

# SDG 6 CLEAN WATER AND SANITATION

Murthy Kaja | Yashwanth Madaka | Kowshik Sola | Nithin Reddy

Department of Computer Science



Stony Brook University

BILLIONS OF PEOPLE STILL LACK  
ACCESS TO SAFE DRINKING WATER,  
SANITATION AND HYGIENE

IN 2020



2 BILLION PEOPLE

26%

LACK  
SAFELY MANAGED  
DRINKING WATER



BETWEEN 1970 AND 2015,  
NATURAL WETLANDS  
SHRANK BY 35% —↓—

3 x THE RATE OF FOREST LOSS



129 COUNTRIES ARE **NOT ON TRACK** TO HAVE  
SUSTAINABLY MANAGED WATER RESOURCES BY 2030

CURRENT RATE OF PROGRESS NEEDS TO **DOUBLE**

6 CLEAN WATER  
AND SANITATION



## GOALS



Target 6.1. Safe and Affordable drinking water.

Indicator: Proportion of population using safely managed drinking water.



Target 6.3. Improve water quality, waste water treatment and safe reuse.

Indicator: Proportion of bodies of water with good ambient water quality.

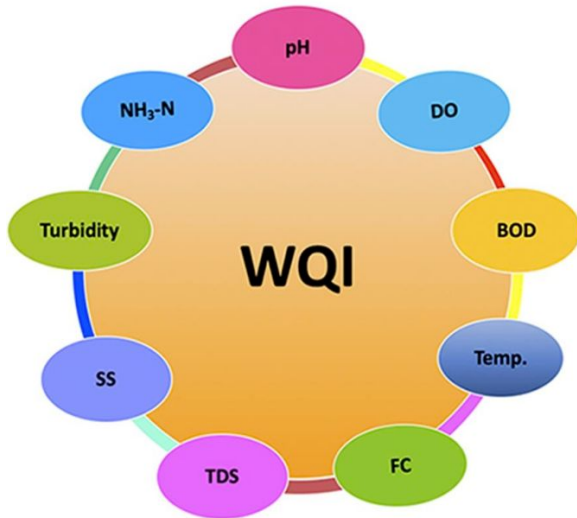
## WHY BIG DATA?



Around 20,000 stations generating tons of data every month.

# Background

“A Review Of Water Quality Index Models And Their Use For Assessing Surface Water Quality”

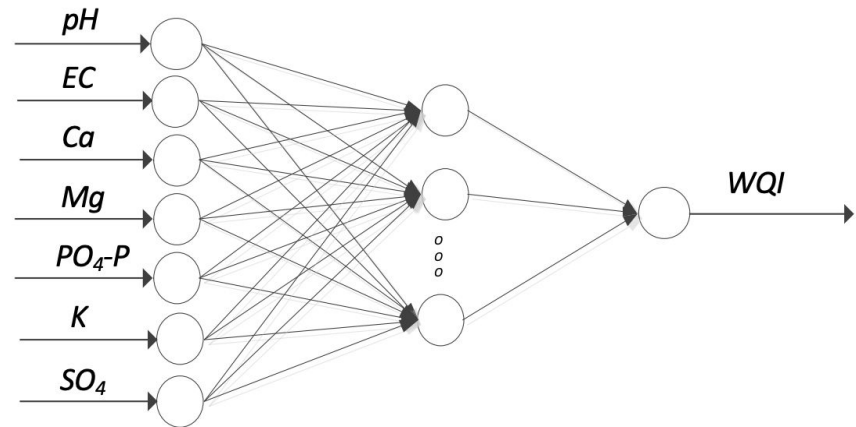


<https://www.sciencedirect.com/science/article/pii/S1470160X20311572>

“Short-term water quality variable prediction using a hybrid CNN–LSTM deep learning model”

<https://link.springer.com/article/10.1007/s00477-020-01776-2>

“Forecasting Water Quality Index in Groundwater Using Artificial Neural Network”



<https://www.mdpi.com/1996-1073/14/18/5875/pdf>

“Provide decision makers and stakeholders with quantitative information to facilitate sustainable management of ongoing and emerging environmental problems”

# Data

Icons taken from: <https://www.flaticon.com/>

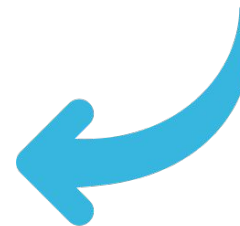
**GEMStat** 



European  
Environment  
Agency



**AGGREGATE**



(Time Series Data, 14 GB)

Station ID	Location	Time	pH	BOD	Temp	TSS	Nitrate	.....
IND1002	India	11-01-1995	....	....	....	....	....	....
NOR341	Norway	12-02-1997	....	....	....	....	....	....
EGY742	Egypt	01-03-2003	...	...	...	...	...	...

Name	From	To
Europe	1941	2018
Africa	1960	2020
India	1971	2008

**Table.** Years of data for Europe, Africa, India

# Methods

## 1. DATA PREPROCESSING

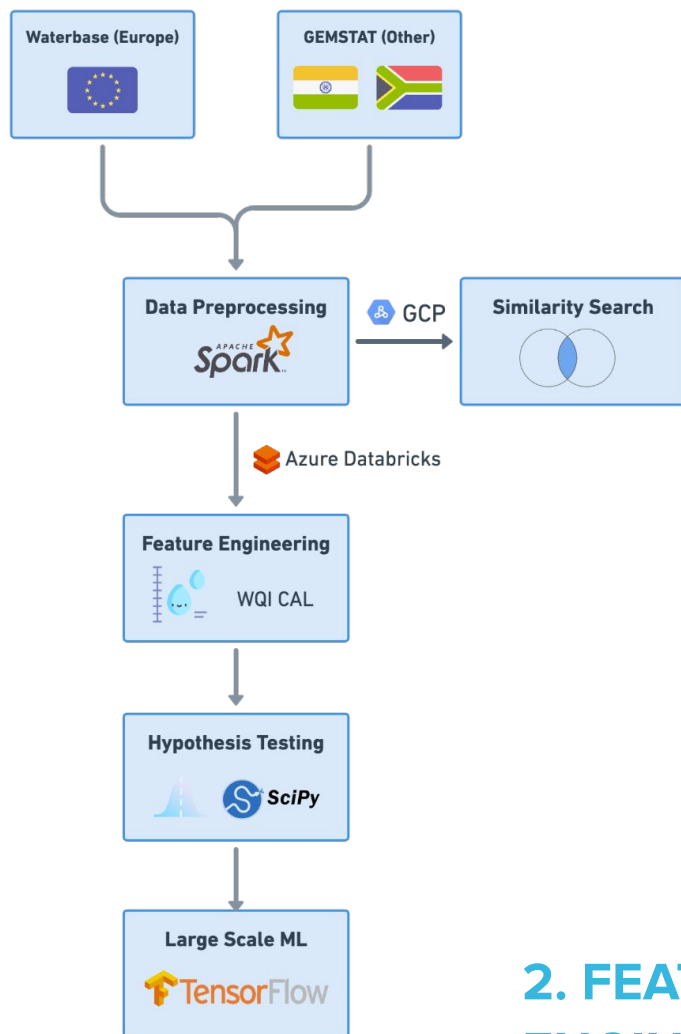


Figure. Data Pipeline

### A. SPARK DATAFRAME

1. Remove NULLS

### B. PANDAS DATAFRAME

1. Use GeoPy to extract Locations

### A. SPARK DATAFRAME

1. Convert to RDD

Structured Chemicals Dataset

### A. SPARK DATAFRAME

1. Extract Date
2. Filter Required Chemicals

### B. HIVE QUERY LANGUAGE

1. Convert Chemical values into unit of measure

### C. SPARK RDD (MAP & REDUCE BY KEY)

1. Transform rows into columns for Chemicals
2. Replace Missing Values with Median

LAT, LONG Dataset

Location Data with Station ID

Chemicals Data with Station ID

SPARK RDD INNER JOIN

Location Data, Chemicals Data with Station ID

## 2. FEATURE ENGINEERING

WQI Calculation Ref : [Link](#)

$$\text{WQI} = \text{TEMP} * (\text{BOD} + \text{TSS} + \text{DO} + \text{COND})$$

DO = Dissolved Oxygen Index  
COND = Conductivity Index

Temp = water temperature index, BOD = Biological Oxygen Index, TSS = Total suspended Index

# Methods

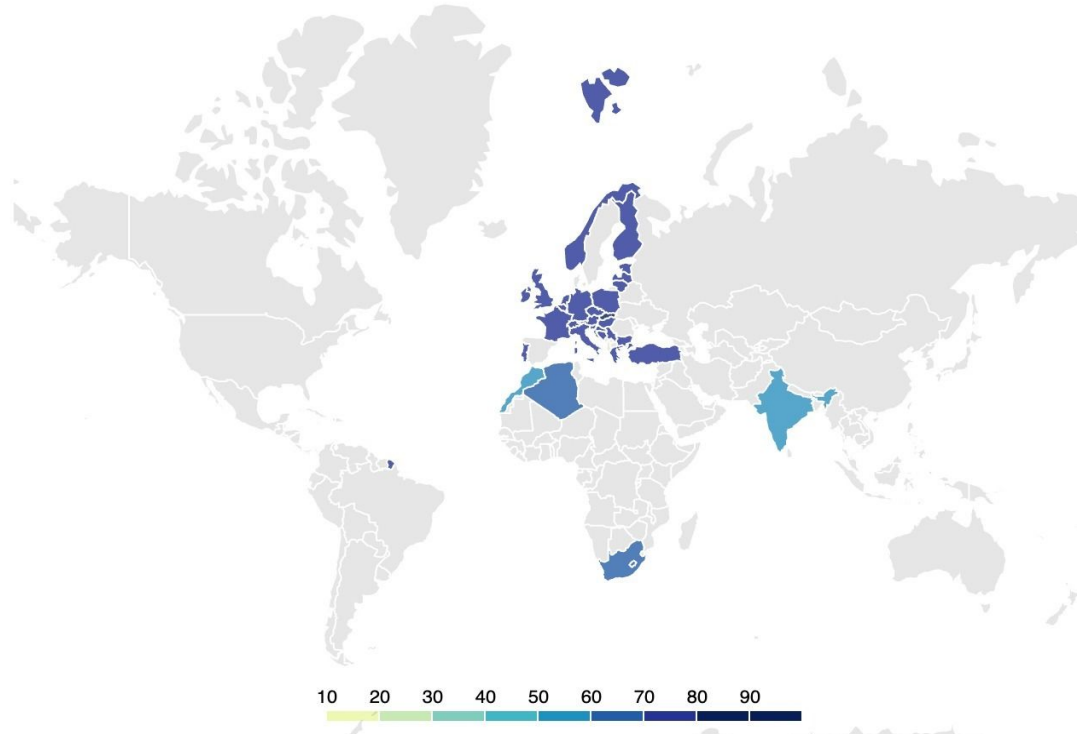


Figure. Heat Map of WQI ranges for Africa, Europe and India

## 5. LARGE SCALE ML

### LSTM Model for Time Series

Batch\_Size: 16

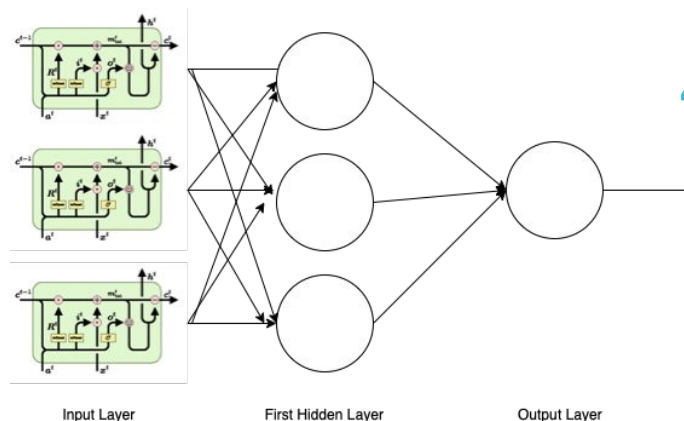
Epochs: 50:

Input Layer: 100 LSTM Units

Hidden Layer: 50 FCN,

Output Layer: 1 FCN,

Split: 80% Train, 20% Test



## 3. SIMILARITY SEARCH

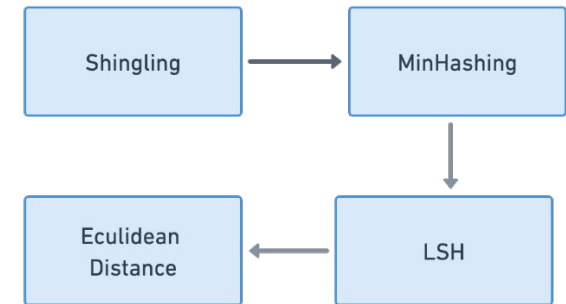


Figure. Similarity Search Pipeline

### Factors

**Economical:** GDP, GNI

**Environmental:** CO2 Emissions, Air Pollution

**Societal:** Population, HDI, Life Expectancy

**All Chemicals:** pH, TSS, BOD, COND, TEMP

## 4. HYPOTHESIS TESTING

WQI VS Chemicals

$$P Value_{Bonferroni} = P Value * m$$

(where  $m$  is Number of countries)

# Evaluations / Results

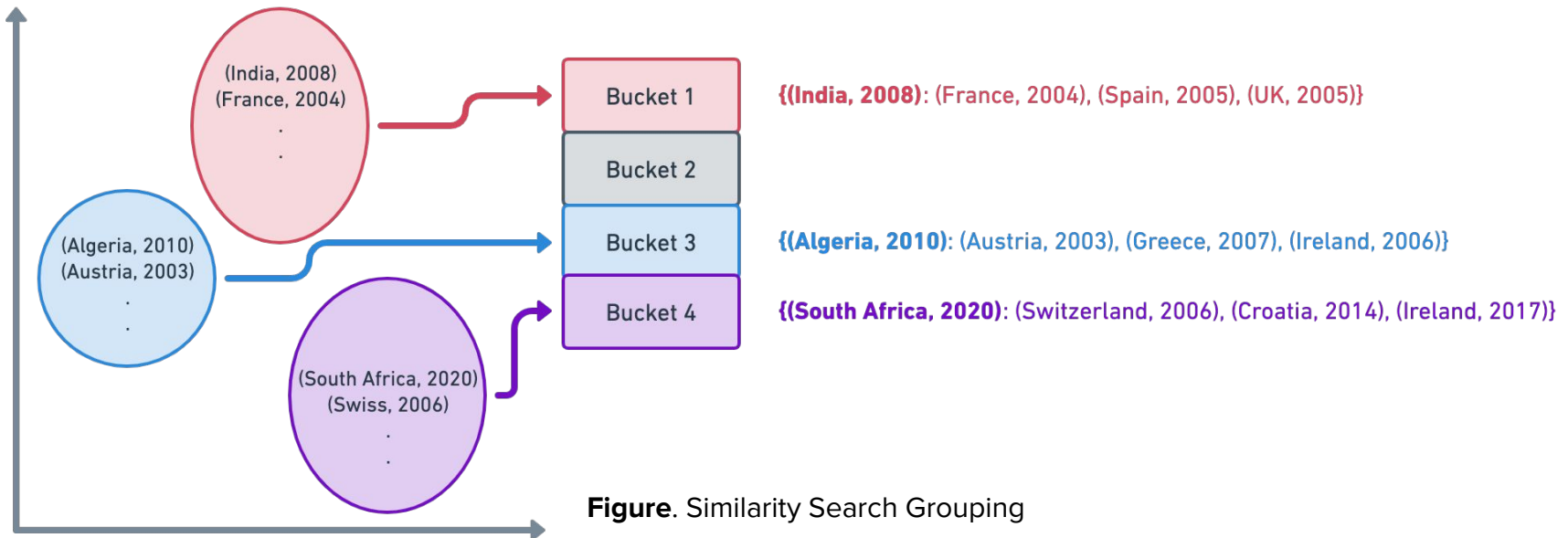
Features	Cosine Similarity	P-Value	Bonferroni corrected P-Value
Iron	0.005157232	0.00011953	0.006932954
Nitrate	-0.074133163	2.9222E-10	1.69490E-08
Chloride	-0.103132809	0.00078071	0.045281067
Sodium	-0.044942523	2.6782E-10	1.55334E-08
Dissolved_oxygen	0.606046465	1.0456E-09	6.06456E-08
Water_temperature	-0.171881863	6.4731E-05	0.003754424
Total_suspended_solids	-0.006519594	7.7272E-05	0.004481761
Conductivity	0.017448662	0.00029797	0.017281976
Phosphate	0.062399443	9.9192E-05	0.005753125
pH	0.07648714	6.3983E-10	3.71104E-08
Non_ionised_ammonia	-0.099626672	0.02091294	1.2129504

**Table.** Hypothesis Testing Results

## Bonferroni corrected P-value

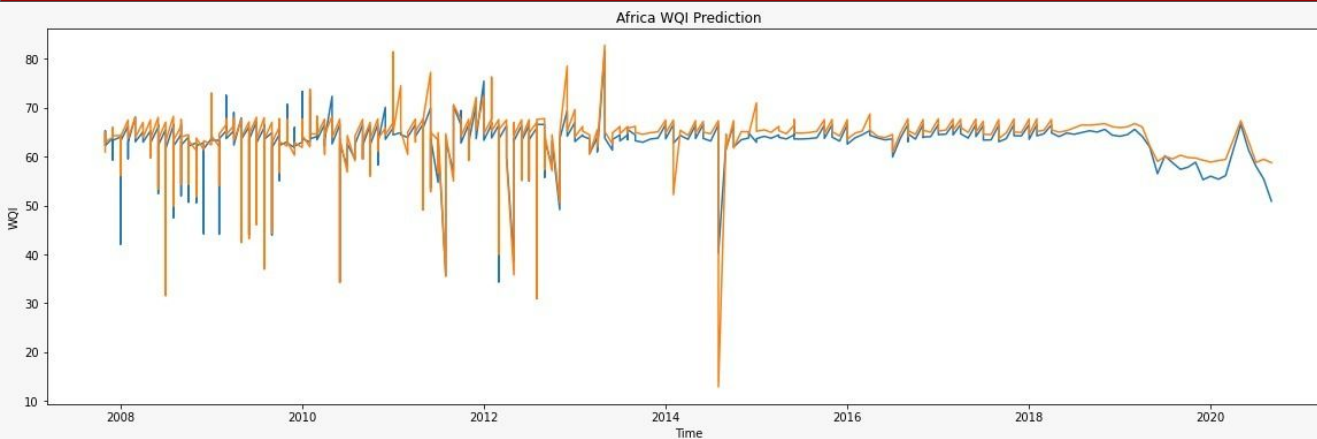
✗ Non\_ionised\_ammonia > 0.05

✓ Others < 0.05

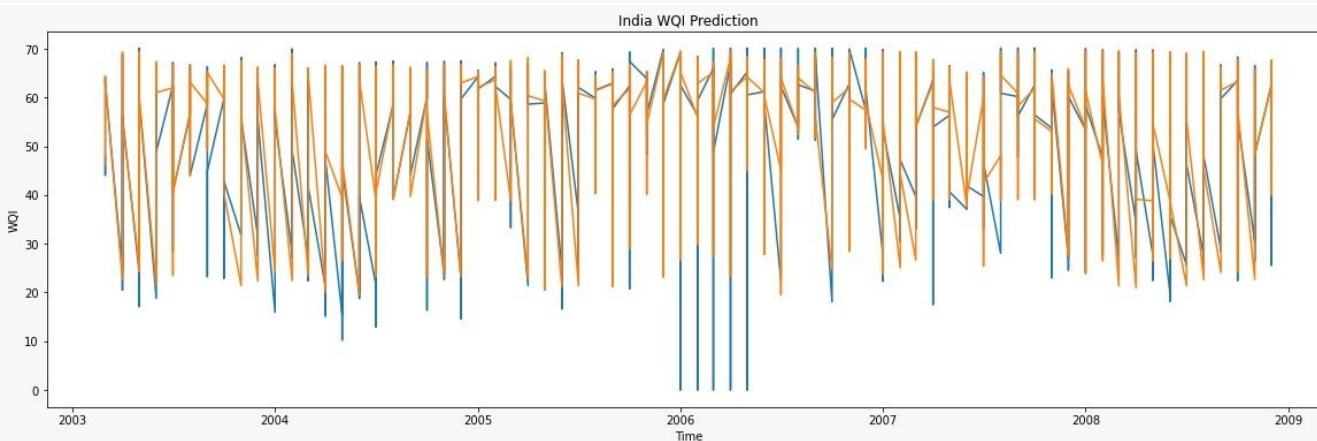


**Figure.** Similarity Search Grouping

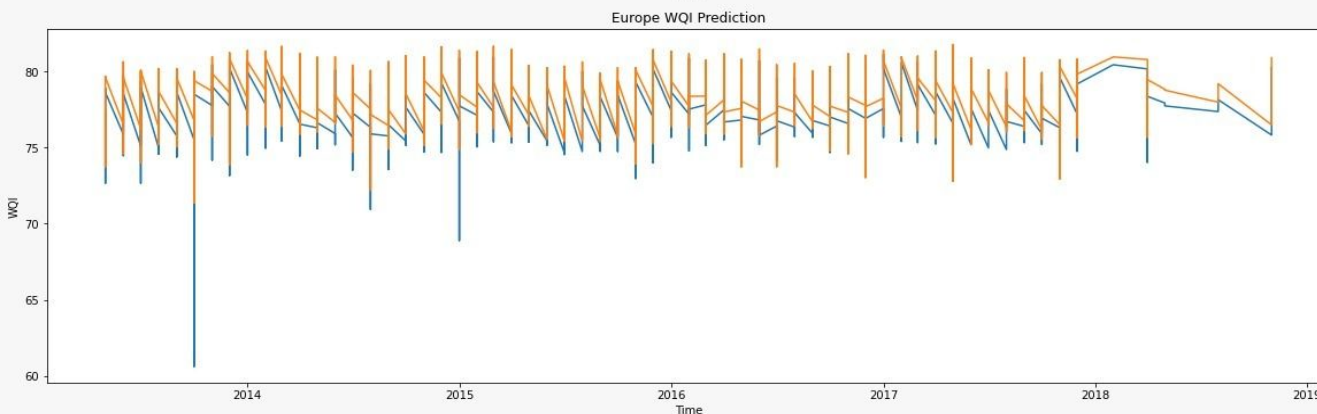
# Evaluations / Results



**AFRICA - LSTM**  
Prediction (Orange) vs Original (Blue)  
From 2008 to 2020  
MAE: 0.84046



**INDIA - LSTM**  
Prediction (Orange) vs Original (Blue)  
From 2003 to 2009  
MAE: 1.05017



**EUROPE - LSTM**  
Prediction (Orange) vs Original (Blue)  
From 2014 to 2019  
MAE: 0.88669



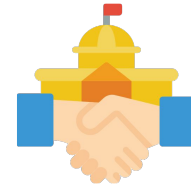
# Conclusion



**Efficient Water Quality Monitoring using our LSTM Model compared to conventional methods**



**Inform Authorities/Governing bodies to take prior action to mitigate damage to Water Quality**



**Inform the Government/Governing body to follow the footsteps of Developed Countries**

**TARGET**

**6·1**



**SAFE AND AFFORDABLE DRINKING WATER**



**TARGET**

**6·3**



**IMPROVE WATER QUALITY, WASTEWATER TREATMENT AND SAFE REUSE**