Appendix A. Cohort Definition Example

751

758

766

767

769

771

772

773

774

781

782

783

785

787

We conducted a retrospective cohort study in the 752 All of Us Registered Tier, restricting to participants 753 with both OMOP-standardized EHRs and connected 754 Fitbit data. All time references are anchored to a 755 participant-specific prediction ("index") time defined 756 below.

A.1. Source population and phenotyping

From 8,477 All of Us participants with valid Fitbit and EHR linkage, we identified 2,148 partici-760 pants with any evidence of major depressive disorder (MDD) using the union of (i) OMOP concept IDs 762 for MDD, (ii) PheCode mappings for depressive disorders, and (iii) depression-specific medications. Phe-764 notyping was performed strictly prior to index. 765

A.2. Incident case definition and index time

Incident cases required the first MDD evidence to occur at least 180 days after Fitbit tracking began, and 768 evidence of EHR activity both before and after Fitbit start. For positives, the index time was set to one 770 day prior to the first MDD flag, yielding 348 incident cases. For negatives, the index time was set to one day prior to the last observed EHR record.

A.3. Control pool and analytic cohort

The control pool comprised participants with (i) no 775 MDD evidence at any time, (ii) ≥ 180 days of valid 776 Fitbit data, and (iii) EHR activity both before and 777 after Fitbit start, producing 2,701 negatives. These 778 criteria formed an initial analytic cohort of 3,049 par-779 ticipants (11.4% positives).

A.4. Wearable/EHR quality filters and final dataset

We applied pre-specified quality control (QC) on wearable data density/regularity and device sanity checks, and excluded participants failing minimal EHR coverage around index. After QC, 2,430 participants remained (12.6% positives).

A.5. Holdout protocol and leakage safeguards

To avoid leakage when comparing to an external 789 CEHR model trained on other AoU participants, we restricted the held-out test set to participants not 791

present in CEHR training (487 participants; 14.8% positives) and used the remainder for model development (1,943; 12.1% positives). All feature extraction windows (EHR and wearable) were strictly pre-index (no look-ahead).

794

795

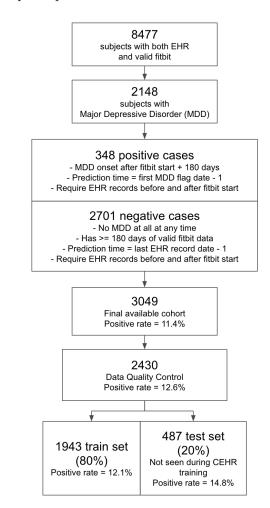
796

797

799

800

Summary of counts. $8,477 \rightarrow 3,049$ (phenotyping and design constraints) $\rightarrow 2,430$ (data quality control) split into 1,943 development and 487 held-out test participants.



Supplementary Figure 1: Example cohort definition steps and numbers for Major Depressive Disorder

Appendix B. More results

Abbreviations used in the tables: To optimize space utilization in the performance comparison tables, we employ standardized abbreviations for health outcomes and methodological components. Health outcomes are abbreviated as follows: HLD (Hyperlipidemia), HTN (Hypertension), GERD (Gastroesophageal Reflux Disease), GAD (Generalized Anxiety Disorder), MDD (Major Depressive Disorder), OB (Obesity), SA (Sleep Apnea), T2D (Type 2 Diabetes), HF (Heart Failure), and AF (Atrial Fibrillation). Methodological abbreviations include CEHR (Clinical Element-based Health Records), OMOP (Observational Medical Outcomes Partnership), and TS Embed (Time Series Embedding). These abbreviations enable comprehensive presentation of performance metrics across all outcomes and fusion approaches while maintaining table readability and fitting within standard page constraints.

Supplementary Table 1: Performance (AUROC) of best feature fusion between EHR and wearable data across health outcomes against EHR baselines.

	HLD	HTN	GERD	GAD	MDD	ОВ	SA	T2D	HF	AF
OMOP	$0.648_{\pm 0.030}$	$0.702_{\pm 0.035}$	$0.662_{\pm0.034}$	$0.733_{\pm 0.027}$	$0.713_{\pm0.036}$	$0.691_{\pm0.034}$	$0.681_{\pm 0.042}$	$0.751_{\pm 0.064}$	$0.708_{\pm 0.055}$	$0.680_{\pm 0.057}$
CEHR	$0.691_{\pm 0.035}$	$0.722_{\pm 0.030}$	$0.669_{\pm0.033}$	$0.792_{\pm0.030}$	$0.767_{\pm 0.025}$	$0.694_{\pm0.036}$	$0.746_{\pm 0.035}$	$0.735_{\pm 0.055}$	$0.862_{\pm 0.030}$	$0.620_{\pm 0.049}$
Feature Fusion	$0.808^{(1)}_{\pm 0.027}$	$0.829_{\pm 0.029}^{(1)}$	$0.785_{\pm 0.034}^{(1)}$	$0.848_{\pm 0.028}^{(3)}$	$0.825_{\pm 0.023}^{(4)}$	$0.805_{\pm 0.032}^{(4)}$	$0.816_{\pm 0.030}^{(1)}$	$0.873_{\pm 0.038}^{(4)}$	$0.881_{\pm 0.037}^{(4)}$	$0.794_{\pm 0.053}^{(1)}$
Significance	***	***	***	***	***	***	***	***	***	***

Notes: The baseline uses two types of EHR encodings: OMOP and CEHR. Significance levels: *** p<0.001, ** p<0.05. Best fusion methods for each outcome: (1) [$\mathbf{x}_{\text{EHR-OMOP}}$, $\mathbf{x}_{\text{wear-sum}}$] (Concat); (3) [$\mathbf{x}_{\text{EHR-CEHR}}$, $\mathbf{x}_{\text{wear-sum}}$] (Concat); (4) [$\mathbf{x}_{\text{EHR-CEHR}}$, $\mathbf{x}_{\text{wear-sum}}$] (Weighted). See next table.

Supplementary Table 2: Performance comparison across all baseline models and feature fusion methods.

Supplementary Table 2. I criormance comparison across an baseline models and leavant rasion methods.										
Encoding & Fusion Methods	HLD	HTN	GERD	GAD	MDD	ОВ	$\mathbf{S}\mathbf{A}$	T2D	HF	\mathbf{AF}
X _{EHR-OMOP}	$0.648_{\pm0.030}$	$0.702_{\pm 0.035}$	$0.662_{\pm0.034}$	$0.733_{\pm0.027}$	$0.713_{\pm 0.036}$	$0.691_{\pm 0.034}$	$0.681_{\pm 0.042}$	$0.751_{\pm 0.064}$	$0.708_{\pm 0.055}$	$0.680_{\pm 0.057}$
$\mathbf{x}_{\mathrm{EHR-CEHR}}$	$0.691_{\pm 0.034}$	$0.721_{\pm 0.030}$	$0.669_{\pm0.033}$	$0.792_{\pm 0.030}$	$0.767_{\pm 0.025}$	$0.694_{\pm 0.036}$	$0.746_{\pm 0.035}$	$0.735_{\pm 0.055}$	$0.862_{\pm 0.030}$	$0.620_{\pm 0.049}$
$\mathbf{X}_{ ext{wear-sum}}$	$0.804_{\pm0.029}$	$0.807_{\pm0.031}$	$0.767_{\pm 0.033}$	$0.833_{\pm 0.028}$	$0.797_{\pm 0.026}$	$0.762_{\pm 0.038}$	$0.741_{\pm 0.036}$	$0.790_{\pm 0.058}$	$0.778_{\pm 0.059}$	$0.791_{\pm 0.053}$
$\mathbf{X}_{ ext{wear-ts}}$	$0.753_{\pm0.034}$	$0.812_{\pm 0.026}$	$0.792_{\pm 0.029}$	$0.745_{\pm 0.030}$	$0.671_{\pm 0.034}$	$0.695_{\pm0.041}$	$0.687_{\pm 0.046}$	$0.774_{\pm 0.062}$	$0.742 _{\pm 0.068}$	$0.764_{\pm 0.058}$
(1) $[\mathbf{x}_{\text{EHR-OMOP}}, \mathbf{x}_{\text{wear-sum}}]$ (Concat)	$0.808_{\pm0.027}$	$0.829_{\pm 0.029}$	$0.785_{\pm0.034}$	$0.842_{\pm 0.025}$	$0.815_{\pm0.020}$	$0.793_{\pm0.032}$	$0.816_{\pm0.030}$	$0.816_{\pm 0.045}$	$0.831_{\pm 0.046}$	$0.794_{\pm 0.053}$
(2) $[\mathbf{x}_{EHR\text{-}OMOP}, \mathbf{x}_{wear\text{-}ts}]$ (Concat)	$0.708_{\pm 0.035}$	$0.782_{\pm 0.026}$	$0.741_{\pm 0.031}$	$0.805_{\pm0.026}$	$0.751_{\pm 0.029}$	$0.752_{\pm 0.035}$	$0.754_{\pm 0.032}$	$0.578_{\pm 0.078}$	$0.768_{\pm 0.066}$	$0.763_{\pm 0.040}$
(3) $[\mathbf{x}_{EHR-CEHR}, \mathbf{x}_{wear-sum}]$ (Concat)	$0.797_{\pm0.031}$	$0.824_{\pm 0.024}$	$0.780_{\pm 0.030}$	$0.848_{\pm 0.028}$	$0.823_{\pm 0.023}$	$0.804_{\pm 0.032}$	$0.797_{\pm 0.028}$	$0.819_{\pm 0.045}$	$0.880_{\pm 0.036}$	$0.783_{\pm 0.047}$
(4) $[\mathbf{x}_{\text{EHR-CEHR}}, \mathbf{x}_{\text{wear-sum}}]$ (Weighted)	$0.798_{\pm0.031}$	$0.818_{\pm 0.026}$	$0.778_{\pm 0.030}$	$0.848_{\pm 0.028}$	$0.825_{\pm 0.023}$	$0.805_{\pm 0.032}$	$0.801_{\pm 0.029}$	$0.873_{\pm 0.038}$	$0.881_{\pm 0.037}$	$0.785_{\pm 0.047}$
(5) $[\mathbf{x}_{\text{EHR-CEHR}}, \mathbf{x}_{\text{wear-ts}}]$ (Concat)	$0.789_{\pm 0.036}$	$0.815_{\pm 0.025}$	$0.785_{\pm 0.029}$	$0.805_{\pm0.030}$	$0.781_{\pm 0.024}$	$0.724_{\pm 0.037}$	$0.781_{\pm 0.029}$	$0.785_{\pm 0.061}$	$0.812_{\pm 0.050}$	$0.749_{\pm 0.047}$
(6) $[\mathbf{x}_{EHR-CEHR}, \mathbf{x}_{wear-ts}]$ (Weighted)	$0.789_{\pm 0.036}$	$0.815_{\pm 0.025}$	$0.785_{\pm 0.029}$	$0.805_{\pm0.030}$	$0.781_{\pm 0.025}$	$0.724_{\pm 0.037}$	$0.781_{\pm 0.029}$	$0.785_{\pm 0.061}$	$0.812_{\pm 0.050}$	$0.749_{\pm 0.047}$

As shown in Supplementary Table 1, the integration of wearable data produced strong improvements for chronic cardiovascular and metabolic conditions. For hyperlipidemia, the [$\mathbf{x}_{\text{EHR-OMOP}}$, $\mathbf{x}_{\text{wear-sum}}$] (Concat) approach achieved an AUROC of 0.808, representing a +0.117 improvement over the CEHR-only baseline of 0.691 ($p=1.6\times10^{-53}$). Similarly, hypertension prediction improved from a CEHR baseline of 0.721 to 0.829 using the same [$\mathbf{x}_{\text{EHR-OMOP}}$, $\mathbf{x}_{\text{wear-sum}}$] (Concat), yielding a +0.107 improvement ($p=9.3\times10^{-52}$). Obesity prediction demonstrated remarkable improvement through [$\mathbf{x}_{\text{EHR-CEHR}}$, $\mathbf{x}_{\text{wear-sum}}$] (Weighted), rising from a CEHR baseline of 0.694 to 0.805 (+0.111, $p=3.3\times10^{-61}$). Type II diabetes achieved the highest final AUROC of 0.873 through [$\mathbf{x}_{\text{EHR-CEHR}}$, $\mathbf{x}_{\text{wear-sum}}$] (Weighted), improving from an OMOP baseline of 0.751 (+0.122, $p=2.0\times10^{-31}$).

Gastroesophageal reflux disease showed substantial improvement as well, with [$\mathbf{x}_{\text{EHR-OMOP}}$, $\mathbf{x}_{\text{wear-sum}}$] (Concat) increasing performance from a baseline CEHR AUROC of 0.669 to 0.785, representing a +0.116 gain ($p = 2.4 \times 10^{-50}$).

Aouwearable

Mental health conditions also benefited significantly from wearable integration. Generalized anxiety disorder prediction improved from a strong CEHR baseline of 0.792 to 0.848 using [$\mathbf{x}_{\text{EHR-CEHR}}$, $\mathbf{x}_{\text{wear-sum}}$] (Concat), achieving a +0.056 improvement ($p = 6.3 \times 10^{-54}$), while major depressive disorder showed similar gains from 0.767 to 0.825 with [$\mathbf{x}_{\text{EHR-CEHR}}$, $\mathbf{x}_{\text{wear-sum}}$] (Weighted) (+0.058, $p = 4.3 \times 10^{-48}$).

Sleep apnea showed consistent improvement, rising from a CEHR baseline of 0.746 to 0.816 using $[\mathbf{x}_{\text{EHR-OMOP}}, \mathbf{x}_{\text{wear-sum}}]$ (Concat) (+0.070, $p = 3.6 \times 10^{-31}$).

Even outcomes with already strong baseline performance showed meaningful improvements. Heart failure, which had the highest baseline CEHR performance at 0.862, still achieved a statistically significant improvement to 0.881 using [$\mathbf{x}_{\text{EHR-CEHR}}$, $\mathbf{x}_{\text{wear-sum}}$] (Weighted) (+0.019, $p = 7.5 \times 10^{-12}$). Atrial fibrillation demonstrated substantial improvement, rising from an OMOP baseline of 0.680 to 0.794 with [$\mathbf{x}_{\text{EHR-OMOP}}$, $\mathbf{x}_{\text{wear-sum}}$] (Concat), representing a +0.114 gain ($p = 1.4 \times 10^{-26}$).