

Table 1: Predictive accuracy table, new dataset

Duration	Name	Evaluation metric				
		RMSE [l/s/km2]	CRPS [l/s/km2]	MAE [l/s/km2]	MRE [%]	MAPE [%]
0 hour	floodGAM	<b>149.7</b>	<b>69.7*</b>	<b>98.8*</b>	<b>21.3*</b>	<b>21.3*</b>
	RFFA_2018	158.1	80.3	110.4	29.5	26.5
	XGBoost	-	-	98.8	-	-
1 hour	floodGAM	<b>150.7</b>	<b>69.5*</b>	<b>98.6*</b>	<b>21.3*</b>	<b>21.2*</b>
	RFFA_2018	156.9	79.2	109.3	29.0	26.2
	XGBoost	-	-	95.0	-	-
6 hour	floodGAM	<b>132.7</b>	<b>63.1*</b>	<b>90.3</b>	<b>20.2*</b>	<b>20.4*</b>
	RFFA_2018	136.4	69.4	97.3	25.4	24.0
	XGBoost	-	-	88.1	-	-
12 hour	floodGAM	<b>109.6</b>	<b>54.6*</b>	<b>77.5</b>	<b>18.9*</b>	<b>19.5*</b>
	RFFA_2018	113.2	58.5	82.6	21.9	21.2
	XGBoost	-	-	76.9	-	-
18 hour	floodGAM	<b>100.0</b>	<b>49.8</b>	<b>70.4</b>	<b>18.3</b>	<b>18.7</b>
	RFFA_2018	100.4	51.7	72.9	20.0	19.5
	XGBoost	-	-	69.2	-	-
24 hours	floodGAM	89.8	<b>45.8</b>	<b>65.0</b>	<b>17.7</b>	18.2
	RFFA_2018	<b>89.4</b>	46.5	65.3	18.5	18.2
	XGBoost	-	-	67.7	-	-
36 hour	floodGAM	74.5	39.4	55.4	<b>16.7</b>	17.3
	RFFA_2018	<b>72.8</b>	<b>38.7</b>	<b>54.2</b>	16.8	<b>16.3</b>
	XGBoost	-	-	52.5	-	-
48 hour	floodGAM	65.0	34.9	48.7	16.0	16.4
	RFFA_2018	<b>63.5</b>	<b>34.2</b>	<b>48.4</b>	<b>15.9</b>	<b>15.6</b>
	XGBoost	-	-	47.2	-	-
72 hour	floodGAM	54.4	29.9	41.6	15.5	15.6
	RFFA_2018	<b>52.7</b>	<b>29.0</b>	<b>41.0</b>	<b>14.9</b>	<b>15.1</b>
	XGBoost	-	-	39.9	-	-
7 days	floodGAM	42.9	23.4	32.4	15.1	<b>15.3</b>
	RFFA_2018	<b>41.4</b>	<b>23.2</b>	<b>31.9</b>	<b>14.6</b>	15.4
	XGBoost	-	-	33.3	-	-
14 days	floodGAM	38.4	<b>19.5</b>	<b>26.7</b>	<b>14.8</b>	<b>14.9*</b>
	RFFA_2018	<b>37.1</b>	20.0	26.9	15.3	16.0
	XGBoost	-	-	25.8	-	-
30 days	floodGAM	33.7	<b>16.0</b>	<b>21.9</b>	<b>14.8*</b>	<b>14.6*</b>
	RFFA_2018	<b>32.7</b>	17.0	23.1	16.3	17.2
	XGBoost	-	-	20.7	-	-

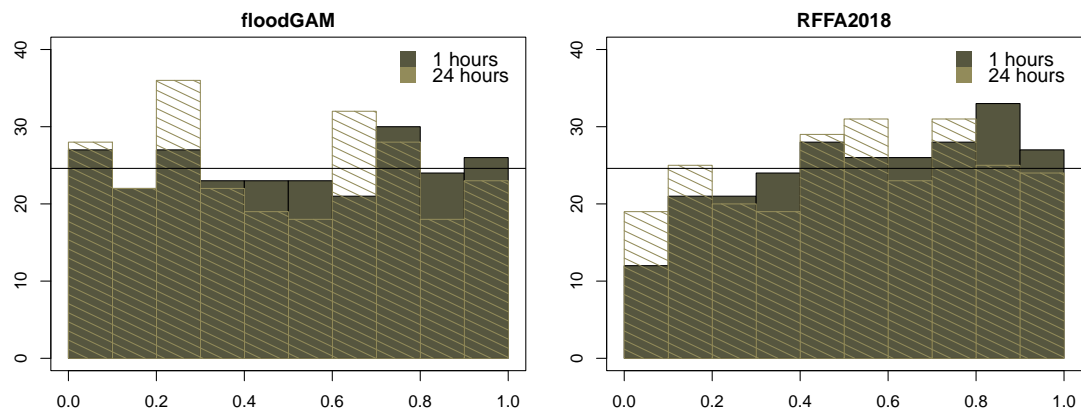


Figure 1: PIT histograms, with new dataset

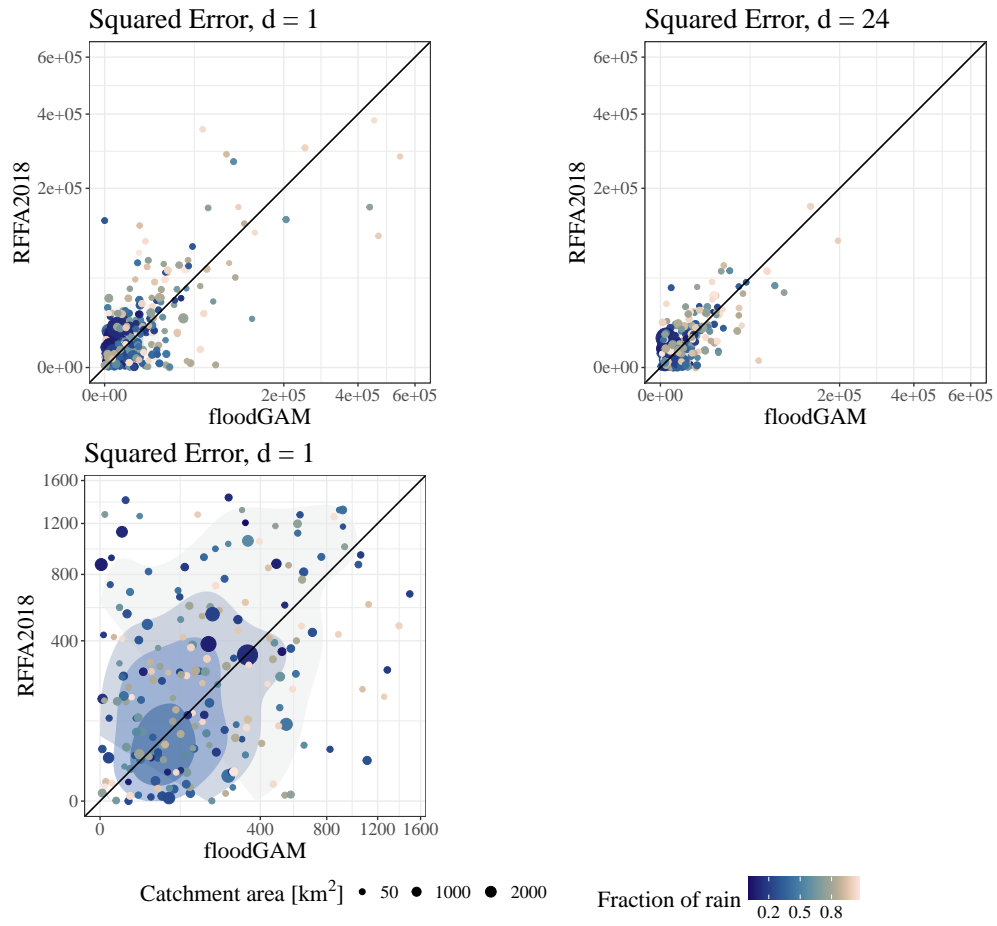


Figure 2: squared error dotplots, new dataset

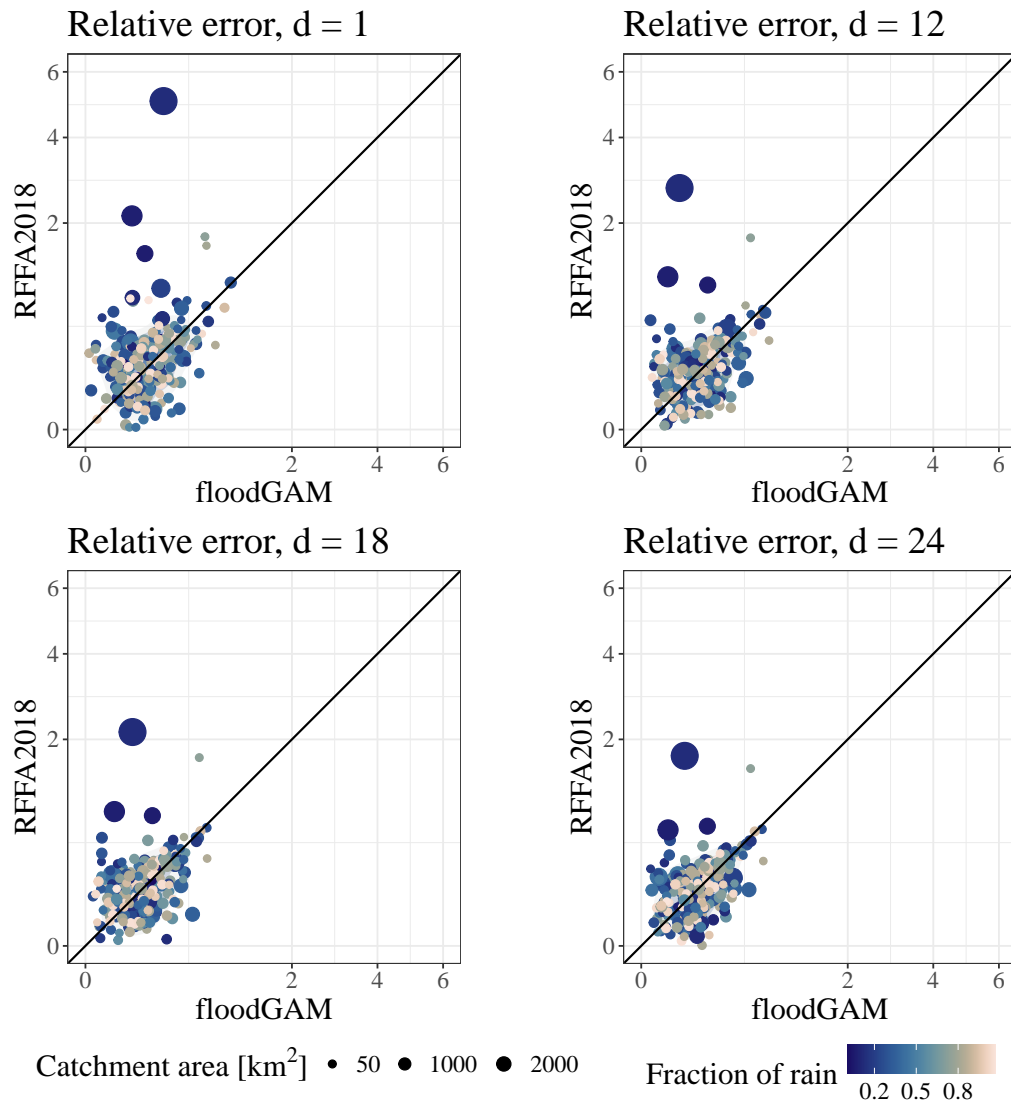


Figure 3: We still see underestimation from RFFA2018 at large, snowmelt driven catchments