

# Both datasets exploratory

## Contents

<b>Objective</b>	<b>1</b>
<b>TPM summaries</b>	<b>1</b>
<b>Filtering out genes</b>	<b>4</b>
<b>GSEA analysis</b>	<b>5</b>
Hallmark analysis . . . . .	5
Curated analysis . . . . .	6
<b>Further analysis</b>	<b>6</b>
Hallmark pathways . . . . .	6
Curated pathways . . . . .	23

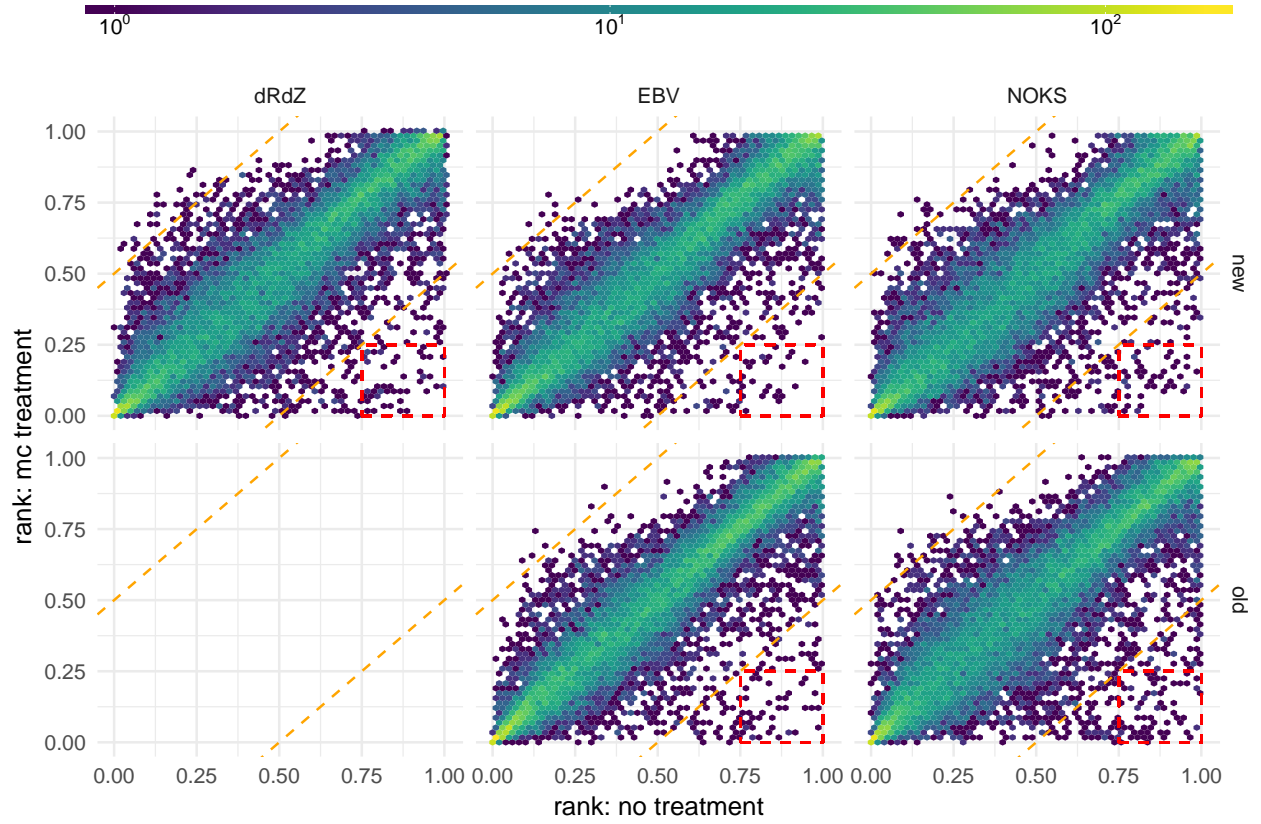
## Objective

The idea of this document is to draft what we need to do to integrate the GSEA pathway analysis with the diff. expression analysis we did with DESeq2. We propose the following pipeline:

1. We fit a model with a formula of the type `Count_ij ~ Treat_i + Cell_j + Interaction_ij`
2. For every gene, test an hypothesis if there is a treatment effect in the cell specific expression.
3. For every gene, summarize the TPM level of each treatment, cell and interaction.
4. For every cell, take the genes that are differentially expressed and are among the top  $K$  most expressed genes for both treatments (the assumption here is to avoid cases where the `log2FC` value is extreme due to a very low quantity of reads in either treatment).
5. Using the `t.stat` values as a signal-2-noise metric, do a pathway analysis with GSEA.

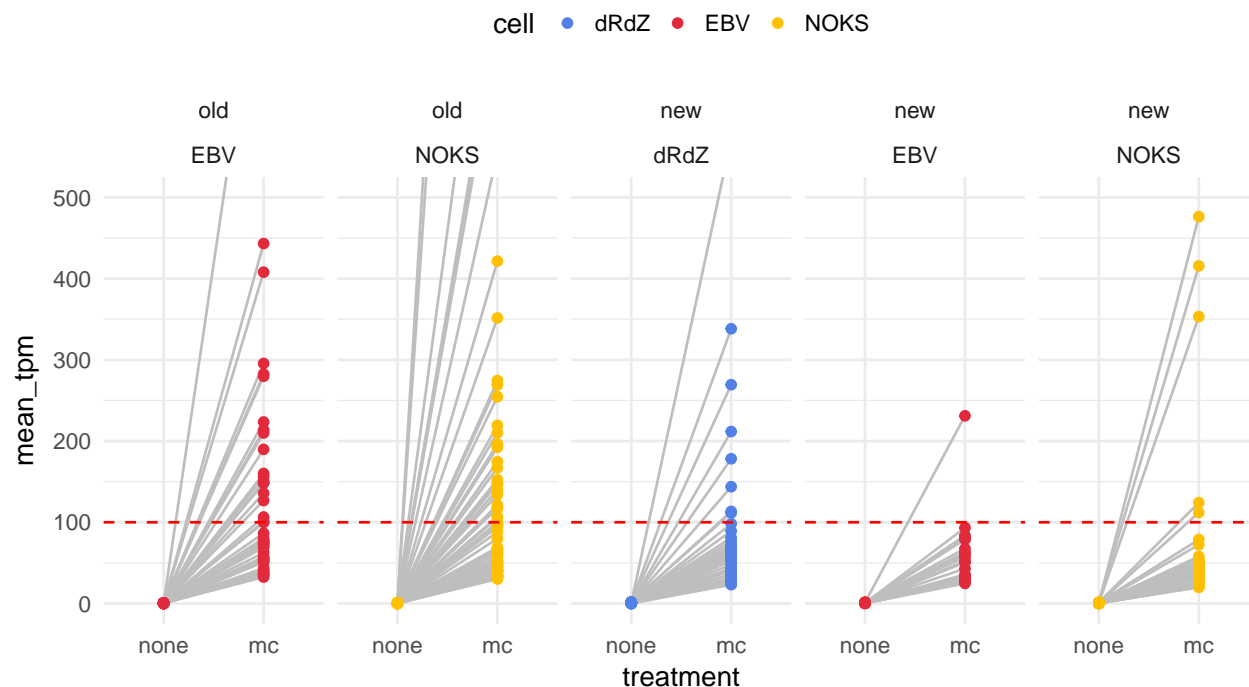
## TPM summaries

Both data batches are grouped, so we are going to rank the genes by the mean TPM.

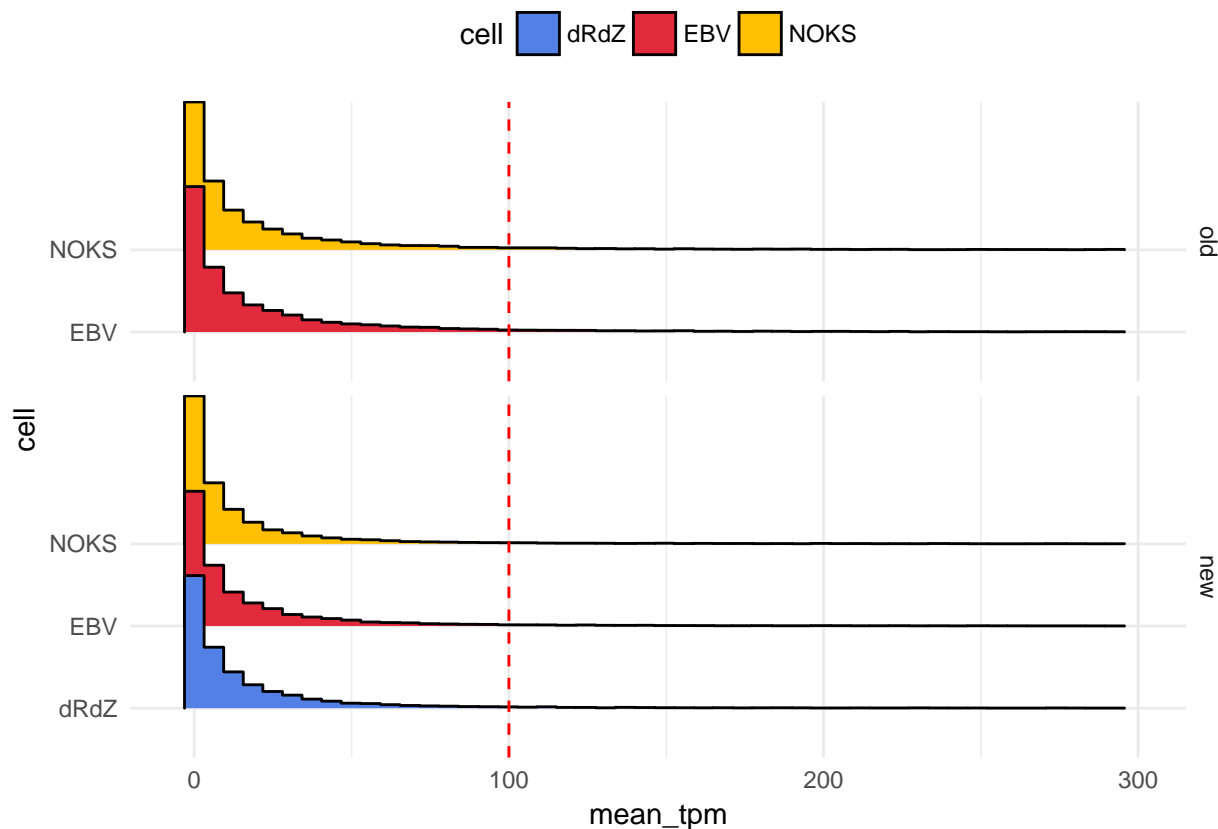


We further explore the genes in the red-square there are the genes that go from being unexpressed without treatment to be very expressed with treatment. In total, there are:

batch	cell	genes_in_square
new	dRdZ	53
new	EBV	25
new	NOKS	47
old	EBV	41
old	NOKS	73



In the figure above, it is shown that as expected when no treatment applied, those genes are not expressed and very likely there is no signal if we observe the tracks. On the other hand, the average tpm for the genes after applying the mc treatment, we can notice that there is some expression. To provide context, we observe the mean TPM distribution:



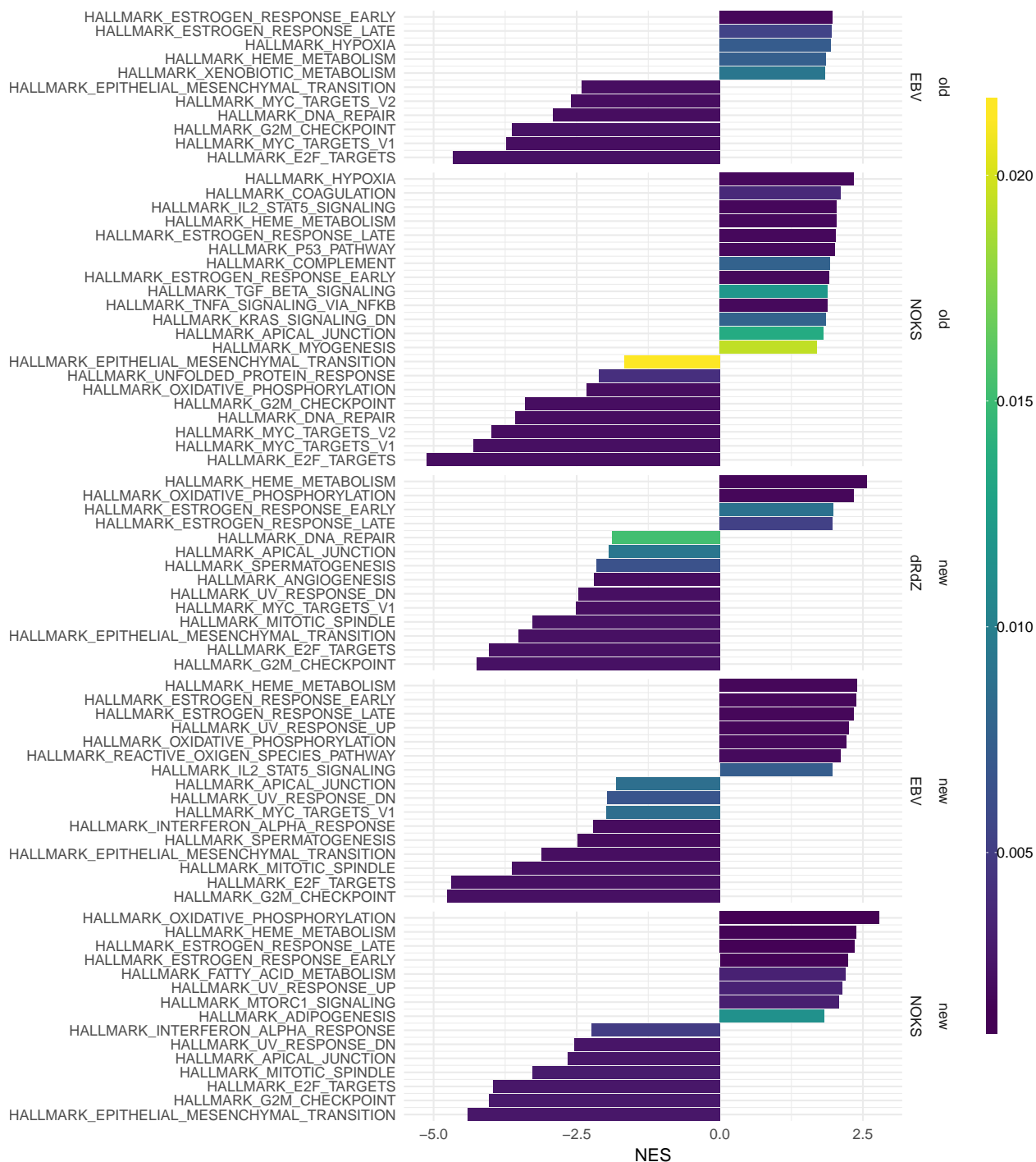
## Filtering out genes

We filtered out genes for the GSEA analysis, by considering only the genes inside the orange lines that are differentially expressed (defined as genes with  $\text{adj. p.value} \leq 0.01$ ):

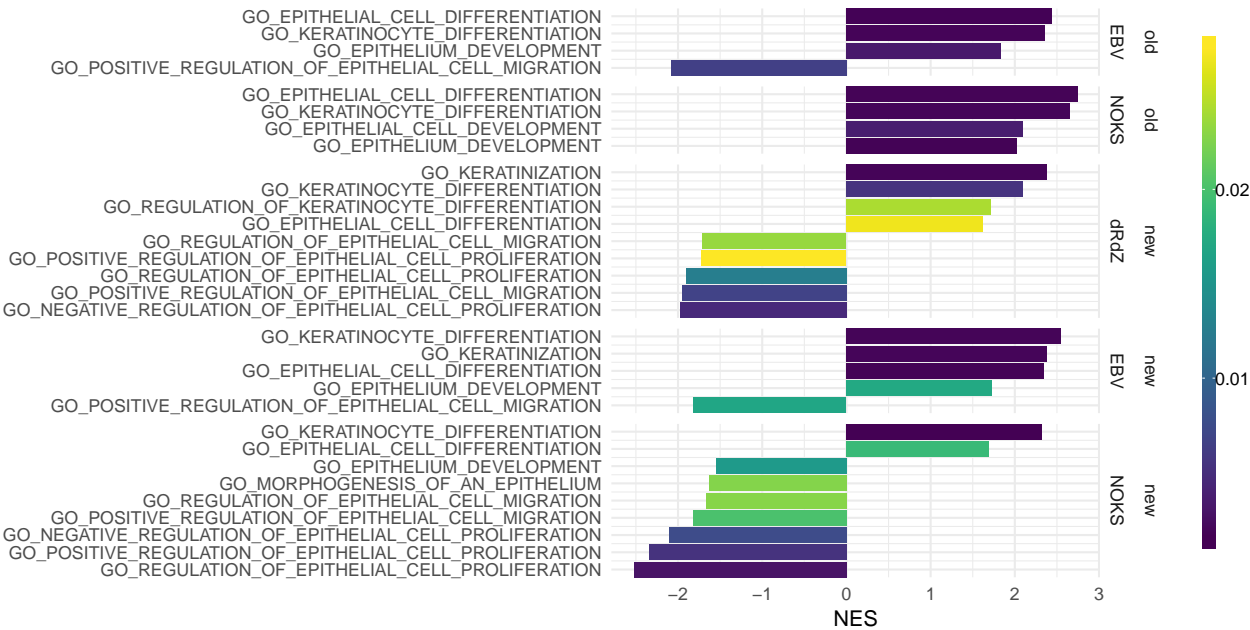
batch	cell	total_genes	diff_genes
old	EBV	16057	2302
old	NOKS	16057	4230
new	EBV	15876	1624
new	NOKS	15876	2445
new	dRdZ	15876	2025

# GSEA analysis

## Hallmark analysis

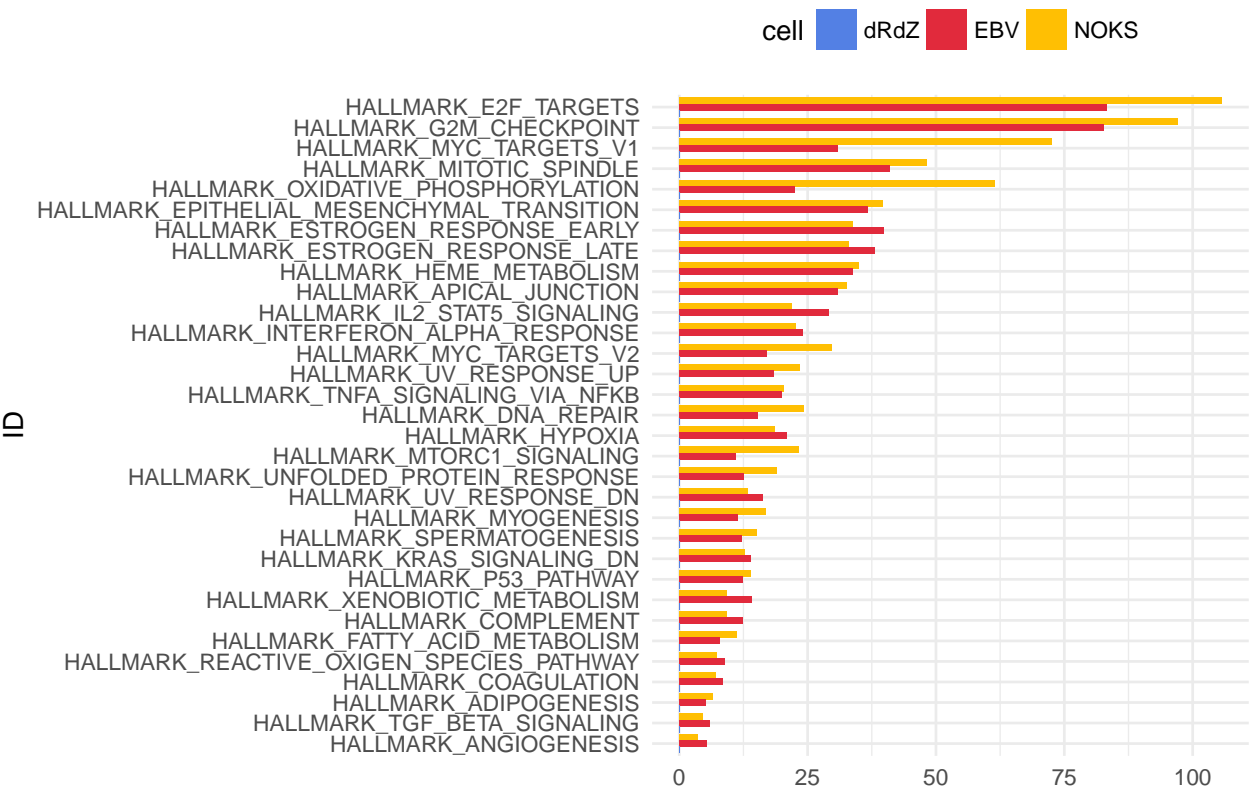


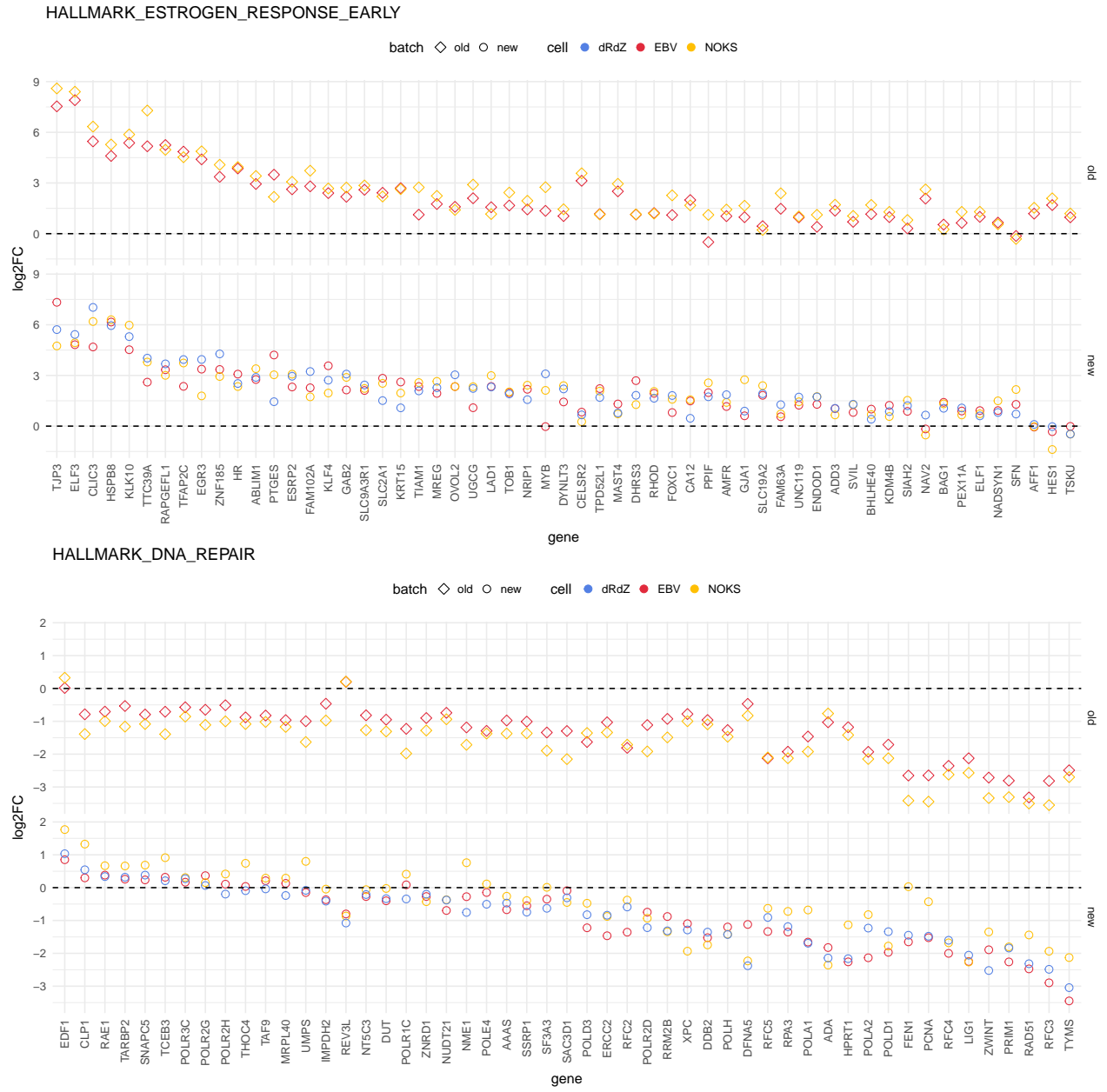
Curated analysis



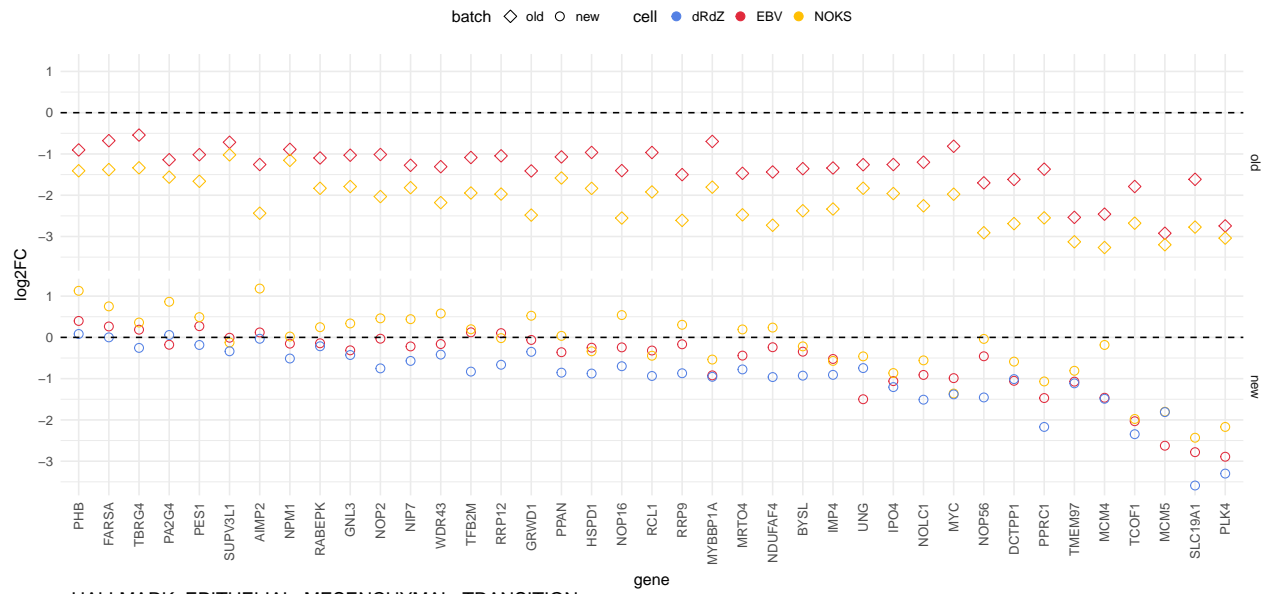
Further analysis

Hallmark pathways

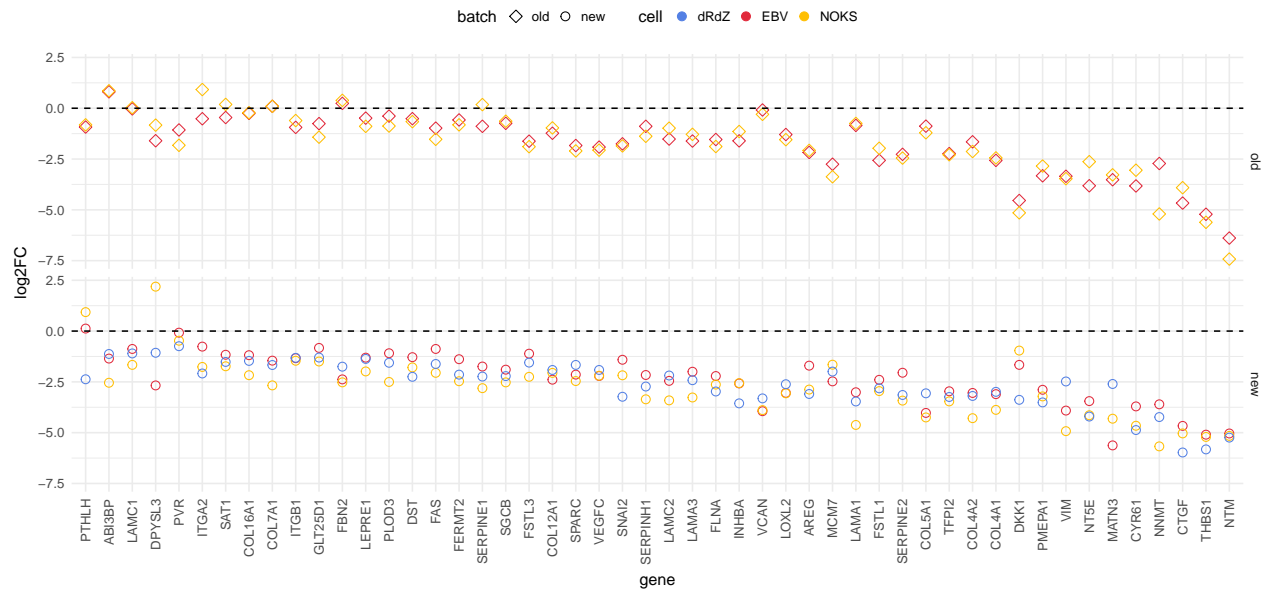




# HALLMARK\_MYC\_TARGETS\_V2

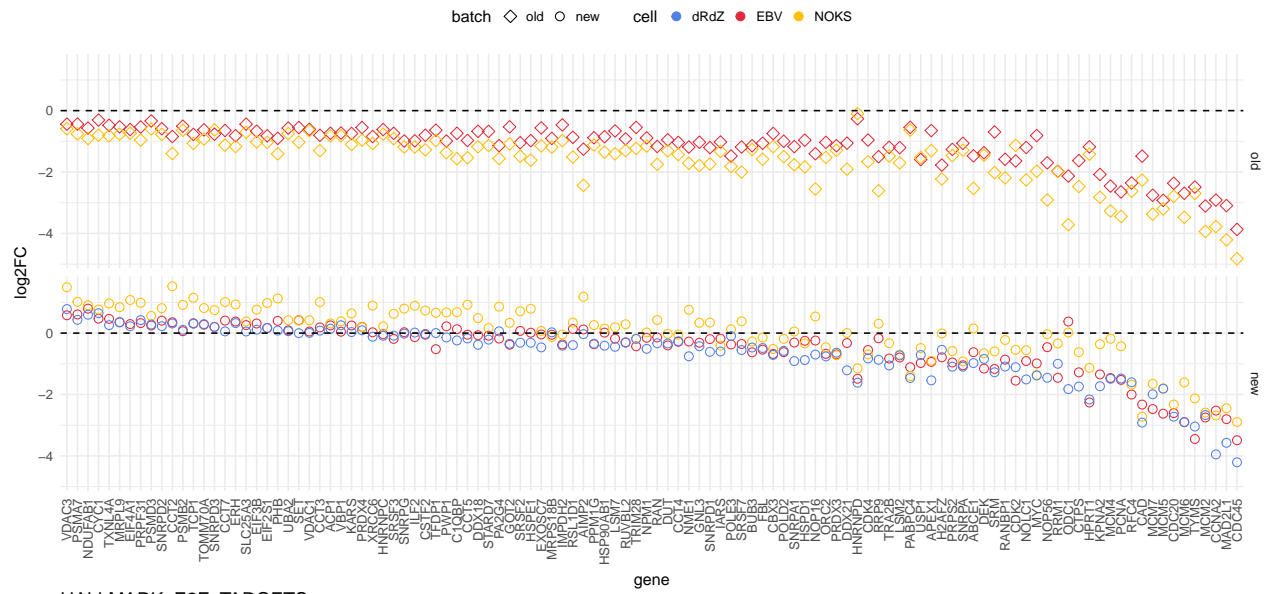


# HALLMARK\_EPITHELIAL\_MESENCHYMAL\_TRANSITION

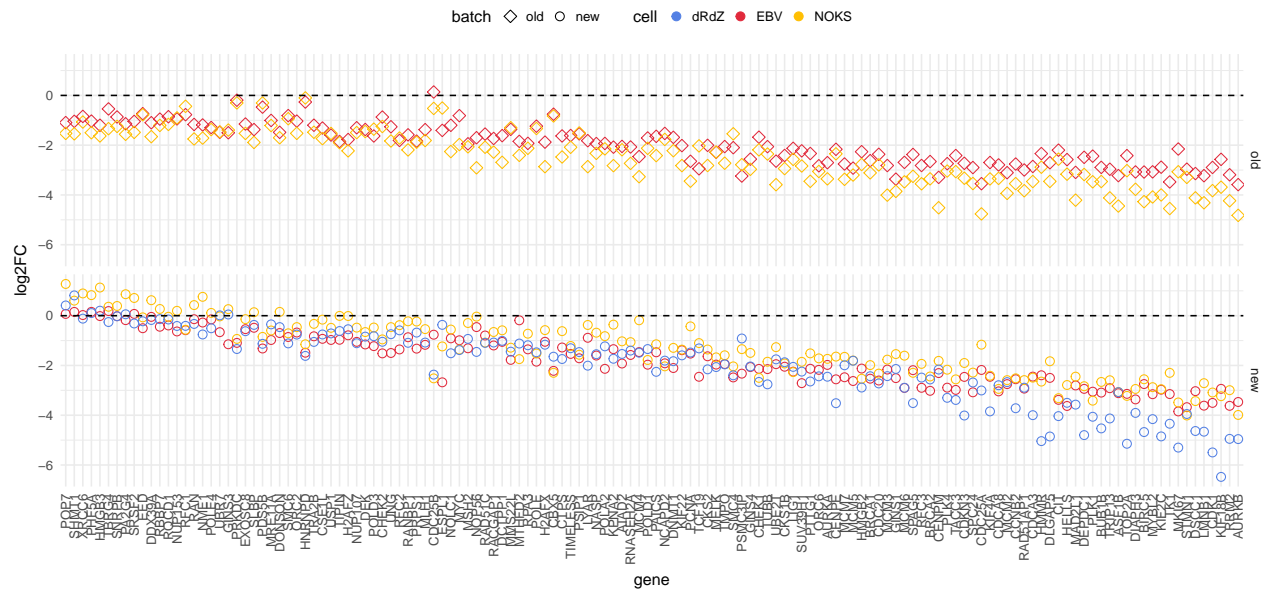




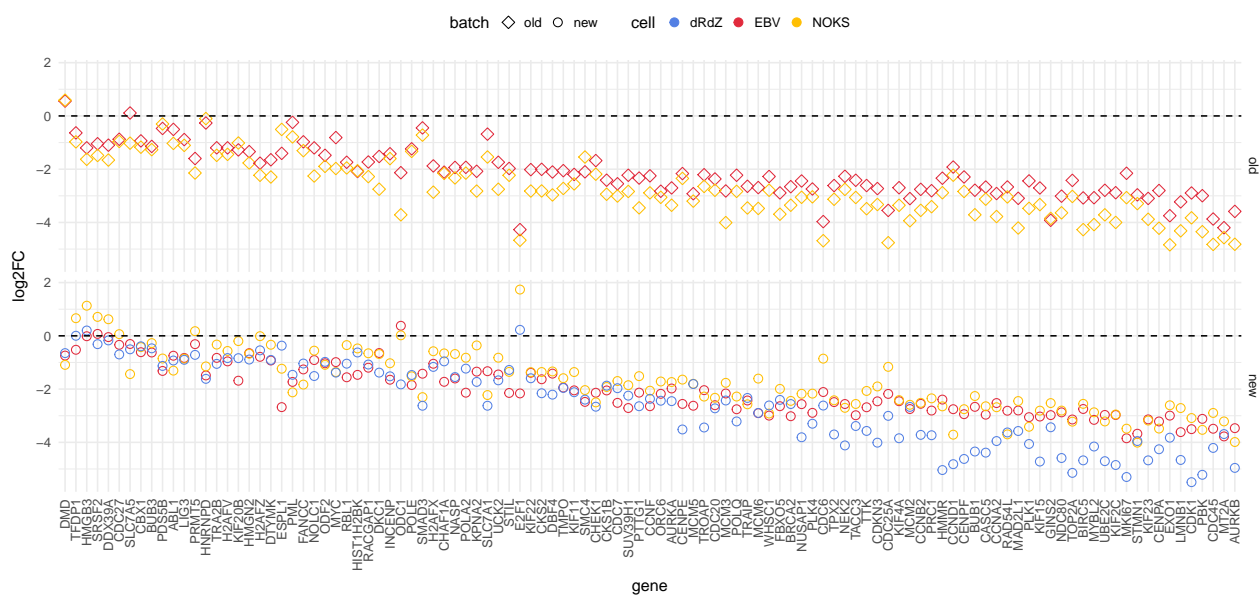
### HALLMARK\_MYC\_TARGETS\_V1



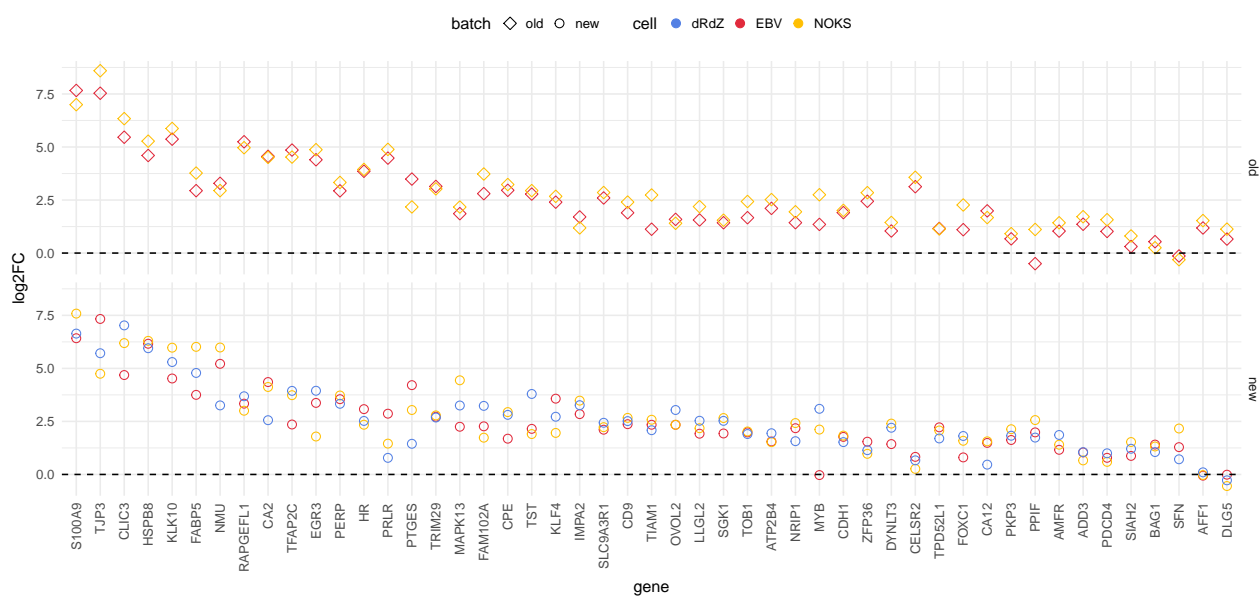
### HALLMARK\_E2F\_TARGETS

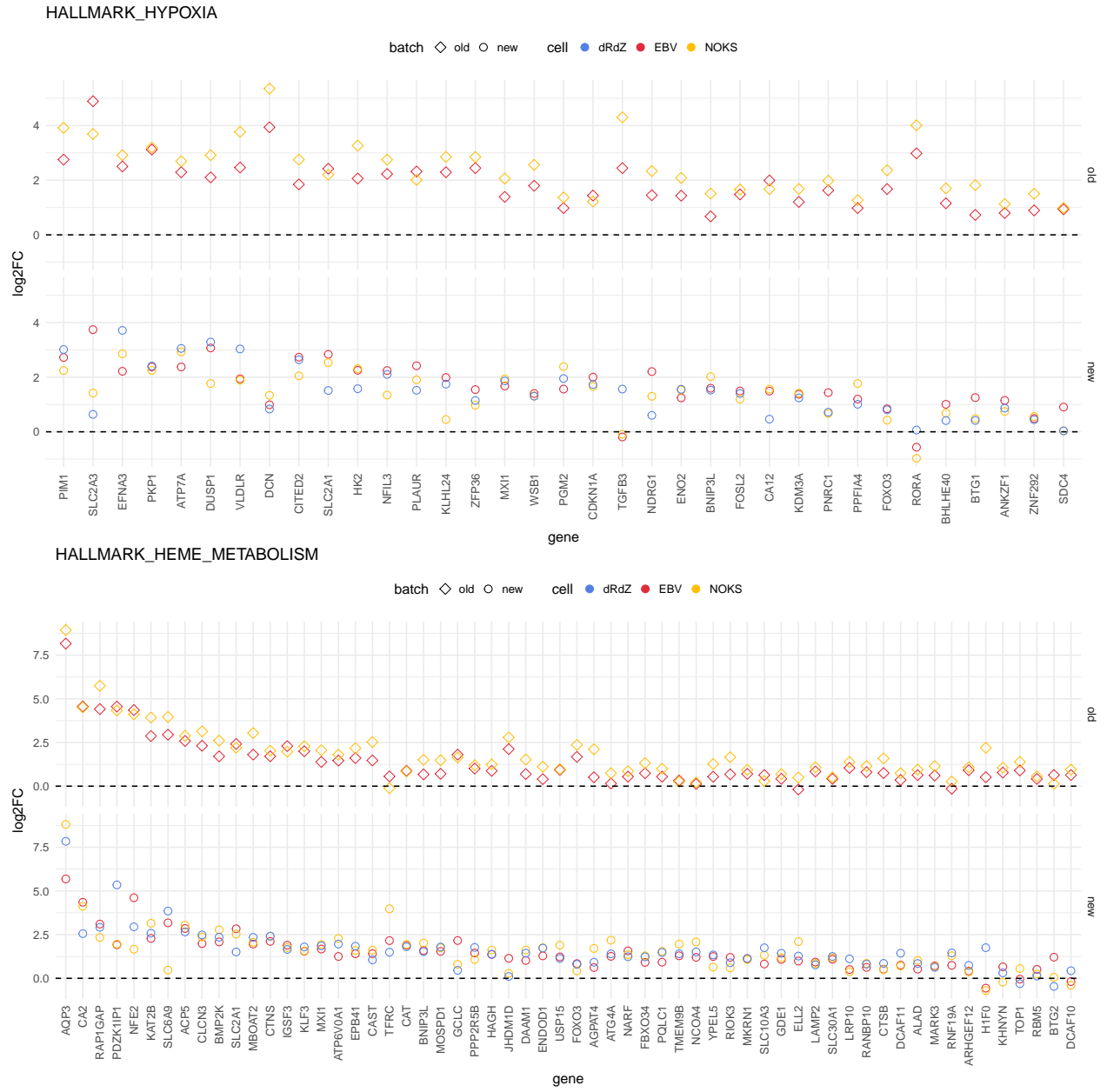


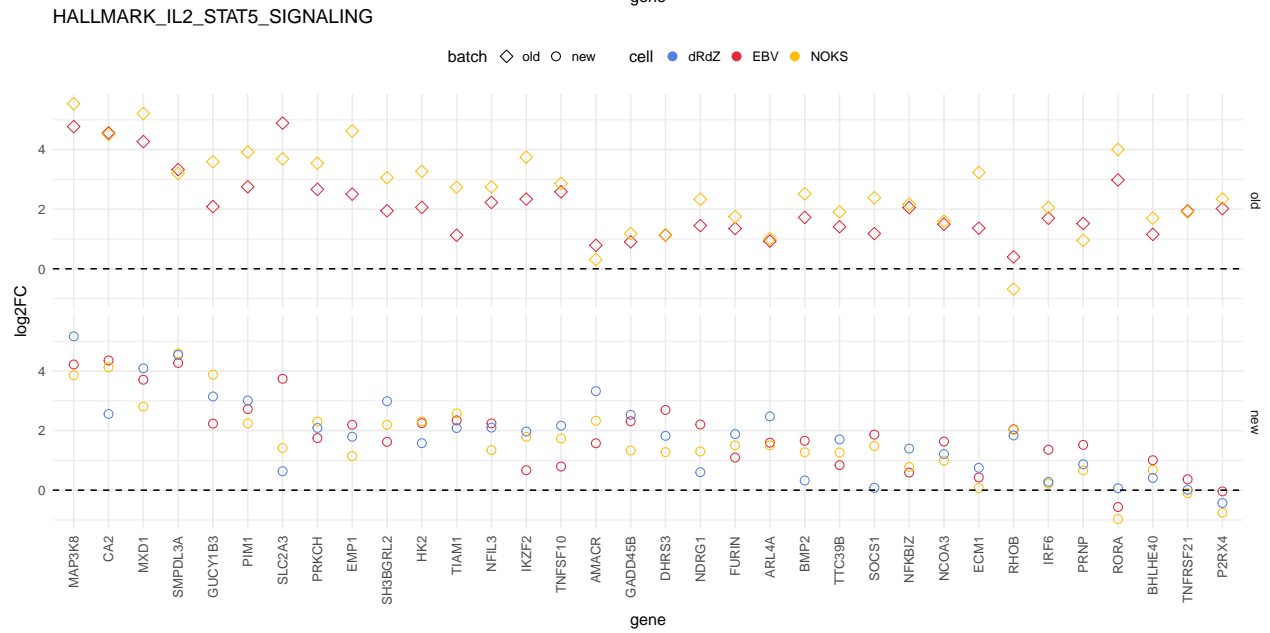
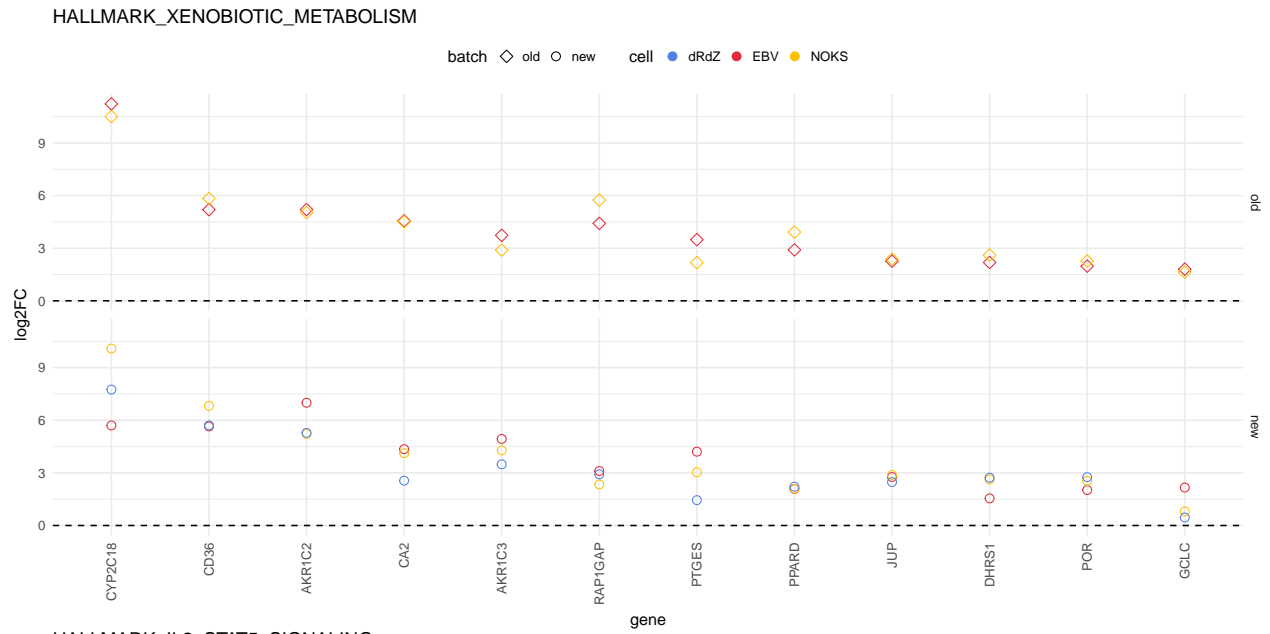
## HALLMARK\_G2M\_CHECKPOINT



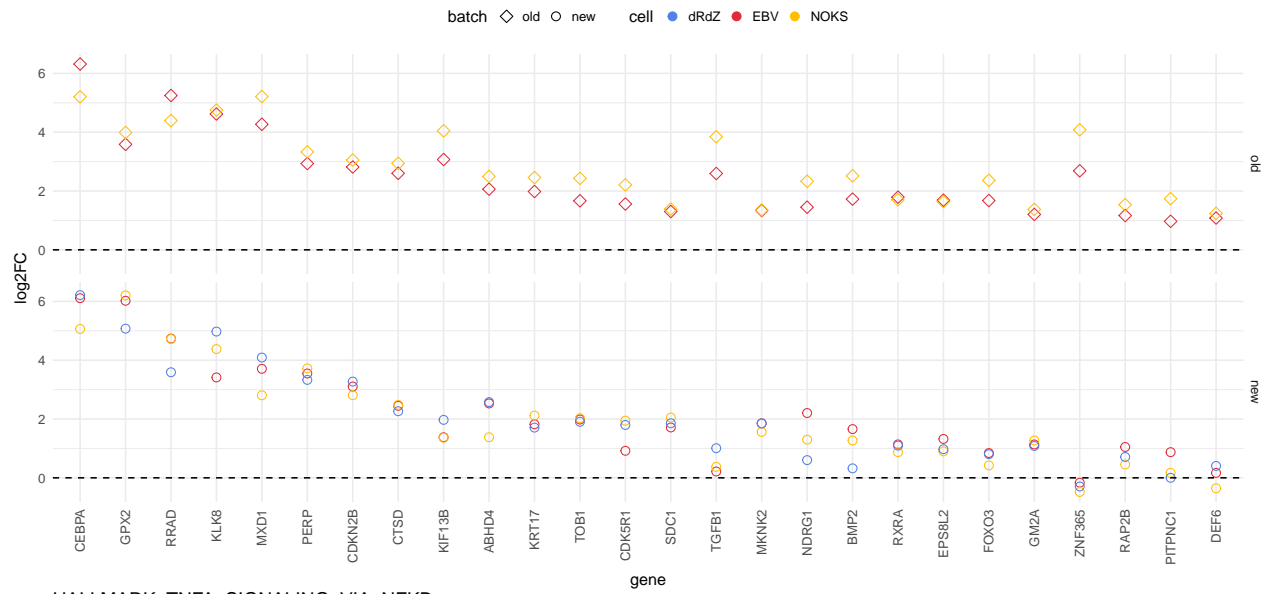
## HALLMARK\_ESTROGEN\_RESPONSE\_LATE



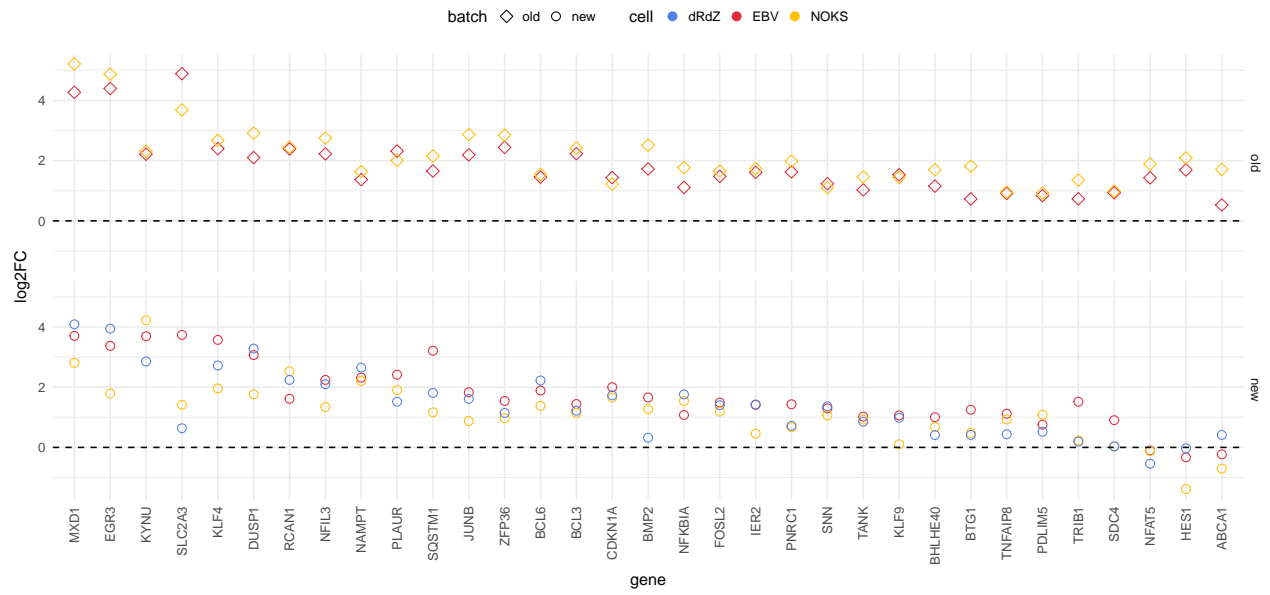




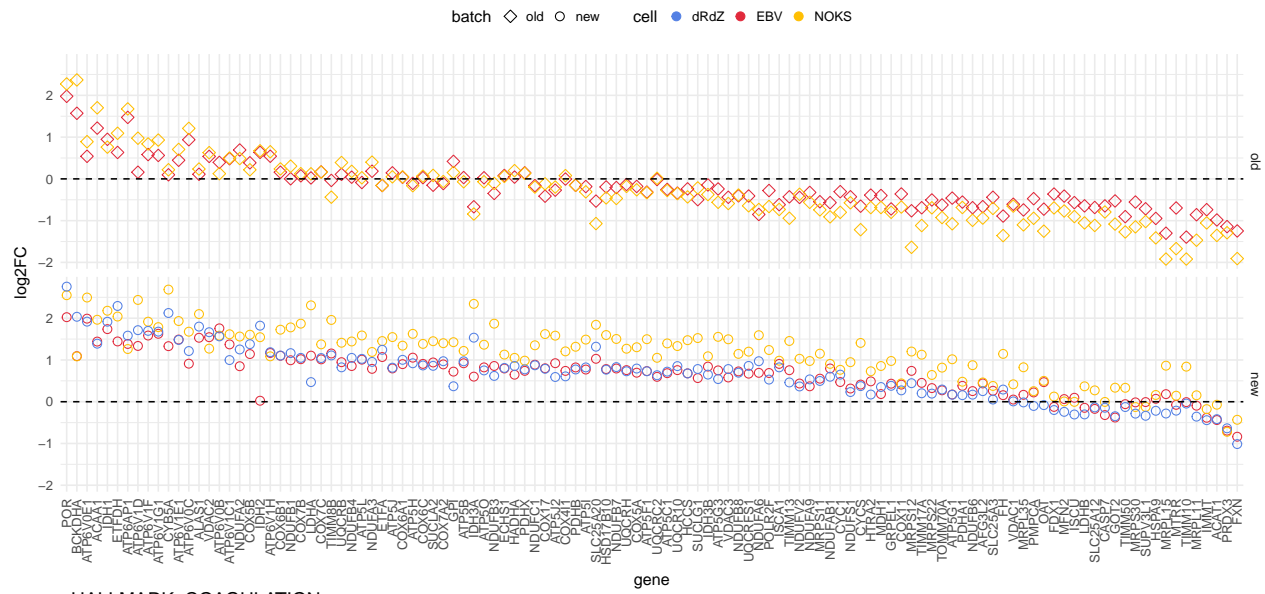
# HALLMARK\_P53\_PATHWAY



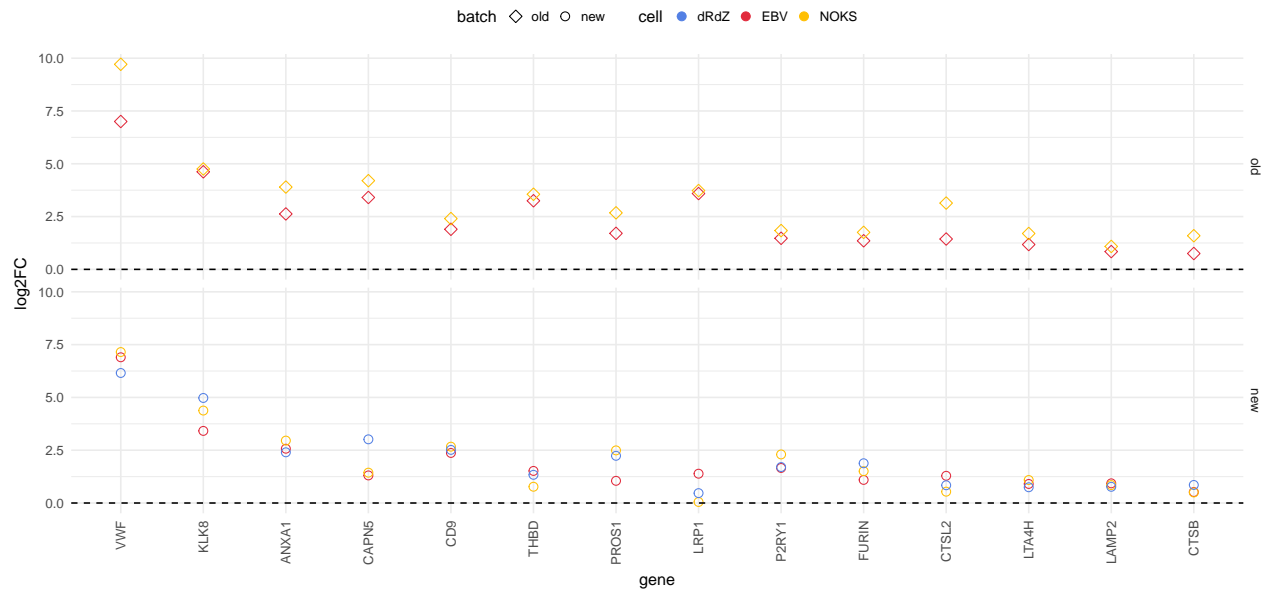
# HALLMARK\_TNFA\_SIGNALING\_VIA\_NFKB

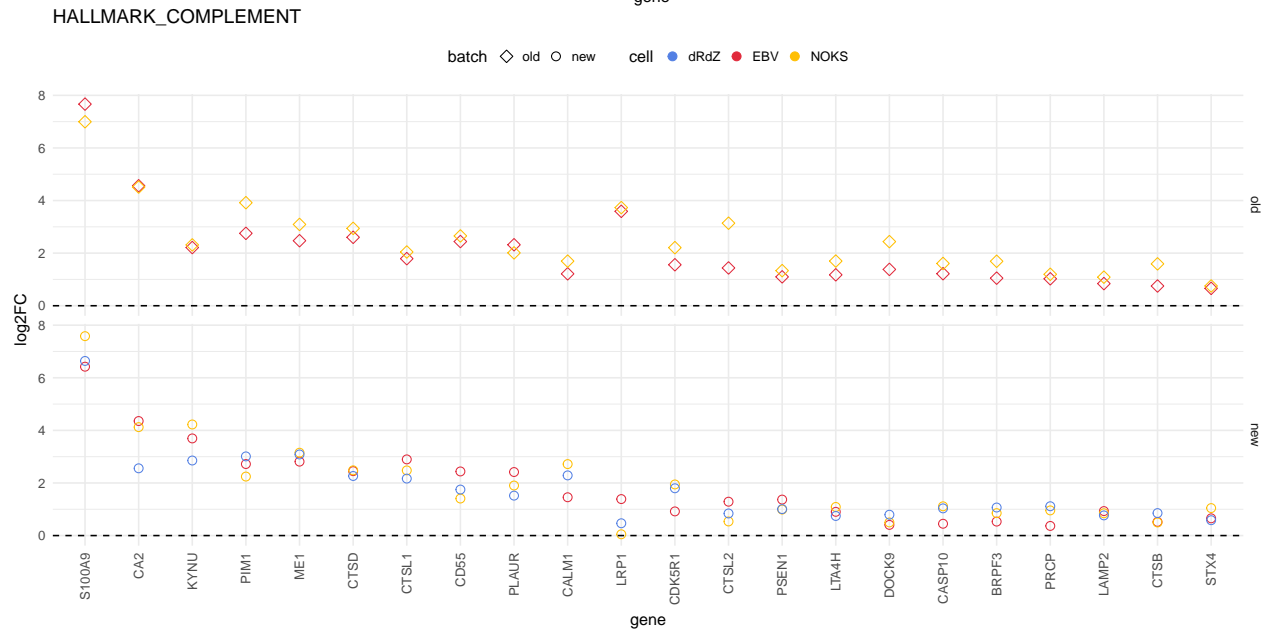
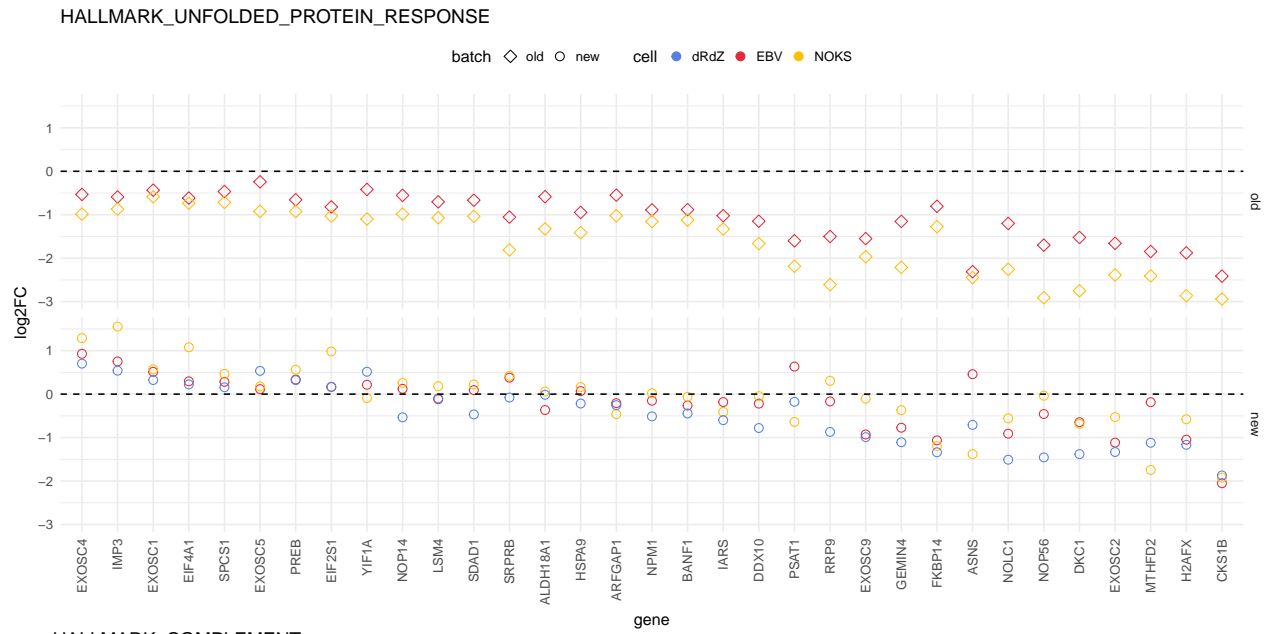


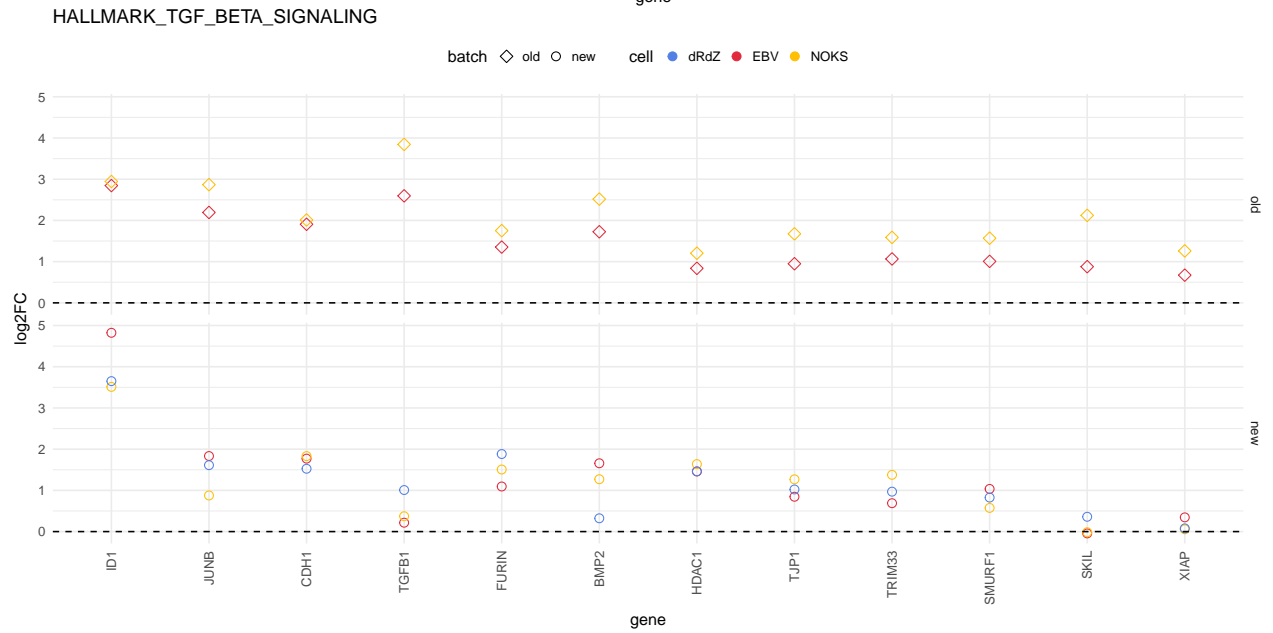
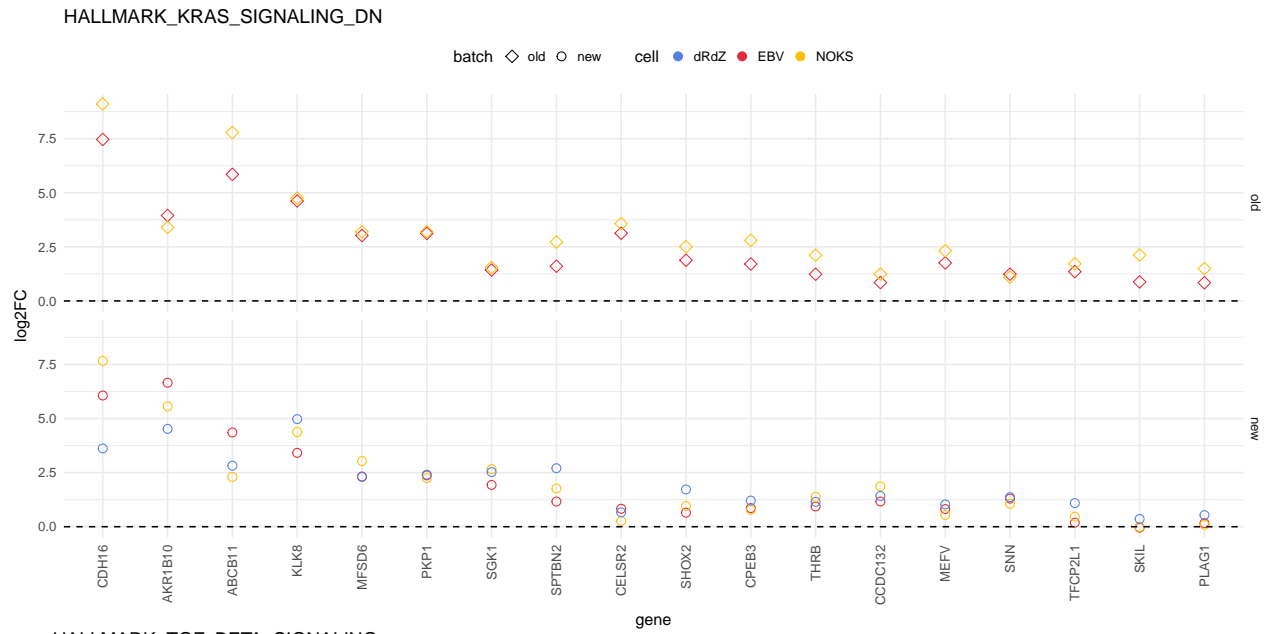
# HALLMARK\_OXIDATIVE\_PHOSPHORYLATION



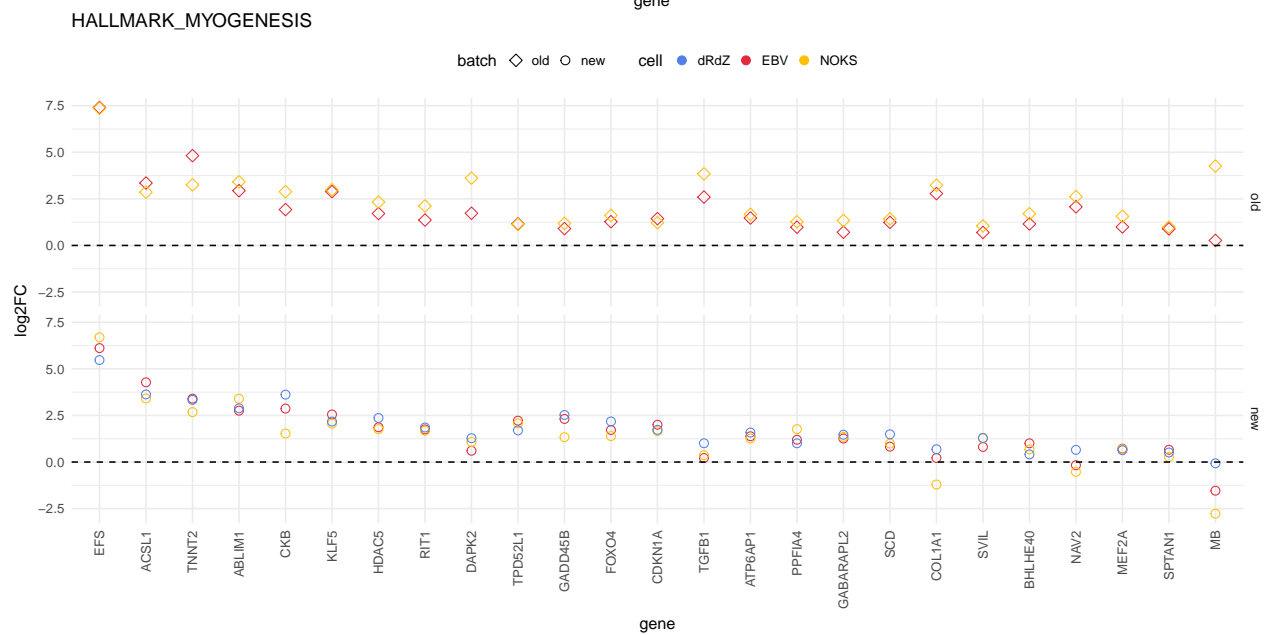
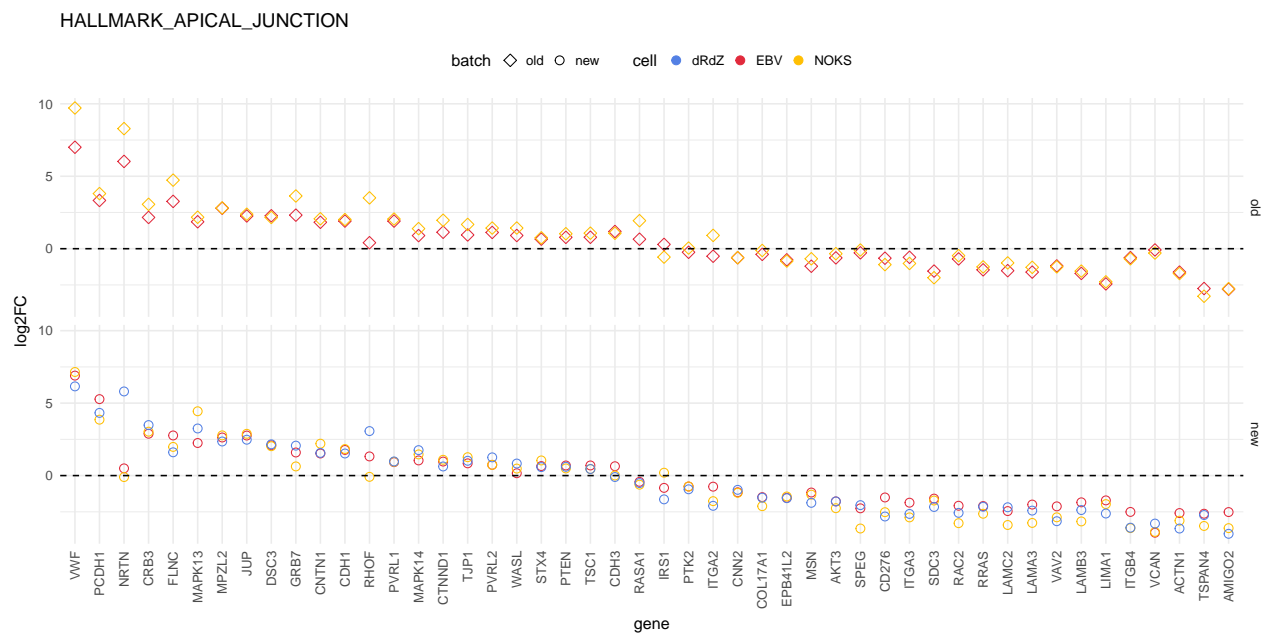
# HALLMARK\_COAGULATION

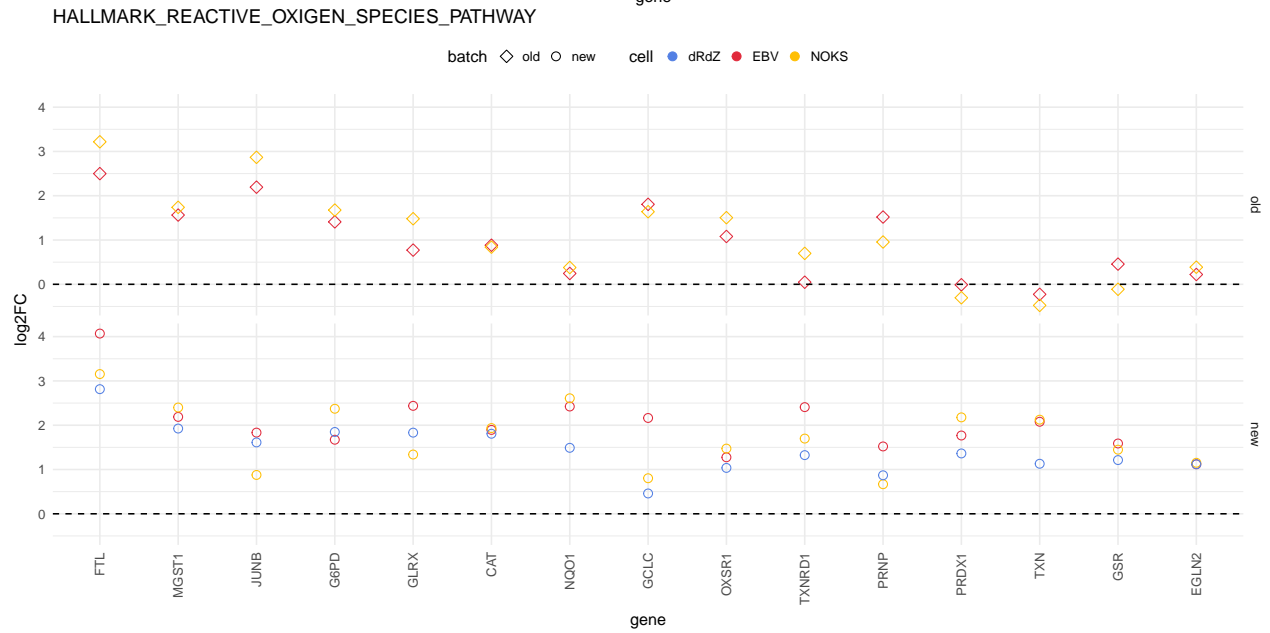
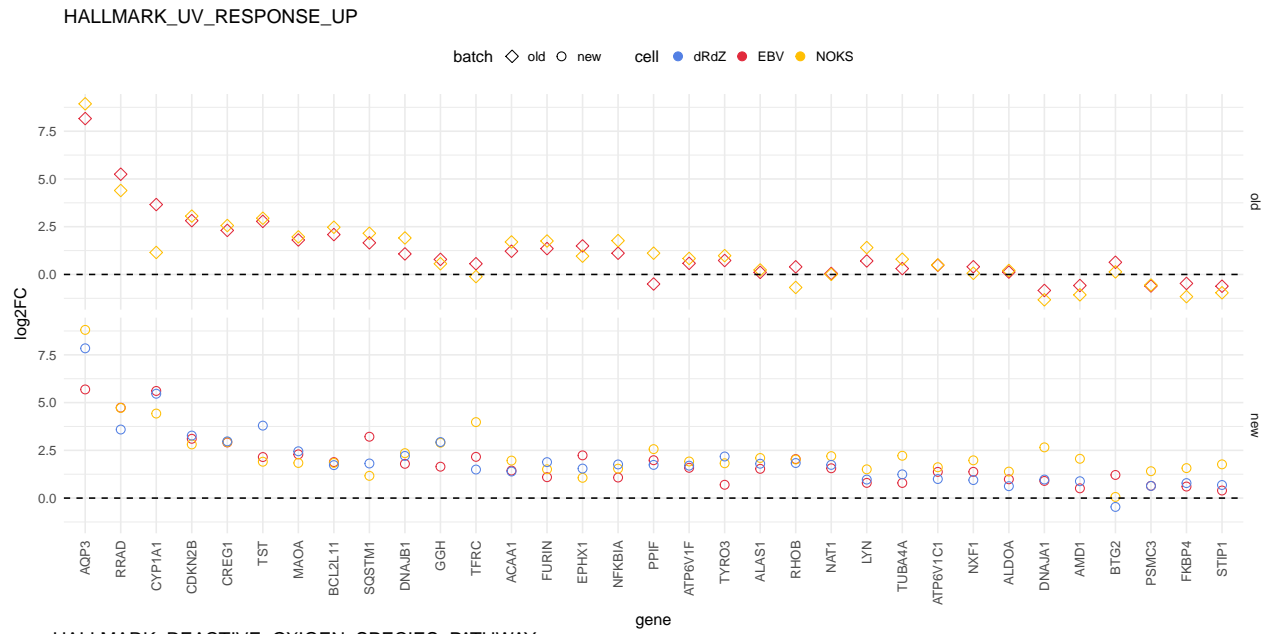


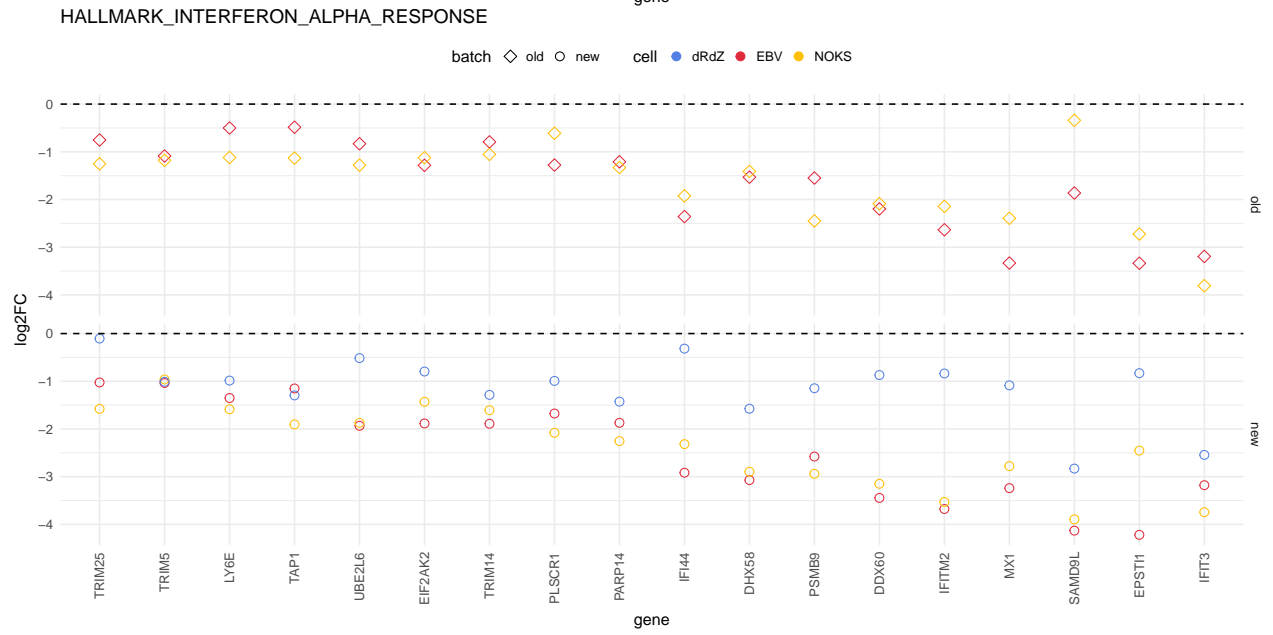
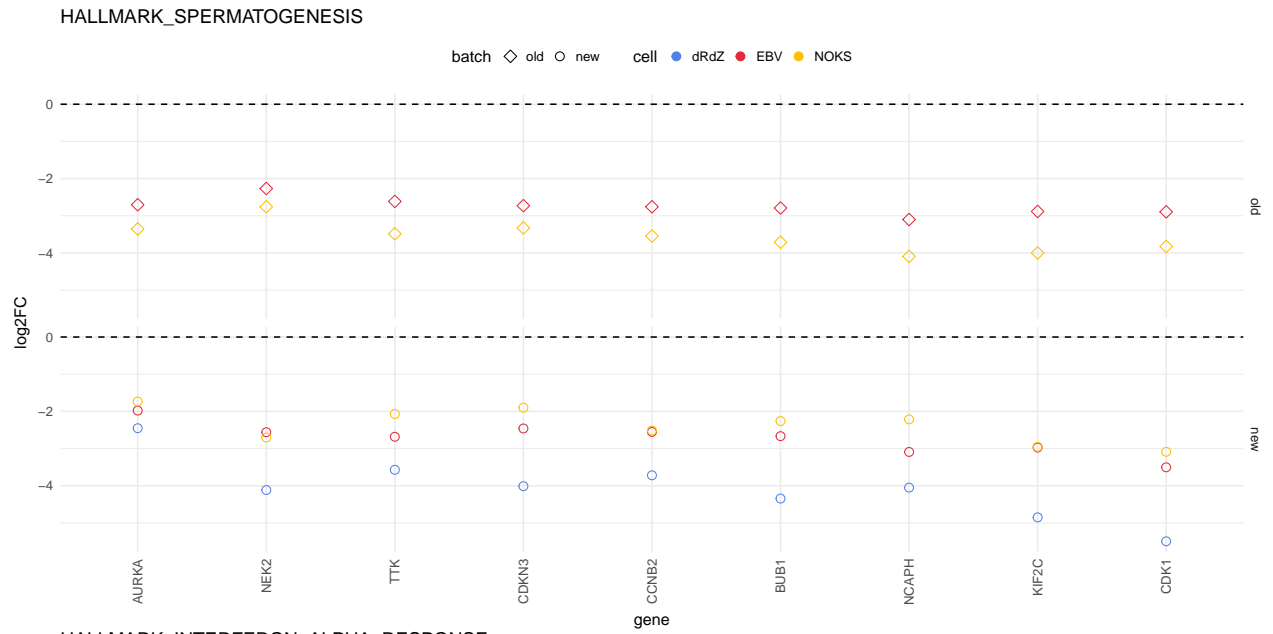


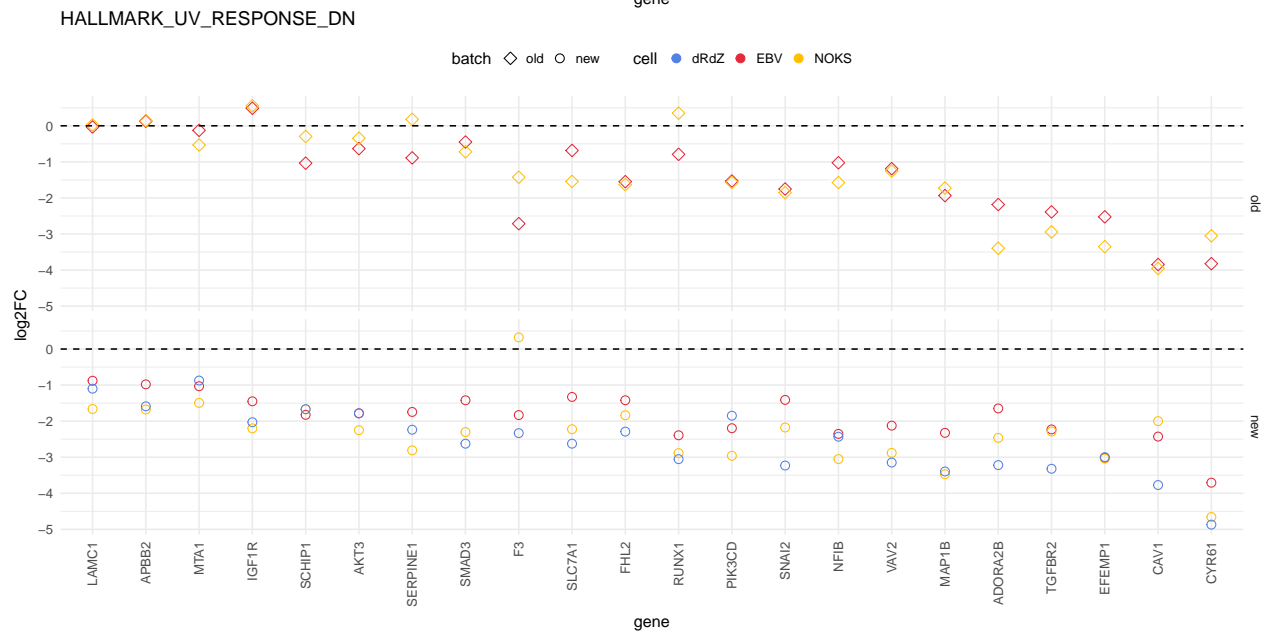
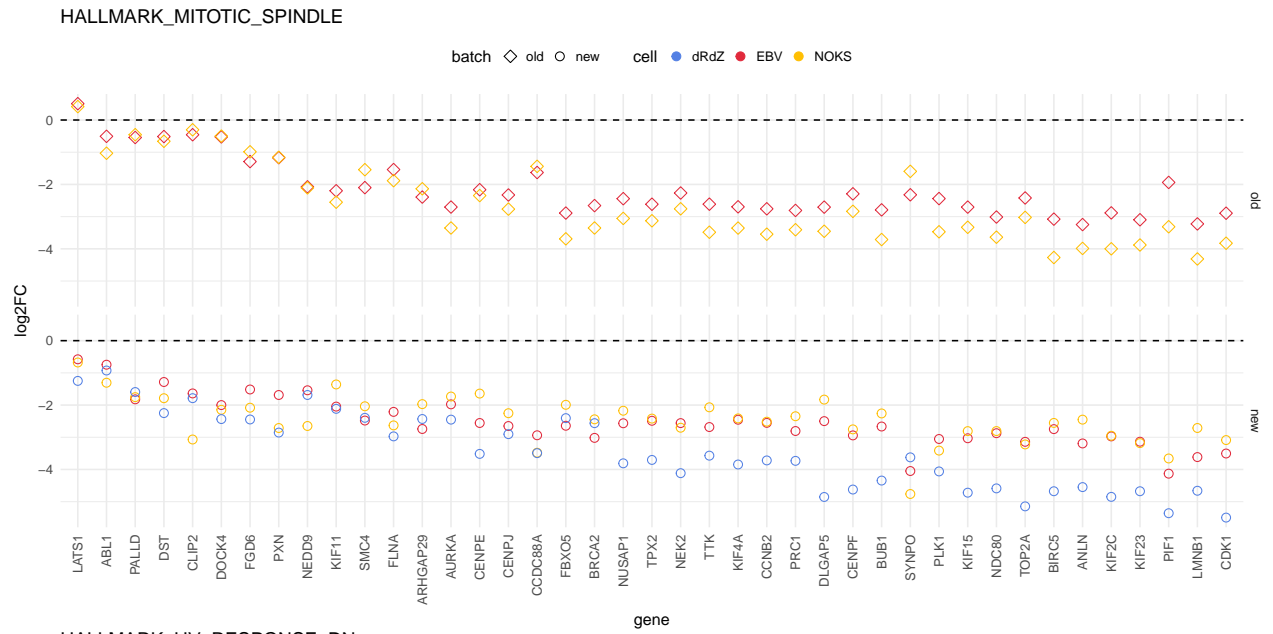


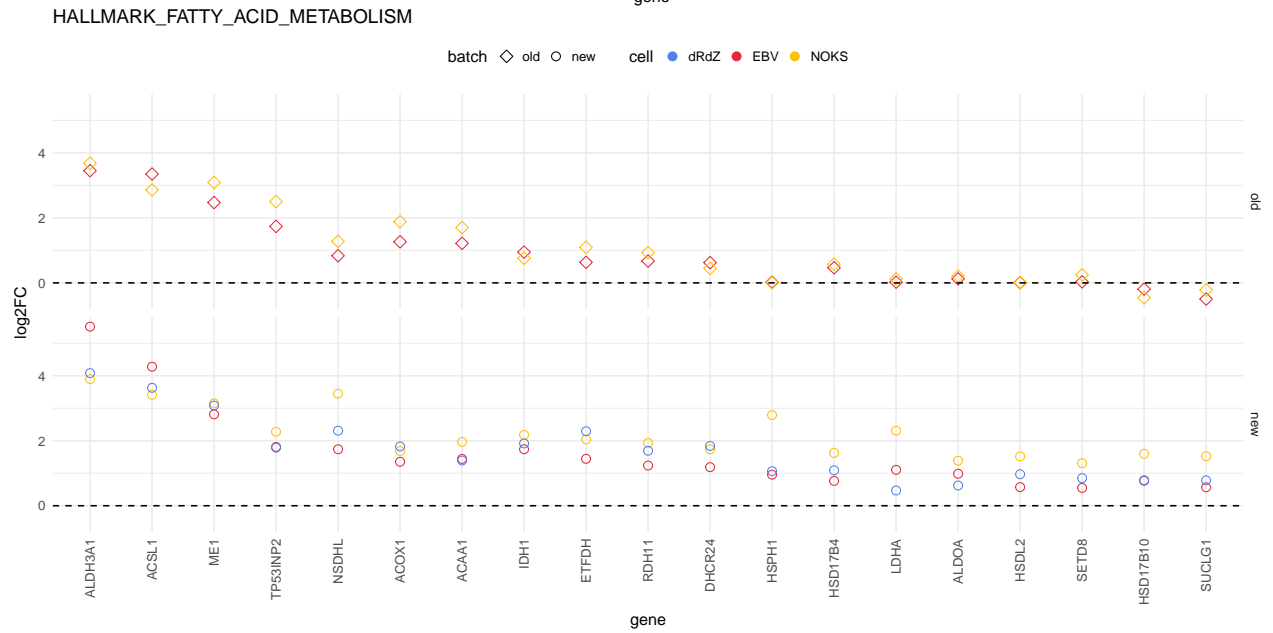
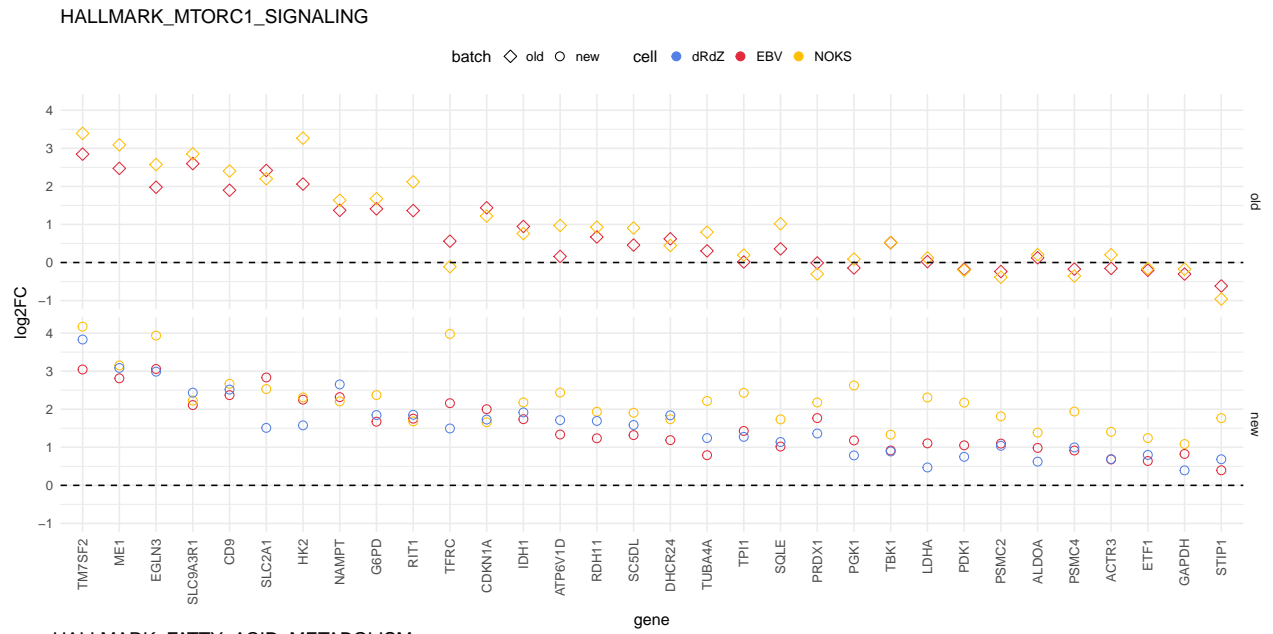




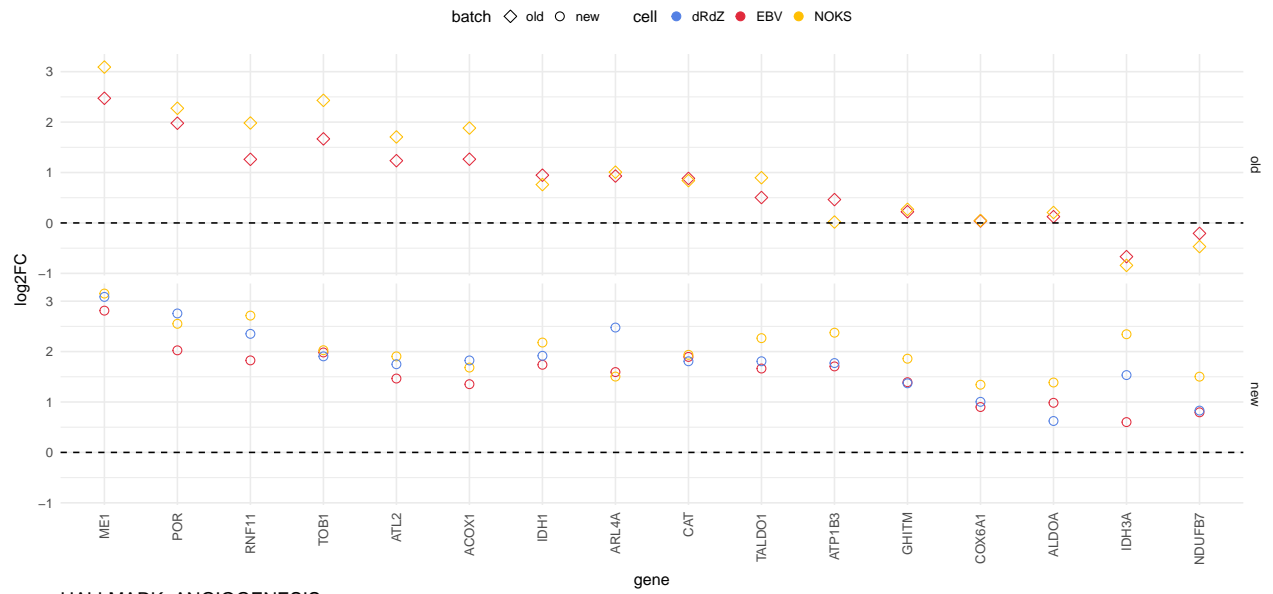




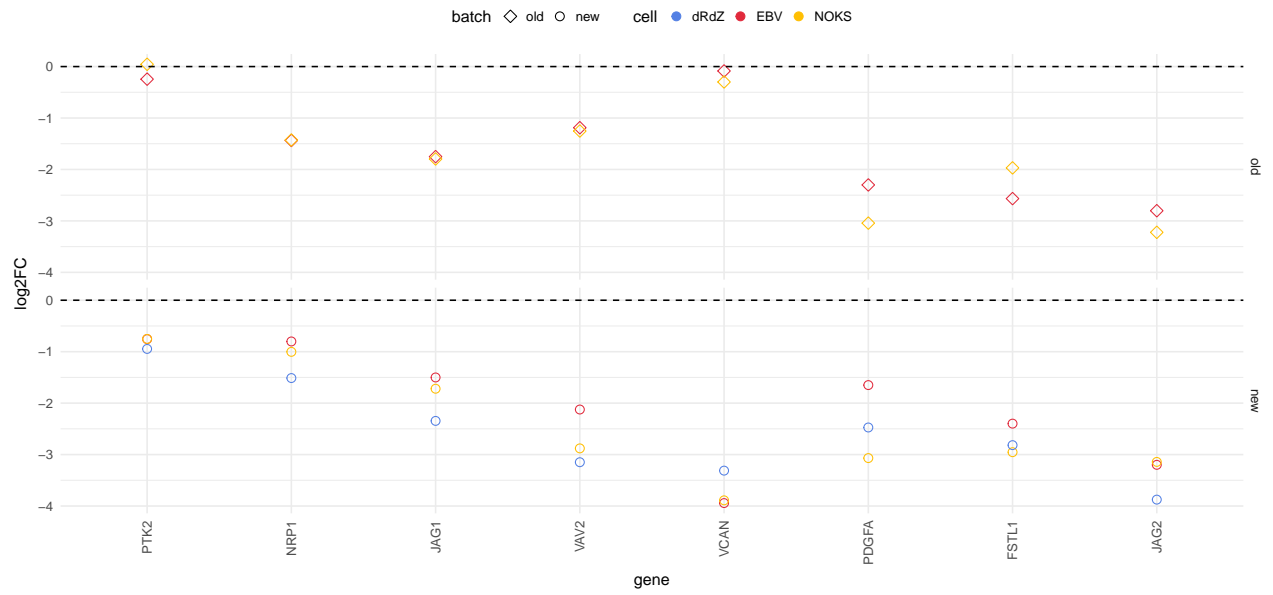




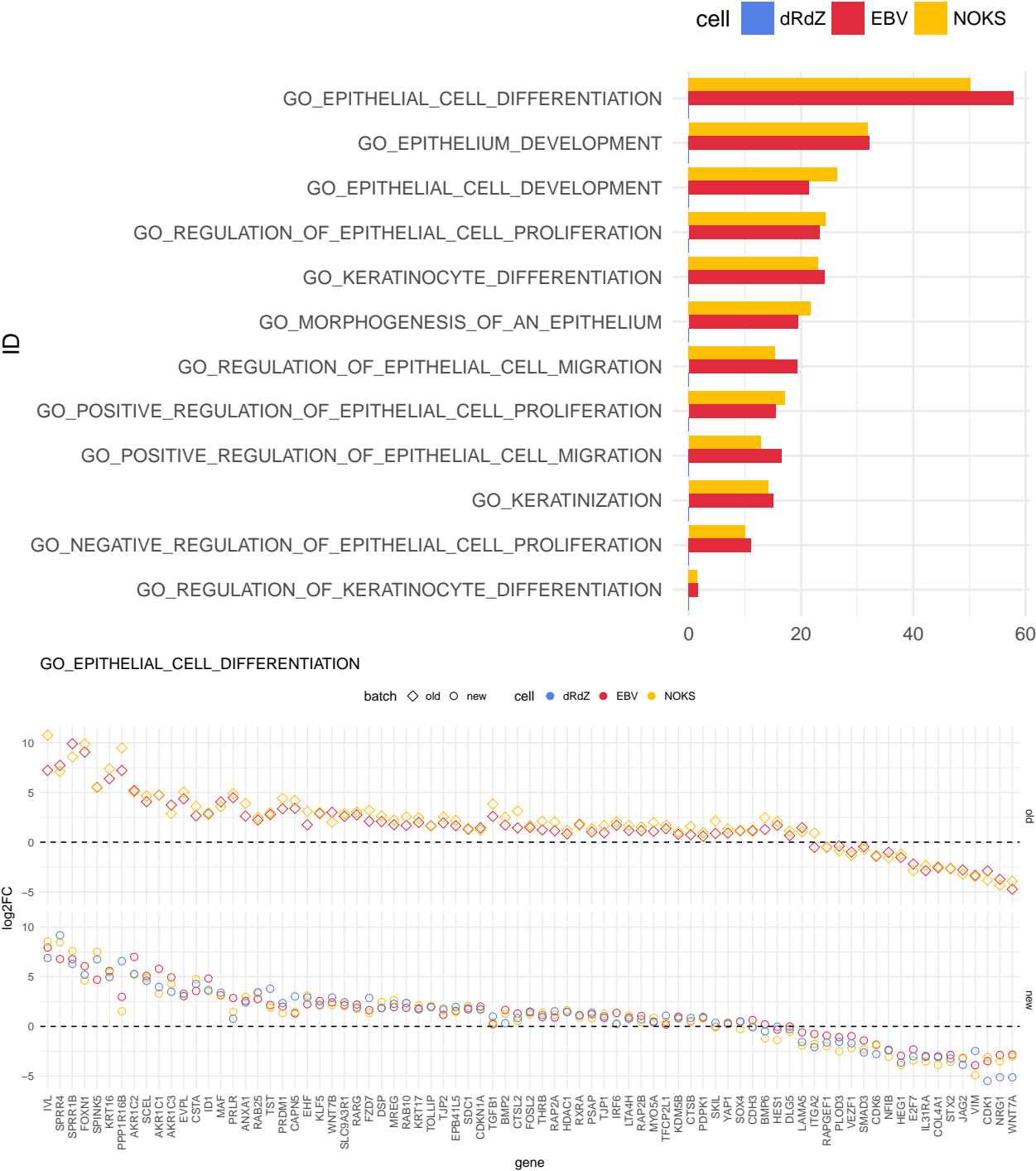
### HALLMARK\_ADIPOGENESIS

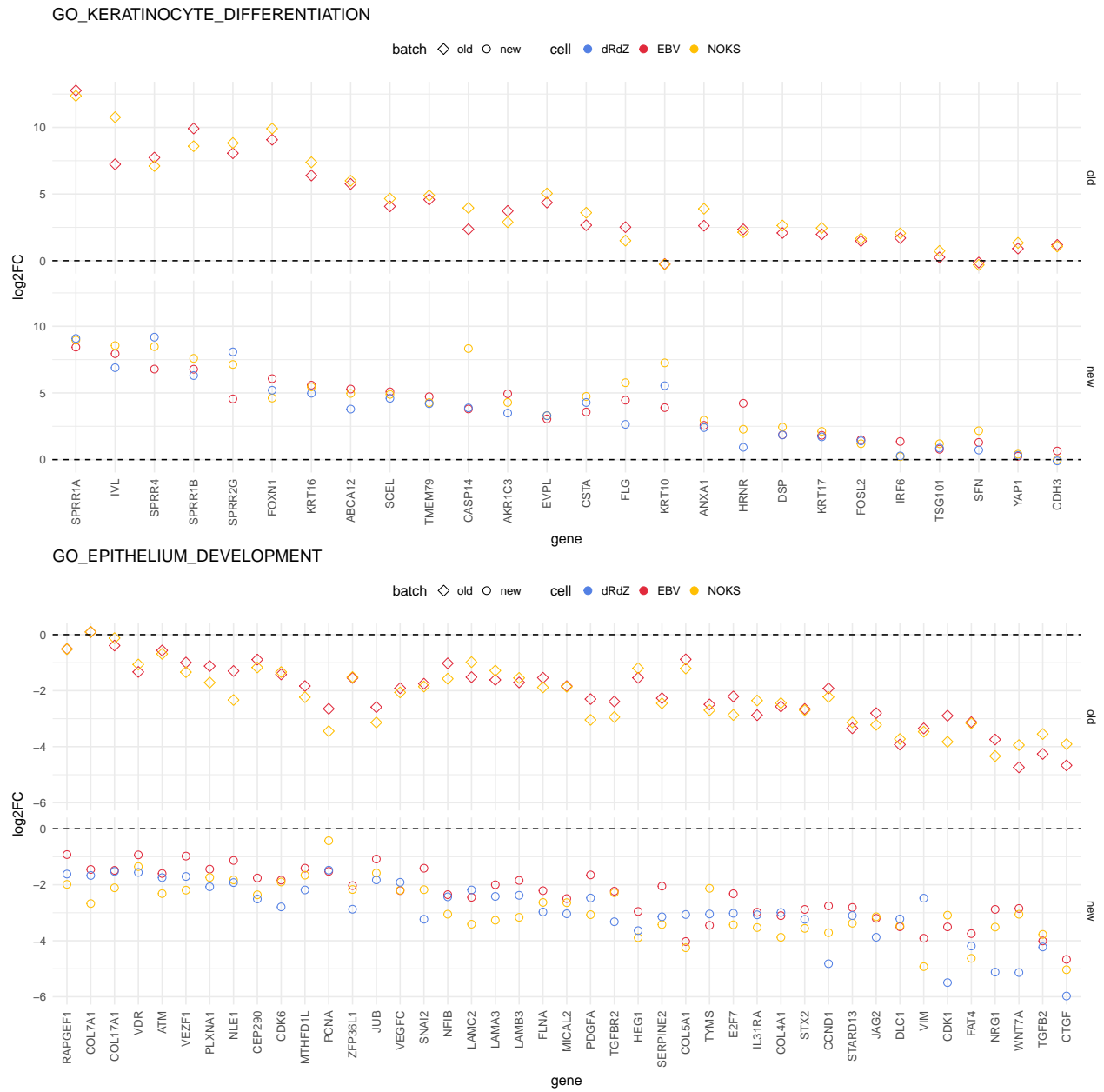


### HALLMARK\_ANGIOGENESIS

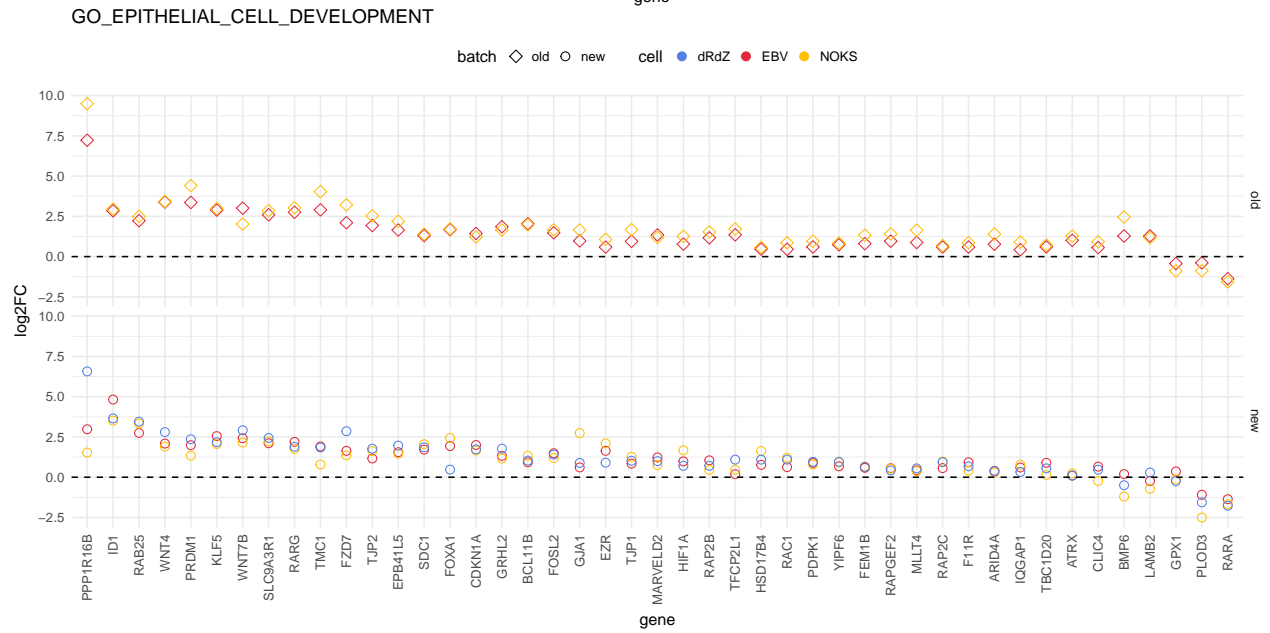
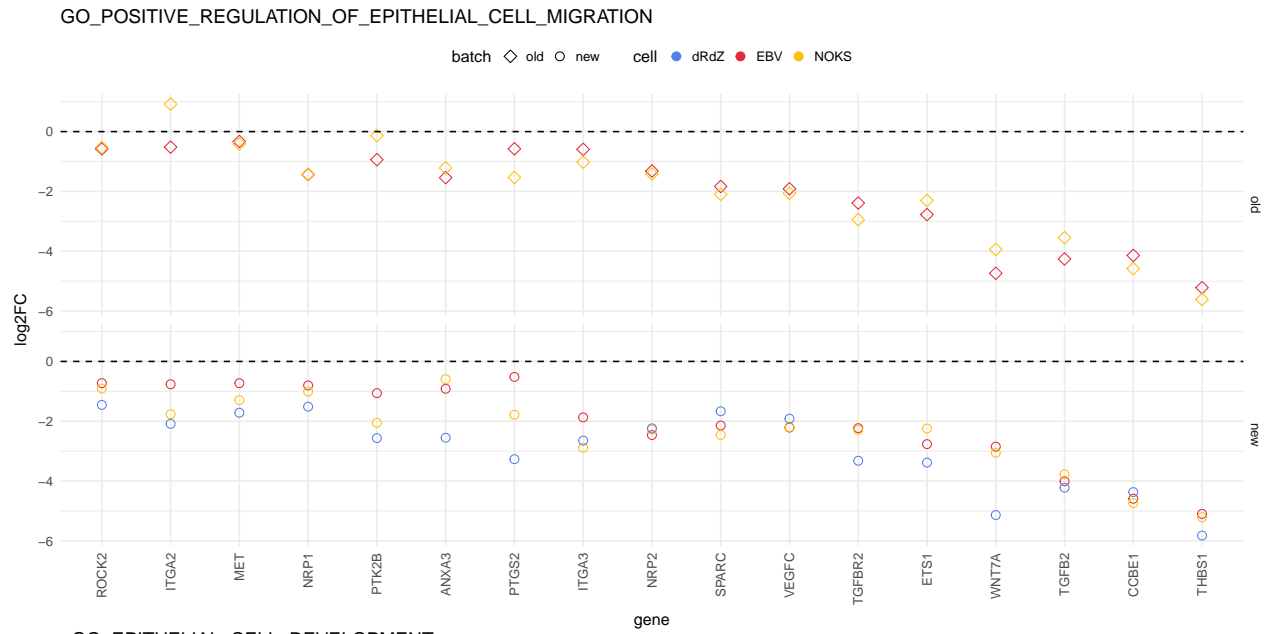


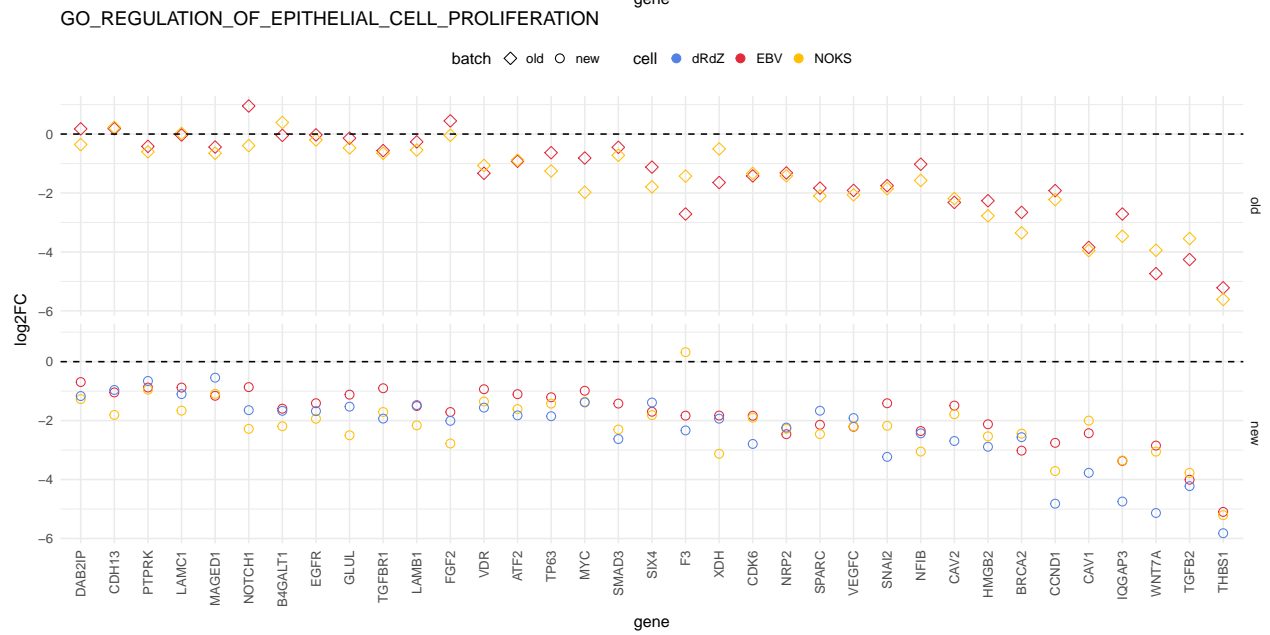
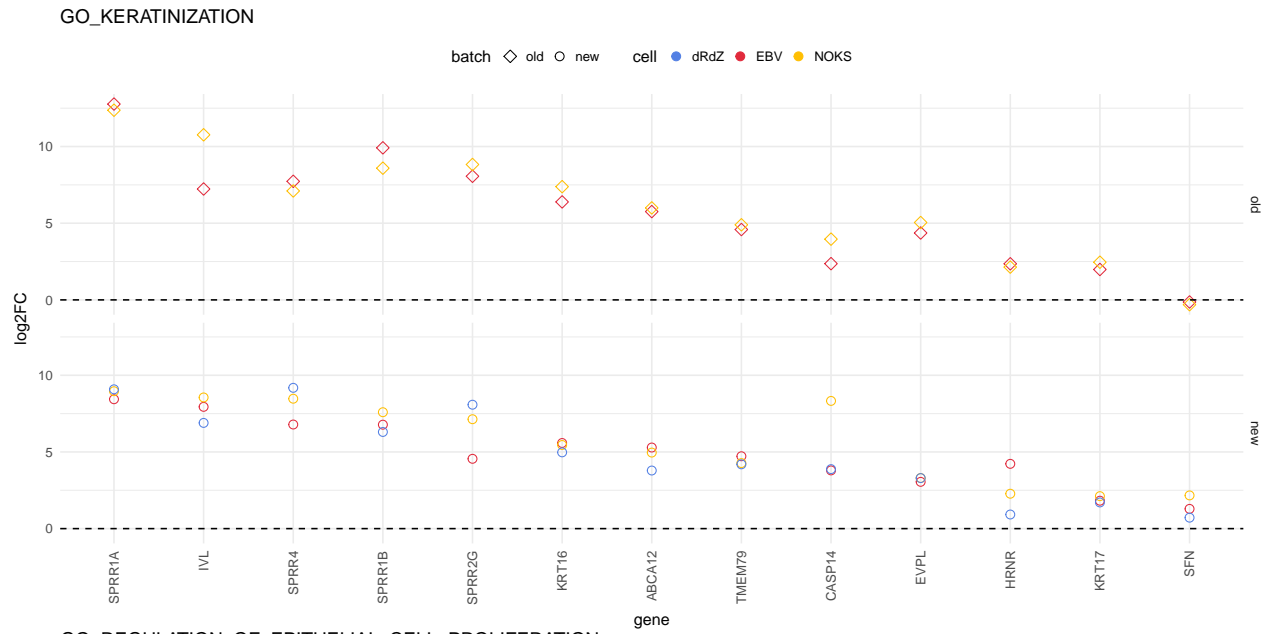
Curated pathways

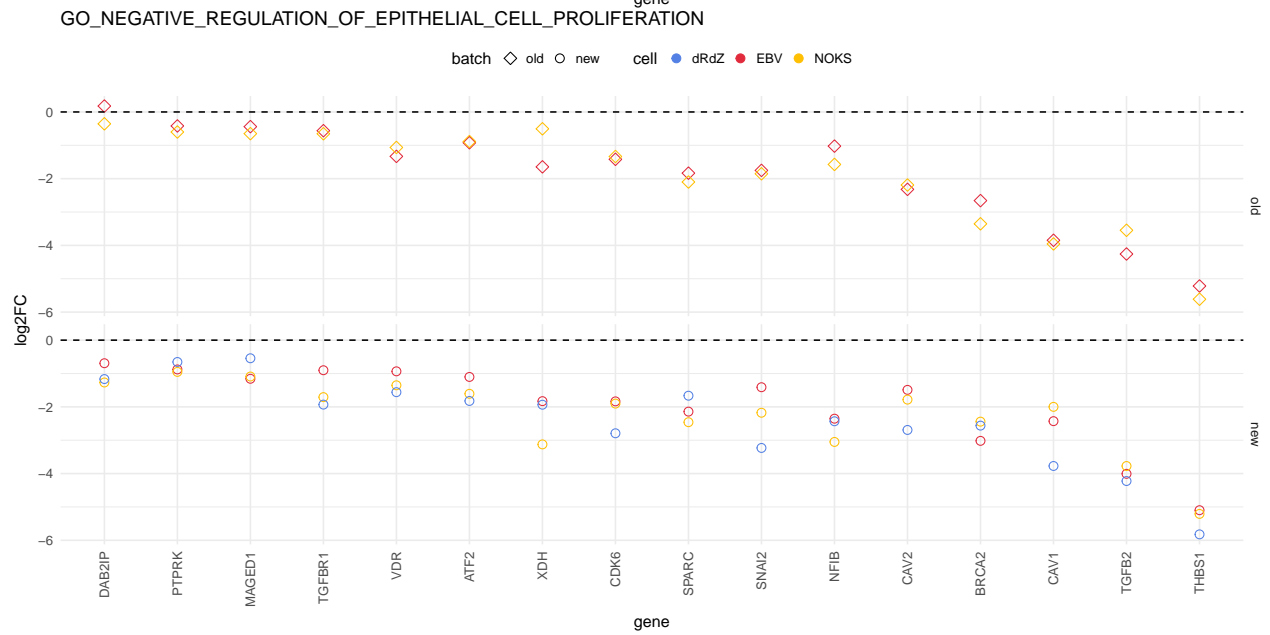
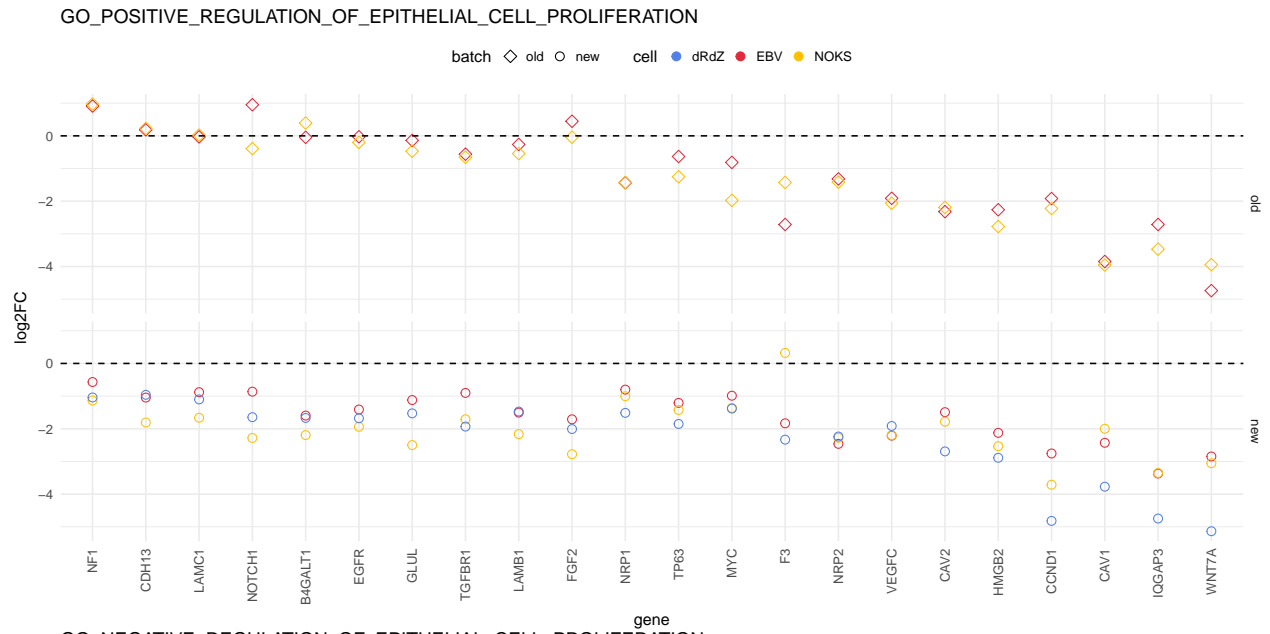




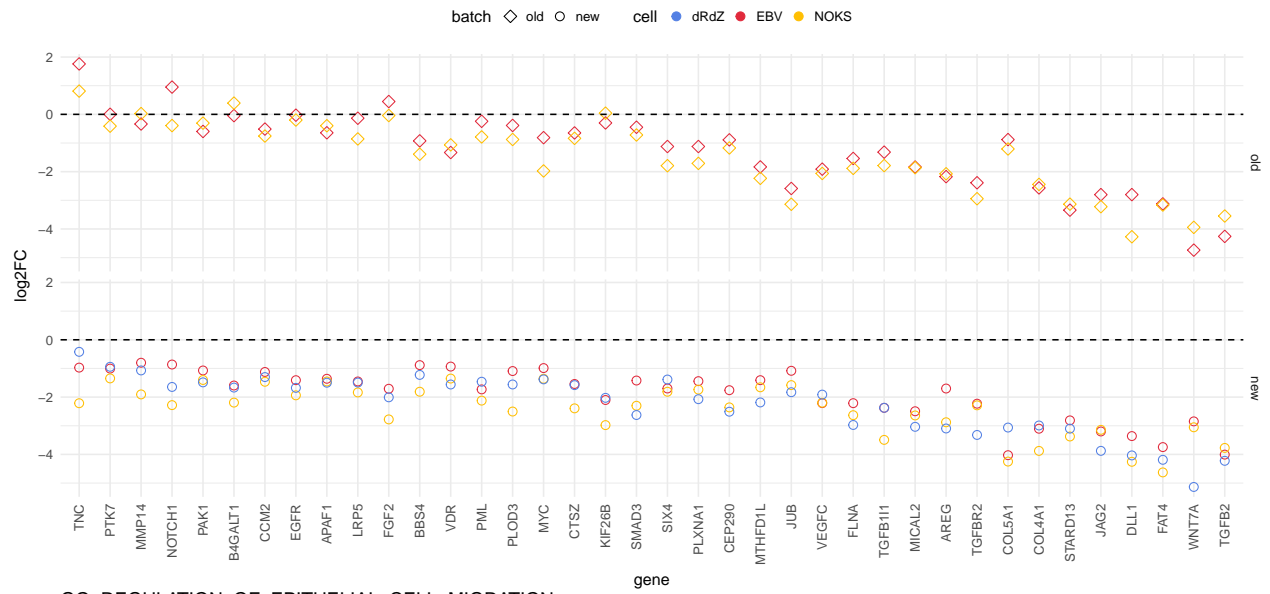








# GO\_MORPHOGENESIS\_OF\_AN\_EPITHELIUM



# GO\_REGULATION\_OF\_EPITHELIAL\_CELL\_MIGRATION

