

# GigaScience

## Proteome-wide association study and functional validation identify novel protein markers for pancreatic ductal adenocarcinoma --Manuscript Draft--

<b>Manuscript Number:</b>	GIGA-D-23-00321	
<b>Full Title:</b>	Proteome-wide association study and functional validation identify novel protein markers for pancreatic ductal adenocarcinoma	
<b>Article Type:</b>	Research	
<b>Funding Information:</b>	NIH/NCI (9 SC1 GM135050-05)	Dr. Yong Wu
	University of Hawaii Cancer Center and V Foundation V Scholar Award	Not applicable
	National Cancer Institute (R01CA263494)	Dr. Chong Wu Associated Professor Lang Wu
	National Human Genome Research Institute (U54HG013243)	Associated Professor Lang Wu
	National Cancer Institute (R00CA218892)	Associated Professor Lang Wu
	National Institute on Minority Health and Health Disparities (U54MD007598)	Not applicable
	NIH/NCI (1U54CA14393)	Not applicable
	NIH/NCI (U56 CA101599-01)	Not applicable
	Department-of-Defense Breast Cancer Research Program (BC043180)	Dr. Jaydutt V. Vadgama
	NIH/NCATS (CTSI UL1TR000124)	Dr. Jaydutt V. Vadgama
	Accelerating Excellence in Translational Science Pilot Grants (G0812D05)	Dr. Yong Wu
	NIH/NCI (SC1CA200517)	Dr. Yong Wu
	VA Merit Award (1 I01 CX001822-01A2)	Dr. Qizhi Yao
	National Cancer Institute (NCI), US National Institutes of Health (NIH) (HHSN261200800001E)	Not applicable
	NIH/NCI (K07 CA140790)	Not applicable
	the American Society of Clinical Oncology Conquer Cancer Foundation	Not applicable
	the Howard Hughes Medical Institute	Not applicable
	the Lustgarten Foundation	Not applicable
	he Robert T. and Judith B. Hale Fund for Pancreatic Cancer Research	Not applicable
	Promises for Purple	Not applicable
	NCI (R01CA154823)	Not applicable
	National Institutes of Health to The Johns Hopkins University (HHSN2682011000111)	Not applicable
	National Institute for Health Research (NIHR)	Not applicable
	NIHR BioResource	Not applicable

	NIHR Cambridge Biomedical Research Centre (BRC-1215-20014)	Not applicable
	NIHR Blood and Transplant Research Unit in Donor Health and Genomics (NIHR BTRU-2014-10024)	Not applicable
	UK Medical Research Council (MR/L003120/1)	Not applicable
	British Heart Foundation (SP/09/002; RG/13/13/30194; RG/18/13/33946)	Not applicable
	NIHR Cambridge BRC (BRC-1215-20014)	Not applicable
<b>Abstract:</b>	<p><b>Abstract</b>  Pancreatic ductal adenocarcinoma (PDAC) remains a lethal malignancy, largely due to the paucity of reliable biomarkers for early detection and therapeutic targeting. Existing blood protein biomarkers for PDAC often suffer from replicability issues, arising from inherent limitations such as unmeasured confounding factors in conventional epidemiologic study designs. To circumvent these limitations, we use genetic instruments to identify proteins with genetically predicted levels to be associated with PDAC risk. Leveraging genome and plasma proteome data from the INTERVAL study, we established and validated models to predict protein levels using genetic variants. By examining 8,275 PDAC cases and 6,723 controls, we identified 40 associated proteins, of which 16 are novel. Functionally validating these candidates by focusing on two selected novel protein-encoding genes, GOLMA1 and B4GALT1, we demonstrated their pivotal roles in driving PDAC cell proliferation, migration, and invasion. Furthermore, we also identified potential drug repurposing opportunities for treating PDAC.</p> <p><b>Significance:</b>  PDAC is a notoriously difficult-to-treat malignancy, and our limited understanding of causal protein markers hampers progress in developing effective early detection strategies and treatments. Our study identifies novel causal proteins using genetic instruments and subsequently functionally validates selected novel proteins. This dual approach enhances our understanding of PDAC etiology and potentially opens new avenues for therapeutic interventions.</p> <p><b>Keywords:</b> Biomarkers, protein, genetics, pancreatic cancer, risk</p>	
<b>Corresponding Author:</b>	Lang Wu University of Hawai'i at Manoa Honolulu, UNITED STATES	
<b>Corresponding Author Secondary Information:</b>		
<b>Corresponding Author's Institution:</b>	University of Hawai'i at Manoa	
<b>Corresponding Author's Secondary Institution:</b>		
<b>First Author:</b>	Jingjing Zhu	
<b>First Author Secondary Information:</b>		
<b>Order of Authors:</b>	Jingjing Zhu	
	Ke Wu	
	Shuai Liu	
	Alexandra Masca	
	Hua Zhong	
	Tai Yang	
	Dalia H Ghoneim	
	Praveen Surendran	
	Tanxin Liu	

	Qizhi Yao
	Tao Liu
	Sarah Fahle
	Adam Butterworth
	Md Ashad Alam
	Jaydutt V. Vadgama
	Youping Deng
	Hong-Wen Deng
	Chong Wu
	Yong Wu
	Lang Wu
<b>Order of Authors Secondary Information:</b>	
<b>Additional Information:</b>	
<b>Question</b>	<b>Response</b>
Are you submitting this manuscript to a special series or article collection?	No
<p><b>Experimental design and statistics</b></p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
<p><b>Resources</b></p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our <a href="#">Minimum</a></p>	Yes

<a href="#">Standards Reporting Checklist?</a>	
<p><b>Availability of data and materials</b></p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in <a href="#">publicly available repositories</a> (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our <a href="#">Minimum Standards Reporting Checklist?</a></p>	Yes

1 **Proteome-wide association study and functional validation identify novel protein markers**  
2 **for pancreatic ductal adenocarcinoma**

3  
4 Jingjing Zhu<sup>1\*</sup>, Ke Wu<sup>2\*</sup>, Shuai Liu<sup>1\*</sup>, Alexandra Masca<sup>1</sup>, Hua Zhong<sup>1</sup>, Tai Yang<sup>3</sup>, Dalia H  
5 Ghoneim<sup>1</sup>, Praveen Surendran<sup>4</sup>, Tanxin Liu<sup>5</sup>, Qizhi Yao<sup>6,7</sup>, Tao Liu<sup>8</sup>, Sarah Fahle<sup>4</sup>, Adam  
6 Butterworth<sup>4,9</sup>, Md Ashad Alam<sup>10</sup>, Jaydutt V. Vadgama<sup>2</sup>, Youping Deng<sup>11</sup>, Hong-Wen Deng<sup>10</sup>,  
7 Chong Wu<sup>12#</sup>, Yong Wu<sup>2#</sup>, Lang Wu<sup>1#</sup>

8  
9 1. Cancer Epidemiology Division, Population Sciences in the Pacific Program, University of  
10 Hawaii Cancer Center, University of Hawaii at Manoa, Honolulu, HI, USA

11 2. Division of Cancer Research and Training, Department of Internal Medicine, Charles R. Drew  
12 University of Medicine and Science, David Geffen UCLA School of Medicine and UCLA  
13 Jonsson Comprehensive Cancer Center, Los Angeles, CA 90095, USA

14 3. Department of Biostatistics, University of Michigan - Ann Arbor, MI, USA

15 4. MRC/BHF Cardiovascular Epidemiology Unit, Department of Public Health and Primary  
16 Care, University of Cambridge, Cambridge, UK

17 5. Department of Epidemiology, Johns Hopkins Bloomberg School of Public Health, Baltimore,  
18 MD, USA

19 6. Division of Surgical Oncology, Michael E. DeBakey Department of Surgery, Baylor College  
20 of Medicine, Houston, Texas

21 7. Center for Translational Research on Inflammatory Diseases (CTRID), Michael E. DeBakey  
22 VA Medical Center, Houston, Texas

23 8. Biological Sciences Division, Pacific Northwest National Laboratory, Richland, WA 99354,  
24 USA

25 9. NIHR Blood and Transplant Research Unit in Donor Health and Genomics, Department of  
26 Public Health and Primary Care, University of Cambridge, Cambridge, UK

27 10. Tulane Center for Biomedical Informatics and Genomics, Division of Biomedical  
28 Informatics and Genomics, Deming Department of Medicine, Tulane University, 1440 Canal  
29 Street, New Orleans, 70112, LA, USA

30 11. Department of Quantitative Health Sciences, John A. Burns School of Medicine, University  
31 of Hawaii at Manoa, Honolulu, HI, USA.

32 12. Department of Biostatistics, The University of Texas MD Anderson Cancer Center, Houston,  
33 TX, USA

34 \* these authors contributed equally to this work and are co-first authors

35 # these authors jointly supervised this work and are co-senior authors

36

37 **Running title:** Predicted protein biomarkers for pancreatic cancer

38

39 **Abbreviations list:**

40 Pancreatic ductal adenocarcinoma (PDAC)

41 protein quantitative trait loci (pQTL)

42 Genome-wide association studies (GWAS)

43 the Pancreatic Cancer Cohort Consortium (PanScan)

44 the Pancreatic Cancer Case-Control Consortium (PanC4)

45 quality control (QC)

46 Hardy-Weinberg equilibrium (HWE)

47 false discovery rate (FDR)

48

49 **Corresponding to:** Lang Wu, Cancer Epidemiology Division, Population Sciences in the Pacific  
50 Program, University of Hawaii Cancer Center, University of Hawaii at Manoa, Honolulu, HI,  
51 96813, USA. Email: [lwu@cc.hawaii.edu](mailto:lwu@cc.hawaii.edu). Phone: (808)564-5965; or Yong Wu, Department of  
52 Internal Medicine, Charles Drew University of Medicine and Science, Los Angeles, CA 90059,  
53 USA. Email: [yongwu@cdrewu.edu](mailto:yongwu@cdrewu.edu); or Chong Wu, Department of Biostatistics, The University  
54 of Texas MD Anderson Cancer Center, Houston, TX, USA. Email: [cwu18@mdanderson.org](mailto:cwu18@mdanderson.org)

55

56 **Competing financial interests**

57 L.W. provided consulting service to Pupil Bio Inc. and received honorarium. No potential  
58 conflicts of interest were disclosed by the other authors.

59

60 **Author contributions**

61 L.W. conceived the study. Y.W. designed the functional experiments and supervised the *in vitro*

62 functional work. C.W. and J.Z. contributed to the study design and/or prediction model building.

63 S.L. performed model building and statistical analyses. D.H.G. contributed to statistical analyses.

64 K.W. conducted *in vitro* functional work. J.Z. performed the drug repurposing curation. M.A.A.  
65 performed molecular docking analysis. H.Z. and S. L. contributed to the bioinformatics and  
66 pathway analyses. L.W., J.Z., K.W., Y.W., A.M., H.Z., and T.Y. wrote the first version of  
67 manuscript. D.H.G., P.S., T.L., E.P., Q.Y., T.L., S.F., J.V.V., H-W. D., Y.D., H.Z., S.L., and  
68 A.B. contributed to manuscript revision and/or INTERVAL data management. All authors have  
69 reviewed and approved the final manuscript.

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84 **Abstract**

85 Pancreatic ductal adenocarcinoma (PDAC) remains a lethal malignancy, largely due to the  
86 paucity of reliable biomarkers for early detection and therapeutic targeting. Existing blood  
87 protein biomarkers for PDAC often suffer from replicability issues, arising from inherent  
88 limitations such as unmeasured confounding factors in conventional epidemiologic study  
89 designs. To circumvent these limitations, we use genetic instruments to identify proteins with  
90 genetically predicted levels to be associated with PDAC risk. Leveraging genome and plasma  
91 proteome data from the INTERVAL study, we established and validated models to predict  
92 protein levels using genetic variants. By examining 8,275 PDAC cases and 6,723 controls, we  
93 identified 40 associated proteins, of which 16 are novel. Functionally validating these candidates  
94 by focusing on two selected novel protein-encoding genes, *GOLM1* and *B4GALT1*, we  
95 demonstrated their pivotal roles in driving PDAC cell proliferation, migration, and invasion.  
96 Furthermore, we also identified potential drug repurposing opportunities for treating PDAC.

97 **Significance:**

98 PDAC is a notoriously difficult-to-treat malignancy, and our limited understanding of causal  
99 protein markers hampers progress in developing effective early detection strategies and  
100 treatments. Our study identifies novel causal proteins using genetic instruments and subsequently  
101 functionally validates selected novel proteins. This dual approach enhances our understanding of  
102 PDAC etiology and potentially opens new avenues for therapeutic interventions.

103 **Keywords:** Biomarkers, protein, genetics, pancreatic cancer, risk

104

105



106

**107 Introduction**

108 Pancreatic cancer is the seventh leading cause of cancer deaths in industrialized countries  
109 with pancreatic ductal adenocarcinoma (PDAC) making up over 90% of pancreatic cancer cases  
110 (1). According to GLOBOCAN 2020 cancer statistics, pancreatic cancer is the 14th most  
111 common cancer type with 495,773 new cases in 2020. There are almost the same number of  
112 deaths caused by pancreatic cancer (466,003 deaths) in 2020, accounting for 4.7% of all cancer  
113 related deaths (2). Owing to its often asymptomatic or non-specific symptoms during early  
114 stages, a majority of patients are usually diagnosed in advanced stages. This results in 80-90% of  
115 pancreatic tumors being unresectable upon diagnosis, leading to a dismal prognosis: a mere 9%  
116 five-year survival rate after diagnosis (1). Given these dire statistics, there is an urgent need to  
117 identify effective biomarkers for screening or early detection in high-risk populations. Equally  
118 crucial is the development of improved therapeutic strategies to improve PDAC outcome.

119 Currently, serum cancer antigen (CA) 19-9 is the only diagnostic biomarker for  
120 pancreatic cancer approved by the U.S. FDA. However, elevated levels of CA 19-9 are related to  
121 other conditions, and its performance as a diagnostic tool for pancreatic cancer is far from ideal  
122 (3): it has a poor positive predictive value (0.5-0.9%), along with restricted specificity (82-90%)  
123 and sensitivity (79-81%). Previous studies have also reported several other circulating blood  
124 protein biomarkers that are potentially associated with pancreatic cancer risk, such as CA242,  
125 PIVKA-II, and PAM4 (4-7). However, results from existing studies often involving small sample  
126 sizes and findings are inconsistent. It is well known that the conventional epidemiologic study  
127 design measuring levels of proteins directly may be subject to selection bias and residual or

128 unmeasured confounding, which could also contribute to the inconsistent findings in the existing  
129 literature.

130 An alternative design of using genetic instruments may decrease many limitations of  
131 existing studies, due to the nature of random assortment of alleles from parents to offspring  
132 during gamete formation (8,9). Inspired by transcriptome-wide association study (TWAS), one  
133 may build comprehensive genetic prediction models for each protein to capture the prediction  
134 value of multiple single nucleotide polymorphisms (SNPs). Unlike conventional TWAS type of  
135 methods, which typically focus solely on cis-acting variants, our study enhanced statistical power  
136 by integrating both cis- and trans-acting elements into our genetic prediction models.  
137 Furthermore, as TWAS or PWAS results imply causality under stringent valid instrumental  
138 variable assumptions, we further functionally validated two novel proteins.

139 In the current study, we applied such a study design to identify novel proteins associated  
140 with PDAC risk. To our knowledge, this is the first large-scale proteome wide association study  
141 (PWAS) using comprehensive protein genetic prediction models as instruments to assess the  
142 associations between genetically predicted blood concentrations of proteins and PDAC risk. We  
143 used data for 8,275 cases and 6,723 controls of European descent from the Pancreatic Cancer  
144 Cohort Consortium (PanScan) and the Pancreatic Cancer Case-Control Consortium (PanC4).  
145 Beyond identifying novel proteins, we functionally validated two of them. Moreover, we  
146 generated a list of drugs targeting the identified proteins which may serve as candidates for drug  
147 repurposing of PDAC.

148

## 149 **Methods**

### 150 *Protein genetic prediction model development and validation*

151 We leveraged the genome and plasma proteome data of healthy European subjects  
152 included in the INTERVAL study to establish (subcohort1) and validate (subcohort2) protein  
153 genetic prediction models. The details of the INTERVAL study data have been published  
154 previously (10-14). Briefly, participants were generally healthy. The SOMAscan assay was used  
155 to collect the relative levels of 3,620 plasma proteins or complexes. Quality control (QC) was  
156 performed at both the sample and SOMAmer level. Approximately ~830,000 genetic variants  
157 were measured on the Affymetrix Axiom UK Biobank genotyping array. Standard sample and  
158 variant QC were conducted. SNPs were phased using SHAPEIT3 and imputed using a combined  
159 1000 Genomes Phase 3-UK10K reference panel, which resulted in over 87 million imputed  
160 variants. The SNPs were further filtered using criteria of 1) imputation quality of at least 0.7, 2)  
161 minor allele count of at least 5%, 3) Hardy Weinberg Equilibrium (HWE)  $p \geq 5 \times 10^{-6}$ , (4) missing  
162 rates < 5%, and (5) presenting in the 1000 Genome Project data for European populations.  
163 Overall there were 4,662,360 variants passing these criteria.

164 In subcohort 1 (N=2,481), as described elsewhere (10), protein concentrations were log  
165 transformed and adjusted for age, sex, duration between blood draw and processing, and the top  
166 three principal components. For the rank-inverse normalized residuals of each protein, we  
167 followed the TWAS/FUSION framework to establish prediction models, using nearby variants  
168 (within 100kb) of potentially associated SNPs as candidate predictors (15). A false discovery rate  
169 (FDR) < 0.05 was used to determine potentially associated SNPs in cis regions (within 1 Mb of  
170 the transcriptional start site (TSS) of the gene encoding the target protein of interest) and  $P$ -value  
171  $\leq 5 \times 10^{-8}$  was used to determine potentially associated SNPs in trans regions. We only included  
172 strand unambiguous SNPs. Four methods of best linear unbiased predictor (blup), elastic net,  
173 LASSO, and top1 were used to develop the models. For each protein of interest, the model

174 showing the most significant cross-validation  $P$ -value among those developed using the four  
175 methods was selected. For protein prediction models with  $R^2 \geq 0.01$ , external validation was  
176 conducted using genetic and protein data of subcohort 2 (N=820). Briefly, predicted protein  
177 expression levels were estimated by applying the developed protein prediction models to the  
178 genetic data, which were further compared with the measured levels for each protein of interest.  
179 Proteins with a model prediction  $R^2$  of  $\geq 0.01$  in subcohort1 and a correlation coefficient of  $\geq 0.1$   
180 in subcohort2 were selected for association analysis with PDAC risk.

181

### 182 *Examine associations of genetically predicted protein levels with PDAC risk*

183 To investigate the associations between genetically predicted circulating protein levels  
184 and PDAC risk, the validated protein genetic prediction models were applied to the summary  
185 statistics from a large GWAS of PDAC risk. In the present work, we used data from GWAS  
186 conducted in the PanScan and PanC4 consortia downloaded from the database of Genotypes and  
187 Phenotypes (dbGaP), including 8,275 PDAC cases and 6,723 controls of European ancestry.  
188 Detailed information on this dataset has been included elsewhere (16-18). Briefly, four GWAS  
189 studies, namely, PanScan I, PanScan II, PanScan III, and PanC4, were genotyped using the  
190 Illumina HumanHap550, 610-Quad, OmniExpress, and OmniExpressExome arrays, respectively.  
191 Standard QC procedures were performed according to the consortia guidelines (17). Study  
192 participants who were related to each other, had sex discordance, had genetic ancestry other than  
193 Europeans, had a low call rate (less than 98% and 94% in PanC4 and PanScan, respectively), or  
194 had missing information on age or sex were excluded. Duplicated SNPs, and those with a high  
195 missing call rate (at least 2% and 6% in PanC4 and PanScan, respectively) or with violations of  
196 Hardy-Weinberg equilibrium (HWE) ( $P < 1 \times 10^{-4}$  and  $P < 1 \times 10^{-7}$  in PanC4 and PanScan,

197 respectively), were also removed. Regarding SNP data from PanC4, those with minor allele  
198 frequency  $< 0.005$ , with more than two discordant calls in duplicate samples, with more than one  
199 Mendelian error in HapMap control trios, and those with sex difference in allele frequency  $> 0.2$   
200 or in heterozygosity  $> 0.3$  for autosomes/XY in European descendants were further removed. We  
201 performed genotype imputation using Minimac3 after prephasing with SHAPEIT from a  
202 reference panel of the Haplotype Reference Consortium (r1.1 2016) (19,20). We retained  
203 imputed SNPs with an imputation quality of  $\geq 0.3$ . The associations between individual genetic  
204 variants and PDAC risk were further estimated adjusting for age, sex and top principal  
205 components. The TWAS/FUSION framework was used to assess the protein-PDAC risk  
206 associations, by leveraging correlations between variants included in the prediction models based  
207 on the phase 3, 1000 Genomes Project data for European populations (15). We used the false  
208 discovery rate (FDR) corrected  $P$ -value threshold of  $\leq 0.05$  to determine significant associations  
209 between genetically predicted protein concentrations and risk of PDAC.

210

### 211 *Somatic variants of genes encoding associated proteins*

212 For each of the genes encoding the proteins that are identified to be associated with PDAC  
213 risk, we evaluated potentially deleterious somatic level mutations in 150 PDAC patients included  
214 in The Cancer Genome Atlas (TCGA). The potentially deleterious somatic variants include  
215 missense mutations, splice site mutations, nonstop mutations, nonsense mutations, frameshift  
216 mutations, in-frame mutations and translation start site mutations.

217 The somatic level genetic changes were called using MuTect2  
218 (doi: <https://doi.org/10.1101/861054>) and deposited to the TCGA data portal. The enrichment of  
219 proportion of assessed genes containing such somatic level genetic events compared with the

220 proportion of all protein-coding genes across the genome was evaluated using socscistatistics  
221 online website (<https://www.socscistatistics.com/tests/ztest/default2.aspx>).

### 222 *Ingenuity Pathway Analysis (IPA) and Protein-Protein Interaction (PPI) analysis*

223 To further assess whether genes encoding the identified PDAC associated proteins are  
224 enriched in specific pathways, molecular and cellular functions, and networks, we performed the  
225 enrichment analysis using Ingenuity Pathway Analysis (IPA) software (21). The "enrichment"  
226 score (Fisher exact test  $P$  value) that measures overlap of observed and predicted regulated gene  
227 sets was generated for each of the tested gene sets. The most significant pathways and functions  
228 with an enrichment  $P$  value less than 0.05 were reported. We also built protein-protein  
229 interaction (PPI) network using STRING database version 11.5 (<https://string-db.org/>) with  
230 0.400 confidence level (22). The STRING database integrates different curated databases  
231 containing information on known and predicted functional protein–protein associations.

### 232 *Drug repurposing analysis*

233 For the identified proteins, we further assessed whether there is any evidence supporting  
234 their potential roles in PDAC by using the OpenTargets (23). Focusing on those showing a  
235 potential relevance, we further mined evidence of their targeting drugs using the DrugBank (24)  
236 database. We also conducted molecular docking analysis for the identified proteins and  
237 corresponding candidate drug agents (25). Specifically, we downloaded the 3D structure of  
238 targeted proteins from Protein Data Bank (PDB) (26) with source code 1CPB, 3CDZ, 1IGR,  
239 3DFK, 5NO06, and drug agents from the PubChem database (27). We further worked out  
240 molecular docking between each of the proteins and the corresponding meta-drug agents to  
241 calculate the binding affinity scores (kcal/mol) for each pair of proteins and drugs.

242

**243 *In vitro functional validation of genes encoding selected associated novel proteins*****244 Cell Lines and Culture Condition**

245 Human pancreatic cancer cell lines PANC-1 and SU.86.86 were obtained from ATCC  
246 (American Type Culture Collection). All cells were cultured in vitro in DMEM (Dulbecco's  
247 modified eagle medium) high glucose medium (Gibco, Novato, CA, United States) supplemented  
248 with 10% (v/v) fetal bovine serum (FBS) (Gibco). Cells were incubated at 37°C with 5% CO<sub>2</sub>.

249

**250 Western blotting**

251 Post 72-hour silencing, we processed control, B4GALT1-silenced, and GOLM1-silenced  
252 cells for Western blotting. Cells were lysed using RIPA buffer, and equal protein amounts were  
253 separated on 10% or 12% SDS polyacrylamide gels, then transferred onto PVDF membranes. To  
254 prevent non-specific antibody binding, membranes were blocked with 5% milk in TBS with 0.1%  
255 Tween for an hour. They were then probed with anti-B4GALT1, anti-GOLM1, and anti-GAPDH  
256 antibodies, followed by their respective HRP-conjugated secondary antibodies. Signal detection  
257 was performed using Pierce™ ECL Western Blotting Substrate and images were captured and  
258 analyzed using Odyssey FC and ImageStudio Software.

259

**260 Quantitative Real-Time PCR (qPCR)**

261 Total RNA was extracted from cells using TRNzol reagent according to the manufacturer's  
262 protocol. The concentration of RNA was determined using a UV spectrophotometer.  
263 Subsequently, 2 mg of total RNA was reverse transcribed into cDNA using the iScript™ cDNA  
264 Synthesis Kit. qPCR analysis was performed on the CFX96™ Real-Time PCR Detection System

265 using the iTaq™ Universal SYBR® Green Supermix. The aim was to detect the expression levels  
266 of three genes: B4GALT1, GOLM1, and GAPDH mRNAs. Specific primer pairs were used for  
267 each gene. For B4GALT1, the forward sequence was GTATTTTGGAGGTGTCTCTGCTC and  
268 the reverse sequence was GGGCGAGATATAGACATGCCTC. For GOLM1, the forward  
269 sequence was ATCACACAGGTGAGAGGCTCA and the reverse sequence was  
270 ACTTCCTCTCCAGGTTGGTCTG. For the housekeeping gene GAPDH, the forward sequence  
271 was GTCTCCTCTGACTTCAACAGCG and the reverse sequence was  
272 ACCACCCTGTTGCTGTAGCCAA. During the qPCR analysis, melting curves were generated  
273 to detect primer-dimer formation and confirm the specificity of the gene-specific peaks for each  
274 target. To ensure accurate quantification, the expression data were normalized to the amount of  
275 GAPDH mRNA expressed.

276

### 277 **Transfection of siRNA**

278 The transfection of small-interfering RNA (siRNA) was performed using specific human  
279 siRNAs targeting GOLM1 (SASI\_Hs01\_00223155), B4GALT1 (SASI\_Hs01\_00080445), and the  
280 MISSION siRNA universal negative control, all of which were obtained from Sigma-Aldrich (St.  
281 Louis, MO). Cells were seeded in 6-well plates at a density of  $1.5 \times 10^5$  cells per well and  
282 subsequently transfected with the siRNAs at a concentration of 40 nM. The transfection procedure  
283 utilized the lipofectamine 2000 reagent (Invitrogen, Carlsbad, CA, United States) following the  
284 manufacturer's recommended guidelines. Gene silencing at both mRNA and protein levels was  
285 typically observed 72 h post-transfection. As such, the cells were collected and subjected to assays  
286 at the 72-hour time point to assess the efficacy of gene silencing.

287



**288 Cell Proliferation Assay**

289 To observe cell proliferation, cells were transfected with Mock siRNA, siGOLM11 and  
290 siB4GAL1 (40 nM). At 24 h after transfection, the cells were trypsinized and seeded into 96-well  
291 plates (Corning, NY, United States) at a density of 5000 cells/well in 200  $\mu$ l media. The plates  
292 were incubated in a 37°C humidified incubator. On each day for [3-(4,5-dimethylthiazol-2-yl)-5-  
293 (3-carboxymethoxyphenyl)-2-(4-sulfophenyl)-2H-tetrazolium] (MTS) assay.

294

**295 *In vitro* invasion assay**

296 Cell invasion was assessed following transfection with Mock siRNA, siGOLM11, and  
297 siB4GAL1 (40 nM). A modified Boyden chamber method was employed. Matrigel (BD  
298 Biosciences) was coated on the upper chamber of Transwell inserts (Corning, 8  $\mu$ m pore size) at a  
299 concentration of 300  $\mu$ g/ml, allowing gel formation for 2 hours at 37°C. Cells ( $5 \times 10^4$ ) were then  
300 suspended in 200  $\mu$ l of serum-free medium and added to the upper chamber. The lower chamber  
301 contained 600  $\mu$ l of medium with 10% FBS, acting as a chemoattractant. Following 24 hours of  
302 incubation at 37°C, non-invading cells on the upper membrane surface were gently removed using  
303 a cotton swab. Cells that invaded the lower membrane surface were fixed with 4%  
304 paraformaldehyde and stained with 0.1% crystal violet. Invasion was quantified by counting the  
305 stained cells on the underside of the membrane using a light microscope (10 random fields at 200x  
306 magnification). All experiments were performed in triplicate to ensure robustness of the findings.

307

**308 Wound Scratch assay**

309 After 24 hours of transfection with Mock siRNA, siGOLM11, and siB4GAL1, PANC-1  
310 and SU.86.86 cells were cultured in a 96-well plate to form a monolayer. Using BioTek's

311 AutoScratch™ Wound Making Tool, straight scratches were carefully created on the cell  
312 monolayer to mimic wounds, following the equipment manual's instructions. Time-lapse images  
313 of the scratches were captured at specific intervals (e.g., 0 hours, 12 hours, 24 hours, etc.) using  
314 the Cytation™ 5 Cell Imaging Multi-Mode Reader. Subsequently, image analysis software was  
315 employed to quantify the closure of the wounds at each time point. Statistical analysis was  
316 performed to compare the wound closure rates at different time points, and the results were  
317 presented graphically.

318

## 319 **Results**

320 The overall workflow of this study is shown in **Figure 1**. Of the proteins assessed, we  
321 were able to develop prediction models for 1,864 proteins with a prediction performance  
322  $R^2 \geq 0.01$ . In the external validation step, 1,389 of them further demonstrated a correlation  
323 coefficient of  $\geq 0.1$  for predicted expression and measured expression levels. Of such proteins,  
324 we observed significant associations between genetically predicted expression levels of 40  
325 proteins and PDAC risk at a false discovery rate (FDR)  $p$ -value of  $\leq 0.05$  (**Figure 2, Tables 1**  
326 **and 2**). Of the associated proteins, 16 are novel ones that have not been reported in previous  
327 studies (**Table 1**). Positive associations were observed for 10 of these proteins, and inverse  
328 associations were observed for six proteins (**Table 1**). The other 24 associated proteins have  
329 been previously reported in our study using pQTL as instruments (28) (**Table 2**). These include  
330 10 that demonstrated positive associations and 14 that showed inverse associations.

331 For the other proteins that were reported in our previous study using pQTL as instruments  
332 (28), while did not show a significant association after FDR correction in the current study  
333 (**Supplementary Table 1**), except for sTie-2, the directions of effect were consistent in the

334 current study compared with those in the published work. Among them, for eight proteins, their  
335 associations were at  $P < 0.05$  in the current work using protein genetic prediction models as  
336 instruments (**Supplementary Table 1**).

337 Based on a comparison of exome-sequencing data of tumor tissue and tumor-adjacent  
338 normal tissue obtained from 150 TCGA PDAC patients, the somatic level changes of potentially  
339 functional variants/mutations were observed in at least one patient for 10 of the 39 genes encoding  
340 identified associated proteins (**Supplementary Table 2**). This proportion ( $10/39=25.64\%$ ) is  
341 significantly higher (enrichment  $P$  value  $< 0.00001$ ) than the overall observed proportion of  
342 potentially functional changes across the genes encoding the proteins tested for association  
343 analyses ( $95/1,218 = 7.80\%$ ; here 1,218 represents the number of the genes available in TCGA  
344 analysis as part of the genes encoding the 1,389 assessed proteins).

345 According to the IPA analysis, several cancer-related functions were enriched for the  
346 genes encoding our identified proteins (**Supplementary Table 3**). The top canonical pathways  
347 identified included IL-15 production ( $P=2.21 \times 10^{-3}$ ), Heparan Sulfate Biosynthesis (Late Stages)  
348 ( $P=2.97 \times 10^{-3}$ ), Heparan Sulfate Biosynthesis ( $P=3.99 \times 10^{-3}$ ), Sperm Motility ( $P=7.73 \times 10^{-3}$ ), and  
349 Dermatan Sulfate Biosynthesis (Late Stages) ( $P=0.01$ ) (**Figure 3**). Among the related networks,  
350 the top network was cell-to-cell signaling and interaction, cardiovascular system development  
351 and function, organismal development (**Supplementary Figure 1**), followed by cancer,  
352 organismal injury and abnormalities, respiratory disease, free radical scavenging, cell death and  
353 survival, organismal injury and abnormalities, carbohydrate metabolism, small molecule  
354 biochemistry, cell cycle, and cancer, cell-to-cell signaling and interaction, cellular assembly and  
355 organization. Interactions among identified proteins were investigated based on STRING

356 database (**Figure 3**). In the network, KDR was predicted to interact with IGF1R, NOTCH1,  
357 MET, SEMA6A, ENG, SELP, and SELE.

358         Based on interrogation using the OpenTargets and DrugBank database, ten of the  
359 identified proteins are supported to be relevant to PDAC (overall score >0 in OpenTargets) and  
360 are targets of existing drugs approved to be used to treat human conditions (**Table 3**). Our work  
361 indicates potential drug repurposing opportunities of these drug targets to other indications. The  
362 scores of molecular docking between each of the proteins and the corresponding meta-drug  
363 agents were included in **Table 3**.

364         Among the 16 novel associated proteins, analysis of TGCA data also revealed potential  
365 relevance of B4GT1 and GOLM1 with tumor development (data not shown). Consequently, these  
366 two proteins were selected as the targets for experimental validation to further investigate their  
367 potential roles in PDAC development. Two gene-specific siRNAs (siGOML1 and siB4GAL1)  
368 were employed for post-transcriptional gene silencing of *GOML1* and *B4GAL1*, resulting in the  
369 knockdown of these two genes. As depicted in **Figure 4A**, qPCR analysis demonstrated a  
370 significant reduction in the mRNA expression of *GOML1* and *B4GAL1* in PANC-1 and SU.86.86  
371 cells at 72 hours after transfection with siGOML1 or siB4GAL1 (40 nM) when compared with the  
372 untreated control group ( $P < 0.05$ ). No significant difference was observed between the negative  
373 control group (NC, Mock-siRNA transfection) and the control groups (**Figure 4A**). This trend was  
374 also consistent in the western blot analysis (**Figure 4B**) in comparison with the qPCR assay,  
375 indicating that siGOML1 and siB4GAL1 effectively reduce the expression of *GOML1* and  
376 *B4GAL1* at both mRNA and protein levels in PANC-1 and SU.86.86 cells.

377         To assess the biological impact of *GOLM1* and *B4GAL1* silencing in PANC-1 and  
378 SU.86.86 cells, cell proliferation was examined using the MTS assay over a span of five

379 consecutive days. As shown in **Figures 4C** and **4D**, transfection of siGOML1 and siB4GAL1  
380 inhibited cell proliferation in both PANC-1 and SU.86.86 cells compared with the control  
381 (untransfected) and NC (Mock-siRNA transfected) groups. Furthermore, a wound healing assay  
382 demonstrated that at 12- and 24-hours post-scratch treatment, the open wound area in *GOLML1*  
383 and *B4GAL1* siRNA-transfected cells was significantly larger than that in mock siRNA-transfected  
384 or untransfected cells (**Figure 4D, 4E**), implying that knockdown of *GOLML1* and *B4GAL1* in  
385 PANC-1 and SU.86.86 cells effectively inhibited cell migration *in vitro*. To investigate whether  
386 the down-regulation of *GOLML1* and *B4GAL1* affects the invasive capabilities of PANC-1 and  
387 SU.86.86 cells, a transwell analysis was performed. The results revealed a significant inhibition of  
388 cell invasion in PANC-1 and SU.86.86 cells upon *GOLML1* or *B4GAL1* silencing. The number of  
389 siGOML1 or siB4GAL1-transfected cells invading through the membrane was markedly lower  
390 than that of control-siRNA transfected cells (**Fig. 4F**,  $P < 0.05$ ). Together, our findings suggest  
391 that GOLM1 and B4GT1 play crucial roles in PDAC cell proliferation, migration, and invasion,  
392 and their suppression could potentially serve as a therapeutic strategy for PDAC.

393

## 394 **Discussion**

395 This is the first PWAS study using comprehensive protein genetic prediction models to  
396 assess the associations between genetically predicted circulating protein concentrations and  
397 PDAC risk. Overall, we identified 40 proteins that were significantly associated with PDAC risk  
398 after FDR correction, including 16 novel proteins that have not been previously reported. Our  
399 results suggest new knowledge on the genetics and etiology of PDAC, and the newly identified  
400 proteins could serve as candidate blood biomarkers for risk assessment of PDAC, a highly fatal

401 malignancy. We also identified potential drug repurposing opportunities targeting the identified  
402 proteins which warrant further investigations.

403         In previous studies, blood concentrations of specific proteins such as CA242, PIVKA-II,  
404 PAM4, S100A6, OPN, RBM6, EphA2, and OPG have been reported to be potentially associated  
405 with PDAC risk (4-7). In the INTERVAL dataset, proteins S100A6 and OPG were captured, and  
406 we were able to develop satisfactory prediction models for their levels in blood (17). We  
407 observed a significant association with the same direction for OPG ( $P$ -value = 0.03, Z-score =  
408 2.23) but not for S100A6 ( $P$ -value=0.93) with PDAC risk. Such inconsistent findings with  
409 previous studies might be explained by potential biases in previous epidemiological studies and  
410 warrant further exploration.

411         In this large study, we identified 16 novel proteins that were associated with PDAC risk.  
412 Previous studies have suggested potential roles for some of the novel proteins in pancreatic  
413 tumorigenesis. Tie1 deficiency is reported to induce endothelial–mesenchymal transition  
414 (EndMT) and promote a motile phenotype (29). EndMT is known to present in human pancreatic  
415 tumors (29). Another study reports that TNF- $\alpha$  that is abundantly present in PDAC, induces  
416 EndMT and acts at least partially through TIE1 regulation in murine pancreatic tumors (30). For  
417 CPB1, immunohistochemistry of tissue microarray from PDAC patients showed that it was  
418 significantly downregulated in pancreatic tumor compared with adjacent normal pancreatic  
419 tissues (31). This aligns with the negative association between genetically predicted levels of  
420 carboxypeptidase B1 and PDAC risk observed in this study. In another study it was reported that  
421 mutations in *CPBI* were associated with pancreatic cancer (32). Regarding GOLM1, one study  
422 supported that long non-coding RNA TP73-AS1 could promote pancreatic cancer progression  
423 through GOLM1 upregulation by competitively binding to miR-128-3p (33). Further

424 investigations are warranted to clarify roles of the identified proteins in pancreatic cancer  
425 development.

426         Based on drug repurposing analyses, we prioritized several drugs that may serve as  
427 promising candidates for treating PDAC, such as Crizotinib, Cabozantinib, Brigatinib,  
428 Capmatinib, Tepotinib, and Tivozanib targeting Met. Previous research has supported potential  
429 link between these drugs and PDAC. For example, earlier research found that Crizotinib and  
430 Cabozantinib could decrease PDAC cell line viability *in vitro* (34). Cabozantinib together with  
431 photodynamic therapy had been shown to achieve local control and decrease in tumor metastases  
432 in preclinical PDAC models (35). A translational mathematical modeling study revealed that  
433 Tepotinib at a dose selection of 500 mg once daily could be effective for PDAC (36). Further  
434 work is needed to assess potential efficacy of these drug candidates in PDAC treatment.

435         There are several strengths of this study for detecting proteins associated with PDAC  
436 risk. We developed comprehensive protein genetic prediction models as instruments, which not  
437 only potentially minimize biases commonly encountered in conventional observational study  
438 design, but also bring improved statistical power compared with the design of only using pQTLs  
439 as instruments. However, several limitations of this study need to be recognized when  
440 interpreting our findings. First, our results may still be susceptible to potential pleiotropic effects  
441 and may not necessarily infer causality. Similar to the design of transcriptome-wide association  
442 study (TWAS), our PWAS should be useful for prioritizing causal proteins; however we cannot  
443 completely exclude the possibility of false positive findings for some of the identified  
444 associations (37). Several likely reasons may induce these, such as correlated protein expression  
445 across participants, correlated genetically predicted protein expression, as well as shared genetic  
446 variants (37). Future functional investigation will better characterize whether the identified

447 proteins play a causal role in PDAC development. Second, since in this work the genetically  
448 regulated components of plasma protein levels were studied but not the overall measured levels,  
449 the utility of the identified proteins as risk biomarkers for PDAC remains unclear. Additional  
450 work for measuring circulating protein levels in pre-diagnostic blood samples are needed to  
451 evaluate the prediction role of these proteins in PDAC risk. Third, for our current model  
452 development design, the candidate predictors for each protein of interest merely rely on the  
453 potentially associated SNPs at a specific statistical threshold. A small proportion of proteins were  
454 excluded for downstream model construction because of the lack of such SNPs. Future work  
455 considering additional potential predictors beyond such statistics-based selection would be  
456 needed to improve the ability to evaluate additional proteins. Fourth, previous work has  
457 supported that covariates of smoking and body mass index are related to blood protein levels  
458 (38,39). In the current study using INTERVAL resources, we were not able to adjust for these  
459 covariates during model construction. Further study is thus needed to validate our results. Lastly,  
460 the current study largely focuses on Europeans for both protein genetic prediction model  
461 development and downstream association analyses with PDAC risk. Future research is warranted  
462 to study proteins associated with PDAC risk in other non-European ancestries.

463 Our TGCA data analysis has revealed potential relevance of B4GT1 and GOLM1 in  
464 tumorigenesis and tumor progression. B4GT1 (Beta-1,4-Galactosyl transferase 1) is an enzyme  
465 primarily responsible for catalyzing the galactose transfer to specific receptor molecules within  
466 organisms (40). Its significance lies in its involvement in various essential biological processes,  
467 such as intercellular communication and cell adhesion. Furthermore, alterations in the expression  
468 level of B4GT1 have been observed in certain cancers, suggesting its potential implication in tumor  
469 initiation and development (41). This intriguing finding has led us to select B4GT1 as a priority



470 target for further exploration of its role in PDAC using experimental techniques. Similarly, our  
471 attention was drawn to GOLM1 (Golgi Membrane Protein 1), a membrane protein predominantly  
472 located in the Golgi apparatus, which plays a pivotal role in cellular secretion and transport  
473 processes. Recent investigations have demonstrated an upregulation of GOLM1 expression in  
474 multiple cancer types, including liver cancer, lung cancer, and pancreatic cancer. Such evidence  
475 strongly suggests that GOLM1 might exert a significant influence on the onset and progression of  
476 these malignancies (42). Consequently, we selected GOLM1 as an additional focus for verification  
477 to gain deeper insights into its involvement in PDAC. By utilizing RNAi technology to silence  
478 these genes, our experimental results corroborated the critical roles of GOLM1 and B4GT1 in  
479 driving PDAC cell proliferation, migration, and invasion. Subduing these genes holds promise as  
480 a potential therapeutic approach for PDAC treatment.

481         In summary, using protein genetic prediction models, we identified 16 novel protein  
482 biomarker candidates for which the genetically predicted circulating levels were significantly  
483 associated with PDAC risk. Future work is needed to better characterize the potential roles of  
484 these proteins in the etiology of PDAC development, assess the predictive role of such markers  
485 in risk assessment of PDAC, and evaluate whether the potential drug repurposing opportunities  
486 we identified may improve PDAC outcomes.

487

#### 488 **Roles of the funders**

489 The funders had no role in the design and conduct of the study; collection, management,  
490 analysis, and interpretation of the data; preparation, review, or approval of the manuscript; and  
491 decision to submit the manuscript for publication.

492

493 **Any prior presentations**

494 A preliminary of this work was presented at the American Association for Cancer Research  
495 Annual Meeting 2021.

496

497 **Acknowledgements**

498 The pancreatic cancer genetic datasets used for the association analyses described in this  
499 manuscript were obtained from dbGaP at <http://www.ncbi.nlm.nih.gov/sites/entrez?db=gap>  
500 through dbGaP accession phs000206.v5.p3 and phs000648.v1.p1. The authors also would like to  
501 thank all the individuals for their participation in the parent studies and all the researchers,  
502 clinicians, technicians and administrative staff for their contribution to the studies. This study is  
503 supported by the University of Hawaii Cancer Center and V Foundation V Scholar Award. Lang  
504 Wu and Chong Wu are supported by NCI R01CA263494. Lang Wu is also supported by  
505 NHGRI/NIMHD U54HG013243 and NCI R00CA218892. This work was supported in part by  
506 NIH-NIMHD U54MD007598, NIH/NCI1U54CA14393, U56 CA101599-01; Department-of-  
507 Defense Breast Cancer Research Program grant BC043180, NIH/NCATS CTSI UL1TR000124  
508 to J.V. Vadgama; Accelerating Excellence in Translational Science Pilot Grants G0812D05,  
509 NIH/NCI SC1CA200517 and 9 SC1 GM135050-05 to Y. Wu. Qizhi Yao is supported by VA  
510 Merit Award 1 I01 CX001822-01A2 (PI: Yao). The PanScan study was funded in whole or in  
511 part with federal funds from the National Cancer Institute (NCI), US National Institutes of  
512 Health (NIH) under contract number HHSN261200800001E. Additional support was received  
513 from NIH/NCI K07 CA140790, the American Society of Clinical Oncology Conquer Cancer  
514 Foundation, the Howard Hughes Medical Institute, the Lustgarten Foundation, the Robert T. and  
515 Judith B. Hale Fund for Pancreatic Cancer Research and Promises for Purple. A full list of

516 acknowledgments for each participating study is provided in the Supplementary Note of the  
517 manuscript with PubMed ID: 25086665. For the PanC4 GWAS study, the patients and controls  
518 were derived from the following PANC4 studies: Johns Hopkins National Familial Pancreas  
519 Tumor Registry, Mayo Clinic Biospecimen Resource for Pancreas Research, Ontario Pancreas  
520 Cancer Study (OPCS), Yale University, MD Anderson Case Control Study, Queensland  
521 Pancreatic Cancer Study, University of California San Francisco Molecular Epidemiology of  
522 Pancreatic Cancer Study, International Agency of Cancer Research and Memorial Sloan  
523 Kettering Cancer Center. This work is supported by NCI R01CA154823 Genotyping services  
524 were provided by the Center for Inherited Disease Research (CIDR). CIDR is fully funded  
525 through a federal contract from the National Institutes of Health to The Johns Hopkins  
526 University, contract number HHSN2682011000111. The content is solely the responsibility of  
527 the authors and does not necessarily represent the official views of the National Institutes of  
528 Health. Participants in the INTERVAL randomised controlled trial were recruited with the active  
529 collaboration of NHS Blood and Transplant England ([www.nhsbt.nhs.uk](http://www.nhsbt.nhs.uk)), which has supported  
530 field work and other elements of the trial. DNA extraction and genotyping were co-funded by the  
531 National Institute for Health Research (NIHR), the NIHR BioResource  
532 (<http://bioresource.nihr.ac.uk>) and the NIHR Cambridge Biomedical Research Centre (BRC-  
533 1215-20014) [\*]. The academic coordinating centre for INTERVAL was supported by core  
534 funding from the: NIHR Blood and Transplant Research Unit in Donor Health and Genomics  
535 (NIHR BTRU-2014-10024), UK Medical Research Council (MR/L003120/1), British Heart  
536 Foundation (SP/09/002; RG/13/13/30194; RG/18/13/33946) and NIHR Cambridge BRC (BRC-  
537 1215-20014) [\*]. A complete list of the investigators and contributors to the INTERVAL trial is

538 provided in reference [\*\*]. The academic coordinating centre would like to thank blood donor  
539 centre staff and blood donors for participating in the INTERVAL trial.

540

541 \*The views expressed are those of the author(s) and not necessarily those of the NIHR or the  
542 Department of Health and Social Care.

543

544 \*\*Di Angelantonio E, Thompson SG, Kaptoge SK, Moore C, Walker M, Armitage J, Ouwehand  
545 WH, Roberts DJ, Danesh J, INTERVAL Trial Group. Efficiency and safety of varying the  
546 frequency of whole blood donation (INTERVAL): a randomised trial of 45 000 donors. *Lancet*.  
547 2017 Nov 25;390(10110):2360-2371.

548

#### 549 **Data sharing statement**

550 The pancreatic cancer genetic datasets used for the association analyses described in this  
551 manuscript can be obtained from dbGaP at <http://www.ncbi.nlm.nih.gov/sites/entrez?db=gap>  
552 through dbGaP accession phs000206.v5.p3 and phs000648.v1.p1. The INTERVAL individual-  
553 level genotype and protein data, and full summary association results from the genetic analysis,  
554 are available through the European Genotype Archive (accession number EGAS00001002555).

555 Summary association results are also publicly available at

556 <http://www.phpc.cam.ac.uk/ceu/proteins/>, through PhenoScanner

557 (<http://www.phenoscaner.medschl.cam.ac.uk>) and from the NHGRI-EBI GWAS Catalog

558 (<https://www.ebi.ac.uk/gwas/downloads/summary-statistics>).

559

560

561

562

563 **Figure legends**564 **Figure 1.** The overall design of this study.

565 **Figure 2.** Manhattan plot of 40 identified proteins associated with PDAC risk. Proteins with blue  
566 color represent those identified in our previous work using pQTL as instruments, and proteins with  
567 red color represent novel ones identified in the current study.

568 **Figure 3.** PPI network and canonical pathways of 40 identified proteins associated with PDAC  
569 risk. Network nodes represent proteins; edge thickness is proportional to the evidence for the PPI;  
570 and dashed lines represent the interaction among clusters. The enrichment of canonical pathways  
571 was determined using IPA software.

572 **Figure 4.** The analysis of cell proliferation, migration and invasion on PANC-1 and SU.86.86  
573 cells with siB4GLAT1 and siGOLM1 transfection. The quantitative real-time PCR (qPCR) assay  
574 and the western blot assay (A) were used to investigate the RNAi effect of siB4GLAT1 and  
575 siGOLM1 (40 nM, 72 h) in PANC-1 and SU.86.86 cells. GAPDH were used as an internal  
576 control for qPCR analyses and western blot analyses, respectively (B,C) The effect of  
577 transfection with siB4GLAT1 and siGOLM1 (40 nM) on cell proliferation. The cells were  
578 detected by MTS [3-(4,5-dimethylthiazol-2-yl)-5-(3-carboxymethoxyphenyl)-2-(4-sulfophenyl)-  
579 2H-tetrazolium] assay on each day for 5 consecutive days. (D,E) Silencing of *B4GLAT1* and  
580 *GOLM1* inhibited migration of PANC-1 and SU.86.86 cells. Representative images of wound  
581 scratch assay performed to evaluate the motility of cells after silencing *B4GLAT1* and *GOLM1*.  
582 After transfection, a scratch was made on cells monolayer and was monitored with microscopy  
583 every 12 hours (0, 12, and 24 h). Bar graphs show normalized wound area, calculated using Gen

584 5. Representative images of invasion assay. Data are represented as mean  $\pm$  SD from triplicate  
585 samples, where  $*p < 0.01$  compared to the control. (F) Effect of siB4GLAT1 and siGOLM1  
586 transfection on the invasion of PANC-1 and SU.86.86 cells. After siB4GLAT1 and siGOLM1  
587 transfection for 48 h, invasive ability of PANC-1 and SU.86.86 cells was identified by transwell  
588 assay.  $**P < 0.01$  compared with the control cells;  $##P < 0.01$  compared with the mock cells;  
589 data are expressed as the mean  $\pm$  SD,  $n = 3$ .

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605 **References**

- 606 1. Sung H, Ferlay J, Siegel RL, Laversanne M, Soerjomataram I, Jemal A, *et al.* Global  
607 Cancer Statistics 2020: GLOBOCAN Estimates of Incidence and Mortality Worldwide  
608 for 36 Cancers in 185 Countries. *CA Cancer J Clin* **2021**;71(3):209-49 doi  
609 10.3322/caac.21660.
- 610 2. Rawla P, Sunkara T, Gaduputi V. Epidemiology of Pancreatic Cancer: Global Trends,  
611 Etiology and Risk Factors. *World J Oncol* **2019**;10(1):10-27 doi 10.14740/wjon1166.
- 612 3. Ballehaninna UK, Chamberlain RS. The clinical utility of serum CA 19-9 in the  
613 diagnosis, prognosis and management of pancreatic adenocarcinoma: An evidence based  
614 appraisal. *J Gastrointest Oncol* **2012**;3(2):105-19 doi 10.3978/j.issn.2078-6891.2011.021.
- 615 4. Tartaglione S, Pecorella I, Zarrillo SR, Granato T, Viggiani V, Manganaro L, *et al.*  
616 Protein Induced by Vitamin K Absence II (PIVKA-II) as a potential serological  
617 biomarker in pancreatic cancer: a pilot study. *Biochem Med (Zagreb)* **2019**;29(2):020707  
618 doi 10.11613/BM.2019.020707.
- 619 5. Duan B, Hu X, Fan M, Xiong X, Han L, Wang Z, *et al.* RNA-Binding Motif Protein 6 is  
620 a Candidate Serum Biomarker for Pancreatic Cancer. *Proteomics Clin Appl*  
621 **2019**;13(5):e1900048 doi 10.1002/prca.201900048.
- 622 6. Koshikawa N, Minegishi T, Kiyokawa H, Seiki M. Specific detection of soluble EphA2  
623 fragments in blood as a new biomarker for pancreatic cancer. *Cell Death Dis*  
624 **2017**;8(10):e3134 doi 10.1038/cddis.2017.545.
- 625 7. Loosen SH, Neumann UP, Trautwein C, Roderburg C, Luedde T. Current and future  
626 biomarkers for pancreatic adenocarcinoma. *Tumour Biol* **2017**;39(6):1010428317692231  
627 doi 10.1177/1010428317692231.
- 628 8. Didelez V, Sheehan N. Mendelian randomization as an instrumental variable approach to  
629 causal inference. *Stat Methods Med Res* **2007**;16(4):309-30 doi  
630 10.1177/0962280206077743.
- 631 9. Lawlor DA, Harbord RM, Sterne JA, Timpson N, Davey Smith G. Mendelian  
632 randomization: using genes as instruments for making causal inferences in epidemiology.  
633 *Stat Med* **2008**;27(8):1133-63 doi 10.1002/sim.3034.
- 634 10. Sun BB, Maranville JC, Peters JE, Stacey D, Staley JR, Blackshaw J, *et al.* Genomic atlas  
635 of the human plasma proteome. *Nature* **2018**;558(7708):73-9 doi 10.1038/s41586-018-  
636 0175-2.
- 637 11. Wu L, Shu X, Bao J, Guo X, Kote-Jarai Z, Haiman CA, *et al.* Analysis of Over 140,000  
638 European Descendants Identifies Genetically Predicted Blood Protein Biomarkers  
639 Associated with Prostate Cancer Risk. *Cancer Res* **2019**;79(18):4592-8 doi  
640 10.1158/0008-5472.CAN-18-3997.
- 641 12. Zhu J, Wu C, Wu L. Associations Between Genetically Predicted Protein Levels and  
642 COVID-19 Severity. *J Infect Dis* **2021**;223(1):19-22 doi 10.1093/infdis/jiaa660.
- 643 13. Zhu J, O'Mara TA, Liu D, Setiawan VW, Glubb D, Spurdle AB, *et al.* Associations  
644 between Genetically Predicted Circulating Protein Concentrations and Endometrial  
645 Cancer Risk. *Cancers (Basel)* **2021**;13(9) doi 10.3390/cancers13092088.

- 646 14. Shu X, Bao J, Wu L, Long J, Shu XO, Guo X, *et al.* Evaluation of associations between  
647 genetically predicted circulating protein biomarkers and breast cancer risk. *Int J Cancer*  
648 **2020**;146(8):2130-8 doi 10.1002/ijc.32542.
- 649 15. Gusev A, Ko A, Shi H, Bhatia G, Chung W, Penninx BW, *et al.* Integrative approaches  
650 for large-scale transcriptome-wide association studies. *Nat Genet* **2016**;48(3):245-52 doi  
651 10.1038/ng.3506.
- 652 16. Liu D, Zhou D, Sun Y, Zhu J, Ghoneim D, Wu C, *et al.* A Transcriptome-Wide  
653 Association Study Identifies Candidate Susceptibility Genes for Pancreatic Cancer Risk.  
654 *Cancer Res* **2020**;80(20):4346-54 doi 10.1158/0008-5472.CAN-20-1353.
- 655 17. Klein AP, Wolpin BM, Risch HA, Stolzenberg-Solomon RZ, Mocchi E, Zhang M, *et al.*  
656 Genome-wide meta-analysis identifies five new susceptibility loci for pancreatic cancer.  
657 *Nat Commun* **2018**;9(1):556 doi 10.1038/s41467-018-02942-5.
- 658 18. Zhu J, Yang Y, Kisiel JB, Mahoney DW, Michaud DS, Guo X, *et al.* Integrating Genome  
659 and Methylome Data to Identify Candidate DNA Methylation Biomarkers for Pancreatic  
660 Cancer Risk. *Cancer Epidemiol Biomarkers Prev* **2021**;30(11):2079-87 doi  
661 10.1158/1055-9965.EPI-21-0400.
- 662 19. McCarthy S, Das S, Kretzschmar W, Delaneau O, Wood AR, Teumer A, *et al.* A  
663 reference panel of 64,976 haplotypes for genotype imputation. *Nat Genet*  
664 **2016**;48(10):1279-83 doi 10.1038/ng.3643.
- 665 20. Howie BN, Donnelly P, Marchini J. A flexible and accurate genotype imputation method  
666 for the next generation of genome-wide association studies. *PLoS Genet*  
667 **2009**;5(6):e1000529 doi 10.1371/journal.pgen.1000529.
- 668 21. Kramer A, Green J, Pollard J, Jr., Tugendreich S. Causal analysis approaches in  
669 Ingenuity Pathway Analysis. *Bioinformatics* **2014**;30(4):523-30 doi  
670 10.1093/bioinformatics/btt703.
- 671 22. Szklarczyk D, Gable AL, Nastou KC, Lyon D, Kirsch R, Pyysalo S, *et al.* The STRING  
672 database in 2021: customizable protein-protein networks, and functional characterization  
673 of user-uploaded gene/measurement sets. *Nucleic Acids Res* **2021**;49(D1):D605-D12 doi  
674 10.1093/nar/gkaa1074.
- 675 23. Koscielny G, An P, Carvalho-Silva D, Cham JA, Fumis L, Gasparyan R, *et al.* Open  
676 Targets: a platform for therapeutic target identification and validation. *Nucleic Acids Res*  
677 **2017**;45(D1):D985-D94 doi 10.1093/nar/gkw1055.
- 678 24. Wishart DS, Knox C, Guo AC, Shrivastava S, Hassanali M, Stothard P, *et al.* DrugBank:  
679 a comprehensive resource for in silico drug discovery and exploration. *Nucleic Acids Res*  
680 **2006**;34(Database issue):D668-72 doi 10.1093/nar/gkj067.
- 681 25. Alam MA SH, Deng H-W. A robust kernel machine regression towards biomarker  
682 selection in multi-omics datasets of osteoporosis for drug discovery. In: University T,  
683 editor. arXiv2022.
- 684 26. Kim S, Chen J, Cheng T, Gindulyte A, He J, He S, *et al.* PubChem 2019 update:  
685 improved access to chemical data. *Nucleic Acids Res* **2019**;47(D1):D1102-D9 doi  
686 10.1093/nar/gky1033.
- 687 27. Kim SY, Jeong HH, Kim J, Moon JH, Sohn KA. Robust pathway-based multi-omics data  
688 integration using directed random walks for survival prediction in multiple cancer  
689 studies. *Biol Direct* **2019**;14(1):8 doi 10.1186/s13062-019-0239-8.



- 690 28. Zhu J, Shu X, Guo X, Liu D, Bao J, Milne RL, *et al.* Associations between Genetically  
691 Predicted Blood Protein Biomarkers and Pancreatic Cancer Risk. *Cancer Epidemiol*  
692 *Biomarkers Prev* **2020**;29(7):1501-8 doi 10.1158/1055-9965.EPI-20-0091.
- 693 29. Garcia J, Sandi MJ, Cordelier P, Binetruy B, Pouyssegur J, Iovanna JL, *et al.* Tie1  
694 deficiency induces endothelial-mesenchymal transition. *EMBO Rep* **2012**;13(5):431-9  
695 doi 10.1038/embor.2012.29.
- 696 30. Adjuto-Saccone M, Soubeyran P, Garcia J, Audebert S, Camoin L, Rubis M, *et al.* TNF-  
697 alpha induces endothelial-mesenchymal transition promoting stromal development of  
698 pancreatic adenocarcinoma. *Cell Death Dis* **2021**;12(7):649 doi 10.1038/s41419-021-  
699 03920-4.
- 700 31. Song Y, Wang Q, Wang D, Junqiang L, Yang J, Li H, *et al.* Label-Free Quantitative  
701 Proteomics Unravels Carboxypeptidases as the Novel Biomarker in Pancreatic Ductal  
702 Adenocarcinoma. *Transl Oncol* **2018**;11(3):691-9 doi 10.1016/j.tranon.2018.03.005.
- 703 32. Tamura K, Yu J, Hata T, Suenaga M, Shindo K, Abe T, *et al.* Mutations in the pancreatic  
704 secretory enzymes CPA1 and CPB1 are associated with pancreatic cancer. *Proc Natl*  
705 *Acad Sci U S A* **2018**;115(18):4767-72 doi 10.1073/pnas.1720588115.
- 706 33. Wang B, Sun X, Huang KJ, Zhou LS, Qiu ZJ. Long non-coding RNA TP73-AS1  
707 promotes pancreatic cancer growth and metastasis through miRNA-128-3p/GOLM1 axis.  
708 *World J Gastroenterol* **2021**;27(17):1993-2014 doi 10.3748/wjg.v27.i17.1993.
- 709 34. Escorcía FE, Houghton JL, Abdel-Atti D, Pereira PR, Cho A, Gutsche NT, *et al.*  
710 ImmunoPET Predicts Response to Met-targeted Radioligand Therapy in Models of  
711 Pancreatic Cancer Resistant to Met Kinase Inhibitors. *Theranostics* **2020**;10(1):151-65  
712 doi 10.7150/thno.37098.
- 713 35. Broekgaarden M, Alkhateeb A, Bano S, Bulin AL, Obaid G, Rizvi I, *et al.* Cabozantinib  
714 Inhibits Photodynamic Therapy-Induced Auto- and Paracrine MET Signaling in  
715 Heterotypic Pancreatic Microtumors. *Cancers (Basel)* **2020**;12(6) doi  
716 10.3390/cancers12061401.
- 717 36. Xiong W, Friese-Hamim M, Johne A, Stroh C, Klevesath M, Falchook GS, *et al.*  
718 Translational pharmacokinetic-pharmacodynamic modeling of preclinical and clinical  
719 data of the oral MET inhibitor tepotinib to determine the recommended phase II dose.  
720 *CPT Pharmacometrics Syst Pharmacol* **2021**;10(5):428-40 doi 10.1002/psp4.12602.
- 721 37. Wainberg M, Sinnott-Armstrong N, Mancuso N, Barbeira AN, Knowles DA, Golan D, *et*  
722 *al.* Opportunities and challenges for transcriptome-wide association studies. *Nature*  
723 *genetics* **2019**;51(4):592-9.
- 724 38. Madhuvanathi M, Lathadevi GV. Serum Proteins Alteration in Association with Body  
725 Mass Index in Human Volunteers. *J Clin Diagn Res* **2016**;10(6):CC05-7 doi  
726 10.7860/JCDR/2016/18278.8047.
- 727 39. Gallus S, Lugo A, Suatoni P, Taverna F, Bertocchi E, Boffi R, *et al.* Effect of Tobacco  
728 Smoking Cessation on C-Reactive Protein Levels in A Cohort of Low-Dose Computed  
729 Tomography Screening Participants. *Sci Rep* **2018**;8(1):12908 doi 10.1038/s41598-018-  
730 29867-9.
- 731 40. Morokuma D, Xu J, Hino M, Mon H, Merzaban JS, Takahashi M, *et al.* Expression and  
732 Characterization of Human beta-1, 4-Galactosyltransferase 1 (beta4GalT1) Using  
733 Silk Worm-Baculovirus Expression System. *Mol Biotechnol* **2017**;59(4-5):151-8 doi  
734 10.1007/s12033-017-0003-1.

- 735 41. Cui Y, Li J, Zhang P, Yin D, Wang Z, Dai J, *et al.* B4GALT1 promotes immune escape  
736 by regulating the expression of PD-L1 at multiple levels in lung adenocarcinoma. *J Exp*  
737 *Clin Cancer Res* **2023**;42(1):146 doi 10.1186/s13046-023-02711-3.
- 738 42. Liu Y, Hu X, Liu S, Zhou S, Chen Z, Jin H. Golgi Phosphoprotein 73: The Driver of  
739 Epithelial-Mesenchymal Transition in Cancer. *Front Oncol* **2021**;11:783860 doi  
740 10.3389/fonc.2021.783860.
- 741

**Table 1.** Novel proteins with genetically predicted concentrations in plasma to be associated with pancreatic cancer risk

Protein	SOMAmer ID	Protein full name	Protein-encoding gene	Region for protein encoding gene	Prediction model method	Number of Predicting SNPs	Number of Predicting SNPs-Cis*	Number of Predicting SNPs-Trans	Model internal cross validation R <sup>2</sup>	Model external validation R <sup>2</sup>	Z-value <sup>a</sup>	P-value <sup>a</sup>	FDR P-value <sup>b</sup>
IL-23 R	IL23R.5088.175.3	Interleukin-23 receptor	<i>IL23R</i>	1p31.3	elastic net	24	24	0	0.04	0.04	3.55	3.80×10 <sup>-4</sup>	0.02
sTie-1	TIE1.2844.53.2	Tyrosine-Protein Kinase Receptor Tie-1, Soluble	<i>TIE1</i>	1p34.2	lasso	18	7	11	0.22	0.28	5.67	1.46×10 <sup>-8</sup>	1.22×10 <sup>-6</sup>
FA20B	FAM20B.7198.197.3	Glycosaminoglycan Xylosylkinase	<i>FAM20B</i>	1q25.2	lasso	8	5	3	0.02	0.04	5.30	1.17×10 <sup>-7</sup>	7.82×10 <sup>-6</sup>
FAM3D	FAM3D.13102.1.3	Protein FAM3D	<i>FAM3D</i>	3p14.2	elastic net	58	16	42	0.37	0.36	6.10	1.07×10 <sup>-9</sup>	1.02×10 <sup>-7</sup>
Carboxypeptidase B1	CPB1.6356.3.3	Carboxypeptidase B	<i>CPB1</i>	3q24	lasso	7	3	4	0.04	0.03	-4.55	5.38×10 <sup>-6</sup>	3.00×10 <sup>-4</sup>
RAP	LRPAP1.3640.14.3	alpha-2-macroglobulin receptor-associated protein	<i>LRPAP1</i>	4p16.3	elastic net	168	23	145	0.27	0.22	3.21	0.001	0.04
Semaphorin-6A	SEMA6A.7945.10.3	Semaphorin-6A	<i>SEMA6A</i>	5q23.1	elastic net	66	44	22	0.05	0.05	-3.57	3.54×10 <sup>-4</sup>	0.02
B4GT1	B4GALT1.13381.49.3	Beta-1,4-galactosyltransferase 1	<i>B4GALT1</i>	9p21.1	elastic net	39	16	23	0.08	0.10	4.65	3.29×10 <sup>-6</sup>	1.96×10 <sup>-4</sup>
GOLM1	GOLM1.8983.7.3	Golgi Membrane Protein 1	<i>GOLM1</i>	9q21.33	lasso	10	0	10	0.14	0.17	8.07	7.12×10 <sup>-16</sup>	2.14×10 <sup>-13</sup>
QSOX2	QSOX2.8397.147.3	Sulfhydryl oxidase 2	<i>QSOX2</i>	9q34.3	elastic net	28	10	18	0.40	0.40	7.98	1.44×10 <sup>-15</sup>	2.75×10 <sup>-13</sup>
KIN17	KIN.14643.27.3	DNA/RNA-binding protein KIN17	<i>KIN</i>	10p14	elastic net	29	0	29	0.05	0.07	-5.52	3.31×10 <sup>-8</sup>	2.60×10 <sup>-6</sup>
ISLR2	ISLR2.13124.20.3	Immunoglobulin superfamily containing leucine-rich repeat protein 2	<i>ISLR2</i>	15q24.1	elastic net	77	32	45	0.14	0.13	-3.45	5.65×10 <sup>-4</sup>	0.02
DPEP2	DPEP2.8327.26.3	Dipeptidase 2	<i>DPEP2</i>	16q22.1	elastic net	36	0	36	0.06	0.05	-4.01	5.97×10 <sup>-5</sup>	0.003
Chymotrypsin	CTRB1.5671.1.3	Chymotrypsinogen B	<i>CTRB1</i>	16q23.1	elastic net	85	69	16	0.23	0.24	-4.32	1.59×10 <sup>-5</sup>	8.50×10 <sup>-4</sup>

Laminin	LAMA1.LAMB1.LAMC1. 2728.62.2	Laminin	<i>LAMA1</i> , <i>LAMB1</i> , <i>LAMC1</i>	18p11.31, 7q31.1, 1q25.3	elastic net	62	14	48	0.08	0.05	3.88	1.06×10 <sup>-4</sup>	0.005
TPST2	TPST2.8024.64.3	Protein-Tyrosine Sulfotransferase 2	<i>TPST2</i>	22q12.1	elastic net	52	28	24	0.07	0.08	5.88	4.16×10 <sup>-9</sup>	3.71×10 <sup>-7</sup>

\* SNPs within 1MB of the protein-encoding gene

a Associations between genetically predicted protein levels and PDAC risk after adjustment for age, sex, and top 10 principle components.

b FDR *P*-value: false discovery rate (FDR) adjusted *P*-value; associations with a FDR  $p \leq 0.05$  considered statistically significant

**Table 2.** Previously reported proteins with genetically predicted concentrations in plasma to be associated with pancreatic cancer risk

Protein	SOMAmer ID	Protein full name	Protein-encoding gene	Region for protein encoding gene	Prediction model method	Number of Predicting SNPs	Number of Predicting SNPs-Cis*	Number of Predicting SNPs-Trans	Model internal cross validation R <sup>2</sup>	Model external validation R <sup>2</sup>	Z-value <sup>a</sup>	P-value <sup>a</sup>	FDR P-value <sup>b</sup>
sE-Selectin	SELE.3470.1.2	E-selectin	<i>SELE</i>	1q24.2	lasso	6	0	6	0.39	0.44	-7.88	3.33×10 <sup>-15</sup>	5.47×10 <sup>-13</sup>
P-Selectin	SELP.4154.57.2	P-Selectin	<i>SELP</i>	1q24.2	lasso	11	7	4	0.26	0.27	-3.77	1.66×10 <sup>-4</sup>	0.008
LMA2L	LMAN2L.8013.9.3	VIP36-like protein	<i>LMAN2L</i>	2q11.2	top1	1	1	0	0.03	0.02	3.35	8.01×10 <sup>-4</sup>	0.03
Alkaline phosphatase, intestine	ALPI.10463.23.3	Intestinal-type alkaline phosphatase	<i>ALPI</i>	2q37.1	lasso	8	0	8	0.03	0.06	-6.79	1.09×10 <sup>-11</sup>	1.21×10 <sup>-9</sup>
VEGF sR2	KDR.3651.50.5	Vascular endothelial growth factor receptor 2	<i>KDR</i>	4q12	elastic net	56	18	38	0.18	0.12	-6.21	5.22×10 <sup>-10</sup>	5.37×10 <sup>-8</sup>
ADH1B	ADH1B.9834.62.3	Alcohol dehydrogenase 1B	<i>ADH1B</i>	4q23	lasso	6	0	6	0.08	0.03	3.21	0.001	0.04
LIF sR	LIFR.5837.49.3	Leukemia inhibitory factor receptor	<i>LIFR</i>	5p13.1	top1	1	0	1	0.03	0.02	-7.39	1.42×10 <sup>-13</sup>	1.73×10 <sup>-11</sup>
gp130, soluble	IL6ST.2620.4.2	Interleukin-6 receptor subunit beta	<i>IL6ST</i>	5q11.2	elastic net	51	21	30	0.06	0.05	-3.69	2.22×10 <sup>-4</sup>	0.01
GP116	ADGRF5.6409.57.3	Adhesion G protein-coupled receptor F5	<i>ADGRF5</i>	6p12.3	lasso	22	15	7	0.46	0.43	-4.65	3.37×10 <sup>-6</sup>	1.96×10 <sup>-4</sup>
CD36 ANTIGEN	CD36.2973.15.2	Platelet glycoprotein 4	<i>CD36</i>	7q21.11	top1	1	0	1	0.03	0.05	3.31	9.25×10 <sup>-4</sup>	0.03
Met	MET.2837.3.2	Hepatocyte growth factor receptor	<i>MET</i>	7q31	blup	1,668	603	1,065	0.07	0.04	-5.06	4.27×10 <sup>-7</sup>	2.72×10 <sup>-5</sup>
STOM	STOM.8261.51.3	Erythrocyte band 7 integral membrane protein	<i>STOM</i>	9q33.2	lasso	5	0	5	0.11	0.05	3.31	9.18×10 <sup>-4</sup>	0.03
BGAT	ABO.9253.52.3	Histo-blood group ABO system transferase	<i>ABO</i>	9q34.2	blup	2,473	2,347	126	0.72	0.72	9.18	4.20×10 <sup>-20</sup>	5.62×10 <sup>-17</sup>
Notch 1	NOTCH1.5107.7.2	Neurogenic locus notch homolog protein 1	<i>NOTCH1</i>	9q34.3	top1	1	0	1	0.01	0.02	3.29	9.97×10 <sup>-4</sup>	0.04

Endoglin	ENG.4908.6.1	Endoglin	<i>ENG</i>	9q34.11	top1	1	0	1	0.01	0.01	-8.04	$8.93 \times 10^{-16}$	$2.14 \times 10^{-13}$
ST4S6	CHST15.4469.78.2	Carbohydrate sulfotransferase 15	<i>CHST15</i>	10q26.13	lasso	5	1	4	0.05	0.03	-8.62	$6.46 \times 10^{-18}$	$4.32 \times 10^{-15}$
	CHST15.14097.86.3				lasso	9	2	7	0.04	0.02	-8.03	$9.60 \times 10^{-16}$	$2.14 \times 10^{-13}$
CHSTB	CHST11.7779.86.3	Carbohydrate sulfotransferase 11	<i>CHST11</i>	12q23.3	elastic net	69	46	23	0.11	0.07	3.52	$4.25 \times 10^{-4}$	0.02
THSD1	THSD1.5621.64.3	Thrombospondin type-1 domain-containing protein 1	<i>THSD1</i>	13q14.3	elastic net	44	27	17	0.04	0.03	-5.34	$9.41 \times 10^{-8}$	$6.62 \times 10^{-6}$
GLCE	GLCE.7808.5.3	D-glucuronyl C5-epimerase	<i>GLCE</i>	15q23	lasso	11	6	5	0.36	0.34	4.18	$2.94 \times 10^{-5}$	0.002
IGF-1 sR	IGF1R.4232.19.2	Insulin-like growth factor 1 receptor	<i>IGF1R</i>	15q26.3	top1	1	0	1	0.01	0.02	-7.39	$1.42 \times 10^{-13}$	$1.73 \times 10^{-11}$
Desmoglein-2	DSG2.9484.75.3	Desmoglein-2	<i>DSG2</i>	18q12.1	elastic net	66	44	22	0.04	0.06	5.34	$9.18 \times 10^{-8}$	$6.62 \times 10^{-6}$
DC-SIGN	CD209.3029.52.2	CD209 Antigen	<i>CD209</i>	19p13.2	elastic net	58	26	32	0.39	0.38	8.52	$1.62 \times 10^{-17}$	$7.22 \times 10^{-15}$
IR	INSR.3448.13.2	Insulin receptor	<i>INSR</i>	19p13.2	lasso	7	0	7	0.09	0.12	-7.53	$4.98 \times 10^{-14}$	$7.40 \times 10^{-12}$

\* SNPs within 1MB of the protein-encoding gene

a Associations between genetically predicted protein levels and PDAC risk after adjustment for age, sex, and top 10 principle components.

b FDR *P*-value: false discovery rate (FDR) adjusted *P*-value; associations with a FDR  $p \leq 0.05$  considered statistically significant

**Table 3.** Drug repurposing opportunities

Protein	Protein full name	Protein-encoding gene	OpenTargets information (overall score)	Drugbank ID	Drug name	Molecular action	Molecular docking score*
sTie-1	Tyrosine-Protein Kinase Receptor Tie-1, Soluble	<i>TIE1</i>	0.006	DB12010	Fostamatinib	inhibitor	-6.1
Carboxypeptidase B1	Carboxypeptidase B	<i>CPB1</i>	0.159	DB04272	Citric acid	NA	-3.9
Chymotrypsin	Chymotrypsinogen B	<i>CTRB1</i>	0.078	DB06692	Aprotinin	NA	MDNA
sE-Selectin	E-selectin	<i>SELE</i>	0.023	DB01136	Carvedilol	inhibitor	-6.9
P-Selectin	P-Selectin	<i>SELP</i>	0.008	DB01109	Heparin	inhibitor	-4.9
				DB08813	Nadroparin	inhibitor	-4.9
				DB06779	Dalteparin	inhibitor	-4.9
				DB15271	Crizanlizumab	inhibitor	3DSNA
VEGF sR2	Vascular endothelial growth factor receptor 2	<i>KDR</i>	0.367	DB06589	Pazopanib	inhibitor	-6.3
				DB08896	Regorafenib	inhibitor	-6.5
				DB09079	Nintedanib	inhibitor	-5.8
				DB14840	Ripretinib	inhibitor	-6.6
				DB00398	Sorafenib	antagonist	-6.6
				DB01268	Sunitinib	inhibitor	-5.6
				DB06595	Midostaurin	antagonist inhibitor	-5.1
				DB06626	Axitinib	inhibitor	-6.0
				DB08875	Cabozantinib	antagonist	<b>-7.0</b>
				DB08901	Ponatinib	inhibitor	-6.9
DB09078	Lenvatinib	inhibitor	-6.1				

				DB05578	Ramucirumab	antagonist	3DSNA
				DB12010	Fostamatinib	inhibitor	-5.3
				DB12147	Erdafitinib	substrate	-5.5
				DB15822	Pralsetinib	inhibitor	-6.9
				DB11800	Tivozanib	inhibitor	-6.4
ADH1B	Alcohol dehydrogenase 1B	<i>ADH1B</i>	0.001	DB00898	Ethanol	substrate	-2.8
				DB09462	Glycerin	NA	-3.7
				DB00157	NADH	substrate	<b>-9.6</b>
				DB01213	Fomepizole	inhibitor	-3.9
Met	Hepatocyte growth factor receptor	<i>MET</i>	0.304	DB08865	Crizotinib	inhibitor	<b>-8.1</b>
				DB08875	Cabozantinib	antagonist	<b>-8</b>
				DB12267	Brigatinib	inhibitor	<b>-8.2</b>
				DB12010	Fostamatinib	inhibitor	-6.7
				DB11791	Capmatinib	inhibitor	<b>-8.7</b>
				DB15133	Tepotinib	inhibitor	<b>-8.3</b>
				DB11800	Tivozanib	inhibitor	<b>-8.2</b>
				DB16695	Amivantamab	antagonist antibody	3DSNA
IGF-I sR	Insulin-like growth factor 1 receptor	<i>IGF1R</i>	0.099	DB00071	Insulin pork	NA	MDNA
				DB00046	Insulin lispro	activator	MDNA
				DB01307	Insulin detemir	activator	MDNA
				DB00047	Insulin glargine	activator	MDNA
				DB01306	Insulin aspart	activator	MDNA
				DB01309	Insulin glulisine	activator	MDNA
				DB09564	Insulin degludec	activator	MDNA



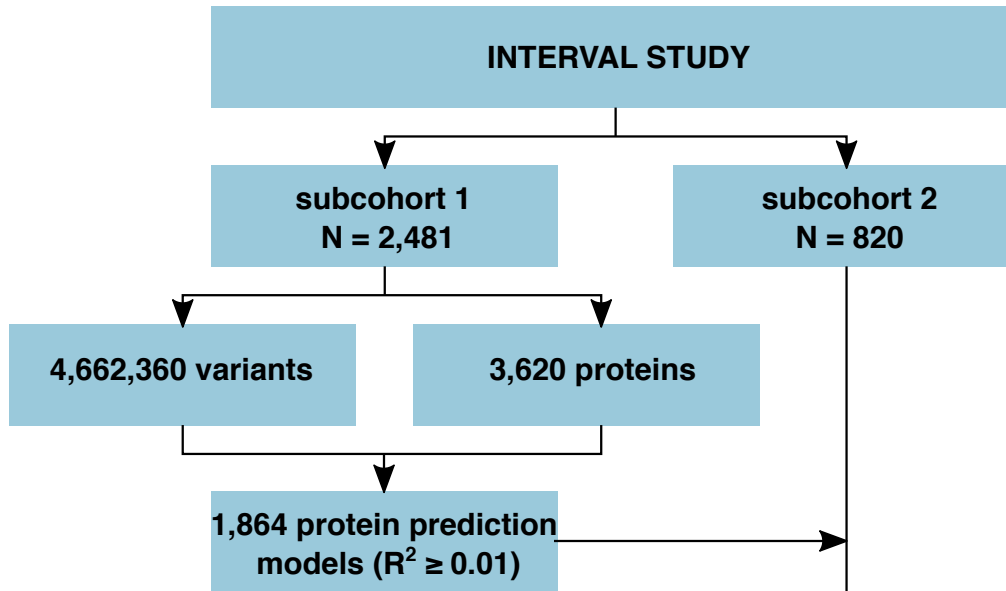
				DB14751	Mecasermin rinfabate	agonist	MDNA
				DB09456	Insulin beef	activator	MDNA
				DB08804	Nandrolone decanoate	inducer	-5.8
				DB01277	Mecasermin	agonist	3DSNA
				DB00030	Insulin human	activator	MDNA
				DB06343	Teprotumumab	binder, antibody	3DSNA
				DB12267	Brigatinib	inhibitor	-5.7
				DB00047	Insulin glargine	agonist	MDNA
				DB00071	Insulin pork	binder	MDNA
				DB01307	Insulin detemir	agonist	MDNA
				DB00046	Insulin lispro	agonist	MDNA
				DB01306	Insulin aspart	agonist	MDNA
				DB01309	Insulin glulisine	agonist	MDNA
				DB09564	Insulin degludec	agonist	MDNA
				DB09129	Chromic chloride	activator	MDNA
				DB14751	Mecasermin rinfabate	NA	MDNA
				DB09456	Insulin beef	agonist	MDNA
				DB00030	Insulin human	agonist	MDNA
				DB01277	Mecasermin	NA	3DSNA
				DB12267	Brigatinib	binding	<b>-8.4</b>
IR	Insulin receptor	<i>INSR</i>	0.013	DB12010	Fostamatinib	inhibitor	<b>-7.5</b>

\* a score of  $\leq -7$  represents a good interaction between the protein and corresponding drug agent and is bolded.

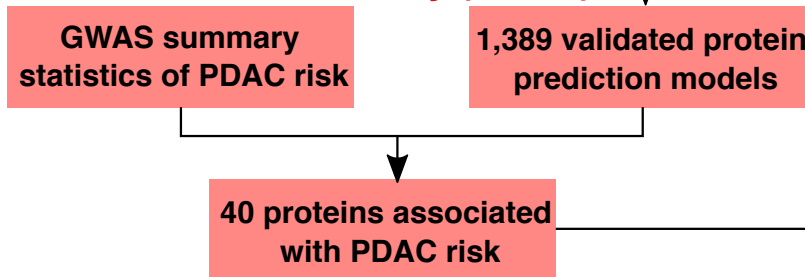
MDNA: Molecular docking not applicable

3DSNA: 3D structure not available.

### ① Establish protein prediction models



### ② Proteome-wide association study (PWAS)



### ③ Downstream analysis

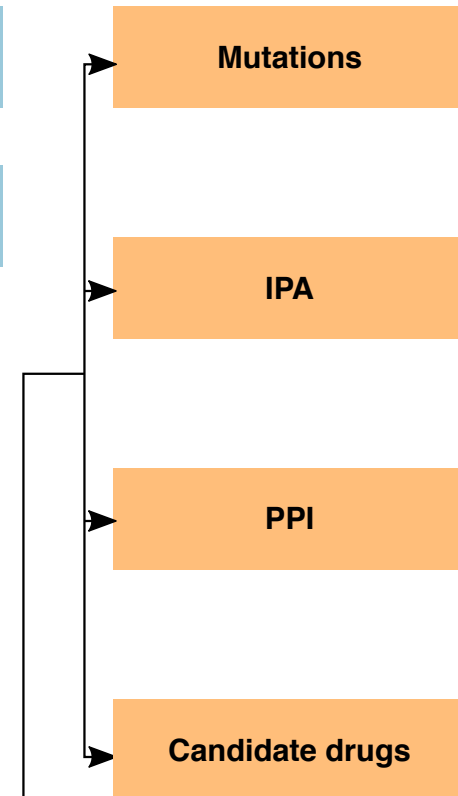


Figure 2

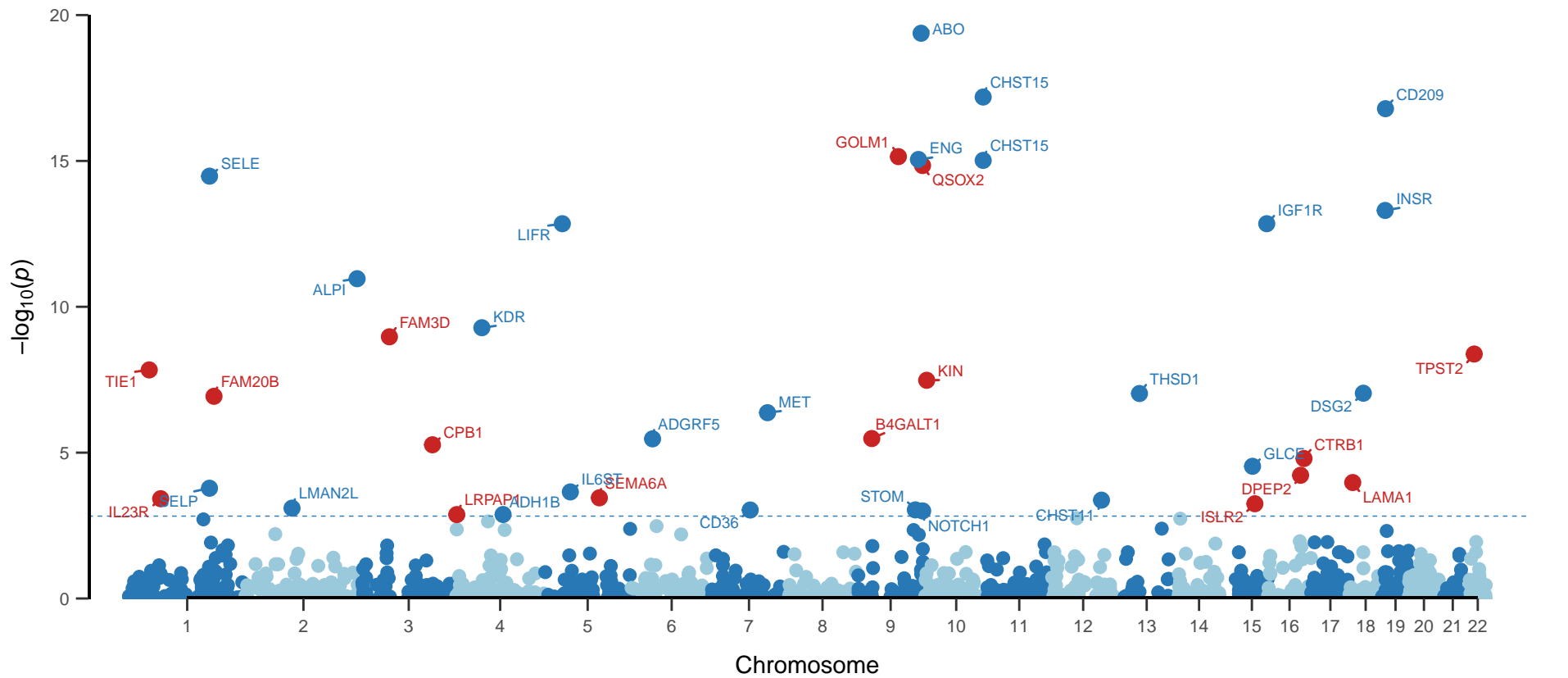
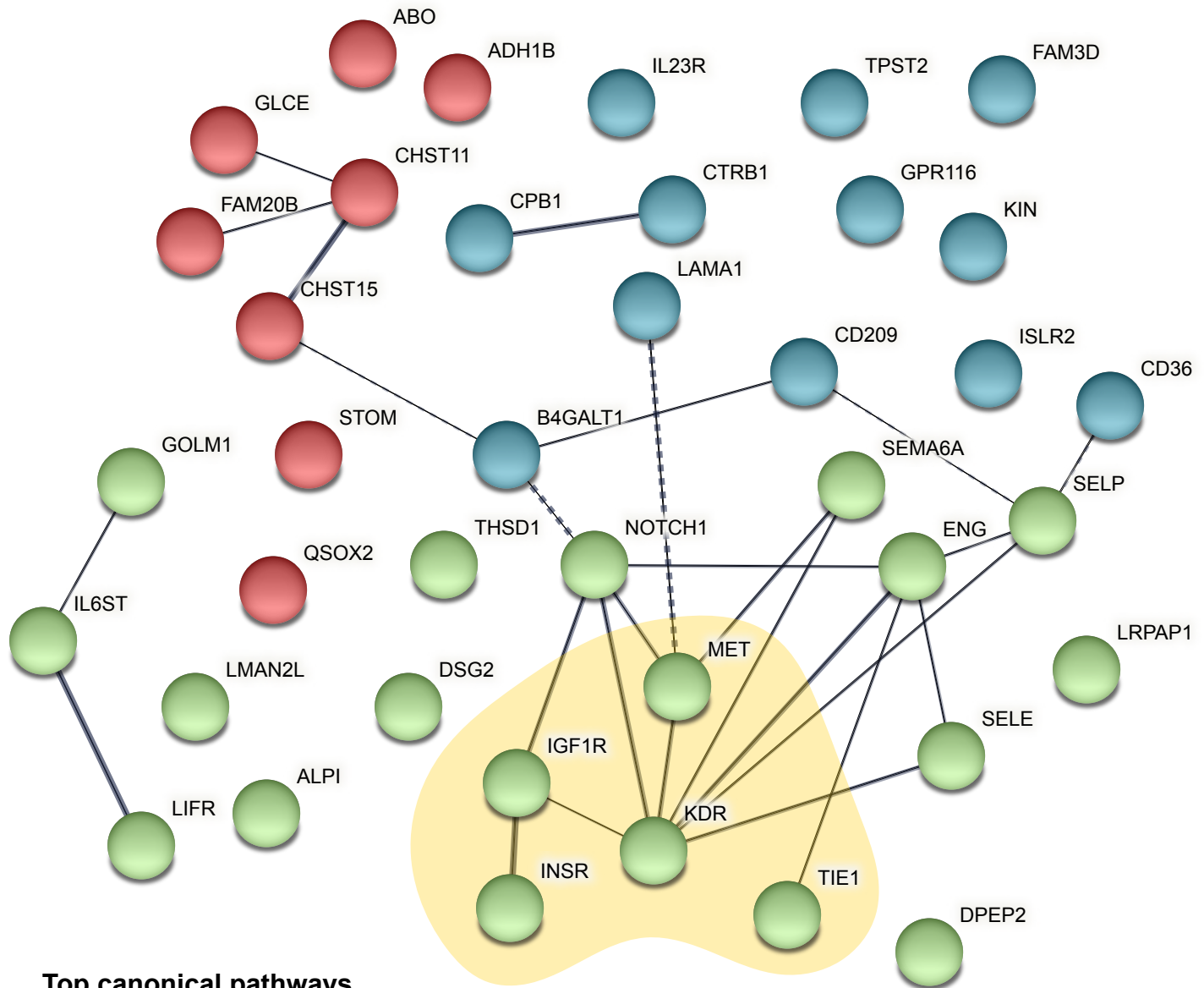
[Click here to access/download;Figure;Figure 2.pdf](#)

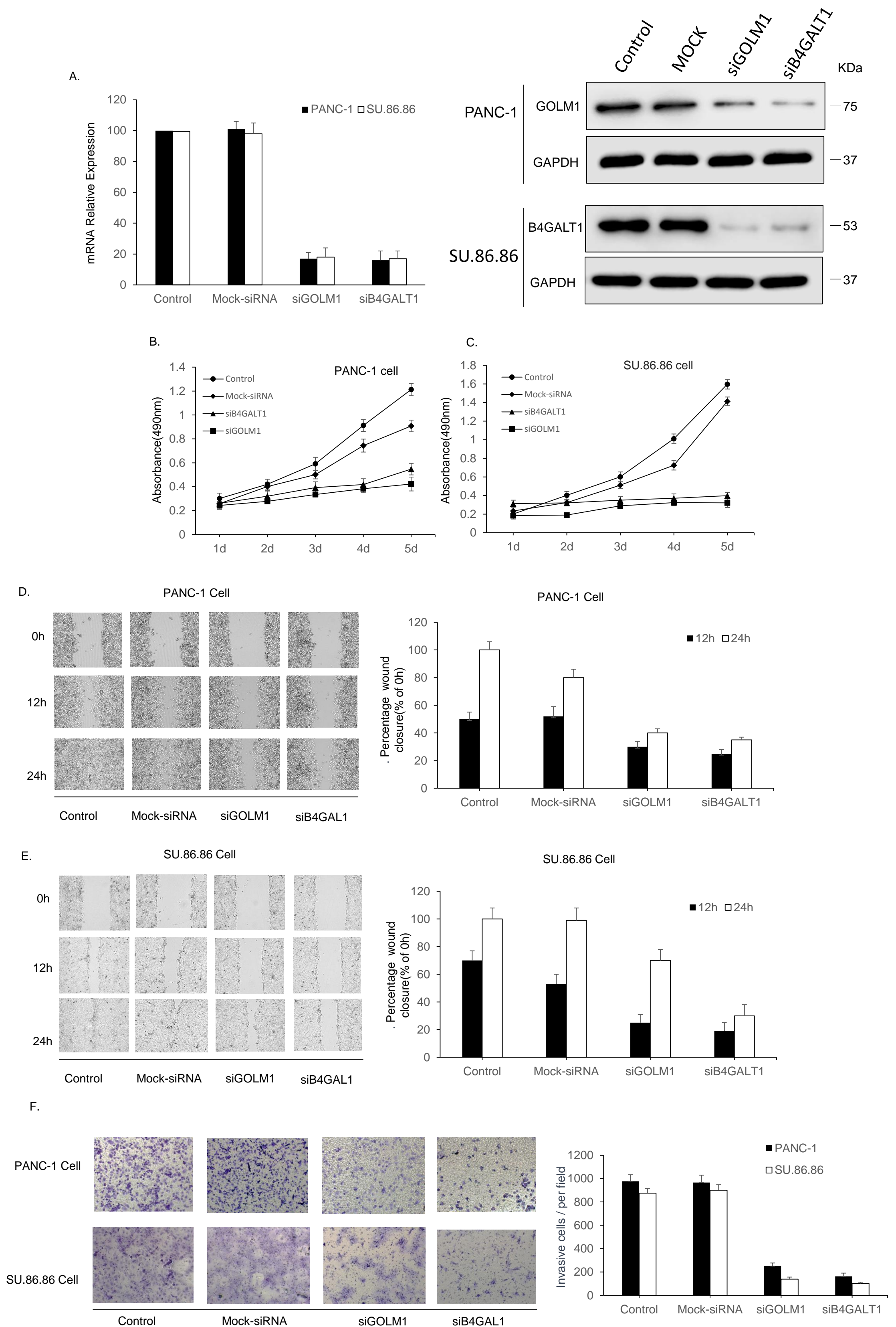
Figure 3

[Click here to access/download;Figure;Figure 3.pdf](#)



**Top canonical pathways**

■ IL-15 Production,  $p\text{-value} = 2.19 \times 10^{-3}$





Click here to access/download  
**Supplementary Material**  
Supplementary File.docx



UNIVERSITY OF HAWAI'I  
**CANCER CENTER**

October 23, 2023

Dr. Scott Edmunds  
Editor-in-Chief, *GigaScience*

**Proteome-wide association study and functional validation identify novel protein markers for pancreatic ductal adenocarcinoma**

Dear Dr. Edmunds:

My colleagues and I would like to submit this manuscript for publication consideration in *GigaScience*.

As you know, pancreatic ductal adenocarcinoma (PDAC) is a lethal malignancy with few known risk factors and biomarkers. Identifying biomarkers is critical for understanding the pathogenesis of this deadly cancer and developing novel therapeutic approaches. Several blood protein biomarkers have been reported to be linked to PDAC in previous studies, however, findings are often inconsistent, potentially due to common biases existing in the conventional epidemiologic study design. One alternative study design is to use genetic instruments to identify proteins whose genetically predicted levels in blood are associated with PDAC risk. This is a design similar to the popular transcriptome-wide association study (TWAS), but focusing on protein expression levels, a novel design that is rarely explored. It is challenging to construct satisfactory genetic prediction models for protein expression levels, because there are far more pQTLs in trans regions than in cis regions.

In this study, we applied a highly novel study to develop comprehensive protein genetic prediction models by considering both cis- and trans-acting elements as instruments for identifying novel PDAC related proteins. We leveraged genome and plasma proteome data of 2,481 healthy European descendants included in the INTERVAL study to establish such prediction models. We selected models with a prediction performance of  $>0.01$  in both internal and external validation for association analyses with PDAC risk, by analyzing 8,275 cases and 6,723 controls of European descent from the Pancreatic Cancer Cohort Consortium and the Pancreatic Cancer Case-Control Consortium.

We identified significant associations between predicted concentrations of 40 proteins and PDAC risk at a false discovery rate of  $< 0.05$ , including 16 novel proteins. For 29 of the genes encoding identified proteins, somatic level potentially functional mutations were detected in PDAC patients in The Cancer Genome Atlas. Relevant protein-encoding genes were also significantly enriched in several cancer-related pathways. We further identified drugs targeting



the identified proteins, which may serve as candidates for drug repurposing for treating PDAC. We also silenced two of the novel protein-encoding genes and observed critical roles of *GOLM1* and *B4GALT1* in driving PDAC cell proliferation, migration, and invasion, by testing two independent cell lines. Our functional characterization further supported critical roles of identified novel proteins in pancreatic tumorigenesis.

We believe that our manuscript should be of great interest to the scientific community served by *GigaScience*. In particular, our study could serve as an excellent model for future research that integrates large genomics and proteomics data to understand the genetics and biology of diseases in the post-GWAS era. We hope that you will find our work interesting and would be willing to consider it for publication in your journal.

Sincerely,

Lang Wu, Ph.D.  
Director, Pacific Center for Genome Research  
Associate Professor, University of Hawaii Cancer Center  
University of Hawaii  
Email: [lwu@cc.hawaii.edu](mailto:lwu@cc.hawaii.edu)