

The Influence of DNA Sequence-Derived Features across the 'omics scales

Gregory Parkes, Mahesan Niranjan

Email: g.m.parkes@soton.ac.uk

Overview

➤ What are sequence-derived features?

➤ What are 'omics or multi-'omics?

➤ Exposure to examples and challenges with SDFs

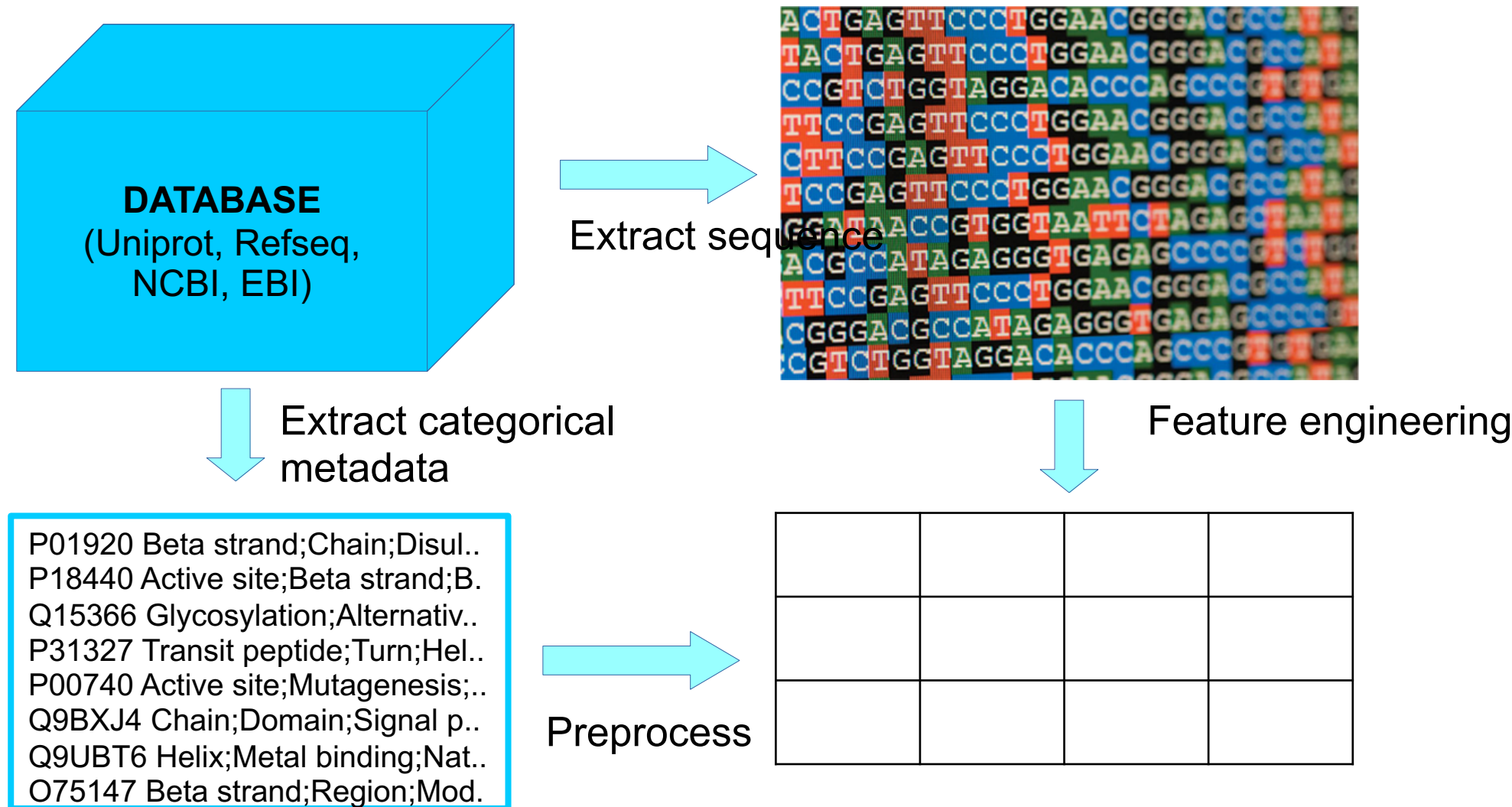
➤ Use of sequence-derived features in yeast for protein prediction

➤ Use of sequence-derived features in human cell cycle study

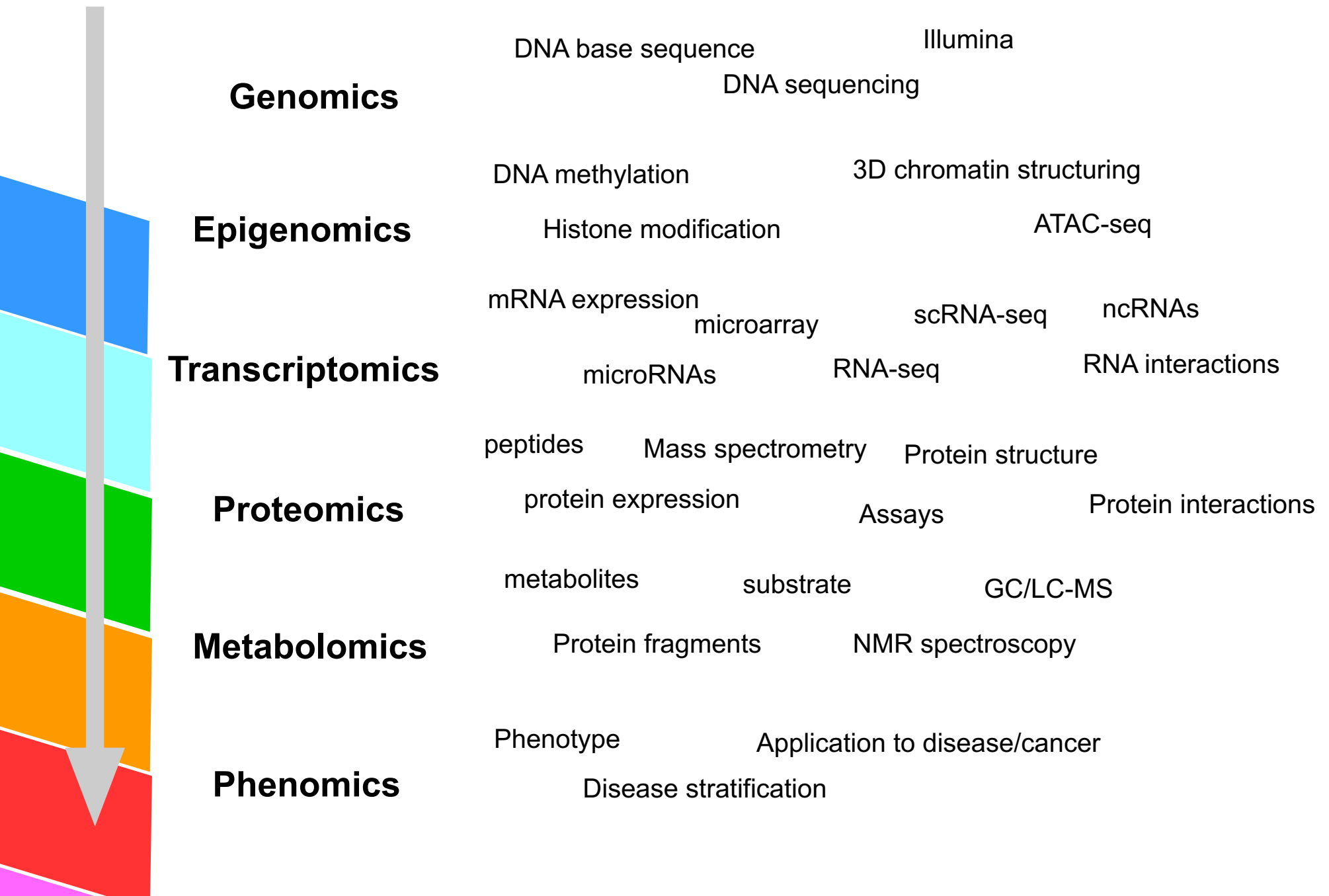
➤ Expanding the number of SDFs, use as proxy

➤ SDFs contribute useful 'static' information into model

What are sequence-derived features?



Analysis and **integration** of multiple levels of expression



Inherent challenges with using data from **different expression levels**

DNA



SEQUENCE

ATGTGGGCTTATAAATGTGCGGTACCAGCCCCCTGTCAATGAGTGGCCCTATACCTCCATCGGCTGATC

STABLE

GC%, Epigenetic markings, Chromosome, Locus

RNA

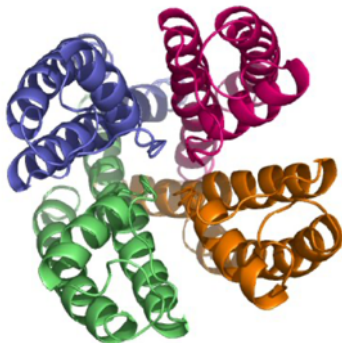
SEQUENCE

AUGACCUGACCUUCACUCUGGCGGCACAAAUGGUCCAUAUGCGCCAUCGUACGUGCGGCUAGCU

UNSTABLE (10h t_{1/2})

mRNA concentration, Codon bias, sequence length,*

Protein



SEQUENCE

MKKNRRSRWL VATAGKRSPSFL LIVAARERW**STOP**

UNSTABLE (20-40h t_{1/2})

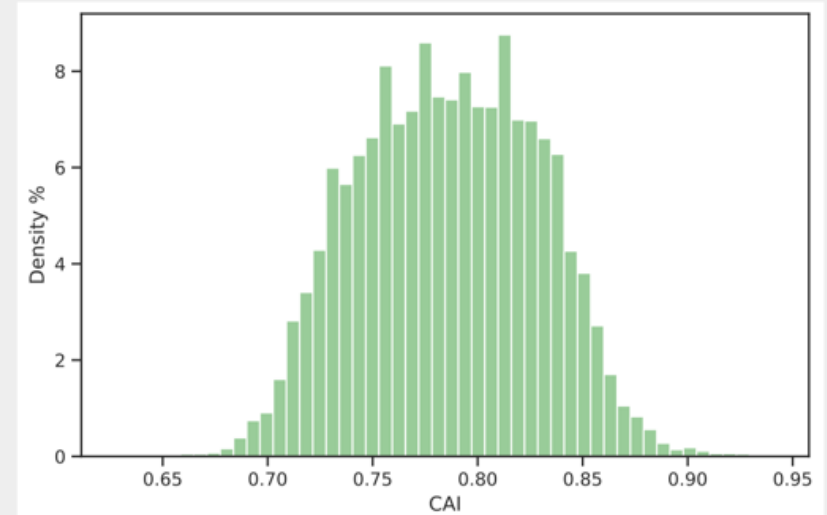
Translation concentration, Protein* concentration, a*

Example: Making a **sequence-derived feature**; RCB & MFE

Relative Codon Usage/Bias [1]

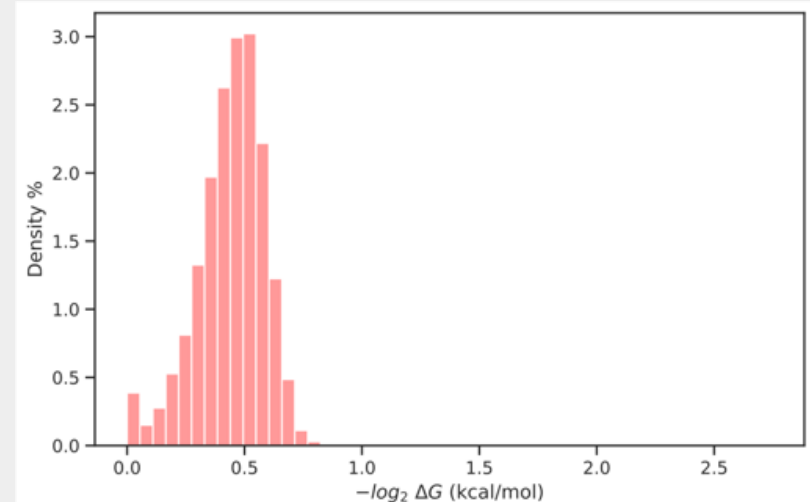
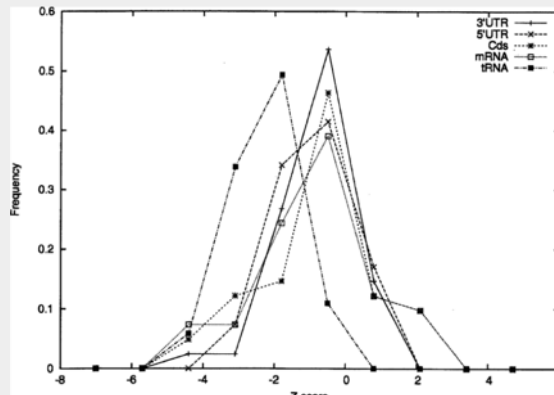
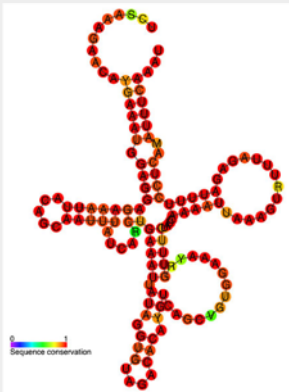
$$RCB_{xyz} = \frac{f(x, y, z)}{f_1(x)f_2(y)f_3(z)},$$

$$RCBS = \left(\prod_{l=1}^L RCB_{xyz}(l) \right)^{1/L} - 1$$



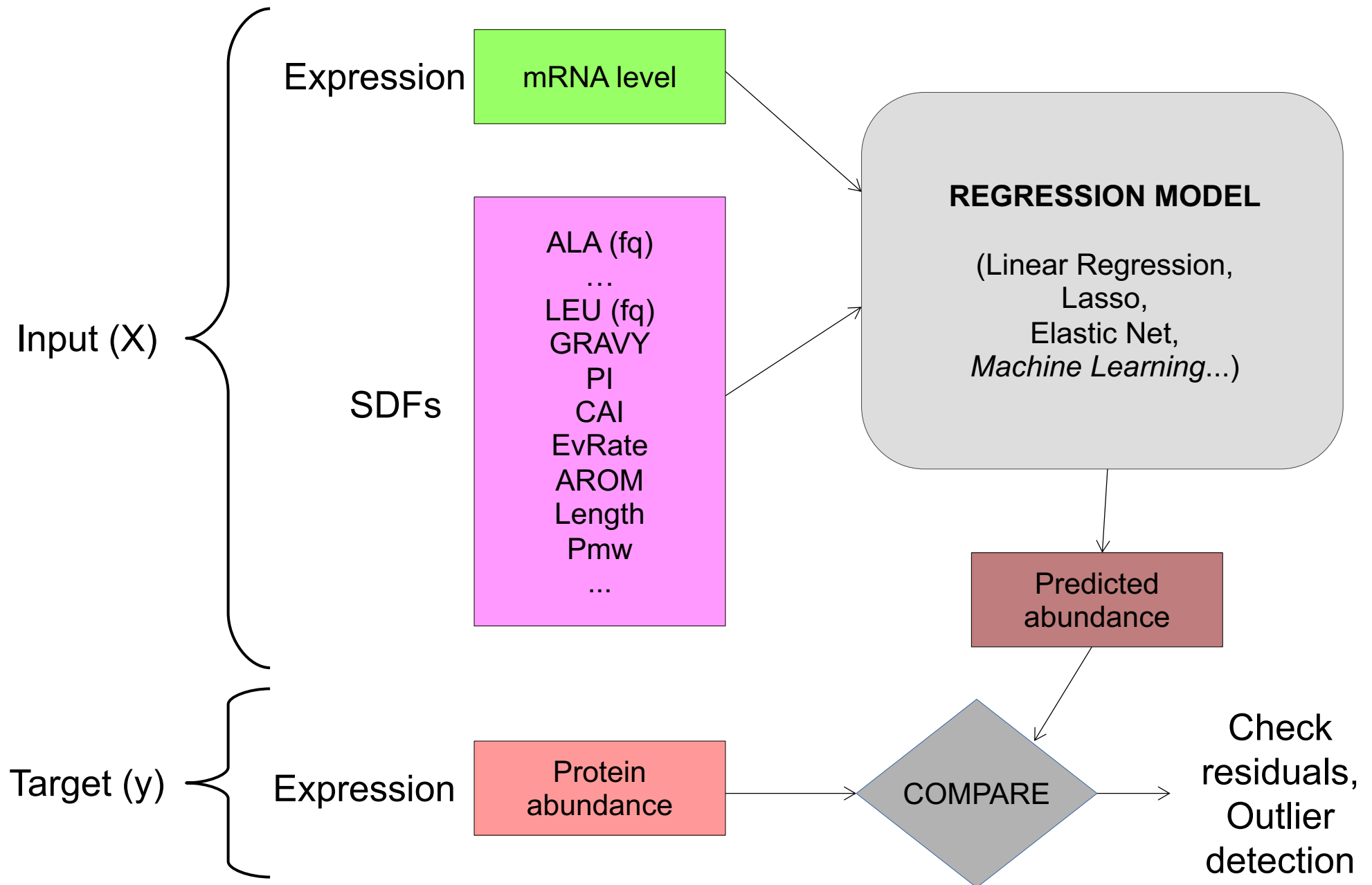
[1] Fox, J. M., & Erill, I. (2010). Relative codon adaptation: a generic codon bias index for prediction of gene expression. DNA research : an international journal for rapid publi

Minimum Free Energy (MFE) [2]

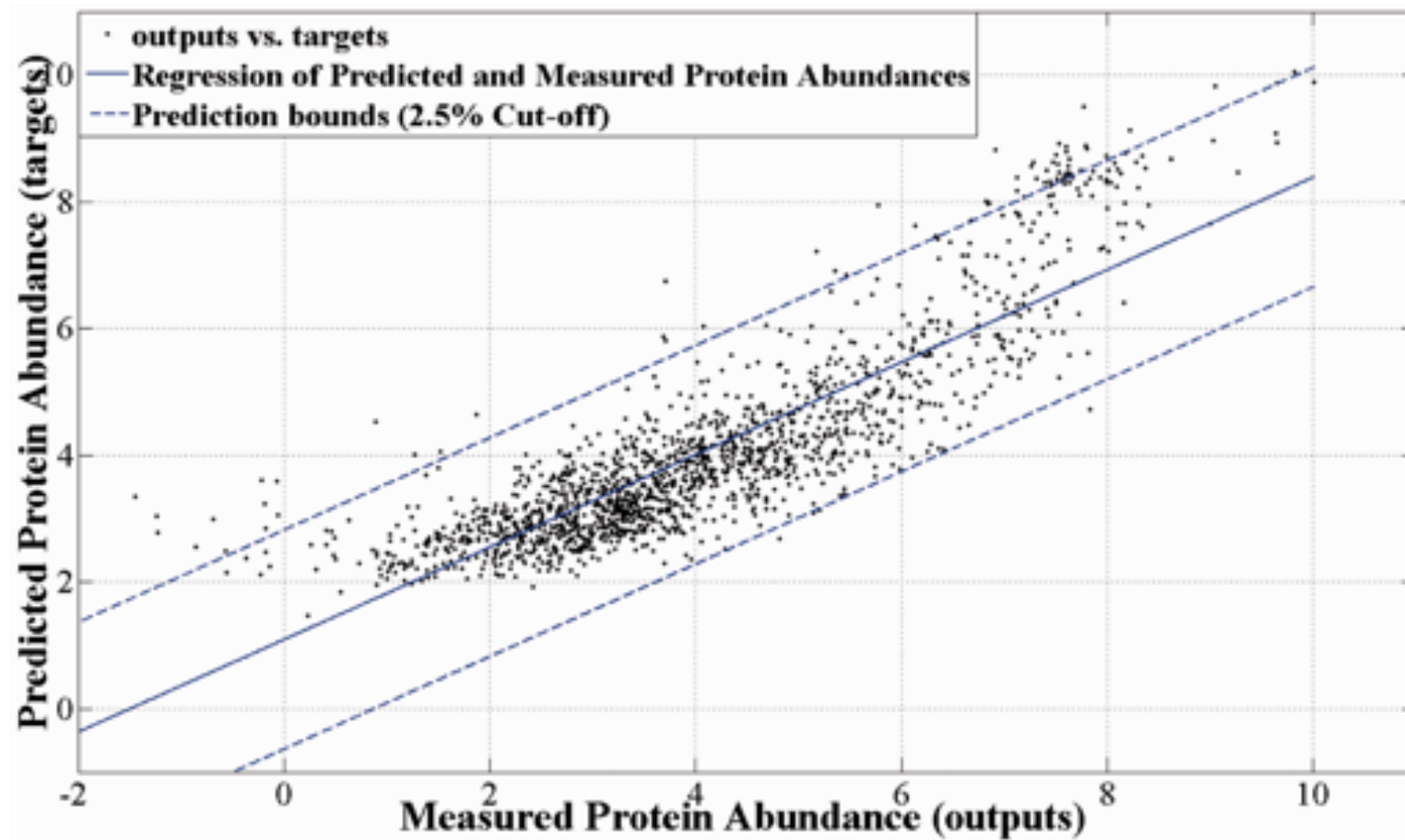


[2] Ringnér M, Krogh M (2005) Folding Free Energies of 5'-UTRs Impact Post-Transcriptional Regulation on a Genomic Scale in Yeast. PLOS Computational Biology 1(7): e72.

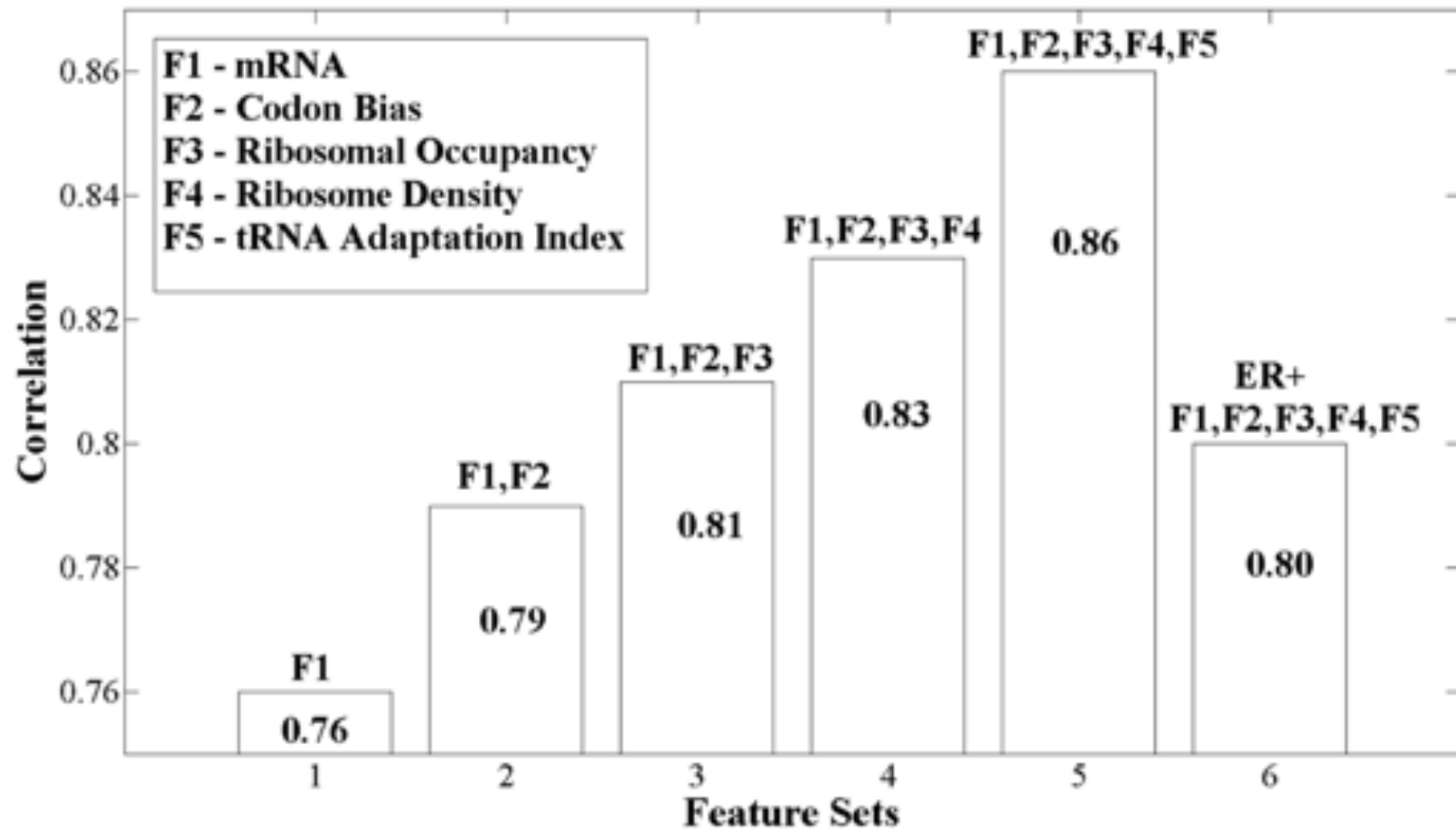
Case study 1: Use of SDFs in modelling protein abundance in yeast



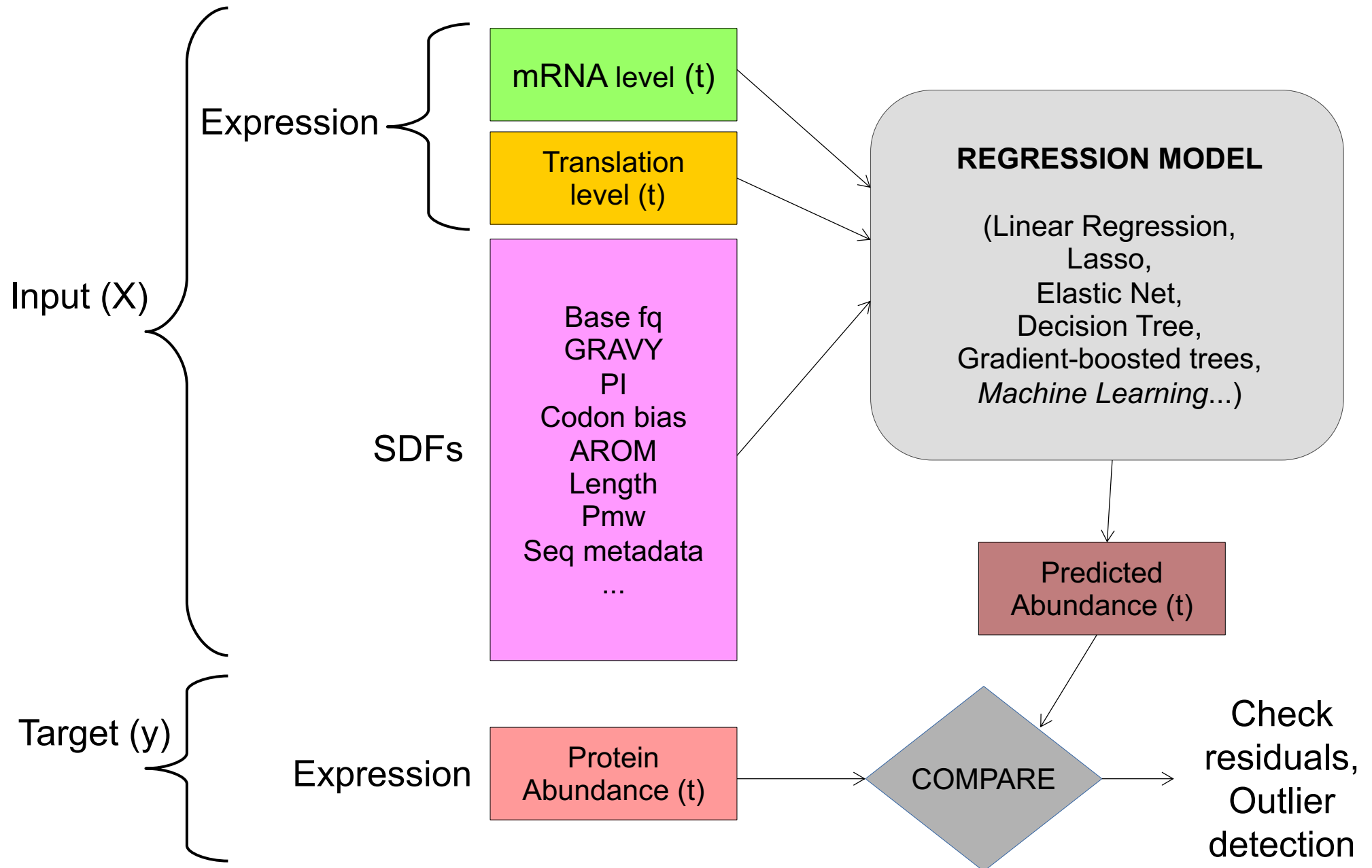
Linear models of mRNA proxy for protein level



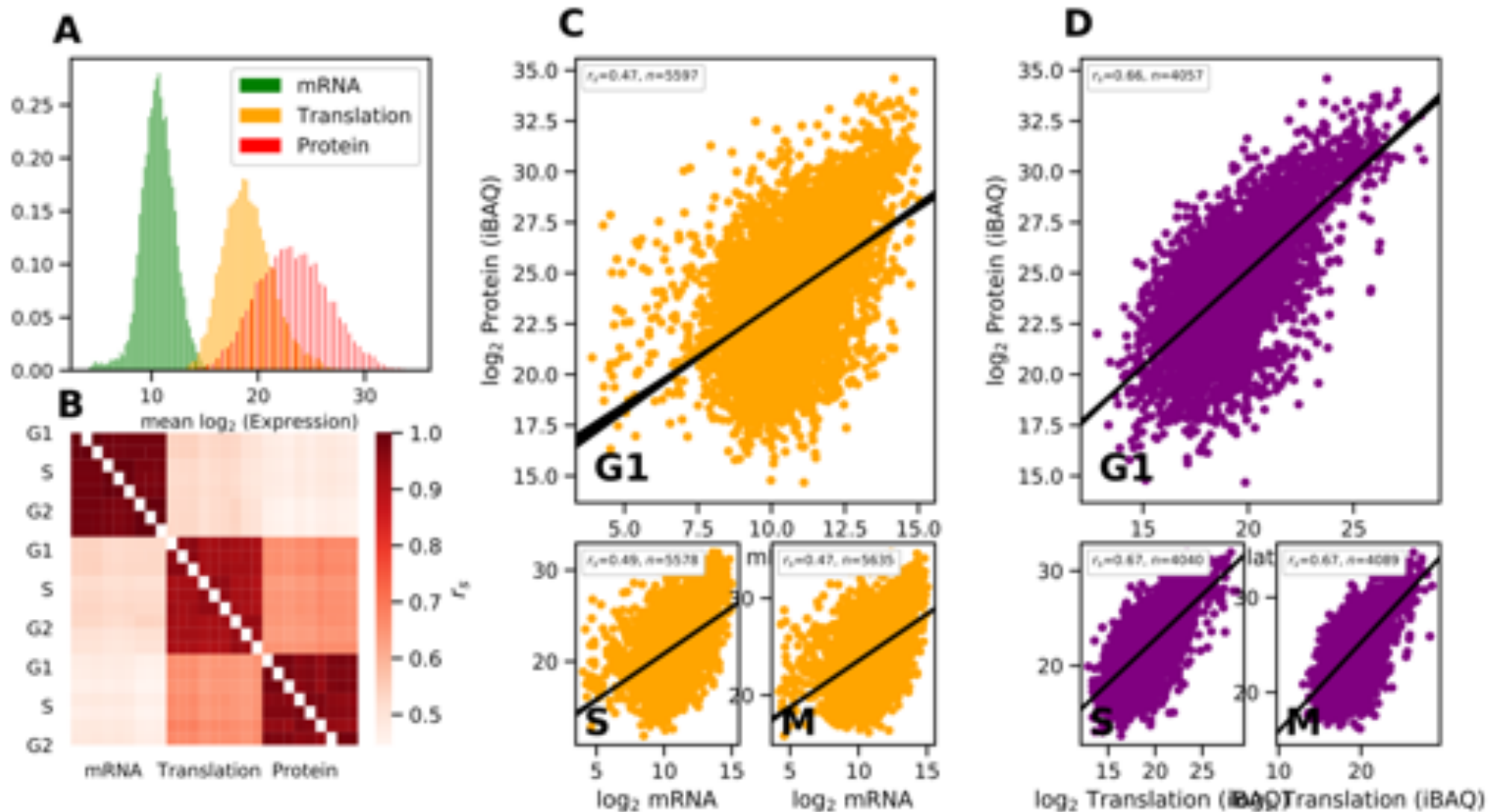
Use of sequence-derived features to additive model



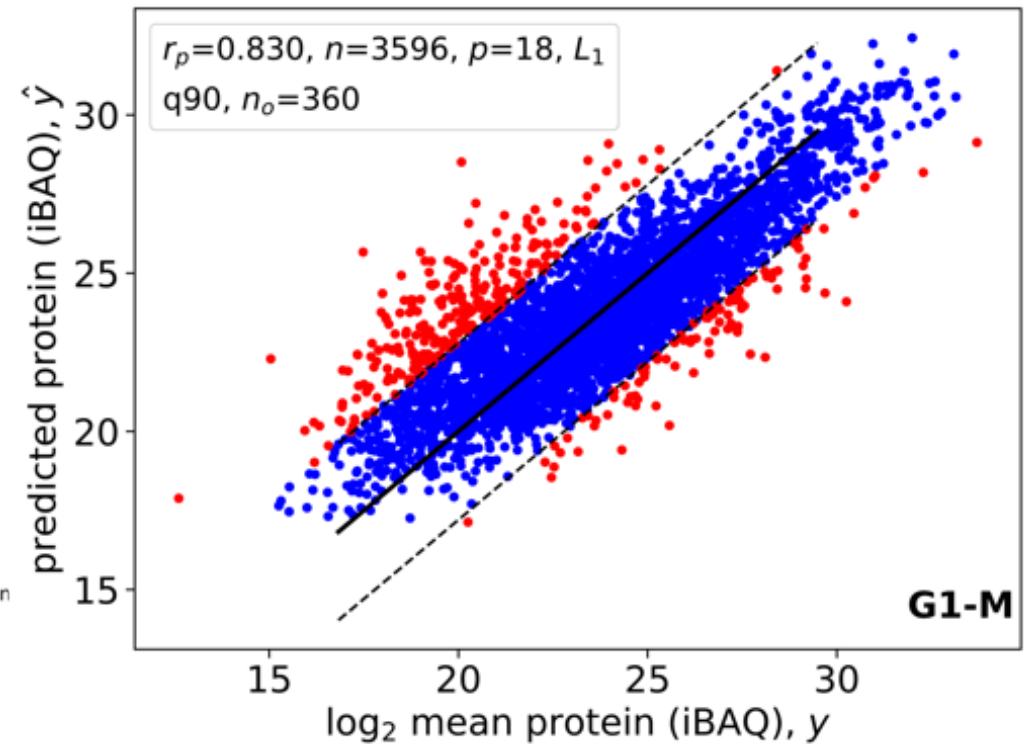
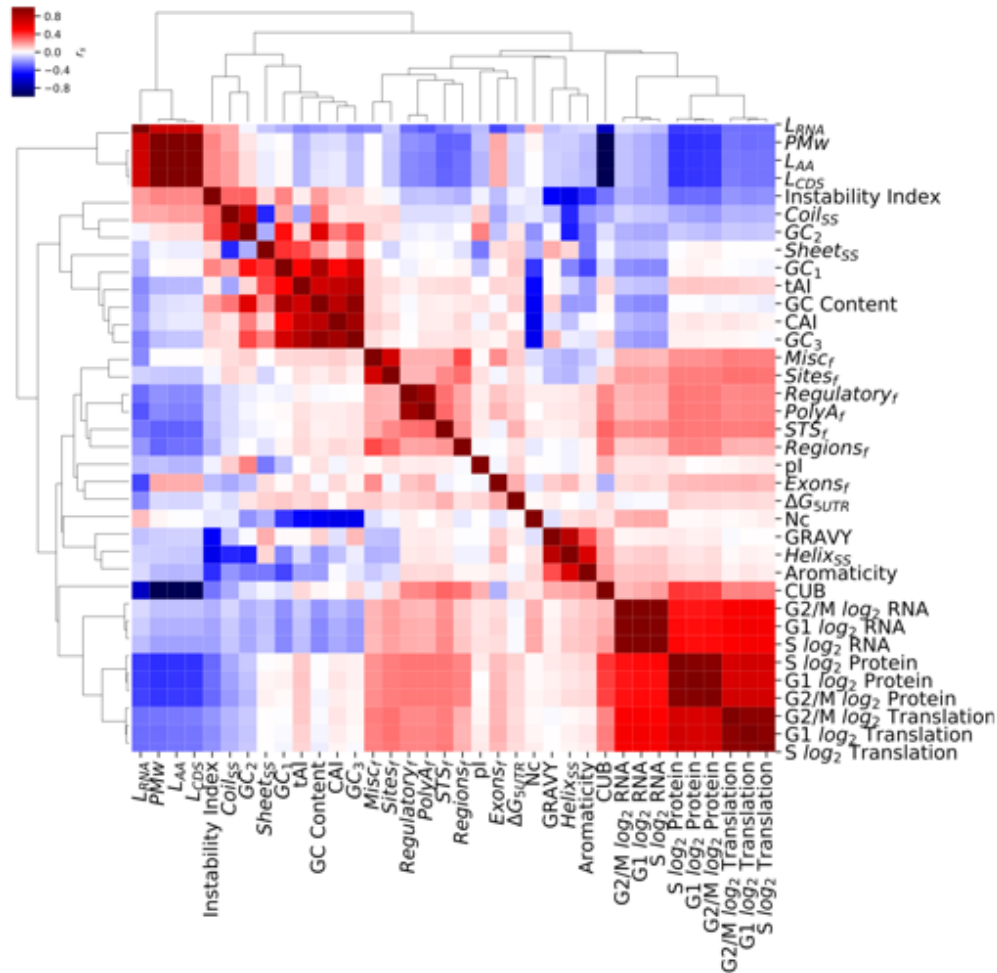
Case study 2: Use of SDFs in modelling protein abundance in human cell cycle



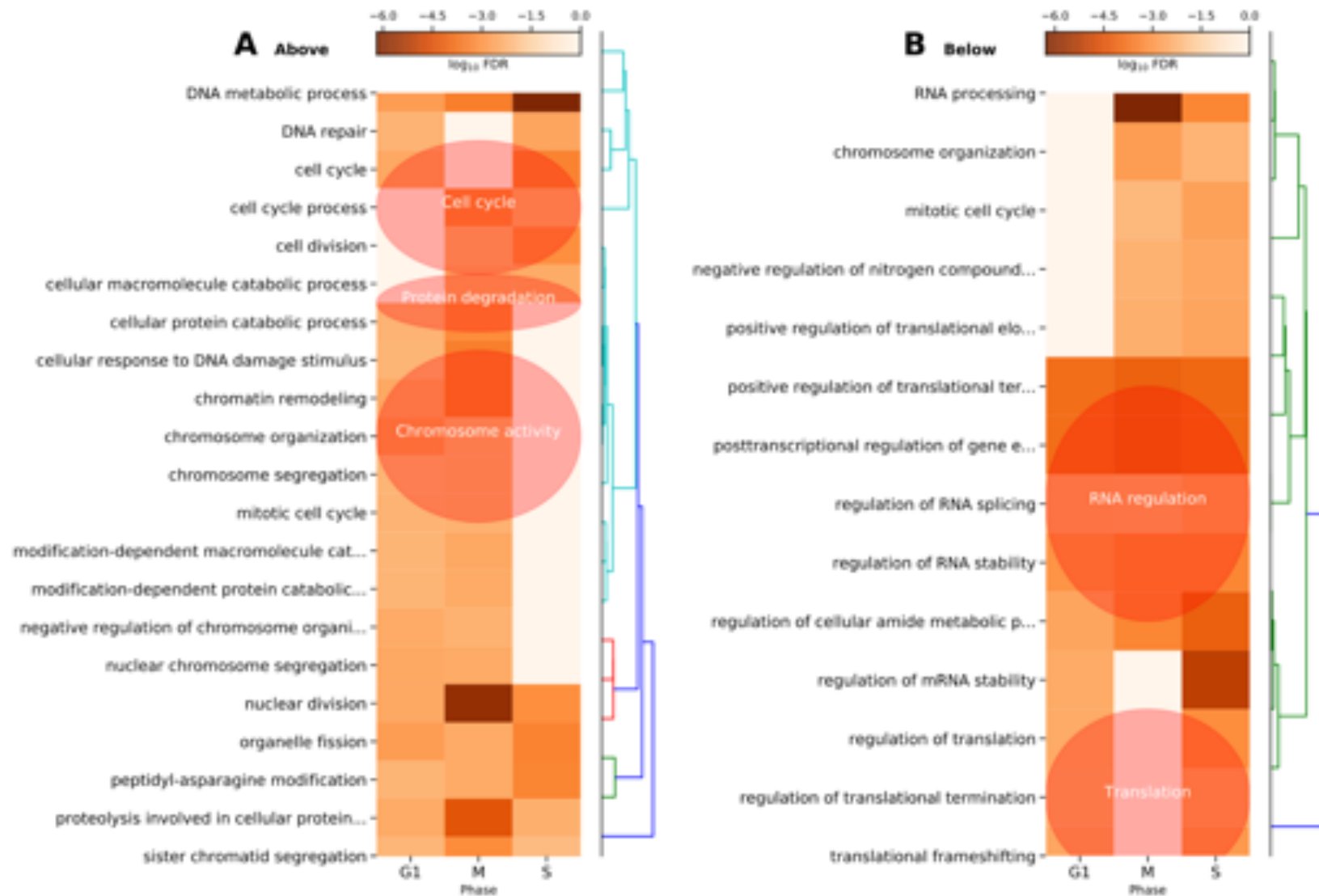
Linear models of relationship between mRNA, translation and protein



Inclusion of SDF demonstrates increase in model prediction

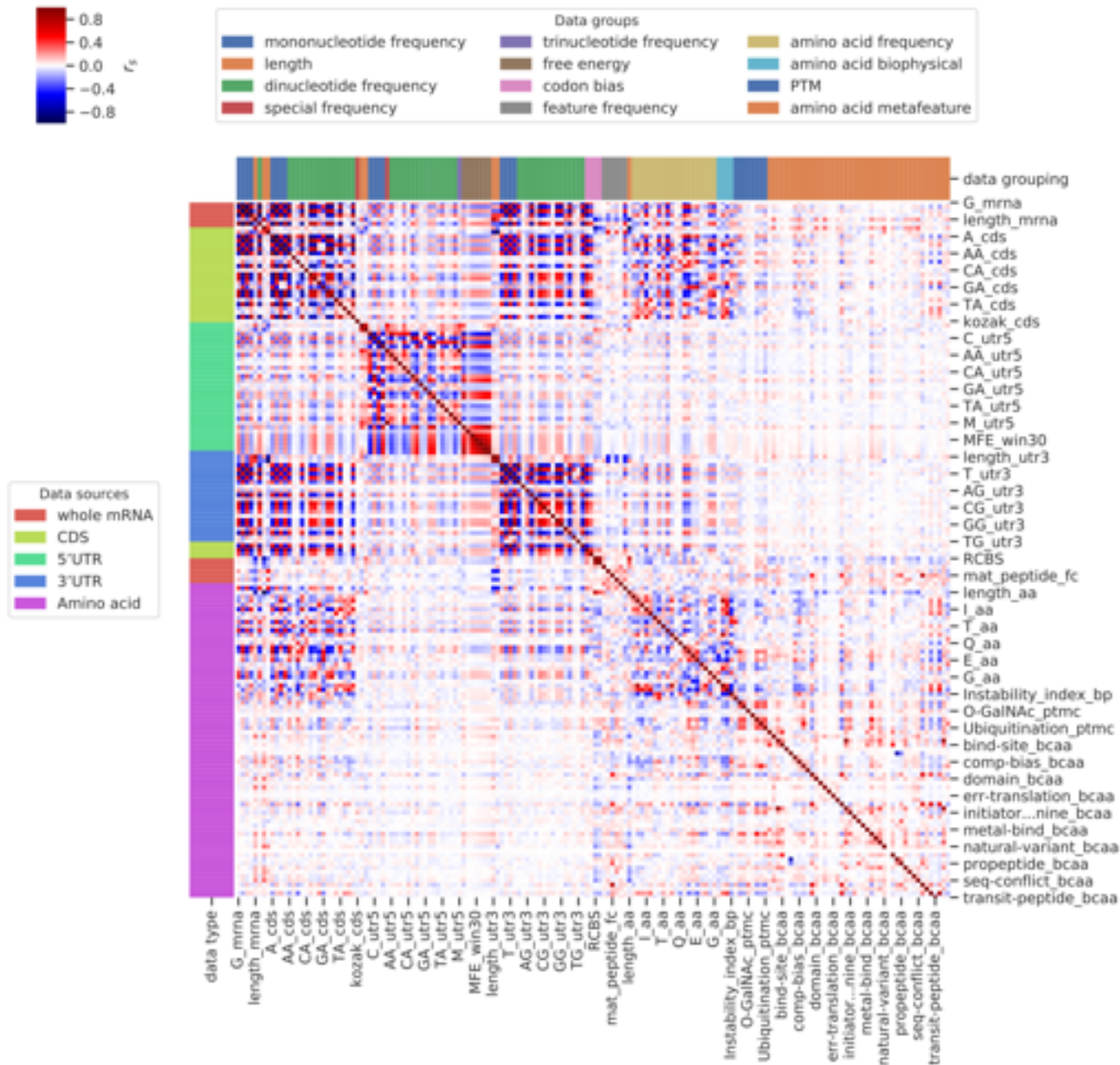


Outliers to SDF have functionally predictable GO terms

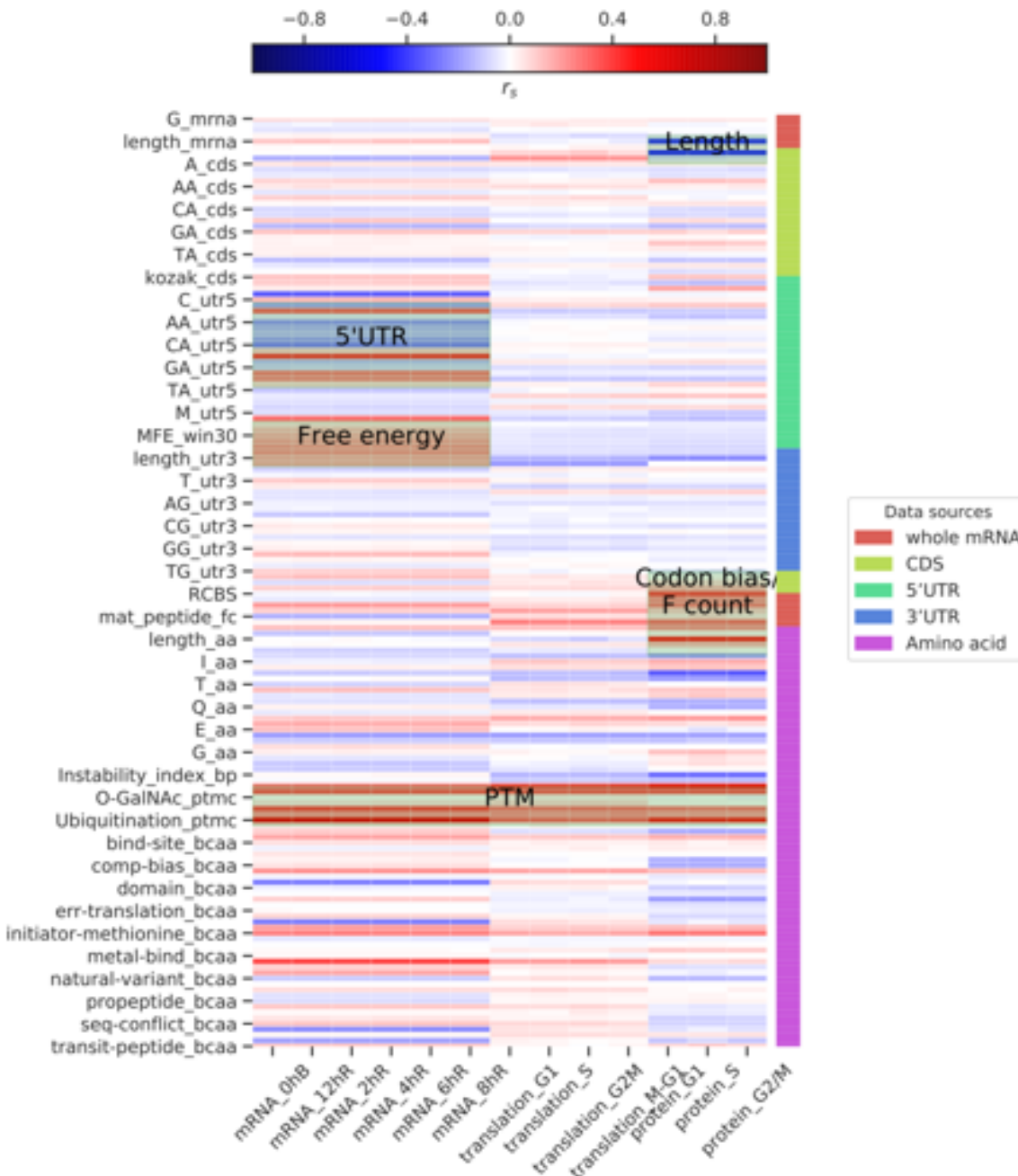


Parkes, Gregory & Niranjan, Mahesan. (2019). Uncovering Extensive Post-Translation Regulation During Human Cell Cycle Progress

Case study 3: Correlation matrix of SDFs reveals data-source subgroups

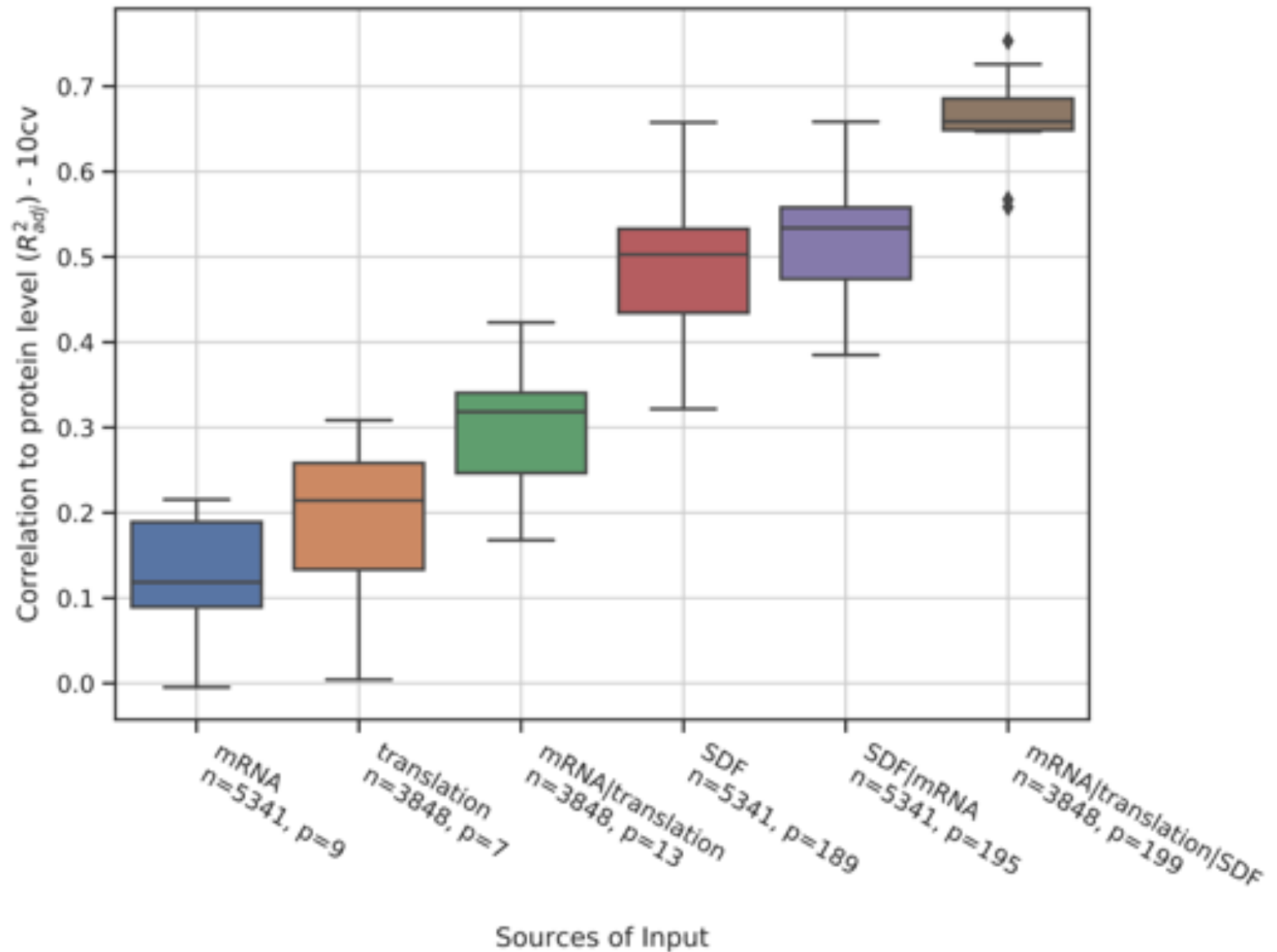


Correlation to multi-'omic features reveals SDF importance



Feature	Target / Correlation R-sq
Acetylation (PTM)	mRNA (0.06), Translation (~0.13) Protein (~0.19)
Ubiquitination (PTM)	mRNA (~0.3) , Translation (~0.09), Protein (~0.19)
Length (CDS, mRNA, AA)	Protein (~0.15-0.17)
Modified residue (bcAA)	mRNA (~0.13), Translation (~0.03)
Proline (AA)	Translation (0.01), Protein (0.04)
RCBS	Protein (~0.13)
Initiator-Methionine (bcAA)	mRNA (0.06), Translation (0.02), Protein (0.07)

Can SDFs act as a reliable proxy for mRNA abundance?



Summary

- What are sequence-derived features?
- What are 'omics or multi-'omics?
- Exposure to examples and challenges with SDFs
- Use of sequence-derived features in yeast for protein prediction
- Use of sequence-derived features in human cell cycle study
- Expanding the number of SDFs, use as proxy
- SDFs contribute useful 'static' information into models

