

# Data integration for predictive modelling

CMI-PB model submission

Nicky Thrupp

# Model concept

Stochasticity in individual measures

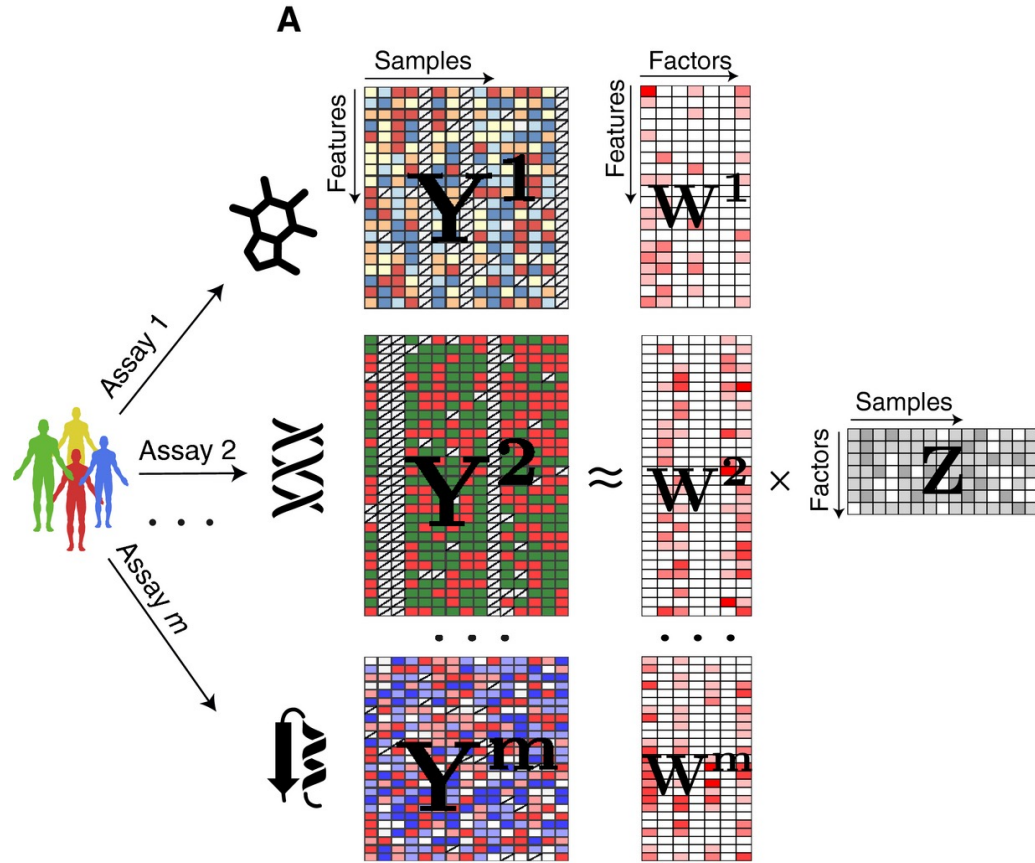
Measures between experiments may suffer from technical bias

How to maximise signals coming from multiple sources

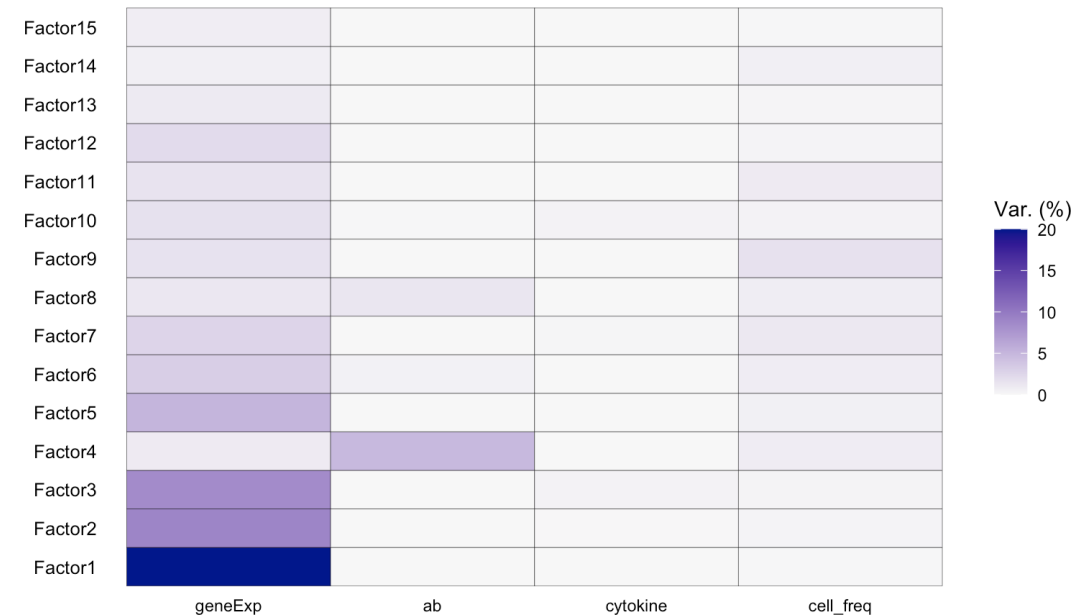
**→ can we address this with data integration**

- Leverage measures that are similar across modalities
- Remove noise
- Biologically interpretable

# MOFA for data integration



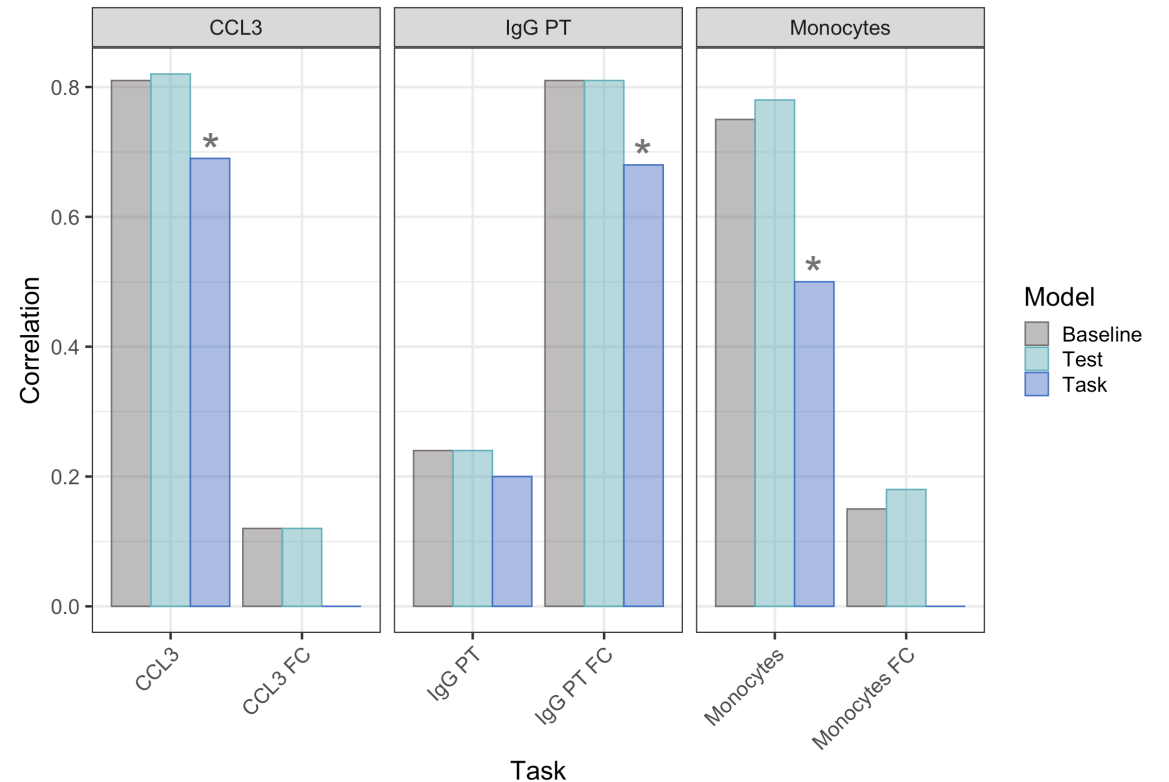
Matrix factorization :  $Y$  = observed measures  
 $Z$  = inferred latent variables  
 $W$  = feature weights



Variance decomposition

# Results

- 2 models:
  - 2020 train, 2021 test
  - 2020+21 train / test
- Lasso regression
- Prediction features:
  - Baseline values
  - Demographic information
  - Top MOFA features



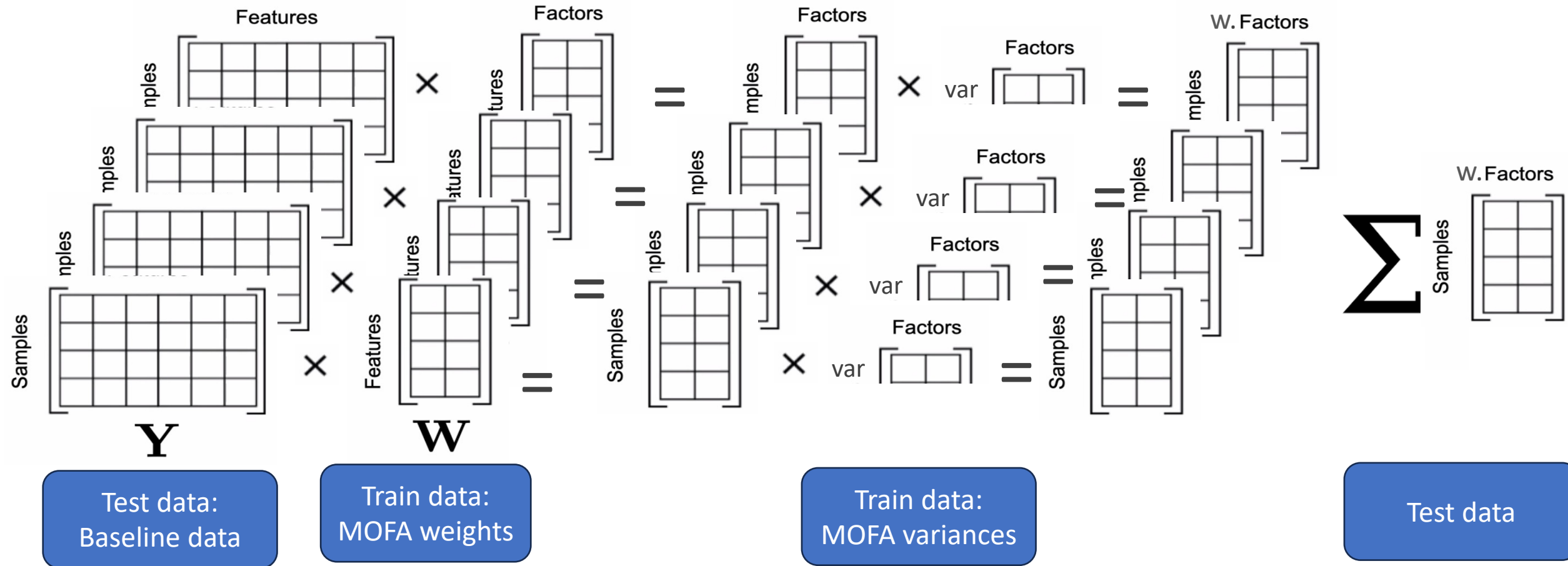
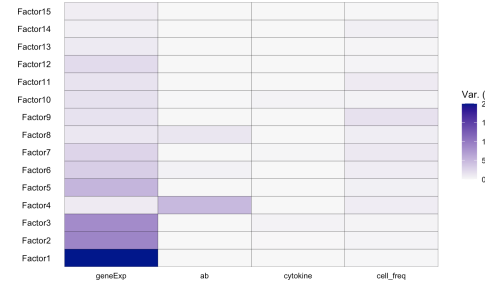
# Conclusions

- Models predicted CCL3 expression and IgG PT fold change well
- MOFA factors improved predictions by small margins
- MOFA model could still be improved:
  - Signals dominated by gene expression modality
  - Too many factors
  - Newer implementation takes groupings (including time courses) into account
- Very similar results when training 2020 and 2020+2021

Thanks !

# Training data:

- clinical data
- baseline assay data (ab, cytokine, cell freq, gene expr)
- MOFA factors:



# LASSO regression in R

```
library(glmnet)
```

```
model<-cv.glmnet(x=as.matrix(predictors.rmNA),  
                 lambda = NULL,  
                 task, family='gaussian',  
                 alpha=alpha,  
                 nfold=nrow(predictors.rmNA) ,  
                 type.measure="mse")
```

Cross-validation using leave-one-out

Hyperparameter tuning of lambda parameter

Additional tuning of alpha parameter (lasso -> elastic -> ridge)