# nature research

# Peer Review Information

**Journal:** Nature Genetics
**Manuscript Title:** Genetics of 35 blood and urine biomarkers in the UK Biobank
**Corresponding author name(s):** Dr. Manuel Rivas

## Editorial Notes:

*EA delete any non-applicable rows, and then delete this instruction.*

## Reviewer Comments & Decisions:

**Decision Letter, initial version:**

**Date:** 12th Aug 20 11:40:08
**Last Sent:** 12th Aug 20 11:40:08
**Triggered By:** Catherine Potenski
**From:** Catherine.Potenski@us.nature.com
**To:** mrivas@stanford.edu
**Subject:** Decision on Nature Genetics submission NG-A55370-T
**Message:** 12th Aug 2020

Dear Dr Rivas,

Your Article, "Genetics of 35 blood and urine biomarkers in the UK Biobank" has now been seen by 3 referees. You will see from their comments below that while they find your work of interest, some important points are raised. We are interested in the possibility of publishing your study in Nature Genetics, but would like to consider your response to these concerns in the form of a revised manuscript before we make a final decision on publication.

1

As you will see, Reviewer #1 thinks that the study is comprehensive and likely will be well cited. There are a series of questions, requests for clarification and comments that should be addressed. Reviewer #2 echoes some of the comments from the previous round about exactly what the novel finding are. Additionally, this reviewer has some questions about the fine-mapping. Reviewer #3 is satisfied and has one minor point. All reviewers are mainly supportive. The points raised are important, but seem addressable.

We therefore invite you to revise your manuscript taking into account all reviewer and editor comments. Please highlight all changes in the manuscript text file. At this stage we will need you to upload a copy of the manuscript in MS Word .docx or similar editable format.

We are committed to providing a fair and constructive peer-review process. Do not hesitate to contact us if there are specific requests from the reviewers that you believe are technically impossible or unlikely to yield a meaningful outcome.

When revising your manuscript:

*1) Include a "Response to referees" document detailing, point-by-point, how you addressed each referee comment. If no action was taken to address a point, you must provide a compelling argument. This response will be sent back to the referees along with the revised manuscript.

*2) If you have not done so already please begin to revise your manuscript so that it conforms to our Article format instructions, available <a href="http://www.nature.com/ng/authors/article_types/index.html">here</a>. Refer also to any guidelines provided in this letter.

*3) Include a revised version of any required Reporting Summary:
https://www.nature.com/documents/nr-reporting-summary.pdf
It will be available to referees (and, potentially, statisticians) to aid in their evaluation if the manuscript goes back for peer review.
A revised checklist is essential for re-review of the paper.

Please be aware of our <a href="https://www.nature.com/nature-research/editorial-policies/image-integrity">guidelines on digital image standards.</a>

Please use the link below to submit your revised manuscript and related files:

[REDACTED]


<strong>Note:</strong> This URL links to your confidential home page and associated information about manuscripts you may have submitted, or that you are reviewing for us. If you wish to forward this email to co-authors, please delete the link to your homepage.

We hope to receive your revised manuscript within four to eight weeks. If you cannot send it within this time, please let us know.

Please do not hesitate to contact me if you have any questions or would like to discuss these revisions further.

Nature Genetics is committed to improving transparency in authorship. As part of our efforts in this direction, we are now requesting that all authors identified as 'corresponding author' on published papers create and link their Open Researcher and Contributor Identifier (ORCID) with their account on the Manuscript Tracking System (MTS), prior to acceptance. ORCID helps the scientific community achieve unambiguous attribution of all scholarly contributions. You can create and link your ORCID from the home page of the MTS by clicking on 'Modify my Springer Nature account'. For more information please visit please visit <a href="http://www.springernature.com/orcid">www.springernature.com/orcid</a>.

We look forward to seeing the revised manuscript and thank you for the opportunity to review your work.

Thank you very much.

All the best,

Catherine

--

Catherine Potenski, PhD
Chief Editor
Nature Genetics
1 NY Plaza, 47th Fl.
New York, NY 10004
catherine.potenski@us.nature.com
https://orcid.org/0000-0002-4843-7071

Reviewers' Comments:

Reviewer #1:
Remarks to the Author:
Sinnott-Armstrong and colleagues present and updated manuscript focused on the genetic determinants of 35 blood or urine biomarkers in > 300K unrelated participants of the UK Biobank. This manuscript reflects a large and important effort, and will be of interest to the general genetics community.

My comments reflect some persistent gaps in clarity that I think are important to close.

Major:

Authors describe biomarker phenotype distributions to start the paper, but then

describe correction for statin medications — is this example relevant to the remaining biomarkers or used only for cholesterol? If it was just these, it doesn't seem like it warrants high profile discussion in the main text. Or if it is felt to be important, then why were other biomarkers (e.g. A1c, triglycerides) similarly affected by medications also adjusted?

Authors describe correction of each biomarker using a residualized model for multiple potential confounders. But I don't have a clear sense of whether this was necessary or helpful as compared to just using the raw biomarkers values and correcting for age and sex. The only place this is addressed is in the paragraph saying 'we recover previously estimated phenotype correlations' but here again I don't know if this was meaningfully better than just adjusting for age and sex — if there was significant improvement, this seems important since it would influence how many other might approach analyzing UKBB biomarkers.

Authors describe fasting glucose as a phenotype - was this actually correct? Most UKBB participants were not fasting.

Is it expected that lp(a) SNP heritability is so low only 0.6% given that is known to be highly heritable?

ST6 row 54 - is it really true that a frameshift variant in SLC22A2 is associated with lp(a) or is this some sort of coding error/typo?

Page 10, line 299 - is isn expected that variants with the same effect on rate would have opposite effects on risk of grout?

With respect to polygenic score component, there are two interesting questions at hand: (a) does the BASIL algorithm which allows for >1M genotypes as predictors enable a better polygenic predictor than a model using just the GWAS summary statistics? and (b) does the prediction of disease outcomes using a combination of biomarker PRS outperform one based just on the disease?

Here, I remain quite confused. Authors profile renal failure as a illustrative example of the PRS. First, I am unable to understand how renal failure was defined, nor is this a standard clinical term with obvious meaning. Second, how is it possible that the snpnet PRS had OR < 1 and AUC < 0.5 (i.e. worse than chance)? Showing improvement over a model that is worse than chance seems like a strange thing to highlight.

Beyond these clarity issues, some of the biomarkers are renal-related biomarkers used to define kidney diseases — so a question is whether use of OTHER biomarkers is adding something or you are best off just combining those that are related?

What were the biomarker PRSs that were used to improve the heart failure and cirrhosis scores?

Among the multi-PRS scores, is it possible to understand/display which PRS contribute the most?


With respect to whether the multi-PRS improve beyond existing PRS for a given

disease, I am also having trouble interpreting. In the GRS46K for CAD, I am not able to understand the approach for the four scores summarized. For the gallstones, is it the case that the multiPRS actually makes the results worse?

Overall, the fundamental question of whether adding PRS of biomarkers to PRS based on disease states meaningfully improves prediction remains unanswered in the revision. Many diseases such as T2D have previously validated PRS available, so uncertain why only CAD was analyzed. Moreover, HR close to 1.5/SD have been reported for many diseases, so the values in Fig 5D gray bars seem very low.

Similarly in ST22, HR of 1.16 for a MI PRS baseline model seems very low.

Does multiPRS approach improve or worsen transportability ocross racial groups?

Can authors give generalizable learnings about when incorporating biomarker PRS might be helpful? For example, it seems like when the disease GWAS is very large, it may be less useful, but more useful when biomarkers are strongly associated with disease, discovery GWAS smaller, etc.

I am not able to individually vet each of the proposed causal inferences, but many of them do not seem particularly biologically plausible. e.g. lowering TG would decrease risk of anorexia?

The list of phenotypes used in PheWAS seems arbitrary/duplicative/not reflective of important conditions. For example, authors describe 'removed disease outcomes that were likely to be duplicated' but then have separate rows fro 'DVT diagnosed by doctor' and 'blood clot or DVT diagnosed by doctor.' Similarly with 'bipolar and major depression status any' and 'bipolar and major depression.' For the 'status any' version, what does that mean?

I wonder if merging diagnosis codes into concepts PheCODES as is available using codings described by the Vanderbilt group, or just using an all by all and not suggesting there was curation would be reasonable alternatives, but best would really be to more effectively curate into disease groupings. Is 'dentures' really an important phenotype?

Minor:
First sentence: I don't think it's true that biomarkers are the 'primary clinical tools' fo adverse health conditions

Authors describe that 90% of variance in testosterone was explained by covariates, but the majority of this is presumably just sex.

Semantics throughout paper are not well-reflective of terms used in clinical practice, suggest review by one of the clinician co-authors. Representative examples include ('kidney problems', 'renal failure', 'quantitative disease symptoms','t2d' in Figure)

Semantics such as 'McCarthy Lab Tools' should be removed

LOR/SD difficult to interpret — any reason to not report odds ratio per SD?

Lipoprotein A is referred to as lipoprotein(a) within clinical research and clinical trial paradigms

Semantics lack precision and warrant proofreading. A representative example is Supp Figure 7A, entitle 'Portability of individual biomarkers,' but I think this may refer to the relative variance explained of polygenic scores for these biomarkers?

Reviewer #2:
Remarks to the Author:
1. What is novel among the many discoveries? I see that some of this information is listed in the Supplementary Tables, but it would be helpful for the reader to also appreciate novelty in the main text. In particular, known/novel information for the different variants and genes listed in the section entitled "Biomarkers associated variants prioritize therapeutic targets" would be very informative. This had already been raised in the previous review.

2. What is the FINEMAP posterior inclusion probability (PIP) for the potential "functional" variants highlighted in section "Biomarkers associated variants prioritize therapeutic targets"?

3. For the FINEMAP analyses, how were "loci" defined? Based on physical distance around independently associated variants? Did the authors merge overlapping loci? More info needed in the Methods section.

4. Figure 3.
Panel B. Consider adding % variance explained by fine-mapped variants for each of the horizontal bar.
Panel D. I cannot see the palest color.

5. In STables 14 A-B-C, what is the meaning of the color code for the different cells?

Reviewer #3:
Remarks to the Author:
I am comfortable with the changes, with one slight exception. I think the HNF1B association with renal failure etc. should be explored relative to diabetes. The response to reviewer indicated the association was substantially attenuated (p of 5x10-7 to 0.035) when T2D cases were removed. I think an additional sentence in the manuscript suggesting that the renal failure may be happening primarily in diabetics would help researchers trying to follow-up this finding.

Nice work.

**Author Rebuttal to Initial comments**

# nature research

Medical School Office Building
1265 Welch Road, 3rd floor
Stanford, CA 94305-5464
mrivas@stanford.edu
http://rivaslab.stanford.edu

2020/10/22

Dear Dr. Catherine Potenski,

Please find the review response and revision regarding our manuscript "**Genetics of 35 blood and urine biomarkers in the UK Biobank**" (NG-A55370-T).

We thank the reviewers for their supportive and constructive comments and their time. In the revised manuscript, we have addressed all the remaining questions and clarifications that were not adequately addressed in earlier versions of the manuscript. Following their suggestions, we have also highlighted the novel findings. We believe that those changes have significantly improved the manuscript.

Manuel A. Rivas,
Assistant Professor of Biomedical Data Science

Our responses to the reviewers' individual comments below are in blue font, the comments from the reviewer are copied in black, and quoted texts from the updated manuscript are shown in gray with a vertical bar (examples are shown below):

This is an example of the reviewer's comments
This is an example of our response.
> This is an example of quoted texts from the updated manuscript

Referee #1:

Sinnott-Armstrong and colleagues present an updated manuscript focused on the genetic determinants of 35 blood or urine biomarkers in > 300K unrelated participants of the UK Biobank. This manuscript reflects a large and important effort, and will be of interest to the general genetics community.

My comments reflect some persistent gaps in clarity that I think are important to close.

Thank you very much for taking the time to review the manuscript and for providing detailed feedback. We are confident that your comments have improved the clarity of the manuscript.
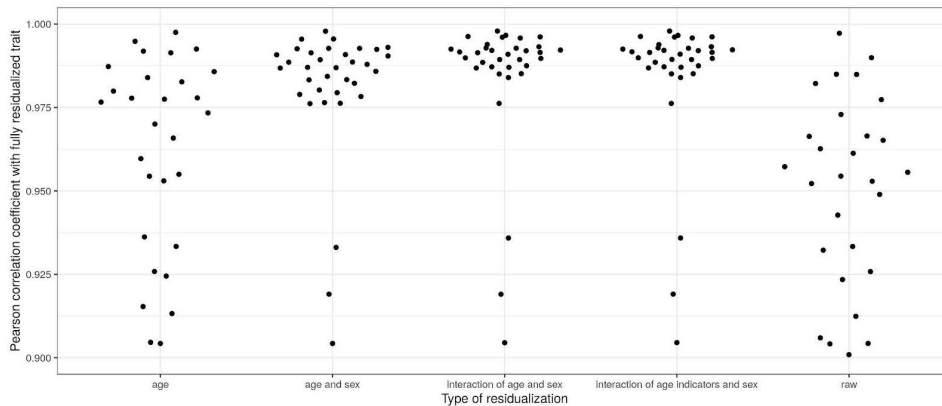
Major:

Authors describe biomarker phenotype distributions to start the paper, but then describe correction for statin medications — is this example relevant to the remaining biomarkers or used only for cholesterol? If it was just these, it doesn't seem like it warrants high profile discussion in the main text. Or if it is felt to be important, then why were other biomarkers (e.g. A1c, triglycerides) similarly affected by medications also adjusted?

Thank you for this comment. We performed a comprehensive analysis of the statin medications, and through this analysis, found that age, sex, and BMI explain the majority of differences in biomarker levels between assessments and that statins do not have an independent effect except for LDL, total cholesterol, and apolipoprotein B levels. As such, we have only adjusted those three biomarkers for downstream analysis. As you note, this procedure is standard and we agree that having it be so high profile is unwarranted. We have kept the description of this adjustment in the Methods section.

Authors describe correction of each biomarker using a residualized model for multiple potential confounders. But I don't have a clear sense of whether this was necessary or helpful as compared to just using the raw biomarkers values and correcting for age and sex. The only place this is addressed is in the paragraph saying 'we recover previously estimated phenotype correlations' but here again I don't know if this was meaningfully better than just adjusting for age and sex — if there was significant improvement, this seems important since it would influence how many other might approach analyzing UKBB biomarkers.

Thank you for noting this. We have performed a number of more specific regressions regarding the residualized phenotypes (in **Supplementary Table 4B**) and the phenotype-PRS correlations (in **Supplementary Table 17B**). We hope that these two clarify the importance of additional covariates. In general, most traits are quite correlated between the totally un-residualized values and the fully adjusted values (~0.95), suggesting that in many cases no covariates of any kind are required. The improvement for individual covariates is substantial in a small number of cases, such as for testosterone and sex, and generally the correlation is around 0.98 for the age- and sex-only residuals, with only Lp(a) being statistically indistinguishable with the fully residualized traits. This suggests that there are meaningful benefits to including additional covariates, but that most results will remain interpretable without them.



The three outlier traits with lower correlations in the more adjustments are Vitamin D (for which month of assessment is quite important), Phosphate (for which assessment centre is quite important), and Glucose (for which fasting time is quite important).

Authors describe fasting glucose as a phenotype - was this actually correct? Most UKBB participants were not fasting.
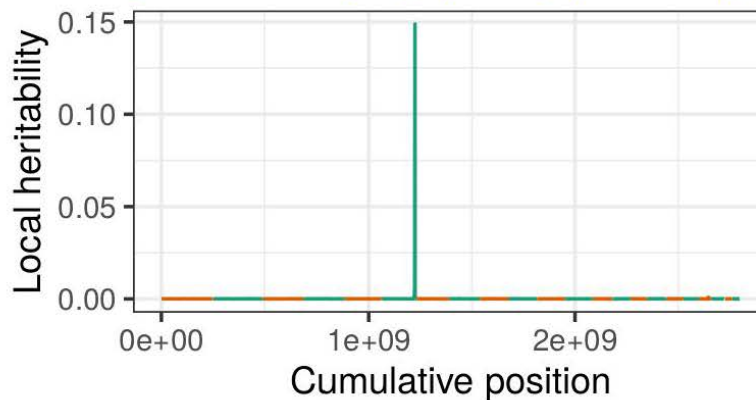
We apologize for the confusion with this and thank you for pointing it out. Indeed, we performed analyses of both random glucose -- for which all individuals, regardless of fasting time, were included -- as well as the ~17,000 individuals who were self-reported fasting and had valid glucose measurements. However, due to the power limitations, we focus primarily on the random glucose measurements for our study. We have clarified that in the text and cite the new wording below:

Page 20 line 596

> *Phenotype and covariate quality control excluded rheumatoid factor and oestradiol from further analyses, and fasting glucose (available for 17,439 self-reported fasting individuals) was used as a phenotype-level quality control for the glucose measurements -- throughout the text, "glucose" refers to glucose levels adjusted for fasting time rather than the GWAS among only fasting individuals (self-report of more than 7 and less than 24 hours of fasting, n = 17,439) unless otherwise noted*

Is it expected that lp(a) SNP heritability is so low only 0.6% given that is known to be highly heritable?

We too found this puzzling, and that was one of the motivations for running HESS. In fact, when examining the local heritability of lp(a), it is very specific to the lp(a) locus:

As such, we believe that the LD Score regression accurately captures the limited polygenic component of the SNP-heritability of lp(a). On the other hand, lp(a)'s heritability under HESS's fixed effects model is 0.237, consistent with the strong association observed at the lp(a) locus itself.

We added a brief explanation to clarify this in the main text (page 6 line 132):

*Estimates were lower in LD Score regression than HESS for traits with low polygenicity (e.g. Lipoprotein A, $h^2_{ldsc}$ = 0.6% and $h^2_{HESS}$ = 24%), as LD Score regression primarily estimates polygenic heritability[24].*

ST6 row 54 - is it really true that a frameshift variant in SLC22A2 is associated with lp(a) or is this some sort of coding error/typo?

Thank you very much for clarifying the reported association between rs8177505 and lp(a). The consequence of the variant is annotated as frameshift in dbSNP (https://www.ncbi.nlm.nih.gov/snp/rs8177505). By performing the following additional analyses, we tested whether the identified association signal can be explained by other genetic variants.

First, we performed a qualitative assessment of the genotyping quality of the variant. We generated and inspected intensity plots with ScatterShot using the "UKB—All Participants" module (http://mccarthy.well.ox.ac.uk/).

**Figure**. The intensity plot of rs8177505, the frameshift variant in *SLC22A2*.

Second, to investigate the possibility that the identified frameshift allele is just a tagging variant of other causal variants in the region, we performed a conditional analysis. Specifically, we ran an association with variants (MAF > 0.1% and INFO > 0.3 in White British individuals) within 200Kb of rs8177505 both unadjusted (beta = 1.44208, p = $9.46191 \times 10^{-233}$) and adjusted for the genotype at rs8177505 ("rs8177505"); for all the >99% posterior probability fine-mapped SNPs from the study ("Finemapped" -- beta = 0.213102, p = $1.1 \times 10^{-8}$); for the list of all variants and stepwise independent variants from an existing study ("NEJMAll" and "NEJMStepwise" -- PMID 20032323, beta = 0.157889 and 0.175442, p = $1.3 \times 10^{-5}$ and $1.3 \times 10^{-6}$); and for the list of independent and copy-number associated variants from another existing study ("TOPMed" and "TOPMedKIV" -- PMID 29973585, beta = 0.27945 and 0.19539, p = $3.3 \times 10^{-16}$ and $1.1 \times 10^{-6}$). The resulting plots are as follows:

**Supplementary Figure 14. Extended comparison of conditional effects at rs8177505 and the *LPA* locus.**
The putative LpA-associated variant rs8177505 is shown in red, including the marginal effect "unconditional" (rs8177505 beta = 1.44208, p = 9.46191 x $10^{-233}$); the effect sizes conditioned on the genotype at rs8177505 ("rs8177505"); the effect sizes conditioned on all the >99% posterior probability fine-mapped SNPs from the study ("Finemapped" -- rs8177505 beta = 0.213102, p = 1.1 x $10^{-8}$); for the list of all variants and stepwise independent

variants from a previous study ("NEJMAll" and "NEJMStepwise" (Clarke et al. 2009) -- rs8177505 beta = 0.157889 and 0.175442, p = 1.3 x $10^{-5}$ and 1.3 x $10^{-6}$); and for the list of independent and Kringle IV repeat associated variants from another study ("TOPMed" and "TOPMedKIV" (Zekavat et al. 2018) -- rs8177505 beta = 0.27945 and 0.19539, p = 3.3 x $10^{-16}$ and 1.1 x $10^{-6}$).

These plots suggest that rs8177505 is not driving the majority of the association signal at the *LPA* locus. However, the conditional p-values are still marginal, and could indicate that rs8177505 plays a minor role in regulation, which we cannot exclude in this study and requires further analyses. (In addition to the potential for a spurious association, we further note that this variant might have an effect independent of it being a coding variant in *SLC22A2*, such as a non-coding regulatory effect on *LPA*, and/or a linked variant that is not captured by our list of fine-mapped associations or previous studies.)

We added a sentence to the discussion to this effect (page 18 line 496):
*In addition, the large and complex linkage present at some loci, including notably the LPA locus, might result in spurious fine-mapped and rare coding variant associations, though conditional analyses of e.g. a rare coding variant in SLC22A2 are inconclusive (***Supplementary Figure 14***).*

Thank you for this feedback and critical evaluation of our findings.

Page 10, line 299 - is expected that variants with the same effect on rate would have opposite effects on risk of grout?

Thank you for clarifying the reported association. In an earlier version of the main text in the manuscript, we reported that two intronic variants in *ABCG2* both have lowering effects on Urate but have opposite effects on Gout. As you correctly pointed out, there was a mistake. The two variants have the opposite effects for both Urate and Gout (**Supplementary Tables 14A**, **13C**).

| chr:pos:ref:alt | rsID | trait | BETA | SE | prob | log10bf |
|---|---|---|---|---|---|---|
| 4:89030920:C:G | rs2622621 | Urate | 0.083 | 2.65e-3 | 1 | 13.4 |
| 4:89074405:T:C | rs13109944 | Urate | -0.085 | 2.52e-3 | 0.999 | 8.00 |

| chr:pos:ref:alt | rsID | trait | OR | LOG(OR)_SE | P |
|---|---|---|---|---|---|
| 4:89030920:C:G | rs2622621 | Gout | 1.384 | 0.019 | 2.83x$10^{-67}$ |
| 4:89074405:T:C | rs13109944 | Gout | 0.717 | 0.019 | 8.18x$10^{-71}$ |

We have updated the texts in the manuscript.

Results section, page 11, line 313:
*An allelic series of two intronic variants in ABCG2 identified for their associations with increasing and lowering urate levels have risk-increasing ($p = 2.8 \times 10^{-67}$, OR = 1.38 [95% CI: 1.33, 1.44]) and protective (OR = 0.72 [95% CI: 0.69, 0.74]) associations with gout, respectively. Both of these associations with gout are also replicated in FinnGen R2 ($p = 6.3 \times 10^{-6}$, OR = 1.25 [95% CI: 1.13, 1.37] and $p = 8.4 \times 10^{-5}$, OR = 0.84 [95% CI: 0.78, 0.92]). Those two variants ($r^2 = 0.47$) have moderately low linkage with a known common protein-altering variant in ABCG2 ($r^2 = 0.22$ and 0.11 in UKB White British for rs2231142 [Q141K]) which contributes to risk of gout[49].*
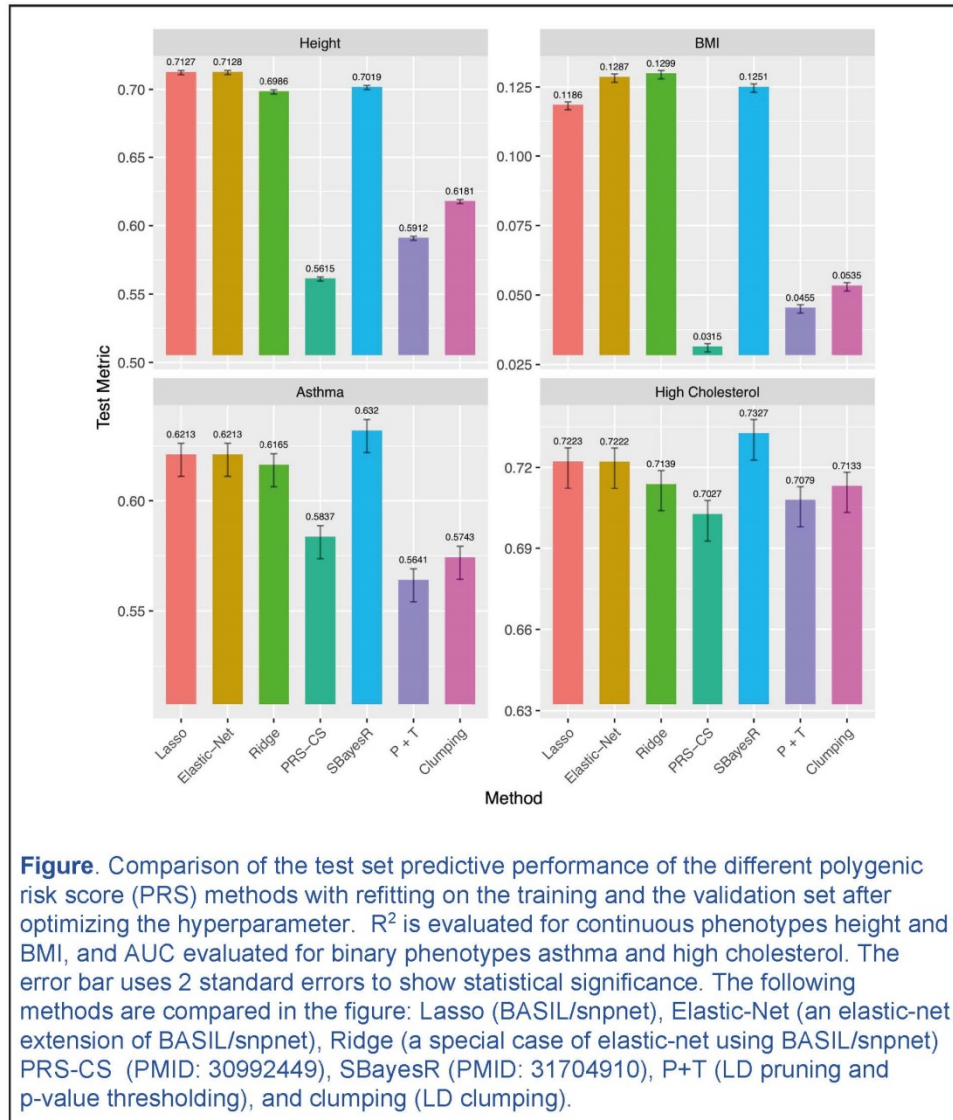
We are sorry for our mistake. Thank you very much for catching this error.

With respect to polygenic score component, there are two interesting questions at hand: (a) does the BASIL algorithm which allows for >1M genotypes as predictors enable a better polygenic predictor than a model using just the GWAS summary statistics?

Thank you very much for clarifying the advantages of the BASIL algorithm and its implementation in R snpnet package, which are described in another manuscript (J. Qian, et al., bioRxiv, 630079, 2019, doi: https://doi.org/10.1101/630079).

In the BASIL/snpnet manuscript, we compared the predictive performance of BASIL/snpnet against other modern PRS methods that takes GWAS summary statistics as input, including SBayesR and PRS-CS, using two quantitative traits (standing height and body mass index) and two binary disease outcomes (high cholesterol and asthma) from UK Biobank. From those benchmarking analysis, we found that BASIL/snpnet shows a competitive predictive performance. For standing height and body mass index, BASIL/snpnet turned out to be the best performing model. We included a figure from the BASIL/snpnet manuscript below. The full BASIL manuscript is also attached as a "related manuscript file" in this revised submission.

**Figure**. Comparison of the test set predictive performance of the different polygenic risk score (PRS) methods with refitting on the training and the validation set after optimizing the hyperparameter. $R^2$ is evaluated for continuous phenotypes height and BMI, and AUC evaluated for binary phenotypes asthma and high cholesterol. The error bar uses 2 standard errors to show statistical significance. The following methods are compared in the figure: Lasso (BASIL/snpnet), Elastic-Net (an elastic-net extension of BASIL/snpnet), Ridge (a special case of elastic-net using BASIL/snpnet) PRS-CS (PMID: 30992449), SBayesR (PMID: 31704910), P+T (LD pruning and p-value thresholding), and clumping (LD clumping).

and (b) does the prediction of disease outcomes using a combination of biomarker PRS outperform one based just on the disease?
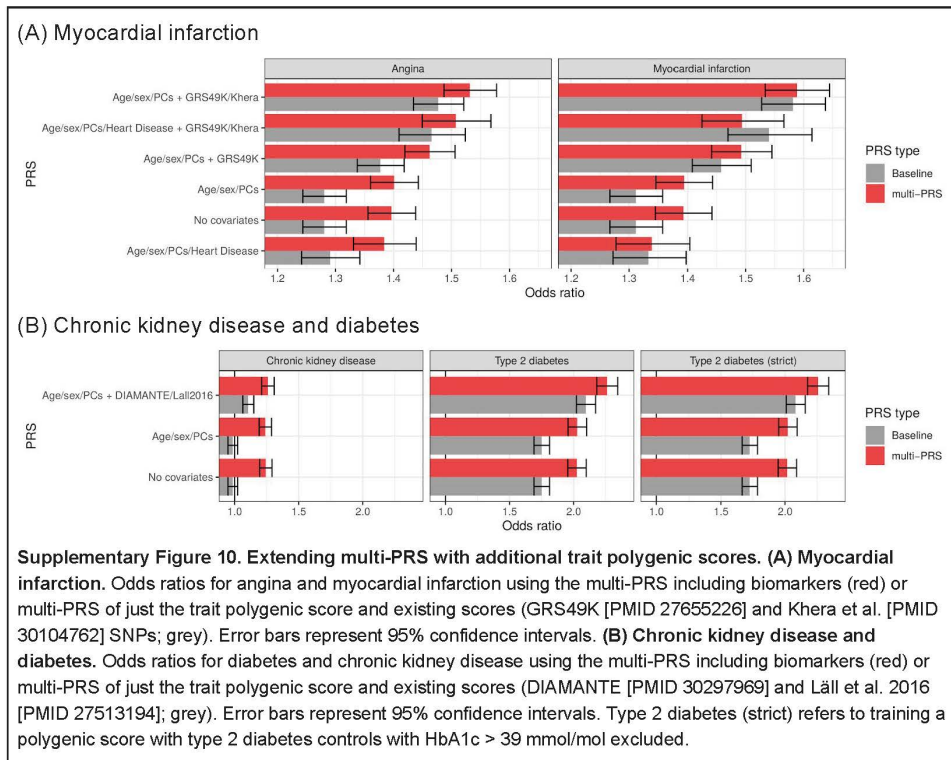
We are sorry that this point was not clearly communicated in the earlier version of the manuscript and revision. Using several disease outcomes as examples, we showed the

improved performance of diseases with multi-PRS compared against trait-specific PRS. For performance evaluation of multi-PRS, we used several different trait-specific PRSs, including the one trained on the same training set in UK Biobank and publicly available PRSs from published literature.

For instance, for chronic kidney disease while adjusting for age, sex, genotyping array, and genotype PCs (**Supplementary Table 20**) we observe a minimally informative PRS for the baseline model which just includes the UK Biobank training set-derived PRS (AUC = 0.495; log odds = -0.014, p = 0.045, 95% CI -0.053 - 0.024) and a somewhat informative PRS for the multi-PRS model which additionally fits the 35 biomarker traits (AUC = 0.561; log odds = 0.215, p = 9.57e-29, 95% CI 0.176 - 0.254). This suggests that including the multi-PRS components from the biomarker traits into the model does improve accuracy substantially, which we confirm in FinnGen (C-index with covariates 0.665; hazard ratio 0.99, p = 0.46, 95% CI 0.95-1.02 for just the chronic kidney disease PRS and C-index with covariates 0.667, hazard ratio 1.12, p = 2.1e-10, 95% CI 1.08-1.16 also including the 35 biomarkers). While existing scores for chronic kidney disease are not available, type 2 diabetes is a major known risk factor for chronic kidney disease. To illustrate that, we fit models which re-weight between the DIAMANTE type 2 diabetes PRS, the Lall 2016 type 2 diabetes PRS, and the training set UKBB PRS. The baseline models which do not include the 35 biomarkers and only include these three PRS does substantially better than just the chronic kidney disease score (AUC = 0.53; log odds = 0.101, p = 1.6e-7, 95% CI 0.062 - 0.140), and further adding the 35 biomarkers does even better (AUC 0.565; log odds = 0.231, p = 9.11e-33, 95% CI 0.192 - 0.269). However, adding the two diabetes PRS does not improve above the model including just the 35 biomarkers and the baseline chronic kidney disease PRS above. Taken together, this suggests that at least for chronic kidney disease, addition of biomarker scores does aid prediction above and beyond the inclusion of related disease traits (type 2 diabetes) or simply training a PRS on the same individuals for chronic kidney disease alone.

To further clarify this, we have added an additional pair of figures, one for myocardial infarction and angina, and the other for type 2 diabetes and chronic kidney disease:

**(A) Myocardial infarction**

**(B) Chronic kidney disease and diabetes**

**Supplementary Figure 10. Extending multi-PRS with additional trait polygenic scores. (A) Myocardial infarction.** Odds ratios for angina and myocardial infarction using the multi-PRS including biomarkers (red) or multi-PRS of just the trait polygenic score and existing scores (GRS49K [PMID 27655226] and Khera et al. [PMID 30104762] SNPs; grey). Error bars represent 95% confidence intervals. **(B) Chronic kidney disease and diabetes.** Odds ratios for diabetes and chronic kidney disease using the multi-PRS including biomarkers (red) or multi-PRS of just the trait polygenic score and existing scores (DIAMANTE [PMID 30297969] and Läll et al. 2016 [PMID 27513194]; grey). Error bars represent 95% confidence intervals. Type 2 diabetes (strict) refers to training a polygenic score with type 2 diabetes controls with HbA1c > 39 mmol/mol excluded.
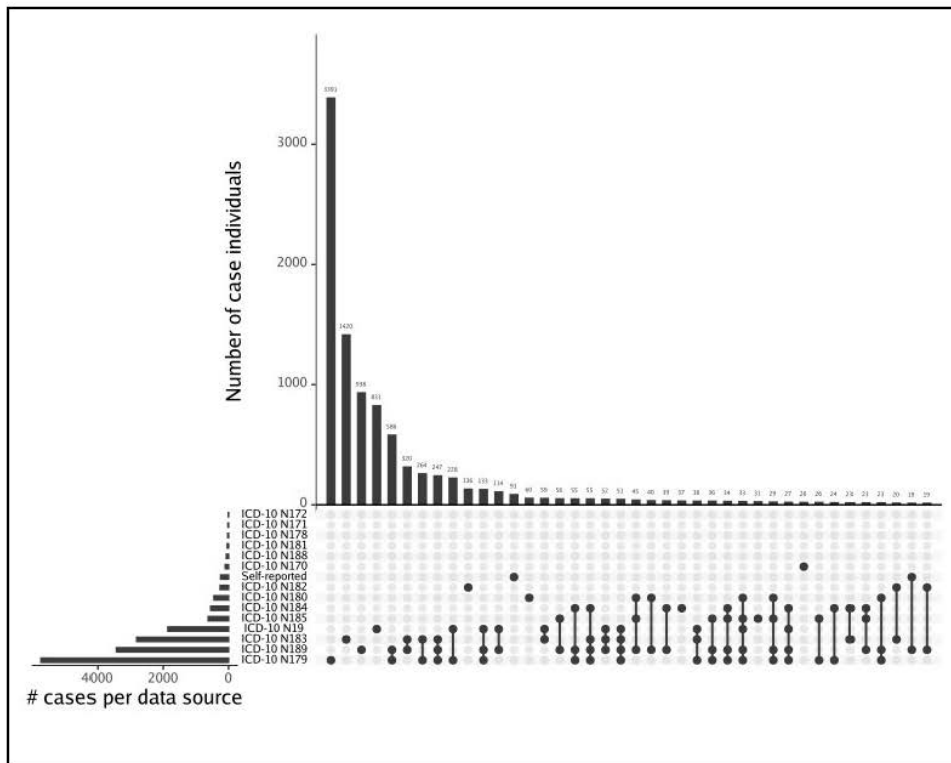
Diabetes, chronic kidney disease, and angina showed consistent improvement over baseline when including multi-PRS terms, while myocardial infarction did not, consistent with the observed AUC improvement (**Supplementary Table 23**) with or without inclusion of pre-existing PRSs. In all cases, the improvements observed were smaller than when excluding these pre-existing polygenic scores, but remain noticeable with the exception of models including family history of heart disease status.

Here, I remain quite confused. Authors profile chronic kidney disease as a illustrative example of the PRS. First, I am unable to understand how renal failure was defined, nor is this a standard clinical term with obvious meaning.

Thank you for clarifying the definition of the disease endpoint used in the study. In our study, the biomarker measurements in the UK Biobank dataset were not directly used in the definition of chronic kidney disease.

In our analysis with UK Biobank, the chronic kidney disease phenotype was defined using a combination of the ICD codes from hospital inpatient records as well as self-reported disease ascertainment status using a procedure described in a previous publication (DeBoever et al "Assessing Digital Phenotyping to Enhance Genetic Studies of Human Diseases" PMID: 32275883). Specifically, we used self-reported chronic kidney disease (coded as "1192" in UKB Data coding ID 6) and ICD-10 code (N17 ["Acute kidney failure"], N18 ["Chronic kidney disease (CKD)"], N19 ["Unspecified kidney failure"], and its sub-concepts) from hospital inpatient data to define this phenotype. Of note, the biomarker measurements were not directly used in the definition of this disease endpoint. We have included a visualization of the breakdown of the data sources below, which is also available as a new phenotyping data source breakdown page in Global Biobank Engine (https://gbe.stanford.edu/RIVAS_HG19/coding_breakdown/HC294).

**Supplementary Figure 8A**. The breakdown of the data sources used for the definition of chronic kidney disease in UK Biobank. The combination of self-reported chronic kidney disease (coded as "1192" in UKB Data coding ID 6) and ICD-10 code (N17 ["Acute kidney failure"], N18 ["Chronic kidney disease (CKD)"], N19 ["Unspecified kidney failure"], and its sub-concepts) from hospital inpatient data are used for the chronic kidney disease definition in UK Biobank. The most common 40 combinations of phenotyping sources are shown in the plot.

In FinnGen, our external validation cohort, disease endpoints including chronic kidney disease are defined based on the ICD codes. A detailed description of the case definition in FinnGen is given in **Supplementary Table 21**. For your convenience, we have quoted the subset of the table showing the definition of chronic kidney disease below.

**Supplementary Table 21.** Description of case definition in FinnGen derived from ICD codes and registry data.

| | Additional definitions | ICD-10 | ICD-9 | Cause of death ICD-10 | Cause of death ICD-9 |
|---|---|---|---|---|---|
| **Chronic kidney disease** | | | | | |
| Acute renal failure | | N17 | 584 | N17 | 584 |
| Chronic kidney diseases | Eligibility for medication reimbursement for anemia in chronic kidney disease | N18 | 585 | N18 | 585 |
| Unspecified kidney failure | | N19 | | N19 | |
| Dialysis | Eligibility for medication reimbursement for uremia requiring dialysis, with ICD-10 N03 or N18 | Z992\|Y841 | | Z992\|Y841 | |

While addressing this comment we found the data sources for the disease endpoints used in our PheWAS analysis (including this renal/kidney failure) were not clearly communicated in the earlier version of the manuscript. We have updated our list of PheWAS disease endpoint phenotypes and added the "Source of the phenotype" column (**Supplementary Table 12**). In case the phenotype was defined using a combination of ICD codes in hospital inpatient records and self-reported disease status, we have provided a link to the phenotyping data source breakdown page in Global

Biobank Engine. We have updated the corresponding Methods section in the manuscript.

Methods section, Page 26, line 810:

*We curated a list of 166 disease outcomes from previously-reported binary phenotypes in Global Biobank Engine (GBE)[27,77,96]. Specifically, we selected disease outcomes with at least 700 cases in white British grouping and removed disease outcomes that were likely to be duplicated (**Supplementary Table 12**). Those phenotypes include non-cancer disease-outcome endpoints derived from a combination of the ICD codes from hospital inpatient records as well as self-reported disease ascertainment status[96], family history phenotype (UK Biobank data category 100034), cancer phenotypes derived from a combination of the UK cancer registry data and questionnaire data[77], and additional set of phenotypes derived from the following data fields in UK Biobank: 2247, 2463, 2834, 3591, 6148, 6149, 6152, 6153, 20126, 20406, 20483, and 21068. The data source of the disease endpoint definitions for each phenotype is described in "Source of the phenotype" column in **Supplementary Table 12**.*

Thank you very much for your question.

Second, how is it possible that the snpnet PRS had OR < 1 and AUC < 0.5 (i.e. worse than chance)? Showing improvement over a model that is worse than chance seems like a strange thing to highlight.

We first randomly split the white British subset in UK Biobank into training (70%), validation (10%), and test (20%) set. The snpnet PRS models were trained on the training set with a hyperparameter optimization using the predictive performance evaluated on the validation set. The predictive performance reported in the manuscript was evaluated on the hold-out test set. We also used the same hold-out test set for the evaluation of the multi-PRS models.

Because we used a hold-out test set, it is possible to have predictive performance not different than chance.

To see whether the "worse than chance" performance is significantly different from the null (AUC = 0.5), we performed an additional statistical test.

- AUC: 0.4953
- The 95% confidence interval (based on bootstrap): [0.4801 - 0.5105]

- P-value (Wilcoxon's rank sum test): 0.5425

This means that our baseline model (snpnt PRS) has a predictive performance equivalent to the random classifier of AUC = 0.5.
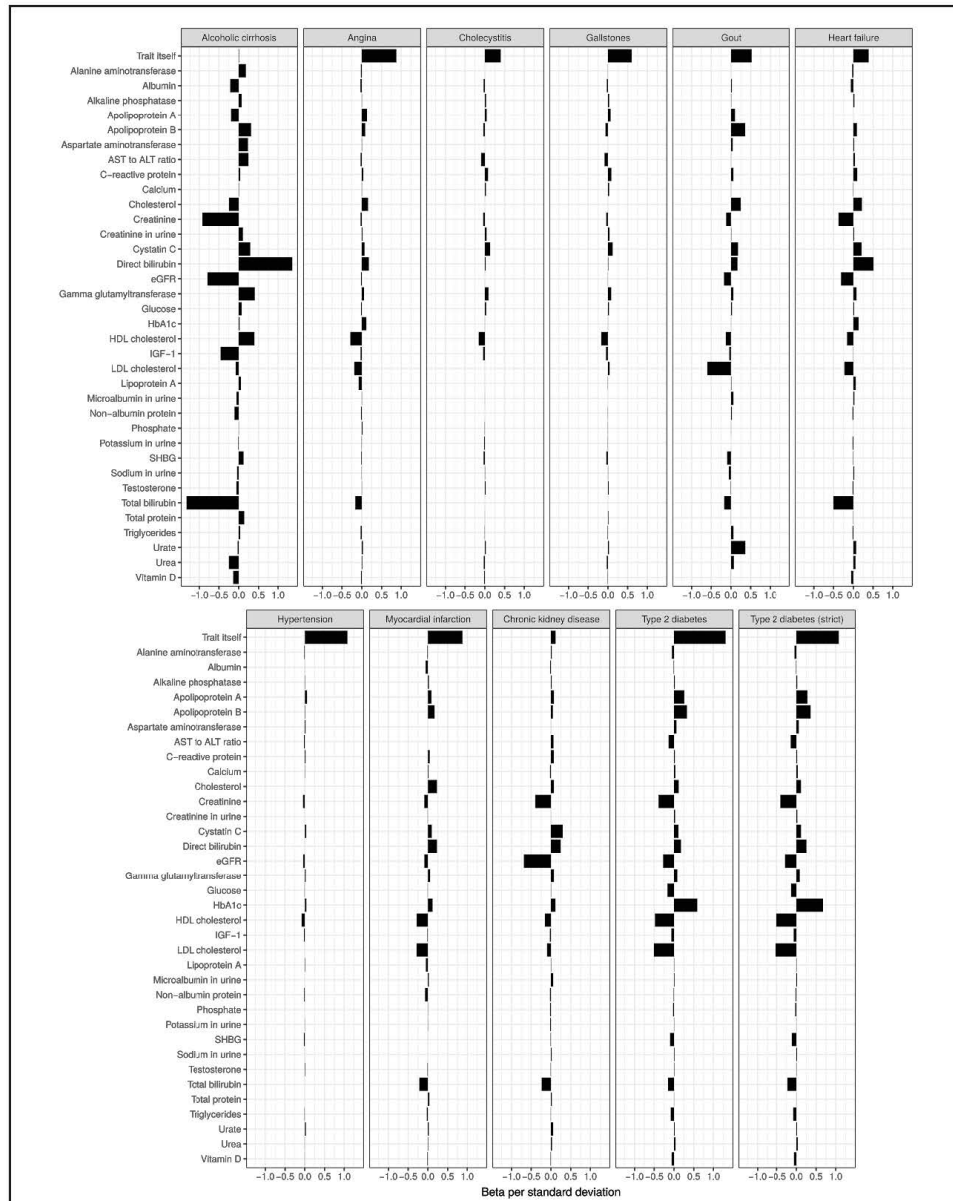
Beyond these clarity issues, some of the biomarkers are renal-related biomarkers used to define kidney diseases — so a question is whether use of OTHER biomarkers is adding something or you are best off just combining those that are related?

What were the biomarker PRSs that were used to improve the heart failure and cirrhosis scores?

Among the multi-PRS scores, is it possible to understand/display which PRS contribute the most?

Thank you very much for your suggestion. In our multi-PRS framework, we indeed provided the PRSs for all of the 35 biomarkers and applied a $L_1$ penalized regression that simultaneously performs variable selection (in this case, selecting the biomarker PRSs that are informative for the disease outcome prediction) and coefficient estimation for each of the selected PRS models.

Following your suggestion, we visualized the coefficients in the penalized regression models (multi-PRS models) to help people interpret the set of biomarker PRSs that have non-zero contribution to the disease prediction in our multi-PRS models. This figure is now added as a new **Supplementary Figure 11**.

**Supplementary Figure 11. Visualization of multi-PRS weights.** Betas per standard deviation fit for each of the multi-PRSs which show, for each biomarker and the trait baseline score, the beta (log odds) of the given outcome

for each standard deviation change in that score. Type 2 diabetes (strict) refers to training a polygenic score with type 2 diabetes controls with HbA1c > 39 mmol/mol excluded.

For chronic kidney disease, we identified eGFR, creatinine, cystatin C, and bilirubin as the main contributors to the polygenic score. For heart failure, in addition to the trait itself, we identified creatinine, bilirubin, total and LDL cholesterol, cystatin C, and eGFR as important contributors. For alcoholic cirrhosis, we identified bilirubin, GGT, eGFR, HDL cholesterol, and IGF-1 levels. Thank you for suggesting these analyses, which we believe substantively improve this work. We added this figure and a corresponding sentence in the main text about this analysis (page 16 line 433):

> *Improved predictions relied on relevant and variable biomarkers across these traits (**Supplementary Figure 11**), including eGFR, creatinine, cystatin C, and bilirubin for CKD; creatinine, bilirubin, total and LDL cholesterol, cystatin C, and eGFR for heart failure; and bilirubin, GGT, eGFR, HDL cholesterol, and IGF-1 for alcoholic cirrhosis.*
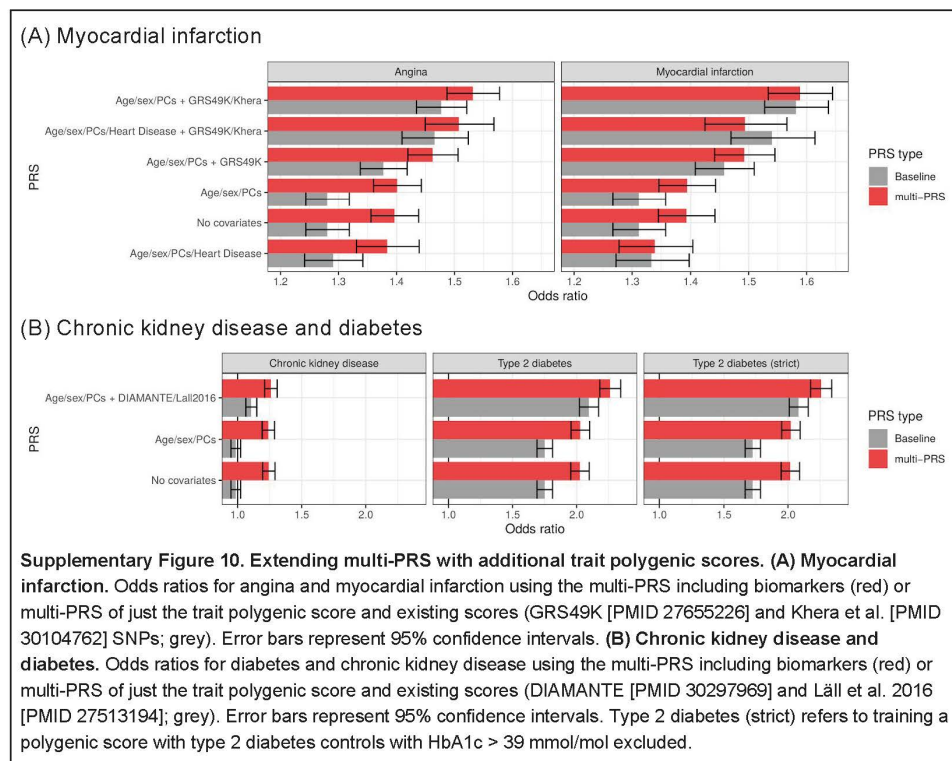
With respect to whether the multi-PRS improve beyond existing PRS for a given disease, I am also having trouble interpreting. In the GRS46K for CAD, I am not able to understand the approach for the four scores summarized. For the gallstones, is it the case that the multiPRS actually makes the results worse?

Thank you for this feedback. We have reformulated these results to clarify the improvements we observe using a new table format. For instance, for chronic kidney disease while adjusting for age, sex, genotyping array, and genotype PCs (**Supplementary Table 20**) we observe a minimally informative PRS for the baseline model which just includes the UK Biobank training set-derived PRS (AUC = 0.495; log odds = -0.014, p = 0.045, 95% CI -0.053 - 0.024) and a somewhat informative PRS for the multi-PRS model which additionally fits the 35 biomarker traits (AUC = 0.561; log odds = 0.215, p = 9.57e-29, 95% CI 0.176 - 0.254). This suggests that including the multi-PRS components from the biomarker traits into the model does improve accuracy substantially, which we confirm in FinnGen (C-index with covariates 0.665; hazard ratio 0.99, p = 0.46, 95% CI 0.95-1.02 for just the chronic kidney disease PRS and C-index with covariates 0.667, hazard ratio 1.12, p = 2.1e-10, 95% CI 1.08-1.16 also including the 35 biomarkers). While existing scores for chronic kidney disease are not available, type 2 diabetes is a major risk factor for chronic kidney disease. To illustrate that, we fit models which re-weight between the DIAMANTE type 2 diabetes PRS, the Lall 2016 type 2 diabetes PRS, and the training set UKBB PRS. The baseline models which do not include biomarkers and only include these three PRS does substantially better than just the chronic kidney disease score (AUC = 0.53; log odds = 0.101, p = 1.6e-7, 95%

CI 0.062 - 0.140), and adding the 35 biomarkers does even better (AUC 0.565; log odds = 0.231, p = 9.11e-33, 95% CI 0.192 - 0.269). However, adding the two diabetes PRS does not improve above the model including just the 35 biomarkers and the baseline chronic kidney disease PRS above. Taken together, this suggests that at least for chronic kidney disease, addition of biomarker scores does aid prediction above and beyond the inclusion of related disease traits (type 2 diabetes) or simply training a PRS on the same individuals for chronic kidney disease alone.

Following your comment, we added **Supplementary Figure 10**.



**Supplementary Figure 10. Extending multi-PRS with additional trait polygenic scores. (A) Myocardial infarction.** Odds ratios for angina and myocardial infarction using the multi-PRS including biomarkers (red) or multi-PRS of just the trait polygenic score and existing scores (GRS49K [PMID 27655226] and Khera et al. [PMID 30104762] SNPs; grey). Error bars represent 95% confidence intervals. **(B) Chronic kidney disease and diabetes.** Odds ratios for diabetes and chronic kidney disease using the multi-PRS including biomarkers (red) or multi-PRS of just the trait polygenic score and existing scores (DIAMANTE [PMID 30297969] and Läll et al. 2016 [PMID 27513194]; grey). Error bars represent 95% confidence intervals. Type 2 diabetes (strict) refers to training a polygenic score with type 2 diabetes controls with HbA1c > 39 mmol/mol excluded.
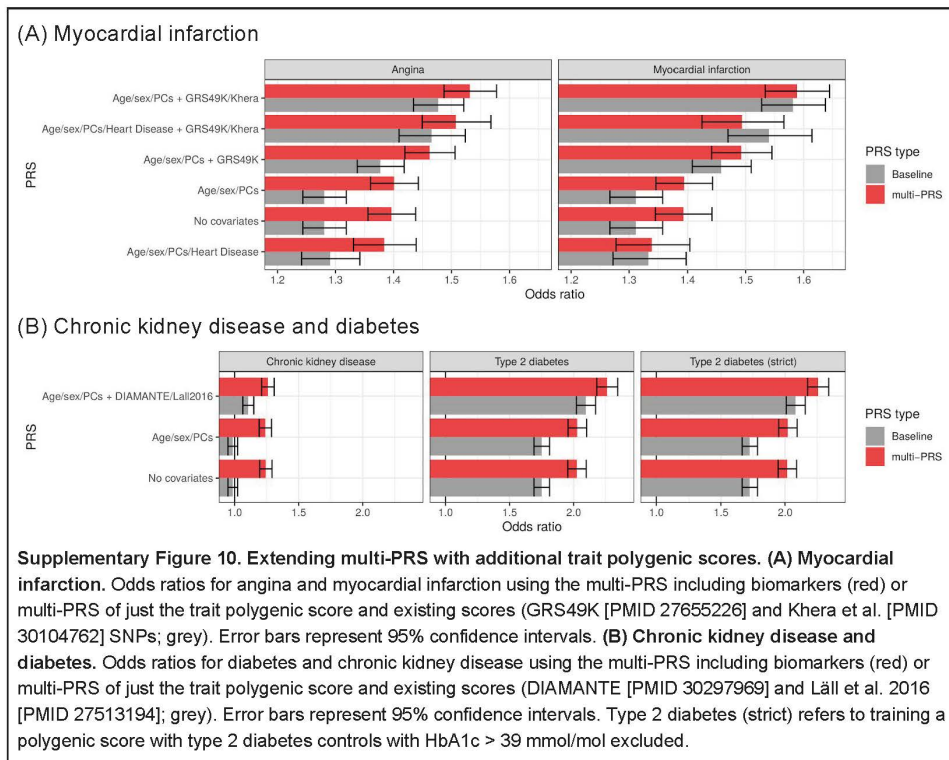
We hope that these plots, and the corresponding AUC values in the **Supplementary Table 20**, help clarify these findings. In short, adding additional external PRSs helps, but does not explain all of the improvement in all traits. For Joshi et al [PMID 27094239] in gallstones, indeed adding the Joshi et al SNPs as a score did not substantially help in

prediction over the baseline UKBB results, and multi-PRS models including the biomarkers did substantially better in both cases.

Overall, the fundamental question of whether adding PRS of biomarkers to PRS based on disease states meaningfully improves prediction remains unanswered in the revision. Many diseases such as T2D have previously validated PRS available, so uncertain why only CAD was analyzed. Moreover, HR close to 1.5/SD have been reported for many diseases, so the values in Fig 5D gray bars seem very low.

Thank you very much for the great question. While we are aware that there are validated PRSs derived from disease-focused study with large-number of cases, our initial analysis focused on the performance comparison of single-trait PRS and multi-PRS within UK Biobank (using the individuals in the hold-out test set), which is shown in **Figure 5D**. We think the relatively low performance (compared to the best reported performance) can be explained by the difference in power resulting from the limited case counts in the population-based cohort.

Following your suggestion, we now included the previously validated PRSs in our extended comparison analysis, which is presented in **Supplementary Figure 10**. In these figures, for a fair comparison across different approaches, we used the same set of individuals from UK Biobank (our test set in the PRS analysis) for evaluation and quantified the odds ratios in prevalent cases in UK Biobank for angina, MI, chronic kidney disease, and diabetes:

(A) Myocardial infarction

(B) Chronic kidney disease and diabetes

**Supplementary Figure 10. Extending multi-PRS with additional trait polygenic scores. (A) Myocardial infarction.** Odds ratios for angina and myocardial infarction using the multi-PRS including biomarkers (red) or multi-PRS of just the trait polygenic score and existing scores (GRS49K [PMID 27655226] and Khera et al. [PMID 30104762] SNPs; grey). Error bars represent 95% confidence intervals. **(B) Chronic kidney disease and diabetes.** Odds ratios for diabetes and chronic kidney disease using the multi-PRS including biomarkers (red) or multi-PRS of just the trait polygenic score and existing scores (DIAMANTE [PMID 30297969] and Läll et al. 2016 [PMID 27513194]; grey). Error bars represent 95% confidence intervals. Type 2 diabetes (strict) refers to training a polygenic score with type 2 diabetes controls with HbA1c > 39 mmol/mol excluded.

In particular, we note that this result is not surprising, since the best PRS should (to a first approximation) come from disease-specific efforts with high case counts (such as these existing PRS based on GWAS meta-analysis). However, we believe these results suggest the inclusion of risk factors for a given disease in the models can improve prediction accuracy, even when including GWAS meta-analyses. Thank you for this feedback and we hope this addition helps to clarify these models.

Similarly in ST22, HR of 1.16 for a MI PRS baseline model seems very low.

Thank you very much for your noting this. We think that this is because we are using a relatively small number of case individuals in UK Biobank. Specifically, we used only 8,459 MI cases in our training set (compared to 63,746 in CARDIoGRAMplusC4D), to fit the PRS models. This enabled us to have a fair comparison of different PRS models using the individuals with outcomes and the corresponding biomarker scores. As noted

above, we have expanded upon the results to include existing scores and see substantially increased HRs as a result.

Does multiPRS approach improve or worsen transportability across racial groups?

This is an interesting question. We have added **Supplementary Table 23**, showing the performance of the polygenic scores per standard deviation across populations in UK Biobank. In short, the performance is improved but the portability (ratio of log odds in White British versus comparison population) is not substantially different. A more quantitative comparison of relative portability, such as including ancestry proportion effects and matching of individuals between populations, is outside the scope of this manuscript, but we thank the reviewer for suggesting this addition which gives a preliminary understanding of this finding.

Can authors give generalizable learnings about when incorporating biomarker PRS might be helpful? For example, it seems like when the disease GWAS is very large, it may be less useful, but more useful when biomarkers are strongly associated with disease, discovery GWAS smaller, etc.

Thank you for bringing up this important point. We hope that the response above about the weights being spread across multiple PRS is useful in answering your question to some extent. We have been training additional multi-PRS models (e.g. including existing scores for T2D) and have included additional anecdotal evidence that could be useful for other researchers this effect in the discussion:

Page 17 line 479:

> In addition to the discovery of multiple individual loci and candidate causal variants, we can also draw some general conclusions across the traits evaluated with our multi-PRS models. Traits and diseases were predicted best when they had individually predictive biomarkers and a complex aetiology (e.g. chronic kidney disease), but underpowered genetic studies. We believe that, barring study-level heterogeneity in case definitions (e.g. early-onset versus late-onset disease), a large number of disease cases is typically most useful in developing well-powered models, as it helps both with the baseline polygenic score and fitting of the multi-PRS components. Further exploration of the conditions where multi-PRS models perform particularly well is an area of future study.

Thank you for this suggestion and we hope that this is useful for the field going forward.

I am not able to individually vet each of the proposed causal inferences, but many of them do not seem particularly biologically plausible. e.g. lowering TG would decrease risk of anorexia?

We agree that these causal inference claims are approximated using genetic instruments and thus only capture part of the possible mechanisms at play. In particular, pleiotropic effects are problematic for Mendelian Randomization based approaches, though our methodology includes a number of filters and decisions that are meant to minimize the number of spurious associations discovered.

Indeed, there is substantial existing literature on the connection between abnormal lipid metabolism and anorexia. Anorexic individuals tend to have higher fasting lipid levels, including triglycerides (PMIDs 12920792, 16791856, 23114953). We further note that a previous genetic analysis identified variants in *EPHX2* consistent with our observed effect (PMID 23999524). Thus, while we do not have conclusive evidence, this adds to a growing body of work with samples of matched anorexic and non-anorexic individuals with lipid panels over time, and could help with plans to develop interventional approaches we could help clarify the causal role of cholesterol metabolism in the development and progression of anorexia nervosa. We added a small section to this effect, shown below (page 12 line 363):

> *We also find the first Mendelian Randomization evidence of a causal association between triglycerides and anorexia nervosa (1.67 OR/SD, FDR-adjusted p = 0.042, supporting both previous rare variant studies and epidemiological literature.*

The list of phenotypes used in PheWAS seems arbitrary/duplicative/not reflective of important conditions. For example, authors describe 'removed disease outcomes that were likely to be duplicated' but then have separate rows fro 'DVT diagnosed by doctor' and 'blood clot or DVT diagnosed by doctor.' Similarly with 'bipolar and major depression status any' and 'bipolar and major depression.' For the 'status any' version, what does that mean?

I wonder if merging diagnosis codes into concepts PheCODES as is available using codings described by the Vanderbilt group, or just using an all by all and not suggesting there was curation would be reasonable alternatives, but best would really be to more effectively curate into disease groupings. Is 'dentures' really an important phenotype?

Thank you very much for raising an important question regarding the set of disease endpoints used in our single-variant and PRS PheWAS analyses.
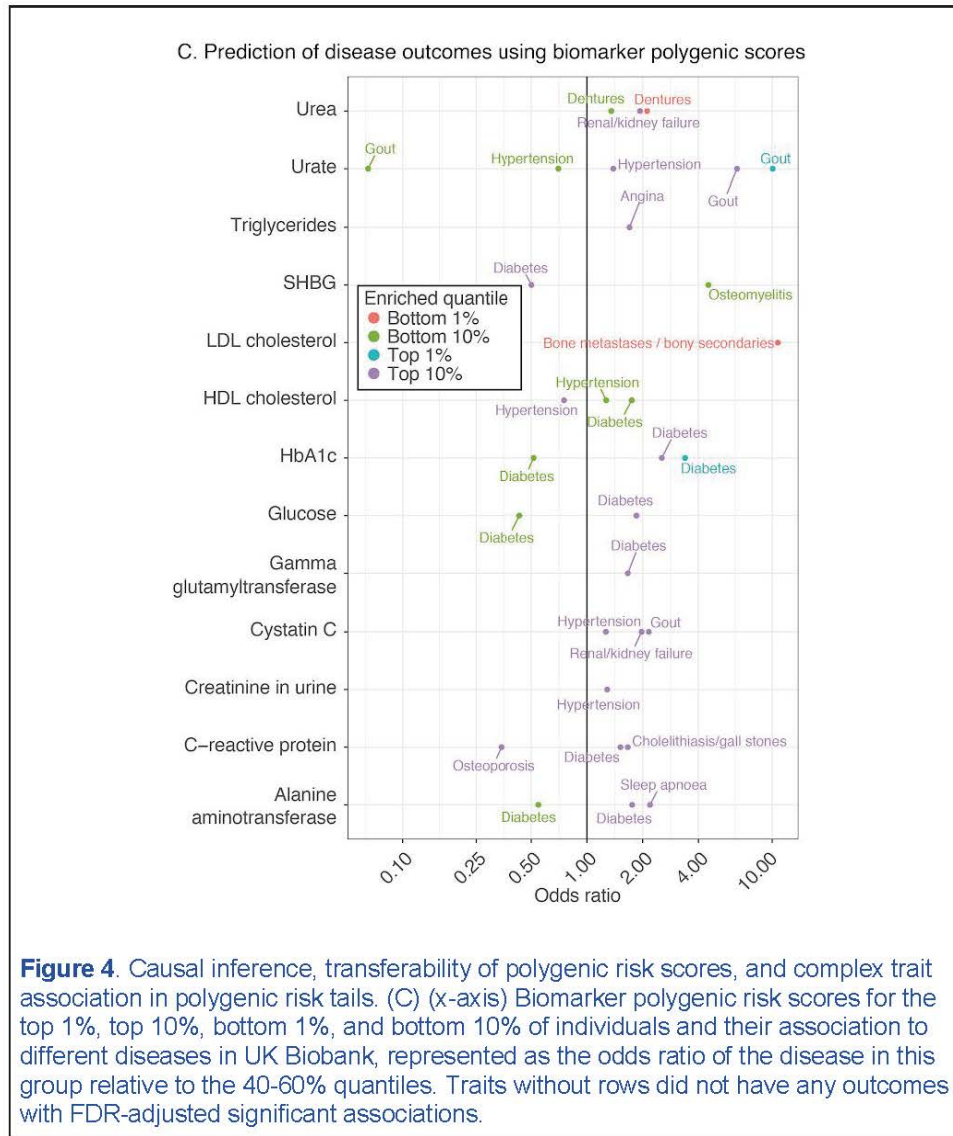
First of all, we do not claim the importance of the selected phenotypes -- our selection criteria of the disease endpoint was due to the case counts in white British unrelated individuals in UK Biobank. With this approach, we aim to perform a comprehensive PheWAS analysis rather than focusing on a predefined set of important and/or relevant phenotypes.

Thank you for suggesting the use of an external ontology and/or concepts, like phecodes. Our disease endpoint definition is based on a combination of ICD codes in hospital inpatient records and self-reported disease ascertainment status. As we showed in our previous publication (PMID: 32275883), our digital phenotyping strategy with combined sources often improves the number of cases and increases the power of association analysis. This made it difficult to map our disease endpoints to external ontology/concepts, like phecodes.

Having said that, we do agree with your concern about duplicated entries in our list of "curated" phenotypes. We sincerely apologize our previous curation effort was incomplete. We have performed one more round of curation and updated the list to 166 phenotypes.

Thank you for commenting on specific phenotypes. We have considered and incorporated your advice when performing the phenotypic curation.

- Deep vein thrombosis (DVT) phenotypes: we agree that there were duplicates. We kept "blood clot or DVT diagnosed by doctor" and removed "DVT diagnosed by doctor"
- Bipolar and major depression phenotypes: there were also duplicates and that has been fixed in the revised supplementary table (**Supplementary Table 12**).
- Dentures - we do not claim the importance of the selected phenotypes. To maximize the discovery from population-based cohorts, we have included all the available phenotypes based on the case counts. Our PRS-PheWAS results, indeed, identified dentures as one of the enriched phenotypes for Urea PRS. We confirmed there are no other duplicated phenotypes for dentures. We have therefore decided to keep this one.

**Figure 4**. Causal inference, transferability of polygenic risk scores, and complex trait association in polygenic risk tails. (C) (x-axis) Biomarker polygenic risk scores for the top 1%, top 10%, bottom 1%, and bottom 10% of individuals and their association to different diseases in UK Biobank, represented as the odds ratio of the disease in this group relative to the 40-60% quantiles. Traits without rows did not have any outcomes with FDR-adjusted significant associations.

Also, as you correctly pointed out, the labeling of "Bipolar and major depression: status any " was a poor phenotype label and we have replaced it with "Bipolar and major depression status: any of the following: bipolar I disorder, bipolar II disorder, probable

single or recurrent major depression episode(s)". This specific phenotype was defined from UKB Data Field 20126 (http://biobank.ctsu.ox.ac.uk/crystal/field.cgi?id=20126). We included individuals in case if the individual answered this questionnaire anything but "No Bipolar or Depression". We have clarified the source Data Field in UK Biobank in the updated **Supplementary Table 12**.

As you can see above, our disease endpoint definitions were not clearly communicated in the earlier version of the manuscript. We have updated our list of PheWAS disease endpoint phenotypes and added the "Source of the phenotype" column (**Supplementary Table 12**). In case the phenotype was defined using a combination of ICD codes in hospital inpatient records and self-reported disease status, we have provided a link to the phenotyping data source breakdown page in Global Biobank Engine.

We have updated the corresponding sections in the manuscript (the PheWAS analysis for the coding variant associations, the PheWAS analysis for the fine-mapped variant associations, and the PRS-PheWAS analysis), figure 4C, and Supplementary Tables 12, 13A-C, and 19.

Results section, Page 7, line 152:

*To further detect whether the variants we found associations with biomarkers impact other diseases or clinically relevant phenotypes, we performed a phenome-wide association analysis (PheWAS) across 166 traits in UK Biobank, compared with literature, and sought replication in FinnGen R2 (Supplementary Tables 12-13, Methods). We found a total of 57 associations (33 and 24 for increasing and decreasing the disease risks, respectively) across 26 disease outcomes for 2 PTVs and 31 PAVs (p < 1 × 10^{-7}), of which 31 associations were previously reported and 26 were novel (Supplementary Tables 13A, Methods).*

Results section, Page 10, line 306:

*We conducted PheWAS of the fine-mapped imputed variants across 166 UK Biobank phenotypes and identified 14 and 263 coding and non-coding associations, of which 109 were not previously reported in literature (p < 10^{-7}, Supplementary Tables 12, 13B-C, Methods).*

Methods section, Page 26, line 810:

*We curated a list of 166 disease outcomes from previously-reported binary phenotypes in Global Biobank Engine (GBE)[27,77,96]. Specifically, we selected disease outcomes with at least 700 cases in white British grouping and removed disease outcomes that were likely to be duplicated (**Supplementary Table 12**). Those phenotypes include non-cancer disease-outcome endpoints derived from a combination of the ICD codes from hospital inpatient records as well as self-reported disease ascertainment status[96], family history phenotype (UK Biobank data category 100034), cancer phenotypes derived from a combination of the UK cancer registry data and questionnaire data[77], and additional set of phenotypes derived from the following data fields in UK Biobank: 2247, 2463, 2834, 3591, 6148, 6149, 6152, 6153, 20126, 20406, 20483, and 21068. The data source of the disease endpoint definitions for each phenotype is described in "Source of the phenotype" column in **Supplementary Table 12**.*

Thank you very much for clarifying the definitions and pointing out the duplicates.

Minor:

First sentence: I don't think it's true that biomarkers are the 'primary clinical tools' fo adverse health conditions

Thank you. We have updated the text.

Introduction, page 3, line 50:

*Serum and urine biomarkers are frequently measured to diagnose and monitor chronic disease conditions, and to assess treatment response.*

Authors describe that 90% of variance in testosterone was explained by covariates, but the majority of this is presumably just sex.

Thank you for this feedback. We agree that the majority of the variance explained in the case of testosterone is indeed by sex. We have added an additional supplemental table which details the variance explained by age, by age and sex, etc in evaluating the biomarkers of interest. From this we can easily read off that indeed sex explains 90.1% of the variance in Testosterone levels.

Semantics throughout paper are not well-reflective of terms used in clinical practice, suggest review by one of the clinician co-authors. Representative examples include ('kidney problems', 'renal failure', 'quantitative disease symptoms','t2d' in Figure)

Thank you very much for pointing this out. With help from clinical experts, we have revised the wording in the manuscript. For example, we now consistently use standard terminology, such as "kidney diseases" and "type 2 diabetes" throughout the text and the figures. Thank you very much.

Semantics such as 'McCarthy Lab Tools' should be removed

Thank you; we have removed this comment and replaced with the following and a citation to these tools (page 28 line 906):

HRC filtering and pre-checking SNPs was applied before imputation[111].

111. Robertson, N., Rayner, N. W. & McCarthy, M. I. ScatterShot web application for UK Biobank, McCarthy Group Tools. https://www.well.ox.ac.uk/~wrayner/tools/

LOR/SD difficult to interpret — any reason to not report odds ratio per SD?

Thank you for this feedback -- we have replaced the log odds with odds ratios in the causal inference section.

Lipoprotein A is referred to as lipoprotein(a) within clinical research and clinical trial paradigms

Thank you for this feedback. We followed documentation in UK Biobank (https://www.ukbiobank.ac.uk/wp-content/uploads/2018/11/BCM023_ukb_biomarker_panel_website_v1.0-Aug-2015-edit-2018.pdf) and used "Lipoprotein A" in our earlier versions of the manuscript.

We appreciate your point and updated the text so that the manuscript is more accessible to clinically focused audiences.

Results section, page 6, line 128:

We found that both LD Score regression and HESS indicate common SNPs explain substantial heritability of some but not all biomarkers (0.6% [Lipoprotein A, also referred

*to as lipoprotein(a)] to 23.9% [IGF-1] using LD Score regression and 3.2%*
*[Microalbumin in urine] to 57% [Total bilirubin] using HESS across the studied*
*continuous phenotypes, **Supplementary Tables 11A-B**).*

Semantics lack precision and warrant proofreading. A representative example is Supp Figure 7A, entitle 'Portability of individual biomarkers,' but I think this may refer to the relative variance explained of polygenic scores for these biomarkers?

Thank you for this comment and we apologize for the ambiguity. We have updated the title of this figure to, "Variance explained for polygenic scores of individual biomarkers across populations." We appreciate the attention to these important details.

Referee #2:

1. What is novel among the many discoveries? I see that some of this information is listed in the Supplementary Tables, but it would be helpful for the reader to also appreciate novelty in the main text. In particular, known/novel information for the different variants and genes listed in the section entitled "Biomarkers associated variants prioritize therapeutic targets" would be very informative. This had already been raised in the previous review.

Thank you very much for encouraging us to emphasize the novelty of the identified associations. Following your suggestion, we have updated the manuscript to improve the clarity of the identified novel association.

Specifically, we present the results from the biomarker GWAS discovery as well as the PheWAS associations to 166 disease endpoints in "Biomarkers associated variants prioritize therapeutic targets" section.

For GWAS, we compared the identified associations with prior studies (listed in **Supplementary Table 5**) and found that **55/58** protein-truncating variant associations and **1,079/1,323** protein-altering variant associated were not reported in the comparison studies.

This is clarified in the main text:

Results section, page 6, line 141:

*We found 58 (43 rare, minor allele frequency [MAF] < 1%, and 55 not reported in comparison study, Methods) PTV and 1,323 (306 rare, 1,079 not reported in comparison studies) protein-altering variant associations outside the major histocompatibility complex (MHC) region (hg19 chr6:25,477,797-36,448,354; meta-analyzed $p < 5 \times 10^{-9}$).*

Methods section, page 24, line 767:

*Using the same set of comparison studies, we also checked whether the protein-truncating and protein-altering associations were previously reported for a given trait by calling the association reported if the p-value of the variant is less than $1 \times 10^{-6}$ in any comparison study for a given trait.*

Furthermore, following your suggestion, we have updated Figure 2 and added '*' symbol to indicate the novel discovery. We hope the readers will have a better understanding on which protein-truncating or protein-altering associations presented in the manuscript were novel. Thank you very much for your suggestion.

**Figure 2.** Genetics of 35 biomarkers. (main panel) Fuji plot of lab phenotypes across the six categories provided by UK Biobank and genetic variant associations shown for LD independent variants with meta-analysis P < 5x10⁻⁹. Large-effect protein-truncating and protein-altering variants (labeled when abs(Beta) >= 0.1 standard deviation [SD]) annotated with the category of association displayed (colored fill boxes) and highlighted if the loci was not previously reported in the comparison studies (**Methods**). Pleiotropic association and trait-specific association are shown by different sized circles.

For the PheWAS associations, we report in a total of 61 associations. To report the novelty of those findings, we queried the NHGRI-EBI GWAS catalog and Open Target Genetics and found that 26 associations were novel. This novelty is reported in column U ("Is_Novel?" column) in **Supplementary Table 13A**. Here is the corresponding text in the "Biomarkers associated variants prioritize therapeutic targets" section mentioning the novelty of PheWAS associations.

Results section, Page 7, line 156:

*We found a total of 57 associations (33 and 24 for increasing and decreasing the disease risks, respectively) across 26 disease outcomes for 2 PTVs and 31 PAVs ($p < 1 \times 10^{-7}$), of which 31 associations were previously reported and 26 were novel (**Supplementary Tables 13A**, **Methods**).*

2. What is the FINEMAP posterior inclusion probability (PIP) for the potential "functional" variants highlighted in section "Biomarkers associated variants prioritize therapeutic targets"?

In our non-synonymous variant association analysis presented in "Biomarkers associated variants prioritize therapeutic targets" section, we employed a marginal effects analysis on the directly-genotyped variants. However, you raise a good point that this could be integrated with the following fine-mapped section on imputed variants. The coding variants evaluated for the prioritization are mostly rare variants which do not overlap with the FINEMAP analysis, and as such, we felt they were better treated as distinct sections. We do believe that this analysis could prove useful in future iterations of this work as fine-mapping techniques continue to improve in rare variants, through the use of e.g. whole-genome sequencing that UK Biobank plans to undergo on all participants, and thank the reviewer for the suggestion.

3. For the FINEMAP analyses, how were "loci" defined? Based on physical distance around independently associated variants? Did the authors merge overlapping loci? More info needed in the Methods section.

Thank you for this feedback. We have added additional details to the methods section (page 25, line 772):

*Independent loci were defined by clumping White British GWAS summary statistics (see section "Derivation of independent loci"). For each putative SNP, we defined*

*distance-independent regions by collating all variants in linkage disequilibrium with the following plink command:*

```
plink1.9 --clump-p1 1e-3 --clump-p2 1e-3 --clump-r2 0.0001
--clump-kb 10000 --clump-field P-value --clump-snp-field
MarkerName
```

*In this way, we defined the individual loci contributing to the fine-mapping. We identified putative causal SNPs in each locus by using the FINEMAP software.*

The loci were defined based on linkage to the lead SNP, and lead SNPs were defined based on a centiMorgan cutoff to provide ample "buffer" of independence to ensure that they were minimally overlapping. We note that even non-overlapping loci have the potential to share signals due to linkage disequilibrium, and further analysis of fine-mapping results should consider this and other methods, such as conditional fine mapping, to aid in interpretation.

4. Figure 3.
Panel B. Consider adding % variance explained by fine-mapped variants for each of the horizontal bar.

We added a new table, **Supplementary Table 14B**, which contains the variance explained by fine-mapped variants, which ranges from 0% for Potassium in urine to 48% for lipoprotein(a). We added a corresponding line to the main text (page 10 line 272):

*These explain between 0% (Potassium in urine) and 48% (Lipoprotein A) of the residual trait variance (**Supplementary Table 14B**).*
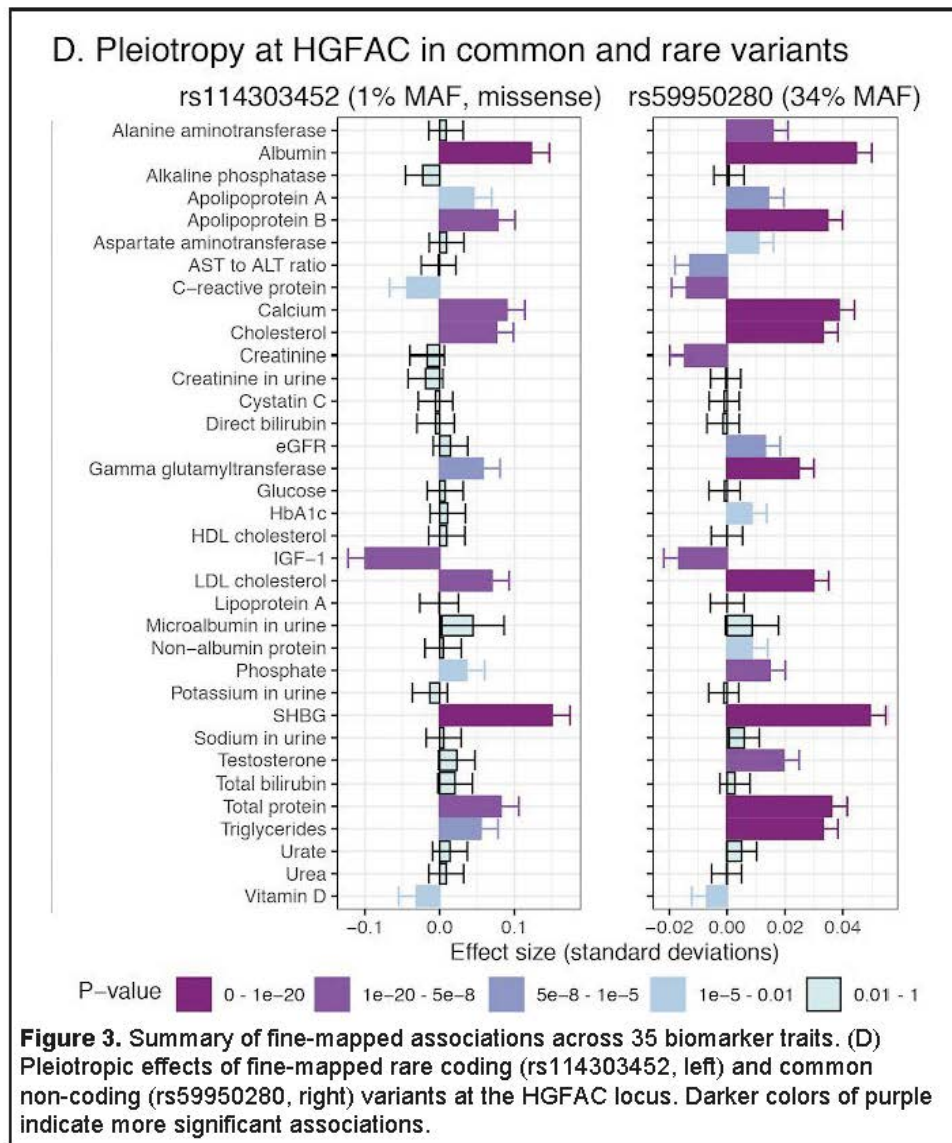
As well as to the figure:

**B. Number of fine-mapped variants per trait**

Figure 3. Summary of fine-mapped associations across 35 biomarker traits. (B) Breakdown of the number of fine-mapped associations with posterior probability greater than 0.95 or 0.99 across all biomarkers. Orange, posterior greater than 0.99, green, posterior between 0.95 and 0.99. Total variance explained for each trait is shown and in **Supplementary Table 14B**.

Thank you for this suggestion.

Panel D. I cannot see the palest color.

Thank you very much for your feedback. We have updated the color to increase the readability:



**Figure 3.** Summary of fine-mapped associations across 35 biomarker traits. (D) Pleiotropic effects of fine-mapped rare coding (rs114303452, left) and common non-coding (rs59950280, right) variants at the HGFAC locus. Darker colors of purple indicate more significant associations.

5. In STables 14 A-B-C, what is the meaning of the color code for the different cells?

Thank you for catching this error and clarifying this. The colors in those tables do not represent meaningful information and should have been removed in the previous revision. We have removed those colors in the updated supplementary tables (which are now **Supplementary Tables 13A-C**).

Referee #3:

I am comfortable with the changes, with one slight exception.

Thank you very much.

I think the HNF1B association with renal failure etc. should be explored relative to diabetes. The response to reviewer indicated the association was substantially attenuated (p of 5x10-7 to 0.035) when T2D cases were removed. I think an additional sentence in the manuscript suggesting that the renal failure may be happening primarily in diabetics would help researchers trying to follow-up this finding.

Following your suggestion, we have added an additional **Supplementary Table 10B** regarding the *HNF1B* association. Indeed, as you note, the odds ratio is substantially larger among diabetes cases, though we felt that the results were sufficiently nuanced that a more comprehensive table might be the best approach to allow readers to draw conclusions on interpretation of the specific patient populations evaluated. We have added this into the main text to reflect the diabetes status information (page 9 line 246):

*Consistent with its developmental role and clinical associations with MODY5, the rare CNVs overlapping HNF1B associate with chronic kidney disease in UK Biobank (p = 1 × 10$^{-7}$; OR = 4.94, SE = 0.30;* **Supplementary Figure 6A**$)^{34,35}$ *in a diabetes-dependent fashion (**Supplementary Table 10B**).*

Nice work.

Thank you very much for your time and feedback on our manuscript.

# nature research

| | |
|---|---|
| **Decision Letter, first revision:** | |

**Date:** 23rd Oct 20 20:27:03

**Last Sent:** 23rd Oct 20 20:27:03

**Triggered By:** Catherine Potenski

**From:** Catherine.Potenski@us.nature.com

**To:** mrivas@stanford.edu

**CC:** nasa@stanford.edu,ytanigaw@stanford.edu,davidama@stanford.edu,nina.mars@helsinki.fi,christian.benner@helsinki.fi,magu@stanford.edu,guhan@stanford.edu,wainberg@stanford.edu,hanna.m.ollila@helsinki.fi,tuomo.kiiskinen@helsinki.fi,aki.havulinna@fimm.fi,jamesp@broadinstitute.org,junyangq@stanford.edu,annashch@stanford.edu,finngen-scientificcommittee@helsinki.fi,frodrigu@stanford.edu,tassimes@stanford.edu,vineeta@stanford.edu,tibs@stanford.edu,hastie@stanford.edu,samuli.ripatti@helsinki.fi,pritch@stanford.edu,mark.daly@helsinki.fi

**Subject:** Decision on Nature Genetics submission NG-A55370R

**Message:** Our ref: NG-A55370R

23rd Oct 2020

Dear Dr. Rivas,

Thank you for submitting your revised manuscript "Genetics of 35 blood and urine biomarkers in the UK Biobank" (NG-A55370R). It has now been seen by original referee #1 and the comments are below. The reviewer finds that the paper has improved in revision, and therefore we will be happy in principle to publish it in Nature Genetics, pending minor revisions to comply with our editorial and formatting guidelines.

\*\* Note that we will send you a checklist detailing these editorial and formatting requirements in about a week. Please do not finalize your revisions or upload the final materials until you receive this additional information.\*\*

In recognition of the time and expertise our reviewers provide to Nature Genetics's editorial process, we would like to formally acknowledge their contribution to the external peer review of your manuscript entitled "Genetics of 35 blood and urine biomarkers in the UK Biobank". For those reviewers who give their assent, we will be publishing their names alongside the published article.

While we prepare these instructions, we encourage the Corresponding Author to begin to review and collect the following:

-- Confirmation from all authors that the manuscript correctly states their names, institutional affiliations, funding IDs, consortium membership and roles, author or collaborator status, and author contributions.

43

-- Declarations of any financial and non-financial competing interests from any author. For the sake of transparency and to help readers form their own judgment of potential bias, the Nature Research Journals require authors to declare any financial and non-financial competing interests in relation to the work described in the submitted manuscript. This declaration must be complete, including author initials, in the final manuscript text.

If you have any questions as you begin to prepare your submission please feel free to contact our Editorial offices at genetics@us.nature.com. We are happy to assist you.

Thank you again for your interest in Nature Genetics.

Congratulations on the paper!

All the best,

Catherine

---

Catherine Potenski, PhD
Chief Editor
Nature Genetics
1 NY Plaza, 47th Fl.
New York, NY 10004
catherine.potenski@us.nature.com
https://orcid.org/0000-0002-4843-7071


Reviewer #1 (Remarks to the Author):

Manuscript significantly improved!




<b>ORCID</b>

Nature Genetics is committed to improving transparency in authorship. As part of our efforts in this direction, we are now requesting that all authors identified as 'corresponding author' create and link their Open Researcher and Contributor Identifier (ORCID) with their account on the Manuscript Tracking System (MTS) prior to acceptance. ORCID helps the scientific community achieve unambiguous attribution of all scholarly contributions. For more information please visit http://www.springernature.com/orcid

For all corresponding authors listed on the manuscript, please follow the instructions in the link below to link your ORCID to your account on our MTS before submitting the final version of the manuscript. If you do not yet have an ORCID you will be able to create one in minutes.
https://www.springernature.com/gp/researchers/orcid/orcid-for-nature-research

# nature research

IMPORTANT: All authors identified as 'corresponding author' on the manuscript must follow these instructions. Non-corresponding authors do not have to link their ORCIDs but are encouraged to do so. Please note that it will not be possible to add/modify ORCIDs at proof. Thus, if they wish to have their ORCID added to the paper they must also follow the above procedure prior to acceptance.

To support ORCID's aims, we only allow a single ORCID identifier to be attached to one account. If you have any issues attaching an ORCID identifier to your MTS account, please contact the <a href="http://platformsupport.nature.com/">Platform Support Helpdesk</a>.

26th Oct 2020

Dear Dr. Rivas,

Thank you for your patience as we've prepared the guidelines for final submission of your Nature Genetics manuscript, "Genetics of 35 blood and urine biomarkers in the UK Biobank" (NG-A55370R). Please follow the instructions provided here and in the attached files.

When you upload your final materials, please:

A) Fill out and upload the attached \*\*\*Publishing Policy Worksheet For Authors\*\*\*, which contains information on how to comply with our legal guidelines for publication and links to the files that you will need to upload prior to final acceptance. You must initial the relevant portions of this checklist, sign it and return it with your final files. We will not be able to proceed further without these files:
a) Reporting Summary (required)
https://www.nature.com/documents/nr-reporting-summary.pdf
b) License to Publish (LTP required for Original Research) or Copyright Assignment (required if commissioned by Journal)
c) Color Fee Form (required for Original Research)
https://www.nature.com/documents/nr-rj-colour-figure-form.pdf
d) Competing Interests Statement (if applicable)
e) Author Approval List (if applicable)
f) Third Party Rights Table (if applicable, either Third Party Rights for Original Research or Third Party Rights if Commissioned by Journal)
g) Institutional Open Access Waiver (if applicable)
h) Inventory of Supplementary Information

B) Include a tracked-changes Word file of your revised article.

C) Include a point-by-point response to the points below:

General formatting:

1. Article: Our standard word limit is 4,000 words for the Introduction, Results and Discussion. Your current manuscript is 5,381 words. Please shorten this to 4,000 words, moving excess text to a Supplementary Note if necessary.
2. Please ensure that sections are in the following order within the same manuscript file: Title, Authors, Affiliations, Abstract, Introduction, Results (with subheadings), Discussion, Acknowledgements, Author Contributions, References for main text, Figure Legends for main text, Tables, Online Methods, Data Availability Statement, Code Availability statement (if applicable), Methods-only references.
3. Currently, you have 2 sets of equally contributing authors. Nature Research journals allow one set of up to six co-authors to be specified as having contributed equally to the work or having jointly supervised the work. Other equal contributions are best described in author contribution statements. Please ensure that there is only one set of "equally contributing" authors.
4. You need to include a numbered affiliation for the FinnGen consortium. Please see instructions below for consortium author formatting. In order for the FinnGen consortium to remain in the main author list, there need to be consortium members who can be considered as authors, with their names and affiliations listed at the end of the article.
5. Online Methods do not have a strict limit, but we suggest 3,000 words as a target. Your methods section is currently 4,992 words. Please retain ~3,000 words in the main article file and move the remaining text to a Supplementary Note. Please note that we cannot accommodate the bulleted formatting found in some methods sections in the online methods section. Please either reformat these sections or ensure that they are moved to the Supplementary Note PDF.
6. Your abstract must be fewer than 150 words and should not include citations.
7. Results sub-headings should be 59 characters or fewer including spaces.
8. Please remove the figures from the main article file and present the figure legends in the main article file after the references (not embedded within the text).
9. We ask that you make all summary statistics publicly available.
10. If it is stated that data are available upon reasonable request, this must be explained and justified. Please be more specific about what "other data" are available upon request and why this is so. We strongly encourage that all data be made available.
11. There is no defined limit for the number of references allowed in the main text. An additional 20 references can be included in the Online Methods. Only papers that have been published or accepted by a named publication or recognized preprint server should be in the numbered list. Published conference abstracts, numbered patents and research data sets that have been assigned a digital object identifier may be included in the reference list.
12. Unpublished meeting abstracts, personal communications and manuscripts under consideration (and not formally accepted) may be cited only internally within the text and should not be added to the reference list. Please provide the names of the first five authors of unpublished data. If you cite personal communications or unpublished data of any individuals who are not authors of your manuscript, you must supply copies of written (including email) permission from the primary investigator of each group cited.

13. All references must be cited in numerical order. Place Methods-only references after the Methods section and continue the numbering of the main reference list (i.e., do not start at 1). Ensure the reference list is up to date for the final submission.
14. Equations and symbols that will be set apart from the text must be in an editable format. Do not use embedded images for equations or symbols.
15. Genes must be clearly distinguished from gene products (e.g., "gene Abc encodes a protein kinase," not "gene Abc is a protein kinase"). For genes, provide database-approved official symbols (for human genes throughout the paper use http://varnomen.hgvs.org/). For the relevant species, use NCBI Gene: http://www.ncbi.nlm.nih.gov/gene. Italicize gene symbols and functionally defined locus symbols; do not use italics for proteins, noncoding gene products and spelled-out gene names.
16. For descriptions of variants, use HGVS notation according to the guidelines at http://varnomen.hgvs.org/. Include the accession code for the corresponding reference sequence at first mention of a variant.

Figures and Tables:

17. All figures and tables, including Extended Data, must be cited in the text in numerical order.
Please correct the following: main Figures and Supplementary Tables are cited out of order.
18. Figure legends should be concise and fewer than 250 words. Begin with a brief title and then describe what is presented in the figure and detail all relevant statistical information (as described and declared in the supplied checklist), avoiding inappropriate methodological detail.
19. Please remove the boxes from around the text in Fig.1.
20. Please ensure that Fig. 2 fits within the figure template and can accommodate minimum font size (4.5pt).
21. Background gridlines should be removed from the graphs.
22. Shadings or symbols in graphs must be defined in some fashion. We prefer that you use a key within the image; do not include colored symbols in the legend.
23. All relevant figures must have a definition for any error bars.
24. Red/green color contrasts can confuse our colorblind readers; please consider recoloring relevant figures, if possible.
25. Graph axes should start at zero and not be altered in scale to exaggerate effects. A 'broken' graph can be used if absolutely necessary due to sizing constraints, but the break must be visually evident and should not impinge on any data points.
26. All bar graphs should be converted to a dot-plot format or to a box-and-whisker format to show data distribution. All box-plot elements (center line, limits, whiskers, points) should be defined.
27. When submitting the revised version of your manuscript, please pay close attention to our href="https://www.nature.com/nature-research/editorial-policies/image-integrity">Digital Image Integrity Guidelines.</a> and to the following points below:
- That unprocessed scans are clearly labelled and match the gels and western blots presented in figures.
- That control panels for gels and western blots are appropriately described as loading on sample processing controls
- All images in the paper are checked for duplication of panels and for splicing of gel lanes.

Finally, please ensure that you retain unprocessed data and metadata files after publication, ideally archiving data in perpetuity, as these may be requested during the peer review and production process or after publication if any issues arise.

Statistics and Reproducibility:

a. The Methods must include a statistics section where you describe the statistical tests used. The supplied checklist provides details of which statistics need to be in the Figure legends, and which assumptions and analytical procedures need to be supplied in the Methods. For all statistics (including error bars), provide the EXACT numbers used to calculate the statistics (reporting individual values rather than a range if n varied among experiments) AND define type of replicates (e.g., cell cultures, technical replicates). Please avoid use of the ambiguous term "biological replicates"; instead state what constituted the replicates (e.g., cell cultures, independent experiments, etc.). For all representative results, indicate number of times experiments were repeated, number of images collected, etc. Indicate statistical tests used, whether the test was one- or two-tailed, exact values (NOT for example: <0.05) for both significant and non-significant P values where relevant, F values and degrees of freedom for all ANOVAs, and t-values and degrees of freedom for t-tests.

b. <b>Reporting Guidelines</b>– Attached you will find an annotated version of the Reporting Summary you submitted, along with a Word document indicating revisions that need to be made in compliance with our reproducibility requirements. These documents detail any changes that will need to be made to the text, and particularly the main and supplementary figure legends, including (but not limited to) details regarding sample sizes, replication, scale and error bars, and statistics. Please use these documents as a guide when preparing your revision and submit an updated Reporting Summary with your revised manuscript. The Reporting Summary will be published as supplementary material when your manuscript is published.

\*\*please note that in a few days we will send you detailed comments on your reproducibility checklist. You may have to modify some of the reporting in the manuscript at that time.\*\*

Supplementary Information:

All Supplementary Information must be submitted in accordance with the instructions in the attached Inventory of Supporting Information, and should fit into one of three categories:

1. EXTENDED DATA: Extended Data are an integral part of the paper and only data that directly contribute to the main message should be presented. These figures will be integrated into the full-text HTML version of your paper and will be appended to the online PDF. There is a limit of 10 Extended Data figures, and each must be referred to in the main text. Each Extended Data figure should be of the same quality as the main figures, and should be supplied at a size that will allow both the figure and legend to be presented on a single legal-sized page. Each figure should be submitted as an individual .jpg, .tif or .eps file with a maximum size of 10 MB each. All Extended Data figure legends must be provided in the attached Inventory of Accessory Information,

not in the figure files themselves.

You have 14 Supplementary Figures. These are currently in the combined Supp. Note PDF and some are multi-page. You can select up to 10 single-page figures to present as Extended Data Figures. The remainder can stay as Supp. Figures.

2. SUPPLEMENTARY INFORMATION: Supplementary Information is material that is essential background to the study but which is not practical to include in the printed version of the paper (for example, video files, large data sets and calculations). Each item must be referred to in the main manuscript and detailed in the attached Inventory of Accessory Information. Tables containing large data sets should be in Excel format, with the table number and title included within the body of the table. All textual information and any additional Supplementary Figures (which should be presented with the legends directly below each figure) should be provided as a single, combined PDF. Please note that we cannot accept resupplies of Supplementary Information after the paper has been formally accepted unless there has been a critical scientific error.

All Extended Data must be called you in your manuscript and cited as Extended Data 1, Extended Data 2, etc. Additional Supplementary Figures (if permitted) and other items are not required to be called out in your manuscript text, but should be numerically numbered, starting at one, as Supplementary Figure 1, not SI1, etc.

You have 26 Supplementary Tables in a single excel sheet, with the titles on each individual tab. Please remove the Supp. Table legends from the combined Supp. Note and include these within the excel tabs themselves.

Supplementary Data 1-4 are provided as figshare links; if these are not files associated with the manuscript, then they cannot be labeled "Supplementary Data". The figshare links can appear in the DAS, but the files should be renamed. The legends should also be removed or rewritten in the Supp. Note.

The Supp. Note contains the FinnGen consortium members. If you wish to include FinnGen in the main author list, some/all of these authors need to be at the end of the main article file (please see below for more details).

3. SOURCE DATA: We encourage you to provide source data for your figures whenever possible. Full-length, unprocessed gels and blots must be provided as source data for any relevant figures, and should be provided as individual PDF files for each figure containing all supporting blots and/or gels with the linked figure noted directly in the file. Statistics source data should be provided in Excel format, one file for each relevant figure, with the linked figure noted directly in the file. For imaging source data, we encourage deposition to a relevant repository, such as figshare (https://figshare.com/) or the Image Data Resource (https://idr.openmicroscopy.org).


TRANSPARENT PEER REVIEW

Nature Genetics offers a transparent peer review option for new original research manuscripts submitted from 20 May 2020. We encourage increased transparency in peer review by publishing the reviewer comments, author rebuttal letters and editorial

decision letters if the authors agree. Such peer review material is made available as a supplementary peer review file. <b>Please state in the cover letter 'I wish to participate in transparent peer review' if you want to opt in, or 'I do not wish to participate in transparent peer review' if you don't.</b> Failure to state your preference will result in delays in accepting your manuscript for publication.

Please note: we allow redactions to authors' rebuttal and reviewer comments in the interest of confidentiality. If you are concerned about the release of confidential data, please let us know specifically what information you would like to have removed. Please note that we cannot incorporate redactions for any other reasons. Reviewer names will be published in the peer review files if the reviewer signed the comments to authors, or if reviewers explicitly agree to release their name. For more information, please refer to our <a href="https://www.nature.com/documents/nr-transparent-peer-review.pdf" target="new">FAQ page</a>.

Authorship and Consortia

c. CONSORTIA: If a consortium has at least one member who qualifies as an author, the name of the consortium has to be listed in the author list. In addition, names and affiliations of all members of a consortium who qualify as an author must be listed in the paper either in the author list and in the consortium, or only in a consortium list at the end of the manuscript. You can find our authorship criteria here: https://www.nature.com/authors/policies/authorship.html. Names and affiliations of all consortium members (those who qualify as authors and those who do not) can be included in the Supplementary information (optional). When submitting your revised manuscript via the online submission system, the consortium name should be entered as an author, together with the contact details of a nominated consortium representative. See https://www.nature.com/documents/nr-consortia-formatting.pdf for further consortia formatting guidelines, which should be adhered to prior to acceptance.

d. PROTOCOL EXCHANGE: Nature Research journals <a href="https://www.nature.com/nature-research/editorial-policies/reporting-standards#protocols" target="new">encourage authors to share their step-by-step experimental protocols</a> on a protocol sharing platform of their choice. Nature Research's Protocol Exchange is a free-to-use and open resource for protocols; protocols deposited in Protocol Exchange are citable and can be linked from the published article. More details can found at <a href="https://www.nature.com/protocolexchange/about" target="new">www.nature.com/protocolexchange/about</a>.

Open Researcher and Contributor Identifier
Nature Genetics is committed to improving transparency in authorship. As part of our efforts in this direction, we are now requesting that all authors identified as 'corresponding author' on published papers create and link their Open Researcher and Contributor Identifier (ORCID) with their account on the Manuscript Tracking System (MTS), prior to acceptance. ORCID helps the scientific community achieve

unambiguous attribution of all scholarly contributions. For more information please visit <a href="http://www.springernature.com/orcid">www.springernature.com/orcid</a>.

Before resubmitting the final version of the manuscript, if you are listed as a corresponding author on the manuscript, please follow the steps below to link your account on our MTS with your ORCID. If you don't have an ORCID yet, you will be able to create one in minutes. If you are not listed as a corresponding author, please ensure that the corresponding author(s) comply.

1. From the home page of the <a href="https://mts-ng.nature.com/cgi-bin/main.plex">MTS</a>) click on '<b>Modify my Springer Nature account</b>' under '<b>General tasks</b>'.
2. In the '<b>Personal profile</b>' tab, click on '<b>ORCID Create/link an Open Researcher Contributor ID(ORCID)</b>'. This will re-direct you to the ORCID website.
3a. If you already have an ORCID account, enter your ORCID email and password and click on '<b>Authorize</b>' to link your ORCID with your account on the MTS.
3b. If you don't yet have an ORCID, you can easily create one by providing the required information and then click on '<b>Authorize</b>'. This will link your newly created ORCID with your account on the MTS.

<b>IMPORTANT:</b> All authors identified as 'corresponding authors' on the manuscript must follow these instructions. Non-corresponding authors do not have to link their ORCIDs, but please note that it will not be possible to add/modify ORCIDs at proof. Thus, if they wish to have their ORCID added to the paper, they must also follow the above procedure prior to acceptance.

To support ORCID's aims, we only allow a single ORCID identifier to be attached to one account. If you have any issues attaching an ORCID identifier to your Manuscript Tracking System account, please contact the at <a href="http://platformsupport.nature.com/">Platform Support Helpdesk</a>.


Please use the following link for uploading these materials:
[REDACTED]


If you have any further questions, please feel free to contact me- I'd be more than happy to help.

Thank you very much.

All the best,

Catherine

---

Catherine Potenski, PhD
Chief Editor
Nature Genetics

1 NY Plaza, 47th Fl.
New York, NY 10004
catherine.potenski@us.nature.com
https://orcid.org/0000-0002-4843-7071

Reviewer #1:
Remarks to the Author:
Manuscript significantly improved!

## Author Rebuttal, first revision:

## Final Decision Letter:

| | |
|---|---|
| **Date:** | 1st Dec 20 13:30:43 |
| **Last Sent:** | 1st Dec 20 13:30:43 |
| **Triggered By:** | Catherine Potenski |
| **From:** | Catherine.Potenski@us.nature.com |
| **To:** | mrivas@stanford.edu |
| **CC:** | rjsproduction@springernature.com |
| **Subject:** | Decision on NG-A55370R1 Rivas |
| **Message:** | In reply please quote: NG-A55370R1 Rivas |

1st Dec 2020

Dear Dr. Rivas,

I am delighted to say that your manuscript "Genetics of 35 blood and urine biomarkers in the UK Biobank" has been accepted for publication in an upcoming issue of Nature Genetics.

Prior to setting your manuscript, we may make minor changes to enhance the lucidity of the text and with reference to our house style. We therefore ask that you examine the proofs most carefully to ensure that we have not inadvertently altered the sense of your text in any way.

Once your manuscript is typeset you will receive a link to your electronic proof via email within 20 working days, with a request to make any corrections within 48 hours. If you have queries at any point during the production process then please contact the production team at rjsproduction@springernature.com. Once your paper has been scheduled for online publication, the Nature press office will be in touch to confirm the

details.

Your paper will be published online after we receive your corrections and will appear in print in the next available issue. You can find out your date of online publication by contacting the Nature Press Office (press@nature.com) after sending your e-proof corrections. Now is the time to inform your Public Relations or Press Office about your paper, as they might be interested in promoting its publication. This will allow them time to prepare an accurate and satisfactory press release. Include your manuscript tracking number (NG-A55370R1) and the name of the journal, which they will need when they contact our Press Office.

Before your paper is published online, we shall be distributing a press release to news organizations worldwide, which may very well include details of your work. We are happy for your institution or funding agency to prepare its own press release, but it must mention the embargo date and Nature Genetics. Our Press Office may contact you closer to the time of publication, but if you or your Press Office have any enquiries in the meantime, please contact press@nature.com.

Acceptance is conditional on the data in the manuscript not being published elsewhere, or announced in the print or electronic media, until the embargo/publication date. These restrictions are not intended to deter you from presenting your data at academic meetings and conferences, but any enquiries from the media about papers not yet scheduled for publication should be referred to us.

The Author's Accepted Manuscript (the accepted version of the manuscript as submitted by the author) may only be posted 6 months after the paper is published, consistent with our <a href="http://www.nature.com/authors/policies/license.html">self-archiving embargo</a>. Please note that the Author's Accepted Manuscript may not be released under a Creative Commons license. For Nature Research Terms of Reuse of archived manuscripts please see: <a href="http://www.nature.com/authors/policies/license.html#terms">http://www.nature.com/authors/policies/license.html#terms</a>
If you have posted a preprint on any preprint server, please ensure that the preprint details are updated with a publication reference, including the DOI and a URL to the published version of the article on the journal website.

To assist our authors in disseminating their research to the broader community, our SharedIt initiative provides you with a unique shareable link that will allow anyone (with or without a subscription) to read the published article. Recipients of the link with a subscription will also be able to download and print the PDF.

As soon as your article is published, you will receive an automated email with your shareable link."

You can now use a single sign-on for all your accounts, view the status of all your manuscript submissions and reviews, access usage statistics for your published articles and download a record of your refereeing activity for the Nature journals.

Please note that we encourage the authors to self-archive their manuscript (the accepted version before copy editing) in their institutional repository, and in their funders' archives, six months after publication. Nature Research recognizes the efforts of funding

bodies to increase access to the research they fund, and strongly encourages authors to participate in such efforts. For information about our editorial policy, including license agreement and author copyright, please visit www.nature.com/ng/about/ed_policies/index.html

An online order form for reprints of your paper is available at <a href="https://www.nature.com/reprints/author-reprints.html">https://www.nature.com/reprints/author-reprints.html</a>. Please let your coauthors and your institutions' public affairs office know that they are also welcome to order reprints by this method.

If you have not already done so, we invite you to upload the step-by-step protocols used in this manuscript to the Protocols Exchange, part of our on-line web resource, natureprotocols.com. If you complete the upload by the time you receive your manuscript proofs, we can insert links in your article that lead directly to the protocol details. Your protocol will be made freely available upon publication of your paper. By participating in natureprotocols.com, you are enabling researchers to more readily reproduce or adapt the methodology you use. Natureprotocols.com is fully searchable, providing your protocols and paper with increased utility and visibility. Please submit your protocol to http://www.nature.com/protocolexchange/. After entering your nature.com username and password you will need to enter your manuscript number (NG-A55370R1). Further information can be found at http://www.natureprotocols.com.


Congratulations on the paper!

All the best,

Catherine

--

Catherine Potenski, PhD
Chief Editor
Nature Genetics
1 NY Plaza, 47th Fl.
New York, NY 10004
catherine.potenski@us.nature.com
https://orcid.org/0000-0002-4843-7071

Click here if you would like to recommend Nature Genetics to your librarian
http://www.nature.com/subscriptions/recommend.html#forms


** Visit the Springer Nature Editorial and Publishing website at <a href="http://editorial-jobs.springernature.com?utm_source=ejP_NGen_email&utm_medium=ejP_NGen_email&utm_campaign=ejp_NGen">www.springernature.com/editorial-and-publishing-jobs</a> for more information about our career opportunities. If you have any questions please click <a href="mailto:editorial.publishing.jobs@springernature.com">here</a>.**