

Contents

Supplementary Fig. 1. Temporal trends of the lab-confirmed monthly new SARS-CoV-2 cases per 10,000 patients in the INSIGHT and OneFlorida+ cohorts, March 2020 to November 2021.	6
Supplementary Fig. 2. Stratified analysis of adjusted excess burden (adjusted excess cumulative incidence per 1,000 patients) of post-acute sequelae of SARS-CoV-2 infection (PASC) over different subgroups, the OneFlorida+ cohort, from March 2020 to November 2021. Subgroups were stratified by acute severity status, age groups, gender, race groups, and baseline pre-existing conditions. Different color panels represent different organ systems, including (from top to bottom): the nervous system, skin, respiratory system, circulatory system, and other signs. CAD, coronary artery disease; CKD, chronic kidney disease; CPD, chronic pulmonary disease; T2D, diabetes type 2; Healthy: no documented pre-existing conditions and no PASC-like symptoms at baseline. Two ICD-10 diagnosis codes B948 (sequelae of other specified infectious and parasitic diseases) and U099 (post-COVID-19 condition, unspecified) were also used to compare general post-acute sequelae of SARS-CoV-2 infection in different groups. The conditions with their aHRs' P-value < 8.39×10^{-5} (the Bonferroni-corrected significance threshold) were highlighted in red squares. The fraction of the subgroup population was shown at the top.....	7
Supplementary Fig. 3. Comparison of Adjusted hazard ratio of identified incident medications for PASC in the OneFlorida+ cohort versus INSIGHT. The 95% confidence intervals (CI) were reported. The aHR's P-value < 8.39×10^{-5} (the Bonferroni-corrected threshold) was used for selecting significant medications. The conditions also replicated in the INSIGHT were marked by ‡ symbols. The color panels represent different organ systems, including (from top to bottom): the respiratory system, endocrine and metabolic, and other general signs. The aHR and its P-value were calculated by the Cox proportional hazard model and the Wald Chi-Square test.	8
Supplementary Fig. 4. Adjusted cumulative incidence (per 1,000 patients) of post-acute sequelae of SARS-CoV-2 infection (PASC) over different subgroups, the INSIGHT cohort, SARS-CoV-2 infected patients, from March 2020 to November 2021. Subgroups were stratified by their acute severity status, age groups, gender, race groups, and baseline pre-existing conditions. Different color panels represent different organ systems. CAD, coronary artery disease; CKD, chronic kidney disease; CPD, chronic pulmonary disease; T2D, diabetes type 2; Healthy: no documented pre-existing conditions and no PASC-like symptoms at baseline. Two ICD-10 diagnosis codes B948 (sequelae of other specified infectious and parasitic diseases) and U099 (post-COVID-19 condition, unspecified) were also used to compare general post-acute sequelae of SARS-CoV-2 infection in different groups. The PASC conditions also identified in OneFlorida+ were marked by ‡ symbols. The conditions with their aHRs' P-value < 8.39×10^{-5} (the Bonferroni-corrected significance threshold) were highlighted in red squares. The fraction of the subgroup population was shown at the top.	10

Supplementary Fig. 5. Adjusted cumulative incidence (per 1,000 patients) of post-acute sequelae of SARS-CoV-2 infection (PASC) over different subgroups, the OneFlorida+ cohort, SARS-CoV-2 infected patients, from March 2020 to November 2021. Subgroups were stratified by their acute severity status, age groups, gender, race groups, and baseline pre-existing conditions. Different color panels represent different organ systems, including (from top to bottom): the nervous system, skin, respiratory system, circulatory system, and general signs. CAD, coronary artery disease; CKD, chronic kidney disease; CPD, chronic pulmonary disease; T2D, diabetes type 2; Healthy: no documented pre-existing conditions and no PASC-like symptoms at baseline. Two ICD-10 diagnosis codes B948 (sequelae of other specified infectious and parasitic diseases) and U099 (post-COVID-19 condition, unspecified) were also used to compare general post-acute sequelae of SARS-CoV-2 infection in different groups. The conditions with their aHRs' P-value < 8.39×10^{-5} (the Bonferroni-corrected significance threshold) were highlighted in red squares. The fraction of the subgroup population was shown at the top.11

Supplementary Fig. 6. Stratified analysis of adjusted excess burden (adjusted excess cumulative incidence per 1,000 patients) of post-acute sequelae of SARS-CoV-2 infection (PASC) over different periods, from March 2020 to November 2021. The ancestral strain wave was defined from March 1, 2020, to September 30, 2020, the Delta wave was defined from June 1, 2021, and November 30, 2021; and from October 1, 2020, to May 31, 2021, were Alpha and other variants. Different color panels represent different organ systems, including (from top to bottom): the nervous system, skin, respiratory system, circulatory system, and general signs. CAD, coronary artery disease; CKD, chronic kidney disease; CPD, chronic pulmonary disease; T2D, diabetes type 2; Healthy: no documented pre-existing conditions and no PASC-like symptoms at baseline. Two ICD-10 diagnosis codes B948 (sequelae of other specified infectious and parasitic diseases) and U099 (post-COVID-19 condition, unspecified) were also used to compare general post-acute sequelae of SARS-CoV-2 infection in different groups. The conditions with their aHRs' P-value < 8.39×10^{-5} (the Bonferroni-corrected significance threshold) were highlighted in red squares.13

Supplementary Fig. 7. Additional post-acute sequelae of SARS-CoV-2 risks in the INSIGHT cohort versus in the OneFlorida+ cohort by using the Benjamini-Yekutieli method to control for false discovery rate, from March 2020 to November 2021. Incidence risk measured by adjusted hazard ratios and 95% confidence intervals (CI) were reported in the main panel. The adjusted cumulative incidences (CIF) per 1,000 patients in both the SARS-CoV-2 positive group and the negative group were also reported. The PASC conditions identified in both datasets were marked by ‡ symbols. The PASC outcomes were ascertained from day 30 after the SARS-CoV-2 infection and all the adjusted risk measures were computed 180 days after the SARS-CoV-2 infection. The threshold of the false discovery rate in the Benjamini-Yekutieli method is 0.05. GERD, gastroesophageal reflux disease.14

Supplementary Fig. 8. The post-acute sequelae of SARS-CoV-2 risks when adjusting for additional baseline vaccination status or not, the INSIGHT and OneFlorida+ cohorts, from March 2020 to November 2021. The baseline vaccination covariates include fully vaccinated, partially vaccinated, and no evidence of vaccination. We defined the fully

vaccinated as two shots of mRNA vaccine (Pfizer or Moderna) or one shot of J&J, see <https://www.cdc.gov/coronavirus/2019-ncov/vaccines/stay-up-to-date.html>. Adjusted hazard ratios (aHR) were compared when adjusting for these additional covariates (aHR w/ vax) versus not adjusting for these covariates in our primary analysis (aHR w/o vax). The PASC conditions identified in both datasets were marked by ‡ symbols. The 95% confidence intervals (CI) of aHR were reported. The color panels represent different organ systems, including (from top to bottom): the nervous system, skin, respiratory system, circulatory system, endocrine and metabolic, digestive system, genitourinary system, and other general signs. The PASC outcomes were ascertained from day 30 after the SARS-CoV-2 infection and all the adjusted risk measures were computed 180 days after the SARS-CoV-2 infection.

.....16

Supplementary Fig. 9. The post-acute sequelae of SARS-CoV-2 risks when adjusting for index day through Cubic B-Spline versus Categorical variables, the INSIGHT and OneFlorida+ cohorts, from March 2020 to November 2021.

The index days since March 1st, 2020, were modeled by a) Cubic B-spline with five knots (namely, 7 spline basis functions of polynomial order 3) and b) categorizing days into five periods (March 2020 – June 2020, July 2020 – October 2020, November 2020 - February 2021, March 2021 – June 2021, July 2021 – November 2021). Adjusted hazard ratios (aHR) were compared when using the cubic B-spline method (aHR spline) versus the categorial method in the primary analysis (aHR category). The PASC conditions identified in both datasets were marked by ‡ symbols. The 95% confidence intervals (CI) of aHR were reported. The color panels represent different organ systems, including (from top to bottom): the nervous system, skin, respiratory system, circulatory system, endocrine and metabolic, digestive system, genitourinary system, and other general signs. The PASC outcomes were ascertained from day 30 after the SARS-CoV-2 infection and all the adjusted risk measures were computed 180 days after the SARS-CoV-2 infection.

.....18

Supplementary Fig. 10. The post-acute sequelae of SARS-CoV-2 risks when using non-linear PS modeling versus linear PS modeling, the INSIGHT and OneFlorida+ cohorts, from March 2020 to November 2021.

Adjusted hazard ratios (aHR) were adjusted for a nonlinear PS modeling using the gradient boosting machine with decision tree base learners (aHR nonlinear) versus the linear PS modeling using regularized logistic regression in our primary analysis (aHR lin). The PASC conditions identified in both datasets were marked by ‡ symbols. The color panels represent different organ systems, including (from top to bottom): the nervous system, skin, respiratory system, circulatory system, endocrine and metabolic, digestive system, genitourinary system, and other general signs. The PASC outcomes were ascertained from day 30 after the SARS-CoV-2 infection and all the adjusted risk measures were computed 180 days after the SARS-CoV-2 infection. The 95% confidence intervals (CI) of aHR were reported. PS, propensity score.

.....20

Supplementary Fig. 11. The post-acute sequelae of SARS-CoV-2 risks when downsampling the INSIGHT cohort into 22,341 SARS-CoV-2 positive patients and 177,010 negative patients versus using original data, the INSIGHT and OneFlorida+ cohorts, from March 2020 to November 2021.

The number of SARS-CoV2 positive (or negative) patients in the INSIGHT after downsampling is the same as the number of positive

(or negative) patients in the OneFlorida+. Adjusted hazard ratios (aHR) and 95% confidence intervals (CI) were compared using the downsampled cohorts (aHR sample) versus the original cohorts in our primary analysis (aHR all) for the INSIGHT cohort. The PASC conditions identified in both datasets were marked by ‡ symbols. The color panels represent different organ systems, including (from top to bottom): the nervous system, skin, respiratory system, circulatory system, endocrine and metabolic, digestive system, genitourinary system, and other general signs. The PASC outcomes were ascertained from day 30 after the SARS-CoV-2 infection and all the adjusted risk measures were computed 180 days after the SARS-CoV-2 infection.....22

Supplementary Fig. 12. Adjusted hazard ratios of post-acute sequelae of SARS-CoV-2 infection (PASC) over different subgroups, the INSIGHT cohort, from March 2020 to November 2021. Subgroups were stratified by their acute severity status, age groups, gender, race groups, and baseline pre-existing conditions. Different color panels represent different organ systems, including (from top to bottom): the nervous system, skin, respiratory system, circulatory system, blood-forming organs, endocrine and metabolic, digestive system, genitourinary system, and general signs. CAD, coronary artery disease; CKD, chronic kidney disease; CPD, chronic pulmonary disease; T2D, diabetes type 2; Healthy: no documented pre-existing conditions and no PASC-like symptoms at baseline. The PASC conditions also identified in OneFlorida+ were marked by ‡ symbols. The conditions with their aHRs' P-value < 8.39×10^{-5} (the Bonferroni-corrected significance threshold) were highlighted in red squares. The fraction of the subgroup population was shown at the top.24

Supplementary Fig. 13. Adjusted hazard ratios of post-acute sequelae of SARS-CoV-2 infection (PASC) over different subgroups, the OneFlorida+ cohort, from March 2020 to November 2021. Subgroups were stratified by their acute severity status, age groups, gender, race groups, and baseline pre-existing conditions. Different color panels represent different organ systems, including (from top to bottom): the nervous system, skin, respiratory system, circulatory system, and general signs. CAD, coronary artery disease; CKD, chronic kidney disease; CPD, chronic pulmonary disease; T2D, diabetes type 2; Healthy: no documented pre-existing conditions and no PASC-like symptoms at baseline. The conditions with their aHRs' P-value < 8.39×10^{-5} (the Bonferroni-corrected significance threshold) were highlighted in red squares. The fraction of the subgroup population was shown at the top.25

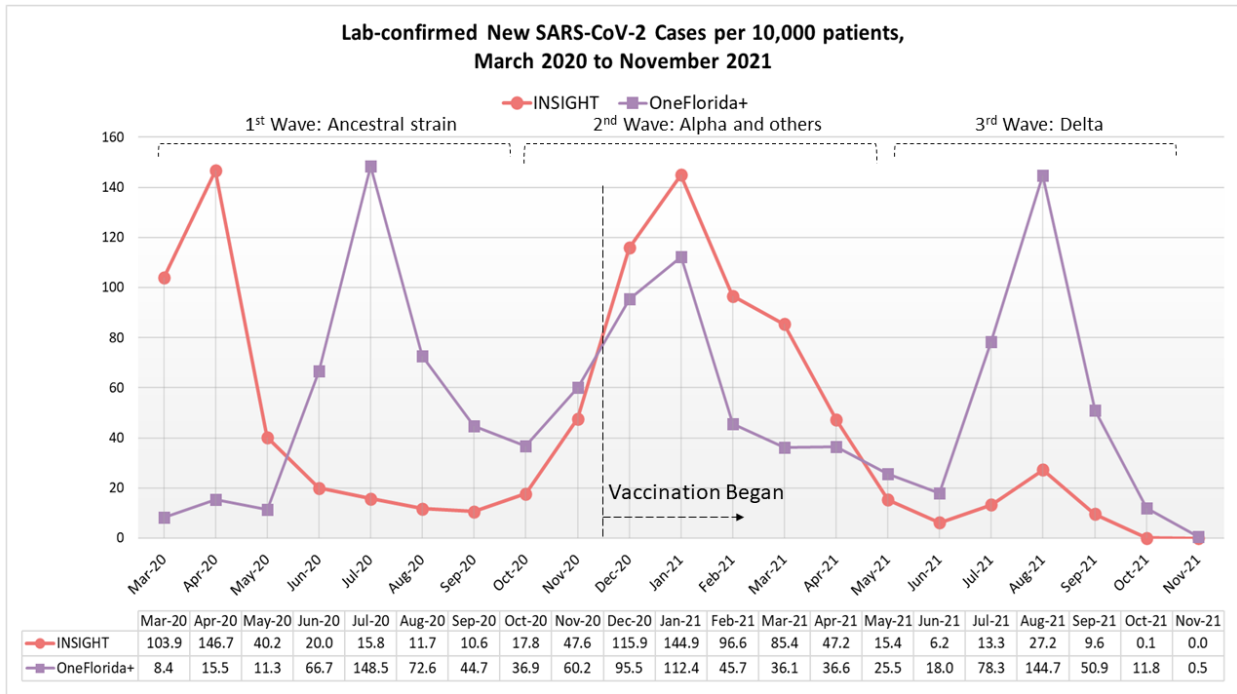
Supplementary Table 1. Results of negative outcome control in both the INSIGHT and OneFlorida+ cohorts, March 2020–November 2021.26

Supplementary Table 2. Specifications of the high-throughput screening framework for identifying potential Post-Acute Sequelae of SARS-CoV-2 infected (PASC) using the INSIGHT Electronic Health Records in New York City and OneFlorida+ in Florida (March 2020 – November 2021).....28

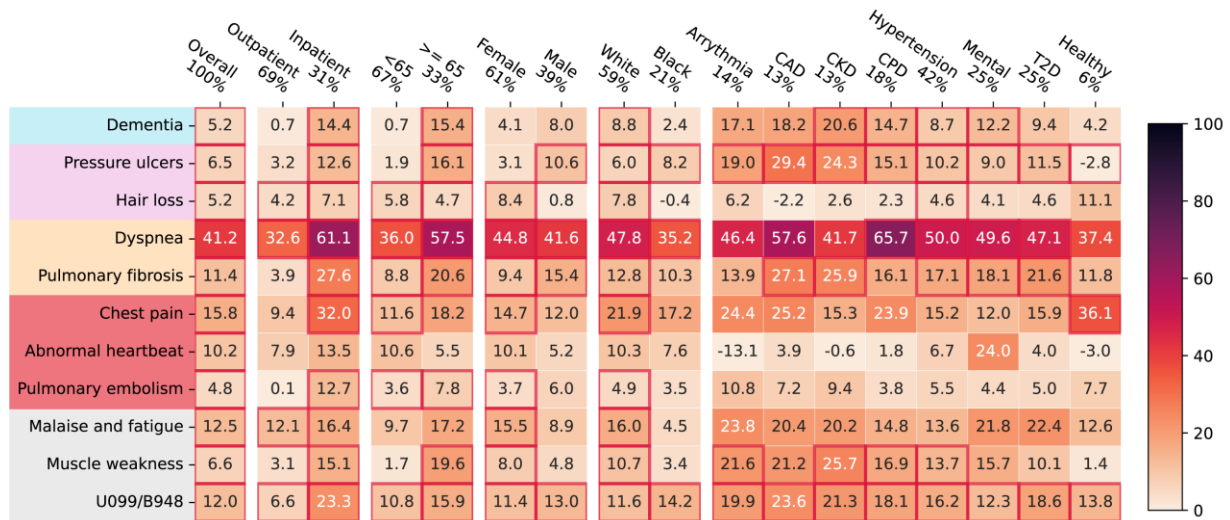
Supplementary Table 3. Cross-validation algorithm tailored for our machine learning-based propensity score calculation for each emulated trial.....31

Supplementary Table 4. Baseline SARS-CoV-2 vaccination status recorded in the two EHR databases^a32

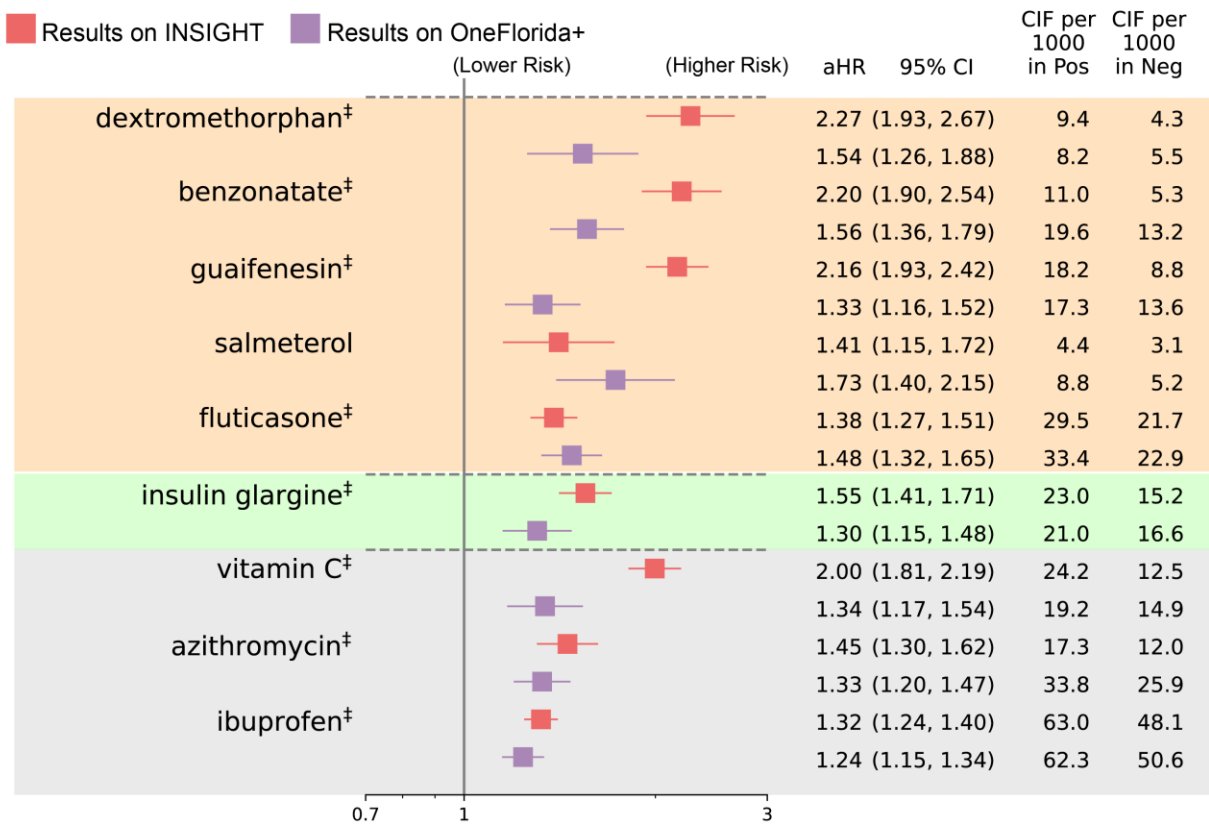
Supplementary Data 1. COVID-19 Phenotyping Lab LOINC codes and Diagnosis ICD10 codes (spreadsheet).....	33
Supplementary Data 2. PASC Adult Diagnostic List for Screening (spreadsheet).....	33
Supplementary Data 3. Baseline population characteristics with more comorbidities information, INSIGHT and OneFlorida+ cohorts, March 2020 to November 2021 (spreadsheet).	33
Supplementary Data 4. Characteristics of PASC-Specific Cohorts on INSIGHT, NYC, March 2020 to November 2021 (spreadsheet).	33
Supplementary Data 5. Characteristics of PASC-Specific Cohorts on OneFlorida, Florida, March 2020 to November 2021 (spreadsheet).....	33



Supplementary Fig. 1. Temporal trends of the lab-confirmed monthly new SARS-CoV-2 cases per 10,000 patients in the INSIGHT and OneFlorida+ cohorts, March 2020 to November 2021.



Supplementary Fig. 2. Stratified analysis of adjusted excess burden (adjusted excess cumulative incidence per 1,000 patients) of post-acute sequelae of SARS-CoV-2 infection (PASC) over different subgroups, the OneFlorida+ cohort, from March 2020 to November 2021. Subgroups were stratified by acute severity status, age groups, gender, race groups, and baseline pre-existing conditions. Different color panels represent different organ systems, including (from top to bottom): the nervous system, skin, respiratory system, circulatory system, and other signs. CAD, coronary artery disease; CKD, chronic kidney disease; CPD, chronic pulmonary disease; T2D, diabetes type 2; Healthy: no documented pre-existing conditions and no PASC-like symptoms at baseline. Two ICD-10 diagnosis codes B948 (sequelae of other specified infectious and parasitic diseases) and U099 (post-COVID-19 condition, unspecified) were also used to compare general post-acute sequelae of SARS-CoV-2 infection in different groups. The conditions with their aHRs' P-value $< 8.39 \times 10^{-5}$ (the Bonferroni-corrected significance threshold) were highlighted in red squares. The fraction of the subgroup population was shown at the top.

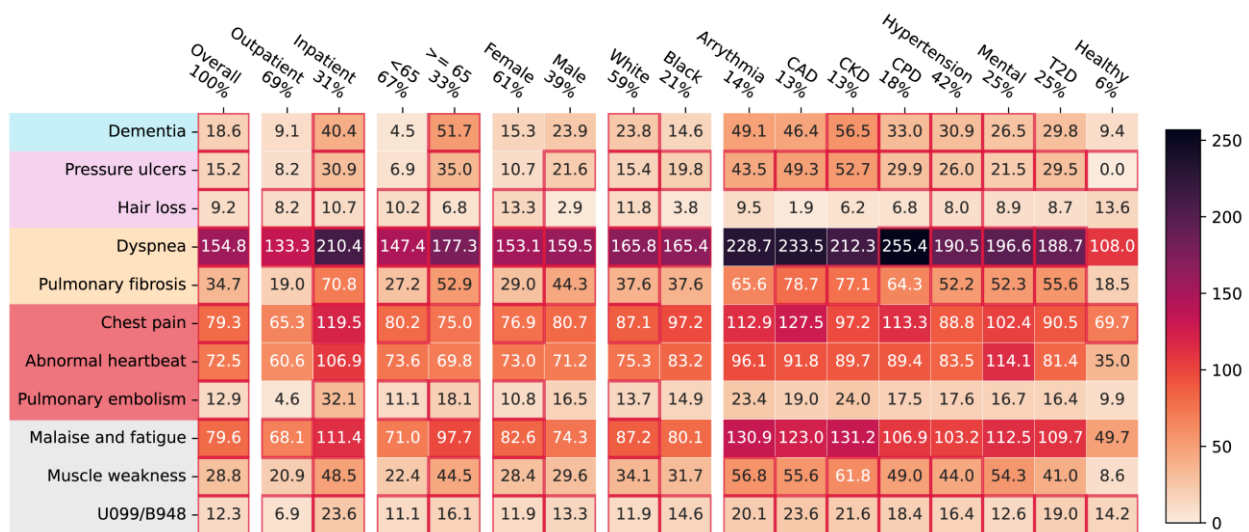


Supplementary Fig. 3. Comparison of Adjusted hazard ratio of identified incident medications for PASC in the OneFlorida+ cohort versus INSIGHT. The 95% confidence intervals (CI) were reported. The aHR's P-value $< 8.39 \times 10^{-5}$ (the Bonferroni-corrected threshold) was used for selecting significant medications. The conditions also replicated in the INSIGHT were marked by ‡ symbols. The color panels represent different organ systems, including (from top to bottom): the respiratory system, endocrine and metabolic, and other general signs. The aHR and its P-value were calculated by the Cox proportional hazard model and the Wald Chi-Square test.

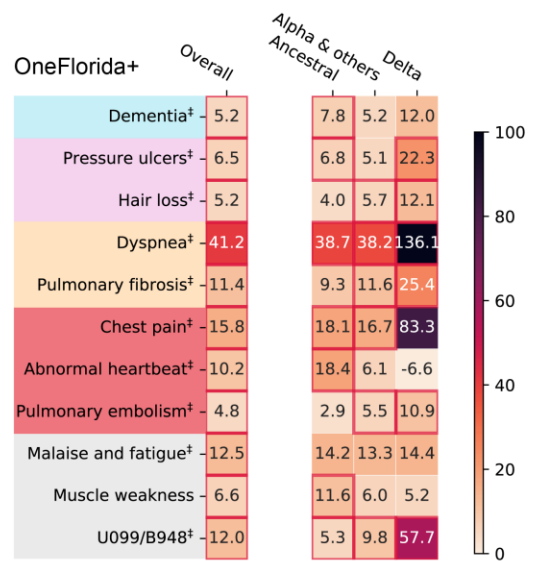
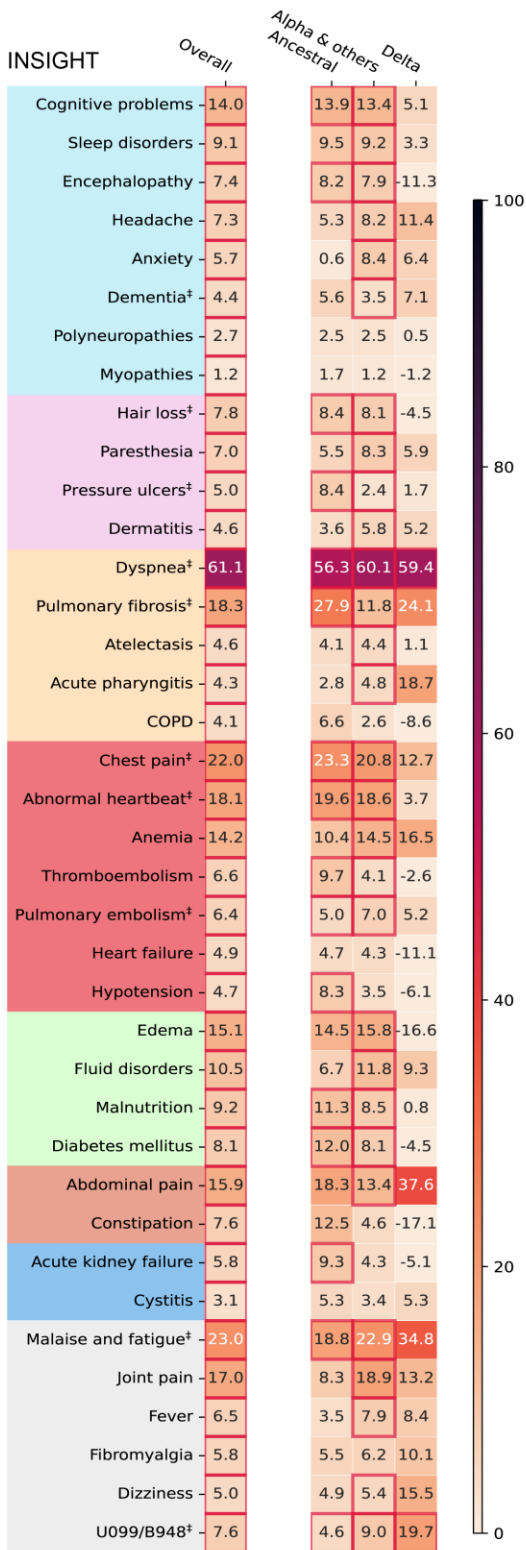
	Overall 100%	Outpatient 72%	Inpatient 28%	<65 66%	>= 65 34%	Female 60%	Male 40%	White 42%	Black 19%	Arrhythmia 11%	CAD 10%	CKD 8%	CPD 12%	Hypertension 30%	Mental 12%	T2D 17%	Healthy 9%
Cognitive problems	45.4	31.5	73.2	31.5	75.8	42.2	50.0	46.6	50.8	84.3	79.2	84.6	66.2	67.6	76.6	63.2	20.5
Sleep disorders	42.2	33.9	58.2	39.6	49.3	37.0	50.1	46.3	41.4	69.8	65.5	62.0	65.3	59.2	73.2	65.0	25.5
Encephalopathy	24.3	15.5	41.7	17.8	37.9	23.1	25.7	24.1	29.6	55.8	53.7	58.3	44.1	39.1	49.6	40.2	8.8
Headache	35.7	35.7	36.0	41.5	24.2	42.8	25.7	29.3	40.1	27.9	31.1	20.6	39.7	30.9	49.0	32.8	24.5
Anxiety	39.9	35.4	47.3	41.9	36.4	43.9	34.4	44.3	28.2	44.8	40.4	40.5	60.7	42.1	80.7	42.3	28.7
Dementia†	14.3	6.3	30.1	3.4	37.6	12.3	17.5	15.3	13.6	36.2	35.6	50.1	23.1	26.6	29.6	27.6	3.9
Polyneuropathies	10.9	8.5	15.3	10.1	12.5	9.5	12.9	11.5	13.1	22.3	15.5	20.7	16.3	15.6	14.8	16.0	2.0
Myopathies	3.0	1.2	6.1	2.9	2.9	2.4	3.9	3.3	2.0	5.7	4.2	6.7	5.1	3.5	5.0	5.1	0.8
Hair loss†	15.0	13.7	16.3	16.0	12.5	22.6	3.7	14.9	8.0	9.6	5.5	9.1	15.1	13.9	17.1	13.8	10.1
Paresthesia	49.1	47.8	52.4	50.7	44.7	51.5	46.1	45.0	55.9	51.0	44.4	58.2	72.5	55.8	69.0	58.5	31.2
Pressure ulcers†	10.4	3.2	23.6	4.0	23.6	8.2	13.3	11.3	12.5	33.6	24.0	39.4	21.6	20.7	19.6	23.4	2.1
Dermatitis	26.6	28.0	24.3	30.9	18.0	28.7	23.8	25.0	26.8	20.9	22.9	21.1	37.2	23.7	25.8	23.3	18.2
Dyspnea†	147.6	116.9	206.1	129.1	194.1	138.7	159.5	146.9	152.0	209.4	203.8	212.0	215.5	181.3	172.0	176.8	114.2
Pulmonary fibrosis†	31.0	17.7	54.3	23.6	47.0	27.5	36.0	35.1	29.8	56.7	46.3	63.1	56.0	44.2	38.6	45.8	15.6
Atelectasis	21.4	13.0	38.6	13.6	39.0	17.6	27.1	24.4	23.6	60.9	49.5	59.8	37.9	37.3	32.6	37.5	4.7
Acute pharyngitis	15.1	15.4	12.9	17.7	9.9	15.3	15.1	13.2	15.1	12.1	13.6	10.3	20.1	13.1	18.2	11.9	17.2
COPD	18.8	10.6	34.0	9.4	40.0	15.4	23.9	21.1	20.3	47.4	44.0	42.4	51.1	32.4	32.4	28.2	6.4
Chest pain†	63.8	58.7	73.3	68.7	52.8	64.4	62.7	51.3	81.1	70.2	83.8	65.6	87.6	70.8	87.7	69.4	43.1
Abnormal heartbeat†	65.5	53.8	87.4	62.2	72.6	61.3	71.1	70.2	68.0	125.2	92.9	96.3	84.1	79.1	85.2	78.8	38.1
Anemia	58.4	39.5	99.6	43.9	91.3	57.0	59.4	52.8	75.6	114.6	113.3	177.7	83.2	89.7	77.7	92.9	22.0
Thromboembolism	16.9	9.8	29.8	11.9	28.1	14.3	20.8	17.0	20.2	37.8	30.4	36.7	22.1	25.9	21.3	26.3	7.8
Pulmonary embolism†	11.6	5.3	23.3	9.6	15.9	11.2	12.4	12.7	15.3	19.3	10.6	12.2	17.0	15.1	14.1	15.4	7.6
Heart failure	26.0	14.1	51.3	11.3	58.1	19.7	34.7	27.2	30.3	107.9	94.1	92.7	46.6	47.2	31.2	45.0	9.5
Hypotension	18.1	10.0	35.5	10.0	35.9	14.5	23.1	21.7	18.3	53.5	48.2	61.1	33.9	33.3	34.6	35.3	2.7
Edema	84.7	64.0	125.8	66.0	128.4	81.6	89.3	86.9	85.1	146.3	134.7	151.7	130.8	119.4	129.2	124.4	35.4
Fluid disorders	45.4	24.5	88.8	26.9	86.8	38.4	55.3	45.8	53.1	130.1	114.4	160.5	73.3	88.4	75.4	88.8	14.5
Malnutrition	25.7	12.5	51.6	14.3	48.3	19.7	34.1	28.1	28.5	62.9	47.0	61.3	40.5	39.1	37.6	37.5	6.4
Diabetes mellitus	40.7	30.0	64.4	27.5	75.4	35.2	48.9	28.7	51.1	65.5	79.6	96.7	44.6	56.7	38.9	135.7	36.6
Abdominal pain	106.6	96.2	125.8	103.8	113.2	109.0	103.4	98.4	116.8	135.5	108.9	136.2	132.8	115.9	143.0	119.0	68.9
Constipation	47.2	35.3	69.9	39.3	63.9	47.3	47.4	48.7	46.1	83.4	70.2	80.9	65.7	67.1	73.6	69.2	14.5
Acute kidney failure	29.6	14.3	60.2	16.5	58.3	22.2	39.8	28.0	40.2	79.9	74.2	130.8	46.6	58.1	50.5	60.2	8.0
Cystitis	12.9	9.7	18.6	8.9	21.0	14.8	9.9	13.0	14.2	21.5	20.6	24.8	15.3	16.3	14.4	17.9	8.2
Malaise and fatigue†	61.1	46.6	85.6	49.3	87.1	59.2	63.1	68.7	55.4	95.9	89.7	94.7	85.8	72.4	77.2	71.4	27.1
Joint pain	110.1	106.8	116.6	108.5	114.9	113.7	105.6	96.4	126.4	112.1	120.9	124.8	149.9	126.5	159.7	134.8	73.2
Fever	20.4	14.8	30.9	17.4	25.7	17.0	25.7	21.9	21.1	31.3	31.1	37.7	27.3	25.2	29.9	25.5	8.4
Fibromyalgia	38.5	33.7	48.3	34.5	48.2	38.7	38.4	35.8	46.7	54.1	53.6	61.7	62.2	52.0	57.1	58.8	17.5
Dizziness	26.9	24.8	31.8	22.7	36.0	29.7	22.9	25.8	24.1	38.0	42.4	34.7	37.0	37.6	44.8	35.2	15.8
U099/B948†	7.8	4.0	14.4	6.2	11.1	7.3	8.4	9.2	6.8	19.1	16.8	17.3	17.0	12.1	10.7	12.2	1.7



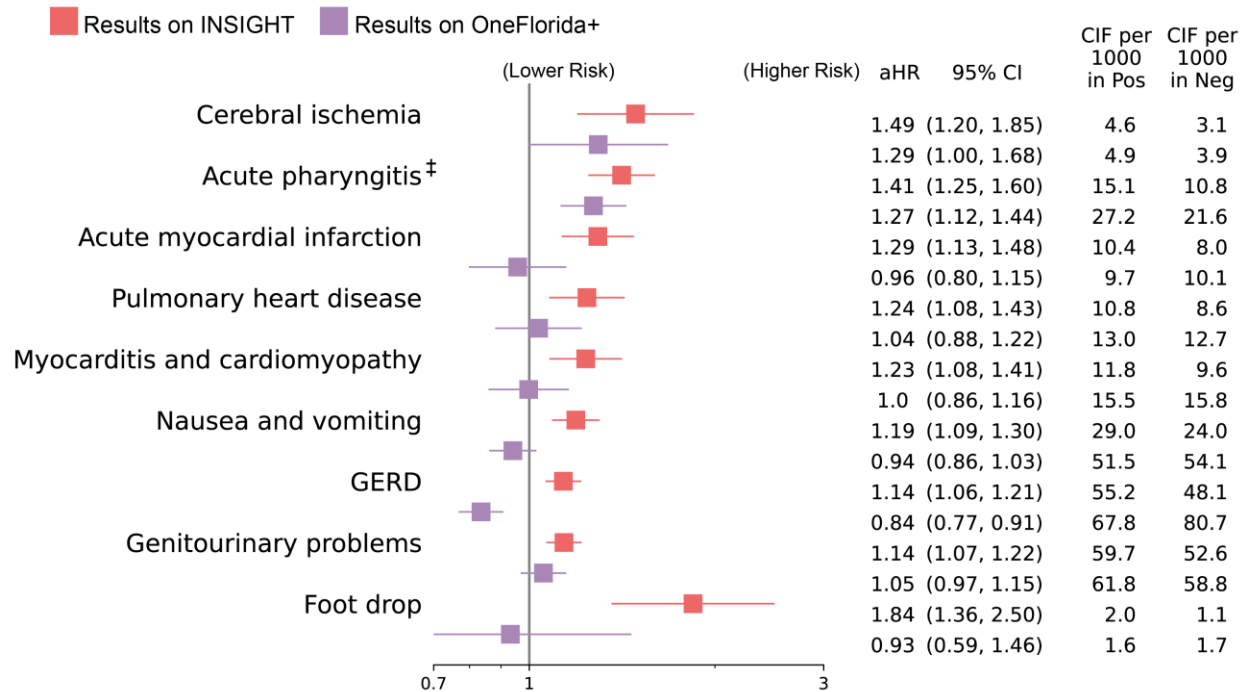
Supplementary Fig. 4. Adjusted cumulative incidence (per 1,000 patients) of post-acute sequelae of SARS-CoV-2 infection (PASC) over different subgroups, the INSIGHT cohort, SARS-CoV-2 infected patients, from March 2020 to November 2021. Subgroups were stratified by their acute severity status, age groups, gender, race groups, and baseline pre-existing conditions. Different color panels represent different organ systems. CAD, coronary artery disease; CKD, chronic kidney disease; CPD, chronic pulmonary disease; T2D, diabetes type 2; Healthy: no documented pre-existing conditions and no PASC-like symptoms at baseline. Two ICD-10 diagnosis codes B948 (sequelae of other specified infectious and parasitic diseases) and U099 (post-COVID-19 condition, unspecified) were also used to compare general post-acute sequelae of SARS-CoV-2 infection in different groups. The PASC conditions also identified in OneFlorida+ were marked by ‡ symbols. The conditions with their aHRs' P-value $< 8.39 \times 10^{-5}$ (the Bonferroni-corrected significance threshold) were highlighted in red squares. The fraction of the subgroup population was shown at the top.



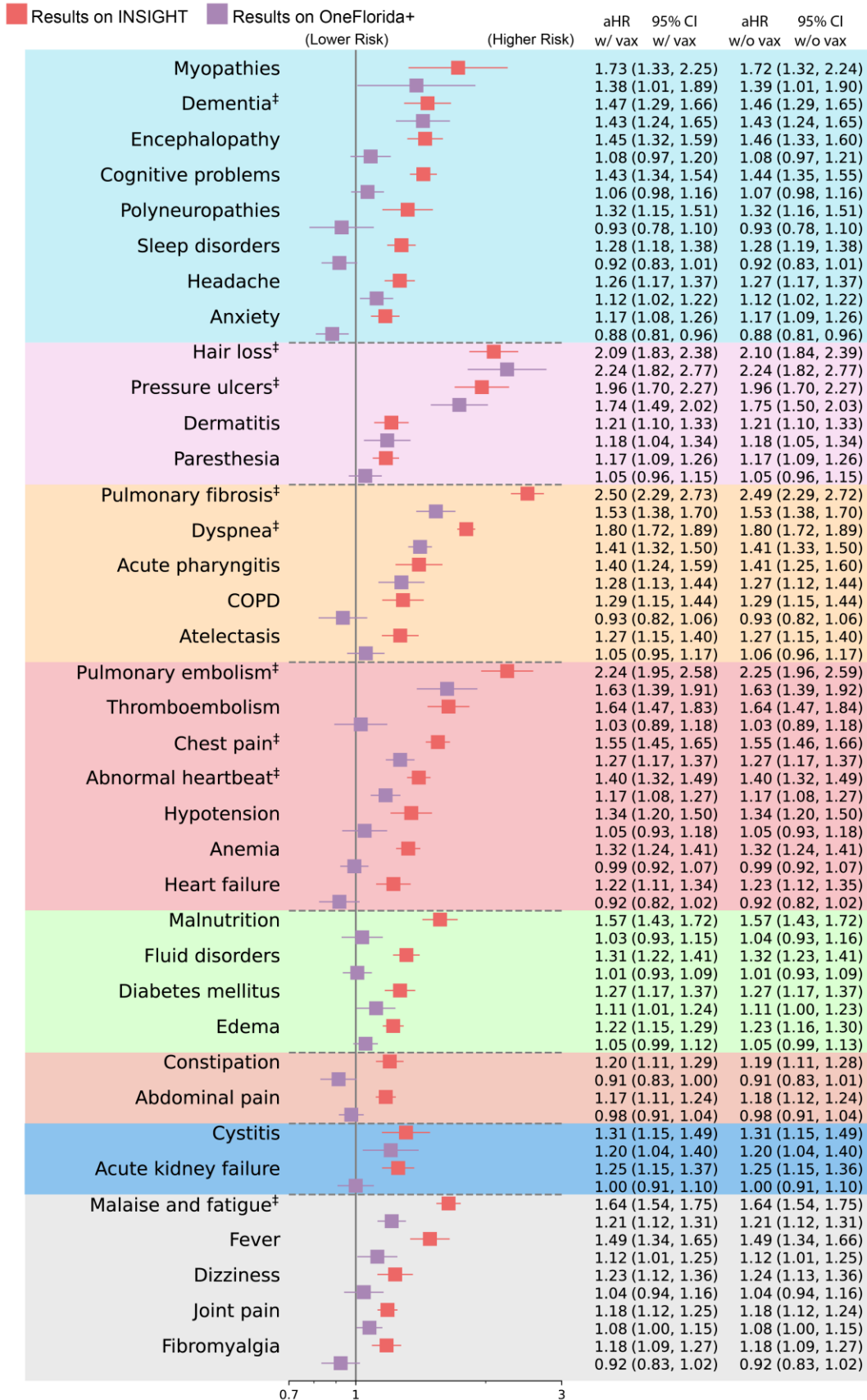
Supplementary Fig. 5. Adjusted cumulative incidence (per 1,000 patients) of post-acute sequelae of SARS-CoV-2 infection (PASC) over different subgroups, the OneFlorida+ cohort, SARS-CoV-2 infected patients, from March 2020 to November 2021. Subgroups were stratified by their acute severity status, age groups, gender, race groups, and baseline pre-existing conditions. Different color panels represent different organ systems, including (from top to bottom): the nervous system, skin, respiratory system, circulatory system, and general signs. CAD, coronary artery disease; CKD, chronic kidney disease; CPD, chronic pulmonary disease; T2D, diabetes type 2; Healthy: no documented pre-existing conditions and no PASC-like symptoms at baseline. Two ICD-10 diagnosis codes B948 (sequelae of other specified infectious and parasitic diseases) and U099 (post-COVID-19 condition, unspecified) were also used to compare general post-acute sequelae of SARS-CoV-2 infection in different groups. The conditions with their aHRs' P-value < 8.39×10^{-5} (the Bonferroni-corrected significance threshold) were highlighted in red squares. The fraction of the subgroup population was shown at the top.



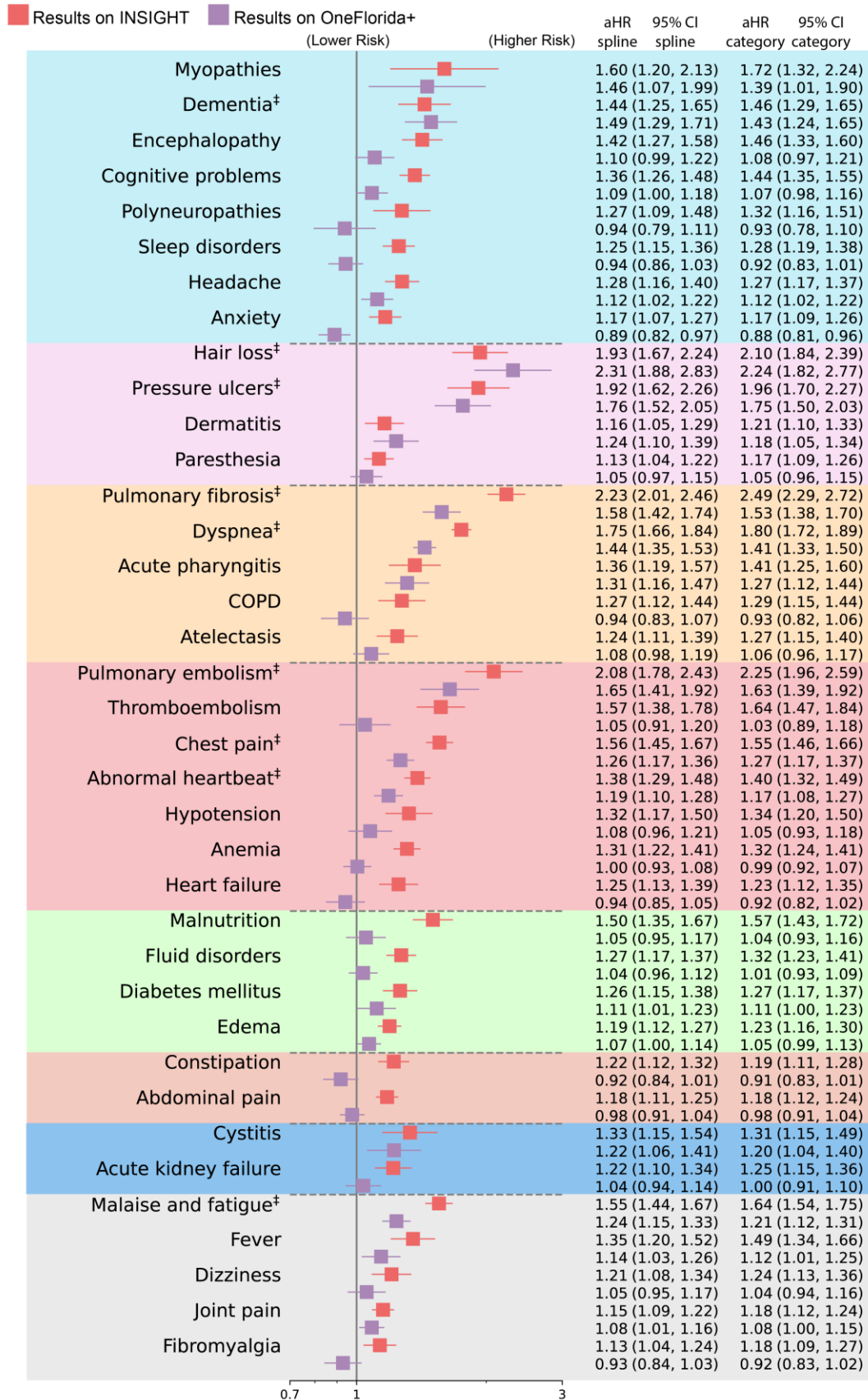
Supplementary Fig. 6. Stratified analysis of adjusted excess burden (adjusted excess cumulative incidence per 1,000 patients) of post-acute sequelae of SARS-CoV-2 infection (PASC) over different periods, from March 2020 to November 2021. The ancestral strain wave was defined from March 1, 2020, to September 30, 2020, the Delta wave was defined from June 1, 2021, and November 30, 2021; and from October 1, 2020, to May 31, 2021, were Alpha and other variants. Different color panels represent different organ systems, including (from top to bottom): the nervous system, skin, respiratory system, circulatory system, and general signs. CAD, coronary artery disease; CKD, chronic kidney disease; CPD, chronic pulmonary disease; T2D, diabetes type 2; Healthy: no documented pre-existing conditions and no PASC-like symptoms at baseline. Two ICD-10 diagnosis codes B948 (sequelae of other specified infectious and parasitic diseases) and U099 (post-COVID-19 condition, unspecified) were also used to compare general post-acute sequelae of SARS-CoV-2 infection in different groups. The conditions with their aHRs' P-value < 8.39×10^{-5} (the Bonferroni-corrected significance threshold) were highlighted in red squares.



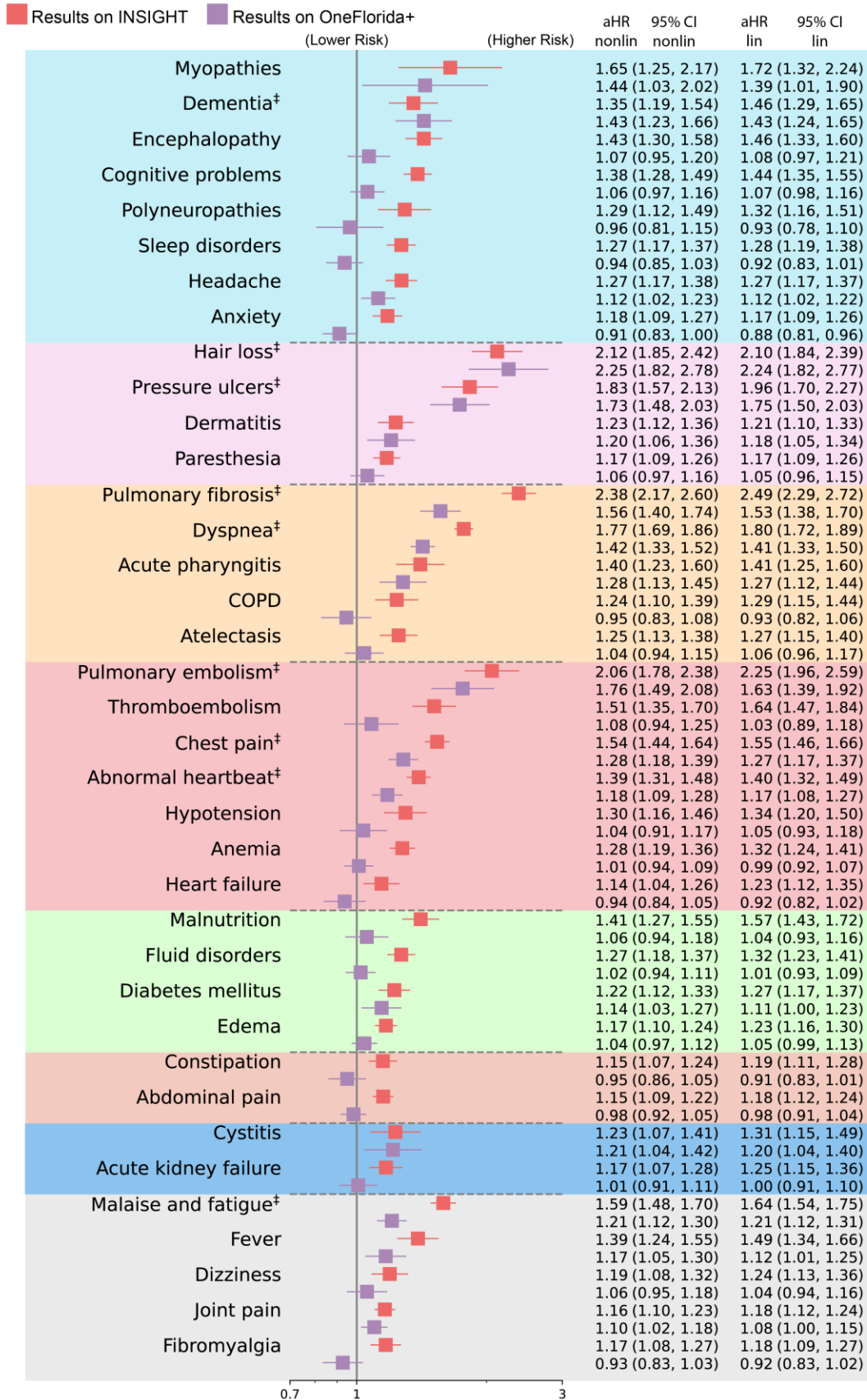
Supplementary Fig. 7. Additional post-acute sequelae of SARS-CoV-2 risks in the INSIGHT cohort versus in the OneFlorida+ cohort by using the Benjamini-Yekutieli method to control for false discovery rate, from March 2020 to November 2021. Incidence risk measured by adjusted hazard ratios and 95% confidence intervals (CI) were reported in the main panel. The adjusted cumulative incidences (CIF) per 1,000 patients in both the SARS-CoV-2 positive group and the negative group were also reported. The PASC conditions identified in both datasets were marked by ‡ symbols. The PASC outcomes were ascertained from day 30 after the SARS-CoV-2 infection and all the adjusted risk measures were computed 180 days after the SARS-CoV-2 infection. The threshold of the false discovery rate in the Benjamini-Yekutieli method is 0.05. GERD, gastroesophageal reflux disease.



Supplementary Fig. 8. The post-acute sequelae of SARS-CoV-2 risks when adjusting for additional baseline vaccination status or not, the INSIGHT and OneFlorida+ cohorts, from March 2020 to November 2021. The baseline vaccination covariates include fully vaccinated, partially vaccinated, and no evidence of vaccination. We defined the fully vaccinated as two shots of mRNA vaccine (Pfizer or Moderna) or one shot of J&J, see <https://www.cdc.gov/coronavirus/2019-ncov/vaccines/stay-up-to-date.html>. Adjusted hazard ratios (aHR) were compared when adjusting for these additional covariates (aHR w/ vax) versus not adjusting for these covariates in our primary analysis (aHR w/o vax). The PASC conditions identified in both datasets were marked by ‡ symbols. The 95% confidence intervals (CI) of aHR were reported. The color panels represent different organ systems, including (from top to bottom): the nervous system, skin, respiratory system, circulatory system, endocrine and metabolic, digestive system, genitourinary system, and other general signs. The PASC outcomes were ascertained from day 30 after the SARS-CoV-2 infection and all the adjusted risk measures were computed 180 days after the SARS-CoV-2 infection.

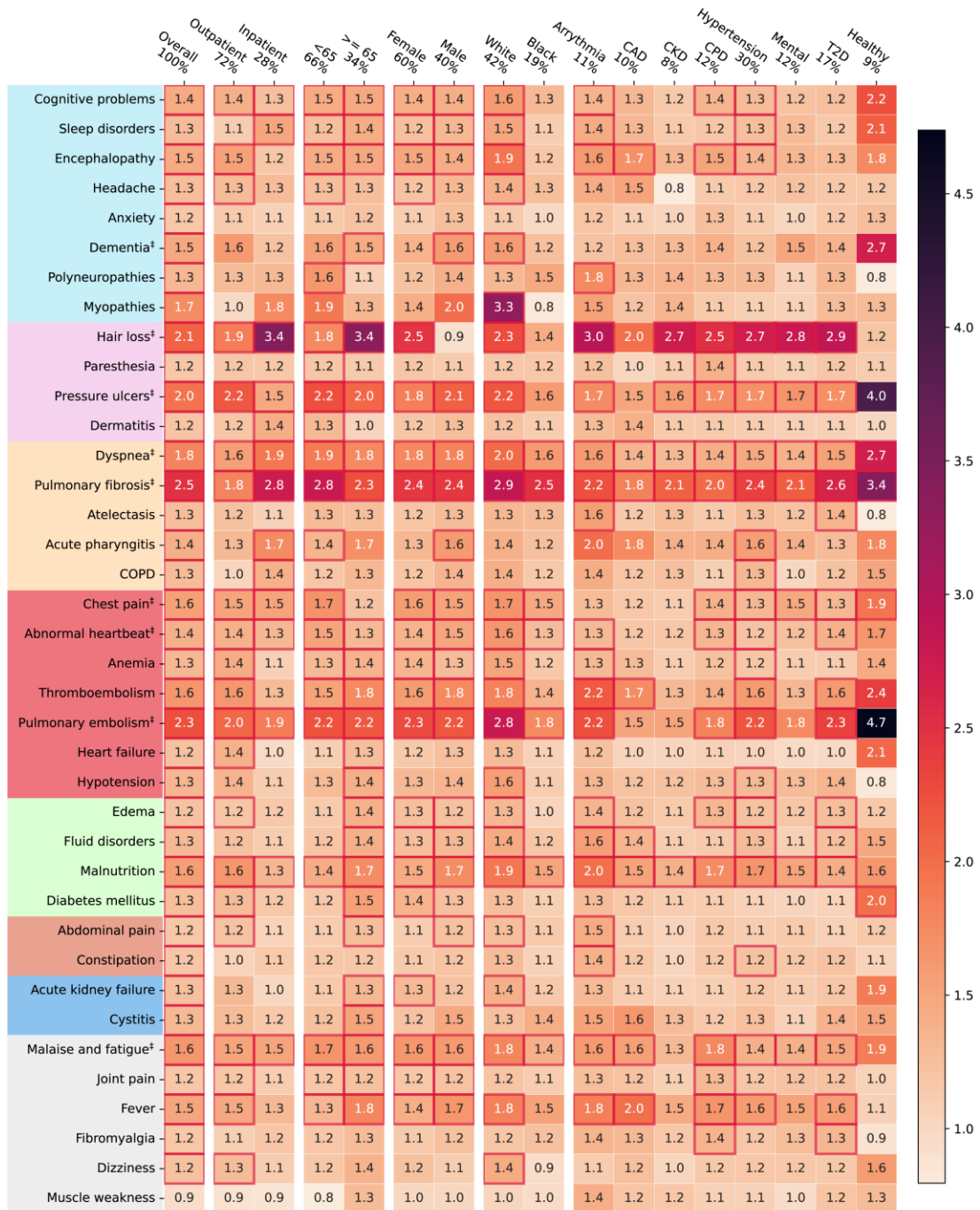


Supplementary Fig. 9. The post-acute sequelae of SARS-CoV-2 risks when adjusting for index day through Cubic B-Spline versus Categorical variables, the INSIGHT and OneFlorida+ cohorts, from March 2020 to November 2021. The index days since March 1st, 2020, were modeled by a) Cubic B-spline with five knots (namely, 7 spline basis functions of polynomial order 3) and b) categorizing days into five periods (March 2020 – June 2020, July 2020 – October 2020, November 2020 - February 2021, March 2021 – June 2021, July 2021 – November 2021). Adjusted hazard ratios (aHR) were compared when using the cubic B-spline method (aHR spline) versus the categorical method in the primary analysis (aHR category). The PASC conditions identified in both datasets were marked by ‡ symbols. The 95% confidence intervals (CI) of aHR were reported. The color panels represent different organ systems, including (from top to bottom): the nervous system, skin, respiratory system, circulatory system, endocrine and metabolic, digestive system, genitourinary system, and other general signs. The PASC outcomes were ascertained from day 30 after the SARS-CoV-2 infection and all the adjusted risk measures were computed 180 days after the SARS-CoV-2 infection.

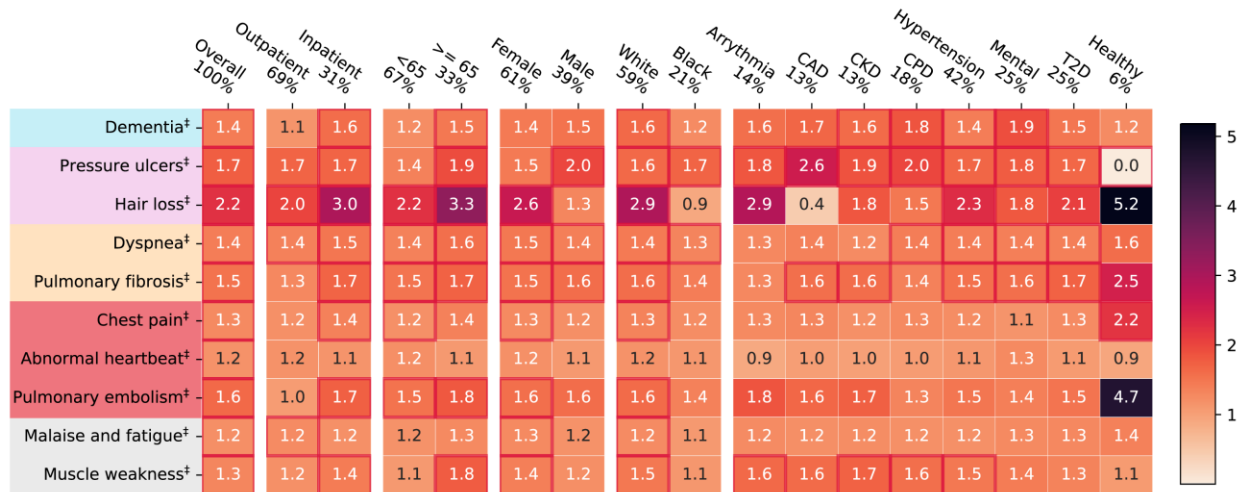


Supplementary Fig. 10. The post-acute sequelae of SARS-CoV-2 risks when using non-linear PS modeling versus linear PS modeling, the INSIGHT and OneFlorida+ cohorts, from March 2020 to November 2021. Adjusted hazard ratios (aHR) were adjusted for a nonlinear PS modeling using the gradient boosting machine with decision tree base learners (aHR nonlinear) versus the linear PS modeling using regularized logistic regression in our primary analysis (aHR lin). The PASC conditions identified in both datasets were marked by ‡ symbols. The color panels represent different organ systems, including (from top to bottom): the nervous system, skin, respiratory system, circulatory system, endocrine and metabolic, digestive system, genitourinary system, and other general signs. The PASC outcomes were ascertained from day 30 after the SARS-CoV-2 infection and all the adjusted risk measures were computed 180 days after the SARS-CoV-2 infection. The 95% confidence intervals (CI) of aHR were reported. PS, propensity score.

Supplementary Fig. 11. The post-acute sequelae of SARS-CoV-2 risks when downsampling the INSIGHT cohort into 22,341 SARS-CoV-2 positive patients and 177,010 negative patients versus using original data, the INSIGHT and OneFlorida+ cohorts, from March 2020 to November 2021. The number of SARS-CoV2 positive (or negative) patients in the INSIGHT after downsampling is the same as the number of positive (or negative) patients in the OneFlorida+. Adjusted hazard ratios (aHR) and 95% confidence intervals (CI) were compared using the downsampled cohorts (aHR sample) versus the original cohorts in our primary analysis (aHR all) for the INSIGHT cohort. The PASC conditions identified in both datasets were marked by ‡ symbols. The color panels represent different organ systems, including (from top to bottom): the nervous system, skin, respiratory system, circulatory system, endocrine and metabolic, digestive system, genitourinary system, and other general signs. The PASC outcomes were ascertained from day 30 after the SARS-CoV-2 infection and all the adjusted risk measures were computed 180 days after the SARS-CoV-2 infection.



Supplementary Fig. 12. Adjusted hazard ratios of post-acute sequelae of SARS-CoV-2 infection (PASC) over different subgroups, the INSIGHT cohort, from March 2020 to November 2021. Subgroups were stratified by their acute severity status, age groups, gender, race groups, and baseline pre-existing conditions. Different color panels represent different organ systems, including (from top to bottom): the nervous system, skin, respiratory system, circulatory system, blood-forming organs, endocrine and metabolic, digestive system, genitourinary system, and general signs. CAD, coronary artery disease; CKD, chronic kidney disease; CPD, chronic pulmonary disease; T2D, diabetes type 2; Healthy: no documented pre-existing conditions and no PASC-like symptoms at baseline. The PASC conditions also identified in OneFlorida+ were marked by ‡ symbols. The conditions with their aHRs' P-value < 8.39×10^{-5} (the Bonferroni-corrected significance threshold) were highlighted in red squares. The fraction of the subgroup population was shown at the top.



Supplementary Fig. 13. Adjusted hazard ratios of post-acute sequelae of SARS-CoV-2 infection (PASC) over different subgroups, the OneFlorida+ cohort, from March 2020 to November 2021. Subgroups were stratified by their acute severity status, age groups, gender, race groups, and baseline pre-existing conditions. Different color panels represent different organ systems, including (from top to bottom): the nervous system, skin, respiratory system, circulatory system, and general signs. CAD, coronary artery disease; CKD, chronic kidney disease; CPD, chronic pulmonary disease; T2D, diabetes type 2; Healthy: no documented pre-existing conditions and no PASC-like symptoms at baseline. The conditions with their aHRs' P-value < 8.39×10^{-5} (the Bonferroni-corrected significance threshold) were highlighted in red squares. The fraction of the subgroup population was shown at the top.

Supplementary Table 1. Results of negative outcome control in both the INSIGHT and OneFlorida+ cohorts, March 2020–November 2021.

Results of negative outcome control in both the INSIGHT and OneFlorida+ cohorts, March 2020–November 2021.

Negative Outcomes	Adjusted Hazard Ratio (95% Confidence Interval)^a	aHR's P-value^b	No. of SARS-CoV-2 Positive patients	No. of SARS-CoV-2 Negative patients^c	Number of events in the case group	Number of events in the control group	No. of unbalanced covariates	No. of unbalanced covariates after re-weighting^a
INSIGHT								
Accidental Injuries	1.02 (0.84, 1.23)	0.87	33,898	33,898	347	248	31	0
Benign Neoplasms	0.93 (0.84, 1.03)	0.14	30,461	30,461	1030	1228	31	0
OneFlorida+								
Accidental Injuries	0.99 (0.88, 1.12)	0.92	18,967	18,967	661	599	22	0
Benign neoplasms	0.93 (0.82, 1.05)	0.23	18,729	18,729	585	693	21	0

- a. Outcomes were ascertained from day 30 after the SARS-CoV-2 infection. The adjusted hazard ratios were computed 180 days after the SARS-CoV-2 infection by adjusting high-dimensional baseline variables the same as in the screening sequelae using inverse probability of treatment weighting as discussed in the Method section. All the baseline covariates were balanced in terms of standardized mean difference < 0.1.
- b. An aHR close to 1 with a large P-value (> 0.05) indicates no significant association between SARS-CoV-2 infection and outcomes was found. The aHR and its P-value were calculated by the Cox proportional hazard model and the Wald Chi-Square test.
- c. The number of patients in the control group was randomly selected from all SARS-CoV-2 negative patients with the same number as patients in the case group.

Supplementary Table 2. Specifications of the high-throughput screening framework for identifying potential Post-Acute Sequelae of SARS-CoV-2 infected (PASC) using the INSIGHT Electronic Health Records in New York City and OneFlorida+ in Florida (March 2020 – November 2021).

Protocol component	Specifications of hypothetical target trials	High-throughput emulation
Eligibility criteria	<ul style="list-style-type: none"> • Age ≥ 20 years at index, and no upper age limit, between March 1, 2020, and November 30, 2021. • Patients without any positive SARS-CoV-2 polymerase-chain-reaction (PCR) or Antigen test, or COVID-19 diagnoses before the index • Use of the INSIGHT health care system, defined as at least one encounter with any diagnoses, in the past 3 years to 7 days before the index. • Use of the INSIGHT health care system, defined as at least one encounter with any diagnoses, 31 days to 180 days after index in the follow-up period, indicates being alive after the potential COVID-19 acute phase and at least 31 days of potential follow-up. • Known documented SARS-CoV-2 PCR/Antigen lab positive or negative results • (High-throughput scenario for PASC) No history of target post-acute sequelae of COVID-19 in the past 3 years to 7 days before the index. • (High-throughput scenario for PASC medications) No history of target medication usage in the past 1 year to 7 days before the index. 	<p>Same as for the hypothetical target trials. We identified the SARS-CoV-2 infected PCR/Antigen tests using the laboratory results table, and the COVID-19 diagnoses using ICD-10 codes in the patients' diagnosis table in the INSIGHT EHR system or OneFlorida+ EHR system, following the PCORnet data model</p>

Exposure strategies	<ul style="list-style-type: none"> • Exposure group: Infection of SARS-CoV-2 and the SARS-CoV-2 PCR/Antigen tested positive • Control group: No infection of SARS-CoV-2, and the SARS-CoV-2 PCR/Antigen tests kept negative in the follow-up period 	<ul style="list-style-type: none"> • Same as for the hypothetical target trial. • SARS-CoV-2 infected cases were defined as patients who had any SARS-CoV-2 positive PCR/Antigen test, and the baseline date, or the index date, is defined as the date of the first documented positive SARS-CoV-2 PCR/Antigen test. • Non-infected controls were defined as patients whose SARS-CoV-2 PCR/Antigen tests were all negative, and there were no COVID-19-related diagnoses documented at any time. The baseline date, or the index date, is defined as the date of the first documented SARS-CoV-2 PCR/Antigen test.
Group assignment	Individuals are randomly assigned to an exposure strategy at baseline and are aware of the assigned exposure strategy.	We classified patients into different exposure groups according to their baseline eligibility criteria and exposure strategies. We assumed that the exposure group and the control group were exchangeable by adjusting for baseline covariates, including age, gender, race, ethnicity, social-economic status, hospital utilization history, period of infection, baseline comorbidities, history of prescriptions, etc.
High-throughput Outcomes	<ul style="list-style-type: none"> • 137 potential post-acute sequelae of COVID-19 (PASC) • 459 categories of drugs due to post-acute sequelae of COVID-19 <p>Hospitalization due to PASC ICU admission due to PASC Death due to PASC (defined as death after 31 days of SARS-CoV-2 infection)</p>	Same as for the hypothetical target trials.
Follow-up	We followed each patient from his/her baseline day until the day of the outcome of interest, death, 180 days after baseline, or the end of the study period (November 30, 2021), whichever happens first.	Same as for the hypothetical target trial.

Outcome contrasts	Excess risk of newly-onset post-acute sequelae of COVID-19 against baseline incidence. The PASC outcomes were ascertained from day 30 after the SARS-CoV-2 infection and all the adjusted risk measures were computed 180 days after the SARS-CoV-2 infection.	Same as for the hypothetical target trial.
High-throughput trials for screening PASC	For each target PASC outcome, we conducted a corresponding hypothetical target trial among which patients were free of the target PASC outcome at baseline and had no history of target outcome before baseline. The number of potential PASCs is large and thus the number of trials is large.	We emulated 137 trials for potential PASC diagnosis outcomes and 459 trials for potential PASC medication outcomes. For each emulated trial, the exposure group consisted of eligible SARS-CoV-2 infected patients without a history of target outcome at baseline. The control group was built from SARS-CoV-2-negative patients who were also without any history of target outcome at baseline. The ratio of the number of patients in the exposure group and the control group was 1:5.
Statistical analysis	Cumulative incidence (risk) curves, excess burden, and hazard ratio (HR) between the two exposure groups were estimated at 6 months. Subgroup analyses by baseline age, race, gender, and severity of acute phase of COVID-19, and different waves. Sensitivity analyses.	Same as for the hypothetical target trial. <ul style="list-style-type: none"> • We used inverse probability of treatment weighting (IPTW) to adjust for high-dimensional baseline covariates. We used L2-norm logistic regression for propensity score calculation, and the best PS model was selected in 10-fold cross-validations for both goodness-of-balance and goodness-of-fit. • Adjusted Cox proportional hazard model, Kaplan-Meier estimator, and Aalen-Johansen estimator were used. • The p-value was corrected by the Bonferroni method for multiple tests of diagnoses and medications separately. • Potential PASCs were selected according to the number of patients with target outcome in the real-world data (>100), the corrected P-value (<0.05/596), and the hazard ratio (>1)

Abbreviations: PCORnet, the National Patient-Centered Clinical Research Network; PASC, Post-Acute Sequelae of COVID-19; PCR, polymerase chain reaction; aHR, adjusted hazard ratio; IPTW, inverse probability of treatment weighting

Supplementary Table 3. Cross-validation algorithm tailored for our machine learning-based propensity score calculation for each emulated trial.

Table S2. Cross-validation algorithm for the ML-based PS model training, selection, and evaluation

Input:

(\mathbf{X}, \mathbf{T}) : n patients' covariates and exposure assignment where $\mathbf{T} \in \{0, 1\}^n$;

\mathbf{F}_θ : a set of machine learning-based propensity score (PS) models;

Output:

f_{best} : the best PS model learned from (\mathbf{X}, \mathbf{T})

Goodness-of-balance performance, and goodness-of-fit performance

1. **for** each f_θ in \mathbf{F}_θ **do**:
 2. randomly splitting (\mathbf{X}, \mathbf{T}) into K ($K = 10$ in our experiments) equal-sized subsets
 3. **for** each $(\mathbf{X}_i, \mathbf{T}_i)$ subset in the K subsets **do**:
 4. training f_θ on the remaining $K - 1$ subsets $(\mathbf{X}_{K-i}, \mathbf{T}_{K-i})$ by optimizing binary cross entropy loss $L(\mathbf{T}, f_\theta(\mathbf{X}))$
 5. computing trimmed and stabilized IPTW \mathbf{w} by using f_θ on (\mathbf{X}, \mathbf{T})
 6. computing re-weighted SMD_i on (\mathbf{X}, \mathbf{T}) by using \mathbf{w}
 7. computing the number of unbalanced features $n_{\text{unbalance-}i}$ after IPTW
 8. computing the AUC_i of f_θ on the testing set $(\mathbf{X}_i, \mathbf{T}_i)$
 9. computing f_θ 's average goodness-of-balance performance $n_{\text{unbalance-}\theta} = E_{i \sim K}[n_{\text{unbalance-}i}]$ and goodness-of-fit performance $\text{AUC}_\theta = E_{i \sim K}[\text{AUC}_i]$ over K folds
 10. updating best-selected model $f_{\text{best}} := f_\theta$, the best performance $n_{\text{unbalance-best}}$ and AUC_{best} if $n_{\text{unbalance-}\theta}$ is smaller than the current minimum $n_{\text{unbalance-best}}$, or $n_{\text{unbalance-}\theta}$ is equal to the current minimum $n_{\text{unbalance-best}}$ but the AUC_θ is larger than the current maximum AUC_{best}
 11. re-training f_{best} on the whole dataset (\mathbf{X}, \mathbf{T})
 12. re-computing stabilized IPTW \mathbf{w}_{best} by using learned f_{best} on (\mathbf{X}, \mathbf{T})
 13. re-computing re-weighted SMD_{best} on (\mathbf{X}, \mathbf{T}) by using \mathbf{w}_{best}
 14. re-computing the number of unbalanced features $n_{\text{unbalance-best}}$ after IPTW
 15. **return** f_{best} , $n_{\text{unbalance-best}}$ and AUC_{best}
-

Supplementary Table 4. Baseline SARS-CoV-2 vaccination status recorded in the two EHR databases^a

INSIGHT (March 2020 - November 2021)						
	Overall	%	SARS-CoV-2 Positive	%	SARS-CoV-2 Negative	%
N	361401	100.0%	35275	100.0%	326126	100.0%
Fully vaccinated ^b	8293	2.3%	358	1.0%	7935	2.4%
Partially vaccinated ^c	6703	1.9%	474	1.3%	6229	1.9%
No evidence	346406	95.9%	34444	97.6%	311962	95.7%
INSIGHT (December 2020 - November 2021)^d						
N	154676	100.0%	20301	100.0%	134375	100.0%
Fully vaccinated ^b	8293	5.4%	358	1.8%	7935	5.9%
Partially vaccinated ^c	6699	4.3%	474	2.3%	6225	4.6%
No evidence	139685	90.3%	19470	95.9%	120215	89.5%
OneFlorida+ (March 2020 - November 2021)						
	Overall	%	SARS-CoV-2 Positive	%	SARS-CoV-2 Negative	%
N	199351	100.0%	22341	100.0%	177010	100.0%
Fully vaccinated ^b	1161	0.6%	98	0.4%	1063	0.6%
Partially vaccinated ^c	499	0.3%	67	0.3%	432	0.2%
No evidence	197691	99.2%	22176	99.3%	175515	99.2%
OneFlorida+ (December 2020 - November 2021)^d						
N	88021	100.0%	13074	100.0%	74947	100.0%
Fully vaccinated ^b	1161	1.3%	98	0.7%	1063	1.4%
Partially vaccinated ^c	499	0.6%	67	0.5%	432	0.6%
No evidence	86361	98.1%	12909	98.7%	73452	98.0%

a. Vaccination capture at baseline. b. The fully vaccinated status is defined as two shots of mRNA vaccine (Pfizer, or Moderna) or one shot of J&J before baseline according to <https://www.cdc.gov/coronavirus/2019-ncov/vaccines/stay-up-to-date.html>. c. The partially vaccinated status is defined as having any recorded vaccination events but not meeting the fully vaccinated criteria. d. The earliest available vaccine began in early December 2020, and nearly half of the study patients got infected before the vaccine was available. Besides, only 9.7% population in INSIGHT had any baseline vaccination capture after December 1st, 2020, and 1.8% population for OneFlorida+

Supplementary Data 1. COVID-19 Phenotyping Lab LOINC codes and Diagnosis ICD10 codes (spreadsheet).

Supplementary Data 2. PASC Adult Diagnostic List for Screening (spreadsheet).

Supplementary Data 3. Baseline population characteristics with more comorbidities information, INSIGHT and OneFlorida+ cohorts, March 2020 to November 2021 (spreadsheet).

Supplementary Data 4. Characteristics of PASC-Specific Cohorts on INSIGHT, NYC, March 2020 to November 2021 (spreadsheet).

Supplementary Data 5. Characteristics of PASC-Specific Cohorts on OneFlorida, Florida, March 2020 to November 2021 (spreadsheet).