# Calibrated prediction intervals for polygenic scores across diverse contexts

In the format provided by the authors and unedited

# Supplementary Information
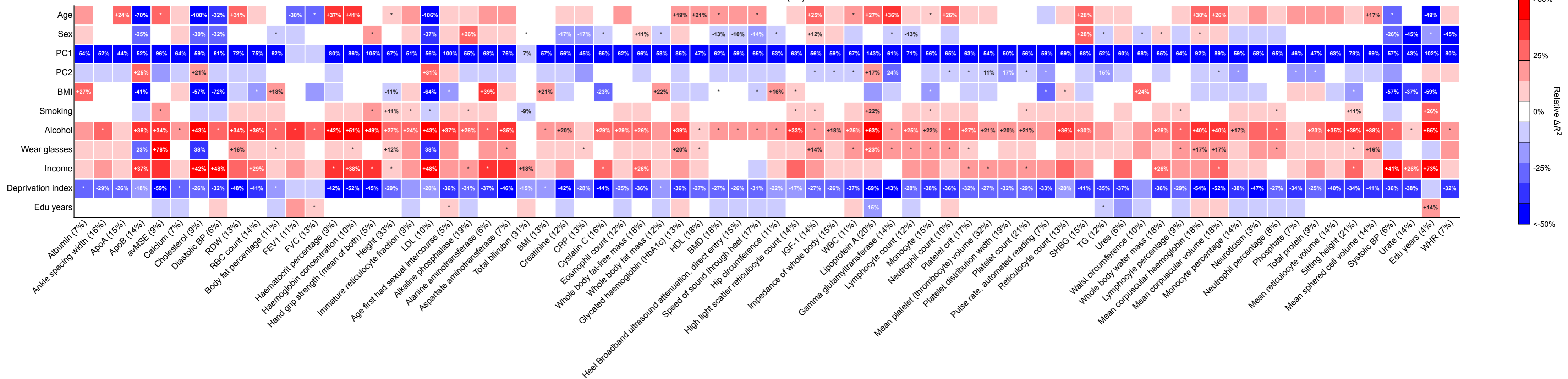
Calibrated prediction intervals for polygenic scores across diverse contexts
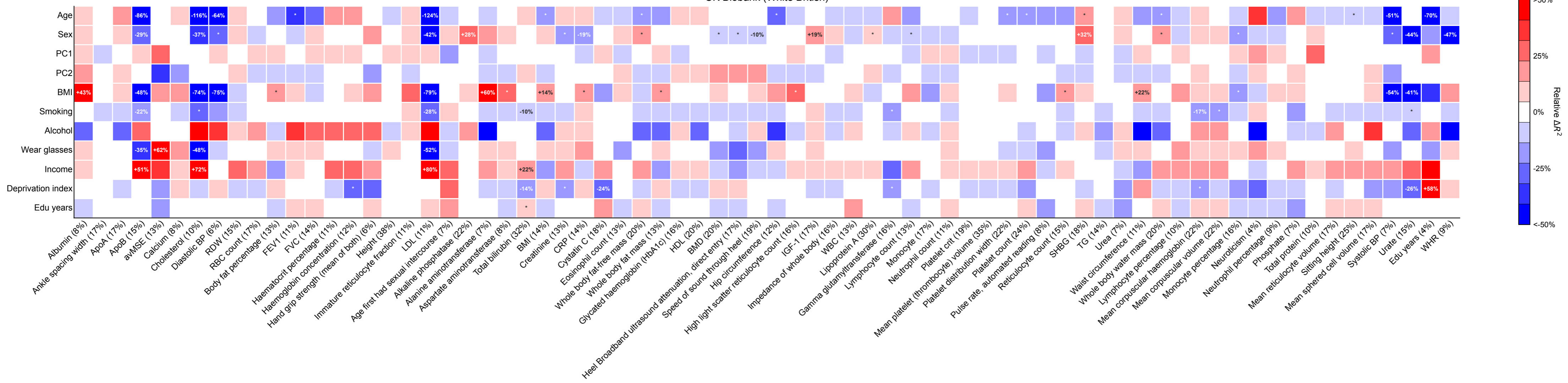
## Supplementary Figures

(see next page)

**Supplementary Figure 1:  Relative $\Delta R^2$ for "white British" and all individuals in UK Biobank.**  Numerical values of relative $\Delta R^2$ are displayed for trait-context pairs with statistically significant differences (multiple testing correction for all 72×11 trait-context pairs in this figure; two-sided $p < 0.05 / (72 \times 11)$). '*' are displayed for context-trait pairs with nominally significant differences (multiple testing correction for 11 contexts; two-sided $p < 0.05 / 11$). See Fig. 2 caption for additional details. Numerical results are reported in Supplementary Table 2.

UK Biobank (All)
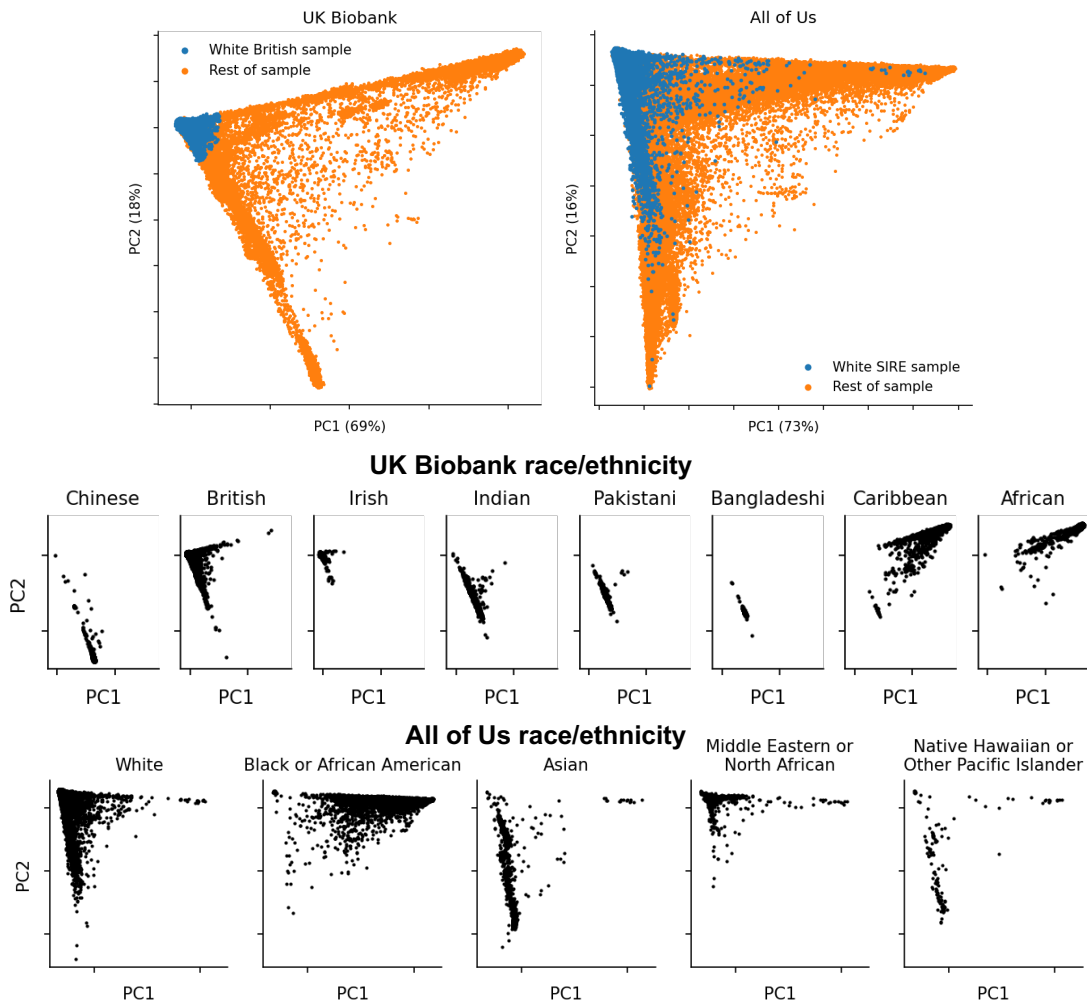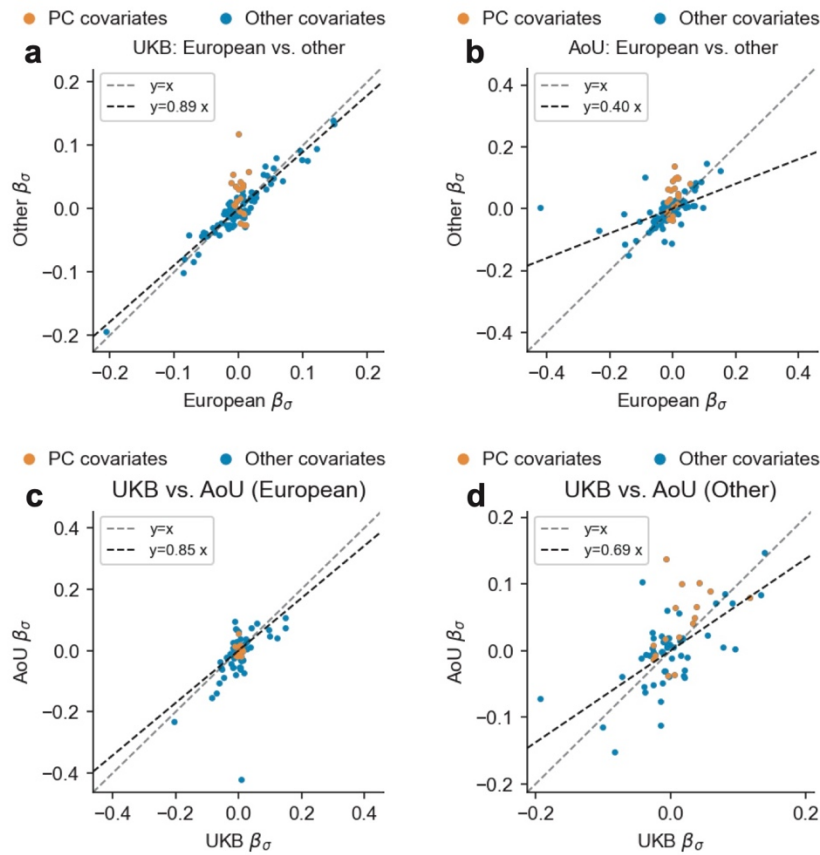
UK Biobank (White British)

**Supplementary Figure 2: Estimated β$_\sigma$ for "white British" and all individuals in UK Biobank.** We show estimated $\beta_\sigma$ in CalPred model. Numerical results are reported in Supplementary Table 2.

**Supplementary Figure 3: Principal components calculated in UK Biobank and All of Us.** PC1 and PC2 were calculated across all individuals separately in UK Biobank and All of Us. We show self-reported race/ethnicity in UK Biobank and All of Us to help interpret PC1/PC2. We show the proportion of variance explained by PC1 and PC2 out of top ten PCs calculated in each dataset in x-axis and y-axis labels.

**Supplementary Figure 4: Comparison of fitted parameters across populations and biobanks.** We compare estimated $\beta_\sigma$ across populations **(a-b)** and biobanks **(c-d)**. Each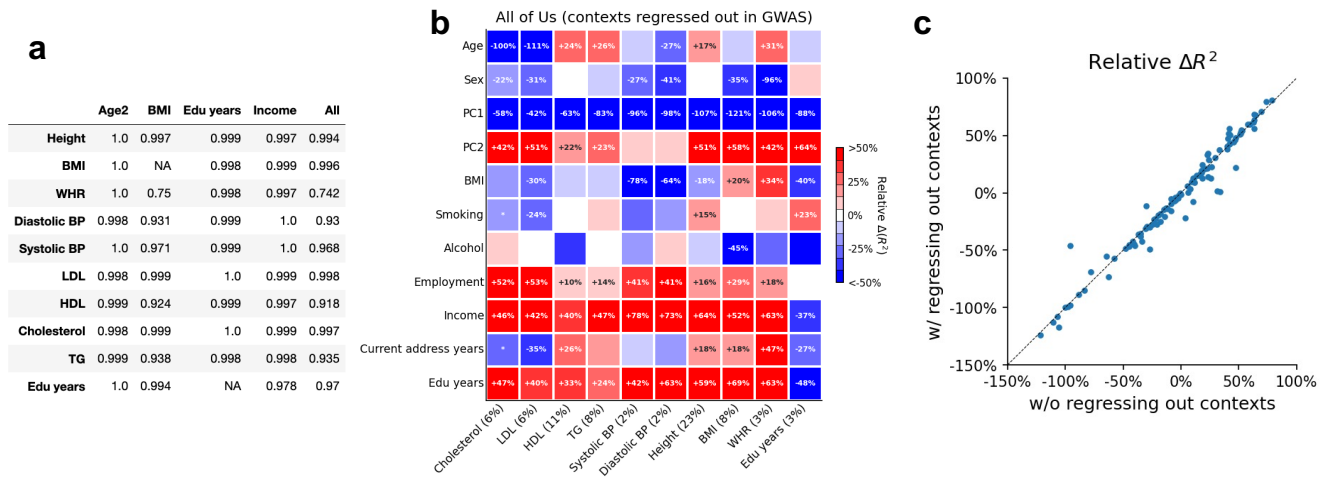 dot denotes a trait-context pair. We separately annotate genetic principal component contexts and other contexts, because PC contexts only have small variations within Europeans ("white British" in UK Biobank or "white SIRE" in All of Us), therefore are not comparable between European and other populations. The regression slope is calculated across all estimated $\beta_\sigma$. For **(c-d)**, we include traits that are shared across biobanks. Overall, we find that $\beta_\sigma$ are highly consistent across populations and biobanks.

**Supplementary Figure 5: $R^2$ and phenotypic variance across context groups for example traits.** We calculate $R^2$ between PGS and covariate-adjusted phenotypes as well as phenotypic variance by PC1 quintile, age quintile and sex for LDL, height and BMI across all individuals in UK Biobank. We determined that variable $R^2$ across contexts were not solely driven by differences of phenotype variance in context strata. The relationship between $R^2$ and phenotypic variance depends on the specific trait-context being studied. For example, height $R^2$ varies across age quintiles while the phenotypic variance remains relatively constant; BMI $R^2$ stays relatively constant across age quintiles while the phenotypic variance varies.

**a**

|  | Age2 | BMI | Edu years | Income | All |
|---|---|---|---|---|---|
| **Height** | 1.0 | 0.997 | 0.999 | 0.997 | 0.994 |
| **BMI** | 1.0 | NA | 0.998 | 0.999 | 0.996 |
| **WHR** | 1.0 | 0.75 | 0.998 | 0.997 | 0.742 |
| **Diastolic BP** | 0.998 | 0.931 | 0.999 | 1.0 | 0.93 |
| **Systolic BP** | 1.0 | 0.971 | 0.999 | 1.0 | 0.968 |
| **LDL** | 0.998 | 0.999 | 1.0 | 0.999 | 0.998 |
| **HDL** | 0.999 | 0.924 | 0.999 | 0.997 | 0.918 |
| **Cholesterol** | 0.998 | 0.999 | 1.0 | 0.999 | 0.997 |
| **TG** | 0.999 | 0.938 | 0.998 | 0.998 | 0.935 |
| **Edu years** | 1.0 | 0.994 | NA | 0.978 | 0.97 |

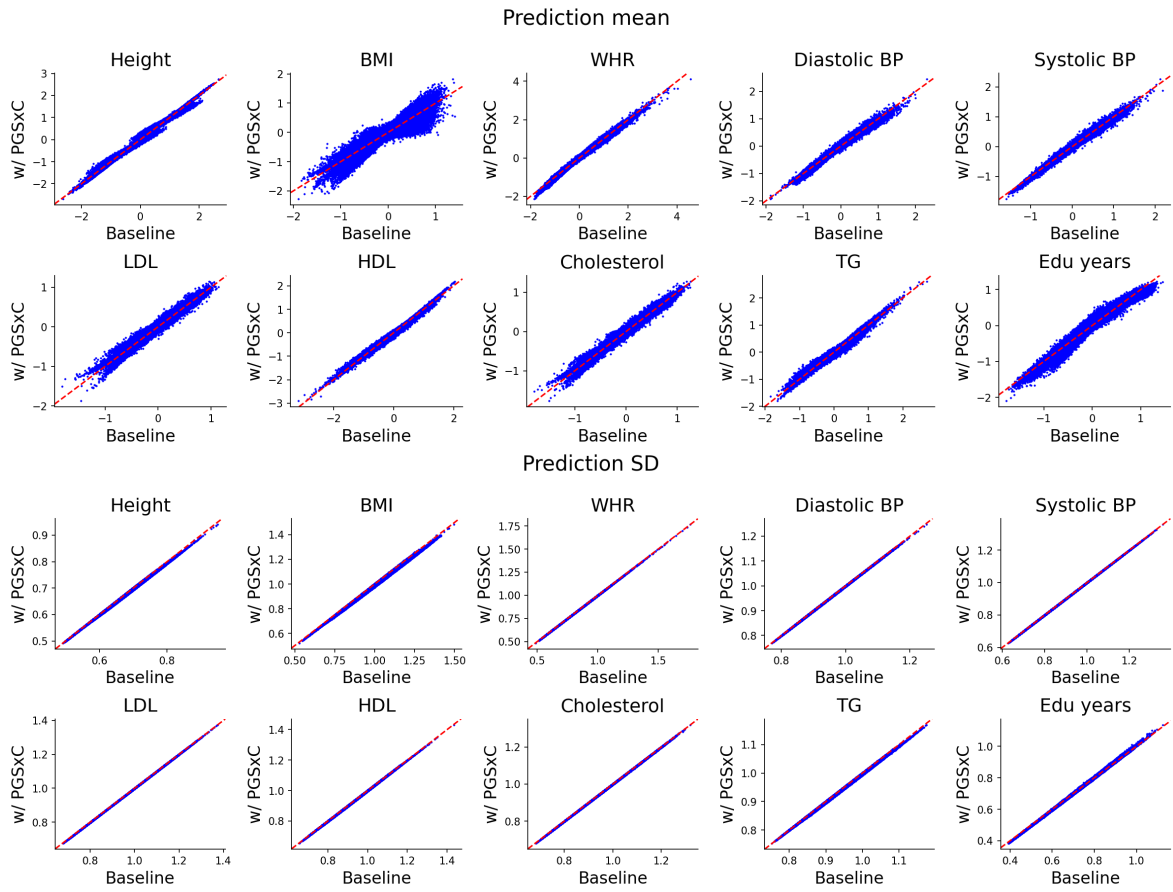**Supplementary Figure 6: (a)** Correlation of PGS evaluated across All of Us individuals derived from UK Biobank GWAS when (1) regressing out age, sex, top 16 PCs and (2) separately regressing out age$^2$, BMI, education years, income, or all of these contexts in addition to age, sex, top 16 PCs. By comparing results of separately regressing out each context with results regressing out all contexts, we determined that regressing out BMI had the most impact to context-specific accuracy. An exception is WHR regressed out of BMI, where PGS for WHR is only weakly correlated with PGS for WHR regressed out of BMI (correlation=0.75). This is expected because WHR regressed out of BMI (WHRadjBMI) has been routinely used as a standard derived phenotype as a proxy of abdominal adiposity and it has been shown that WHRadjBMI had moderate genetic overlap with WHR (Locke 2015 Nature[1]). **(b)** Context-specificity patterns for PGS derived from GWAS regressing out contexts of age$^2$, BMI, education years, income. These context-specificity patterns were similar to those in Fig. 3 without regressing out these contexts. **(c)** Consistency of relative ΔR$^2$ for all context-trait pairs, with or without regressing out additional contexts when performing GWAS.
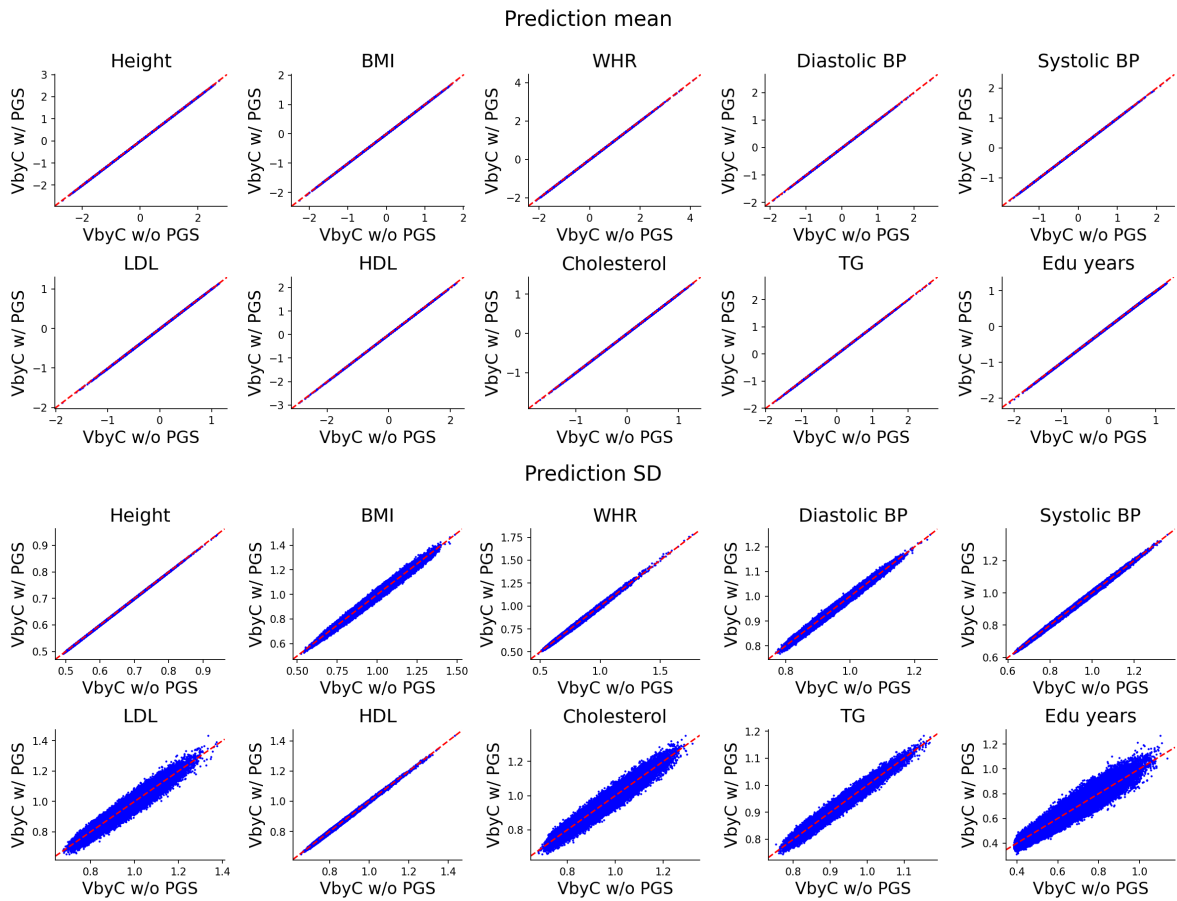
**Supplementary Figure 7: Simulations with varying phenotype-prediction slopes.** We simulated variable slopes (in regressing phenotypes against point predictions) in addition to variable prediction variances in Fig. 5. We applied CalPred with or without fitting the slope parameters (denoted by 'with slope' / 'without slope') and evaluated the coverage of prediction intervals. We simulate point predictions as $y \sim \mathcal{N}(\hat{y} \times (1 + \sum \beta_{\eta,c} \times c), \exp(\beta_{\sigma,0} + \sum_c \beta_{\sigma,c} \times c))$, where true $\beta_{\eta,\mathrm{age}} = 0, \beta_{\eta,\mathrm{sex}} = 0.3, \beta_{\eta,\mathrm{PC1}} = -0.05$. Other simulation settings were the same as in Fig. 5. **(a)** Example data with variable slope. Because $\beta_{\eta,\mathrm{sex}} = 0.3$, slopes between phenotype and point predictions were different between male and female. When interaction parameters $\beta_{\eta}$ were not fitted in the calibration, an overall fitted line was obtained. Consequently, at the extreme of point predictions, prediction intervals would have lower-than-expected coverages because of the shift in the mean predictions. **(b)** Coverage of prediction intervals evaluated across all 5,000 target individuals. CalPred produced prediction intervals with expected coverage level at the overall level regardless of whether interaction term is fitted or not. **(c)** Coverage was evaluated across 5% individuals at the tail distribution of point predictions (left tail~2.5%, right tail ~2.5%). When applied without fitted slope, prediction intervals had lower-than-expected coverage because of biased point predictions. When applied with fitted slope, prediction intervals resumed well-calibration at extremes of the distribution. **(d)** Numerical results of parameter estimation. We report the parameter values and 95% confidence intervals of estimated parameters across 100 simulations. 'Constant slope' column denotes simulations in Fig. 5. 'Variable slope (with slope)' / 'Variable slope (without slope)' denotes simulations in **(a-c)** with / without fitting interaction term. We determined that when model is correctly specified ('Constant slope' and 'Variable slope (with slope)'), parameter estimation was unbiased. When model is mis-specified ('Variable slope (without slope)'), parameter estimation remained robust. For (b) and (c), each box plot contains results across 100 simulations (each box contains n=100 points). For box plots, the center corresponds to the median; the box represents the first and third quartiles of the points; the whiskers represent the minimum and maximum points located within 1.5× interquartile ranges from the first and third quartiles, respectively.

**Supplementary Figure 8: Data likelihood comparison of different modeling strategies in ten quantitative traits in All of Us. (a)** We evaluated four modeling strategies with regard to their log-likelihood increase (ΔlogL) over baseline model (using PGS, age, sex, age*sex, age$^2$, and top 10 PCs, BMI, smoking status, drinking, employment, income, current address years, education and constant prediction variance across individuals): (1) "PGSxC": modeling PGSxC interaction in addition to baseline model; (2) "VbyC": modeling context-specific variance which does not include PGS; (3) "PGSxC+VbyC": modeling both PGSxC interaction and VbyC; (4) "PGSxC+VbyC(w/PGS)": modeling PGSxC interaction and VbyC where VbyC also contains PGS. We found limited improvement of model fitting when including PGS in VbyC and therefore did not include PGS in the main analysis because model interpretation is more straightforward when prediction variance is solely a function of contexts. However, it is technically valid to include PGS in VbyC as genetic contexts may modify prediction precisions. **(b)** The summation of ΔlogL(PGSxC) and ΔlogL(VbyC) is approximately equal to ΔlogL(PGSxC+VbyC), indicating that PGSxC and VbyC components had independent contribution to data modeling. **(c)** Prediction $R^2$ between prediction and phenotype for "context only": using contexts excluding PGS; "context+PGS" using contexts and PGS; "context+PGS+PGSxC" using contexts, PGS and PGSxC interaction terms. Results for "context+PGS+PGSxC+VbyC" were not shown as including VbyC terms did not impact prediction $R^2$.

**Supplementary Figure 9: Consistency of prediction mean and SD with/without PGSxC interaction terms.** Evaluation was performed in ten quantitative traits across All of Us individuals. PGSxC had moderate impact to prediction mean and had no impact to prediction SD, suggesting the VbyC and PGSxC captured independent aspects of trait modeling.

**Supplementary Figure 10: consistency of prediction mean and SD with/without PGS in modeling variance by contexts (VbyC).** Evaluation was performed across ten quantitative traits in all individuals from All of Us. We determined including PGS in VbyC had no impact to prediction mean and had minor impact to the prediction SD.

**Supplementary Figure 11: correspondence between quantile normalized and raw LDL levels for UK Biobank (a) and All of Us (b).** We plot distribution of LDL in raw unit of mg/dL in upper panels, and plot the correspondence between rank-based inverse normal transformed LDL and original LDL measurement in unit of mg/dL in lower panels. We note that distribution of LDL is right-skewed in both UK Biobank and All of Us. Therefore, the same unit increase in normalized scale at different LDL levels correspond to different amount of change in the original LDL measurement.

**Supplementary Figure 12: Results for other populations and traits in UK Biobank.** We plot results for LDL and height in "white British", and height in all individuals. **(top panel)** prediction $R^2$ between phenotype and point predictions. **(middle panel)** coverage of generic vs. context-specific 90% prediction intervals. **(bottom panel)** average length of generic vs. context-specific 90% prediction intervals in each context. See Fig. 6 caption for more details. For box plots, the center corresponds to the median; the box represents the first and third quartiles of the points; the whiskers represent the minimum and maximum points located within 1.5× interquartile ranges from the first and third quartiles, respectively. Each box contains results across 30 random samples with each sample of 5000 training and 5000 target individuals (n=30 points for each box).

**Supplementary Figure 13: Results for other traits in All of Us.** We plot results for LDL and height in "white SIRE" and height in all individuals. **(top panel)** prediction $R^2$ between phenotype and point predictions. **(middle panel)** coverage of generic vs. context-specific 90% prediction intervals. **(bottom panel)** average length of generic vs. context-specific 90% prediction intervals in each context. See Fig. 6 caption for more details. For box plots, the center corresponds to the median; the box represents the first and third quartiles of the points; the whiskers represent the minimum and maximum points located within 1.5× interquartile ranges from the first and third quartiles, respectively. Each box contains results across 30 random samples with each sample of 5000 training and 5000 target individuals (n=30 points for each box).

**Supplementary Figure 14: Comparison of model-predicted and observed SD levels incorporating all contexts.** We display model-predicted versus observed prediction SDs across traits incorporating effects of all contexts in All of Us. We observed high consistency between expected model-predicted SD and observed prediction SD (standard deviation of differences between observed and predicted phenotypes). For each trait, we also show the overall prediction SD produced by a model without modeling VbyC (vertical black lines). Here we did not show results for education years which was coded with only four discrete levels in our analysis and was not suitable to calculate prediction SD designed for continuous values.



**Supplementary Figure 15: Prediction interval length in predicting LDL as a function of PC1 and age.** Average length of context-specific 90% prediction intervals stratified by both five quintiles of PC1 and five quintiles of age. By contrasting individuals with youngest age and smallest PC1 (leftmost) quintiles versus those of oldest age and largest PC1 quintiles (rightmost), we find larger differences of prediction interval length across these subgroups compared to results obtained when single context (either age or PC1) is considered as in Fig. 6. By considering the contribution of both age and PC1 (two largest contributors to context-specific accuracy), we detected larger differences for individuals with youngest age and smallest PC1 quintiles (more similar to European) versus those of oldest age and largest PC1 quintiles (less similar to European) (27.3 vs. 36.3 mg/dL, 33% difference). Each box plot contains data across 30 random samples with each sample of 5,000 training individuals and 5,000 target individuals (n=30 points for each box plot)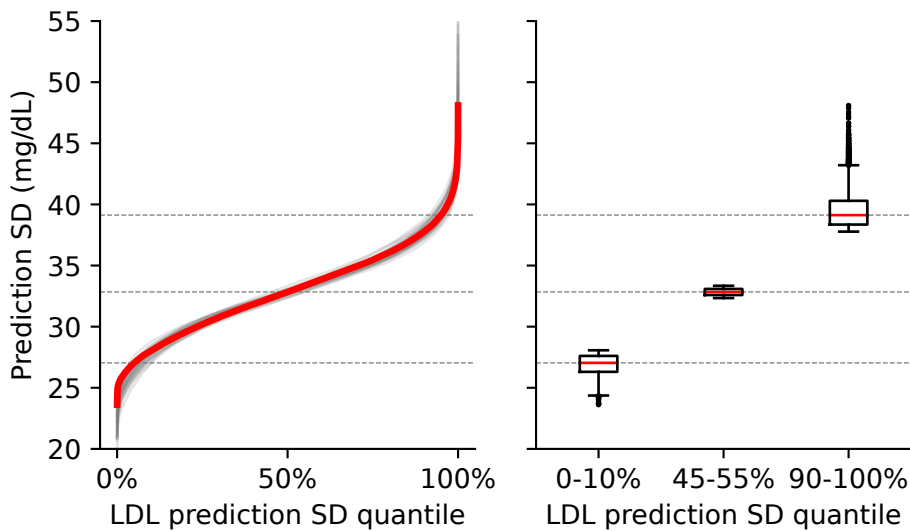, the center corresponds to the median; the box represents the first and third quartiles of the points; the whiskers represent the minimum and maximum points located within 1.5× interquartile range from the first and third quartiles, respectively.

**Supplementary Figure 16: Variation of prediction SD accounting for all contexts in All of Us. (a)** Ordered LDL prediction SD in unit of mg/dL. Gray lines denote prediction SD obtained with random sample of 5,000 training and applied to 5,000 testing individuals. Red line denote prediction SD obtained from all individuals. **(b)** Box plots of results in (a) from individuals of LDL prediction SD quantile of 0-10%, 45-55%, 90-100% (n=96K individuals in total); the center corresponds to the median; the box represents the first and third quartiles of the points; the whiskers represent the minimum and maximum points located within 1.5× interquartile range from the first and third quartiles, respectively.



**Supplementary Figure 17: Variation of prediction standard deviation (SD) accounting for all contexts across prediction SD quintiles.** Relative difference of prediction SD between top and bottom prediction SD quintiles (80-100% vs. 0-20%) for all traits in UK Biobank **(a)** and All of Us **(b)**. Traits are ranked by prediction SD. See Fig. 7 for additional details.

**Supplementary Figure 18: Comparison of different types of models in predicting disease traits in All of Us.** We evaluated three modeling strategies in their **(a)** incremental data likelihood and **(b)** area under the curve (AUC) compared to baseline models. "Baseline": a baseline model with prediction mean modeled using age, sex, age*sex, $age^2$, BMI, and top 10 PCs; "Context": additionally include contexts of smoking status, drinking, employment, income, current address years and education; "Context+PGSxC": additionally include PGSxC interaction terms; "Context+PGSxC+VbyC": additionally include VbyC terms.

(see next four pages)

**Supplementary Figure 19: Calibration of predicted probability.** Calibration curve of observed proportions versus predicted probabilities for different context groups in prediction models for all individuals in All of Us. We fit four models using different sets of predictors in regression across all individuals: "Baseline": using PGS, age, sex, age*sex, $age^2$, BMI, and top 10 PCs; "Baseline+C": also including smoking status, drinking, employment, income, current address years and education; "Baseline+C+PGSxC": also including corresponding PGSxC interaction; "Baseline+C+PGSxC; VbyC": additionally model variance by contexts within a liability threshold model. For binary contexts, we show the calibration curve separately for two groups. For continuous context variables, we divided individuals into five quintiles and compared the calibration between top and bottom quintile groups. Error bars denote observed disease proportions and their 95% confidence intervals for each predicted probability bin (n="number of individuals shown in parenthesis" / 20 bins). See Fig. 8 caption for additional details.

Prostate_Cancer (Baseline)

Prostate_Cancer (Baseline+C)

Prostate_Cancer (Baseline+C+PGSxC)

Prostate_Cancer (Baseline+C+PGSxC; VbyC)

Breast_Cancer (Baseline)

| Overall | Age | Edu years | Income | Current address years | Smoking |
|---|---|---|---|---|---|
| AUC=0.776 (N=93733) | AUC=0.703 (N=18797)<br>AUC=0.638 (N=17116) | AUC=0.753 (N=8475)<br>AUC=0.759 (N=40493) | AUC=0.780 (N=23161)<br>AUC=0.735 (N=18026) | AUC=0.829 (N=30190)<br>AUC=0.694 (N=17664) | AUC=0.779 (N=83575)<br>AUC=0.740 (N=10158) |

| Employment | Alcohol | BMI | PC1 | PC2 |
|---|---|---|---|---|
| AUC=0.760 (N=50405)<br>AUC=0.788 (N=43328) | AUC=0.765 (N=81111)<br>AUC=0.817 (N=356) | AUC=0.782 (N=18791)<br>AUC=0.776 (N=18573) | AUC=0.737 (N=18754)<br>AUC=0.747 (N=18747) | AUC=0.794 (N=18747)<br>AUC=0.749 (N=18746) |

Breast_Cancer (Baseline+C)

| Overall | Age | Edu years | Income | Current address years | Smoking |
|---|---|---|---|---|---|
| AUC=0.783 (N=93733) | AUC=0.700 (N=18797)<br>AUC=0.645 (N=17116) | AUC=0.761 (N=8475)<br>AUC=0.762 (N=40493) | AUC=0.786 (N=23161)<br>AUC=0.737 (N=18026) | AUC=0.835 (N=30190)<br>AUC=0.697 (N=17664) | AUC=0.786 (N=83575)<br>AUC=0.750 (N=10158) |

| Employment | Alcohol | BMI | PC1 | PC2 |
|---|---|---|---|---|
| AUC=0.767 (N=50405)<br>AUC=0.794 (N=43328) | AUC=0.772 (N=81111)<br>AUC=0.824 (N=356) | AUC=0.789 (N=18791)<br>AUC=0.782 (N=18573) | AUC=0.747 (N=18754)<br>AUC=0.762 (N=18747) | AUC=0.797 (N=18747)<br>AUC=0.759 (N=18746) |

Breast_Cancer (Baseline+C+PGSxC)

| Overall | Age | Edu years | Income | Current address years | Smoking |
|---|---|---|---|---|---|
| AUC=0.784 (N=93733) | AUC=0.702 (N=18797)<br>AUC=0.647 (N=17116) | AUC=0.765 (N=8475)<br>AUC=0.763 (N=40493) | AUC=0.787 (N=23161)<br>AUC=0.739 (N=18026) | AUC=0.836 (N=30190)<br>AUC=0.699 (N=17664) | AUC=0.787 (N=83575)<br>AUC=0.750 (N=10158) |

| Employment | Alcohol | BMI | PC1 | PC2 |
|---|---|---|---|---|
| AUC=0.768 (N=50405)<br>AUC=0.795 (N=43328) | AUC=0.773 (N=81111)<br>AUC=0.829 (N=356) | AUC=0.790 (N=18791)<br>AUC=0.784 (N=18573) | AUC=0.747 (N=18754)<br>AUC=0.765 (N=18747) | AUC=0.799 (N=18747)<br>AUC=0.760 (N=18746) |

Breast_Cancer (Baseline+C+PGSxC; VbyC)

| Overall | Age | Edu years | Income | Current address years | Smoking |
|---|---|---|---|---|---|
| AUC=0.785 (N=93733) | AUC=0.708 (N=18797)<br>AUC=0.647 (N=17116) | AUC=0.765 (N=8475)<br>AUC=0.764 (N=40493) | AUC=0.788 (N=23161)<br>AUC=0.741 (N=18026) | AUC=0.837 (N=30190)<br>AUC=0.698 (N=17664) | AUC=0.788 (N=83575)<br>AUC=0.751 (N=10158) |

| Employment | Alcohol | BMI | PC1 | PC2 |
|---|---|---|---|---|
| AUC=0.769 (N=50405)<br>AUC=0.796 (N=43328) | AUC=0.774 (N=81111)<br>AUC=0.829 (N=356) | AUC=0.791 (N=18791)<br>AUC=0.784 (N=18573) | AUC=0.748 (N=18754)<br>AUC=0.765 (N=18747) | AUC=0.799 (N=18747)<br>AUC=0.760 (N=18746) |

## T2D (Baseline)

**Overall** — AUC=0.850 (N=92879)

**Age** — AUC=0.859 (N=19077); AUC=0.798 (N=18017)

**Sex** — AUC=0.857 (N=58825); AUC=0.835 (N=34054)

**Edu years** — AUC=0.823 (N=8251); AUC=0.865 (N=40506)

**Income** — AUC=0.832 (N=22115); AUC=0.871 (N=10626)

**Current address years** — AUC=0.870 (N=28733); AUC=0.827 (N=33520)

**Smoking** — AUC=0.855 (N=79268); AUC=0.819 (N=13611)

**Employment** — AUC=0.828 (N=50857); AUC=0.866 (N=42022)

**Alcohol** — AUC=0.848 (N=77276); AUC=0.821 (N=740)

**BMI** — AUC=0.854 (N=18767); AUC=0.765 (N=18517)

**PC1** — AUC=0.843 (N=18579); AUC=0.811 (N=18576)

**PC2** — AUC=0.863 (N=18576); AUC=0.846 (N=18575)

## T2D (Baseline+C)

**Overall** — AUC=0.856 (N=92879)

**Age** — AUC=0.863 (N=19077); AUC=0.806 (N=18017)

**Sex** — AUC=0.863 (N=58825); AUC=0.841 (N=34054)

**Edu years** — AUC=0.823 (N=8251); AUC=0.869 (N=40506)

**Income** — AUC=0.836 (N=22115); AUC=0.873 (N=10626)

**Current address years** — AUC=0.873 (N=28733); AUC=0.835 (N=33520)

**Smoking** — AUC=0.859 (N=79268); AUC=0.826 (N=13611)

**Employment** — AUC=0.833 (N=50857); AUC=0.869 (N=42022)

**Alcohol** — AUC=0.854 (N=77276); AUC=0.821 (N=740)

**BMI** — AUC=0.865 (N=18767); AUC=0.773 (N=18517)

**PC1** — AUC=0.854 (N=18579); AUC=0.812 (N=18576)

**PC2** — AUC=0.867 (N=18576); AUC=0.857 (N=18575)

## T2D (Baseline+C+PGSxC)

**Overall** — AUC=0.858 (N=92879)

**Age** — AUC=0.862 (N=19077); AUC=0.809 (N=18017)

**Sex** — AUC=0.865 (N=58825); AUC=0.843 (N=34054)

**Edu years** — AUC=0.829 (N=8251); AUC=0.870 (N=40506)

**Income** — AUC=0.838 (N=22115); AUC=0.874 (N=10626)

**Current address years** — AUC=0.875 (N=28733); AUC=0.837 (N=33520)

**Smoking** — AUC=0.862 (N=79268); AUC=0.829 (N=13611)

**Employment** — AUC=0.837 (N=50857); AUC=0.870 (N=42022)

**Alcohol** — AUC=0.856 (N=77276); AUC=0.825 (N=740)

**BMI** — AUC=0.871 (N=18767); AUC=0.774 (N=18517)

**PC1** — AUC=0.855 (N=18579); AUC=0.817 (N=18576)

**PC2** — AUC=0.868 (N=18576); AUC=0.859 (N=18575)

## T2D (Baseline+C+PGSxC; VbyC)

**Overall** — AUC=0.859 (N=92879)

**Age** — AUC=0.870 (N=19077); AUC=0.809 (N=18017)

**Sex** — AUC=0.867 (N=58825); AUC=0.844 (N=34054)

**Edu years** — AUC=0.831 (N=8251); AUC=0.871 (N=40506)

**Income** — AUC=0.840 (N=22115); AUC=0.875 (N=10626)

**Current address years** — AUC=0.877 (N=28733); AUC=0.838 (N=33520)

**Smoking** — AUC=0.863 (N=79268); AUC=0.829 (N=13611)

**Employment** — AUC=0.838 (N=50857); AUC=0.871 (N=42022)

**Alcohol** — AUC=0.857 (N=77276); AUC=0.821 (N=740)

**BMI** — AUC=0.873 (N=18767); AUC=0.777 (N=18517)

**PC1** — AUC=0.856 (N=18579); AUC=0.819 (N=18576)

**PC2** — AUC=0.869 (N=18576); AUC=0.860 (N=18575)

CAD (Baseline)

| Overall | Age | Sex | Edu years | Income | Current address years |

AUC=0.813 (N=170474)

Age: AUC=0.694 (N=36007), AUC=0.679 (N=32411)
Sex: AUC=0.799 (N=107485), AUC=0.794 (N=62989)
Edu years: AUC=0.797 (N=15261), AUC=0.826 (N=74373)
Income: AUC=0.791 (N=41301), AUC=0.822 (N=18902)
Current address years: AUC=0.847 (N=53941), AUC=0.745 (N=33963)

Smoking: AUC=0.813 (N=145972), AUC=0.770 (N=24502)
Employment: AUC=0.768 (N=93998), AUC=0.840 (N=76476)
Alcohol: AUC=0.806 (N=142043), AUC=0.818 (N=1399)
BMI: AUC=0.827 (N=34776), AUC=0.795 (N=33676)
PC1: AUC=0.801 (N=34108), AUC=0.759 (N=34095)
PC2: AUC=0.842 (N=34095), AUC=0.800 (N=34094)

CAD (Baseline+C)

AUC=0.819 (N=170474)

Age: AUC=0.712 (N=36007), AUC=0.685 (N=32411)
Sex: AUC=0.808 (N=107485), AUC=0.798 (N=62989)
Edu years: AUC=0.801 (N=15261), AUC=0.828 (N=74373)
Income: AUC=0.795 (N=41301), AUC=0.822 (N=18902)
Current address years: AUC=0.851 (N=53941), AUC=0.751 (N=33963)

Smoking: AUC=0.817 (N=145972), AUC=0.775 (N=24502)
Employment: AUC=0.773 (N=93998), AUC=0.842 (N=76476)
Alcohol: AUC=0.812 (N=142043), AUC=0.820 (N=1399)
BMI: AUC=0.832 (N=34776), AUC=0.802 (N=33676)
PC1: AUC=0.808 (N=34108), AUC=0.762 (N=34095)
PC2: AUC=0.847 (N=34095), AUC=0.808 (N=34094)

CAD (Baseline+C+PGSxC)

AUC=0.820 (N=170474)

Age: AUC=0.713 (N=36007), AUC=0.685 (N=32411)
Sex: AUC=0.809 (N=107485), AUC=0.800 (N=62989)
Edu years: AUC=0.804 (N=15261), AUC=0.829 (N=74373)
Income: AUC=0.798 (N=41301), AUC=0.822 (N=18902)
Current address years: AUC=0.852 (N=53941), AUC=0.752 (N=33963)

Smoking: AUC=0.818 (N=145972), AUC=0.777 (N=24502)
Employment: AUC=0.774 (N=93998), AUC=0.843 (N=76476)
Alcohol: AUC=0.813 (N=142043), AUC=0.821 (N=1399)
BMI: AUC=0.833 (N=34776), AUC=0.804 (N=33676)
PC1: AUC=0.808 (N=34108), AUC=0.766 (N=34095)
PC2: AUC=0.848 (N=34095), AUC=0.808 (N=34094)

CAD (Baseline+C+PGSxC; VbyC)

AUC=0.821 (N=170474)

Age: AUC=0.718 (N=36007), AUC=0.687 (N=32411)
Sex: AUC=0.810 (N=107485), AUC=0.802 (N=62989)
Edu years: AUC=0.805 (N=15261), AUC=0.830 (N=74373)
Income: AUC=0.800 (N=41301), AUC=0.824 (N=18902)
Current address years: AUC=0.854 (N=53941), AUC=0.753 (N=33963)

Smoking: AUC=0.820 (N=145972), AUC=0.778 (N=24502)
Employment: AUC=0.775 (N=93998), AUC=0.843 (N=76476)
Alcohol: AUC=0.814 (N=142043), AUC=0.823 (N=1399)
BMI: AUC=0.835 (N=34776), AUC=0.805 (N=33676)
PC1: AUC=0.809 (N=34108), AUC=0.771 (N=34095)
PC2: AUC=0.849 (N=34095), AUC=0.809 (N=34094)

# Supplementary Tables

**Supplementary Table 1: Trait information in UK Biobank.** We report trait names, sample size used in training PGS weights, estimated heritability in the training sample, prediction $R^2$ (between PGS and covariate-adjusted phenotypes) and sample sizes used in testing populations, separately for white British and all populations. These traits are selected for their sufficient predictive power and/or biological importance.

**Supplementary Table 2: Numerical results of relative $\Delta R^2$ and estimated $\beta_\sigma$ in UK Biobank.** We report trait, context, relative $\Delta R^2$ differences between groups, and z-score for the significance of $R^2$ differences between groups. We also report $\beta_\sigma$ estimates and their standard errors.

**Supplementary Table 3: Trait information in All of Us.** We report trait names, prediction $R^2$ (between PGS and covariate-adjusted phenotypes) and sample size used in testing populations, separately for self-reported race of white and all populations.

**Supplementary Table 4: Numerical results of relative $\Delta R^2$ and estimated $\beta_\sigma$ in All of Us.** We report trait, context, relative $\Delta R^2$ differences between groups, and z-score for the significance of $R^2$ differences between groups. We also report $\beta_\sigma$ estimates and corresponding standard errors.

# Supplementary Note

## Comparison to other PGS calibration methods

The observation that PGS distribution differs across genetic ancestry continuum[2] motivates methods that regress out effects of variables representing genetic ancestry from PGS distribution to facilitate comparison across individuals locating at different positions in genetic ancestry continuum[3,4]. However, such approaches may unintentionally remove true biological differences of PGS distribution across genetic ancestry continuum (e.g., African Americans have reduced neutrophil count explained by the large effect of a single Duffy-null SNP[5]) as they do not consider phenotype value distribution in calibration procedure; in addition, these approaches do not represent different standard errors in PGS predictions for individuals across genetic ancestry continuum. Our method CalPred leverages a set of calibration data to adjust point predictions across contexts according to true phenotype distribution. Compared to other data-driven calibration methods[6], our approach provides a framework to incorporate context information.

## Comparison between single-context and combined-context analyses

Here, we clarify differences between single-context and combined-context analyses. Overall, single-context analyses compared $R^2$ between covariate-adjusted phenotypes and PGS, while combined-context analyses assessed how residual phenotypic variance increases by contexts in a regression model. Technically, single-context analyses require discretization of continuous context variable into pre-specified groups to calculate $R^2$ differences across groups, while combined-context analyses use a regression model that naturally accommodate continuous context variable. Moreover, single-context analyses only consider one context variable at a time, while combined-context analyses consider multiple context variables jointly and can evaluate effect of each context variable conditional on other contexts.

To illustrate their differences, we consider a model where phenotype $y$ is a function of predicted phenotype $\hat{y}$ multiplied by a slope $s$ and environmental noise $e$: $y = s \cdot \hat{y} + e$, and we discuss how $R^2(y, \hat{y})$ varies as a function of both $s$ and $\mathrm{Var}[e]$. The slope $s$ can vary across contexts (e.g., phenotype-PGS regression slope $s$ can vary across sex (see ref.[7]). Here prediction $R^2$ is

$$R^2(y, \hat{y}) = R^2(s \cdot \hat{y} + e, \hat{y}) = \frac{s^2 \cdot \mathrm{Var}[\hat{y}]}{s^2 \cdot \mathrm{Var}[\hat{y}] + \mathrm{Var}[e]}.$$

Holding $\mathrm{Var}[\hat{y}]$ as constant, $R^2(y, \hat{y})$ increases with increasing $s$ and decreasing $\mathrm{Var}[e]$ – variable $R^2$ can result from variable slope $s$ when residual variance $\mathrm{Var}[e]$ and prediction interval length are constant. Therefore, the distinction between single- and combined- context analyses is that single-context analyses models $R^2$ while combined-context analyses models $\mathrm{Var}[e]$.

Next, we performed case studies on two context-trait pairs (sex-BMI, sex-WHR) where we observed large differences between single-context and combined-context analyses. We focused on evaluation within white British individuals from UK Biobank to minimize confounding due to differences across ancestry groups. First, for sex-BMI, we observed almost no cross-sex difference in $R^2_{PGS}$ between covariate-adjusted phenotype and PGS, while we observed a significant $\beta_\sigma(\mathrm{sex})$. We determined that this was because both $\mathrm{Var}[\hat{y}]$ and $\mathrm{Var}[e]$ varied proportionally across sex, while the ratio between the two $R^2(y, \hat{y})$ did not vary substantially across sex.

| Sex | $R^2$ | $R^2_{PGS}$ | SD(pred) | SD(resid) | AVG(y) | SD(y) |
|---|---|---|---|---|---|---|
| Female | 0.159 | 0.137 | 0.423 | 0.992 | -0.12 | 1.082 |
| Male | 0.146 | 0.136 | 0.343 | 0.809 | 0.139 | 0.875 |

**Supplementary Note Table 1.** Statistics of PGS-based predictions of BMI across sex in UK Biobank. $R^2$ denotes squared Pearson correlation between measured phenotype and predicted phenotype incorporating PGS, age, sex and top 10 PCs. $R^2_{PGS}$ denotes squared Pearson correlation between covariate-adjusted phenotype (adjusted for age, sex and top 10 PCs) and PGS. SD(pred), SD(resid), SD(y) denote standard deviation of predicted phenotype, residual phenotype and phenotype, respectively. AVG(y) denotes average of phenotype.

Second, for sex-WHR, we observed lower $R^2_{PGS}$ between covariate-adjusted phenotype and PGS in male compared to female in single-context analyses, while the sign of $\beta_\sigma$ indicated lower prediction residual error in male compared to female, which seemed to be inconsistent with results from $R^2_{PGS}$. This was explained by that

$\beta_\sigma$ captured the residual error level in trait prediction, and the overall $R^2$ between prediction (with contribution from including PGS and other contexts) and phenotype is higher in male compared to female. This discrepancy between $R^2$ and $R^2_{PGS}$ was because BMI – the context with the highest prediction power to WHR had a much higher effect in male compared to female.
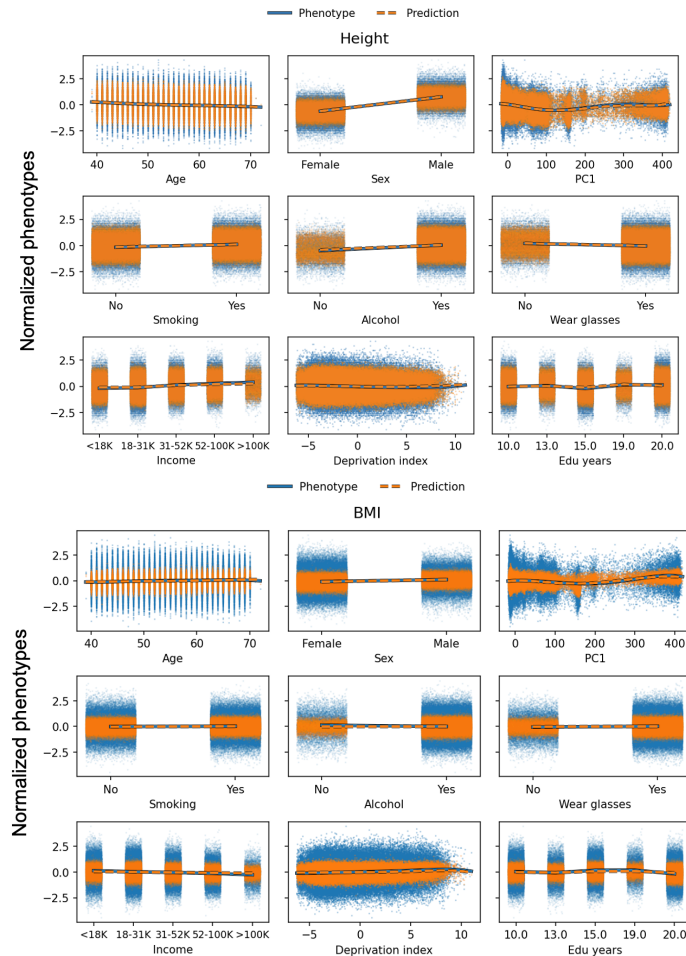
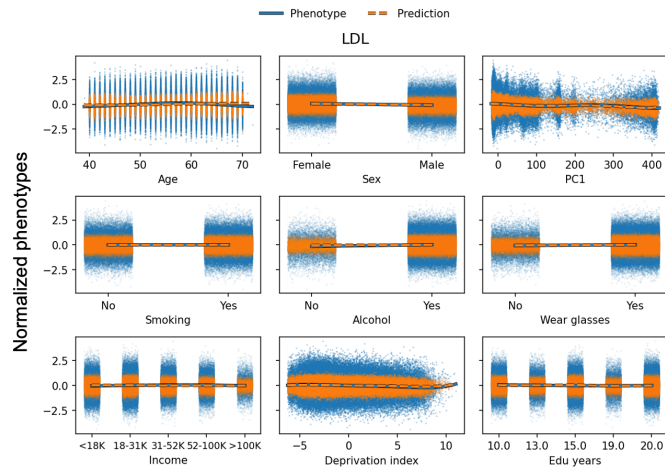| Sex | $R^2$ | $R^2_{PGS}$ | SD(pred) | SD(resid) | AVG(y) | SD(y) |
|---|---|---|---|---|---|---|
| Female | 0.31 | 0.104 | 0.45 | 0.656 | -0.603 | 0.79 |
| Male | 0.438 | 0.064 | 0.48 | 0.541 | 0.702 | 0.721 |

**Supplementary Note Table 2.** Statistics of PGS-based predictions of WHR across sex in UK Biobank. $R^2$ denotes squared Pearson correlation between phenotype and predicted phenotype incorporating PGS, age, sex and top 10 PCs. $R^2_{PGS}$ denotes squared Pearson correlation between covariate-adjusted phenotype (adjusted for age, sex and top 10 PCs) and PGS. SD(pred), SD(resid), SD(y) denote standard deviation of predicted phenotype, residual phenotype and phenotype, respectively. AVG(y) denotes average of phenotype.

Overall, the commonly-used metric of $R^2_{PGS}$ between covariate-adjusted phenotype and PGS and $\beta_\sigma$ in this work are related but different quantities for PGS-based predictions. The extent of concordance of the two metrics depends on the trait. Notably, $\beta_\sigma$ we estimated here directly relates to individual-specific residual noise level when PGS and other contexts are jointly used to predict quantitative trait values.

**Investigating source of variable accuracy due to bias versus conditional variance**
With $y_i = \mathcal{N}\big(\mu(\mathbf{c}_i), \sigma^2(\mathbf{c}_i)\big)$, CalPred models the variable conditional variance through $\sigma^2(\mathbf{c}_i)$. Therefore, CalPred performance relies on the unbiasedness of prediction mean $\mu(\mathbf{c}_i)$ to properly model the conditional variance term. We investigated the bias of prediction mean by comparing distribution of phenotype values versus prediction mean across contexts. We performed the comparison for three example traits of height, BMI and LDL across all individuals in UK Biobank.

Across three example traits, we observed that the prediction mean tracked well with the phenotypic mean, indicating that prediction biases were small across context groups and that CalPred model assumption was valid. Meanwhile, more fine-grained modeling of prediction factors that capture more variation, for example, by modeling interactions between prediction factors will benefit CalPred model via more precise, and shorter, prediction intervals.

## References

1. Locke, A. E. *et al.* Genetic studies of body mass index yield new insights for obesity biology. *Nature* **518**, 197–206 (2015).

2. Martin, A. R. *et al.* Human demographic history impacts genetic risk prediction across diverse populations. *Am. J. Hum. Genet.* **107**, 788–789 (2020).

3. Hao, L. *et al.* Development of a clinical polygenic risk score assay and reporting workflow. *Nat. Med.* **28**, 1006–1013 (2022).

4. Khera, A. V. *et al.* Whole-genome sequencing to characterize monogenic and polygenic contributions in patients hospitalized with early-onset myocardial infarction. *Circulation* **139**, 1593–1602 (2019).

5. Reich, D. *et al.* Reduced neutrophil count in people of African descent is due to a regulatory variant in the Duffy antigen receptor for chemokines gene. *PLoS Genet.* **5**, e1000360 (2009).

6. Sun, J. *et al.* Translating polygenic risk scores for clinical use by estimating the confidence bounds of risk prediction. *Nat. Commun.* **12**, 5276 (2021).

7. Zhu, C., Ming, M. J., Cole, J. M., Kirkpatrick, M. & Harpak, A. Amplification is the primary mode of gene-by-sex interaction in complex human traits. (2022) doi:10.1101/2022.05.06.490973.