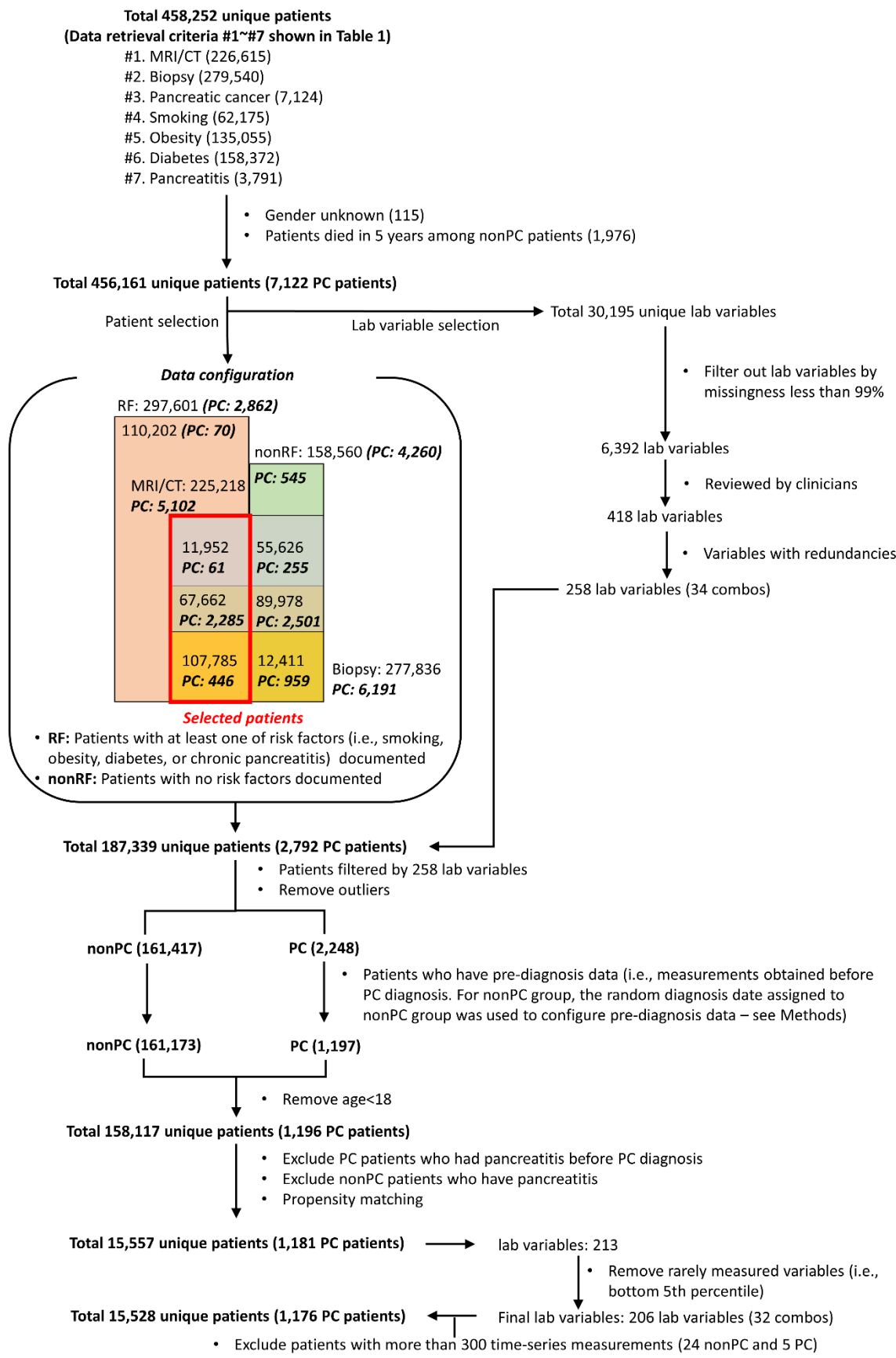


**Patterns, Volume 4**

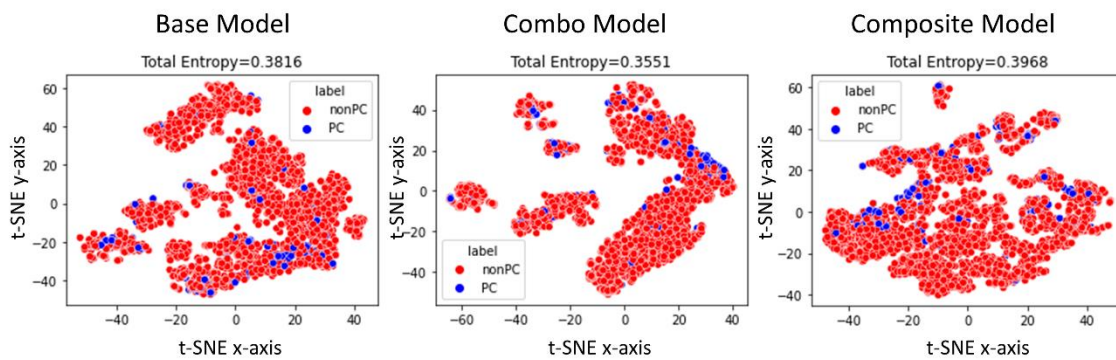
**Supplemental information**

**Structured deep embedding model to generate  
composite clinical indices from electronic health  
records for early detection of pancreatic cancer**

**Jiheum Park, Michael G. Artin, Kate E. Lee, Benjamin L. May, Michael Park, Chin Hur, and Nicholas P. Tatonetti**



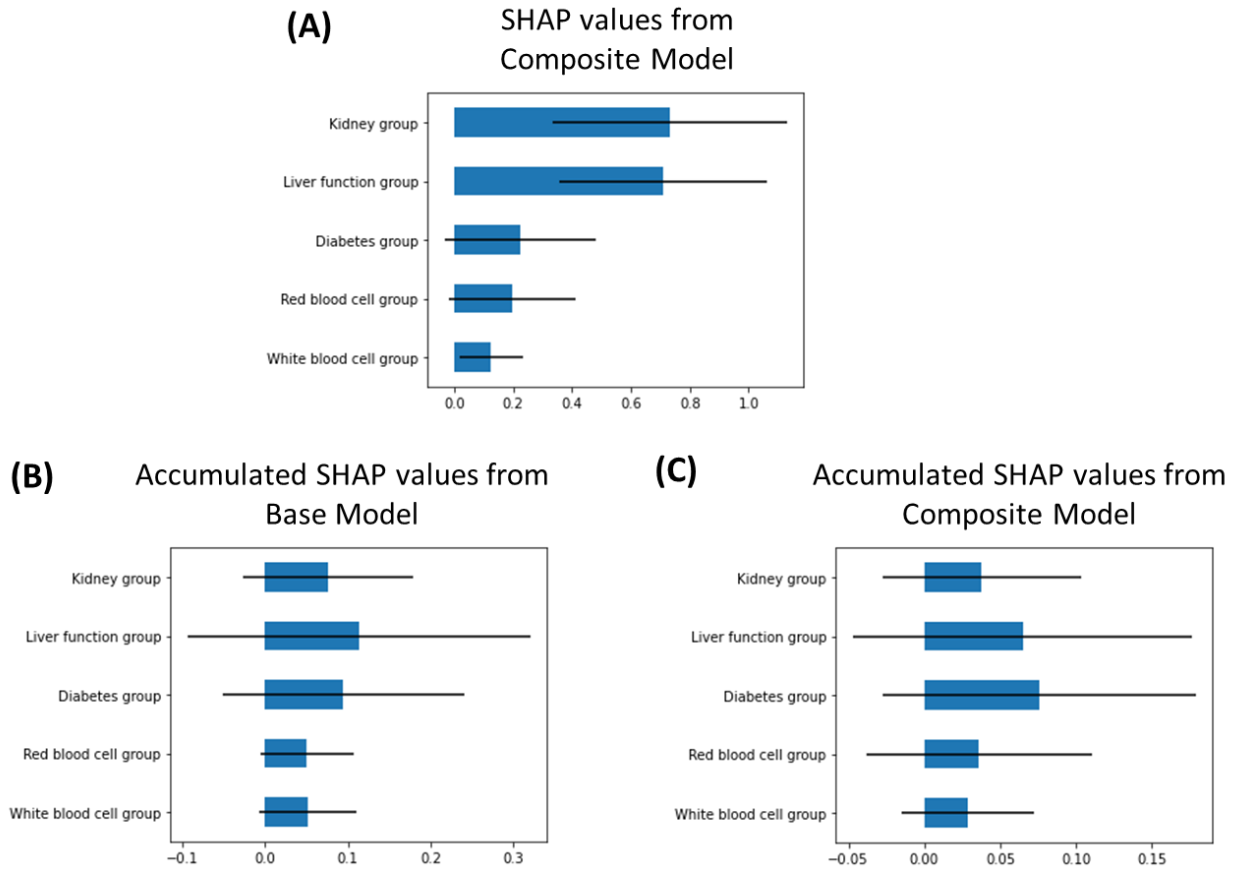
**Fig S1. Data preprocessing flow chart.** We obtained 458,252 patient samples with 30,195 lab variables from New York-Presbyterian Hospital (NYP)/Columbia University Medical Center (CUMC) EHR data. We focused on high-risk population for PC (i.e., red box: selected patients), composed of the patient group who has one of the four risk factors (i.e., smoking, obesity, diabetes, or chronic pancreatitis) documented and also received either imaging or biopsy. This selected patients' data processed into the final dataset is composed of 206 lab variables in total and 15,528 patients where 1,176 are PC patients.



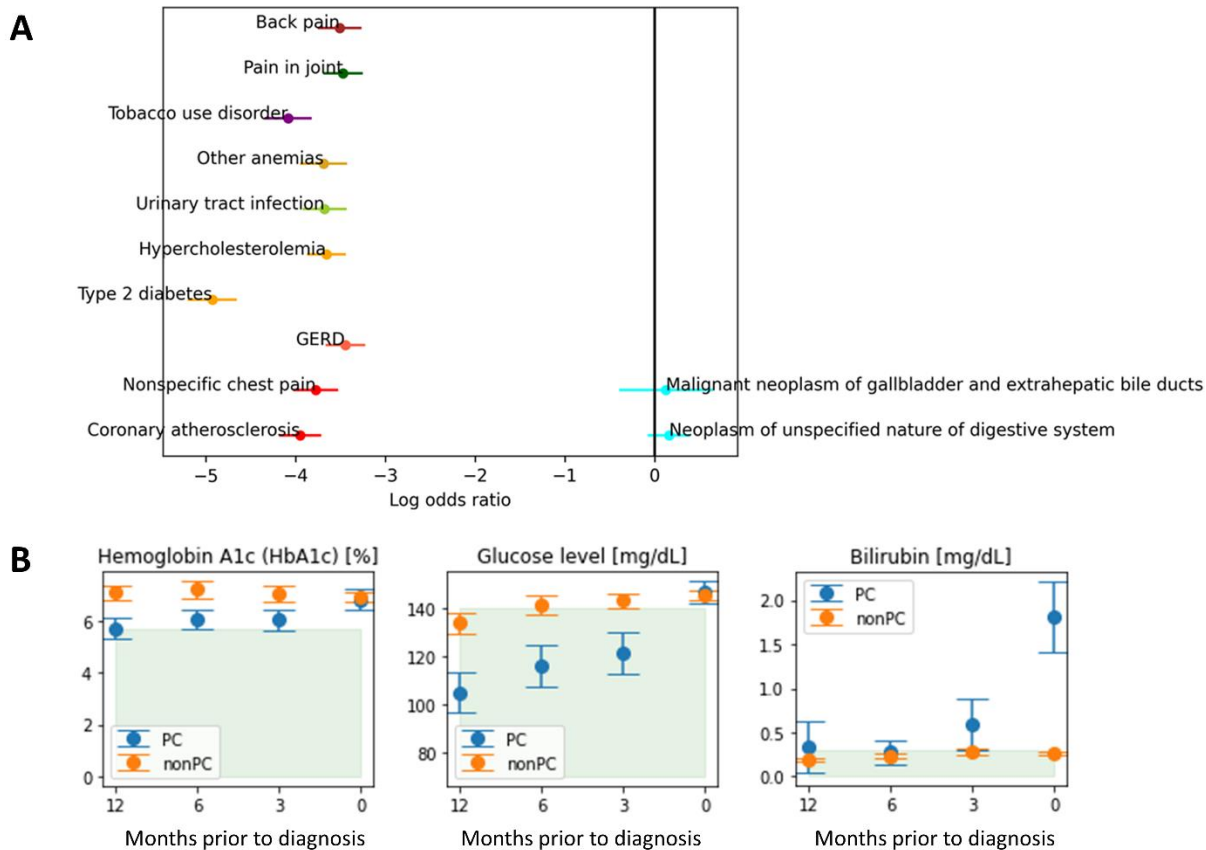
**Fig S2. Cluster analysis.** We evaluated clusters created by the resulting embeddings from each model (i.e., 206 embeddings from the Base Model, 32 embeddings from the Combo Model, 5 embeddings from Composite Modelg1, 3 embeddings from Composite Modelg2, and 7 embeddings from Composite Modelg3).



**Fig S3. Correlation matrix of 32 embeddings from Combo Model.** We filtered the correlation matrix by absolute value of correlation coefficients greater than 0.3 (A) and 0.4 (B) to come up with Composite Model<sub>g2</sub> and Composite Model<sub>g3</sub> respectively. We then bundled combo variables with ones that were correlated to each other (Table 1). The remaining combo variables that were not correlated with any other ones were bundled into “comp3” and “comp7” in grouping strategy 2 and 3 respectively (Table 1).



**Fig S4. Accumulated feature importance test.** (A) SHAP values of 5 composite indices from the Composite Model. Accumulated SHAP values of 5 composite variables from (B) the Base Model and (C) the Composite Model by grouping the SHAP values of 206 individual variables into 5 composite indices.



**Fig S5. Phenome-wide association study (PheWAS) results.** (A) Log Odds Ratio (LOR) plot where top 10 PheCodes resulted in negative LOR and all PheCodes resulted in positive LOR are shown. (B) Temporal changes in time at 0, 3, 6, and 12 months prior to diagnosis. Green shade area indicates normal ranges.



**Table S1. Baseline characteristics.** This table shows brief baseline characteristics for the final dataset used in the analysis. A full demographics include 7 categories of race, 8 categories of ethnicity, 66 categories of language, and 103 categories of zip codes, which are not shown in this table.

		<b>PC/nonPC</b>	
<b>Total</b>		<b>1176 (8%)/14,352 (92%)</b>	
Risk factors	Smoking	Yes	215 (18%)/2,670 (19%)
		Not documented	961 (82%)/11,682 (81%)
	Obesity	Yes	235 (20%)/2,944 (21%)
		Not documented	941 (80%)/11,408 (79%)
	Diabetes	Yes	880 (75%)/11,098 (77%)
		Not documented	296 (25%)/3,254a (23%)
Demographics	Race	White	543 (46%)/6,284 (44%)
		Asian	36 (3%)/368 (3%)
		African American	144 (12%)/1,882(13%)
		Other Combinations not described	103 (10%)/1,451 (10%)
		Unknown	344 (29%)/4,288 (30%)
	Ethnicity	Caucasian	21 (2%)/280 (2%)
		Hispanic	9 (1%)/68 (1%)
		Not Hispanic	240 (20%)/2,453 (17%)
		African American	124 (11%)/1,517 (10%)
		Unknown	778 (66%)/9,981 (70%)
	Sex	Male	631 (54%)/7,644 (53%)
		Female	545 (46%)/6,708 (47%)
	Zip code	Starts with 0 (MA, NH, ME, VT, CT, NJ)	186 (16%)/1,904 (13%)
		Starts with 1 (NY, PA)	958 (82%)/12,029 (85%)
		Starts with 3 (GA, FL, AL, TN, MS)	20 (2%)/215 (2%)

Language	English	637 (55%)/7,440 (53%)
	Spanish	103 (9%)/1,520 (11%)
	Other	311 (27%)/4,109 (29%)
	Unknown	105 (8%)/1,035 (7%)

Age 73.9 (CI95%=73.2-74.6)/74.5 (CI95%=74.3-74.7)

**Table S3. Performance comparison of model results** We performed 10 repetitive experiments for each model by randomly splitting the dataset into train set (80%) and test set (20%), and presented mean AUROC and AUPRC with 95% confidence intervals.

Prediction model	Train set		Test set	
	AUROC	AUPRC	AUROC	AUPRC
Base Model	$0.873 \pm 0.004$	$0.473 \pm 0.010$	$0.846 \pm 0.008$	$0.410 \pm 0.020$
Combo Model	$0.888 \pm 0.005$	$0.524 \pm 0.017$	$0.855 \pm 0.010$	$0.436 \pm 0.022$
Composite Model <sub>g1</sub>	$0.893 \pm 0.004$	$0.538 \pm 0.009$	$0.858 \pm 0.009$	$0.435 \pm 0.033$
Composite Model <sub>g2</sub>	$0.893 \pm 0.005$	$0.539 \pm 0.020$	$0.859 \pm 0.008$	$0.444 \pm 0.025$
Composite Model <sub>g3</sub>	$0.888 \pm 0.006$	$0.523 \pm 0.018$	$0.854 \pm 0.011$	$0.432 \pm 0.029$