

Life Is On



DRIVEN DATA



## Rinse over Run

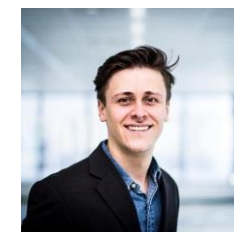
Team Fatima Yamaha



Private Score: 0.2658 (2<sup>nd</sup> Place)

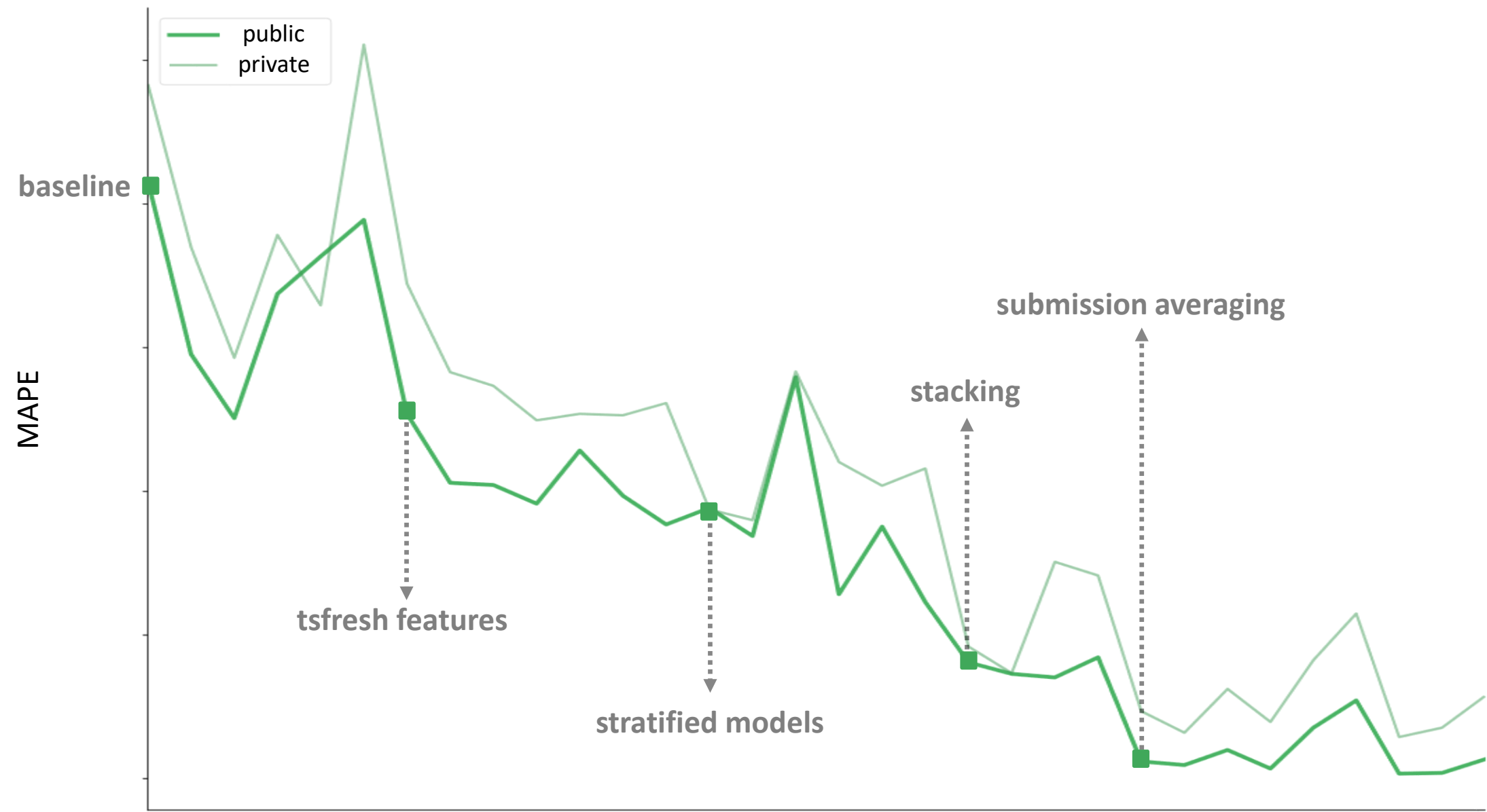


Gilles Vandewiele  
PhD student

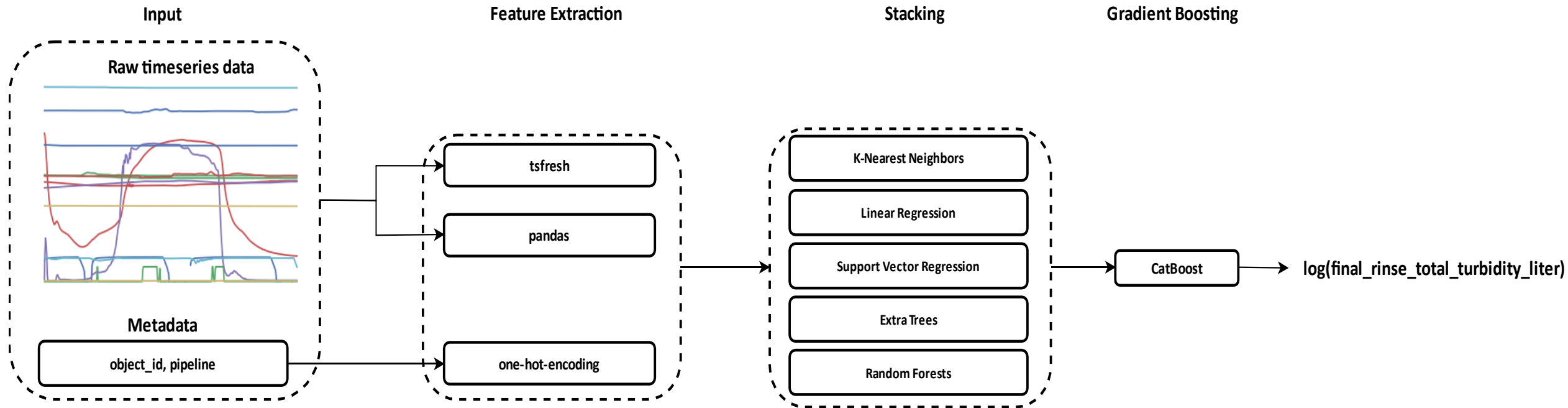


Thomas Mortier  
PhD student

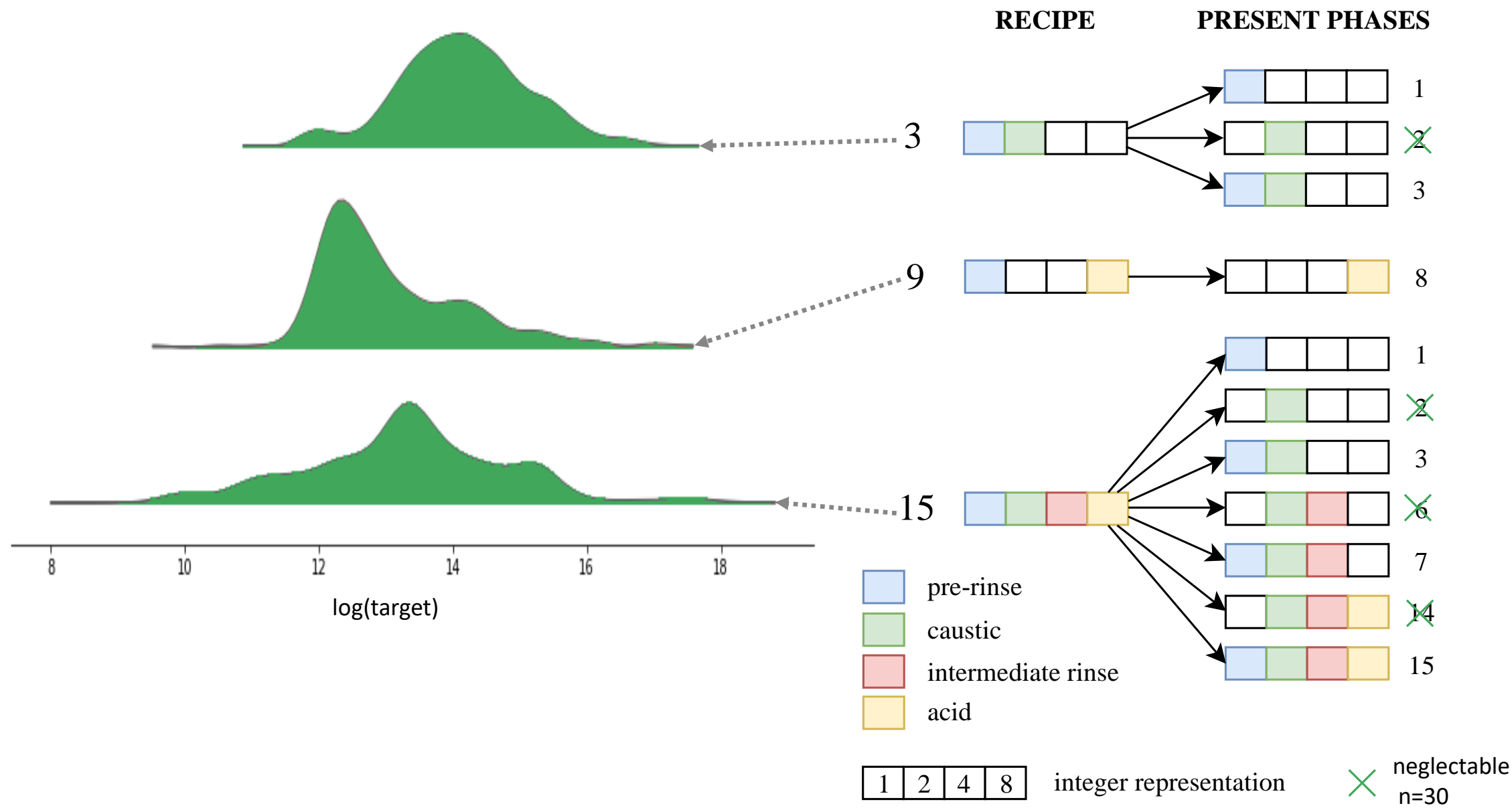
Starting with a simple baseline, using useful features and smart stratification, and ending up with a complex high performant ensemble, secured us the second spot



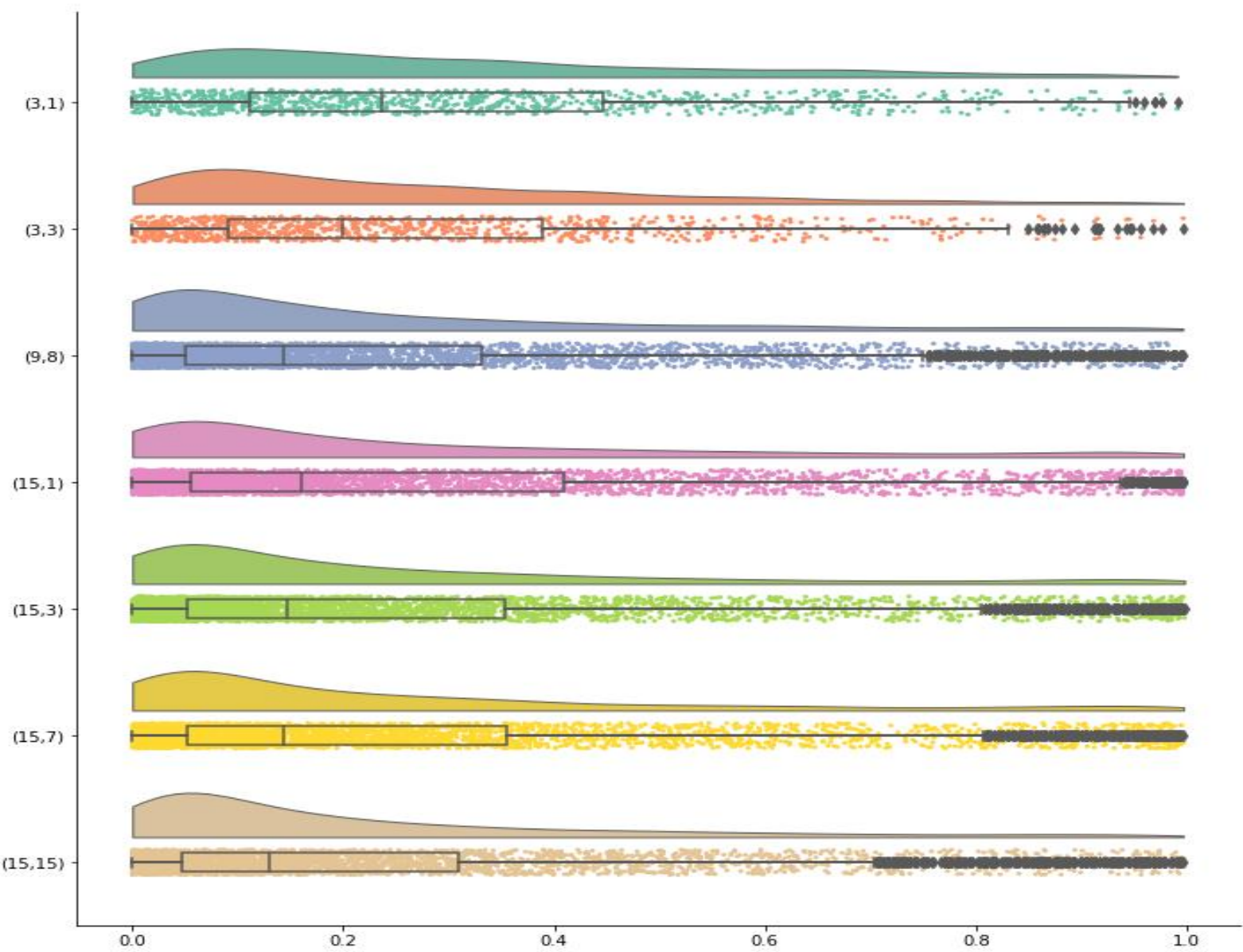
# An overview of our final pipeline



Stratifying models allow to predict after every cleaning phase and explains differences in variance in the outcome of the three cleaning recipes

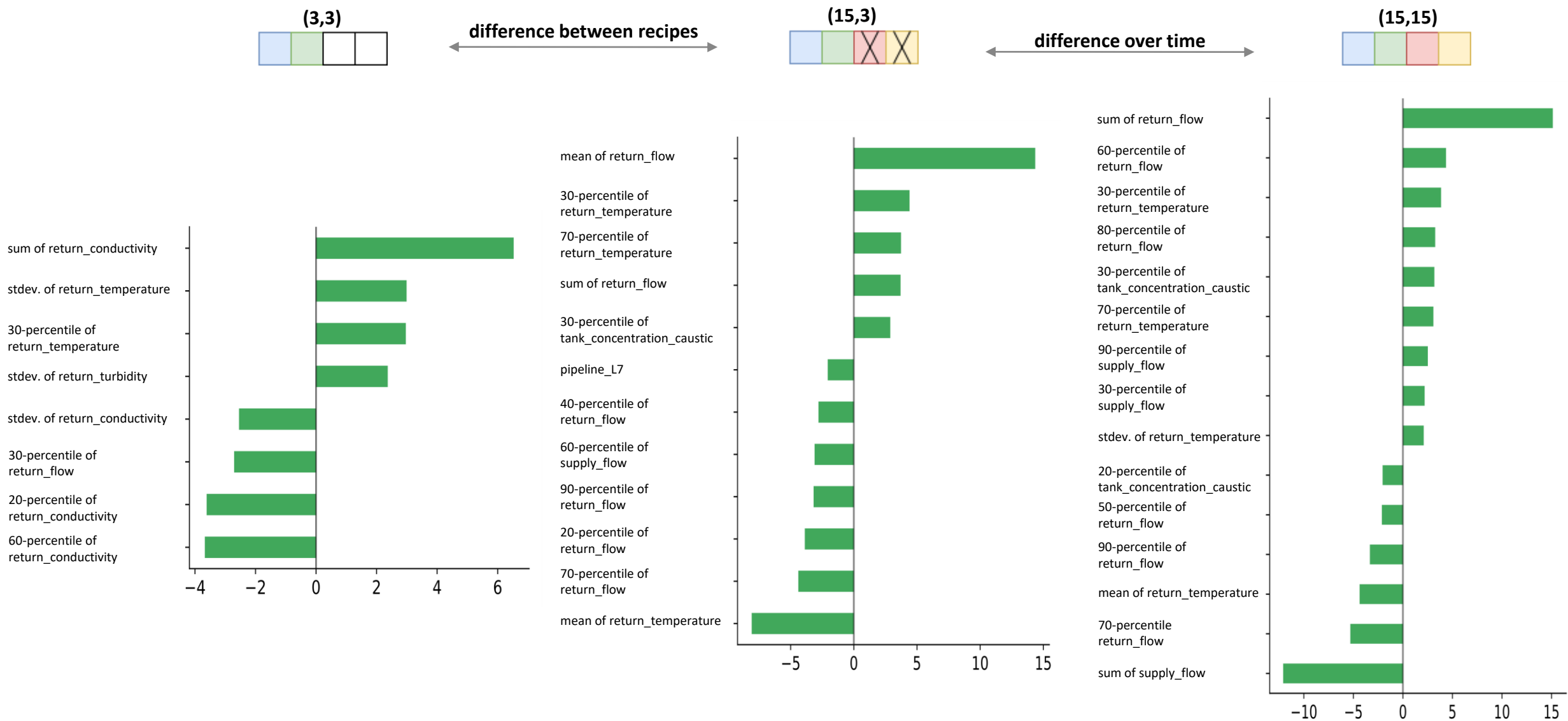


A higher number of processes and using more phases per processes decreases the five-fold cross-validation error, and reveals very high variance, hence, the choice for ensemble techniques  
(\* ) the fraction of test processes      (\*\*) model without stacking



Group	MAPE ( $\mu \pm \sigma$ )	Weight*
(3,1)	$0.35 \pm 0.36$	0.022
(3,3)	$0.29 \pm 0.32$	0.17
(9,8)	$0.38 \pm 0.33$	0.04
(15,1)	$0.30 \pm 0.37$	0.07
(15,3)	$0.28 \pm 0.36$	0.23
(15,7)	$0.27 \pm 0.34$	0.22
(15,15)	$0.25 \pm 0.37$	0.22
Total: **	$0.29 \pm 0.36$	

By using logistic regression, together with a p-value correction by means of the Benjamini-Hochberg correction (FDR,  $\alpha = 0.05$ ), we identified positive and negative significant features which explain high targets



# Shapley values allow to highlight most impactful features for the final rinse outcome

## INTERPRETATION:

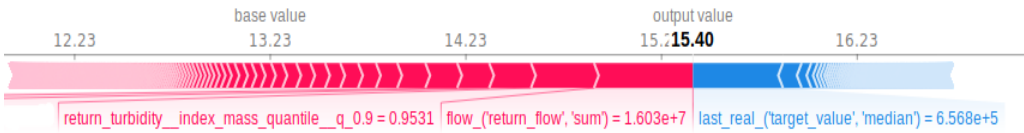
- dots along the x-axis are the different impacts on the model (left = large negative impact; right = large positive impact)
- color of dot impact feature magnitude (blue = lowest value, red = highest value)
- e.g. for (3, 3), a higher number of peaks in supply\_flow increases the prediction



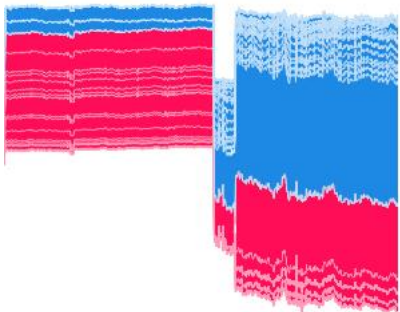


# Shapley values can be generated for each individual prediction and can be grouped together to find clusters of similar predictions and model behavior

individual Shapley plot



(15,3)



~ 300 predictions around 2.25 million liters due to high non-linearity in supply\_flow

(3,3)

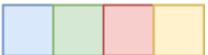


~ 50 predictions around 750,000 liters due to objects 205-209

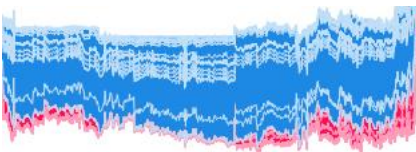
~ 300 predictions around 1.25 million liters due to low non-linearity and high variance in supply\_flow

rotate & cluster

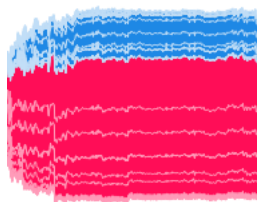
(15,15)



~ 600 predictions around 800,000 liters due to relatively large sum of return\_flow



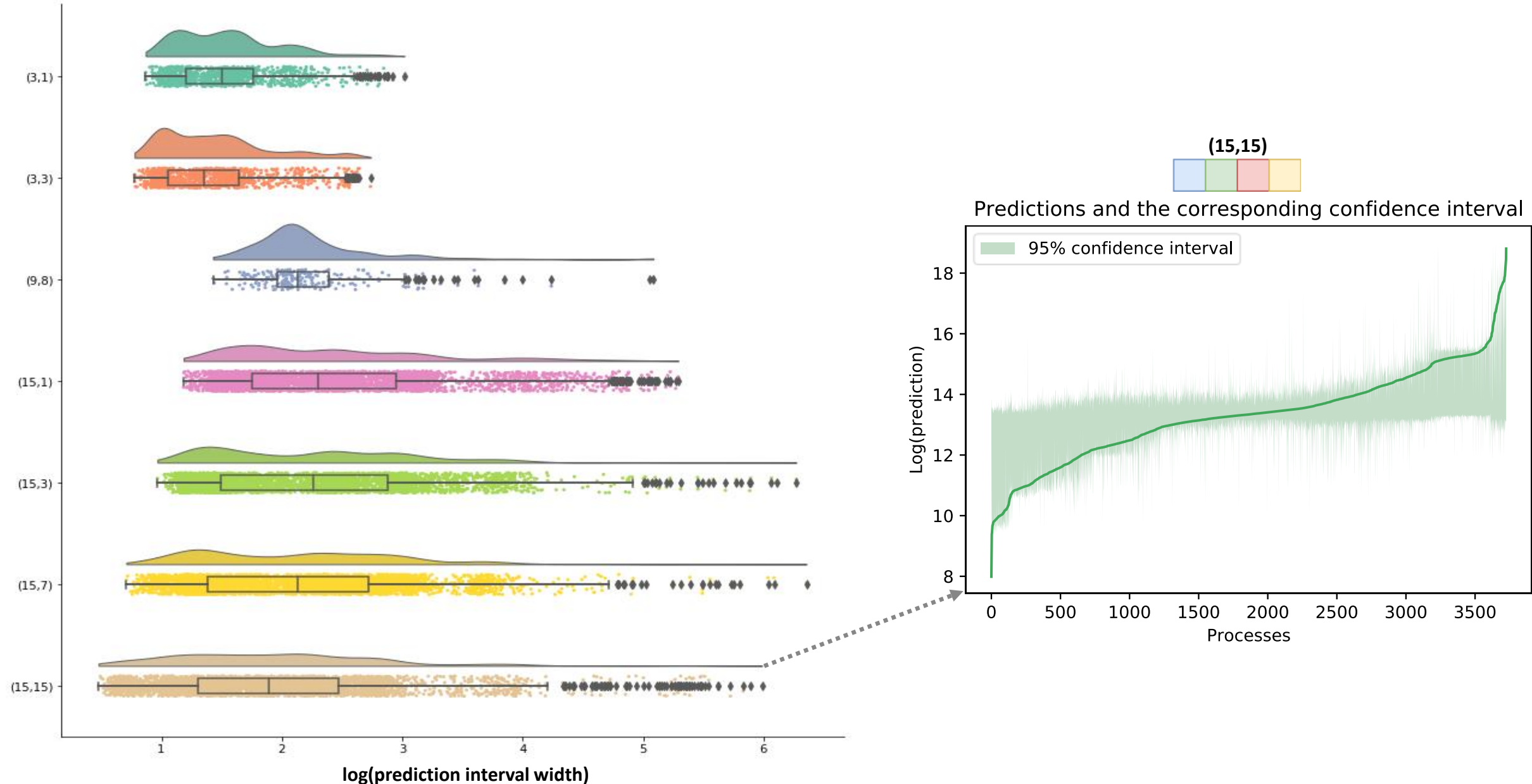
~ 400 predictions around 150,000 liters due to low sum of return\_flow and high variance in supply\_flow



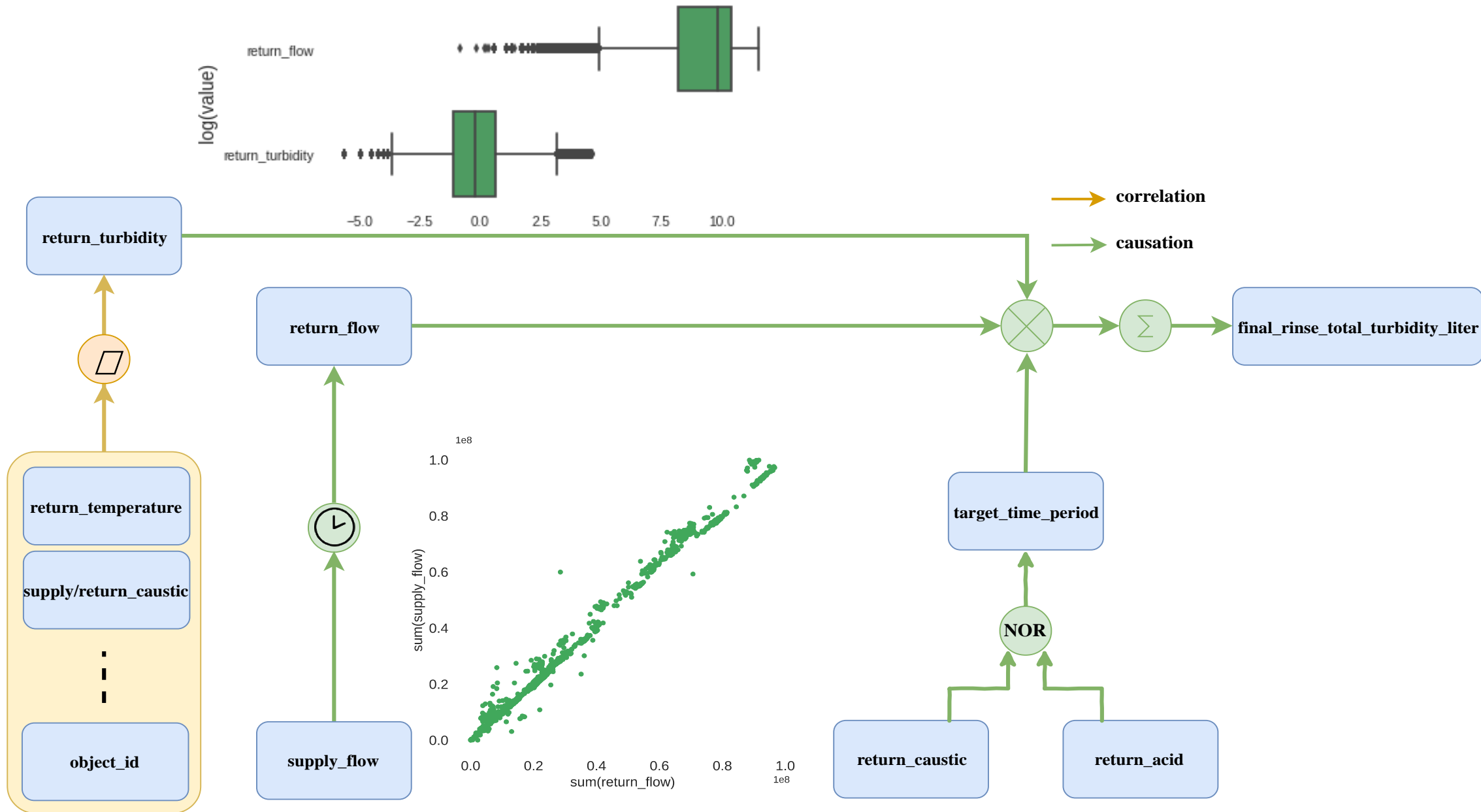
~ 350 predictions around 2 million liters due to much higher return\_flow and supply\_flow sums



Training two gradient boosters with quantile loss ( $1 - \frac{\alpha}{2}$  and  $\frac{\alpha}{2}$ ) allows to construct  $(1 - \alpha)$ -prediction intervals with  $\alpha = 0.05$



An aggregated target causes the less controllable return turbidity to be masked by return flow, which is of a larger order of magnitude while less informative



Predicting the turbidity outcome of the final rinse seems a harder task than predicting the flow outcome, however, interesting interactions appear

return_turbidity	
Group	MAPE ( $\mu \pm \sigma$ )
(3,3)	$0.33 \pm 0.40$
(9,8)	$0.54 \pm 0.49$
(15,15)	$0.33 \pm 0.46$
Total:	<b><math>0.34 \pm 0.45</math></b>

return_flow	
Group	MAPE ( $\mu \pm \sigma$ )
(3,3)	$0.11 \pm 0.15$
(9,8)	$0.16 \pm 0.17$
(15,15)	$0.10 \pm 0.15$
Total:	<b><math>0.10 \pm 0.15</math></b>

Group	MAPE ( $\mu \pm \sigma$ )
(3,3)	$0.10 \pm 0.16$
(9,8)	$0.19 \pm 0.21$
(15,15)	$0.11 \pm 0.17$
Total:	<b><math>0.11 \pm 0.17</math></b>

use only 22  
flow features

