

## Context

\* Emulsifiers stabilise & optimise biological formulations

\* Batch to batch variation affects QA, develop model to quantify emulsifiers in blends.

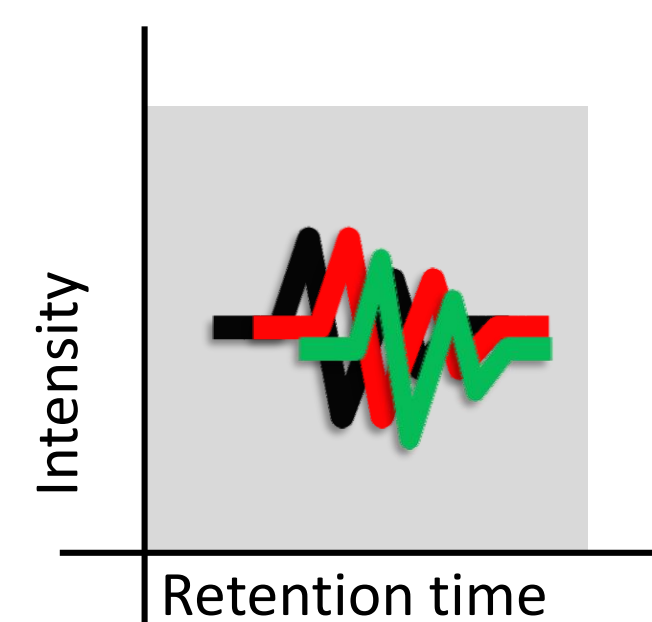
\* Quantitation of emulsifiers by LCMS :- co-elution, overlapping peaks, matrix effects; multivariate methods for resolution needed.

## Data description & pre-processing

Each formulation is a mix of **three emulsifiers/co-formulants (Y1, Y2, Y3)** at the following concentrations

- \* Low (0.06mg/ml)
- \* Mid (0.15mg/ml)
- \* High (0.2399mg/ml)

**Objective** :- Create individual models to predict level of each component in any given blend.



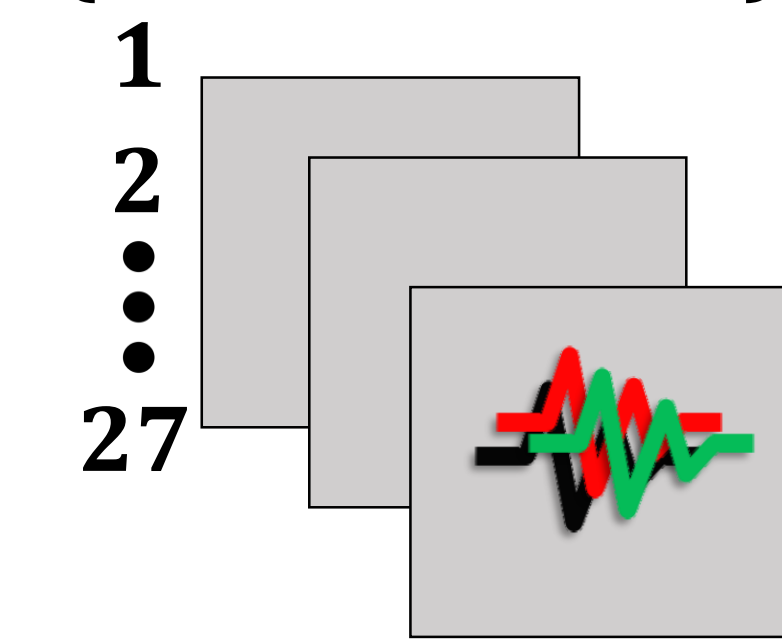
n = 27 unique blends

3 replicates

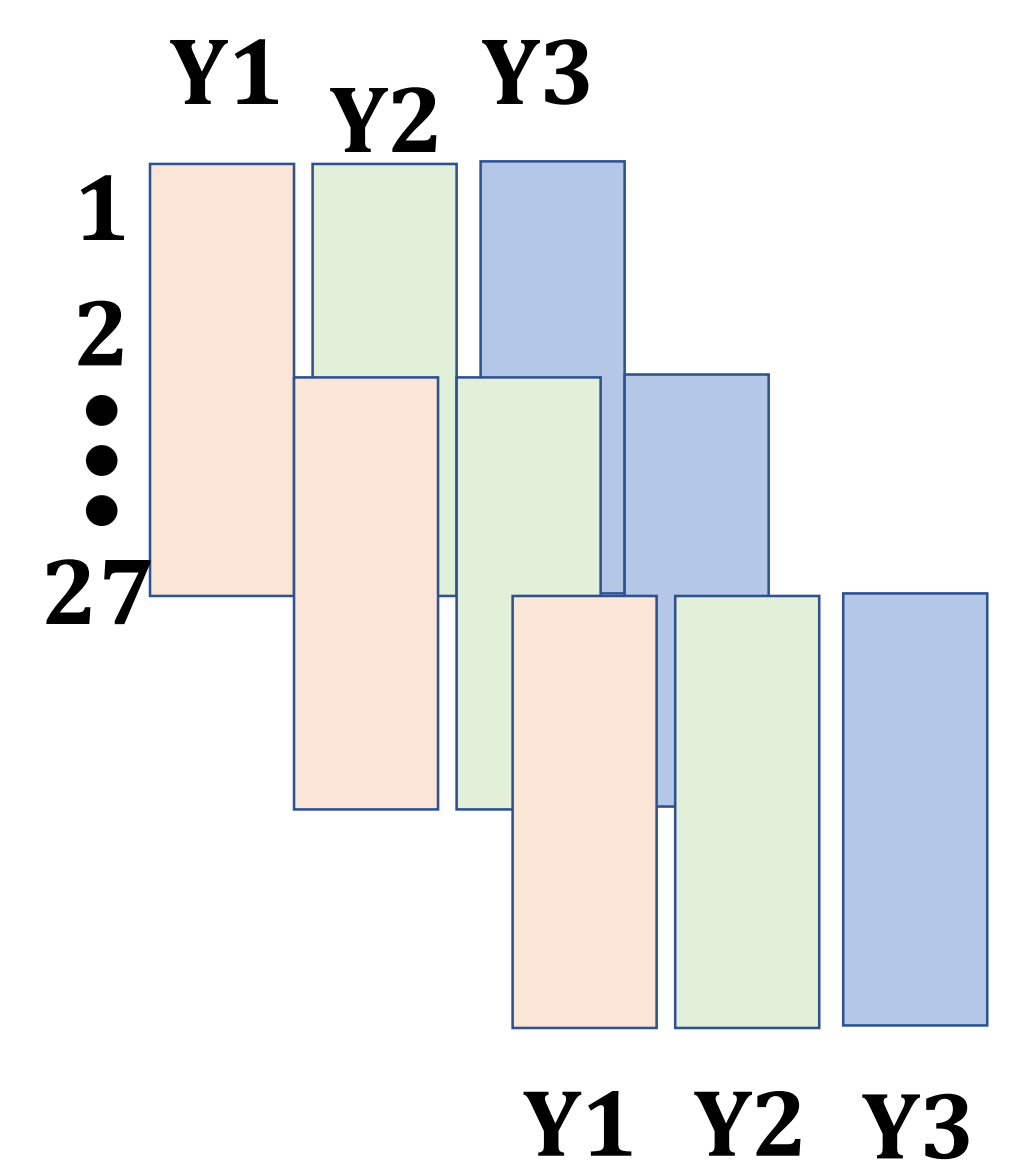
3 concentration values

30-minute LC run (Agilent 1290 Infinity II)

**X - block**  
(n = 27 \* 3 = 81)

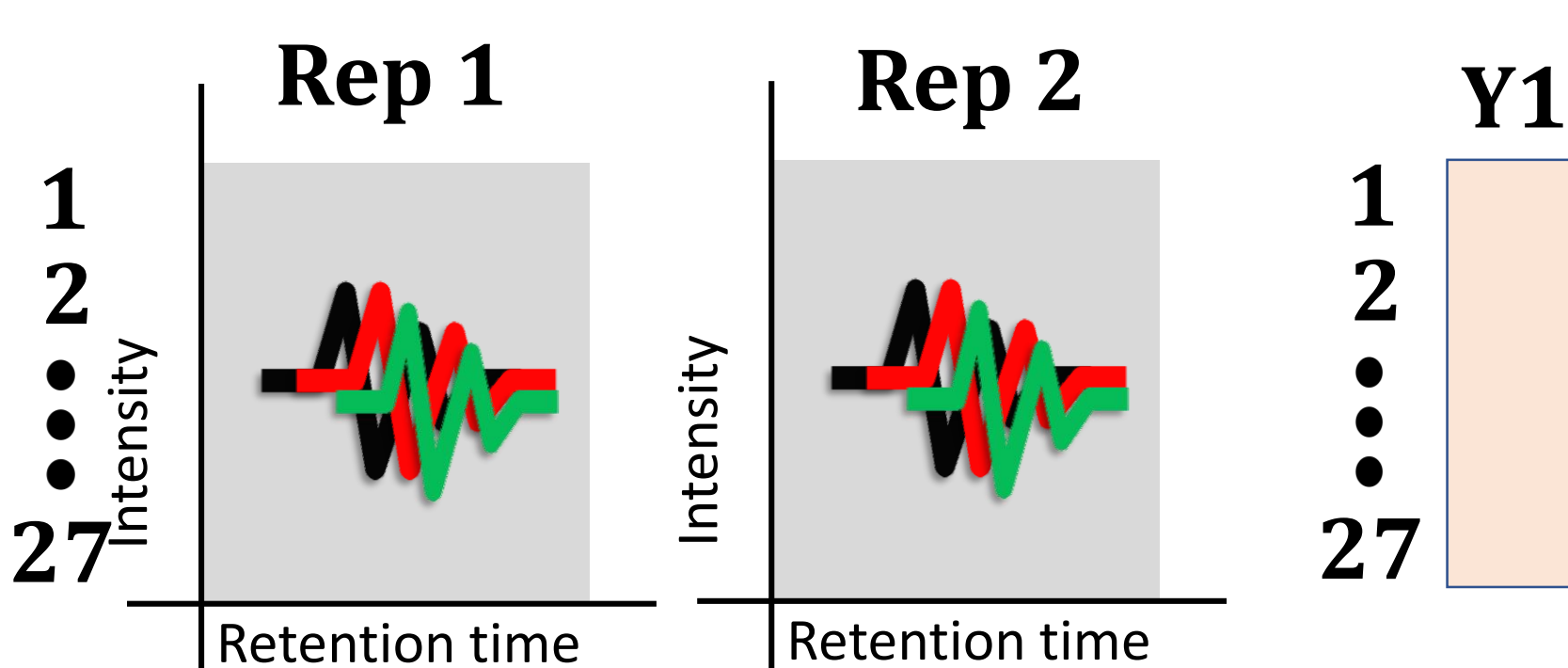


(a) Truncate X - block Rt  
16 - 27min (manual  
vassel)



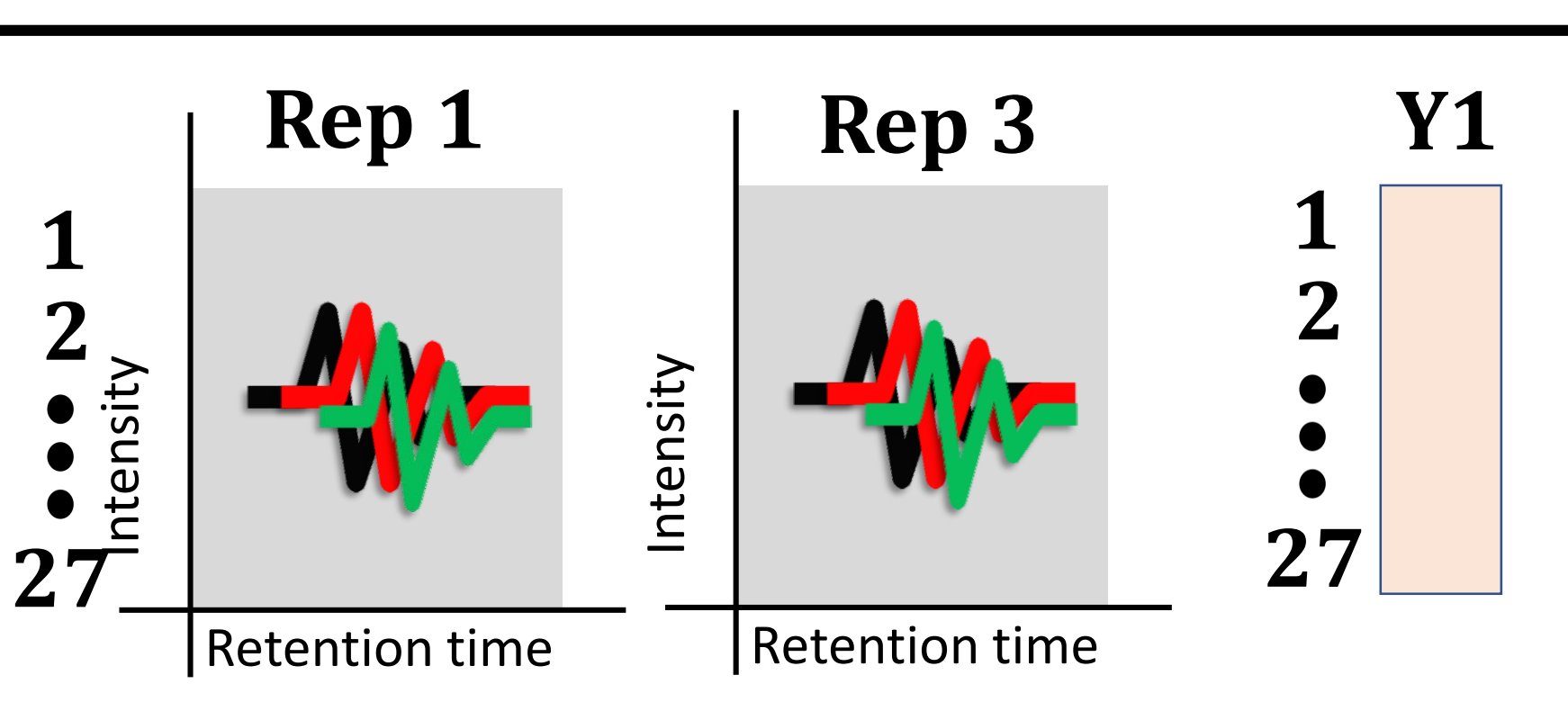
(b) Autoscale X & Y

## Modelling I (Calibration)



\* Create PLSR1 model using n = 27 \* 2 as calibration set for Y1

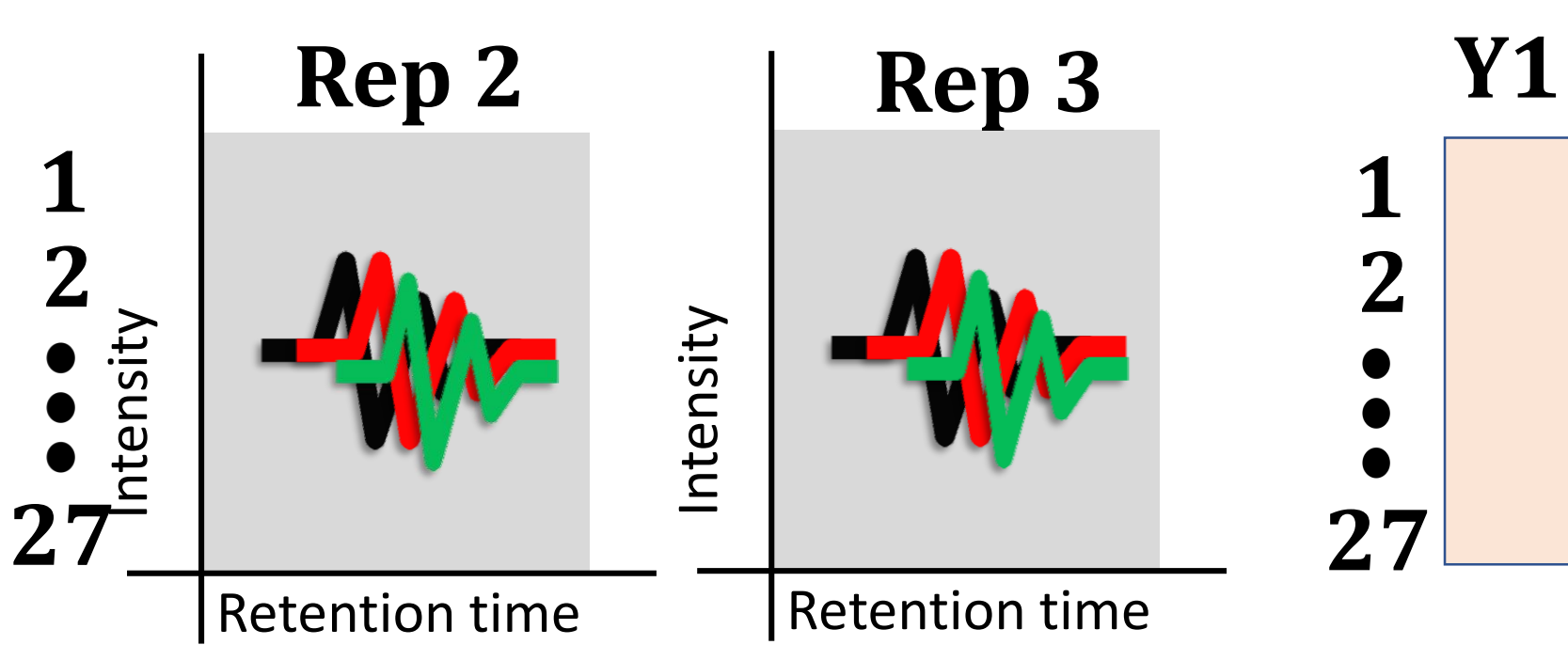
\* 10-fold cross validation



\* Select no. of latent variables based on performance indicators :- **lowest RMSE<sub>CV</sub>**

**highest R<sup>2</sup><sub>CV</sub>**

\* Check if variable selection needed and/or stable



\* Evaluate variability in model performance with respect to cycling of calibration replicates to *select optimal model*

## Modelling II (Validation)

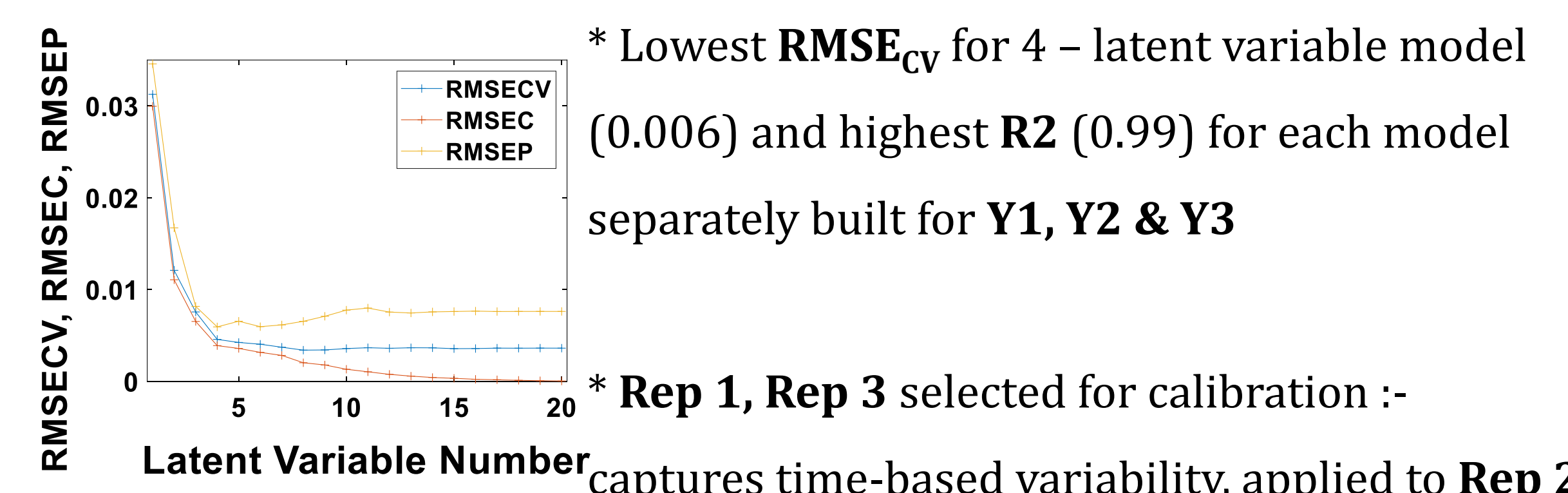
\* Apply model on validation set (n = 27), i.e. replicate 2

\* Evaluate model performance ;select model with **lowest RMSE<sub>p</sub>**, **highest R<sup>2</sup><sub>p</sub>**

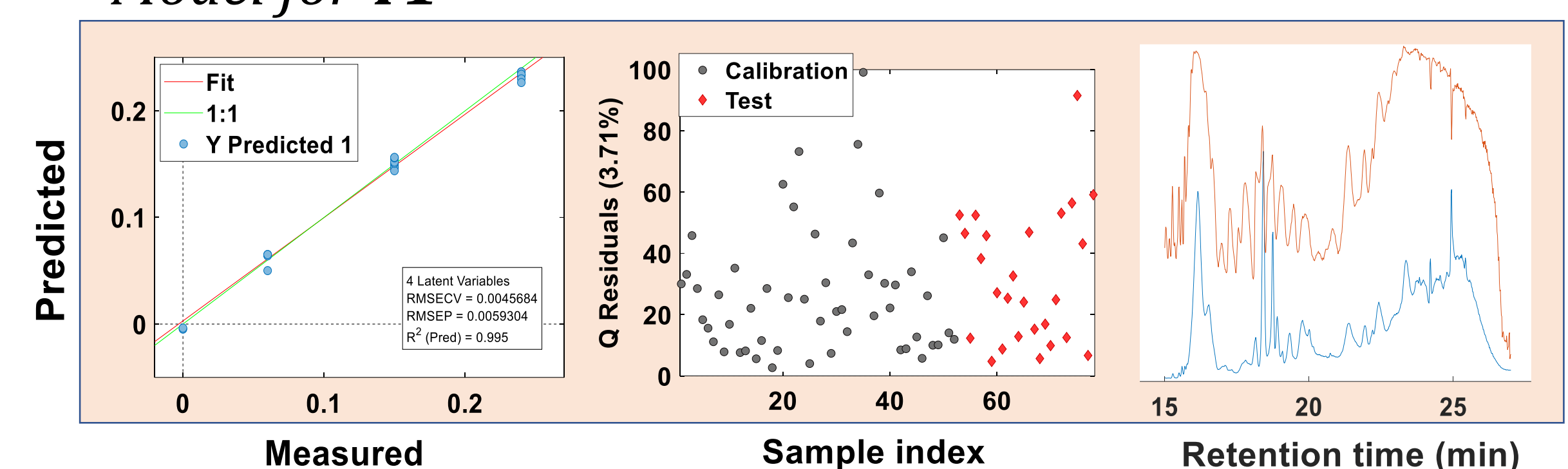
\* Inspect calibration-validation subspace using **Q residuals** to determine if optimal model reached ; consider **exclusion of datapoints** if needed ; consider inclusion of latent variables if needed

\* Repeat modelling with **Y2 and Y3** and compare

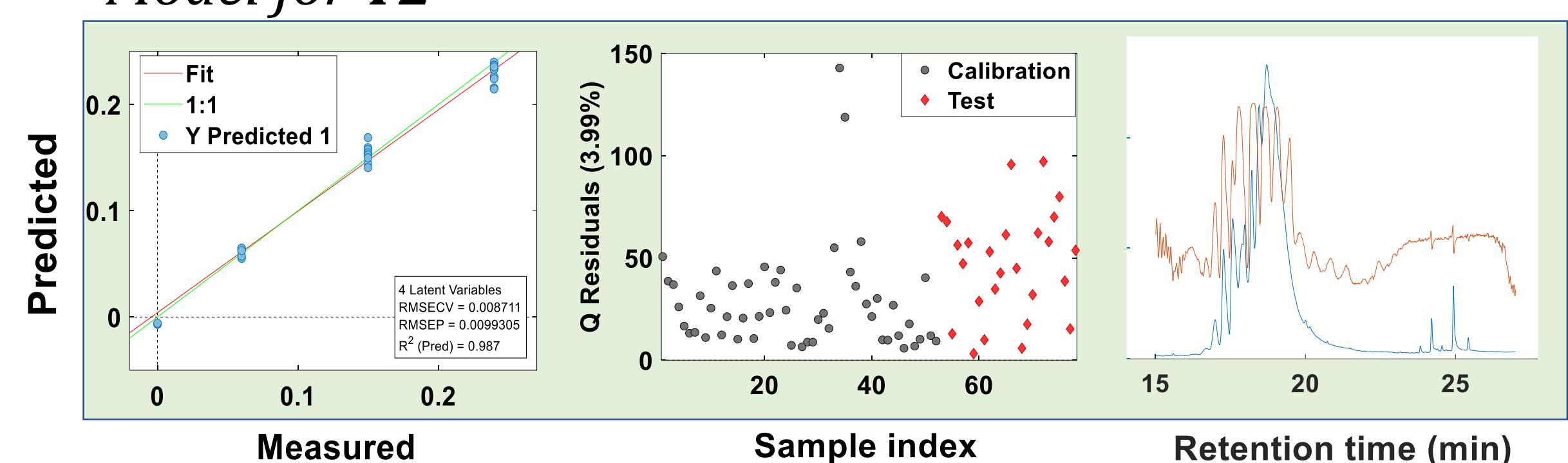
## Results & Discussion



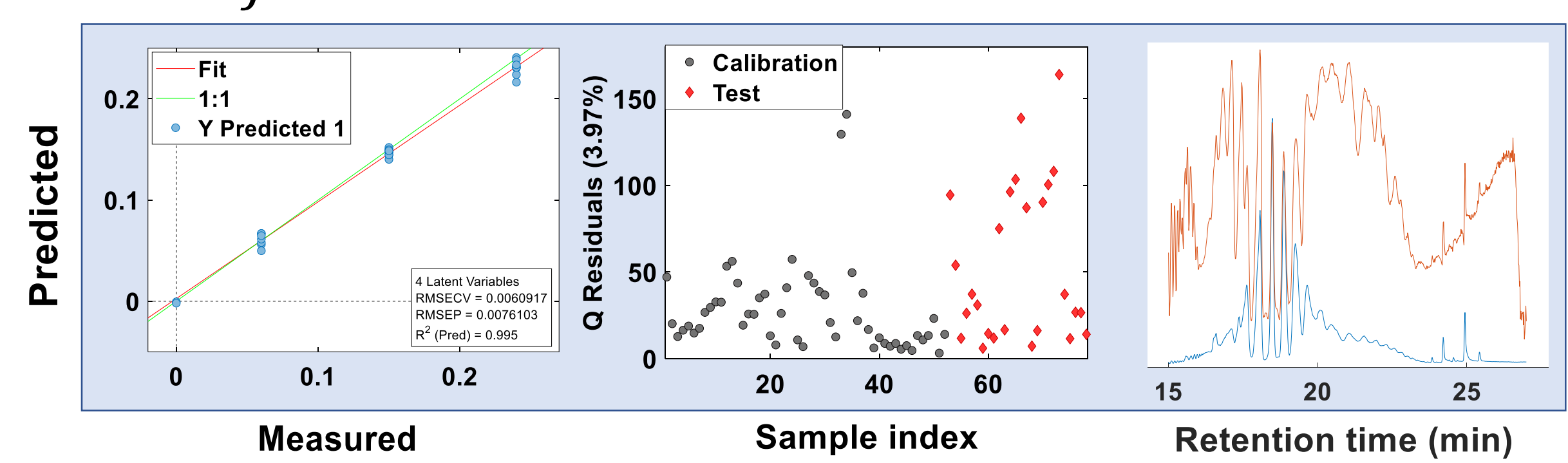
**Model for Y1**



**Model for Y2**



**Model for Y3**



\* **Q residual** subspace indicates inclusion of more latent variables could improve model specificity, at the cost of overfitting > more conservative 4 latent variable model preferred to keep predictive performance consistent

\* Overlay of 'pure' emulsifier chromatogram with 4-latent variable model regression vector indicates key chromatographic features agree more or less

## Acknowledgements

The authors would like to thank Enterprise Ireland for funding support and to SFdS, GC, Univ. Liege for organising Chimiometrie 2020. All analyses performed using MatLab™ 2019a and PLS Toolbox 8.7.1.