

Learning Landscape Features from Streamflow with Autoencoders

Alberto Bassi^{1,2}, Marvin Höge², Antonietta Mira^{3,4}, Fabrizio Fenicia², and Carlo Albert²

¹Department of Physics, ETH Zurich, Switzerland

²Swiss Federal Institute for Aquatic Science and Technology (EAWAG), Dübendorf, Switzerland

³Euler Institute, Università della Svizzera italiana, Lugano, Switzerland

⁴Insubria University, Como, Italy

Correspondence: Alberto Bassi (abassi@ethz.ch)

Abstract. Recent successes with Machine Learning (ML) models in catchment hydrology have highlighted their ability to extract crucial information from catchment properties pertinent to the rainfall-runoff relationship. In this study, we aim to identify a minimal set of catchment signatures in streamflow that, when combined with meteorological drivers, enable an accurate reconstruction of the entire streamflow time series. To achieve this, we utilize an Explicit Noise Conditional Autoencoder (ENCA), which, assuming an optimal architecture, separates the influences of meteorological drivers and catchment properties on streamflow. The ENCA architecture feeds meteorological forcing and climate attributes to the decoder in order to incentivize the encoder to only learn features that are related to landscape properties minimally related to climate. By isolating the effect of meteorology, these hydrological features can thus be interpreted as landscape fingerprints. The optimal number of features is found by means of an intrinsic dimension estimator. We train our model on the hydro-meteorologic time series data of 568 catchments of the continental United States from the CAMELS dataset. We compare the reconstruction accuracy with models that take as input a subset of static catchment attributes (both climate and landscape attributes) along with the meteorological forcing variables. Our results suggest that available landscape attributes can be summarized by only two relevant learnt features (or signatures), while at least a third one is needed for about a dozen difficult to predict catchments in the central US, mainly characterized by high aridity index. The principal components of the learnt features strongly correlate with the baseflow index and aridity indicators, which is consistent with the idea that these indicators capture the variability of catchment hydrological response. The correlation analysis further indicates that soil-related and vegetation attributes are of importance.

1 Introduction

Hydrological signatures encompass descriptive statistics derived from meteorological and streamflow time series. They serve various purposes in hydrology, such as hydrological model calibration or evaluation (Fenicia et al., 2018; Kiraz et al., 2023), process identification (McMillan, 2020a), and ecological characterization (Olden and Poff, 2003). Alongside with catchment attributes (distinguished here in landscape and climate attributes), they are also used for catchment classification and regionalization studies (Wagener et al., 2007).

Streamflow signatures, i.e. hydrological signatures solely based on streamflow, hold significant importance as they relate to the variable one aims to predict and understand (Gnann et al., 2021). Hydrologists have developed diverse signatures reflecting

25 different aspects of streamflow dynamics. Examples include those linked to the flow duration curve (e.g., slope of various
segments), the baseflow index, or the flashiness index. Numerous other such signatures exist. For instance, Olden and Poff
(2003) compiled a list of 171 indices from prior work, reflecting aspects associated to magnitude, frequency, duration, timing,
and rate of change of flow events. As streamflow depends on meteorological forcing and landscape attributes, streamflow
signatures generally contain information from both sources. In particular for predictions in ungauged basins, it is vital to be
30 able to disentangle them.

One way of doing so is through hydrological models, which condense catchment attributes into model parameters (Wagner
et al., 2003). Previous research indicates that observed hydrographs can be represented by a handful of model parameters
(Jakeman and Hornberger, 1993). For instance, the GR4J model (Perrin et al., 2003), resulting from a continuous refinement
process aimed at balancing model complexity and performance, has only four parameters. However, these analyses are based
35 on predefined model assumptions.

Model parameters can, in principle, directly be estimated from streamflow signatures. The Approximate Bayesian Computa-
tion (ABC) technique (Albert et al., 2015) has recently been used to infer model parameters from streamflow signatures - which
in this context are called summary statistics - bypassing the need to directly compare the complete time series (Fenicia et al.,
2018). If these summary statistics contained all the information necessary to estimate model parameters they would be termed
40 as sufficient. Sufficiency is therefore not an inherent property of the summary statistics but depends on the specific hydrological
model and on the parameters that need to be inferred. For ABC to converge efficiently, we also want the summary statistics
to be minimal, i.e., while they should ideally encode all the parameter related information available in the streamflow, they
should encode no other information, neither from the forcing nor from the noise that is used for the simulations (Albert et al.,
2022). Such minimally sufficient summary statistics could thus be considered the relevant fingerprints of landscape features
45 on the streamflow. Of course, this holds true only if the model is capable of encoding all the information in such features that
is relevant for the input-output relationship. However, recent studies show that purely data-driven models outperform process-
based models in prediction accuracy (Kratzert et al., 2019; Mohammadi, 2021), because they suggest information in catchment
attributes previously not utilized for streamflow prediction.

Our goal is to employ Machine Learning (ML) techniques to identify a minimal set of streamflow features enabling accurate
50 streamflow predictions when combined with meteorological forcing. Thus, our aim is to eliminate all forcing-related informa-
tion from the streamflow, distilling features solely from the catchments themselves. We approach this objective from a purely
data-driven perspective.

To identify minimal sets of streamflow features, we employ an Explicit Noise Conditional Autoencoder (ENCA) (Albert
et al., 2022), where the bare noise utilized by the stochastic model simulator is fed into the decoder together with the learnt
55 summary statistics. This way, the encoder is encouraged to encode only those features containing information on the model
parameters while disregarding the noise. Albert et al. (2022) applied ENCA to infer parameters of simple one-dimensional
stochastic maps, showing that the learnt features allow for an excellent approximation of the true posterior. In our case, instead
of noise, we input meteorological forcing into the decoder. By feeding also climate attributes to the decoder, we encourage
the encoder to exclusively encode landscape-related information within the streamflow. Moreover, since we make use of uni-

60 directional LSTM (see Appendix C), conditioning ENCA also on climate attributes could help the decoder to obtain future information about the climate that it would not be normally able to retrieve.

In order to reduce the computational costs and learn a minimal set of catchment features, the dimension of the latent space is chosen according to the estimation of its Intrinsic Dimension (ID) (Facco et al., 2017; Allegra et al., 2020; Denti et al., 2022). In particular, we employ the ID estimator GRIDE (Denti et al., 2022), which is robust to noise. Learnt features will then be
65 compared with known catchment attributes (both from the landscape and the climate) and hydrological signatures to provide a hydrological interpretation and guide knowledge domain experts towards the pertinent information necessary for streamflow prediction.

We apply our approach to the US-CAMELS dataset (Newman et al., 2015), covering several hundred catchments over the continental US. LSTMs (Long Short-Term Memory networks) have emerged as state-of-the-art models for streamflow data-
70 driven predictions. In the study of Kratzert et al. (2019), LSTMs validated on unseen catchments, enriched with static landscape and climate attributes from Addor et al. (2017), outperformed conceptual models. First investigations towards mechanistic interpretation of the LSTM states, e.g. linking hidden states to the dynamics of soil moisture, demonstrated the potential of eliciting physics from data-driven models (Lees et al., 2022). Here, by linking learnt features to known catchment attributes, we explore a further aspect in this broader field of explainable AI or interpretable ML (Molnar, 2024; Molnar et al., 2020).

75 Our specific objectives are: (i) find the minimal number of dominant streamflow-features stemming from the landscape; (ii) relate them to known landscape and climate attributes as well as established hydrological signatures. This will not only allow us to determine how many features are required for streamflow prediction, but also to answer the question whether there is missing information in known catchment attributes.

A similar attempt of learning signatures has recently been made by Botterill and McMillan (2023). In pursuit of an inter-
80 pretable latent space on a continental scale, they employed a convolutional encoder to compress high-dimensional information derived from meteorological forcing and streamflow data. This approach was aimed at learning hydrological signatures (McMillan, 2020b) within the US-CAMELS dataset. Their approach differs from ours in three aspects: (i) they used a traditional conceptual model as a decoder whereas we use an LSTM architecture which has been shown to be superior to conceptual models when provided with catchment properties; (ii) they fed both streamflow and meteorological forcing into the encoder
85 whereas we feed in only streamflow data in an attempt to separate landscape- from forcing-information; (iii) they did not attempt to find a minimal number of signatures sufficient for streamflow prediction, whereas this is a primary objective of our work.

It is important to note that our main objective is not to beat state-of-the-art models regarding their predictive performance (Kratzert et al., 2021; Klotz et al., 2022). Our goal is rather to investigate the information content in streamflow. However, we
90 believe our research will provide valuable insights into the most critical features for streamflow prediction.

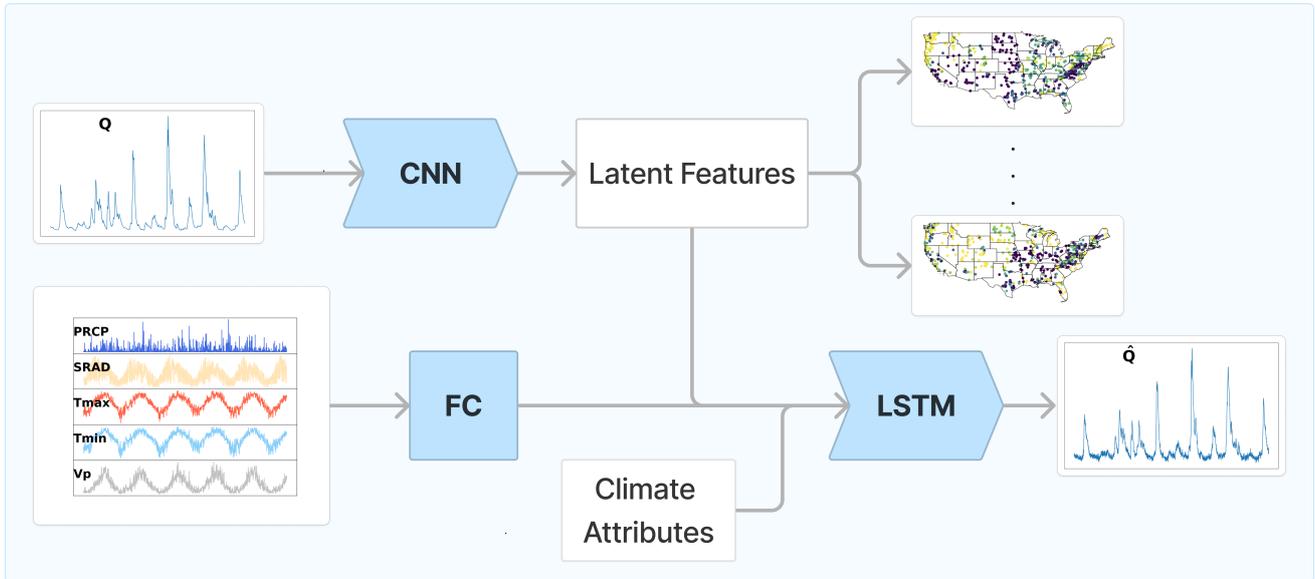


Figure 1. Explicit Noise Conditional Autoencoder used in this study. For the hyper-parameters and the implementation details, the reader is referred to Appendix C. The neural network architectures employed are Convolutional Neural Networks (CNN), a Fully Connected (FC) layer and LSTM. The observed and simulated streamflows are denoted with Q and \hat{Q} , respectively. The meteorological forcing variables are denoted with PRCP (precipitation), SRAD (solar radiation), T_{max} (maximum temperature), T_{min} (minimum temperature) and V_p (vapour pressure).

2 Models and Methods

2.1 Data

We employ the Catchments Attributes and Meteorology for Large-sample Studies (CAMELS) (Newman et al., 2015), which consists of 671 catchments in the contiguous United States (CONUS), ranging in size from 4 to $25 \cdot 10^3$ km². For this study, we select those 568 catchments out of 671 whose data span continuously on a daily basis the time period from 1 October 1980 to 30 September 2010, corresponding to 30 hydrological years. The first 15 years are used for training and the last 15 for testing. Along with the streamflow time series and the meteorological forcing variables, US-CAMELS also provides information about catchments static attributes (Addor et al., 2017), encompassing both landscape (vegetation, soil, topography and geology) and climate. Streamflow data is retrieved from the U.S. Geological Survey gauges, while the meteorological forcing comes from the extended North America Land Data Assimilation Systems (NLDAS) (Kratzert, 2019) and includes maximum and minimum daily temperature, solar radiation, vapour pressure and precipitation.

2.2 Explicit Noise Conditional Autoencoder

We use the Explicit Noise Conditional Autoencoders (ENCA) (Albert et al., 2022), where the streamflow is fed to a convolutional encoder. ENCA has been designed to distill sufficient summary statistics which contain minimal noise information from the output of stochastic models. Here (Figure 1 - for the detailed architectures the reader is referred to Appendix C), the noise is substituted by all the variables we are not interested in, namely the meteorological forcing. The convolutional encoder is thus followed by a LSTM decoder that takes as input 15 hydrological years of meteorological forcing, i.e. 5478 time points, and nine climate attributes (reported in Table 1). The LSTM capacity is limited by the dimension of the input layer. In order to enlarge the available capacity and capture more complex patterns from the meteorological forcing, the meteorological time series are first fed to a single linear layer with 1350 output units. The output of this linear layer is then concatenated with the output of the encoder and fed to the LSTM decoder. This way the decoder sees tensors of size $(BS, 5478, 1350 + N)$, where BS is the batch size (batches are selected across different catchments) and N is the latent space dimension. We opted for such an architecture in order to extract as much static information related to the streamflow as possible.

We expect to be able to compress almost all streamflow information not already contained in the forcing into a low (N)-dimensional feature vector¹. Because they should be largely devoid of forcing information, we call these features the *relevant landscape features* and explain in subsection 2.4 how we fix their number. In principle, these features are hydrological signatures (since they are functions of the streamflow) which contain as little information as possible about the forcing and the climate. However, this can only be achieved if the decoder is capable of utilizing all the available information in the meteorological drivers, which is never fully realized in practice. Also the sufficiency of the learnt features is not guaranteed a priori and depends on the encoder architecture employed, here a convolutional network.

Comparing relevant landscape features with known static catchment attributes in terms of their capacity for streamflow reconstruction will allow us to find out whether static catchment attributes lack information that is crucial for streamflow prediction. For this comparison, we use an LSTM model augmented with catchment attributes (Addor et al., 2017) in the input, stemming from both the landscape and the climate. We refer to this model as Catchment Attributes Augmented Model (CAAM). This model differs from Figure 1 solely by the fact that the latent features are substituted by known landscape attributes. Following Kratzert et al. (2019), CAAM is fed with 27 catchment attributes (reported in Table 1), which are representative of climate, topography, geology, soil and vegetation.

In order to mitigate numerical instability, it is crucial to standardize the catchment attributes or latent features before feeding the LSTM. In CAAM, we standardize the catchment attributes with the mean and standard deviation computed over all the considered catchments. This is not possible for ENCA, since the mean and standard deviation of the latent features are not known a priori. Therefore, we standardize the latent features by means of a batch normalization layer. This way, we ensure that the magnitude of the LSTM input is comparable between CAAM and ENCA.

¹We refer to the ENCA model with latent space dimension N as ENCA- N .

2.3 Training and Testing

We use the first 15 years of data for training and the last 15 for testing. Training is performed by maximizing the Nash-Sutcliffe Efficiency (NSE) (Nash and Sutcliffe, 1970), defined as:

$$\text{NSE} = 1 - \frac{\sum_{t=1}^T (q_{\text{sim},t} - q_{\text{obs},t})^2}{\sum_{t=1}^T (q_{\text{obs},t} - \mu_{\text{obs}})^2}, \quad (1)$$

where $q_{\text{obs},t}$ and $q_{\text{sim},t}$ are, respectively, the observed and predicted streamflow expressed in *mm/day* at day t , and μ_{obs} is the average of the observed streamflow. We notice that maximizing the NSE is equivalent to minimizing the Mean Square Error (MSE) between data and prediction. Each model is trained with the Adam optimizer (Kingma and Ba, 2015), with learning rate equal to 10^{-5} for 10,000 epochs. The batch size is set to 64 and the first 270 days of the predicted streamflow are excluded when computing the loss.

We also report the three components in which NSE can be decomposed, see Eq. 4 in the main text of Gupta et al. (2009). These components are the linear correlation coefficient (R), the bias normalized by the observed streamflow standard deviation (BIAS) and standard deviation ratio (STDEV). The linear correlation coefficient is related to timing, whereas STDEV measures the streamflow variability and it is defined as

$$\text{STDEV} = \frac{\sigma_{\text{sim}}}{\sigma_{\text{obs}}}, \quad (2)$$

where σ_{sim} and σ_{obs} are the standard deviations of the simulated and observed streamflows, respectively. Finally, the BIAS is related to volume errors and it is defined as

$$\text{BIAS} = \frac{\mu_{\text{sim}} - \mu_{\text{obs}}}{\sigma_{\text{obs}}}. \quad (3)$$

Each algorithm is affected by noise, due to the random initialization of the neural network parameters. To minimize this effect we run each model with four random restarts, each one providing the streamflow prediction in the whole testing period. We compute the evaluation metrics on the predicted streamflow after averaging the streamflow over these four random restarts.

2.4 Intrinsic Dimension Estimation

The selection of the encoder latent space dimension, specifically the number of relevant features, is informed by the ID estimator GRIDE (Denti et al., 2022). We utilize the GRIDE paths, which involve estimating the ID at several distance scales at which the data is analyzed (for an in-depth discussion on ID, refer to Appendix A). The ID intuitively measures the dimension of the manifold where the data resides, which may be lower than the dimension of the embedding space. Most ID estimators depend on calculating the distance scale between data points, and the estimated ID itself can vary with this distance scale. Figure 2, derived from Denti et al. (2022), shows points on a one-dimensional line with added noise embedded in a three-dimensional space. When the distance scale is too small, the data points appear to fill the space uniformly, making the manifold seem three-dimensional. However, as the distance scale increases and the noise is bypassed, the estimated ID decreases until the correct value of one is achieved.

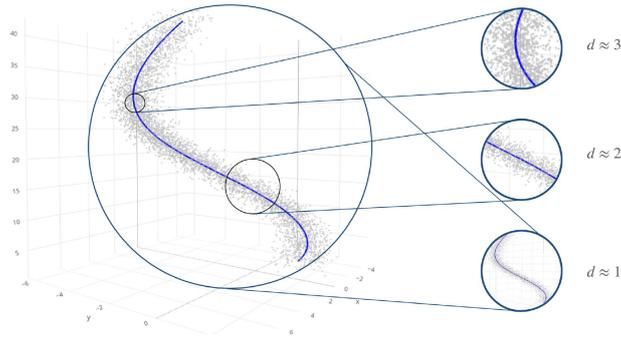


Figure 2. Unidimensional line with noise, embedded in three dimensions. The estimated intrinsic dimension depends on the distance scale.

To identify the dimension of the latent space of ENCA we proceed with the following methodology. First, we train an ENCA- N with a relatively large number of latent features N . Since we fed 27 catchment attributes to the reference model (CAAM), we use a 27-dimensional latent space in order to have a fair comparison in terms of model capacity. We refer to this model as ENCA-27. The exact dimension of the latent space we start with does not matter much, as long as it is larger than the expected number of relevant landscape features. Then, we estimate the ID of the latent space, and train another ENCA with the number of latent features equal to the estimated ID and, in turn, estimate its ID to check if the dimension of the latent space can be further reduced. We thus use the ID as a guide to progressively diminish the dimension of the latent space of the autoencoder. However, in the end we train ENCA for several dimensions of the latent space and evaluate the information content of the learnt features in terms of their ability to reconstruct streamflows. In Figure A1 we report GRIDE paths for different models trained in this work.

3 Latent Space Interpretation

The relevant features are first projected using a Principal Component Analysis (PCA), since in general the autoencoder latent representation is in arbitrary coordinates. Doing so, we ensure a fair comparison among different random restarts, since we change the basis of each latent space by ordering the new coordinates according to the explained variance. Finally, in order to interpret the relevant landscape features, we report the absolute Spearman correlation (Zar, 2005) matrix among the learnt features, static catchment attributes and hydrological signatures, which are reported in Table 1.

Meteorological Forcing Variables

PRCP	Average daily precipitation (<i>mm/day</i>).
SRAD	Surface incident solar radiation (W/m^2).
Tmax	Maximum daily atmosphere temperature ($^{\circ}C$).
Tmin	Minimum daily atmosphere temperature ($^{\circ}C$).
Vp	Nearly surface daily vapour pressure average (Pa).

Climate Attributes

Prec Mean	Mean daily precipitation.
PET Mean	Mean daily potential evapotranspiration.
Prec Seasonality	Seasonality of precipitation estimated by using sinusoidal waves.
Fraction Snow	Fraction of precipitation falling on days with $T < 0^{\circ}C$.
Aridity Index	Ratio between the mean PET and mean precipitation.
High Prec Frequency	Frequency of days with $\leq 5x$ mean daily precipitation.
High Prec Duration	Mean duration of high precipitation events.
Low Prec Frequency	Frequency of days with ≤ 1 mm/day of precipitation.
Low Prec Duration	Mean duration of dry periods.

Hydrological Signatures

Q Mean	Mean daily streamflow (mm/day).
Streamflow Ratio	Ratio between the mean daily streamflow and mean daily precipitation.
Slope FDC	Slope of the flow duration curve.
Baseflow Index	Ratio between the average daily baseflow and streamflow.
Stream ELAS	Streamflow precipitation elasticity.
Q5	5 % flow quantile (mm/day).
Q95	95 % flow quantile (mm/day).
High Q Frequency	Frequency of high-flow days (> 9 times the median daily flow).
High Q Duration	Mean duration of high flow events (number of consecutive days > 9 times the median daily flow).
Low Q Frequency	Frequency of low flow days ($< 0.2 x$ the mean daily flow).
Low Q Duration	Mean duration of low-flow events (number of consecutive days < 0.2 times the mean daily flow).
HFD Mean	Mean half-low-date (date on which the cumulative streamflow since October the 1st reached half of the annual streamflow).

Landscape Attributes

Topographic Attributes

Elevation Mean	Mean elevation of the catchment.
Slope Mean	Mean slope of the catchment.
Area Catchment	Area of the catchment.

Geological Attributes

Carbonate Rocks Fraction	Carbonate sedimentary rocks fraction area in the catchment.
Geological Permeability	Surface permeability (in log10 scale).

Soil Attributes

Soil Depth (Pelletier)	Depth to bedrock (maximum 50 m).
Soil Depth (STATSGO)	Soil depth (maximum 1.5 m).
Soil Porosity	Volumetric porosity.
Soil Conductivity	Saturated hydraulic conductivity.
Max Water Content	Maximum water content of the soil.
Sand Fraction	Fraction of sand in the soil.
Silt Fraction	Fraction of silt in the soil.
Clay Fraction	Fraction of clay in the soil.

Vegetation Attributes

Fraction Forest	Fraction of the catchment covered by forest.
LAI Max	Maximum monthly mean of leaf area index.
LAI Diff	Difference between the max. and the min. monthly mean of the leaf area index.
GVF Max	Maximum monthly mean of the green vegetation fraction.
GVF Diff	Difference between the max. and min. monthly mean of the green vegetation fraction.

Table 1. Meteorological forcing variables, climate and landscape (topographic, geological, soil and vegetation) attributes and hydrological signatures compared in this study. Both climate and landscape attributes are fed to CAAM. ENCA models are conditioned on climate attributes too.

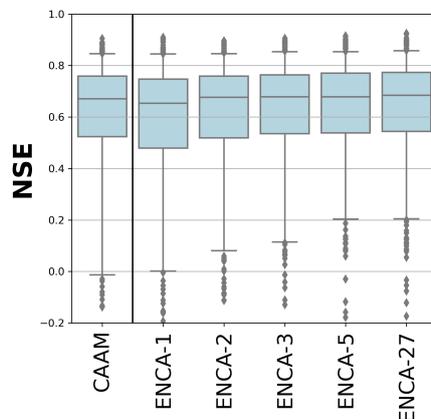


Figure 3. NSE values for the considered models in the test period. The boxes are delimited by the 25 % and the 75 % quantiles, while the whiskers indicate to the 5 % and 95 % quantiles. CAAM performance is similar to ENCA-2.

Figure 3 depicts the boxplot of the test NSE values for the considered models. We report and discuss the associated statistics, the boxplots of the NSE components (R, BIAS, STDEV) and the correlations among them in Appendix B.

In terms of NSE, we observe a performance improvement from ENCA-1 to ENCA-2 in the bulk of the distribution, and a further minor improvement from ENCA-2 to ENCA-3. The NSE improvement between ENCA-3 and ENCA-5 is minor and is mainly related to outliers while ENCA-5 and ENCA-27 distributions are almost identical.

In general, Figure 3 and the related statistics (Figure B1) show that increasing the number of latent features improves the prediction accuracy of the considered metrics. Even though it is difficult to set a cut-off dimension, we can state that: i) with more than five latent features we do not observe a performance improvement anymore, meaning that 5 features are a sufficient set of summary statistics of the streamflow (which, however, can still depend on the chosen encoder architecture). ii) Overall, CAAM performance is most similar to ENCA-2. We therefore argue that known catchment attributes (selected in this study) account for two relevant landscape features that appear to be sufficient for most catchments, while at least a third one is needed to resolve specific catchments.

To study which catchments are most affected when using the latent features of the ENCA models in place of the known catchment attributes in CAAM, we report (Figure 4) the NSE difference between CAAM and ENCA-2 (left panels) or ENCA-3 (right panels), respectively. While, for most catchments, switching from ENCA-2 to ENCA-3 does not result in a high performance gain, we see a clear improvement on about a dozen or so catchments, mostly located in the central CONUS. This corroborates the hypothesis that the known catchment attributes account for two relevant landscape features and the improvement due to the third one is related to only few catchments that are particularly difficult to predict. It is interesting to

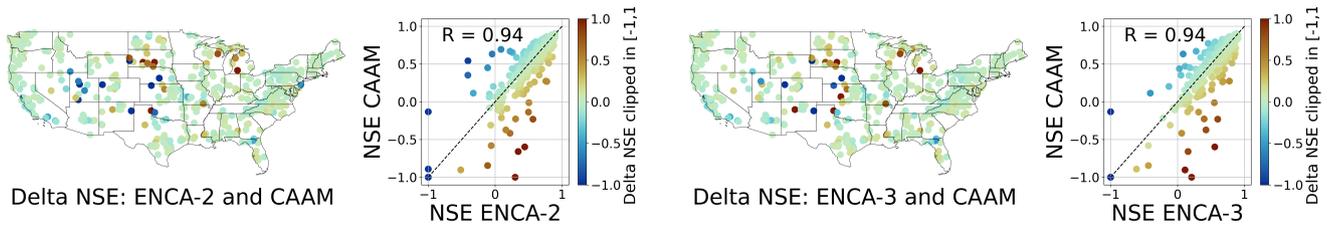


Figure 4. Test NSE of ENCA-2 (left panel) and ENCA-3 (right panel) versus CAAM, color-coded with the NSE difference per catchment clipped in $[-1,1]$. Red means ENCA performs better, blue means CAAM performs better.

count the number of catchments for which CAAM’s NSE is negative (i.e. predictions that are worse than average streamflow) but ENCA’s NSE is positive. This number is 17 for ENCA-2 and 15 for ENCA-3. On the other hand, the number of catchments for which CAAM’s NSE is positive but ENCA’s negative decreases from eight (ENCA-2) to only two (ENCA-3).

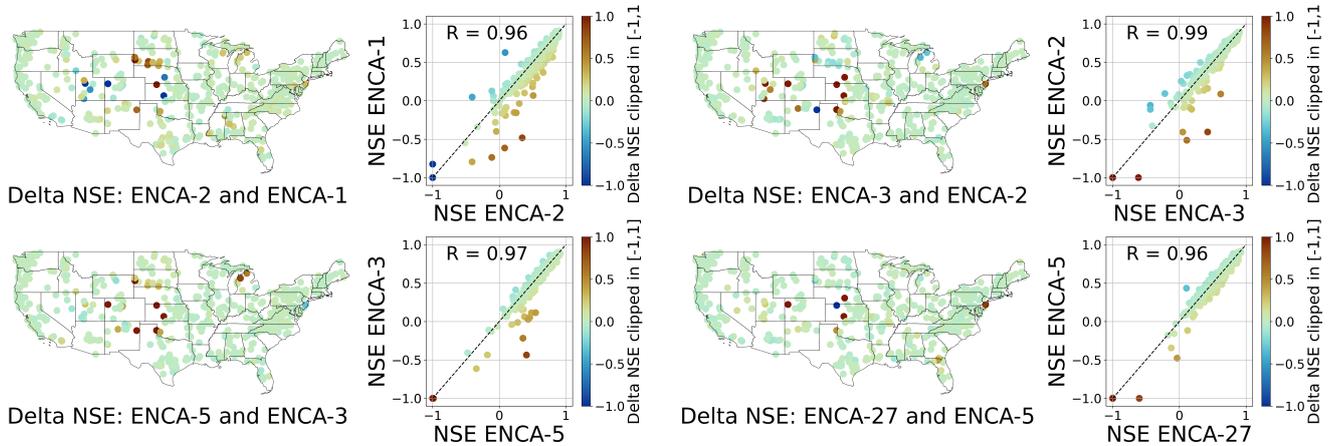


Figure 5. Test NSE of ENCA models with different number of latent features, color-coded according to the NSE difference per catchment clipped in $[-1,1]$. The improvement of increasing the number of latent features is significant from ENCA-2 to ENCA-3 and marginal for more complex models. The improved catchments are mainly located in the central CONUS, dominated by arid climate conditions.

In order to evaluate the impact of additional features, we compare the performances of ENCA models differing in their number of latent variables (Figure 5). The results corroborate our earlier findings that two features are sufficient to cover most of the catchments, and additional features provide information about relatively few, difficult to predict, catchments mostly located in the central CONUS and dominated by arid climate conditions. While the number of such catchments informed by the third feature is relatively high, additional features only have a minor effect. Indeed, adding a third feature turns seven catchments from negative NSE to positive NSE and only leads to marginal deterioration of three other catchments. Additional features have much less dramatic effects.

Note that test NSE obtained in this work are good, but still far from state-of-the-art approaches on the same dataset. For comparison, the global model of Kratzert et al. (2019) (augmented with the same catchment attributes reported in Table 1) achieves a median NSE of 0.74, while in this work the best model achieves a median NSE of 0.68. Later approaches (with multiple input forcings) achieve even higher performance of 0.82 (Kratzert et al., 2021). One might be tempted to attribute this to over-fitting due to the very large number of parameters of our architecture (about 3 million). However, the test MSE loss curves (Figure C1) do not support this hypothesis. Also the very long sequences fed to the LSTMs might affect their performance, as they are known to suffer from vanishing/exploding gradients. While we use time-series of length 5478, state-of-the-art approaches use lengths of 270 (Kratzert et al., 2019) and 365 (Kratzert et al., 2021).

4.2 Interpretation of the Relevant Feature Principal Components

Figure 6 shows the absolute Spearman correlation matrix between the principal components of the identified three relevant features, the known streamflow signatures and catchment attributes across different random restarts of the model. The relevant features share information with catchment attributes and hydrological signatures. For instance, feature one carries information about basic hydrological attributes like baseflow index and low flow frequency. Moreover, feature one is (weakly) correlated with soil-related attributes like soil porosity and conductivity, sand, silt and clay fraction. Feature two is correlated with climatic indicators, such as the aridity index, the mean precipitation, high and low precipitation frequency, but also with hydrological signatures like mean streamflow and the 95% quantile of the flow duration curve. We point out that even though the encoder is explicitly designed to learn non-climate landscape features, we can still observe a correlation between latent features and climate attributes. This correlation can be due to collinearities between landscape and climate attributes. In this case, the collinear attributes are those related to vegetation, like the fraction of forests and the maximum GVF (Green Vegetation Fraction), which are obviously correlated with climate. For instance, from Figure D1 we can observe that the aridity index is highly correlated with the mean precipitation (0.88) and the fraction of forest (0.74), while these last two attributes are fairly correlated between each other (0.67).

Finally, feature three is mostly correlated with high and low flow duration and frequency, signatures relating to the extremes of streamflow. Interestingly, this principal component does not hold much information about neither landscape nor climate attributes, indicating that it encodes catchment information that has not yet been considered or that is not related to any discernible catchment feature. Since the third feature mainly conveys information about certain dry and hard to predict catchments, the latter might very well be the case.

The discussed principal components, however, do not share the same amount of explained variance: feature one accounts for about 60%, feature two for about 30% and feature three for about 10%. In Appendix E we report the correlation matrices for the other ENCA models, where we can verify that the principal components carry the same information for streamflow prediction consistently across different models.

The geospatial distribution of feature importance is shown in Figure 7, where a non-trivial distribution of the features appears, highlighting that the different features have different information content for different regions: feature one dominates in the less to non-arid eastern CONUS, while feature two is mainly dominant in the western part. Feature three does not show such

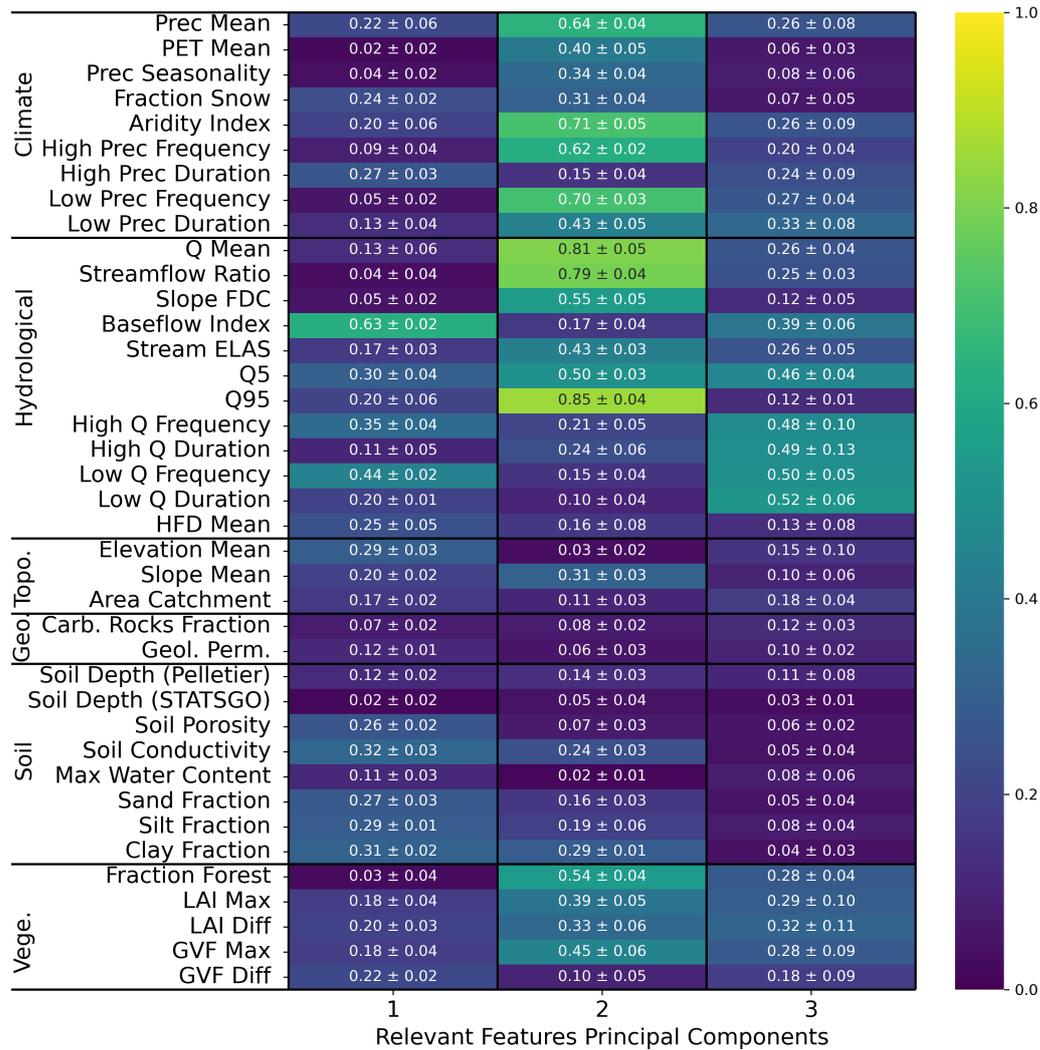


Figure 6. Average (plus minus the standard deviation) of the absolute Spearman correlation of relevant features principal components of ENCA-3 with respect catchment attributes and hydrological signatures across four different random model restarts. The colors refer to the average.

a clear spatial representation. Overall, the potential of delineating geospatial relations is another indicator that the encoder has learnt from the landscape signal in the data.

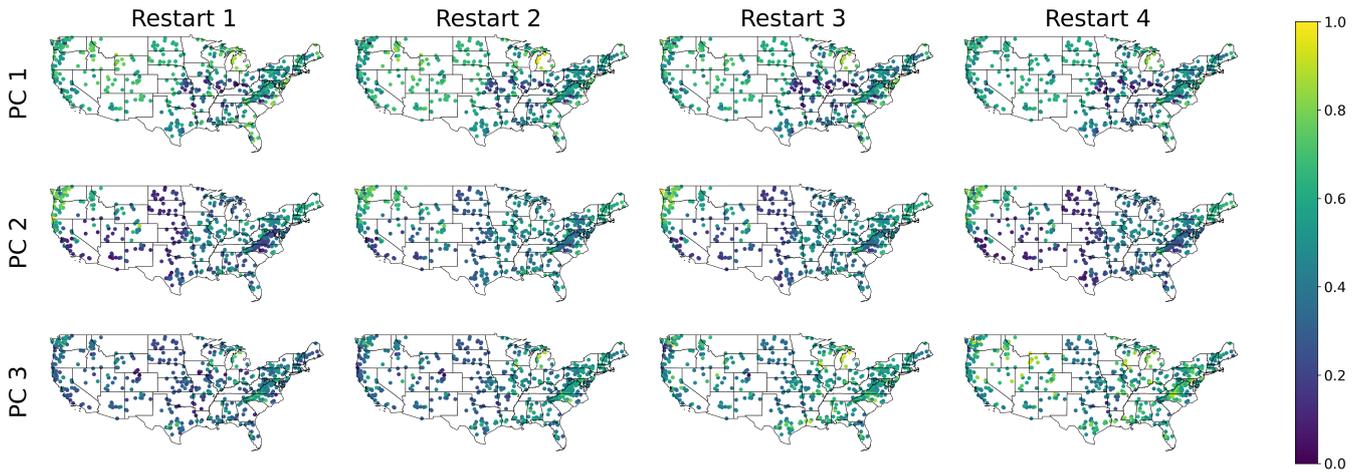


Figure 7. Principal Components (PCs) of the relevant features of ENCA-3 in the CONUS. For a better comparison across different restarts, PCs are normalized in the interval $[0, 1]$ and their sign is adjusted such that the PCs of the first catchment in the dataset are always positive.

245 5 Conclusions and Outlook

We employed a conditional autoencoder to distill a minimal set of streamflow features (signatures) necessary for streamflow reconstruction in conjunction with meteorological data. These features are minimally related to climate and can be interpreted as landscape fingerprints on the streamflow. We compared these features with known catchment attributes in terms of their capacity for streamflow reconstruction. The primary conclusions we highlight in this study are:

- 250 – For all the metrics considered, ENCAs (Explicit Noise Conditional Autoencoders) perform better than the reference attributes enhanced model (CAAM) when the number of latent features is greater than two. In fact, two features seem to be sufficient for most catchments, while a relatively small number of catchments, mostly located in the central CONUS, require a third one. Including more than three features, however, only leads to marginal improvements. We therefore conjecture that most of the information contained in the static attributes used for CAAM, insofar as it is relevant for streamflow prediction, can be reduced to two independent features. The third latent feature, however, seems to encode information that is not fully contained in those static attributes.
- 255
- 260 – The correlation between attributes and importance of the relevant features (see Figure 6) suggests an ordering of the information contained in the features for accurately predicting discharge: first, basic hydrological attributes like baseflow and soil-related attributes, followed by the average streamflow and the 95% flow quantile (correlated to climate due to collinearities with vegetation-related attributes) and, third, specifics on the high and low flow, i.e. the extremes of the hydrograph. Looking back at Figure 4, this last feature appears to encode the information that is needed to exceed the model performance that is only based on the 27 static attributes (CAAM).

In summary, our research reveals a significant reduction in the dimensionality of the streamflow time series, at least in relation to the calibration metric used, i.e. the NSE. In principle, using different calibration metrics can modify the type and the number of learnt features. This comparison lies beyond the scope of this work and could be an interesting avenue of exploration. Despite the plethora of hydrological signatures and catchment attributes at our disposal, only a small subset proves essential for NSE-accurate streamflow prediction. This finding echoes established results from prior studies (Jakeman and Hornberger, 1993; Edijanto et al., 1999; Perrin et al., 2003), suggesting that hydrological systems might be effectively modelled using only a limited set of parameters. The low dimensionality of the relevant catchment information opens up the opportunity for a better understanding of its nature, suggesting some future research directions:

- A promising approach could be the adoption of NeuralODEs (Höge et al., 2022), which offer a high level of interpretability due to their low number of states. This combination of a few states and a few features may help to decipher not only the nature of the relevant catchment information but also how it influences streamflow.
- Preliminary analysis (not shown in this paper) has revealed that the known static catchment attributes live on a low-dimensional manifold, which is in line with our finding that only two independent features seem to capture most of the information that is relevant for streamflow. While the correlation-based analysis presented in this paper gives some clues as to how these features can be interpreted, more sophisticated types of analysis like those based on Information Imbalance (Glielmo et al., 2022) might allow for a more precise understanding of their physical nature.

Code and data availability. The US-CAMELS dataset, as well as the catchment attributes, is available at the site <https://ral.ucar.edu/solutions/products/camels>. The extended NLDAS forcing dataset is available at <https://doi.org/10.4211/hs.0a68bfd7ddf642a8be9041d60f40868c>. All the code used for this work is publicly available at <https://doi.org/10.5281/zenodo.13132951>.

Appendix A: The Intrinsic Dimension

In order to estimate the ID, we apply the GRIDE estimator (Denti et al., 2022). Given sample points, $\mathbf{x}_i \in \mathbb{R}^D$, for $i = 1, \dots, M$, and a distance measure, $r : \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}^+$, GRIDE assumes that points in a given neighbourhood are counted with a Poisson point process with intensity ρ , which is constant at least up to the scale of the diameter of the considered neighbourhood. Let $r_{i,l}$ be the distance between the point \mathbf{x}_i and its l -th nearest neighbour, and define $\mu_{i,n_2,n_1} = \frac{r_{i,n_2}}{r_{i,n_1}}$, where n_1 and n_2 (with $0 \leq n_1 \leq n_2 \leq M$) are given integers. The distribution of μ_{i,n_2,n_1} can be computed in closed form and depends only on the ID of the data while, crucially, does not depend on ρ , as long as ρ is constant in the considered neighbourhood of i whose diameter is set by the distance between i and its n_2 -th nearest neighbour (Denti et al., 2022).

In order to correctly identify the ID of a dataset, a scale-independent analysis is essential. We therefore make use of GRIDE paths, the evolution of the ID estimate as a function of n_2 , which can be interpreted as the scale at which we look at the data. We set $n_1 = n_2/2$, as usually done in the literature. As a function of n_2 , the ID is first expected to increase, due to the noise present at small distance scales, and then to reach a plateau corresponding to the correct ID.

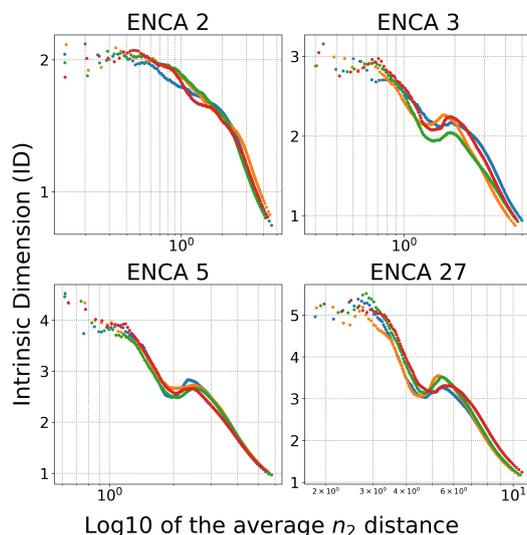


Figure A1. GRIDE evolution plot for the different ENCA models employed for the four random restarts of the models.

Figure A1 shows the GRIDE path for different ENCA models trained in this work. In particular, the ID estimate of the latent space of ENCA-27 decreases after showing a plateau around five, then reaches a minimum around three, then increases again and finally collapses to low values at larger distance scale. The plateau at five motivates us to train an ENCA-5 and study its ID. The local minimum of the GRIDE path of the latent space of ENCA-5 is consistent with an ID of three. We can deduce that, for most of the catchments, the ID is three, while for some it can be higher. However, the fact that the GRIDE path of the latent space of ENCA-27 shows two plateaus around five and three can be an indicator of the existence of two or more manifolds with different IDs. From Figure 3 we see that, indeed, three features seem to capture most of the relevant information, which

is in line with the GRIDE path for ENCA-5.

Appendix B: Metric Comparisons

We report some summary statistics of the NSE, R, BIAS and STDEV values across the 568 catchments considered in this study.

In terms of the linear correlation coefficient R and the BIAS (left and central panels of Figure B1), we observe a similar pattern to what we observed for the NSE distribution (Figure 3). By increasing the number of latent features to three, the bulk distribution improves significantly. At the same time, further increasing the number of latent features improves only the BIAS outliers.

Regarding the STDEV (right panel of Figure B1), it is clear that both CAAM and ENCA models tend to underestimate streamflow variability, a known issue associated with using NSE as the objective function (Gupta et al., 2009). However, the differences between CAAM and ENCA models are less pronounced, and unlike the observations for NSE, R, and BIAS, a clear performance hierarchy cannot be established.

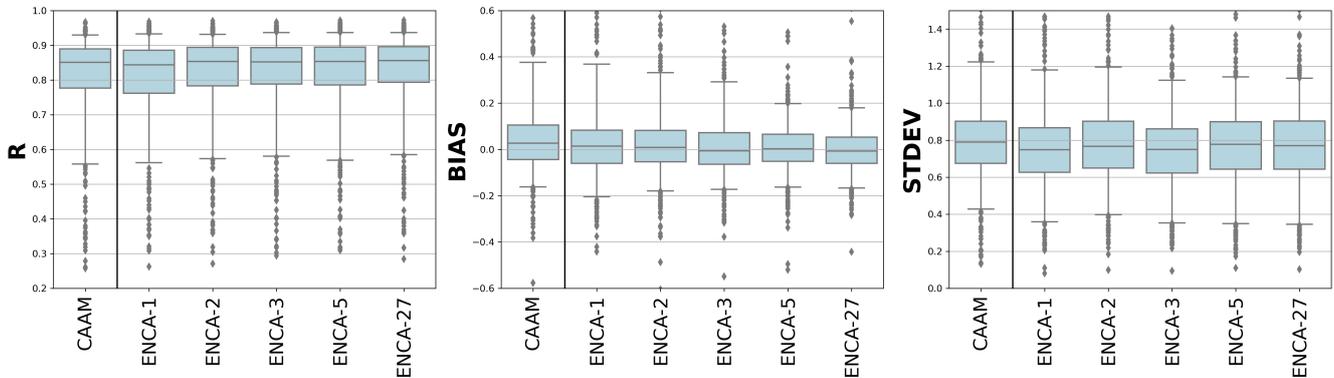


Figure B1. R, BIAS and STDEV values for the considered models in the test period. The boxes are delimited by the 25 % and the 75 % quantiles, while the whiskers indicate to the 5 % and 95 % quantiles. For these metrics as well, CAAM performance is similar to ENCA-2.

Table B1 demonstrates that increasing the number of latent features in ENCA generally enhances performance. However, the improvement is modest between ENCA-3 and ENCA-5 and nearly negligible between ENCA-5 and ENCA-27. The performance of CAAM is more comparable to ENCA-2 overall, suggesting that the third latent feature is crucial for better predictions in certain catchments, while five features appear to be the maximum number our encoder can effectively learn. Furthermore, as shown in Figure 5, it is clear that almost all catchments that improved from ENCA-3 to ENCA-5 still exhibit very poor performance (NSE below -1.0).

From Table B2, we can observe that the NSE is strongly correlated with the linear coefficient R and fairly correlated with STDEV and the correlation increases with the performance (see Figure 3). At the same time, the NSE is anti-correlated with the BIAS, but this correlation is low. While the correlation between R and BIAS and between BIAS and STDEV is weak, the correlation between R and STDEV is intermediate and increases with performance. These results agree with the findings of

(Gupta et al., 2009), who showed that with optimal BIAS values, optimal NSE values are reached when R is correlated with
325 STDEV.

	CAAM	ENCA-1	ENCA-2	ENCA-3	ENCA-5	ENCA-27	
NSE	Mean	0.52	0.47	0.50	0.53	0.58	0.59
	Min	-17.66	-15.08	-17.74	-14.77	-11.03	-6.62
	Q5	-0.02	-0.00	0.07	0.11	0.19	0.20
	Q25	0.52	0.48	0.52	0.53	0.54	0.54
	Median	0.67	0.65	0.68	0.68	0.68	0.68
	Q75	0.76	0.75	0.76	0.76	0.77	0.77
	Q95	0.85	0.85	0.85	0.86	0.85	0.86
	Max	0.90	0.91	0.90	0.91	0.92	0.93
	# < 0.0	31	29	22	18	14	13
R	Mean	0.81	0.80	0.82	0.82	0.82	0.82
	Min	0.11	0.12	0.14	0.15	0.13	0.11
	Q5	0.55	0.55	0.57	0.58	0.57	0.58
	Q25	0.78	0.76	0.78	0.79	0.79	0.79
	Median	0.85	0.84	0.85	0.85	0.85	0.86
	Q75	0.89	0.89	0.89	0.89	0.89	0.90
	Q95	0.93	0.93	0.93	0.94	0.94	0.94
	Max	0.97	0.97	0.97	0.97	0.97	0.97
	BIAS	Mean	0.05	0.05	0.05	0.02	0.02
Min		-1.13	-0.86	-0.86	-1.05	-0.52	-0.63
Q5		-0.16	-0.21	-0.18	-0.17	-0.16	-0.17
Q25		-0.04	-0.06	-0.05	-0.06	-0.05	-0.06
Median		0.03	0.01	0.01	-0.01	0.00	-0.01
Q75		0.11	0.08	0.08	0.07	0.07	0.05
Q95		0.40	0.40	0.34	0.30	0.20	0.18
Max		1.70	2.87	2.53	2.03	1.70	1.28
STDEV		Mean	0.82	0.77	0.80	0.76	0.78
	Min	0.13	0.08	0.10	0.09	0.11	0.10
	Q5	0.42	0.35	0.40	0.35	0.35	0.34
	Q25	0.67	0.63	0.65	0.62	0.64	0.64
	Median	0.79	0.75	0.77	0.75	0.78	0.77
	Q75	0.90	0.87	0.90	0.86	0.90	0.90
	Q95	1.23	1.18	1.20	1.13	1.14	1.13
	Max	4.26	3.64	4.12	3.88	3.55	2.88

Table B1. Metrics comparison for different models computed in the test period. We report the mean, the minimum, the 5% quantile, the 25% quantile, the median, the 75% quantile, the 95% quantile and the maximum values. Additionally, we report the number of catchments whose predicted NSE values are lower than zero.

	CAAM	ENCA-1	ENCA-2	ENCA-3	ENCA-5	ENCA-27
NSE - R	0.90	0.89	0.89	0.88	0.91	0.92
NSE - BIAS	-0.33	-0.30	-0.33	-0.32	-0.10	-0.12
NSE - STDEV	0.28	0.29	0.29	0.40	0.48	0.50
R - BIAS	-0.23	-0.33	-0.34	-0.30	-0.21	-0.20
R - STDEV	0.45	0.39	0.38	0.43	0.47	0.48
BIAS - STDEV	0.24	0.21	0.19	0.28	0.32	0.33

Table B2. Linear correlation coefficient between different metric in the models analyzed in this work. Catchments whose NSE prediction is lower than zero are excluded.

Appendix C: Neural Networks Details and Training Losses

We report the architecture details of the encoder (Table C1) of the ENCA models used in this work.

For the encoder, we chose a single layer uni-directional LSTM, followed by a dropout layer (with probability 0.4) and a fully connected layer that maps the LSTM output layer to the predicted output. For the LSTM, we chose a hidden size of 256 (number of memory cells), an initial forget bias of 5.

Input	Layer name	Hyper-parameters	Output
	streamflow	input	(BS, 5478, 1)
streamflow	Conv 1	7, 8, BN, Leakyrelu, DR(0.4)	(BS, 5472, 8)
Conv 1	Avgpool 1	4 (BS, 1368, 8)	
Avgpool 1	Conv 2	5, 16, BN, Leakyrelu, DR(0.4)	(BS, 1364, 16)
Conv 2	Avgpool 2	4	(BS, 341, 16)
Avgpool 2	Conv 3	2, 32, BN, Leakyrelu, DR(0.4)	(BS, 340, 32)
Conv 3	Avgpool 3	4	(BS, 85, 32)
Avgpool 3	Flatten	N/A	(BS, 2720)
Flatten	Linear	BN, Leakyrelu, DR(0.4)	(BS, 512)
Linear	Output	BN	(BS, N)

Table C1. The Convolutional encoder architecture used in this study. Batch Normalization (BN) and Dropout (DR) with probability 0.4 are added among layers. A last BN layer is applied to the decoder output in order to standardize the latent features. N is the number of latent features. The batch size is indicated with BS.

In Figure C1 we report the training and test losses of some models employed. We observe that all the models employed are about to reach a plateau, where the test loss does not decrease anymore. Though convergence is not perfectly reached due to computational limitations, the fact that the test loss is almost at the reachable minimum is an indicator that the models are not overfitting the dataset.

335 Additionally, we report the mean and standard deviations of the latent features of ENCA-5 (Table C2). We can appreciate a small amount of bias, even if the encoder succeeds in preserving the standard deviation of the latent features close to one. We found a similar behaviour in the latent features of other ENCA models (not reported).

Mean	0.23	0.27	-0.07	-0.25	0.28
Standard Deviation	1.15	1.1	0.92	0.67	0.99

Table C2. Mean and standard deviation of the latent features of ENCA-5.

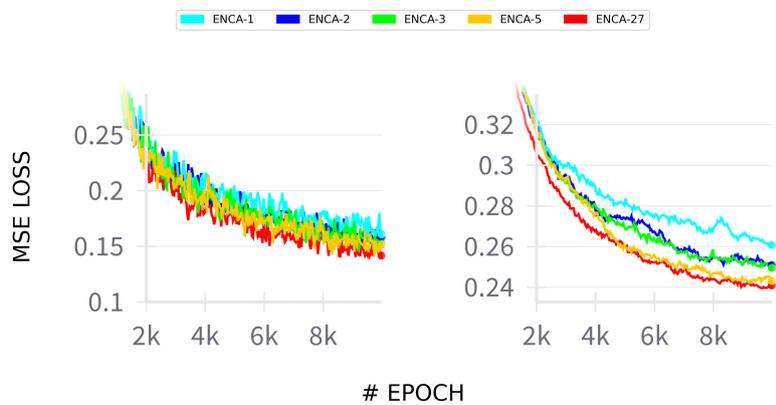


Figure C1. Training MSE loss (left panel) and test loss (right panel) during training for different ENCA models. The curves shown are obtained by averaging the losses across different random restarts. Loss variability across random restarts (not shown) is negligible.

Appendix E: Effect of Random Restart

340 We ascertain that the random restart does not affect much the prediction accuracy (Figure E1). Apart from some catchments, most of them show a consistent behaviour across different random seeds.

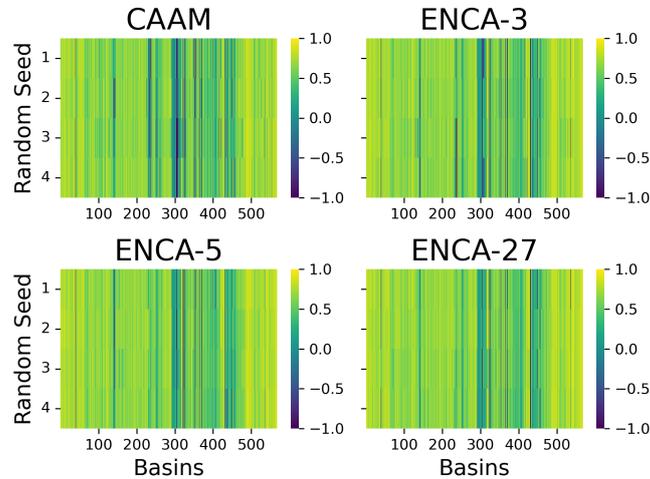


Figure E1. Test NSE values of four random restarts for CAAM (upper left), ENCA-3 (upper right), ENCA-5 (lower left) and ENCA-27 (lower right) for the 568 catchments considered in this study. NSE values are clipped in the interval $[-1, 1]$. We do not observe much performance variability across different random restarts.

We report the correlation matrix between the principal components of the learnt features of ENCA-5 (Figure E2) and ENCA-27 (Figure E3) for different random restarts. We notice a consistency across random restarts and different models. Moreover, the correlation becomes weaker and weaker with the fourth component, indicating that three features carry most of the information related to streamflow prediction.

345

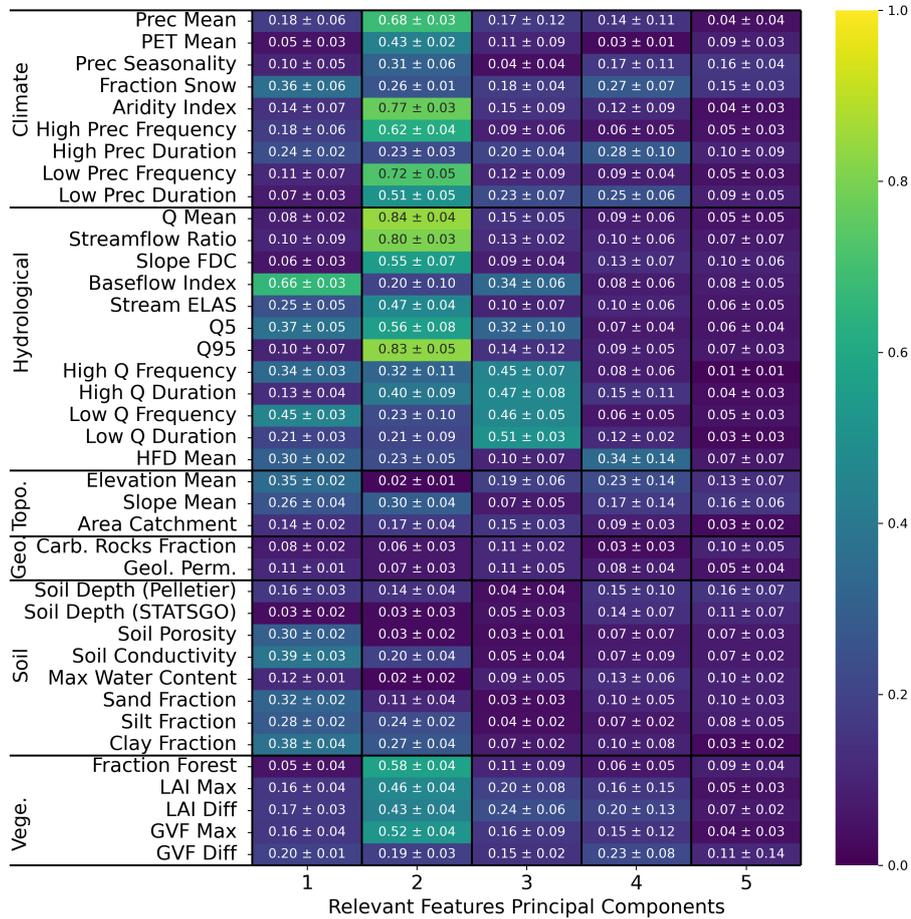


Figure E2. Average (plus minus the standard deviation) of the absolute Spearman correlation of relevant features principal components of ENCA-5 with respect catchment attributes and hydrological signatures across four different random model restarts. The colors refer to the average.

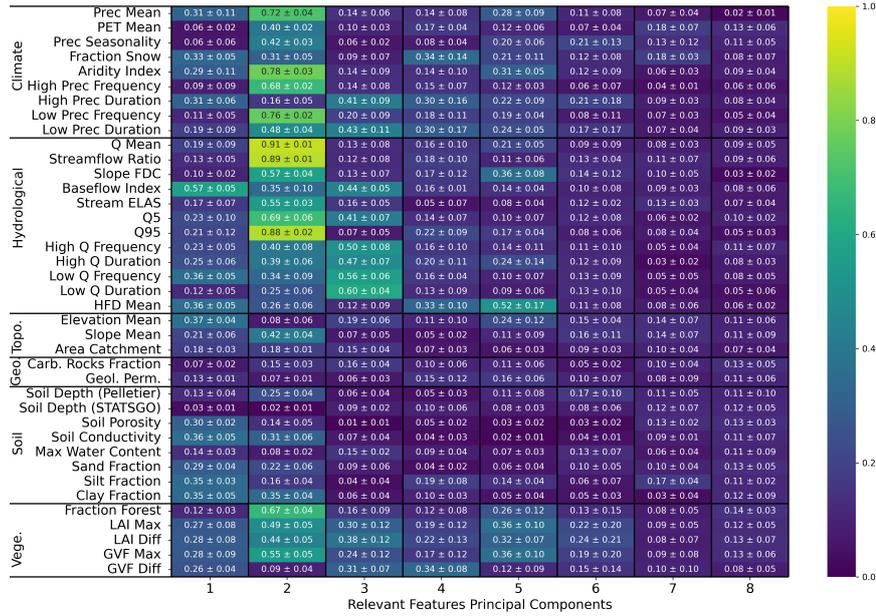


Figure E3. Average (plus minus the standard deviation) of the absolute Spearman correlation of relevant features principal components of ENCA-27 with respect catchment attributes and hydrological signatures across four different random model restarts. The colors refer to the average. For a better visualization, we report only the first eight Principal Components.

Author contributions. CA had the original idea and AB and CA developed the conceptualization and methodology of the study. The idea of using intrinsic dimension is of AM. AB developed the software and conducted all model simulations and their formal analysis. Results were discussed and interpreted between MH, CA, AM, FF and AB. The visualizations and the original draft of the manuscript were prepared by AB, and revisions and editing were provided by MH, CA, AM and FF. Funding was acquired by AM and CA. All authors have read and
350 agreed to the current version of the paper.

Competing interests. At least one of the (co-)authors is a member of the editorial board of Hydrology and Earth System Sciences. The authors also have no other competing interests to declare.

Acknowledgements. This research has been partly founded by the SNSF (Swiss National Science Foundation) grant 200021_208249. A special thanks goes to Antonio Di Noia (Università della Svizzera italiana, ETH Zurich) for his insights and the fruitful discussions. We also
355 thank Fernando Perez Cruz (ETH Zurich), Andreas Scheidegger (Eawag), Marco Baity-Jesi (Eawag) and Dmitri Kavetski (The University of Adelaide) for the insightful discussions.

References

- Addor, N., Newman, A. J., Mizukami, N., and Clark, M. P.: The CAMELS data set: catchment attributes and meteorology for large-sample studies, *Hydrology and Earth System Sciences*, 21, 5293–5313, <https://doi.org/10.5194/hess-21-5293-2017>, 2017.
- 360 Albert, C., Künsch, H., and Scheidegger, A.: A Simulated Annealing Approach to Approximate Bayes Computations, *Statistics and Computing*, 25, 1217–1232, <https://doi.org/https://doi.org/10.1007/s11222-014-9507-8>, 2015.
- Albert, C., Ulzega, S., Ozdemir, F., Perez-Cruz, F., and Mira, A.: Learning Summary Statistics for Bayesian Inference with Autoencoders, *SciPost Phys. Core*, 5, 043, <https://doi.org/10.21468/SciPostPhysCore.5.3.043>, 2022.
- Allegra, M., Facco, E., Denti, F., Laio, A., and Mira, A.: Data segmentation based on the local intrinsic dimension, *Scientific Reports*, 10, 16 449, <https://doi.org/10.1038/s41598-020-72222-0>, 2020.
- 365 Botterill, T. E. and McMillan, H. K.: Using Machine Learning to Identify Hydrologic Signatures With an Encoder–Decoder Framework, *Water Resources Research*, 59, e2022WR033 091, <https://doi.org/https://doi.org/10.1029/2022WR033091>, 2023.
- Denti, F., Doimo, D., Laio, A., and Mira, A.: The generalized ratios intrinsic dimension estimator, *Scientific Reports*, 12, 20 005, <https://doi.org/10.1038/s41598-022-20991-1>, 2022.
- 370 Edijanto, N. D. O. N., Yang, X., Makhlof, Z., and Michel, C.: GR3J: a daily watershed model with three free parameters, *Hydrological Sciences Journal*, 44, 263–277, <https://doi.org/10.1080/02626669909492221>, 1999.
- Facco, E., d’Errico, M., Rodriguez, A., and Laio, A.: Estimating the intrinsic dimension of datasets by a minimal neighborhood information, *Scientific Reports*, 7, 12 140, <https://doi.org/10.1038/s41598-017-11873-y>, 2017.
- Fenicia, F., Kavetski, D., Reichert, P., and Albert, C.: Signature-Domain Calibration of Hydrological Models Using Approximate Bayesian Computation: Empirical Analysis of Fundamental Properties, *Water Resources Research*, 54, 3958–3987, <https://doi.org/https://doi.org/10.1002/2017WR021616>, 2018.
- 375 Glielmo, A., Zeni, C., Cheng, B., Csányi, G., and Laio, A.: Ranking the information content of distance measures, *PNAS Nexus*, 1, pgac039, <https://doi.org/10.1093/pnasnexus/pgac039>, 2022.
- Gnann, S. J., McMillan, H. K., Woods, R. A., and Howden, N. J. K.: Including Regional Knowledge Improves Baseflow Signature Predictions in Large Sample Hydrology, *Water Resources Research*, 57, e2020WR028 354, <https://doi.org/https://doi.org/10.1029/2020WR028354>, e2020WR028354 10.1029/2020WR028354, 2021.
- 380 Gupta, H. V., Kling, H., Yilmaz, K. K., and Martinez, G. F.: Decomposition of the mean squared error and NSE performance criteria: Implications for improving hydrological modelling, *Journal of Hydrology*, 377, 80–91, <https://doi.org/https://doi.org/10.1016/j.jhydrol.2009.08.003>, 2009.
- 385 Höge, M., Scheidegger, A., Baity-Jesi, M., Albert, C., and Fenicia, F.: Improving hydrologic models for predictions and process understanding using neural ODEs, *Hydrology and Earth System Sciences*, 26, 5085–5102, <https://doi.org/10.5194/hess-26-5085-2022>, 2022.
- Jakeman, A. J. and Hornberger, G. M.: How much complexity is warranted in a rainfall-runoff model?, *Water Resources Research*, 29, 2637–2649, <https://doi.org/https://doi.org/10.1029/93WR00877>, 1993.
- Kingma, D. P. and Ba, J.: Adam: A Method for Stochastic Optimization, in: 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings, edited by Bengio, Y. and LeCun, Y., <http://arxiv.org/abs/1412.6980>, 2015.
- 390

- Kiraz, M., Coxon, G., and Wagener, T.: A Signature-Based Hydrologic Efficiency Metric for Model Calibration and Evaluation in Gauged and Ungauged Catchments, *Water Resources Research*, 59, e2023WR035321, <https://doi.org/https://doi.org/10.1029/2023WR035321>, e2023WR035321 2023WR035321, 2023.
- 395 Klotz, D., Kratzert, F., Gauch, M., Keefe Sampson, A., Brandstetter, J., Klambauer, G., Hochreiter, S., and Nearing, G.: Uncertainty estimation with deep learning for rainfall–runoff modeling, *Hydrology and Earth System Sciences*, 26, 1673–1693, <https://doi.org/10.5194/hess-26-1673-2022>, 2022.
- Kratzert, F.: CAMELS Extended NLDAS Forcing Data, <https://doi.org/10.4211/hs.0a68bfd7ddf642a8be9041d60f40868c>, 2019.
- Kratzert, F., Klotz, D., Herrnegger, M., Sampson, A. K., Hochreiter, S., and Nearing, G. S.: Toward Improved Predictions in Ungauged Basins: Exploiting the Power of Machine Learning, *Water Resources Research*, 55, 11344–11354, <https://doi.org/https://doi.org/10.1029/2019WR026065>, 2019.
- 400 Kratzert, F., Klotz, D., Hochreiter, S., and Nearing, G. S.: A note on leveraging synergy in multiple meteorological data sets with deep learning for rainfall–runoff modeling, *Hydrology and Earth System Sciences*, 25, 2685–2703, <https://doi.org/10.5194/hess-25-2685-2021>, 2021.
- 405 Lees, T., Reece, S., Kratzert, F., Klotz, D., Gauch, M., De Bruijn, J., Kumar Sahu, R., Greve, P., Slater, L., and Dadson, S. J.: Hydrological concept formation inside long short-term memory (LSTM) networks, *Hydrology and Earth System Sciences*, 26, 3079–3101, <https://doi.org/10.5194/hess-26-3079-2022>, 2022.
- McMillan, H.: Linking hydrologic signatures to hydrologic processes: A review, *Hydrological Processes*, 34, 1393–1409, <https://doi.org/https://doi.org/10.1002/hyp.13632>, 2020a.
- 410 McMillan, H.: A review of hydrologic signatures and their applications, *Wiley Interdisciplinary Reviews: Water*, <https://doi.org/10.1002/wat2.1499>, 2020b.
- Mohammadi, B.: A review on the applications of machine learning for runoff modeling, *Sustainable Water Resources Management*, 7, 98, <https://doi.org/10.1007/s40899-021-00584-y>, 2021.
- Molnar, C.: *Interpretable Machine Learning*, open source online book, 2 edn., <https://christophm.github.io/interpretable-ml-book>, 2024.
- 415 Molnar, C., Casalicchio, G., and Bischl, B.: Interpretable machine learning—a brief history, state-of-the-art and challenges, in: *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pp. 417–431, Springer, https://doi.org/10.1007/978-3-030-65965-3_28, 2020.
- Nash, J. and Sutcliffe, J.: River flow forecasting through conceptual models part I — A discussion of principles, *Journal of Hydrology*, 10, 282–290, [https://doi.org/https://doi.org/10.1016/0022-1694\(70\)90255-6](https://doi.org/https://doi.org/10.1016/0022-1694(70)90255-6), 1970.
- 420 Newman, A. J., Clark, M. P., Sampson, K., Wood, A., Hay, L. E., Bock, A., Viger, R. J., Blodgett, D., Brekke, L., Arnold, J. R., Hopson, T., and Duan, Q.: Development of a large-sample watershed-scale hydrometeorological data set for the contiguous USA: data set characteristics and assessment of regional variability in hydrologic model performance, *Hydrology and Earth System Sciences*, 19, 209–223, <https://doi.org/10.5194/hess-19-209-2015>, 2015.
- Olden, J. D. and Poff, N. L.: Redundancy and the choice of hydrologic indices for characterizing streamflow regimes, *River Research and Applications*, 19, 101–121, <https://doi.org/https://doi.org/10.1002/rra.700>, 2003.
- 425 Perrin, C., Michel, C., and Andréassian, V.: Improvement of a parsimonious model for streamflow simulation, *Journal of Hydrology*, 279, 275–289, [https://doi.org/https://doi.org/10.1016/S0022-1694\(03\)00225-7](https://doi.org/https://doi.org/10.1016/S0022-1694(03)00225-7), 2003.
- Wagener, T., McIntyre, N., Lees, M. J., Wheeler, H. S., and Gupta, H. V.: Towards reduced uncertainty in conceptual rainfall-runoff modelling: dynamic identifiability analysis, *Hydrological Processes*, 17, 455–476, <https://doi.org/https://doi.org/10.1002/hyp.1135>, 2003.

- 430 Wagener, T., Sivapalan, M., Troch, P., and Woods, R.: Catchment Classification and Hydrologic Similarity, *Geography Compass*, 1, 901–931, <https://doi.org/https://doi.org/10.1111/j.1749-8198.2007.00039.x>, 2007.
- Zar, J. H.: Spearman Rank Correlation, John Wiley & Sons, Ltd, <https://doi.org/https://doi.org/10.1002/0470011815.b2a15150>, 2005.