue-dotted lines.

Fig. 4 (a) demonstrates a graphical representation of COVID-19 infection in Alabama, a province of US using the without-TL approach. The x-axis represents the number of days and the y-axis represents infectious spread. The days ranging from 0 to 30 describe a span of 30 days

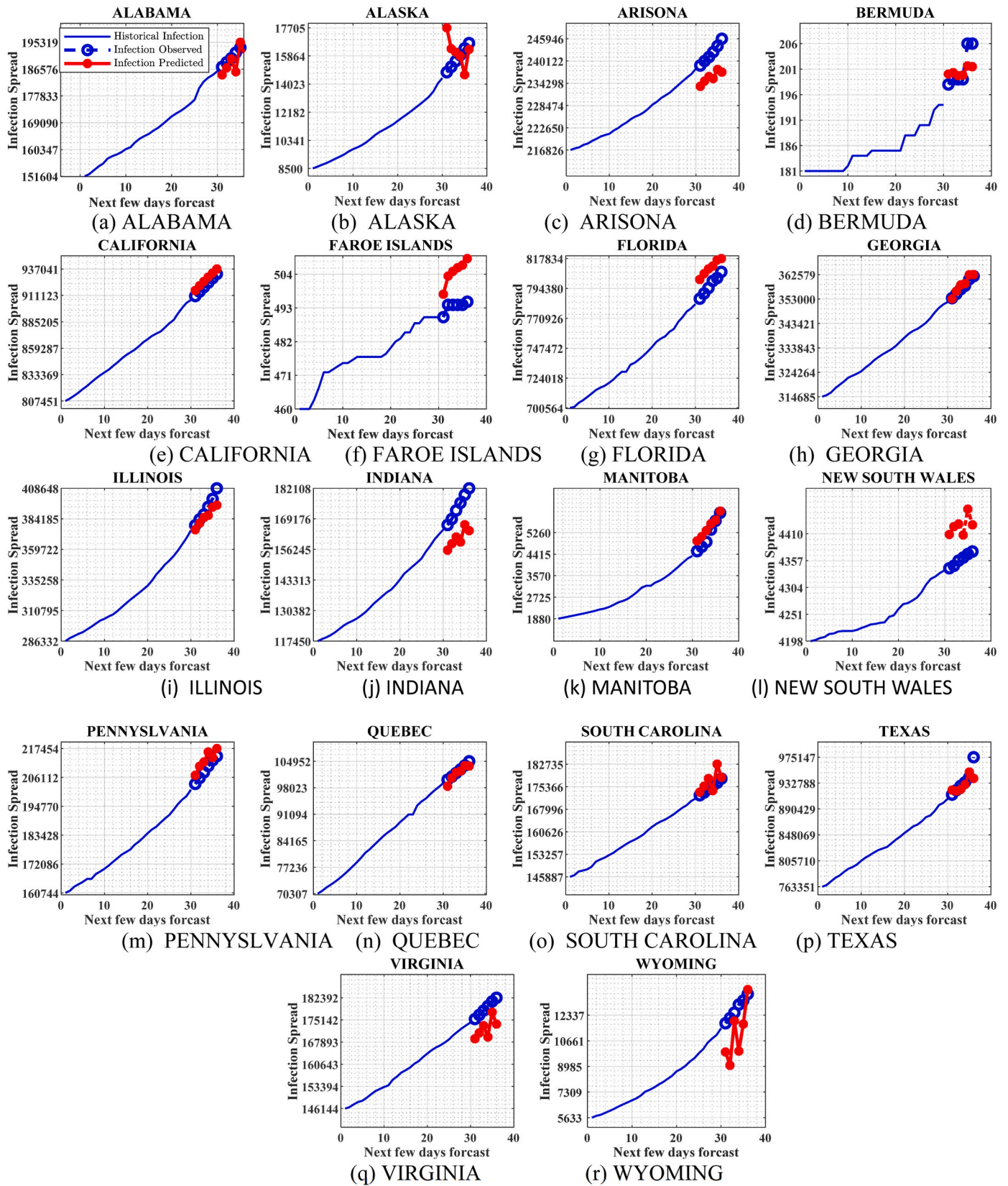


Fig. 5. COVID-19 Infection forecast using multi-source deep- TL approach.

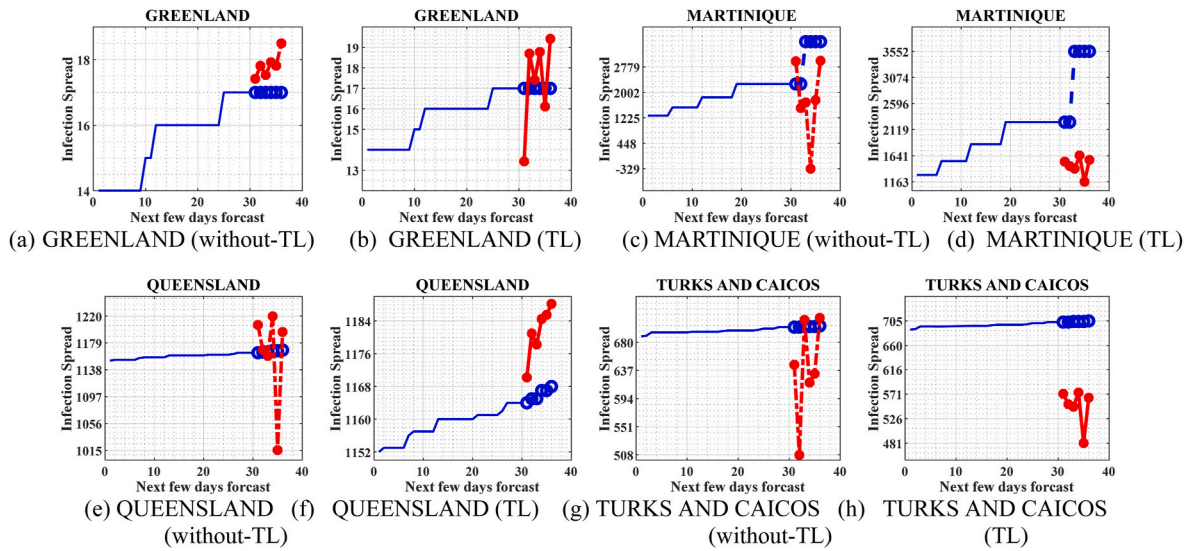


Fig. 6. COVID-19 Infection forecast depicting weakness of the proposed approach.

demonstrating historical infection observed, whereas days from 31 onwards describe infection prediction in the upcoming 6 days as plotted with the red-dotted line, and actual infection observed is depicted in a blue-dotted line. However, the blue dotted line and red dotted line are far apart. Thus, significant variations can be visualized from the figure between actual cases and forecasted ones which resulted in 5.9% of error. Similar observations can be visualized for other provinces when trained using a simple-LSTM regression network without incorporating the TL approach as represented in Fig. 4 (b)–(r). These figures with infection prediction using simple LSTM networks without performing TL depicts large fluctuating outcomes resulting in huge prediction error.

Fig. 5 describes COVID-19 infection using the multi-source deep-TL methodology. These figures have 3 different lines for visualization. The solid blue line depicts historical Infection in the past 30 days, solid red lines demonstrate infection forecasting using the multi-source deep-TL model in the upcoming 6 days, and actual suspects in the next 6 days is represented by blue-dotted lines. Fig. 5 (a) represents forecast modeling in Alabama with the use of the multi-source deep-TL model. The observed infection and predictions are overlapping with each other as it can be visualized from blue dotted points and solid red lines. Thus, these predictions demonstrate a good match of predicted infection against actual infection which resulted in a minimal error of 1.4%. Thus, we can say that the proposed model is more robust and precise for COVID-19 infection forecasting. Similar forecasting can be visualized for other provinces using our proposed multi-source deep-TL approach as depicted in Fig. 5(b)–(r). Hence, it is observed that the proposed multi-source deep-TL methodology has resulted precise forecasting in comparison to the without-TL approach. Also, the without-TL approach attained irregular predictions whereas the proposed approach obtained overlapping predictions with actual observed values.

However, there are few cases where the proposed algorithm doesn't perform better. In these cases, the without-TL approach has obtained better accuracy for predicting future forecasts of infectious spread. For certain provinces such as Greenland, Martinique, Queensland and Turks and Caicos Islands, the without-TL approach has resulted in better forecasting as depicted in Fig. 6 (a), (c), (e), and (g). The red-dotted lines are closer to blue-dotted points in comparison to other figures with captions of TL for the same provinces. In case of Greenland 4.86% error using without-TL approach is obtained, which got significantly increased with the deep-TL approach. Similar results are obtained using  $R^2$  statistical measure. The deterioration of  $-0.86$  in  $R^2$  is observed for Greenland province using the proposed TL model over classical (without-TL) model. The Queensland province attained 4.03% MAPE

using without-TL approach whereas TL approach achieved 10.12% error. Similarly, a decline of  $-108.7$  in  $R^2$  evaluation is detected for the province Queensland. Thus, the proposed approach resulted in a steep decline in performance. For such cases, the LSTM approach (without-TL) has resulted in better forecasting accuracy than the multi-source deep-TL technique.

## 5. Discussion

In this section, the importance of feature selection and the prominent reason behind the formation of the proposed MSDTL approach is described. It consists of three sub-sections, the first sub-section discusses the significance of the feature-set that has been deployed in this paper; the second sub-section provides a comparative analysis with some existing state-of-art algorithms, whereas the last sub-section focuses on the performance evaluation of the proposed approach.

### A. Importance of Features

This section discusses the significance of all the features i.e., {COVID-19 confirmed cases, mortality rate, number of recoveries, population, area, and GDP}, and emphasizes on the reason behind the selection of these parameters using a practical approach. An experimentation has been performed by considering features independently on the proposed MSDTL approach, then forecasting COVID-19 infections, and recording MAPE values. Therefore, four-different experiments were performed, and the results are recorded in Table 5. Initially, we started by considering COVID-19 statistics as a feature-set  $F_1$  i.e., the parameters of interest are.

$$F_1 = \{\text{'confirmed cases'}, \text{'deaths'}, \text{'recoveries'}\}$$

Which are trained on the MSDTL model, future infections are forecasted and model's performance is recorded in terms of MAPE values. Thus, these results are presented in column two of Table 5. Later, the feature 'GDP' is added to the original feature set i.e.,

$$F_2 = \{\text{'confirmed cases'}, \text{'deaths'}, \text{'recoveries'}, \text{'GDP'}\}$$

and training is performed. The column three depicts the MAPE score when executed using the feature set  $F_2$ . Further, another experiment has been performed which consists of feature set  $F_3$  i.e.,

$$F_3 = \{\text{'confirmed cases'}, \text{'deaths'}, \text{'recoveries'}, \text{'population'}, \text{'area'}\}$$

**Table 5**  
Significance of feature selection in proposed approach.

| Province Name                    | COVID        | COVID-GDP | COVID-Population-Area MAPE | ALL PARAMETERS MAPE |
|----------------------------------|--------------|-----------|----------------------------|---------------------|
| Alabama                          | 37.66        | 97.62     | 62.04                      | <b>1.40</b>         |
| Alaska                           | 84.92        | 88.19     | 78.39                      | <b>7.42</b>         |
| Alberta                          | 60.32        | 87.85     | 73.67                      | <b>30.26</b>        |
| Arizona                          | 32.49        | 98.05     | 49.65                      | <b>2.58</b>         |
| Arkansas                         | 49.76        | 96.24     | 67.78                      | <b>2.96</b>         |
| Aruba                            | 1819.15      | 72.86     | 59.87                      | <b>8.62</b>         |
| Australian Capital Territory     | 19.68        | 1.79      | 15.90                      | <b>0.43</b>         |
| Bermuda                          | 19.29        | 40.35     | 20.07                      | <b>1.13</b>         |
| Bonaire, Sint Eustatius and Saba | 63.63        | 10.75     | 20.10                      | <b>14.70</b>        |
| British Columbia                 | 29.32        | 84.53     | 80.05                      | <b>18.56</b>        |
| California                       | 51.70        | 99.46     | 99.99                      | <b>0.58</b>         |
| Cayman Islands                   | 12.14        | 30.92     | 15.48                      | <b>2.48</b>         |
| Channel Islands                  | 36.13        | 50.00     | 30.27                      | <b>4.89</b>         |
| Colorado                         | 42.92        | 96.53     | 102.27                     | <b>50.72</b>        |
| Curacao                          | 107.59       | 86.47     | 45.53                      | <b>55.54</b>        |
| Faroe islands                    | 23.04        | 39.65     | 32.54                      | <b>3.61</b>         |
| Florida                          | 65.05        | 99.36     | 48.34                      | <b>1.71</b>         |
| French Guiana                    | 15.11        | 62.74     | 48.56                      | <b>2.77</b>         |
| French Polynesia                 | 63.56        | 95.11     | 87.49                      | <b>70.54</b>        |
| Georgia                          | 34.91        | 98.64     | 57.60                      | <b>0.29</b>         |
| Gibraltar                        | 38.07        | 71.06     | 46.78                      | <b>5.75</b>         |
| Greenland                        | 2551.63      | 17.34     | 23.17                      | <b>10.56</b>        |
| Guadeloupe                       | 17.73        | 81.29     | 73.73                      | <b>19.78</b>        |
| Hawaii                           | 14.57        | 18.77     | 7.35                       | <b>28.60</b>        |
| Idaho                            | 55.14        | 94.31     | 74.60                      | <b>13.09</b>        |
| Illinois                         | 65.29        | 98.88     | 67.75                      | <b>1.51</b>         |
| Indiana                          | 57.56        | 97.75     | 72.00                      | <b>7.57</b>         |
| Iowa                             | 72.64        | 96.91     | 72.69                      | <b>3.24</b>         |
| Isle of Man                      | 122.60       | 5.46      | 7.31                       | <b>5.11</b>         |
| Kansas                           | 54.11        | 95.77     | 74.98                      | <b>42.12</b>        |
| Kentucky                         | 34.97        | 96.37     | 75.37                      | <b>49.65</b>        |
| Manitoba                         | 135.59       | 91.75     | 90.69                      | <b>53.97</b>        |
| Martinique                       | 40.64        | 82.70     | 70.52                      | <b>50.86</b>        |
| Mayotte                          | 16.90        | 46.44     | 41.15                      | <b>18.96</b>        |
| New Brunswick                    | 27.80        | 46.86     | 26.93                      | <b>25.09</b>        |
| Newfoundland and Labrador        | 9.55         | 10.68     | 15.36                      | <b>12.30</b>        |
| New South Wales                  | <b>5.56</b>  | 31.35     | 26.02                      | 23.28               |
| New York                         | 41.16        | 98.76     | 30.06                      | <b>0.71</b>         |
| Nova Scotia                      | 18.69        | 14.97     | 13.01                      | <b>5.82</b>         |
| Ohio                             | 61.09        | 98.12     | 70.13                      | <b>7.10</b>         |
| Ontario                          | 82.37        | 92.81     | 65.04                      | <b>1.81</b>         |
| Oregon                           | 30.53        | 91.33     | 69.82                      | <b>16.54</b>        |
| Pennsylvania                     | 58.28        | 97.73     | 56.78                      | <b>1.73</b>         |
| Puerto Rico                      | 235.30       | 71.76     | <b>33.97</b>               | 36.49               |
| Quebec                           | 65.77        | 94.40     | 55.04                      | <b>0.62</b>         |
| Queensland                       | 103.77       | 11.66     | 15.23                      | <b>4.06</b>         |
| Reunion                          | 128.22       | 79.13     | 71.69                      | <b>60.09</b>        |
| Saskatchewan                     | 95.79        | 79.52     | 77.20                      | <b>22.98</b>        |
| Sint Maarten                     | <b>15.23</b> | 60.08     | 33.88                      | <b>17.92</b>        |
| South Australia                  | <b>15.38</b> | 19.94     | 20.83                      | <b>16.66</b>        |
| South Carolina                   | 35.01        | 97.02     | 58.15                      | <b>1.38</b>         |
| South Dakota                     | 39.84        | 92.45     | 76.53                      | <b>24.99</b>        |
| St Martin                        | <b>25.85</b> | 40.00     | 48.27                      | 27.31               |
| Texas                            | 57.25        | 99.46     | 58.15                      | <b>1.08</b>         |
| Turks and Caicos Islands         | 17.35        | 42.21     | 36.09                      | <b>22.06</b>        |
| Victoria                         | 199.56       | 76.72     | 53.79                      | <b>21.57</b>        |
| Virginia                         | 25.93        | 96.94     | 56.05                      | <b>3.68</b>         |
| Washington                       | 21.27        | 20.83     | 62.16                      | <b>4.07</b>         |
| West Virginia                    | 37.39        | 30.23     | 78.57                      | <b>17.86</b>        |
| Western Australia                | 17.79        | 20.83     | 19.29                      | <b>8.91</b>         |
| Wisconsin                        | 57.66        | 98.31     | 82.03                      | <b>42.13</b>        |
| Wyoming                          | 184.58       | 88.52     | 85.37                      | <b>44.17</b>        |

and in the similar way, MAPE values are recorded which are depicted in column four of Table 5. Finally, the performance of all the three-different feature-set is compared with F<sub>4</sub> which consists of all these features and that is used entirely in this paper.

$$F_4 = \{ \text{'confirmed cases', 'deaths', 'recoveries', 'population', 'area', 'GDP'} \}$$

The column five of Table 5 presents the performance of MSDTL using feature-set F<sub>4</sub>. Thus, it can be clearly visualized from Table 5 that, when COVID-19 infections are forecasted using F<sub>4</sub> i.e., consisting of all the features, then highly precise model is resulted. When forecasting is achieved using feature set F<sub>1</sub>, F<sub>2</sub> or F<sub>3</sub>, then a biased model is trained and high magnitude of MAPE values are resulted. For every province, performance of the model using the feature-set generating minimal error values is highlighted in bold. Therefore, we can conclude that the features such as-population, area and GDP are alone not useful because they are unable to capture the non-linearity of the datasets', but when all these features are combined, their impact is profound, resulting in magnanimous improvement in prediction precision.

### B. Comparison with other state-of-art algorithms

A comparative analysis has been performed to analyze the performance of the proposed MSDTL approach with some other existing state-of-art algorithms which are proficient in handling sequential data. The first one is the deep learning model i.e., Recurrent Neural Network (RNN), and another one is the VARMAX model, which is a statistical modeling technique.

A recurrent Neural Network is a special type of Artificial Neural Network [52] designed to work for sequential data or time-series models. It works on the principle of memorization of historical information, which helps in storing states of previous inputs while generating the next output of the sequence. RNN is trained with a same number of neurons, trainable weights, activation, optimization, and loss function as that of LSTM to conduct an unbiased comparison. The RNN model is trained in the similar behavior as it is done for LSTM i.e., by incorporating transfer-learning, and the performance is evaluated in terms of MAPE metric which is depicted mathematically in Table 6.

Later, the performance of MSDTL approach has been compared with a statistical VARMAX model. Vector Autoregression Moving Average with Exogenous variable (VARMAX) model [53] is a generalized statistical time-series model. It is a multivariate version of classical univariate ARIMA model. The mathematical representation is as follows:

V: It denotes a Vector, which means that it is composed of two or more features (multivariate parameters).

AR: It denotes Auto Regression term, which means that output of the model is based on linear combinations of inputs. It is characterized by 'p' parameter.

$$X_t = C_0 + C_1X_{t-1} + C_2X_{t-2} + \dots + C_pX_{t-p} \tag{4}$$

MA: It denotes Moving Average term, which is calculated as the difference between observed and the predicted values. In this case, the output of the model is based on linear combinations of residual errors which are being generated at each timestep. It is characterized by 'q' parameter.

$$X_t = C_0 + C_1\varepsilon_{t-1} + C_2\varepsilon_{t-2} + \dots + C_q\varepsilon_{t-q} \tag{5}$$

X: It denotes exogenous variable, which are used to model the primary variables which are of more interest. It is a type of variable whose value is determined outside the model and is imposed here.

These variables are not directly modelled but used in terms of weighted input to the model. Whereas there is other type of variables known as endogenous variable, which is determined by the computation of the model. Thus, in this paper, the parameters: 'Population', 'Area' and 'GDP' are considered as exogenous variables, since their value is independent of time. Their value remains same for a particular province. However, the parameters: 'Confirmed cases', 'Deaths', and 'Recoveries' are considered as endogenous variables, since their values are changing constantly and are a prime source of interest. Thus, the multisource domain dataset has been collected and trained using transfer learning. The performance has been estimated using MAPE as an evaluation

**Table 6**  
Comparative analysis of proposed approach with state-of-art algorithms.

| Province Name                    | RNN         | VARMAX       | MSDTL        |
|----------------------------------|-------------|--------------|--------------|
| Alabama                          | 99.74       | 87.48        | <b>1.40</b>  |
| Alaska                           | 98.74       | 92.27        | <b>7.30</b>  |
| Alberta                          | 99.14       | 89.74        | <b>30.26</b> |
| Arizona                          | 99.78       | 47.28        | <b>2.58</b>  |
| Arkansas                         | 99.63       | 82.23        | <b>2.96</b>  |
| Aruba                            | 93.50       | 4957.97      | <b>8.62</b>  |
| Australian Capital Territory     | 160.80      | 20.71        | <b>0.43</b>  |
| Bermuda                          | 30.15       | 72.59        | <b>1.13</b>  |
| Bonaire, Sint Eustatius and Saba | 85.78       | <b>11.73</b> | 14.70        |
| British Columbia                 | 98.45       | 93.07        | <b>18.56</b> |
| California                       | 99.93       | 1000.67      | <b>1.00</b>  |
| Cayman Islands                   | 12.74       | 40.03        | <b>2.48</b>  |
| Channel Islands                  | 74.77       | 79.41        | <b>4.89</b>  |
| Colorado                         | 99.69       | 482.94       | <b>50.72</b> |
| Curacao                          | 82.54       | 94.69        | <b>55.54</b> |
| Faroe islands                    | 40.07       | 32.04        | <b>0.57</b>  |
| Florida                          | 99.92       | 34.20        | <b>3.58</b>  |
| French Guiana                    | 96.89       | 22.72        | <b>2.77</b>  |
| French Polynesia                 | 97.50       | 95.54        | <b>70.54</b> |
| Georgia                          | 99.84       | 36.21        | <b>1.65</b>  |
| Gibraltar                        | 67.16       | 80.82        | <b>5.75</b>  |
| Greenland                        | 1569.28     | 41.13        | <b>8.37</b>  |
| Guadeloupe                       | 96.18       | 85.00        | <b>19.78</b> |
| Hawaii                           | 60.70       | 38.20        | <b>28.60</b> |
| Idaho                            | 99.51       | 87.46        | <b>13.09</b> |
| Illinois                         | 99.88       | 83.95        | <b>8.37</b>  |
| Indiana                          | 99.78       | 87.87        | <b>1.50</b>  |
| Iowa                             | 99.72       | 86.48        | <b>3.24</b>  |
| Isle of Man                      | 17.14       | 27.77        | <b>5.11</b>  |
| Kansas                           | 99.63       | 88.21        | <b>42.12</b> |
| Kentucky                         | 99.67       | 88.70        | <b>49.65</b> |
| Manitoba                         | 97.37       | 96.50        | <b>4.21</b>  |
| Martinique                       | 93.54       | 90.99        | <b>50.86</b> |
| Mayotte                          | 93.60       | 67.45        | <b>18.96</b> |
| New Brunswick                    | 30.02       | 50.47        | <b>25.09</b> |
| Newfoundland and Labrador        | <b>8.06</b> | 34.81        | 12.30        |
| New South Wales                  | 92.60       | 22.64        | <b>1.20</b>  |
| New York                         | 99.86       | 40.49        | <b>0.71</b>  |
| Nova Scotia                      | 75.11       | 37.78        | <b>5.82</b>  |
| Ohio                             | 99.81       | 88.45        | <b>7.10</b>  |
| Ontario                          | 99.45       | 83.53        | <b>1.81</b>  |
| Oregon                           | 99.34       | 85.88        | <b>16.54</b> |
| Pennsylvania                     | 99.78       | 84.75        | <b>1.72</b>  |
| Puerto Rico                      | 95.09       | 76.20        | <b>36.49</b> |
| Quebec                           | 99.54       | 74.58        | <b>0.62</b>  |
| Queensland                       | 74.02       | 22.38        | <b>10.16</b> |
| Reunion                          | 95.68       | 78.11        | <b>60.09</b> |
| Saskatchewan                     | 94.78       | 92.38        | <b>22.98</b> |
| Sint Maarten                     | 70.52       | 30.17        | <b>17.92</b> |
| South Australia                  | 43.18       | 34.67        | <b>16.66</b> |
| South Carolina                   | 99.71       | 83.41        | <b>1.38</b>  |
| South Dakota                     | 99.43       | 89.90        | <b>24.99</b> |
| St Martin                        | 55.51       | 84.72        | <b>27.31</b> |
| Texas                            | 99.94       | 75.54        | <b>1.07</b>  |
| Turks and Caicos Islands         | 59.06       | 35.08        | <b>22.06</b> |
| Victoria                         | 98.11       | 48.35        | <b>21.57</b> |
| Virginia                         | 99.71       | 39.71        | <b>3.68</b>  |
| Washington                       | 99.61       | 81.33        | <b>4.07</b>  |
| West Virginia                    | 99.08       | 91.43        | <b>17.86</b> |
| Western Australia                | 61.96       | 31.618       | <b>8.91</b>  |
| Wisconsin                        | 99.84       | 89.88        | <b>42.13</b> |
| Wyoming                          | 98.68       | 93.91        | <b>4.17</b>  |

metric which is represented in tabular format in [Table 6](#).

[Table 6](#) presents a comparative analysis of the proposed MSDTL approach with other existing state-of-art algorithms in terms of MAPE evaluation metrics. It consists of four columns: column one presents the name of all provinces which are considered for this experiment, column two depicts the empirical results of the MAPE metric on the RNN approach, column three represents the evidence for VARMAX model, and finally, the last column describes the performance of the proposed MSDTL approach in terms of MAPE measure.

For every province, the best resulting model i.e., the model

generating minimum error value is highlighted in bold. Thus, it can be clearly visualized that the proposed MSDTL approach is providing very accurate infection forecasts and resulting in minimal error for 60 out of 62 provinces. Conversely, it can be analyzed that the RNN model is resulting in very inaccurate forecasting of COVID-19 infection, and a high magnitude of MAPE values is obtained. However, RNN suffers from a vanishing gradient issue and the model is unable to learn relevant information of input states for a longer time, thus the proposed approach involves the usage of an enhanced version of this algorithm i.e., LSTM which provides the characteristics of learning long-term dependencies with the help of three different gates. Similar behavior is observed for a statistical model that, being a time-series model still VARMAX model is not able to perform well on the sequential dataset. Thus, it resulted in very poor predictions and generated large MAPE score values. Therefore, based on the above study, the proposed MSDTL approach is utilized for efficient forecasting of COVID-19 infections, and it is further recommended for future use.

### C. Performance Evaluation of the proposed approach

The proposed approach has been evaluated with the COVID-19 data sets from sixty-two provinces from a wide range of countries around the world. The empirical evidence shows that an improvement of up to 96% is attained using the proposed-TL approach. The maximum similarity up to 0.84 has been observed using the R<sup>2</sup>-test as discussed in detail in section 4 of this paper. However, degradation in the performance of the proposed approach is observed for very few provinces. [Table 7](#) presents the rate of infectious spread and its relationship with the performance of the proposed methodology. It consists of four columns, the first column presents the province name, the second one calculates the gradient of historical infections, the third column shows the improvement in MAPE score, and the fourth column shows the improvement of R<sup>2</sup>. The columns 3 and 4 are respectively taken from [Tables 3 and 4](#) for better analysis. The infections gradient may be mathematically calculated as follows:

$$\text{Gradient} = \text{Increase in historical infections in last 30 days}/30 \quad (6)$$

For the province of California, Alabama, Florida, Manitoba, Quebec, South Carolina, Virginia and Wyoming, the magnitude of gradient is 3304.83,1144.46, 2714.96, 82.3, 964.24,853.8, 937.7, and 194.8 respectively. A higher gradient signifies that rate of infections per day is

**Table 7**  
Performance evaluation of the proposed approach.

| Province Name            | Rate of infectious spread | Improvement in MAPE score (%) | Improvement in R <sup>2</sup> |
|--------------------------|---------------------------|-------------------------------|-------------------------------|
| Alabama                  | 1144.46                   | 76.2                          | 18.3                          |
| Alaska                   | 197.23                    | 77.6                          | 77.12                         |
| Arizona                  | 711.23                    | 93.01                         | 1378.42                       |
| Bermuda                  | 2.43                      | 81.35                         | 12                            |
| California               | 3304.83                   | 72.67                         | 20.26                         |
| Faroe Islands            | 2.1                       | 74.66                         | 19.9                          |
| Florida                  | 2714.96                   | 96.39                         | 12023.36                      |
| Georgia                  | 1239.86                   | 95.44                         | 2177.84                       |
| Illinois                 | 2952.3                    | 20.28                         | 17.94                         |
| Indiana                  | 1571.03                   | 97.73                         | 468.78                        |
| Manitoba                 | 82.3                      | 21.60                         | 28.08                         |
| New South Wales          | 44.66                     | 21.56                         | 60                            |
| Pennsylvania             | 1350.13                   | 53.51                         | 21.05                         |
| Quebec                   | 964.26                    | 84.91                         | 13.6                          |
| South Carolina           | 853.8                     | 93.86                         | 1326.19                       |
| Texas                    | 4797.26                   | 92.24                         | 43.31                         |
| Virginia                 | 937.7                     | 83.93                         | 309.21                        |
| Wyoming                  | 194.8                     | 69.71                         | 65.54                         |
| Greenland                | 0.1                       | -72.22                        | -0.86                         |
| Martinique               | 0.9                       | -4.69                         | -1.78                         |
| Queensland               | 0.4                       | -152                          | -108.7                        |
| Turks and Caicos Islands | 0.46                      | -115                          | -34.01                        |

very high i.e., infections are spreading at a very faster pace each day. A similar trend is followed for all 58 provinces. However, there are few provinces where an increase in the infections (gradient) is even less than 1, i.e., infections are not even increasing linearly. For Example, Greenland's Martinique, Queensland, and Turks and Caicos Island observed gradients of 0.1, 0.9, 0.4, and 0.46, respectively.

From Table 7, it can be seen that if the magnitude of the gradient is significantly higher for a province, then the proposed approach is quite efficient in forecasting COVID-19 Infections. For Alabama, the magnitude of the gradient is 1144.46, which is comparatively very large. So, the improvement in MAPE score and  $R^2$  values are also better (with improvements of 76.2% and 18.3, respectively) than the without-TL approach. Similarly, for Manitoba, the magnitude of the gradient is 82.3. Thus, the prediction preciseness of the proposed approach is also effectual i.e., improvement in MAPE score and  $R^2$  value of 21.60% and 28.08 respectively as compared to without-TL. This trend is observed for as many as 58 provinces (out of a total of 62 provinces), which confirms the increased efficacy of the proposed model. The increase of infectious spread in the past 30 days is minimal for 4 provinces out of a total of 62 provinces i.e., the magnitude of the gradient of infectious spread in the past 30 days is very small. In such cases, enhancement of the prediction preciseness from the proposed approach is absent as compared to the without-TL approach. For Example, the gradient for Greenland is 0.1 (which is very minuscule).

So, the proposed approach has resulted in poor performance i.e., deterioration in the MAPE score and  $R^2$  values ( $-72.22\%$  and  $-0.86$ , respectively) as compared to without-TL. Similarly, for Martinique, the small gradient of 0.9, deteriorated the performance of MAPE score  $R^2$  values by  $-4.69\%$  and  $-1.78$ , respectively. Similar behavior is observed for Queensland and Turks and Caicos Islands, as visualized in Table 7. For these provinces, the proposed approach has been found as not so suitable.

Therefore, our major observations on the performance of the proposed multi-source deep-TL approach are as follows:

- a) If historical infections of 30 days are variably increasing with respect to time i.e., a stationary pattern (gradient of infectious spread) is not observed, then our proposed model will result in efficient forecasting of COVID-19 infection. It can be visualized from Table 7, that when the gradient value is higher in magnitude, the proposed-TL approach is providing better efficacy of the model. For example, the gradient of infectious spread is higher for Florida and New South Wales (2714.96 and 44.66, respectively), thus the proposed approach attained precise predictions. Hence, it can be visualized from Fig. 4 (g), and Fig 4(l) for Florida and New South Wales in which the without-TL approach has resulted in imprecise predictions (i.e., very high MAPE score 99.31, 1.53 respectively) whereas with-TL approach in Fig. 5 (g), Fig 5 (l) has provided more precise forecasting (better MAPE score: 3.58 and 1.20, respectively). Such trends are also observed for 58 provinces from a total of 62 provinces.
- b) The proposed model of multi-source deep-TL methodology has provided more stable outcomes with a lower rate of fluctuations in future infection forecasting in comparison with the without-TL approach. For example, it can be visualized from Fig. 4(b) and 4(h) for Alaska and Georgia in which the without-TL approach depicts very large fluctuations (i.e., very poor  $R^2$  values: 81.21 and  $-2177$ , respectively) whereas the with-TL approach (Fig. 5(b) and (h)) has resulted in more stable outcomes (i.e., better  $R^2$  values: 4.09 and 0.84, respectively). Since the infection gradient for Alaska and Georgia is higher in magnitude (197.23 and 1239.86, respectively), the performance of the proposed approach has resulted in more stable and accurate results. Such trends of stable forecasts (better  $R^2$  values) are also observed for as many as 58 provinces out of a total of 62 provinces. Thus, the proposed technique has provided more steady forecasts of the expected number of daily cases per province in the next 6 days.

- c) If gradient of infectious spread is very low, it means certain level of stationarity is observed in past infections. That is, if increase in the confirmed cases witnesses little variation over the time, then the proposed approach may not result in efficient performance. For example, it can be visualized from Fig. 6(b) and (f) for the province Greenland and Queensland etc., in which with-TL resulted in poor performance. That is, deterioration of MAPE score and  $R^2$  values  $-72.22\%$  and  $-0.86$  for Greenland and  $-152\%$  and  $-108.7$  for Queensland, respectively.
- d) When the gradient of infectious spread is very minimal, the proposed model has resulted in high fluctuating outcomes (less robust) in comparison to without-TL approach. For example, the proposed approach depicts very large fluctuations for Martinique, and Turks and Caicos Island (see Figs. 6(d), Fig 4(h)) whereas the without-TL approach (Fig. 6(c) and (g)) has resulted in more stable outcomes (deterioration of MAPE score and  $R^2$  values by  $-4.69\%$  and  $-1.78$  for Martinique, and  $-115\%$  and  $-34.01$  for Turks and Caicos Island, respectively). Since the infection gradient for Martinique and Turks and Caicos Island is low in magnitude (0.9 and 0.46, respectively), the performance of the proposed approach degrades.

Thus, it can be concluded as whenever the gradient of infection spread is variably increasing, then the proposed predictive modeling results in more stable and precise forecasting than the without-TL approach. However, when the gradient is very low i.e., a certain level of stationarity is observed in past infections, then the without-TL approach performs better. In such a scenario, the proposed approach results in high fluctuations and imprecise forecasting.

## 6. Conclusion and future works

This paper has proposed a novel approach for COVID-19 infections forecasting using a multi-source deep transfer learning methodology. The proposed approach is also novel as for the first time it considers parameters such as population density and economic conditions for COVID-19 efficient forecasting since COVID-19 is a contagious disease. It has been proved experimentally, by investigating different features independently and recording their performance. Later, when all the features are combined, better performance resulted. Apart from these novel parameters, the COVID-19 infection spread in published literature is found to be dependent on the rate of infections, fatality, and recovery over a period. This time-dependent task can effectively be solved using the well-known LSTM model. However, a comparative analysis has been performed with other existing state-of-art algorithms such as RNN and VARMAX, to validate the performance of LSTM over them. The availability of the COVID-19 province-specific dataset is very little, thus LSTM suffers from low variance and high bias. So, to tide over these limitations of LSTM, the proposed approach has utilized transfer learning to enhance the LSTM accuracy in COVID-19 forecasting by employing a multi-source domain dataset for initial training.

The performance of the proposed approach has been evaluated using two different statistical measures, which are: MAPE and coefficient of determination. The empirical evidences show that, a maximum improvement up to 96% (with the proposed TL approach) with the average accuracy ranging from 65 to 85% have been observed in terms of MAPE score. The maximum similarity up to 0.84 (with the proposed TL approach) has been observed between actual infections and forecasted infectious in terms of R-squared test, whereas the without-TL LSTM model resulted in very poor similarity. The proposed approach provides more stable and precise forecasting, whenever the gradient is high, as compared to without-TL. However, when stationarity is observed, then without-TL approach performs better than the proposed approach. Thus, keeping in mind the strength and weakness of the proposed approach, a stable and more accurate transmission forecast is a crucial requirement for administrators for analyzing the current scenario to deal with the threats of COVID-19.

Towards the end, we may firmly infer that novel technique for COVID-19 infection forecasting, using multisource deep transfer learning is extremely effective. Henceforth, we would improve our approach to forecast the COVID-19 infection evolution of developing nations where coverage of COVID-19 testing is inadequate. We will incorporate more parameters which are related to infectious spread, to provide better infection forecasting in future. Also, the proposed approach of MSDTL using LSTM will be further evolved so that its prediction precision is enhanced including the cases of stationarity.

The proposed approach can be extended in modelling the prediction of any contagious disease precisely, since future infections will be dependent on density, population, and time. The proposed approach may also target many other real-world applications with data scarcity and data insufficiency. The proposed approach can be applied for the prediction of stock prices for a company that is recently listed in a country. Since the company is a recent entrant to that country, the data is insufficient to precisely forecast its stock prices. Thus, in this case, multi-source deep transfer learning can be employed and related data from other countries having similar economic conditions can be incorporated, and accurate forecasts of these stocks can be modelled. Similarly, the proposed approach can also be applied in weather forecasting if some unusual weather phenomenon is observed. In such cases, data pertaining to this unusual phenomenon will be rare. But this limitation can be optimally tackled by the proposed multi-source deep TL, which applies deep TL on related multi source datasets for precise forecasting of such unusual weather phenomena. For example, recently in some parts of the world, sudden heat waves were observed, and at some places, the recorded ground temperature was nearly 50° Celsius, which was very unusual. Thus, in such scenarios, multi-source deep transfer learning can be employed and similar data from different parts of the regions/world can be captured which may help in appropriate weather forecasting. Such a study may therefore aim for better prediction by employing multi-source deep transfer learning.

#### Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

#### Acknowledgments

Authors are thankful to the editors and the reviewers for all their valuable comments which have helped a lot in improving the work. First author gratefully acknowledges the financial assistance obtained from South Asian University (SAU), New Delhi, India in the form of a master's scholarship. Second author is grateful to Department of Science and Technology, Government of India, New Delhi for the financial support in the form of INSPIRE Fellowship for his PhD research. All authors are also thankful to the SAU, New Delhi for providing the infrastructural facilities to conduct this research through the Computational Intelligence research lab.

#### References

- <https://www.who.int/health-topics/severe-acute-respiratory-syndrome-> (Accessed January 2021).
- <https://covid19.who.int/>. (Accessed January 2021).
- H. Swapnarekha, Himansu Sekhar Behera, Janmenjoy Nayak, Bighnaraj Naik, Role of intelligent computing in COVID-19 prognosis: a state-of-the-art review, *Chaos, Solit. Fractals* 138 (2020), 109947.
- Maryam M. Najafabadi, Flavio Villanustre, M. Taghi, Khoshgoftaar, naeem seliya, randall wald, and edin muharemagic. "Deep learning applications and challenges in big data analytics.", *Journal of big data* 2 (1) (2015) 1–21.
- Sinno Jialin Pan, Qiang Yang, A survey on transfer learning, *IEEE Trans. Knowl. Data Eng.* 22 (10) (2009) 1345–1359.
- Nadajat Kadi, Mounia Khelfaoui, Population density, a factor in the spread of COVID-19 in Algeria: statistic study, *Bull. Natl. Res. Cent.* 44 (1) (2020) 1–7.
- [https://www.who.int/water\\_sanitation\\_health/emergencies/qa/emergencie\\_s\\_qa9/en/](https://www.who.int/water_sanitation_health/emergencies/qa/emergencie_s_qa9/en/).
- Charlie H. Zhang, Gary G. Schwartz, Spatial disparities in coronavirus incidence and mortality in the United States: an ecological analysis as of May 2020, *J. Rural Health* 36 (3) (2020) 433–445.
- Patrick GT. Walker, Charles Whittaker, Oliver J. Watson, Marc Baguelin, Winkill Peter, Arran Hamlet, A. Bimandra, Djafaara, et al., The impact of COVID-19 and strategies for mitigation and suppression in low-and middle-income countries, *Science* 369 (6502) 413–422.
- Domenico Benvenuto, Marta Giovanetti, Lazzaro Vassallo, Silvia Angeletti, Massimo Ciccozzi, Application of the ARIMA model on the COVID-2019 epidemic dataset, *Data Brief* 29 (2020), 105340.
- Trishan Panch, Heather Mattie, Leo Anthony Celi, The "inconvenient truth" about AI in healthcare, *NPJ Digit. Med.* 2 (1) (2019) 1–3.
- Knight, M. Gwenan, J. Nila, Dharan, J.Fox Gregory, Stennis Natalie, Alice Zwerling, Renuka Khurana, David W. Dowdy, Bridging the gap between evidence and policy for infectious diseases: how models can aid public health decision-making, *Int. J. Infect. Dis.* 42 (2016) 17–23.
- Alex Sherstinsky, Fundamentals of recurrent neural network (RNN) and long short-term memory (LSTM) network, *Phys. Nonlinear Phenom.* 404 (2020), 132306.
- K. Singh, S. Shastri, A.S. Bhadwal, P. Kour, et al., Implementation of exponential smoothing for forecasting time series data, *Int. J. Sci. Res. Comput. Sci. Appl. Manag. Stud.* 8 (1) (2019) 1–8.
- Z. Zhao, K. Nehil-Puleoa, Y. Zhao, How well can we forecast the COVID-19 pandemic with curve fitting and recurrent neural networks? medRxiv preprint (2020) <https://doi.org/10.1101/2020.05.14.20102541>.
- Jerome T. Connor, R. Douglas Martin, Les E. Atlas, Recurrent neural networks and robust time series prediction, *IEEE Trans. Neural Network.* 5 (2) (1994) 240–254.
- Sepp Hochreiter, Jürgen Schmidhuber, Long short-term memory, *Neural Comput.* 9 (8) (1997) 1735–1780.
- David W. Berger, Kyle F. Herkenhoff, Mongey Simon, An Seir Infectious Disease Model with Testing and Conditional Quarantine. No. W26901, National Bureau of Economic Research, 2020.
- Alberto Godio, Francesca Pace, Andrea Vergnano, SEIR modeling of the Italian epidemic of SARS-CoV-2 using computational swarm intelligence, *Int. J. Environ. Res. Publ. Health* 17 (10) (2020) 3535.
- Ghanbari, A., R. Khordad, and Mostafa Ghaderi-Zefrehei. "Mathematical prediction of the spreading rate of COVID-19 using entropy-based thermodynamic model." *Indian J. Phys.*: 1-7.
- Lixiang Li, Zihang Yang, Zhongkai Dang, Meng Cui, Jingze Huang, Haotian Meng, Deyu Wang, et al., Propagation analysis and prediction of the COVID-19, *Infect. Dis. Model.* 5 (2020) 282–292.
- Shi Zhao, Qianyin Lin, Jinjun Ran, Salihu S. Musa, Guangpu Yang, Weiming Wang, Yijun Lou, et al., Preliminary estimation of the basic reproduction number of novel coronavirus (2019-nCoV) in China, from 2019 to 2020: a data-driven analysis in the early phase of the outbreak, *Int. J. Infect. Dis.* 92 (2020) 214–217.
- Abir, Farhan Fuad, Khalid Alyafei, Muhammad EH. Chowdhury, Amith Khandaker, Rashid Ahmed, Md Shafayet Hossain, Sakib Mahmud, et al., PCovNet: a presymptomatic COVID-19 detection framework using deep learning model using wearables data, *Comput. Biol. Med.* (2022), 105682.
- Trishan Panch, Heather Mattie, Leo Anthony Celi, The "inconvenient truth" about AI in healthcare, *NPJ Digit. Med.* 2 (1) (2019) 1–3.
- Dianbo Liu, Leonardo Clemente, Canelle Poirier, Xiyu Ding, Matteo Chinazzi, Jessica T. Davis, Alessandro Vespignani, Mauricio Santillana, A machine learning methodology for real-time forecasting of the 2019-2020 COVID-19 outbreak using Internet searches, news alerts, and estimates from mechanistic models, arXiv preprint arXiv:2004 (2020), 04019.
- Yao Zhang, Heng Xue, Mixue Wang, He Nan, Zhibin Lv, Ligang Cui, Lung ultrasound findings in patients with coronavirus disease (COVID-19), *Am. J. Roentgenol.* 216 (1) (2021) 80–84.
- Babady, N. Esther, Tracy McMillen, Krupa Jani, Agnes Viale, Elizabeth V. Robilotti, Anoshe Aslam, Maureen Diver, et al., Performance of severe acute respiratory syndrome coronavirus 2 real-time RT-PCR tests on oral rinses and saliva samples, *J. Mol. Diagn.* 23 (1) (2021) 3–9.
- Ali Narin, Ceren Kaya, Ziyet Pamuk, Automatic detection of coronavirus disease (covid-19) using x-ray images and deep convolutional neural networks, arXiv preprint arXiv:2003 (2020), 10849.
- Lei Qin, Qiang Sun, Yidan Wang, Ke-Fei Wu, Mingchih Chen, Ben-Chang Shia, Szu-Yuan Wu, Prediction of number of cases of 2019 novel coronavirus (COVID-19) using social media search index, *Int. J. Environ. Res. Publ. Health* 17 (7) (2020) 2365.
- Mei-Ling Huang, Yu-Chieh Liao, A lightweight CNN-based network on COVID-19 detection using X-ray and CT images, *Comput. Biol. Med.* (2022), 105604.
- Ahmad, Muhammad, Saima Sadiq, Ala Saleh Alluhaidan, Muhammad Umer, Saleem Ullah, and Michele Nappi. "Industry 4.0 technologies and their applications in fighting COVID-19 pandemic using deep learning techniques." *Comput. Biol. Med.* 145 (2022): 105418.
- K.E. Arun Kumar, Dinesh V. Kalaga, Ch Mohan Sai Kumar, Govinda Chilkoor, Masahiro Kawaji, M. Timothy, Brenza, Forecasting the dynamics of cumulative COVID-19 cases (confirmed, recovered and deaths) for top-16 countries using statistical machine learning models: auto-Regressive Integrated Moving Average (ARIMA) and Seasonal Auto-Regressive Integrated Moving Average (SARIMA), *Appl. Soft Comput.* 103 (2021), 107161.
- Andres Hernandez-Matamoros, Hamido Fujita, Toshitaka Hayashi, Hector Perez-Meana, Forecasting of COVID19 per regions using ARIMA models and polynomial functions, *Appl. Soft Comput.* 96 (2020), 106610.



- [34] Chiou-Jye Huang, Yung-Hsiang Chen, Yuxuan Ma, Ping-Huan Kuo, Multiple-input deep convolutional neural network model for covid-19 forecasting in China, medRxiv (2020).
- [35] Parul Arora, Himanshu Kumar, Bijaya Ketan Panigrahi, Prediction and analysis of COVID-19 positive cases using deep learning models: a descriptive case study of India, *Chaos, Solit. Fractals* 139 (2020), 110017.
- [36] Chimmula, Vinay Kumar Reddy, Lei Zhang, Time series forecasting of COVID-19 transmission in Canada using LSTM networks, *Chaos, Solit. Fractals* 135 (2020), 109864.
- [37] Sandeep Kumar, Pranab K. Muhuri, A novel GDP prediction technique based on transfer learning using CO2 emission dataset, *Appl. Energy* 253 (2019), 113476.
- [38] Mohamed Loey, Gunasekaran Manogaran, Mohamed Hamed N. Taha, Nour Eldeen M. Khalifa, A hybrid deep transfer learning model with machine learning methods for face mask detection in the era of the COVID-19 pandemic, *Measurement* 167 (2020), 108288.
- [39] Yogesh Gautam, Transfer Learning for COVID-19 Cases and Deaths Forecast Using LSTM Network, *ISA transactions*, 2021.
- [40] Hang Su, Zhao Dong, Hela Elmannaï, , Ali Asghar Heidari, Sami Bourouis, Zongda Wu, Zhennao Cai, Wenyong Gui, Mayun Chen, Multilevel threshold image segmentation for COVID-19 chest radiography: a framework using horizontal and vertical multiverse optimization, *Comput. Biol. Med.* (2022), 105618.
- [41] Armando Carteni, Luigi Di Francesco, Maria Martino, The role of transport accessibility within the spread of the Coronavirus pandemic in Italy, *Saf. Sci.* 133 (2021), 104999.
- [42] Yu Wu, Wenzhan Jing, Jue Liu, Qiuyue Ma, Jie Yuan, Yaping Wang, Min Du, Min Liu, Effects of temperature and humidity on the daily new cases and new deaths of COVID-19 in 166 countries, *Sci. Total Environ.* 729 (2020), 139051.
- [43] <https://coronavirus.jhu.edu/map.html>-. (Accessed January 2021).
- [44] <https://ourworldindata.org/grapher/tests-of-covid-19-per-thousand-people-vs-gdp-per-capita> -. (Accessed January 2021).
- [45] <https://images.app.goo.gl/MGqc2kEwvpv4ojNaZA->. (Accessed January 2021).
- [46] Aytac Altan, Seçkin Karasu, Recognition of COVID-19 disease from X-ray images by hybrid model consisting of 2D curvelet transform, chaotic salp swarm algorithm and deep learning technique, *Chaos, Solit. Fractals* 140 (2020), 110071.
- [47] Seçkin Karasu, Aytac Altan, Stelios Bekiros, Wasim Ahmad, A new forecasting model with wrapper-based feature selection approach using multi-objective optimization technique for chaotic crude oil time series, *Energy* 212 (2020), 118750.
- [48] Aytac Altan, Seçkin Karasu, Enrico Zio, A new hybrid model for wind speed forecasting combining long short-term memory neural network, decomposition methods and grey wolf optimizer, *Appl. Soft Comput.* 100 (2021), 106996.
- [49] İsmail Kirbaş, Adnan Sözen, Azim Doğuş Tuncer, Fikret Şınasi Kazancıoğlu, Comparative analysis and forecasting of COVID-19 cases in various European countries with ARIMA, NARNN and LSTM approaches, *Chaos, Solit. Fractals* 138 (2020), 110015.
- [50] Peipei Wang, Xinqi Zheng, Gang Ai, Dongya Liu, Bangren Zhu, Time series prediction for the epidemic trends of COVID-19 using the improved LSTM deep learning method: case studies in Russia, Peru and Iran, *Chaos, Solit. Fractals* 140 (2020), 110214.
- [51] Michael T. Rosenstein, Zvika Marx, Leslie Pack Kaelbling, Thomas G. Dietterich, To transfer or not to transfer, in: *NIPS 2005 Workshop on Transfer Learning*, vol. 898, 2005, pp. 1–4.
- [52] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, Yoshua Bengio, Learning phrase representations using RNN encoder-decoder for statistical machine translation, arXiv preprint arXiv:1406 (2014) 1078.
- [53] G. Peter Zhang, Time series forecasting using a hybrid ARIMA and neural network model, *Neurocomputing* 50 (2003) 159–175.
- [54] William Ogilvy Kermack, Anderson G. McKendrick, A contribution to the mathematical theory of epidemics, *Proc. R. Soc. Lond. - Ser. A. Contain. Pap. a Math. Phys. Character* 115 (772) (1927) 700–721.
- [55] Christopher Y. Shen, Logistic growth modelling of COVID-19 proliferation in China and its international implications, *Int. J. Infect. Dis.* 96 (2020) 582–589.
- [56] Xiaolei Zhang, Renjun Ma, Lin Wang, Predicting turning point, duration and attack rate of COVID-19 outbreaks in major Western countries, *Chaos, Solit. Fractals* 135 (2020), 109829.
- [57] Jin Wang, Mathematical models for COVID-19: applications, limitations, and potentials, *J. Publ. Health Epidemiol.* 4 (2020).
- [58] K.E. ArunKumar, Dinesh V. Kalaga, Ch Mohan Sai Kumar, Masahiro Kawaji, Timothy M. Brenza, Comparative analysis of Gated Recurrent Units (GRU), long Short-Term memory (LSTM) cells, autoregressive Integrated moving average (ARIMA), seasonal autoregressive Integrated moving average (SARIMA) for forecasting COVID-19 trends, *Alex. Eng. J.* 61 (10) (2022) 7585–7603.
- [59] Sweeti Sah, B. Surendiran, R. Dhanalakshmi, Sachi nandan mohanty, fayadh alenezi, and kemal polat. "Forecasting COVID-19 pandemic using prophet, ARIMA, and hybrid stacked LSTM-GRU models in India, *Comput. Math. Methods Med.* (2022) 2022.
- [60] Aishwarya Jakka, Forecasting COVID-19 cases in India using machine learning models, in: *2020 International Conference on Smart Technologies in Computing, Electrical and Electronics (ICSTCEE)*, IEEE, 2020, pp. 466–471.
- [61] Li, Cheng-Fan, Yi-Duo Xu, Xue-Hai Ding, Jun-Juan Zhao, Rui-Qi Du, Li-Zhong Wu, and Wen-Ping Sun. "MultiR-net: a novel joint learning network for COVID-19 segmentation and classification." *Comput. Biol. Med.* 144 (2022): 105340.
- [62] Jin, Weiqiu, Shuqing Dong, Chengqing Yu, and Qingquan Luo. "A data-driven hybrid ensemble AI model for COVID-19 infection forecast using multiple neural networks and reinforced learning." *Comput. Biol. Med.* 146 (2022): 105560.