

Quantifying and Visualising Model Performance

About this talk

Some background

Why model performance data is so useful

Measures

Process

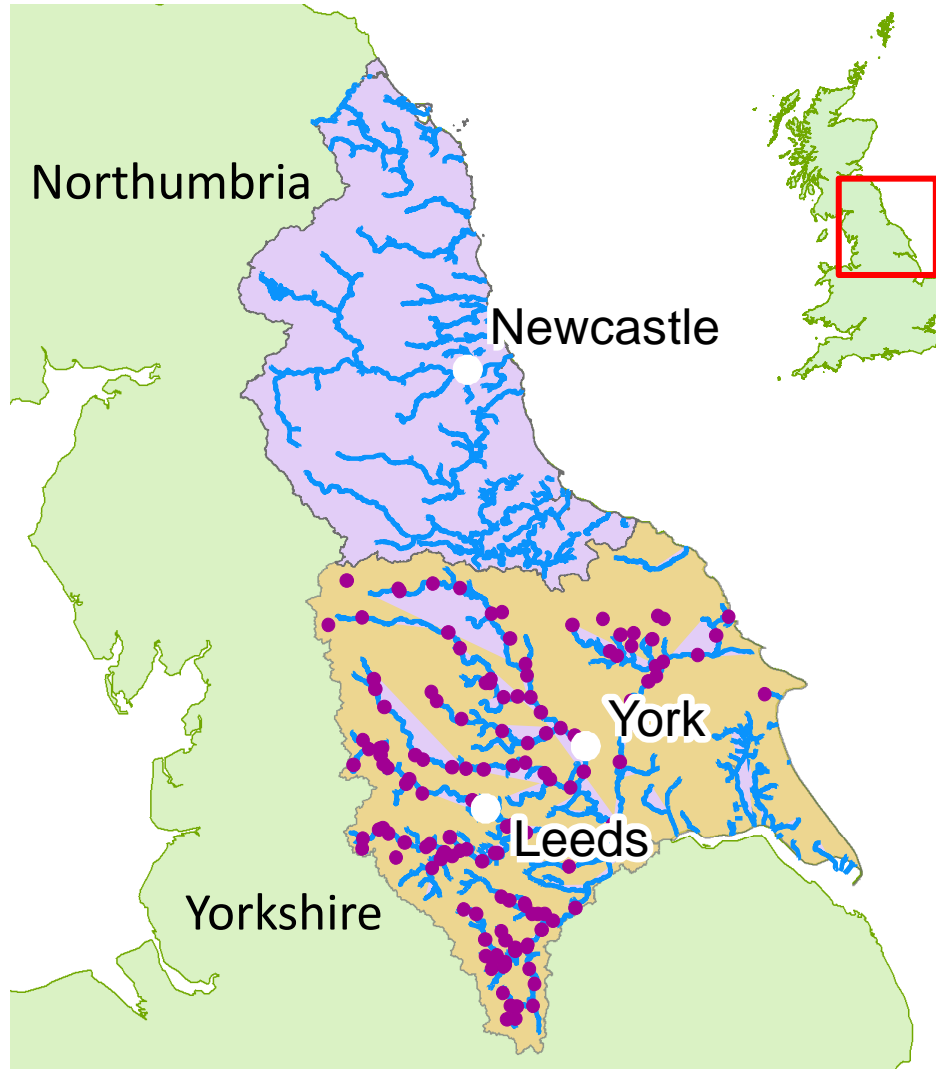
The future



Some background

Quantifying and visualising performance

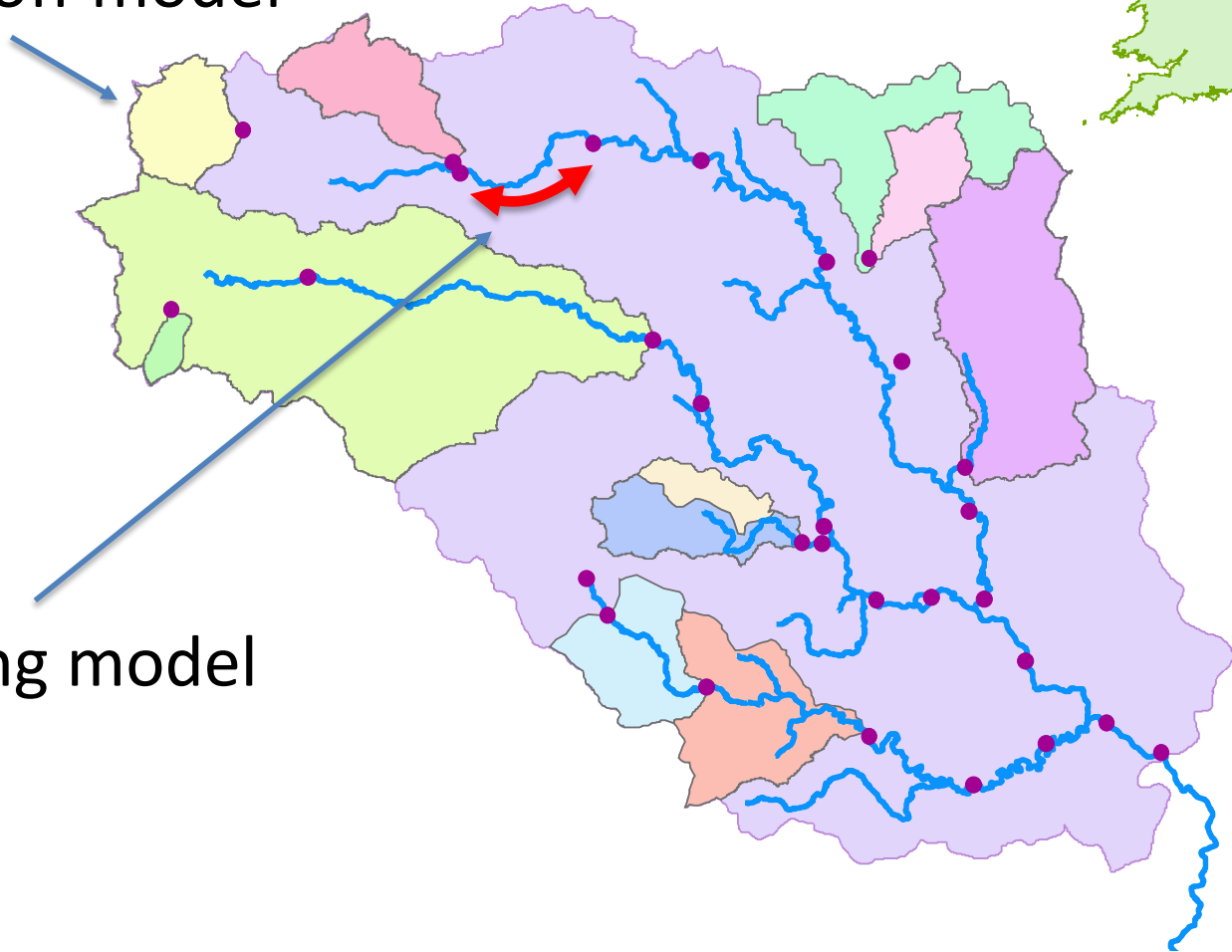
Northumbria and Yorkshire



Ouse: example network

Rainfall runoff model

Flow routing model

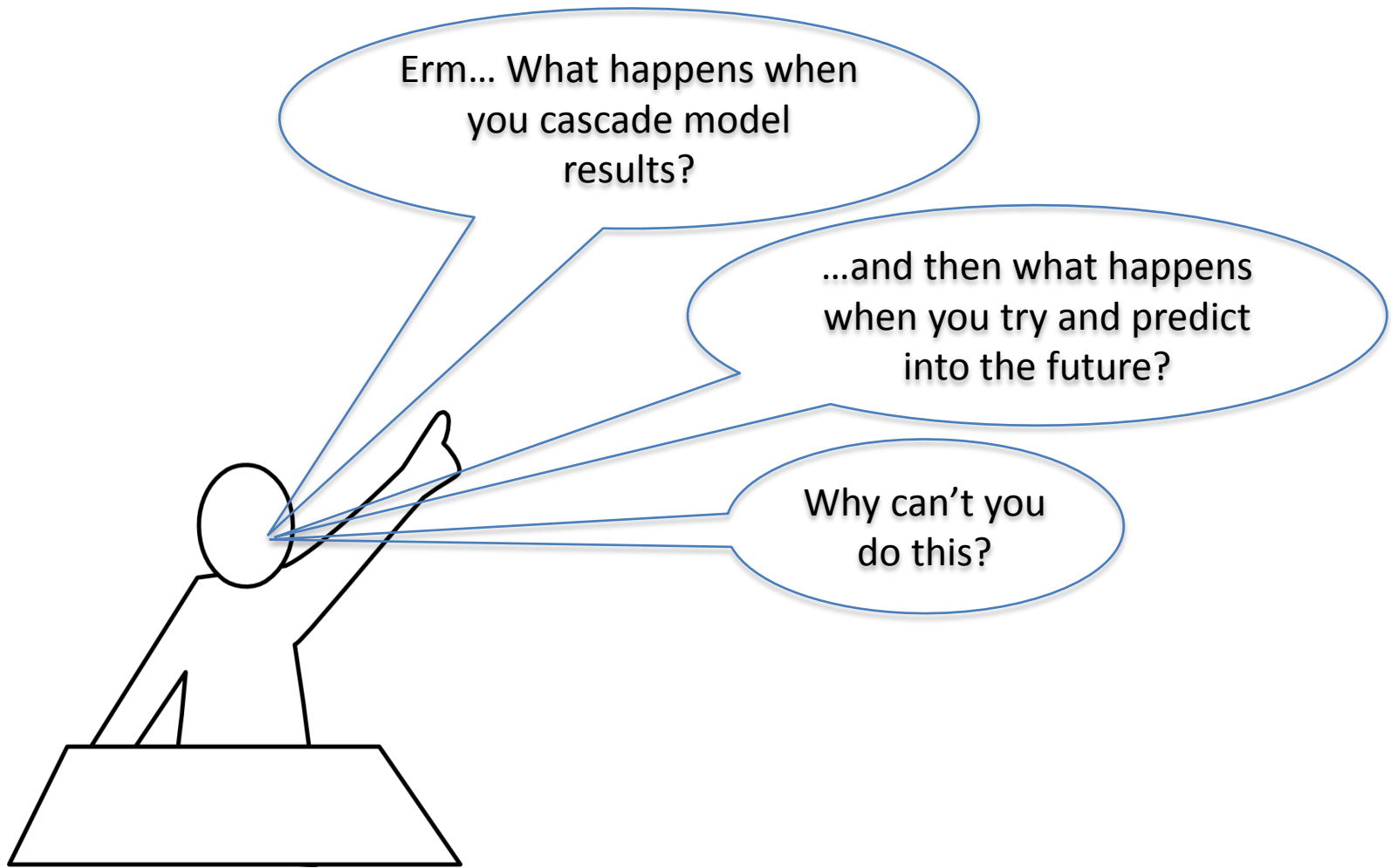


Model development for RFFS

- Develop individual models
- Calibrate error model
- Supplier integrates
- Hope for the best



Awkward questions



Early attempt

Figure 6.6 : No error correction but foreknowledge of rainfall

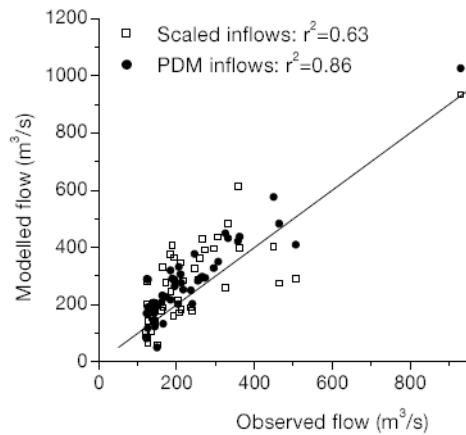
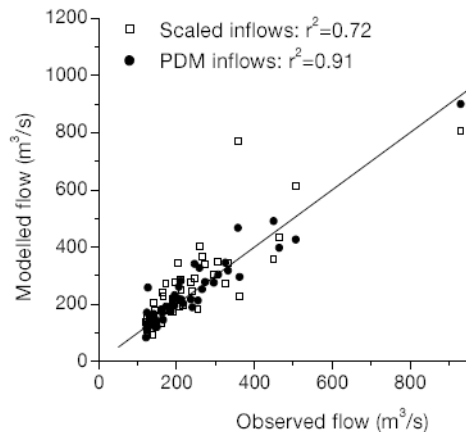
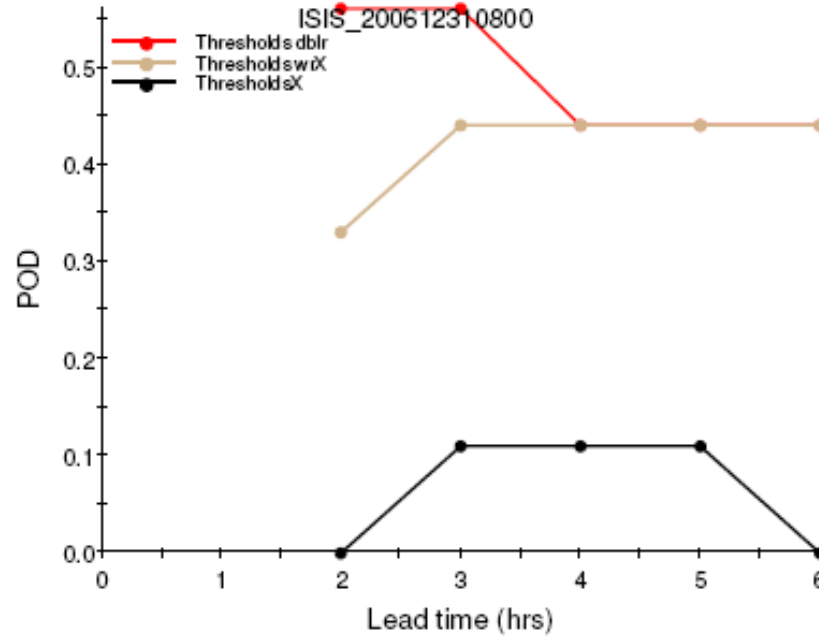


Figure 6.8 : Data available till 3 hours before peak



Real time performance measures for Colliers



POD table

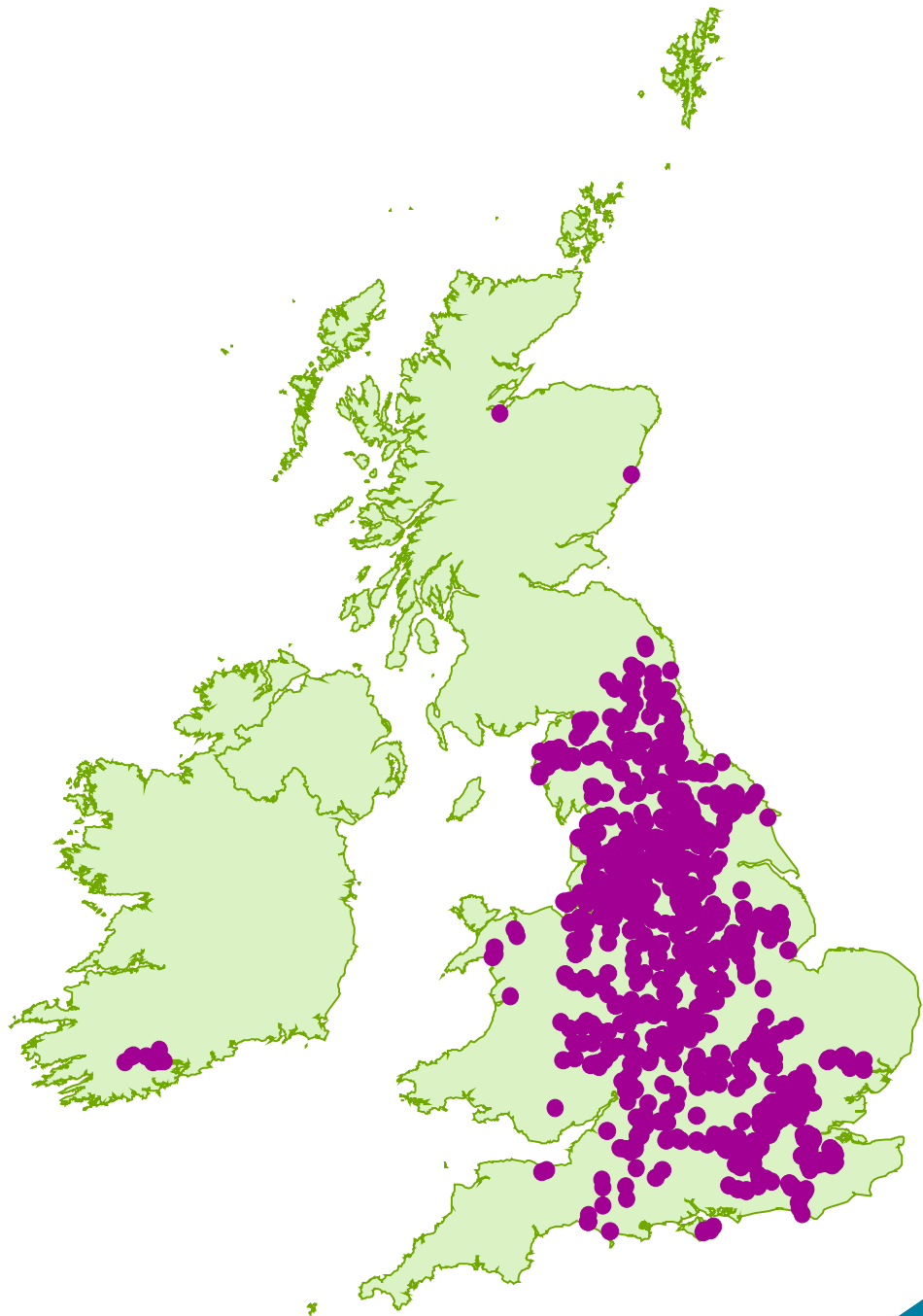
Scenario	N	POD at increasing lead time (hrs)				
		t-2	t-3	t-4	t-5	t-6
Thresholds dblr	9	0.56	0.56	0.44	0.44	0.44
Thresholds wrX	9	0.33	0.44	0.44	0.44	0.44
ThresholdsX	9	0.00	0.11	0.11	0.11	0.00



And by 2016...

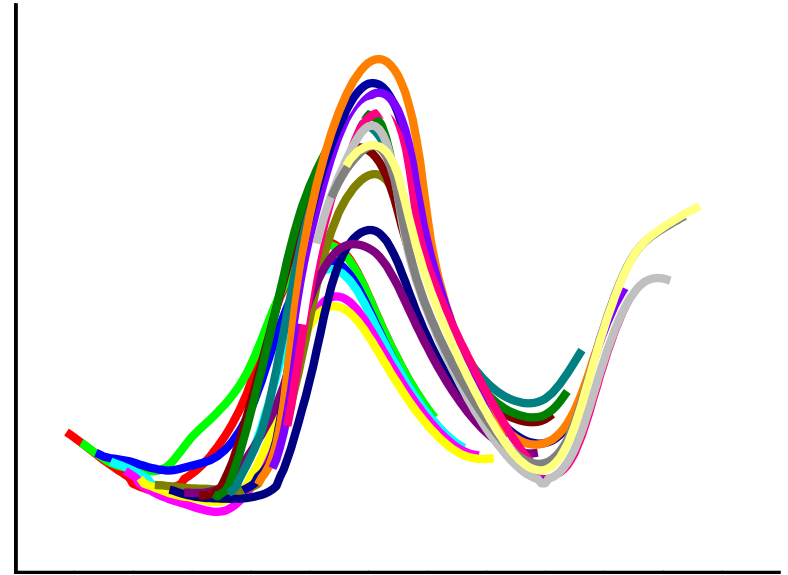
>950 locations
tested

Well developed
software, skills
and methods



Four levels of testing

- Individual component
- Cascaded components
- With data assimilation
- With forecast rainfall

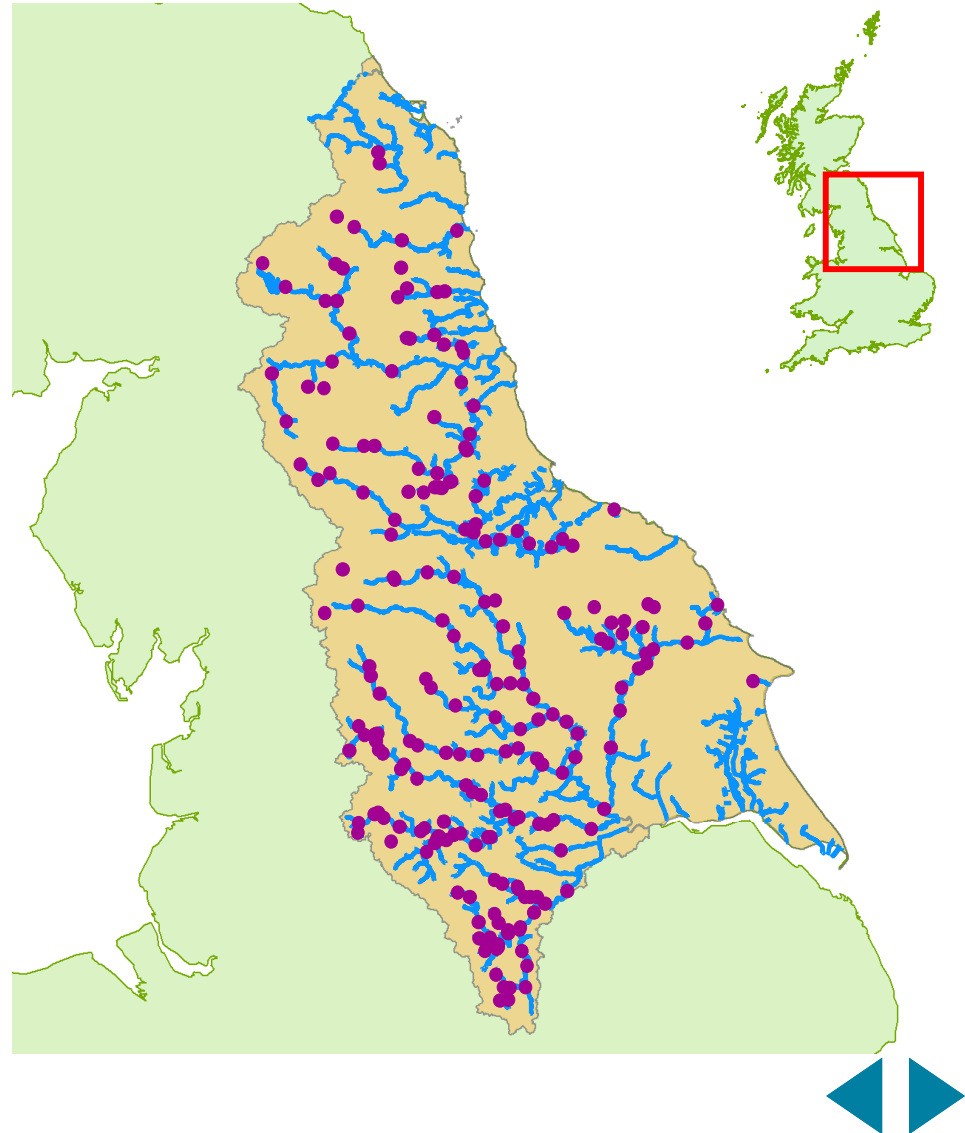


Why performance data are useful

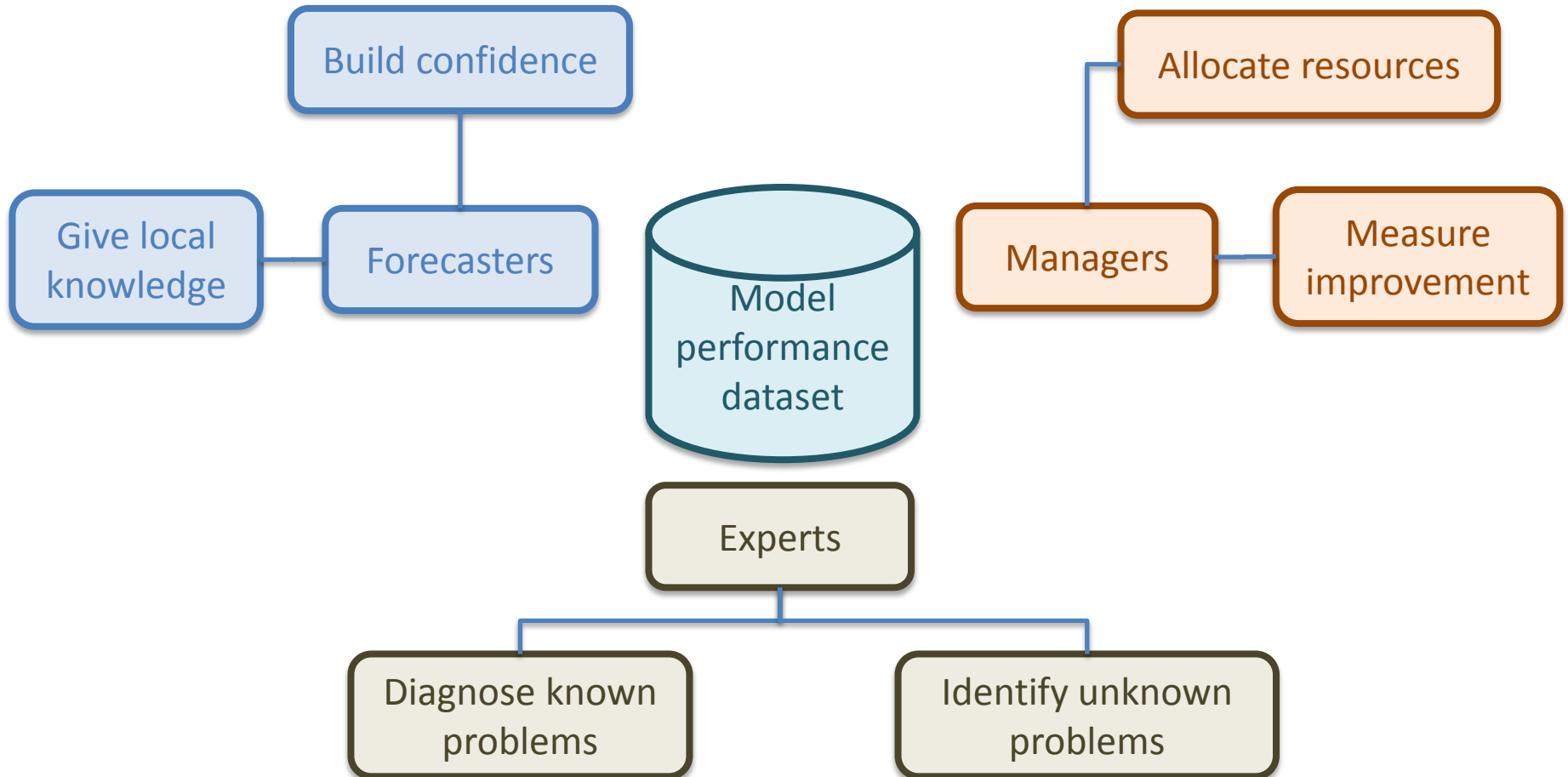
Quantifying and visualising performance

Leeds forecast centre

- Operating since early 1990s
- 210 forecast locations
- 274 models



Uses for performance data

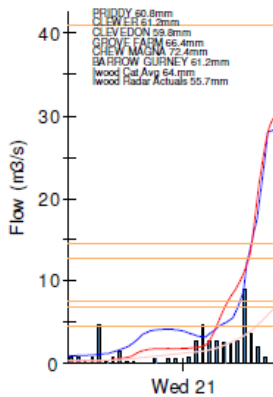


How do you quantify performance?

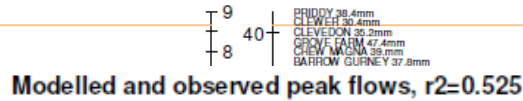
Quantifying and visualising performance

How do you quantify performance?

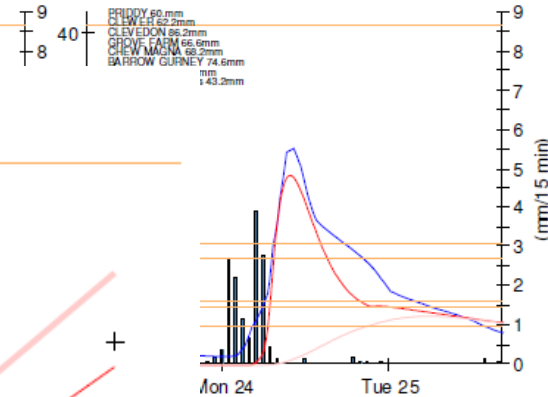
1. 21 Nov 2012 13:15 (56-64-72mm)



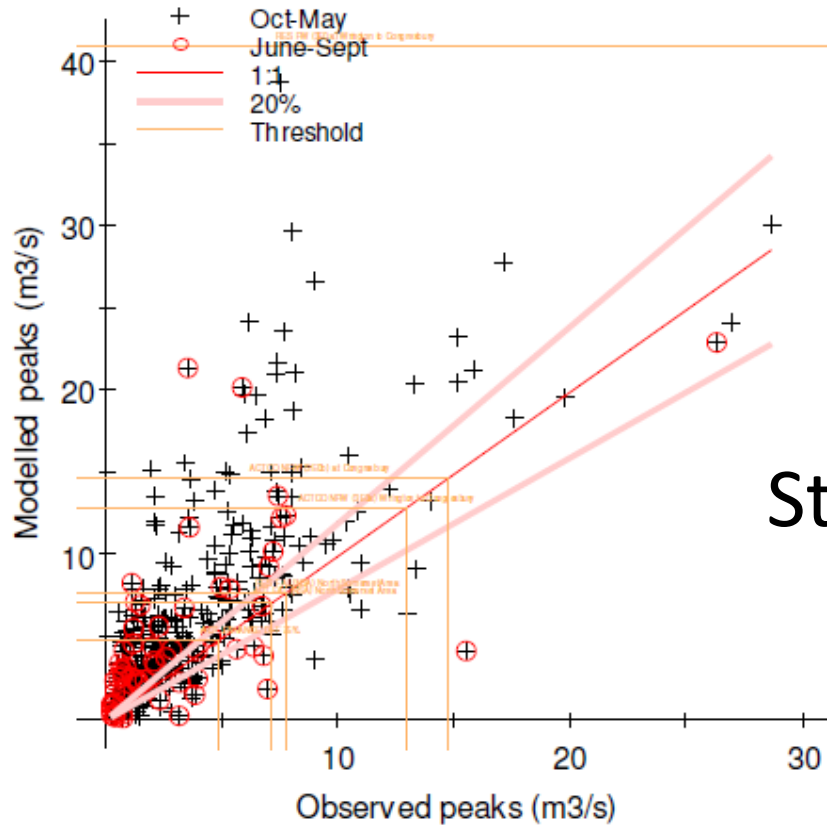
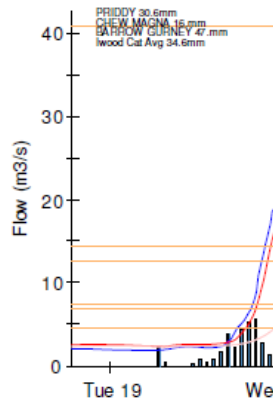
2. 25 Nov 2012 04:00 (30-35-47mm)



3. 24 Sep 2012 10:15 (43-65-86mm)



4. 20 Jan 1999 00:15 (1)



Still Good?



Types of simulation

- Real time
- Simulation



Questions forecasters are asked



Real time measures

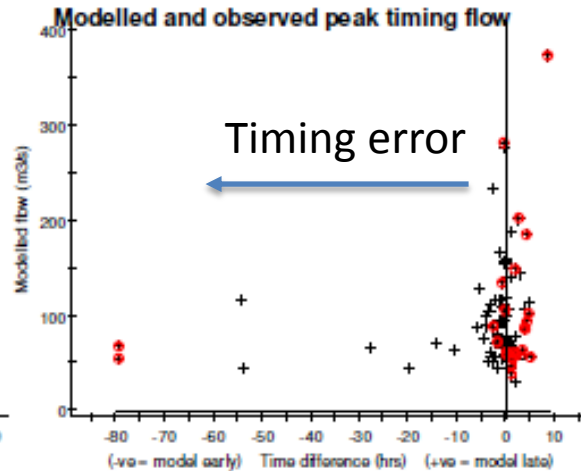
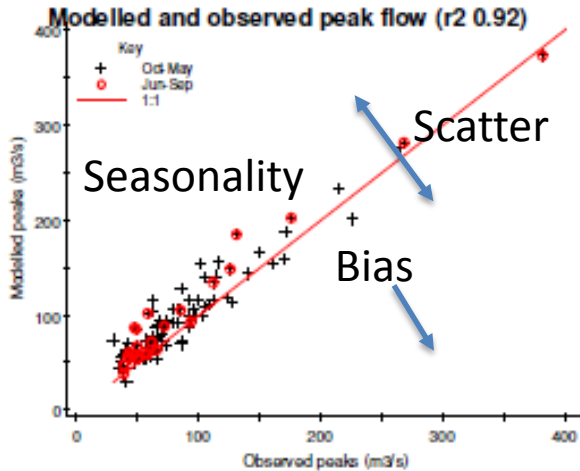
- Probability of Detection (POD)
- False Alarm Rate
- Peak error

Categorising performance

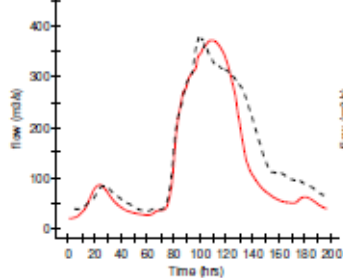
Grade ²	Description ¹	Probability of Detection (POD)	False Alarm Rate (FAR)
A+	Exceeds targets	$POD \geq 0.8$	$FAR \leq 0.2$
A	Meets targets	$0.8 > POD \geq 0.7$	$0.2 < FAR \leq 0.3$
B	Meets targets with tolerance	$POD \geq \text{target with } \pm 0.2\text{m tolerance}$	$FAR \leq \text{target with } \pm 0.2\text{m tolerance}$
C	Does not meet target	$0.7 > POD \geq 0.5$	$0.3 < FAR \leq 0.5$
D	Significantly below target	$0.5 > POD \geq 0.3$	$0.5 < FAR \leq 0.7$
E ²	Poor	$POD < 0.3$	$FAR > 0.7$



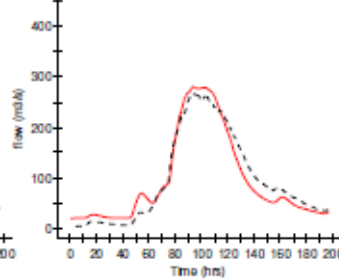
Simulation measures



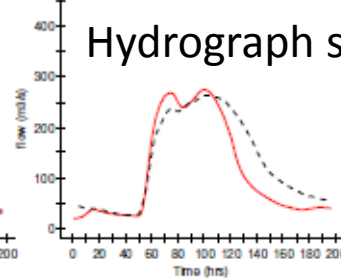
1. 26 Jun 2007 08:00



2. 16 Jun 2007 02:00

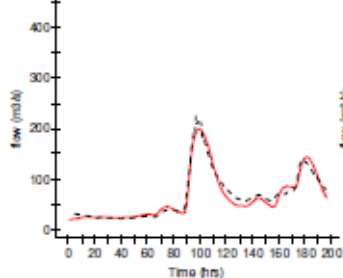


3. 08 Nov 2000 00:15

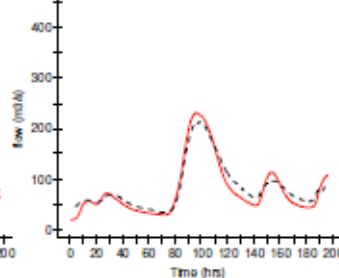


Hydrograph shape

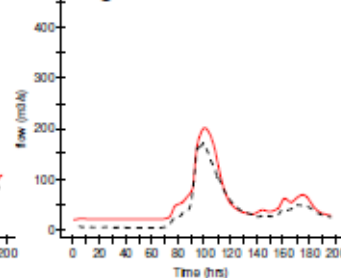
4. 18 Jan 2007 20:15



5. 30 Dec 2002 15:15



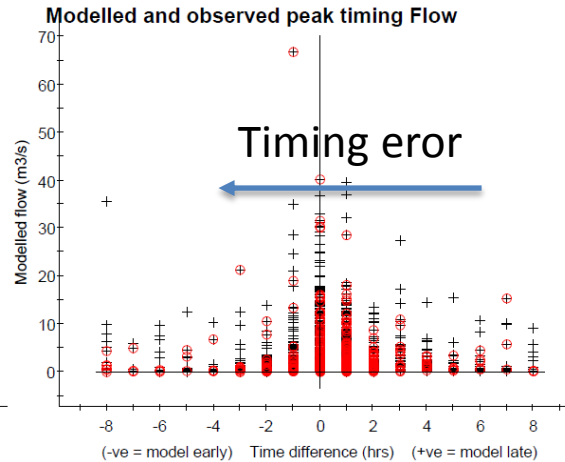
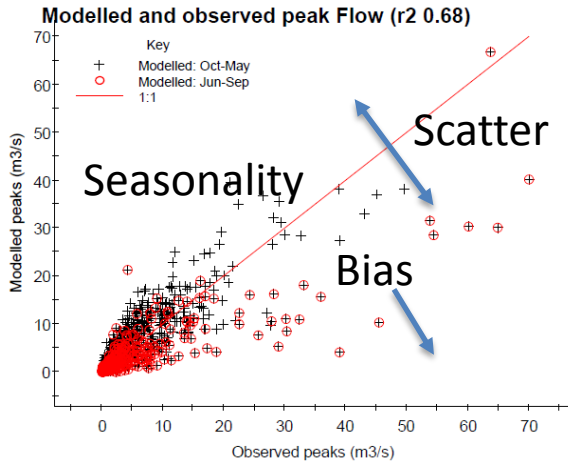
6. 10 Aug 2004 18:15



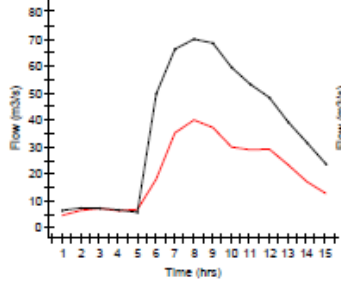
Variability	bias - absolute	bias - significance	shape	Timing bias - abs.	Timing bias - sign.	Seasonality	Overall Score
1.19	2.62	0.71	1.56	0.58	0.86	2.53	1.36



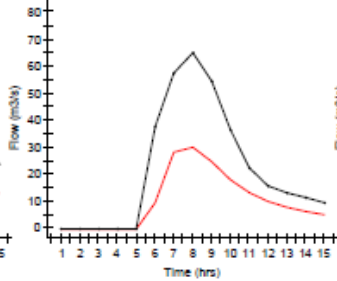
Simulation measures



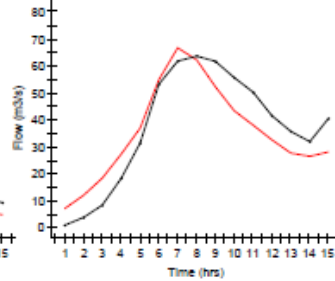
1. 06/09/2008 15:00:00



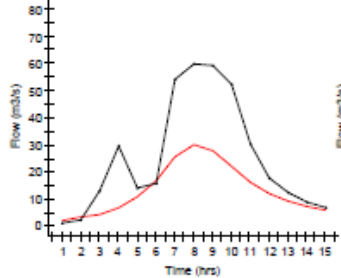
2. 28/06/2012 12:00:00



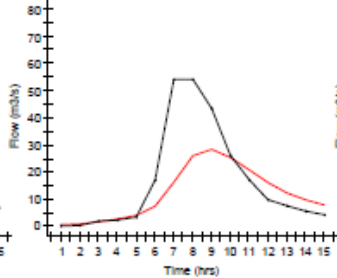
3. 20/07/2007 16:00:00



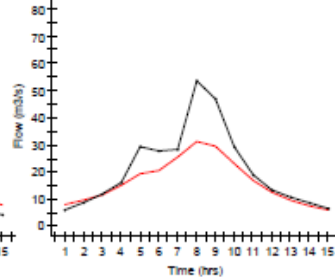
4. 15/06/2007 05:00:00



5. 08/08/1999 21:00:00



6. 24/09/2012 11:00:00



Variability	bias -absolute	bias -significance	shape	Timing bias - abs.	Timing bias - sign.	Seasonality	Overall Score
2.79	3.49	3.54	1.38	2.61	4.00	4.00	2.97

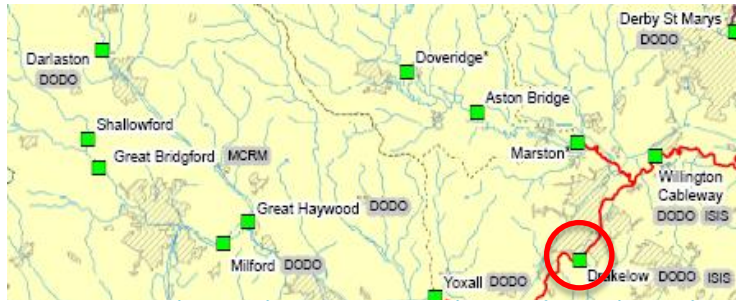


Visualisation

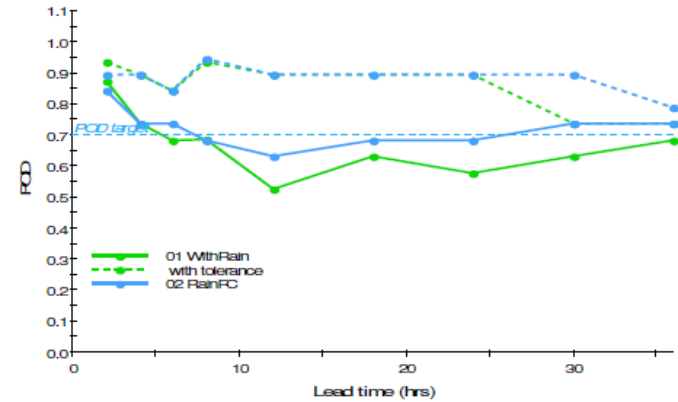
Location ID: 4019
 Network: Middle Trent
 Model Type: ISIS

Real time forecast performance gra

Scenario	Thresh (m)	N
01 WithRain	2.66	19
	2.88	
	3.25	
	3.45	
	3.8	
02 RainFC	2.66	19
	2.7	
	2.88	
	3.25	
	3.45	



Real time performance measures for 40192.66H.rated.forecast



Model scores

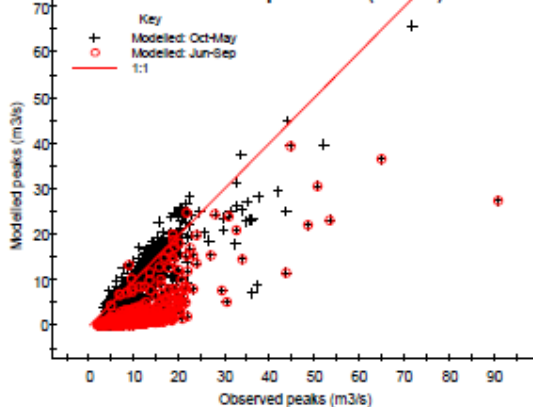
	Variability	bias - absolute	bias - significance	shape	Timing bias - abs.	Timing bias - sign.	Seasonality	Overall Score
Modelled	3.26	3.37	3.17	1.46	2.15	3.17	4.00	2.93

Table; threshold tested = 2.66m

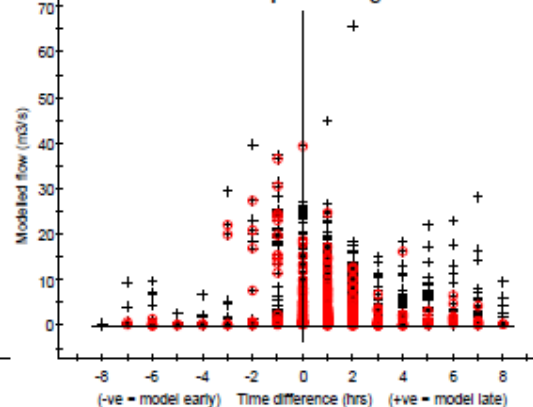
Scenario	N	POD at increasing lead time (hrs)								
		t-2	t-4	t-6	t-8	t-12	t-18	t-24	t-30	t-36
WithRain	19	0.88	0.74	0.68	0.69	0.53	0.63	0.58	0.63	0.68
RainFC	19	0.84	0.74	0.74	0.68	0.63	0.68	0.68	0.74	0.74

the number of observed crossings used to calculate POD
 olute POD values are tabulated. POD values with a tolerance of 0.2 are also plotted as dashed lines, i.e. if a
 ulation gets within a specified distance of the threshold, then it counts as a hit

Modelled and observed peak Flow (r2 0.56)



Modelled and observed peak timing Flow



1.1	0.6	0.2	1.5	2.4	3.8	1.7	1.5
1.3	2.5	0.5	0.6	2.3	2.2	2.8	1.5
2.1	4.0	2.1	1.0	1.5	1.1	4.0	2.1
1.5	0.2	0.1	1.2	0.1	0.1	2.2	1.0
1.6	3.3	2.1	2.2	2.2	2.7	3.6	2.3
1.5	2.4	0.8	1.8	3.4	3.8	4.0	2.3
1.0	4.0	2.3	1.0	0.7	1.0	2.0	1.4
3.3	3.4	3.2	1.5	2.2	3.2	4.0	2.9
3.0	3.0	2.5	1.6	1.9	3.8	1.1	2.4
2.2	0.7	0.4	1.2	1.8	3.9	4.0	2.1
1.9	2.9	2.4	2.7	1.4	3.5	2.9	2.5
1.0	2.7	0.6	0.8	1.8	1.9	2.2	1.3
2.6	2.3	1.8	2.4	1.1	4.0	2.3	2.4
1.6	1.1	0.5	1.8	0.8	0.7	2.3	1.4
1.3	1.4	0.6	2.1	2.4	4.0	1.6	1.8
2.7	0.4	0.7	1.5	0.5	2.3	4.0	2.1
2.5	1.7	2.1	1.9	0.2	0.5	3.5	2.0
2.8	3.5	3.5	1.4	2.6	4.0	4.0	3.0
1.1	4.0	2.4	1.4	1.5	2.9	2.7	2.0

[Explanation of performance grades](#)

[Explanation of model scores](#)



Visualisation

- Forecasters



Real time forecast performance grades

Scenario	Thresh.	N	Lead time (hrs) [grey cell = indicative catchment response time]									
Jan 04 - Jul 14	(m)		2	4	6	8	12	18	24	30	36	
01 WithRain	2.66	19	B	B	A	A	A+	A+	B	B	B	
	2.7	16	A+	A	A+	A+	A+	B	B	B	B	
	2.88	10	A+	A+	A	A+	A+	A+	B	B	B	
	3.25	5	B	B	B	B	B	B	B	B	E	
	3.45	2	B	B	B	B	B	B	B	B	C	
	3.8	0	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	
	Peaks		A	A	A	A	A	B	B	B	A	
Jan 04 - Jul 14	(m)											
02 RainFC	2.66	19	B	B	A	A	A+	A+	A	B	B	
	2.7	16	A+	A	A+	A+	A+	A+	B	B	B	
	2.88	10	A+	A+	A	A+	A+	A+	A	B	B	
	3.25	5	B	B	B	B	B	B	B	E	E	
	3.45	2	B	B	B	B	B	B	B	B	C	
	3.8	0	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	
	Peaks		A	A	A	A	A	B	B	B	A	



Visualisation

- Managers



See footnote for an explanation of table column headings

Forecast location performance commentary



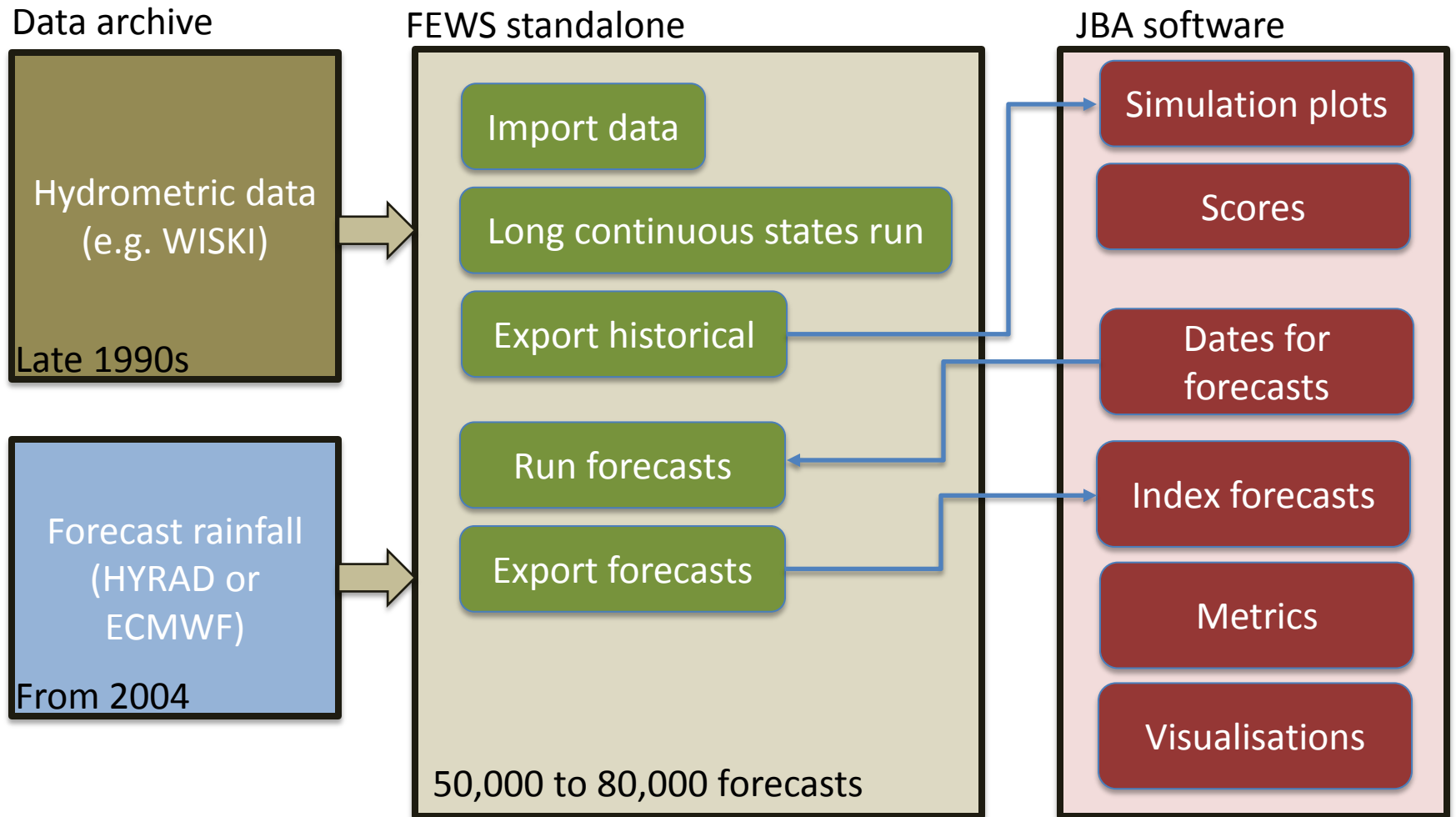
ID	Site	River	Type	Response time	Obs Rain		Forecast Rain		Intrinsic c	Variability	Bias abs	Shape	Timing bias abs	Seasonality	Overall Score	Overview	Recommendations
					B	B	D	2.2									
2516	Halesowen	Stour	MCRM	2	B	B	D	2.2	3.1	3.4	0.7	3.6	2.8		MCRM simulating the headwaters of the Upper Stour. Catchment urbanised	Consider modelling the urban response separately to improve the model's performance in smaller events particularly.	
2083	Stourbridge	Stour	DODO	2.31	B	B	E	3.0	2.3	1.5	1.2	2.9	2.5		DODO routed flows from Halesowen, with significant urban lateral inputs from an MCRM	For lateral inflows, consider modelling the urban response separately to improve the model's performance in smaller events particularly. Calibrating the model's wavespeed may improve performance above 20m ³ /s.	
2641	Wightwick	Stour	MCRM	2	B	B	E	2.4	1.0	3.4	1.2	0.4	2.2		Small MCRM simulating the Upper Smeston. Fast response	Consider recalibration to improve the hydrograph shape (this will also require a reduction in volume of runoff)	
2706	Wombourne	Stour	MCRM	2	A	D	E	3.0	2.0	3.7	2.0	4.0	3.2		Small MCRM simulating the Wom Brook. Fast response, with an obviously early urban runoff peak	Consider modelling the urban response separately to improve the model's performance in smaller events particularly.	
2067	Swindon	Stour	DODO	4	A+	B	C	1.8	1.4	1.4	2.2	0.2	1.9		DODO routes flows from Wombourne and Wightwick with additional lateral inflow from an MCRM. Observed hydrographs quite 'flat topped', indicating storage or bypassing	Understand the cause of the 'flat topping' in the observed series and decide whether it needs to be incorporated into the DODO model	
2084	Stourton	Stour	DODO	8.5	A+	B	B	1.2	3.2	1.1	1.0	3.6	1.8		DODO routing flows from Swindon and Stourbridge with some lateral inflow from an MCRM	Consider correcting the peak bias and timing bias by adjusting the DODO's inputs and wavespeed (although low priority)	
															DODO routing flows from Stourton with additional lateral flow from MCRM	Investigate the source of the bias and then correct.	



How can I get my hands on this stuff?

Quantifying and visualising performance

What's the process?



What specialist skills do I need?

- FEWS configuration
- Software for reporting/processing
 - but could use FEWS performance module



The future

Quantifying and visualising performance

In Future?

- Testing now the norm
- National results database
- Changing measures
- Make it easier to re-run forecasts (archive)
- Other uses for results (e.g. QR)



Thanks for listening!

Quantifying and visualising performance