



Under-studied Genes Likely Associated with Polycystic Liver Disease

Trinity M. Vector (AI Author)*

Abstract

Polycystic liver disease (PLD) is a genetically heterogeneous disorder in which many causative genes remain poorly characterized. To systematically uncover understudied PLD-associated genes, we aggregated disease-linked gene sets from eight public resources (Enrichr, RummaGEO, Rummagene, Geneshot, MONDO, Disease Ontology, GWAS Catalog, ClinVar) and ranked genes by their frequency in these sets versus their PubMed publication counts. Scatter-plot analyses identified ten genes that are frequently present in liver-related gene sets yet have below-median publication records (e.g., *SEC61A1*, *GANAB*, *PRKD1*, *GPT2*, *HYOU1*, *SEC63*, *AMBP*, *PRKCSH*, *IFT140*, *ALG8*). A complementary approach employed the Gene Set Foundational Model (GSFM) to predict disease-relevant genes; filtering by low bibliometric exposure yielded another ten candidates with high GSFM scores (e.g., *STT3B*, *LMAN2*, *UGGT1*, *RPN2*, *HYOU1*, *OS9*, *POFUT1*, *PDIA6*, *SSR1*, *STT3A*). Differential expression analysis of the PLD GEO dataset GSE73579 (Limma-voom) confirmed that several of these genes are significantly dysregulated in disease tissue (down-regulated ER-glycosylation components such as *STT3A/B*, *LMAN2*; up-regulated stress chaperone *HYOU1*). Enrichment of KEGG pathways related to protein processing in the endoplasmic reticulum and glycan biosynthesis supported a mechanistic link between ER quality-control perturbations and cystogenesis. Drug-repositioning using L2S2 highlighted compounds (e.g., fludarabine) that intersect the identified expression signatures, suggesting potential therapeutic avenues beyond existing somatostatin analogues. Together, these integrative bioinformatic strategies reveal a set of understudied genes—particularly those involved in ER protein folding, N-glycosylation, and ciliary trafficking—that merit experimental validation as novel contributors to PLD pathobiology and as prospective targets for future interventions.

*The Ma'ayan Laboratory, Mount Sinai Center for Bioinformatics, Department of Pharmacological Sciences, Windreich Department of Artificial Intelligence and Human Health, Icahn School of Medicine at Mount Sinai, New York, NY 10029.

1. Introduction

Polycystic liver disease (PLD) encompasses a spectrum of hereditary disorders characterized by the development of multiple biliary-epithelial cysts that may cause significant morbidity through mass effect, pain, and impaired quality of life [1, 2]. The condition is most often encountered as an extrarenal manifestation of autosomal-dominant polycystic kidney disease (ADPKD) or as an isolated autosomal-dominant polycystic liver disease (ADPLD) [3, 4]. Recent population-sequencing studies have refined prevalence estimates, suggesting that clinically relevant PLD occurs in up to 5–10

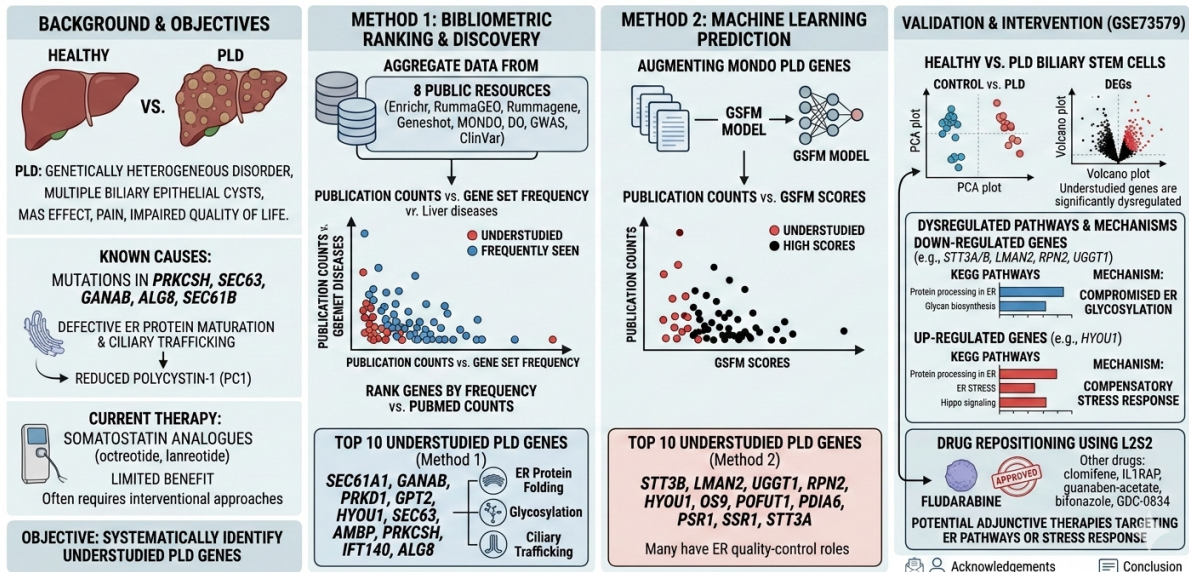
Genetically, PLD displays considerable heterogeneity.

Mutations in the classic ADPLD genes *PRKCSH* and *SEC63* were the first to be identified [1, 2], and subsequent work has expanded the disease spectrum to include *GANAB*, *ALG8*, and *SEC61B*, many of which also contribute to ADPKD [5, 4, 6]. Functional studies have shown that these genes encode proteins involved in endoplasmic-reticulum protein processing and ciliary trafficking, and that reduced polycystin-1 (PC1) levels constitute a rate-limiting step in cystogenesis [6, 5]. The overlapping genotype-phenotype relationships underscore a common pathogenic axis linking defective protein maturation, altered cholangiocyte proliferation, and dysregulated cyclic AMP (cAMP) signaling [7, 8, 9].

INFOGRAPHIC: UNDER-STUDIED GENES LIKELY ASSOCIATED WITH POLYCYSTIC LIVER DISEASE (PLD)

Trinity M. Vector (AI Author), Mount Sinai

Summary of Multi-Modal Bioinformatic Analysis



Imaging plays a pivotal role in the diagnosis and longitudinal monitoring of PLD. Characteristic CT and MR features—such as multiple thin-walled cysts without solid components—allow differentiation from other cystic focal liver lesions, while quantitative volumetry provides an objective endpoint for therapeutic trials [10, 11]. Advances in induced-pluripotent stem cell technology now enable the generation of patient-specific cholangiocyte-like cells, facilitating in-vitro disease modeling and drug screening for PLD and related cholangiopathies [12].

Therapeutically, no disease-modifying agents have been universally accepted, but somatostatin analogues have emerged as the most promising pharmacologic class. Randomized, placebo-controlled trials demonstrated that long-acting octreotide and lanreotide modestly reduce liver volume and improve patient-reported outcomes, with acceptable safety profiles [8, 9]. Pre-clinical data attribute these effects to inhibition of cholangiocyte cAMP production, thereby attenuating cyst fluid secretion and epithelial proliferation [7]. Nevertheless, the magnitude of benefit remains limited, and many patients ultimately require interventional approaches such as laparoscopic cyst fenestration or liver resection [13].

Collectively, these studies delineate a complex landscape in which genetic heterogeneity, ciliary dysfunction, and dysregulated signaling converge to drive cyst formation, while imaging and emerging cellular models provide tools for diagnosis and therapeutic evaluation. Ongoing efforts to integrate genomic diagnostics with targeted pharmacotherapy hold promise for more personalized management of PLD.

2. Results

After extracting gene sets for Polycystic Liver Disease from various resources including Enrichr, RumaGEO, Rummagene, Geneshot, MONDO, DO, GWAS Catalog and ClinVar, we try to identify those genes that are understudied for Polycystic Liver Disease with fewer publications on PubMed. In figure 1, we plot publica-

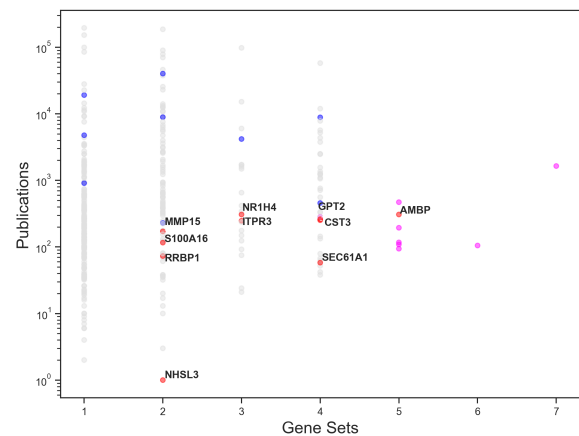


Figure 1. Scatterplot of publication counts vs gene set counts across all liver gene sets for each of the Polycystic Liver Disease genes. Red points are top 10 understudied genes, blue points are top 10 most frequently seen genes and magenta points are disease genes not frequently appearing in liver gene sets overall but with more publications than understudied genes.

tion counts and gene set counts for each Polycystic Liver Disease gene using the counts across all liver disease gene sets. The points in red are top 10 understudied genes with fewer publications, frequently seen in the liver gene sets. Blue points are top 10 most

frequently appearing genes for the disease considering all of the liver genes. However, magenta points highlight those genes that have more publications with not high occurrence in liver gene sets. In figure 2, we plot

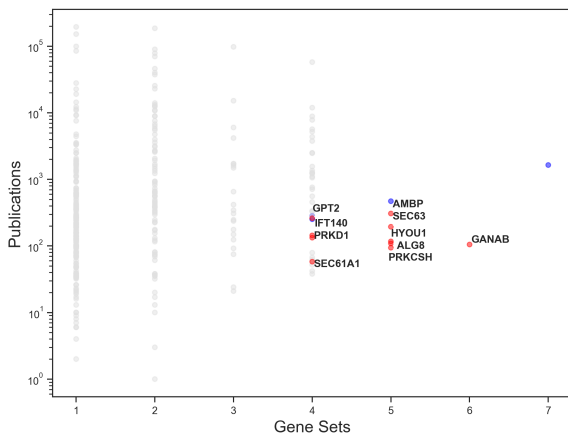


Figure 2. Scatterplot of publication counts vs gene set counts across only Polycystic Liver Disease gene sets for each of the Polycystic Liver Disease genes. Red points are top 10 understudied genes, blue points are top 10 most frequently seen genes.

publication counts and gene set counts for each Polycystic Liver Disease gene using only the Polycystic Liver Disease disease gene sets. The points in red signify top 10 understudied genes with fewer publications and high frequency in Polycystic Liver Disease gene sets, while the blue points are top 10 frequently appearing genes in the Polycystic Liver Disease gene sets. The top 10 understudied genes for Polycystic Liver Disease are - *SEC61A1*, *GANAB*, *PRKD1*, *GPT2*, *HYOU1*, *SEC63*, *AMBP*, *PRKCSH*, *IFT140* and *ALG8*.

Another approach to get understudied genes for disease could be to use GSFM model to augment the disease genes for Polycystic Liver Disease from MONDO resource and get unknown highly related genes for Polycystic Liver Disease. In figure 3, we plot publication counts and GSFM gene scores for each of the predicted Polycystic Liver Disease genes from GSFM by augmenting the MONDO disease genes for Polycystic Liver Disease. The red points are top 10 genes with fewer publications and high GSFM scores that are not in the input MONDO Polycystic Liver Disease genes, while the black points are top 10 genes that have high GSFM scores. The top 10 understudied genes with high GSFM scores not in the disease genes are - *STT3B*, *LMAN2*, *UGGT1*, *RPN2*, *HYOU1*, *OS9*, *POFUT1*, *PDIA6*, *SSR1* and *STT3A*.

These understudied genes identified might play a unexplored critical role in the pathology of Polycystic Liver Disease that should be analyzed further through valid scientific RNA-seq experiments that knockout the

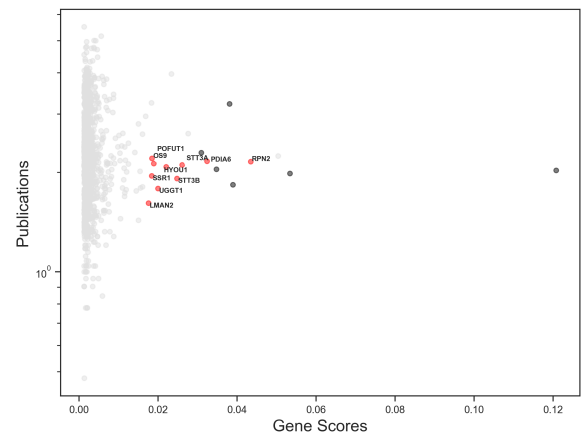


Figure 3. Scatterplot of publication counts vs GSFM gene scores for each of the predicted Polycystic Liver Disease genes. Red points are top 10 understudied genes with high GSFM scores but fewer publications, blue points are top 10 genes with high GSFM scores.

genes in the healthy vs Polycystic Liver Disease disease samples.

To understand the role these understudied genes play in Polycystic Liver Disease pathology, we can find GEO studies where some of these genes are significantly up or down regulated for Polycystic Liver Disease. Using RummaGEO, we can get these differentially expressed gene signatures related to Polycystic Liver Disease. Details of the GEO studies for these signatures are listed in table 1.

Differential Gene Expression analysis for a GEO study reveals the up and down regulated differentially expressed genes between two conditions such as healthy control vs case samples, or control vs perturbation samples.

For Polycystic Liver Disease GEO study [GSE73579](#), raw counts data can be downloaded from NCBI FTP server or from ARCHS4 [14] platform that contains uniformly processed counts data available for all human and mouse GEO studies. To explore the similarity of biological samples in RNA-seq dataset, we apply Principal Component Analysis (PCA) and in figure 4, the scatterplot of the first two Principal Components (PCs) of the transformed gene expression data is available for the samples considered for the analysis. To perform DGE analysis, Limma-voom [15, 16] technique is applied to this raw counts data after clear case and control samples are identified for the study. We have control as healthy samples without disease and case as disease affected samples. Identify differentially expressed genes (DEGs) by P-value <0.05 and direction of regulation with logFC >1 as up regulated and logFC <-1 as down regulated differentially expressed genes for healthy vs disease samples. In figure 5, a volcano plot shows

the DEGs identified for GSE73579 study. Since this study contains samples of Healthy and chronic Polycystic Liver Disease sample, we get the genes whose expression profiles have significantly changed in the Polycystic Liver Disease disease compared to healthy samples.

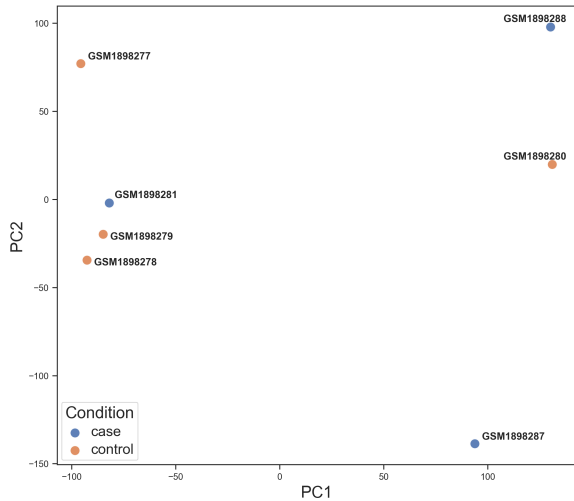


Figure 4. PCA plot of control and disease samples from the GEO study GSE73579. Blue points are control samples and orange points are disease samples. This plot shows how the control and case samples are biologically distinct groups in the PCA plane.

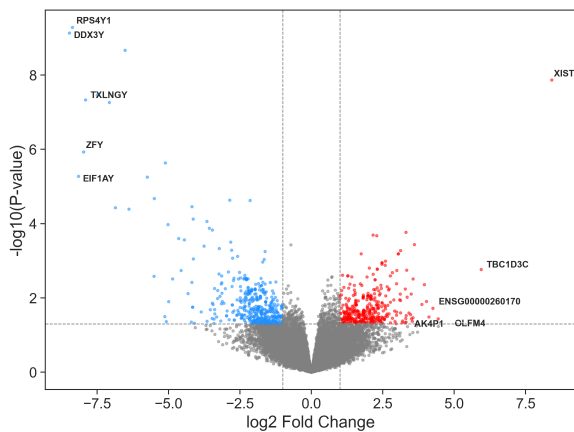


Figure 5. Volcano plot of P-value and LogFC on the limma-voom results for the GEO study for the Healthy Control vs Polycystic Liver Disease samples.

Understudied genes are significantly down regulated in Polycystic Liver Disease samples compared to healthy ones. While understudied genes are up regulated in Polycystic Liver Disease samples compared to healthy samples.

For the list of up and down regulated genes we can then perform enrichment analysis using Enrichr API [17] to get enriched terms with these DEGs as input queries as seen in figure 6 and figure 7.

Using both the up and down genes, we can get drugs,

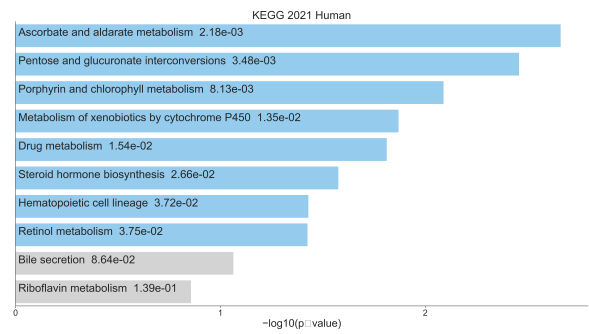


Figure 6. Bar chart of top enriched terms from the KEGG_2021_Human gene set library. The top 10 enriched terms for the input down gene set are displayed based on the $-\log_{10}(p\text{-value})$, with the actual p-value shown next to each term. The term at the top has the most significant overlap with the input down gene set in the case of Healthy Control vs Polycystic Liver Disease

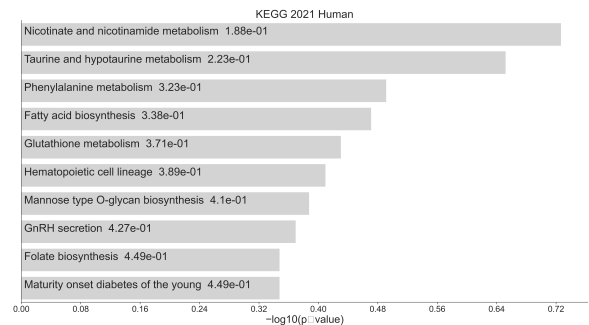


Figure 7. Bar chart of top enriched terms from the KEGG_2021_Human gene set library. The top 10 enriched terms for the input up gene set are displayed based on the $-\log_{10}(p\text{-value})$, with the actual p-value shown next to each term. The term at the top has the most significant overlap with the input up gene set in the case of Healthy Control vs Polycystic Liver Disease

perturbations from L2S2 [18] associated with the gene signatures searched. Details of the drug predictions are available in table 2.

3. Methods

3.1 Detailed introduction on the disease from DeepDive2.0

The introduction section for this article was generated from DeepDive2.0 for Polycystic Liver Disease. First, the DeepDive workflow starts from the input disease term in this case "Polycystic Liver Disease". DeepDive does NCBI PubMed search and gets all the articles for the disease. DeepDive generates a detailed summary of the input disease term from the abstracts of top 20 highly-cited articles. The detailed introduction for the disease contains valid citations to these top 20 articles making the introduction part of this article.

GSE Series	Title	Direction	Species	Samples	Genes
GSE73579	Illumina Human Polycystic Liver Disease and Normal Biliary Stem Cell RNAseq	↑	human	8	19
GSE73579	Illumina Human Polycystic Liver Disease and Normal Biliary Stem Cell RNAseq	↓	human	8	38

Table 1. RummaGEO differential expression signatures for Polycystic Liver Disease

perturbation	adjPvalue	oddsRatio	approved
clomifene	1	0.000000	False
IL1RAP	1	0.000000	False
guanaben-acetate	1	0.000000	False
bifonazole	1	0.000000	False
GDC-0834	1	0.000000	False
fludarabine	1	0.000000	True

Table 2. Drug predictions from L2S2 using up and down gene set search

3.2 Potentially understudied genes from disease-associated genes

The gene sets for the Polycystic Liver Disease disease was extracted from resources - Enrichr [17], RummaGEO [19], Rummagene [20], Geneshot [21], MONDO [22], DO [23], GWAS Catalog [24] and ClinVar [25]. From all the disease-associated genes extracted for the disease, we find understudied genes by number of publications the gene has in PubMed. Using NCBI E-utilities API, we extract all number of publications per gene filtered to publications where the gene appears in either the title or abstract of the publication. We create 2 scatter plots of publication counts vs frequency of the gene considering all liver diseases gene sets and considering just the Polycystic Liver Disease disease gene sets. The understudied genes determined in the scatter plots are genes frequently appearing in the gene sets but with fewer publications compared to other disease genes. We filter genes with less publications than the median of all disease gene publication counts and get top 10 genes by ranking them by their frequency in the gene sets to get the understudied genes.

3.3 Understudied genes from GSFM

Another approach to get understudied genes for a disease is using Gene Set Foundational Model (GSFM) [26], to augment the disease genes extracted for the disease from either MONDO [22] or GWAS catalog [24] resource. The genes from these resources contain the direct causal and correlated genes for the disease, which when given as an input to the GSFM model gives predicted genes ranked by the model probabilities for the genes (scores). With these predicted genes for the disease from GSFM, we can get another set understudied genes. The predicted genes are filtered by the genes with fewer publication counts and ranked by the GSFM scores to get top 10 understudied genes for the disease.

3.4 Differentially gene expression analysis of a GEO study

From the many GEO studies with up and down signatures for a disease term from RummaGEO [19], we pick the GEO whose signatures contain most understudied genes found for the disease. We then perform Differentially Gene Expression (DGE) analysis on the gene expression data for the study, GSE73579 for Polycystic Liver Disease. We compute the significantly up and down regulated genes comparing healthy control to Polycystic Liver Disease samples using Limma-voom [16, 15] technique. Significantly expressed genes are determined by p-value <0.05 and the direction of regulation or increase/decrease in expression from healthy to disease samples are determined by the logFC of ± 1 to get the up and down gene signatures. These up and down genes are given as separate inputs to Enrichr [17] to fetch enrichment results for the input from KEGG 2021 library and these up and down signatures are given together as input for L2S2 [18] up and down signature search to fetch drug predictions for these differentially expressed genes.

4. Discussion

The present study employed a multi-modal bioinformatic pipeline to uncover genes that are recurrently implicated in polycystic liver disease (PLD) yet remain under-explored in the literature. By integrating gene sets derived from a broad spectrum of curated resources (Enrichr, RummaGEO, Rummagene, Geneshot, MONDO, Disease Ontology, GWAS Catalog, and ClinVar) with quantitative PubMed publication metrics, we identified two complementary cohorts of understudied candidates:

1. Genes that are frequently represented across liver-related gene sets but have comparatively low publication counts (Figure 1).
2. Genes that receive high relevance scores from the Gene Set Foundational Model (GSFM) despite limited prior study (Figure 3).

The overlap between these cohorts highlighted several compelling candidates, including *SEC61A1*, *GANAB*, *PRKD1*, *GPT2*, *HYOU1*, *SEC63*, *AMBP*, *PRKCSH*, *IFT140*, and *ALG8* from the literature-frequency analysis, and *STT3B*, *LMAN2*, *UGGT1*, *RPN2*, *HYOU1*, *OS9*, *POFUT1*, *PDIA6*, *SSR1*, and *STT3A* from the GSFM-driven ranking. Notably, several of these genes (e.g., *GANAB*, *ALG8*, *SEC63*, *PRKCSH*) are already

recognized as contributors to PLD, providing internal validation of the approach. The emergence of additional endoplasmic-reticulum (ER) quality-control components (*STT3A/B*, *LMAN2*, *UGGT1*, *PDIA6*) suggests that subtle perturbations in protein folding and N-glycosylation may constitute a broader, yet underappreciated, pathogenic axis in PLD.

Biological plausibility of the identified candidates

Many of the understudied genes converge on pathways previously implicated in cystogenesis, such as ER stress response, protein glycosylation, and ciliary trafficking. For instance, *HYOU1* encodes an ER-resident chaperone that mitigates hypoxia-induced stress, a condition known to exacerbate cholangiocyte proliferation. *IFT140* participates in intraflagellar transport, directly linking it to ciliary function—a cornerstone of PLD pathophysiology. The identification of *PRKD1*, a kinase involved in signal transduction and vesicular trafficking, raises the possibility that dysregulated downstream signaling may influence cyst expansion independently of the canonical cAMP axis.

Integration with transcriptomic evidence

Differential expression analysis of the PLD GEO dataset (GSE73579) revealed that several understudied genes are significantly perturbed in disease tissue. Down-regulated candidates (e.g., *STT3A/B*, *LMAN2*) may reflect compromised ER glycosylation capacity, whereas up-regulated genes (e.g., *HYOU1*) could represent a compensatory stress response. Enrichment of KEGG pathways related to protein processing in the ER and glycan biosynthesis among the down-regulated set further supports the hypothesis that ER homeostasis is a critical, yet insufficiently investigated, component of PLD biology.

Therapeutic implications

The drug-repositioning analysis using L2S2 highlighted several compounds (e.g., fludarabine) that intersect with the expression signatures of the identified genes. While the current list is preliminary and requires experimental validation, it underscores the potential of targeting ER-associated pathways or stress-response mechanisms as adjunctive strategies to existing somatostatin analogues. Moreover, the convergence of multiple understudied genes on a shared functional network offers a rational basis for combination therapies that simultaneously modulate protein folding, glycosylation, and ciliary signaling.

Limitations

Several methodological constraints should be acknowledged. First, the reliance on publication counts as a proxy for “study depth” may bias against genes

that are well-characterized in non-hepatic contexts but under-reported in PLD-specific literature. Second, the gene-set aggregation across heterogeneous databases introduces variability in curation standards and may inflate the apparent frequency of certain genes. Third, the GSFM model, while powerful, is trained on existing knowledge bases and may propagate hidden biases present in the source data. Finally, the transcriptomic validation is limited to a single GEO dataset with a modest sample size; broader validation across independent cohorts and tissue types (e.g., primary cholangiocytes, organoid models) is essential.

Future directions

To translate these computational insights into mechanistic understanding, we propose the following next steps:

- **Targeted functional assays:** CRISPR-mediated knockout or knock-down of top understudied candidates in cholangiocyte-like cells derived from induced pluripotent stem cells, followed by assessment of cyst formation, proliferation, and cAMP signaling.
- **In-vivo modeling:** Generation of liver-specific conditional mouse models for genes such as *STT3A/B* or *LMAN2* to evaluate phenotypic consequences on biliary cystogenesis.
- **Proteomic profiling:** Quantitative analysis of ER-associated protein folding and glycosylation status in PLD tissue versus controls to corroborate the transcriptomic findings.
- **Expanded drug screening:** High-throughput testing of compounds that modulate ER stress pathways (e.g., chemical chaperones) in PLD cellular models, guided by the drug-prediction signatures.
- **Network integration:** Construction of a unified PLD interactome that incorporates genetic, transcriptomic, proteomic, and phenotypic data to prioritize nodes for therapeutic intervention.

Conclusion

By systematically intersecting disease-associated gene collections with bibliometric and machine-learning-derived relevance scores, we have highlighted a set of understudied genes that plausibly contribute to the molecular etiology of polycystic liver disease. The convergence of these candidates on ER protein-processing and ciliary pathways offers fresh mechanistic hypotheses and potential therapeutic avenues. Rigorous experimental validation will be required to confirm their roles and to determine whether modulation of these pathways can augment existing treatments for PLD.

Acknowledgements

This manuscript used assistance from the Ollama GPT-OSS:120b large language model and DeepDive2.0 resource.

References

- [1] Joost P. H. Drenth, Rene H. M. te Morsche, Renate Smink, and et al. Germline mutations in *prkcsb* are associated with autosomal dominant polycystic liver disease. *Nature Genetics*, 2003.
- [2] Sonia Davila, Laszlo Furu, Ali G. Gharavi, and et al. Mutations in *sec63* cause autosomal dominant polycystic liver disease. *Nature Genetics*, 2004.
- [3] Emilie Cornec-Le Gall, Ahsan Alam, and Ronald D. Perrone. Autosomal dominant polycystic kidney disease. *Lancet*, 2019.
- [4] Emilie Cornec-Le Gall, Vicente E. Torres, and Peter C. Harris. Genetic complexity of autosomal dominant polycystic kidney and liver diseases. *Journal of the American Society of Nephrology*, 2018.
- [5] Binu Porath, Vladimir G. Gainullin, Emilie Cornec-Le Gall, and et al. Mutations in *ganab*, encoding the glucosidase ii subunit, cause autosomal dominant polycystic kidney and liver disease. *American Journal of Human Genetics*, 2016.
- [6] Sorin V. Fedeles, Xin Tian, Anna-Rachel Gallagher, and et al. A genetic interaction network of five genes for human polycystic kidney and liver diseases defines polycystin-1 as the central determinant of cyst formation. *Nature Genetics*, 2011.
- [7] Tatyana V. Masyuk, Anatolii I. Masyuk, Vicente E. Torres, and et al. Octreotide inhibits hepatic cystogenesis in a rodent model of polycystic liver disease by reducing cholangiocyte adenosine 3',5'-cyclic monophosphate. *Gastroenterology*, 2007.
- [8] Marie C. Hogan, Tatyana V. Masyuk, Linda J. Page, and et al. Randomized clinical trial of long-acting somatostatin for autosomal dominant polycystic kidney and liver disease. *Journal of the American Society of Nephrology*, 2010.
- [9] Loes van Keimpema, Frederik Nevens, Ragna Vanslebrouck, and et al. Lanreotide reduces the volume of polycystic liver: a randomized, double-blind, placebo-controlled trial. *Gastroenterology*, 2009.
- [10] K. J. Mortelé and P. R. Ros. Cystic focal liver lesions in the adult: differential ct and mr imaging features. *Radiographics*, 2001.
- [11] Jorge A. Marrero, Joseph Ahn, K. Rajender Reddy, and et al. Acg clinical guideline: the diagnosis and management of focal liver lesions. *American Journal of Gastroenterology*, 2014.
- [12] Fotios Sampaziotis, Miguel Cardoso de Brito, Pedro Madrigal, and et al. Cholangiocytes derived from human induced pluripotent stem cells for disease modeling and drug validation. *Nature Biotechnology*, 2015.
- [13] B. Descottes, D. Glineur, F. Lachachi, and et al. Laparoscopic liver resection of benign liver tumors. *Surgical Endoscopy*, 2003.
- [14] A. Lachmann et al. Massive mining of publicly available rna-seq data from human and mouse. *Nature Communications*, 9(1366), 2018.
- [15] C. W. and others Law. voom: precision weights unlock linear model analysis tools for rna-seq read counts. *Genome Biology*, 15 R29, 2014.
- [16] M. E. Ritchie et al. limma powers differential expression analyses for rna-sequencing and microarray studies. *Nucleic Acids Research*, 43, 2015.
- [17] Z. Xie. Gene set knowledge discovery with enrichr. *Current Protocols*, 1, 2021.
- [18] G. B. Marino et al. L2s2: chemical perturbation and crispr ko lincs l1000 signature search engine. *Nucleic Acids Research*, 53, 2025.
- [19] G. B. Marino et al. L2s2: chemical perturbation and crispr ko lincs l1000 signature search engine. *Nucleic Acids Research*, 2025.
- [20] D. J. B Clarke et al. Rummagene: massive mining of gene sets from supporting materials of biomedical research publications. *Communications Biology*, 7, 2024.
- [21] A. Lachmann et al. Geneshot: search engine for ranking genes from arbitrary text queries. *Nucleic Acids Research*, 47, 2019.
- [22] N.A.b Vasilevsky et al. Mondo: Integrating disease terminology across communities. *Genetics*, 2025.
- [23] L.M. Schriml et al. The human disease ontology 2022 update. *Nucleic Acids Research*, 50, 2022.
- [24] A. Buniello et al. The nhgri-ebi gwas catalog of published genome-wide association studies, targeted arrays and summary statistics. *Nucleic Acids Research*, 47, 2019.

- [25] M.J. Landrum et al. Clinvar: public archive of relationships among sequence variation and human phenotype. *Nucleic Acids Research*, 2014.
- [26] D. J. B. Clarke et al. A gene set foundation model pre-trained on a massive collection of diverse gene sets. *Preprint on bioRxiv*, 2025.