An aerial photograph of a city, likely Seoul, showing a large river (Han River) with several bridges. The city is densely packed with high-rise apartment buildings. The image is partially obscured by a black silhouette of a person's head and shoulders on the right side.

# Enhancing Streamflow Prediction in Vulnerable Regions through Probabilistic Deep Learning and Satellite-Derived Data

Fatemeh Ghobadi<sup>1</sup>, Amir Saman Tayerani Charmchi<sup>1</sup>,  
JungMin Lee<sup>2</sup>, Myeong In Kim<sup>2</sup>, and Kichan jung<sup>1</sup>

<sup>1</sup> Department of Digital Convergence, Onpoom Co., LTD, Seoul, South Korea

<sup>2</sup> Land and Housing Research Institute, Daejeon, 34047, Korea



# Contents

1. [Introduction](#)
2. [Materials and methods](#)
  - [Proposed framework](#)
3. [Results](#)
4. [Conclusion](#)



# Problem statement

## Streamflow forecasting incorporates various parameters

Meteorological data, Hydrological models, and Historical data;

## Importance of streamflow forecasting

Flood and water resource management, environmental monitoring, climate change adaptation, and informed decision making;

## Challenges in poorly gauged basins

Data scarcity, climate variability and change, human intervention, and inadequate infrastructure;

## Innovative approaches to overcome data limitations and improve forecasting accuracy

Geo-spatiotemporal models, mesoscale data, Attention-based networks





# Improve forecasting accuracy

## Geo-spatiotemporal features

Meteorological data (such as temperature, precipitation, humidity, and wind speed);

Geographic locations, time stamps, and related attributes

## Data Preprocessing

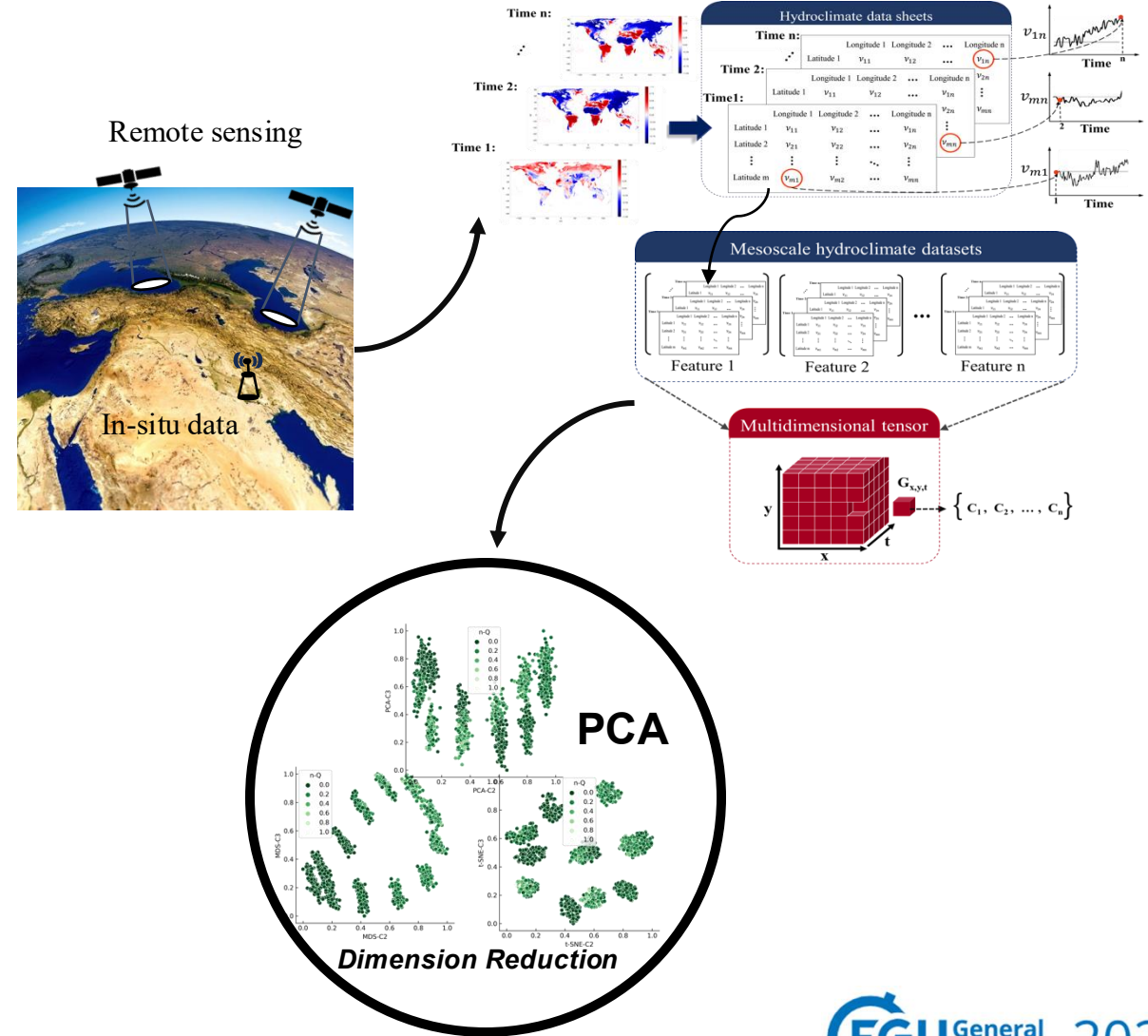
Cleaning the data, normalizing the data, and dealing with missing data

## Feature engineering

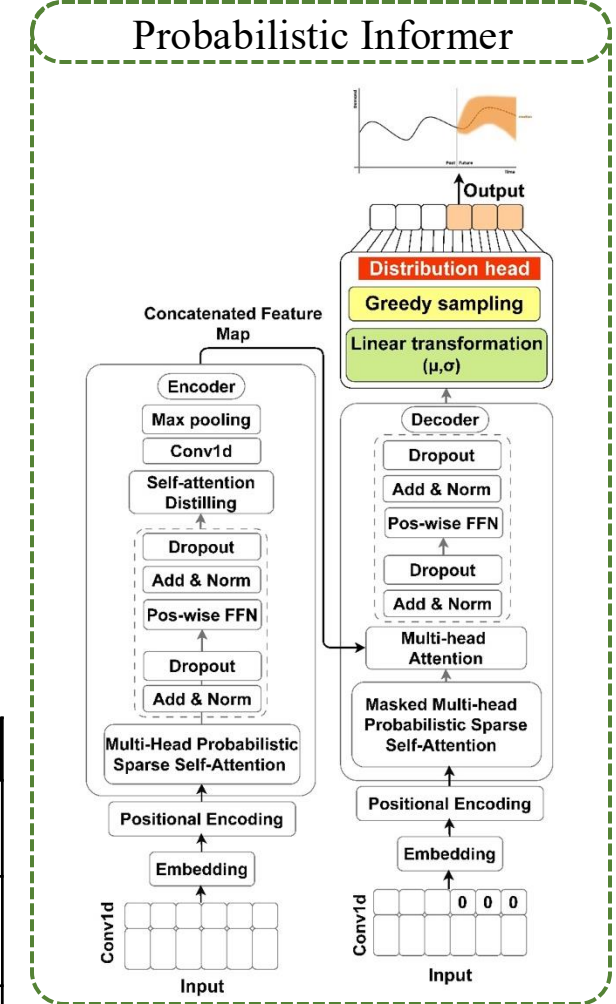
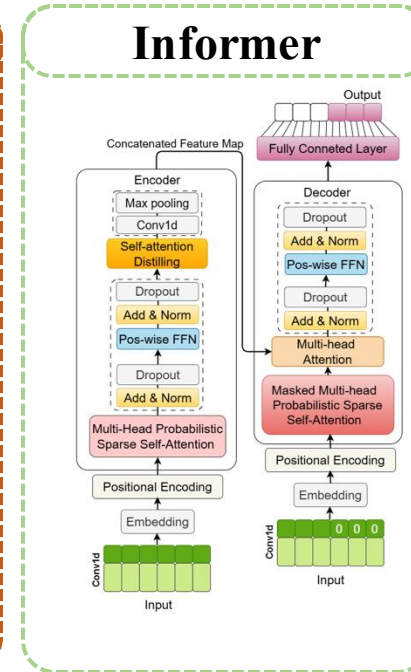
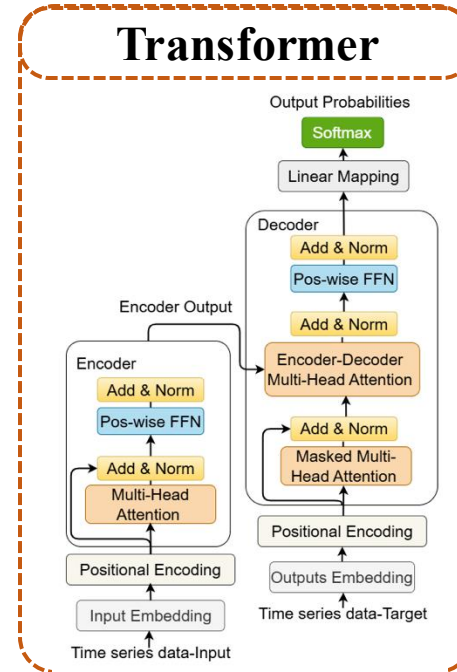
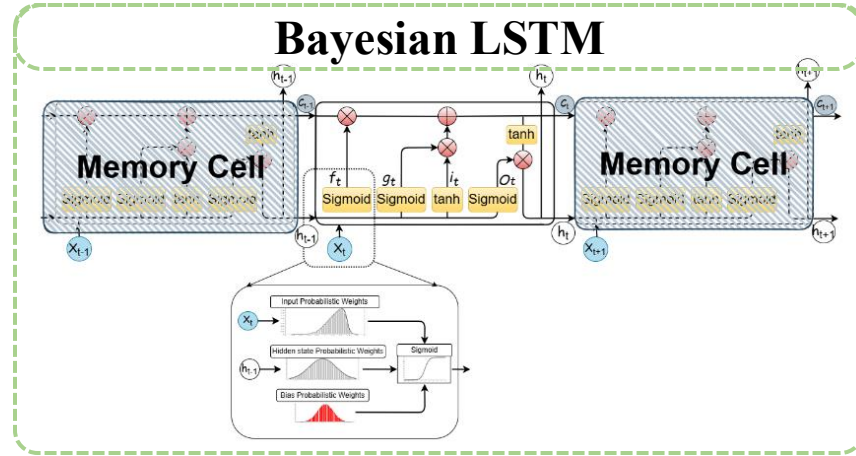
Cascade dimensionality reduction for feature extraction, Cross-correlation analysis, Feature selection

## Model development

Advanced algorithms, optimization, and validation for geo-spatiotemporal modeling.

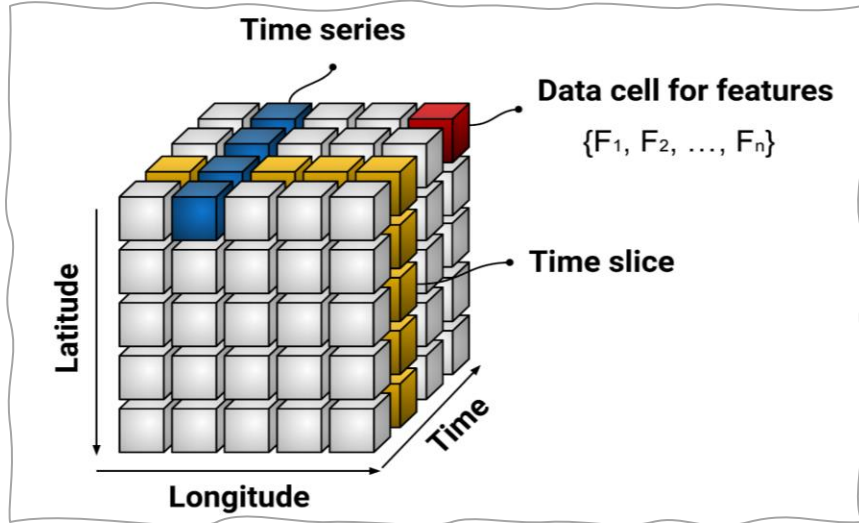


## Attention-based architecture

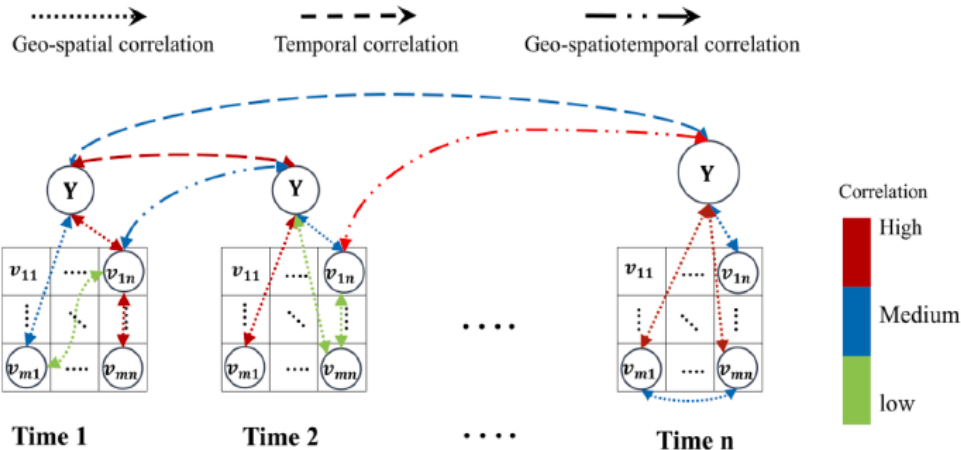
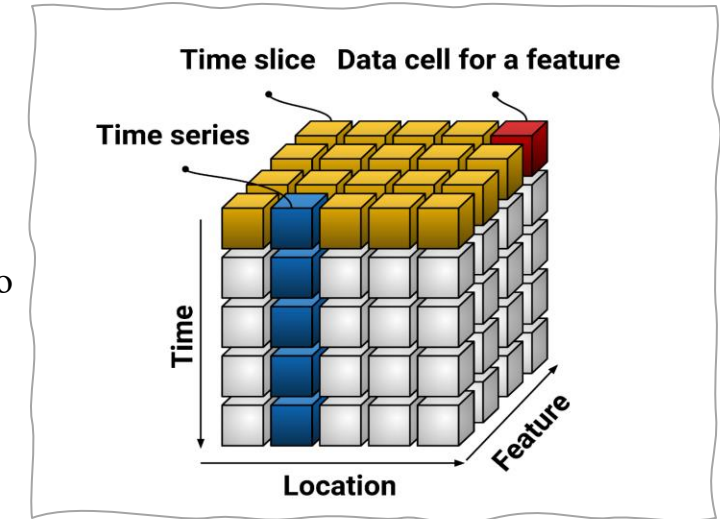


## Three attention-based networks

Model	Input Encoding	Attention Style	Strengths	Limitations
<b>SageFormer</b>	GNN (label encoding)	Spatial + Temporal	Graph-based spatiotemporal learning	Needs graph structure, more complex
<b>Informer</b>	CNN + Time features	ProbSparse Attention	Efficient, long-sequence forecasting	Less spatial modeling, heuristic pruning
<b>Transformer</b>	Linear + Positional	Full Self-Attention	General, deep temporal dependencies	$O(L^2)$ cost, no feature relation modeling



Convert high-dimensional geo-spatiotemporal data into lower-dimensional representations (Spatiotemporal Space ) that still capture essential variability.

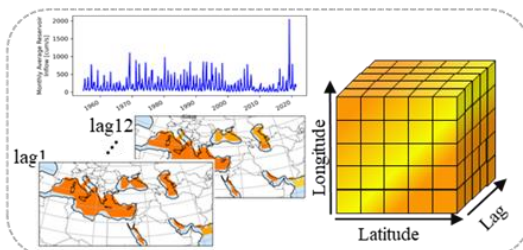


Visualization of interdependencies among geo-spatiotemporal data, highlighting the complex relationships and interactions between dimensions

## Why High Dimensionality is Problematic?

- **Computational Complexity:**  
High-dimensional datasets significantly increase computation cost and processing time.
- **Overfitting Risk:**  
With many features (spatial grid points over time), machine learning models can become prone to overfitting, degrading prediction accuracy.
- **Redundant Information:**  
Satellite-derived data often contains redundant information (spatial correlation) across adjacent grid points.
- **Curse of Dimensionality:**  
As dimensions grow, data points become sparse, causing performance deterioration and reducing statistical significance

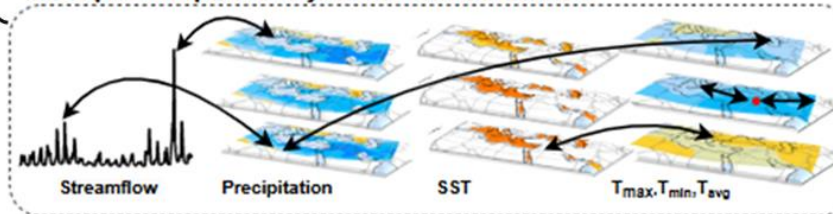




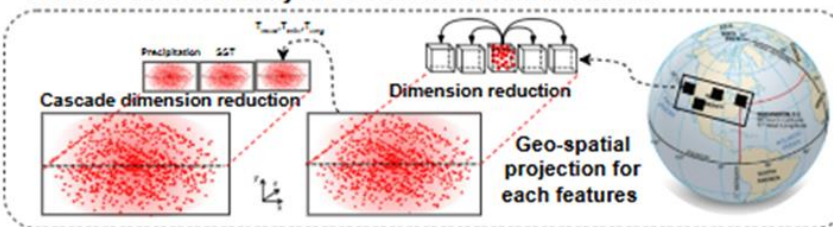
Data fusion for daily flood forecasting



Geo-spatiotemporal analysis for feature selection

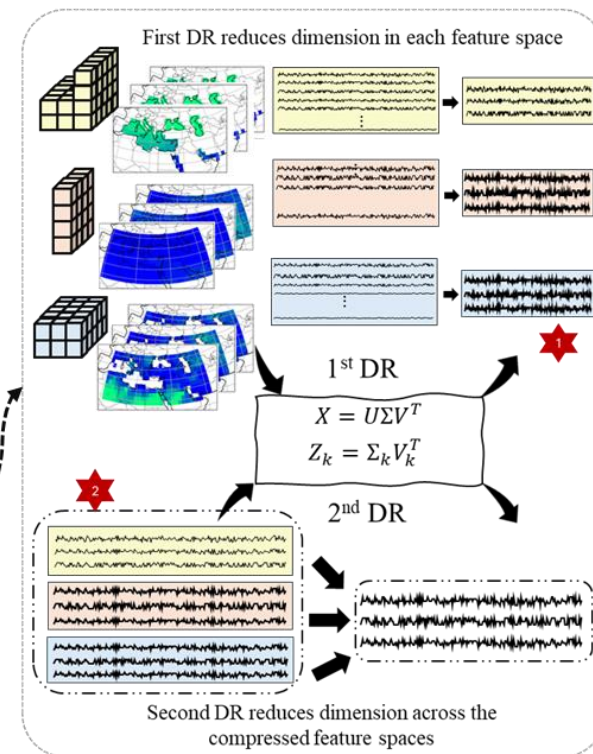
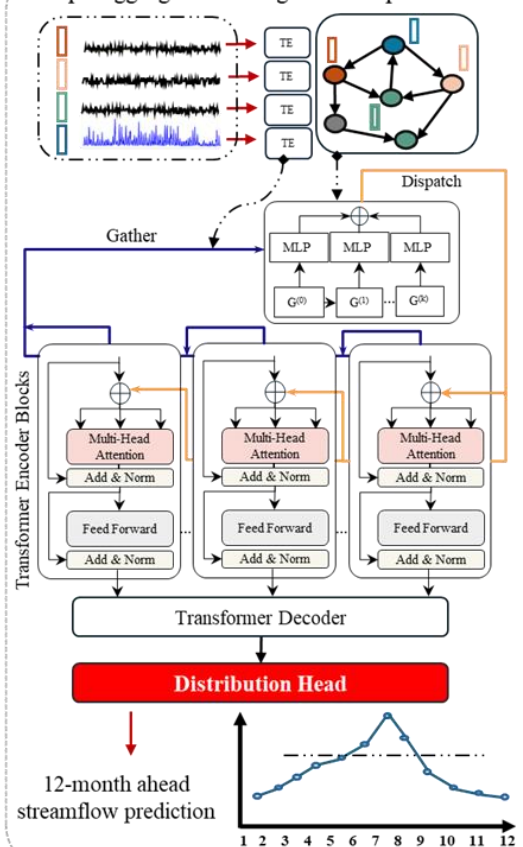


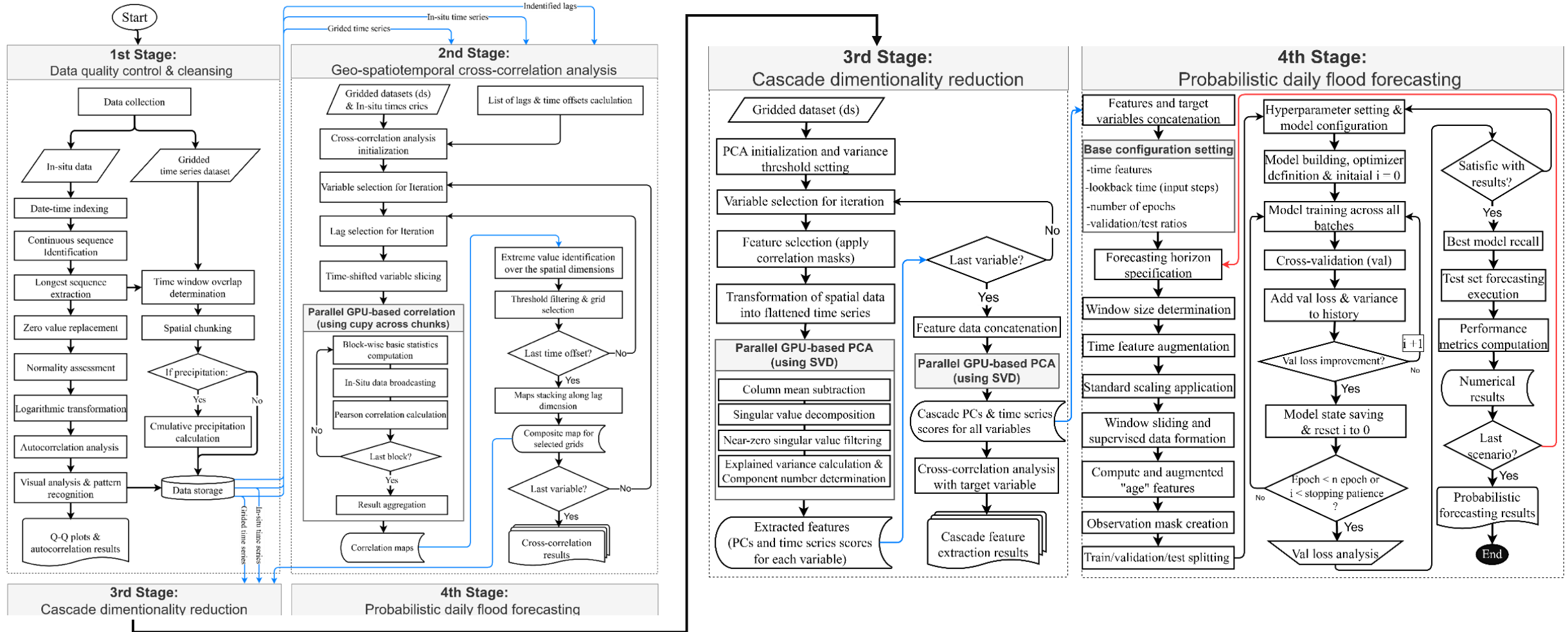
Cascade dimensionality reduction for feature extraction



Probabilistic deep learning for multi-step ahead prediction

Graph aggregation using multi-hope GNN





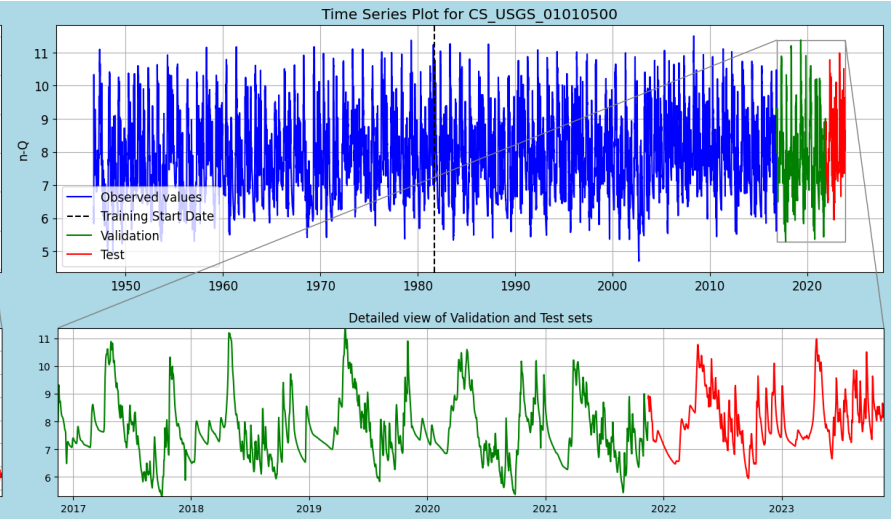
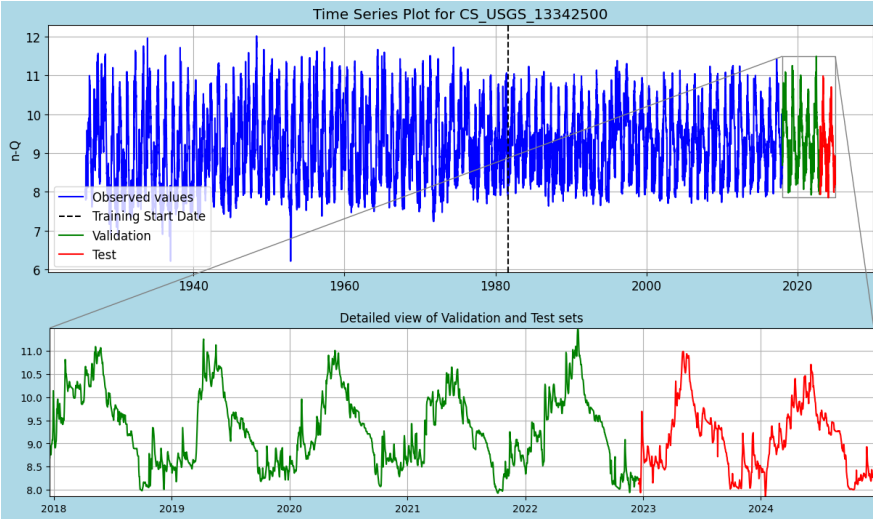


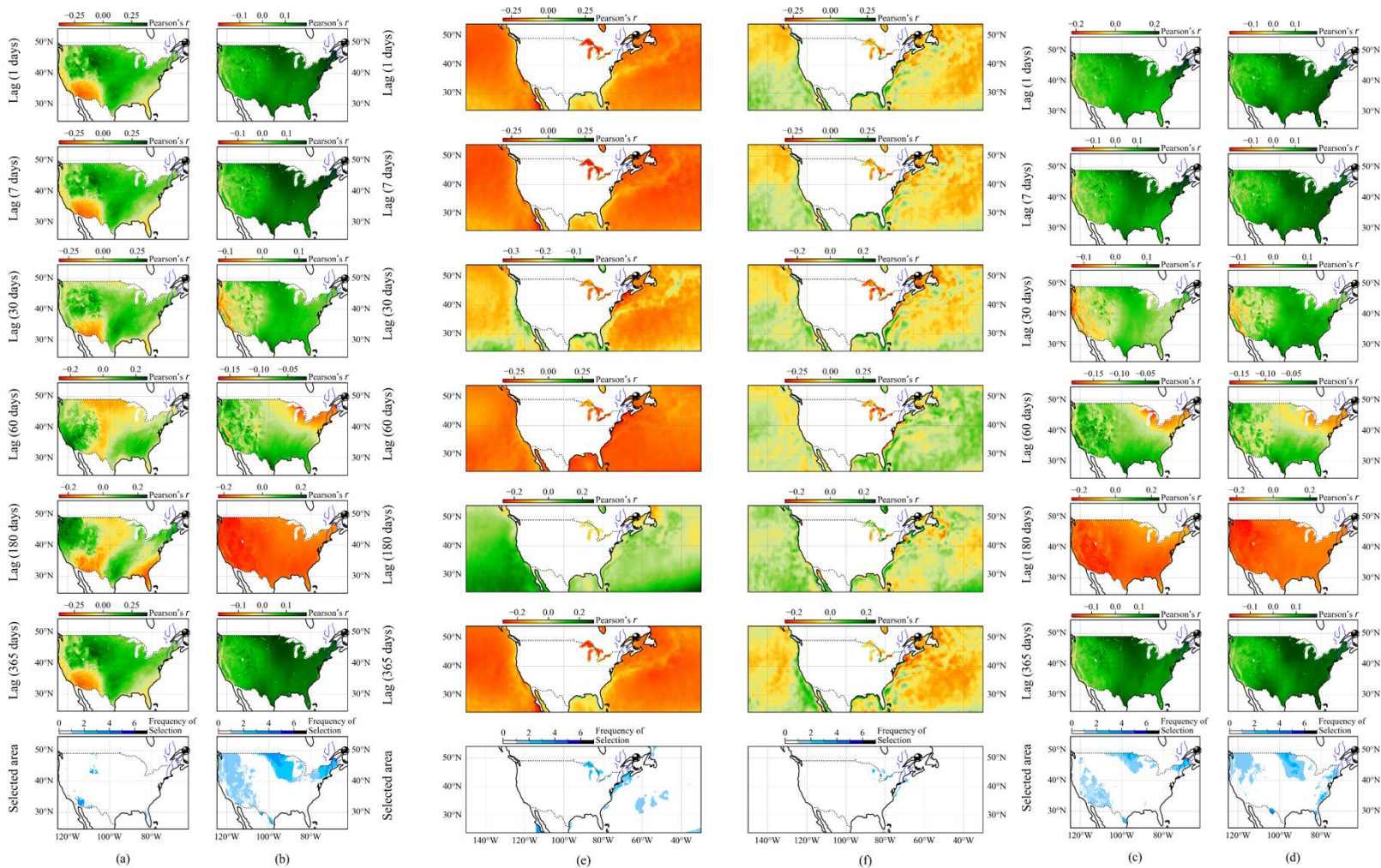


Mesoscale hydroclimate datasets

- Optimum Interpolation Sea Surface Temperature
- NOAA nClimGrid-Daily Version 1

In Situ timeseries





Number of selected grids for each feature in each case study, along with the number of features extracted for each variable for initial and cascaded PCA

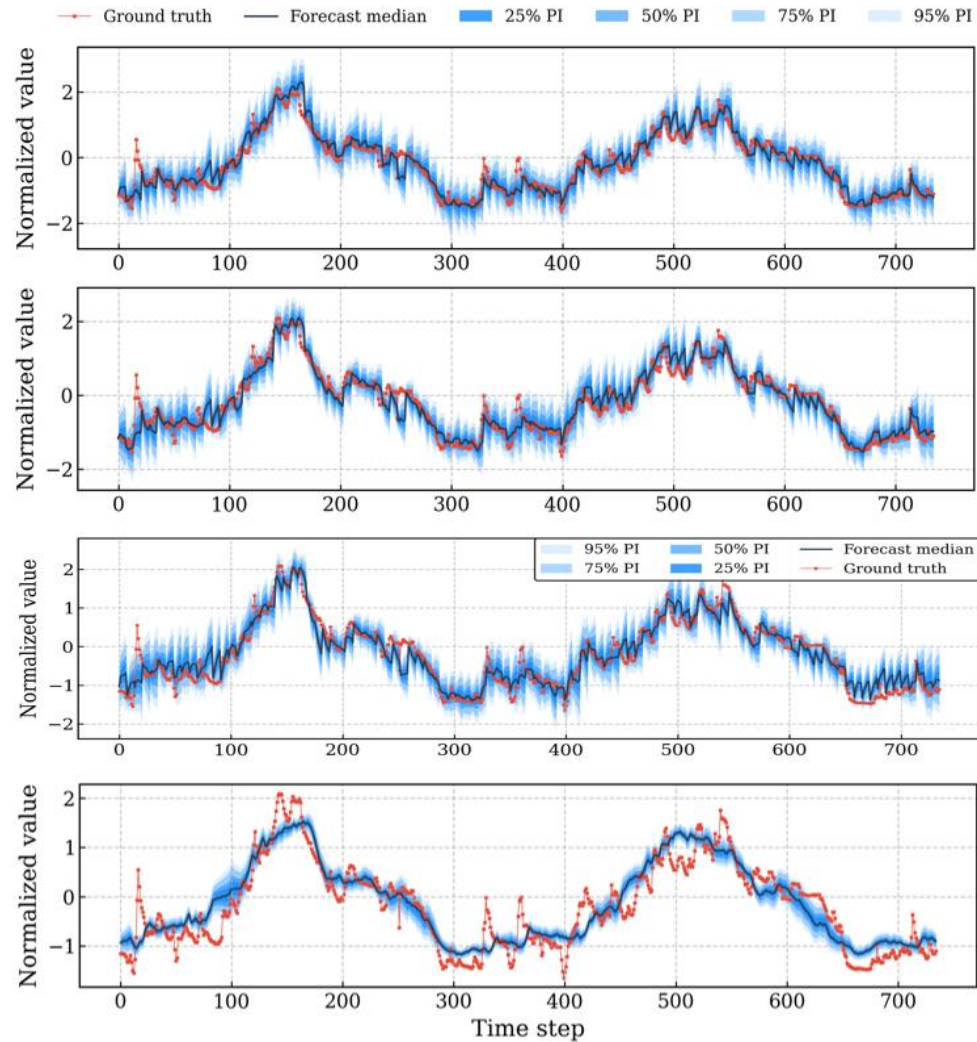
Case study	Features	No. of selected grids	No. of principal components	
			Initial PCA	Cascaded PCA
USGS-1010500	Precip	120	3	
	T <sub>avg</sub>	63	3	
	T <sub>min</sub>	100	3	21
	T <sub>max</sub>	50	3	
	SST	230	3	
	Anomaly SST	19	19	
USGS-13342500	Precip	400	3	
	T <sub>avg</sub>	100	4	
	T <sub>min</sub>	200	4	22
	T <sub>max</sub>	120	5	
	SST	150	5	
	Anomaly SST	130	15	

Geo-spatiotemporal correlation between streamflow and (a) Precip, (b) Tavg, (c) Tmin, (d) Tmax, (e) sea surface temperature (SST), and (f) anomaly SST.

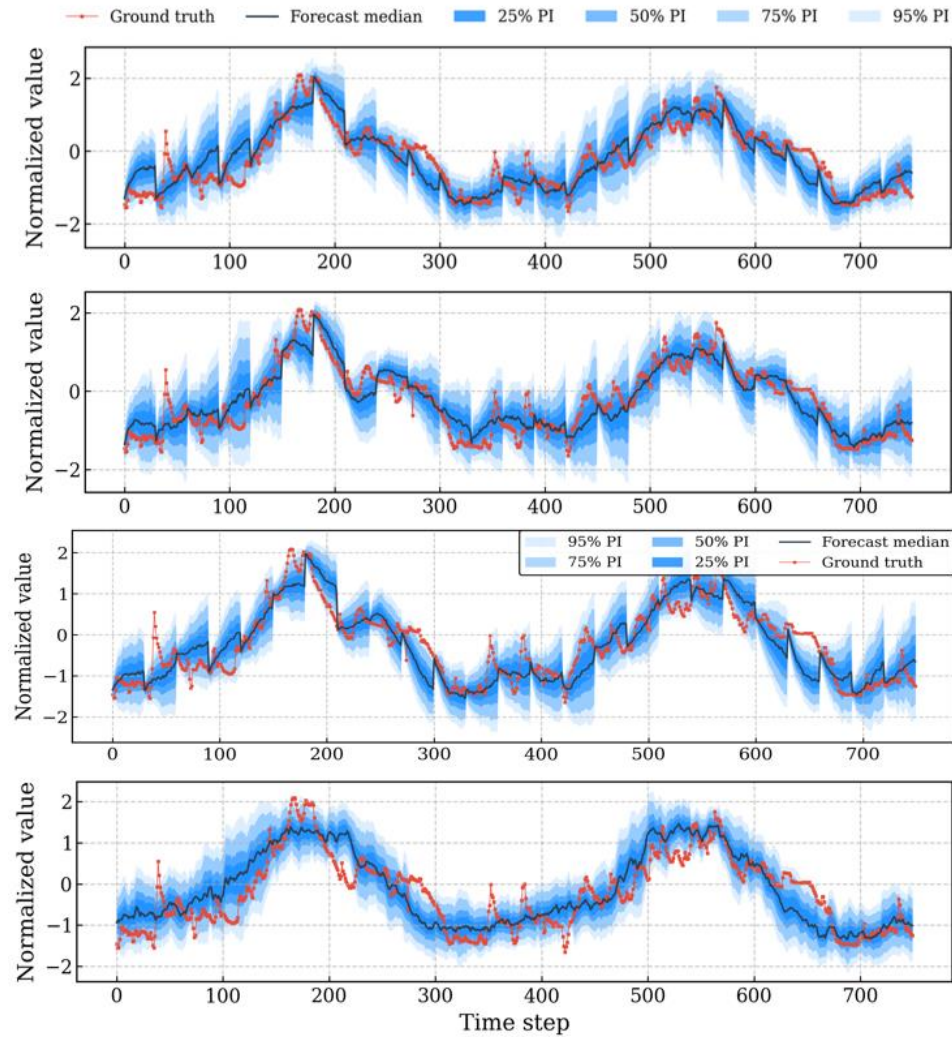
The red areas represent a negative correlation, and the green areas represent a positive correlation.



# Two-year test dataset, Case study 1

[Introduction](#)[Method](#)[Results](#)

7-day forecasting



30-day forecasting

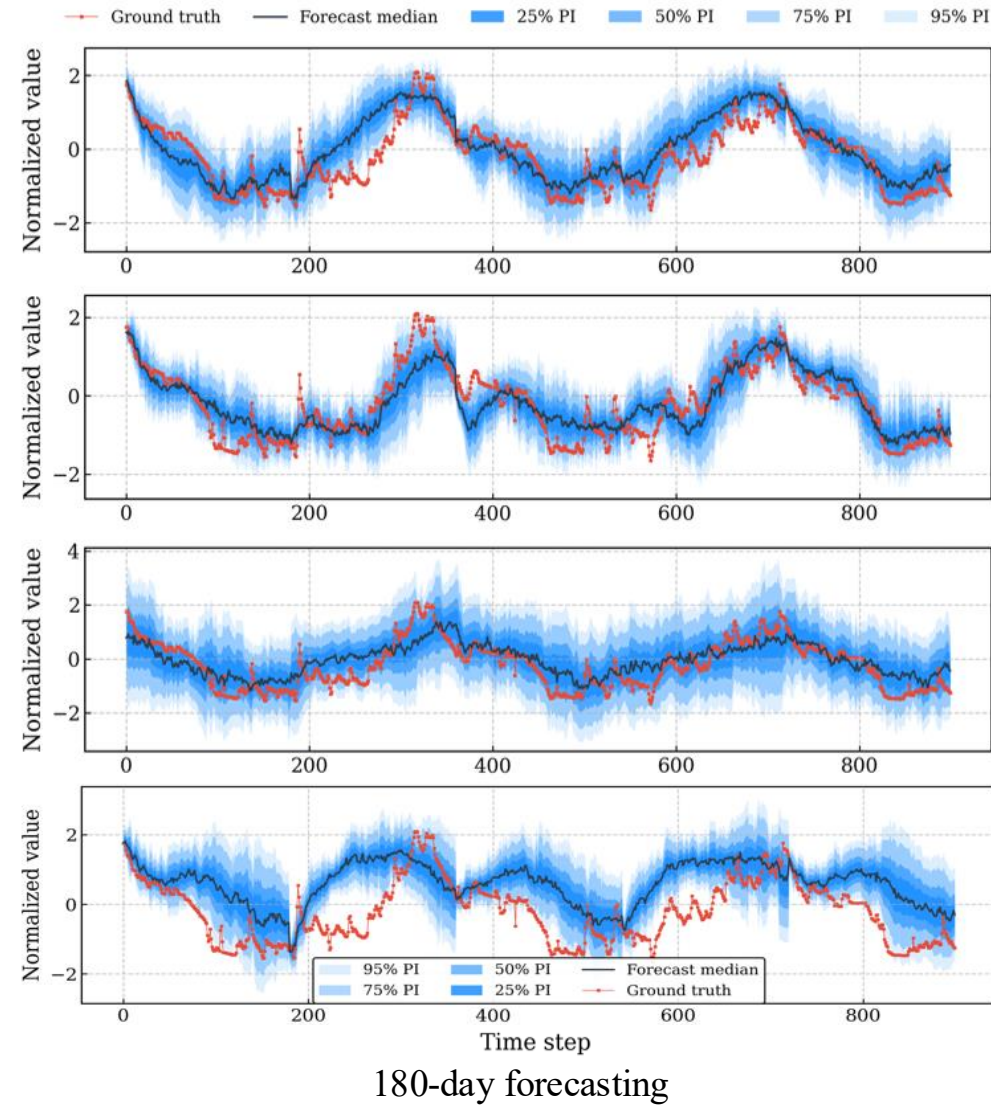
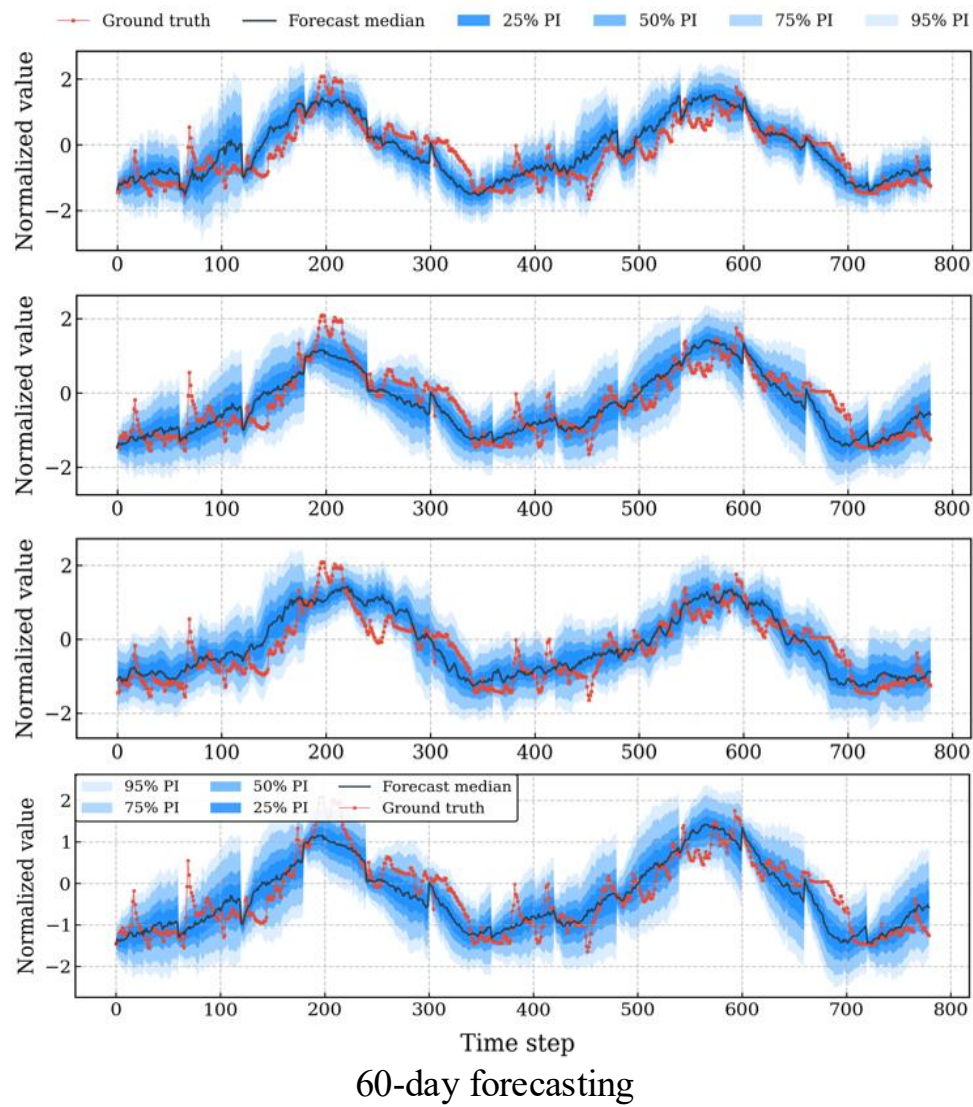
SageFormer

Transformer

Informer

PLSTM

# Two-year test dataset, Case study 1

[Introduction](#)[Method](#)[Results](#)

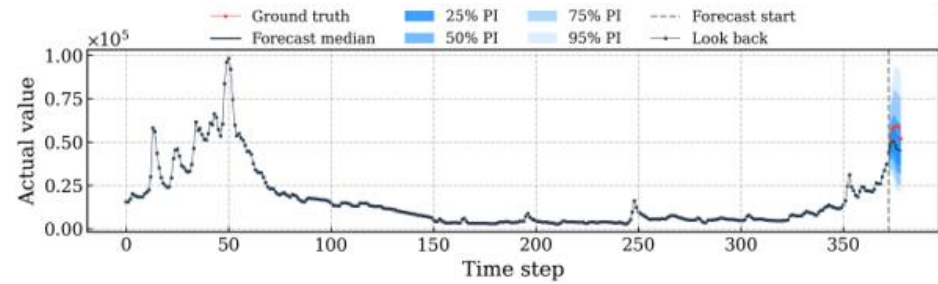
SageFormer

Transformer

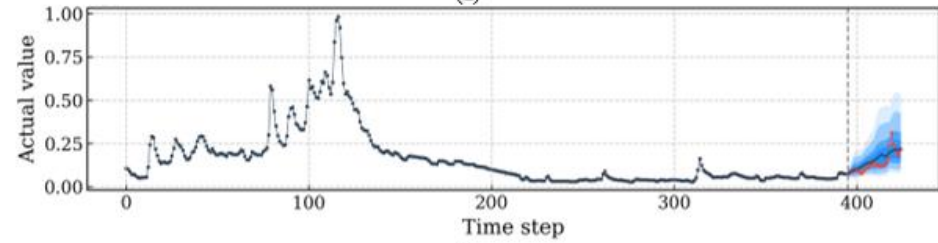
Informer

PLSTM

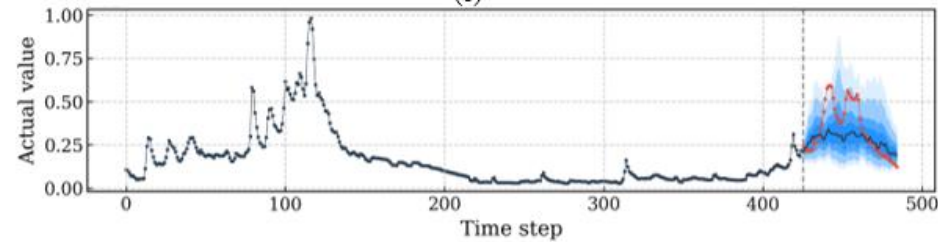




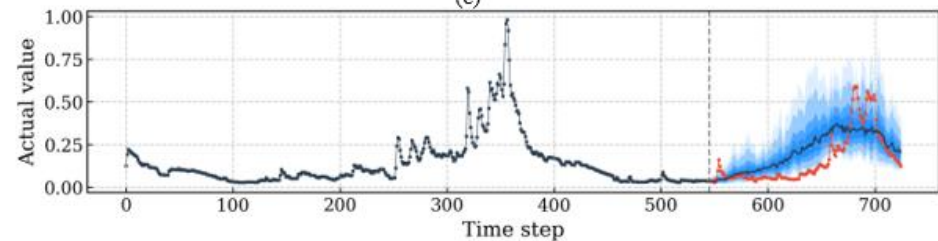
(a)



(b)



(c)



(d)

(a)

(b)

(c)

(d)

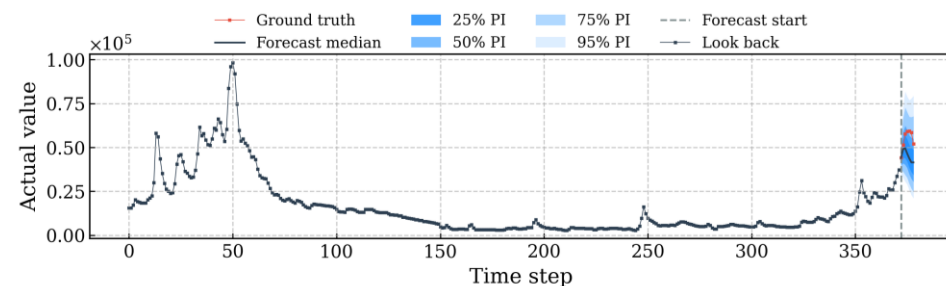
Back-tested probabilistic daily flood forecasting results for separately trained models of **SageFormer network** with different look-back periods corresponding to forecast horizons:

7-day forecast (372-day look-back)

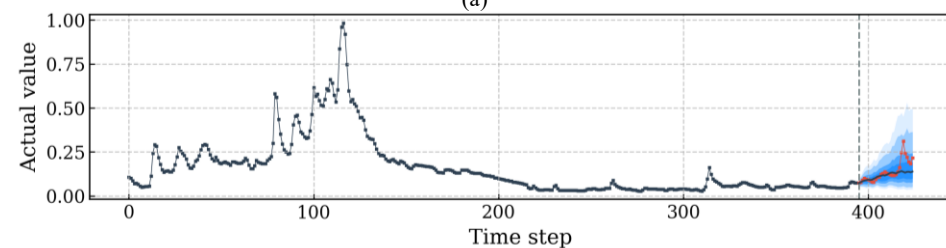
30-day forecast (395-day look-back)

60-day forecast (425-day look-back)

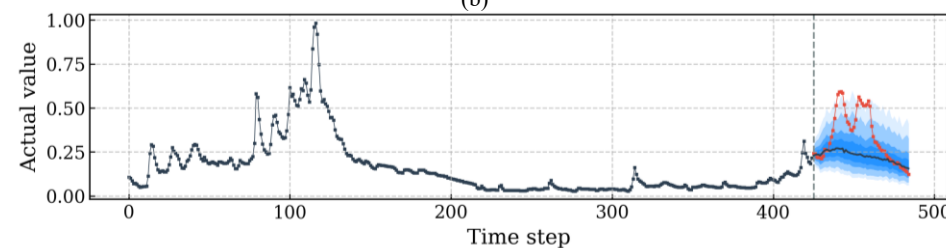
180-day forecast (545-day look-back)



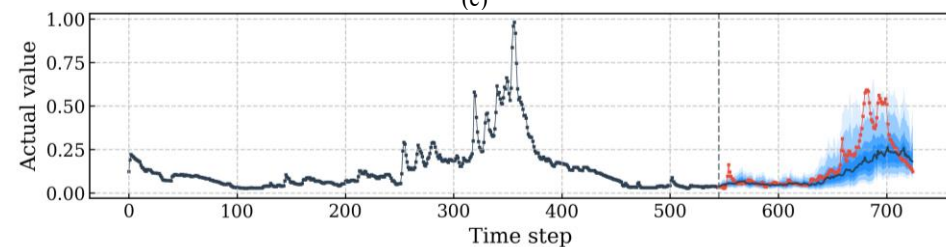
(a)



(b)



(c)



(d)

(a)

(b)

(c)

(d)

Back-tested probabilistic daily flood forecasting results for separately trained models of **Transformer network** with different look-back periods corresponding to forecast horizons:

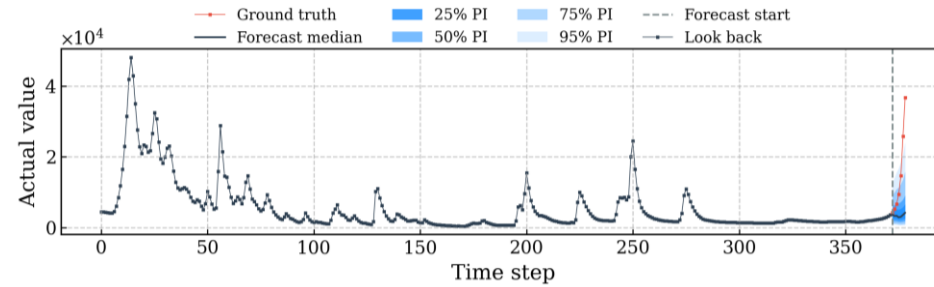
7-day forecast (372-day look-back)

30-day forecast (395-day look-back)

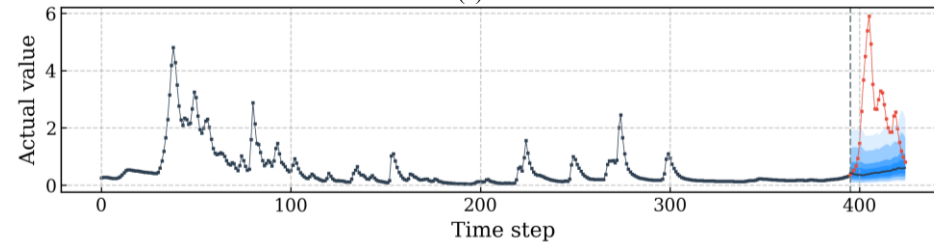
60-day forecast (425-day look-back)

180-day forecast (545-day look-back)

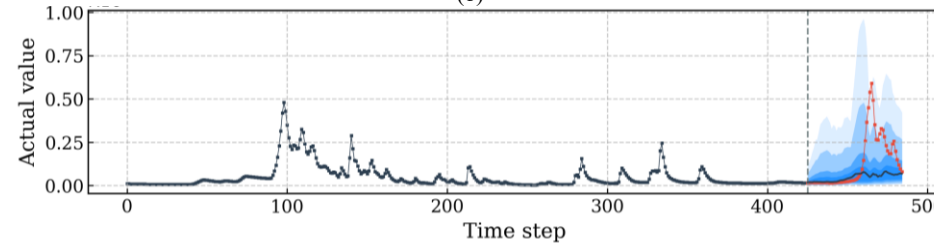




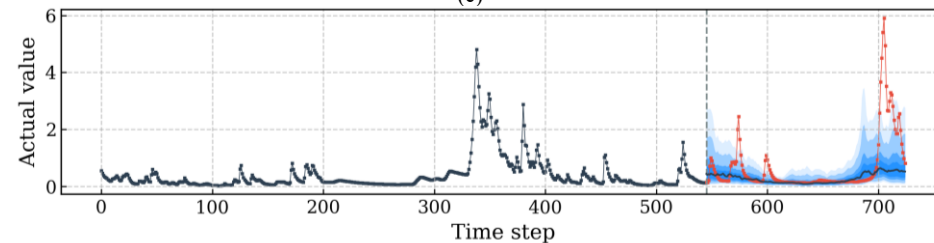
(a)



(b)



(c)



(d)

(a)

(b)

(c)

(d)

Back-tested probabilistic daily flood forecasting results for separately trained models of **PLSTM network** with different look-back periods corresponding to forecast horizons:

7-day forecast (372-day look-back)

30-day forecast (395-day look-back)

60-day forecast (425-day look-back)

180-day forecast (545-day look-back)

# Conclusion

- **SageFormer** exhibited superior performance in terms of forecast **sharpness** (lowest MPIW) and **overall accuracy** (lowest CRPS) across all forecasting horizons, emphasizing its capability to reliably predict flood events even months in advance.
- Informer offered significant computational advantages, achieving accuracy comparable to Transformer but with reduced complexity, making it particularly suitable for resource-constrained operational contexts.
- Attention-based models effectively calibrated their uncertainty estimates, as indicated by PICP values closely matching nominal prediction intervals (75–95%), thus offering **valuable decision-making support during extreme hydrological events**.
- PLSTM, despite generating wide prediction intervals, consistently underperformed in capturing critical peaks and failed to provide precise uncertainty quantification, demonstrating inherent limitations of recurrent architectures in extended probabilistic forecasting scenarios.
- SageFormer, Transformer and Informer maintained stable performance across increasing forecast horizons, likely benefiting from their autoregressive architectures and progressively enriched historical context, underscoring the importance of attention mechanisms for managing long-range hydrological dependencies.



# 감사합니다

# THANK YOU FOR YOUR ATTENTION

Fatemeh Ghobadi    ✉ [Ghobadi@khu.ac.kr](mailto:Ghobadi@khu.ac.kr)