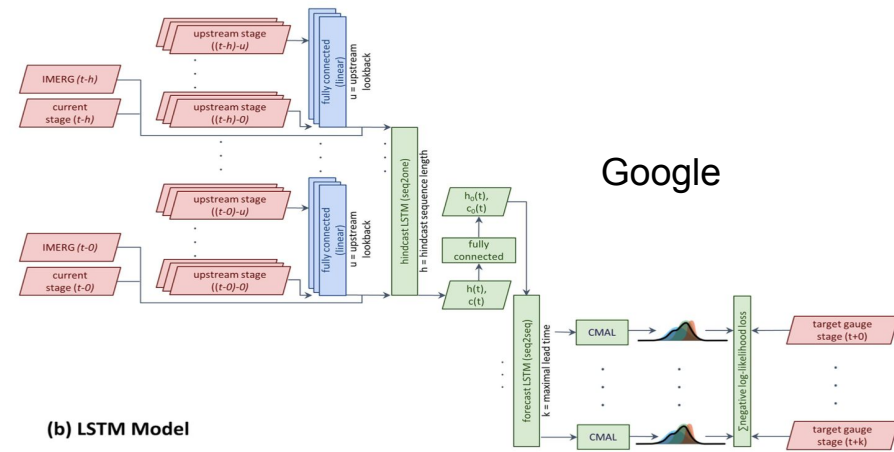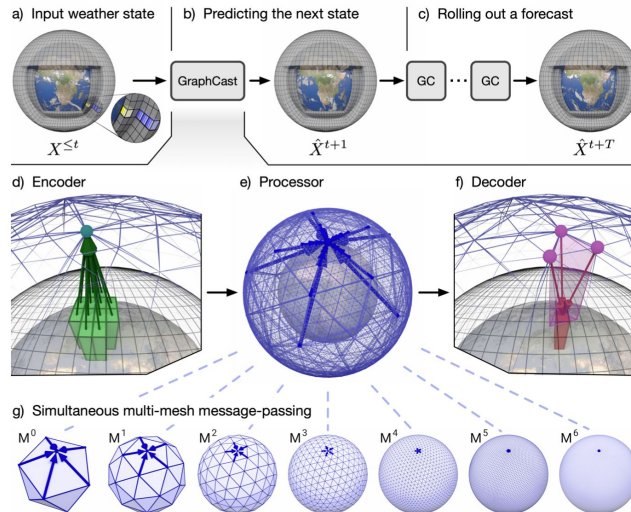- Google's forecasting system consists of four subsystems: data validation, stage forecasting, inundation modeling, and alert distribution;  stage forecasting is modeled with the long short-term memory (LSTM) networks and the linear models
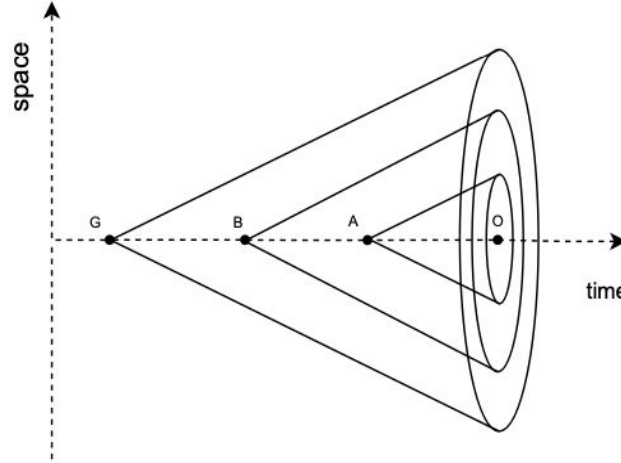
- GraphCast capitalises on GNN's ability to model arbitrary sparse interactions by introducing internal multi-mesh representation, which has homogeneous spatial resolution over the globe, and allows long-range interactions within few message-passing steps

- In certain domains, incorporating domain-specific knowledge into the input features can improve performance. For example, in computer vision tasks, handcrafted features like edges or texture descriptors can provide additional information to the neural network
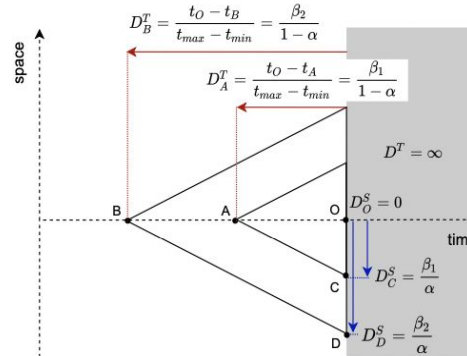


Google

(b) LSTM Model



GraphCast

- Most predictive spatio-temporal machine learning models transform the original problem into a multiple regression task, where the target variable is the future value of the series and the predictors are previous past values of the series up to a certain time window

- Assumes that the future conditions depend on recently observed conditions in the same location

- Certain spatio-temporal variables of interest depend not only on the recent past conditions at the same location, but also on recent past conditions of nearby locations

- The "spatial radius of influence" may decrease for older observations such that observations that are more proximate in both time and location have a larger influence on the outcome variable than temporally and spatially distance observations.
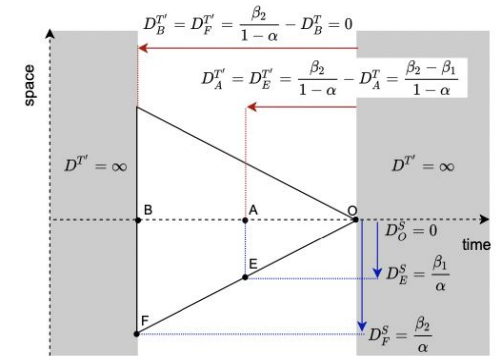
$$D_A^T = \begin{cases} \frac{t_O - t_A}{t_{max} - t_{min}}, & \text{if } t_A < t_O \\ \infty, & \text{otherwise} \end{cases}$$

$$D_A = \alpha \cdot D_A^S + (1 - \alpha) \cdot D_A^T \le \beta$$

$$\mathcal{N}_o^\beta = \{k \in \mathcal{D} : D_{o,k} < \beta\}$$



$$D_B^T = \frac{t_O - t_B}{t_{max} - t_{min}} = \frac{\beta_2}{1 - \alpha}$$

$$D_A^T = \frac{t_O - t_A}{t_{max} - t_{min}} = \frac{\beta_1}{1 - \alpha}$$

$D^T = \infty$

$D_O^S = 0$

$D_C^S = \frac{\beta_1}{\alpha}$

$D_D^S = \frac{\beta_2}{\alpha}$

(a) Original cone

$$D_B^{T'} = D_F^{T'} = \frac{\beta_2}{1 - \alpha} - D_B^T = 0$$

$$D_A^{T'} = D_E^{T'} = \frac{\beta_2}{1 - \alpha} - D_A^T = \frac{\beta_2 - \beta_1}{1 - \alpha}$$

$D^{T'} = \infty$

$D^{T'} = \infty$

$D_O^S = 0$

$D_E^S = \frac{\beta_1}{\alpha}$

$D_F^S = \frac{\beta_2}{\alpha}$

(b) Reversed cone

- Commonly used performance estimation procedures such as cross-validation (CV) and out-of-sample (OOS) validation face challenges due to the implicit dependence between observations in spatiotemporal datasets

- Standard cross-validation leads to over-optimistic estimates in spatio-temporal settings

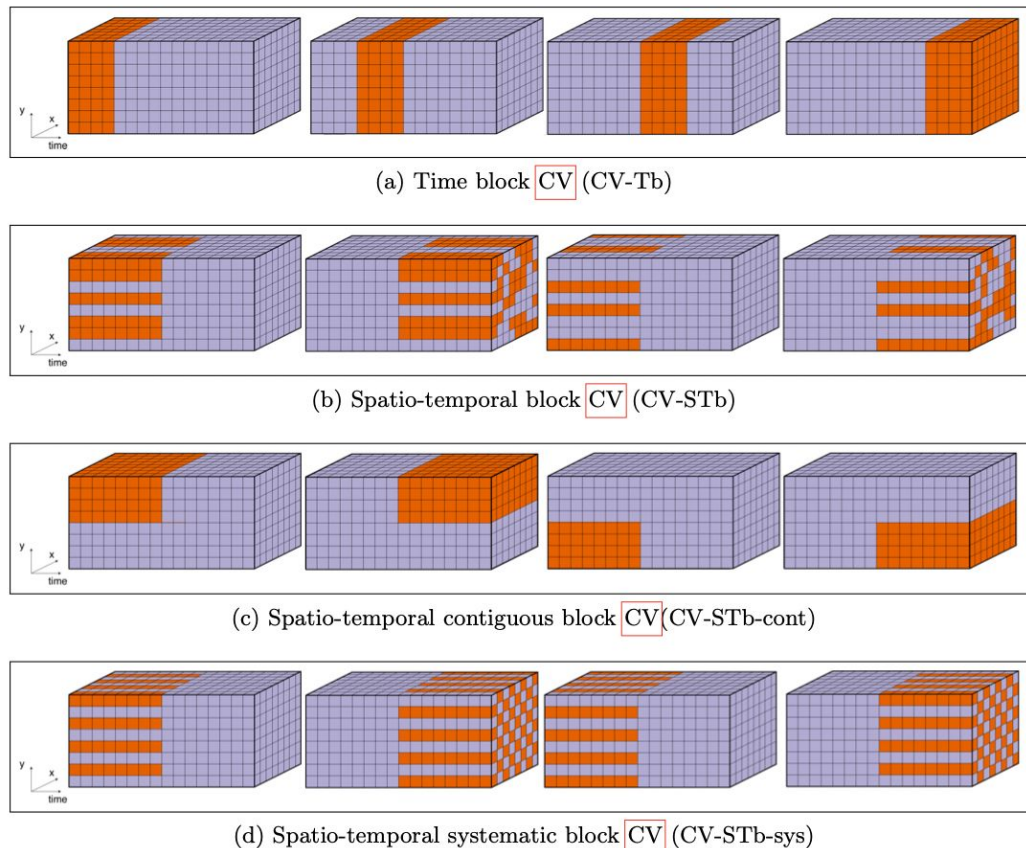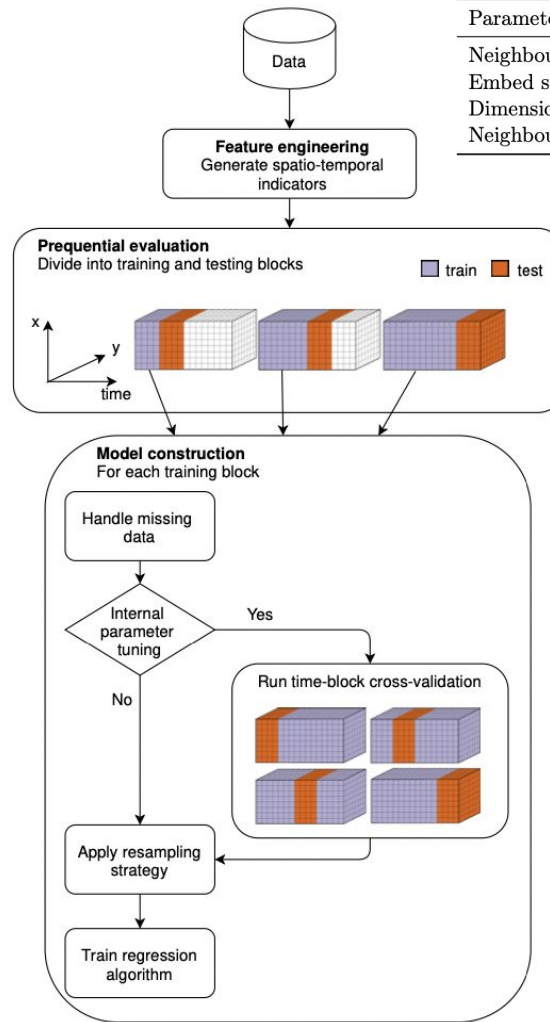- Blocking data in space and/or in time is useful in mitigating CV's bias to underestimate error



(a) Time block CV (CV-Tb)

(b) Spatio-temporal block CV (CV-STb)

(c) Spatio-temporal contiguous block CV (CV-STb-cont)

(d) Spatio-temporal systematic block CV (CV-STb-sys)

Figure 3.3: Block cross-validation methods that have prequential equivalents. Folds used for training in lighter lilac; folds used for testing in dark orange. Time flows left to right

- Identify methods for optimizing parameters of spatio-temporal cone shape, $\alpha$ and $\beta$, as well as orientation, determined by $D_a{}^\mathsf{T}$, for a given data set (e.g. producing a grid search algorithm or using numerical optimization, where possible)

- Apply block cross-validation to evaluate parameter tuning methods on a variety of spatio-temporal data sets using different modeling approaches (linear regression, SVM, regression trees, random forests, CNN, LSTM, GNN)

- If sufficient time remains, develop and publish a package in R to select optimal spatio-temporal neighbourhoods for different spatial data types (raster, point, polygon, and line)

| Parameter | Search space |
|---|---|
| Neighbourhood type | cone, reversed |
| Embed size (K) | 4, 8 |
| Dimension weight $\alpha$ | 0.1, 0.25, 0.5, 0.75, 0.9 |
| Neighbourhood radii $\beta$ | $\{0.01, 0.02, 0.03\}$, $\{0.02, 0.03, 0.04\}$ |

Data

**Feature engineering**
Generate spatio-temporal indicators

**Prequential evaluation**
Divide into training and testing blocks

train   test

x
y
time

**Model construction**
For each training block

Handle missing data

Internal parameter tuning

Yes

No

Run time-block cross-validation

Apply resampling strategy

Train regression algorithm

- So far have gathered list of key terms (mostly for my own reference) and summarised 43 relevant papers in literature review

scientific **data**

**OPEN**

**DATA DESCRIPTOR**

## Caravan - A global community dataset for large-sample hydrology

Frederik Kratzert [1 ✉], Grey Nearing[2], Nans Addor[3,4], Tyler Erickson[5], Martin Gauch [6], Oren Gilon[7], Lukas Gudmundsson [8], Avinatan Hassidim[7], Daniel Klotz[6], Sella Nevo[7], Guy Shalev[7] & Yossi Matias [7]

High-quality datasets are essential to support hydrological science and modeling. Several CAMELS (Catchment Attributes and Meteorology for Large-sample Studies) datasets exist for specific countries or regions, however these datasets lack standardization, which makes global studies difficult. This paper introduces a dataset called *Caravan* (a series of CAMELS) that standardizes and aggregates seven existing large-sample hydrology datasets. Caravan includes meteorological forcing data, streamflow data, and static catchment attributes (e.g., geophysical, sociological, climatological) for 6830 catchments. Most importantly, Caravan is both a dataset and open-source software that allows members of the hydrology community to extend the dataset to new locations by extracting forcing data and catchment attributes in the cloud. Our vision is for Caravan to democratize the creation and use of globally-standardized large-sample hydrology datasets. Caravan is a truly global open-source community resource.

- Next step is data collection, methods that are used are highly dependent on the type of data that is available; more detailed with high space/time resolution is preferable

- Domain expertise required

**CTDC**

Home  About ▼  Download ▼  Visualise  Map  Search  Related Resources ▼

⌂ / Home / Download / The Global Synthetic Dataset

### THE GLOBAL SYNTHETIC DATASET

**News:** *In December 2022, IOM released* The Global Victim-Perpetrator Synthetic Dataset *produced using an updated version of* Synthetic Data Showcase *with added support for differential privacy. The resulting dataset describes victim-perpetrator relations. It is CTDC's second synthetic dataset, and the first to provide the guarantee of differential privacy.*

Microsoft Research has worked with IOM to develop a new algorithm to derive "synthetic data" from CTDC's sensitive victim case data. Rather than systematically redacting cases, which results in a substantial amount of data being suppressed, the algorithm generates a synthetic dataset that accurately preserves the statistical properties and relationships in the original data. Representative data on all of CTDC's victim of trafficking cases are now available as a downloadable data file thanks to the new algorithm.