

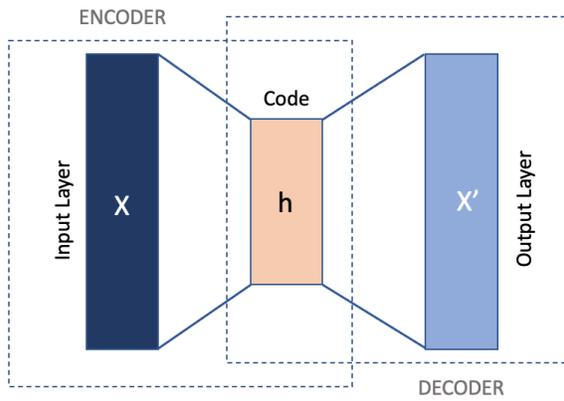
Unsupervised Deep Learning of Electrocardiograms Enables Scalable Human Disease Profiling

Supplementary Materials

Table of Contents

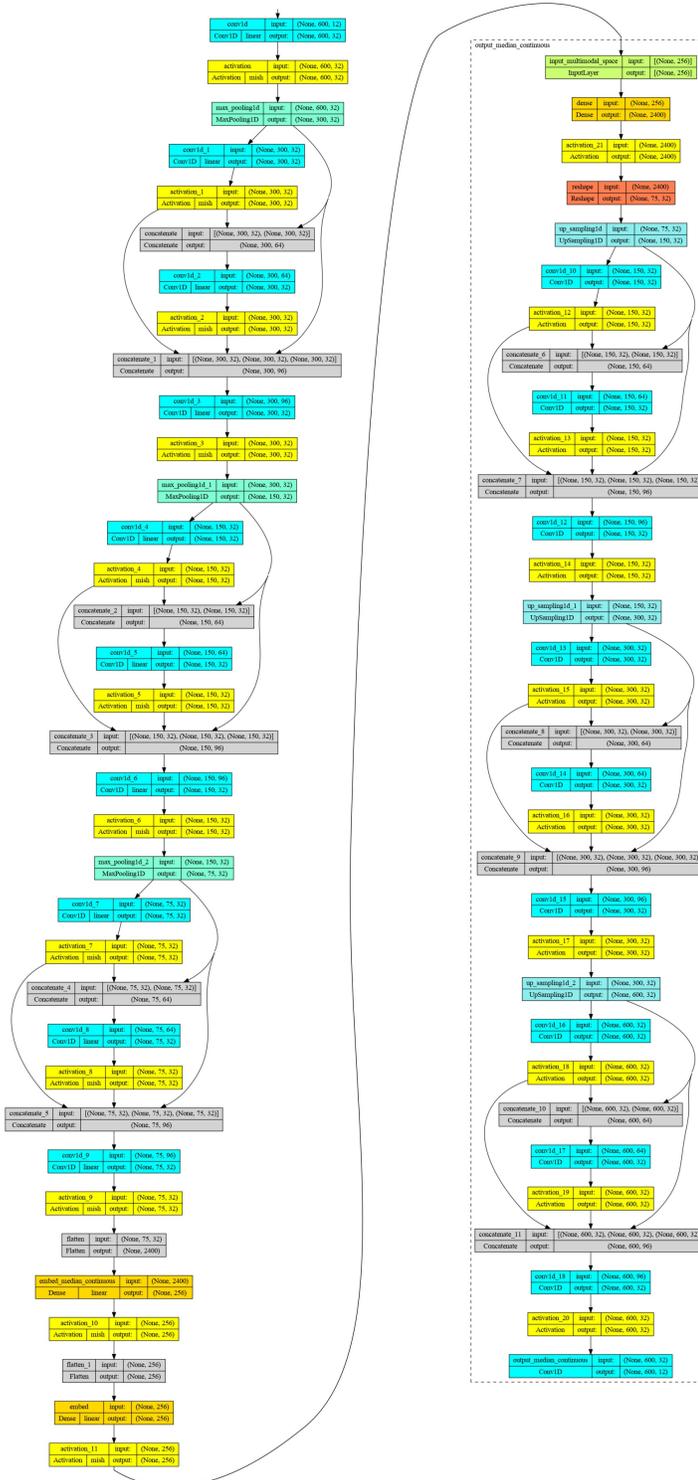
Supplementary Figure 1. General autoencoder schematic.	1
Supplementary Figure 2. Autoencoder architecture.	2
Supplementary Figure 3. Example of median waveforms generated from an original 12-lead electrocardiogram.	3
Supplementary Figure 4. Latent space representations of age, sex, and body mass index.	4
Supplementary Figure 5. Cluster centroids and phenotype vector derivation.	5
Supplementary Figure 6. Latent space 12-lead electrocardiogram phenome-wide association study perturbation tests.	6
Supplementary Figure 7. Forest plots demonstrating associations for the top 3 phecodes in each of the datasets for both 12-lead and single-lead ECG latent space PheWAS analyses.	7
Supplementary Figure 8. Improvements in Phecode discrimination using ECG vector component scores.	8
Supplementary Figure 9. Phenotype vector correlation matrix.	9
Supplementary Figure 10. Vector component score illustration.	10
Supplementary Tables 1-10.	11
Supplementary Table 11. Top associations by effect size across disease category for incident disease.	12
Supplementary Tables 12-13.	11

Supplementary Figure 1. General autoencoder schematic.



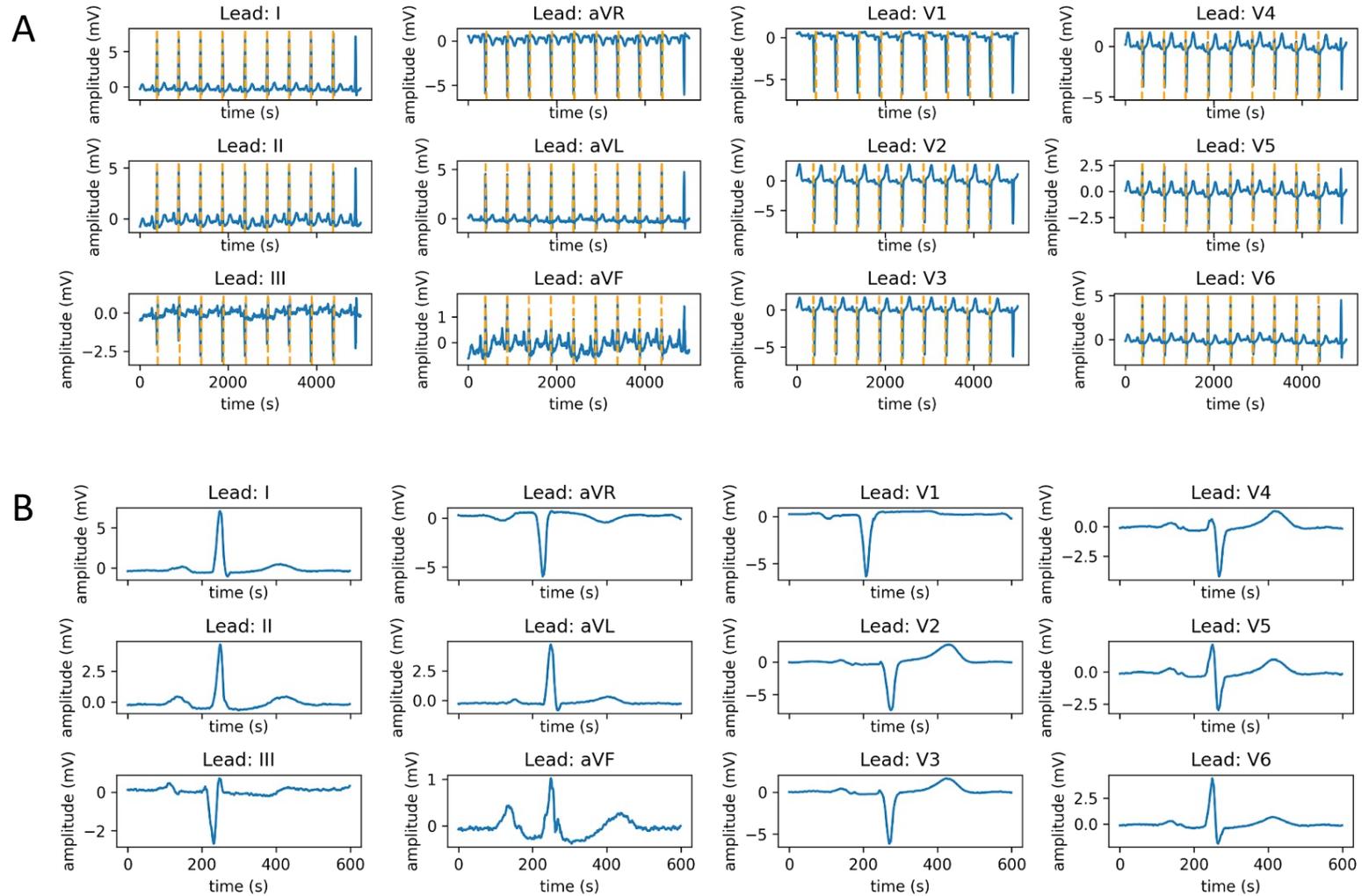
General autoencoder structure. The input, X , is encoded into a multi-dimensional latent space, h . This “code” is then used to reconstruct the input as accurately as possible at the output layer, X' .

Supplementary Figure 2. Autoencoder architecture.



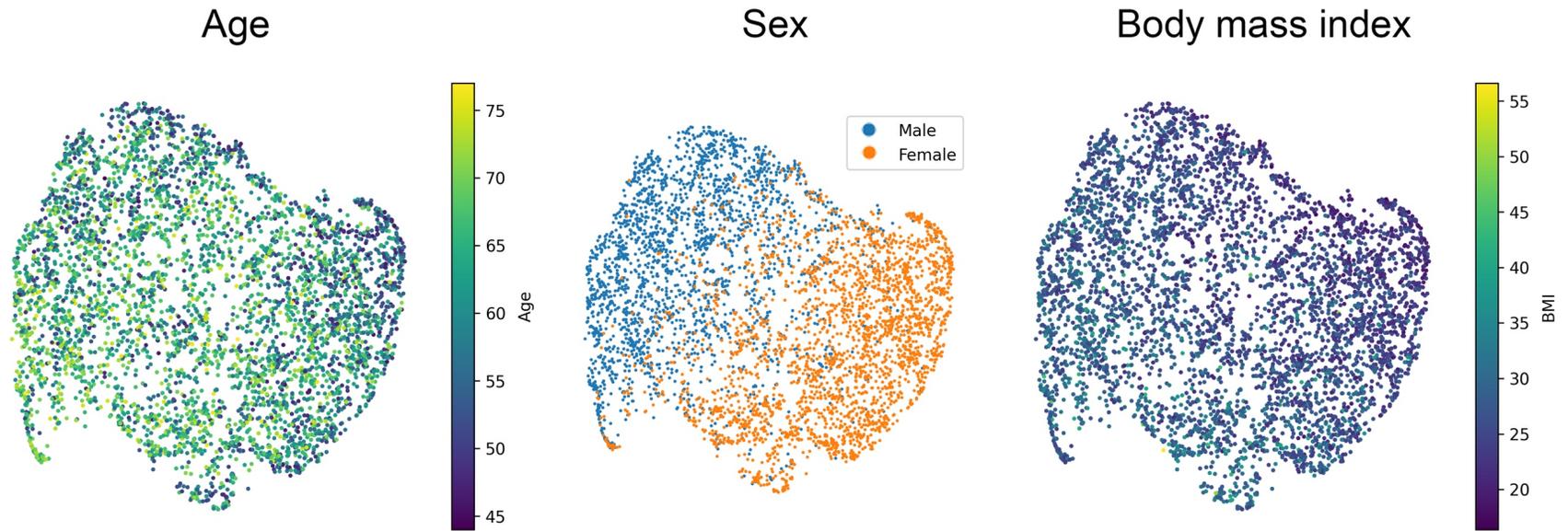
Depicted is a summary of the architecture of the convolutional autoencoder using in the current study. The encoder (left) takes 10 seconds of continuous 12-lead ECG waveform data as input to generate a lower dimensional latent space. The decoder (right) takes the latent space as input to reconstruct the 12-lead ECG. Box colors correspond to the following layer types: green = input, cyan = convolution, yellow = activation, teal = pooling, red = reshape, and orange = dense.

Supplementary Figure 3. Example of median waveforms generated from an original 12-lead electrocardiogram.



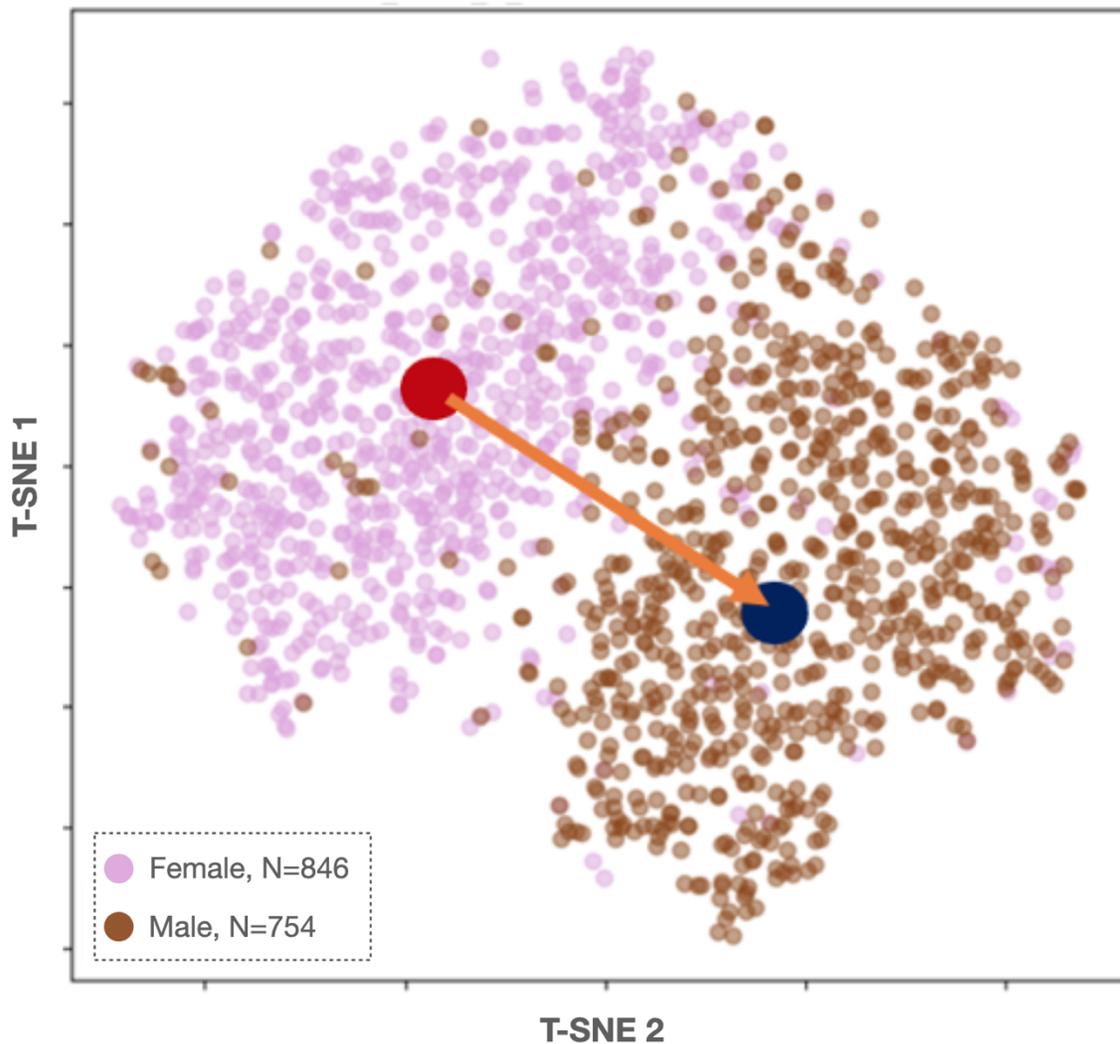
Panel A represents an exemplar original 12-lead tracing, and panel B represents the corresponding median waveform for each lead.

Supplementary Figure 4. Latent space representations of age, sex, and body mass index



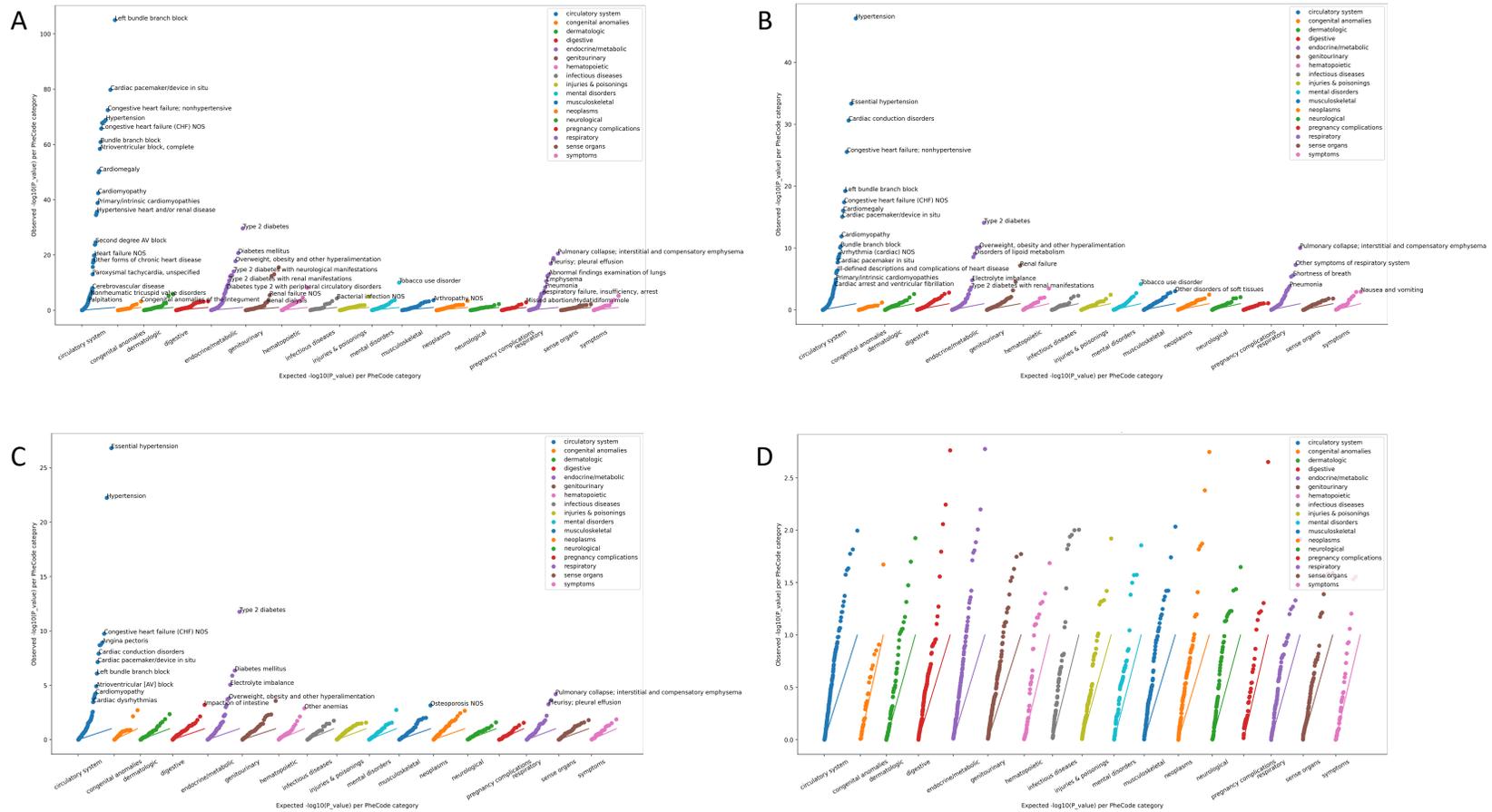
Figures show uniform manifold approximation and projections (UMAPs) depicting depictions of the ECG latent space stratified by age (left), sex (middle), and body mass index (right). Maps were derived using ECG encodings in latent space of XX.

Supplementary Figure 5. Cluster centroids and phenotype vector derivation.



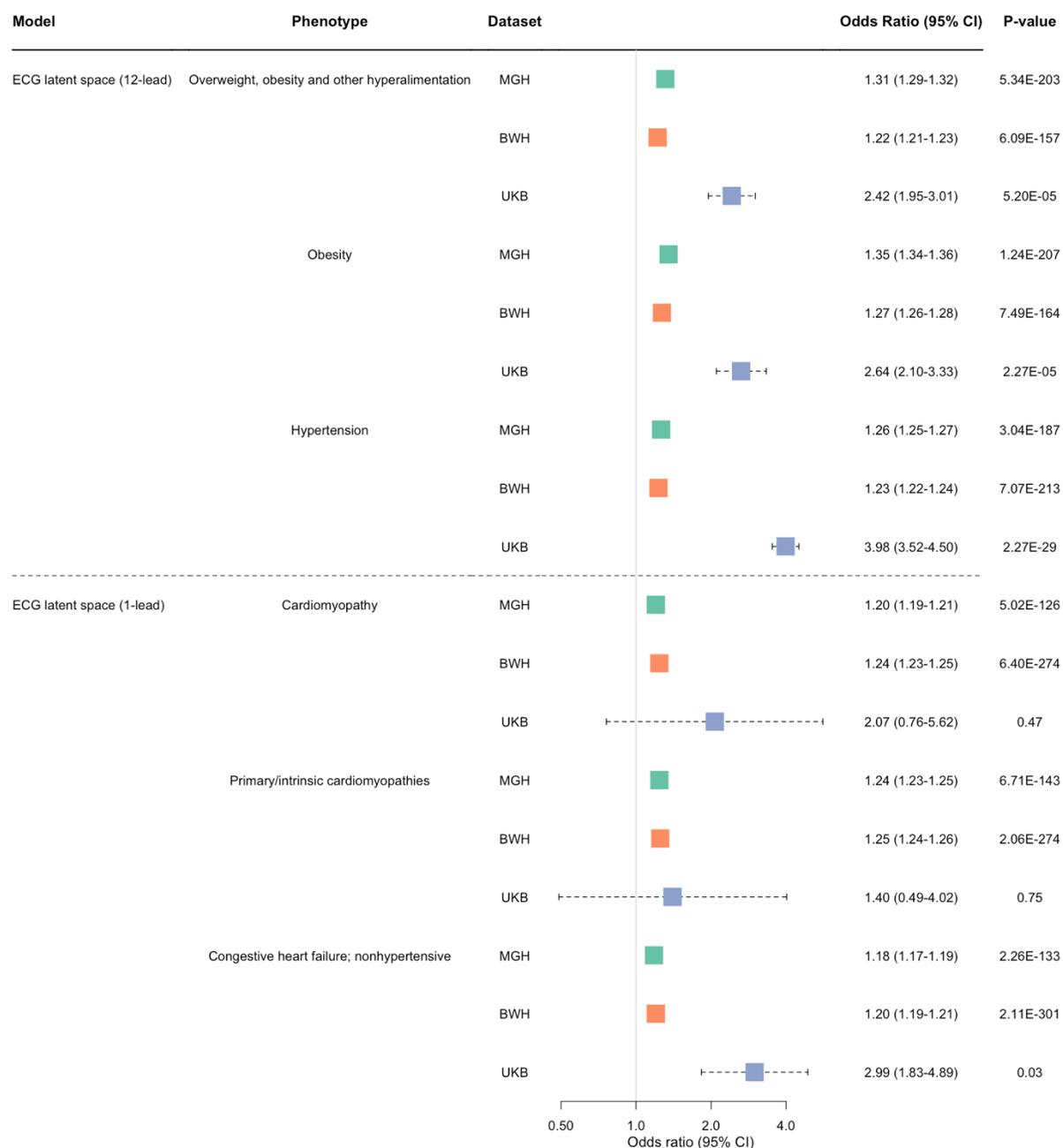
Illustrative example of phenotype vector derivation. The figure shows a two-dimensional t-Distributed Stochastic Neighbor Embedding (t-SNE) scatter plot and demonstrates how sex affects the distribution of electrocardiogram (ECG) encodings in latent space of 1600 UKB participants. ECG encodings from female participants (small pink dots) and male participants (small brown dots) form distinct phenotypic clusters. The larger red and blue dots mark the centroids of these clusters, respectively. The orange arrow represents the phenotype vector for sex.

Supplementary Figure 6. Latent space 12-lead electrocardiogram phenome-wide association study perturbation tests.



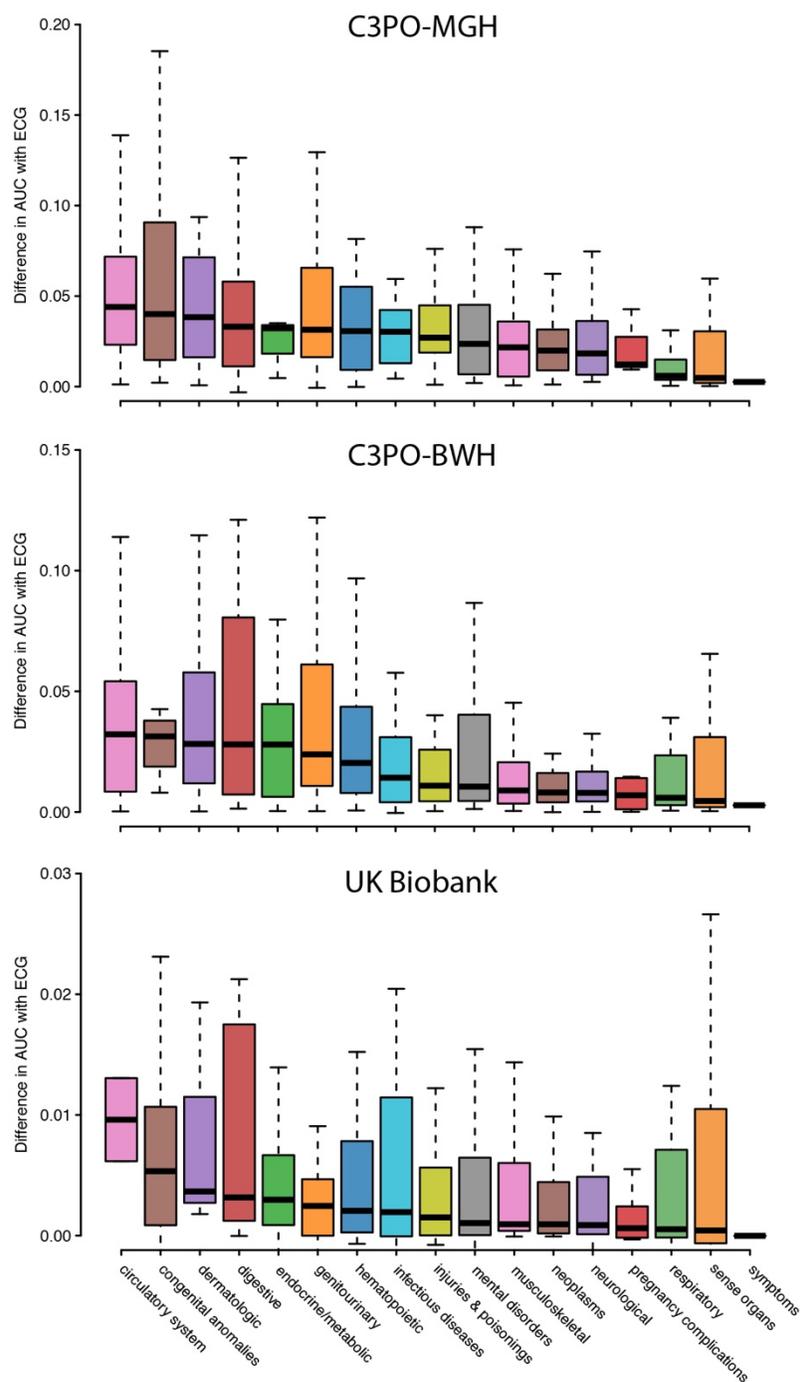
Panel A represents the results of the phenome-wide association study results in the Massachusetts General Hospital Community Care Cohort Project sample at baseline. Each phecode tested for association is represented as a single point on the plot. The x-axis represents the phenotype category and the y-axis represents the $-\log_{10}(p\text{-value})$ for the association test. The remaining panels represent the results in which the phecode labels are randomly reclassified with (B) 10%, (C) 20%, and D (100%) reclassification.

Supplementary Figure 7. Forest plots demonstrating associations for the top 3 phecodes in each of the datasets for both 12-lead and single-lead ECG latent space PheWAS analyses.



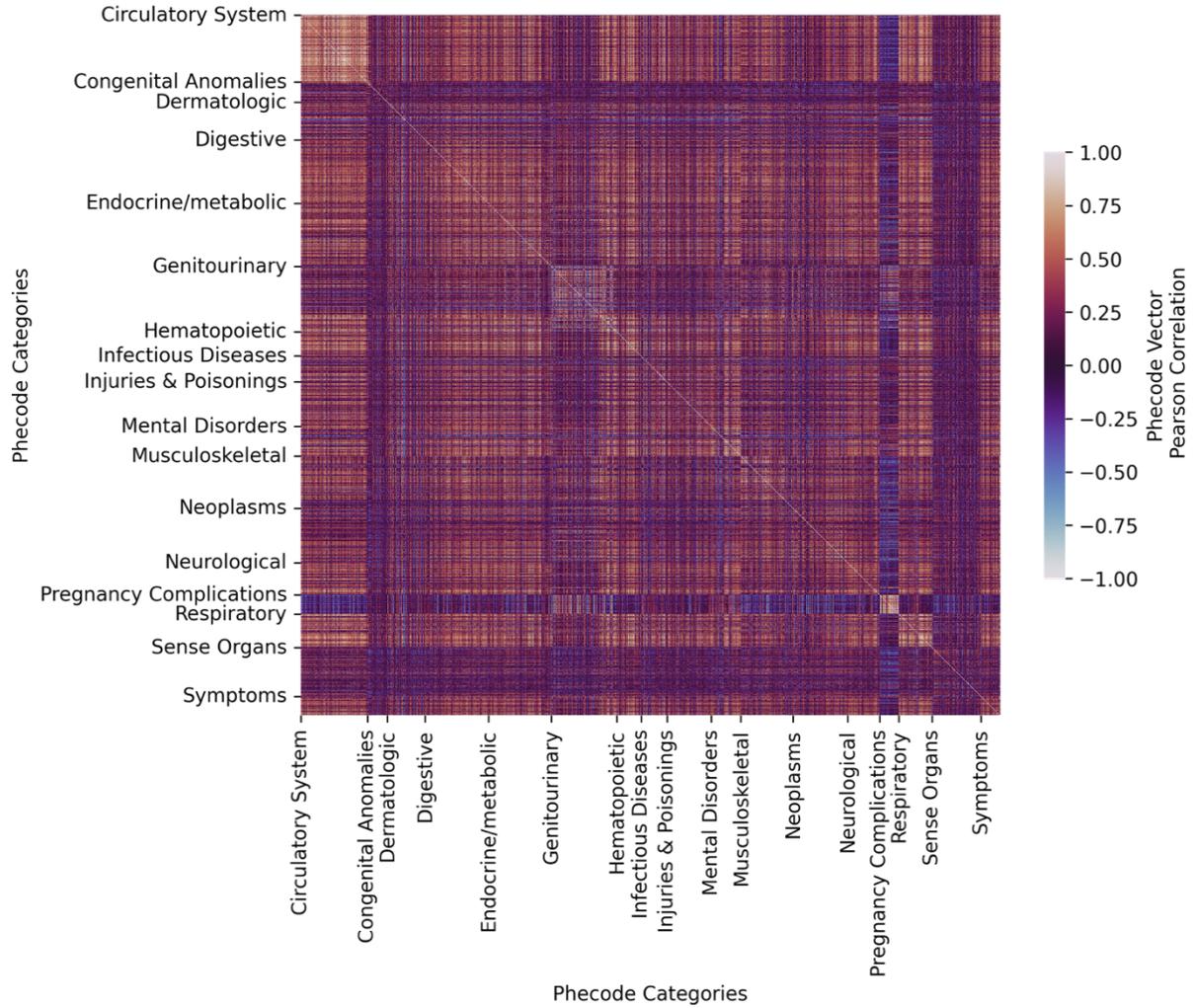
Depicted is a forest plot summarizing the top three strongest Phecode associations with the autoencoder model using 12-lead ECG (top) and single-lead ECG (bottom) in the primary meta-analysis. Points depict study-specific odds ratios (per 1-point increase in vector component score) and error bars 95% confidence intervals.

Supplementary Figure 8. Improvements in Phecode discrimination using ECG vector component scores



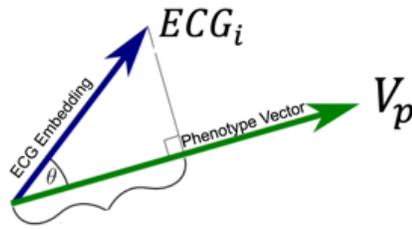
Plots depict the difference in the area under the receiver operator characteristic curve (AUC) between models with and without the ECG vector component scores, stratified by Phecode category and dataset. Only conditions meeting Bonferroni-corrected phenome-wide significance in the primary meta-analysis ($p < 3.1 \times 10^{-5}$) are included. Multivariable logistic regression models were fit with adjustment for age, sex, race, and vector component score. In the MGH-C3PO and BWH-C3PO datasets, models were further adjusted for ECG acquisition date and zero padding as described in the methods.

Supplementary Figure 9. Phenotype vector correlation matrix.



Correlation matrix for phenotype vectors based on the 12-lead electrocardiogram model, organized by disease category. The color scale represents Pearson correlation coefficients.

Supplementary Figure 10. Vector component score illustration



Depicted is an illustration of the calculation of vector component scores. Each ECG encoding (“ECG embedding”, ECG_i) projects onto each phenotype vector, V_p . The projected component is calculated from the angle between the ECG encoding and the phenotype vector, scaled by the length of the ECG. Thus, the projected component signifies the latent space position of a single individual along a single phenotype vector and therefore represents a disease-specific “vector component score”.

Supplementary Tables 1–10.
See separate data file.

Supplementary Table 11. Top associations by effect size across disease category for incident disease

Disease grouping*	N events	Odds ratio (95% CI)	p
<i>Circulatory system</i>			
Cardiac defibrillator in situ	166	1.80 (1.65-1.95)	3.70x10 ⁻⁴³
Right bundle branch block	468	1.77 (1.70-1.84)	7.47x10 ⁻¹⁸³
Bundle branch block	835	1.69 (1.63-1.75)	4.22x10 ⁻¹⁶²
Left bundle branch block	409	1.64 (1.57-1.72)	5.06x10 ⁻¹¹²
Paroxysmal ventricular tachycardia	495	1.61 (1.53-1.70)	2.58x10 ⁻⁶⁵
<i>Congenital anomalies</i>			
Cardiac congenital anomalies	1051	1.33 (1.28-1.39)	1.46x10 ⁻⁴⁰
Cardiac and circulatory congenital anomalies	1307	1.32 (1.27-1.38)	2.04x10 ⁻⁴¹
Valvular heart disease/ heart chambers	285	1.27 (1.18-1.38)	1.74x10 ⁻⁰⁹
Cardiac shunt/ heart septal defect	494	1.20 (1.13-1.28)	2.36x10 ⁻⁰⁸
<i>Dermatologic</i>			
Chronic ulcer of leg or foot	1147	1.24 (1.20-1.29)	5.41x10 ⁻³⁷
Chronic ulcer of unspecified site	532	1.24 (1.17-1.31)	1.13x10 ⁻¹³
Decubitus ulcer	650	1.24 (1.17-1.30)	9.49x10 ⁻¹⁵
Chronic ulcer of skin	1856	1.24 (1.20-1.27)	1.95x10 ⁻⁴⁹
Cellulitis and abscess of trunk	804	1.17 (1.12-1.22)	3.22x10 ⁻¹²
<i>Digestive</i>			
Hemorrhage from gastrointestinal ulcer	319	1.34 (1.22-1.47)	7.84x10 ⁻¹⁰
Portal hypertension	123	1.31 (1.18-1.46)	8.84x10 ⁻⁰⁷
Gastric ulcer	470	1.29 (1.20-1.38)	1.50x10 ⁻¹¹
Liver abscess and sequelae of chronic liver disease	267	1.28 (1.19-1.38)	9.40x10 ⁻¹¹
Splenomegaly	398	1.25 (1.18-1.31)	4.99x10 ⁻¹⁶
<i>Endocrine/Metabolic</i>			
Cachexia	253	1.40 (1.27-1.54)	1.43x10 ⁻¹¹
Dysmetabolic syndrome X	367	1.37 (1.27-1.47)	7.44x10 ⁻¹⁸
Other disorders of pancreatic internal secretion	341	1.35 (1.25-1.46)	2.42x10 ⁻¹³
Disorders of magnesium metabolism	747	1.34 (1.27-1.41)	2.26x10 ⁻²⁸
Acid-base balance disorder	1320	1.33 (1.29-1.38)	1.26x10 ⁻⁶⁰
<i>Genitourinary</i>			
Renal dialysis	613	1.33 (1.27-1.39)	1.04x10 ⁻³⁵
End stage renal disease	363	1.31 (1.23-1.40)	8.49x10 ⁻¹⁸
Nephritis and nephropathy in diseases classified elsewhere	184	1.31 (1.21-1.42)	4.25x10 ⁻¹¹
Acute renal failure	2959	1.29 (1.26-1.31)	2.71x10 ⁻¹¹⁹
Chronic Kidney Disease, Stage III	1621	1.26 (1.22-1.30)	1.00x10 ⁻⁴⁰
<i>Hematopoietic</i>			
Anemia in chronic kidney disease	453	1.26 (1.19-1.33)	1.07x10 ⁻¹⁵
Acute posthemorrhagic anemia	1306	1.25 (1.21-1.30)	4.48x10 ⁻³⁴
Hemorrhagic disorder due to intrinsic circulating anticoagulants	170	1.25 (1.14-1.37)	3.74x10 ⁻⁰⁶
Secondary thrombocytopenia	547	1.24 (1.17-1.31)	1.40x10 ⁻¹⁴
Encounter for long-term (current) use of anticoagulants	2661	1.24 (1.21-1.27)	7.91x10 ⁻⁶³
<i>Infectious diseases</i>			
Septicemia	1848	1.31 (1.28-1.35)	5.20x10 ⁻⁷⁸
Bacteremia	905	1.27 (1.22-1.33)	3.17x10 ⁻²⁸
Infection/inflammation of internal prosthetic device; implant; and graft	601	1.26 (1.20-1.32)	1.34x10 ⁻²⁰
Gram positive septicemia	173	1.26 (1.14-1.39)	5.80x10 ⁻⁰⁶
Methicillin resistant Staphylococcus aureus	340	1.24 (1.14-1.35)	3.60x10 ⁻⁰⁷
<i>Injuries and poisonings</i>			

Non-healing surgical wound	181	1.34 (1.22-1.48)	5.08x10 ⁻⁰⁹
Traumatic amputation	124	1.33 (1.17-1.52)	2.19x10 ⁻⁰⁵
Certain early complications of trauma or procedure	463	1.31 (1.22-1.40)	2.84x10 ⁻¹⁵
Septic shock	561	1.30 (1.23-1.37)	2.82x10 ⁻²³
Sepsis	1452	1.30 (1.25-1.34)	1.37x10 ⁻⁵⁶
<i>Mental disorders</i>			
Delirium due to conditions classified elsewhere	1193	1.26 (1.22-1.31)	7.90x10 ⁻³⁷
Altered mental status	2465	1.24 (1.21-1.27)	5.92x10 ⁻⁵³
Alcoholic liver damage	176	1.23 (1.13-1.34)	3.81x10 ⁻⁰⁶
Psychosis	540	1.22 (1.15-1.29)	1.70x10 ⁻¹²
Alteration of consciousness	1357	1.18 (1.14-1.23)	4.87x10 ⁻¹⁹
<i>Musculoskeletal</i>			
Acute osteomyelitis	209	1.30 (1.20-1.41)	6.41x10 ⁻¹¹
Unspecified osteomyelitis	349	1.28 (1.21-1.35)	9.71x10 ⁻¹⁸
Osteomyelitis	446	1.26 (1.20-1.33)	2.50x10 ⁻¹⁸
Chronic osteomyelitis	177	1.24 (1.13-1.37)	5.86x10 ⁻⁰⁶
Osteomyelitis, periostitis, and other infections involving bone	474	1.24 (1.18-1.3)	6.46x10 ⁻¹⁷
<i>Neoplasms</i>			
Cancer of connective tissue	240	1.23 (1.12-1.34)	8.73x10 ⁻⁰⁶
Cancer of bone and connective tissue	388	1.18 (1.11-1.26)	1.36x10 ⁻⁰⁷
Cancer within the respiratory system	1045	1.17 (1.12-1.23)	4.84x10 ⁻¹³
Cancer of bronchus; lung	961	1.17 (1.12-1.22)	4.32x10 ⁻¹³
Secondary malignant neoplasm of liver	374	1.17 (1.10-1.24)	6.45x10 ⁻⁰⁷
<i>Neurological</i>			
Encephalopathy, not elsewhere classified	754	1.24 (1.18-1.31)	8.10x10 ⁻¹⁶
Chronic pain syndrome	662	1.20 (1.13-1.27)	1.25x10 ⁻⁰⁹
Convulsions	1054	1.16 (1.12-1.21)	1.04x10 ⁻¹³
Epilepsy, recurrent seizures, convulsions	1207	1.15 (1.11-1.20)	2.37x10 ⁻¹⁴
Other conditions of brain	1712	1.15 (1.11-1.19)	3.71x10 ⁻¹⁶
<i>Respiratory</i>			
Pseudomonal pneumonia	103	1.62 (1.43-1.84)	2.79x10 ⁻¹⁴
Pulmonary congestion and hypostasis	1958	1.35 (1.31-1.39)	5.26x10 ⁻¹⁰⁷
Respiratory failure	1134	1.35 (1.30-1.40)	2.16x10 ⁻⁵²
Emphysema	759	1.32 (1.25-1.39)	8.90x10 ⁻²⁶
Pneumonitis due to inhalation of food or vomitus	692	1.31 (1.25-1.37)	2.57x10 ⁻²⁸
<i>Sense organs</i>			
Other nondiabetic retinopathy	282	1.18 (1.10-1.27)	3.52x10 ⁻⁰⁶
Blindness and low vision	769	1.11 (1.06-1.17)	3.52x10 ⁻⁰⁵
Presbyopia	1118	1.07 (1.04-1.10)	1.41x10 ⁻⁰⁶
Dry eyes	1579	1.06 (1.03-1.09)	3.18x10 ⁻⁰⁵
Senile cataract	3649	1.05 (1.03-1.07)	5.34x10 ⁻¹⁰
<i>Symptoms</i>			
Cardiogenic shock	313	1.45 (1.34-1.56)	1.31x10 ⁻²²
Shock	747	1.38 (1.32-1.44)	2.71x10 ⁻⁴⁴
Gangrene	209	1.37 (1.25-1.50)	5.46x10 ⁻¹²
Debility unspecified	223	1.28 (1.15-1.43)	4.77x10 ⁻⁰⁶
Nonspecific abnormal findings on radiological and other examination of other intrathoracic organs (echocardiogram, etc)	1632	1.14 (1.10-1.18)	1.41x10 ⁻¹³
*Displayed are significant associations with the top 5 largest effect sizes within each disease category. In cases where there are fewer than 5 significant associations, all significant associations are shown.			

Supplementary Tables 12–13.
See separate data file.