# When to Make Mountains out of Molehills: The Pros and Cons of Simple and Complex Model Calibration Procedures

Katie A. Smith[1], Lucy J. Barker[1], Shaun Harrigan[1,2], Christel Prudhomme[2,1,3], Jamie Hannaford[1], Maliko Tanguy[1] and Simon Parry[1,3]

[1]Centre for Ecology and Hydrology, Wallingford, England, 0X10 8BB, [2]ECMWF, Reading, England, RG2 9AX, [3]Department of Geography Loughborough University, England, LE11 3TU

**Centre for Ecology & Hydrology**
NATURAL ENVIRONMENT RESEARCH COUNCIL
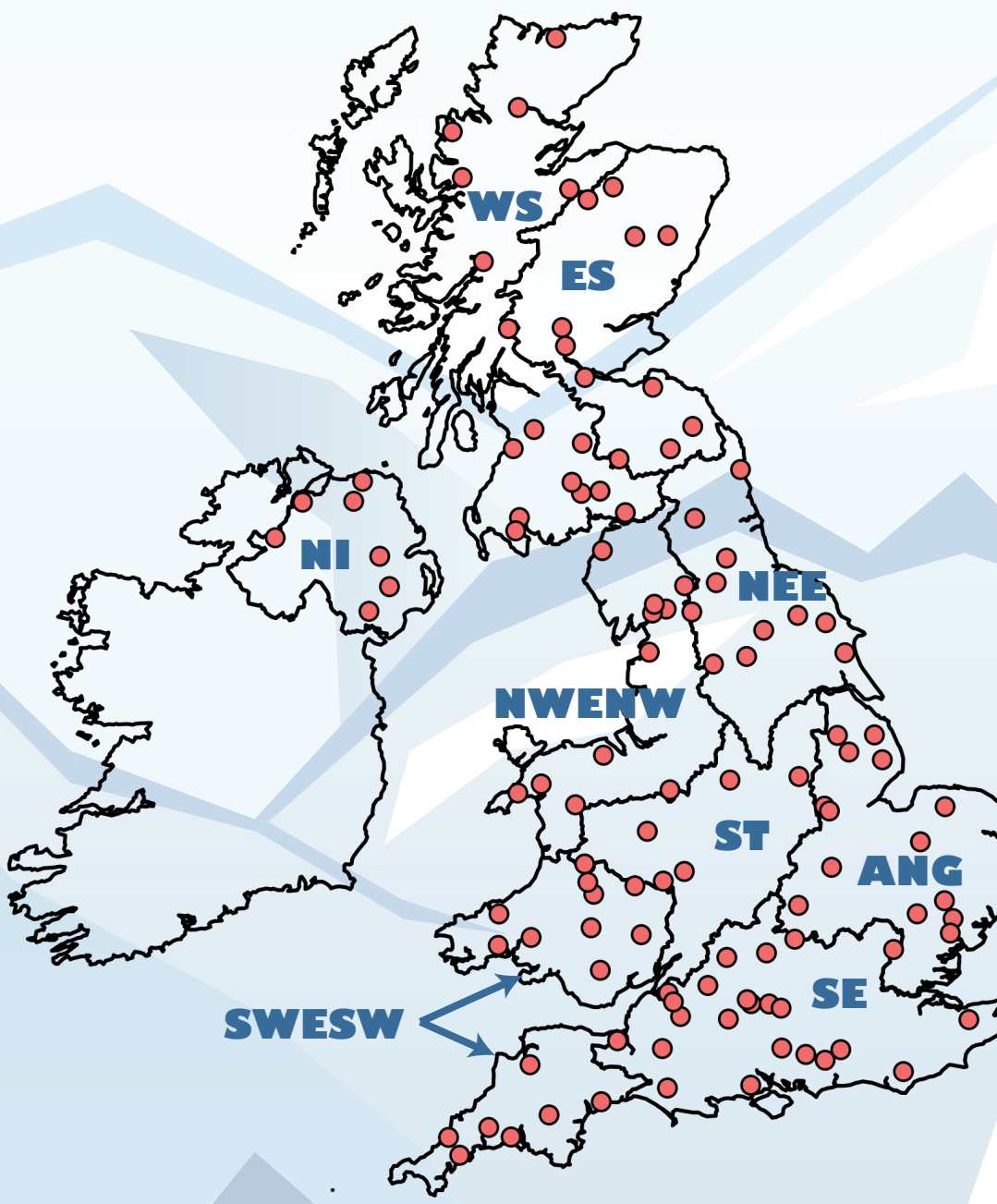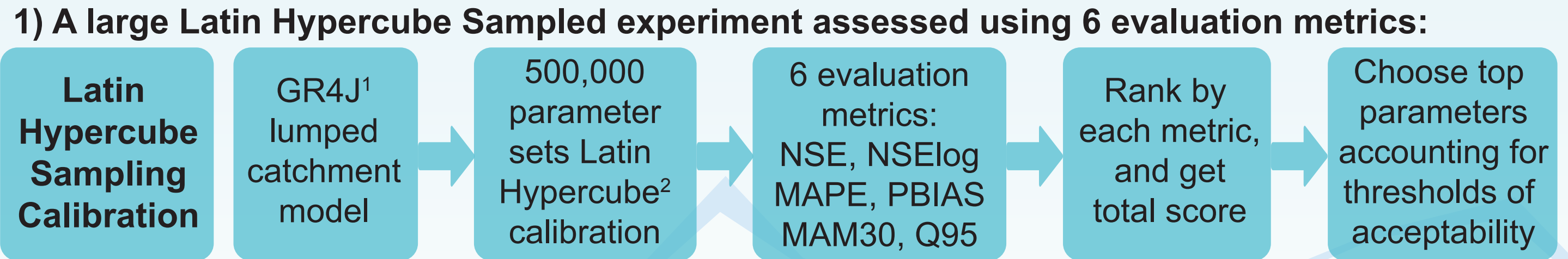enquiries@ceh.ac.uk
www.ceh.ac.uk

## Introduction

♦ Environmental models are used to simulate processes that cannot otherwise be examined.

♦ Models are made up of parameters that represent physical or conceptual variables.

♦ The values assigned to these parameters are determined through model calibration.

♦ Calibration techniques vary and require evaluation metrics to compare model outputs with observations.

♦ Models often include quick to run inbuilt automatic calibration functions, but these apply a limited selection of evaluation metrics.

♦ More complex methods that randomly sample the parameter space can also be used.

So, why would we use these more complex methods, when we can use simple methods?

### Why would we make mountains out of molehills?

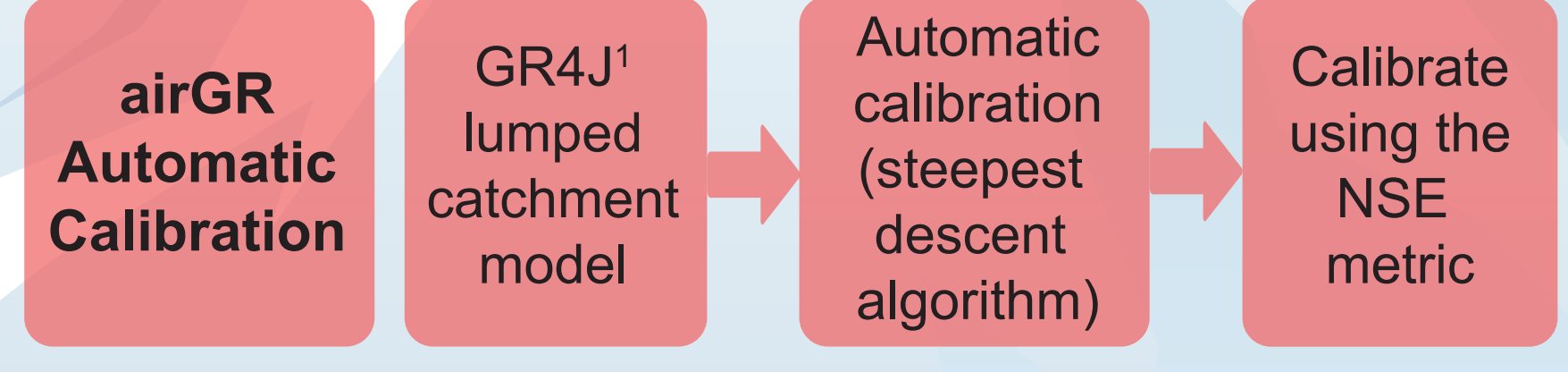This poster explores the pros and cons of each method applied to drought simulation in the UK.

## Methods

Two methods of model calibration have been compared using the GR4J lumped catchment model for 115 near-natural catchments across the UK over the period 1982-2014:
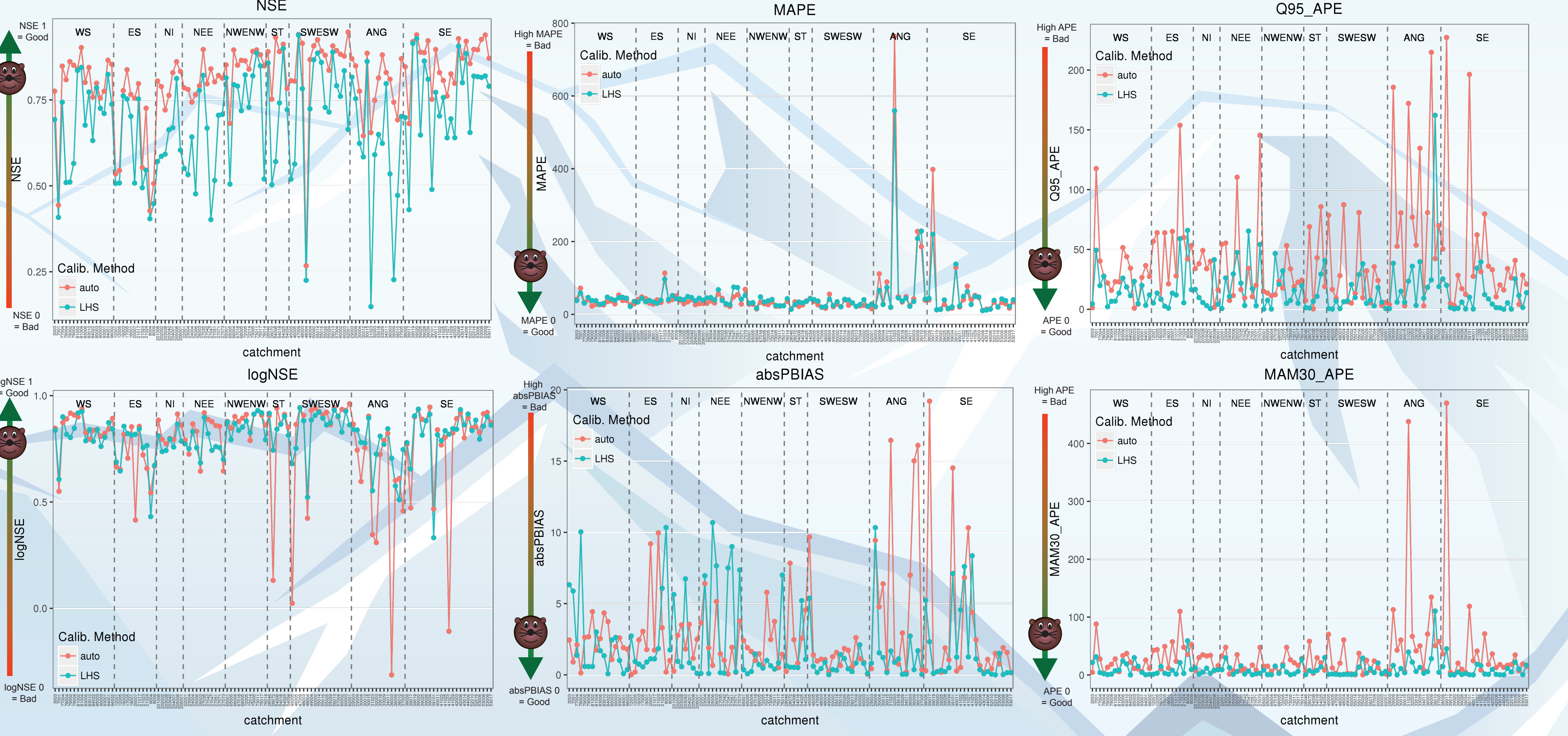
**1) A large Latin Hypercube Sampled experiment assessed using 6 evaluation metrics:**

Latin Hypercube Sampling Calibration → GR4J[1] lumped catchment model → 500,000 parameter sets Latin Hypercube[2] calibration → 6 evaluation metrics: NSE, NSElog MAPE, PBIAS MAM30, Q95 → Rank by each metric, and get total score → Choose top parameters accounting for thresholds of acceptability

6 evaluation metrics:
♦ Nash Sutcliffe (NSE),
♦ NSE on log flows (logNSE),
♦ Mean Absolute Percent Error (MAPE),
♦ Percent Bias (PBIAS),
♦ Absolute Percent Error (APE) in Mean Annual Minima on 30 day flow accumulations (MAM30_APE),
♦ APE on the flow exceeded 95 percent of the time (Q95_APE).

**2) The GR4J inbuilt automatic calibration function using the Nash Sutcliffe Efficiency (NSE)**

airGR Automatic Calibration → GR4J[1] lumped catchment model → Automatic calibration (steepest descent algorithm) → Calibrate using the NSE metric

Black boundaries indicate UK hydroclimatic regions, pink dots show gauging station locations for the 115 near-natural catchments included in this study.

## Comparing Metric Scores

Here, the two methods of model calibration are compared. The automatic calibration used just NSE, whilst the Latin Hypercube Sampling (LHS) method used 6 metrics. These plots show that for NSE, the automatic calibration produced better results across the majority of the 115 catchments. However when other metrics were considered, the LHS calibration method gave comparable, or better results. This is particularly evident when looking at the drought metrics Q95 and MAM30.



## Drought Events

This section explores the model calibrations during the 1975-1976 drought event in the UK.
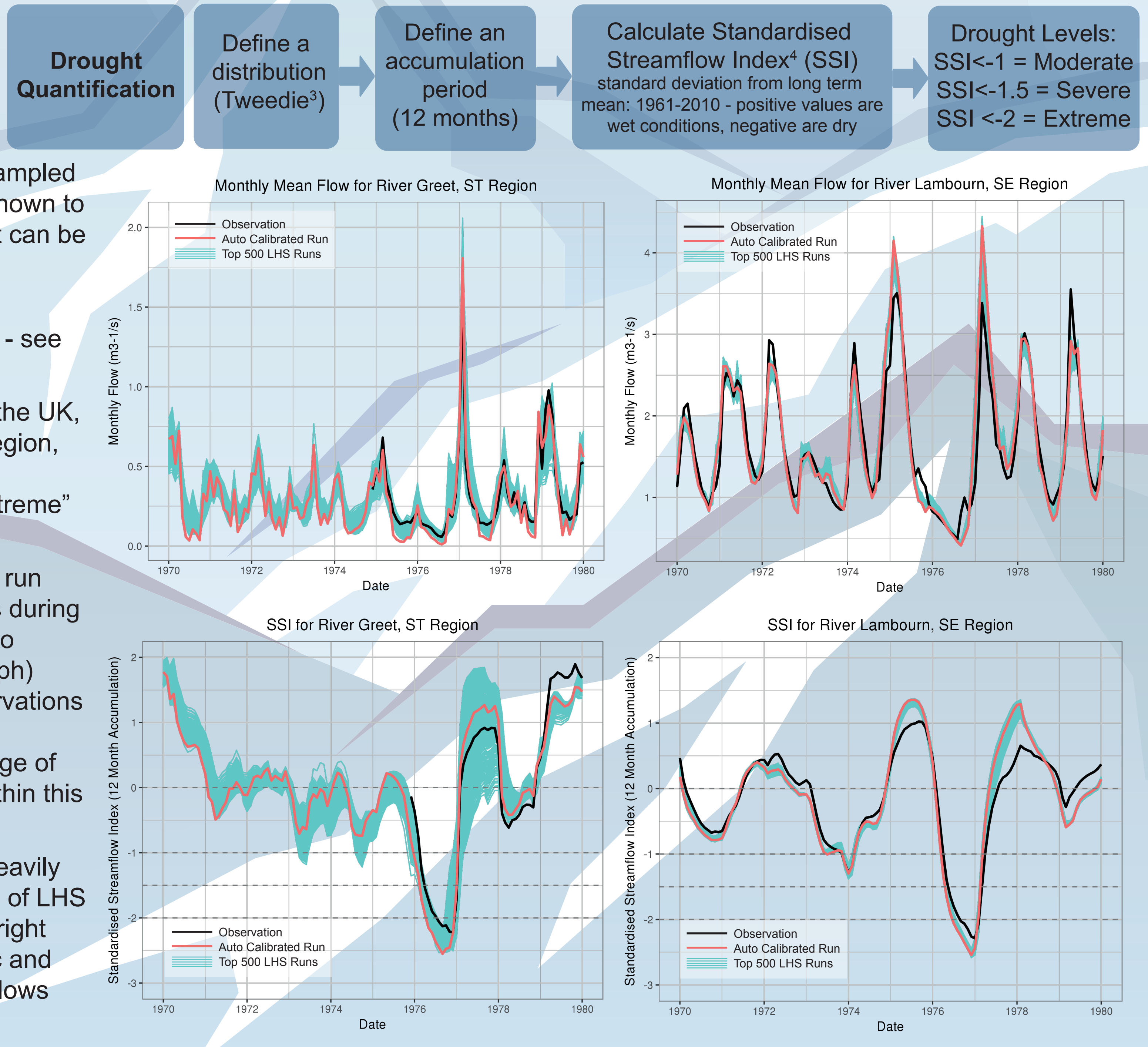
Here, the top 500 Latin Hypercube Sampled (LHS) model parameterisations are shown to represent the probabilistic results that can be explored using this approach.

Drought can be quantified using the Standardised Streamflow Index (SSI) - see the diagram above on the right.

These plots show two catchments in the UK, the River Greet in the Severn Trent region, and the River Lambourn in Southern England. Both rivers experienced "extreme" drought (SSI <-2) in this period.

In the Greet, the automatic calibrated run siginifcantly underestimates low flows during 1975-1980 (top left graph). This is also shown in the SSI plot (bottom left graph) where the SSI is lower than the observations in 1976-1977. Although the 500 LHS calibrated runs show quite a wide range of model results, the observations sit within this range.

For the Lambourn, where flows are heavily dominated by groundwater, the range of LHS model results is much narrower (see right hand graphs), and both the automatic and the LHS runs underestimate the low flows during the drought event.

**Drought Quantification** → Define a distribution (Tweedie[3]) → Define an accumulation period (12 months) → Calculate Standardised Streamflow Index[4] (SSI) standard deviation from long term mean: 1961-2010 - positive values are wet conditions, negative are dry → Drought Levels: SSI<-1 = Moderate SSI<-1.5 = Severe SSI <-2 = Extreme



Monthly Mean Flow for River Greet, ST Region



Monthly Mean Flow for River Lambourn, SE Region



SSI for River Greet, ST Region



SSI for River Lambourn, SE Region

## Key Messages

The table below summarises some of the key differences between the automatic and the Latin Hypercube Sampled (LHS) calibration techniques.

Automatic calibration algorithms are very efficient, and can yield good model results.

The LHS calibration technique is a far more versatile approach which allows the modeller to calibrate the model to the one, or many, evaluation metrics they are interested in.

However, this approach requires a significant amount of time and computational resource when modelling multiple catchments.

| | Automatic Calibration | Latin Hypercube Calibration |
|---|---|---|
| **Speed** | Quick (2 Mins) | Slow (2 Weeks) |
| **Metrics Available** | Limited by developer (usually only 2 or 3) | Any |
| **Application** | Included metrics are usually mid to high flow | Any |
| **Multi-Objective Optimisation** | No | Yes |
| **Output Type** | Deterministic (1 answer) | Probabilistic or Deterministic (top result) |

Sometimes, a model may perform badly however you calibrate it!

References:
1 Perrin, C. et al (2003) Improvement of a parsimonious model for streamflow simulation. Journal of Hydrology 279, pp 275-289;
2 Cheng and Druzdel (2000) Latin Hypercube Sampling in Bayesian Networks. Uncertain Reasoning. American Association for Artificial Intelligence;
3 Svensson, C. et al (2017) Statistical distributions for monthly aggregations of precipitation and streamflow in drought indicator applications. Water Resources Research, 53(2) pp 999-1018;
4 Vicente-Serrano, S. M. et al (2012) Accurate computation of a streamflow drought index. Journal of Hydrologic Engineering 17(2), pp 318-332;

IMPETUS IMPROVING PREDICTIONS OF UK DROUGHT

Historic Droughts

NERC SCIENCE OF THE ENVIRONMENT