

Hypercap CC NLP Analysis

Table of contents

1	Workbook for MIMIC Hypercapnia Presenting Chief Concern Analysis	1
1.1	Environment Gate	1
1.2	Load Data	2
1.3	Descriptive Checks	8
1.4	ICD And Inclusion Categories	9
1.5	Symptom Composition By Hypercapnia Definition	13
1.6	Symptom Distribution By Ascertainment Overlap	18
1.7	ICD Diagnostic Performance (ICD as predictor)	22
1.8	Ascertainment overlap UpSet	24
1.9	PDF-ready long tables	26
1.10	Association Model	31
1.11	Export Verification	32

1 Workbook for MIMIC Hypercapnia Presenting Chief Concern Analysis

This notebook is a deterministic analysis workflow for the NLP-augmented hypercapnia cohort workbook.

1.1 Environment Gate

Fail fast if required packages are missing. Use `uv sync` to repair the environment.

```
import importlib.util
```

```
required_packages = [  
    "numpy",
```

```

    "pandas",
    "matplotlib",
    "seaborn",
    "statsmodels",
    "upsetplot",
    "openpyxl",
]
missing = [pkg for pkg in required_packages if importlib.util.find_spec(pkg)
↪ is None]
if missing:
    raise ModuleNotFoundError(
        "Missing required packages: "
        + ", ".join(missing)
        + ". Run `uv sync` from the repository root and rerun the notebook."
    )
print("Environment check passed.")

```

Environment check passed.

1.2 Load Data

Use a single canonical workbook path under MIMIC tabular data.

```

import json
import os
import sys
from pathlib import Path

import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
import statsmodels.api as sm
from statsmodels.stats.proportion import proportion_confint
from upsetplot import UpSet, from_indicators

WORK_DIR = Path(os.getenv("WORK_DIR", Path.cwd())).expanduser().resolve()
SRC_DIR = WORK_DIR / "src"
if SRC_DIR.exists() and str(SRC_DIR) not in sys.path:
    sys.path.insert(0, str(SRC_DIR))
CANONICAL_NLP_FILENAME = "MIMICIV all with CC_with_NLP.xlsx"

```

```

def resolve_analysis_input_path(work_dir: Path, input_filename: str | None =
    ↪ None) -> Path:
    filename = input_filename or CANONICAL_NLP_FILENAME
    input_path = (work_dir / "MIMIC tabular data" /
    ↪ filename).expanduser().resolve()
    if not input_path.exists():
        raise FileNotFoundError(
            "Expected analysis input workbook was not found at "
            f"{input_path}. Run the classifier notebook first or set "
            "ANALYSIS_INPUT_FILENAME."
        )
    return input_path

def ensure_required_columns(df: pd.DataFrame, required: list[str]) -> None:
    missing = sorted(set(required).difference(df.columns))
    if missing:
        raise KeyError(f"Missing required columns: {missing}")

def to_binary_flag(series: pd.Series) -> pd.Series:
    numeric = pd.to_numeric(series, errors="coerce").fillna(0)
    return (numeric > 0).astype(int)

def _binary_or_zero(df: pd.DataFrame, column: str) -> pd.Series:
    if column in df.columns:
        return to_binary_flag(df[column])
    return pd.Series(0, index=df.index, dtype="int64")

def classify_icd_category_vectorized(df: pd.DataFrame) -> pd.Series:
    j9602 = _binary_or_zero(df, "ICD10_J9602")
    j9612 = _binary_or_zero(df, "ICD10_J9612")
    j9622 = _binary_or_zero(df, "ICD10_J9622")
    j9692 = _binary_or_zero(df, "ICD10_J9692")
    e662 = _binary_or_zero(df, "ICD10_E662")
    icd9_27803 = _binary_or_zero(df, "ICD9_27803")

    category = np.select(
        [
            j9602.eq(1),
            j9612.eq(1),

```

```

        j9622.eq(1),
        j9692.eq(1),
        e662.eq(1) | icd9_27803.eq(1),
    ],
    [
        "Acute RF with hypoxia",
        "Acute RF with hypercapnia",
        "Acute RF with hypoxia & hypercapnia",
        "Respiratory failure, unspecified",
        "Obesity hypoventilation syndrome",
    ],
    default="Other / None",
)
return pd.Series(category, index=df.index, name="icd_category")

def classify_inclusion_type_vectorized(any_icd: pd.Series, gas_any:
    ↪ pd.Series) -> pd.Series:
    any_icd_bin = to_binary_flag(any_icd)
    gas_any_bin = to_binary_flag(gas_any)
    labels = np.select(
        [
            any_icd_bin.eq(1) & gas_any_bin.eq(1),
            any_icd_bin.eq(1) & gas_any_bin.eq(0),
            any_icd_bin.eq(0) & gas_any_bin.eq(1),
        ],
        ["Both", "ICD_only", "Gas_only"],
        default="Neither",
    )
    return pd.Series(labels, index=any_icd.index, name="inclusion_type")

def binary_crosstab_yes_no(df: pd.DataFrame, row_col: str, flag_col: str) ->
    ↪ pd.DataFrame:
    ensure_required_columns(df, [row_col, flag_col])
    tab = pd.crosstab(df[row_col], to_binary_flag(df[flag_col]),
    ↪ margins=False, dropna=False)
    tab = tab.reindex(columns=[0, 1], fill_value=0)
    tab.columns = ["No", "Yes"]
    row_totals = tab.sum(axis=1).replace(0, np.nan)
    tab["Percent_yes"] = (tab["Yes"] / row_totals * 100).round(1).fillna(0)
    return tab

```

```

def symptom_distribution_by_overlap(
    df: pd.DataFrame,
    group_col: str,
    symptom_col: str,
    top_k: int = 10,
) -> tuple[pd.DataFrame, pd.DataFrame]:
    ensure_required_columns(df, [group_col, symptom_col])
    tmp = df.dropna(subset=[group_col, symptom_col]).copy()
    if tmp.empty:
        return pd.DataFrame(columns=[group_col, "symptom_group", "N",
            ↪ "Percent"]), pd.DataFrame()
    top_symptoms =
    ↪ tmp[symptom_col].value_counts(dropna=False).head(top_k).index
    tmp["symptom_group"] =
    ↪ tmp[symptom_col].where(tmp[symptom_col].isin(top_symptoms), "Other")
    counts = (
        tmp.groupby([group_col, "symptom_group"], dropna=False)
        .size()
        .reset_index(name="N")
    )
    counts["Percent"] = (
        counts.groupby(group_col)["N"].transform(lambda x: x / x.sum() *
    ↪ 100).round(1)
    )
    pivot = counts.pivot_table(
        index="symptom_group",
        columns=group_col,
        values="Percent",
        fill_value=0,
    ).round(1)
    return counts, pivot

def classify_gas_source_overlap(
    abg_series: pd.Series,
    vbg_series: pd.Series,
    other_series: pd.Series,
) -> pd.Series:
    abg = to_binary_flag(abg_series)
    vbg = to_binary_flag(vbg_series)
    other = to_binary_flag(other_series)
    labels = np.select(

```

```

[
    abg.eq(1) & vbg.eq(1) & other.eq(1),
    abg.eq(1) & vbg.eq(1) & other.eq(0),
    abg.eq(1) & vbg.eq(0) & other.eq(1),
    abg.eq(0) & vbg.eq(1) & other.eq(1),
    abg.eq(1) & vbg.eq(0) & other.eq(0),
    abg.eq(0) & vbg.eq(1) & other.eq(0),
    abg.eq(0) & vbg.eq(0) & other.eq(1),
],
[
    "ABG+VBG+OTHER",
    "ABG+VBG",
    "ABG+OTHER",
    "VBG+OTHER",
    "ABG-only",
    "VBG-only",
    "OTHER-only",
],
default="No-gas",
)
return pd.Series(labels, index=abg_series.index,
    ↪ name="gas_source_overlap")

def render_latex_longtable(
    table_df: pd.DataFrame,
    *,
    caption: str,
    label: str,
    landscape: bool = False,
    index: bool = True,
) -> str:
    latex_text = table_df.to_latex(
        index=index,
        escape=False,
        longtable=True,
        caption=caption,
        label=label,
    )
    if landscape:
        latex_text = "\\begin{landscape}\\n" + latex_text +
    ↪ "\\n\\end{landscape}\\n"
    return latex_text

```

```

from hypercap_cc_nlp.pipeline_audit import collect_run_manifest

ANALYSIS_INPUT_FILENAME = os.getenv("ANALYSIS_INPUT_FILENAME")
ANALYSIS_INPUT_PATH = resolve_analysis_input_path(
    WORK_DIR,
    ANALYSIS_INPUT_FILENAME if ANALYSIS_INPUT_FILENAME else None,
)
OUTPUT_DIR = WORK_DIR

HYPERCAP_CRITERIA = [
    "any_hypercap_icd",
    "abg_hypercap_threshold",
    "vbg_hypercap_threshold",
    "other_hypercap_threshold",
    "pco2_threshold_any",
]

SYMPTOM_COL = "RFV1_name"

df = pd.read_excel(ANALYSIS_INPUT_PATH, engine="openpyxl")
required_analysis_cols = sorted({SYMPTOM_COL, *HYPERCAP_CRITERIA})
try:
    ensure_required_columns(df, required_analysis_cols)
except KeyError as exc:
    raise KeyError(
        "Analysis input schema mismatch. Run 'Hypercap CC NLP Classifier.qmd'"
        ↪ " "
        f"to regenerate '{CANONICAL_NLP_FILENAME}' before running analysis."
    ) from exc

for column in HYPERCAP_CRITERIA:
    df[column] = to_binary_flag(df[column])

print(
    f"Loaded {ANALYSIS_INPUT_PATH.name}: {df.shape[0]:,} rows x"
    ↪ {df.shape[1]:,} columns"
)
print(f"Analysis input path: {ANALYSIS_INPUT_PATH}")

Loaded MIMICIV all with CC_with_NLP.xlsx: 41,322 rows x 266 columns
Analysis input path: /Users/blocke/Box Sync/Residency Personal Files/Scholarly Work/Locke Res

```

1.3 Descriptive Checks

Compute core cohort summaries with guarded column checks.

```
gender_candidates = [col for col in df.columns if
    ↪ col.lower().startswith("gender")]
if not gender_candidates:
    raise KeyError("No gender-like column found. Expected a column starting
    ↪ with 'gender'.")
gender_col = gender_candidates[0]

gender_summary = (
    df[gender_col]
    .value_counts(dropna=False)
    .rename_axis(gender_col)
    .to_frame("Count")
)
gender_summary["Percent"] = (gender_summary["Count"] / len(df) *
    ↪ 100).round(1)

age_summary = pd.Series(
    {
        "Mean": round(float(df["age"].mean()), 2),
        "SD": round(float(df["age"].std()), 2),
        "Q1": round(float(df["age"].quantile(0.25)), 2),
        "Q3": round(float(df["age"].quantile(0.75)), 2),
    },
    name="Age (years)",
)

prevalence_label_map = {
    "any_hypercap_icd": "Hypercapnic RF ICD (any)",
    "abg_hypercap_threshold": "ABG hypercapnia threshold",
    "vbg_hypercap_threshold": "VBG hypercapnia threshold",
    "other_hypercap_threshold": "PCO2 OTHER threshold",
    "pco2_threshold_any": "PCO2 threshold any source",
}
cohort_n = int(len(df))
hypercap_prevalence = (
    pd.DataFrame(
        {
```



```

        "Definition": [prevalence_label_map[col] for col in
            ↪ HYPERCAP_CRITERIA],
        "Column": HYPERCAP_CRITERIA,
        "Count": [int(df[col].sum()) for col in HYPERCAP_CRITERIA],
        "Denominator_N": [cohort_n for _ in HYPERCAP_CRITERIA],
        "Percent": [round(float(df[col].mean() * 100), 1) for col in
            ↪ HYPERCAP_CRITERIA],
    }
)
.set_index("Definition")
.sort_values("Count", ascending=False)
)

display(gender_summary)
display(age_summary.to_frame())
display(hypercap_prevalence)

```

	Count	Percent
gender		
M	21865	52.9
F	19457	47.1

	Age (years)
Mean	64.78
SD	17.78
Q1	54.00
Q3	78.00

	Column	Count	Denominator_N	Percent
Definition				
PCO2 threshold any source	pco2_threshold_any	41080	41322	99.4
PCO2 OTHER threshold	other_hypercap_threshold	36518	41322	88.4
VBG hypercapnia threshold	vbg_hypercap_threshold	17548	41322	42.5
ABG hypercapnia threshold	abg_hypercap_threshold	10749	41322	26.0
Hypercapnic RF ICD (any)	any_hypercap_icd	1983	41322	4.8

1.4 ICD And Inclusion Categories

Use vectorized helper functions to avoid row-wise `apply(axis=1)`.

```

df["icd_category"] = classify_icd_category_vectorized(df)
df["inclusion_type"] = classify_inclusion_type_vectorized(
    df["any_hypercap_icd"],
    df["pco2_threshold_any"],
)

icd_category_summary = (
    df["icd_category"]
    .value_counts(dropna=False)
    .rename_axis("ICD Category")
    .to_frame("Count")
)
icd_category_summary["Percent"] = (icd_category_summary["Count"] / len(df) *
    ↪ 100).round(1)
icd_category_summary["Denominator_N"] = int(len(df))

inclusion_summary = (
    df["inclusion_type"]
    .value_counts(dropna=False)
    .rename_axis("Inclusion Type")
    .to_frame("Count")
)
inclusion_summary["Percent"] = (inclusion_summary["Count"] / len(df) *
    ↪ 100).round(1)
inclusion_summary["Denominator_N"] = int(len(df))

icd_positive_df = df.loc[df["any_hypercap_icd"].eq(1)].copy()
icd_positive_n = int(len(icd_positive_df))
icd_positive_breakdown = pd.DataFrame(
    {
        "Definition": [
            "ABG threshold positive",
            "VBG threshold positive",
            "PCO2 OTHER threshold positive",
            "Any gas threshold positive",
        ],
        "Count": [
            int(icd_positive_df["abg_hypercap_threshold"].sum()),
            int(icd_positive_df["vbg_hypercap_threshold"].sum()),
            int(icd_positive_df["other_hypercap_threshold"].sum()),
            int(icd_positive_df["pco2_threshold_any"].sum()),
        ],
    }
)

```

```

)
if icd_positive_n > 0:
    icd_positive_breakdown["Percent"] = (
        icd_positive_breakdown["Count"] / icd_positive_n * 100
    ).round(1)
else:
    icd_positive_breakdown["Percent"] = 0.0
icd_positive_breakdown["Denominator_N"] = icd_positive_n

icd_positive_category_summary = (
    icd_positive_df["icd_category"]
    .value_counts(dropna=False)
    .rename_axis("ICD Category (ICD-positive subset)")
    .to_frame("Count")
)
if icd_positive_n > 0:
    icd_positive_category_summary["Percent"] = (
        icd_positive_category_summary["Count"] / icd_positive_n * 100
    ).round(1)
else:
    icd_positive_category_summary["Percent"] = 0.0
icd_positive_category_summary["Denominator_N"] = icd_positive_n

display(icd_category_summary)
display(inclusion_summary)
display(icd_positive_category_summary)
display(icd_positive_breakdown)

```

	Count	Percent	Denominator_N
ICD Category			
Other / None	39339	95.2	41322
Acute RF with hypoxia	793	1.9	41322
Obesity hypoventilation syndrome	524	1.3	41322
Acute RF with hypoxia & hypercapnia	386	0.9	41322
Respiratory failure, unspecified	187	0.5	41322
Acute RF with hypercapnia	93	0.2	41322

	Count	Percent	Denominator_N
Inclusion Type			
Gas_only	39339	95.2	41322
Both	1741	4.2	41322

Inclusion Type	Count	Percent	Denominator_N
ICD_only	242	0.6	41322

ICD Category (ICD-positive subset)	Count	Percent	Denominator_N
Acute RF with hypoxia	793	40.0	1983
Obesity hypoventilation syndrome	524	26.4	1983
Acute RF with hypoxia & hypercapnia	386	19.5	1983
Respiratory failure, unspecified	187	9.4	1983
Acute RF with hypercapnia	93	4.7	1983

Definition	Count	Percent	Denominator_N
0 ABG threshold positive	1012	51.0	1983
1 VBG threshold positive	1482	74.7	1983
2 PCO2 OTHER threshold positive	1596	80.5	1983
3 Any gas threshold positive	1741	87.8	1983

```

symptom_work_df = df.copy()
symptom_text =
    ↪ symptom_work_df[SYMPTOM_COL].fillna("").astype(str).str.strip()
symptom_work_df["symptom_missing_flag"] = symptom_text.eq("")
top_symptom_labels = symptom_text.loc[~symptom_work_df[
    ↪ "symptom_missing_flag"]].value_counts().head(10).index
symptom_work_df["symptom_group"] = symptom_text.where(
    symptom_text.isin(top_symptom_labels),
    "Other",
)
symptom_work_df.loc[symptom_work_df["symptom_missing_flag"], "symptom_group"]
    ↪ = "No symptom recorded"

crosstab_tables = {}
for definition in HYPERCAP_CRITERIA:
    definition_table = binary_crosstab_yes_no(symptom_work_df,
    ↪ "symptom_group", definition)
    crosstab_tables[definition] = definition_table.sort_values("Percent_yes",
    ↪ ascending=False)

display(crosstab_tables["pco2_threshold_any"].head(10))

```

```
symptom_non_null =
  ↳ symptom_work_df.loc[~symptom_work_df["symptom_missing_flag"]].copy()
```

	No	Yes	Percent_yes
symptom_group			
Injuries & adverse effects	7	2920	99.8
Diseases (patient-stated)	9	3315	99.7
Symptom – Digestive	17	5765	99.7
Symptom – Nervous	17	6562	99.7
Symptom – Genitourinary	10	2586	99.6
Symptom – Eye/Ear	5	1022	99.5
Symptom – Circulatory	27	4159	99.4
Other	23	3208	99.3
Symptom – General	18	2570	99.3
Symptom – Skin/Hair/Nails	10	978	99.0

1.5 Symptom Composition By Hypercapnia Definition

Generate counts, percentages, and clipped Wald 95% confidence intervals; export stable tables for downstream reporting.

```
definition_long_df = symptom_non_null.melt(
    id_vars=["symptom_group"],
    value_vars=HYPERCAP_CRITERIA,
    var_name="Hypercapnia_Definition",
    value_name="Positive",
)
definition_positive_df =
  ↳ definition_long_df.loc[definition_long_df["Positive"].eq(1)].copy()

definition_counts_df = (
    definition_positive_df.groupby(["Hypercapnia_Definition",
  ↳ "symptom_group"], dropna=False)
    .size()
    .reset_index(name="Count")
)
definition_counts_df["Total"] = definition_counts_df.groupby(
  ↳ "Hypercapnia_Definition")["Count"].transform("sum")
definition_counts_df["Percent"] = definition_counts_df["Count"] /
  ↳ definition_counts_df["Total"] * 100
```

```

p_hat = (definition_counts_df["Percent"] / 100).clip(0, 1)
n_obs = definition_counts_df["Total"].replace(0, np.nan)
se = np.sqrt((p_hat * (1 - p_hat)) / n_obs).fillna(0)
definition_counts_df["CI_lower"] = ((p_hat - 1.96 * se).clip(0, 1) *
    ↪ 100).round(2)
definition_counts_df["CI_upper"] = ((p_hat + 1.96 * se).clip(0, 1) *
    ↪ 100).round(2)
definition_counts_df["Percent"] = definition_counts_df["Percent"].round(2)

definition_counts_df = definition_counts_df.sort_values(
    ["Hypercapnia_Definition", "Count"],
    ascending=[True, False],
)

definition_pivot_df = definition_counts_df.pivot_table(
    index="symptom_group",
    columns="Hypercapnia_Definition",
    values="Percent",
    fill_value=0,
).round(2)

definition_output_path = OUTPUT_DIR /
    ↪ "Symptom_Composition_by_Hypercapnia_Definition.xlsx"
pivot_output_path = OUTPUT_DIR / "Symptom_Composition_Pivot_ChartReady.xlsx"
definition_counts_df.to_excel(definition_output_path, index=False)
definition_pivot_df.to_excel(pivot_output_path)

display(definition_counts_df.head(12))
print(f"Exported: {definition_output_path}")
print(f"Exported: {pivot_output_path}")

```

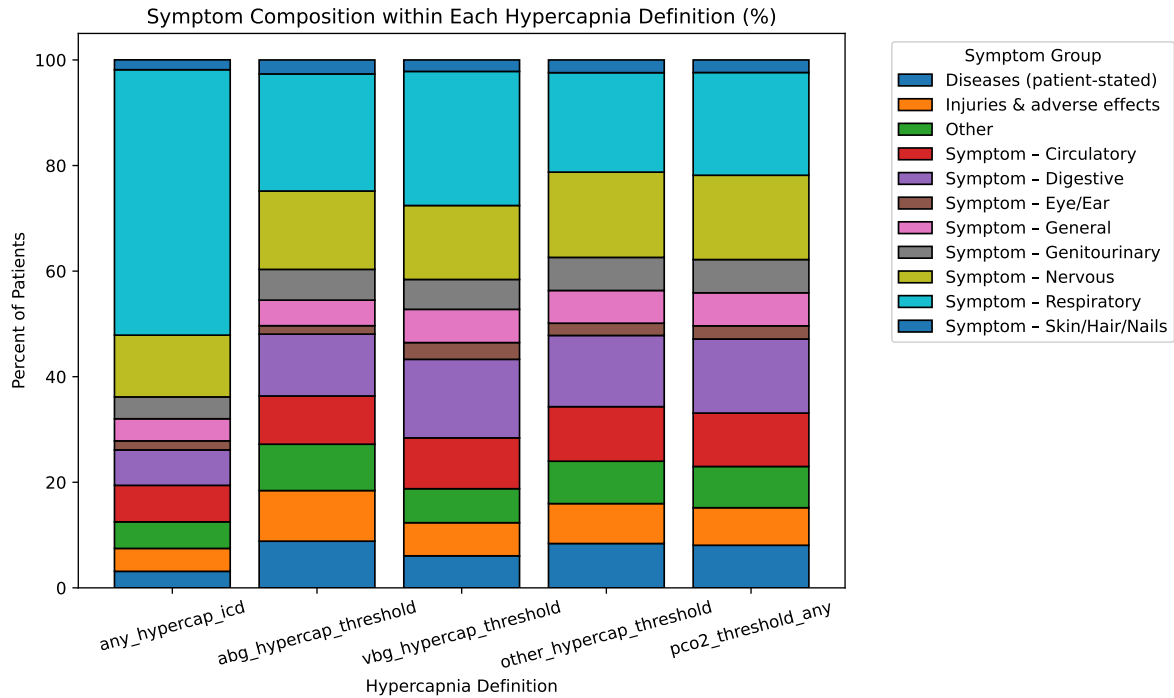
	Hypercapnia_Definition	symptom_group	Count	Total	Percent	CI_lower	CI_upper
9	abg_hypercap_threshold	Symptom – Respiratory	2384	10749	22.18	21.39	22.96
8	abg_hypercap_threshold	Symptom – Nervous	1595	10749	14.84	14.17	15.51
4	abg_hypercap_threshold	Symptom – Digestive	1259	10749	11.71	11.10	12.32
1	abg_hypercap_threshold	Injuries & adverse effects	1031	10749	9.59	9.03	10.15
3	abg_hypercap_threshold	Symptom – Circulatory	982	10749	9.14	8.59	9.68
0	abg_hypercap_threshold	Diseases (patient-stated)	950	10749	8.84	8.30	9.37
2	abg_hypercap_threshold	Other	944	10749	8.78	8.25	9.32
7	abg_hypercap_threshold	Symptom – Genitourinary	627	10749	5.83	5.39	6.28
6	abg_hypercap_threshold	Symptom – General	520	10749	4.84	4.43	5.24

	Hypercapnia_Definition	symptom_group	Count	Total	Percent	CI_lower	CI_upper
10	abg_hypercap_threshold	Symptom – Skin/Hair/Nails	285	10749	2.65	2.35	2.96
5	abg_hypercap_threshold	Symptom – Eye/Ear	172	10749	1.60	1.36	1.84
20	any_hypercap_icd	Symptom – Respiratory	997	1983	50.28	48.08	52.48

Exported: /Users/blocke/Box Sync/Residency Personal Files/Scholarly Work/Locke Research Project/CC-NLP/Symptom_Composition_by_Hypercapnia_Definition.xlsx
Exported: /Users/blocke/Box Sync/Residency Personal Files/Scholarly Work/Locke Research Project/CC-NLP/Symptom_Composition_Pivot_ChartReady.xlsx

```
composition_plot_df = definition_pivot_df.T.loc[HYPERCAP_CRITERIA]

ax = composition_plot_df.plot(
    kind="bar",
    stacked=True,
    figsize=(10, 6),
    width=0.8,
    edgecolor="black",
)
ax.set_title("Symptom Composition within Each Hypercapnia Definition (%)")
ax.set_xlabel("Hypercapnia Definition")
ax.set_ylabel("Percent of Patients")
ax.tick_params(axis="x", labelrotation=15)
ax.legend(title="Symptom Group", bbox_to_anchor=(1.05, 1), loc="upper left")
plt.tight_layout()
plt.show()
```



```

top_for_ci = (
    definition_counts_df.groupby("symptom_group")["Count"]
    .sum()
    .sort_values(ascending=False)
    .head(5)
    .index
)

ci_plot_df = definition_counts_df.loc[definition_counts_df["symptom_group"]
    ↪ ].isin(top_for_ci)].copy()
symptom_order = list(top_for_ci)
definition_order = HYPERCAP_CRITERIA

x = np.arange(len(symptom_order))
width = 0.18

fig, ax = plt.subplots(figsize=(11, 6))
for idx, definition in enumerate(definition_order):
    subset = (
        ci_plot_df.loc[ci_plot_df["Hypercapnia_Definition"].eq(definition)]
        .set_index("symptom_group")
        .reindex(symptom_order)
        .fillna(0)
    )

```



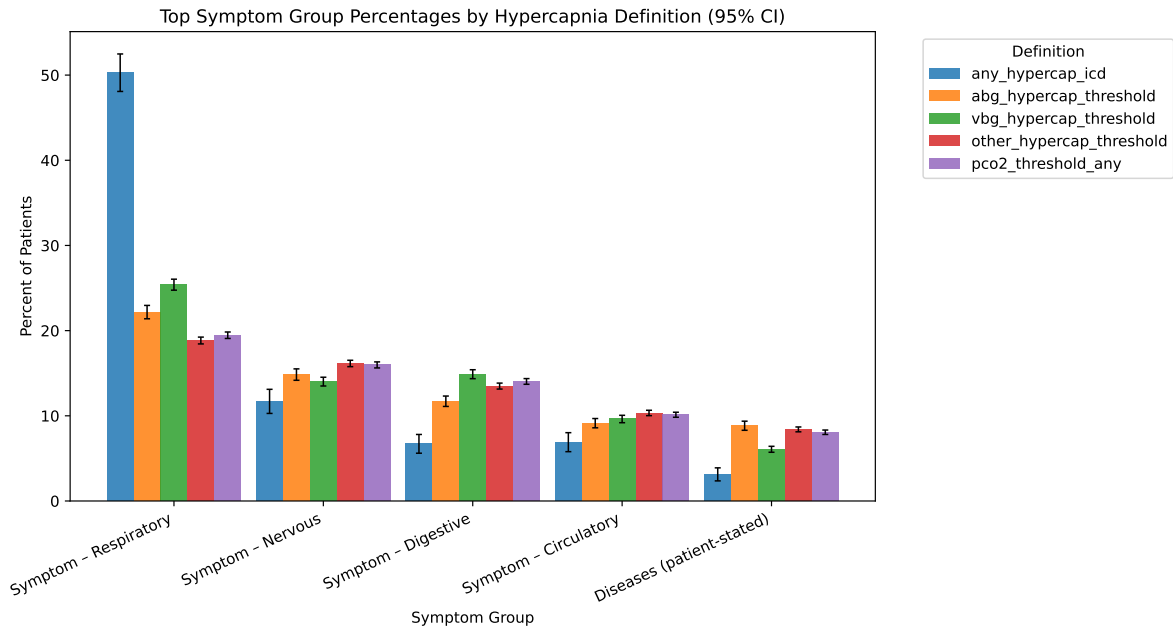
```

)
x_pos = x + (idx - (len(definition_order) - 1) / 2) * width
y = subset["Percent"].to_numpy()
lower = subset["CI_lower"].to_numpy()
upper = subset["CI_upper"].to_numpy()

ax.bar(x_pos, y, width=width, label=definition, alpha=0.85)
ax.errorbar(
    x_pos,
    y,
    yerr=[y - lower, upper - y],
    fmt="none",
    ecolor="black",
    elinewidth=1,
    capsize=2,
)

ax.set_xticks(x)
ax.set_xticklabels(symptom_order, rotation=25, ha="right")
ax.set_ylabel("Percent of Patients")
ax.set_xlabel("Symptom Group")
ax.set_title("Top Symptom Group Percentages by Hypercapnia Definition (95%
    ↪ CI)")
ax.legend(title="Definition", bbox_to_anchor=(1.05, 1), loc="upper left")
plt.tight_layout()
plt.show()

```



1.6 Symptom Distribution By Ascertainment Overlap

```

overlap_required = [
    SYMPTOM_COL,
    "abg_hypercap_threshold",
    "vbg_hypercap_threshold",
    "other_hypercap_threshold",
    "any_hypercap_icd",
    "pco2_threshold_any",
]
ensure_required_columns(df, overlap_required)

abg_flag = to_binary_flag(df["abg_hypercap_threshold"])
vbg_flag = to_binary_flag(df["vbg_hypercap_threshold"])
other_flag = to_binary_flag(df["other_hypercap_threshold"])
icd_flag = to_binary_flag(df["any_hypercap_icd"])
gas_flag = to_binary_flag(df["pco2_threshold_any"])

gas_source_labels = classify_gas_source_overlap(abg_flag, vbg_flag,
    ↪ other_flag)
abg_vbg_labels = np.select(
    [
        abg_flag.eq(1) & vbg_flag.eq(1),

```

```

        abg_flag.eq(1) & vbg_flag.eq(0),
        abg_flag.eq(0) & vbg_flag.eq(1),
    ],
    ["ABG+VBG", "ABG-only", "VBG-only"],
    default="Neither",
)

icd_gas_labels = np.select(
    [
        icd_flag.eq(1) & gas_flag.eq(1),
        icd_flag.eq(1) & gas_flag.eq(0),
        icd_flag.eq(0) & gas_flag.eq(1),
    ],
    ["ICD+Gas", "ICD-only", "Gas-only"],
    default="Neither",
)

overlap_df = df.copy()
overlap_df["gas_source_overlap"] = gas_source_labels
overlap_df["abg_vbg_overlap"] = abg_vbg_labels
overlap_df["icd_gas_overlap"] = icd_gas_labels

gas_positive_df = overlap_df.loc[abg_flag.eq(1) | vbg_flag.eq(1) |
    ⇨ other_flag.eq(1)].copy()
abg_vbg_positive_df = overlap_df.loc[abg_flag.eq(1) | vbg_flag.eq(1)].copy()
abg_vbg_counts_df, abg_vbg_pivot_df = symptom_distribution_by_overlap(
    abg_vbg_positive_df,
    group_col="abg_vbg_overlap",
    symptom_col=SYMPTOM_COL,
    top_k=10,
)
gas_source_counts_df, gas_source_pivot_df = symptom_distribution_by_overlap(
    gas_positive_df,
    group_col="gas_source_overlap",
    symptom_col=SYMPTOM_COL,
    top_k=10,
)
icd_gas_counts_df, icd_gas_pivot_df = symptom_distribution_by_overlap(
    overlap_df,
    group_col="icd_gas_overlap",
    symptom_col=SYMPTOM_COL,
    top_k=10,
)

```

```

gas_source_output_path = OUTPUT_DIR /
↳ "Symptom_Composition_by_ABG_VBG_Overlap.xlsx"
gas_source_expanded_output_path = OUTPUT_DIR /
↳ "Symptom_Composition_by_Gas_Source_Overlap.xlsx"
icd_gas_output_path = OUTPUT_DIR /
↳ "Symptom_Composition_by_ICD_Gas_Overlap.xlsx"
abg_vbg_pivot_df.to_excel(gas_source_output_path)
gas_source_pivot_df.to_excel(gas_source_expanded_output_path)
icd_gas_pivot_df.to_excel(icd_gas_output_path)

print("Symptom distribution by ABG/VBG overlap (legacy output):")
display(abg_vbg_pivot_df.head(15))
print("Symptom distribution by ABG/VBG/OTHER overlap (expanded output):")
display(gas_source_pivot_df.head(15))
print("Symptom distribution by ICD/Gas overlap:")
display(icd_gas_pivot_df.head(15))
print(f"Exported: {gas_source_output_path}")
print(f"Exported: {gas_source_expanded_output_path}")
print(f"Exported: {icd_gas_output_path}")

```

Symptom distribution by ABG/VBG overlap (legacy output):

abg_vbg_overlap symptom_group	ABG+VBG	ABG-only	VBG-only
Diseases (patient-stated)	6.4	10.8	5.9
Injuries & adverse effects	9.2	9.9	5.2
Other	7.0	10.2	6.2
Symptom – Circulatory	9.1	9.2	9.8
Symptom – Digestive	11.7	11.8	16.1
Symptom – Eye/Ear	1.9	1.4	3.7
Symptom – General	5.0	4.7	6.8
Symptom – Genitourinary	5.3	6.3	5.8
Symptom – Nervous	12.7	16.6	14.5
Symptom – Respiratory	29.5	16.1	23.8
Symptom – Skin/Hair/Nails	2.3	2.9	2.1

Symptom distribution by ABG/VBG/OTHER overlap (expanded output):

gas_source_overlap symptom_group	ABG+OTHER	ABG+VBG	ABG+VBG+OTHER	OTHER-only	VBG+
Diseases (patient-stated)	10.8	0.0	6.4	9.1	6.2
Injuries & adverse effects	9.9	0.0	9.2	7.0	6.1
Other	10.2	0.0	7.0	8.4	6.3
Symptom – Circulatory	9.2	100.0	9.0	10.9	10.6
Symptom – Digestive	11.8	0.0	11.7	13.9	14.8
Symptom – Eye/Ear	1.4	0.0	1.9	2.2	3.6
Symptom – General	4.7	0.0	5.0	6.7	6.9
Symptom – Genitourinary	6.3	0.0	5.3	6.9	5.4
Symptom – Nervous	16.6	0.0	12.7	17.7	14.5
Symptom – Respiratory	16.1	0.0	29.6	14.7	23.5
Symptom – Skin/Hair/Nails	2.9	0.0	2.3	2.4	2.1

Symptom distribution by ICD/Gas overlap:

icd_gas_overlap symptom_group	Gas-only	ICD+Gas	ICD-only
Diseases (patient-stated)	8.3	3.0	3.7
Injuries & adverse effects	7.2	4.5	2.9
Other	8.0	4.4	9.5
Symptom – Circulatory	10.3	6.3	11.2
Symptom – Digestive	14.4	6.7	7.0
Symptom – Eye/Ear	2.5	1.7	2.1
Symptom – General	6.4	3.7	7.4
Symptom – Genitourinary	6.4	4.1	4.1
Symptom – Nervous	16.1	12.3	7.0
Symptom – Respiratory	18.0	51.6	40.9
Symptom – Skin/Hair/Nails	2.4	1.6	4.1

Exported: /Users/blocke/Box Sync/Residency Personal Files/Scholarly Work/Locke Research Proj
CC-NLP/Symptom_Composition_by_ABG_VBG_Overlap.xlsx
Exported: /Users/blocke/Box Sync/Residency Personal Files/Scholarly Work/Locke Research Proj
CC-NLP/Symptom_Composition_by_Gas_Source_Overlap.xlsx
Exported: /Users/blocke/Box Sync/Residency Personal Files/Scholarly Work/Locke Research Proj
CC-NLP/Symptom_Composition_by_ICD_Gas_Overlap.xlsx

1.7 ICD Diagnostic Performance (ICD as predictor)

```
performance_targets = [
    ("pco2_threshold_any", "Gas any"),
    ("abg_hypercap_threshold", "ABG threshold"),
    ("vbg_hypercap_threshold", "VBG threshold"),
    ("other_hypercap_threshold", "PCO2 OTHER threshold"),
]

icd_positive = to_binary_flag(df["any_hypercap_icd"])
performance_rows = []
for target_col, target_label in performance_targets:
    target_positive = to_binary_flag(df[target_col])
    tp = int(((icd_positive == 1) & (target_positive == 1)).sum())
    fp = int(((icd_positive == 1) & (target_positive == 0)).sum())
    fn = int(((icd_positive == 0) & (target_positive == 1)).sum())
    tn = int(((icd_positive == 0) & (target_positive == 0)).sum())

    sens_denom = tp + fn
    ppv_denom = tp + fp
    sensitivity = float(tp / sens_denom) if sens_denom else np.nan
    ppv = float(tp / ppv_denom) if ppv_denom else np.nan
    sens_ci = (
        proportion_confint(tp, sens_denom, alpha=0.05, method="wilson")
        if sens_denom
        else (np.nan, np.nan)
    )
    ppv_ci = (
        proportion_confint(tp, ppv_denom, alpha=0.05, method="wilson")
        if ppv_denom
        else (np.nan, np.nan)
    )
    performance_rows.append(
        {
            "Target": target_label,
            "Target_Column": target_col,
            "TP": tp,
            "FP": fp,
            "FN": fn,
            "TN": tn,
            "Sensitivity": sensitivity,
            "Sensitivity_CI_Lower": sens_ci[0],
            "Sensitivity_CI_Upper": sens_ci[1],
        }
    )
```

```

        "PPV": ppv,
        "PPV_CI_Lower": ppv_ci[0],
        "PPV_CI_Upper": ppv_ci[1],
    }
)

icd_performance_df = pd.DataFrame(performance_rows)
icd_performance_df[[
    "Sensitivity",
    "Sensitivity_CI_Lower",
    "Sensitivity_CI_Upper",
    "PPV",
    "PPV_CI_Lower",
    "PPV_CI_Upper",
]] = icd_performance_df[[
    "Sensitivity",
    "Sensitivity_CI_Lower",
    "Sensitivity_CI_Upper",
    "PPV",
    "PPV_CI_Lower",
    "PPV_CI_Upper",
]].clip(lower=0.0, upper=1.0)

icd_subset_output_path = OUTPUT_DIR / "ICD_Positive_Subset_Breakdown.xlsx"
icd_performance_output_path = OUTPUT_DIR / "ICD_vs_Gas_Performance.xlsx"
with pd.ExcelWriter(icd_subset_output_path, engine="openpyxl") as writer:
    icd_positive_breakdown.to_excel(writer, index=False,
    ↪ sheet_name="Gas_criteria")
    icd_positive_category_summary.reset_index().to_excel(
        writer, index=False, sheet_name="ICD_categories"
    )
icd_performance_df.to_excel(icd_performance_output_path, index=False)

display(icd_positive_breakdown)
display(icd_performance_df)
print(f"Exported: {icd_subset_output_path}")
print(f"Exported: {icd_performance_output_path}")

```

	Definition	Count	Percent	Denominator_N
0	ABG threshold positive	1012	51.0	1983
1	VBG threshold positive	1482	74.7	1983
2	PCO2 OTHER threshold positive	1596	80.5	1983

	Definition	Count	Percent	Denominator_N
3	Any gas threshold positive	1741	87.8	1983

	Target	Target_Column	TP	FP	FN	TN	Sensitivity	Sensitivity
0	Gas any	pco2_threshold_any	1741	242	39339	0	0.042381	0.04047
1	ABG threshold	abg_hypercap_threshold	1012	971	9737	29602	0.094148	0.08877
2	VBG threshold	vbg_hypercap_threshold	1482	501	16066	23273	0.084454	0.08043
3	PCO2 OTHER threshold	other_hypercap_threshold	1596	387	34922	4417	0.043704	0.04165

Exported: /Users/blocke/Box Sync/Residency Personal Files/Scholarly Work/Locke Research Proj
CC-NLP/ICD_Positive_Subset_Breakdown.xlsx
Exported: /Users/blocke/Box Sync/Residency Personal Files/Scholarly Work/Locke Research Proj
CC-NLP/ICD_vs_Gas_Performance.xlsx

1.8 Ascertainment overlap UpSet

```
ascertainment_flags = pd.DataFrame(
    {
        "ICD": to_binary_flag(df["any_hypercap_icd"]).astype(bool),
        "ABG": to_binary_flag(df["abg_hypercap_threshold"]).astype(bool),
        "VBG": to_binary_flag(df["vbg_hypercap_threshold"]).astype(bool),
        "OTHER": to_binary_flag(df["other_hypercap_threshold"]).astype(bool),
    }
)

upset_series = from_indicators(ascertainment_flags.columns.tolist(),
    ↪ ascertainment_flags)
plt.figure(figsize=(12, 7))
upset_plot = UpSet(
    upset_series,
    subset_size="count",
    show_counts=True,
    sort_by="cardinality",
)
upset_plot.plot()
plt.suptitle("Ascertainment Overlap (ICD / ABG / VBG / OTHER)")
plt.tight_layout()

upset_output_path = OUTPUT_DIR / "Ascertainment_Overlap_UpSet.png"
```



```

plt.savefig(upset_output_path, dpi=300, bbox_inches="tight")
plt.show()

intersection_counts = (
    ascertainment_flags.groupby(["ICD", "ABG", "VBG", "OTHER"], dropna=False)
    .size()
    .reset_index(name="Count")
    .sort_values("Count", ascending=False)
    .reset_index(drop=True)
)

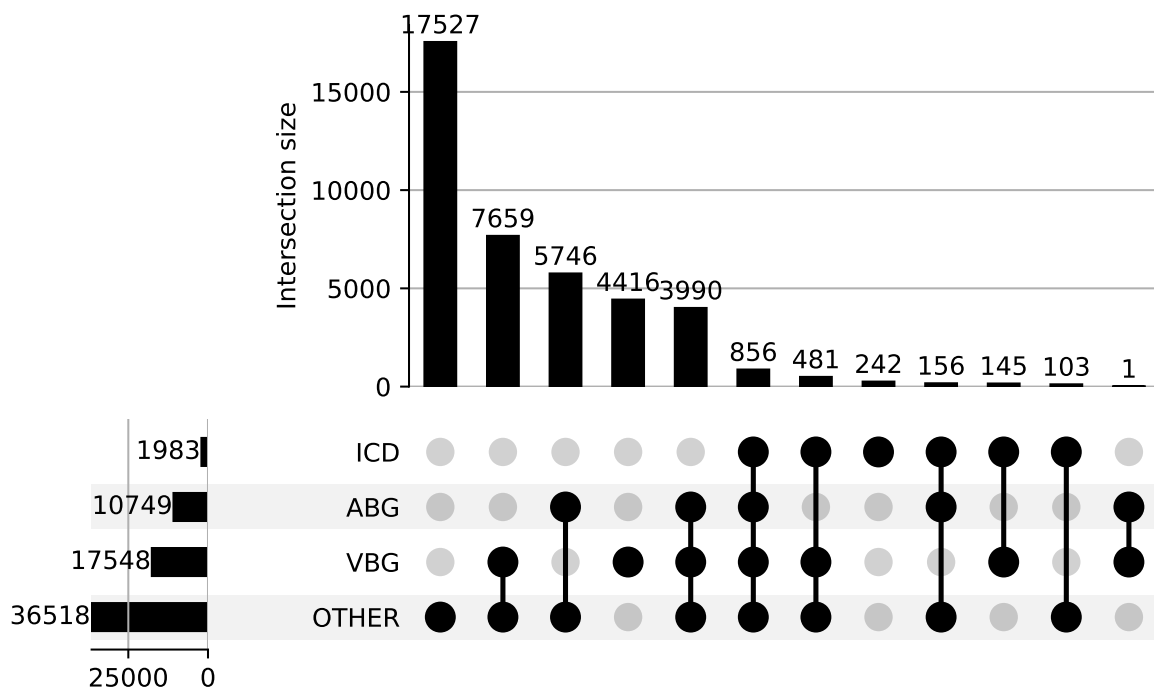
intersection_output_path = OUTPUT_DIR /
    ↪ "Ascertainment_Overlap_Intersections.xlsx"
intersection_counts.to_excel(intersection_output_path, index=False)

display(intersection_counts.head(20))
print(f"Exported: {upset_output_path}")
print(f"Exported: {intersection_output_path}")

```

<Figure size 3600x2100 with 0 Axes>

Ascertainment Overlap (ICD / ABG / VBG / OTHER)



	ICD	ABG	VBG	OTHER	Count
0	False	False	False	True	17527
1	False	False	True	True	7659
2	False	True	False	True	5746
3	False	False	True	False	4416
4	False	True	True	True	3990
5	True	True	True	True	856
6	True	False	True	True	481
7	True	False	False	False	242
8	True	True	False	True	156
9	True	False	True	False	145
10	True	False	False	True	103
11	False	True	True	False	1

Exported: /Users/blocke/Box Sync/Residency Personal Files/Scholarly Work/Locke Research Proj
CC-NLP/Ascertainment_Overlap_UpSet.png
Exported: /Users/blocke/Box Sync/Residency Personal Files/Scholarly Work/Locke Research Proj
CC-NLP/Ascertainment_Overlap_Intersections.xlsx

1.9 PDF-ready long tables

```
print(
    render_latex_longtable(
        hypercap_prevalence.reset_index(),
        caption=f"Hypercapnia prevalence summary (denominator = full cohort
↪ N={cohort_n:,}) .",
        label="tab:prevalence_summary",
        index=False,
    )
)
print(
    render_latex_longtable(
        icd_category_summary.reset_index(),
        caption=f"ICD category composition (denominator = full cohort
↪ N={len(df):,}) .",
        label="tab:icd_category",
        index=False,
    )
)
print(
    render_latex_longtable(
```

```

        inclusion_summary.reset_index(),
        caption=f"Inclusion source composition (denominator = full cohort
↪ N={len(df):,}).",
        label="tab:inclusion_type",
        index=False,
    )
)
print(
    render_latex_longtable(
        icd_positive_breakdown,
        caption=f"Among ICD-positive encounters, which gas criteria are also
↪ met (denominator = ICD-positive N={icd_positive_n:,}).",
        label="tab:icd_positive_breakdown",
        index=False,
    )
)
print(
    render_latex_longtable(
        icd_positive_category_summary.reset_index(),
        caption=f"Among ICD-positive encounters, ICD category distribution
↪ (denominator = ICD-positive N={icd_positive_n:,}).",
        label="tab:icd_positive_categories",
        index=False,
    )
)
print(
    render_latex_longtable(
        icd_performance_df,
        caption="ICD diagnostic performance vs gas-confirmed hypercapnia
↪ definitions (Wilson 95% CI).",
        label="tab:icd_performance",
        landscape=True,
        index=False,
    )
)

\begin{longtable}{llrrr}
\caption{Hypercapnia prevalence summary (denominator = full cohort N=41,322).} \label{tab:pr
\toprule
Definition & Column & Count & Denominator_N & Percent \\
\midrule
\endfirsthead
\caption[]{}{Hypercapnia prevalence summary (denominator = full cohort N=41,322).} \\

```

```

\toprule
Definition & Column & Count & Denominator_N & Percent \\
\midrule
\endhead
\midrule
\multicolumn{5}{r}{Continued on next page} \\
\midrule
\endfoot
\bottomrule
\endlastfoot
PCO2 threshold any source & pco2_threshold_any & 41080 & 41322 & 99.400000 \\
PCO2 OTHER threshold & other_hypercap_threshold & 36518 & 41322 & 88.400000 \\
VBG hypercapnia threshold & vbg_hypercap_threshold & 17548 & 41322 & 42.500000 \\
ABG hypercapnia threshold & abg_hypercap_threshold & 10749 & 41322 & 26.000000 \\
Hypercapnic RF ICD (any) & any_hypercap_icd & 1983 & 41322 & 4.800000 \\
\end{longtable}

\begin{longtable}{lrrrr}
\caption{ICD category composition (denominator = full cohort N=41,322).} \label{tab:icd_category}
\toprule
ICD Category & Count & Percent & Denominator_N \\
\midrule
\endfirsthead
\caption[]{}{ICD category composition (denominator = full cohort N=41,322).} \\
\toprule
ICD Category & Count & Percent & Denominator_N \\
\midrule
\endhead
\midrule
\multicolumn{4}{r}{Continued on next page} \\
\midrule
\endfoot
\bottomrule
\endlastfoot
Other / None & 39339 & 95.200000 & 41322 \\
Acute RF with hypoxia & 793 & 1.900000 & 41322 \\
Obesity hypoventilation syndrome & 524 & 1.300000 & 41322 \\
Acute RF with hypoxia & hypercapnia & 386 & 0.900000 & 41322 \\
Respiratory failure, unspecified & 187 & 0.500000 & 41322 \\
Acute RF with hypercapnia & 93 & 0.200000 & 41322 \\
\end{longtable}

\begin{longtable}{lrrrr}

```

```

\caption{Inclusion source composition (denominator = full cohort N=41,322).} \label{tab:incl}
\toprule
Inclusion Type & Count & Percent & Denominator_N \\
\midrule
\endfirsthead
\caption[]{}{Inclusion source composition (denominator = full cohort N=41,322).} \\
\toprule
Inclusion Type & Count & Percent & Denominator_N \\
\midrule
\endhead
\midrule
\multicolumn{4}{r}{Continued on next page} \\
\midrule
\endfoot
\bottomrule
\endlastfoot
Gas_only & 39339 & 95.200000 & 41322 \\
Both & 1741 & 4.200000 & 41322 \\
ICD_only & 242 & 0.600000 & 41322 \\
\end{longtable}

```

```

\begin{longtable}{lrrr}
\caption{Among ICD-positive encounters, which gas criteria are also met (denominator = ICD-
positive N=1,983).} \label{tab:icd_positive_breakdown} \\
\toprule
Definition & Count & Percent & Denominator_N \\
\midrule
\endfirsthead
\caption[]{}{Among ICD-positive encounters, which gas criteria are also met (denominator = ICD-
positive N=1,983).} \\
\toprule
Definition & Count & Percent & Denominator_N \\
\midrule
\endhead
\midrule
\multicolumn{4}{r}{Continued on next page} \\
\midrule
\endfoot
\bottomrule
\endlastfoot
ABG threshold positive & 1012 & 51.000000 & 1983 \\
VBG threshold positive & 1482 & 74.700000 & 1983 \\
PCO2 OTHER threshold positive & 1596 & 80.500000 & 1983 \\

```

```
Any gas threshold positive & 1741 & 87.800000 & 1983 \\
\end{longtable}
```

```
\begin{longtable}{lrrrr}
\caption{Among ICD-positive encounters, ICD category distribution (denominator = ICD-
positive N=1,983).} \label{tab:icd_positive_categories} \\
\toprule
ICD Category (ICD-positive subset) & Count & Percent & Denominator_N \\
\midrule
\endfirsthead
\caption[]{}{Among ICD-positive encounters, ICD category distribution (denominator = ICD-
positive N=1,983).} \\
\toprule
ICD Category (ICD-positive subset) & Count & Percent & Denominator_N \\
\midrule
\endhead
\midrule
\multicolumn{4}{r}{Continued on next page} \\
\midrule
\endfoot
\bottomrule
\endlastfoot
Acute RF with hypoxia & 793 & 40.000000 & 1983 \\
Obesity hypoventilation syndrome & 524 & 26.400000 & 1983 \\
Acute RF with hypoxia & hypercapnia & 386 & 19.500000 & 1983 \\
Respiratory failure, unspecified & 187 & 9.400000 & 1983 \\
Acute RF with hypercapnia & 93 & 4.700000 & 1983 \\
\end{longtable}
```

```
\begin{landscape}\n\begin{longtable}{llrrrrrrrrrrrr}
\caption{ICD diagnostic performance vs gas-confirmed hypercapnia definitions (Wilson 95% CI)}
\toprule
Target & Target_Column & TP & FP & FN & TN & Sensitivity & Sensitivity_CI_Lower & Sensitivity_CI_Upper & Specificity & Specificity_CI_Lower & Specificity_CI_Upper \\
\midrule
\endfirsthead
\caption[]{}{ICD diagnostic performance vs gas-confirmed hypercapnia definitions (Wilson 95% CI)}
\toprule
Target & Target_Column & TP & FP & FN & TN & Sensitivity & Sensitivity_CI_Lower & Sensitivity_CI_Upper & Specificity & Specificity_CI_Lower & Specificity_CI_Upper \\
\midrule
\endhead
\midrule
\multicolumn{12}{r}{Continued on next page} \\
\midrule
```

```

\endfoot
\bottomrule
\endlastfoot
Gas any & pco2_threshold_any & 1741 & 242 & 39339 & 0 & 0.042381 & 0.040475 & 0.044372 & 0.8
ABG threshold & abg_hypercap_threshold & 1012 & 971 & 9737 & 29602 & 0.094148 & 0.088772 & 0
VBG threshold & vbg_hypercap_threshold & 1482 & 501 & 16066 & 23273 & 0.084454 & 0.080430 & 0
PCO2 OTHER threshold & other_hypercap_threshold & 1596 & 387 & 34922 & 4417 & 0.043704 & 0.0
\end{longtable}
\n\end{landscape}\n

```

1.10 Association Model

Logistic regression of respiratory symptom flag on hypercapnia definitions.

```

model_df = df.dropna(subset=[SYMPTOM_COL]).copy()
model_df["is_respiratory"] = model_df[SYMPTOM_COL].astype(str).str.contains(
    r"\brespir", case=False, na=False
).astype(int)

design_matrix = sm.add_constant(model_df[HYPERCAP_CRITERIA],
    ↪ has_constant="add")
outcome = model_df["is_respiratory"]
logit_result = sm.Logit(outcome, design_matrix,
    ↪ missing="drop").fit(dis=False)

or_table = pd.DataFrame(
    {
        "OR": np.exp(logit_result.params),
        "CI_lo": np.exp(logit_result.conf_int()[0]),
        "CI_hi": np.exp(logit_result.conf_int()[1]),
        "p": logit_result.pvalues,
    }
).round(3)

display(or_table.loc[HYPERCAP_CRITERIA])

```

	OR	CI_lo	CI_hi	p
any_hypercap_icd	3.823	3.453	4.231	0.000
abg_hypercap_threshold	1.110	1.047	1.177	0.000
vbg_hypercap_threshold	1.666	1.576	1.761	0.000
other_hypercap_threshold	0.914	0.842	0.993	0.033
pco2_threshold_any	1.018	0.760	1.362	0.906

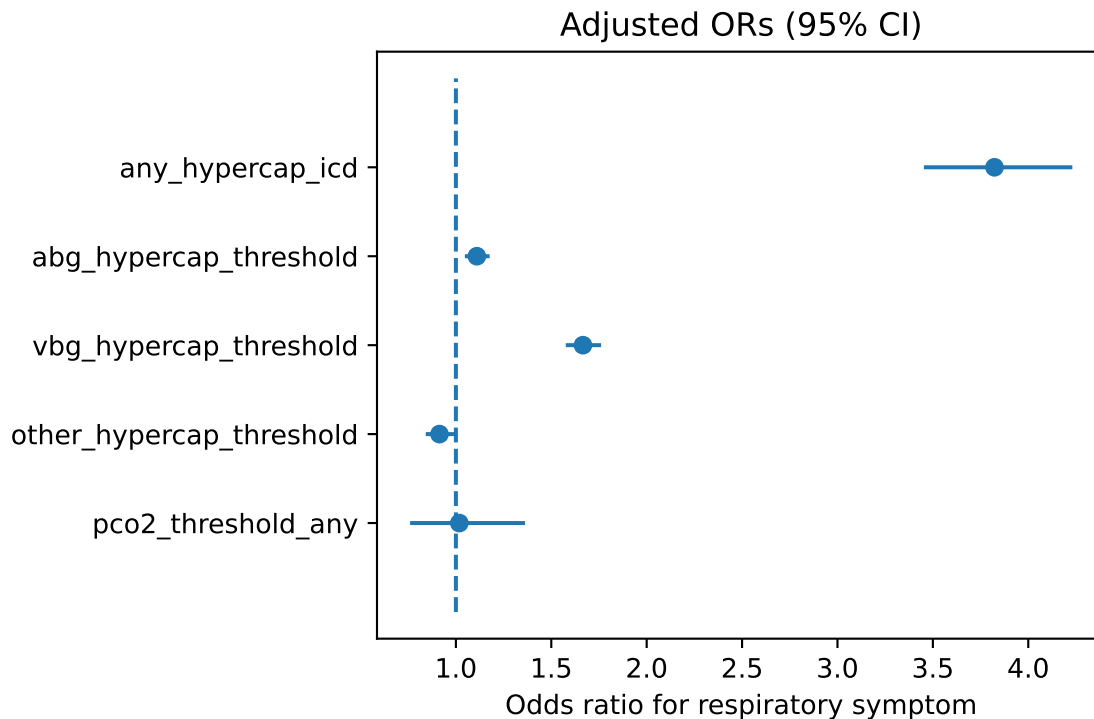
OR	CI_lo	CI_hi	p
----	-------	-------	---

```

or_plot_df = or_table.loc[HYPERCAP_CRITERIA]
y_positions = np.arange(len(or_plot_df))[:-1]

plt.figure(figsize=(6, 4))
plt.hlines(y=y_positions, xmin=or_plot_df["CI_lo"], xmax=or_plot_df["CI_hi"],
↪ linewidth=1.5)
plt.plot(or_plot_df["OR"], y_positions, "o")
plt.vlines(1, ymin=-1, ymax=len(or_plot_df), linestyle="dashed")
plt.yticks(y_positions, or_plot_df.index)
plt.xlabel("Odds ratio for respiratory symptom")
plt.title("Adjusted ORs (95% CI)")
plt.tight_layout()
plt.show()

```



1.11 Export Verification

```

expected_outputs = [
    definition_output_path,

```



```

    pivot_output_path,
    gas_source_output_path,
    gas_source_expanded_output_path,
    icd_gas_output_path,
    icd_subset_output_path,
    icd_performance_output_path,
    upset_output_path,
    intersection_output_path,
]

verification_rows = []
for output_path in expected_outputs:
    verification_rows.append(
        {
            "path": str(output_path),
            "exists": output_path.exists(),
            "size_bytes": output_path.stat().st_size if output_path.exists()
                ↪ else 0,
        }
    )

output_verification = pd.DataFrame(verification_rows)
display(output_verification)

```

	path	exists	size_bytes
0	/Users/blocke/Box Sync/Residency Personal File...	True	7367
1	/Users/blocke/Box Sync/Residency Personal File...	True	5560
2	/Users/blocke/Box Sync/Residency Personal File...	True	5373
3	/Users/blocke/Box Sync/Residency Personal File...	True	5531
4	/Users/blocke/Box Sync/Residency Personal File...	True	5366
5	/Users/blocke/Box Sync/Residency Personal File...	True	5855
6	/Users/blocke/Box Sync/Residency Personal File...	True	5581
7	/Users/blocke/Box Sync/Residency Personal File...	True	133494
8	/Users/blocke/Box Sync/Residency Personal File...	True	5225

```

from datetime import datetime

prior_runs_dir = WORK_DIR / "MIMIC tabular data" / "prior runs"
prior_runs_dir.mkdir(parents=True, exist_ok=True)
run_date = datetime.now().strftime("%Y-%m-%d")

```

```

analysis_manifest = collect_run_manifest(
    WORK_DIR,
    run_id=f"analysis_{datetime.now().strftime('%Y%m%d_%H%M%S')}",
)
analysis_manifest["stage"] = "analysis"
analysis_manifest["analysis_input_path"] = str(ANALYSIS_INPUT_PATH)
analysis_manifest["outputs"] = {
    "definition_output_path": str(definition_output_path),
    "pivot_output_path": str(pivot_output_path),
    "abg_vbg_overlap_output_path": str(gas_source_output_path),
    "gas_source_overlap_output_path": str(gas_source_expanded_output_path),
    "icd_gas_overlap_output_path": str(icd_gas_output_path),
    "icd_subset_output_path": str(icd_subset_output_path),
    "icd_performance_output_path": str(icd_performance_output_path),
    "upset_output_path": str(upset_output_path),
    "intersection_output_path": str(intersection_output_path),
}
analysis_manifest["output_verification"] = verification_rows
analysis_manifest_path = prior_runs_dir / f"{run_date}"
↪ analysis_run_manifest.json"
analysis_manifest_path.write_text(json.dumps(analysis_manifest, indent=2))
print(f"Wrote: {analysis_manifest_path}")

```

Wrote: /Users/blocke/Box Sync/Residency Personal Files/Scholarly Work/Locke Research Project
 CC-NLP/MIMIC tabular data/prior runs/2026-02-17 analysis_run_manifest.json