# Bayesian Machine Learning approach for modelling gene expression time series.

# Contents

# Abstract

Cellular stress responses involve rapid and dynamic changes in gene expression, which are best captured through high-resolution time-series data. High-throughput RNA sequencing generates count data often exhibiting overdispersion due to biological variability and sequencing noise. Accurately modelling temporal gene expression changes, especially during stress responses, requires models which can account for complex temporal patterns and count-based variability.

In this study, we compared two time-series differential expression tools – GPcounts, which utilises a negative-binomial Gaussian Process framework, and maSigPro, which applies polynomial regression – to analyse RNA-seq data from mouse liver under fasting and *ad libitum* conditions. A model dataset of 80 mice (40 fasting; 40 *ad libitum*: equal male/female) was sampled every 4 hours from 8 to 44 hours, for 10 time points with two biological replicates each. Differentially expressed genes were identified for both model- and sex-dependent conditions, then classified as shared or unique across tools. GO enrichment analysis and expression trajectory plots were used to examine biological relevance.

Both GPcounts and maSigPro identified core metabolic genes involved in fatty acid metabolism. However, whilst maSigPro-unique genes were predominantly enriched for fatty acid metabolic pathways, GPcounts-unique genes, especially in males, were enriched for cell-cycle processes, potentially indicating robust repair mechanisms or compensatory proliferation in response to fasting.

Our findings demonstrate that model choice significantly affects differential gene expression identification: GPcounts excels at capturing non-linear, proliferative gene

expression but is computationally intensive, whilst maSigPro offers a faster, more accessible approach for detecting broad metabolic trends. Future studies should integrate kernel extensions and shrinkage methods to further refine these models. Additionally, future extensions for GPcounts could include periodic kernel implementation to identify rhythmic circadian gene expression. Functional validation, like CRISPR knockouts, will be essential for confirming the biological relevance of uniquely identified genes.

# Acknowledgements

# Introduction

High-throughput RNA sequencing (RNA-seq) generates count data to quantify gene expression across conditions; however, such counts often exhibit overdispersion, where the variance exceeds the mean [1]. This can arise from both biological variability and technical factors such as sequencing noise [2]. Accurately capturing dynamic changes in gene expression over time requires models which accommodate both complex temporal patterns and count-based variability. Supervised machine learning techniques like Gaussian Processes (GPs) paired with a negative binomial likelihood offer a flexible, non-parametric framework that can capture such dynamics, providing a tool for identifying differentially expressed genes (DEGs) and modelling non-linear temporal trajectories in RNA-seq datasets [3-6].

Robust analysis of time-series RNA-seq data is challenging due to issues with normalisation, scalability and reproducibility [7]. Traditional models, such as Poisson regression and linear regression, make simplified assumptions about variance: Poisson regression assumes equal mean and variance, while linear regression assumes normally distributed errors with constant variance. However, RNA-seq data often exhibits overdispersion, which these models fail to capture [8]. Negative binomial GPs offer a robust Bayesian approach by modelling count variability, missing data and temporal uncertainty by updating prior beliefs based on observed data. This framework is particularly suited to modelling temporal gene expression as its dispersion parameter explicitly accounts for overdispersion [9].

GPcounts is a specialised Python package designed for temporal and spatial counts analysis, implementing a negative binomial likelihood with GPs, variational inference and log-link transformations [10]. Here, it has been applied to study differential gene expression in mouse liver under fasting and *ad libitum* (unrestricted access) feeding conditions.

In mouse liver, fasting triggers transcriptional reprogramming, upregulating key metabolic pathways like glycogenolysis and lipolysis [11,12]. In previous animal studies, intermittent fasting has been shown to increase stress resistance, extend lifespan, and improve inflammatory and oxidative stress responses, in addition to reducing Alzheimer's-associated β-amyloid accumulation [13-17]. Sexual asymmetry in transcriptional responses to fasting has also been observed. Sex-specific hormonal regulation, such as increased glucocorticosteroid circulation, has been implicated in greater amino acid catabolism from muscle in males, whilst greater oestrogen levels are associated with higher lipid metabolism in females [18-21].

Given that maSigPro and GPcounts differ in their approaches to modelling count data: maSigPro uses polynomial regression, assuming normally distributed residuals, whereas GPcounts employs GP regression with a negative binomial likelihood, we expect to observe divergence in model sensitivity and the genes detected. We hypothesise that GPcounts will identify additional biologically relevant genes, whilst maSigPro will efficiently detect broader metabolic trends. Both models were evaluated using a high-resolution RNA-seq time-series dataset comprising liver samples from 80 mice (40 fasting, 40 *ad libitum*; 40 males, 40 females) collected across 10 time points at 4-hour intervals, with two biological replicates per time point, spanning over 20,000 genes. This study aims to compare model performance and investigate sex-dependent differences in temporal gene expression.

# Methods

**Data**

RNA-seq count data from 80 mouse liver samples (40 fasting, 40 *ad libitum*; 40 males, 40 females) were provided by Dr. Fustin's laboratory at the University of Manchester. Mice were kept in a standard 12:12 hour light-dark cycle for the first 10 days to synchronise their internal clock with external cues before being moved to constant darkness on day 11 until the experiment's conclusion. Samples were collected at 4-hour intervals from 8 to 44 hours, with two biological replicates per collection point providing a high-resolution dataset of over 20,000 genes. Sacrifice and sample collection were conducted in darkness to preserve endogenous circadian rhythms.

Initial filtering steps were applied to refine the RNA-seq data. For inclusion into the combined-sex dataset, at least 40 non-zero observations were required; for male- or female-specific datasets, a minimum of 20 non-zero observations were required. Data was log-transformed for normalisation. A total of 20,545 genes were retained after filtering.

**Linear Regression**

Linear regression models the relationship between a dependent variable and one or more independent variables. Here, we modelled gene expression as a function over time, under different conditions (fasting vs *ad libitum*) and across sexes. In its simplest form, linear regression is defined as:

$$y = mx + c$$

Where $y$ (dependent variable) is related to $x$ (independent variable) through $m$ (gradient) and $c$ (intercept). However, gene expression over time is rarely linear and often follows complex, non-linear dynamics. This simplistic model is inappropriate for capturing time-series gene expression data, requiring flexible models that can accurately capture time-series gene expression profiles.

**maSigPro**

Differential expression analysis was performed using the maSigPro R package [v1.78.0] [22]. maSigPro applies a two-step regression strategy to identify genes with significant temporal expression changes between groups.

Firstly, a general regression model was fitted using dummy variables to account for categorical effects (Fasting vs *ad libitum*). The model is adjusted using the least squares method, minimising squared differences between observed and predicted values. Second, a fourth-degree polynomial regression was used to model non-linear expression trajectories over time:

$$y = b_0 + b_1 x_1 + b_2 x_1^2 + \cdots + b_n x_1^n$$

Where $y$, representing gene expression, is modelled as a function of time $x_1$, through a quartic polynomial. The coefficients $b$ define the magnitude and direction of temporal changes in expression, shaping the curve to best capture the gene's dynamic behaviour over time [30].

The p.vector() function was used with degree = 4, to define the polynomial, assuming Gaussian error by default. Significance was assessed using the F-statistic with p-values corrected using the Benjamini-Hochberg method [31]. Genes with FDR < 0.05 and $R^2 \geq$

0.6 were considered significant. Stepwise regression refinement was performed using the T.fit() function with α = 0.05, returning p-values for each coefficient.

## GPcounts

GPcounts [v2.0.0] models gene expression using GPs with a negative binomial likelihood, allowing for non-parametric modelling of non-linear temporal dynamics while accounting for overdispersion in RNA-seq data. GPs define the distribution over possible curves, using a kernel function to determine similarity between gene expression values at nearby time points [3].

The latent function $f(x)$ represents the underlying expression pattern, where $k(x, x')$ represents the kernel function and 0 represents the mean function, assuming no prior bias about the function's average behaviour.

$$f(x) \sim GP\big(0, k(x, x')\big)$$

This was linked to expected gene expression via a logarithmic link, ensuring expected count values remain positive, as required for sequencing data. This transformation also helps stabilise variance.

$$\mu = exp(f(x))$$

Observed counts are assumed to follow a negative binomial distribution, accounting for variance, where $\theta$ represents the dispersion parameter.

$$y \sim NB(\mu, \theta)$$

A two-sample log likelihood ratio test was used to compare conditions. The null hypothesis assumed a shared trajectory between conditions, whereas the alternative allowed separate trajectories.

Genes with LLR > 0 and FDR < 0.05 (Benjamini-Hochberg correction) were considered significant to retain consistency with maSigPro [31]. LLRs were compared using a chi-squared distribution with 1 degree of freedom.

Pre-filtering required ≥ 15 observations per gene. Log-normalised counts were rounded to the closest integer to minimise Cholesky decomposition failures. For computational efficiency, safe mode was disabled, allowing for a single fit per gene. Ten inducing points were manually set at 4-hour intervals from 8 to 44 hours. Analysis was conducted using the University of Manchester HPC (CSF3) using 1TB cores.

**Model Effect analysis**

Differential expression analysis was performed between fasting and *ad libitum* conditions separately for male, female and combined-sex datasets. For each group, differential expression was assessed using maSigPro and GPcounts. The intersection of significant genes identified by both methods was recorded, generating three gene lists (maSigPro, GPcounts, shared) per dataset.

Genes were classified as consistent if they were identified by both methods and unique if only identified by one. Intersections between lists were calculated using Python's set () function. Venn diagrams were generated using Matplotlib [v3.10.1] to visualise overlaps within sex group (two-circle venn, maSigPro vs GPcounts) and across all groups (maSigPro/GPcounts: male, female and combined, where GPcounts explicitly models sex as a covariate). A summary table of significant genes was also produced [32].

**Sex-Dependent analysis**

Male and female DEGs were compared using three analyses, GPcounts, maSigPro and the intersection of both, to assess sex-dependent responses. All analyses were performed on datasets combining fasting and *ad libitum* samples. Genes identified as significant in both sexes were labelled as shared, while those significant in only one sex were labelled unique. Using Python's set () function, overlaps were visualised as two-circle Venn diagrams using Matplotlib [v3.10.1].

**GO enrichment**

Go enrichment was performed locally in R using clusterProfiler [v4.14.6], AnnotationDbi [v1.68.0] and org.MM.eg.db [v3.20.0] [33-35]. Gene sets for both model and sex comparisons were tested separately using the enrichGO() function with the following parameters:

- OrgDb = org.Mm.eg.db
- ont = "BP" (Biological Processes)
- keyType = "SYMBOL"
- pvalueCutoff = 1
- qvalueCutoff = 1

The top 30 GO terms (by adjusted p-value) were visualised in heatmaps using pheatmap [v1.0.12] [36].

## Model- and Sex- Specific Gene Visualisation

DEGs identified by both maSigPro (min obs. 15, FDR <0.05, $R^2 \geq 0.6$) and GPcounts (min obs. 15, FDR <0.05), LLR >0) between fasting and *ad libitum* conditions were visualised using GPcounts' in-built plot () function based on Matplotlib framework [v3.10.1].

For model-specific visualisation, six genes previously identified through GO enrichment analysis were selected: three associated with fatty acid metabolism and three with cell-cycle regulation. For sex-specific visualisation, three genes were selected: one shared between sexes, one male-specific, and one female-specific. Data was plotted using the GPcounts plot () function.

## Reproducibility and Computational Challenges

GPcounts encountered reproducibility challenges due to GPflow-related memory errors on the CSF. Datasets were split into ten subsets using Pandas [v2.2.3] and processed independently on the CSF [38]. Certain errors (e.g., AttributeError: 'GP_nb_zinb' object has no attribute 'fit_model') were consistently encountered and were not reproducible locally using the same packages, suggesting environment-specific issues. Genes which consistently encountered this error were removed from the final dataset, resulting in a smaller gene set analysed than maSigPro. Code for splitting and merging results was written in Python using Pandas [v2.2.3].

# Results

**Model differences in identification of Significantly Differentially Expressed Genes (DEGs) between fasting and non-fasting conditions.**

To identify DEGs under fasting and *ad libitum* conditions, both maSigPro and GPcounts were applied to male, female and combined mouse liver datasets. maSigPro was initially applied to 20,545 genes across all datasets. Genes with fewer than 15 observations were filtered out, and the Benjamini-Hochberg method was applied (FDR < 0.05). 10,395 male, 12,780 female and 14,652 combined-sex genes were initially considered differentially expressed. Following stepwise regression and an $R^2$ threshold of ≥ 0.6, this was reduced to 4,227, 4,542 and 2,197 DEGs, respectively (Fig. 1d).

GPcounts, using the same input data and observation threshold (≥15), identified 2,364 (male), 3,509 (female), and 6392 (combined) DEGs after applying a log-likelihood ratio (LLR > 0) and correcting for false discovery rate using the Benjamini-Hochberg procedure (Fig. 1d).

To compare both models, consistent and unique DEGs for each condition were compared (Fig. 1):
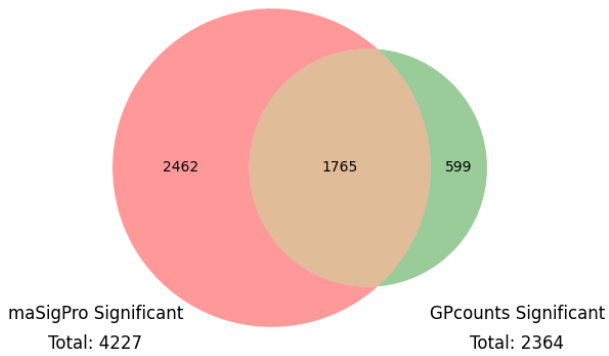
- Consistently identified DEGs by both methods: 1,765 (male), 2,223 (female), 1,942 (combined) (Fig. 1a).
- GPcounts-unique DEGs: 599 (male), 1,286 (female), 4,350 (combined) (Fig. 1b).
- maSigPro-unique DEGs:  2,462 (male), 2,319 (female), 255 (combined) (Fig. 1c).

A notable crossover of DEGs between models and sexes was observed, suggesting these genes are significant in fasting-induced stress responses. The higher number of DEGs identified by GPcounts for the combined dataset, and in maSigPro for the sex-specific datasets, suggests the models vary in how they handle replicates, variance and model complexity. The difference between male-unique DEGs in GPcounts vs maSigPro (599 vs 2,462) and combined sex unique DEGs (4350 vs 255) is particularly notable (Fig. 1a, c).

A large proportion of genes not identified in male or female datasets using GPcounts were detected in the combined sex dataset (2690) (Fig. 1e). Almost all genes detected in the combined sex dataset for maSigPro and GPcounts were also detected in the male- and female-only datasets (Fig. 1f).
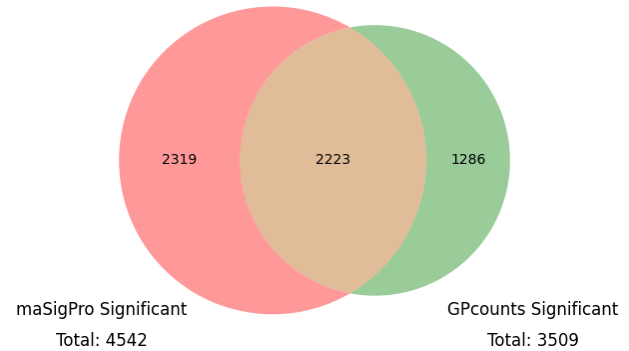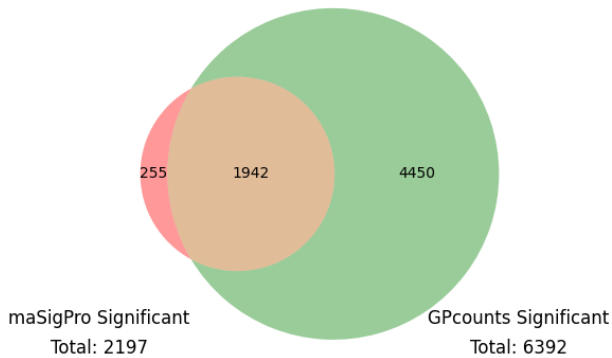
**a)**

Male maSigPro and GPcounts Significant Genes
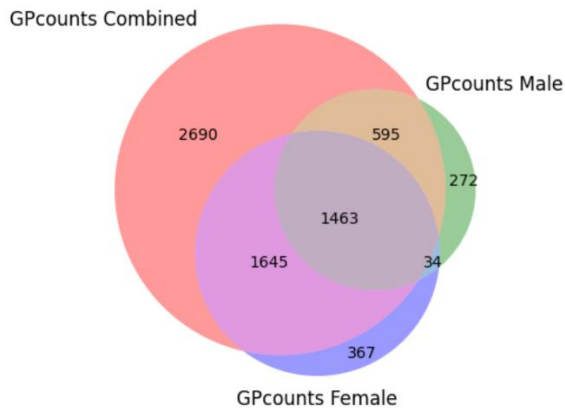


**b)**

Female maSigPro and GPcounts Significant Genes
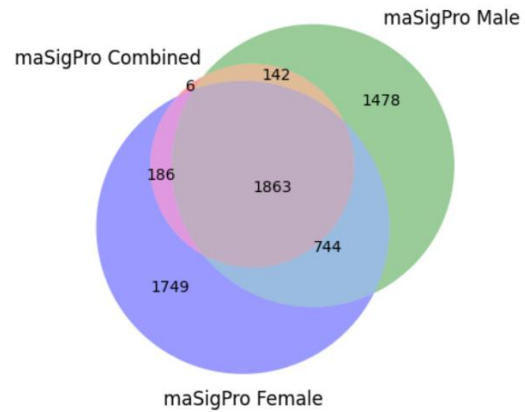


**c)**

Combined maSigPro and GPcounts Significant Genes



**d)**

| Significant Genes | Male | Female | Combined |
|---|---|---|---|
| Identified by maSigPro | 4227 | 4542 | 2197 |
| Identified by GPcounts | 2364 | 3509 | 6392 |
| Consistently Identified | 1765 | 2223 | 1942 |
| Unique to GPcounts | 599 | 1286 | 4450 |
| Unique to maSigPro | 2462 | 2319 | 255 |

**e)**



**f)**



**Fig. 1. Model-dependent overlap between differentially expressed genes identified in single-cell mouse liver datasets.** GPcounts ((min obs. 15, FDR <0.05), LLR >0) and maSigPro (min obs. 15, FDR <0.05, rsq = 0.6) were used to measure significantly differentially expressed genes between fasting and *ad libitum* conditions, sampled at 4-hour intervals from 8 to 44 hours. (a-c) Significantly differentially expressed genes identified by GPcounts and maSigPro in (a) male, (b) female, and (c) combined sexes. (d) Combined data table. (e) Comparison of significant differentially expressed genes identified by GPcounts in male, female and combined datasets. (f)Comparison of significant differentially genes identified by maSigPro in male, female and combined datasets.

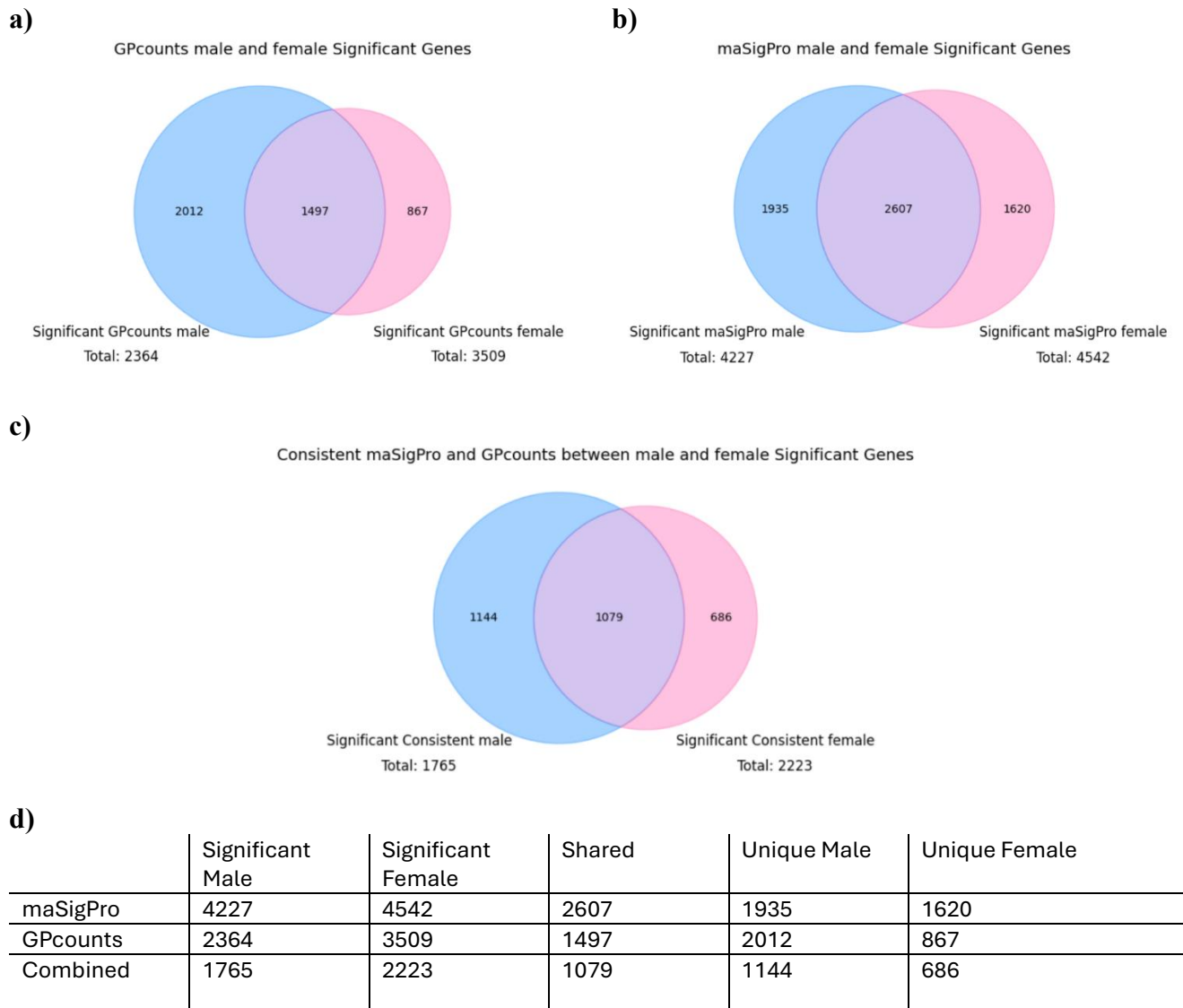**Model differences in sex-dependent Differential Gene Expression**

To assess sex-specific differences in gene expression under fasting and *ad libitum* conditions, we analysed DEGs in male and female mice using both maSigPro (minimum observations ≥15, FDR < 0.05, $R^2$ ≥ 0.6) and GPcounts (minimum observations ≥15, FDR < 0.05, LLR > 0). Analysis was conducted separately by model to evaluate how each statistical model detects sex-dependent expression changes. Given the distinct hormonal and metabolic responses between sexes, as well as differences in statistical models, divergence between the DEG sets was expected.

For each model, DEGS were classified as unique to males, unique to females, or shared between sexes (Fig. 2):

- Sex-Dependent DEGs identified by GPcounts: 2,012 (unique male), 867 (unique female), 1,497 (male and female) (Fig. 1a).
- Sex-Dependent DEGs identified by maSigPro: 1,935 (unique male), 1,620 (unique female), 2,607 (male and female) (Fig. 1b).
- Sex-Dependent DEGs identified by both GPcounts and maSigPro: 1,114 (unique male), 686 (unique female), 1079 (male and female) (Fig. 1c).

Both models identified more male-specific DEGs compared to female-specific DEGs, this trend was more pronounced in GPcounts, where only 867 genes were unique to females versus 2,012 in males (Fig. 2. a). A similar pattern was observed in maSigPro, although the gap was narrower, with 1,620 in female's vs 1,935 in males (Fig. 2. b). The number of shared DEGs between sexes was also higher in maSigPro (2,607) than in GPcounts (1,497), indicating strong overlap between male and female expression profiles in the polynomial-based model.

The set of DEGs consistently identified across both models and sexes showed a similar trend, with 1,114 genes unique to males, 686 unique to females, and 1,079 shared (Fig. 2. c). This suggests a reproducible bias towards identifying more male-specific transcriptional responses under these conditions, regardless of the statistical model applied.

**a)**

GPcounts male and female Significant Genes

| | | |
|---|---|---|
| 2012 | 1497 | 867 |

Significant GPcounts male
Total: 2364

Significant GPcounts female
Total: 3509

**b)**

maSigPro male and female Significant Genes

| | | |
|---|---|---|
| 1935 | 2607 | 1620 |

Significant maSigPro male
Total: 4227

Significant maSigPro female
Total: 4542

**c)**

Consistent maSigPro and GPcounts between male and female Significant Genes

| | | |
|---|---|---|
| 1144 | 1079 | 686 |

Significant Consistent male
Total: 1765

Significant Consistent female
Total: 2223

**d)**

| | Significant Male | Significant Female | Shared | Unique Male | Unique Female |
|---|---|---|---|---|---|
| maSigPro | 4227 | 4542 | 2607 | 1935 | 1620 |
| GPcounts | 2364 | 3509 | 1497 | 2012 | 867 |
| Combined | 1765 | 2223 | 1079 | 1144 | 686 |

**Fig. 2. Sex dependent overlap between differentially expressed genes identified in single-cell mouse liver datasets.** GPcounts ((min obs. 15, FDR <0.05), LLR >0) and maSigPro (min obs. 15, FDR <0.05, rsq = 0.6) were used to measure significantly differentially expressed genes between fasting and *ad libitum* conditions in male and female mice sampled at 4-hour intervals from 8 to 44 hours. (a-c) Significantly differentially expressed genes identified by GPcounts and maSigPro between sexes, (a) male-female GPcounts comparison, (b) male-female maSigPro comparison, (c) male-female combined maSigPro GPcounts comparison. (d) Combined data table.

## Model differences in fasting versus non-fasting pathways

To explore the biological relevance of genes identified by maSigPro and GPcounts, GO enrichment analysis was conducted on the consistent and unique gene sets using Biological Process ontology [33–35]. Genes involved in core metabolic processes - particularly fatty acid metabolism, lipid breakdown, and energy derivation- were consistently enriched across both models (Fig. 3a). These shared genes showed strong enrichment, with -$\log_{10}(p)$ values commonly exceeding 15. Genes related to nucleoside, purine, and ribose-dependent metabolic processes were also enriched, though with lower significance.

Unique gene sets identified by GPcounts in both and female datasets were strongly enriched for cell-cycle processes, like chromosome segregation and cell-cycle transitions, with -$\log_{10}(p)$ values commonly exceeding 15 (Fig. 3a). This was more pronounced in females, with high-$\log_{10}(p)$ values observed for nuclear division, chromatid and chromosome segregation. In comparison, maSigPro-unique genes were more broadly distributed across biological processes with less pronounced enrichment (Fig. 3a). Enrichment in males was associated with small-molecule catabolic processes, while in females, these genes showed enrichment in energy derivation and small-molecule metabolic processes.

Together, these findings suggest consistently identified genes are linked to core metabolic processes like fatty acid metabolism, whilst model-specific genes highlight distinct biological features. Specifically, GPcounts-unique genes emphasise cell-cycle regulation, and maSigPro-unique genes are enriched in general catabolism.

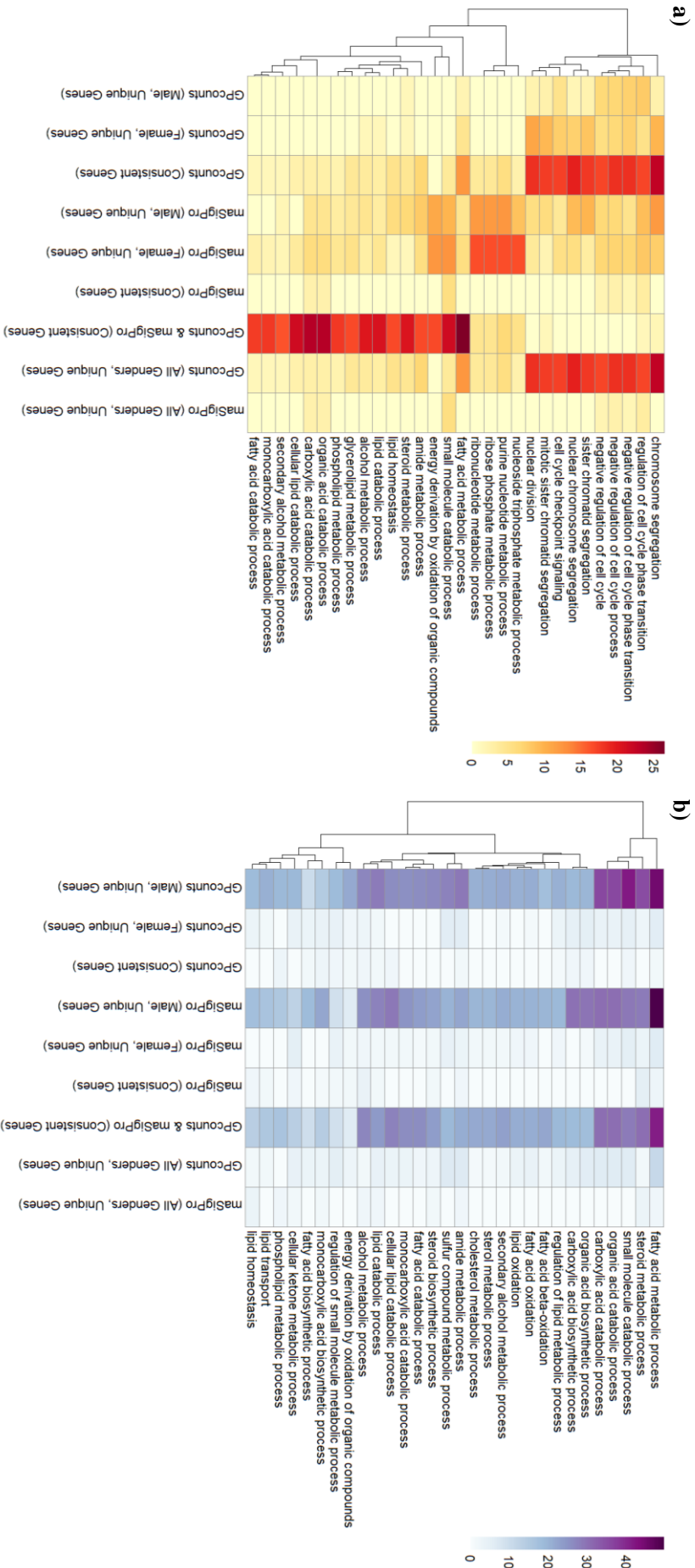## Model differences in sex dependent GO enrichment

To explore the functional consequences of sex-specific gene expression differences, GO enrichment analysis was performed on the unique and shared DEG sets identified in male and female mice using both maSigPro and GPcounts [33 – 35].

Across all gene sets, the most strongly enriched biological processes involve key metabolic and catabolic processes – fatty acid metabolism, steroid metabolism, and small molecule degradation. These processes were consistently upregulated in both sexes (Fig. 3b). The combined gene set, comprising DEGs identified in both males and females, exhibited strong enrichment with -$\log_{10}(p)$ values exceeding 40. Related processes such as lipid catabolism were also significantly enriched.

In contrast, sex-specific gene sets showed more selective enrichment. Female-unique genes, identified by either model, showed weaker enrichment, with low to moderate significance for pathways for fatty acid metabolism and sulphur compound metabolic processes. No -$\log_{10}(p)$ values for female-unique genes exceeded 20.

Male-unique genes, particularly those identified by GPcounts, showed exceptionally strong enrichment for fatty acid metabolism and small molecule catabolism ($-\log_{10}(p) > 40$), alongside broader enrichment across metabolic processes like steroid biosynthesis, amide metabolism, and lipid catabolism (Fig. 3b).

Combined, these findings suggest a shared trend of metabolic regulation across both sexes in response to fasting and *ad libitum* conditions. The exceptionally strong enrichment and broader pathway representation in male-unique genes, identified by both models, indicate potential sex-specific transcriptional regulation, especially within fatty acid metabolic networks.

**Fig 3. GO enrichment heatmaps of significantly differentially expressed genes (adjusted p-value < 0.05) identified by GPcounts and maSigPro in male and female fasting and *ad libitum* mouse liver samples.** Enrichment was performed using the enrichGO() function from the clusterProfiler package using org.Mm.eg.db mouse genome database. The top 30 enriched biological processes were selected based on their -log10 (adjusted p-value). Colour value reflects significance level, with darker shades indicating stronger enrichment and smaller P-values. a) Model dependent GO enrichment for genes uniquely identified in male and female datasets by GPcounts or maSigPro, as well as those consistently found by both models. b) Sex dependent GO enrichment for genes identified as significantly differentially expressed between male and female datasets. Here the different sexes are compared using the same models, again highlighting consistent and unique gene sets.

## Comparison of maSigPro and GPcounts models

To illustrate model differences, expression profiles of selected genes were visualised under fasting and *ad libitum* conditions (Fig. 4). Differences in sensitivity between maSigPro and GPcounts were apparent across genes linked to fatty acid metabolism (Acox1, Ppat, HAGH), and cell-cycle regulation (E2F1, RIOK2, Pola1).

E2F1 (Fig. 4a), a regulator of S-phase cyclin transcription and DNA replication, repair and apoptosis, was identified as significant by both methods. Dysregulation has been linked to higher rates of cell proliferation [39]. Expression gradually declined from a peak of ~200 at 8 hours to below 100 by 44 hours, with downregulation observed under fasting conditions.

RIOK2 (*RIO Kinase 2*; Fig. 4c), an ATPase involved in 40s ribosomal subunit maturation and preventing premature translation initiation, was identified as significant by maSigPro [40-41]. Expression followed a sinusoidal pattern between ~1500 and ~2500 counts. A sharp peak observed between 36-44 hours may contribute to its classification as significant, potentially reflecting sensitivity to local noise. Notably, the alternative model shows a sinusoidal fluctuation with a period of around 30 hours, suggesting potential intradian regulation.

Pola1 (*DNA polymerase α 1, catalytic subunit*; Fig. 4e) encodes the catalytic subunit of DNA polymerase α, key for initiating DNA replication during cell division and repair [42]. Expression steadily declines from ~400 at 8 hours to ~200 by 44 hours. Pola1 was only identified as significant by GPcounts, potentially reflecting its better handling of low-count, variable data using a negative binomial likelihood.

Acox1 (*acyl-CoA oxidase*; Fig. 4b) encodes the first enzyme in the fatty acid β-oxidation pathway, catalysing desaturation which leads to fatty acid catabolism and ATP production [43, 44]. Identified as significant by both models, the alternative model showed a sinusoidal expression pattern. Expression rose strongly from ~150,000 at 8 hours to ~350,000 at 20, followed by a decline and trough to ~200,000 at ~34 hours, roughly 26 hours later. This was followed by a secondary rise at 44 hours.

Ppat (*phosphoribosyl pyrophosphate amidotransferase*; Fig. 4d), associated with purine biosynthesis, cell metabolism, and proliferation, exhibited a sinusoidal, rhythmic expression pattern between ~2000 (8 and 36 hours) and ~1000 (20 and 44 hours) counts [45 – 47]. It was classified as significant by maSigPro but not by GPcounts, suggesting potential differences in model sensitivity to local noise.

HAGH (*hydroxyacylglutathione hydrolase*; Fig. 4f), a mitochondrial enzyme involved in glutathione metabolism, was identified exclusively by GPcounts [47-48]. Its expression showed a smooth decline and partial recovery pattern, which appears weakly sinusoidal. Variance between conditions became more pronounced between 20 and 40

hours, potentially explaining the difference in model classification. Noise identification using a negative-binomial distribution may enable the identification of HAGH as significantly differentially expressed.

These examples illustrate how statistical model assumptions can lead to the identification of distinct sets of differentially expressed genes, even with identical temporal trends.
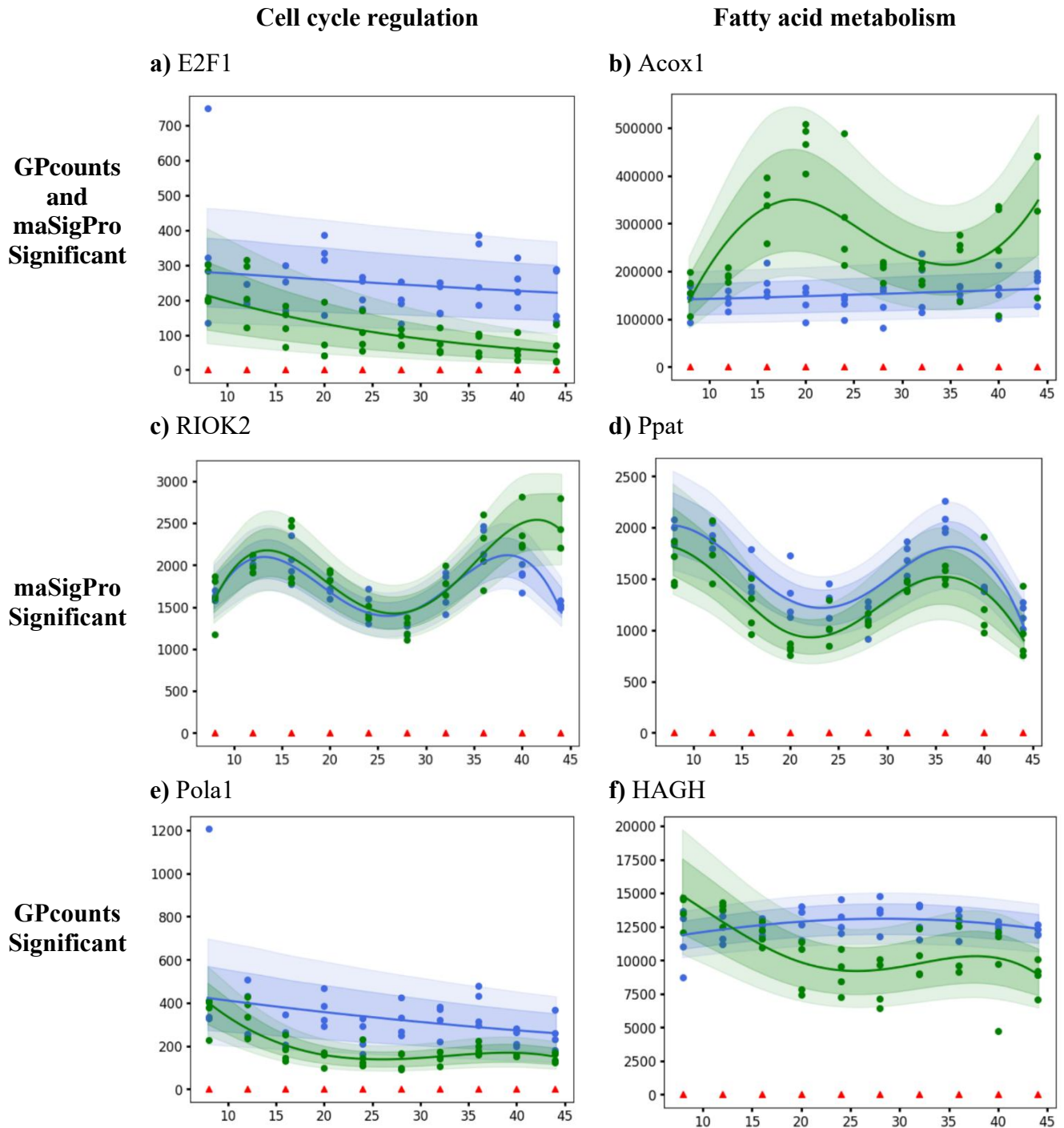
**Sex-dependent gene visualisation**

Expression profiles of three representative genes were visualised to illustrate sex-dependent regulation (Fig. 5). While fatty acid metabolism pathways were enriched in both sexes, minor gene-level differences were observed.

Etfdh, upregulated in both sexes, encodes ETFDH (*electron transfer flavoprotein dehydrogenase*; Fig. 5a), a mitochondrial electron transfer system enzyme, transferring electrons from flavin-containing dehydrogenases to the respiratory chain [49, 50]. Defects can cause acyl-CoA-dehydrogenase deficiency and hypoglycaemia [51]. A sinusoidal pattern is observed atop a general upward trend in expression with a peak at 25 hours, consistent with dynamic but progressively increasing gene activity.
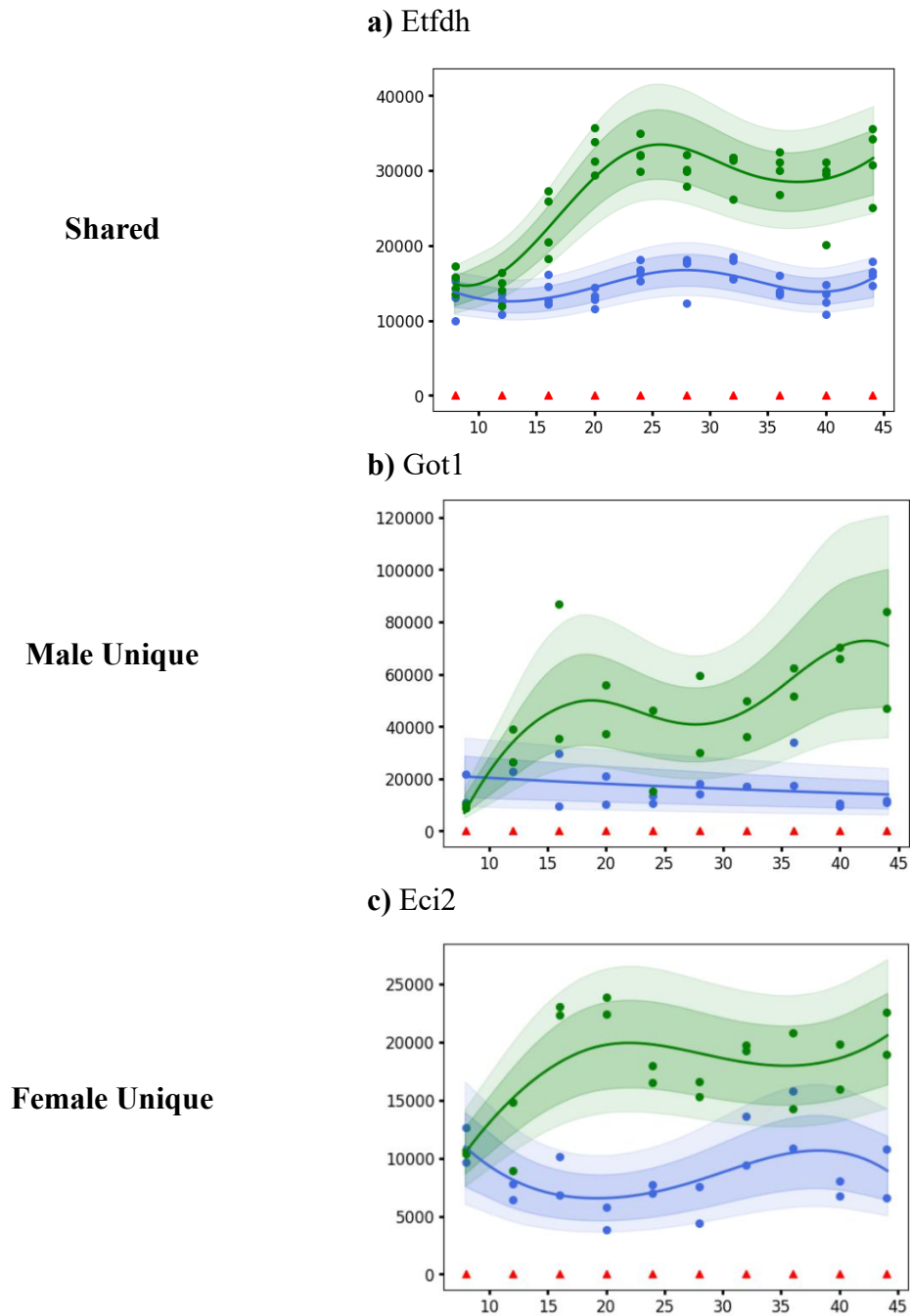
GOT1 (*glutamic-oxaloacetic transaminase* 1: Fig. 5b) encodes a cytosolic enzyme involved in amino acid metabolism, energy production, glucocorticoid responses, and steroid metabolism [52, 53]. Identified exclusively in males, the significant increase in counts (~10,000 to ~80,000) suggests sex-dependent induction, potentially linked to enhanced amino acid catabolism, glycogenolysis, and steroid metabolism under fasting conditions [12, 20]. A sinusoidal pattern was observed atop the upward trend, characterised by a smaller peak at 18 hours, a decline, and a larger peak at 42 hours.

Eci2 (*Enoyl-CoA Delta Isomerase 2*; Fig. 5c) was identified exclusively in females. Eci2 encodes a key mitochondrial enzyme involved in β-oxidation of unsaturated fatty acids. Expression rose from ~10,000 to ~20,000, peaking at 20 hours, with a second increase beginning near 44 hours. After the initial rise, gentle oscillation suggests a transition from sustained upregulation to rhythmic control. Combined, this suggests Eci2 has a role in enhancing fatty-acid metabolism and potentially reflects sex-specific lipolysis [12].

These examples illustrate how sex-dependent gene expression differences modulate distinct aspects of metabolic regulation, while still converging on fatty-acid metabolism.

**Fig. 4. Model-specific and shared significant genes in cell cycle regulation and fatty acid metabolism: maSigPro and GPcounts.** Genes significantly differentially expressed between fasting and *ad libitum* conditions in the combined dataset using maSigPro (min obs. 15, FDR <0.05, rsq = 0.6) and GPcounts (min obs. 15, FDR <0.05), LLR >0). Genes uniquely identified by one model or shared by both were visualised. RNA-seq counts are plotted over time (x-axis: 4-44 hours in 4-hour increments), with the y-axis representing expression levels. The blue line shows the null model fit, the green line, the alternative model. Shaded areas represent 95% confidence intervals. Red triangles indicate inducing points and specific time points of sample collection. (a, c, e) Genes identified in cell-cycle regulation, (b, d, f) Genes identified in fatty acid metabolism.

**a)** Etfdh

**Shared**

**b)** Got1

**Male Unique**

**c)** Eci2

**Female Unique**

**Fig 5. Sex-specific and shared significant genes in fatty acid metabolism.** Significantly differentially expressed genes between fasting and *ad libitum* conditions in shared, male and female datasets. Genes were considered unique if identified by both maSigPro (min obs. 15, FDR <0.05, rsq = 0.6) and GPcounts (min obs. 15, FDR <0.05), LLR >0). RNA-seq counts are plotted over time (x-axis: 4-44 hours in 4-hour increments), with the y-axis representing expression levels. The blue line shows the null model fit, the green line, the alternative model. Shaded areas represent 95% confidence intervals. Red triangles indicate inducing points and specific time points of sample collection. (a) Etfdh, a metabolic gene identified in both male and female datasets. (b) Got1 a steroid-associated gene identified exclusively in male datasets. (c) Eci2 a lipid-metabolism associated gene identified exclusively in female datasets.

# Discussion and Conclusion

This study systematically compared two statistical frameworks for time-series RNA-seq analysis: maSigPro, which models gene expression over time using polynomial regression, and GPcounts, which applies a Bayesian GP regression with a negative binomial likelihood to explicitly account for variance [2, 10, 22]. Both models were evaluated using mouse liver gene expression profiles under fasting and *ad libitum* feeding conditions.

Accurate modelling of temporal gene expression dynamics is critical for identifying regulatory pathways responsive to metabolic stress. Both maSigPro and GPcounts identified a shared core set of genes involved in metabolic regulation; however, they differed in their sensitivity to genes associated with cell-cycle regulation and sex-specific expression patterns (Fig. 4). These differences reflect the distinct statistical assumptions and modelling strategies underlying each approach.

**Flexibility and robustness**

GPcounts captures subtle, non-linear patterns in noisy time-series data by modelling gene expression using GPs with a negative binomial likelihood, making it well suited for RNA-seq datasets where variance exceeds the mean [6, 7, 10]. This approach explicitly accounts for overdispersion and enables the modelling of smooth, non-linear expression dynamics without imposing a predefined curve [4]. Features like sparse inference, a safe mode to mitigate matrix decomposition errors and options to manually specify inducing points enhance its robustness. Together, these capabilities allow GPcounts to distinguish subtle, biologically relevant signals in noisy data (Fig. 4c-f), providing a distinct advantage over fixed regression models when gene expression deviates from simple polynomial behaviour [5].

However, the flexibility of GP models increases the risk of overfitting, especially in smaller datasets or when biological variation is masked by noise [56, 57]. These risks must be weighed against the benefits when selecting the appropriate statistical model for time-series analysis.

In comparison, maSigPro applied a two-step regression strategy based on fourth-degree polynomials, assuming Gaussian noise [22]. This parametric approach is more computationally efficient but less flexible for capturing complex or non-stationary expression dynamics, potentially overlooking biologically relevant processes (Fig. 3a). Originally developed for microarray data, maSigPro remains accessible for researchers with limited computational backgrounds. Its R-based interface allows for local, rapid analysis without requiring high-performance computing, specialised libraries or

extensive programming experience. This makes it well-suited for exploratory analyses and hypothesis generation in academic labs lacking dedicated bioinformatic support.

## GPcounts' current computational demands limit reproducibility and scalability

Despite its modelling flexibility, GPcounts computational intensity presents practical limitations. Fitting thousands of individual GPs to large RNA-seq datasets demands significant time and memory, often requiring high-performance (HPC) systems with sufficient RAM and CPU availability to avoid memory failures during model fitting. Additionally, this may limit the application of GPcounts to larger datasets such as exon count data.

Here, GPcounts required substantial memory and time even with sparse inference enabled, manually defined inducing points, and safe mode disabled. On the HPC system, several genes failed due to GPflow-related attribute errors and memory issues. These failures were not reproducible locally, suggesting potential environment-specific incompatibilities. Although the proportion of excluded genes was small (0.11%, reducing the dataset to 20,521 genes for females and 20,506 for males), it raises concerns about tool reliability in high-throughput settings. Failures may result to overlooked DEGs, particularly for genes near detection thresholds.

Even with HPC resources, achieving full reproducibility across runs remains difficult. Hyperparameter optimisation for GP regression models relies on stochastic optimisation methods, which can converge to different local minima depending on the random initialisation [58]. Consequently, small differences in starting conditions, random seeds, or different hardware can produce slight variations in model fits and selected genes between runs. Whilst such variations are typically minor and are unlikely to alter major biological conclusions, they may complicate the reproducibility of exact gene lists, particularly for borderline differential or rhythmic genes.

To mitigate these issues, the dataset was split into ten subsets and processed separately. While this allowed the pipeline to run, it increased workflow complexity and introduced risks of human error during downstream consolidation. Technical workarounds, like manually defining inducing points, reduce GPcounts' accessibility for users without strong Python or HPC experience.

In contrast, maSigPro's parametric approach is substantially less computationally demanding and more deterministic, once a polynomial degree is specified, results are highly reproducible across different runs, systems and users. This enhances maSigPro's appeal for studies prioritising exact reproducibility, rapid turnaround or broader accessibility in computationally limited environments.

## GPyTorch migration could improve GPcounts' reliability and accessibility

Given these computational challenges, an improvement would be to re-implement GPcounts core modelling framework using GPyTorch, a GP library built on Pytorch [57]. Unlike TensorFlow, which is optimised for large-scale operations, PyTorch is often preferred in research settings for its flexibility and ease of debugging [58, 59]. This may alleviate the observed runtime inconsistencies and attribute errors on the HPC, especially under parallel or diverse environments [60]. GPyTorch also supports scalable variational inference and fast kernel approximations, offering a practical foundation for future development [57]. In principle, this would support running large datasets with full inference while maintaining safe-mode protections to minimise Cholesky decomposition failures or memory and attribute-based errors.

## Shrinkage may mitigate overfitting in low-expression genes

GPcounts' flexibility can lead to overfitting, where a model captures random noise rather than generalisable trends, reducing predictive accuracy. This is common in smaller or noisy datasets, common in RNA-seq experiments with biological replicates or limited time points [61]. Although GPcounts addresses this with overdispersion-aware negative binomial modelling, residual noise in low-expression or highly variable genes remains challenging.

Shrinkage, which pulls extreme or unstable parameter estimates toward group or global averages, could help mitigate overfitting by stabilising variance without distorting biological signals [62]. For GPcounts, shrinkage could be implemented by averaging fitted values that fall beyond a user-defined number of standard deviations from the model, condition-specific group, or global dataset mean. Similar strategies have been successfully employed in edgeR and DESeq2 to stabilise dispersion estimates [28, 63].

## Tool accessibility

While GPcounts provides greater modelling flexibility and sensitivity than maSigPro, its technical demands, including Python coding, dependency management, and HPC configuration, make it less accessible to researchers without computational expertise [64, 65]. maSigPro, in comparison, is user-friendly, fast, and easily implemented in R on standard desktop environments.

In academia, where individual labs computational resources are often limited, simpler tools with accessible pipelines are preferred. From practical experience in both industry and academia, tools using R with simpler pipelines are often preferred due to their limited computational demands on the user. There is a growing need for computational

tools that balance statistical sophistication with usability to ensure broader adoption in biology.

## maSigPro analysis is limited by its regression framework

Originally developed in 2006 for microarray data, where gene expression is continuous and approximately Gaussian, maSigPro has been updated over time [22, 66]. Its default approach applies polynomial regression, assuming a Gaussian error structure. Count-based RNA-seq support was later introduced by adapting this framework to a generalised linear model using a negative binomial distribution [67]. However, this extension remains a retrofit rather than a full redesign, as the core framework still relies heavily on polynomial regression, which is limited for complex RNA-seq datasets.

Visualising gene expression is crucial for assessing model performance, helping identify well-fitted profiles and poorly fitted genes where models may overfit, underfit, or miss biological trends. maSigPro includes built-in plotting functions (see.Genes() and PlotGroups ()), allowing gene expression to be plotted individually or by cluster [10,22]. Plotting regression curves acts as an additional quality control step alongside quantitative thresholds, especially for borderline cases. However, maSigPro's functions sometimes fail to plot regression curves for borderline genes, typically those with large fluctuations, low variance, or outliers, making it challenging to visually assess fit quality. This issue was not present when using GPcounts, even for genes with poor fits.

Another limitation of maSigPro is the rigidity introduced by polynomial regression. A polynomial of degree $n$ can have at most ($n$ – 1) turning points (where the graph of the function changes its slope) [68]. In maSigPro, the polynomial degree is specified before fitting data [22]. Consequently, a quartic polynomial can model at most three bends, a significant limitation for capturing gene expression profiles with multiple oscillations, like ultradian or noisy circadian rhythms [69]. Using higher-degree polynomials also introduces the risk of overfitting, especially at time-series boundaries where oscillations can become exaggerated [61]. This rigid framework restricts maSigPro's ability to model complex temporal patterns, increasing the risk of overfitting noise or underfitting biological rhythms.

The application of a fixed polynomial degree across all genes' risks missing dynamic ultradian oscillations or mischaracterising circadian trends. This likely contributes to false positives and negatives in gene selection and may explain the lack of cell-cycle associated genes identified by maSigPro (Fig. 3a). Combined with the use of an $R^2$ threshold in maSigPro, this could limit the number of discoveries compared to applying only an FDR threshold, as used in GPcounts. maSigPro's criteria may cause it to miss weaker but biologically relevant signals, like sex-specific genetic variations.

In contrast, GPcounts was purpose-built to handle RNA-seq data, integrating a negative binomial likelihood into a Bayesian GP model. This allows it to explicitly account for overdispersion and non-linear temporal dynamics, offering a statistically rigorous approach better suited for complex time-series count data. Re-implementing GPcounts using GPyTorch could further improve its usability and scalability. In practice, maSigPro may remain the preferred model for fast, exploratory analysis, whilst GPcounts is better suited for detailed studies, particularly when computational resources are available.

**GPcounts may reveal compensatory proliferative responses to metabolic stress**

GPcounts identified DEGs enriched for cell-cycle processes, suggesting extended fasting may trigger adaptive regenerative programs in hepatocytes, causing transient cell-cycle re-entry to maintain liver function (Fig. 3a, Fig. 4e) [70-73]. Unlike maSigPro's fixed polynomial fits, GPcounts' flexible GP model with a negative binomial likelihood may better capture these subtle, transient activation patterns, enhancing detection of cell-cycle related DEGs during fasting. Functional validation is required to confirm whether these gene expression changes reflect compensatory proliferation.

**Future directions**

### 1. Advancing Temporal Modelling with Periodic Kernels

Mice were kept under a 12:12 light-dark cycle for 10 days, followed by constant darkness on day 11 and throughout sample collection, to preserve endogenous circadian rhythms. Under a standard light-dark cycle, light inputs to the suprachiasmatic nucleus regulate and coordinate most circadian rhythms in the body [74]. Removal of light cues can cause circadian rhythms to shift or become "free-running", leading to period deviations [75-77]. This may explain the shifts in period length observed in metabolic genes like Acox1, identified as significant by both maSigPro and GPcounts (Fig. 4b).

Future work could extend GPcounts by incorporating periodic kernels into the GP framework to capture ultradian (< 12h), intradian (> 30h), and free-running circadian rhythms (~24h) [24, 25, 69]. Similar models like CircaCompare have previously been implemented in R [78]. Applying this analysis across male/female and fasting/*ad libitum* datasets could reveal model- or sex-dependent circadian regulation. For example, fasting induced upregulation of steroid and gluconeogenesis-related genes in male mice, leading to greater amino acid metabolism [18-20, 79, 80].

Visualising kernel-derived metrics using histograms or density plots could reveal subtle cycle patterns missed by traditional DEG analysis, particularly when expression changes are modest but periodic. A wide spread of posterior periods, especially those

slightly above or below 24 hours may reflect free-running circadian clocks as light-dependent suprachiasmatic nucleus regulation is removed [81].

### 2. Functional Validation

To validate the biological relevance of DEGs, CRISPR/Cas9 knockout models could be employed [82]. Candidate genes should be prioritised based on statistical significance, functional novelty, or identification by both models, balancing scientific value with ethical, financial, and time constraints.

Examples include Ppat, involved in purine biosynthesis, FBP2, which plays a role in liver gluconeogenesis, and MAP9, a mitotic spindle regulator implicated in cell-cycle developmental arrest (Fig. 4d) [19, 45, 46, 83]. These genes may play important roles in the liver's response to metabolic stress. GPcounts uniquely identified DEGs were enriched for mitotic and cell-cycle processes, suggesting a stress-induced regenerative response (Fig. 3a-b). Comparing knockout and control mice under fasting could test whether loss of these genes affects survival, regeneration, or metabolic adaptation.

### 3. Bioinformatic Follow-Up

After experimental validation, a second round of transcriptomic profiling could reveal downstream effects. Enrichment analysis (e.g., using enrichGO), network construction (e.g., using WGCNA), and temporal clustering could identify altered pathways and co-expression modules [84]. Bayesian trajectory inference and hierarchical modelling could integrate prior biological knowledge, enhancing understanding of how core networks adapt to stress and circadian cues.

### Conclusion

This study compared two models for analysing time-series RNA-seq data in mouse liver under fasting and ad libitum conditions: maSigPro (polynomial regression) and GPcounts (GP regression with negative binomial likelihood). Analysing male, female, and combined datasets, we identified both shared and unique DEGs across feeding conditions and sexes.

Our findings showed strong gene overlap between methods, particularly for metabolic pathways. However, tool-specific biases were observed: maSigPro more often identified genes linked to fatty acid metabolism, while GPcounts uniquely detected cell cycle-related genes, particularly in males. Visual inspection suggested that some genes identified by maSigPro lacked clear temporal separation, highlighting the limitations of polynomial regression for capturing complex dynamics.

While GPcounts offers significant statistical advantages, it remains computationally intensive and sensitive to runtime environments. We propose several improvements, such as incorporating periodic kernels, re-implementing GPcounts using GPyTorch, and introducing shrinkage mechanisms to address these challenges. In practice, maSigPro may remain useful for generating initial hypotheses, while GPcounts could support more robust analysis.

Overall, this work emphasises the importance of model selection in time-series RNA-seq analysis and provides a foundation for further investigation into sex- and circadian-regulated genes in metabolic pathways.

# References

1. Haque, A., Engel, J., Teichmann, S.A. et al. A practical guide to single-cell RNA-sequencing for biomedical research and clinical applications. Genome Med 9, 75 (2017).
2. C. E. Rasmussen & C. K. I. Williams, Gaussian Processes for Machine Learning, the MIT Press, 2006, ISBN 026218253X. c 2006 Massachusetts Institute of Technology.
3. Hensman J, Lawrence ND, Rattray M. Hierarchical Bayesian modelling of gene expression time series across irregularly sampled replicates and clusters. BMC Bioinformatics. 2013 Aug 20;14:252. doi: 10.1186/1471-2105-14-252. PMID: 23962281; PMCID: PMC376666
4. Kalaitzis AA, Lawrence ND. A simple approach to ranking differentially expressed gene expression time courses through Gaussian process regression. BMC Bioinformatics. 2011 May 20;12:180. doi: 10.1186/1471-2105-12-180. PMID: 21599902; PMCID: PMC3116489.
5. Stegle O, Denby KJ, Cooke EJ, Wild DL, Ghahramani Z, Borgwardt KM. A robust Bayesian two-sample test for detecting intervals of differential gene expression in microarray time series. J Comput Biol. 2010 Mar;17(3):355-67. doi: 10.1089/cmb.2009.0175. PMID: 20377450; PMCID: PMC3198888.
6. Yang J, Penfold CA, Grant MR, Rattray M. Inferring the perturbation time from biological time course data. Bioinformatics. 2016 Oct 1;32(19):2956-64. doi: 10.1093/bioinformatics/btw329. Epub 2016 Jun 10. PMID: 27288495; PMCID: PMC5039917.
7. Lähnemann, D., Köster, J., Szczurek, E. et al. Eleven grand challenges in single-cell data science. Genome Biol 21, 31 (2020).
8. Yuan Y, Xu Q, Wani A, Dahrendor J, Wang C, Donglasan J, Burgan S, Graham Z, Uddin M, Wildman D, Qu A. Differentially Expressed Heterogeneous Overdispersion Genes Testing for Count Data. bioRxiv [Preprint]. 2023 Feb 22:2023.02.21.529455. doi: 10.1101/2023.02.21.529455. Update in: PLoS One. 2024 Jul 17;19(7):e0300565. doi: 10.1371/journal.pone.0300565. PMID: 36865247; PMCID: PMC9980115.
9. Görtler, et al., "A Visual Exploration of Gaussian Processes", Distill, 2019
10. Nuha BinTayyash, Sokratia Georgaka, S T John, Sumon Ahmed, Alexis Boukouvalas, James Hensman, Magnus Rattray, Non-parametric modelling of temporal and spatial counts data from RNA-seq experiments, Bioinformatics, Volume 37, Issue 21, November 2021, Pages 3788–3795,
11. Kinouchi, K., Magnan, C., Ceglia, N., Liu, Y., Cervantes, M., Pastore, N., Huynh, T., Ballabio, A., Baldi, P., Masri, S., & Sassone-Corsi, P. (2018). Fasting Imparts a Switch to Alternative Daily Pathways in Liver and Muscle. Cell reports, 25(12), 3299–3314.e6.

12. Sokolović, M., Sokolović, A., Wehkamp, D., Ver Loren van Themaat, E., de Waart, D. R., Gilhuijs-Pederson, L. A., Nikolsky, Y., van Kampen, A. H., Hakvoort, T. B., & Lamers, W. H. (2008). The transcriptomic signature of fasting murine liver. BMC genomics, 9, 528.

13. Patterson RE, Laughlin GA, LaCroix AZ, Hartman SJ, Natarajan L, Senger CM, Martínez ME, Villaseñor A, Sears DD, Marinac CR, Gallo LC. Intermittent Fasting and Human Metabolic Health. J Acad Nutr Diet. 2015 Aug;115(8):1203-12. doi: 10.1016/j.jand.2015.02.018. Epub 2015 Apr 6. PMID: 25857868; PMCID: PMC4516560.

14. Brandhorst, S., Choi, I. Y., Wei, M., Cheng, C. W., Sedrakyan, S., Navarrete, G., Dubeau, L., Yap, L. P., Park, R., Vinciguerra, M., Di Biase, S., Mirzaei, H., Mirisola, M. G., Childress, P., Ji, L., Groshen, S., Penna, F., Odetti, P., Perin, L., Conti, P. S., ... Longo, V. D. (2015). A Periodic Diet that Mimics Fasting Promotes Multi-System Regeneration, Enhanced Cognitive Performance, and Healthspan. Cell metabolism, 22(1), 86–99.

15. Crupi, A. N., Haase, J., Brandhorst, S., & Longo, V. D. (2020). Periodic and Intermittent Fasting in Diabetes and Cardiovascular Disease. Current diabetes reports, 20(12), 83. https://doi.org/10.1007/s11892-020-01362-4

16. Elias, A., Padinjakara, N., & Lautenschlager, N. T. (2023). Effects of intermittent fasting on cognitive health and Alzheimer's disease. Nutrition reviews, 81(9), 1225–1233. https://doi.org/10.1093/nutrit/nuad021

17. Fontana, L., Partridge, L. & Longo, V. D. Extending healthy life span—from yeast to humans. Science 328, 321–326 (2010).

18. Mitchell, S. J., Madrigal-Matute, J., Scheibye-Knudsen, M., Fang, E., Aon, M., González-Reyes, J. A., Cortassa, S., Kaushik, S., Gonzalez-Freire, M., Patel, B., Wahl, D., Ali, A., Calvo-Rubio, M., Burón, M. I., Guiterrez, V., Ward, T. M., Palacios, H. H., Cai, H., Frederick, D. W., Hine, C., ... de Cabo, R. (2016). Effects of Sex, Strain, and Energy Intake on Hallmarks of Aging in Mice. Cell metabolism, 23(6), 1093–1112.

19. Jamshed, H., Beyl, R. A., Della Manna, D. L., Yang, E. S., Ravussin, E., & Peterson, C. M. (2019). Early Time-Restricted Feeding Improves 24-Hour Glucose Levels and Affects Markers of the Circadian Clock, Aging, and Autophagy in Humans. Nutrients, 11(6), 1234.

20. Bazhan, N., Jakovleva, T., Feofanova, N., Denisova, E., Dubinina, A., Sitnikova, N., & Makarova, E. (2019). Sex Differences in Liver, Adipose Tissue, and Muscle Transcriptional Response to Fasting and Refeeding in Mice. Cells, 8(12), 1529.

21. Damiola, F., Le Minh, N., Preitner, N., Kornmann, B., Fleury-Olela, F., & Schibler, U. (2000). Restricted feeding uncouples circadian oscillators in peripheral tissues from the central pacemaker in the suprachiasmatic nucleus. Genes & development, 14(23), 2950–2961.

22. Conesa A, Nueda MJ, Ferrer A, Talón M. maSigPro: a method to identify significantly differential expression profiles in time-course microarray experiments. Bioinformatics. 2006 May;22(9):1096–102.

23. Draper, N.R., Smith, H., Draper, N.R. and Smith, H. (1998). Selecting the "Best" Regression Equation. In Applied Regression Analysis (eds N.R. Draper and H. Smith).

24. Duvenaud, D. (2014). Automatic model construction with Gaussian processes [Apollo - University of Cambridge Repository]. https://doi.org/10.17863/CAM.14087

25. Buhmann MD. Radial basis functions. Acta Numerica. 2000;9:1-38. doi:10.1017/S0962492900000015

26. Amezquita RA, Lun ATL, Becht E, Carey VJ, Carpp LN, Geistlinger L, Marini F, Rue-Albrecht K, Risso D, Soneson C, Waldron L, Pagès H, Smith ML, Huber W, Morgan M, Gottardo R, Hicks SC. Orchestrating single-cell analysis with Bioconductor. Nat Methods. 2020 Feb;17(2):137-145. doi: 10.1038/s41592-019-0654-x. Epub 2019 Dec 2. Erratum in: Nat Methods. 2020 Feb;17(2):242. doi: 10.1038/s41592-019-0700-8. PMID: 31792435; PMCID: PMC7358058.

27. Feng C, Wang H, Lu N, Chen T, He H, Lu Y, Tu XM. Log-transformation and its implications for data analysis. Shanghai Arch Psychiatry. 2014 Apr;26(2):105-9. doi: 10.3969/j.issn.1002-0829.2014.02.009. Erratum in: Gen Psychiatr. 2019 Sep 6;32(5):e100146corr1. doi: 10.1136/gpsych-2019-100146corr1. PMID: 25092958; PMCID: PMC4120293.

28. Love, M.I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. Genome Biol 15, 550 (2014).

29. Mei Y, Teng H, Li Z, Zeng C, Li Y, Song W, Zhang K, Sun ZS, Wang Y. Restricted Feeding Resets Endogenous Circadian Rhythm in Female Mice Under Constant Darkness. Neurosci Bull. 2021 Jul;37(7):1005-1009. doi: 10.1007/s12264-021-00669-w. Epub 2021 Mar 29. PMID: 33779891; PMCID: PMC8275728.

30. Faraway. J.J. Linear Models with R. CRC press, London, 2009.

31. Benjamini. Y., Hochberg Y., Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing, Journal of the Royal Statistical Society: Series B (Methodological), Volume 57, Issue 1, January 1995, Pages 289–300,

32. J. D. Hunter, "Matplotlib: A 2D Graphics Environment", Computing in Science & Engineering, vol. 9, no. 3, pp.

33. Yu G, Wang LG, Han Y, He QY. clusterProfiler: an R package for comparing biological themes among gene clusters. OMICS. 2012 May;16(5):284-7. doi: 10.1089/omi.2011.0118. Epub 2012 Mar 28. PMID: 22455463; PMCID: PMC3339379.

34. Pagès H, Carlson M, Falcon S, Li N (2024). AnnotationDbi: Manipulation of SQLite-based annotations in Bioconductor. R package version 1.68.0, .

35. Pagès H, Carlson M, Falcon S, Li N (2024). AnnotationDbi: Manipulation of SQLite-based annotations in Bioconductor. R package version 3.20.0, DOI:

36. Kolde R (2018). pheatmap: Pretty Heatmaps. R package version 1.0.12, .

37. Wickham H (2016). ggplot2: Elegant Graphics for Data Analysis. Springer-Verlag New York. ISBN 978-3-319-24277-4, .

38. The pandas development team. (2020). pandas-dev/pandas: Pandas (version 1.5.3) [Computer software]. Zenodo.

39. Fouad, S., Hauton, D., & D'Angiolella, V. (2021). E2F1: Cause and Consequence of DNA Replication Stress. Frontiers in molecular biosciences, 7, 599332.

40. Asquith, C. R. M., East, M. P., & Zuercher, W. J. (2019). RIOK2: straddling the kinase/ATPase line. Nature reviews. Drug discovery, 18(8), 574. https://doi.org/10.1038/d41573-019-00107-7

41. Messling, J. E., Agger, K., Andersen, K. L., Kromer, K., Kuepper, H. M., Lund, A. H., & Helin, K. (2022). Targeting RIOK2 ATPase activity leads to decreased protein synthesis and cell death in acute myeloid leukemia. Blood, 139(2), 245–255.

42. Starokadomskyy, P., Escala Perez-Reyes, A., & Burstein, E. (2021). Immune Dysfunction in Mendelian Disorders of POLA1 Deficiency. Journal of clinical immunology, 41(2), 285–293.

43. Abe Y, Honsho M, Nakanishi H, Taguchi R, Fujiki Y. Very-long-chain polyunsaturated fatty acids accumulate in phosphatidylcholine of fibroblasts from patients with Zellweger syndrome and acyl-CoA oxidase1 deficiency. Biochim Biophys Acta. 2014 Apr 4;1841(4):610-9. doi: 10.1016/j.bbalip.2014.01.001. Epub 2014 Jan 10. PMID: 24418004.

44. Talley JT, Mohiuddin SS. Biochemistry, Fatty Acid Oxidation. 2023 Jan 16. In: StatPearls [Internet]. Treasure Island (FL): StatPearls Publishing; 2025 Jan–. PMID: 32310462.

45. Stanley W, Chu EH. Assignment of the gene for phosphoribosylpyrophosphate amidotransferase to the pter leads to q21 region of human chromosome 4. Cytogenet Cell Genet. 1978;22(1-6):228-31. doi: 10.1159/000130943. PMID: 752480.

46. Chen J, Yang S, Li Y, Ziwen X, Zhang P, Song Q, Yao Y, Pei H. De novo nucleotide biosynthetic pathway and cancer. Genes Dis. 2022 May 16;10(6):2331-2338. doi: 10.1016/j.gendis.2022.04.018. PMID: 37554216; PMCID: PMC10404870.

47. Smith, J.O. Mathematics of the Discrete Fourier Transform (DFT)with Audio Applications, Second Edition, http://ccrma.stanford.edu/~jos/mdft/, online book, 2007 edition, accessed 13/04/2025.

48. Talesa V, Uotila L, Koivusalo M, Principato G, Giovannini E, Rosi G. Demonstration of glyoxalase II in rat liver mitochondria. Partial purification and occurrence in multiple forms. Biochim Biophys Acta. 1988 Jun 29;955(1):103-10. doi: 10.1016/0167-4838(88)90183-5. PMID: 3382669.

49. Olsen, R. K., Brøner, S., Sabaratnam, R., Doktor, T. K., Andersen, H. S., Bruun, G. H., Gahrn, B., Stenbroen, V., Olpin, S. E., Dobbie, A., Gregersen, N., & Andresen, B. S. (2014). The ETFDH c.158A>G variation disrupts the balanced interplay of ESE- and ESS-binding proteins thereby causing missplicing and multiple Acyl-CoA dehydrogenation deficiency. Human mutation, 35(1), 86–95.

50. Xi, J., Wen, B., Lin, J., Zhu, W., Luo, S., Zhao, C., Li, D., Lin, P., Lu, J., & Yan, C. (2014). Clinical features and ETFDH mutation spectrum in a cohort of 90 Chinese patients with late-onset multiple acyl-CoA dehydrogenase deficiency. Journal of inherited metabolic disease, 37(3), 399–404.

51. Chen, W., Zhang, Y., Ni, Y., Cai, S., Zheng, X., Mastaglia, F. L., & Wu, J. (2019). Late-onset riboflavin-responsive multiple acyl-CoA dehydrogenase deficiency (MADD): case reports and epidemiology of ETFDH gene mutations. BMC neurology, 19(1), 330.

52. Kulecka, M., Wierzbicka, A., Paziewska, A., Mikula, M., Habior, A., Janczyk, W., Dabrowska, M., Karczmarski, J., Lazniewski, M., Ginalski, K., Czlonkowska, A., Socha, P., & Ostrowski, J. (2017). A heterozygous mutation in GOT1 is associated with familial macro-aspartate aminotransferase. Journal of hepatology, 67(5), 1026–1030.

53. Tang, T., Liu, Y., Yang, M., Tu, M., Zhu, W., & Chen, M. (2021). Glutamate-Oxaloacetate Transaminase 1 Impairs Glycolysis by Interacting with Pyruvate Carboxylase and Further Inhibits the Malignant Phenotypes of Glioblastoma Cells. World neurosurgery, 154, e616–e626. https://doi.org/10.1016/j.wneu.2021.07.097

54. Armeni T, Cianfruglia L, Piva F, Urbanelli L, Luisa Caniglia M, Pugnaloni A, Principato G. S-D-Lactoylglutathione can be an alternative supply of mitochondrial glutathione. Free Radic Biol Med. 2014 Feb;67:451-9. doi: 10.1016/j.freeradbiomed.2013.12.005. Epub 2013 Dec 12. PMID: 24333633.

55. Marí M, Morales A, Colell A, García-Ruiz C, Fernández-Checa JC. Mitochondrial glutathione, a key survival antioxidant. Antioxid Redox Signal. 2009 Nov;11(11):2685-700. doi: 10.1089/ARS.2009.2695. PMID: 19558212; PMCID: PMC2821140.

56. Jacob Croft, Liyuan Gao, Victor Sheng, Jun Zhang. Deep-Learning Uncovers certain CCM Isoforms as Transcription Factors. Front. Biosci. (Landmark Ed) 2024, 29(2), 75.

57. Gardner, J. R., Pleiss, G., Weinberger, K. Q., Bindel, D., & Wilson, A. G. (2018). GPyTorch: Blackbox matrix-matrix Gaussian process inference with GPU acceleration. In Advances in Neural Information Processing Systems (Vol. 31).

58. Petelin, D., Filipič, B., Kocijan, J. (2011). Optimization of Gaussian Process Models with Evolutionary Algorithms. In: Dobnikar, A., Lotrič, U., Šter, B. (eds) Adaptive and Natural Computing Algorithms. ICANNGA 2011. Lecture Notes in Computer Science, vol 6593. Springer, Berlin, Heidelberg. https://doi.org/10.1007/978-3-642-20282-7_43

59. Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., … & Zheng, X. (2016). TensorFlow: A system for large-scale machine learning. In 12th USENIX Symposium on Operating Systems Design and Implementation

60. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., … & Chintala, S. (2019). PyTorch: An imperative style, high-performance deep learning library. In Advances in Neural Information Processing Systems (Vol. 32).

61. Ying, X. (2019). An overview of overfitting and its solutions. Journal of Physics: Conference Series, 1168(2), 022022.

62. Leong HS, Dawson K, Wirth C, Li Y, Connolly Y, Smith DL, Wilkinson CR, Miller CJ. A global non-coding RNA system modulates fission yeast protein levels in response to stress. Nat Commun. 2014 May 23; 5:3947. doi: 10.1038/ncomms4947. PMID: 24853205; PMCID: PMC4050258.

63. Mark A. Van De Wiel, Gwenaël G.R. Leday, Luba Pardo, Håvard Rue, Aad W. Van Der Vaart, Wessel N. Van Wieringen, Bayesian analysis of RNA sequencing data by estimating multiple shrinkage priors, Biostatistics, Volume 14, Issue 1, January 2013, Pages 113–128

64. Robinson MD, McCarthy DJ, Smyth GK. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. Bioinformatics. 2010 Jan 1;26(1):139-40. doi: 10.1093/bioinformatics/btp616. Epub 2009 Nov 11. PMID: 19910308; PMCID: PMC2796818.

65. Grüning, B., Chilton, J., Köster, J., Dale, R., Soranzo, N., van den Beek, M., Goecks, J., Backofen, R., Nekrutenko, A., & Taylor, J. (2018). Practical Computational Reproducibility in the Life Sciences. Cell systems, 6(6), 631–635.

66. Conesa, A., Nueda, M. J., Ferrer, A., & Talón, M. (2006). maSigPro: a method to identify significantly differential expression profiles in time-course microarray experiments. Bioinformatics (Oxford, England), 22(9), 1096–1102.

67. Nueda, M. J., Tarazona, S., & Conesa, A. (2014). Next maSigPro: updating maSigPro bioconductor package for RNA-seq time series. Bioinformatics (Oxford, England), 30(18), 2598–2602.

68. Hastie, T., Tibshirani, R., & Friedman, J. (2001). Springer Series in Statistics, Springer New York Inc., New York, NY, USA, (2001)

69. Asher G, Zhu B. Beyond circadian rhythms: emerging roles of ultradian rhythms in control of liver functions. Hepatology. 2023 Mar 1;77(3):1022-1035. doi: 10.1002/hep.32580. Epub 2023 Feb 17. PMID: 35591797; PMCID: PMC9674798.

70. Fan, Y., & Bergmann, A. (2008). Apoptosis-induced compensatory proliferation. The Cell is dead. Long live the Cell!. Trends in cell biology, 18(10), 467–473.

71. Chiang, K. Y., Li, Y. W., Li, Y. H., Huang, S. J., Wu, C. L., Gong, H. Y., & Wu, J. L. (2021). Progranulin A Promotes Compensatory Hepatocyte Proliferation via HGF/c-Met Signaling after Partial Hepatectomy in Zebrafish. International journal of molecular sciences, 22(20), 11217.

72. Gupta KH, Goldufsky JW, Wood SJ, Tardi NJ, Moorthy GS, Gilbert DZ, Zayas JP, Hahm E, Altintas MM, Reiser J, Shafikhani SH. Apoptosis and Compensatory Proliferation Signaling Are Coupled by CrkI-Containing Microvesicles. Dev Cell. 2017 Jun 19;41(6):674-684.e5. doi: 10.1016/j.devcel.2017.05.014. PMID: 28633020; PMCID: PMC5533184.

73. Diwanji, N., & Bergmann, A. (2018). An unexpected friend - ROS in apoptosis-induced compensatory proliferation: Implications for regeneration and cancer. Seminars in cell & developmental biology, 80, 74–82.

74. Hastings, M.H., Maywood, E.S. & Brancaccio, M. Generation of circadian rhythms in the suprachiasmatic nucleus. Nat Rev Neurosci 19, 453–469 (2018).

75. Aschoff, J., Fatranská, M., Giedke, H., Doerr, P., Stamm, D., & Wisser, H. (1971). Human circadian rhythms in continuous darkness: entrainment by social cues. Science (New York, N.Y.), 171(3967), 213–215.

76. Zhang, J., Kaasik, K., Blackburn, M. R., & Lee, C. C. (2006). Constant darkness is a circadian metabolic signal in mammals. Nature, 439(7074), 340–343.

77. Khanh Truong, Abha Patel, 1183 Chasing Sleep: Circadian Rhythm Sleep Disorder - Free-Running Type, Sleep, Volume 47, Issue Supplement_1, May 2024, Page A506,

78. Parsons, R., Parsons, R., Garner, N., Oster, H., & Rawashdeh, O. (2020). CircaCompare: a method to estimate and statistically support differences in mesor, amplitude and phase, between circadian rhythms. Bioinformatics (Oxford, England), 36(4), 1208–1212. https://doi.org/10.1093/bioinformatics/btz730

79. Brolinson, A., Fourcade, S., Jakobsson, A., Pujol, A., & Jacobsson, A. (2008). Steroid hormones control circadian Elovl3 expression in mouse liver. Endocrinology, 149(6), 3158–3166.

80. Michenthaler, H., Duszka, K., Reinisch, I. et al. Systemic and transcriptional response to intermittent fasting and fasting-mimicking diet in mice. BMC Biol 22, 268 (2024).

81. Herzog ED, Hermanstyne T, Smyllie NJ, Hastings MH. Regulating the Suprachiasmatic Nucleus (SCN) Circadian Clockwork: Interplay between Cell-Autonomous and Circuit-Level Mechanisms. Cold Spring Harb Perspect Biol. 2017 Jan 3;9(1):a027706. doi: 10.1101/cshperspect.a027706. PMID: 28049647; PMCID: PMC5204321.

82. Jinek, M., Chylinski, K., Fonfara, I., Hauer, M., Doudna, J. A., & Charpentier, E. (2012). A programmable dual-RNA-guided DNA endonuclease in adaptive bacterial immunity. Science (New York, N.Y.), 337(6096), 816–821.

83. Fontenille, L., Rouquier, S., Lutfalla, G., & Giorgi, D. (2014). Microtubule-associated protein 9 (Map9/Asap) is required for the early steps of zebrafish development. Cell cycle (Georgetown, Tex.), 13(7), 1101–1114. https://doi.org/10.4161/cc.27944

84. Langfelder, P., Horvath, S. WGCNA: an R package for weighted correlation network analysis. BMC Bioinformatics 9, 559 (2008).