# AI LAB

## 3.1 - SUPERVISED LEARNING
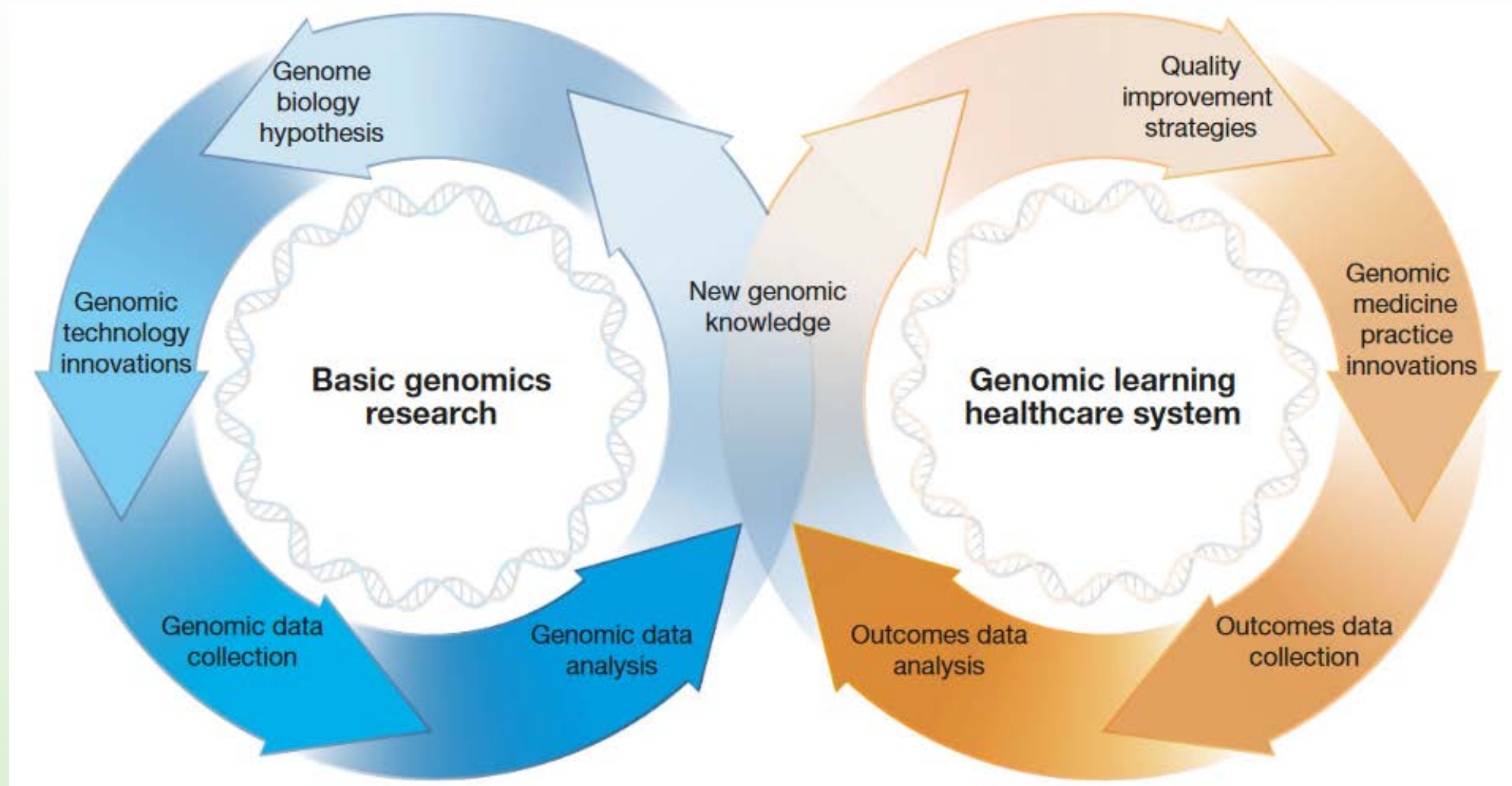
FRANCESCA M. BUFFA

# LAB STRUCTURE

- INTRODUCTION
- THE DATA
- THE AI-LAB CHALLENGE – PART 1
- PART 1 - SHARING AND DISCUSSION

- UNSUPERVISED LEARNING EXAMPLES
- THE AI-LAB CHALLENGE – PART 2
- PART 2 - SHARING AND DISCUSSION

- SUPERVISED LEARNING EXAMPLES
- THE AI-LAB CHALLENGE – PART 3
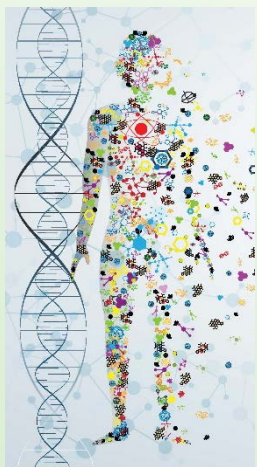- LARGE PROJECTS AND DATABASES

- THE AI-LAB CHALLENGE PARTS 1-3, SHARING AND DISCUSSION
- DATA INTERPRETATION
- DISCUSS AND PREPARE WORKSHOP PRESENTATIONS

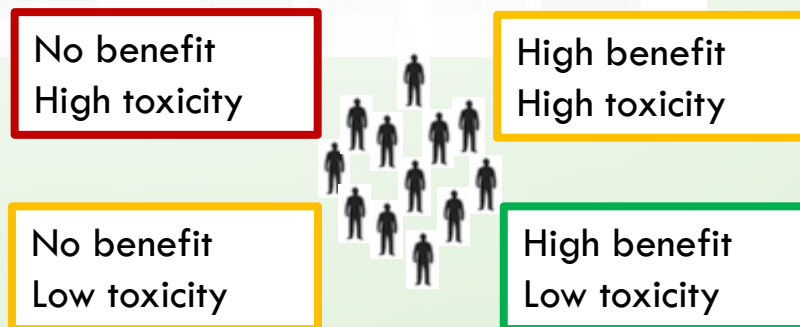# VIRTUOUS CYCLES IN HUMAN GENOMICS RESEARCH AND CLINICAL CARE
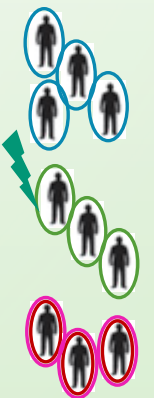
Precision medicine paradigm

# A 70-GENE SIGNATURE FOR RISK OF METASTASIS
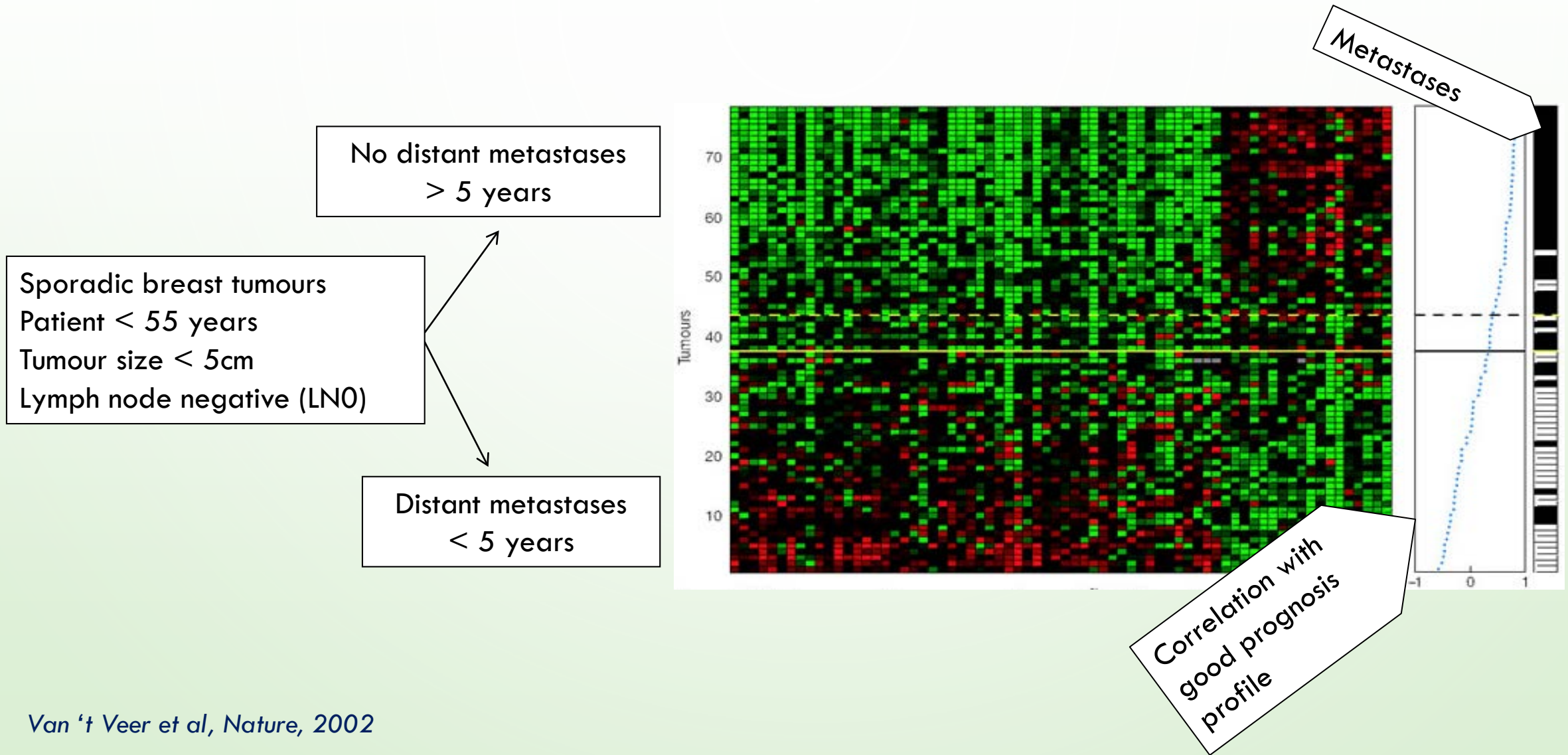


No distant metastases
> 5 years

Sporadic breast tumours
Patient < 55 years
Tumour size < 5cm
Lymph node negative (LN0)

Distant metastases
< 5 years

Metastases

Correlation with good prognosis profile

*Van 't Veer et al, Nature, 2002*
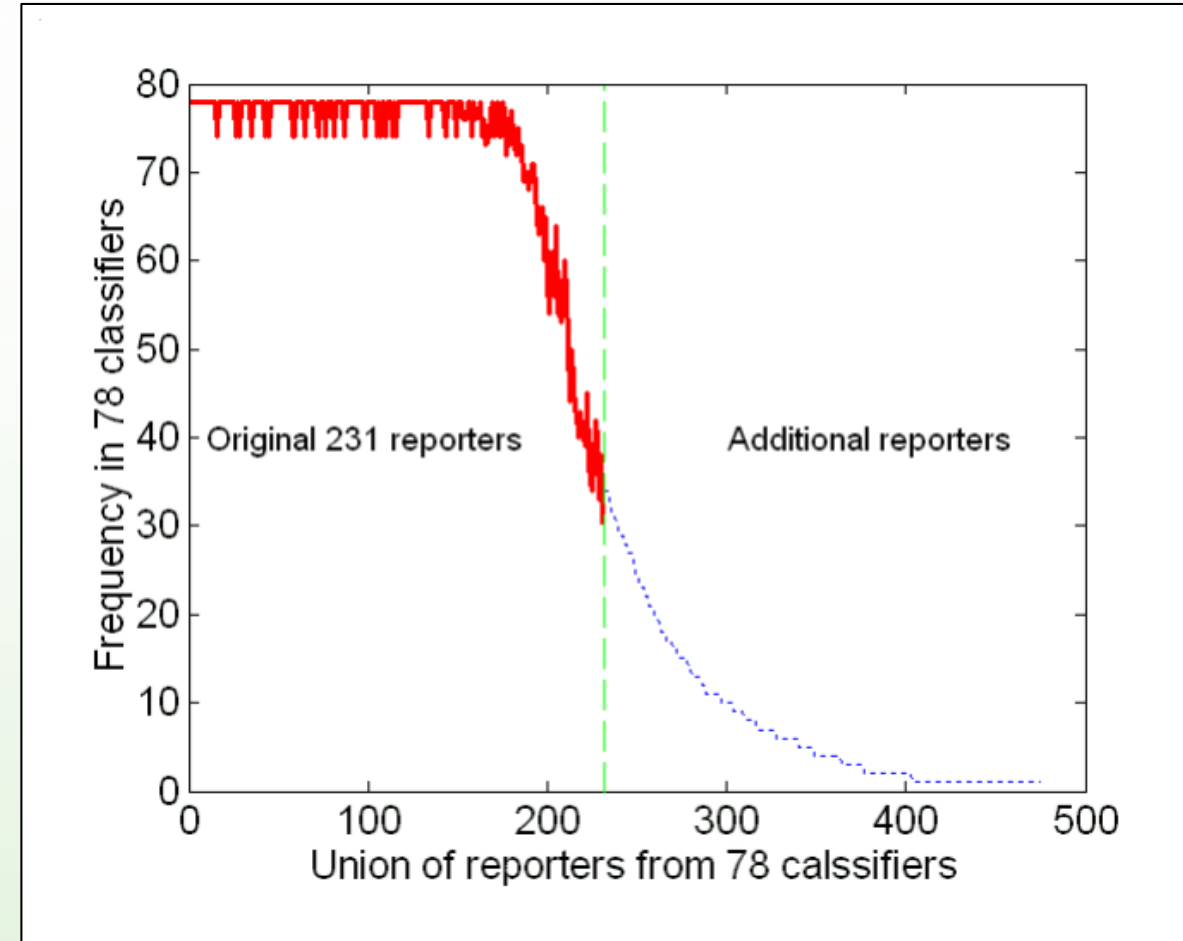
# A 70-GENE SIGNATURE FOR RISK OF METASTASIS

Correlation between the prognostic category (metastasis vs. no-metastasis) and the logarithmic expression ratio across all 78 samples for each individual gene in 5,000 significantly expressed genes.

Permutation to evaluate significance.
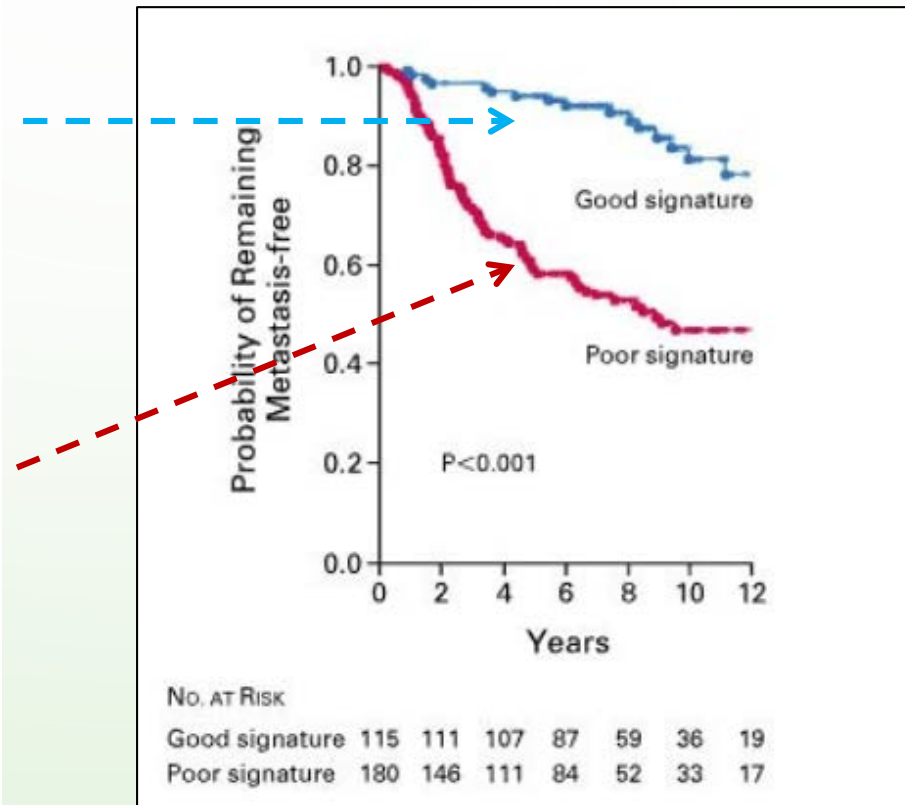
231 genes corr coeff > 0.3 or < − 0.3 selected (p= 0.3%)

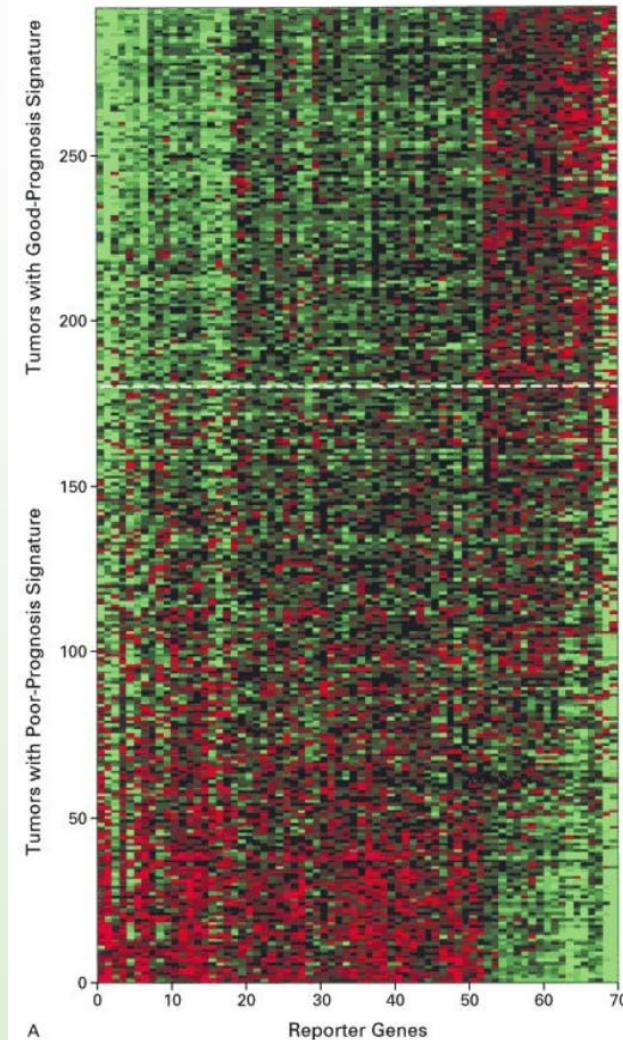Leave-one-out: (1) leave one sample out, (2) define reporters based on the remaining 77 samples among the set of ~5000 significant genes, (3) use the reporters to predict the outcome of the one sample that was left out in step (1), (4) repeat steps (1)-(3) exhaustively for all 78 samples.

Select top 70

# 70-GENE SIGNATURE IS PREDICTIVE OF SURVIVAL AND RISK OF METASTASES [VAN DE VIJVER ET AL, N ENGL J MED, 2002]

# Performance of gene expression signatures in published studies

| Study reference | Cancer type | Clinical endpoint | Sample size | Number of events (%) | Number of channels (type) | Number of genes after filtration* |
|---|---|---|---|---|---|---|
| 2 | Non-Hodgkin lymphoma | Survival | 240 | 138 (58%) | 2 (Lymphochip) | 6693 |
| 3 | Acute lymphocytic leukaemia | Relapse-free survival | 233 | 32 (14%) | 1 (Affymetrix) | 12 236 |
| 4 | Breast cancer | 5-year metastasis-free survival | 97 | 46 (47%) | 2 (Agilent) | 4948 |
| 5 | Lung adenocarcinoma | Survival | 86 | 24 (28%) | 1 (Affymetrix) | 6532 |
| 6,7 | Lung adenocarcinoma | 4-year survival | 62† | 31 (50%) | 1 (Affymetrix) | 5403 |
| 8 | Medulloblastoma | Survival | 60 | 21 (35%) | 1 (Affymetrix) | 6778 |
| 9 | Hepatocellular carcinoma | 1-year recurrence-free survival | 60 | 20 (33%) | 1 (Affymetrix) | 4861 |

*For the data of van 't Veer and colleagues,[4] the same filter was used as in the original publication. For other studies, genes with little variation in expression were excluded. †Only patients with clinical follow-up of at least 4 years after surgical resection were analysed.[7]

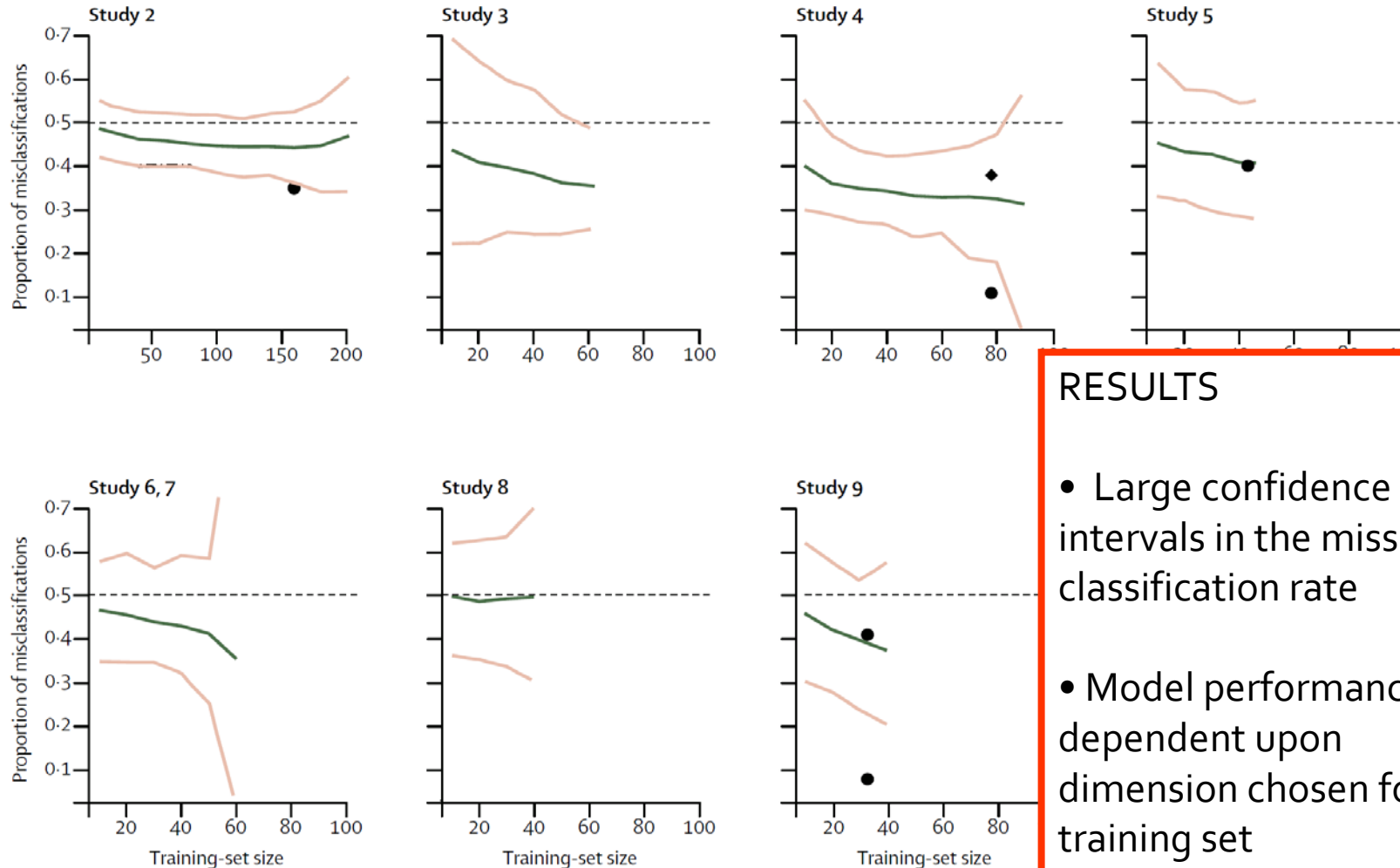**Table:** Description of eligible studies ordered by sample size

*Michiels et al, Lancet 2005*

Resampling strategy:
1) Sample a set of patients for training at random, and leave the rest for validation
2) Identify a gene expression signature (GES) for the clinical endpoint in the training set
3) Predict for the patients in validation set and estimate the proportion of misclassifications
4) Iterate on multiple random sets to study stability and performance of the GES

# Performance of gene expression signatures in published studies

Misclassification rate from 500 random training-validation sets vs. training-set size (mean and 95% CIs)



RESULTS

• Large confidence intervals in the miss-classification rate

• Model performance dependent upon dimension chosen for the training set

# Performance of gene expression signatures in published studies

Misclassification rate from 500 random training-validation sets vs. training-set size (mean and 95% CIs)
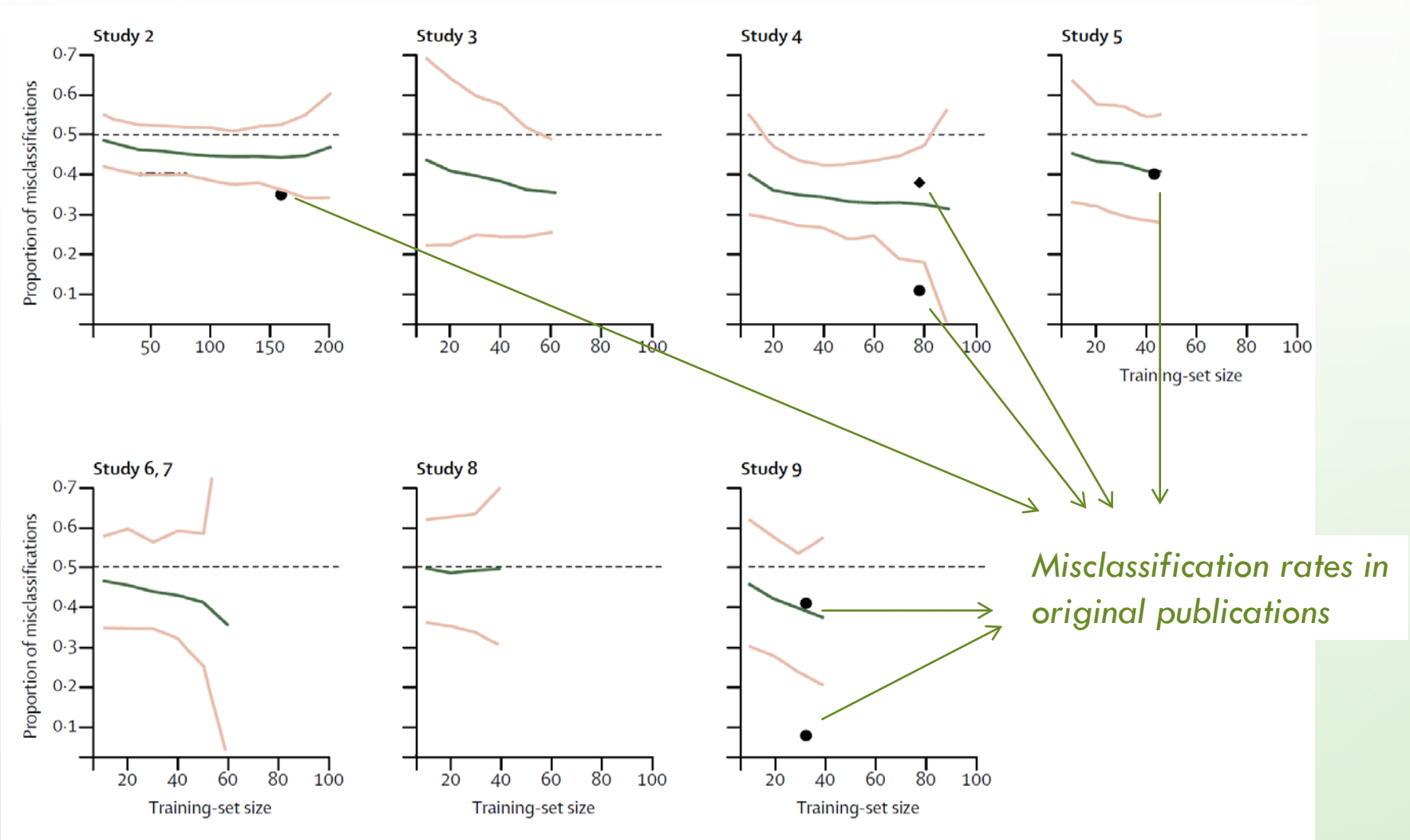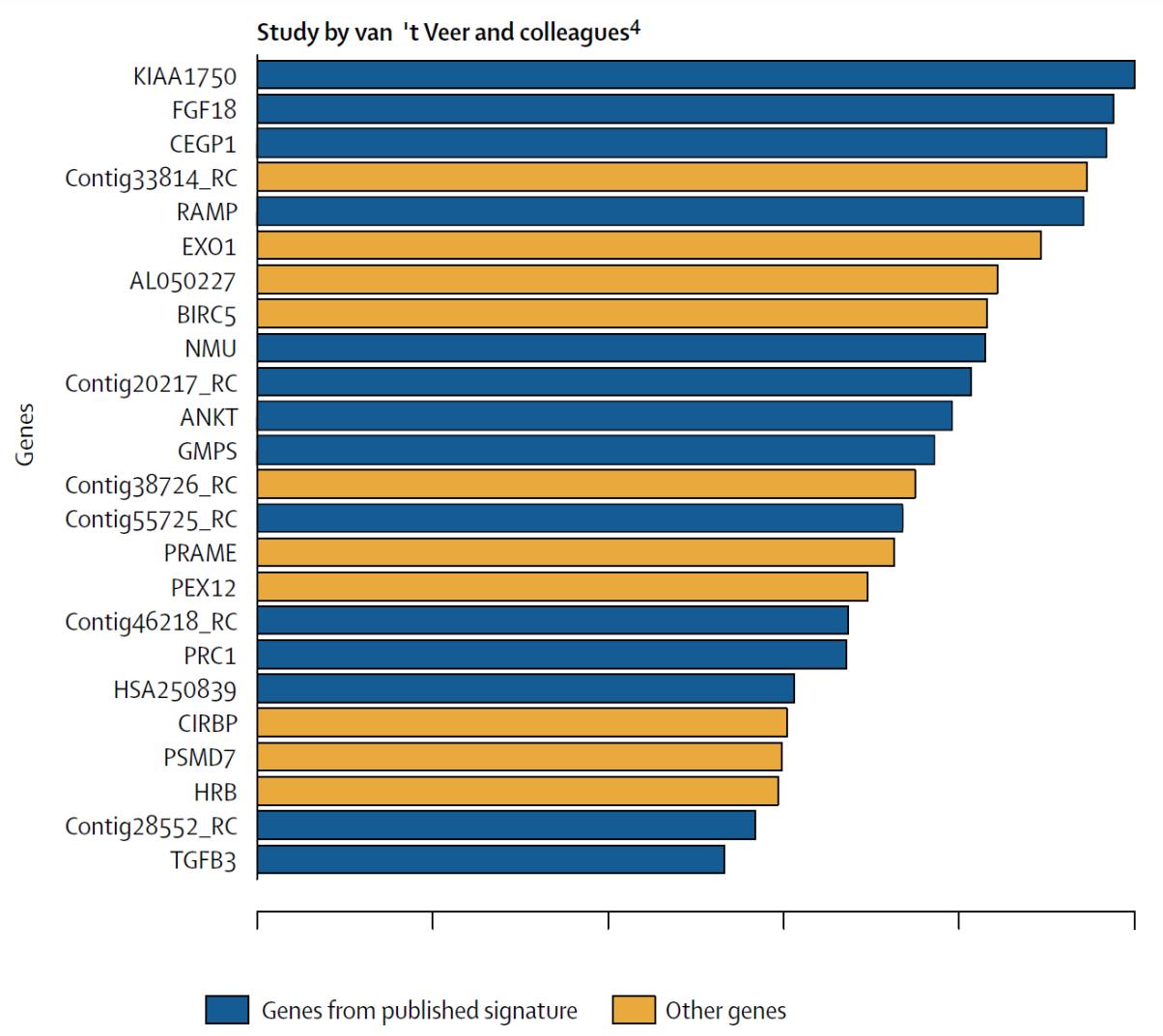


*Misclassification rates in original publications*
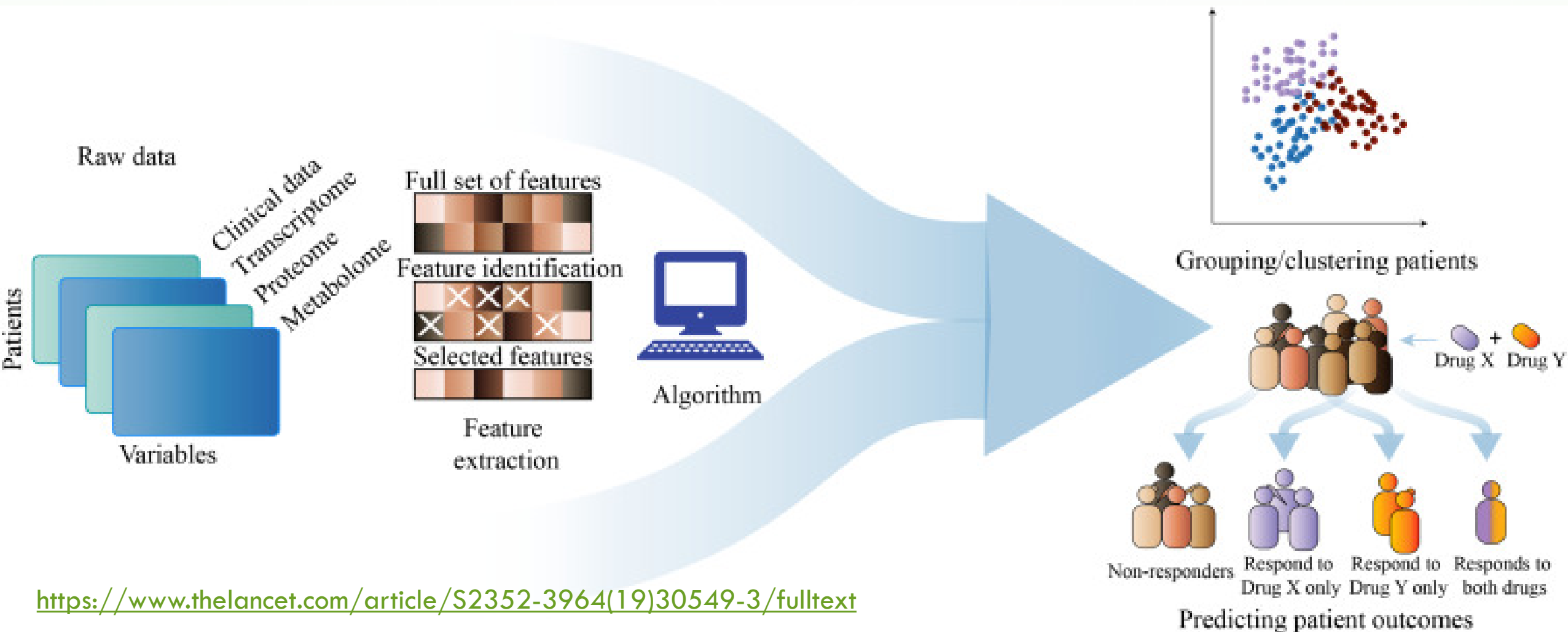
# LARGE VARIABILITY IN SELECTED GENES

• *Several possible models with similar correlation with outcome*

• *Very little overlap in gene content between models*

24 genes within the common set were selected in at least 50% of the random simulations



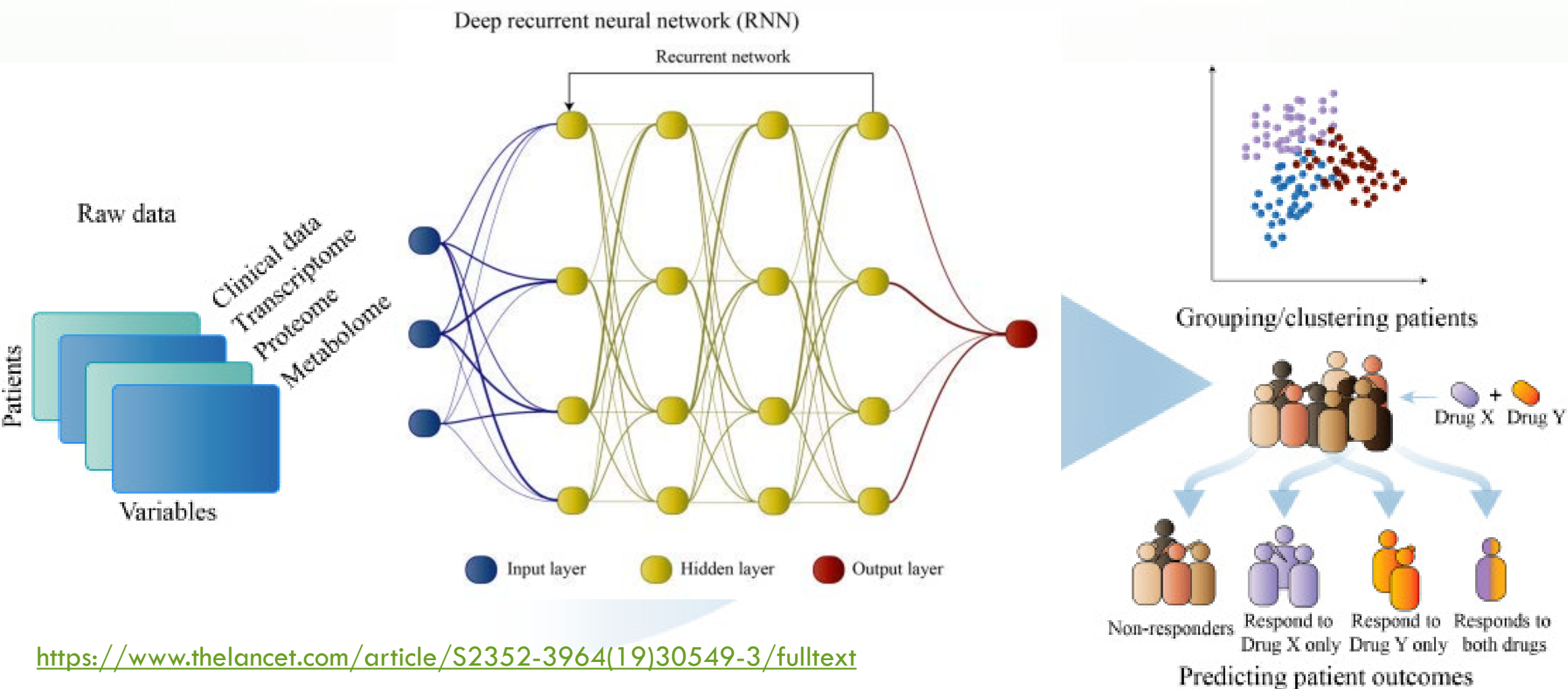Study by van 't Veer and colleagues[4]

Genes from published signature | Other genes

# (DEEP) LEARNING TO IDENTIFY PREDICTIVE BIOMARKERS

# (DEEP) LEARNING TO IDENTIFY PREDICTIVE BIOMARKERS



https://www.thelancet.com/article/S2352-3964(19)30549-3/fulltext

# Multi-omic machine learning predictor of breast cancer therapy response

Nature, 623–629 (2022)



https://www.nature.com/articles/s41586-021-04278-5

Clinical work flow and data acquisition

Prediction of response to neoadjuvant therapies

**Methods**
- clinical, digital pathology, genomic and transcriptomic profiles of pre-treatment biopsies of breast tumours from 168 patients treated with chemotherapy with or without HER2 (encoded by *ERBB2*)-targeted therapy before surgery.
- Pathology end points: complete response or residual disease
- Multi-omic features in these diagnostic biopsies

**Multi-omic machine learning predictor of breast cancer therapy response**
Nature, 623–629 (2022)



a

**Methods**

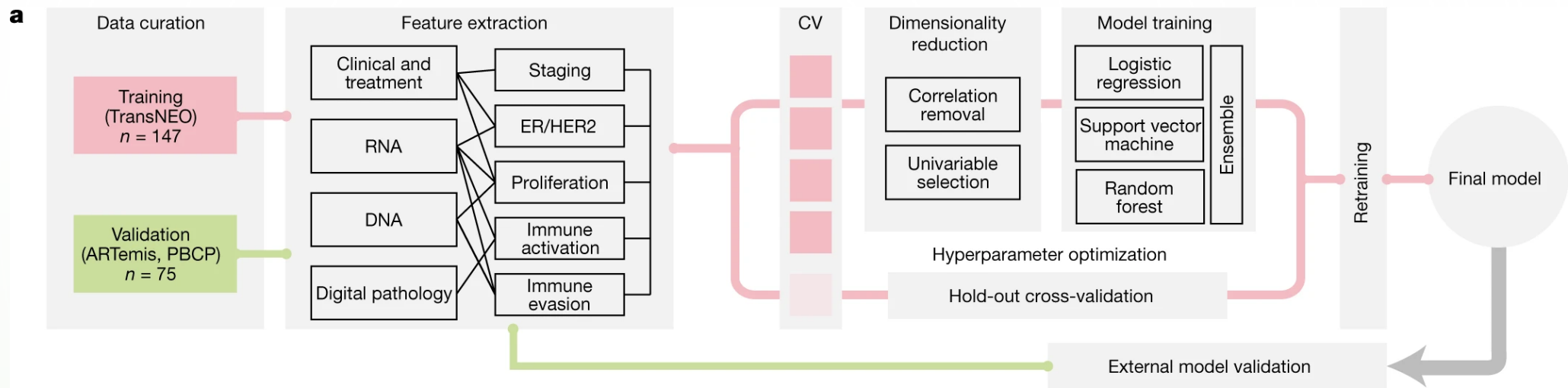-   Six pCR prediction models including different feature combinations were derived using:
    (1) clinical features only, and adding (2) DNA, (3) RNA, (4) DNA and RNA, (5) DNA, RNA and digital pathology, and (6) DNA, RNA, digital pathology and treatment.

-   The models were based on a multi-step predictor pipeline. Features were first filtered by univariable selection and collinearity reduction, and then fed into an unweighted ensemble classifier.

-   Each ensemble consisted of three algorithms acting in parallel: logistic regression with elastic net regularization, a support vector machine and a random forest. The three algorithm scores were then averaged to form the predictor.

-   A fivefold cross-validation scheme was used to optimize model hyperparameters

To maximise the robustness of the predictions, the models contain two levels of averaging. Firstly, predictions are obtained by averaging three classifier pipelines, as follows:

$$Prob(attaining\ pCR) = \frac{1}{3} \times (Pipeline_{LR}^{HER2+} + Pipeline_{SVC}^{HER2+} + Pipeline_{RF}^{HER2+}),$$

where

$$Pipeline_{Classifier} = \{Coll.Reduction\ (0.8) \Rightarrow Univ.selection \Rightarrow Classifier\},$$

with the corresponding hyperparameters listed in Tables 1 and 2. In particular, model hyperparameters that were *fixed a priori* are listed in Table 1, while model hyperparameters that were *optimised* for each of the classifiers using a 5-fold cross-validation setup are listed in Table 3. Once all hyperparameters are set, the model is re-trained on the entire training cohort, and subsequently frozen.

Secondly, to account for possible biases in the optimisation due to the particular cross validation splitting used, we repeated the process explained above 5 times, with 5 different cross-validation splitting seeds (integers from 1 to 5). As a result, we obtain 5 alternative optimised models. The final predictions are the average of the 5.

From Suppl Methods:
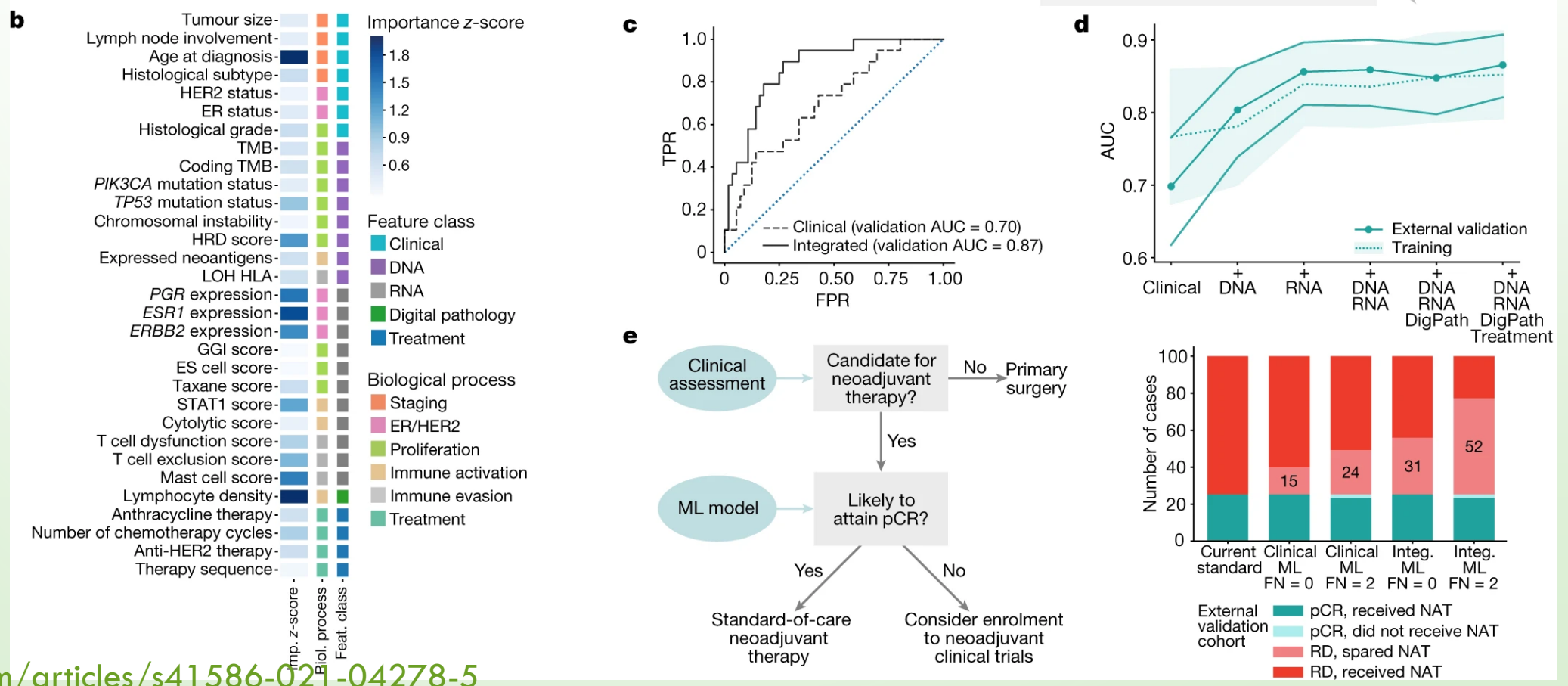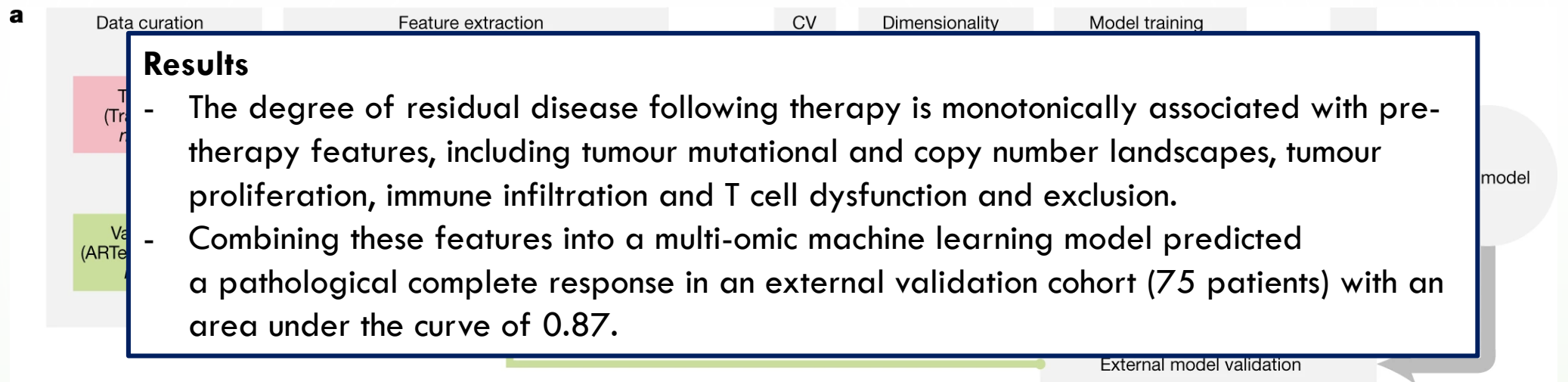https://static-content.springer.com/esm/art%3A10.1038%2Fs41586-021-04278-5/MediaObjects/41586_2021_4278_MOESM1_ESM.pdf
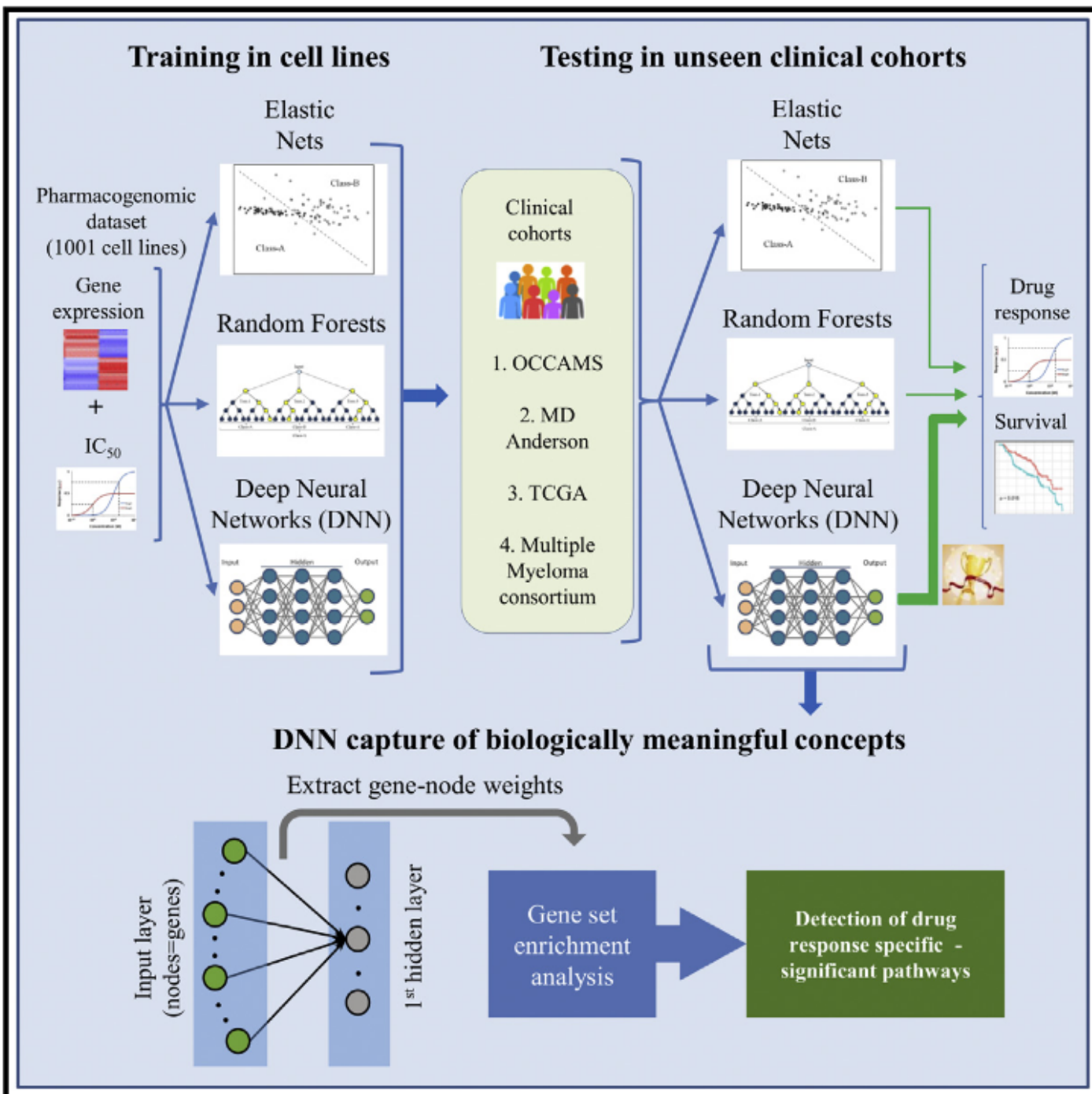Code used:
https://github.com/cclab-brca/neoadjuvant-therapy-response-predictor/blob/master/R/06%20-%20ML%20predictor.R

**Multi-omic machine learning predictor of breast cancer therapy response**
Nature, 623–629 (2022)

**Results**
- The degree of residual disease following therapy is monotonically associated with pre-therapy features, including tumour mutational and copy number landscapes, tumour proliferation, immune infiltration and T cell dysfunction and exclusion.
- Combining these features into a multi-omic machine learning model predicted a pathological complete response in an external validation cohort (75 patients) with an area under the curve of 0.87.

# A Deep Learning Framework for Predicting Response to Therapy in Cancer

**Highlights**
- A machine learning (ML) workflow is designed to predict drug response in cancer patients
- Deep neural networks (DNNs) surpass current ML algorithms in drug response prediction
- DNNs predict drug response and survival in various large clinical cohorts
- DNNs capture intricate biological interactions linked to specific drug response pathways