

1 **Using synthetic datasets to better understand and explain health outcomes associated with**
2 **common single nucleotide polymorphisms**

3
4 Thomas R. Wood^{1,2} and Nathan Owens³

5
6 ¹Division of Neonatology, Department of Pediatrics, University of Washington, Seattle, WA

7 ²Institute for Human and Machine Cognition, Pensacola, FL

8 ³Independent Researcher, Seattle, WA

9
10 **Corresponding author:**

11 Dr. Thomas Wood, BM BCh, PhD

12 Division of Neonatology

13 Department of Pediatrics

14 University of Washington Medical Center

15 Box 356320, HSB RR-535

16 Seattle, Washington 98195

17 E-mail: tommyrw@uw.edu

18
19 **Key words:** Genetics, single nucleotide polymorphisms, risk, obesity, methylation, blood glucose

1 **ABSTRACT**

2 Due to decreasing costs and a move towards “personalised medicine”, the use of direct-to-consumer genetic
3 analyses is increasing. Both consumers and healthcare practitioners must therefore be able to understand
4 the true disease risks associated with common genetic single nucleotide polymorphisms (SNPs). However,
5 most population studies of common SNPs only provide average (+/- error) phenotypic or risk descriptions for
6 a given genotype, which hides the true heterogeneity of the population and reduces the ability of an individual
7 to determine how they themselves might truly be effected. Here, we describe the use of synthetic datasets
8 generated from descriptive phenotypic data published on common SNPs associated with obesity, elevated
9 fasting blood glucose, and methylation status. Using both simple statistical theory and full graphical
10 representation of the generated data, we show that single common SNPs are associated with a less than
11 10% likelihood of effecting final phenotype, even in homozygotes. The significant heterogeneity in the data,
12 as well as the baseline disease risk of Western populations suggests that most disease risk is dominated by
13 the effect of the modern environment.

14

1 INTRODUCTION

2 Due to decreasing costs and a move towards “personalised medicine”, the use of direct-to-consumer
3 (DTC) genetic analyses and third party interpretation services is increasing.¹ Though whole genome
4 sequencing (WGS) is also increasing in popularity, most DTC products involve the analysis of common single
5 nucleotide polymorphism (SNPs). These SNPs are then reported, either by the testing company or a third
6 party tool that analyses the data, with specific disease risks based on published population data such as that
7 from genome-wide association studies (GWAS). These risk predictions are generally based on population
8 average outcomes, with the heterogeneity of a given phenotype or disease risk infrequently reported. In fact,
9 most GWAS studies tend to only report descriptive data (e.g. mean and standard error) for a given phenotype
10 (such as body mass index, BMI, or fasting blood glucose) within a risk genotype. By only comparing or
11 providing group averages based on genotype, the consumer is likely to overestimate the disease risk
12 associated with a given SNP. Presenting only simplified descriptive data, either graphically or numerically,
13 for a given genotype gives the impression that each SNP has consistent penetrance with respect to the
14 phenotype in question, which is known to not be the case.² Therefore, the interpretation of disease risk based
15 on SNPs by those not involved in the original studies and without access to the original data is almost
16 impossible.

17 More important than the mean population effects of a given SNP or combination of SNPs that
18 influence a common phenotype is the likelihood of a physiologically-relevant effect in a given individual. This
19 includes the likelihood that there is no overall effect of genotype, particularly compared to common
20 environmental factors that drive chronic disease risk in high income countries such as diet, sleep, and
21 exercise. In order to allow for healthcare practitioners or self-interested parties to better understand the
22 likelihood of a given phenotype being altered by a specific genotype, we developed a method by which
23 synthetic datasets could be generated and analysed. This is largely possible due to the fact that the effects
24 of SNPs on measurable phenotypes are generally considered to follow a normal distribution, with the number
25 of alleles or weighted genetic scores being linearly associated with the target phenotype. Using this approach,
26 the significant heterogeneity of population data can be better understood, particularly with respect to how a
27 given individual may or may not display phenotypic changes based on the presence of common genotypes.

28 METHODS

29 *Selection of representative SNPs*
30

1 One of the authors (TRW) performed at-home SNP analysis using a 23andMe DTC kit (23andMe, Mountain
2 View, CA). The data were run through a third-party analysis tool (FoundMyFitness Genetic Report) to identify
3 SNPs most commonly reported to be associated with differential disease risk. Individual studies and meta-
4 analyses of per allele effects for common SNPs most strongly-associated with risk of type 2 diabetes
5 (Melatonin Receptor 1B, MTNR1B rs10830963), obesity (Fat mass and obesity-associated protein, FTO
6 rs9939609), and altered methylation and nutrient handling resulting in elevated homocysteine levels
7 (Methylenetetrahydrofolate Reductase, MTHFR rs1801131 and rs1801133) were identified from the third
8 party tool output as well as the online SNP wiki SNPedia.com. Due to the significant effect of ethnicity on
9 SNP disease penetrance, example population data that were likely to most closely match the Anglo-
10 Scandinavian background of the author were used, including data from deCODE (Iceland) and the Northern
11 Finland Birth Cohort (NFBC) that were included in large multi-population GWAS studies.^{3,4} According to
12 recently-published methods suggested by Pontzer *et al.*, published hunter gatherer data for fasting glucose
13 were used to provide an estimate of the effect of the Western environment on fasting glucose and diabetes
14 risk compared to a published genetic risk score.⁵

16 *Generation of synthetic datasets*

17 Published per allele or per genetic risk score means were used to construct synthetic datasets for a given
18 phenotype. All publications assumed data were normally distributed and that per allele/genetic risk score
19 effects were linear. If data were expressed as mean with standard error (SE) or 95% confidence interval (CI),
20 the standard deviation (SD) was calculated using the number (N) of participants in each group, where
21 $SD=SE*\sqrt{N}$ and $SD=\sqrt{N}*(width\ of\ 95\%\ CI)/3.92$. When the descriptive data were not included in the
22 publication, as was the case for genetic risk scores associated with obesity and fasting blood glucose,^{3,4} they
23 were estimated from published graphs by extracting images and determining the number of pixels in each
24 column and error bar relative to the scale bars on the axes. In all cases, enough data was included in the
25 manuscript body to confirm that at least one of the estimated values was correctly determined using this
26 method (such as total number of participants, or mean values in the highest or lowest genetic risk groups).
27 For each genotype and gene, 1,000 synthetic individuals were randomly generated to re-create a normally-
28 distributed dataset with the same mean and SD characteristics as those in the associated publication.
29 Numbers were generated using Python 3.7 and the NumPy (1.17.0) and Pandas (0.25.0) libraries. The
30 necessary code is available on GitHub (<https://github.com/root-causing-health/SNPGaussianDistGenerator>).

1 Visual inspection of the data (Prism version 8, GraphPad Software, San Diego, CA) confirmed that they were
2 normally distributed.

4 *Statistical analysis*

5 Each synthetic dataset was graphically represented using a violin plot to show the full distribution of the data.
6 Percent chance of a null effect from a risk allele was calculated by determining the percent overlap of the
7 normal distribution of the wild type phenotype with that of a risk genotype using statistics.NormalDist in
8 Python 3.8 Beta. The percent likelihood of the phenotype in a risk allele group being at or below the mean
9 value of the “wild type” was also calculated, and linear regression analysis was performed to determine the
10 percent contribution of risk alleles to a given phenotype. Similar analyses were performed using published
11 multi-SNP genetic risk scores for type 2 diabetes and obesity.^{3,4}

13 *Alternative methods*

14 To encourage attempts to perform similar analyses, a number of free online tools can be used that do not
15 require significant technical skills. After calculating mean and SD as described above, free gaussian random
16 number generators such as from Random.org (<https://www.random.org/gaussian-distributions/>) can be used
17 to generate synthetic datasets. Though the Box-Muller transform used by this tool is unlikely to produce a
18 truly normal distribution,⁶ this is also unlikely to meaningfully affect the outcome. Similar online tools can be
19 used to determine the likelihood of being at, above, or below, a given point in a normal distribution to
20 determine null effects of a given SNP or risk score (http://onlinestatbook.com/2/calculators/normal_dist.html).
21 Finally, free online graphing software can be used to visually represent the datasets for visual examination
22 of variability and overlap (<https://plot.ly/>), and perform linear regression analyses
23 (<https://www.graphpad.com/quickcalcs/linear1/>).

25 **RESULTS AND DISCUSSION**

26 *FTO rs9939609 (A:T) and risk of being overweight*

27 Published meta-analyses suggest an increase in body mass index (BMI) of 0.3 kg/m² per FTO rs9939609 A
28 allele.⁷ From this meta-analysis, data from the NFBC at 31 years of age (n=4,435) were used as a graphical
29 example (**Figure 1A**).⁸ Mean (SD) BMI across the three genotypes was 24.12 (3.87) kg/m², 24.43 (3.94)
30 kg/m², and 24.82 (3.95) kg/m² for TT, AT, and AA respectively. In this population the risk of being overweight

1 (BMI >25 kg/m²) was 41%, 44%, and 48%, resulting in an absolute 7% increase in risk in the TT genotype.
2 BMI at or below the TT genotype was 47% in those with the AT genotype, and 43% in those with the TT
3 genotype. The likelihood of null effect (percent overlap in BMI distribution of those with AT and AA genotypes
4 compared to TT) was 96.8% and 92.8%, respectively. Therefore, only 3.2% of AT and 7.2% of AA genotypes
5 would be expected to display any increase in BMI due to FTO genotype relative to TT. Linear regression
6 found a significant association between number of A copies and BMI ($p=0.001$, $R_2=0.0035$), suggesting that
7 only around 0.4% of the variability in BMI is determined by FTO genotype (**Figure 1B**).

9 *Genetic BMI risk score*

10 Willer *et al.* established a BMI genetic score using eight validated SNPs associated with BMI, weighted to
11 effect size (with FTO rs9939609 given the largest weighting).⁴ This score was applied to the European
12 Prospective Investigation of Cancer (EPIC) Norfolk cohort, where the top 1.2% of people (risk score >12) had
13 an average BMI of 1.46 kg/m² greater than those in the bottom 1.4% (risk score <4). However, the majority
14 of participants had risk scores in the middle of the range (6-10), with large variability across the whole range
15 of scores (**Figure 2A**). In the highest genetic risk groups (genetic scores of 11, 12, and >12), the likelihood
16 of null effect was at least 80% (**Table 1**). The likelihood of null effect in the most common genetic score (score
17 of 8, 18.4% of participants) was 88.1%. This suggests that regardless of an individual's genetic score, there
18 is less than a 20% chance that they will display any increase in BMI due to their score relative to those 1.4%
19 of individuals with the lowest genetic risk. Across the entire range of scores, linear regression found a
20 significant association between risk score and BMI ($p<0.001$, $R_2=0.018$), suggesting that only around 2% of
21 BMI is determined by the eight SNPs most significantly associated with BMI (**Figure 2B**).

23 *MTNR1B rs10830963 (C:G) and fasting blood glucose*

24 Of the common SNPs associated with increased blood sugar, rs10830963 (C:G) has one of the largest effect
25 sizes, with each G copy associated with around a 1.3 mg/dl increase in fasting blood glucose.⁹ Data from the
26 deCODE cohort ($n=6,240$) were used as a graphical example (**Figure 3A**).⁹ Mean (SD) fasting blood glucose
27 across the three genotypes was 95.2 (12.8) mg/dl, 97.0 (12.8) mg/dl, and 97.9 (12.8) mg/dl for CC, CG, and
28 GG respectively. The likelihood of null effect was 94.4% in those with the CG genotype, and 91.6% in those
29 with the GG genotype. Linear regression found a significant association between number of G copies and

1 fasting blood glucose ($p < 0.001$, $R^2 = 0.01$), with around 1% of the variability in blood glucose being determined
2 by MTNR1B rs10830963 genotype (**Figure 3B**).

3 4 *Genetic type 2 diabetes risk score*

5 Similar to the approach of Willer *et al.*, Dupuis *et al.* published a genetic risk score for elevated fasting blood
6 glucose and risk of type 2 diabetes,³ including MTNR1B and 15 other loci. This score was applied to the
7 Framingham cohort, where the top 3.1% of people (risk score > 22) had an average fasting blood glucose ~ 6
8 mg/dl greater than those in the bottom 4.2% (risk score < 13). Similar to the obesity risk score, significant
9 heterogeneity in blood glucose levels was seen across the range of scores (**Figure 4A**). The likelihood of
10 null effect in the most common genetic score (score of 18, 14.3% of participants) was 84.5% (**Table 2**). In
11 those with the highest genetic risk score (scores 21, 22, and > 22), the risk of prediabetic level blood glucose
12 (> 100 mg/dl) was double that of those in the lowest risk group. However, even in these groups the likelihood
13 of a given genetic score being associated with blood sugar outside of the distribution of those in the lowest
14 risk group was only 25.5-27.7%, suggesting that fewer than 30% of people with the highest genetic risk of
15 prediabetes experience that risk as a disease phenotype. Across the entire range of scores, linear regression
16 found a significant association between risk score and fasting glucose ($p < 0.001$, $R^2 = 0.049$), suggesting that
17 around 5% of fasting glucose is determined by the 16 SNPs most significantly associated with type 2 diabetes
18 risk (**Figure 4B**). By comparison to the Framingham cohort, where mean (SD) fasting blood glucose was
19 92.5 (8.7) mg/dl in the lowest genetic risk group, free living hunter gathers from Tukisenta and Kitava
20 reportedly have fasting blood glucose of around 75 (8) and 65 (14) mg/dl, respectively (**Figure 4C**).^{10,11} Based
21 on these data, the Tukisentans would have a 98.6% likelihood of having a blood sugar below the mean of
22 those in the Framingham cohort with the lowest genetic risk score, with a 97.5% likelihood in the Kitavans,
23 and normal distributions that display only 19.5% and 27.3% overlap with the lowest risk Framingham group.
24 This translates to a 0.09% and 0.05% risk of prediabetic fasting blood glucose, respectively. Therefore, even
25 in the lowest risk genetic group in the Framingham cohort, the relative risk of prediabetic fasting blood sugar
26 levels (19.4%) is around 200-400 times higher than in hunter gatherer populations.

27 28 *MTHFR rs1801131 (A:C) and rs1801133 (C:T) and homocysteine*

29 Two common polymorphisms in the MTHFR gene, which alter *in vitro* enzyme activity and are associated
30 with reduced capacity to produce 5-methyltetrahydrofolate, are frequently discussed in the popular and

1 alternative health fields with regard to the methyl cycle and associated changes in detoxification, cellular
2 repair, and detoxification pathways. In 1998, van der Put *et al.* described *in vitro* MTHFR activity of the most
3 common combinations of alleles at rs1801131 and rs1801133, as well as homocysteine levels in the same
4 participants.¹² In the most common genotypes, excluding 1298AA/677TT, which account for around 88% of
5 the population on average, MTHFR function across five genotypes varies from 100% to 47.7% (**Table 3**).
6 However, even in those with 47.7% function (1298AC/677CT) there is an 82.1% chance of null effect
7 compared to 1298AA/677CC “wild type” with 100% function (**Table 3**). Across these common mutations,
8 MTHFR function only explains around 1% of the variability in homocysteine levels ($p < 0.001$, $R_2 = 0.01$; **Figure**
9 **5A**). The addition of 1298AA/677TT, which has around 12% prevalence in the population and is associated
10 with a 75.2% loss of MTHFR function, increases the explanation of variance to 7% (**Figure 5B**); however,
11 the synthetic dataset included 6.9% negative values due to the large SD in this population. This suggests
12 significant heterogeneity of homocysteine in those with the 677TT/1298AA genotype, which is not normally
13 distributed. Indeed, though the percent chance of non-significant difference in homocysteine levels compared
14 to 1298AA/677CC was only 35% in those with 1298AA/677TT, this includes a large proportion of the
15 distribution in homocysteine levels that would be below that of the “wild type” due to the very large SD in the
16 1298AA/677TT group; 31.3% would be predicted to have homocysteine levels below the mean of
17 1298AA/677CC.

18 19 **DISCUSSION**

20 The increasing prevalence of DTC genetic analyses is resulting in more and more healthcare providers being
21 asked to interpret SNP-based disease risk by their patients, or attempting to incorporate these analyses into
22 personalised treatment approaches. Here we demonstrate that, by using simple statistical theory and
23 synthetic datasets generated based on published population phenotypic data from well-characterised SNPs,
24 the likelihood of any given genotype resulting in a meaningful difference in phenotype is relatively small. For
25 individual common SNPs determined to have large effect sizes, such as *FTO* rs9939609 on BMI and
26 *MTNR1B* rs10830963 on fasting glucose, even those with two alleles have a less than 10% chance of
27 displaying a difference in phenotype due to significant population variability. Additionally, baseline disease
28 risks suggest that the vast majority of health outcomes associated with common SNPs are dominated by the
29 environment.

1 The best-characterised SNP associated with risk of overweight and obesity is FTO rs9939609, with
2 an average per A allele increase in BMI of 0.3 kg/m².⁷ However, an average population effect is less useful
3 to an individual than the likelihood that they are going to be affected in the first place. For a single FTO A
4 allele, this likelihood is around 3%, increasing to 7% in individuals with two A alleles, with 0.4% of overall BMI
5 explained by FTO genotype. Though it may be the SNP most well associated with increases in BMI, the vast
6 majority of individuals are unlikely to have their BMI meaningfully affected by their FTO SNPs. Importantly,
7 even this negligible effect of FTO on BMI is largely dominated by the environment, with recent analyses
8 suggesting that FTO rs9939609 genotype was not associated with BMI in those born before 1942.¹³ Similarly,
9 analyses of both FTO rs9939609 SNPs and composite obesity genetic risk scores suggest that those who
10 partake in regular movement or exercise (~1h of moderate-vigorous physical activity per day) have similar
11 BMIs regardless of genetics.^{14,15} In the well-characterised EPIC Norfolk cohort, the risk of being overweight
12 was above 50% regardless of a genetic score consisting of the eight SNPs most tightly-associated with BMI,
13 again suggesting a significant environmental component. Considering that current Centers for Disease
14 Control and Prevention data suggests that 39.8% of the adult population in the United States is obese,¹⁶ the
15 degree to which common genetic SNPs contribute to BMI may be statistically significant but borderline
16 physiologically irrelevant compared to the impact of the environment.

17 Similar results to those seen with genetic obesity risk were found when analysing genetic risk of
18 elevated fasting blood glucose and type 2 diabetes. Of the SNPs associated with increased fasting blood
19 glucose, MTNR1B SNP rs10830963 (C:G) has one of the largest effect sizes, with each G copy associated
20 with around a 1.3 mg/dl increase in fasting blood glucose.⁹ In our analysis, only 5.6% of individuals with a
21 single G copy would be expected to experience an increase in fasting blood sugar relative to those with the
22 CC genotype, increasing to 8.2% in homozygotes. Using the genetic risk score developed by Dupuis *et al.* is
23 more predictive, with more than a doubling of risk of prediabetes in those with the highest genetic risk score
24 compared to those with the lowest genetic risk. However, linear regression analysis suggested that only
25 around 5% of fasting blood glucose is determined by genetic risk. This is just very similar to the proportion of
26 explained variance that Dupuis *et al.* state in their original manuscript,³ which provides some support for the
27 use of synthetic datasets when variance and absolute numbers are not provided in the published literature.
28 More importantly, however, it the way in which this information is placed into the context of the consumer
29 using DTC genetic analysis to assess disease risk. For instance, the variance in fasting blood glucose (~5%)
30 attributed to the loci included in the genetic risk score is smaller than the variance in reproducibility of

1 commonly-used hand held at home glucometers used to monitor blood glucose in individuals with diabetes.
2 Any effect of genetic risk is also largely a reflection of a slight amplification of the risk associated with the
3 Western environment. Compared to hunter gatherer populations,^{5,10,11} fasting glucose is around 25-30 mg/dl
4 higher even in the lowest genetic risk group, and the risk of prediabetes is 200-400 times higher. Indeed, in
5 a recent analysis of the Bolivian Tsimane, prevalence of type 2 diabetes was 0%,¹⁷ on top of which any
6 increase in genetic risk would be essentially meaningless. Therefore, the presence of any prediabetes
7 appears to simply be a reflection of disease risk in the US as a whole, where more than 80% are thought to
8 have suboptimal metabolic health, including more than 50% with fasting glucose >100 mg/dl.¹⁸ Based on
9 multiple lines of evidence, close to 100% of the disease risk associated with elevated fasting blood glucose
10 in the Western world can be attributed to the modern environment.

11 The concept of methylation capacity and its association with long-term health has recently gained a
12 lot of interest in the alternative health community and popular press. As a result, DTC testing of common
13 SNPs in the MTHFR and other related genes is being used to estimate an individual's capacity to (re)generate
14 methylfolate in order to guide disease risk or nutrient supplementation. One potential biomarker of methyl
15 cycle function, including MTHFR activity, is homocysteine, which is associated with and increased risk of
16 cardiovascular disease, dementia, and all-cause mortality when elevated.¹⁹⁻²¹ Though there are multiple
17 pathways for the metabolism of homocysteine, one is dependent on methylfolate, and homocysteine levels
18 are often used as a proxy for the status of the folate cycle.²² Importantly, SNPs resulting in decreased *in vitro*
19 MTHFR function are common. The "wild type" genotype 677CC/1298AA associated with 100% MTHFR
20 function is only found in around 15% of the population,²³ which makes some degree of reduced MTHFR
21 function a more representative "normal" state. In addition to this, the degree of MTHFR function appears to
22 be only loosely associated with homocysteine levels. For instance, only 1% of homocysteine was accounted
23 for by the five rs1801131 (A1298C) and rs1801133 (C677T) combinations that encompass 47.7-100% mean
24 MTHFR activity. This suggests significant redundancy in the system that is unlikely to be able to inform any
25 interventions based solely on genotype. Additionally, homocysteine levels are more likely to be determined
26 by factors not associated with direct enzyme function, as those with the 1298AC/677CC genotype have
27 higher MTHFR activity than 1298AA/677CT (83.2% versus 66.8% relative enzyme function), but also had
28 higher mean homocysteine levels (13.6 $\mu\text{mol/L}$ versus 12.8 $\mu\text{mol/L}$).¹² The non-linearity of the association
29 between MTHFR and homocysteine levels is typified by the 1298AA/677TT genotype, who have around 75%
30 loss of enzyme function and 50% higher mean homocysteine levels but, importantly, display a high degree

1 of variability and values that do not appear to be normally distributed. Therefore, any specific
2 recommendations to this group must be based in phenotypic measurements, including individual
3 homocysteine levels and nutrient status. Indeed, though MTHFR is associated with the folate cycle, ensuring
4 adequate B6 and B12 may be at least as important with respect to homocysteine levels.²⁴ Homocysteine in
5 677TT carriers can also be significantly reduced with a small amount of supplemental riboflavin.²⁵ This again
6 suggests that phenotypic measurements and ensuring adequate environmental/nutrient status has a much
7 greater impact than does knowledge of genotype. However, it must be cautioned that, as yet, reducing
8 homocysteine with nutritional supplements has not yet been shown to result robustly improve health
9 outcomes, though there may be a small reduction in stroke risk.²⁶

10 This study does have some limitations. The approach used relies on the use of both simulated and
11 statistically-ideal normal distributions based on published descriptive data rather than the data itself.
12 However, where the methods could be tested against known data, such as the degree to which the glucose
13 risk score explains glucose variability, the results were very similar to the original analyses. Importantly, if
14 this approach fails to accurately recreate datasets similar to those in the published literature, then it is likely
15 that those datasets were not normally-distributed and the original analyses were therefore inappropriate. This
16 is probably the case for homocysteine levels in individuals with the MTHFR 1298AA/677TT genotype based
17 on the widely-cited study by van der Put *et al.*¹² Though all the SNPs analysed here have low penetrance,
18 they were specifically chosen because they are well-characterized in multiple populations and commonly
19 included in third party DTC analyses of consumer genetic data. Though we have only highlighted a few SNPs,
20 the techniques applied here could be used by any practitioner or interested individual to better understand
21 their disease or outcome risk based on common genetic SNPs.

22 Even though there is inherent error in our approach, it is clear that using population means to
23 determine genetic risk and make recommendations based on genetics, as is very common in the DTC market,
24 is likely to be highly-flawed due to inherent phenotypic variability. This includes variability in risk based on
25 common factors such as socioeconomic status and ethnicity. For instance, FTO genotypes are associated
26 with increased BMI in Caucasians, but not in those of African origin.⁷ For the risk of both obesity and
27 prediabetes or type 2 diabetes, particularly, the effect of the environment (diet, exercise, nutrient status) is
28 likely to dominate the phenotype such that knowing about an individual's SNPs associated with risk will have
29 little benefit. A focus on genetic risk may indeed be detrimental due to the fact that i) thinking that you have
30 a risk SNP can have an effect on physiology regardless of whether you have that SNP,²⁷ ii) the majority of

1 people have average genetic risk for a given phenotype, iii) DTC genetics testing still includes significant
2 variability and error,²⁸ iv) there is little to no evidence that specific interventions for a given common SNP
3 have any effect on health outcomes, v) communicating genetic risk does not appear to alter health
4 behaviours,²⁹ and vi) though statistically significant, the final effect of most SNPs on phenotype could often
5 be considered physiologically irrelevant. These risks have generally been acknowledged by the scientific
6 community performing genetic research, but the over-interpretation of risk by third-parties relying on
7 published population averages remains a significant worry, likely due to misinterpretation of the nature of the
8 data.

10 **CONCLUSIONS**

11 Using simple statistical techniques, either with Python code or freely-available online tools, we have
12 outlined a method by which healthcare providers and third-party genetic analysis tools can more accurately
13 analyse genetic disease risk. Importantly, it is worth noting that the widely-characterized and cited SNPs for
14 obesity, type 2 diabetes, and methylation status appear to have negligible overall effects on phenotype
15 compared to the dominant effect of the environment.

17 **ACKNOWLEDGEMENTS**

18 T.R.W is supported by start-up funds from the University of Washington Department of Pediatrics.

20 **AUTHOR CONTRIBUTIONS**

21 T.R.W developed the concept, performed the statistical analyses, and drafted the manuscript. N.O performed
22 the number generation, assisted with the statistical analyses, and edited the manuscript. Both authors
23 approved the final version.

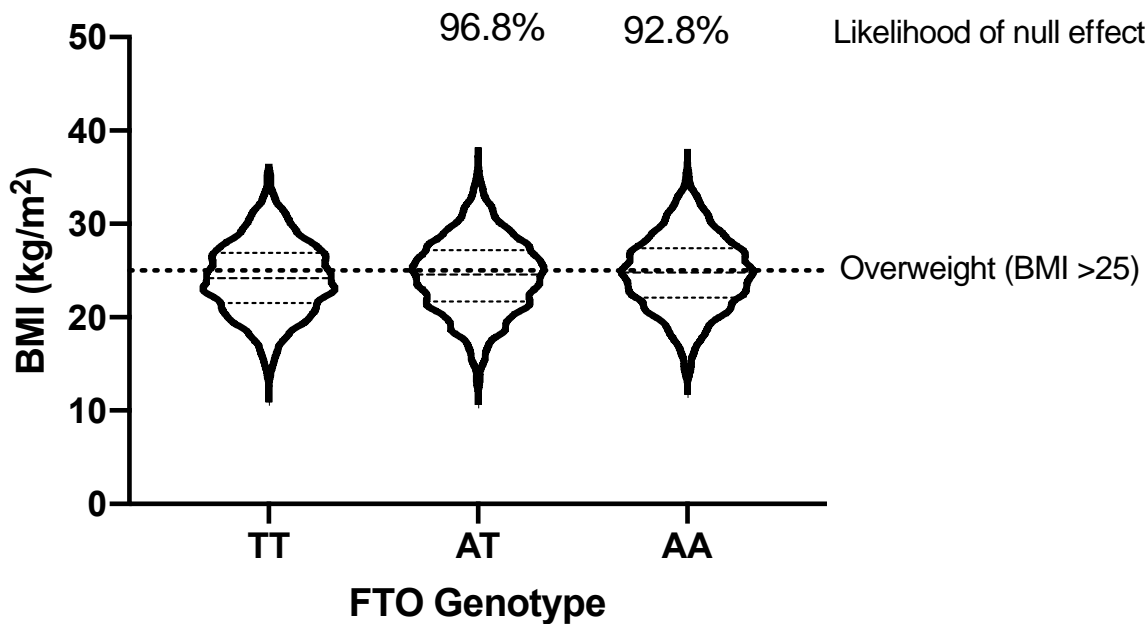
25 **COMPETING INTERESTS**

26 T.R.W and N.O declare that they have no competing interests.

28 **DATA AVAILABILITY**

29 All data was randomly-generated based on published descriptive measures from population studies.

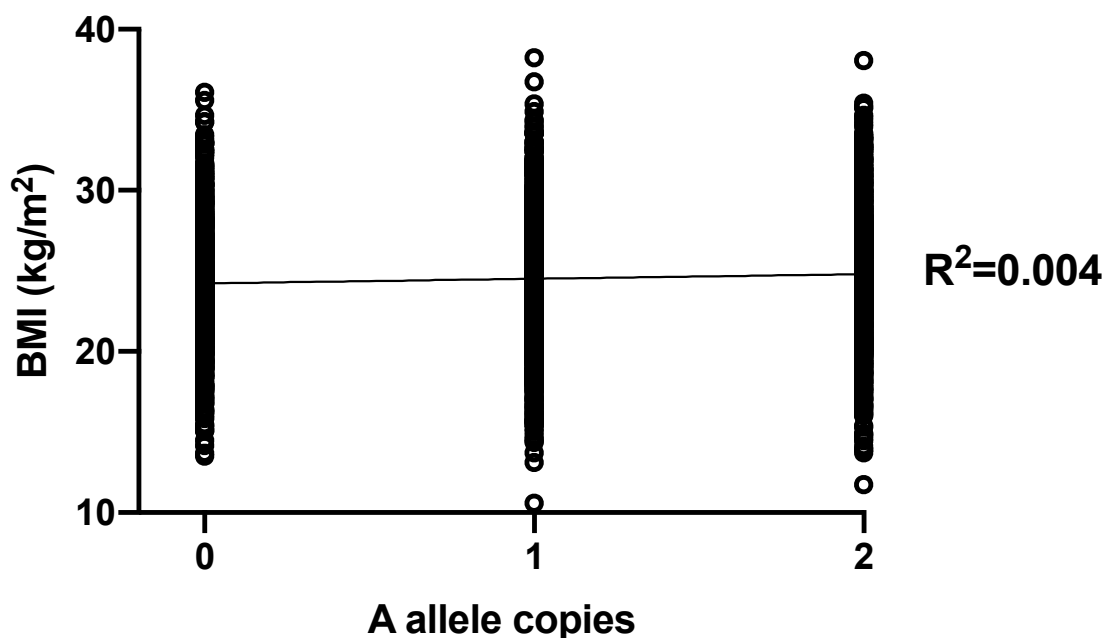
1 FIGURES AND FIGURE LEGENDS



2

3 **Figure 1A. Effect of FTO rs9939609 genotype on BMI in the NFBC cohort.** Violin plot displaying 1,000
4 synthetic BMI datapoints per FTO rs9939609 genotype, based on published population mean and SD values
5 from the NFBC cohort.⁸ Percent overlap between the AT and AA normal distributions with that of the “wild
6 type” (TT) genotype are displayed as a measure of the likelihood of these risk genotypes having no overall
7 effect on BMI.

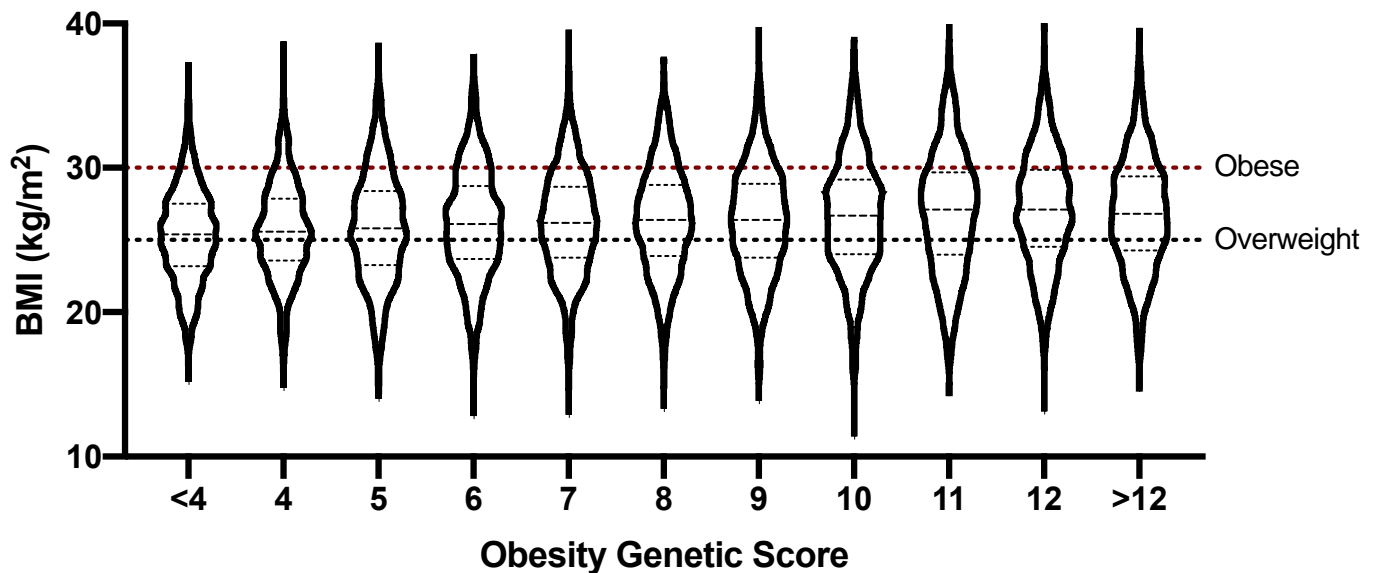
8



9

10 **Figure 1B. Linear regression of FTO rs9939609 A alleles versus BMI.** Linear regression of 1,000 synthetic
11 BMI datapoints per FTO rs9939609 A allele copy. There was a significant association between number of A
12 copies and BMI ($p=0.001$, $R_2=0.0035$), suggesting that only around 0.4% of the variability in BMI is
13 determined by FTO genotype.

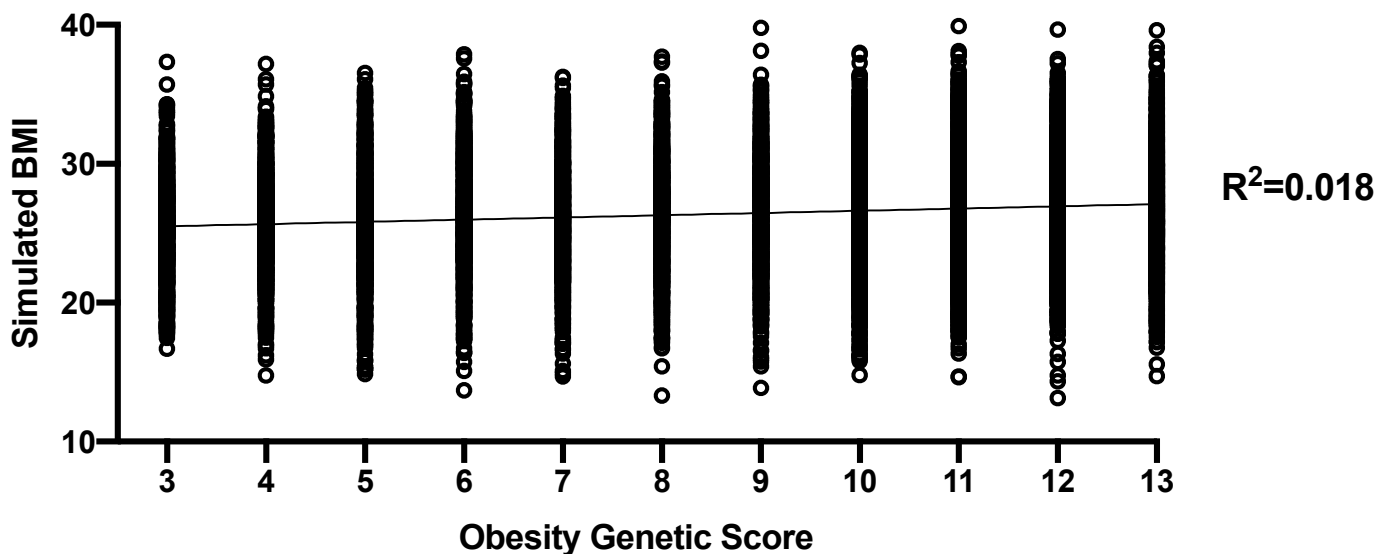
14



1

2 **Figure 2A. Effect of BMI genetic score on BMI in the EPIC Norfolk cohort.** Violin plot displaying 1,000
3 synthetic BMI datapoints per group of BMI genetic risk score, as developed by Willer *et al.*,⁴ using population
4 mean and SD values from the EPIC Norfolk cohort. Significant variability is seen across the entire range of
5 genetic scores, with more than 50% of individuals being overweight regardless of genotype.

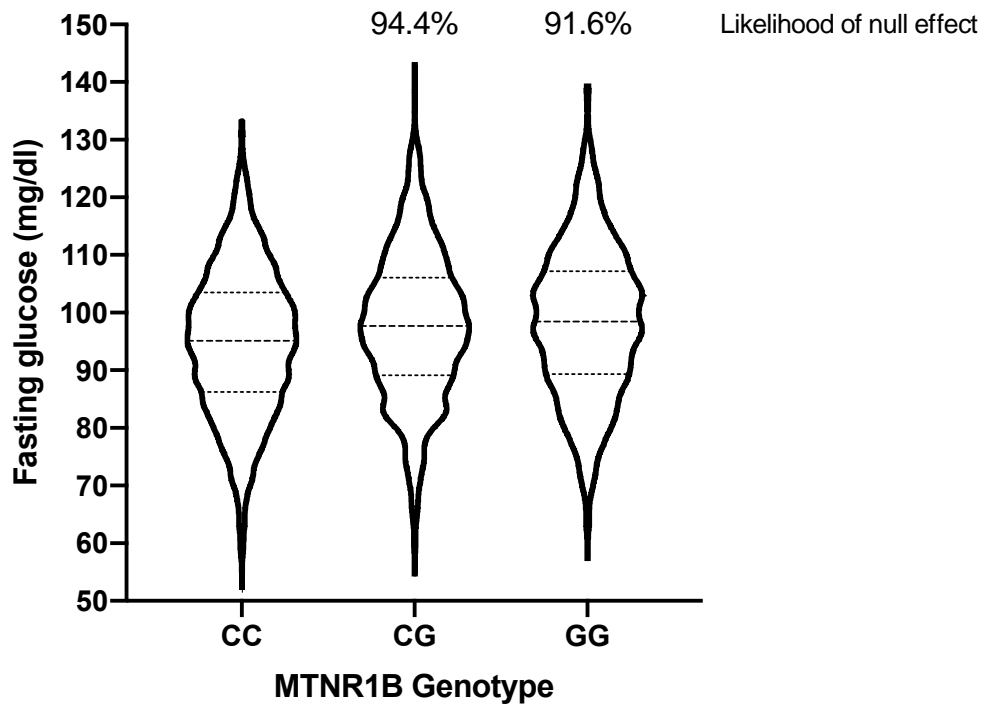
6



7

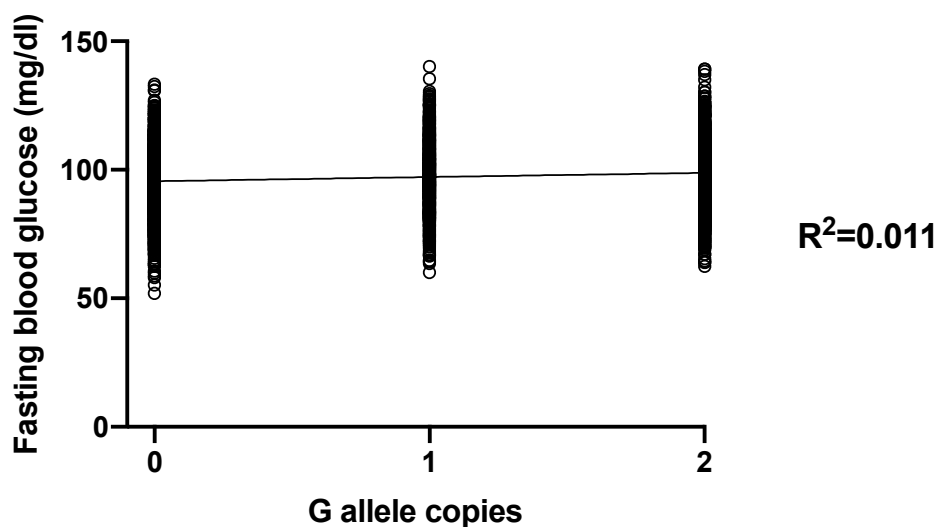
8 **Figure 2B. Linear regression of genetic BMI risk score versus BMI.** Linear regression of 1,000 synthetic
9 BMI datapoints per group of BMI genetic risk score, as developed by Willer *et al.*,⁴ There was a significant
10 association between risk score and BMI ($p<0.001$, $R^2=0.018$), with around 2% of BMI determined by the eight
11 SNPs most significantly associated with BMI.

12



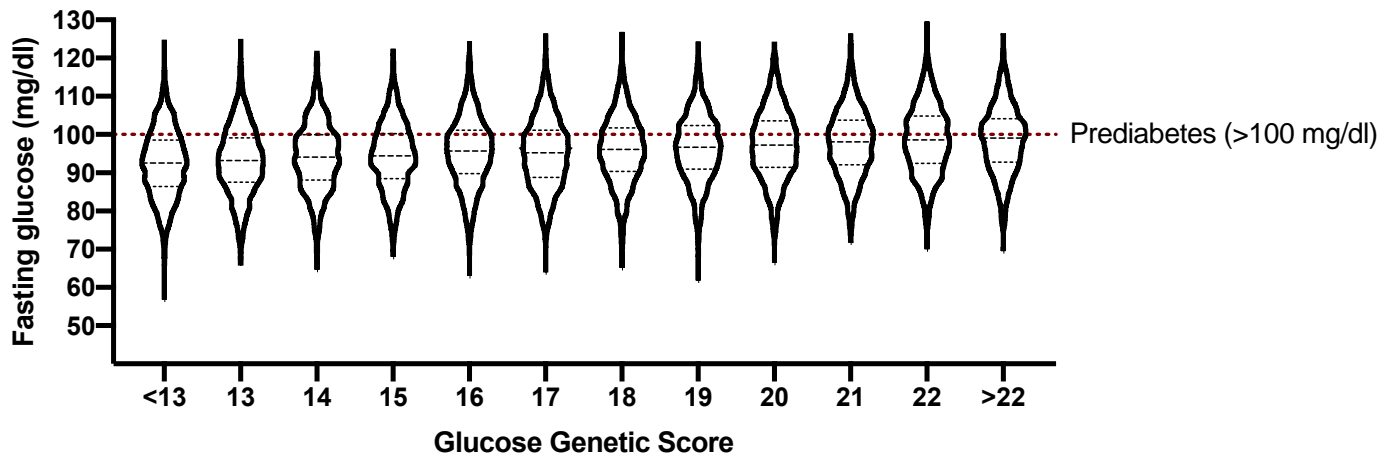
1

2 **Figure 3A. Effect of MTNR1B rs10830963 genotype on fasting glucose in the deCODE cohort.** Violin
3 plot displaying 1,000 synthetic glucose datapoints per MTNR1B rs9939609 genotype, based on published
4 population mean and SD values from the deCODE cohort.⁹ Percent overlap between the CG and GG normal
5 distributions with that of the “wild type” (CC) genotype are displayed as a measure of the likelihood of these
6 risk genotypes having no overall effect on fasting blood glucose.



7

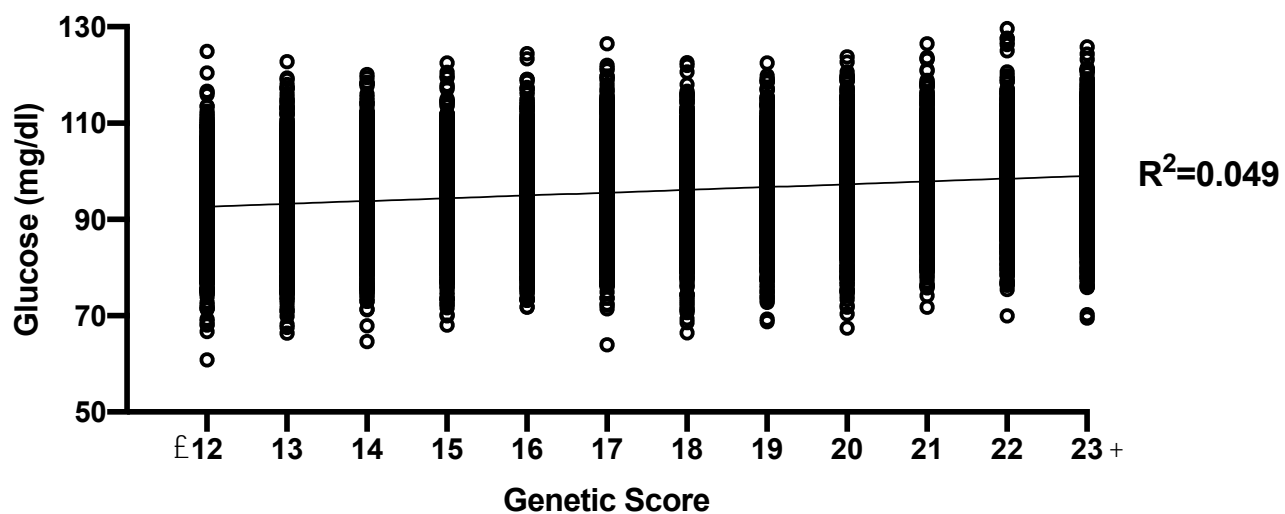
8 **Figure 3B. Linear regression of MTNR1B rs10830963 genotype versus fasting glucose.** Linear
9 regression of 1,000 synthetic BMI datapoints per MTNR1B rs9939609 G allele copy. There was a significant
10 association between number of G copies and fasting glucose ($p<0.001$, $R_2=0.011$), suggesting that only
11 around 1% of the variability in fasting blood glucose is determined by MTNR1B genotype.



1

2 **Figure 4A. Effect of glucose genetic score on fasting glucose in the Framingham cohort.** Violin plot
3 displaying 1,000 synthetic glucose datapoints per group of glucose genetic risk score, as developed by
4 Dupuis *et al.*,³ using population mean and SD values from the Framingham cohort. Significant variability is
5 seen across the entire range of genetic scores. Risk of prediabetes (fasting glucose >100 mg/dl) increases
6 from 19.4% to 43.5% from the lowest to highest risk group, with the most common genetic risk profiles (scores
7 of 16-19, ~50% of individuals) having around a 30% risk of prediabetes.

8



9

10 **Figure 4B. Linear regression of genetic glucose risk score versus fasting glucose.** Linear regression
11 of 1,000 synthetic BMI datapoints per group of glucose genetic risk score, as developed by Dupuis *et al.*³
12 There was a significant association between risk score and fasting blood glucose ($p<0.001$, $R^2=0.049$), with
13 around 5% of fasting glucose variability determined by the 16 SNPs most significantly associated with glucose
14 homeostasis.

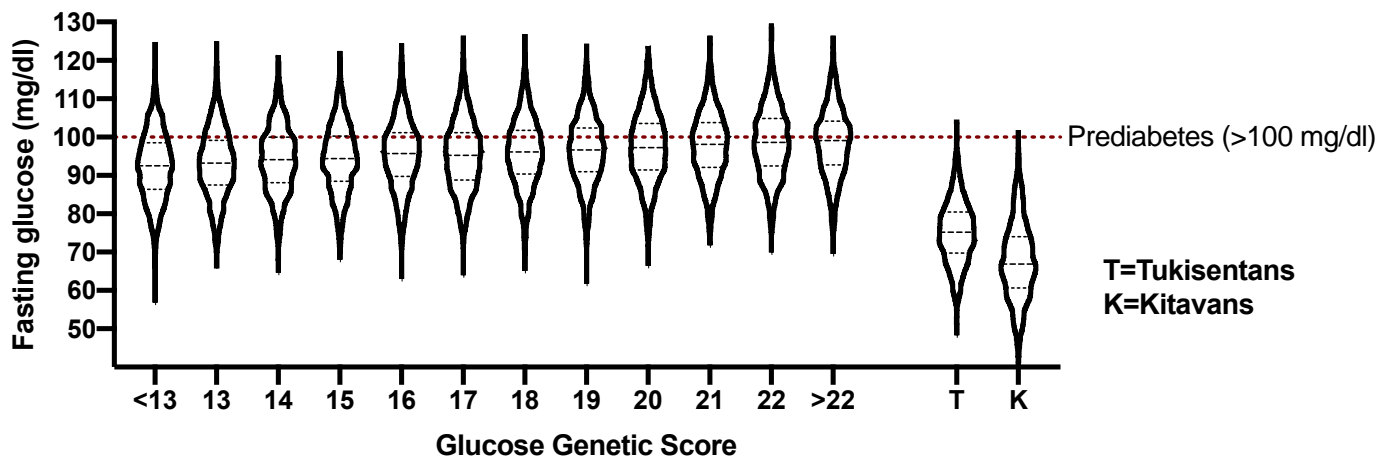


Figure 4C. Comparison between fasting glucose in the Framingham cohort and in hunter gatherer populations. Violin plot displaying 1,000 synthetic glucose datapoints per group of glucose genetic risk score using population mean and SD values from the Framingham cohort,³ as well as using data from two hunter gatherer cohorts, the Tukisentans and Kitavans.^{10,11} The Tukisentans would have a 98.6% likelihood of having a blood sugar below the mean of those in the Framingham cohort with the best genetic score, with a 97.5% likelihood in the Kitavans. They display normal distributions with only 19.5% and 27.3% overlap with the lowest risk Framingham group and 0.09% and 0.05% risk of prediabetic fasting blood glucose, respectively. Even in the lowest risk genetic group in the Framingham cohort, the estimated prevalence of prediabetic fasting blood sugar levels (19.4%) is around 200-400 times higher than in hunter gatherer populations.

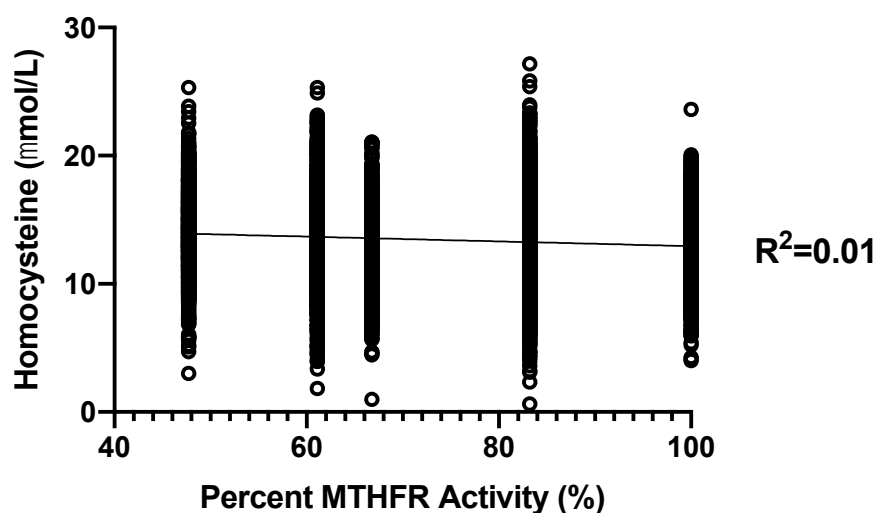
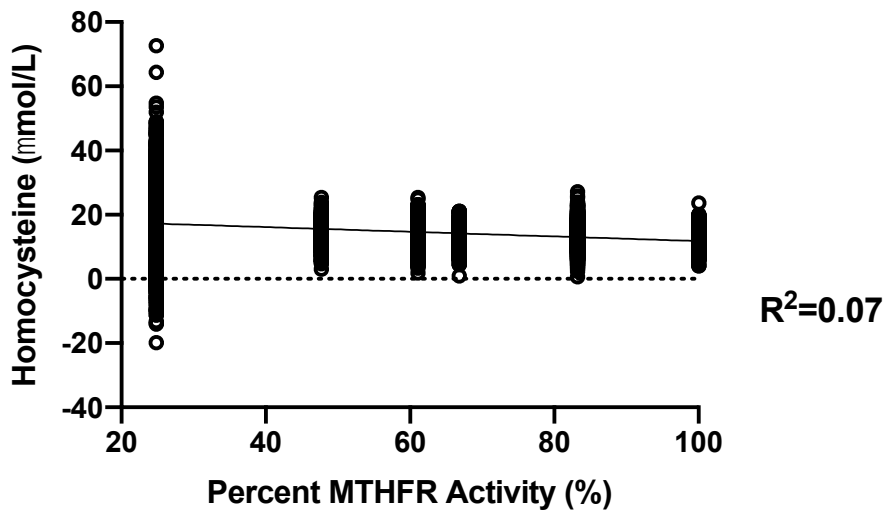


Figure 5A. Linear regression of MTHFR activity versus homocysteine for the most common genotypes. Linear regression of 1,000 synthetic homocysteine datapoints per combination of common rs1801131 (A1298C) and rs1801133 (C677T) SNPs by *in vitro* MTHFR activity, excluding 1298AA/677TT. There was a significant association between MTHFR function and homocysteine ($p < 0.001$, $R^2 = 0.01$), suggesting that only around 1% of the variability in homocysteine is determined by MTHFR activity across these genotypes.



1

2

3

4

5

6

7

8

9

Figure 5B. Linear regression of MTHFR activity versus homocysteine including 1298AA/677TT. Linear regression of 1,000 synthetic homocysteine datapoints per combination of rs1801131 (A1298C) and rs1801133 (C677T) SNPs by *in vitro* MTHFR activity. There was a significant association between MTHFR function and homocysteine ($p < 0.001$, $R^2 = 0.07$); however, the large SD (66% of the mean) in those with 1298AA/677TT resulted in 6.9% of predicted homocysteine levels being negative. This suggests that homocysteine in those with 1298AA/677TT is highly-variable, non-normally distributed, and that the effects of MTHFR activity on homocysteine levels are non-linear.

1 **TABLES**

2 **Table 1. Effect of BMI genetic score on risk of overweight and obesity.** BMI genetic risk score, as
 3 developed by Willer *et al.*,⁴ and risk of being overweight or obese, using population mean and SD values
 4 from the EPIC Norfolk cohort. Genetic scores of 6-10 cover around 75% of the population. The likelihood of
 5 null effect of each score was determined as the percent overlap of its normal distribution with that of the
 6 lowest risk group (score <4). Even in the highest risk groups (11, 12, <12) percent overlap was at least 80%,
 7 with only 12-17% of those with a genetic score of 6-10 predicted to have BMI affected by their genotype.

Genetic BMI Score	Prevalence (%)	Mean (SD) BMI	Overweight (%)	Obese (%)	Distribution overlap with score <4 (%)
<4	1.4	25.4 (3.1)	55.1	7.1	
4	3.4	25.7 (3.4)	58.2	10.2	95.3
5	7.2	25.9 (3.8)	59.3	14.2	89.1
6	12.9	26.2 (3.7)	62.6	15.6	88.2
7	17.4	26.2 (3.6)	63.1	14.4	89.3
8	18.1	26.3 (3.7)	63.9	15.5	88.1
9	15.8	26.5 (3.7)	65.7	17.3	85.9
10	10.6	26.6 (3.9)	66.0	19.1	83.4
11	7.7	26.8 (4.2)	66.7	22.1	80.9
12	2.8	27.0 (4.0)	69.2	22.6	80.0
>12	1.2	26.8 (3.8)	68.1	20.2	81.7

8

9 **Table 2. Effect of glucose genetic score on risk of prediabetes.** Glucose genetic risk score, as developed
 10 by Dupuis *et al.*,³ and risk of having prediabetes, using population mean and SD values from the Framingham
 11 cohort. Genetic scores of 16-19 cover around 52% of the population, and have around 30% prevalence of
 12 prediabetes. The likelihood of null effect of each score was determined as the percent overlap of its normal
 13 distribution with that of the lowest risk group (score <13). In those with the highest genetic risk scores (21,
 14 22, and >22), the risk of prediabetic blood glucose (>100mg/dl) levels was double that of the lowest risk
 15 group. However, even in these groups the likelihood of a given genetic score being associated with blood
 16 sugar outside of the distribution of those in the lowest risk group was only 25.5-27.7%, suggesting that fewer
 17 than 30% of people with the highest genetic risk of prediabetes experience that risk as a disease phenotype.

Genetic Glucose Score	Prevalence (%)	Mean (SD) fasting glucose (mg/dl)	Prediabetes (%)	Distribution overlap with score <13 (%)
<13	4.2	92.5 (8.7)	19.4	
13	5.0	93.6 (8.8)	23.4	94.8
14	8.2	94.2 (8.6)	25.0	92.3
15	9.8	94.3 (8.8)	25.9	91.7
16	13.0	95.2 (8.9)	29.5	87.7
17	13.8	95.4 (8.7)	29.9	86.7
18	14.3	95.9 (8.9)	32.3	84.5
19	11.5	96.5 (8.7)	34.4	81.9
20	8.5	97.3 (8.8)	38.0	78.3
21	5.4	98.1 (8.6)	41.3	74.5
22	3.2	98.6 (8.7)	43.6	72.7
>22	3.1	98.6 (8.6)	43.5	72.3

18

1 **Table 3. Effect of common MTHFR SNPs on average function and homocysteine.** Average MTHFR
 2 function and homocysteine levels by rs1801131 (A1298C) and rs1801133 (C677T) SNP combination, as
 3 published by van der Put *et al.*¹² Genotypes are listed in order of *in vitro* MTHFR enzyme function, with
 4 estimated population prevalence taken from Brown *et al.*²³ The likelihood of null effect of each combination
 5 of MTHFR SNPs was determined as the percent overlap of its normal distribution with that with the “wild type”
 6 genotype (1298AA/677CC). Significant non-linearity between degree of MTHFR function and homocysteine
 7 is seen, with the most common genotypes, representing close to 88% of the population and 47.7%-83.2%
 8 enzyme function displaying 81-95% likelihood of null effect on resulting homocysteine levels.

Genotype		Estimated	Relative	Homocysteine	Homocysteine Distribution
1298	677	Prevalence (%)	Mean (%)	Mean (SD)	overlap with 1298AA/677CC (%)
AA	CC	15.3	100	12.9 (2.8)	
AC	CC	20.8	83.2	13.6 (4.0)	81.5
AA	CT	22.8	66.8	12.8 (3.1)	94.9
CC	CC	8.8	61.1	13.9 (3.9)	81.0
AC	CT	19.8	47.7	14.2 (3.1)	82.1
AA	TT	12.2	24.8	19 (2.5)	35.0

9

10

1 REFERENCES

- 2 1 Guerrini, C. J., Wagner, J. K., Nelson, S. C., Javitt, G. H. & McGuire, A. L. Who's on third?
3 Regulation of third-party genetic interpretation services. *Genet Med*, doi:10.1038/s41436-019-0627-
4 6 (2019).
- 5 2 Bush, W. S. & Moore, J. H. Chapter 11: Genome-wide association studies. *PLoS Comput Biol* **8**,
6 e1002822, doi:10.1371/journal.pcbi.1002822 (2012).
- 7 3 Dupuis, J. *et al.* New genetic loci implicated in fasting glucose homeostasis and their impact on type
8 2 diabetes risk. *Nature genetics* **42**, 105-116, doi:10.1038/ng.520 (2010).
- 9 4 Willer, C. J. *et al.* Six new loci associated with body mass index highlight a neuronal influence on
10 body weight regulation. *Nature genetics* **41**, 25-34, doi:10.1038/ng.287 (2009).
- 11 5 Pontzer, H., Wood, B. M. & Raichlen, D. A. Hunter-gatherers as models in public health. *Obes Rev*
12 **19 Suppl 1**, 24-35, doi:10.1111/obr.12785 (2018).
- 13 6 Johnston, D. *Random Number Generators—Principles and Practices: A Guide for Engineers and*
14 *Programmers* (De|G Press, 2018).
- 15 7 Qi, Q. *et al.* FTO genetic variants, dietary intake and body mass index: insights from 177,330
16 individuals. *Hum Mol Genet* **23**, 6961-6972, doi:10.1093/hmg/ddu411 (2014).
- 17 8 Kaakinen, M. *et al.* Life-course analysis of a fat mass and obesity-associated (FTO) gene variant
18 and body mass index in the Northern Finland Birth Cohort 1966 using structural equation modeling.
19 *Am J Epidemiol* **172**, 653-665, doi:10.1093/aje/kwq178 (2010).
- 20 9 Prokopenko, I. *et al.* Variants in MTNR1B influence fasting glucose levels. *Nature genetics* **41**, 77-
21 81, doi:10.1038/ng.290 (2009).
- 22 10 Lindeberg, S., Eliasson, M., Lindahl, B. & Ahren, B. Low serum insulin in traditional Pacific
23 Islanders--the Kitava Study. *Metabolism: clinical and experimental* **48**, 1216-1219,
24 doi:10.1016/s0026-0495(99)90258-5 (1999).
- 25 11 Sinnett, P. F. & Whyte, H. M. Epidemiological studies in a total highland population, Tuisenta, New
26 Guinea. Cardiovascular disease and relevant clinical, electrocardiographic, radiological and
27 biochemical findings. *J Chronic Dis* **26**, 265-290, doi:10.1016/0021-9681(73)90031-3 (1973).
- 28 12 van der Put, N. M. *et al.* A second common mutation in the methylenetetrahydrofolate reductase
29 gene: an additional risk factor for neural-tube defects? *Am J Hum Genet* **62**, 1044-1051,
30 doi:10.1086/301825 (1998).
- 31 13 Rosenquist, J. N. *et al.* Cohort of birth modifies the association between FTO genotype and BMI.
32 *Proc Natl Acad Sci U S A* **112**, 354-359, doi:10.1073/pnas.1411893111 (2015).
- 33 14 Vimalaswaran, K. S. *et al.* Physical activity attenuates the body mass index-increasing influence of
34 genetic variation in the FTO gene. *The American journal of clinical nutrition* **90**, 425-428,
35 doi:10.3945/ajcn.2009.27652 (2009).
- 36 15 Li, S. *et al.* Physical activity attenuates the genetic predisposition to obesity in 20,000 men and
37 women from EPIC-Norfolk prospective population study. *PLoS Med* **7**,
38 doi:10.1371/journal.pmed.1000332 (2010).
- 39 16 CentersforDiseaseControl. *Adult Obesity Facts*, <<https://www.cdc.gov/obesity/data/adult.html>> (
40 17 Kaplan, H. *et al.* Coronary atherosclerosis in indigenous South American Tsimane: a cross-sectional
41 cohort study. *Lancet (London, England)* **389**, 1730-1739, doi:10.1016/S0140-6736(17)30752-3
42 (2017).
- 43 18 Araujo, J., Cai, J. & Stevens, J. Prevalence of Optimal Metabolic Health in American Adults:
44 National Health and Nutrition Examination Survey 2009-2016. *Metab Syndr Relat Disord* **17**, 46-52,
45 doi:10.1089/met.2018.0105 (2019).
- 46 19 Fan, R., Zhang, A. & Zhong, F. Association between Homocysteine Levels and All-cause Mortality:
47 A Dose-Response Meta-Analysis of Prospective Studies. *Scientific reports* **7**, 4769,
48 doi:10.1038/s41598-017-05205-3 (2017).
- 49 20 Ganguly, P. & Alam, S. F. Role of homocysteine in the development of cardiovascular disease.
50 *Nutrition journal* **14**, 6, doi:10.1186/1475-2891-14-6 (2015).
- 51 21 Smith, A. D. *et al.* Homocysteine and Dementia: An International Consensus Statement. *J*
52 *Alzheimers Dis* **62**, 561-570, doi:10.3233/JAD-171042 (2018).
- 53 22 Blom, H. J. & Smulders, Y. Overview of homocysteine and folate metabolism. With special
54 references to cardiovascular disease and neural tube defects. *J Inherit Metab Dis* **34**, 75-81,
55 doi:10.1007/s10545-010-9177-4 (2011).
- 56 23 Brown, N. M. *et al.* Detection of 677CT/1298AC "double variant" chromosomes: implications for
57 interpretation of MTHFR genotyping results. *Genet Med* **7**, 278-282,
58 doi:10.109701.GIM.0000159904.92850.D5 (2005).

- 1 24 Moll, S. & Varga, E. A. Homocysteine and MTHFR Mutations. *Circulation* **132**, e6-9,
2 doi:10.1161/CIRCULATIONAHA.114.013311 (2015).
- 3 25 McNulty, H. *et al.* Riboflavin lowers homocysteine in individuals homozygous for the MTHFR 677C-
4 >T polymorphism. *Circulation* **113**, 74-80, doi:10.1161/CIRCULATIONAHA.105.580332 (2006).
- 5 26 Marti-Carvajal, A. J., Sola, I., Lathyris, D. & Dayer, M. Homocysteine-lowering interventions for
6 preventing cardiovascular events. *The Cochrane database of systematic reviews* **8**, CD006612,
7 doi:10.1002/14651858.CD006612.pub5 (2017).
- 8 27 Turnwald, B. P. *et al.* Learning one's genetic risk changes physiology independent of actual genetic
9 risk. *Nat Hum Behav* **3**, 48-56, doi:10.1038/s41562-018-0483-4 (2019).
- 10 28 Tandy-Connor, S. *et al.* False-positive results released by direct-to-consumer genetic tests highlight
11 the importance of clinical confirmation testing for appropriate patient care. *Genet Med* **20**, 1515-
12 1521, doi:10.1038/gim.2018.38 (2018).
- 13 29 Hollands, G. J. *et al.* The impact of communicating genetic risks of disease on risk-reducing health
14 behaviour: systematic review with meta-analysis. *BMJ* **352**, i1102, doi:10.1136/bmj.i1102 (2016).
- 15