# AN LSTM APPROACH TOWARDS DIABETES DETECTION USING MACHINE LEARNING

Srinivas Kalluri, Dr.D.V.Divakara Rao, Dr.K.V.Satyanarayana

MTech Student of Raghu Engineering College, Associate Professor in Raghu Engineering College
Department of Computer Science and Engineering
Raghu Engineering College, Visakhapatnam, India

**Abstract:** Diabetes Mellitus (DML) is a condition incited by unregulated diabetes that may prompt multi-organ dysfunctionality in patients. On interpretation of developments in artificial intelligence and manufactured consciousness, which empowers the early identification and analysis of diabetes through a machine level interaction which is more worthwhile than a manual conclusion. As of now, numerous articles are distributed on programmed diabetes location, finding, and self-administration by means of artificial intelligence and man-made reasoning methods. This paper conveys an examination of the recognition, finding, and self-administration strategies of diabetes from six distinct aspects viz., datasets of DML pre-preparing techniques, include extraction techniques, artificial intelligence based distinguishing proof, grouping, and determination of DML, machine-made intelligence based canny diabetes collaborator and execution measures. It likewise talks about the finishes of the past investigation and the significance of the consequences of the examination. Additionally, three momentum research issues in the field of diabetic condition discovery and analysis and self-administration and personalization are recorded. This paper gives an itemized outline of diabetes detection and self-administration strategies which demonstrate important to the monitor patient's condition.

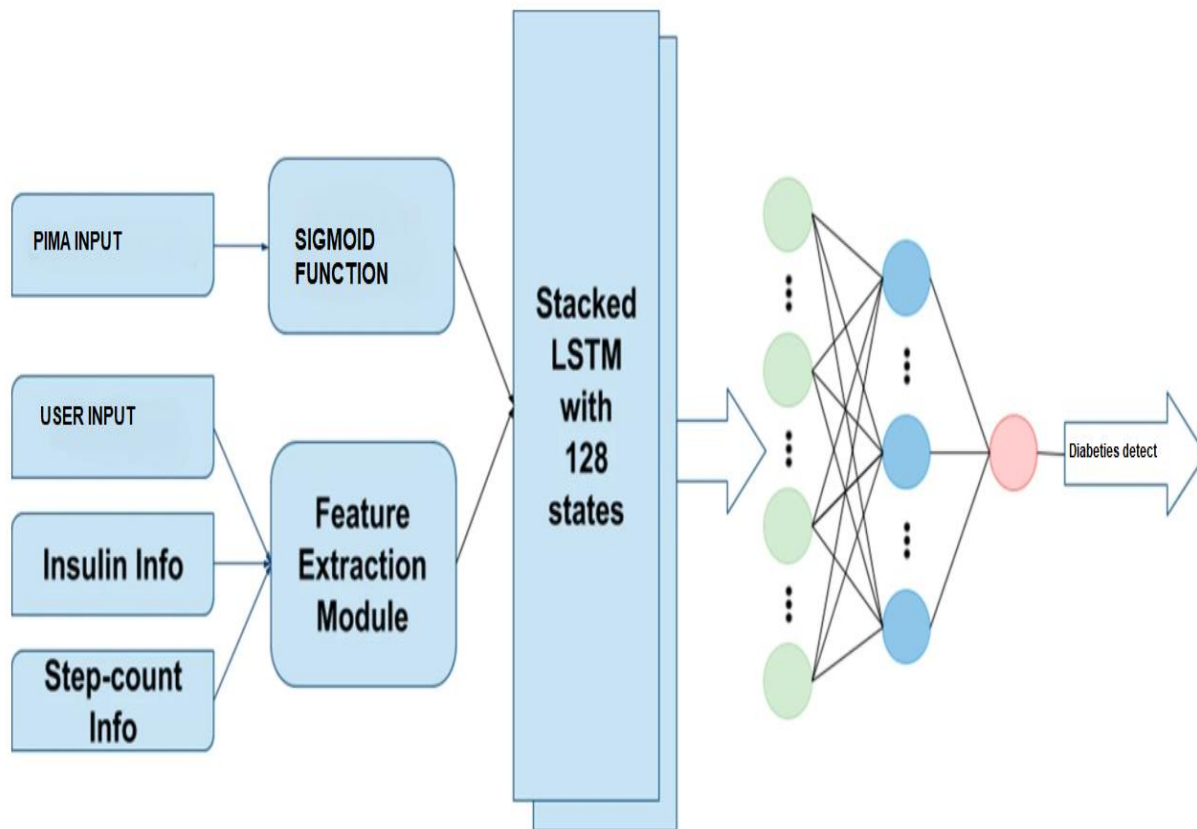*Index Terms* - **Diabetes Mellitus, Machine learning, Neural networks.**

## I. INTRODUCTION

Universally, an expected 422 million of population are experiencing diabetes mellitus (DML), as indicated by the Health Report on Diabetes. Be that as it may, an expected 30–80 percent of diabetes cases lay dormant. Diabetes without clinical consideration is fundamentally connected to genuine complexities, which can add an impressive weight to the general wellbeing framework. The predominance of diabetes is relied upon to increment quickly later on because of the commonness of weight, maturing of the populace, and other cardiovascular danger factors. Confusions of diabetes mellitus like cardiovascular infection, kidney harm, etc, ought to be forestalled in beginning phase. Be that as it may, diabetes is generally asymptomatic. Individuals with undiscovered diabetes are bound to be determined to have complexities than the individuals who know about their diabetes status. In spite of the fact that fasting glucose plasma, the oral glucose resistance test, and hemoglobin are grounded determinants in diabetes determination, they are lacking to give intrusive screen tests to an enormous populace. Hazard evaluating frameworks for patients with undiscovered diabetes has been created. Built up a self-evaluation score for diabetes hazard in Korean grown-ups suggested by lee et al and Zhou etal. They proposed a diabetes evaluating model for moderately aged provincial Chinese population. An expectational hazard score for individuals at high danger of diabetes in south Asia by Aekplakorn et al. The model to anticipate the three-year occurrence of type 2 diabetes in a Japanese populace by Aekplakorn et al. A diabetes hazard score for screening undiscovered diabetes and approved it utilizing Chinese grown-ups by Gao et al. Proposal's framework are utilized to forestall diabetes through changes in way of life and mediation with drug medicines. Be that as it may, research utilizing AI innovation to create evaluating apparatuses for undiscovered diabetes has been lacking. Past investigations have presented prescient models for infections like diabetic retinopathy, skin malignancy, lung illness, cardiovascular breakdown, persistent kidney sickness, etc utilizing computerized methods. These examinations that utilization machine learning procedures to cause significant advances in tackling issues to have opposed the best endeavors of the man-made brainpower local area much of the time. Albeit past investigations have created prescient models dependent on calculations made by machine models, it is hazy whether these models can be appropriately utilized for assessing undiscovered diabetes. Thus, the target of the current investigation was to build up a model (DLM) for patients with undiscovered diabetes.

**II. PROPOSED SYSTEM**

**2.1 System Architecture**

In this Methodology, patient data is collected from PIMA india dataset extracted from Kaggle repository. Data preprocessing is done with PIMA dataset to remove missing values and noisy information. After this process user enter the patient's data from online. The model is trained to predict Diabetes from the patient's input data through LSTM as shown below figure shown in next page.



**2.2. Data Understanding and Sources of Data**

The dataset taken for the training is PIMA medical Indian dataset extracted from Kaggle repository. PIMA datasets contains the pregnancies, glucose levels, blood pressure levels, skin thickness, insulin, BMI, diabetes pedigree function and age as its parameter list.

PIMA Dataset details

| S.No. | Parameter | Parameter Description |
|-------|-----------|----------------------|
| 1 | Pregnancies | Number of times pregnant |
| 2 | Glucose | Plasma glucose concentration a 2 hours in an oral glucose tolerance test |
| 3 | Blood Pressure | Diastolic blood pressure (mm Hg) |
| 4 | Skin thickness | Triceps skin fold thickness (mm) |
| 5 | Insulin | 2-Hour serum insulin (mu U/ml) |
| 6 | BMI(Body Mass Index) | Body mass index |
| 7 | Diabetes Pedigree Function | Diabetes pedigree function |
| 8 | Age | Age (years) |

Random forest classifies all the parameters into different classes of valuation parameters. The valuation parameters are taken as input and fed into LSTM classifier. It is a recurrent network constructed on the basis of valuation parameters. Where each valuation parameter is fed into the network. It can be defined as

$F(x) \rightarrow h$

Where h is a valuation parameters property. Y and v are the series functions. Which denotes the parameters time variation with reference to variations of changing health constants. F(x) will be used as a training parameter that defines levels of the disease. By

this F(x) the system determines level of diabetes that has to be trained to the classifier. By this F(x) parameter the training model is built. The next phase constitutes of testing part. Each and every parameter will be fed into the classifier by the user. This can be defined as

$$h= E[M(D,t+1)] \geq \frac{F(D,t)}{\bar{F}(t)} N(D,t)\{1 - \epsilon(D,t)\}$$

Where F(D,t) is the fitness of schema D, $\bar{F}$(t) is the average fitness of the population, and $\epsilon$(D,t) is a term which reflects the potential for genetic operators to destroy instances of schema S.

From this, each parameter will be compared with the trained model. The values are pitted against the trained model where each parameter will be evaluated as per the diabetes classifier. Each parameter is then forwarded to Naive – Bayes classifier which is defined as The distance d(u, v') between v(x, y) and v'(x', y') is given on the zw-coordinate system v(z, w) and v'(z', w') by

$$d(v, v') = |x - x'| + |y - y'| = max(|z - z'|,|w - w'|).$$

Before explaining our algorithm, we define some terminology. For a Steiner point v, (weighted) path lengths from v to all leaves in the subtree rooted by v are equal in outspreading nodes. This length is called the path length from the Steiner point v and is denoted by c(v). This definition is consistent with the weight c(v) for leaves because we can consider the (weighted) path length from a leaf v as its weight. Note that c(v) = c(v') for all points v' on prospective segment $l_v$. Thus, we define path length c($l_v$) from a prospective segment $l_v$ by c(v). We consider the prospective segment $l_{vi}$ for each leaf $v_i$ as a segment with length 0 containing only the leaf $v_i$, and, therefore, c($l_{vi}$) = c($v_i$).From the above function it can be said that all the disease parameters are probabilistically determined for prediction of diabetes function.

**2.3. LSTM :** Long short-term memory (LSTM) units are a special type of building units for RNN. It can analyse, classify and predict temporal data sequences of time lags of any size. A typical
LSTM network is made up of memory, input, output and forget gates. The memory in LSTM can remember values overarbitrary time intervals. Each of the three gates is a form of neuron (which computes an activation function of a weighted sum). More than that, these gates control the passage of values in LSTM layers; hence these special neurons are named as gates. By long short-term, the fact underlined is that LSTM's memory can really last for large time duration. LSTM tackles the issue of exploding and vanishing gradient problem which is an important issue while training traditional RNNs.

## III. Methodology

### 3.1Data Ingestion:

The LSTM  Long Short-Term Memory proposed by Hochreiter and Schmidhuber can efficiently deal with the big time dependencies in the classification, and its assembly. The focal point of LSTM lies in its memory unit, and related realities is communicated in reverse through the memory unit. Hypothetically, the memory unit can switch measurements during the entire assortment spread strategy so the insights at the past time can be utilized to are anticipating the yield at the later time, so it could resolve the quick time-frame memory issue of the regular intermittent neural local area. Moreover, all through the regressive switch of realities in the memory unit, the LSTM gives or erases data in the memory unit through 3 entryways. Through these entryways they are Pregnancies, Glucose, Blood Pressure, Skin Thickness, Insulin, BMI, Diabetes Pedigree Function and Age. These parameters are taken from the PIMA diabetes dataset. Each of these parameters are feed to the memory unit of the algorithm. These parameters are individually feed to the memory cell of the LSTM.

### 3.2 Data Preprocessing:

These doors of the LSTM can be viewed as unique neural organizations, which might be prepared to naturally contemplate what data to keep or neglect. The data arriving from the data Ingestion phase will be fed to the doors. The arrangement of LSTM preparing information is as per the following. To start with, the LSTM will utilize the "neglect door" to choose which data must be eliminated. The elimination data is omitted according to the value ie if the person is having glucose as 126 mg/dL, Insulin as 120mg/dL and Diabetes Pedigree Function above 6.5% then the values are allowed through the door rest of the values are omitted from the dataset. The entry is planned somewhere in the range of 0 and 1 via the Sigmoid trademark these are defined to be the values that calculated by the equation

$$V=Max \left\{ \frac{(l_1 - x_1^0)a'}{a_2}, \frac{(x_2^0 - u_2)a'}{a_1} \right\} \leq t \leq Min \left\{ \frac{(u_1 - x_1^0)a'}{a_2}, \frac{(x_2^0 - l_2)a'}{a_1} \right\}.$$

Here 'V' is the values of the the max and min of the function.  The design to "1" signifies to keep the data; in some other case, it strategy to ignore measurements. The "input entryway" is utilized to choose which insights wants to be exceptional. The Sigmoid capacity is utilized to decide if it wants to be held. Then, at that point, the tan h trademark maps the information cost to [−1, 1], along these lines creating another memory unit realm and adding it to the valid memory unit. The data in the memory unit will be of the parameter such as Pregnancies, Glucose, Blood Pressure, Skin Thickness, Insulin, BMI, Diabetes Pedigree Function and Age are stored in the memory cell with its corresponding parameters. Then, to refresh the cost of the memory unit, first increase the memory unit via the disregard door, dispose of the records that desires to be neglected, and afterward transfer the information data acquired from the enter entryway to achieve the shiny new memory unit cost. At last, the "yield door" goes to a choice which memory unit records to yield, this is, resultant model arriving from the calculation part . The computation plan of LSTM is as per the following:

**3.3 Modelling:**

**Step 1:** The data from the PIMA dataset is input to the Entry of the LSTM

**Step 2:** The data is omitted from the input door according to the base conditions ie the peak values of Pregnancies, Glucose, Blood Pressure, Skin Thickness, Insulin, BMI, Diabetes Pedigree Function and Age.

**Step 3:** The data then will be stored in the forwarded to the temporary memory unit of the classifiers by calculating the sigmoid function.

**Step 4:** After the step3 each data block from the parameter will be stored in the memory cell of the classifier.

**Step 5:** The consortium of the memory cells constitute of the LSTM Diabetes model.

**4.4 Detection Phase**

**Step 1:** The data from the user is given as input to the model of LSTM

**Step 2:** The data values of Pregnancies, Glucose, Blood Pressure, Skin Thickness, Insulin, BMI, Diabetes Pedigree Function and Age are provided to the flask based web interface of the system.

**Step 3:** The data then will be forwarded to the temporary memory unit model.

**Step 4:** Each value of the data is then forwarded to respective unit of the memory cell of the LSTM Model.

**Step 5:** The value comparison is processed inside memory cells of the LSTM Diabetes model.

**Step 6:** After value comparison is processed inside memory cells then the values matching the values of the memory cell are allowed to pass from the memory.

**Step 7:** The value that is allowed to pass from the memory will reach the output gate of the LSTM classifier.

**Step 8:** The values that are passed to output gate according to the memory node comparison will be displayed on the output of the webpage.

**Step 9:** The stage of the diabetes is displayed as output will be in accordance of node memory values.

**Algorithm**

*Begin*

1. Input datasets . $K := \{lv1, lv2, \ldots, lvn\}$ input node(gate).
2. Extract disease metadata and classifier. $R(S) = \{(z, w) \mid z_{min} \leq z \leq z_{max}, w_{min} \leq w \leq w_{max}\}$ the omitted parameters.
3. Determine the base conditions $z_{max} = max_i(z_i + c(v_i))$ the values to store in the memory node.
4. Set upper limit and lower limit for disease according metadata and related data zmax.
5. Train the classifier $K(x_i, x_j) = exp\left[-\frac{(x - xj) \cdot (xi - xj)}{2\sigma^2}\right]$ the values stored in memory node.
6. End training creation of model.
7. Input disease parameter from user $\{(x_1, y_1), (x_2, y_2), (x_3, y_3), (x_4, y_4) \ldots \ldots, (x_n, y_n)\}$ from webpage.
8. Compare with model where K is memory node $K(x_i, x_j) = (\gamma x_i^T x_j + r)^d$ , $\gamma > 0$.
9. Detect the diabetes.

*End*

## IV. RESULTS AND DISCUSSION

**4.1 Results of Descriptive Statics of Study Variables**

PIMA datasets parameters were trained and fed into the Algorithm for the classification part of each parameter. This parameters were arranged in various class levels of hierarchy. This data was trained in the LSTM classifier where the diabetes conditions were trained. This trained classifier is existing in the form of a model. A web interface provided to the user to input all the disease parameters and check the diabetes level and to see the person is diabetic or not. The accuracy found in this system was found to be 92% and was found to be efficient in the most cases. Below Fig 4.1 shows Result of the Daiabetes Detection after the user enter the patient's data. The result given as "Time for Dessert" means Patient is not Diabetic. Figure 4.2 shows the correlation of the parameters, skin and thickness are highly correlated.

## Diabetes Detection

Number of pregnancies
`2`

Plasma glucose concentration in oral glucose tolerance test
`142`

Diastolic blood pressure (mm Hg)
`85`

Triceps skin fold thickness (mm)
`24`

2-Hour serum insulin (mu U/ml)
`385`

Body mass index
`34`

Diabetes Pedigree Function
`1.2`

Age
`49`

[Send]

Time for dessert

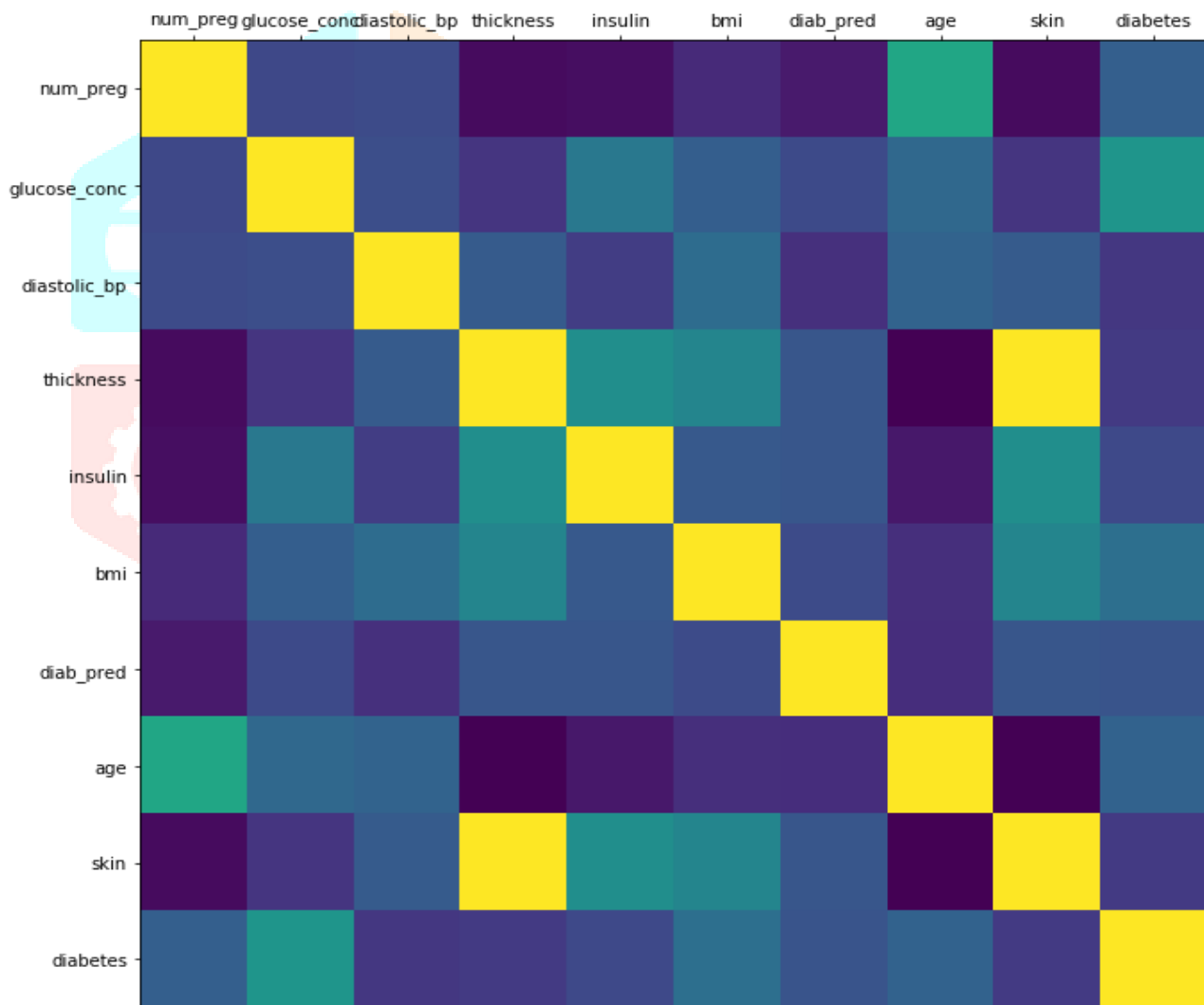Fig 4.1 Result of the Daiabetes Detection after the user enter the patient's data.



Fig 4.2  Parameter Correlation

## V. CONCLUSION

A lion's share of human populace has been influenced by diabetes. Diabetes can't be relieved; it must be monitored. Uncontrolled diabetes can prompt difficulties, prompting a few other constant sicknesses. Hence, early recognition, proficient treatment and appropriate administration of diabetes is vital. Diabetes causes nerve issues which can influence heart and along these lines pulse as well. Here, in our work, we use HRV information (removed from ECG motions toward) distinguish diabetes with a high exactness. To the extent our insight, this is the primary work to utilize profound learning in recognizing diabetes utilizing HRV information. The precision of 95.1% accomplished utilizing CNN-LSTM network with 5-crease cross-approval is the most noteworthy exactness gotten so far in the robotized recognition of diabetes utilizing HRV. There is no necessity of unequivocal component extraction and utilization of customary classifiers. Our technique is non-intrusive and reproducible. Our framework can help the clinicians to analyze diabetes precisely. As clarified under Results segment, further improvement in exactness can be investigated by taking care of into the proposed design huge measured info dataset contrasted with the dataset size utilized in this work.

## VI. ACKNOWLEDGMENT

## VII. REFERENCES

[1] E. Daskalaki, A. Prountzou, P. Diem, and S. G. Mougiakakou, "Real-Time Adaptive Models for the Personalized Prediction of Glycemic Profile in Type 1 Diabetes Patients," Diabetes Technol. Ther., vol. 14, no. 2, pp. 168–174, 2012.

[2] R. Bunescu, N. Struble, C. Marling, J. Shubrook, and F. Schwartz, "Blood Glucose Level Prediction Using Physiological Models and Support Vector Regression," 2013 12th Int. Conf. Mach. Learn. Appl., pp. 135–140, 2013.

[3] E. I. Georga, V. C. Protopappas, and D. I. Fotiadis, "Glucose Prediction in Type 1 and Type 2 Diabetic Patients Using Data Driven Techniques," Knowledge-Oriented Appl. Data Min., pp. 277–296, 2011.

[4] H. N. Mhaskar, S. V. Pereverzyev, and M. D. van der Walt, "A Deep Learning Approach to Diabetic Blood Glucose Prediction," Front. Appl. Math. Stat., vol. 3, no. July, pp. 1–11, 2017.

[5] Y. Bengio, P. Simard, and P. Frasconi, "Learning Long-Term Dependencies with Graident Descent is Difficult," Saudi Med J, vol. 33, pp. 3–8, 2012.

[6] S. Park, S. Min, H.-S. Choi, and S. Yoon, "Deep Recurrent Neural Network-Based Identification of Precursor microRNAs," Nips, no. Nips, 2017.

[7] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," Neural Comput., vol. 9, no. 8, pp. 1735–1780, 1997.

[8] M. Roondiwala, H. Patel, and S. Varma, "Predicting Stock Prices Using LSTM," vol. 6, no. 4, pp. 2015–2017, 2017.

[9] F. Altché, A. De, and L. Fortelle, "An LSTM Network for Highway Trajectory Prediction."

[10] Q. Zhang, H. Wang, J. Dong, G. Zhong, and X. Sun, "Prediction of Sea Surface Temperature using Long Short-Term Memory," pp. 1–5, 2017.

[11] S. Mirshekarian, R. Bunescu, C. Marling, and F. Schwartz, "Using LSTMs to learn physiological models of blood glucose behavior," Proc. Annu. Int. Conf. IEEE Eng. Med. Biol. Soc. EMBS, pp. 2887– 2891, 2017.

[12] P. Su, X. Ding, Y. Zhang, F. Miao, and N. Zhao, "Learning to Predict Blood Pressure with Deep Bidirectional LSTM Network," pp. 1–19, 2017.

[13] C. Olah, "Understanding LSTM Networks," 2015. [Online]. Available: http://colah.github.io/posts/2015-08-Understanding-LSTMs/.

[14] M. Schuster and K. K. Paliwal, "Bidirectional recurrent neural networks," IEEE Trans. Signal Process., vol. 45, no. 11, pp. 2673– 2681, 1997.

[15] B. Lia et al., "Carbohydrate Estimation Supported by the GoCARB system in Individuals With Type 1 Diabetes: A Randomized Prospective Pilot Study," Diabetes Care, vol. 40, no. 2, pp. e6–e7, 2017.

[16] M. Anthimopoulos et al., "Computer vision-based carbohydrate estimation for type 1 patients with diabetes using smartphones," J. Diabetes Sci. Technol., vol. 9, no. 3, pp. 507–515, 2015.

[17] E. Daskalaki, "Towards the External Artificial Pancreas : Design and Development of a Personalized Control System for Glucose Regulation in Individuals with Type 1 Diabetes," University of Bern, 2013.