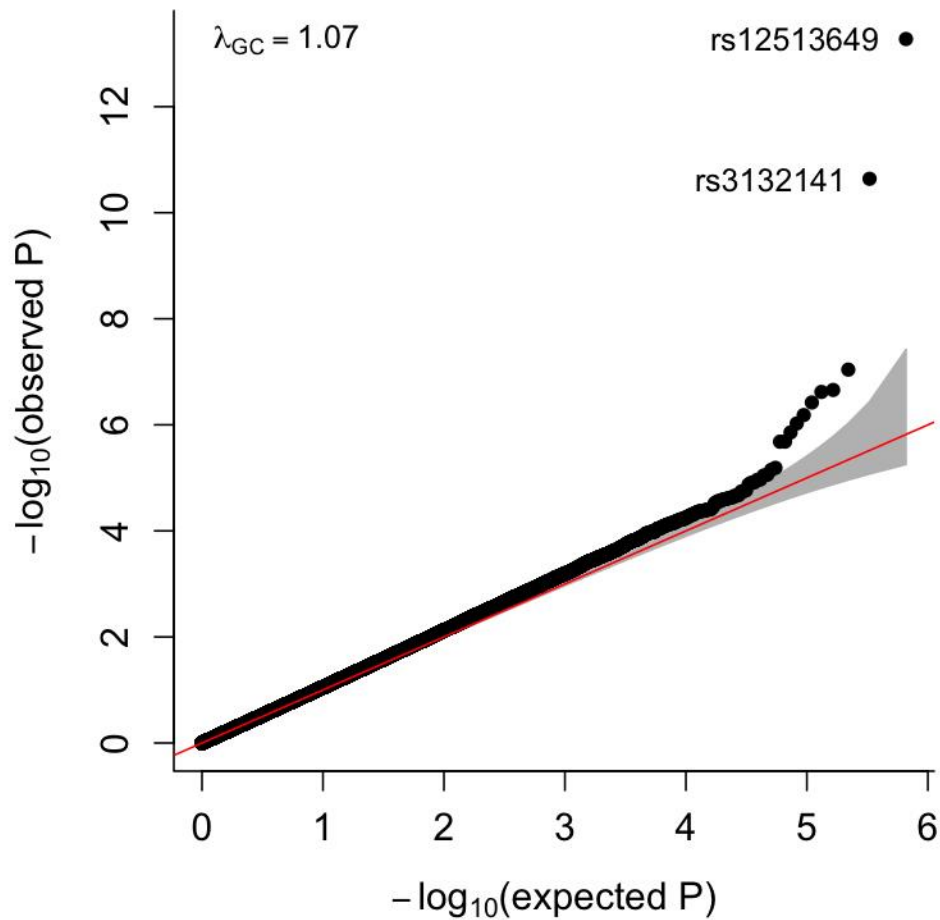**Supplementary Figure 1**
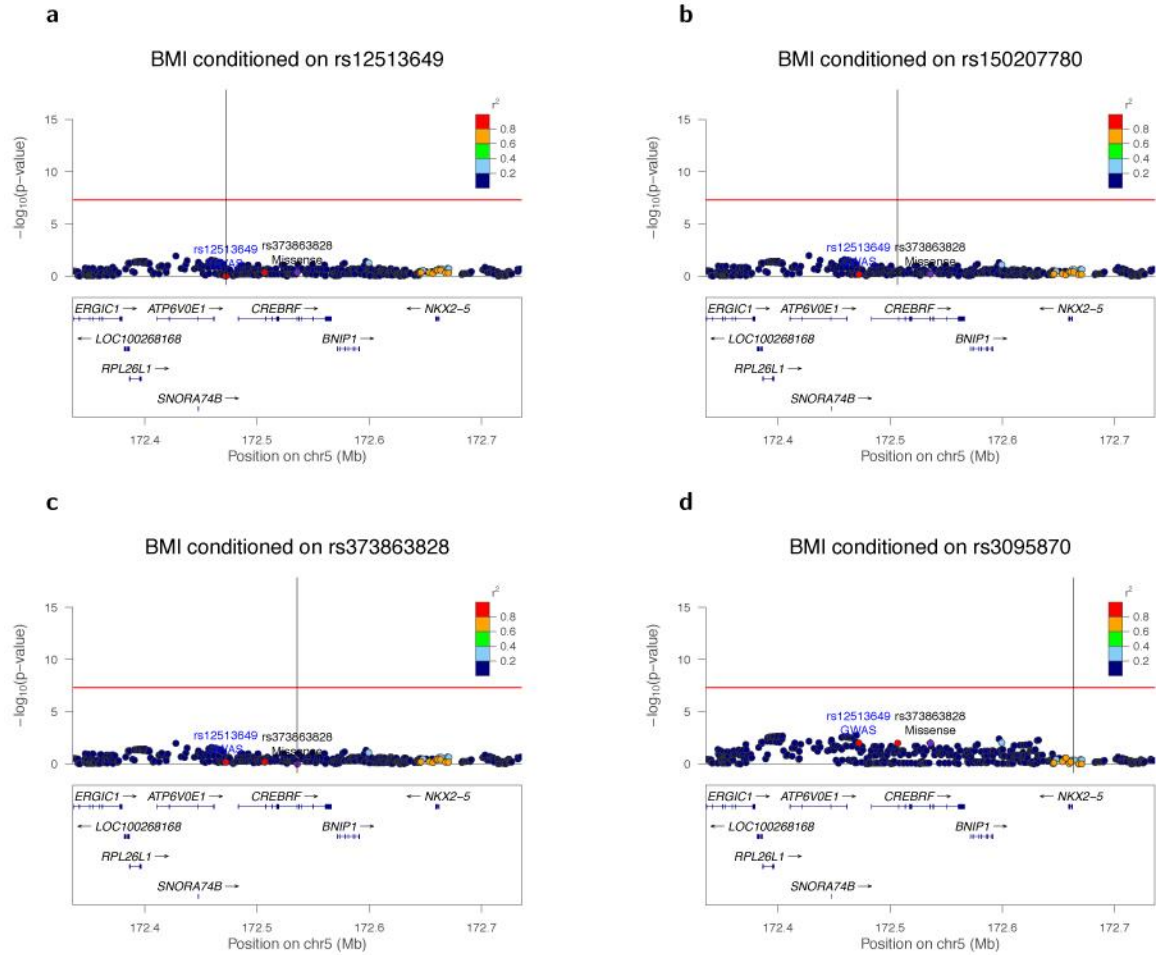
**Principal-components analyses.**

(**a**) Scatterplot of the first three principal components from the principal-components analysis of the Samoan and HapMap phase 3 populations. Continental population abbreviations: SAM, Samoans ($n$ = 250); EUR, Europeans ($n$ = 253); AFR, Africans ($n$ = 511); EAS, East Asians ($n$ = 255); SAS, South Asians ($n$ = 88); AMR, admixed Americans ($n$ = 77). **Supplementary Video 1** shows a rotating animation of this figure. (**b**) Scatterplots of the first six principal components from the principal-components analysis of the Samoans alone ($n$ = 3,094) plotted against each other.

**Supplementary Figure 2**

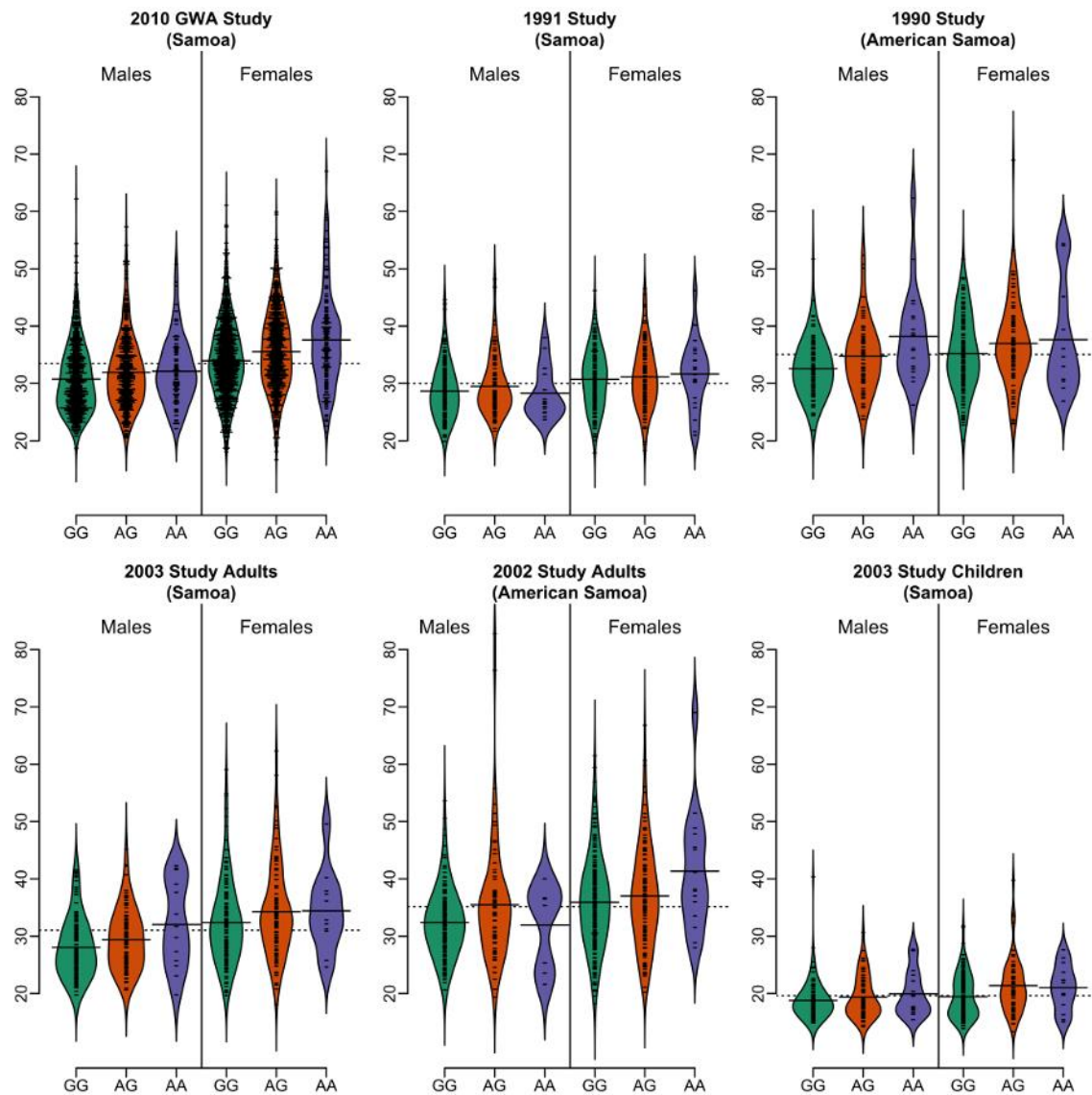**Quantile–quantile plot for the BMI GWAS.**

A quantile–quantile (QQ) plot of the observed –$\log_{10}$ ($P$ values) from **Figure 1a** for association of BMI in the discovery sample versus – $\log_{10}$ ($P$ values) as expected under no association. The second most significant variant, rs3132141, lies between *BNIP1* and *NKX2-5* and is 184.5 kb from the most significant variant, rs12513649. $n$ = 3,072 Samoans.

**Supplementary Figure 3**

**Conditional associations of targeted sequencing genotypes with BMI.**
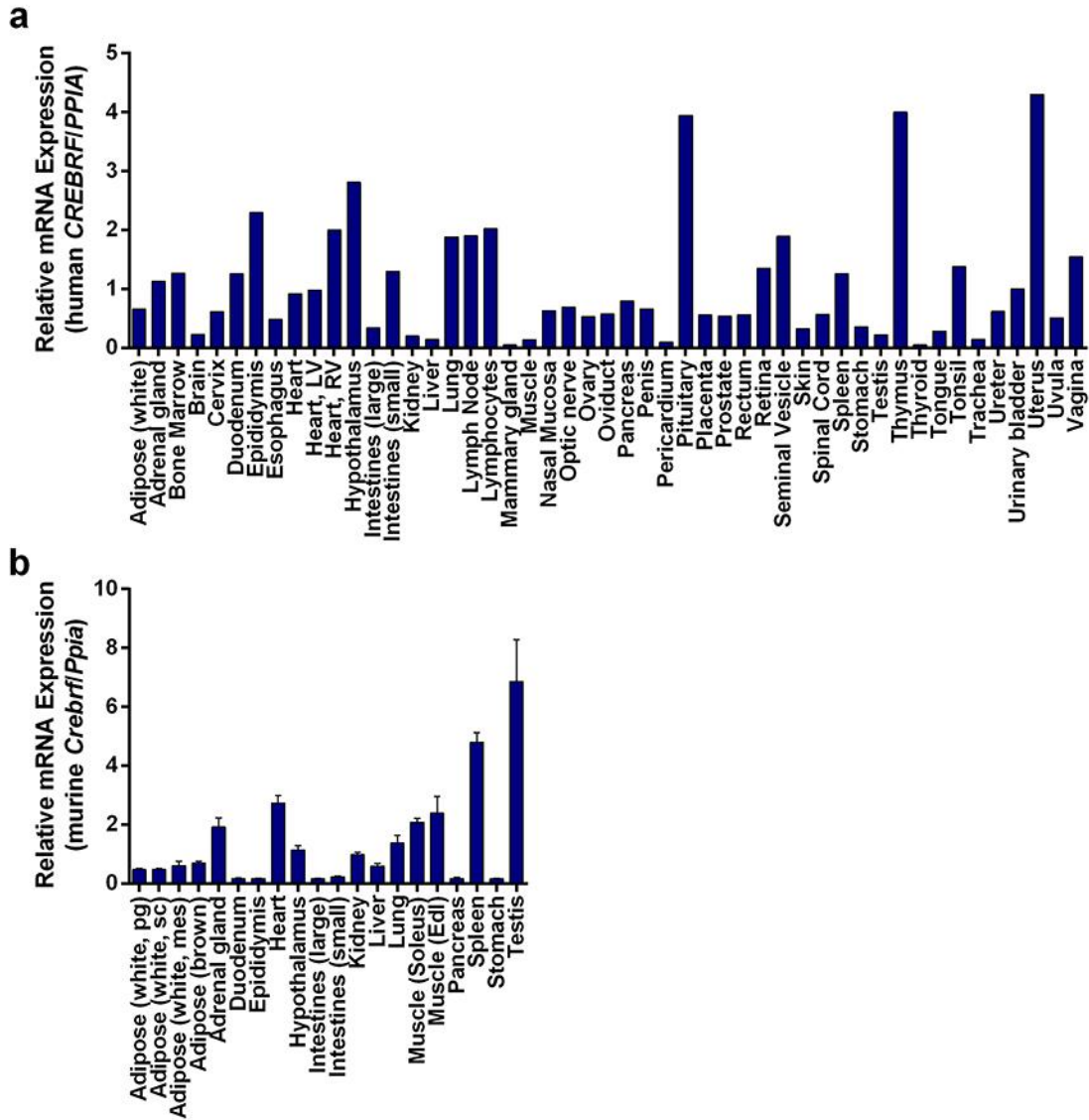
(**a–d**) Associations between SNPs in the targeted sequencing regions and BMI conditioned on rs12513649 (**a**), rs150207780 (**b**), rs373863828 (**c**), and rs3095870 (**d**). The red line in each plot corresponds to a $P$ value of $5 \times 10^{-8}$. $n = 3,072$ Samoans.

**Supplementary Figure 4**

**Beanplots of BMI in GWAS and replication samples stratified by missense variant rs373863828 genotype, sex, and nation.**
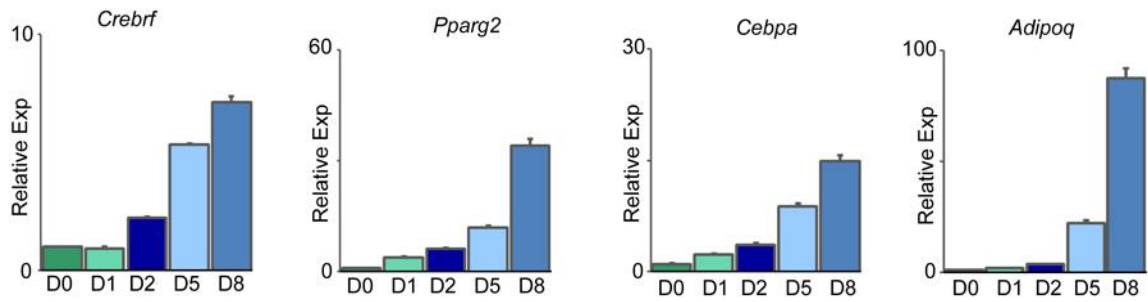
Each bean consists of a mirrored density curve containing a one-dimensional scatterplot of the individual data. The heavy dark line shows the average within each group, and the dotted line indicates the overall average. Plots were drawn using the R beanplot package[33]. Sample sizes are as indicated in **Supplementary Table 1**.

**Supplementary Figure 5**
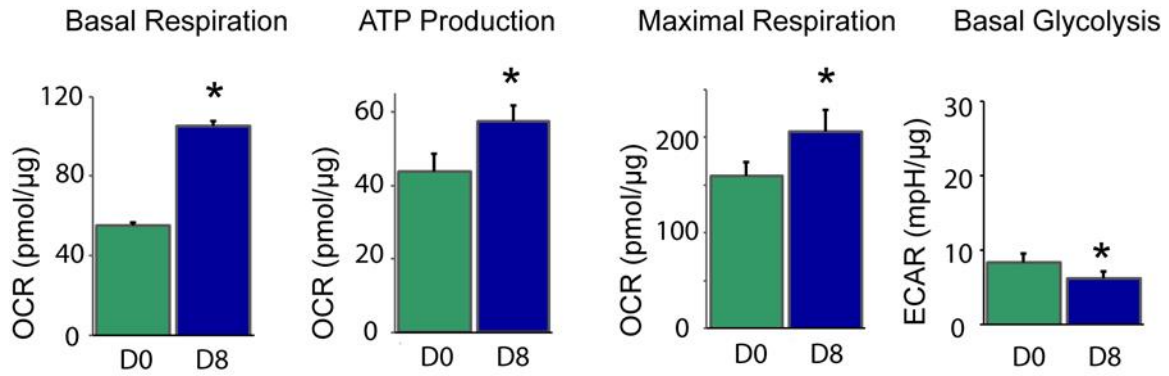
**Expression of *CREBRF* in human and mouse tissues.**

(**a**) Human *CREBRF* mRNA expression was determined in multiple human tissues using Human cDNA Arrays from Origene (*n* = 1/tissue; nutritional status not known). (**b**) Mouse *Crebrf* mRNA expression was determined in mouse tissues obtained from 10-week-old, littermate-matched, *ad libitum*–fed, male C56BL/6J mice (*n* = 6/group). Expression was normalized to the endogenous control gene peptidylprolyl isomerase A/cyclophilin A (*PPIA* for human; *Ppia* for mouse). Values represent relative expression and are expressed as means plus s.e.m. No statistical comparisons were performed. pg, perigonadal; sc, inguinal subcutaneous; mes, mesenteric. These data support the presence/absence of CREBRF in specific tissues but should be used with caution when assessing relative expression, particularly in humans where precise conditions at the time of tissue collection are not known. Gene expression can be compared to additional *in silico* resources including the BGTEx and BioGPS portals (see URLs).

**Supplementary Figure 6**

**Expression of mouse *Crebrf* relative to key adipogenic genes during adipocyte differentiation.**

3T3-L1 cells were treated with a hormonal differentiation cocktail at 2 d after confluence (day 0, D0), and RNA samples were collected at the indicated time points. mRNA expression relative to the β-actin (*Actb*) reference gene was determined using quantitative RT–PCR, with day 0 expression values set at 1. Values are given as means $\pm$ s.e.m. ($n$ = 8). A representative of five independent experiments is shown.

**Supplementary Figure 7**

**Bioenergetic profile changes during adipocyte differentiation.**

3T3-L1 cells were treated with a hormonal differentiation cocktail at 2 d after confluence (day 0, D0), and key bioenergetic variables were determined on the basis of oxygen consumption rate (OCR) and extracellular acidification rate (ECAR) measurements normalized to protein content. Values are given as means $\pm$ s.e.m. ($n = 6$). *$P < 0.01$ compared to day 0 (two-tailed $t$ test with unequal variances). As the results were consistent with previously published data[24,25], the experiment was performed once.

**a**

**b**

**Supplementary Figure 8**

**iHS and nS$_L$ scores in an 800-kb region centered on the missense variant rs373863828 ($n$ = 626 non-closely related Samoans).**

(**a**) iHS scores versus physical position. (**b**) nS$_L$ scores versus physical position. In both **a** and **b**, the blue dot indicates the score at the missense variant rs373863828 and the yellow dot indicates the score at the discovery variant rs12513649; the dotted horizontal line indicates the score at the missense variant rs373863828.

Supplementary Note for

# A thrifty variant in *CREBRF* strongly influences body mass index in Samoans

Ryan L. Minster, Nicola L. Hawley, Chi-Ting Su, Guangyun Sun, Erin E. Kershaw, Hong Cheng, Olive D. Buhule, Jerome Li, Muagututi'a Sefuiva Reupena, Satupa'itea Viali, John Tuitele, Take Naseri, Zsolt Urban, Ranjan Deka, Daniel E. Weeks, Stephen T. McGarvey

## 1) Participants

The discovery sample of 3,072 phenotyped and genotyped individuals was drawn from 3,475 men and women ($n$ = 1,437 men), ages 24.5 to <65 years who reported Samoan ancestry (based on having four Samoan grandparents). Recruitment took place between February and July 2010 in 33 villages from the islands of 'Upolu and Savai'i of Samoa[1]. A population-based design was employed and consenting participants completed interviews targeting lifestyle factors related to cardiometabolic health (health history, socio-economic position, dietary intake, and physical activity) and anthropometric measurements (height, weight, blood pressure, body composition), and gave fasting whole blood samples for biochemical and genetic assays. A description of the prevalence of non-communicable diseases and associated risk factors is provided in Hawley *et al.*[1]

The replication sample consists of individuals from two samples of Samoans studied in 1990–95 and in 2002–03. The 1990–95 study sample derives from a longitudinal study of adiposity and cardiovascular disease risk factors among Samoan adults from American Samoa and Samoa. Although there is substantial economic disparity between the two polities, the Samoans from both territories form a single socio-cultural unit with frequent exchange of mates. Genetically they represent a single homogenous population[2,3]. Participants were between 25–55 years of age at baseline and reported that they were of Samoan ancestry. Detailed descriptions of the sampling and recruitment were reported previously[4-6]. Briefly, participants were recruited from 46 villages and worksites in American Samoa in 1990 and nine villages in (then Western) Samoa in 1991. All participants were free of self-reported history of heart disease, hypertension, or diabetes at baseline. There were 413 and 607 genotyped and phenotyped individuals available from American Samoa in 1990 and from Samoa in 1991, respectively (**Supplementary Table 1**). Due to lack of genome-wide marker data on these samples, we were unable to infer relatedness, and so these were treated as unrelated in the analyses.

The 2002–03 family study sample includes adults and children recruited as part of an extended family-based genetic linkage analysis of cardiometabolic traits[7-11]. Probands and relatives were unselected for obesity or related phenotypes. The recruitment process and criteria used for inclusion in this study are described in detail previously[7,9]. All individuals self-reported as having Samoan ancestry[10]. There were 590 adults, 18–89 years, from 2002 in American Samoa; and 493 adults, 19–82 years, and 409 children ages 5–<18 years, from 2003 in Samoa, available with genotypes and phenotypes (**Supplementary Table 1**). The analyses of these samples were adjusted for relatedness using kinships derived from the known family structures (which had been verified to be consistent with relatedness estimates derived using genome-wide microsatellite markers)[10].

## 2) Genotyping

DNA was extracted from whole blood as previously reported[1,7]. In the discovery sample, genotyping was attempted on 3,298 DNA samples (including 3,194 participants, 34 duplicates and 70 positive controls) across 909,622 probes using a Genome-Wide Human SNP Array 6.0 (Affymetrix, California, USA). Genotyping of the discovery samples was performed on 96 well plates, each plate containing two reference samples: 1) REF103 provided by Affymetrix, and 2) a Coriell DNA sample, NA15510, and a negative control. A duplicate sample from the same plate was introduced in each plate with blinded IDs for the laboratory personnel. The samples were not randomized and were processed in the order collected in the field. Laboratory personnel were blind to the sample phenotypes.

Extensive quality control was conducted based on a pipeline developed by Laurie et al.[12] including assessment of probe and sample quality (probes and samples excluded with missingness rates > 5%), sex validation, investigation of genotyping batch effects, assessment of cryptic relatedness and population substructure, and duplicate sample and duplicate probe discordance. Of the 3,194 samples attempted for genotyping, 4 were dropped due to high genotyping missingness, 3 due to discrepancy between reported and apparent genetic gender, 7 due to apparent sex chromosome aneuploidy, 9 due to chromosomal abnormalities such as deletions and duplications, 2 due to apparent sample admixture, and 50 due to poor cluster resolution across the genome. An additional 25 participants were excluded due to self-reported pregnancy. After quality control, 3,119 samples genotyped for 894,139 unique autosomal and X-linked markers were available to conduct genome-wide association studies. Since much of the analysis software used does not properly handle twins, we also excluded 3 more individuals: the least genotyped member of three monozygotic twin pairs. There were 19 participants missing BMI. Complete phenotype and genotype data were available for up to 3,072 participants.

To test for possible overlap between the samples from our three studies, we used 116 single-nucleotide polymorphisms (SNPs) genotyped in common across all our samples. These 116 SNPs, including rs12513649, were chosen based on their association signals for a whole suite of traits in the discovery sample. At loci with multiple significant SNPs, the peak SNP was chosen as representative of that locus. At loci (defined as 1 Mbp windows) with different peak SNPs for different phenotypes, the SNP with the smallest $P$ value among the associated phenotypes was genotyped as representative of that locus. These SNPs spanned all autosomal chromosomes and the X chromosome, and were at least 1 Mb away from each other and not in linkage disequilibrium with each other ($r^2$<0.3 for all but one pair of adjacent markers; $r^2$=0.73 between rs4932738 and rs7252689 on chromosome 19). Genotyping of these variants in the replication sample was performed using custom-designed TaqMan OpenArray Real-Time PCR assays (Applied Biosystems). SNPs that could not be genotyped using OpenArray assays, including rs12513649, were genotyped individually using TaqMan SNP Genotyping assays (Applied Biosystems). For replication genotyping, in each 384 well plate ($n$ = 8), 4 duplicates from the same plate with blinded ID were included; each plate also contained 8 negative controls and 8 Coriell samples (NA15510). The quality of genotype clustering for each SNP was verified and corrected manually.

2

### 3) Statistical Analysis

BMI was log-transformed to approximate normality. Residuals were generated by linear regression against age, age$^2$, sex and the interactions between age and sex. We tested for association between autosomal marker genotypes and the BMI residuals while using the empirical kinship matrix to adjust for population substructure and subject relatedness. Note that population substructure is accounted for in our analyses by inclusion of the empirical kinship model in the analysis models, because, as Hofmann[13] states "explicitly modeling the pairwise relatedness between all individuals captures both population structure and kinship". The tests were conducted using a score test as implemented in the mmscore function in GenABEL[14]. The statistics between X chromosome genotypes and BMI residuals were calculated in GenABEL without adjusting for the empirical kinship estimates. Following analysis, 230,554 SNPs with a minor allele frequency < 0.01 (including 23,612 monomorphic SNPs) and then 4,093 SNPs with HWE test $P$<0.00001 were filtered out, resulting in 659,492 autosomal and X-linked SNPs used for analyses. Inflation due to population stratification and cryptic relatedness was assessed by estimating $\lambda_{GC}$ using the lower 90% of the $P$ value distribution[15].

The replication sample was divided into three groups for separate analysis: the 1990–1995 study participants, the 2002–03 family study adults (age≥18), and the 2002–03 family study children (age<18). For the purposes of the meta-analyses, we did not further subdivide the studies by nation; doing so would have broken up pedigrees in the family study that span both nations. For consistency, we therefore did not subdivide the 1990–1995 study by nation either. All samples, including those from the discovery sample, were examined using the 116 SNPs typed in common across our samples for genetic identity that might have arisen through recruitment into multiple studies over the two decades that they span. One sample of each pair that had an estimated identity-by-descent > 0.9 as estimated in PLINK[16] were removed from analysis. In total, we removed 72 samples, preferentially from the 1990–1995 study. For the participants, both adults and children, from the 2002–03 family study, kinship coefficients were calculated from the recorded pedigrees using the 'kinship2' package[17] in R[18]. Replication association analyses were performed using GenABEL[19] for each group, using the kinship coefficients to adjust for relatedness in the family sample. We do not have sufficient marker data to infer relatedness and adjust for it in the 1990-1995 study, so they were treated as unrelated. We used the same covariates in the replication regression models that we used in the GWAS analyses, with an additional variable indicating whether subjects were from Samoa or American Samoa.

The genetic ancestry of our discovery sample, where every individual self-reported having four Samoan grandparents, was assayed via principal components analyses using PC-AiR[20]. We conducted two principal components analyses. Firstly, to examine the relationship of the Samoans against other continent populations, we compared the genotypes of a randomly chosen subset of 250 Samoans against genotypes from individuals comprising HapMap Phase 3. Genotype management was performed using PLINK[16]. HapMap Phase 3 genotypes[21] were merged with the genotypes from the Samoan discovery sample. SNPs with a minor allele frequency < 0.05, with a missingness rate > 0.1, and located within regions problematic for the calculation of principal components analysis (the major histocompatibility locus on 6p21, the region near *LCT* on 2q21, and common inversion regions on 8p23 and 17q21) were dropped. Markers were further pruned down to every fourth marker. The PC-AiR algorithm was applied to the remaining 111,438 markers: the PCs were estimated in the unrelated subjects as determined by the KING-robust kinship coefficient estimator[22] and extended to relatives in the

3

dataset based on their genetic similarity. The first three principal components from this analysis are shown in **Supplementary Figure 1a**. Secondly, to examine the potential for population stratification within the Samoans, we calculated principal components within the Samoan participants in our sample. SNPs were again removed based on the same minor allele frequencies, missingness rates and location within problematic regions as above. Markers were pruned based on linkage disequilibrium down to a set of independent SNPs, and the PC-AiR algorithm was applied to the remaining 72,586 markers. The first six principal components from this analysis are shown in **Supplementary Figure 1b**. Note that the between-population 'distances' shown in **Supplementary Figure 1a** should be interpreted with caution, as we did not correct for how SNPs were selected to be on the Affymetrix genotyping array[23]. Correcting for SNP ascertainment bias in a well-calibrated manner requires not only sophisticated and careful modeling of the ascertainment process but also requires sequencing data (which we do not yet have) to validate that the correction method works correctly[24].

Prior to meta-analysis, we performed quality control of the summary statistics using EasyQC[25] to check for strand and allele frequency consistency. Meta-analysis of the adult samples was performed using METAL[26] to generate two replication $P$ values: one for the adult replication samples and one for the adult replication samples and discovery sample together (**Table 1**). The $P$ value–based method was used with sample sizes as weights and the genomic control correction turned off. We assessed heterogeneity across all the cohorts by calculating both Cochran's $Q$ and the $I^2$ statistic[27-29].

## 4) Targeted Sequencing

We selected a 1.5 Mbp segment (NC_000005.09:171583933_173083933) for targeted sequencing centered on the Samoan linkage disequilibrium block containing rs12513649. Sequencing was performed on 96 discovery sample participants optimally chosen using INFOSTIP[30]. The sample size of 96 was chosen due to fiscal constraints, and was estimated to recover 94% of the information had we been able to sequence everyone. Baits were derived using SureDesign (Agilent Technologies), with additional baits derived based on blat analysis. DNA libraries were prepared using SureSelect (Agilent Technologies) and sequenced using 100 bp paired-end runs on an Illumina HiSeq 2500 with the goal that at least 95% of the targeted region achieves a coverage depth of 20× or greater. Mean bait coverage was 81×. Samples were processed using BWA, GATK3 (QD<2.0, MQ<40.0, FS>60.0, MQRankSum<-12.5, ReadPosRankSum<-8.0), and HaplotypeCaller with hard cutoffs. This resulted in 99.6% concordance to VeraCode array calls, and 98.35% of single nucleotide variants were in dbSNP 138.

## 5) Bayesian Fine-Mapping

For fine-mapping using the imputed variants, we selected 160 variants with minor allele frequency ≥ 0.05 on either side of the missense variant rs373863828. These 321 SNPs spanned from 172368674 to 172670745 on chromosome 5, including from the GWAS variant rs12513649 on the left to the variants with significant $P$ values near *NKX2-5* on the right (**Figure 1b**). We then used the PAINTOR program[31] to estimate posterior probabilities of causality for each variant in the region, based on Z scores derived from the ProbABEL estimates described above and the linkage disequilibrium correlation matrix as estimated by the R package 'snpStats'[32]. We used the default maximal number of causal variants of 2 and the default number of maximum iterations of 10. We also used PAINTOR to incorporate prior information

4

about coding and regulatory DNA regions using the genome segmentation data derived by the ENCODE project[33]. This annotation segments the genome into seven classes: 1) CTCF enriched element, 2) Predicted enhancer, 3) Predicted promoter flanking region, 4) Predicted repressed or low activity region, 5) Predicted transcribed region, 6) Predicted promoter region including transcription start site, and 7) Predicted weak enhancer or open chromatin cis regulatory region. PAINTOR was run using these segmentations in each of the six ENCODE cell lines, and then the most significant annotation (a predicted transcribed region in the HepG2 liver carcinoma cell line) was used when estimating the posterior probabilities. The 'combined' ENCODE genomic segmentation annotation was downloaded from the Ensembl Encode ftp site (see URLs).

## 6) Cell Culture and Transfection

The mouse preadipocyte cell line 3T3-L1 was obtained from ATCC. No genetic authentication has been performed. However, the phenotype of the cells is consistent with previous publications. Cells were maintained in Dulbecco's modified Eagle's medium (DMEM, Gibco) supplemented with 10% newborn calf serum (NCS, Sigma), 100 units/mL penicillin and 100 µg/mL stremptomycin (Sigma), 3.7 g/L $NaHCO_3$, 4.77 g/L HEPES in a 37°C with 5% $CO_2$ humidified incubator. 3T3-L1 preadipocytes were transfected with plasmids containing eGFP-only negative control, wild-type human *CREBRF*, or the p.Arg457Gln variant using Lipofectamine 2000 (ThemoFisher Scientific) in triplicates. Transfected cells were kept under selection with 500 µg/mL Geneticin (G418, ThemoFisher Scientific) for 3 weeks to generate stable cell lines. Mycoplasma testing was performed by PCR and DAPI staining. All cells used in this study tested negative.

## 7) Adipocyte Differentiation

The differentiation of 3T3-L1 to adipocytes was carried out as described previously[34]. Differentiation was induced 2 days post confluence with a differentiation cocktail including 3-isobutyl-1-methylxanthine (IMBX, 0.5 mM; Sigma), dexamethasone (0.25 µM; Sigma), human insulin (1 µg/mL; Sigma) in basic media with 10% fetal bovine serum (FBS). After 2 days, the media was replaced with maintenance media with 10% FBS and 1 µg/mL human insulin. After further 2 days, the maintenance media was replaced with growth media containing 10% FBS and was changed every other day for up to 10 days. Geneticin (500 µg/mL) selection was maintained throughout the differentiation protocol for stable transfected cells.

## 8) Oil Red O Plate Assay

Oil Red O staining detects intracellular triglyceride accumulation[35]. Cells were seeded in 96-well cell culture plates at 10,000 cells/well with 8 technical replicates. At endpoints of interest, cells were fixed with 4% paraformaldehyde for 15 min. To obtain a working solution, an Oil Red O stock (0.3% in isopropanol) was diluted with water 24:16 v/v. After fixation, cells were rinsed with phosphate-buffered saline (PBS) and incubated with Oil Red O working solution for 15 min (30 µL per well). The wells were washed with PBS three times. Then, 100 µL isopropanol was added in each well to elute the dye and the absorbance was measured at 560 nm. Cells containing media only served as blanks. Blank values were subtracted from experimental samples. Cells in a parallel plate were lysed using CelLytic M (Sigma) and the protein concentration was measured using the Bradford assay[36] (Bio-Rad). Absorbance data were normalized to protein concentration and expressed in $OD_{560}$/µg units.

5

### 9) Oil Red O Staining and Microscopy

To visualize lipid accumulation, cells were cultured on coverslips. Eight days after confluence the media was removed and the cells were washed twice with PBS. Fixation in 4% paraformaldehyde for 10 min at room temperature was followed by staining with Oil Red O working solution for 30 min at room temperature. The Oil Red O solution was aspirated and the cells were rinsed 6 times in distilled water. The cells were counterstained with hematoxylin for 5 minutes at room temperature followed by rinsing 6 times with distilled water. The coverslips were mounted with glycerol-gelatin media and images were captured using a DM5000 (Leica Microsystems) photomicroscope.

### 10) Triglyceride Assay

Cells were harvested 8 days after confluence and the PicoProbe Triglyceride Quantification Assay Kit (Abcam, ab178780) was used to measure the level of triglycerides in cell lysate. The triglyceride level (pmol) was normalized to the amount of protein measured by the Bradford method[36] in each lysate sample.

### 11) Quantitative RT-PCR

Total RNA was harvested using an RNeasy Mini Kit (Qiagen) and cDNA was generated using the Superscirpt III Reverse Transcriptase (ThemoFisher Scientific). Quantitative RT-PCR analysis used SYBR Green PCR Master Mix (BioRad) with specific primers (**Supplementary Table 3**). Samples were run on a QuantStudio 12 Flex Real Time PCR System (ThemoFisher Scientific). The efficiency of the qPCR assays was determined using a template dilution series and was found to be $\geq 0.9$. The results were analyzed using ExpressionSuite Software v1.0.4 either using the $\Delta\Delta$Ct method[37], or by calculating the $2^{e^{*}\Delta Ct}$ value, where $e$ is PCR efficiency and $\Delta$Ct is the threshold cycle difference between the target gene and $\beta$-actin (*Actb*) as a reference gene.

**URLs.** Ensembl Encode ftp site
ftp.ebi.ac.uk/pub/software/ensembl/encode/supplementary/integration_data_jan2011/byDataType/segmentations/jan2011/hub

### REFERENCES
1.  Hawley, N.L. *et al.* Prevalence of adiposity and associated cardiometabolic risk factors in the Samoan genome-wide association study. *Am J Hum Biol* **26**, 491-501 (2014).
2.  McGarvey, S.T. Cardiovascular disease (CVD) risk factors in Samoa and American Samoa, 1990-95. *Pac Health Dialog* **8**, 157-62 (2001).
3.  Deka, R. *et al.* Genetic characterization of American and Western Samoans. *Hum Biol* **66**, 805-22. (1994).
4.  McGarvey, S.T., Levinson, P.D., Bausserman, L., Galanis, D.J. & Hornick, C.A. Population-change in adult obesity and blood-lipids in American-Samoa from 1976-1978 to 1990. *American Journal of Human Biology* **5**, 17-30 (1993).
5.  Chin-Hong, P.V. & McGarvey, S.T. Lifestyle incongruity and adult blood pressure in Western Samoa. *Psychosom Med* **58**, 131-7 (1996).

6.    Galanis, D.J., McGarvey, S.T., Quested, C., Sio, B. & Afele-Fa'amuli, S. Dietary intake of modernizing Samoans: Implications for risk of cardiovascular disease. *Journal of the American Dietetic Association* **99**, 184-190 (1999).

7.    Dai, F. *et al.* Genome-wide scan for adiposity-related phenotypes in adults from American Samoa. *Int J Obes (Lond)* **31**, 1832-42 (2007).

8.    Åberg, K. *et al.* A genome-wide linkage scan identifies multiple chromosomal regions influencing serum lipid levels in the population on the Samoan islands. *J Lipid Res* **49**, 2169-2178 (2008).

9.    Dai, F. *et al.* A whole genome linkage scan identifies multiple chromosomal regions influencing adiposity-related traits among Samoans. *Ann Hum Genet* **72**, 780-792 (2008).

10.   Åberg, K. *et al.* Susceptibility loci for adiposity phenotypes on 8p, 9p, and 16q in American Samoa and Samoa. *Obesity (Silver Spring)* **17**, 518-24 (2009).

11.   Åberg, K. *et al.* Suggestive linkage detected for blood pressure related traits on 2q and 22q in the population on the Samoan islands. *BMC Med Genet* **10**, 107 (2009).

12.   Laurie, C.C. *et al.* Quality control and quality assurance in genotypic data for genome-wide association studies. *Genet Epidemiol* **34**, 591-602 (2010).

13.   Hoffman, G.E. Correcting for population structure and kinship using the linear mixed model: theory and extensions. *PLoS One* **8**, e75707 (2013).

14.   Chen, W.M. & Abecasis, G.R. Family-based association tests for genomewide association scans. *Am J Hum Genet* **81**, 913-26 (2007).

15.   Devlin, B. & Roeder, K. Genomic control for association studies. *Biometrics* **55**, 997-1004 (1999).

16.   Purcell, S. *et al.* PLINK: a tool set for whole-genome association and population-based linkage analyses. *American Journal of Human Genetics* **81**, 559-75 (2007).

17.   Therneau, T.M. & Sinnwell, J. kinship2: Pedigree Functions. R package. (2015).

18.   R Development Core Team. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria.  ISBN 3-900051-07-0. (2004).

19.   Aulchenko, Y.S., Ripke, S., Isaacs, A. & van Duijn, C.M. GenABEL: an R library for genome-wide association analysis. *Bioinformatics* **23**, 1294-6 (2007).

20.   Conomos, M.P., Miller, M.B. & Thornton, T.A. Robust inference of population structure for ancestry prediction and correction of stratification in the presence of relatedness. *Genet Epidemiol* **39**, 276-93 (2015).

21.   International HapMap Consortium *et al.* Integrating common and rare genetic variation in diverse human populations. *Nature* **467**, 52-8 (2010).

22.   Manichaikul, A. *et al.* Robust relationship inference in genome-wide association studies. *Bioinformatics* **26**, 2867-73 (2010).

23.   Albrechtsen, A., Nielsen, F.C. & Nielsen, R. Ascertainment biases in SNP chips affect measures of population divergence. *Mol Biol Evol* **27**, 2534-47 (2010).

24.   Wollstein, A. *et al.* Demographic history of Oceania inferred from genome-wide data. *Curr Biol* **20**, 1983-92 (2010).

25.   Winkler, T.W. *et al.* Quality control and conduct of genome-wide association meta-analyses. *Nat Protoc* **9**, 1192-212 (2014).

26.   Willer, C.J., Li, Y. & Abecasis, G.R. METAL: fast and efficient meta-analysis of genomewide association scans. *Bioinformatics* **26**, 2190-1 (2010).

27.   Cochran, W.G. The comparison of percentages in matched samples. *Biometrika* **37**, 256-66 (1950).

28.   Higgins, J.P. & Thompson, S.G. Quantifying heterogeneity in a meta-analysis. *Stat Med* **21**, 1539-58 (2002).

7

29. Higgins, J.P., Thompson, S.G., Deeks, J.J. & Altman, D.G. Measuring inconsistency in meta-analyses. *BMJ* **327**, 557-60 (2003).
30. Gusev, A. *et al.* Low-pass genome-wide sequencing and variant inference using identity-by-descent in an isolated human population. *Genetics* **190**, 679-89 (2012).
31. Kichaev, G. *et al.* Integrating functional data to prioritize causal variants in statistical fine-mapping studies. *PLoS Genet* **10**, e1004722 (2014).
32. Clayton, D. snpStats: SnpMatrix and XSnpMatrix classes and methods. R package. 1.20.0 edn Vol. 1.20.0 (2015).
33. Encode Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57-74 (2012).
34. Zebisch, K., Voigt, V., Wabitsch, M. & Brandsch, M. Protocol for effective differentiation of 3T3-L1 cells to adipocytes. *Anal Biochem* **425**, 88-90 (2012).
35. Ramirez-Zacarias, J.L., Castro-Munozledo, F. & Kuri-Harcuch, W. Quantitation of adipose conversion and triglycerides by staining intracytoplasmic lipids with Oil red O. *Histochemistry* **97**, 493-7 (1992).
36. Bradford, M.M. A rapid and sensitive method for the quantitation of microgram quantities of protein utilizing the principle of protein-dye binding. *Anal Biochem* **72**, 248-54 (1976).
37. Livak, K.J. & Schmittgen, T.D. Analysis of relative gene expression data using real-time quantitative PCR and the $2^{-\Delta\Delta C}T$ Method. *Methods* **25**, 402-8 (2001).
38. Cole, T.J., Bellizzi, M.C., Flegal, K.M. & Dietz, W.H. Establishing a standard definition for child overweight and obesity worldwide: international survey. *BMJ* **320**, 1240-3 (2000).
39. Kampstra, P. Beanplot: A boxplot alternative for visual comparison of distributions. *Journal of Statistical Software* **28**, 1-9 (2008).
40. Wilson-Fritch, L. *et al.* Mitochondrial biogenesis and remodeling during adipogenesis and in response to the insulin sensitizer rosiglitazone. *Mol Cell Biol* **23**, 1085-94 (2003).
41. Keuper, M. *et al.* Spare mitochondrial respiratory capacity permits human adipocytes to maintain ATP homeostasis under hypoglycemic conditions. *FASEB J* **28**, 761-70 (2014).

8

## Supplementary Table 1: Characteristics of genotyped individuals from the Samoan Studies

| | 2010 GWA Study (Samoa) | | | | 1991 Study (Samoa) | | | | 1990 Study (American Samoa) | | | |
| | Men (n = 1235) | | Women (n = 1837) | | Men (n = 291) | | Women (n = 316) | | Men (n = 188) | | Women (n = 225) | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Age (years) | 45.4 | (11.4) | 44.7 | (11.1) | 38.6 | (9.1) | 39.1 | (9.1) | 40.4 | (9.9) | 39.2 | (10.5) |
| **Adiposity traits** | | | | | | | | | | | | |
| BMI (kg/m$^2$) | 31.3 | (5.9) | 34.9 | (6.8) | 28.9 | (4.9) | 30.9 | (5.3) | 33.9 | (5.9) | 36.0 | (7.0) |
| Body fat (%) | 24.0 | (11.8) | 37.2 | (11.8) | — | — | — | — | — | — | — | — |
| Abdominal circ. (cm) | 102.1 | (15.0) | 108.3 | (14.5) | 93.1 | (12.9) | 97.8 | (13.7) | 106.8 | (14.6) | 109.6 | (15.3) |
| Hip circ. (cm) | 105.7 | (10.2) | 114.5 | (12.6) | 100.0 | (9.0) | 107.3 | (10.5) | 111.0 | (11.8) | 118.2 | (13.8) |
| Abdominal–hip ratio | 0.962 | (0.07) | 0.945 | (0.07) | 0.928 | (0.07) | 0.909 | (0.07) | 0.961 | (0.06) | 0.928 | (0.08) |
| Obesity (> 32 kg/m$^2$) | 509 | 41% | 1195 | 65% | 65 | 22% | 129 | 41% | 114 | 61% | 162 | 72% |
| **Metabolic traits** | | | | | | | | | | | | |
| Fasting glucose (mg/dL)[†] | 89.6 | (14.4) | 88.0 | (13.6) | 85.3 | (12.3) | 84.7 | (11.0) | 94.4 | (10.8) | 93.1 | (10.7) |
| Fasting insulin (µU/mL)[†] | 12.5 | (13.7) | 16.2 | (14.4) | 10.4 | (11.3) | 12.2 | (10.8) | 19.8 | (24.9) | 21.4 | (17.2) |
| HOMA-IR[†] | 2.9 | (3.6) | 3.6 | (3.6) | 2.3 | (3.0) | 2.6 | (2.6) | 4.9 | (7.4) | 5.1 | (4.5) |
| Adiponectin (µg/mL) | 4.9 | (2.5) | 6.1 | (3.1) | — | — | — | — | — | — | — | — |
| Leptin (ng/mL) | 7.7 | (7.0) | 25.5 | (13.8) | 4.3 | (4.4) | 17.0 | (9.2) | 10.1 | (22.6) | 25.7 | (11.7) |
| Diabetes | 185 | 16% | 293 | 17% | 9 | 3% | 12 | 3% | 25 | 13% | 17 | 8% |
| Hypertension | 441 | 36% | 583 | 32% | 60 | 21% | 41 | 13% | 57 | 30% | 53 | 24% |
| **Serum lipid levels** | | | | | | | | | | | | |
| Total cholesterol (mg/dL) | 200.3 | (38.7) | 199.2 | (36.1) | 204.4 | (37.2) | 209.6 | (35.1) | 202.1 | (39.4) | 196.0 | (36.8) |
| Triglycerides (mg/dL) | 139.4 | (112.9) | 115.2 | (80.6) | 91.5 | (52.7) | 81.2 | (38.4) | 162.8 | (117.4) | 103.6 | (48.1) |
| HDL (mg/dL) | 43.7 | (11.2) | 46.5 | (10.8) | 40.5 | (11.6) | 43.3 | (10.4) | 36.0 | (7.6) | 38.3 | (8.1) |
| LDL (mg/dL) | 129.6 | (35.3) | 129.9 | (32.7) | 145.4 | (36.0) | 150.1 | (30.9) | 134.8 | (35.7) | 137.1 | (34.3) |

| | 2003 Study Adults (Samoa) | | | | 2002 Study Adults (American Samoa) | | | | 2003 Study Children (Samoa) | | | |
| | Men (n = 245) | | Women (n = 248) | | Men (n = 254) | | Women (n =336) | | Boys (n = 189) | | Girls (n = 220) | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Age (years) | 40.9 | (16.3) | 44.0 | (17.0) | 43.0 | (16.5) | 43.0 | (16.0) | 11.3 | (3.5) | 11.6 | (3.5) |
| **Adiposity traits** | | | | | | | | | | | | |
| BMI (kg/m$^2$) | 28.8 | (5.4) | 33.2 | (7.7) | 33.4 | (7.6) | 36.5 | (8.4) | 19.1 | (3.5) | 20.1 | (4.2) |
| Body fat (%) | 28.1 | (7.3) | 39.4 | (6.8) | 33.5 | (6.8) | 41.6 | (6.3) | 16.2 | (5.3) | 22.6 | (7.5) |
| Abdominal circ. (cm) | 95.5 | (14.9) | 107.0 | (16.5) | 107.5 | (16.4) | 111.0 | (16.5) | 67.0 | (9.7) | 70.4 | (12.3) |
| Hip circ. (cm) | 103.3 | (9.7) | 114.8 | (14.2) | 113.5 | (15.7) | 123.2 | (16.2) | 77.0 | (12.2) | 82.1 | (14.4) |
| Abdominal–hip ratio | 0.921 | (0.08) | 0.931 | (0.08) | 0.947 | (0.07) | 0.902 | (0.08) | 0.873 | (0.05) | 0.859 | (0.05) |
| Obesity (> 32 kg/m$^2$) | 59 | 24% | 130 | 52% | 138 | 54% | 229 | 68% | 5* | 3%* | 13* | 6%* |
| **Metabolic traits** | | | | | | | | | | | | |
| Fasting glucose (mg/dL)[†] | 88.6 | (11.4) | 89.8 | (12.2) | 88.1 | (14.9) | 86.8 | (15.8) | 83.4 | (8.4) | 82.4 | (8.4) |
| Fasting insulin (µU/mL)[†] | 7.1 | (9.1) | 10.0 | (9.6) | 12.7 | (13.4) | 14.5 | (17.9) | 5.1 | (5.0) | 8.6 | (10.3) |
| HOMA-iR[†] | 1.7 | (2.4) | 2.4 | (2.6) | 2.9 | (3.3) | 3.3 | (4.7) | 1.1 | (1.1) | 1.8 | (2.5) |
| Adiponectin (µg/mL) | 10.0 | (8.4) | 12.5 | (7.9) | 8.1 | (5.5) | 11.0 | (9.7) | 13.9 | (10.7) | 13.7 | (6.3) |
| Leptin (ng/mL) | 6.4 | (6.9) | 24.5 | (14.1) | 11.4 | (9.7) | 30.0 | (15.8) | 4.0 | (3.9) | 9.8 | (7.7) |
| Diabetes | 19 | 8% | 25 | 10% | 58 | 23% | 65 | 19% | — | — | — | — |
| Hypertension | 68 | 28% | 75 | 30% | 119 | 47% | 117 | 35% | — | — | — | — |
| **Serum lipid levels** | | | | | | | | | | | | |
| Total cholesterol (mg/dL) | 195.8 | (40.4) | 202.3 | (35.9) | 189.5 | (37.9) | 187.2 | (38.6) | 158.9 | (25.0) | 168.3 | (26.9) |
| Triglycerides (mg/dL) | 120.3 | (91.4) | 110.9 | (58.9) | 200.2 | (207.3) | 130.9 | (78.3) | 73.4 | (27.6) | 87.1 | (44.6) |
| HDL (mg/dL) | 46.3 | (11.2) | 47.2 | (10.2) | 38.6 | (8.8) | 42.1 | (8.5) | 49.7 | (11.2) | 49.8 | (11.4) |
| LDL (mg/dL) | 126.1 | (37.7) | 133.0 | (32.2) | 118.2 | (34.5) | 118.9 | (34.3) | 94.8 | (21.0) | 101.5 | (24.7) |

Summary statistics based on those who were both phenotyped and successfully genotyped (for either rs12513649 or rs373863828). Numbers are means and (standard deviations) for all traits except obesity, diabetes and hypertension, which are counts and percentages. Percent body fat and serum adiponectin are not available for the 1990–91 Studies; self-reported diabetes and hypertension were exclusion criteria for the 1990–91 Studies. [†]Non-diabetics only (n = 966 men, n = 1,423 women). *Children were classified as obese per Cole et al.[38]

**Supplementary Table 2: Association of rs373863828 with untransformed adiposity, metabolic, and lipid traits in (a) the discovery sample and (b) the adult replication sample.**

| (a) Discovery sample | All adults | | | | | Men | | | | | Women | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Quantitative Trait | n | β | s.e. | P | Covariates* | n | β | s.e. | P | Covariates* | n | β | s.e. | P | Covariates* |
| **Adiposity traits** | | | | | | | | | | | | | | | |
| BMI (kg/m²) | 3066 | 1.356 | 0.183 | **1.12E-13** | A, A², S, A×S | 1233 | 0.967 | 0.265 | **2.57E-04** | A, A² | 1833 | 1.644 | 0.247 | **2.75E-11** | A, A² |
| Body fat (%) | 2893 | 2.199 | 0.345 | **1.78E-10** | A, A², S, A×S | 1150 | 1.677 | 0.548 | 2.20E-03 | A, A² | 1743 | 2.559 | 0.442 | **6.92E-09** | A, A² |
| Abdominal circ. (cm) | 3057 | 2.842 | 0.404 | **2.05E-12** | A, A², S, A×S, A²×S | 1231 | 2.258 | 0.638 | **3.98E-04** | A, A² | 1826 | 3.235 | 0.520 | **5.01E-10** | A, A² |
| Hip circ. (cm) | 3058 | 2.361 | 0.332 | **1.19E-12** | A, A², S, A²×S | 1230 | 1.769 | 0.462 | **1.30E-04** | A, A² | 1828 | 2.776 | 0.458 | **1.31E-09** | A, A² |
| Abdominal–hip ratio | 3056 | 0.005 | 0.002 | 2.23E-03 | A, A², S, A×S, A²×S | 1230 | 0.005 | 0.003 | 0.051 | A, A² | 1826 | 0.005 | 0.002 | 0.019 | A |
| **Metabolic traits** | | | | | | | | | | | | | | | |
| Fasting glucose (mg/dL)† | 2393 | -1.652 | 0.423 | **9.52E-05** | A, A², S | 970 | -2.448 | 0.687 | **3.62E-04** | A, A² | 1423 | -1.019 | 0.535 | 0.057 | A, A² |
| Fasting insulin (µU/mL)† | 2392 | 1.342 | 0.449 | 0.003 | A, S, A×S | 970 | 0.619 | 0.684 | 0.365 | | 1422 | 1.809 | 0.592 | 2.23E-03 | A |
| HOMA-IR† | 2392 | 0.241 | 0.114 | 0.035 | A, S, A×S | 970 | 0.080 | 0.181 | 0.660 | A, A² | 1422 | 0.355 | 0.146 | 0.015 | A |
| Adiponectin (µg/mL) | 2858 | -0.228 | 0.083 | 0.006 | A, A², S, A×S | 1151 | -0.251 | 0.113 | 0.027 | A, A² | 1707 | -0.235 | 0.116 | 0.043 | A, A² |
| Leptin (ng/mL)‡ | — | — | — | — | | 1151 | 0.719 | 0.326 | 0.027 | A | 1707 | 1.888 | 0.525 | **3.25E-04** | |
| **Metabolic traits adjusted for BMI** | | | | | | | | | | | | | | | |
| Fasting glucose (mg/dL)† | 2383 | -2.248 | 0.417 | **6.89E-08** | A, A², S, B | 964 | -2.833 | 0.682 | **3.24E-05** | A, A², B | 1419 | -1.756 | 0.524 | **8.01E-04** | A, B |
| Fasting insulin (µU/mL)† | 2382 | 0.225 | 0.420 | 0.592 | A, A², S, B, A×S, A²×S | 964 | -0.224 | 0.632 | 0.723 | B | 1418 | 0.513 | 0.557 | 0.357 | A, A², B |
| HOMA-IR† | 2382 | -0.034 | 0.107 | 0.754 | A, B | 964 | -0.130 | 0.170 | 0.444 | B | 1418 | 0.029 | 0.138 | 0.834 | A, B |
| Adiponectin (µg/mL) | 2844 | -0.066 | 0.080 | 0.412 | A, A², S, B, A×S | 1143 | -0.130 | 0.109 | 0.233 | A, A², B | 1701 | -0.042 | 0.111 | 0.707 | A, A², B |
| Leptin (ng/mL)‡ | — | — | — | — | | 1143 | -0.262 | 0.210 | 0.213 | A, A², B | 1701 | -0.516 | 0.366 | 0.159 | A, A², B |
| **Serum lipid levels** | | | | | | | | | | | | | | | |
| Total cholesterol (mg/dL) | 2858 | -3.203 | 1.029 | **1.84E-03** | A, A², S, A×S, A²×S | 1151 | -3.423 | 1.731 | 0.048 | A, A² | 1707 | -3.319 | 1.256 | 0.008 | A, A² |
| Triglycerides (mg/dL) | 2858 | 0.349 | 2.769 | 0.900 | A, S, A×S | 1151 | -5.838 | 5.220 | 0.263 | A | 1707 | 4.676 | 2.981 | 0.117 | A |
| HDL (mg/dL) | 2858 | -0.322 | 0.321 | 0.317 | A, A², S | 1151 | 0.406 | 0.516 | 0.431 | A | 1707 | -0.914 | 0.408 | 0.025 | A |
| LDL (mg/dL) | 2851 | -2.347 | 0.945 | 0.013 | A, A², S, A²×S | 1145 | -2.115 | 1.586 | 0.182 | A, A² | 1706 | -2.647 | 1.155 | 0.022 | A, A² |

| Dichotomous Trait | n | OR | 95% CI | p | Covariates* | n | OR | 95% CI | p | Covariates* | n | OR | 95% CI | p | Covariates* |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Obesity (> 32 kg/m²) | 3066 | 1.305 | (1.159, 1.470) | **1.12E-05** | A, A², S, A×S | 1233 | 1.270 | (1.052, 1.535) | 0.013 | A, A² | 1833 | 1.335 | (1.144, 1.557) | **2.38E-04** | A, A² |
| Diabetes | 2876 | 0.637 | (0.536, 0.758) | **3.86E-07** | A | 1157 | 0.611 | (0.461, 0.811) | **6.31E-04** | A, A² | 1719 | 0.669 | (0.537, 0.833) | **3.40E-04** | A, A² |
| Diabetes adj. for BMI | 2861 | 0.586 | (0.489, 0.702) | **6.68E-09** | A, B | 1149 | 0.623 | (0.495, 0.784) | **5.49E-05** | A, A², B | 1712 | 0.566 | (0.422, 0.760) | **1.50E-04** | A, A², B |
| Hypertension | 3041 | 1.014 | (0.898, 1.145) | 0.818 | A, S | 1226 | 0.923 | (0.760, 1.120) | 0.416 | A | 1815 | 1.087 | (0.930, 1.269) | 0.295 | A |

**Boldface** represents a *P* value < 2.17E-03.

* A = age, A² = age², S = sex, A×S = age × sex interaction, A²×S = age² × sex interaction, B = log(BMI).

† Analysis conducted only in non-diabetics

‡ Leptin was not analyzed in men and women combined because the distributions in each sex were very different.

Abbreviations: s.e., standard error; OR, odds ratio; 95% CI, 95% confidence interval; circ., circumference; adj., adjusted

10

| (b) Replication sample (mega analysis) | All adults | | | | | Men | | | | | Women | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Quantitative Trait | *n* | β | s.e. | *P* | Covariates* | *n* | β | s.e. | *P* | Covariates* | *n* | β | s.e. | *P* | Covariates* |
| **Adiposity traits** | | | | | | | | | | | | | | | |
| BMI (kg/m²) | 2103 | 1.453 | 0.237 | **8.22E-10** | A, A², S, N, C | 978 | 1.501 | 0.306 | **9.54E-07** | A, A², N | 1125 | 1.389 | 0.348 | **6.41E-05** | A, A², N, C |
| Body fat (%) | 880 | 1.335 | 0.392 | **6.58E-04** | A, A², S, A×S, N | 401 | 1.192 | 0.595 | 0.045 | A, A², N | 479 | 1.314 | 0.491 | 0.007 | A, A², N |
| Abdominal circ. (cm) | 2172 | 3.218 | 0.518 | **5.12E-10** | A, A², S, N, C | 1008 | 3.318 | 0.704 | **2.42E-06** | A, A², N, C | 1164 | 3.087 | 0.735 | **2.64E-05** | A, A², N, C |
| Hip circ. (cm) | 2165 | 2.716 | 0.462 | **4.27E-09** | A, A², S, N, C | 1002 | 2.838 | 0.605 | **2.76E-06** | A, A², N, C | 1163 | 2.597 | 0.674 | **1.16E-04** | A, A², N, C |
| Abdominal–hip ratio | 2162 | 0.006 | 0.002 | 0.017 | A, A², S, A×S, A²×S, N, C | 1001 | 0.006 | 0.003 | 0.052 | A, A², N, C | 1161 | 0.005 | 0.004 | 0.126 | A, A² |
| **Metabolic traits** | | | | | | | | | | | | | | | |
| Fasting glucose (mg/dL)[†] | 1948 | -1.541 | 0.463 | **8.84E-04** | A, A², S, N | 901 | -1.508 | 0.668 | 0.024 | A, A², N | 1047 | -1.764 | 0.634 | 0.005 | A, A², N |
| Fasting insulin (μU/mL)[†] | 1947 | 2.500 | 0.565 | **9.55E-06** | A, A², S, A×S, N, C | 900 | 2.595 | 0.838 | **1.96E-03** | N, C | 1047 | 2.174 | 0.740 | 0.003 | A², N, C |
| HOMA-IR[†] | 1947 | 0.572 | 0.150 | **1.43E-04** | A, A², S, A×S, N, C | 900 | 0.663 | 0.228 | 0.004 | A, A², N, C | 1047 | 0.440 | 0.191 | 0.022 | A, A², N, C |
| Adiponectin (μg/mL) | 1079 | -1.078 | 0.426 | 0.011 | A, A², S, A×S, A²×S, N | 497 | -1.153 | 0.529 | 0.029 | A, A², N | 582 | -0.878 | 0.628 | 0.162 | A, A², N |
| Leptin (ng/mL)[‡] | — | — | — | — | | 831 | 2.237 | 0.607 | **2.26E-04** | A, A², N, C | 952 | 2.548 | 0.726 | **4.47E-04** | A, A², N, C |
| **Metabolic traits adjusted for BMI** | | | | | | | | | | | | | | | |
| Fasting glucose (mg/dL)[†] | 1867 | -2.094 | 0.468 | **7.62E-06** | A, A², S, B, N | 866 | -2.137 | 0.672 | **1.47E-03** | A, A², B, N | 1001 | -2.274 | 0.642 | **3.96E-04** | A, A², B, N |
| Fasting insulin (μU/mL)[†] | 1866 | 1.557 | 0.539 | 0.004 | A, A², S, B, A×S, N, C | 865 | 1.874 | 0.781 | 0.016 | B, N, C | 1001 | 1.185 | 0.723 | 0.101 | A, A², B, N, C |
| HOMA-IR[†] | 1866 | 0.358 | 0.147 | 0.015 | A, S, B, A×S, N, C | 865 | 0.475 | 0.220 | 0.031 | B, N, C | 1001 | 0.221 | 0.190 | 0.245 | A, B, N, C |
| Adiponectin (μg/mL) | 1068 | -0.780 | 0.428 | 0.068 | A, A², S, B, A×S, A²×S, N | 491 | -0.928 | 0.527 | 0.078 | A, A², B, N | 577 | -0.439 | 0.633 | 0.488 | A², B |
| Leptin (ng/mL)[‡] | — | — | — | — | | 801 | 0.863 | 0.521 | 0.098 | A, A², B, C | 919 | -0.009 | 0.511 | 0.985 | A, A², B, N, C |
| **Serum lipid levels** | | | | | | | | | | | | | | | |
| Total cholesterol (mg/dL) | 1849 | -1.905 | 1.344 | 0.157 | A, A², S, A×S, A²×S, N | 860 | -0.891 | 1.945 | 0.647 | A, A², N | 989 | -2.525 | 1.812 | 0.163 | A, A², N, C |
| Triglycerides (mg/dL) | 1849 | -4.888 | 4.153 | 0.239 | A, A², S, A×S, A²×S, N, C | 860 | -13.11 | 7.729 | 0.090 | A, A², N, C | 989 | 1.683 | 3.629 | 0.643 | A, A², N, C |
| HDL (mg/dL) | 1834 | -1.097 | 0.391 | 0.005 | A, A², S, N, C | 848 | -1.088 | 0.578 | 0.060 | A, A², N, C | 986 | -0.948 | 0.516 | 0.066 | A, A², N, C |
| LDL (mg/dL) | 1805 | -1.047 | 1.291 | 0.417 | A, A², S, A×S, A²×S, N, C | 825 | 0.156 | 1.951 | 0.936 | A, A², N, C | 980 | -1.857 | 1.671 | 0.266 | A, A², N, C |

| Dichotomous Trait | n | OR | 95% CI | p | Covariates* | n | OR | 95% CI | p | Covariates* | n | OR | 95% CI | p | Covariates* |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Obesity (> 32 kg/m²) | 2103 | 1.441 | (1.227, 1.692) | **8.49E-06** | A, A², S, N, C | 978 | 1.586 | (1.252, 2.009) | **1.35E-04** | A, A², N | 1125 | 1.32 | (1.061, 1.643) | 0.013 | A, A², N, C |
| Diabetes | 2145 | 0.831 | (0.639, 1.081) | 0.168 | A, A², S, A×S, N, C | 1000 | 0.698 | (0.472, 1.031) | 0.071 | A, A², N, C | 1145 | 0.950 | (0.670, 1.348) | 0.774 | A, A², N, C |
| Diabetes adj. for BMI | 2053 | 0.742 | (0.567, 0.969) | 0.029 | A, A², S, B, N, C | 960 | 0.550 | (0.358, 0.845) | 0.006 | A, A², B, N, C | 1093 | 0.915 | (0.646, 1.296) | 0.616 | A, A², B, N, C |
| Hypertension | 2173 | 1.045 | (0.881, 1.240) | 0.613 | A, A², S, A×S, N, C | 1006 | 1.029 | (0.817, 1.296) | 0.809 | A, A², N, C | 1167 | 1.072 | (0.834, 1.377) | 0.587 | A, A², N, C |

**Boldface** represents a *P* value < 2.17E-03.

\* A = age, A² = age², S = sex, A×S = age × sex interaction,

A²×S = age² × sex interaction, B = log(BMI), N = nation, C = study (1990s vs 2000s)

† Analysis conducted only in non-diabetics

‡ Leptin was not analyzed in men and women combined because the distributions in each sex were very different.

Abbreviations: s.e., standard error; OR, odds ratio; 95% CI, 95% confidence interval; circ., circumference; adj., adjusted

11

**Supplementary Table 3: Oligonucleotide primers used in the study**

| Gene | Forward | Reverse |
|---|---|---|
| CREBRF[1] | ATGTATGAACTGGATAGAGAGATG | GTTAGGTCTTCACAGTATGTATCC |
| CREBRF[2] | GAAGACCTGAAGGAGGTGACT | GTTCCACTCAGATGGTCTCAGC |
| Crebrf | GAGGACTTGAAGGAGATGACG | CAGAAGGCCTCAGAATCCTC |
| Pparg2 | CCAGAGCATGGTGCCTTCGCT | CAGCAACCATTGGGTCAG |
| Cebpa | CAAGAACAGCAACGAGTACCG | GTCACTGGTCAACTCCAGCAC |
| Adipoq | TGTTCCTCTTAATCCTGCCCA | CCAACCTGCACAAGTTCCCTT |
| Actb | CCACTGCCGCATCCTCTTCC | CTCGTTGCCAATAGTGATGACCTG |

[1]Used for data in Supplementary Fig 5a.
[2]Used for data in Fig. 2a.

I apologize, but I notice the transcription got corrupted. Let me provide the clean version:

12

Nature Genetics: doi:10.1038/ng.3620