# Temporal Air Quality Forecasting using Hybrid RNN Model.

DharsanGM

*DepartmentofDataScienceandBusinessSystems,*
*SchoolofComputing,*
*CollegeofEngineeringandTechnology,*
*SRMInstituteofScienceandTechnology,*
Kattankulathur,Chennai,603203,India.

PriyadarsiniK

*DepartmentofDataScienceandBusinessSystems,*
*SchoolofComputing,*
*CollegeofEngineeringandTechnology,*
*SRMInstituteofScienceandTechnology,*
Kattankulathur,Chennai,603203,India.

*Abstract*—Inmanynations,poorairqualityhasturnedintoa significant environmental issue. The majority of air pollutantsseriouslyharmpeople'shealthandqualityoflife.Rapi durbanisation,industrialexpansion,andtrafficareallcontrib utingto a significant decline in air quality. In this work, we suggest ahybrid recurrent neural network (RNN) model that combines thebenefits of wavetransform decomposition and LSTM model fortemporalairqualityforecasting.Thetimeseriesisdividedint oitstrend, seasonal, and residual components using the decompositiontechnique, which is then fed into the LSTM model for forecasting.Weexaminedtheeffectivenessofthesuggestedmod elusingdailyairqualitydatacollectedfroma monitoring station in Delhi.OuroutcomeindicatethattheH-LSTMmodeloutperformstraditionaltime-series models and the LSTM model in termsof predicting accuracy, having reduced the mean absolute errorandmeansquarederrorof50%forPM2.5concentrations, performs better in terms of forecasting accuracy than traditionaltime-series models and the Long short - term memory model.Furthermore, the decomposition technique allows us to identifytheunderlyingtrendsandseasonalpatternsinthedata,w hichcanprovideinsightsforairqualitymanagement.Ourresults highlight the potential for further advancements in thisfield andshow the value of using a decomposition technique and LSTMmodelforairqualityforecasting.

*IndexTerms*—
Airquality,wavedecomposition,PM2.5concentration,LSTM

## I. INTRODUCTION

On a global level, air pollution poses a serious hazard tohumanhealththatisgrowinginimportance.TheWorldHealth Organization (WHO) 2018 report states that 92% of peopleworldwidebreathetoxicair,whichisthoughttobethecaus eof7 million annual fatalities. Two respiratory and cardiovasculardisorders,emphysemaandasthma,areamongthe harmfulimpacts of pollution in the air on human health. In addition,ambientairpollutionhas1.8yearsofftheaveragelifetim e.

[9] The World Bank estimates that the economic costs of airpollution, including lost productivity, medical costs, and earlydeaths,total$5trillionannually.By2050,itisanticipatedtha t the number of premature deaths brought on by particulatematter(PM)willhavemorethantripled.Thus,itisesse ntialto identify a precise and practical method to lower

ambient airpollution if we are to guarantee the security and well-being oftheglobalpopulace.

Wecannotseepollutants,andwearenotawareofthereasons whypollutionlevelsarerising.Tocomprehendtheorigins,itisnec essaryto firstgoovertheprinciplesofairpollution[3].PM2.5endangersth eenvironmentglobally.Long-termexposuretotheseparticles,whichcanquicklyreachdeepwit hintherespiratoryorgans,canresultinavarietyofdiseases.Thedi ameterofPM2.5,whichismuchsmallerthanthediameterofahum anhair,is2.5microns.Asaresult,thereis asignificant sourceof exposureand sicknessonthe humanbody.Theyendangerhumanhealth,raisingthelikelihood ofcancerandotherendemicandpandemicdiseases.Theultimate goalistoraisepublicawarenessofthevariablesthataffectenergyc onsumptionandmethodsforcontrollingit.[8]

Inthispaper,wepresentahybridRNNmodelthatcombinest hedecompositiontechniquewithlongshort-termmmemoryfor temporal air quality forecasting (LSTM). While the de-compositiontechniquesareusedtobreakdownthetimeseries into wave transforms, which aid in identifying trends,cyclic patterns, etc., the LSTM cells are renowned for theircapacity to capture long-term dependencies in the data. Theproposed model also incorporates input features, which havebeen demonstrated to be significant predictors of air quality,suchasmeteorologicalvariablesandtimeofday.

Usingactualairqualitydatagatheredfromamonitoringstati oninDelhi,weassessedtheperformanceofthesuggestedmodel. Withameanabsoluteerrorofg/m3forPM2.5concentrations,our resultsshowthatthehybridRNNmodeloutperformsconvention altime-seriesmodelsandotherRNNmodelsintermsofforecastingaccur acy.Inaddition,thehybridRNNmodeloffersusefulinsightsintot hetemporaltrendsinairpollutionthatmaybeutilisedtoguidepoli cyanddecision-makingforthemanagementofairquality.Therestofthispaperisst ructuredasfollows.AsummaryofrelatedresearchemployingR NNmodelsforairqualityforecastingisincludedinSection2.Exp erimentalareaofstudyexplainedinSection3.Theexperimentals etupandproposedmethodologyisexplainedinsection4.The findingsofourmodelevaluationareshowninSection5and6.Sect ion7concludesthepaperandconsiderspossiblefuturelinesofenq uiryforthis fieldofstudy.

## II. RELATEDWORKS

In order to forecast particulate matter (fraction PM2.5) airpollution, a study comparing artificial neural networks (ANNs)andadaptiveneuro-fuzzyinferencesystems(ANFIS)wasconducted. The results are provided in this publication. PM2.5hourly measurement records from the Airbase databases wereused for the trials. The two key statistical measures calculatedwere mean absolute error and root mean square error (MAE).[6]

This study proposes improved models for estimating hourlyair pollution concentrations using meteorological data frompreviousdays.Amulti-tasklearning(MTL)problemwithvariousregularisationtechniq uesisusedtoformulatethepredictionoverthenext24hours.There sultsshowthatthe parameter-reducing formulations and regularisations as-sociatedwithsuccessivehoursoutperformcommonlyusedstand ard regression models and regularisations for MTL whencomparedtotheproposedregularisation.[10]

Thisstudyusesalongshort-termmemory(LSTM)approachtoforeseeO3,PM2.5,NOx,and COconcentrationsinDelhi,theheavilypollutedIndianNational CapitalTerritory.Theresearchadoptsiveseparatesetsofparamet ersandcriteria,suchastrafficdata,levelsofpollution,climaticsta tus,andvehicleemission.TheLSTMmodelsaccuratelyestimate hourlyconcentrationswhilehandlingthecomplexitiesoflong-terminteractionscomingfrombothhumanandnaturalsources.T heresultsofthisstudymaybeusedbythegovernmentandpolicym akerstoplanactionstoreducethenegativeeffectsofdecliningairq uality.[4]

Becausepollutantscanbeharmfultohumanhealth,predictin gairpollutionisahottopic.Traditionalmachinelearningmodelst ypicallyfocusonincreasingpredictionaccuracyoverallwhileig noringaccuracyforpeakvalues.Furthermore,itisimpossibletoc omprehendthesemodels.Theydonotadequatelyexplainhowvar iousdecidingfactorsinteracttoaffectairpollution.Inthisstudy,a newcombinationofinterpretablepredictivemachinelearningm odelisenhancedwithtwonoveltiesforthepredictionofPM2.5.Fi rst,afusionmodelframeworkisdevelopedusingdeepneuralnet worksandanonlinearautoregressivemovingaveragewithexoge nousinputmodel.Second,thishybridmodelincludesmethodsfor automaticallycreatingandchoosingfeatures.Theexperimental findingsshowthat ourmodeloutperformsothermodelsintermsofhighestvaluepred ictionprecisionandmodelinterpretability.Theproposedmodeli llustrateshowtocalculatePM2.5estimatesusingexistingPM2.5 ,climate,andseasondata.Therecommendedapproachalsoprovi desanoveleasilyunderstoodmachinelearningframeworkforte mporalseriesdata,enablingthecreationofprecisepredictivemo delsandthedefenceofintricaterelationshipsbetweenmultimode inputs.[2]

Globally,airpollutionisakeyfactorinprematuremortality.I t's crucial to understand and anticipate air pollution patternsto lessen its caused damage. Complex prediction algorithmsareneededforthis.InordertounderstandPM2.5patter nsthroughout space and time, deep learning models like GCNNand ConvLSTM are used. Remote sensing satellite photos areutilised to monitor atmospheric pollutant issues while time-series multidimensional directed graphs are employed to char-acterise climate aspects. ConvLSTM is

used to foresee PM2.5in Los Angeles region using data from the preceding 10 daysandground-basedPM2.5sensordata.Thespatiotemporaldeeppredictive algorithms outperform earlier studies significantly.[5]

Intheproposedstudy,ahybridapproachisintroducedfor effectively breaking down single, one-dimensional PM2.5time data into various dimensional time data and extractinghiddeninformationfromit.Eachsequenceispredicted bythealgorithmusingtraditionalpredictiontechniques,andthe finalpredictionsareobtainedbyreconstructingtheresults.Three hybridmodelsW-ANN,W-ARIMA,andW-SVMweredevelopedinthestudytoforesee2016PM2.5trendsinf iveChinesecities.Accordingtotheoutputs,combinedmodelsper formbetterthantraditionalARIMA,ANN,andSVMmodelsatfo recastingshort-termPM2.5concentrations.WhenitcomestoforecastingPM2.5 concentrations,theW-ARIMAmodelperformsbetter,especiallywhenitcomestocaptu ringthemutationalpointsthatcouldhelpcreatepollutionwarning s.[1]

A new model for predicting PM2.5 concentrations for eachhour during heavy haze episodes is presented in the proposedwork. For increased forecasting accuracy, this model com-bines the mode decomposition-recombination method with anensemble learning strategy. The data is divided into severalfrequency modes using the FEEMD to lessen the effects ofnoise in the data. Then, to ensure precise extraction of information and computational efficiency, related modes are combined and excessive breakdown is avoided using the sampleentropy (SE). As a forecasting model,SDEM is built which isstack driven to improve feature representation and informationconsumptioncapabilities.Eachbasemodel'sgener alisationperformance is enhanced via K-fold cross-validation, and themeta-model generates higher prediction results by using eachbasemodel'soutputsasnewinputs.Regardingpredictingpr ecision, stability, and class prediction rate , the FEEMD-SE-SDEM model performs better than earlier contrast models,making it a suitable choice for an early air quality warningsystem.[7]

## III. STUDY AREA

### A. DataCollection

Delhi,thecapitalofIndia,islocatedinthenorthernpartof the nation.It is a big, multiethnic city that has a thrivingmodern culture and history. The city, which has a populationofaround20million,isahubforbusiness,politics,and education. Due to the large population, transportation, urbanisation,industrialization, etc. are all growing quickly. These are someof the factors that contribute to Delhi's high pollution levels,whichcanimpairpeople'shealth.Thereare46livemonitor ingstationsforairqualityinDelhi,oneofwhichIhavechosentobei nAshokVihar.ItisaresidentialneighbourhoodinDelhi's northwest. It is roughly at latitude 28.6958°. Delhi hasa hot semi-arid climate, with extremely hot summers and coolwinters.Thetemperaturecanreachashighas45°C(113°F)in the summer months of May and June, while the wintertemperaturescanfallaslowas2°C(35.6°F)inDecembera

nd January. The city experiences monsoon rains from July toSeptember, which provide relief from the hot summer months.Overall, the climate of Delhi is characterized by its extremetemperaturevariationsandseasonalwinds.



Fig.1.GeographicallocationofAshokVihar

## B.DatasetDescription

The dataset collected from Ashok Vihar contain 1450 datawith11featureswhichincludeconcentrationslikePM2.5,PM10, NO2, NOxetc and meteorological factors like windspeed, humidity, temperature, wind direction etc. The data iscollectedfromyear2019January-2022December.

### IV. PROPOSED METHODOLOGY

The proposed methodoly is a technique which illustrates themodel used in this work step by step. The dataset is collectedfromhttps://cpcb.nic.in/.InthatwebsitefromAshokViharliveairqualitymonitoringstationthedataiscollectedforpast4years.TheworkFowisprovidedbelowasshownin[?].The data is extracted in CSV format. Data preprocessing is atechniquewhichisusedtoprepareabetterdatafortrainingthemodel.Thenullvaluesinthedatasetareidentifedandreplacedwith distribution of each feature. For outlier detection a robustoutlierdetectionalgorithmHampelIdentifierisused.

## A.HampelIdentfier

Toexcludeoutliersfromadataset,onetechniqueistoimplementtheHampelidentifier.Itstartsbyfiguringoutthedata'smedian andmedianabsolutedeviation(MAD).Bycomputingthemedian oftheabsolutedeviationfromthemedian,theMADameasureofhedata'svariabilityisestablished.Eachdatapointisthencompar edtothemedian;iftheabsolutedifferenceexceedsacertainnumb er(oftenamultipleoftheMAD),itisdeemedanoutlierandelimina tedfromthedataset.Instatisticalanalysis,datamining,andmachi nelearningapplications,theHampelIdentifierisuseful,especiall ywhenworkingwithdatasetsthathaveerrorsorextremevalues.

## B. Decomposition

Decompositionistheprocessofbreakingdowntimeseriesin towaveformatusingmathematicaltechniques.Itisusedtoidentif ythehiddentrends,patternandcyclesintimeseriesdata.Thereare differenttypesofdecompositiontechnique,EnsembleEmpirical ModeDecomposition(EEMD)isusedinthismodelfordecompsi ngthePM2.5feature.ThedecomposedPM2.5areshownin figure2.Sincemodemixingmakesitchallengingtodistinguishbe tweenseveraloscillatorymodesinthedata,theEEMDapproachi sanimprovementovertheEMDmethod.Thebreakdownbecome

smorediversewhennoiseisaddedtothedata,makingiteasiertoid entifythemanyoscillatorymodes.Thismethodoffersastrongtoo lforunderstandingtheunderlyingprocessesthatgeneratethedata ,makingitvaluable.ThePM2.5dataisbrokedownintoIMFsandr esidual.TheIMFsarehelpfulincapturingthehighertrendsorhigh erfrequencyandresidualusedtocapturethelowertrendsorlowerf requencyinthemodel.


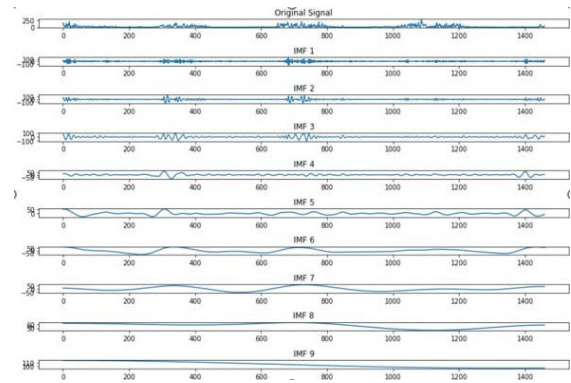
Fig.2.DecompositionofPM2.5

## C.LSTM

TheissueoftypicalRNNs'vanishingorexpandinggradients isresolvedbytheRNNarchitectureapproachknownasLSTM.B ecauseitcanidentifylong-termdependenciesinthedata. LSTM is useful for time series data. LSTMs use memorycells to store information over time and gates made of sigmoidandtanhfunctionstoregulateinformationfow.Thesegat esincreasethemodel'sabilitytorepresentcomplexsequencesby allowingittochooseretainorrejectinformationbasedontheinput data.Overall,LSTMhasbeenextensivelyappliedinmanydiffere ntapplicationsandisastrongtoolformodellingsequentialdata.B eforepassingthedatatothelstmodel,thedataisscaledwithfeature rangeof(o,1).LongShort-TermMemory(LSTM)networksusethreedifferenttypesofgate s:input,output,andforgetgatesasshowninfgure3.InorderforLS TMstomanagelong-termdependenciesinsequentialdata,thesegatesarecrucialbeca usetheyregulatethevowofinformationintoandoutofthememor ycell.
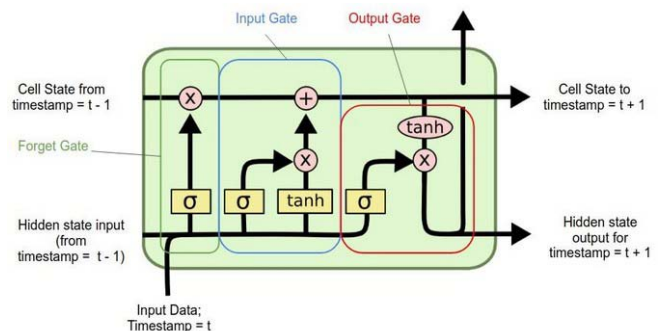


Fig.3.ArchitectureofLSTM

## D.ModellingProcessofHybridLstmApproach

Overall,thisworkoffersahybridLSTM(HI-EEMD-LSTM)modelforthedailyPM2.5forecasting,followingthetech

nicalpathof"decompositionandLSTM."Themodellingprocess isdepictedinfigure4,anditmaybesummedupasfollows:

Step-1:Gethistoricalinformation,suchasPM2.5concentrationsandweatherinformation.

Step-2:ToobtainnnumberofIMFsandoneresiduemode,performtheEEMD

Step-3:Standardizetheresultingdecomposedwavesinthefeaturerangefromthemeteorologicaldatafrom0to1

Step-4:Createa(n+1)LSTMmodelwitheachIMFacquiredfromdecompositionasthetrainingdatafortheweather.

Step-5:ToacquiretheoriginalPM2.5concentrations,dothe inverse EEMD to produce the anticipated PM2.5 concentrations
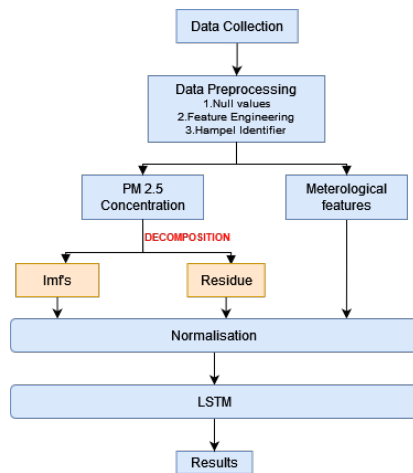
Step-6:Producetheforecastingexercise'sfinaloutput.



Fig.4.ArchitectureofHybridLSTMmodel

## V. EVALUATION METRICS

Many evaluation metrics root mean squared error (RMSE),mean absolute error (MAE), and mean squared error (MSE)areusedtoidentifythecapabilityofthemodel.Theevaluationmetrics of the model is compared with basic LSTM model toidentify how much better the proposed model is executingbetterthantheLSTMmodel

## VI. RESULTS

Thestudy's findingsdemonstratedthatthehybridLSTMmodelgreatlyperformedbetterthantheLSTMmodelsinpredictingPM2.5valuesbasedonmeteorologicaldata.TheMAEofthebestLSTMmodelwas 42.20,whichwasasignificantimprovementoverthebaselineMAEof91.83.Similarly,therootmeansquarederror(RMSE)wasalsosignificantlylowerfortheLSTMmodelscomparedtothebaselinemodels.TheLSTMmodelswereparticularlyeffectiveincapturingthelong-termdependenciesinthedata,whichhelpedtoimproveprediction

ns.Thefigure5andfigure6givesagraphicalrepresentationofboththemodelswithactualandpredictedvaluesThese findingssuggestthatLSTMmodelsmaybeausefultoolforpredictingstockprices,andcouldhaveimportant applications in finance and investment. However, itisimportanttonotethatthereweresomehindrancetothestudy,includingarelativelylittlesamplesizeandalimitedtimeperiodfortheanalysis.Furtherresearchisneededtovalidatethese findingsandexplorethepotentialofLSTMmodelsforotherapplicationsin financeandbeyond.

TABLEI: MODEL PERFORMANCE

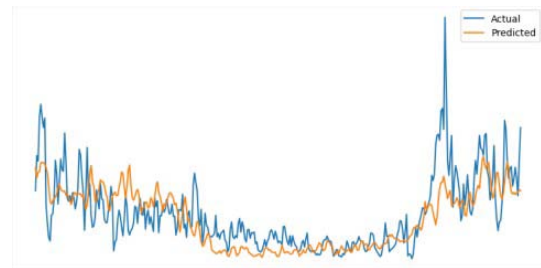| Model | EvaluationMetrics | | |
|---|---|---|---|
| | MAE | MSE | RMSE |
| LSTM | 91.83 | 12161.98 | 110 |
| HI-EEMD-LSTM | 42.2 | 3997.97 | 63.22 |



Fig.5.LSTMmodel

## VII. CONCLUSION

Inthisstudy,aninnovativeHybridLSTMmodelisintroducedforpredictingdailyPM2.5concentrationsbyutilisingmultiple LSTMmodelsindifferentmodes.ThemodelfirstappliesEEMDtotransformthetimeseriesfusedfeaturesintosimplerfeaturesinmono-mode.Thematchingbetween
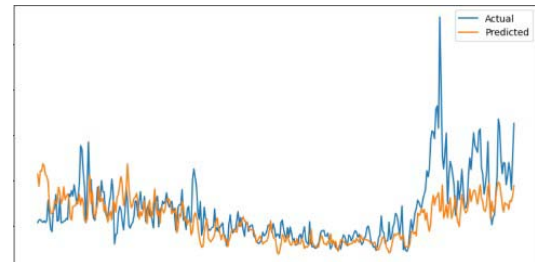


Fig.6.HybridLSTMmodel

meteorologicalconditionsandeachmode'scoefficientsisthenconstructed using LSTM, resulting in the development of agroup of LSTM models for ensemble learning. The predictedmode coefficients are then incorporated into the results usingthe reverse EEMD. Using wave decomposition, the H-LSTMmodelhasincreasedtheaccuracyofPM2.5observation. WhenthesuggestedH-LSTMmodeliscomparedtotheLSTMmodel on the Ashok Vihar dataset, it is evident that the H-LSTMmodeloutperformsbothofthosemethods.

## REFERENCES

[1] YongCheng,HongZhang,ZhenhaiLiu,LongfeiChen,andPingWang, "Hybrid algorithm for short-term forecasting of pm2. 5 in china,"Atmosphericenvironment,vol. 200, pp. 264–279,2019.

[2] YuanlinGu,BaihuaLi,andQinggangMeng, "Hybridinterpretablepredictivemachinelearningmodelforairpollutionprediction," Neurocomputing,vol. 468, pp. 123–136,2022.

[3] ChiouJyeHuang,andPing-HuanKuo,"Adeepcnn-lstmmodelforparticulatematter(pm2.5)forecastinginsmartcities," Sensors,vol. 18, no. 7, p. 2220,2018.

[4] Sitharthan, R., Vimal, S., Verma, A., Karthikeyan, M., Dhanabalan, S. S., Prabaharan, N., ...&Eswaran, T. (2023). Smart microgrid with the internet of things for adequate energy management and analysis.Computers and Electrical Engineering, 106, 108556..

[5] PratyushMuthukumar,EmmanuelCocom,KabirNagrecha,DawnComer, IreneBurga,JeremyTaub,ChisatoFukudaCalvert,JeanneHolm,andMohammadPourhomayoun," Predictingpm2.5atmosphericairpollutionusingdeeplearningwithmeteorologicaldataandground-basedobservationsandremote-sensingsatellitebigdata," AirQuality,Atmosphere&Health,vol. 15, no. 7, pp. 1221–1234,2022.

[6] MihaelaOprea,SandaFlorentinaMihalache,andMarianPopescu, "Acomparativestudyofcomputationalintelligencetechniquesappliedtopm2.5airpollutionforecasting,"In20166thInternationalConferenceonComputersCommunicationsandControl(ICCCC),IEEE, pp.103–108,2016.

[7] Moshika, A., Thirumaran, M., Natarajan, B., Andal, K., Sambasivam, G., &Manoharan, R. (2021).Vulnerability assessment in heterogeneous web environment using probabilistic arithmetic automata. IEEE Access, 9, 74659-74673.

[8] Nur'atiahZaini,LeeWoenEan,AliNajahAhmed,andMarlindaAbdulMalek," Asystematicliteraturereviewofdeeplearningneuralnetworkfortimeseriesairqualityforecasting," EnvironmentalScienceandPollutionResearch,pp.1–33,2021.

[9] GuangjiZheng,HuiLiu,ChengqingYu,YeLi,andZijieCao, "Anewpm2.5forecastingmodelbasedondatapreprocessing,reinforcementlearningandgatedrecurrentunitnetwork," AtmosphericPollutionResearch,vol. 13, no. 7, p. 101475,2022.

[10] Dixian Zhu, ChangjieCai, Tianbao Yang, and Xun Zhou, "A machinelearning approach for air quality prediction: Model regularization andoptimization," Bigdataandcognitivecomputing,vol. 2, no. 1, p. 5,2018.